

Bayesian and parsimony approaches reconstruct informative trees from simulated morphological datasets

Journal:	Biology Letters
Manuscript ID	RSBL-2018-0632.R3
Article Type:	Research
Date Submitted by the Author:	n/a
Complete List of Authors:	Smith, Martin; Durham University, Earth Sciences
Subject:	Evolution < BIOLOGY, Taxonomy and Systematics < BIOLOGY, Palaeontology < BIOLOGY
Categories:	Evolutionary Biology
Keywords:	phylogenetic inference, parsimony analysis, Bayesian phylogenetic methods, information content, equal weights, implied weighting



1 Running head: COMPARING PHYLOGENETIC METHODS

- 2 Bayesian and parsimony approaches reconstruct informative trees from simulated
- 3 morphological datasets
- 4 Martin R. Smith
- 5 Department of Earth Sciences, Lower Mount Joy, Durham University, Durham, DH1 3LE, UK;
- 6 *martin.smith@durham.ac.uk*

Phylogenetic analysis aims to establish the true relationships between taxa. Different analytical methods, however, can reach different conclusions. In order to establish which approach best reconstructs true relationships, previous studies have simulated datasets from known tree topologies, and identified the method that reconstructs the generative tree most accurately. On this basis, researchers have argued that morphological datasets should be analysed by Bayesian approaches, which employ an explicit probabilistic model of evolution, rather than parsimony methods – with implied weights parsimony <u>sometimes</u> identified as particularly inaccurate.

14 Accuracy alone, however, is an inadequate measure of a tree's utility: a fully unresolved 15 tree is perfectly accurate, yet contains no phylogenetic information. The highly resolved trees 16 recovered by implied weights parsimony in fact contain as much useful information as the more 17 accurate, but less resolved, trees recovered by Bayesian methods. By collapsing poorly 18 supported groups, this superior resolution can be traded for accuracy, resulting in trees as 19 accurate as those obtained by a Bayesian approach. In contrast, equally weighted parsimony 20 analysis produces trees that are less resolved and less accurate, leading to less reliable 21 evolutionary conclusions.

Keywords.—phylogenetic inference; parsimony analysis; equal weights; implied weighting;
 Bayesian phylogenetic methods; information content

Smith

24

25 INTRODUCTION

26 Evolutionary history can be reconstructed using parsimony-based or probabilistic approaches. 27 Because models used with molecular datasets generally share a common probabilistic 28 construction, statistical methods can be used to determine the most appropriate model [1]. With 29 morphological datasets, however, it is more difficult to establish whether probabilistic models or 30 parsimony better reconstruct phylogenetic relationships (which are typically unknown). 31 A pragmatic approach to this question is to simulate data from a known tree. With the 32 important caveat that generative trees and simulated morphological datasets may be unrealistic 33 [2,3], probabilistic approaches typically reconstruct the generative tree most accurately (i.e. with 34 least conflict), followed by parsimony under equal and implied weights in turn [4–9]. 35 Previous studies have advocated accuracy as the sole criterion by which to select a 36 method [5–11]. Congreve & Lamsdell [9] (problematically [2]) define the most accurate tree as 37 the one that bears the fewest of incorrect splits. Other authors [5-8,11] use the Robinson-Foulds 38 distance as a proxy for accuracy (even though the RF distance is also influenced by precision; a 39 pair of trees can be made two units more similar by replacing an incorrect partition with a correct 40 one, or by collapsing two incorrect partitions.) Goloboff et al. [2] propose alternative tree 41 similarity metrics as proxies for accuracy.

Accuracy alone, however, is not the only goal when reconstructing trees [11]. No tree shows less conflict than a single polytomy, for a total absence of relationship information guarantees that no relationship is incorrectly resolved. An emphasis on accuracy therefore disadvantages methods that produce highly resolved trees [11] (and *vice versa*). This trade-off has been acknowledged by collapsing some poorly supported groups before calculating accuracy (which even if accuracy is still equated with 'performance') [2,6,8,11]. Naturally [12], methods that yield less resolution are consistently more accurate [2,5,7,8,11].

