# Systematic Biology

OXFORD UNIVERSITY PRESS

## Morphological Datasets Fit a Common Mechanism Much More Poorly than DNA Sequences and Call Into Question the Mkv Model

SCHOLARONE™
Manuscripts

1    **RUNNING-TITLE: Morphological datasets reject Mkv model**

2

3    **Morphological Datasets Fit a Common Mechanism Much More Poorly than DNA**

4    **Sequences and Call Into Question the Mkv Model**

5

6                Pablo A. Goloboff[1*], Michael Pittman[2], Diego Pol[3], and Xing Xu[4]

7    [1] Unidad Ejecutora Lillo (UEL), Consejo Nacional de Investigaciones Científicas y

8    Técnicas (CONICET), S.M. Tucumán, Argentina. E-mail: pablogolo@yahoo.com.ar.

9    [2] Vertebrate Palaeontology Laboratory, Department of Earth Sciences, University of Hong

10   Kong, Pokfulam, Hong Kong

11   [3] Museo Egidio Feruglio, Consejo Nacional de Investigaciones Científicas y Técnicas

12   (CONICET), Trelew, Argentina

13   [4] Key Laboratory of Vertebrate Evolution and Human Origins, Institute of Vertebrate

14   Paleontology and Paleoanthropology, Chinese Academy of Sciences, Beijing, China

15   * Corresponding Author

16

17

18  ABSTRACT

19  The Mkv evolutionary model, based on minor modifications to models of molecular

20  evolution, is being increasingly used to infer phylogenies from discrete morphological data,

21  often producing different results from parsimony.  The critical difference between Mkv and

22  parsimony is the assumption of a "common mechanism" in the Mkv model, with branch

23  lengths determining that probability of change for all characters increases or decreases at

24  the same tree branches by the same exponential factor.  We evaluate whether the

25  assumption of a common mechanism applies to morphology, by testing the implicit

26  prediction that branch lengths calculated from different subsets of characters will be

27  significantly correlated.  Our analysis shows that DNA (38 datasets tested) is often

28  compatible with a common mechanism, but morphology (86 datasets tested) generally is

29  not, showing very disparate branch lengths for different character partitions.  The low

30  levels of branch length correlation demonstrated for morphology (fitting models without a

31  common mechanism) suggest that the Mkv model is too unrealistic and inadequate for the
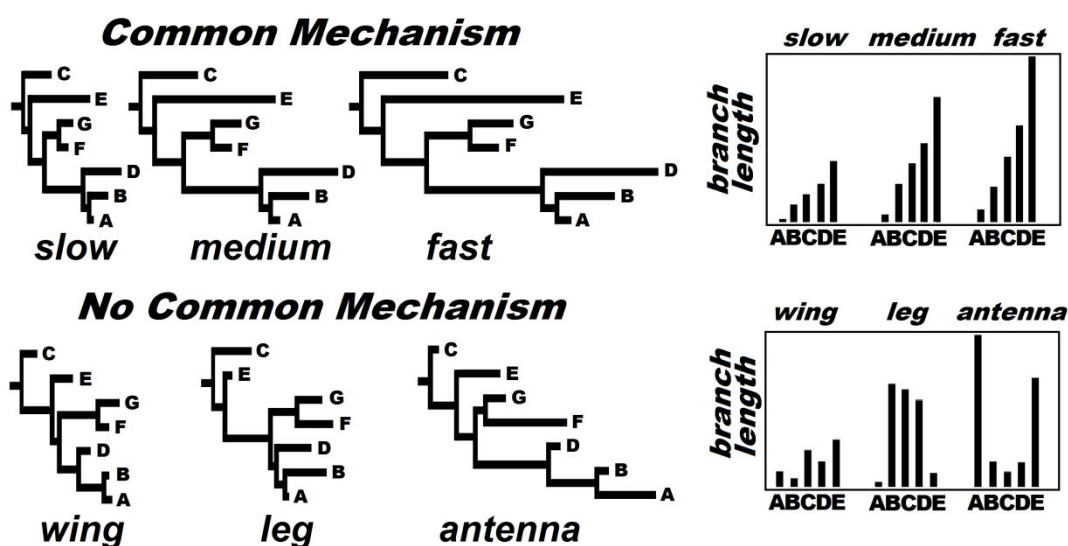
32  analysis of morphological datasets.

33

34  KEYWORDS: phylogenetics, Bayesian analysis, morphological data, Mkv model

35

36     Discrete morphological characters, despite the predominance of molecular datasets,

37     continue playing an important role in inferring phylogenetic trees (e.g. as the sole source of

38     evidence for most fossil taxa). Parsimony (implemented in PAUP*, Swofford 2002, or

39     TNT, Goloboff and Catalano 2016) is widely used for morphological data. The Mkv model

40     (Lewis 2001), based on minor modifications to models of molecular evolution, is being

41     increasingly used for phylogenetic inference (Wright and Hillis 2014, O'Reilly et al. 2016,

42     Puttick et al. 2017), even when it is often acknowledged that morphology and molecules

43     may evolve in very different ways (e.g. Lee 2016, Zhang 2018). The Mkv model is

44     implemented in several major phylogeny programs such as PAUP*, MrBayes (Ronquist et

45     al. 2012), or RAxML (Stamatakis 2014). The Mkv model critically differs from parsimony

46     in assuming a "common mechanism" (CM, Tuffley and Steel 1997), in which the

47     probability of change in different tree branches varies simultaneously for all characters,

48     exponentially depending on the "length" of the branch (expected number of changes per

49     character, the product of time and instantaneous rate, both affecting all characters equally;

50     for details, see Swofford et al. 1996, Felsenstein 2004). This assumption of a CM is in fact

51     what Lewis (2001: 915-916) considered that systematists would likely find most

52     unrealistic. Eliminating this commonality assumption causes parsimony and likelihood to

53     select the same tree (e.g. with the "no-common-mechanism" model of Tuffley and Steel

54     1997, Steel 2013; NCM); if the data have indeed not evolved with common branch lengths

55     (e.g. with heterotachy), parsimony may produce better results than model-based methods

56     that assume homogeneity (Kolaczkowski and Thornton 2004, Goloboff et al. 2017).

57         Although there have been no empirical comparisons of molecules and morphology

58     in terms of their fit to a CM, patterns of change in discrete morphological characters seem

59   not to follow this assumption of commonality. Sets of characters highly variable in a group

60   are often almost invariable in another, where different characters become highly variable

61   instead (e.g. Farris 1983:15, Sereno 2009), suggesting that the Mkv model may be

62   inappropriate for morphological data (Goloboff and Pol 2005, Nyakatura and



63

64   **Figure 1.** Differences between common and no-common mechanisms. Under a common mechanism,
65   there can be slower and faster characters (e.g. with a gamma distribution), but branch lengths
66   (expected changes per character, product of time and instantaneous rate of change for the branch)
67   increase or decrease for all the characters together. This is shown in the barplot diagram, with five
68   branches of the tree, A–D, ordered in increasing length. Without a common mechanism, there can be
69   characters with different overall rates (e.g. wing, leg, and antenna), but the expected changes show no
70   correlation between the different characters. The branch leading to taxon A is intermediate in the first
71   character, shortest in the second, and longest in the third, while the branch leading to taxon B is
72   shortest in the first character, longest in the second, and intermediate in the third, and the shapes of
73   length distributions vary for the three characters.

74

75   Bininda-Emonds 2012). With methods like the discretized gamma distribution (see details

76   in Felsenstein 2004), the Mkv model allows for rate heterogeneity among characters, but

77   this still assumes that the expected changes per character increase or decrease, together, for

78   faster and slower evolving characters, along the same branches of the tree, as illustrated in

79   Figure 1. The patterns of change in morphological characters would seem instead to depart

80　strongly from that CM, both at the level of character partitions, and individual characters.

81　Multiple (unlinked) partitions (Duchene et al. 2014, Lanfear et al. 2017) allow expected

82　changes per character at a branch to change separately in each partition, but are rarely used

83　in morphological datasets and continue requiring both the commonality assumption within

84　each partition and a prior identification of the correct partitions.

