Version dated: April 19, 2018

# Influence of different modes of morphological character correlation on phylogenetic tree inference

THOMAS GUILLERME[1,2,*], AND MARTIN D. BRAZEAU[2,3]

[1]*School of Biological Sciences, University of Queensland, St. Lucia, Queensland, Australia.*

[2]*Imperial College London, Silwood Park Campus, Department of Life Sciences, Buckhurst Road, Ascot SL5 7PY, United Kingdom.*

[3]*Department of Earth Sciences, Natural History Museum, Cromwell Road, London, SW75BD, United Kingdom.*

**\*Corresponding author.** *guillert@tcd.ie*

## Abstract

Phylogenetic analysis algorithms require the assumption of character independence - a condition generally acknowledged to be violated by morphological data. Correlation between characters can originate from intra-organismal features, shared phylogenetic history or forced by particular character-state coding schemes. Although the two first sources can be investigated by biologists *a posteriori* and the third one can be avoided *a priori* with good practices, phylogenetic software do not distinguish between any of them.

In this study, we propose a new metric of raw character difference as a proxy for character correlation. Using thorough simulations, we test the effect of increasing or decreasing character differences on tree topology. Overall, we found an expected positive effect of reducing character correlations on recovering the correct topology. However, this effect is less important for matrices with a small number of taxa (25 in our simulations) where reducing character correlation is not more effective than randomly drawing characters. Furthermore, in bigger matrices (350 characters), there is a strong effect of the inference method with Bayesian trees being consistently less affected by character correlation than maximum parsimony trees.

These results suggest that ignoring the problem of character correlation or independence can often impact topology in phylogenetic analysis. However, encouragingly, they also suggest that, unless correlation is actively maximised or minimised, probabilistic methods can easily accommodate for a random correlation between characters.

(Keywords: Character difference, correlation, topology, Bayesian, maximum parsimony)

# INTRODUCTION

24

25    The last two decades have witnessed a "resurgence" of interest in the use of

26   morphological character data in phylogenetic studies. This owes in large part to the use

27   of fossils to undertake at least partial reconstructions of phylogenetic trees, especially

28   where ancestral states reconstructions or absolute calibrations of divergence times are

29   necessary. Morphological character data are often considered inferior to molecular

30   sequence data, but are often the only source of phylogenetic data for extinct species.

31   While there is a general appreciation of the limits of morphological data, they are

32   frequently dismissed without any empirical investigations into their statistical

33   properties. As morphological data are likely to continue to play an extensive role in

34   phylogenetic analysis, it is essential to understand the circumstances under which

35   morphological data might be expected to "misbehave". This opens up possibilities for

36   predicting problematic datasets and possibly proposing new confidence measures in

37   phylogenetic datasets.

38    The non-independence of large numbers of morphological characters is often

39   cited in anticipation of problems with morphological data. The assumption of character

40   independence is central to phylogenetic inference methods such as maximum

41   likelihood and maximum parsimony (e.g. Joysey and Friday, 1982; Felsenstein, 1985;

42   Lewis, 2001; Felsenstein, 2004). Especially for discrete morphological data, this

43   assumption of independence is probably violated frequently due to the very nature of

44   phylogenetic data: correlations are expected to occur (to some degree) when characters

are depending on each other. Before discussing character correlation further, it is important to understand that it may manifest itself in at least three distinct ways:

- **Intra-organismal dependence:** this is the result of an intrinsic biological link between two characters through development, pleiotropy, and/or biological function. For example the lower and upper molar characters in mammals generally occlude one another. Therefore, one character describing a feature of a lower molar will be expected to be complemented by the surface of the occluding upper molar. Characters of the occlusal surface of two opposing molars will be expected to directly covary. Pleiotropy also results in covariation between different aspects of phenotype. From a phylogenetic perspective, it can be especially pernicious because the relationship between the traits in question may have no obvious link from a morphological or functional comparison alone. Intra-organismal links can be the targets of comparative developmental biology (Goswami and David Polly, 2010; Kelly and Sears, 2010; Stoessel et al., 2013; Goswami et al., 2014) or functional investigations.

- **Evolutionary dependence:** this is the result of sets of characters co-evolving due to selection, likely related to functional links between two traits that help serve an overall lifestyle trait. Unlike the case of intra-organismal dependence, there need not be an intrinsic constraint that causes these traits to covary. For example, in vertebrates, axial elongation can be correlated to limb reduction with snake-like bodies evolving multiple times in numerous tetrapod lineages. This is thought to

4

correspond to adaptations for fossoriality or aquatic lifestyles. Such covariances are generally studied in the context of a given phylogeny, often one derived from molecular data with the morphological traits of interest mapped on it. Many methods have been developed to study these correlations, especially since they can provide us with a lot of information on how specific groups acquired specific characteristics (Russell Lande, 1983; Maddison, 1990; Pagel, 1994; Mark Pagel, 2006; Grabowski and Porto, 2016). However, these methods do not give us a means to objectively control correlations that might adversely affect phylogenetic inference.

- **Coding dependence:** this is the results of researcher methodology for defining or/and coding discrete morphological characters (Brazeau, 2011; Simões et al., 2017). Coding dependence manifests itself in several ways, particularly in coding redundant information. For instance, coding for the same absence in different characters creates state transformations associated with the loss or gain of a particular character. This occurs when a number of multistate characters include two variable feature states (e.g. large, small; red, blue etc.) in conjunction with absence. It is worth noting, however, that these correlations could also be due to the nature of the available data, especially in palaeontology. For example, when only one fragmentary molar is available to describe a specimen, researchers have to "extract" as much phylogenetic information from the available data as possible, potentially inducing correlations. This coding dependency is linked to hierarchical

87    dependency between characters (Wilkinson, 1995; Brazeau et al., 2017). Finally

88    this can also be due to a bias in the amount of characters available. For example,

89    in skulls, because of their complexity, there is a high likelihood of inducing

90    correlation (by effectively reducing structural complexity to discrete characters).

91   Of course, the three sources of dependence have an interaction: characters describing

92   the left and right lower/upper molars will have induced dependence due to the

93   modularity of the molars, their shared history and the duplicated coding. Logical

94   dependence, however, is easily distinguished prior to phylogenetic inference, while the

95   two other ones (intra-organismal and evolutionary) are much harder. However, the

96   development of algorithms and software has not yet caught up with the need to deal

97   with these interdependencies (De Laet, 2015; Brazeau et al., 2017). Intra-organismal

98   dependence requires more detailed, often extremely time-consuming studies (and

99   possibly beyond the limits of available technology). Even after all of the effort is

100  expended, the results might then only be known for a single (model) species.

101  Evolutionary dependence itself requires the resolution of a phylogenetic tree, and is

102  best determined by independent character sets. This is frequently accomplished by

103  mapping morphological traits on molecular phylogenetic trees.

104    These sources of dependence between characters are well studied in biology.

105  Biological and evolutionary dependences are inherent parts to Evo-Devo and

106  macroevolutionary studies and best practices to avoid coding-induced dependences are

107  commonly known and applied. However, eventually, all these characters, whether they

108  are independent or not are analysed through phylogenetic inferences software that are

109  blind to these distinctions. If fact, what the software are confronted with is a two

110  dimensional matrix problem that renders the morphological subtleties described

111  opaque. This introduces a new, less studied, source of character dependence:

112  **Correlation between characters detected by the software**: this is the result how

113  software actually interprets the differences between characters. The vast majority of

114  phylogenetic software ignores both the character's definition and the different states

115  signification (simply treating them as different or similar tokens). Therefore a great

116  number of characters and - traditionally - a few number of tokens can easily lead to

117  dependence between characters. For example, if we consider the following matrix

118  containing four cetartiodactyls - say a pig (e.g. *Sus*), a deer (*Cervus*), a hippo

119  (*Hippopotamus*) and a whale (*Balaenoptera*) - and four binary characters - say (**C1**:

120  presence (1) or absence (0) of an astragalus; **C2**: presence (0) or absence (1) of baleen;

121  **C3**: presence (0) or absence (1) of a left astragalus with a double pulley; **C4**: presence

122  (0) or absence (1) of a right astragalus with a double pulley:

123  In the example in Table 1, the characters **C1** and **C2** are the most likely to be

124  truly independent; characters **C3** and **C4** suffer from a a coding induced dependency;

125  characters **C1** and **C3/C4** have an evolutionary induce dependency and again,

126  characters **C2** and **C3/C4** are likely to be independent. Yet a phylogenetic software will

127  treat all these four characters in exactly the same way: only the sheer difference

128  between the character states tokens will be used in order to infer the tree. Some

7

|  | C1 | C2 | C3 | C4 |
|---|---|---|---|---|
| *Sus* | 1 | 1 | 1 | 1 |
| *Cervus* | 1 | 1 | 1 | 1 |
| *Hippopotamus* | 1 | 1 | 1 | 1 |
| *Balaenoptera* | 0 | 0 | 0 | 0 |

Table 1: Example of a matrix with software induced character correlation. **C1**: presence (1) or absence (0) of an astragal; **C2**: presence (0) or absence (1) of baleens; **C3**: presence (0) or absence (1) of a left astragalus with a double pulley; **C4**: presence (0) or absence (1) of a right astragal with a double pulley.