Comparing Phylogenetic Methods

Smith

49 We should be seeking not the most accurate method, but the method that recovers as 50 much *information* as possible about the true tree, striking a balance between the complementary 51 quantities [12] of accuracy and resolution. For example, a tree that resolves 20 relationships 52 conveys much information about the correct tree, even if one of those relationships is incorrect; a 53 tree that resolves just one relationship conveys less information, even if that single relationship is 54 correct. If two trees are equally accurate, we should prefer the more precise. Here I explore the 55 impact on previous studies of evaluating trees according to their total shared information content, rather than 'accuracy' alone. 56

57 Methods

58 Congreve and Lamsdell [9; CL hereafter] simulated 55-character matrices from a bifurcating 22-59 tip tree using a Markov *k*-state 1 parameter model with rates sampled from a discretized Gamma 60 distribution. Their generative tree is the single most parsimonious tree obtained from a study of 61 Ordovician trilobites; its edges were assigned a unit length.

O'Reilly *et al.* [5; OR hereafter] simulated matrices containing 100, 350 and 1000
 characters from a bifurcating 75-tip tree using a modified HKY85 model; they followed a
 previous simulation study [4] in selecting a single bifurcating tree from a morphological +
 molecular analysis of Lissamphibia.

66 I used TNT [13] to conduct parsimony searches on each of these matrices under equal 67 and implied weights, using the parsimony ratchet and sectorial search heuristics (search options: 68 xmult:hits 20 level 4 chklevel 5 rat10 drift10). I took a strict consensus of all optimal trees obtained under equal weights, and under implied weights [14] at the concavity 69 70 constants used in each respective study (CL: k = 1, 2, 3, 5 and 10; OR: k = 2, 3, 5, 10, 20 and 71 200). For each dataset I generated a further strict consensus of all trees that were optimal under 72 any of the concavity constants, excluding the unreasonable value of k = 1, which inadequately 73 penalises extra steps beyond the first, and thus exhibits undesirable properties of clique analysis 74 [15] (see Supplementary Text).

Comparing Phylogenetic Methods

Page 4 of 16

I also generated majority-rule consensus trees in MrBayes 3.2.2 [16] using an Mk model, with rates distributed according to a gamma parameter. I combined results from four independent runs, each of which employed four Metropolis-coupled Markov chains. After a burn-in period of 4 000 000 generations, the cold chain in each run was sampled every 10 000 generations for 6 000 000 generations. The sampled topologies faithfully reflected the posterior distribution for each dataset (0.999 < PSRF < 1.001; ESS > 400).

81 To explore the relationship between resolution and accuracy, I generated further trees for 82 each analysis by collapsing poorly supported groups. Under the Mk model, I collapsed groups 83 whose posterior probability was < 95%, 90%, 85%, ... 50%. In parsimony analyses, I compared 84 different measures of node support. Under Jackknife and Bootstrap resampling, I collapsed 85 groups with (i) absolute frequency supports of < 0%, 2%, 4% ... 100%; (ii) relative frequency 86 (GC) support of < --100%, --95%, ... 95%, 100%. Under Bremer support, I collapsed groups 87 with Bremer support values less than $1, 2, 3, \dots 20$ with equally weighted trees (TNT command 88 subopt x; bbreak;); under implied weighting, Bremer support values were drawn from a logarithmic distribution (0.73^{0...19}, 2.5×10⁻³ \rightarrow 1×10⁰), reflecting the fractional nature of tree 89 90 scores under implied weights [14].

91 Symmetric difference metrics calculate how much information two trees hold in common 92 [17] — that is, how much information a generated tree contains about the generative tree. 93 Where the generative tree is bifurcating, a particular relationship may be resolved the same way 94 (s) or a different way (d) on each tree, or resolved in the comparison tree only (r) [18,19]. The 95 symmetric difference ('SD', also termed the Robinson-Foulds distance) is given by 2d + r. The 96 symmetric difference is conventionally normalized against the total information present ('TIP') in the two trees, 2d + 2s + r [19]. Undesirably, this assigns a fully unresolved tree an equal score 97 98 to a tree that is perfectly resolved and completely incorrect (Fig. 1a). In the present context, 99 therefore, it is more appropriate to normalize against the maximum information ('MaxI') that 100 could potentially have been resolved, 2 (d + s + r).