85　　　　The present study evaluates, for the first time, the assumption of a CM for

86　morphological datasets.  Bayesian model selection has been applied in some studies to

87　evaluate differences between morphological partitions, but only to assess among-character

88　rate variation (e.g Harrison and Larsson 2015), or the fit of different partitions to alternative

89　rate parameters (Lanfear et al. 2017, Clarke and Middleton 2008), instead of critically

90　evaluating the adequacy of a CM.  **Model selection may be problematic when both**

91　**models compared are incorrect (Yang and Zhu 2018), and can only be used to**

92　**compare two alternative models (instead of testing whether a single model has an**

93　**acceptable fit).  The latter becomes particularly difficult when the alternatives to CM**

94　**are to be sought among phylogenetic methods approaching parsimony:  Tuffley and**

95　**Steel's (1997) NCM is equivalent to parsimony, but (as noted by Holder et al. 2010:**

96　**478; see also Sober 2004) NCM is too highly parameterized to be ever selected, and**

97　**probably not the only way to characterize parsimony –yet no currently available**

98　**implementation emulates parsimony methods with fewer parameters, to enable a**

99　**more meaningful comparison of likelihoods.**

100　　　　Given those difficulties, we use here an approach based on statistical hypothesis

101　testing, to assess the adequacy of the Mkv model for morphological datasets. The CM of

102　the Mkv model predicts that branch lengths for different subsets of data will be correlated,

103 and our test is based on evaluating whether that prediction is met in empirical datasets. The

104 paper begins by outlining the test and its justification, then applies it to morphology and

105 DNA sequences, first to partitions predefined on the basis of contiguity (DNA) or anatomy

106 (morphology), then to randomly defined subpartitions (within predefined partitions, and

107 whole datasets). As these tests show that the vast majority of morphological datasets do

108 not conform to a CM, we then apply a similar test to evaluate the alternative: whether the

109 degrees of correlation between branch lengths in morphological datasets could have been

110 produced by models without a CM. These tests reject a pure NCM, but an alternative

111 model for generating datasets without a CM (which we call the *episodic* model, with

112 character changes restricted to certain parts of the tree; see below) produces a correlation

113 between numbers of character changes in each partition that is well within the values

114 observed in real morphological datasets. Finally, we examine the relative performance of

115 phylogenetic methods on datasets simulated under the episodic model, and show that

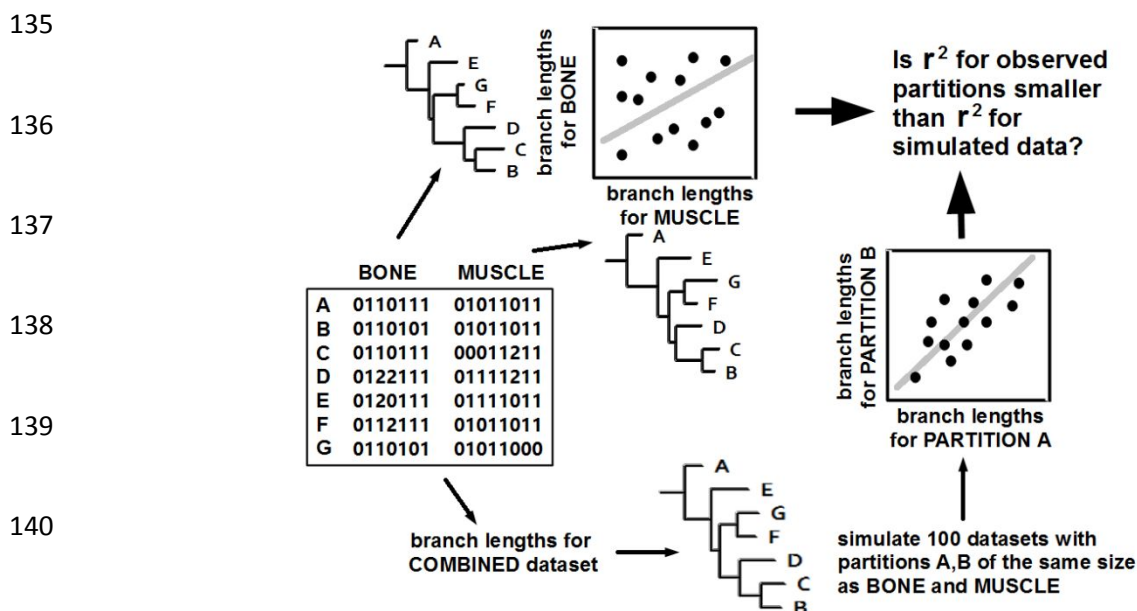116 parsimony tends to perform on par or better than Bayesian analysis.

117

118 METHODS AND MATERIALS

119 *Datasets.*– Source and details in the Supplementary Material. A total of 86 morphological

120 datasets was used, with 26–188 taxa, and 80–4541 characters. For 8 of the morphological

121 datasets, it was possible to define partitions on the basis of anatomy (with 2–8 partitions per

122 dataset, 40–1451 characters per partition). For sequences, a total of 38 datasets for 35

123 different genes, with 60–500 taxa and 305–2218 characters, was examined. These

124 molecular datasets were prepartitioned in 2–5 partitions of 100–500 contiguous positions,

125 depending on the length of the sequences (only one case, with very short sequences, used

126    two partitions of 50 positions).  These molecular partitions were created leaving out the

127    initial 100 positions (which, due to alignment, often contain large proportions of missing

128    entries), except in the shortest sequences (where partitions started at position #50).

129    *Branch length tests.–*  Branches shorter for one partition and longer for the other are

130    evidence of at least heterotachy (or, at most, the complete absence of a CM).  The strength

131    of the observed correlation can be measured with the $r^2$ statistic; the observed $r^2$ was then

132    compared with that for partitions of the same size generated on a model tree with the same

133    branch lengths as the combined dataset; if observed $r^2$ is matched with a low probability,

134    then the CM of the Mkv model can be confidently rejected.  Figure 2 displays the

135

136

137

138

139

140



141
142
143
144    **Figure 2.** General scheme of the test to evaluate significance of heterogeneity between branch lengths
145        for different partitions.
146

147    procedure for testing branch length homogeneity in two partitions.   The only similar

148    evaluation of which we are aware is that of Clarke and Middleton (2008), who compared

149  branch lengths for different morphological partitions; however, they did not evaluate the

150  significance of the differences in branch lengths by reference to a specific model of

151  evolution.

152      Even when the data have been generated by a model with a CM (e.g. Mkv or JC69)

153  the expected homogeneity in branch lengths for the simulated partitions will depend on

154  both the branch lengths of the model tree, as well as the numbers of characters in the two

155  partitions being compared.  To the extent that the branch lengths of the model tree are more

156  dissimilar, the correlation between branch lengths for two sets of characters generated on

157  the same model tree will be stronger; when all branch lengths of the model tree are

158  identical, character changes can be located equiprobably on any tree branch, resulting in

159  very low correlation.  On the other hand, to the extent that there are more characters in the

160  partitions, branch lengths will more accurately converge to the values in the model tree,

161  thus increasing the correlation between the branch lengths for both partitions.  Therefore, a

162  proper test cannot be based solely on the observed value of $r^2$ for the correlation between

163  branch lengths for two partitions: the values of $r^2$ must be compared against the values

164  expected under the specific situation being tested, i.e. using the same numbers of characters

165  of the observed partitions, and a model tree with the same branch lengths as the combined

166  dataset.

167      For completeness, most of the tests were repeated calculating branch lengths with

168  most parsimonious reconstructions (MPR).  In this case, the scripts calculated branch

169  lengths simply as the number of characters in the partition unambiguously changing along

170  the branch, divided by the total number of characters in the partition.

171    *Calculation of Branch Lengths.–* Branch lengths for the results reported were calculated

172    using maximum likelihood, unless noted otherwise.  Taxa with missing entries for all

173    characters in one (or both) partition(s) were pruned from the tree, and the branch lengths

174    were calculated on the resulting reduced tree.  This was necessary only in few comparisons.

175    TNT scripts (Goloboff et al. 2008) automatically created Nexus files and called PAUP*

176    with commands to calculate and save branch lengths in Newick format, then reading back

177    the branch lengths into TNT, for further processing.  For morphological datasets, invariant

178    characters were excluded (for different pairwise comparisons between partitions, some of

179    the variable characters in a partition could become invariant if some taxa with only missing

180    entries in the other partition are deactivated). For morphological data, branch lengths were

181    calculated with default PAUP* options (in the absence of invariant characters, PAUP*

182    defaults to the Mkv model, estimating the proportion of invariant characters automatically).

183    For sequence data, the simplest model (JC69, Jukes and Cantor 1969) was invoked, with

184    *lset nst=1 rates=equal basefr=equal*, which is the closest equivalent to the Mkv model

185    (except for the estimation of invariant characters, which has a minimum effect on branch

186    length proportionality).  Invoking more complex DNA models and adding more parameters

187    to be estimated seemed unnecessary, given that the goal of the analysis is only evaluating

188    the heterogeneity in branch lengths for different partitions.