129 characters will therefore be expected to covary in non-phylogenetic way, and that this

130 phenomenon can reasonably be expected to mislead phylogenetic analysis. Yet the

131 question has never been explored through a thorough simulation framework (although

132 it has been tackled empirically for morphological data Dávalos et al. 2014 or molecular

133 data Zou and Zhang 2016).

134       How does these correlation really affect topology? We expect matrices with a

135 high level of correlation to recover precise but inaccurate topologies but will matrices

136 with low level of correlation (i.e. with high levels of homoplasy) actual cancel out the

137 effects of correlation? Here we formally assess the effect of discrete character's

138 correlation using simulated data. We propose a new distance metric to measure the

139 difference between characters (as a proxy for these three sources of correlation as

140 interpreted by the software) and a protocol to modify discrete morphological matrices

141 to increase/decrease the overall differences or similarities between characters. We

142 found that overall, there is a detectable effect of character correlation on topology

143 where an increase in character dependence results in a decrease in the ability to recover

144 the correct topology. These results, however, vary greatly in magnitude depending on

145 the size of matrix and the inference method used.

# Methods

146

147      To assess the effects of character correlation on the accuracy of phylogenetic

148 inference we generated a series of matrices exhibiting different levels of correlation

149 between some characters (Fig.1 - note that each step is described in more details below):

150   1. **Simulating matrices**: we simulated discrete morphological matrices with 25, 75

151      and 150 taxa for 100, 350 and 1000 characters, hereafter called the "normal"

152      matrices. This step resulted in 9 matrices.

153   2. **Modifying matrices**: we changed the "normal" matrices by modifying the

154      characters in order to maximise or minimise characters differences (hereafter

155      called respectively "maximised" and "minimised" matrices) by removing

156      respectively the least different or most different characters and replacing them

157      randomly by the remaining characters. Our protocol for measuring character

158      difference is detailed below.

9

159    We also randomly duplicated characters from the "normal" matrices without

160    biasing towards maximising or minimising character differences to create

161    randomised matrices (hereafter called the "randomised" matrices - equivalent to a

162    null expectancy). This step resulted in 36 matrices.

163    3. **Inferring topologies**: we inferred the topologies from the "normal",

164    "maximised", "minimised" and "randomised" matrices using both maximum

165    parsimony and Bayesian inference. Hereafter, the resultant topologies are called

166    the "normal", "maximised", "minimised" and "randomised" trees). This step

167    resulted in 72 topologies.

168    4. **Comparing topologies**: finally, we compared the "normal" to the "maximised",

169    "minimised" and "randomised" trees to measure the effect of character

170    correlation on topology.

171 Each step was replicated 35 times and are described below in more detail, along with

172 our proposed definition for measuring the difference between characters.

173                         *Measuring differences between characters*

174 To measure the effect of character correlation as interpreted by the phylogenetic

175 software, we define characters as being entirely correlated if they give the same

176 phylogenetic information. In order to measure this, we propose a new distance metric

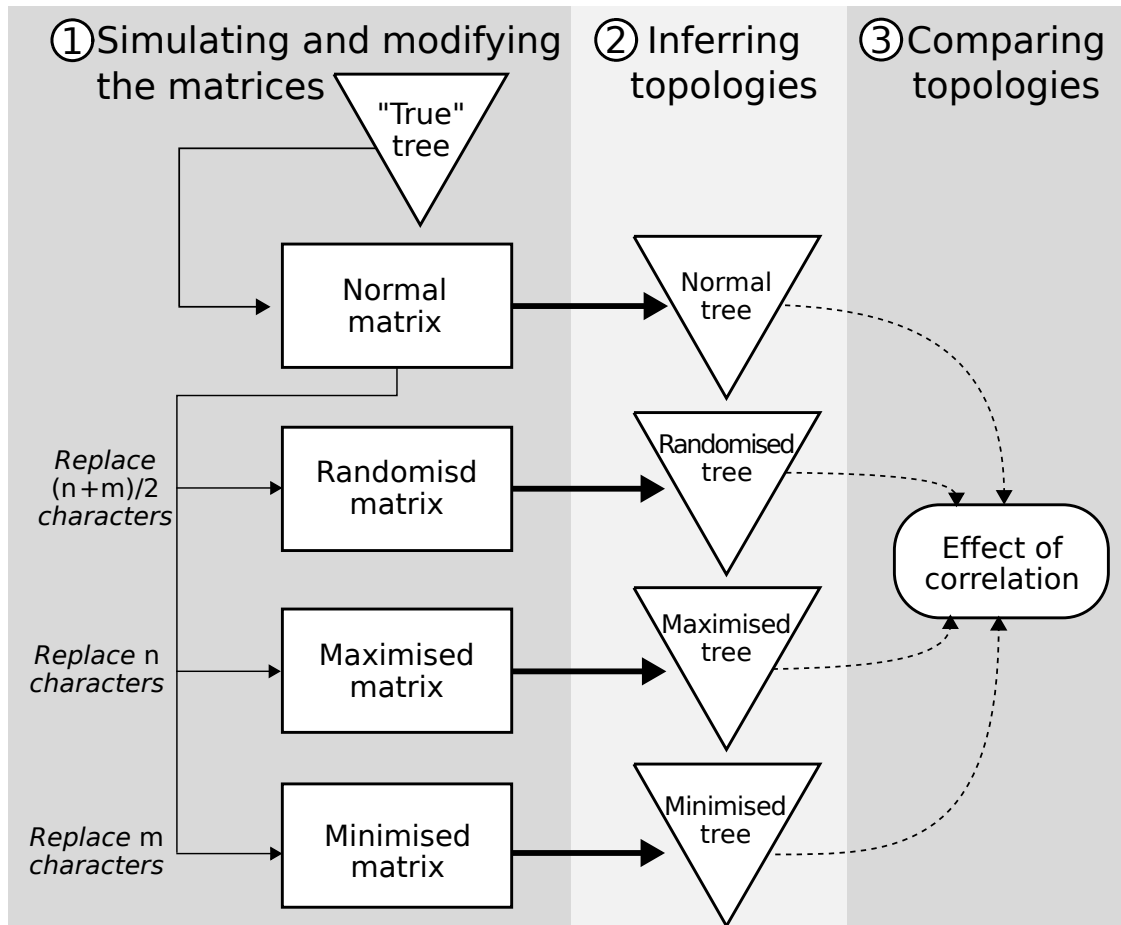177 to measure the difference between two characters:

10

Figure 1: Outline of the simulation protocol: the first step includes both the simulation and the modification of the matrices (thin solid lines); the second step includes tree inference using MP and BPP methods (thick solid lines); the third step includes comparing the resulting tree topologies (dashed lines). *n* and *m* corresponds to the number of characters with a character difference $< 0.25$ and $> 0.75$ respectively.

178 *Character Difference (CD).—*

$$CD_{(x,y)} = 1 - 2\left(\left|\frac{\sum_i^n |x_i - y_i|}{n} - \frac{1}{2}\right|\right) \tag{1}$$

179  Where $n$ is the number of taxa with comparable characters $x$, $y$ and $x_i$, $y_i$ are each

180  character's state for the $i^{th}$ taxon. *CD* is a continuous distance metric bounded between

181  0 and 1 (see the mathematical demonstration in the supplementary material 1). Since

182  we are considering differences as being only Fitch-like (non-additive) and unweighted,

183  we calculated the difference between character states in a qualitative way. Two same

184  character states tokens have a difference of zero and two different ones have a

185  difference of one (e.g. $0 - 0 = 0$ or $1 - 8 = 1$). Additionally, we only consider

186  differences for taxa with shared information (i.e. a Gower distance; Gower, 1971).