Comparing Phylogenetic Methods

Smith

101 The unit of relationship information may be a quartet (a four-taxon statement) [18-20] or 102 a bipartition split [21–23]. (Each clade in a tree corresponds to a bipartition that splits taxa into 103 'members' and 'non-members'.) Partitions offer a simple but incomplete measure of the 104 relationship topological information accommodated in a tree. The trees ((A, (X, B)), (C, D)) and 105 ((A, B), ((C, X), D)) both contain the same information regarding the relationships between (A, 106 B) and (C, D), yet have no partitions in common. As a consequence, the partition difference (= 107 Robinson-Foulds distance) suffers four essential shortcomings [21]. Firstly, it is imprecise; the 108 number of unique values that the metric can take is two fewer than the number of taxa. (Simply 109 put, a precise method can allocate distinct difference values to two trees that an imprecise 110 method would assign an identical score.) Secondly, it is rapidly saturated; relatively small 111 differences can result in the maximum distance value. Thirdly, its value can be counterintuitive; 112 for example, moving a single tip to a particular location can generate a higher difference value 113 than moving both that tip and its immediate neighbour to the same point (Supplementary Text). 114 Fourthly, balanced trees contain proportionally more uneven partitions, and thus attract lower 115 average distances than asymmetric trees (Supplementary Text). 116 Quartets, in contrast, completely represent all topological information within a tree. The 117 quartet dissimilarity measure is precise, does not rapidly reach saturation, generates a meaningful 118 value for random trees, is robust to the placement of wildcard taxa, and consistently increases in 119 value as trees become more different; and every quartet represents an equal quantity of 120 information. I consider it to represent a more useful, meaningful and interpretable indicator of 121 tree similarity. 122 I calculated quartet distances using the tqDist algorithm [24] via the QuartetStatus 123 function in the new R package Quartet [25]. Partition distances were calculated using the

124 Quartet function SplitStatus. To summarise results, *s*, *d*, and *r* were calculated for each

individual tree relative to the generative tree, and the mean of each of parameter was calculated

126 at each resolution in each analysis.

Comparing Phylogenetic Methods

Smith

127 Previous studies (e.g. [5,6]) have plotted unnormalized symmetric difference against 128 resolution. The unnormalized symmetric difference, however, is a function of both resolution 129 and accuracy: a change in resolution (x) necessarily influences the value, and the range of 130 possible values, of the symmetric difference (y). Because the axes are not independent, this is 131 analogous to plotting x against y/x; the inherent correlation between the axes makes it difficult 132 to interpret the relative contributions of x and y to the plotted function. I instead plotted the 133 proportion of quartets or partitions that are the same in both trees (s), different in both trees (d), 134 and only resolved in the generative tree (r) on ternary plots using the Ternary R package [26], 135 oriented such that SD/MaxI decreases vertically, and resolution decreases horizontally (Figure 136 1a). This plotting configuration distinguishes the relative contributions of resolution and 137 accuracy to overall similarity (Figure 1b).

138

Data, scripts and analyses used in this study are archived on GitHub [27,28].

139 Results

Ideally, measures of node support would assign incorrect nodes low support values. With the CL
datasets (55 characters, 22 tips), resampling methods accomplished this more effectively than
Bremer support (Figure 1c,d), a metric that has attracted criticism [29,30]. The groups
contradicted/supported (GC) metric outperformed group frequency (as anticipated by [31]),
whereas bootstrap resampling outperformed the jackknife approach (contra [32]); subsequent
analyses thus employed the bBootstrap GC metric. Differences between methods were not
statistically significant (Supplementary Text).

