

1           On the importance of credibility: Bayesian and likelihood  
2           phylogenetic reconstructions of morphological traits are  
3           concordant. A comment on Puttick *et al.*

4           Joseph W. Brown<sup>1\*</sup>, Caroline Parins-Fukuchi<sup>1\*</sup>, Gregory W. Stull<sup>1</sup>, Oscar M. Vargas<sup>1</sup>, and  
5           Stephen A. Smith<sup>1</sup>

6           <sup>1</sup>Department 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.* (hereafter, PEA) [1] performed a simulation study to compare accuracy among  
11 methods of inferring phylogeny from discrete morphological characters. They report that a  
12 Bayesian implementation of the Mk model [2] was most accurate (but with low resolution), while  
13 a maximum likelihood (ML) implementation of the same model was least accurate. They  
14 conclude by strongly advocating that Bayesian implementations of the Mk model should be the  
15 default method of analysis for such data. While we appreciate the authors' attempt to investigate  
16 the accuracy of alternative methods of analysis, their conclusion is based on an inappropriate  
17 comparison of the ML point estimate, which does not consider confidence, with the Bayesian  
18 consensus, which incorporates estimation credibility into the summary tree. Using simulation, we  
19 demonstrate that ML and Bayesian estimates are concordant when confidence and credibility are  
20 comparably reflected in summary trees, a result expected from statistical theory. We therefore  
21 disagree with the conclusions of PEA and consider their prescription of any default method to be  
22 poorly founded. Instead, we recommend caution and thoughtful consideration of the model or  
23 method being applied to a morphological dataset.

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

25 **Comparing point estimates to consensus summaries**

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

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

### 43 **Confidence and credibility are fundamental to inference**

44 To avoid drawing untenable conclusions, it is *de rigueur* of any statistical analysis to explicitly  
45 assess the robustness of an inference. Non-parametric bootstrap sampling [3] is the overwhelming  
46 standard in phylogenetic confidence estimation. PEA did not assess edge support in their ML  
47 estimates, stating that morphological data do not meet an underlying assumption of the bootstrap  
48 statistical procedure that “phylogenetic signal is distributed randomly among characters,” but  
49 provide no references to support the assertion. Non-parametric bootstrapping has been a staple of  
50 phylogenetic reconstruction for decades, including for the analysis of discrete morphological  
51 characters. Like Bayesian credibility estimation, bootstrapping estimates confidence by assuming  
52 that empirical data are a representative sample from an underlying distribution of characters  
53 evolving independently under a shared process [3]. As PEA note, the assumption of independence  
54 may often be violated. However, this violation is fundamentally problematic to model-based  
55 phylogenetics in general. Contrary to PEA, Bayesian and frequentist approaches to confidence  
56 estimation are similar in the sense that both provide distribution-based summaries of uncertainty,  
57 the sole distinguishing factor of Bayesian approaches is the incorporation of prior densities. Many  
58 of the concerns raised in relation to the bootstrap can thus also be shared by Bayesian approaches  
59 and should not preclude its use more generally. While there are concerns about the use and  
60 interpretation of the bootstrap [4], genetic datasets are routinely bootstrapped. Without additional  
61 information, it may be reasonable to assume that individual characters in a morphological matrix  
62 would be more independent than adjacent sites from the same gene (for which the  
63 interdependence among characters is far better understood). We thus dispute the assertion that  
64 bootstrapping is uniquely problematic for morphological data.

65 While Bayesian approaches estimate credibility intervals during parameter sampling,  
66 confidence assessment is equally fundamental to likelihood analyses. In addition to the bootstrap,  
67 alternatives such as jackknifing and the SH-like test [5] are also implemented in popular software  
68 packages such as RAxML [6], one of the programs used by PEA. ML packages also frequently  
69 offer an option to collapse edges on a MLE tree that fall below some minimum threshold length.  
70 Use of any of these options would enable a more sensible comparison of likelihood and Bayesian  
71 reconstructions.

