

## **A global envelope test to detect early and late bursts of trait evolution.**

**D. J. Murrell<sup>1,2,\*</sup>**

<sup>1</sup>Department of Genetics, Evolution and Environment, University College London, Gower Street,  
London WC1E 6BT, UK.

<sup>2</sup>Centre for Biodiversity and Environment Research, University College London, Gower Street,  
London WC1E 6BT, UK.

\*Corresponding author: [d.murrell@ucl.ac.uk](mailto:d.murrell@ucl.ac.uk)

Running title:

**Key-words:** Disparity through time; adaptive radiation; phenotypic diversification; null model;  
morphological diversification; macroevolution; comparative methods

## 1 **Abstract**

2 The joint analysis of species' evolutionary relatedness and their morphological evolution has  
3 offered much promise in understanding the processes that underpin the generation of biological  
4 diversity. Disparity through time (DTT) is a popular method that estimates the relative trait  
5 disparity within and between subclades at each time point, and compares this to the null hypothesis  
6 that trait values follow an uncorrelated random walk along the time calibrated phylogenetic tree. A  
7 simulation envelope is normally created by calculating, at every time point, the 95% minimum and  
8 95% maximum disparity values from multiple simulations of the null model on the phylogenetic  
9 tree. The null hypothesis is rejected whenever the empirical DTT curve falls outside of this  
10 envelope, and these time periods may then be linked to events that may have sparked non-random  
11 trait evolution. However, this method of envelope construction leads to multiple testing and a poor,  
12 uncontrolled, false positive rate. As a consequence it cannot be recommended. A recently  
13 developed method in spatial statistics is introduced that constructs a confidence envelope by giving  
14 each DTT curve a single ranking value based upon its most extreme disparity value. This method  
15 avoids the pitfalls of multiple testing whilst retaining a visual interpretation. Results using  
16 simulated data show this new test has desirable type 1 properties and is at least as powerful in  
17 correctly rejecting the null hypothesis as the morphological disparity index and node height test that  
18 lack a visual interpretation. Three example datasets are reanalyzed to show how the new test may  
19 lead to different inferences being drawn. Overall the results suggest the new rank envelope test  
20 should be used in null model testing for DTT analyses, and that there is no need to combine the  
21 envelope test with other tests such as has been done previously. Moreover, the rank envelope  
22 method can easily be adopted into recently developed posterior predictive simulation methods.  
23 More generally, the rank envelope test should be adopted when-ever a null model produces a vector  
24 of correlated values and the user wants to determine where the empirical data is different to the null  
25 model.

26

27

## 28 **Introduction**

29 Understanding the joint temporal dynamics of taxonomic and phenotypic diversity can provide  
30 tremendous insights into evolutionary success and its relationship with ecological opportunity,  
31 selective pressures, constraints, biotic interactions and environmental conditions. At the most basic  
32 level evolutionary biologists are often interested in detecting non-random evolution of biological  
33 traits within and across clades of species. Non-random bursts in evolution are often thought to be  
34 associated with events that open up ecological opportunities and enable a rapid increase in  
35 speciation rates and trait evolution, followed by slowdown in both processes as the ecological  
36 niches become filled. The evolutionary theory of adaptive radiation is the special case where the  
37 burst in speciation rate and trait evolution occur early in the clade's history (Schluter, 2000), but  
38 such bursts in trait evolution may occur at other times and can be triggered by other processes such  
39 as major events in the external environment.

40

41 A variety of methods exist to look for the signature of evolutionary bursts and fall into the  
42 categories of null model testing, and model selection (Harmon et al., 2003, Freckleton and Harvey,  
43 2006, Harmon et al., 2010, Slater et al., 2010, Slater and Pennell, 2014). The model selection  
44 approach takes a variety of candidate models (Brownian evolution, early burst, selective peak) and  
45 fits these to the data using maximum-likelihood methods before choosing the model that has the  
46 'best fit' (Harmon et al., 2010, Slater and Pennell, 2014). The null model approach remains more  
47 popular, partly because the methods have been established for longer, and the overall aim is to  
48 investigate if the data can be distinguished from the null model of uncorrelated evolution of trait  
49 values (Harmon et al., 2003, Freckleton and Harvey, 2006).

50

51 One of the more popular null model approaches is to look at morphological traits to see if trait  
52 disparity increases, decreases or stays the same as species accumulate in evolutionary time, and also  
53 see whether this disparity is greater within clades or between clades. Convergent evolution of traits  
54 is implied if morphological disparity is predominantly found within one or more subclades; whereas

55 adaptive radiations are expected to show divergence of traits between subclades, and in this  
56 scenario between clade morphological disparity should be greater than among subclade disparity.  
57 This analysis of between and within clade trait disparity has been championed by the disparity  
58 through time (DTT) approach introduced by Harmon et al. (2003). Here the empirical DTT curve is  
59 compared to the distribution of DTT curves generated on the same phylogenetic tree but under a  
60 specific model of how the trait diversity evolves. Generally the null model is an uncorrelated  
61 random walk, and this is generally referred to as Brownian evolution (ie a Brownian random walk  
62 over time in trait space). The method of comparison is critical in determining whether the empirical  
63 data can be distinguished from the null model. Early analyses used an integral deviation method  
64 called the Morphological Disparity Index (MDI) which sums the deviations of the empirical DTT  
65 curve from the median of the null model simulations (Harmon et al., 2003). The index can then be  
66 compared to the distribution of values produced by the simulation to test whether it is significantly  
67 different from the null model (e.g. Ingram (2015)). Where  $MDI > 0$ , this implies within-clade trait  
68 variation is generally greater than expected under the null model, and  $MDI < 0$  implies between-  
69 clade trait variation is more dominant than expected under the null model, and is suggestive of an  
70 adaptive radiation. The strength of the MDI is that it is a global test and so avoids multiple testing  
71 that can occur in analyses of time series data (see below), but the power of the MDI to detect non-  
72 Brownian bursts in trait evolution may be compromised if short-lived but extreme deviations in one  
73 direction at one point in the time series are balanced by a weak but long-lived deviation in the  
74 opposite direction.

75

76 Since the MDI produces a number, visualization of where any non-random bursts might have  
77 occurred (e.g. early on in the radiation) can only proceed by plotting the empirical DTT curve  
78 against the DTT curves sampled from the null model. However, determining where *statistically*  
79 *significant* local deviations from the null model are occurring in the time series requires another  
80 test. Slater et al. (2010) provided this by introducing an envelope method where the  $(100 - \alpha)\%$

81 upper and lower confidence intervals for the null model are estimated by sampling from the null  
82 model  $n$  times (where typically,  $n > 1000$ ) and then ordering the relative disparity values at each  
83 time point to arrive at an envelope that contains  $(100-2\alpha)\%$  of the simulated relative disparity  
84 values at each time point. This method is also referred to as the pointwise envelope method  
85 (Myllymaki et al., 2017). The observed relative disparity can then be compared to this envelope and  
86 if it falls outside the null model is said to be rejected at that level of significance.  
87  
88 The pointwise envelope method continues to be a popular method of inference. For example, Weber  
89 et al. (2016) used the DTT method to investigate the macroevolution of perfume signalling in  
90 orchid bees, a group known for their chemical sexual communications, and therefore a likely  
91 candidate for rapid diversification in traits. They found strong non-Brownian evolution for perfume  
92 signal, and weaker support for non-Brownian evolution of the labial gland. In both cases disparity  
93 was greater than the median DTT of the null model simulations indicating clades are overlapping in  
94 trait space, a signature of convergent evolution. The visual/graphical interpretation of the DTT with  
95 an envelope test has extra appeal as it can be used to identify time points where the burst of non-  
96 Brownian evolution occurred, enabling correlation with known evolutionary or environmental  
97 events that have triggered the burst. For example, Aristide et al. (2016) investigated brain shape,  
98 encephalization, and log body mass in New World monkeys. Their DTT envelope analyses  
99 supported the conclusion that a burst in evolution of brain shape occurred approximately 17-12 Ma,  
100 was associated with a burst in evolution of body mass that has previously been linked to  
101 diversification of diet and locomotion strategies, and was followed by a slowdown in disparity  
102 changes that persists to the current day. Conversely, Feilich (2016) found disparity in cichlid fin and  
103 body morphology was often greater than expected under the null model of Brownian evolution  
104 indicating most variation in morphology occurred within subclades. Moreover, the observed body  
105 and median fin disparity above the 95% confidence interval produced by the null model simulations

106 coincided with the Cichlinae–Pseudocrenilabrinae split, and a later split caused by the radiation of  
107 the haplochromine cichlids.

