

## Accuracy in near-perfect virus phylogenies

JOEL O. WERTHEIM<sup>1,\*</sup>, MIKE STEEL<sup>2</sup>, AND MICHAEL J. SANDERSON<sup>3</sup>

<sup>1</sup> *Department of Medicine, University of California San Diego, La Jolla, CA 92093, USA*

<sup>2</sup> *Biomathematics Research Center, School of Mathematics and Statistics, University of Canterbury, Christchurch, 8041, New Zealand*

<sup>3</sup> *Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ 85721 USA*

*\*To whom correspondence should be addressed: [jwertheim@health.ucsd.edu](mailto:jwertheim@health.ucsd.edu); [sanderm@email.arizona.edu](mailto:sanderm@email.arizona.edu)*

## ABSTRACT

1 Phylogenetic trees from real-world data often include short edges with very few  
2 substitutions per site, which can lead to partially resolved trees and poor accuracy. Theory  
3 indicates that the number of sites needed to accurately reconstruct a fully resolved tree  
4 grows at a rate proportional to the inverse square of the length of the shortest edge.  
5 However, when inferred trees are partially resolved due to short edges, “accuracy” should  
6 be defined as the rate of discovering false splits (clades on a rooted tree) relative to the  
7 actual number found. Thus, accuracy can be high even if short edges are common.  
8 Specifically, in a “near-perfect” parameter space in which trees are large, the tree length  $\xi$   
9 (the sum of all edge lengths) is small, and rate variation is minimal, the expected false  
10 positive rate is less than  $\xi/3$ ; the exact value depends on tree shape and sequence length.  
11 This expected false positive rate is far below the false negative rate for small  $\xi$  and often  
12 well below 5% even when some assumptions are relaxed. We show this result analytically  
13 for maximum parsimony and explore its extension to maximum likelihood using theory  
14 and simulations. For hypothesis testing, we show that measures of split “support” that rely  
15 on bootstrap resampling consistently imply weaker support than that implied by the false  
16 positive rates in near-perfect trees. The near-perfect parameter space closely fits several  
17 empirical studies of human virus diversification during outbreaks and epidemics, including  
18 Ebolavirus, Zika virus, and SARS-CoV-2, reflecting low substitution rates relative to high  
19 transmission/sampling rates in these viruses.

20 *Key words:* Perfect phylogeny, Homoplasy, Yule-Harding model, Virus, Epidemic,  
21 SARS-CoV-2, Ebolavirus, Mumps virus, HIV, West Nile virus, Zika virus

## INTRODUCTION

A “perfect phylogeny” is an evolutionary tree constructed from discrete character data in which no character state evolves more than once (Gusfield, 1997; Fernandez-Baca and Lagergren, 2003). Homoplasy (Wake et al., 2011) is absent. Real-world datasets rarely allow reconstruction of perfect phylogenies, but algorithms can be modified to search efficiently for “near-perfect” trees when a small amount of homoplasy is present (Fernandez-Baca and Lagergren, 2003; Awasthi et al., 2012). In this paper, we address how best to measure accuracy in such “near-perfect” trees, what factors guarantee accuracy is high, and whether real datasets with such minimal levels of homoplasy even exist.

The concept of perfect and near-perfect phylogenies played a key role in early attempts to understand the connections among phylogenetic tree reconstruction methods, such as maximum likelihood (ML), maximum parsimony (MP), and maximum compatibility. In a landmark paper, Felsenstein (1973) showed that a sufficient condition for ML and MP to infer the same tree was for the expected number of substitutions on edges of the tree to be very small. Then, “[i]f our assumption were true that evolutionary change is improbable during the relevant period of time, most characters should be uniform over the group. A few would show a single change of state during the evolution of the group. But only very rarely would we find more than one change of state, so that few or no characters would show convergence.” This last statement may have been the first hint of a probabilistic description of “near-perfect phylogeny”. This condition can be stated more formally as  $\xi \leq 1$ , where  $\xi$  is the expected number of substitutions per site summed over the entire tree (i.e., the tree length per site). Homoplasy is rare but has a non-zero probability of occurring.