187     We standardised each character by arbitrarily modifying their character state

188  tokens (or symbols) by order of appearance. In other words, we replaced all the

189  occurrences of the first token to be 1, the second to be 2, etc. This procedure was used

190  to treat all the characters are unordered with no assumption on the meaning of the

191  character state (e.g. in a binary character 0 is not necessary ancestral to 1). It also

192  greatly improved the speed of our algorithm implementation to compare the characters.

193  This way, a character `A = {2,2,3,0,0,3}` for six taxa would be standardised as `A' =`

194  `{1,1,2,3,3,2}` (following the *xyz* notation in Felsenstein, 2004, p.13). Note that in

195  terms of phylogenetic signal, both `A` and `A'` are exactly identical (forming three distinct

196  splits in the tree inference process).

197     When the character difference is null (0) it means that characters convey the

12

198  same phylogenetic signal (i.e. characters are entirely correlated). When the character

199  difference is maximal (1) it means it conveys the greatest difference in phylogenetic

200  signal (i.e. characters are uncorrelated). It is important to stress that a character

201  difference of 0 (i.e. the same phylogenetic signal) does not mean the opposite of 1 (i.e.

202  *not* the opposite phylogenetic signal but the most different number of implied splits) .

203  For example with three characters A = $\{0,1,1,1\}$, B = $\{1,0,0,0\}$ and C = $\{0,1,2,3\}$,

204  $CD_{(A,B)} = 0$ and $CD_{(A,C)} = 1$. Because the character is continuous and bounded

205  between $(0,1)$, it can be interpreted as the probability of two characters leading to a

206  different set of splits (i.e. a different phylogenetic signal).

### *Simulating discrete morphological matrices*

208  To simulate the matrices we applied a protocol very similar to Guillerme and Cooper

209  (2016b). First, we generate random birth-death trees with the birth ($\lambda$) and death ($\mu$)

210  parameters sampled from a uniform $(0,1)$ distribution maintaining $\lambda > \mu$ using the

211  `diversitree` R package (v0.9-8; FitzJohn, 2012) and saving the tree after reaching either

212  25, 75 or 150 taxa. For each tree, we arbitrarily set the outgroup to be the first taxon

213  (alphabetically) thus effectively rooting the trees on this taxon. These trees are hereafter

214  called the "true" trees (see distinction below). We then simulated discrete

215  morphological characters on the topology of these trees using the either of the two

216  following models:

217  • The morphological HKY-binary model (O'Reilly et al., 2016) which is an HKY

218      model (Hasegawa et al., 1985) with a random states frequency (sampled from a

219   Dirichlet distribution $Dir(1,1,1,1)$) and using a transition/transvertion rate of 2

220   (Douady et al., 2003) but where the purines (A,G) were changed into state 0 and

221   the pyrimidines (C,T) in state 1. This model has the advantage of not favouring

222   Bayesian inference (since it doesn't use an M$k$ model; O'Reilly et al., 2016, ; see

223   discussion) but the downside of it is it can only generate binary state characters

224   (or 4 states; Puttick et al., 2017).

225   • To generate more than binary states characters, we used the M$k$ model (Lewis,

226   2001). We draw the number of character states with a probability of 0.85 for

227   binary characters and 0.15 for three state characters (Guillerme and Cooper,

228   2016b; Zou and Zhang, 2016). This model assumes a equal transition rate between

229   character states which might seem overly simplistic, excluding other observed

230   transition patterns (e.g. Dollo characters; Dollo, 1893; Wright et al., 2015).

231   Recently however, Wright et al. (2016) have shown that an equal rate transition is

232   still the most present in empirical data.

233   For each character, both models (morphological HKY-binary or M$k$) where chosen

234   randomly and run with an overall evolutionary rate drawn from a gamma distribution

235   ($\beta = 100$ and $\alpha = 5$). These low evolutionary rate values allowed reduction in the

236   number of homoplasic character changes, thus reinforcing the phylogenetic information

237   in the matrices. We re-simulated every invariant characters to obtain a matrix with no

238   invariant characters in order to better approximate real morphological data matrices. To

239   ensure that our simulations were reflecting realistic observed parameters, we only

14

240  selected matrices with Consistency Indices (CI) superior to 0.26 (O'Reilly et al., 2016).

241      For each tree with 25, 75 or 150 taxa we generated matrices with 100, 350 and

242  1000 characters following O'Reilly et al. (2016). The matrices were generated using the

243  dispRity R package (Guillerme, 2016). To estimate the variance of our simulations and

244  assess the effect of our random parameters, we repeated this step 35 times resulting in

245  315 "normal" morphological matrices.

246                      *Modifying the matrices*

247  We calculated the pairwise character differences for each generated matrix using the

248  dispRity R package (Guillerme, 2016). We then modified the matrices to either

249  maximise or minimise the pairwise character differences for each matrix using three

250  different algorithms. For maximising the pairwise differences between characters, we

251  selected the characters that were the most similar to all the others (i.e. with an average

252  character difference $< 0.25$) and replaced them randomly by any of the remaining

253  characters. This operation increased the overall pairwise character difference in the

254  matrix thus making the characters more dissimilar. Conversely, for minimising the

255  pairwise character differences, we selected the most dissimilar characters (i.e. with an

256  average character difference $< 0.75$) and randomly replaced them with the remaining

257  ones. Finally, because this operation effectively changes the weight of characters (i.e.

258  giving the characters $< 0.25$ or $> 0.75$ a weight of 0 and giving the randomly selected

259  remaining characters a weight of +1), we randomly replaced the average number of

260  characters replaced in the character maximisation and minimisation by any other

15

261  characters as a randomised expectation modification (i.e. randomly weighting

262  characters). Each of the three matrices are effectively a bootstrap pseudo-replication of

263  the "normal" matrix with the "randomised" one being a random one and the

264  "maximised" and "minimised" being conditional bootstraps. This step resulted in a

265  total of 1260 matrices (hereafter called the "normal", "maximised", "minimised" and

266  "randomised" matrices - see Fig. 2 for an illustration). The algorithms for the three

267  modifications are available on GitHub

268  (https://github.com/TGuillerme/CharactersCorrelation)

269                              *Inferring topologies*

270  We inferred the topologies with both BPP and MP using MrBayes (v3.2.6; Ronquist

271  et al., 2012) and PAUP* (v4.0a151; Swofford, 2001) respectively. For both methods, we

272  used the arbitrarily chosen outgroup in the simulations to root our trees. The

273  maximum parsimony inference was run using a heuristic search with random sequence

274  addition replicate 100 times with a limit of $5 \times 10^6$ rearrangements per replicates

275  (hsearch addseq=random nreps=100 rearrlimit=5000000 limitperrep=yes).

276        Bayesian inference was run using an M$k$ model with ascertainment bias and four

277  discrete gamma rate categories (M$kv$ 4Γ - lset nst=1 rates=gamma Ngammacat=4) with

278  an variable rate prior an exponential (0.5) shape (prset ratepr=variable

279  Shapepr=Exponential(0.5)). We ran two runs of 6 chains each (2 hot, 4 cold) for a

280  maximum of $1 \times 10^9$ generations with a sampling every 200 generations. We