108

109 However, the pointwise envelope method leads to weaker than expected statistical performance  
110 because multiple tests, one at each time point, are being performed simultaneously. This is an issue  
111 that occurs in many different areas (eg spatial statistics Baddeley et al. (2014)), and generally  
112 whenever the pointwise envelope method is used in conjunction with a non-parametric method that  
113 produces a function as its summary output. Multiple testing leads to an increased type 1 statistical  
114 error rate (an elevated rate of rejection of the null hypothesis when it is true) that is no longer in line  
115 with the significance level being used to generate the confidence intervals of the envelope.

116 Although multiple testing problems may be solved using a Bonferroni correction, it is not  
117 appropriate here because the assumption of independence of tests is violated by the correlation of  
118 disparity values between consecutive time points and also the (often) large number of time points  
119 being simultaneously evaluated (Loosmore and Ford, 2006). Perhaps as a consequence of this many  
120 studies, including those discussed above, have used multiple methods to look for non-Brownian  
121 trait evolution including the MDI and the node height test (Freckleton and Harvey, 2006). However,  
122 the continued use of the pointwise envelope suggests its graphical interpretation is very appealing  
123 and it would therefore be worthwhile to circumvent its multiple testing issues.

124

125 Recently, a new method that avoids the multiple testing problems of the pointwise envelope but  
126 retains the visual interpretation has been developed in spatial statistics (Myllymaki et al., 2017).  
127 Spatial analysis of ecological data often leads to use of a non-parametric summary statistic such as  
128 Ripley's K that plots the tendency to cluster against the radial distance (eg. Law et al. (2009)), and  
129 the problem of pointwise envelopes for inference of non-random patterns is well established  
130 (Loosmore and Ford, 2006, Baddeley et al., 2014). Instead of producing confidence intervals by  
131 ranking each curve at each time point the rank envelope method ranks each curve for its overall

132 extremeness (more details are given below). As shown by Myllymaki et al. (2017), it has good type 1  
133 and type 2 error rates and is recommended for testing point pattern data against the null model of  
134 complete spatial randomness. The rank envelope can be developed and applied to any model that  
135 produces a vector (eg van Veen and Murrell, 2005), however, its performance needs to be tested  
136 since there are many ways of ordering curves based on their ‘extremeness’ and not all methods will  
137 produce desirable results.

138

139 In what follows, the rank envelope test will be developed for DTT null model analyses and its type  
140 1 and type 2 statistical properties compared to the pointwise envelope, MDI, and node height tests.  
141 The pointwise envelope test will be shown to have extremely poor type 1 error rates and should not  
142 be used for inference. In contrast, the rank envelope method will be shown to possess desirable type  
143 1 error rates, and be at least as powerful as the other tests in detecting accelerating or decelerating  
144 rates of trait evolution whilst retaining the useful property of graphical interpretation.

145

## 146 **Methods**

### 147 *Data simulation*

148 Phylogenetic trees were generated within *R* (version 3.3.3) using the *pbtree* function with the  
149 *phytools* R library (version 0.6, Revell (2012)), using the pure-birth (Yule) model. A pure-birth  
150 model is a simple way to ensure the number of tips is constant across simulations, but may be  
151 biologically plausible for clades over relatively short periods of time. These phylogenetic trees were  
152 then used to simulate quantitative trait evolution under a variety of scenarios including the null  
153 model of Brownian evolution. Specifically, trait evolution was simulated using the *fastBM* (in  
154 *phytools*) and *rescale* (in *geiger*, version 2.0.6, Pennell et al. (2014)) functions. The *rescale* function  
155 allows the simulation of early or late burst trait evolution using the early burst option, and a rate  
156 change parameter,  $a$ . When  $a < 0$  the rate of evolution decreases with time, mimicking an early  
157 burst of trait evolution, whereas  $a > 0$  models a late burst. The magnitude of  $a$  determines how

158 quickly this burst of activity fades away or builds up, with large magnitudes delivering a rapid  
159 decay or late increase in evolutionary change (examples are given in Figure S1 of the  
160 supplementary information). The null model of Brownian evolution is simulated under the same  
161 framework, but where  $a = 0$ .

162

### 163 *Disparity Through Time (DTT) Analyses*

164 DTT has proven to be one of the more popular approaches and uses the average pairwise Euclidean  
165 distance between species trait values as a measure of disparity. Following Harmon et al. (2003)  
166 relative disparity is calculated by dividing disparity of each subclade by the disparity of the whole  
167 tree. At each time point (speciation event) the average relative disparity for that time point is  
168 calculated as the mean of the relative disparities for all subclades whose ancestral lineages are  
169 present at that time. Disparity values close to zero indicate that variation in the trait(s) is  
170 predominantly partitioned between subclades rather than within them. Disparity values close to  
171 unity suggest that a clade contains a large amount of that variation, and that clades may overlap in  
172 trait space. By definition, disparity is 1 at the base of the phylogenetic tree, but is 0 at the present  
173 day.

174

175 Two methods that are currently used to search for the signal of bursts in morphological evolution  
176 using DTT are (1) the pointwise envelope test (Slater et al., 2010), and (2) an integral deviation test  
177 known as the Morphological Disparity Index (MDI, Harmon et al. (2003)). As well as these a third  
178 test, the global envelope test, was investigated. The global envelope test retains the same visual  
179 interpretation as the pointwise envelope test whilst avoiding the problem of multiple testing. All  
180 make comparisons of the empirical DTT to the DTT taken from the ensemble of simulations  
181 generated by the null model of Brownian evolution, and all use the same measure of disparity  
182 defined above.

183



## 184 *The Pointwise Envelope Test*

185 The pointwise envelope test is a Monte-Carlo simulation method that aims to produce a confidence  
186 interval, or envelope within which any part of the empirical DTT curve is said to be statistically  
187 indistinguishable from the null model. The method currently implemented in *Geiger* (v2.0.6 )  
188 constructs an envelope by extracting the  $\alpha$ th and  $(100-\alpha)$ th quantile of the DTT at each time point  
189 (speciation event) for all the null model simulations (normally  $\alpha = 2.5$ ). This produces a lower and  
190 upper interval within which  $(100-2\alpha)\%$  of all values for DTT at each speciation event under the null  
191 model. More formally, the envelope is defined by the lower and upper bounding curves

$$192 \quad T_{low}^{(k)}(t) = \min_{i=1,2,\dots,s}^k T_i(t) \quad (1a)$$

$$193 \quad T_{upp}^{(k)}(t) = \max_{i=1,2,\dots,s}^k T_i(t) \quad (1b)$$

194 where  $\min^k$  and  $\max^k$  denotes the  $k$ th smallest and largest values of the DTT across all simulations  
195  $s$ , of the null model at time (speciation event)  $t$ . If the empirical DTT curve falls outside of this  
196 envelope it is interpreted as being evidence for a departure from the null model of Brownian  
197 evolution.