With the CL datasets, there is no significant difference (at p = 0.01) between the MaxInormalized quartet symmetric difference of the best trees generated by the Mk model or implied weights ($k \in \{2, 3, 5, 10\}$) – but the best trees generated by equal weights, implied weights with k = 1, and the consensus of k values are significantly worse than those produced by the other methods (Figure 2a; Supplementary Text). Submitted to Biology Letters

7

Comparing Phylogenetic Methods

Smith

152	Collapsing the least-supported groups initially increases the overall accuracy (as
153	predicted in [2,33]), leading to a slight increase in the overall informativeness of the tree (Figure
154	2a,b). Beyond a GC score of c15, the gain in accuracy no longer offsets the resolution lost;
155	collapsing further groups thus removes 'correct' information and reduces the similarity between
156	the tree and the reference tree. Indeed, the optimal tree is only perfectly resolved in a minority of
157	cases (CL, 18%; OR: < 0.2%). Because a Bayesian approach results in less resolution, its most
158	resolved trees cannot generally be improved by collapsing groups (Figures 1c,d, 2).
159	These results hold even if the (problematic) partition difference metric is employed
160	(Figure 1b), though relatively more groups must be collapsed (those with a GC score of < 10) to
161	maximise this metric. The results do not meaningfully change when datasets with low
162	consistency indices are excluded.
163	Similar results are observed in the OR datasets (Figure 2c-e): at any given level of
164	resolution, the best trees obtained by the Mk model are similar in accuracy to those obtained
165	under implied weights (except with very small values of k), but are more accurate than those
166	obtained using equal weights.
167	These datasets also demonstrate the impact of dataset size on tree quality. With larger
168	ratios of characters to taxa (1000 or 350 characters, 75 tips), all methods produced reasonably
169	accurate, well-resolved trees (Figure 2d-e). With the smallest (100 character) datasets (Figure
170	2c), trees were much more different from the generative tree, and the choice of method
171	influenced results more strongly: the Bayesian approach could obtain substantially less
172	resolution, and implied weights recovered poor trees at low values of k . No existing method can
173	overcome the inherent limitation of a low character to taxon ratio.

174 DISCUSSION

175 When accuracy and resolution are recognized as complementary aspects of information [12],

- 176 parsimony and probabilistic analyses generate equally informative reconstructions of
- 177 evolutionary history in the simulation studies analysed herein. Parsimony results are most

Comparing Phylogenetic Methods

178

Page 8 of 16

179 as Bayesian results if nodes are collapsed until trees exhibit an equal resolution. As an important 180 caveat, parsimony analysis must employ a moderate weighting scheme. At low values of the 181 concavity constant (k < 2, say), implied weights begins to exhibit the undesirable properties of 182 clique analysis, whereas at high values (as $k \to \infty$), it converges to the inferior equally weighted 183 parsimony (Supplementary Text). Each of these extremes yields results that are less accurate 184 and less resolved, making them more different from the generative tree and consequently less 185 informative about evolutionary history; results encountered only under such parameters do not 186 merit biological interpretation.

187 Oute aside from issues with the validity of data simulation protocol [2,3], previous 188 results that favour Bayesian methods over parsimony [5–8,10], or equal weights over implied 189 weights [9], have arisen because accuracy has been considered the sole measure of a method's 190 performance. Future simulation studies should evaluate methods based on normalized tree 191 similarity metrics that reflect the total *information* contained within two trees – a quantity that reflects both resolution and accuracy. In the analyses examined herein, neither Bayesian nor 192 193 parsimony analyses generate consistently superior results. Of course, other factors may 194 influence a researcher's choice of methods: Bayesian models, for instance, can readily integrate 195 non-morphological data [34,35] and allow probabilistic hypothesis testing using Bayes Factors 196 [36]. Such considerations notwithstanding, researchers may wish to explicitly compare the 197 results of both Bayesian and implied weights analyses when conducting phylogenetic analysis; 198 observations common to both approaches and receiving strong node support values are 199 particularly likely to be well supported by underlying data.

200 FIGURE LEGENDS

201 Figure 1. Method selection. (a), normalizing symmetric difference against the total information 202 present in two trees (SD/TIP, dotted dashed lines) scores a completely incorrect bifurcating tree

Comparing Phylogenetic Methods

Smith

203 (all relationships resolved differently; bottom corner) no worse than a polytomy (all relationships 204 unresolved; rightmost corner). Random trees (coloured line) with more relationships resolved 205 receive better scores, as some relationships will by chance be resolved correctly. Normalizing 206 against the maximum possible relationship information (SD/MaxI, solid lines) penalizes 207 misinformation over non-information; random trees with more relationships resolved (which thus 208 contain more misinformation) consequently receive worse scores. (b), four measures of tree 209 quality. (c-df), impact on tree quality when least-supported groups are collapsed: (c-d), counting 210 quartets; (e-fd), counting partitions.