281  automatically stopped the MCMC when the average standard deviation of split
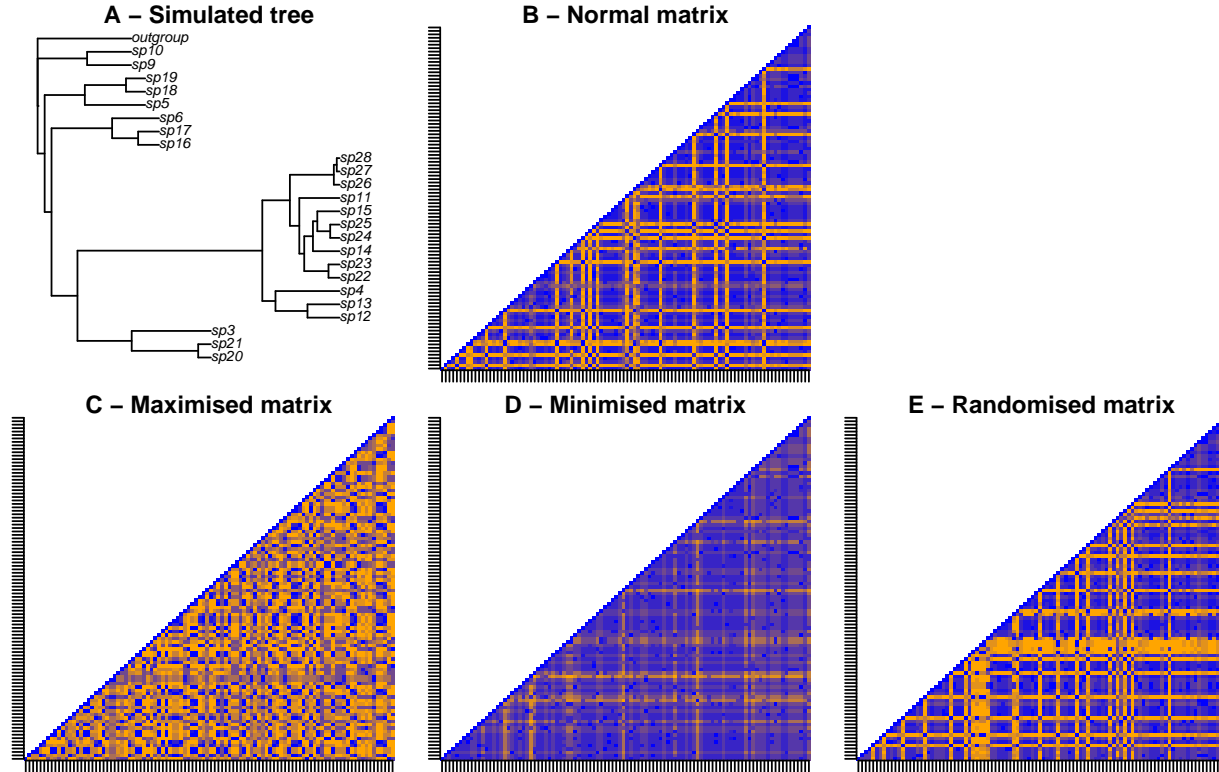
16

Figure 2: Example illustration of the protocol for modifying matrices. The matrices represent the pairwise character differences for 100 characters. Blue colours correspond to low character differences and orange colours correspond to high character differences. **A** - a random Birth-Death tree is simulated and used for generating the "normal" matrix (**B**), characters in this matrix are then removed or duplicated to favour maximised (**C**), minimised (**D**) or randomise character difference (**E**). The differences between the characters is low in **C** (minimised compared to **A**) implying a high correlation between the characters. Conversely, the character differences is high in **D** (maximised compared to **A**) implying a low correlation between the characters.

282  frequencies (ASDSF) between both runs fell below 0.01 (with a diagnosis every $1 \times 10^4$

283  generations - `mcmc nruns=2 Nchains=6 ngen=1000000000 samplefreq=200`

284  `printfreq=2000 diagnfreq=10000 Stoprule=YES stopval=0.01 mcmcdiagn=YES`). Due

285  to cluster hardware requirements an to save some time, when chains didn't converged

286  and the runs exceeded 5GB each, we aborted the MCMC and computed the consensus

287  tree from the unconverged chains. In practice, these few MCMC got stuck at an ASDSF

288  around (but not below) 0.01.

289       A strict majority rule tree was then calculated for both Bayesian an maximum

290  parsimony trees. For the Bayesian consensus trees, the 25% first trees of the posterior

291  tree distribution were excluded as a burnin. The 2880 tree inferences took around one

292  CPU century on the Imperial College High Performance Computing Service (2-3GHz

293  clock rate; ICHPC, 2011).

## *Comparing topologies*

295  We compared the topologies using the same approach as in Guillerme and Cooper

296  (2016b): we measured both the Robinson-Foulds distance (Robinson and Foulds, 1981)

297  and the triplets distance (Dobson, 1975) between the trees inferred from the

298  "maximised", "minimised" and "randomised" matrices and the tree inferred from the

299  "normal" matrix. We explored the effect of character difference on recovering the

300  "normal" topology by comparing the "maximised", "minimised" and "randomised"

301  trees to the "normal" tree (Figs 3 and 4 and supplementary materials 3 Figs 1 and 2).

302  Note that we are not comparing the trees to the "true" tree used to simulate the

matrices. First, in biology, this tree is always unknown. Second, our objective is to

measure the direct effect of character correlation approximated by the difference in

topology between the "normal", "maximised" and "minimised" trees. When measuring

the difference between these trees and the "true" tree, we would also confound the

effect of simulating a birth-death tree and simulating a discrete morphological matrices

from it.

The metric scores where calculated using the `TreeCmp` javascript (Bogdanowicz

et al., 2012). The measurements where then standardised using the Normalised Tree

Similarity metric ($NTS$; i.e. centering the metrics scores using the mean metric score for

1000 pairwise comparisons between random trees with $n$ taxa; Bogdanowicz et al.,

2012; Guillerme and Cooper, 2016b). When the normalised metric has a score of one it

means both trees are identical, when it has a score of zero it means the trees are no

more different than expected by chance and when it has a score $< 0$ the trees are more

different than expected by chance. The normalised score for both metrics thus reflects

two distinct aspects of tree topology: (1) the Normalised Robinson-Foulds ($NTS_{RF}$)

Similarity reflects the conservation of clades (i.e. a score close to 1 indicates that most

clades are identical in both trees); and (2) the Normalised Triplets Similarity ($NTS_{Tr}$)

reflects the position of taxa (i.e. a score close to 1 indicates that most taxa have the same

neighbours in both trees).

Because both $NTS_{RF}$ and $NTS_{Tr}$ metrics are bounded at one. The residuals of

any model based on the $NTS$ scores were not normal thus preventing the use of

324  parametric tests for comparisons (see online material

325  https://rawgit.com/TGuillerme/CharactersCorrelation/master/Analysis/

326  02-EffectCorrelationFullResults.html). Similarly, a non-parametric Wilcoxon rank

327  test (Hollander et al., 2013) would be biased in its p-value calculation due to the

328  presence of equal values in the $NTS$ distributions (e.g. when multiple trees are equal to

329  the "normal" tree). Therefore, we used a combination of the Wilcoxon rank test with a

330  Bonferonni-Holm corrections (to ensure our significant results were robust to Type I

331  error rate inflation; Holm, 1979) and a simple non-parametric metric for measuring the

332  probability of overlap between two distributions, the Bhattacharyya Coefficient ($BC$;

333  Bhattacharyya, 1943; Guillerme and Cooper, 2016b; Soto et al., 2016). Thus, additionally

334  to the Wilcoxon test results, we considered distribution to be significantly similar if they

335  had an overlap probability $> 0.95$ and different if they had an overlap probability

336  $> 0.05$. Comparisons falling between these range can not be designated as strictly

337  similar/different but can still be ranked (e.g. for three distributions A, B, C, if

338  $BC_{(A,B)} = 0.15$ and $BC_{(A,C)} = 0.65$, we cannot consider either distribution significantly

339  different or similar but $B$ still has a lower probability of being similar to $A$ than $C$).

340      The resulting full simulation was 3.5TB big so is not shared here (though the

341  parameters are). However, the resulting consensus trees on which the topological

342  differences are calculated are available at

343  https://figshare.com/s/7a8fde8eaa39a3d3cf56.

