

1 Missing the point (estimate): Bayesian and likelihood phylogenetic
2 reconstructions of morphological characters produce generally
3 concordant inferences. A comment on Puttick *et al.*

4 Joseph W. Brown^{1*}, Caroline Parins-Fukuchi^{1*}, Gregory W. Stull¹, Oscar M. Vargas¹, and
5 Stephen A. Smith¹

6 ¹Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor,
7 Michigan 48109, USA

8 *Equal authorship. Emails: josephwb@umich.edu, cfukuchi@umich.edu

9 **Abstract**

10 Puttick *et al.* [1] performed a simulation study to compare accuracy between methods inferring
11 phylogeny from discrete morphological characters. They report that a Bayesian implementation
12 of the Mk model [2] was most accurate (but with low resolution), while a maximum likelihood
13 (ML) implementation of the same model was least accurate. They conclude by strongly
14 advocating that Bayesian implementations of the Mk model should be the default method of
15 analysis for such data. While we applaud investigations into accuracy and alternative methods of
16 analysis, this conclusion is based on an inappropriate comparison of the ML point estimate with
17 the Bayesian consensus. We revisit these issues through simulation by considering uncertainty in
18 ML reconstructions, and demonstrate that Bayesian and ML estimates are generally concordant
19 when conventional edge support thresholds are considered. We therefore disagree with the
20 conclusions of [1], and consider their prescription of any default method to be unfounded.
21 Instead, we recommend caution and thoughtful consideration of the model or method being
22 applied to a morphological dataset.

23 **Key words:** phylogeny, morphology, paleontology, Bayesian, likelihood

24 **Comparing point estimates to consensus summaries**

25 Puttick *et al.* [1] report that ML tree inference under the Mk model results in higher topological
26 error than Bayesian implementations. However, this result is driven precisely by the comparison
27 of maximum likelihood point estimates (MLE) to Bayesian majority-rule (BMR) consensus trees.
28 MLE topologies are fully resolved, but this stems from the standard binary tree searching
29 algorithms employed and not from an explicit statistical rejection of unresolved nodes. Therefore,
30 individual MLE estimates may contain edges with negligible statistical support. On the other
31 hand, consensus summaries, independent of phylogenetic method, may have reduced resolution

32 as a product of uncertainty arising by summarization across conflicting sampled topologies. Thus,
33 a direct comparison between a consensus tree (i.e., BMR) and a point estimate (i.e., MLE) is
34 inappropriate. BMR topologies of [1] are more accurate because poorly supported conflicted
35 edges were collapsed, while MLE topologies were fully resolved, even if poorly supported. While
36 comparison between MLE and Bayesian maximum *a posteriori* (MAP) trees would be a more
37 appropriate for optimal point estimates, the incorporation of uncertainty is an integral part of all
38 phylogenetic analysis. Therefore, comparison of consensus trees from Bayesian and ML analyses
39 hold more practical utility for systematists. For these reasons, we argue that the results of [1] are
40 an artifact of their comparison between fundamentally incomparable sets of trees.

41 **Support metrics are available for morphological characters**

42 To avoid drawing untenable conclusions, it is *de rigueur* of any phylogenetic analysis to explicitly
43 assess edge support. Systematists often accomplish this via non-parametric bootstrap sampling
44 [3], though other measures exist (see below). Puttick et al. did not assess edge support in their
45 ML estimates, stating that morphological (but not genetic) data do not meet an underlying
46 assumption of the bootstrap statistical procedure that phylogenetic signal is distributed randomly
47 among characters. The authors do not explain the meaning of this statement, and no references
48 are provided to support the assertion. Non-parametric bootstrapping has been a staple of
49 phylogenetic reconstruction for decades, including for the analysis of discrete morphological
50 characters. Bootstrapping works via the assumption that the observed characters are a
51 representative sample from a population of possible characters evolving under the same process,
52 and thus can be resampled to assess confidence in parameters [3]. While morphological matrices
53 typically include only variable characters (i.e., an ascertainment bias), this is an informative
54 subset of the possible characters, and should not be thought of as misleading calculations. Were
55 this otherwise, the original sample would likewise be suspect, as the use of model-based
56 phylogenetic inference (such as Mk) explicitly assumes characters evolve according to the same
57 process. Concerns about the interpretation and use of the bootstrap exist [4], the primary of which
58 involves the assumption that individual characters are statistically independent. However, it is
59 reasonable to assume that individual sites in a morphological matrix would be more independent
60 than adjacent sites from the same gene, and genetic datasets are routinely bootstrapped. We
61 therefore disagree with the claims of [1] that bootstrapping is inappropriate for morphological
62 data, or at least any *more* inappropriate than for genetic data.