198

## 199 *The MDI test*

200 Perhaps the simplest way to avoid multiple testing is to perform a deviation test that sums the  
201 deviations of the empirical DTT from the median DTT of the ensemble of null model simulations.  
202 Known as the Morphological Disparity Index (MDI), negative values indicate the empirical DTT  
203 curve is below the null model median DTT for at least some of the range of time points, again  
204 pointing to the possibility of an early burst in diversity (Harmon et al., 2003). A disadvantage of  
205 this approach is that it is not possible to say where the empirical DTT deviates from the null model  
206 without plotting it against the null model simulations and then performing some sort of envelope  
207 test. Moreover, since the index sums up the deviations from the median of the ensemble of  
208 simulations of the null model, it is theoretically possible for time periods where the empirical DTT  
209 is above the median DTT to be cancelled out by time periods where it is below the median DTT,

210 thus giving an MDI value close to that expected under Brownian evolution. However, most likely  
211 because it avoids the issue of multiple testing, the MDI has been well used (eg Harmon et al.  
212 (2003), Slater et al. (2010), Colombo et al. (2015), Ingram (2015), Jonsson et al. (2015)). In the  
213 results below Monte Carlo simulations were used to generate a distribution of MDI values from the  
214 null model. The empirical MDI value is then compared to this distribution and the null hypothesis  
215 (that the empirical data is from the same distribution as the null model of Brownian evolution) is  
216 rejected if it is smaller than the  $\alpha$ th quantile value or larger than the  $(100-\alpha)$ th quantile. As for the  
217 envelope tests,  $\alpha = 2.5$ .

218

### 219 *Rank Envelope Test*

220 A method that avoids multiple testing and therefore one which should show better type 1 error  
221 properties is the rank envelope test. In this case the whole DTT curve is ranked in terms of its  
222 extremeness (rather than rank each speciation event individually). This so-called ‘extreme rank’  
223 depth measure method was recently introduced by Myllymaki et al. (2017) for a similar problem  
224 that occurs in hypothesis testing for spatial point processes where non-parametric functions (pair  
225 correlation function and Ripley’s K) are used to detect departures from complete spatial  
226 randomness (ie a lack of any correlation between points). Such tests have been popular in ecology  
227 (e.g. Wiegand and Moloney (2004), Flüggé et al. (2012)). The more formal underpinnings of the  
228 test can be found in Myllymaki et al. (2017), but briefly the DTT curves from the  $r$  simulations of  
229 the null model are ordered according to the largest  $k$  for which they are still present in the  $k$ th  
230 envelope as defined by the pointwise envelope (equations 1a and 1b).

$$231 \quad R_i = \max\left\{k: T_{low}^{(k)}(t) \leq T_i(t) \leq T_{upp}^{(k)}(t) \text{ for all } t\right\} \quad (2)$$

232  $R_i$  is therefore the most extreme of the  $t$ -wise ranks of the curve  $T_i(t)$ . In other words, each curve is  
233 ranked relative to the other simulation curves, with its ranking being taken as its most extreme  
234 ranking within the pointwise envelope.

235

236 Since ties are possible (indeed are likely) the set of curves can only be weakly ordered. In order to  
237 generate a  $p$ -value a method for dealing with the ties needs to be used. Following Myllymaki et al.  
238 (2017), a range of  $p$ -values is reported that encompasses the most liberal and most conservative  $p$ -  
239 values respectively defined as

$$240 \quad p_- = \frac{1}{s+1} \sum_{i=1}^{s+1} \mathbf{1}(R_i < R_1) \quad (3a)$$

$$241 \quad p_+ = \frac{1}{s+1} \sum_{i=1}^{s+1} \mathbf{1}(R_i \leq R_1) \quad (3b)$$

242 and where  $R_1$  is the rank of the empirical DTT curve. This does raise the further problem that the  
243 interval defined by  $p_-$  and  $p_+$  could include the significance level  $\alpha$ , leading to an ambiguous result.  
244 However the likelihood of this happening is very small as long as  $s$ , the number of Monte Carlo  
245 simulations of the null model is sufficiently large. Myllymaki et al. (2017) recommend  $s \geq 2500$ ,  
246 and the results below use  $s = 5000$ . From the above DTT curve ordering, it is straightforward to  
247 visualise the global envelope determined by the significance level used, and in so-doing the user  
248 can readily see where the empirical data falls outside of the global envelope, and thus where the  
249 DTT is significantly different to Brownian evolution.

250

251 The extreme depth ranking method used here is but one of a number of possible ways to order the  
252 DTT curves. Some other functional depth orderings are discussed in Myllymaki et al. (2017) and  
253 the reader is directed to their paper for more details, but as will be shown below, the extreme rank  
254 depth measure described above leads to desirable type 1 and type 2 statistical errors for DTT  
255 analyses.

256

### 257 *The Node Height Test*

258 The final test does not use simulations of the null model to compare to the empirical data but  
259 instead relies upon the expectation that trait evolution should slow as niche space become packed.  
260 The node height test (Freckleton and Harvey, 2006) investigates if there is a significant correlation  
261 between the absolute magnitude of the standardized independent contrasts of the trait(s) and the

262 height above the root of the node at which they were being compared to. The height of a node is  
263 defined as the absolute distance between the root and the most recent common ancestor of the pair  
264 from which the contrast is generated. A significant relationship between these indicates that the rate  
265 of trait evolution is changing systematically through the tree with early and late bursts in trait  
266 evolution being diagnosed by the sign of the slope. Since this is an established test, the analyses of  
267 the node height test were performed using the function *nh.test* with the R library *geiger* (version  
268 2.0.6).

269

270 *Software code*

271 All results were obtained using R (version) and the code used to produce the results reported below  
272 can be found at <https://github.com/djmurrell/DTT-Envelope-code>. Details of the current methods to  
273 detect non-Brownian bursts of trait evolution can be found in the *geiger* software library, and the  
274 rank envelopes are produced by modifying code from the *spptest* library (Myllymaki et al., 2017).

275

## 276 **Results**

277 *False positive rates (Type 1 errors)*

278 The false positive rate is investigated by simulating an empirical dataset of trait evolution under the  
279 Brownian null model and testing how frequently each of the four tests described above incorrectly  
280 rejects the null hypothesis. Results for the DTT approach using the pointwise envelope test at the  
281 5% level of significance show a disappointing, but unsurprising high rate of false positives (Figure  
282 1). The multiple testing nature of this method means that the false positive rate is dependent upon  
283 the number of species in the comparison, and in the simulations the rate of false positives ranges  
284 approximately between 0.25 and 0.5 for 10-200 species (Figure 1). That is to say, for comparisons  
285 using more than 100 species the pointwise envelope test is incorrectly rejecting the null hypothesis  
286 of Brownian evolution in approximately 50% of cases. As such it is impossible to recommend this  
287 method for inference of non-Brownian bursts of trait evolution. In comparison, the node height test

288 and the global rank envelope test both return consistent false positive rates that hover around the  
289 significance level used (Figure 1). The MDI shows a high type 1 error rate for phylogenetic trees  
290 with less than 40 species, but thereafter a desirable false positive rate is returned.

291

### 292 *True positive rates (Type 2 errors)*

293 Evolutionary biologists are often most interested in early bursts in trait evolution, since this is  
294 argued to be the hallmark of adaptive radiations (Harmon et al., 2010). However, a number of  
295 studies have also reported late bursts in trait evolution where departures from Brownian evolution  
296 occur only in the recent past (e.g. Koecke et al. (2013), Tran (2014), Pincheira-Donoso et al. (2015),  
297 Feilich (2016)). Simulations for both scenarios confirm that the MDI, the node height, and the rank  
298 envelope test can all successfully detect both early and late bursts in trait evolution (Figure 2).  
299 Convention dictates that at the 5% significance level a desirable test shows a true positive rate of  
300 0.8. The ability of all tests to reach this mark is dependent on the number of species and the strength  
301 of the early or late burst, but other generalities do emerge. Firstly, for early burst trait evolution the  
302 MDI deviation test is less powerful than the node height and global envelope tests, with the global  
303 envelope test generally showing slightly higher power. Secondly, early bursts are slightly easier to  
304 detect than late bursts of trait evolution i.e. the desired true positive rate is reached with fewer  
305 species in early burst models compared to late burst models. The node height and global envelope  
306 tests show similar power to detect late bursts, but the MDI is better able to detect late bursts in trait  
307 evolution with fewer species at the tips. Example simulations for each early/late burst rate with 100  
308 species at the tips of the phylogenetic tree can be found in the supplementary material (Figure S1).