## Results

344

Figure 3: Effect of character difference on recovering the "normal" topology. The y axis represents the Normalised Tree Similarity using Robinson-Fould distance for matrices with 25, 75 and 150 taxa from top to bottom respectively. The x axis represents the different character difference scenarios and tree inference method with the "maximised" character difference in Bayesian (red) and under maximum parsimony (orange), the "minimised" character difference in Bayesian (dark green) and under maximum parsimony (light green) and the "randomised" character difference in Bayesian (dark blue) and under maximum parsimony (light blue) for matrices of 100, 350 and 1000 characters in the panels from left to right.
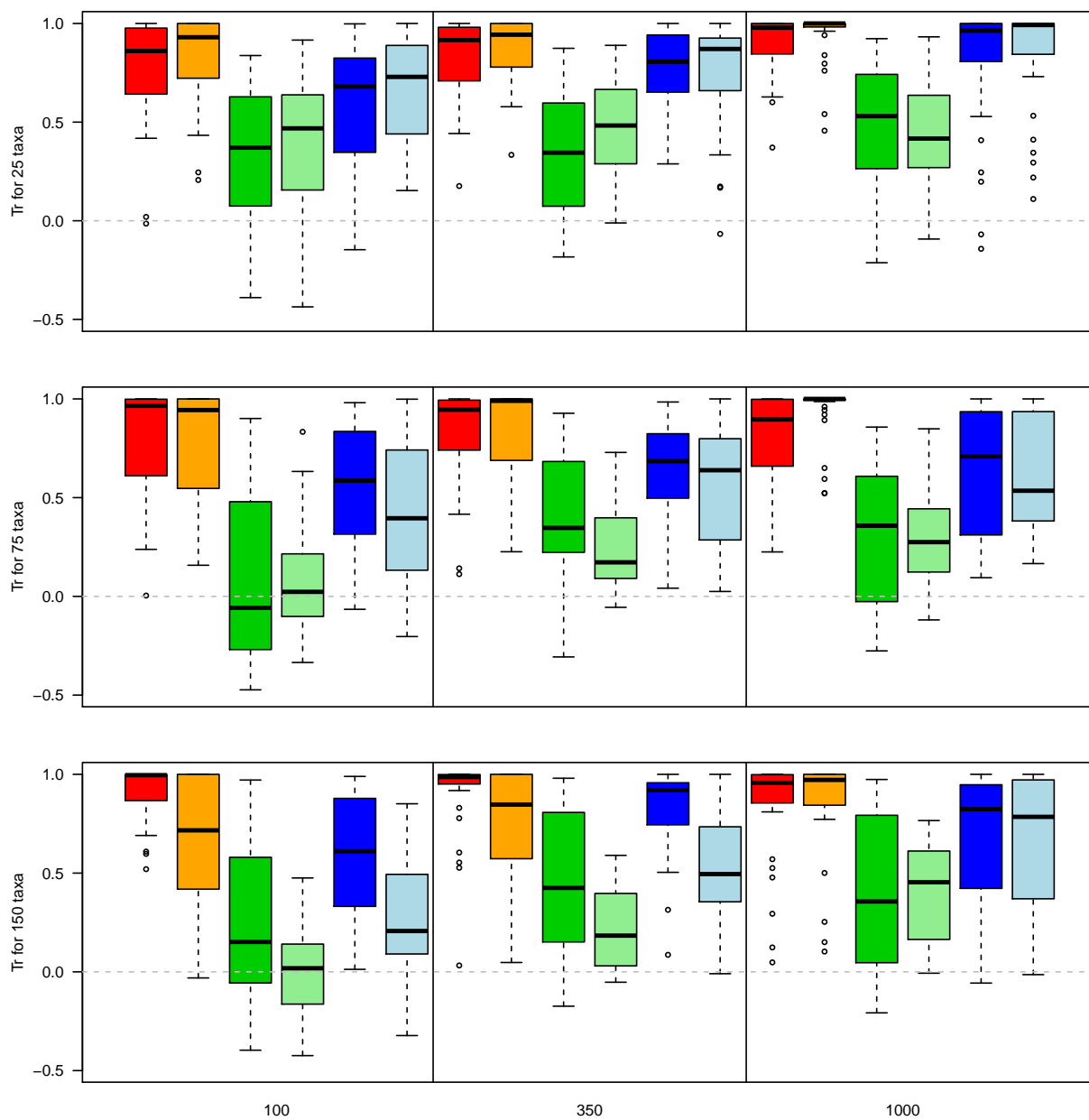
Figure 4: Effect of character difference on recovering the "normal" topology. The axis are identical to figure 3 but y axis represents the Normalised Tree Similarity using Triplets distance.

<p style="text-align:center">345</p>

## *Effect of character differences on topology*

346    The overall amount of character difference in a matrix has an effect of the ability

347    to recover the correct topology when maximising character difference leading to the

348    smallest loss in phylogenetic information (median $NTS_{RF}$ = 0.956 and median $NTS_{Tr}$ =

349    0.839) followed by simply randomising the characters (median $NTS_{RF}$ = 0.762 and

350    median $NTS_{Tr}$ = 0.628) and minimising the character difference (median $NTS_{RF}$ = 0.605

351    and median $NTS_{Tr}$ = 0.303 - see supplementary material 3 for the full summary

352    statistics). There is a significant difference between all scenarios (maximising,

353    minimising and randomising) with the highest probability of overlap being between

354    maximising and randomising the character difference (Bhattacharrya Coefficient of

355    0.873 for the $NTS_{RF}$ and 0.908 for the $NTS_{Tr}$ - Table 2) and the lowest probability

356    between maximising and minimising the character difference (Bhattacharrya coefficient

357    of 0.573 for the $NTS_{RF}$ and 0.614 for the $NTS_{Tr}$ - Table 2)

358    *Number of characters.—* This effect of the character difference is not dependent on the

359    number of characters when looking at clade conservation (i.e. $NTS_{RF}$). The median

360    $NTS_{RF}$ was similar for 100, 350 and 1000 characters (0.730, 0.745, 0.767 respectively -

361    see supplementary materials 3) with a significant difference only between 100 and 1000

362    and 350 and 1000 characters (Table 3). The number of characters affects the character

363    difference more in terms of taxon placement for a low number of characters (median

364    $NTS_{Tr}$ for 100, 350 and 1000 characters equals 0.544, 0.693, 0.799 respectively - see

365    supplementary materials 3) with a significant difference between 100 and 350 or 1000

<p style="text-align:center">23</p>

| metric | test | bhatt.coeff | statistic | p.value |
|--------|------|------------:|----------:|---------|
| RF | maxi:mini | 0.573 | 356436.000 | 0 |
|    | maxi:rand | 0.873 | 287225.000 | 0 |
|    | mini:rand | 0.856 | 95841.500 | 0 |
| Tr | maxi:mini | 0.614 | 358800.000 | 0 |
|    | maxi:rand | 0.908 | 288223.500 | 0 |
|    | mini:rand | 0.858 | 98507.500 | 0 |

Table 2: Difference between the pooled scenarios. Bhatt.coeff is the Bhattacharrya Coefficient (probability of overlap), the statistic and the p.value are from a non-parametric wilcoxon test (with Bonferonni-Holm correction)

characters (Table 3). However, these differences have to be contrasted by a very high probability of overlap between each number of characters and metrics (Bhattacharrya Coefficient always $> 0.95$) suggesting that the significant effects of the number of characters still leads to really similar distributions.

*Number of taxa.—*

Similar to the effect of number of characters on character difference, the number of taxa seems to have only a marginal effect. A low number of taxa (25) resulted in significant differences with both 75 or 150 taxa in both $NTS_{RF}$ and $NTS_{Tr}$ but no differences between 75 and 150 taxa (medians for 25, 75 and 150 taxa equals 0.802, 0.76, 0.763 $NTS_{RF}$ and 0.758, 0.588 and 0.615 $NTS_{Tr}$ respectively - Table 4 and see

| metric | test | bhatt.coeff | statistic | p.value |
|--------|------|-------------|-----------|---------|
| RF | c100:c350 | **0.99** | 190357.500 | 1 |
| | c100:c1000 | **0.98** | 174085.500 | **0.001** |
| | c350:c1000 | **0.984** | 180460.000 | **0.032** |
| Tr | c100:c350 | **0.961** | 166609.500 | **0** |
| | c100:c1000 | **0.956** | 151389.500 | **0** |
| | c350:c1000 | **0.981** | 178793.500 | **0.014** |

Table 3: Difference between the pooled number of characters. Bhatt.coeff is the Bhattacharrya Coefficient (probability of overlap), the statistic and the p.value are from a non-parametric wilcoxon test (with Bonferonni-Holm correciton)

supplementary materials 3). Again, however, the significant differences have to be contrasted with still high probabilities of overlaps for each $NTS_{RF}$ and $NTS_{Tr}$ distributions for every number of taxa (Table 4).

### *Effect of character differences on the inference method*

Regarding the inference method, there is a significant difference in clade conservation between Bayesian and maximum parsimony (Table 5 - median $NTS_{RF}$ of 0.828 and 0.679 respectively) but not in terms of individual taxon placements (Table 5 - median $NTS_{Tr}$ of 0.738 and 0.601 respectively).