## 72 **ML and Bayesian comparisons incorporating uncertainty**

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

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

## 100 **The expected concordance of Bayesian and ML results**

101 Our results reveal much greater congruence between Bayesian and ML estimates than suggested  
102 by PEA. This is to be expected and is reassuring. ML and Bayesian tree construction methods  
103 should yield similar results under the conditions in which they are often employed. While  
104 Bayesian tree reconstruction differs from ML by incorporating prior distributions, the methods  
105 share likelihood functions. In phylogenetics, researchers typically adopt non-informative priors,  
106 with a few exceptions (e.g., priors on divergence time parameters). Arguments can be made for  
107 pseudo-Bayesian approaches when care is taken to ensure that priors used are truly uninformative,  
108 which result in posterior probabilities that mirror the likelihood and are therefore congruent with  
109 ML [9, 10]. If prior distributions are formulated thoughtfully, as with [11] in shaping the Mk  
110 model using hyperpriors to accommodate character change heterogeneity, Bayesian methods can  
111 outperform ML. Alternatively, inappropriate priors can positively mislead [10]. Generally, when  
112 informative prior distributions are known or can be estimated using hierarchical approaches,  
113 Bayesian reconstruction methods may be strongly favoured over ML. It is unclear whether PEA

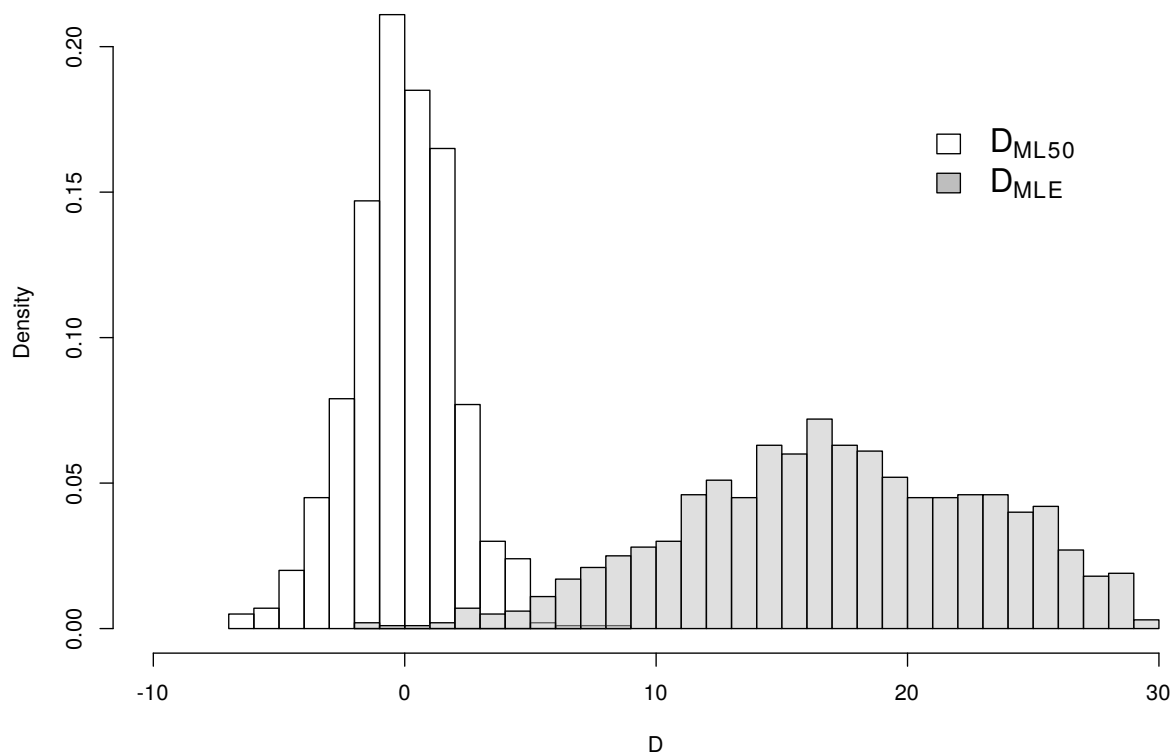


Figure 1: Topological accuracy of ML vs. Bayesian reconstructions for the most discordant comparison identified by PEA (see text).  $D$  measures how much larger ML distances are from the true tree ( $d_{ML}$ ) than are Bayesian distances ( $d_{BMR}$ ). MLE trees are indeed less accurate than BMRs ( $D_{MLE}$ ; mean = 17.63), but when conventional bootstrap thresholds are employed the difference in efficacy disappears ( $D_{ML50}$ ; mean = 0.43).

114 intend to draw the comparisons discussed above as they do not describe any reasons to prefer  
115 Bayesian over ML in principle.

116 Although our results demonstrate general concordance between ML and Bayesian approaches  
117 when uncertainty is represented, further simulation work is needed to determine the extent and  
118 conditions of this concordance. Issues surrounding the application of Bayesian methods are  
119 particularly important in paleontology, where researchers often conduct inference upon very  
120 limited data. In these cases, it may be desirable to construct informative prior distributions when  
121 conducting Bayesian analyses [10]. The questions posed by PEA are sensible in light of current  
122 enthusiasm for statistical morphological phylogenetics. However, the relative performance of the  
123 implementations of the Mk model remain unresolved due to the authors' misaligned treatment of  
124 confidence. This lack of resolution extends to their treatment of parsimony, which is invalid for  
125 the same reason as their ML comparison.