309

### 310 *Data Examples*

311 Having established the rank envelope test possesses desirable type 1 and type 2 statistical error  
312 properties, three datasets were used to illustrate how inference of the rates of morphological  
313 evolution can change depending on whether the pointwise or global envelope test is used. Since the

314 pointwise envelope test is too liberal in its rejection of the null model the expectation should be for  
315 a reduction in support for non-Brownian bursts in morphological evolution.

316

317 The first example uses the morphological and phylogenetic data on Darwin's finches (*Geospiza*)  
318 which is currently found in the *geiger* (version 2.0.6) *R* package. Re-analysis shows support for two  
319 late bursts in culmen length evolution by the pointwise envelope test, both showing diversification  
320 predominantly occurring within clade(s) followed by decreases in disparity caused by increased  
321 diversity between clades (Figure 3a). In contrast, there is no departure from the null model of  
322 Brownian evolution according to the rank envelope test (Figure 3b).

323

324 The second example uses a time-calibrated molecular phylogeny of extant cetaceans and a  
325 morphological dataset on body size from Slater et al. (2010) which is also located within *geiger*  
326 (version 2.0.6). Their analyses used a combination of the node height test (NHT), and DTT tests  
327 using the MDI and the DTT pointwise envelope method. On the basis of the pointwise envelope  
328 results they came to the conclusion that cetaceans do show a burst in evolution of body size and that  
329 this occurred predominantly during the period 6-11 Ma. Reanalysis confirms that the pointwise  
330 envelope approach finds the same burst in trait evolution but that the global envelope test fails to  
331 find any departure from the null model of Brownian evolution at the 5% level of statistical  
332 significance (Figure 3c, d). These results are covered in more detail in the discussion.

333

334 The final example is taken from (Feilich, 2016) who investigated the evolution of  
335 body shape, caudal fin shape, dorsal fin shape, and anal fin shape in 131 African cichlid fishes. Re-  
336 analysing the data for anal fin shape using the pointwise envelope (Figure 3e) confirms the spike in  
337 relative disparity coinciding with the Cichlinae–Pseudocrenilabrinae split 45-75 MYA reported in  
338 the original paper as well as the spike nearer to the present day that coincides with the  
339 haplochromine radiation (Feilich, 2016). In contrast, the rank envelope method finds no discernable

340 difference (at the 5% level of significance) from the null model of Brownian evolution at any point  
341 in the evolutionary timeline (Figure 3f). Re-analysis of body shape, dorsal fin shape, and caudal fin  
342 shape evolution using the rank envelope method does retain support for the peaks in disparity  
343 associated with the haplochromine radiation found by (Feilich, 2016) using the pointwise envelope  
344 method (Figure S2).

345

## 346 **Discussion**

347 The object of this study has been to highlight the unacceptably high false positive rates of the  
348 pointwise envelope test, and offer an alternative solution, the global rank envelope test (Myllymaki  
349 et al., 2017) that does not have the same multiple testing issues, but which still allows identification  
350 of the time period where the non-Brownian trait evolution may have occurred. Envelope tests using  
351 the DTT pointwise envelope method continue to be useful and popular (e.g. Johnson and Omland,  
352 2004, Slater et al., 2010, Dornburg et al., 2011, Blackburn et al., 2013, Ingram, 2015, Arbour and  
353 Lopez-Fernandez, 2016, Aristide et al., 2016, Feilich, 2016, Hlusko et al., 2016, Weber et al., 2016)  
354 as the time periods over which trait evolution has been non-Brownian can often be linked to  
355 specific events that may have triggered the burst of trait evolution (e.g. Slater et al., 2010, Feilich,  
356 2016, Hlusko et al., 2016). Unfortunately, the results presented here (Figures 1- 3) suggest re-  
357 interpretation of some previous analyses may be required using a global envelope test such as that  
358 introduced here rather than a pointwise envelope test.

359

360 A number of methods beyond the pointwise envelope test have been used and there is clear  
361 variability in their ability to detect departures from the null model (Figure 2). The MDI (Harmon et  
362 al., 2003) is a global test as it sums up the difference between the empirical DTT and the average of  
363 the simulations of the null model. However, the MDI can be relatively insensitive to early bursts in  
364 diversity that lead to larger between-clade disparity (Figure 2); can return a high false positive rate  
365 for datasets with less than 40 species (Figure 1); and it is possible for a large difference from the

366 null model in one time period to be cancelled out by smaller deviations in the opposite direction at  
367 other time periods. However, the MDI might be better at detecting late bursts in diversity (Figure  
368 2). This is possibly because late bursts are often associated with long time periods of disparity  
369 concentrated within clades (eg Figure 2a, e) and there is less chance the deviation from the median  
370 of the simulations is cancelled out by negative deviations elsewhere. On the other hand, there is  
371 little to choose between the rank envelope test and node height test based upon the results presented  
372 here (Figure 1, 2). The main advantage of the rank envelope method is that it provides a  
373 visualization of how disparity changes over time and it is easier to see where the burst of trait  
374 evolution may have occurred. Given this, and the slightly better performance when species numbers  
375 are small and/or early/late bursts are weak, the rank envelope method might be preferable.

376

377 Tests using three datasets that are freely available show how inferences and conclusions can change  
378 in quite important ways when the rank envelope test is used instead of the pointwise envelope  
379 (Figure 3). Macroevolutionary investigations have often used multiple methods (eg node height test,  
380 MDI test and DTT plots (Slater et al., 2010); MDI test and DTT plots (Arbour and Lopez-  
381 Fernandez, 2016)), to test for departures from Brownian evolution, but given the results presented  
382 here there is a risk that a mixture of results will be produced. For example, before removing species  
383 considered to be outliers, Slater et al. (2010) found evidence supporting an early burst in cetacean  
384 body size using the pointwise envelope test (replicated here in Figure 3), but neither the MDI test  
385 nor the node height test could find a statistically significant deviation from the null model of  
386 Brownian evolution. These discrepancies are expected given the high false positive rates of the  
387 pointwise envelope test (Figure 1), and re-analysis using the global rank envelope test shows  
388 agreement with the results of the initial node height test, and MDI test before removal of outliers  
389 (Slater et al. (2010); Figure 3).

390



391 The methods investigated here all use the same hypothesis testing approach. That is to say we test  
392 our data against a suitable null model to see if there are detectable departures from the null model.  
393 A different approach is to consider a number of candidate models and ask which model best  
394 describes the data (Johnson and Omland, 2004). The advantage of this model selection approach is  
395 that multiple models are considered simultaneously, but of course there is no guarantee that the best  
396 model, usually determined by some information theoretic criterion, is a ‘good’ descriptor of the  
397 data. The model selection approach has been developed for trait evolution by Harmon et al. (2010)  
398 who used maximum likelihood methods to fit models that could produce Brownian evolution,  
399 increasing or decreasing trait diversification rates, as well as selective peaks where the trait value  
400 has a tendency to return to a medial value. Using the likelihood ratio test, they found the Brownian  
401 evolution and the selective peak (Ornstein-Uhlenbek) models to be the most frequently selected  
402 across 49 clades, implying early bursts in trait evolution are relatively rare. Slater and Pennell  
403 (2014) extended this method by employing a posterior predictive approach instead of the likelihood  
404 ratio test. The posterior predictive approach proceeds by fitting the parameters to the candidate  
405 models using maximum likelihood as in (Harmon et al., 2010), but model selection is based upon  
406 sampling the trait evolution from the fitted models and then comparing the fit of each model to the  
407 observed trait values. Slater and Pennell (2014) developed this method using either the MDI test, or  
408 the node height test and showed both of these posterior predictive methods can have a higher power  
409 to detect early bursts in trait evolution compared to the maximum likelihood ratio approach of  
410 Harmon et al. (2010). Re-analysing the cetacean dataset with these methods led to the conclusion  
411 that an early burst model best described the evolution of whale body size (Slater and Pennell, 2014).  
412 This is not surprising given the rank envelope test clearly shows the empirical DTT curve is close to  
413 falling below the lower confidence interval (Figure 2d). Ultimately, the user needs to choose  
414 between the null model testing and model selection methods, but the rank envelope test developed  
415 here could easily be incorporated into the posterior predictive methods of Slater and Pennell (2014),  
416 since the ranking of the observed DTT curve in the ensemble of simulations from each of the

417 candidate models generates a single metric, the global rank amongst the set of model curves, that  
418 could then be used to compare the models.