189    *Model tree.–* The datasets for calculating the statistical distribution of the correlation

190    between partitions with the same numbers of characters as the observed partitions were

191    simulated using the observed tree as model.  The "observed" tree is the published tree,

192    when available, or a most parsimonious tree for the combined dataset otherwise (in the case

193    of phylogenomic datasets, this is a tree for the dataset combining all the genes). We did
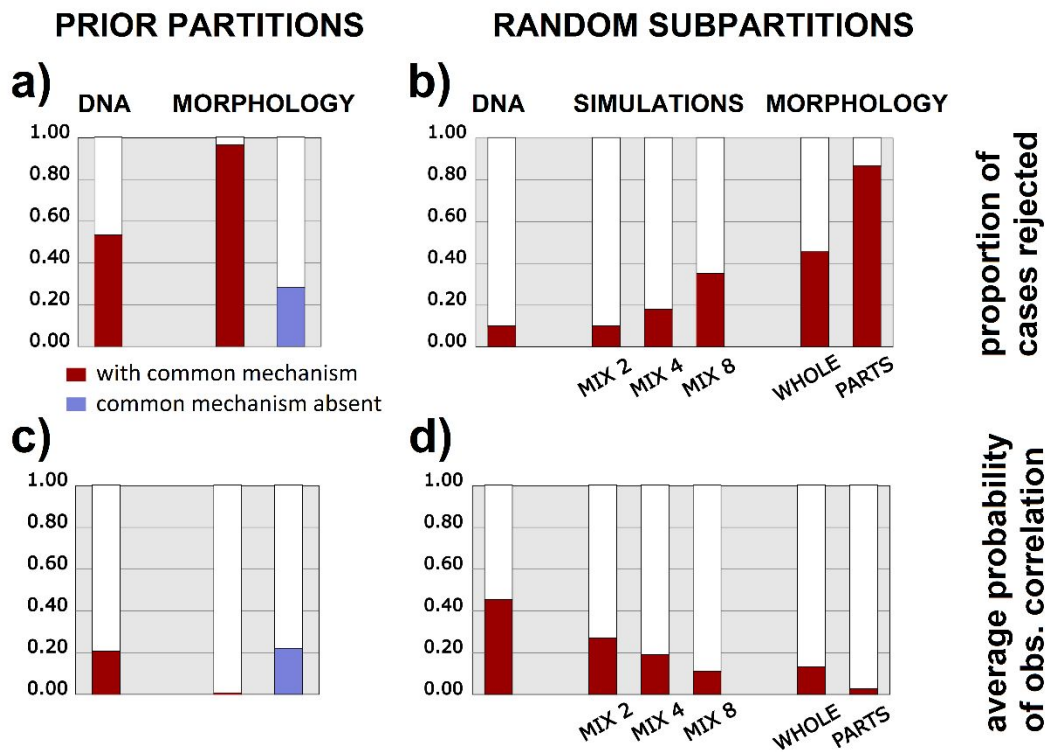
194    experiments to confirm that the test does not strongly depend on the topology of the tree

195    used to calculate branch lengths, so even if the observed tree is slightly different from the

196    correct phylogeny, the results of the test continue being valid (See Supplementary

197    Material). This makes the test radically different from "empirical" comparisons where real

198    datasets are analyzed with different methods of phylogenetic inference and the resulting

199    groupings are evaluated on whether they agree with groupings presumed to be correct prior

200    to the analysis (e.g. Puttick et al.'s 2017 discussion of results for 4 empirical datasets). No

201    presumption of prior knowledge is needed for the present correlation test, which considers

202    only the fit of the model to the dataset, not the accuracy of the trees produced by assuming

203    the model.

204

205    COMMON MODEL TESTED BETWEEN PRE-DEFINED PARTITIONS

206    We first tested 8 large published matrices, containing partitions corresponding to

207    anatomical regions or organ systems with numerous characters (40 characters per partition

208    was considered as the minimum for appropriate testing). Given the different numbers of

209    partitions per dataset, a total of 79 pairwise comparisons were possible. The vast majority

210    of these partitions (Figs. 3, 4) have much more pronounced differences in numbers of

211    character changes along branches than expected under the Mkv model (only 3.8% of cases

212    fail to reject the Mkv model as null model with $\alpha=0.01$; Fig. 3a). The results of a similar

213    test performed on DNA sequences (38 datasets for 35 different genes, with partitions

214    defined by contiguity, 66 possible comparisons) are very different, with a common

215    mechanism accepted for 42.4% of comparisons (Fig. 3a), over ten times more frequently

216    than for morphology.  The average probabilities of observed r$^2$ values under a CM are also

217    much higher for DNA than for morphological datasets (Fig. 3c, 4). Therefore, branch

218    lengths for partitions of DNA sequences are clearly much less heterogeneous than for

219    morphological data.



220

221    **Figure 3.** Proportion of cases where different models are rejected with α=0.01 by branch length tests
222    (a, b), and average probabilities of observed correlation (c, d). Prior partitions (a, c) correspond to
223    characters grouped on the basis of anatomy in the case of morphology, and on the basis of contiguity
224    in the case of DNA. For morphology, the common mechanism model is Mkv; for DNA, its closest
225    equivalent, JC69. The model without a common mechanism is the Episodic model described in the
226    text. Random subpartitionings (b, d) for DNA were tested on a mid-sequence group of positions, on
227    whole datasets simulated with mixtures of 2–8 independent sets of branch lengths (MIX), on the
228    partitions predefined on the basis of anatomy (PARTS), and on whole datasets when no anatomical
229    partitions could be predefined (WHOLE).
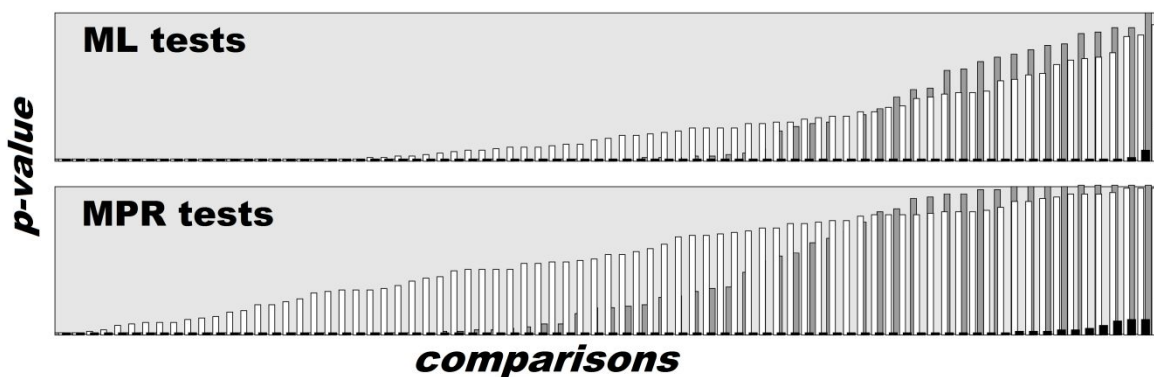
230

231        Our tests evaluate multiple instances, and not all comparisons are fully independent,

232    because some of those imply combinations of partitions.  Appropriate corrections for

233    confidence levels on individual cases would have required making prohibitively time

234    consuming simulations (i.e. with many more replications per test).  Our interest, however,

235    is not in the significance of individual comparisons, but rather in the collective results, and

236    the differences between morphological and molecular datasets.  While corrections for

237    multiple tests might have lowered somewhat the rejection rate of homogeneity, a correction

238    would equally affect the comparisons for morphology and sequences, so that the

239    differences between the degree to which branch length homogeneity is, or is not rejected,

240    by each type of dataset, would have remained equally strong.

241         An important caveat of the test performed here is that the observed branch lengths

242    were calculated using a single rate category (i.e. no gamma parameter).  There are

243    indications (e.g. Marshall et al. 2006; Nguyen et al. 2017) that taking into account among-

244    site rate heterogeneity improves estimations of branch lengths.  A more accurate appraisal

245    would perhaps have analyzed both the observed and simulated datasets allowing rate

246    heterogeneity, simulating data under the same gamma values estimated for the combined

247    dataset (instead of the single rate now used); this would have made evaluations

248    significantly slower, would have required modifications to the functions of TNT that

249    simulate data under a CM, and would have added another layer of complexity (and thus,

250    potential errors) to the estimations.  It seems doubtful, however, that using a gamma

251    correction would have changed much the evaluations.  The test focuses on correlations

252    between branch lengths for two partitions, and the main numerical effect of applying a

253    gamma correction is to alter the absolute values of all branch lengths by roughly the same

254    factor, with only minor modifications to their proportionality.  This is indeed a problem

255    when the interest is in calculating the correct values of branch lengths for each partition

256    (e.g. as in the study of Nguyen et al. 2017), but does not have a strong effect on the values

257   of correlation (changing only the regression slopes).  The best indication that the use of a

258   single rate category did not bias the comparisons in the case of morphology is in the results

259   for DNA sequence data: those analyses did not use, either, a gamma parameter for among-

260   site rate variation, yet they produced a high proportion of cases where correlation between

261   estimated branch lengths was within the range expected under the single-rate model.  This

262   suggests that the effect of a test considering among-site rate variation would have been

263   minor, and that the same differences between DNA and morphological datasets would have

264   been obtained.