211 Figure 2. Status of quartets and bipartitions in trees recovered from simulated datasets.

Points denote the average number of quartets (a, $e\underline{d}-\underline{i}e$) or partitions (b_c) that are the same as the generative tree, resolved differently to the generative tree, or not resolved. Each series indicates the effect of progressively collapsing the least-supported groups in trees generated by analysis of CL (a-bc) and OR datasets (de, g, 100; e, hd, 350; f, ie, 1000 characters) under the specified analytical parameters. The vertical direction corresponds to similarity (i.e. more informative trees); the horizontal direction corresponds to resolution.

218 References

- Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermiin, L. S.
 200 2017 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods*14, 587–589. (doi:10.1038/nmeth.4285)
- 222 2 Goloboff, P. A., Torres, A. & Arias, J. S. 2018 Weighted parsimony outperforms other
 223 methods of phylogenetic inference under models appropriate for morphology. *Cladistics*224 34, 407–437. (doi:10.1111/cla.12205)
- 225 3 Goloboff, P. A., Torres Galvis, A. & Arias, J. S. 2018 Parsimony and model-based
- 226 phylogenetic methods for morphological data: comments on O'Reilly et al. Palaeontology
- 227 **61**, 625–630. (doi:10.1111/pala.12353)

Smith

228	4	Wright, A. M. & Hillis, D. M. 2014 Bayesian analysis using a simple likelihood model
229		outperforms parsimony for estimation of phylogeny from discrete morphological data.
230		PLoS One 9, e109210. (doi:10.1371/journal.pone.0109210)
231	5	O'Reilly, J. E., Puttick, M. N., Parry, L., Tanner, A. R., Tarver, J. E., Fleming, J., Pisani,
232		D. & Donoghue, P. C. J. 2016 Bayesian methods outperform parsimony but at the expense
233		of precision in the estimation of phylogeny from discrete morphological data. Biol. Lett.
234		12, 20160081. (doi:10.1098/rsbl.2016.0081)
235	6	Puttick, M. N., O'Reilly, J. E., Pisani, D. & Donoghue, P. C. J. 2019 Probabilistic
236		methods outperform parsimony in the phylogenetic analysis of data simulated without a
237		probabilistic model. Palaeontology 62, 1–17. (doi:10.1111/pala.12388)
238	7	Puttick, M. N. et al. 2017 Uncertain-tree: discriminating among competing approaches to
239		the phylogenetic analysis of phenotype data. Proc. R. Soc. B 284, 20162290.
240		(doi:10.1098/rspb.2016.2290)
241	8	O'Reilly, J. E., Puttick, M. N., Pisani, D. & Donoghue, P. C. J. 2018 Probabilistic
242		methods surpass parsimony when assessing clade support in phylogenetic analyses of
243		discrete morphological data. Palaeontology 61, 105–118. (doi:10.1111/pala.12330)
244	9	Congreve, C. R. & Lamsdell, J. C. 2016 Implied weighting and its utility in
245		palaeontological datasets: a study using modelled phylogenetic matrices. Palaeontology
246		59 , 447–462. (doi:10.1111/pala.12236)
247	10	Puttick, M. N. et al. 2017 Parsimony and maximum-likelihood phylogenetic analyses of
248		morphology do not generally integrate uncertainty in inferring evolutionary history: a
249		response to Brown et al. Proc. R. Soc. B 284, 20171636. (doi:10.1098/rspb.2017.1636)
250	11	Brown, J. W., Parins-Fukuchi, C., Stull, G. W., Vargas, O. M. & Smith, S. A. 2017
251		Bayesian and likelihood phylogenetic reconstructions of morphological traits are not
252		discordant when taking uncertainty into consideration: a comment on Puttick et al. Proc.
253		R. Soc. B 284, 20170986. (doi:10.1098/rspb.2017.0986)
254	12	Mackay, D. M. 1953 Quantal aspects of scientific information. IEEE Trans. Inf. Theory 1,