### *Combined effects of taxa, characters and correlation on topology*

25

| metric | test | bhatt.coeff | statistic | p.value |
|--------|------|------------:|-----------|---------|
| RF | t25:t75 | **0.976** | 218421.000 | **0.012** |
| | t25:t150 | **0.988** | 220529.000 | **0.004** |
| | t75:t150 | **0.99** | 201037.000 | 1 |
| Tr | t25:t75 | **0.976** | 233282.000 | **0** |
| | t25:t150 | **0.978** | 227288.000 | **0** |
| | t75:t150 | **0.992** | 194201.000 | 1 |

Table 4: Difference between the pooled number of taxa. Bhatt.coeff is the Bhattacharrya Coefficient (probability of overlap), the statistic and the p.value are from a non-parametric wilcoxon test (with Bonferonni-Holm correciton)

When looking at the combined effect of each parameter, the "maximised" and "minimised" scenarios are always significantly different with no high probability of overlap for both $NTS_{RF}$ and $NTS_{Tr}$ (Wilcoxon rank test p.value $< 0.05$ and Bhattacharrya Coefficient $< 0.95$ - see supplementary material 3). The same differences are observed when comparing the "maximised" scenario against the "randomised" one expect for: (1) the Bayesian inference with 25 taxa (with 100, 350 and 1000 characters) and with 75 taxa with 1000 characters for both $NTS_{RF}$ and $NTS_{Tr}$; and (2) the maximum parsimony for 25 taxa (with 350 and 1000) characters for both $NTS_{RF}$ and $NTS_{Tr}$ and 75 taxa with 100 characters for $NTS_{Tr}$. Identically, there was always a significant difference between the "minimised" scenario and the "randomised" one was

| metric | test | bhatt.coeff | statistic | p.value |
|--------|------|-------------|-----------|---------|
| RF | bayesian:parsimony | 0.891 | 579437.500 | **0** |
| Tr | bayesian:parsimony | **0.984** | 470621.500 | 0.084 |

Table 5: Difference between the pooled methods. Bhatt.coeff is the Bhattacharrya Coefficient (probability of overlap), the statistic and the p.value are from a non-parametric wilcoxon test (with Bonferonni-Holm correciton)

395  expect for the matrix of 150 taxa and 100 characters under maximum parsimony for

396  $NTS_{RF}$ and the matrix of 150 and 1000 characters under Bayesian inference for $NTS_{Tr}$.

397  The full list of comparisons and summary statistics are available in the supplementary

398  materials 3.

# Discussion

399

## *Effect of character differences on topology*

400

401  As expected, there is a significant effect of the character difference in the ability to

402  recover the correct topology. The character difference metric can be seen as the inverse

403  of character correlation (see Methods ): a high character difference approximates a low

404  level of character correlation and vice versa. When characters are correlated, one could

405  expect the matrices to convey a strong (but potentially misleading) phylogenetic signal

406  since every character agrees with each other and conversely, when characters are

27

uncorrelated, one could expect them to convey a weaker phylogenetic signal with a high amount of homoplasy. Intuitively, this would lead the "minimised" character difference scenario to lead to incorrect but consistent trees, the "maximised" scenario to lead to poorly resolved once (really homoplasic trees) and the "randomised" scenario to perform the best at recovering the correct topology. Although the expected results appear to be true for a low character difference scenario, increasing the character difference surprisingly improves the ability to recover the "normal" topology both in terms of clade conservation ($NTS_{RF}$) and taxa placement ($NTS_{Tr}$) for both inference methods (especially in bigger matrices; Figs 3 and 4). Furthermore, the trees generated by the "minimised" scenario do not appear better resolved (towards any topology) than the other scenarios (see Supplementary material 3, Figs 3, 4 and 5).

*Number of characters and taxa.*— Because of the nature of our simulation protocol, one could expect that the effect of character correlation would have increased with the number of characters (i.e. the more characters available, the more characters are modified in each scenario). Similarly, one could expect the number of taxa to have an effect of the raw ability to recover the "normal" topology (i.e. the more taxa, the more likely taxa are misplaced by chance).

Although we measured a significant difference between "small" and larger matrices (both in terms of number of taxa and characters; Tables 3 and 4), these differences have to be contrasted with the probability of overlap between the results distributions that are always really high between every matrices sizes (Bhattacharrya

Coefficients > 0.965 for both the different number of characters and taxa). This suggest

that the effect of character correlation on recovering the right topology is independent

of the size of the matrix when pooling the data. For the number of characters, this

suggests that the overall character difference metric is a good proxy for character

correlation as it is independent of the number of characters analysed. Similarly, using

the a Normalised Tree Similarity metric ($NTS$) accounts for the fact that topological

difference is affected by the sheer number of taxa considered (i.e. we corrected for the

expected difference when comparing two random trees with the same number of taxa).

## *Effect of character differences on the inference method*

When considering the pooled effect of the tree inference method, we only detected a

significant difference between the Bayesian and the maximum parsimony trees in terms

of clade conservation but none in terms of taxa placement (both using a Wilcoxon test

and the Bhattacharrya Coefficient; Table 5). The difference in the ability of each method

to recover the "correct" topology has been heavily discussed in the last five years with

some indications that Bayesian inference will outperform parsimony when analysing

discrete morphological characters alone (Wright and Hillis 2014; O'Reilly et al. 2016;

Puttick et al. 2017; although some critics have raised issues with these investigations

Spencer and Wilberg 2013; Goloboff et al. 2017). In this study, it is possible that our

simulation protocol for generating the characters (favouring slightly more M*k*-based

characters rather than HKY ones) could slightly favour Bayesian inference over

maximum parsimony, however, our protocol for selecting matrices (i.e. those with in a

29

*CI* < 0.26 in a quick preliminary parsimony search; O'Reilly et al., 2016) could also favour maximum parsimony analysis. It was however not the purpose of this study to compare the overall performance of both methods but rather to measure the effect of character correlation on each of those methods separately.

The differences in performance of the two methods observed here could be due to the inherent mechanisms of each method. For any given topology *T* that was obtained from the "normal" matrix and a matrix with high homoplasy, both methods will generate score differently: (1) in parsimony, the topology will probably be given a bad optimality score (on that implies many changes along the tree) and the optimality criterion (favouring the minimum score) will likely discard the tree. The tree search will thus likely result in a topology island that will not contain the given topology *T*. (2) in Bayesian inference, the topology will also be given a bad optimality score (i.e. low likelihood) although the high homoplasy can be accommodated in the tree through high evolution rates or/and long branches. The rate and the branch length being two parameters among others, the optimality score (the likelihood) will change less drastically than for using parsimony. Furthermore, in Bayesian inference, a reasonable difference in optimality between two topologies (the acceptance) will not necessarily mean that the given topology *T* will be discarded. This difference in both mechanisms could explain why, on average, Bayesian Inference seems better to recover the "normal" topology than maximum parsimony.b

*Distinction between different character correlations*

30

470 Here we mention three different types of character correlations but evolutionary

471 biologists are mainly interested in the intra-organismal and evolutionary correlations

472 (e.g. in evo-devo Goswami and Janis 2006; or in macroevolution FitzJohn et al. 2014).

473 These two types of correlations can only be studied *a posteriori* with a phylogenetic

474 hypothesis and should not used *a priori* as a criterion to select characters. In other

475 words, intra-organisaml and evolutionary correlation should be studied based on an

476 underlying phylogenetic framework making the correlation induced by data collection

477 (i.e. coding correlation) the only type of correlation that can affect the phylogeny *a*

478 *priori*. This dichotomy thus creates a trade of between: (1) coding fewer characters

479 (stochastically reducing *a priori* correlation) but making the *a posteriori* correlation more

480 dependent on the coding; and (2) coding more characters (increasing *a priori*

481 dependence) but allowing the *a posteriori* correlation being less dependent on the

482 coding correlations.

483      It is important to note that the two other sources of character correlation could

484 also be present in our simulations although they were not explicitly modelled: (1)

485 evolutionary correlation is implied by simulating the characters based on Birth-Death

486 trees; and (2) intra-organismal correlation could also be present in the matrices for

487 those characters randomly simulated but sharing similar evolutionary simulation

488 regimes (i.e. creating "modules" of characters). However, the effect of these sources of

489 correlation was out of the scope of this study and would have required *a posteriori*

490 changes to the matrices which are - when using empirical data - at best bad practice

491 and at worth dishonest.