265       The results obtained when comparing branch lengths for the partitions calculated

266   with MPR are, overall, similar to those obtained with likelihood, with the same difference

267   between DNA and morphological datasets.  The probability of obtaining the observed

268   correlation for the morphological datasets (p-values under the episodic and Mkv models),

269   and for DNA sequences (p-values under JC69) is shown in Figure 4, for each individual

270   comparison.  The similarity in results obtained using two methods as different as MPR and

271   maximum likelihood also suggests that the rejection of a CM in morphology does not

272   strongly depend on method used for calculating branch lengths (including the use of a
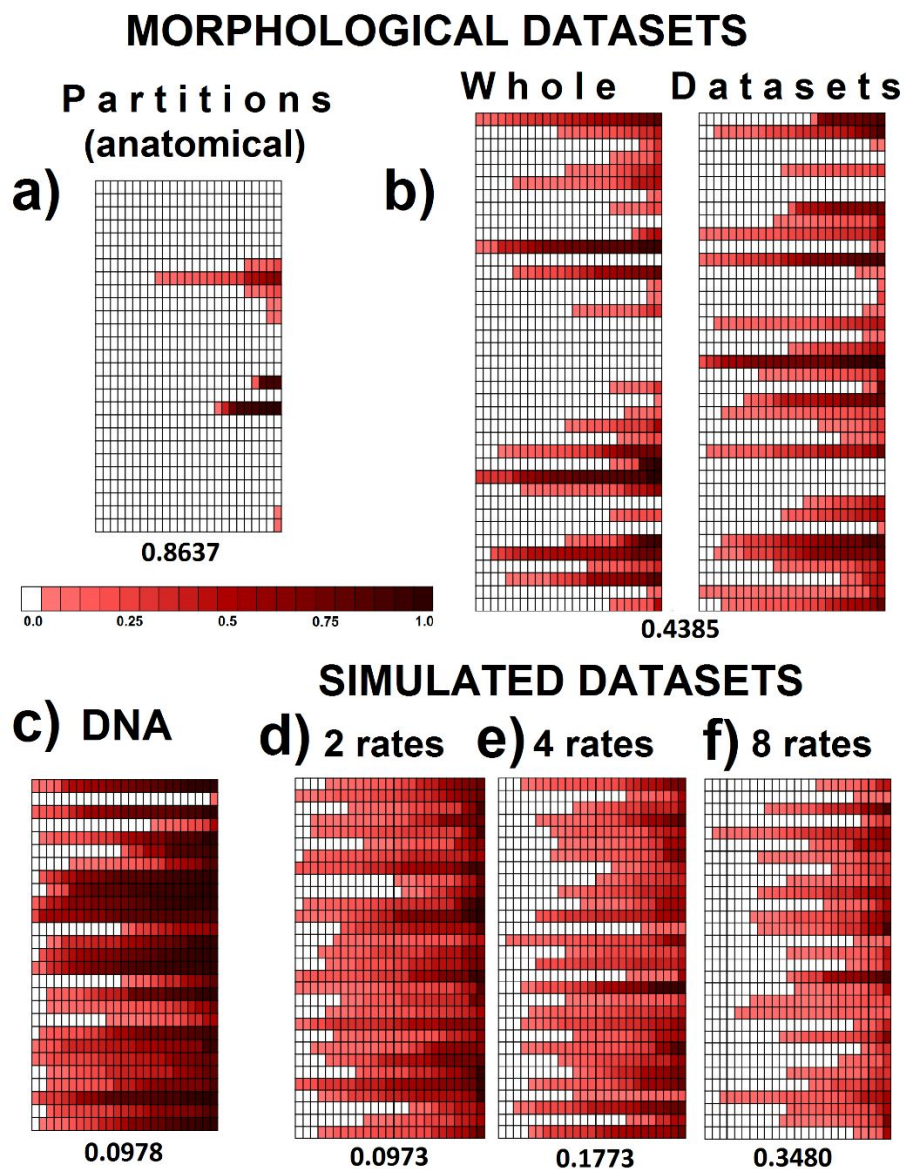
273   gamma parameter).

274

**Figure 4.** Plot showing the probability (P-values) of obtaining a correlation between branch lengths for two partitions as strong as the observed one, under different models for generating data, for morphological partitions defined by anatomy, and DNA partitions defined by contiguity, arranged in increasing order of P-value. Probabilities calculated both with likelihood (ML) and most parsimonious reconstructions (MPR). Black, morphological partitions tested against Mkv model; gray, DNA, tested against JC69; white, morphological partitions tested against episodic model. Gray and black bars are models with a common mechanism, white bars are for a model lacking a common mechanism.

COMMON MODEL TESTED WITHIN PARTITIONS AND ENTIRE DATASETS

Some studies have already demonstrated (with different methods; Clarke and

Middleton 2008, Tarasov and Genier 2015, Lee 2016) heterogeneity in branch lengths for

predefined partitions, so a meaningful evaluation must test whether a CM is in effect *within*

individual partitions. Two subpartitions containing similar proportions of characters

evolving under two completely different sets of branch lengths will have similar mixtures

of rates, combining to provide a common average "rate" for each branch (Kolaczkowski

and Thornton 2004), similar for both subpartitions. The internal heterogeneity of such

mixtures cannot be detected by the present test (or any test we know), unless the correct

partitions are known in advance –seldom the case for morphological data. Some

partitioning schemes will produce the opposite effect, of making datasets generated from a

single set of branch lengths to appear heterogeneous (e.g. by separating the characters in

296    two groups, depending on which half of the tree they have more changes), but those

297    partitionings are unlikely to be obtained at random.  Thus, a conservative test can randomly

298    subpartition characters, comparing the degree of branch length correlation between the

299    random subpartitions with that expected under a CM; mixtures with similar proportions of

300    two (or a few) sets of distinct branch lengths will often appear relatively homogeneous

301    under such a test, for the mixtures will be sampled in roughly similar proportions.

302    Therefore, rejection of branch length homogeneity in a majority of randomly chosen

303    subpartitions is especially meaningful: more than just a few alternative rates, such a result

304    suggests the absence of a CM altogether.

305         For testing random subpartitions, only the partitions with 80 or more characters

306    were considered, dividing in two evenly-sized subpartitions.  The results for random

307    subpartitions are summarized in Figures 3b, 3d.  Figure 5 shows the results of testing each

308    subpartition individually; Figure 6 shows the average results for all the subpartitions of

309    each partition (or dataset).  For 86.4% of cases, random subpartitionings of the

310    anatomically defined partitions produced a heterogeneity beyond ($\alpha$=0.01) expected under

311    the CM of the Mkv model (white boxes in Fig. 5a).  Given that the test based on random

312    subpartitioning requires no prior definition of partitions, an additional set of 78

313    morphological datasets (for which partitions could not be easily defined on the basis of

314    anatomy) were tested as a whole.  A CM was rejected ($\alpha$=0.01) in 43.8% of all

315    subpartitions (white boxes in Fig. 5b).  For molecular datasets, instead, random

316    subpartitionings (for 200 mid-sequence positions) reject a CM in only 9.8% of cases (Fig.

317    5c).  To give these results further context, we simulated datasets (200 characters) under a

318    Mkv model but with independent sets of branch lengths; as expected, the proportion of

# MORPHOLOGICAL DATASETS

### Partitions
### (anatomical)

### Whole   Datasets



0.8637

0.4385

# SIMULATED DATASETS

### c) DNA

### d) 2 rates   e) 4 rates   f) 8 rates

0.0978

0.0973

0.1773

0.3480

319

320 **Figure 5.** Plots of subpartition tests (25 per partition/dataset). Every row corresponds to a dataset or
321 partition, every **individual** box corresponds to a subpartition. **The color of each box indicates the**
322 **probability of obtaining the branch length correlation in the subpartition under a common**
323 **mechanism (white color is p < 0.05, with a darker color as p increases).  The numbers  below**
324 **frames correspond to proportion of subpartitions where common mechanism is rejected with**
325 **α=0.01 (i.e. lower proportions correspond to cases where the common mechanism is less likely to**
326 **have generated the data). (a) Subpartitions of partitions predefined on the basis of anatomy; (b)**
327 **Whole datasets; (c) Molecular datasets; (e-f) Datasets simulated with 2, 4 and 8 independent sets**
328 **of branch lengths.**

329

**Figure 6.** Comparison of proportion of subpartitions with a common mechanism rejected (at α=0.01) per partition (or dataset), for morphology and DNA.

cases where a single CM could be rejected on random subpartitions increased with the

number of independent sets of branch lengths (Figs. 5d–f), reaching up to 34.8% for

mixtures of 8 independent sets (Fig. 5f).  This is still well below the rejection rate for

morphological datasets, suggesting that (on average) morphological characters evolved

with even larger deviations from a single CM.