255		60-80. (doi:10.1109/TIT.1953.1188569)
256	13	Goloboff, P. A., Farris, J. S. & Nixon, K. C. 2008 TNT, a free program for phylogenetic
257		analysis. Cladistics 24, 774–786. (doi:10.1111/j.1096-0031.2008.00217.x)
258	14	Goloboff, P. A. 1993 Estimating character weights during tree search. Cladistics 9, 83-91.
259		(doi:10.1111/j.1096-0031.1993.tb00209.x)
260	15	Wilkinson, M. 1994 Three-taxon statements: when is a parsimony analysis also a clique
261		analysis? Cladistics 10, 221-223. (doi:10.1111/j.1096-0031.1994.tb00174.x)
262	16	Huelsenbeck, J. P. & Ronquist, F. 2001 MRBAYES: Bayesian inference of phylogenetic
263		trees. Bioinformatics 17, 754–755. (doi:10.1093/bioinformatics/17.8.754)
264	17	Gaudesi, M., Squillero, G. & Tonda, A. 2014 Universal information distance for genetic
265		programming. Proc. 2014 Conf. companion Genet. Evol. Comput. companion - GECCO
266		Comp '14, 137-138. (doi:10.1145/2598394.2598440)
267	18	Estabrook, G. F., McMorris, F. R. & Meacham, C. A. 1985 Comparison of undirected
268		phylogenetic trees based on subtrees of four evolutionary units. Syst. Zool. 34, 193–200.
269		(doi:10.2307/sysbio/34.2.193)
270	19	Day, W. H. E. 1986 Analysis of quartet dissimilarity measures between undirected
271		phylogenetic trees. Syst. Biol. 35, 325-333. (doi:10.1093/sysbio/35.3.325)
272	20	Bandelt, H. J. & Dress, A. 1986 Reconstructing the shape of a tree from observed
273		dissimilarity data. Adv. Appl. Math. 7, 309-343. (doi:10.1016/0196-8858(86)90038-2)
274	21	Steel, M. A. & Penny, D. 1993 Distributions of tree comparison metrics—some new
275		results. Syst. Biol. 42, 126-141. (doi:10.1093/sysbio/42.2.126)
276	22	Penny, D. & Hendy, M. 1985 The use of tree comparison metrics. Syst. Zool. 34, 75-82.
277		(doi:10.2307/2413347)
278	23	Robinson, D. F. & Foulds, L. R. 1981 Comparison of phylogenetic trees. Math. Biosci. 53,
279		131–147. (doi:10.1016/0025-5564(81)90043-2)
280	24	Sand, A., Holt, M. K., Johansen, J., Brodal, G. S., Mailund, T. & Pedersen, C. N. S. 2014
281		tqDist: a library for computing the quartet and triplet distances between binary or general

Page 12 of 16

282 trees. Bioinformatics 30, 2079–2080. (doi:10.1093/bioinformatics/btu157) 283 25 Smith, M. R. 2019 Quartet: comparison of phylogenetic trees using quartet and bipartition 284 measures. Compr. R Arch. Netw. (doi:10.5281/zenodo.2536318) 285 26 Smith, M. R. 2017 Ternary: an R package to generate ternary plots. Compr. R Arch. Netw. 286 (doi:10.5281/zenodo.1068997). 287 27 Smith, M. R. 2019 Distance metrics for trees generated by Congreve and Lamsdell (2016). 288 github.com/ms609/CongreveLamsdell2016, doi:10.5281/zenodo.2536874. 289 28 Smith, M. R. 2019 Distance metrics for trees generated by O'Reilly et al. (2016). 290 github.com/ms609/OReillyEtAl2016, doi:10.5281/zenodo.2536935. 291 29 Wilkinson, M., Thorley, J. L. & Upchurch, P. 2000 A chain is no stronger than its weakest 292 link: double decay analysis of phylogenetic hypotheses. Syst. Biol. 49, 754–776. 293 (doi:10.1080/106351500750049815) 294 30 DeBry, R. W. 2001 Improving interpretation of the decay index for DNA sequence data. 295 Syst. Biol. 50, 742–752. (doi:10.1080/106351501753328866) 296 31 Goloboff, P. A., Farris, J. S., Källersjö, M., Oxelman, B., Ramírez, M. J. & Szumik, C. A. 297 2003 Improvements to resampling measures of group support. *Cladistics* **19**, 324–332. 298 (doi:10.1016/S0748-3007(03)00060-4) 299 32 Kopuchian, C. & Ramírez, M. J. 2010 Behaviour of resampling methods under different 300 weighting schemes, measures and variable resampling strengths. *Cladistics* 26, 86–97. 301 (doi:10.1111/j.1096-0031.2009.00269.x) 302 33 Goloboff, P. A. & Szumik, C. A. 2015 Identifying unstable taxa: efficient implementation 303 of triplet-based measures of stability, and comparison with Phyutility and RogueNaRok. 304 *Mol. Phylo. Evol.* **88**, 93–104. (doi:10.1016/j.ympev.2015.04.003) 305 Lee, M. S. Y., Soubrier, J. & Edgecombe, G. D. 2013 Rates of phenotypic and genomic 34 306 evolution during the Cambrian explosion. Curr. Biol. 23, 1889–1895. 307 (doi:10.1016/j.cub.2013.07.055)