## *Limitations*

493 First, simulating evolutionary history is complex. Not only because the models we're

494 using to infer phylogenies are ever improving (e.g. Heath et al., 2014; Wright et al.,

495 2016) but also because generalising morphological evolution across vastly different

496 organisms is probably impossible (see constrasted discussions from Goloboff et al.,

497 2018; O'Reilly et al., 2018). However, we do not compare the "maximised",

498 "minimised" and "randomised" to the "true" tree but rather to the "normal" tree. This

499 allows us to reduce the caveats from our simulations on the effect of character

500 correlation since we only compare the simulation end products to each other (the

501 outputs) rather than to the simulation inputs.

502 Second, measuring and modifying character correlation is difficult. In our

503 simulation protocol we chose to create simulation by duplicating characters in a matrix

504 to maximise or minimise correlation. In biology, this correlation arises from either

505 intra-organismal or evolutionary mechanisms. This could lead to correlations between

506 characters to be more present in some parts of the trees that other (e.g. in the case of

507 inapplicable data  Brazeau et al., 2017). However, because of the number of characters, it

508 is actually complex to actually measure their correlation in a biological sense and is still

509 actively discussed in the literature (Russell Lande, 1983; Maddison, 1990; Pagel, 1994;

510 Mark Pagel, 2006; Goswami and Janis, 2006; Goswami and David Polly, 2010; Goswami

511 et al., 2014; Grabowski and Porto, 2016). Additionally, as discussed in the introduction,

32

512 character correlation can also simply arise by chance due to the discrete coding scheme

513 (i.e. some sets of characters can be highly correlated but effectively describe

514 independent information). Therefore, we made the choice to simplify our simulations

515 by generating character correlation as a stochastic process rather than a biological one.

516　　　　Third, comparing phylogenetic inference methods is not trivial. As mentioned

517 above, both maximum parsimony and Bayesian inference, although aiming (and often

518 achieving) to infer evolutionary history only have similar outputs and vastly differ in

519 how optimality is measured. But there are also difficulties in summarising both

520 methods with consensus trees OReilly and Donoghue (2017). However, we want to

521 point out again that here we're no comparing the methods to each other *per se* but

522 rather how they both, individually, react to an increase or decrease of correlated

523 characters.


524　　　　　　　　　　　　　　　*Potential applications*


525 Effectively, our simulation protocol bootstraps our data "with bias". In the

526 "randomised" scenarios the data is simply randomly bootstrapped simply we

527 randomly remove and resample characters (i.e. giving the weight of 0 to some and $> 1$

528 to other). However, in the "minimised" and "maximised" scenario, the bootstrapping

529 we remove the characters with the lowest/highest overall character difference. For

530 example, in the "maximised" scenario, we randomly remove some characters that are

531 strongly correlated with other and randomly resample from the left characters.

532   It is noteworthy to point that in rather small matrices ($25 \times 100$), there was no

533   significant difference in terms of recovering the right topology when maximising or

534   randomising the character differences. Since many discrete morphological matrices are

535   of similar size (Guillerme and Cooper, 2016a) a simple bootstrap re-sampling (i.e. the

536   equivalent of randomising the character differences in our analysis) will be sufficient to

537   obtain a robust topology (*cf.* actively collecting different characters). In matrices with

538   more taxa, however, the "maximised" scenario resulted in better topological recovery

539   than any other scenarios. Applying this kind of bootstraps that maximises character

540   difference by biasing the random sampling could thus results in better resolved trees.


541                              *Conclusion*


542   Correlation between characters can be induced through three main phenomena:

543   intra-organismal relationships, selection-driven covariation or biases in coding the

544   characters yet only the latter can be improved upon to investigate phylogenetic

545   relationships. Useful best practices guidelines (e.g. Brazeau, 2011; Simões et al., 2017)

546   and algorithms for dealing with different types of character correlations (e.g. for

547   characters hierarchy **?**Brazeau et al., 2017) already exist. However, with the regain of

548   popularity in discrete morphological data and the expansion of dataset size (e.g. Ni

549   et al., 2013; O'Leary et al., 2013, with more than 1000 characters each), we can expect

550   the correlation between characters to increase stochastically. Moreover, because

551   phylogenetic inference software are unable to *a priori* differentiate these difference

34

552 correlations, it is important to understand to what extant topologies can be induced by

553 such bias.

554      We found that character differences as a proxy for character correlation have a

555 strong effect on recovering the "normal" topology: when character correlation was high

556 (low character differences), the topology was always the furthest away from the

557 "normal" topology. Conversely, when correlation between characters was low, the

558 topology was always the closest to the "normal" topology. These results seem

559 independent on the size of the matrix (number of taxa and/or characters) but can be

560 influenced by the phylogenetic inference method used with Bayesian inference faring

561 better in terms of clade conservation, especially in larger matrices.

562      However, in modest size matrices (25 taxa; 100 to 350 characters), the effect of

563 actively choosing to minimise character correlation was not more significant than

564 simply bootstrapping the matrix, suggesting that character correlation is more a

565 problem in large discrete morphological matrices. For such matrices, minimising the

566 character correlation (resampling characters $< 25\%$ different) or maximising it ($> 75\%$)

567 respectively significantly decreased and increased correct topological recovery

568 compared to randomly resample matrices.


## DATA AVAILABILITY, REPEATABILITY AND REPRODUCIBILITY

569

570 The consensus trees are available on figshare at

571 `https://figshare.com/s/7a8fde8eaa39a3d3cf56`. The simulations are fully replicable

572 following the explanations at

573 `https://github.com/TGuillerme/CharactersCorrelation`. The post-simulation

574 analysis, tables and figures (reported in this manuscript) are fully reproducible see

575 (`https://github.com/TGuillerme/CharactersCorrelation`).

# Funding

# Acknowledgements

584 *

585 References

586 Bhattacharyya, A. 1943. On a measure of divergence between two statistical populations

587 defined by their probability distributions. Bulletin of the Calcutta Mathematical

588 Society 35:99–109.

589  Bogdanowicz, D., K. Giaro, and B. Wróbel. 2012. TreeCmp: Comparison of trees in

590     polynomial time. Evolutionary Bioinformatics 8:475–487.

591  Brazeau, M. D. 2011. Problematic character coding methods in morphology and their

592     effects. Biological Journal of the Linnean Society 104:489–498.

593  Brazeau, M. D., T. Guillerme, and M. R. Smith. 2017. Morphological phylogenetic

594     analysis with inapplicable data. bioR$\chi$iv .

595  Dávalos, L. M., P. M. Velazco, O. M. Warsi, P. D. Smits, and N. B. Simmons. 2014.

596     Integrating incomplete fossils by isolating conflicting signal in saturated and

597     non-independent morphological characters. Systematic Biology 63:582–600.

598  De Laet, J. 2015. Parsimony analysis of unaligned sequence data: maximization of

599     homology and minimization of homoplasy, not minimization of operationally

600     defined total cost or minimization of equally weighted transformations. Cladistics

601     31:550–567.

602  Dobson, A. J. 1975. Comparing the shapes of trees vol. 452 of *Lecture Notes in*

603     *Mathematics* Pages 95–100. Springer Berlin Heidelberg.

604  Dollo, L. 1893. Les lois de l'évolution. Bull Soc Belge Geol Pal Hydr Page 164166.

605  Douady, C., F. Delsuc, Y. Boucher, W. Doolittle, and E. Douzery. 2003. Comparison of

606     bayesian and maximum likelihood bootstrap measures of phylogenetic reliability.

607     Molecular Biology and Evolution 20:248–254.

608    Felsenstein, J. 1985. Phylogenies and the comparative method. The American Naturalist

609       125:1–15.

610    Felsenstein, J. 2004. Inferring phylogenies vol. 2. Sinauer Associates Sunderland.

611    FitzJohn, R. G. 2012. Diversitree: comparative phylogenetic analyses of diversification

612       in R. Methods in Ecology and Evolution 3:1084–1092.

613    FitzJohn, R. G., M. W. Pennell, A. E. Zanne, P. F. Stevens, D. C. Tank, and W. K.

614       Cornwell. 2014. How much of the world is woody? Journal of Ecology 102:1266–1272.

615    Goloboff, P. A., A. T. Galvis, J. S. Arias, and A. Smith. 2018. Parsimony and modelbased

616       phylogenetic methods for morphological data: comments on o'reilly etal.

617       Palaeontology 0.

618    Goloboff, P. A., A. Torres, and J. S. Arias. 2017. Weighted parsimony outperforms other

619       methods of phylogenetic inference under models appropriate for morphology.