TESTING MODELS WITHOUT A COMMON MECHANISM

The homogeneity of branch lengths for DNA sequences can be expected from

theoretical considerations and previous empirical work evaluating the CM in sequences

(Huelsenbeck et al. 2008).  The results for morphological datasets, in contrast, strongly

refute the CM (and hence the Mkv model), both in datasets taken as a whole, and most

importantly, within partitions defined on the basis of anatomy.  Whether the data evolve

under a CM is indeed relevant for phylogenetic inference based on morphology:

347    simulations show that Bayesian inference works best when the data evolve homogeneously

348    (Wright and Hillis 2014, O'Reilly et al. 2016, Puttick et al. 2017), but parsimony may work

349    best (Kolaczkowski and Thornton 2004, Goloboff et al. 2017) when they do not.

350          Note that the probability distribution of character patterns of both NCM and the

351    "Ultra-Conserved-Mechanism" (UCM, with a CM and all tree-branches having exactly the

352    same length for all characters) are exactly identical, as shown by Huelsenbeck et al. (2008)

353    and Steel (2011).  Either of those models will have any change equiprobably located (CEL)

354    on any tree branch, which is how Goloboff et al. (2017) generated their data.  Given that

355    parsimony is an appropriate method if the data do evolve under NCM (as shown by Tuffley

356    and Steel 1997, Steel 2011), it follows that so it is under the equivalent (but less strongly

357    parameterized) UCM or CEL, which generate the same probability distributions for

358    character patterns.  Note that CEL is a statement of the product of evolution (i.e. on how

359    character changes will be located on tree branches), more than a statement of process; this

360    product may be achieved by different processes (NCM, UCM, and possibly others).

361          The equiprobability of location of changes on any tree branch in the simulations of

362    Goloboff et al. (2017) is a uniform distribution which (given the difficulties in modelling

363    morphology) can be defended as an initial reference assumption, and produces no branch

364    length correlation between partitions.  Parsimony is then a well-justified method, but the

365    model is also rejected by morphological datasets: branch length correlation between

366    partitions is higher than expected (with 41 of the 79 comparisons between predefined

367    partitions rejecting the model with α=0.05). The generating model can be made more

368    realistic with characters equiprobably changing in every branch but only within a certain

369    region of the tree, thus following the mosaicism proposed by Farris (1983: 15) and

370    Goloboff et al. (2018).  This model (Fig. 7) assumes that, during evolution, the possibility



371

372    **Figure 7.** Episodic model.  Colors indicate regions of the tree where a character (or group of
373    characters) can change; black branches indicate regions where characters cannot change.  Within the
374    colored region, a change has the same probability of being located in any of the branches. The
375    example shows 4 pivots (i.e. points where change becomes possible or impossible); because of the
376    interaction between pivots, different branches of the tree have different numbers of characters
377    (indicated on the rightmost tree) that could possibly change.

378

379    that formerly invariable characters become variable (or viceversa) can be triggered in

380    *episodic* events; the point at which the character becomes variable (or invariable) is a node

381    in the tree acting as a *pivot*.  Morphological characters, by its very hierarchical nature

382    (Maddison 1993, De Laet 2005, Brazeau et al. 2017) and by being subject to selection

383    shifting from stabilizing to directional along time or changes in the developmental

384    constraints, may well be liable to such episodic evolution.  This *episodic* model is

385    reminiscent of the covarion model (Fitch and Markowitz 1970, and successive

386    modifications), differing in that character changes within regions of variability can be

387    equiprobably located at any possible branch, thus lacking a CM and a formal branch-length

388    parameter, consequently being more suited for morphological data.  In the presence of

389    multiple pivots affecting groups of characters, some branches will have larger numbers of

390    synapomorphies for each group, generating a correlation between branch lengths for

391    different partitions, and the correlation observed for empirical partitions is mostly within

392    that expected from the episodic model (with only 27.8% of comparisons rejecting the

393    model when half the characters are affected by pivots and half are not; see Figs. 3a, 4).

394    This does not prove that an episodic model is the best general explanation for

395    morphological patterns of character change, but at least the model is not as widely rejected

396    as the Mkv or NCM models.  More interestingly theoretically, the model shows that trees

397    with some correlation between changes per branch for different partitions can result from

398    models that do not assume a CM.  The episodic model is used here solely to generate data,

399    not to infer trees; likelihood inference assuming that model has not been implemented, and

400    (by analogy to the covarion model) may suffer from identifiability problems (as noted by

401    Gruenheit et al. 2008 for standard covarion models) unless significant restrictions are

402    imposed.

403

404    IMPLICATIONS FOR CHOICE OF PHYLOGENETIC METHODS

405        The episodic model resembles the NCM, UCM, or any other process leading to

406    CEL, except that change is restricted to some parts of the tree.  Given this similarity, we

407    conjecture that only multiple pivots per character could produce inconsistency for

408    parsimony if the model truly generated the data.  In other words, if only one pivot per

409    character occurs in the tree, parsimony can be justified just like the model with changes in

410    each character occurring equiprobably over all the tree.  With a single pivot per character,

411    several tree branches may have more changes by virtue of being intermediate between

412    pivots, but each of those long branches would have changes in different groups of

413    characters (just like the synapomorphies for the long branches leading to e.g. Cetacea and

414    Chiroptera correspond to different characters; Goloboff et al. 2018), so that they would be

415    unlikely to attract.

416        The fact that models with a CM are strongly rejected by morphological data, and

417    some models without a CM are not, is relevant for the choice of phylogenetic method.

418    Previous studies where Bayesian analysis outperformed parsimony (e.g. Wright and Hillis

419    2014, O'Reilly et al. 2016, Puttick et al. 2017) had generated their data with a CM.  In

420    addition, for implied weighting (Goloboff 1993), O'Reilly et al. (2016) and Puttick et al.

421    (2017) chose the worst concavity value (k=2, close to a clique, contrary to

422    recommendations of Goloboff 1995: 99) and did not eliminate poorly supported groups

423    (Fig. 8a).  With a milder concavity and poorly supported groups eliminated, Bayesian

424    analysis with the Mkv model outperforms implied weighting by a much smaller difference

425    (Fig. 8b), but by a difference nonetheless, when the data are generated with a CM.  When

426    the data are generated instead with the half-episodic model, which does not assume a CM,

427    parsimony tends to produce (as in the unrestricted model of Goloboff et al. 2017, and in

428    agreement with expectations) slightly better results than Bayesian analysis (see Fig. 8c). As

429    the number of characters increases (Fig. 8d), **both methods improve their results, but**

430    Bayesian analysis has a slightly poorer performance for every statistic, perhaps as a result

431    of the departure from the CM assumed by the Mkv model becoming more evident (given

432    the large amounts of data).

433

**Figure 8.** Comparison between implied weights and Bayesian analysis, using different methods for
simulating and analyzing data (columns), and four different statistics to evaluate performance (rows).
Proportional error is the number of incorrect groups found, divided by the number of groups in the
inferred tree.  The values of different statistics for Bayesian analysis are plotted against implied
weights parsimony; by plotting the values for implied weighting on the x-axis, and those for BI on the
y-axis, the deviation from the diagonal allows the difference in performance between the two methods
to be easily detected.  Datasets generated with both the Mk model of Lewis (2001) (columns A, B),
and with the half-episodic model (C, D).  Each of 100 points represents the average of 10 simulations
with the same numbers of taxa and characters (to reduce dispersion, for a total of 1,000 simulated
datasets). As the datasets are generated with the half-episodic model (lacking a common mechanism),
the number of characters increases, and a concavity value of k=12 is used for implied weighting
(instead of k=2, the worst performing value, chosen by O'Reilly et al. 2016 and Puttick et al. 2017 for
their comparisons), parsimony outperforms Bayesian analysis **by a smaller margin, but more
consistently**.  The average values for each statistic are indicated in the x-axis for implied weighting,
and on the y-axis for Bayesian analysis.