419

420 In summary, the pointwise envelope test method that has been employed to investigate bursts in  
421 trait evolution shows unacceptably high type 1 statistical errors and should not be used. The rank  
422 envelope test that was introduced by Myllymaki et al. (2017) for spatial point pattern analysis, ranks  
423 the extremeness of the empirical DTT curve against the suite of realisations of an appropriate null  
424 model, and is instead shown to possess the desirable statistical properties. This method is shown to  
425 be at least as powerful as other popular methods for detecting bursts of trait evolution, whilst  
426 retaining the advantages of graphical interpretation of where in time the significant bursts of  
427 evolution occur. Null model testing is a commonly used tool in ecology and evolution (eg Gotelli  
428 and Ulrich (2012)) and the rank envelope test is a flexible method that allows inference without the  
429 problem of multiple testing when the null model returns a functional relationship (ie a vector) rather  
430 than a singular value (ie a scalar quantity).

431

432

### 433 **Acknowledgements**

434 This work was supported by the Engineering and Physical Sciences Research Council grant  
435 EP/N007336/1. The author would like to thank Alex Pigot who provided useful suggestion and  
436 comments on earlier versions of the manuscript.

437

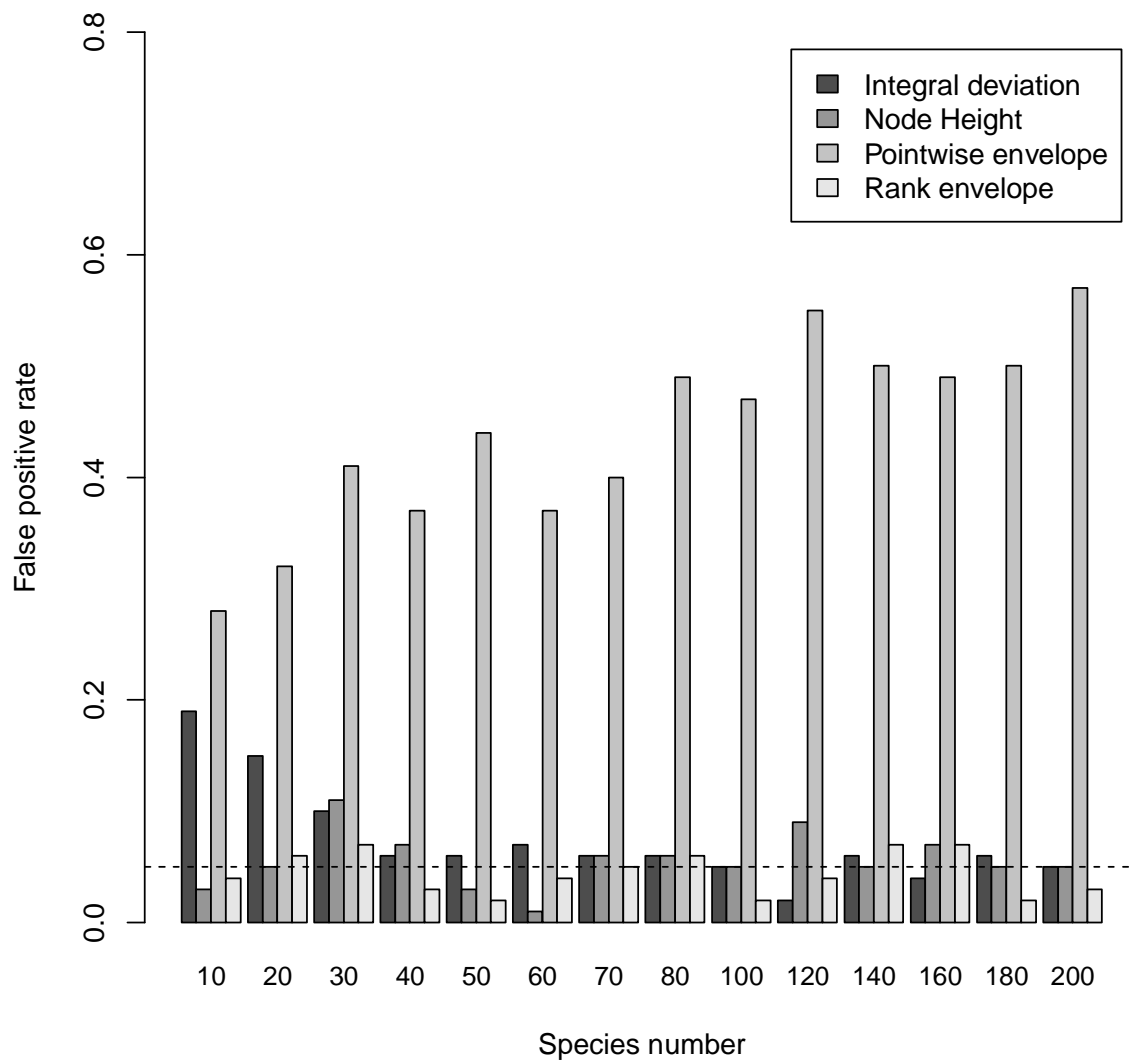
438

439 **References**

- 440 Arbour, J. H. & Lopez-Fernandez, H. (2016) Continental cichlid radiations: functional diversity  
441 reveals the role of changing ecological opportunity in the Neotropics. *Proceedings of the*  
442 *Royal Society B-Biological Sciences*, **283**.
- 443 Aristide, L., dos Reis, S. F., Machado, A. C., Lima, I., Lopes, R. T. & Perez, S. I. (2016) Brain  
444 shape convergence in the adaptive radiation of New World monkeys. *Proceedings of the*  
445 *National Academy of Sciences of the United States of America*, **113**, 2158-2163.
- 446 Baddeley, A., Diggle, P. J., Hardegen, A., Lawrence, T., Milne, R. K. & Nair, G. (2014) On tests of  
447 spatial pattern based on simulation envelopes. *Ecological Monographs*, **84**, 477-489.
- 448 Blackburn, D. C., Siler, C. D., Diesmos, A. C., McGuire, J. A., Cannatella, D. C. & Brown, R. M.  
449 (2013) An adaptive radiation of frogs in a southeast Asian island archipelago. *Evolution*, **67**,  
450 2631-46.
- 451 Colombo, M., Damerau, M., Hanel, R., Salzburger, W. & Matschiner, M. (2015) Diversity and  
452 disparity through time in the adaptive radiation of Antarctic notothenioid fishes. *Journal of*  
453 *Evolutionary Biology*, **28**, 376-394.
- 454 Dornburg, A., Sidlauskas, B., Santini, F., Sorenson, L., Near, T. J. & Alfaro, M. E. (2011) The  
455 influence of an innovative locomotor strategy on the phenotypic diversification of  
456 triggerfish (family: Balistidae). *Evolution*, **65**, 1912-26.
- 457 Feilich, K. L. (2016) Correlated evolution of body and fin morphology in the cichlid fishes.  
458 *Evolution*, **70**, 2247-2267.
- 459 Flügge, A. J., Olhede, S. C. & Murrell, D. J. (2012) The memory of spatial patterns: changes in  
460 local abundance and aggregation in a tropical forest. *Ecology*, **93**, 1540-1549.
- 461 Freckleton, R. P. & Harvey, P. H. (2006) Detecting non-Brownian trait evolution in adaptive  
462 radiations. *PLoS Biol*, **4**, e373.
- 463 Gotelli, N. J. & Ulrich, W. (2012) Statistical challenges in null model analysis. *Oikos*, **121**, 171-  
464 180.