620       Cladistics .

621    Goswami, A. and P. David Polly. 2010. The influence of character correlations of

622       phylogenetic analyses: a case study of the carnivoran cranium. Pages 141–164 *in*

623       Carnivoran Evolution: New Views on Phylogeny, Form, and Function. (A. Goswami

624       and A. Friscia, eds.). Cambridge University Press, Cambridge.

625    Goswami, A. and C. Janis. 2006. Morphological integration in the carnivoran skull.

626       Evolution 60:169–183.

627  Goswami, A., J. Smaers, C. Soligo, and P. Polly. 2014. The macroevolutionary

628    consequences of phenotypic integration: from development to deep time. Phil. Trans.

629    R. Soc. B 369:20130254.

630  Gower, J. C. 1971. A general coefficient of similarity and some of its properties.

631    Biometrics 27:857–871.

632  Grabowski, M. and A. Porto. 2016. How many more? sample size determination in

633    studies of morphological integration and evolvability. Methods in Ecology and

634    Evolution Pages n/a–n/a.

635  Guillerme, T. 2016. disprity: v0.2.

636  Guillerme, T. and N. Cooper. 2016a. Assessment of available anatomical characters for

637    linking living mammals to fossil taxa in phylogenetic analyses. Biology letters

638    12:20151003.

639  Guillerme, T. and N. Cooper. 2016b. Effects of missing data on topological inference

640    using a total evidence approach. Molecular Phylogenetics and Evolution 94, Part

641    A:146 – 158.

642  Hasegawa, M., H. Kishino, and T. A. Yano. 1985. Dating of the human ape splitting by a

643    molecular clock of mitochondrial-DNA. Journal of Molecular Evolution 22:160–174.

644  Heath, T. A., J. P. Huelsenbeck, and T. Stadler. 2014. The fossilized birth–death process

645    for coherent calibration of divergence-time estimates. Proceedings of the National

646    Academy of Sciences 111:E2957–E2966.

647 Hollander, M., D. A. Wolfe, and E. Chicken. 2013. Nonparametric statistical methods.

648     John Wiley & Sons.

649 Holm, S. 1979. A simple sequentially rejective multiple test procedure. Scandinavian

650     journal of statistics Pages 65–70.

651 ICHPC. 2011. Imperial college high performance computing service. http://www.

652     imperial.ac.uk/admin-services/ict/self-service/research-support/hpc/.

653 Joysey, K. A. and A. E. Friday. 1982. Problems of phylogenetic reconstruction. London

654     etc.:[Systematics Assn. Spec. vol. No. 21.] Academic Press 576:561.

655 Kelly, E. M. and K. E. Sears. 2010. Reduced phenotypic covariation in marsupial limbs

656     and the implications for mammalian evolution. Biological Journal of the Linnean

657     Society 102:22–36.

658 Lewis, P. 2001. A likelihood approach to estimating phylogeny from discrete

659     morphological character data. Systematic Biology 50:913–925.

660 Maddison, W. P. 1990. A method for testing the correlated evolution of two binary

661     characters: Are gains or losses concentrated on certain branches of a phylogenetic

662     tree? Evolution 44:539–557.

663 Mark Pagel, A. E. B. J. C. E. J. B. L., Andrew Meade. 2006. Bayesian analysis of

664     correlated evolution of discrete characters by reversiblejump markov chain monte

665     carlo. The American Naturalist 167:808–825.

666   Ni, X., D. Gebo, M. Dagosto, J. Meng, P. Tafforeau, J. Flynn, and K. Beard. 2013. The

667       oldest known primate skeleton and early haplorhine evolution. Nature 498:60–64.

668   O'Leary, M. A., J. I. Bloch, J. J. Flynn, T. J. Gaudin, A. Giallombardo, N. P. Giannini, S. L.

669       Goldberg, B. P. Kraatz, Z.-X. Luo, J. Meng, X. Ni, M. J. Novacek, F. A. Perini, Z. S.

670       Randall, G. W. Rougier, E. J. Sargis, M. T. Silcox, N. B. Simmons, M. Spaulding, P. M.

671       Velazco, M. Weksler, J. R. Wible, and A. L. Cirranello. 2013. The placental mammal

672       ancestor and the post-K-Pg radiation of placentals. Science 339:662–667.

673   O'Reilly, J. E., M. N. Puttick, L. Parry, A. R. Tanner, J. E. Tarver, J. Fleming, D. Pisani,

674       and P. C. J. Donoghue. 2016. Bayesian methods outperform parsimony but at the

675       expense of precision in the estimation of phylogeny from discrete morphological

676       data. Biology Letters 12.

677   O'Reilly, J. E., M. N. Puttick, D. Pisani, P. C. J. Donoghue, and A. Smith. 2018. Empirical

678       realism of simulated data is more important than the model used to generate it: a

679       reply to goloboff etal. Palaeontology 0.

680   OReilly, J. E. and P. C. Donoghue. 2017. The efficacy of consensus tree methods for

681       summarising phylogenetic relationships from a posterior sample of trees estimated

682       from morphological data. Systematic biology .

683   Pagel, M. 1994. Detecting correlated evolution on phylogenies: A general method for

684       the comparative analysis of discrete characters. Proceedings of the Royal Society of

685       London B: Biological Sciences 255:37–45.

41

686  Puttick, M. N., J. E. O'Reilly, A. R. Tanner, J. F. Fleming, J. Clark, L. Holloway,

687  J. Lozano-Fernandez, L. A. Parry, J. E. Tarver, D. Pisani, et al. 2017. Uncertain-tree:

688  discriminating among competing approaches to the phylogenetic analysis of

689  phenotype data. Proceedings of the Royal Society B 284:20162290.

690  Robinson, D. F. and L. R. Foulds. 1981. Comparison of phylogenetic trees. Mathematical

691  Biosciences 53:131–147.

692  Ronquist, F., M. Teslenko, P. van der Mark, D. L. Ayres, A. Darling, S. Hohna, B. Larget,

693  L. Liu, M. A. Suchard, and J. P. Huelsenbeck. 2012. MrBayes 3.2: efficient Bayesian

694  phylogenetic inference and model choice across a large model space. Systematic

695  Biology 61:539–42.

696  Russell Lande, S. J. A. 1983. The measurement of selection on correlated characters.

697  Evolution 37:1210–1226.

698  Simões, T. R., M. W. Caldwell, A. Palci, and R. L. Nydam. 2017. Giant taxon-character

699  matrices: quality of character constructions remains critical regardless of size.

700  Cladistics 33:198–219.

701  Soto, D. X., J. Benito, E. Gacia, E. García-Berthou, and J. Catalan. 2016. Trace metal

702  accumulation as complementary dietary information for the isotopic analysis of

703  complex food webs. Methods in Ecology and Evolution 7:910–918.

704  Spencer, M. R. and E. W. Wilberg. 2013. Efficacy or convenience? Model-based

42

705    approaches to phylogeny estimation using morphological data. Cladistics 29:663–671.

706

707    Stoessel, A., B. M. Kilbourne, and M. S. Fischer. 2013. Morphological integration versus

708        ecological plasticity in the avian pelvic limb skeleton. Journal of Morphology

709        274:483–495.

710    Swofford, D. L. 2001. Paup*: Phylogenetic analysis using parsimony (and other

711        methods) 4.0. b5 .

712    Wilkinson, M. 1995. Coping with abundant missing entries in phylogenetic inference

713        using parsimony. Syst. Biol. 44:501–514.

714    Wright, A. M. and D. M. Hillis. 2014. Bayesian analysis using a simple Likelihood

715        model outperforms parsimony for estimation of phylogeny from discrete

716        morphological data. PLoS ONE 9:e109210.

717    Wright, A. M., G. T. Lloyd, and D. M. Hillis. 2016. Modeling character change

718        heterogeneity in phylogenetic analyses of morphology through the use of priors.

719        Systematic Biology 65:602–611.

720    Wright, A. M., K. M. Lyons, M. C. Brandley, and D. M. Hillis. 2015. Which came first:

721        the lizard or the egg? robustness in phylogenetic reconstruction of ancestral states.

722        Journal of Experimental Zoology Part B: Molecular and Developmental Evolution

723        324:504–516.

724    Zou, Z. and J. Zhang. 2016. Morphological and molecular convergences in mammalian

725    phylogenetics. Nature Communications 7:12758 EP –.