## CONCLUSIONS

451        Our findings provide the first empirical demonstration, in a phylogenetic

452   framework, of the differences in modes of evolution of molecules and morphology.  While

453   models that lack a CM (such as the episodic model) can produce degrees of branch length

454   correlation between partitions that are in line with those observed in real datasets, the CM

455    assumed by the Mkv model is strongly rejected by the morphological datasets.  Of course,

456    as generally acknowledged, a model need not reflect reality perfectly to be a useful aid in

457    estimation, but a model still needs to have *some* basis in reality. If it is accepted that "all

458    models are wrong, but some are useful", then one must also accept that some models are

459    *not* useful.  The extent to which the CM assumed by the Mkv model deviates from reality

460    seems strong enough to suspect the model may well do more harm than good. It is possible

461    that violations of its assumptions rarely mislead Bayesian inference of trees in practice**; our**

462    **simulations show that MrBayes seems rather robust to such violations.  Such**

463    **robustness may well be a result of the mechanics of the Markov chain and subsequent**

464    **tree summarization, more than the result of assuming the Mkv model.  If this is**

465    **correct, MrBayes with the "parsimony" model might well produce (for datasets**

466    **generated without a CM) trees of about the same quality as those produced with the**

467    **Mkv model (a possibility that has not hitherto been examined in detail).**  But one of the

468    advantages claimed for model-based methods is that (by incorporating biological

469    knowledge about evolutionary processes; Huelsenbeck et al. 2011) they allow estimating

470    more than just tree topologies.  Unrealistic assumptions built into phylogenetic models,

471    therefore, can also affect studies of character mapping, dating of nodes on given trees,

472    calculation of probabilities of specific evolutionary events, and even how taxonomists think

473    of characters or diagnose groups. Thus, in light of the evidence against the common

474    mechanism assumption, we strongly advise against the uncritical use of the Mkv model.

475 SUPPLEMENTARY MATERIAL

476 Material and methods, datasets, results, and scripts are available at the Dryad repository,

477 doi:10.5061/dryad.3680n0c.

478

499     REFERENCES

500     Brazeau, M., Guillerme, T., Smith, M. 2017. Morphological phylogenetic analysis with

501     inapplicable data. BioRxiv https://doi.org/10.1101/209775 [Note: this is under review

502     in Syst.Bio, so it will likely be citable as a paper by the time our paper comes out]

503     Clarke, J., Middleton, K. 2008. Mosaicism, Modules, and the evolution of birds: results

504     from a Bayesian approach to the study of morphological evolution using discrete

505     character data. *Syst. Biol.* **57**, 185–201.

506     De Laet, J. 2005. Parsimony and the problem of inapplicables in sequence data. In:

507     *Parsimony, Phylogeny and Genomics (*V. Albert, ed.), Oxford University Press, pp 81–

508     116.

509     Duchene, S., Molak, M., Ho, S. 2014. ClockstaR: choosing the number of relaxed-clock

510     models in molecular phylogenetic analysis. *Bioinformatics* **30**, 1017–1019.

511     Farris, J. 1983. The logical basis of phylogenetic analysis. In: Platnick, N., Funk, V.

512     (Eds.), *Advances in Cladistics II*. Columbia University Press, New York, NY, pp. 7–36.

513     Felsenstein, J. 2004. *Inferring Phylogenies*. Sinauer Associates, Sunderland, MA.

514     Fitch, W., Markowitz, E. 1970. An improved method for determining codon variability

515     in a gene and its application to the rate of fixation of mutations in evolution. *Biochem.*

516     *Genet.* **4**, 579–593.

517     Goloboff, P. 1993. Estimating character weights during tree search. *Cladistics* **9**, 83–

518     91.

519    Goloboff, P. 1995.  Parsimony and weighting: a reply to Turner and Zandee.  *Cladistics*

520    **11**, 91–104.

521    Goloboff, P., Pol, D. 2005.  Parsimony and Bayesian phylogenetics.  In: *Parsimony,*

522    *phylogeny, and genomics* (Victor Albert, ed.).  Oxford University Press, pp. 148-159.

523    Goloboff, P., Farris, J.  2008. Nixon, K. TNT, a free program for phylogenetic analysis.

524    *Cladistics* **24**, 774–786.

525    Goloboff, P., Catalano, S. 2016. TNT version 1.5, including a full implementation of

526    geometric morphometrics.  *Cladistics* **32**, 221–238.

527    Goloboff, P., Torres, A., Arias, S. 2017. Weighted parsimony outperforms other

528    methods of phylogenetic inference under models appropriate for morphology.

529    *Cladistics*, https://doi.org/10.1111/cla.12205.

530    Goloboff, P., Torres, A., Arias, S. 2018.  Parsimony and model-based phylogenetic

531    methods for morphological data:  comments on O'Reilly et al. (2017).  *Palaeontology*

532    doi.org/10.1111/pala.12353.

533    Gruenheit, N., Lockhart, P., Steel, M., Martin, W. 2008.  Difficulties in testing for

534    covarion-like properties of sequences under the confounding influence of changing

535    proportions of variable sites. *Mol. Biol. Evol.* **25**, 1512–1520.

536    Harrison, L., Larsson, H. 2015. Among-character rate variation distributions in

537    phylogenetic analysis of discrete morphological characters.  *Syst. Biol.* **64**, 307–324.

538    **Holder, M., Lewis, P., Swofford, D. 2010. The Akaike Information Criterion will**

539    **not choose the No Common Mechanism model.  *Syst. Biol.* 59, 477–485.**

540  Huelsenbeck, J., Annè, C., Larget, B., Ronquist, F. 2008. A Bayesian perspective on a

541  non-parsimonious parsimony model. *Syst. Biol.* **57**, 406–419.

542  Huelsenbeck, J., Alfaro, M., Suchard, M. 2011. Biologically inspired phylogenetic

543  models strongly outperform the no common mechanism model. *Syst. Biol.* **60**, 225–

544  232.

545  Jukes, T., Cantor, C. 1969. Evolution of protein molecules. In Munro, N. (editor),

546  *Mammalian protein metabolism*. Vol. 3, New York, Academic Press, pp. 21–132.

547  Kolaczkowski, B., Thornton, J. 2004. Performance of maximum parsimony and

548  likelihood phylogenetics when evolution is heterogeneous. *Nature* **431**, 980–984.

549  Lanfear, R., Frandsen, P., Wright, A., Senfeld, T., Calcott, B. 2017. PartitionFinder 2:

550  new methods for selecting partitioned models of evolution for molecular and

551  morphological phylogenetic analyses. *Mol. Biol. Evol.* **34**, 772–773.

552  Lee, M. 2016. Multiple morphological clocks and total evidence tip-dating in

553  mammals. *Biol. Lett.* **12**, 20160033. http://dx.doi.org/10.1098/rsbl.2016.0033.

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

555  morphological character data. *Syst. Biol.* **50**, 913–925.

556  Maddison, W. 1993. Missing data versus missing characters in phylogenetic analysis.

557  *Syst. Biol.* 42, 576–581.

558  Marshall, D., Simon, C., Buckley, T. 2006. Accurate branch length estimation in

559  partitioned Bayesian analyses requires accommodation of among-partition rate variation

560  and attention to branch length priors. *Syst. Biol.* **55**, 993–1003.

561    Nguyen, L.-T., von Haeseler, A., Minh, B. 2017. Complex models of sequence

562    evolution require accurate estimators as exemplified with the invariable site plus

563    gamma model. *Syst. Biol.* **67**, 552-558.

564    Nyakatura, K., Bininda-Emonds, O.  2012. Updating the evolutionary history of

565    Carnivora (Mammalia): a new species-level supertree complete with divergence time

566    estimates.  *BMC Biol.* **10**, 1–31.

567    O'Reilly, J., et al. 2016.  Bayesian methods outperform parsimony but at the expense of

568    precision in the estimation of phylogeny from discrete morphological data. *Biol. Lett.*

569    **12**, 20160081, 1–5. https://doi.org/10.1098/rsbl.2016.0081.

570    Puttick, M., et al. 2017.  Uncertain-tree: discriminating among competing approaches to

571    the phylogenetic analysis of phenotype data. *Proc. R. Soc. B* **284**, 20162290.

572    Sereno, P. 2009.  Comparative cladistics.  *Cladistics*  **26**, 624–659.

573    Swofford, D.  2002.  *PAUP*: Phylogenetic analysis using parsimony (* and other*

574    *methods)*. Version 4. Sunderland (MA): Sinauer Associates.