Felsenstein’s concluding comment on near-perfect phylogenies was skeptical: “Real data is certainly not like this...” (Felsenstein, 1973). Homoplasy has since been viewed as a commonplace feature of phylogenetic datasets (Wake et al., 2011) and, reasonably enough, most phylogenetic theory has been developed with this sentiment as an implicit

49 assumption. However, extensive surveys of genetic diversity in RNA viruses have revealed  
50 that some viral phylogenies, particularly those associated with outbreaks and epidemics,  
51 do exhibit small per site total tree lengths consistent with near-perfect phylogenies (Dudas  
52 and Bedford, 2019). These datasets often comprise full-length viral genomes from RNA  
53 viruses, which are typically 10–30 kb in length and have a substitution rate of around  $10^{-3}$   
54 substitutions/site/year.

55 The potential of these data to yield fully resolved phylogenies has been of particular  
56 interest in epidemiology, because internal nodes in viral trees represent transmission events  
57 (Campbell et al., 2018; Grubaugh et al., 2019; Dudas and Bedford, 2019). This objective  
58 motivates placing a premium on minimizing false negatives (i.e., on deciphering all such  
59 transmission events) and thereby maximizing resolution. Increased phylogenetic resolution  
60 is achievable by analyzing longer genomic fragments from viruses with faster evolutionary  
61 rates (Dudas and Bedford, 2019). However, understanding the false positive rate remains a  
62 key issue in characterizing phylogenetic accuracy (Felsenstein and Kishino, 1993),  
63 particularly in the special case of a poorly resolved tree with few—but  
64 well-supported—clades.

65 Here we explore what assumptions comprise “near-perfect” phylogenies and  
66 decouple the false-positive and false-negative components of accuracy in such trees. In  
67 particular, by focusing on a mathematically tractable case in which tree size is large yet  
68 tree length is small, we will show that the false positive rate can be very good, even when  
69 the false negative rate is not: most of the clades inferred are probably correct, even though  
70 the tree may be only partly resolved. We also survey a set of viral phylogenies that have  
71 many properties of this near-perfect space and estimate their accuracy. Finally, we briefly  
72 consider phylogenetic “support” measures in relation to accuracy in near-perfect data.  
73 Whereas accuracy relates to the overall performance of a tree estimator relative to the true  
74 tree, support relates to the probability of making a mistake in deciding about some aspect  
75 of that tree—typically the presence of a particular split—using a statistically based

76 decision rule such as the bootstrap support value or a posterior probability (Felsenstein,  
77 1985; Felsenstein and Kishino, 1993; Hillis and Bull, 1993; Efron et al., 1996; Susko, 2008,  
78 2009; Alfaro and Holder, 2006; Simmons and Norton, 2014).

79 This paper is organized as follows. “Materials and Methods” are divided into two  
80 parts: first, mathematical theory (with proofs in the Supplement), and second, simulation  
81 protocols, data, and data analysis. “Results” begin with a more expository description of  
82 the theory, illustrated with simulation results, and then describes results from analyses of  
83 robustness and support, and data analyses. Following these is the Discussion.

## 84 MATERIALS AND METHODS I. THEORY

### 85 *Definitions of Accuracy*

86 Given a true unrooted binary tree,  $T$ , and an estimated tree,  $\hat{T}$ , a strict measure of  
87 accuracy is just  $\text{Prob}(\hat{T} = T)$  (Huelsenbeck and Hillis, 1993; Erdős et al., 1999). In large  
88 trees it is useful to measure partial agreement, such as the proportion of nontrivial splits  
89 on  $\hat{T}$  that are also on  $T$ , out of a possible  $n - 3$  (Yang, 1998).

90 A still more nuanced definition of accuracy is useful when either  $T$  or  $\hat{T}$  is only  
91 partially resolved (not binary), that is, when the number of nontrivial splits,  $C(T)$ , is less  
92 than  $n - 3$  (Warnow, 2013). Let  $N_{FP}$  be the number of splits on  $\hat{T}$  but not  $T$  (false  
93 positives), and let  $N_{FN}$  be the number of splits on  $T$  but not  $\hat{T}$  (false negatives). When  
94 both trees are binary,  $N_{FP} = N_{FN}$  (Berry and Gascuel, 1996; Smirnov and Warnow, 2021);  
95 otherwise they can contribute differentially to error. The Robinson–Foulds (RF) distance  
96 (Robinson and Foulds, 1981),  $d_{RF} = N_{FP} + N_{FN}$ , combines both errors in one measure of  
97 overall accuracy. Here we distinguish between these errors explicitly by defining false  
98 positive and negative rates (Smirnov and Warnow, 2021):

$$\begin{aligned} FP_T &= \mathbb{E}[N_{FP}/C(\hat{T})], \\ FN_T &= \mathbb{E}[N_{FN}/C(T)]. \end{aligned} \tag{1}$$