- 465 Harmon, L. J., Losos, J. B., Davies, T. J., Gillespie, R. G., Gittleman, J. L., Jennings, W. B., Kozak,  
466 K. H., McPeck, M. A., Moreno-Roark, F., Near, T. J., Purvis, A., Ricklefs, R. E., Schluter,  
467 D., Schulte, J. A., Seehausen, O., Sidlauskas, B. L., Torres-Carvajal, O., Weir, J. T. &  
468 Mooers, A. O. (2010) Early bursts of body size and shape evolution are rare in comparative  
469 data. *Evolution*, **64**, 2385-2396.
- 470 Harmon, L. J., Schulte, J. A., 2nd, Larson, A. & Losos, J. B. (2003) Tempo and mode of  
471 evolutionary radiation in iguanian lizards. *Science*, **301**, 961-4.
- 472 Hlusko, L. J., Schmitt, C. A., Monson, T. A., Brasil, M. F. & Mahaney, M. C. (2016) The  
473 integration of quantitative genetics, paleontology, and neontology reveals genetic  
474 underpinnings of primate dental evolution. *Proceedings of the National Academy of*  
475 *Sciences of the United States of America*, **113**, 9262-9267.
- 476 Ingram, T. (2015) Diversification of body shape in *Sebastes* rockfishes of the north-east Pacific.  
477 *Biological Journal of the Linnean Society*, **116**, 805-818.
- 478 Johnson, J. B. & Omland, K. S. (2004) Model selection in ecology and evolution. *Trends in*  
479 *Ecology & Evolution*, **19**, 101-108.
- 480 Jonsson, K. A., Lessard, J. P. & Ricklefs, R. E. (2015) The evolution of morphological diversity in  
481 continental assemblages of passerine birds. *Evolution*, **69**, 879-889.
- 482 Koecke, A. V., Muellner-Riehl, A. N., Pennington, T. D., Schorr, G. & Schnitzler, J. (2013) Niche  
483 evolution through time and across continents: The story of neotropical *Cedrela* (Meliaceae)  
484 *American Journal of Botany*, **100**, 1800-1810.
- 485 Law, R., Illian, J., Burslem, D. F., Gratzler, G., Gunatilleke, C. & Gunatilleke, I. (2009) Ecological  
486 information from spatial patterns of plants: insights from point process theory. *Journal of*  
487 *Ecology*, **97**, 616-628.
- 488 Loosmore, N. B. & Ford, E. D. (2006) Statistical inference using the G or K point pattern spatial  
489 statistics. *Ecology*, **87**, 1925-1931.

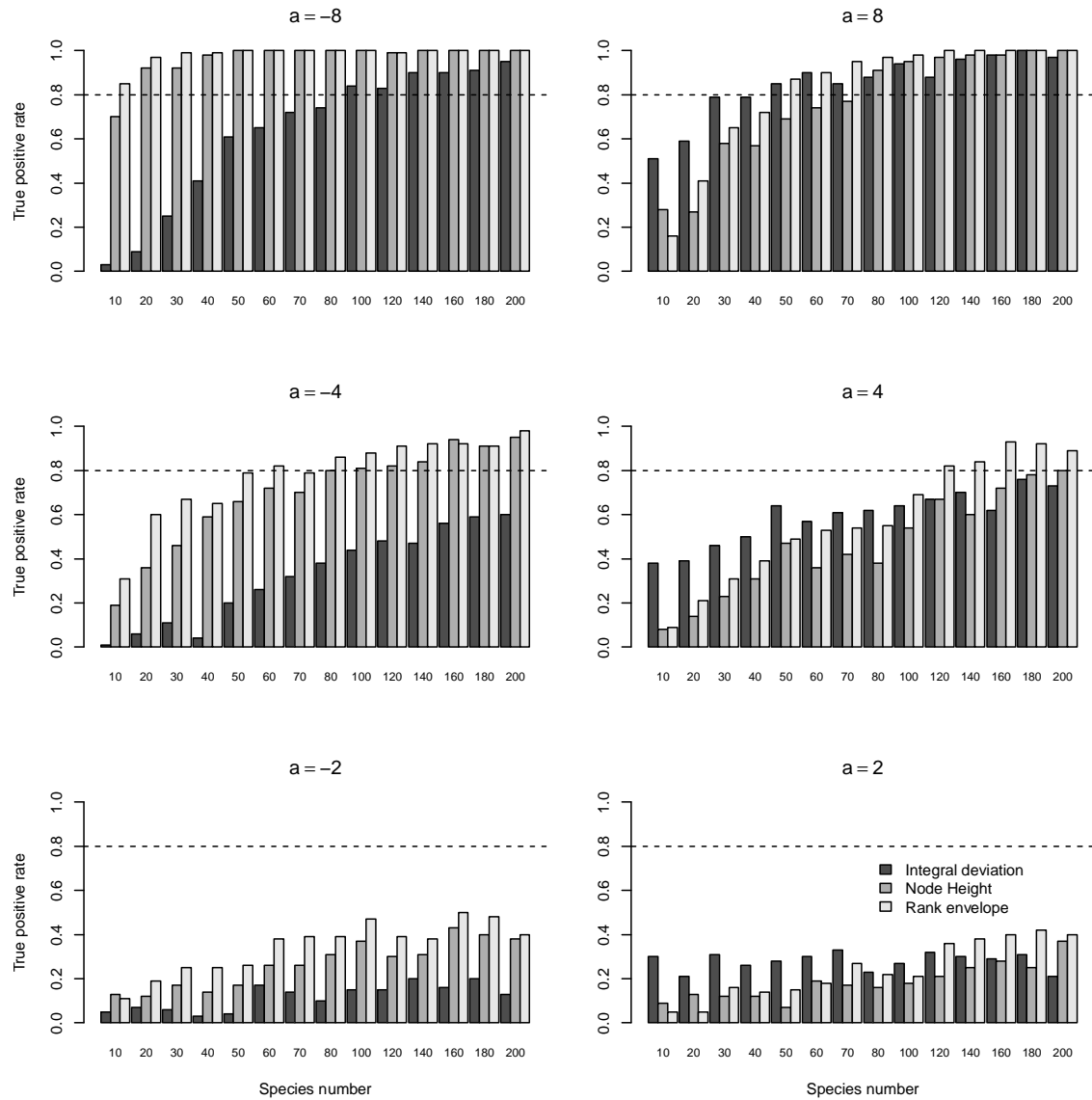
- 490 Myllymaki, M., Mrkvicka, T., Grabarnik, P., Seijo, H. & Hahn, U. (2017) Global envelope tests for  
491 spatial processes. *Journal of the Royal Statistical Society Series B-Statistical Methodology*,  
492 **79**, 381-404.
- 493 Pennell, M. W., Eastman, J. M., Slater, G. J., Brown, J. W., Uyeda, J. C., FitzJohn, R. G., Alfaro,  
494 M. E. & Harmon, L. J. (2014) geiger v2.0: an expanded suite of methods for fitting  
495 macroevolutionary models to phylogenetic trees. *Bioinformatics*, **30**, 2216-2218.
- 496 Pincheira-Donoso, D., Harvey, L. P. & Ruta, M. (2015) What defines an adaptive radiation?  
497 Macroevolutionary diversification dynamics of an exceptionally species-rich continental  
498 lizard radiation. *Bmc Evolutionary Biology*, **15**.
- 499 Revell, L. J. (2012) phytools: an R package for phylogenetic comparative biology (and other  
500 things). *Methods in Ecology and Evolution*, **3**, 217-223.
- 501 Schluter, D. (2000) *The ecology of adaptive radiation*. OUP Oxford.
- 502 Slater, G. J. & Pennell, M. W. (2014) Robust regression and posterior predictive simulation  
503 increase power to detect early bursts of trait evolution. *Syst Biol*, **63**, 293-308.
- 504 Slater, G. J., Price, S. A., Santini, F. & Alfaro, M. E. (2010) Diversity versus disparity and the  
505 radiation of modern cetaceans. *Proc Biol Sci*, **277**, 3097-104.
- 506 Tran, L. A. P. (2014) The role of ecological opportunity in shaping disparate diversification  
507 trajectories in a bicontinental primate radiation. *Proceedings of the Royal Society B-*  
508 *Biological Sciences*, **281**.
- 509 van Veen, F. & Murrell, D. (2005) A simple explanation for universal scaling relations in food  
510 webs. *Ecology*, **86**, 3258-3263.
- 511 Weber, M. G., Mitko, L., Eltz, T. & Ramirez, S. R. (2016) Macroevolution of perfume signalling in  
512 orchid bees. *Ecology Letters*, **19**, 1314-1323.
- 513 Wiegand, T. & Moloney, K. (2004) Rings, circles, and null models for point pattern analysis in  
514 ecology. *Oikos*, **104**, 209-229.
- 515