575    Swofford, D., Olsen, G., Waddell, P., Hillis, D. 1996.  Phylogenetic inference. In:

576    Hillis, D., Moritz, C., Mable, B. (Eds.), *Molecular Systematics*, second ed. Sinauer,

577    Sunderland, MA, pp. 407–514.

578    Ronquist, F., et al. 2012.  MrBayes 3.2: efficient Bayesian phylogenetic inference and

579    model choice across a large model space.  *Syst. Biol.* **61**, 539–542.

580    **Sober, E. 2004. The contest between parsimony and likelihood.  *Syst. Biol.* 53, 644–**

581    **653.**

582 Stamatakis, A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-

583 analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313.

584 Steel, M. 2011. Can we avoid "sin" in the house of "no common mechanism"? *Syst.*

585 *Biol.* **60**, 96–109.

586 Tarasov, S., Génier, F. 2015. Innovative Bayesian and parsimony phylogeny of dung

587 beetles (Coleoptera, Scarabaeidae, Scarabaeinae) enhanced by ontology-based

588 partitioning of morphological characters. *PLOS One* **10**, e0116671.

589 Tuffley, C., Steel, M. 1997. Links between maximum likelihood and maximum

590 parsimony under a simple model of site substitution. *Bull. Math. Biol.* **59**, 581–607.

591 Wright, A., Hillis, D. 2014. Bayesian analysis using a simple likelihood model

592 outperforms parsimony for estimation of phylogeny from discrete morphological data.

593 *PLoS ONE* **9**, e109210. https://doi.org/10.1371/journal.pone.0109210.

594 Yang, Z., Zhu, T. 2018. Bayesian selection of misspecified models is overconfident

595 and may cause spurious posterior probabilities for phylogenetic trees. *Proc. Nat. Acad.*

596 *Sc.* **115**, 1854–1859.

597 Zhang, J. 2018. Neutral Theory and Phenotypic Evolution. *Mol. Biol. Evol.* **35**, 1327–
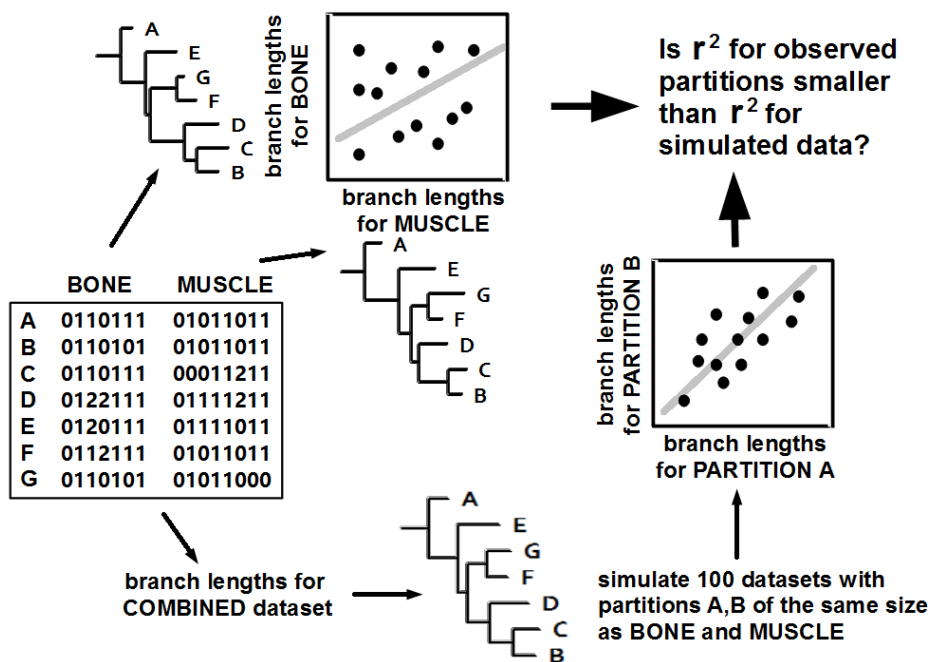
598 1331.

Figure 1

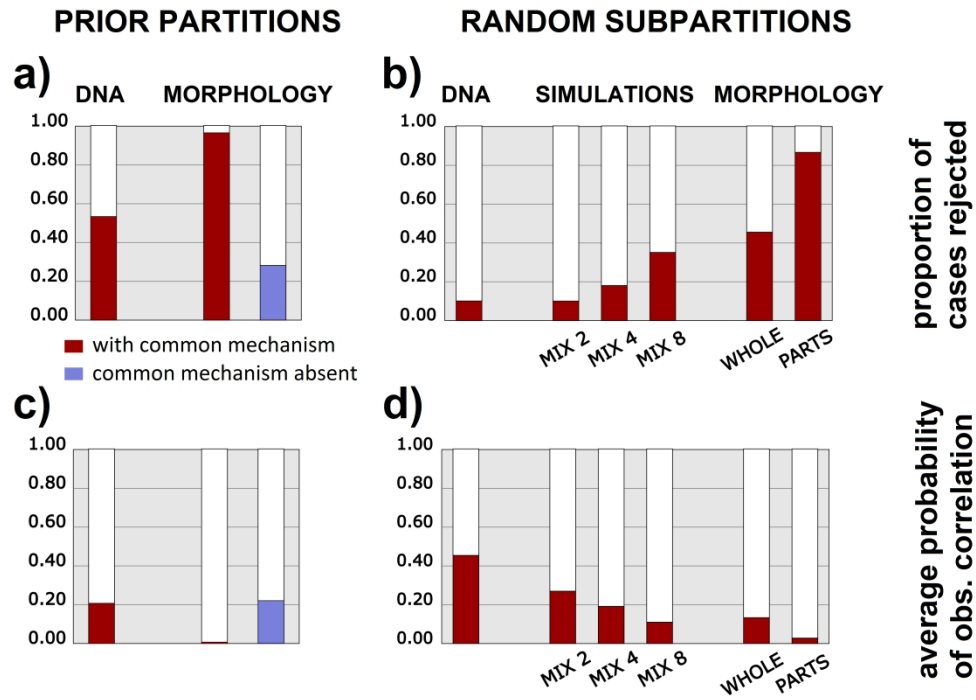757x382mm (96 x 96 DPI)

Figure 2

278x210mm (96 x 96 DPI)