63 There are also other methods researchers can use to assess edge support in a likelihood
64 framework. Jackknifing, unlike bootstrapping, samples without replacement, conditioning on
65 strict subsets of the observed data. More recently, the SH-like test [5] computes support for each
66 internal edge in the MLE tree by considering all nearest neighbour interchanges (NNIs). This test
67 is implemented in several software packages including RAxML [6], one of the programs used by
68 [1]. Alternatively, some ML programs offer an option to collapse edges on a MLE tree that fall
69 below some minimum threshold length. Use of any of these options would enable a fairer
70 comparison of likelihood and Bayesian reconstructions.

71 **ML and Bayesian comparisons incorporating uncertainty**

72 To measure the effect of comparing BMR and MLE trees, we used the simulation code from [1]
73 to generate 1000 character matrices, each of 100 characters on a fully pectinate tree of 32 taxa, as
74 these settings generated the most discordant results. Each matrix was analyzed in both Bayesian
75 and ML frameworks using the Mk+G model [2]. Bayesian reconstructions were performed using
76 MrBayes v3.2.6 [7], using the same settings as [1]: 2 runs, each with 5×10^5 generations,
77 sampling every 50 generations, and discarding the first 25% samples as burnin. As in [1], we
78 summarized each analysis with a BMR consensus tree (i.e. only edges with ≥ 0.5 posterior
79 probability are represented). Likelihood analyses were performed in RAxML v8.2.9 [6]. For each
80 simulated matrix we inferred both the MLE tree and 200 nonparametric bootstrap trees. Accuracy
81 in topological reconstruction was assessed using the Robinson-Foulds (RF) distance [8], which
82 counts the number of unshared bipartitions between trees. We measured the following distances
83 from the true tree: d_{BMR} , the distance to the Bayesian majority-rule consensus; d_{MLE} , the distance
84 to the MLE tree; d_{ML50} , the distance to the MLE tree which has had all edges with $< 50\%$
85 bootstrap support collapsed. Finally, for each matrix we calculate $D_{\text{MLE}} = d_{\text{MLE}} - d_{\text{BMR}}$, and D_{ML50}
86 $= d_{\text{ML50}} - d_{\text{BMR}}$. These paired distances measure the relative efficacy of ML and Bayesian
87 reconstructions: values of D greater than 0 indicate that ML produces less accurate estimates (that
88 is, with a greater RF distance from the true generating tree).

89 As demonstrated by [1], MLE trees are indeed less accurate than BMR trees (Figure 1; D_{MLE}),
90 with MLE trees on average having an RF distance 17.6 units greater than the analogous Bayesian
91 distance. However, when collapsing MLE edges with less than 50% bootstrap support, Bayesian
92 and ML differences are normally distributed around 0 (Figure 1; D_{ML50}), indicating that when
93 standardizing the degree of uncertainty in tree summaries there is no difference in topology
94 reconstruction accuracy. These results support the argument that the original comparisons of MLE
95 and BMR trees are inappropriate. Depending on the level uncertainty involved, an optimal point
96 estimate from a distribution (e.g., MLE or MAP) may be arbitrarily distant from a summary of the
97 same distribution. And so, the differences in MLE vs. BMR are not expected to be consistent.