308 35 Zhang, C., Stadler, T., Klopfstein, S., Heath, T. A. & Ronquist, F. 2016 Total-evidence

Comparing Phylogenetic Methods

- dating under the fossilized birth-death process. *Syst. Biol.* **65**, 228–249.
- 310 (doi:10.1093/sysbio/syv080)
- 311 36 Kass, R. E. & Raftery, A. E. 1995 Bayes factors. J. Am. Stat. Assoc. 90, 773–795.
- 312 (doi:10.1080/01621459.1995.10476572)
- 313

314 ACKNOWLEDGEMENTS

- 315 The TNT software is supported by the Willi Hennig Society. Detailed comments from two
- 316 anonymous referees substantially improved the manuscript.
- 317 AUTHOR CONTRIBUTIONS
- 318 MS is the sole author.
- 319 DATA ACCESSIBILITY
- 320 Data and analyses are available on GitHub and archived with Zenodo. CL:
- 321 <u>https://github.com/ms609/CongreveLamsdell2016</u> [27]; OR:
- 322 <u>https://github.com/ms609/OReillyEtA12016</u> [28].
- 323 Electronic Supplementary Material accompanies the article:
- Supplementary Text (PDF, Rmd): includes a detailed comparison of the suitability of
- 325 tree comparison metrics, and discusses the use of small concavity constants in implied
- 326 weighting. The Rmd file provides the source R code used in analyses, to enable the
- 327 reproduction of research results. These files accompany the online article at [publisher to
- 328 provide URL]
- 329 FUNDING
- 330 None to report.

- 331 COMPETING INTERESTS
- 332 None to report.
- **333** ETHICAL STATEMENT
- 334 No ethical approval was required to conduct this research.

Smith



Figure 1. Method selection. (a), normalizing symmetric difference against the total information present in two trees (SD/TIP, dashed lines) scores a completely incorrect bifurcating tree (all relationships resolved differently; bottom corner) no worse than a polytomy (all relationships unresolved; rightmost corner). Random trees (coloured line) with more relationships resolved receive better scores, as some relationships will by chance be resolved correctly. Normalizing against the maximum possible relationship information (SD/MaxI, solid lines) penalizes misinformation over non-information; random trees with more relationships resolved (which thus contain more misinformation) consequently receive worse scores. (b), four measures of tree quality. (c-f), impact on tree quality when least-supported groups are collapsed: (c-d), counting quartets; (e-f), counting partitions.

174x262mm (300 x 300 DPI)



Figure 2. Status of quartets and bipartitions in trees recovered from simulated datasets. Points denote the average number of quartets (a, d-i) or partitions (b-c) that are the same as the generative tree, resolved differently to the generative tree, or not resolved. Each series indicates the effect of progressively collapsing the least-supported groups in trees generated by analysis of CL (a-c) and OR datasets (d, g, 100; e, h, 350; f, i, 1000 characters) under the specified analytical parameters. The vertical direction corresponds to similarity (i.e. more informative trees); the horizontal direction corresponds to resolution.

175x170mm (300 x 300 DPI)