Figure 3

1280x927mm (96 x 96 DPI)
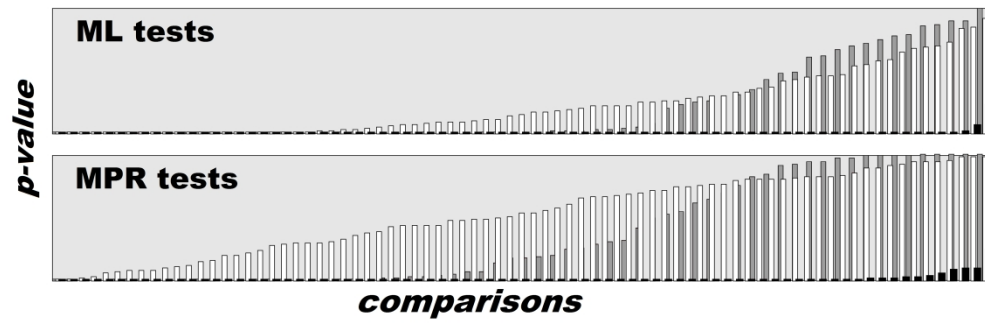
Figure 4
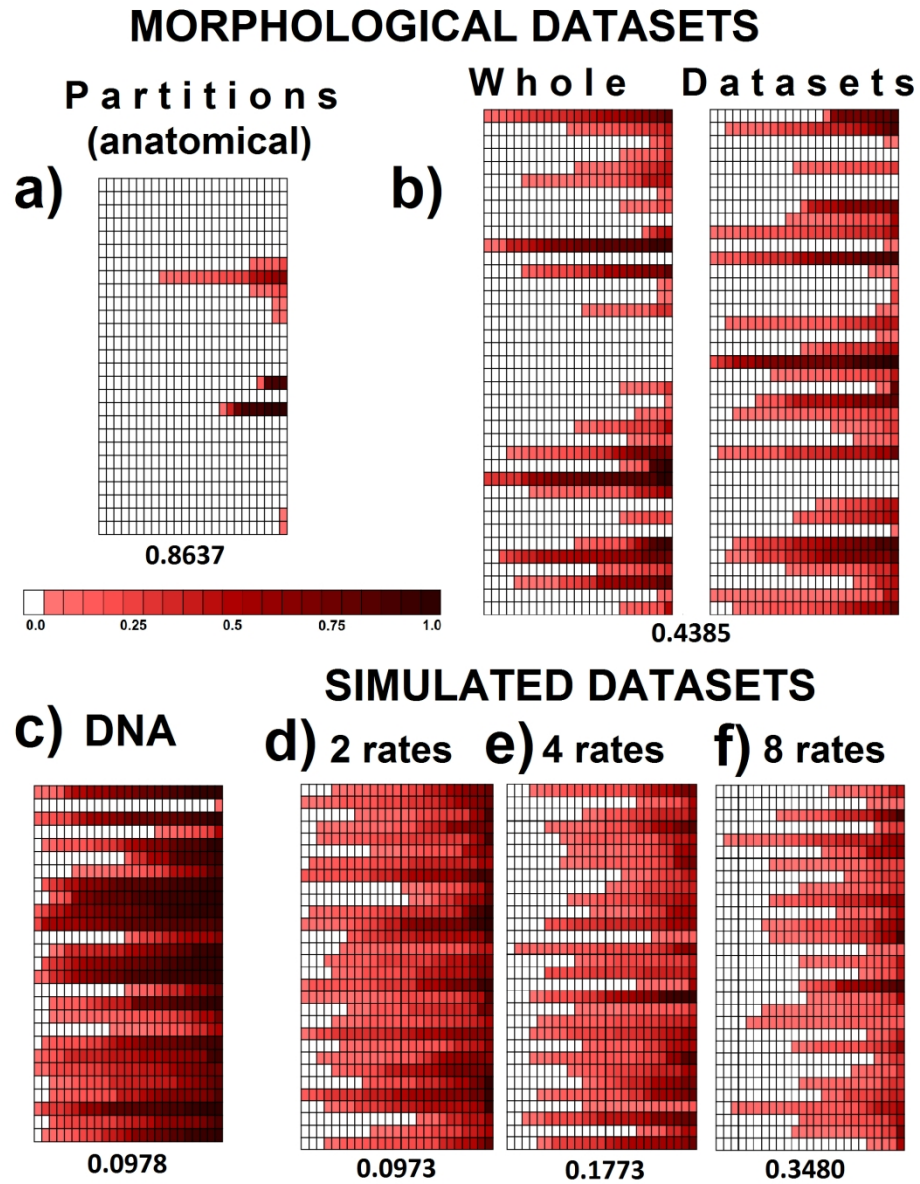
1320x438mm (96 x 96 DPI)
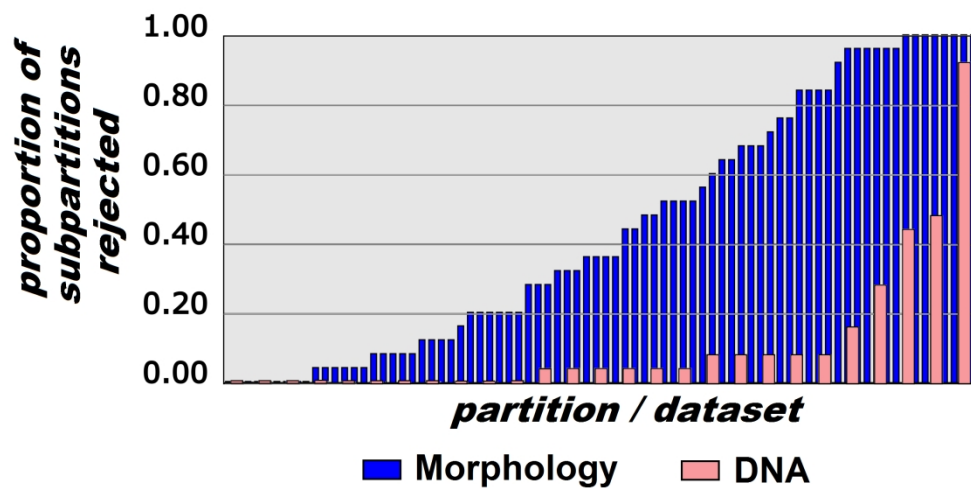
Figure 5

546x702mm (96 x 96 DPI)
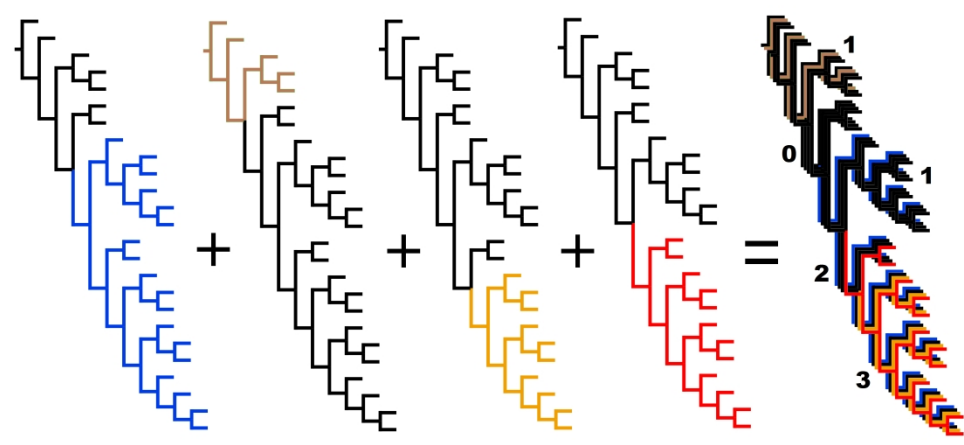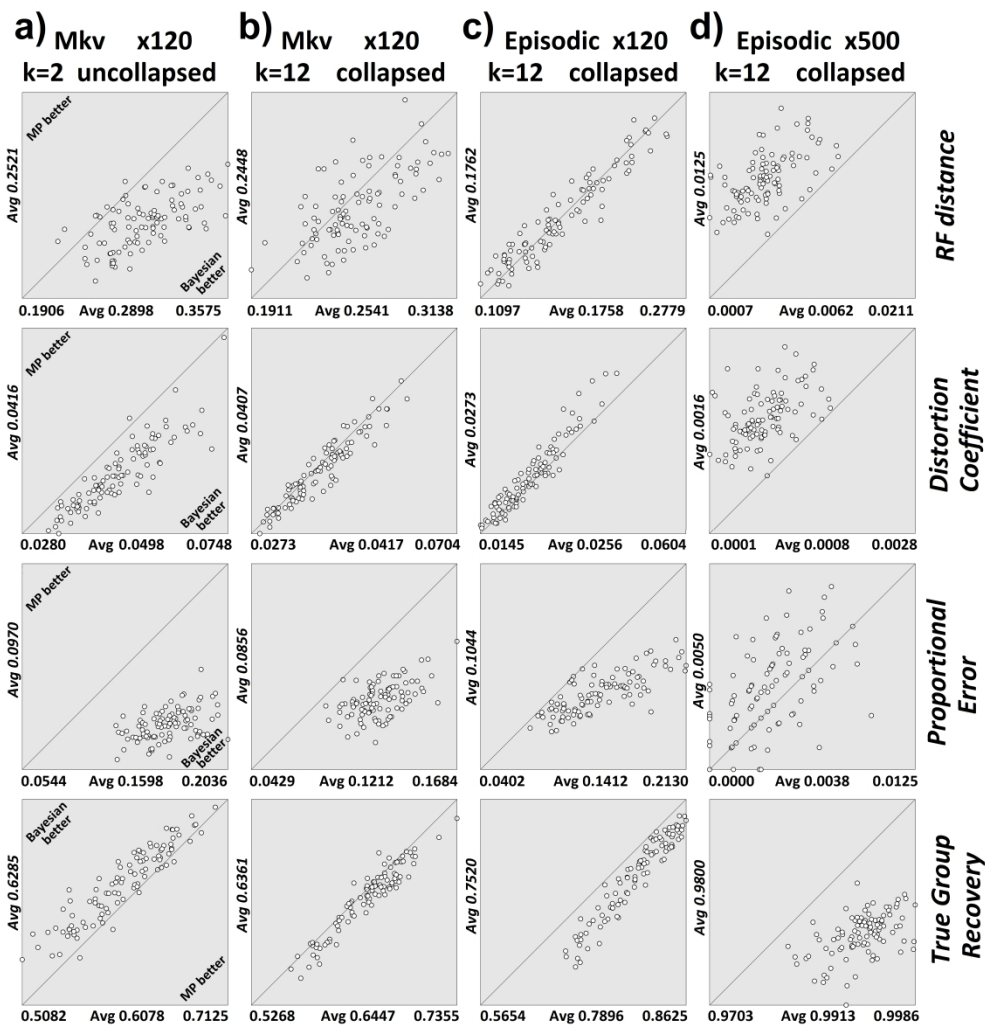
Figure 6

712x358mm (96 x 96 DPI)

Figure 7

1301x618mm (96 x 96 DPI)

Figure 8

1508x1547mm (96 x 96 DPI)