98 **The expected concordance of Bayesian and ML results**

99 Our results reveal much greater congruence between Bayesian and ML estimates than suggested
100 by [1]. This is to be expected. ML and Bayesian tree construction methods should yield similar
101 results under the conditions in which they are often employed. While Bayesian tree
102 reconstruction differs from ML by incorporating prior distributions, the methods share likelihood
103 functions. In phylogenetics, researchers typically adopt non-informative priors, with a few
104 exceptions (e.g., priors on divergence time parameters). Arguments can be made for
105 pseudo-Bayesian approaches when care is taken to ensure that priors used are truly uninformative,
106 which result in posterior probabilities that mirror the likelihood and are therefore congruent with
107 ML [9,10]. If prior distributions are formulated thoughtfully, as with [11] in shaping the Mk
108 model using hyperpriors to accommodate character change heterogeneity, Bayesian methods can
109 outperform ML. Alternatively, inappropriate priors can positively mislead [10]. Generally, when
110 informative prior distributions are known or can be estimated using hierarchical approaches,
111 Bayesian reconstruction methods may be strongly favored over ML. It is unclear whether [1]
112 intend to draw the comparisons discussed above as they do not describe any reasons to prefer

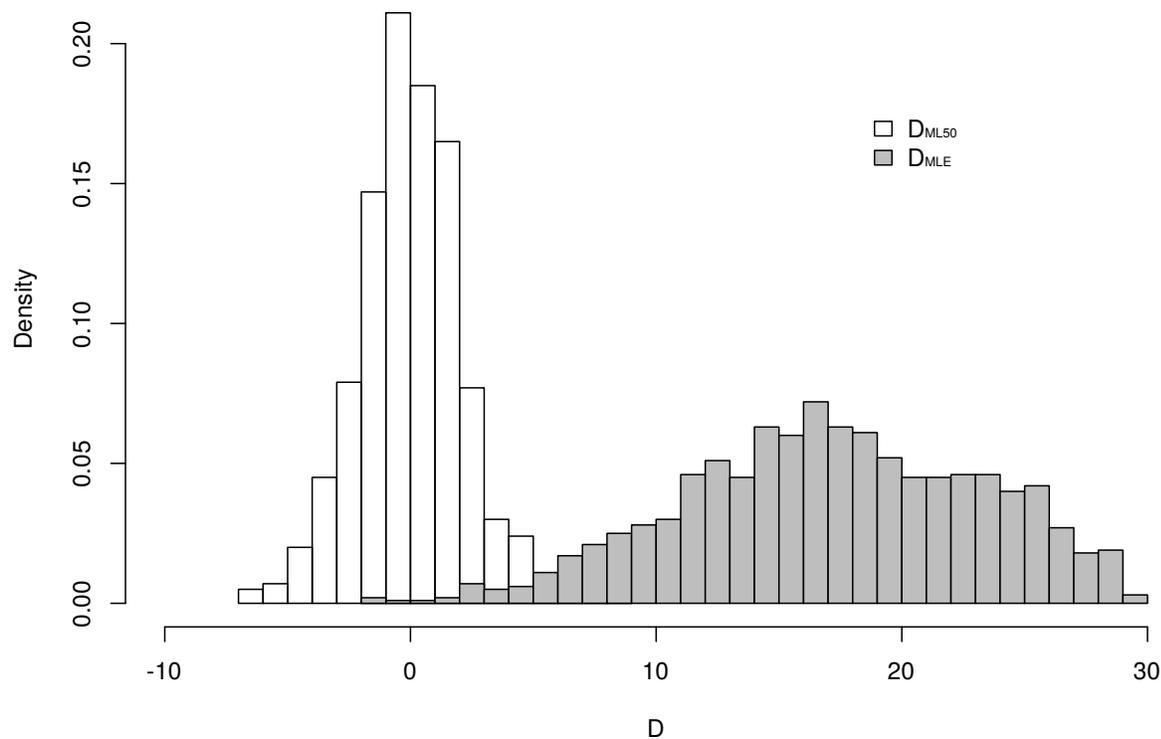


Figure 1: Topological accuracy of ML vs. Bayesian reconstructions. D measures how much larger ML distances are from the true tree (d_{ML}) than Bayesian distances (d_{BMR}). MLE trees are indeed less accurate than BMRs (D_{MLE}), but when conventional bootstrap thresholds are employed (D_{ML50}) the difference in efficacy disappears.