126 We do not advocate any one method for morphological phylogenetic reconstruction. Methods  
127 differ in model (Mk vs. parsimony), inferential paradigm (parsimony vs. ML/Bayesian),

128 assumptions (prior distributions, model adequacy), interpretation, and means to incorporate  
129 uncertainty (ML/parsimony vs. Bayesian). We therefore recommend caution and thoughtful  
130 consideration of the biological question being addressed and then choosing the method that will  
131 best address that question. All inferential approaches possess strengths and weaknesses, and it is  
132 the task of researchers to determine the most appropriate given available data and the questions  
133 under investigation. The excitement of new morphological data sources and new means for  
134 analyzing these data should not overshadow the obligation to apply methods thoughtfully.

## 135 **Data accessibility**

136 Simulated and inferred trees analyzed in this publication can be accessed in the electronic  
137 supplementary material.

## 138 **Authors' contributions**

139 All authors jointly developed the conceptual basis of the manuscript and study; J.W.B. conceived  
140 of and performed the simulations; J.W.B. and C.P.-F. drafted the manuscript; all authors  
141 contributed to the interpretation of results and the writing of the manuscript.

## 142 **Funding**

143 J.W.B. and S.A.S were supported by NSF DEB AVATOL grant 1207915. G.W.S. was supported  
144 by NSF DBI grant 1612032. O.M.V. was supported by NSF grants FESD 1338694 and DEB  
145 1240869.

## 146 **Acknowledgements**

147 We thank M. Puttick for an open and constructive discourse throughout this process. Two  
148 anonymous reviewers provided thoughtful reviews on an earlier draft of this manuscript. We  
149 thank Editor G. Carvalho and Associate Editor P. Makovicky for considering an appeal to an  
150 earlier decision on this manuscript, and P. Donoghue for reversing his original recommendation.  
151 Finally, J.W.B. and C.P.-F. thank Annika Hansen for being a stalwart leading example of objective  
152 criticism. This is paper #1 of the PRUSSIA working group at UM.

## 153 **References**

154 [1] Puttick MN, O'Reilly JE, Tanner AR, Fleming JF, Clark J, Holloway L, et al.  
155 Uncertain-tree: discriminating among competing approaches to the phylogenetic analysis of  
156 phenotype data. *Proceedings of the Royal Society of London B: Biological Sciences*.  
157 2017;284(1846):20162290.

- 158 [2] Lewis PO. A likelihood approach to estimating phylogeny from discrete morphological  
159 character data. *Systematic Biology*. 2001;50(6):913–925.
- 160 [3] Felsenstein J. Confidence limits on phylogenies: an approach using the bootstrap.  
161 *Evolution*. 1985;39(4):783–791.
- 162 [4] Sanderson MJ. Objections to bootstrapping phylogenies: a critique. *Systematic Biology*.  
163 1995;44(3):299–320.
- 164 [5] Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms  
165 and methods to estimate maximum-likelihood phylogenies: assessing the performance of  
166 PhyML 3.0. *Systematic Biology*. 2010;59(3):307–321.
- 167 [6] Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
168 phylogenies. *Bioinformatics*. 2014;30(9):1312–1313.
- 169 [7] Ronquist F, Teslenko M, Van Der Mark P, Ayres DL, Darling A, Höhna S, et al. MrBayes  
170 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space.  
171 *Systematic Biology*. 2012;61(3):539–542.
- 172 [8] Robinson DF, Foulds LR. Comparison of phylogenetic trees. *Mathematical Biosciences*.  
173 1981;53(1-2):131–147.
- 174 [9] Alfaro ME, Holder MT. The posterior and the prior in Bayesian phylogenetics. *Annual*  
175 *Review of Ecology, Evolution, and Systematics*. 2006;37:19–42.
- 176 [10] Gelman A, Carlin JB, Stern HS, Rubin DB. *Bayesian Data Analysis*. vol. 2. Chapman &  
177 Hall/CRC Boca Raton, FL, USA; 2014.
- 178 [11] Wright AM, Lloyd GT, Hillis DM. Modeling character change heterogeneity in  
179 phylogenetic analyses of morphology through the use of priors. *Systematic Biology*.  
180 2016;65(4):602–611.