516

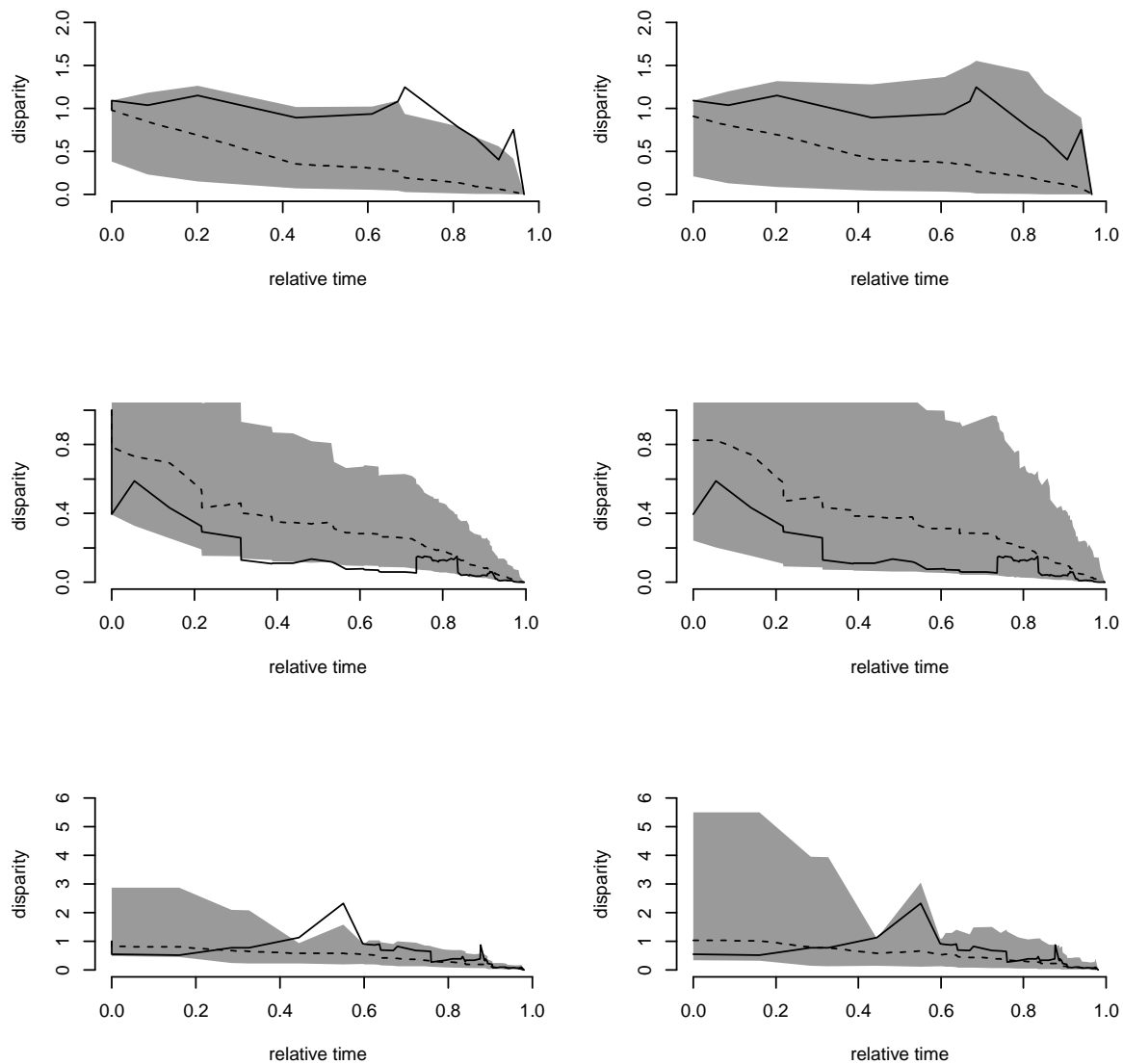
517 **Figure 1.** False positive rates for the four tests for non-random disparity through time (DTT) as a  
518 function of the number of species at the tips of the phylogenetic tree. Rates are estimated from 100  
519 simulated phylogenetic trees for each number of species using a pure birth model to generate the  
520 phylogenetic tree, and assuming Brownian evolution of the trait at each speciation event. All tests  
521 requiring Monte Carlo simulations were run with  $s = 5000$  trait evolution simulations.

522



524 **Figure 2.** True positive rates (statistical power) of the rank envelope test, the MDI test and the node  
525 height test under a range of simulated early and late burst scenarios, and for a range of size of  
526 phylogenetic tree. Rates are estimated from 100 simulated phylogenetic trees for each number of  
527 species using a pure birth model to generate the phylogenetic tree, and assuming trait evolution at  
528 each speciation event speeds ups or slow downs over evolutionary time. All tests requiring Monte  
529 Carlo simulations were run with  $s = 5000$  trait evolution simulations for every phylogenetic tree.  
530 When  $a < 0$ , a burst in trait evolution occurs early in evolutionary time, and late bursts occur when  
531  $a > 0$ . Large magnitudes lead to the bursts occurring over a smaller period of time. Corresponding  
532 example plots of DTT in each scenario are given in the supplementary information.





533

534 **Figure 3.** Comparisons of inference from using the pointwise envelope test (left hand column) and

535 the rank envelope test (right hand column) for trait disparity through time for three showcase

536 datasets. In each panel, the empirical pattern (solid black line) is compared to the median of 5000

537 simulations of the null model of Brownian evolution (broken lines), and the shaded regions

538 correspond to the 95% confidence intervals calculated using the pointwise (left column) and rank