113 Bayesian over ML in principle.

114 Although our results demonstrate general concordance between ML and Bayesian approaches
115 when uncertainty is represented, further simulation work is needed to determine the extent and
116 conditions of this concordance. Issues surrounding the application of Bayesian methods are
117 particularly important in paleontology, where researchers often conduct inference upon very
118 limited data. In these cases, it may be desirable to construct informative prior distributions when
119 conducting Bayesian analysis [10]. The questions posed by [1] are critically important as
120 statistical morphological phylogenetics moves forward. However, their inappropriate comparison
121 between ML and Bayesian approaches leaves the relative performance of the two
122 implementations of the Mk model unresolved.

123 We are not advocating one method over another for morphological phylogenetic
124 reconstruction. Methods differ in model (Mk vs. parsimony), inferential paradigm (parsimony vs.
125 ML/Bayesian), assumptions (prior distributions, model adequacy), interpretation, and means to
126 incorporate uncertainty (ML/parsimony vs. Bayesian). We therefore recommend caution and
127 thoughtful consideration of the biological question being addressed and then choosing the method

128 that will best address that question. All inferential approaches possess strengths and weaknesses,
129 and it is the task of researchers to determine the most appropriate given available data and the
130 questions under investigation. The excitement of new morphological data sources and new means
131 for analyzing these data should not overshadow the obligation to apply methods thoughtfully.

132 **Authors' contributions**

133 J.W.B. conceived the design of the study and performed the analyses; J.W.B. and C.F.-P. drafted
134 the manuscript; all authors contributed to the interpretation of results and the writing of the
135 manuscript.

136 **Acknowledgements**

137 We thank Mark Puttick for sharing datasets and thoughts on bootstrap resampling. We thank
138 members of the Smith laboratory for thoughtful discussions on an earlier draft of this manuscript.
139 J.W.B. and C.F.-P. thank Annika Hansen for being a stalwart leading example of objective
140 criticism. This is paper #1 of the PRUSSIA working group at UM.

141 **References**

- 142 1. Puttick MN *et al.* 2017 Uncertain-tree: discriminating among competing approaches to the
143 phylogenetic analysis of phenotype data. *Proc. R. Soc. B* **284**, 20162290.
- 144 2. Lewis PO. 2001 A likelihood approach to estimating phylogeny from discrete morphological
145 character data. *Syst. Biol.* **50**, 913-925.
- 146 3. Felsenstein J. 1985 Confidence limits on phylogenies: an approach using the bootstrap.
147 *Evolution* **39**, 783-791.
- 148 4. Sanderson MJ. 1995 Objections to bootstrapping phylogenies: a critique. *Syst. Biol.* **44**,
149 299-320.
- 150 5. Guindon S *et al.* 2010 New algorithms and methods to estimate maximum-likelihood
151 phylogenies: assessing the performance of PhyML 3.0. *Syst. Biol.* **59**, 307-321.
- 152 6. Stamatakis A. 2014 RAxML version 8: a tool for phylogenetic analysis and post-analysis of
153 large phylogenies. *Bioinf.* **30**, 1312-1313.
- 154 7. Ronquist F *et al.* 2012 MrBayes 3.2: efficient Bayesian phylogenetic inference and model
155 choice across a large model space. *Syst. Biol.* **61**, 539-542.
- 156 8. Robinson DF, Foulds LR. 1981 Comparison of phylogenetic trees. *Math. Biosci.* **53**, 131-147.
- 157 9. Alfaro ME, Holder MT. 2006 The posterior and the prior in Bayesian phylogenetics. *Ann. Rev.*
158 *Ecol. Evol. Syst.* **37**, 19-42.
- 159 10. Gelman A, Carlin JB, Stern HS, Rubin DB. 2014 Bayesian Data Analysis (Vol. 2). Boca
160 Raton, FL, USA: Chapman & Hall/CRC.
- 161 11. Wright AM, Lloyd GT, Hillis DM. 2016 Modeling character change heterogeneity in
162 phylogenetic analyses of morphology through the use of priors. *Syst. Biol.* **65**, 602-611.