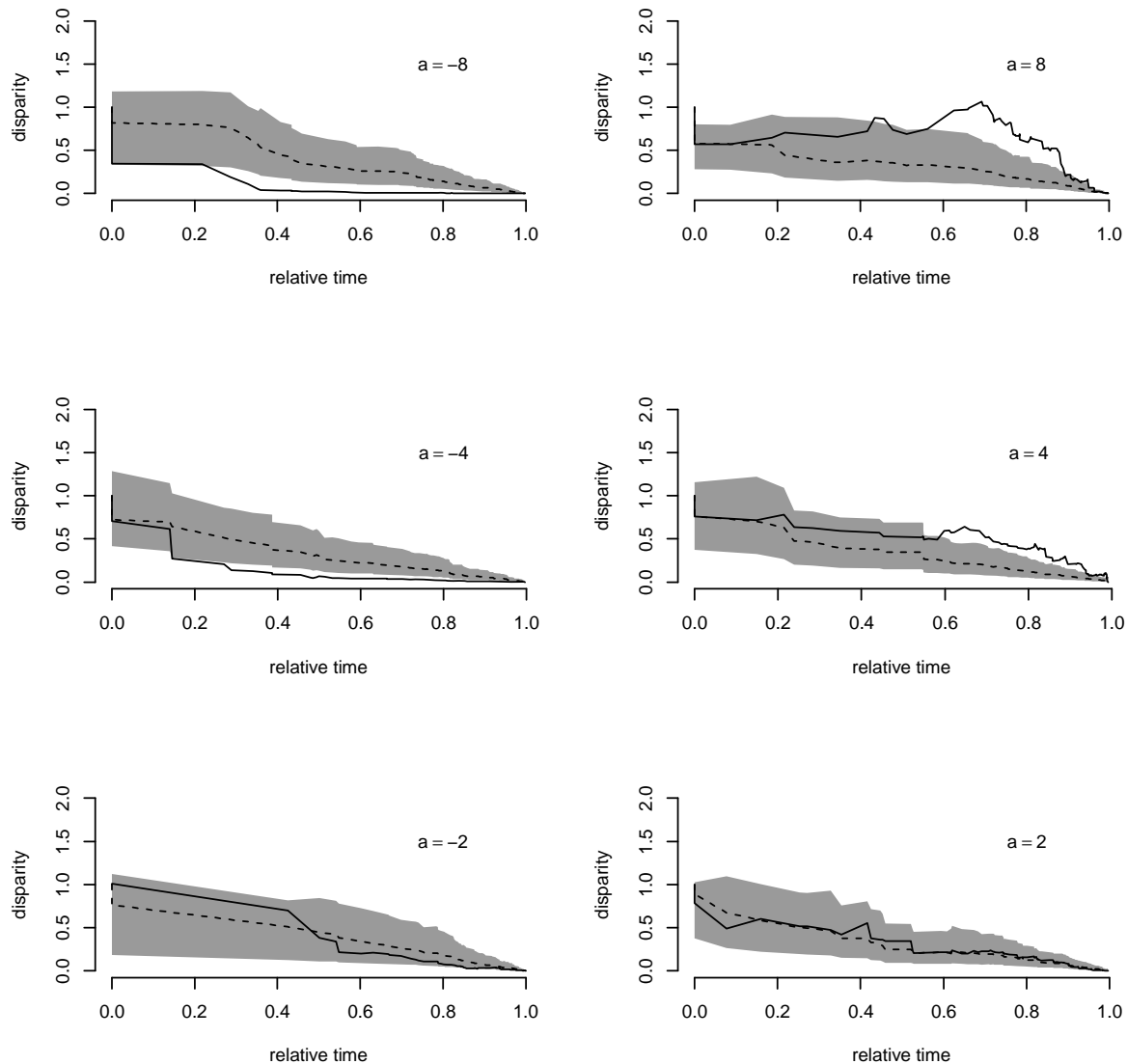
539 envelope (right column) methods. Top row is for Darwin's finches (*Geospiza*) and the evolution of

540 culmen length; middle row is for the evolution of Cetacean body size (Slater et al., 2010); bottom

541 row is for the evolution of anal fin shape in 131 species of African cichlids (Feilich, 2016).

542

## 543 Supplementary Information

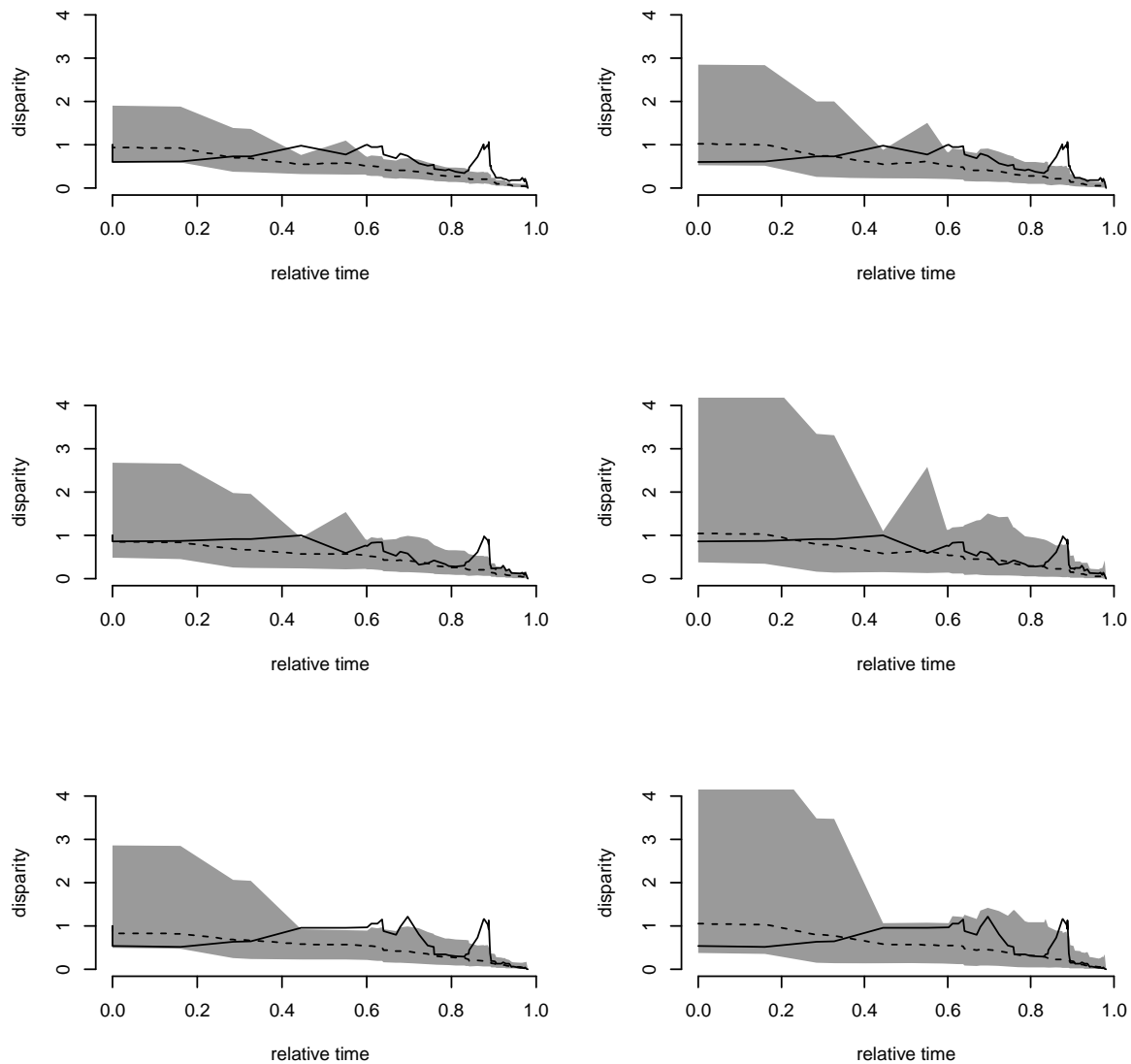


544

545

546 **Figure S1.** Examples of simulated data sets for disparity through time analyses where the rate of  
547 trait evolution either slows down ( $a < 0$ ) or speeds up ( $a > 0$ ) over time (black solid lines). Broken  
548 lines represent the median disparity through time for 2500 simulations of the null model of  
549 Brownian evolution (ie  $a = 0$ ), and the shaded region is the 95% confidence interval for the  
550 simulations based upon the rank envelope method (see main text). Simulations are for simulated  
551 phylogenetic trees with 100 species at the tips.

552



553

554 **Figure S2.** Comparisons of inference from using the (left hand column) pointwise envelope test and  
555 the (right hand column) rank envelope test for trait disparity through time for (top row) body shape;  
556 (middle row) caudal fin shape; and (bottom row) dorsal fin shape, for 131 species of African cichlid  
557 fishes. Data is taken from (Feilich, 2016). In each panel the empirical pattern (solid black line) is  
558 compared to the median of 5000 simulations of the null model of Brownian evolution of the trait  
559 values (broken lines) on the phylogenetic tree, and the shaded regions correspond to the 95%  
560 confidence region.

561