99 Both error rates are expectations over some generating model for the data, described next.

*Evolutionary Model*

Let  $B(n)$  denote the set of unrooted binary phylogenetic trees with leaf set  $[n] = \{1, 2, \dots, n\}$ . Note that a tree  $T \in B(n)$  has  $2n - 3$  edges. Consider a Jukes-Cantor model (JC69; Felsenstein, 2004), with rate parameter  $\lambda$ , in which the probability of a state change between the endpoints of an edge  $e$ , denoted  $p_e$ , is given by  $p_e = p$ , where  $p = \frac{3}{4}(1 - \exp(-4\lambda/3))$ . Assume further that all edges have the same value of  $\lambda$ . Let  $\xi$  denote the expected number of state changes per character in  $T$ . Thus  $\xi = \lambda \cdot (2n - 3)$ .

A *character* refers to the assignment of states to the taxa at a given site of an alignment.

We will say that a character evolves ‘perfectly’ on  $T$  if there is a single change of state across one interior edge (say  $e$ ) and no change of state on any other edge of  $T$ . Thus, a character that evolves perfectly on  $T$  is homoplasy-free, and the two notions are equivalent for binary characters. However, for multi-state characters, the notion of a perfectly evolved character is stronger than that of being merely homoplasy-free. We deal here with this stronger notion for two reasons: firstly, it simplifies the mathematical analysis, and second, the expected proportion of homoplasy-free characters that are not perfectly evolved under the models we consider tends to zero as the number of taxa becomes large.

We will say that a character  $f$  evolves on  $T$  with  $c$  edge changes on  $e_1, \dots, e_c$  if state changes occur on edges  $e_1, \dots, e_c$  and on no other edge of  $T$ . More briefly, we say that  $f$  evolves on  $T$  with  $c$  edge changes if  $f$  evolves with  $c$  edge changes for some set of  $c$  distinct edges of  $T$  (mostly we will deal with the case  $c = 2$ ).

Recall that a *split* refers to a bipartition of the leaf set  $[n]$  into two nonempty subsets (and splits are induced by binary characters). A character that has evolved perfectly on  $T$  produces a split, and these splits (across a set of perfectly evolved characters) are compatible and so form a (generally unresolved/non-binary) tree on leaf set  $[n]$ .

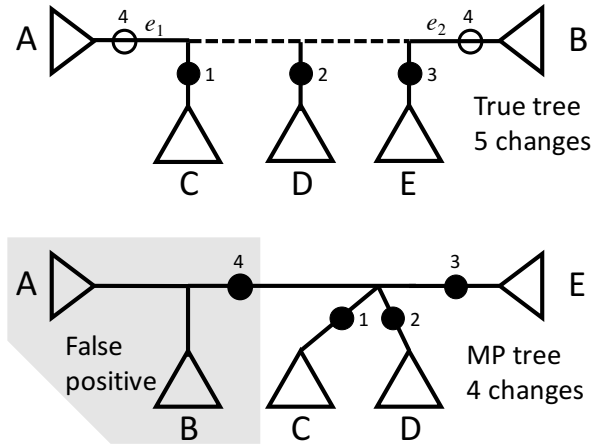


Fig. 1. How a false positive split is inferred by maximum parsimony (MP). On true tree (top) sites 1–3 are binary and “perfect”; that is, they have only a single change (locations marked by black circles), but site 4 is binary and homoplastic, changing twice (open circles), on edges  $e_1$  and  $e_2$ . The dotted line is the path between the two homoplastic changes in site 4. If no perfect sites change along the dotted line path on the true tree, a false positive split is inferred on the MP tree (bottom).

*Probability of False Splits*

125

126 Suppose that  $m$  characters evolve on  $T$  and that, of these  $m$  characters,  $k$  of them  
 127 are perfectly evolved on  $T$  (note that more than one of these characters may correspond to  
 128 the same split of  $T$ ). Next, consider a single additional character  $f$  which has evolved on  $T$   
 129 with 2 edge changes, on  $e_1, e_2$  (there is no restriction that these must be interior edges).  
 130 Under certain conditions, the MP tree for these characters will include a false split (false  
 131 positive)—a split not on  $T$  (Fig. 1). In particular, a false split occurs if no perfect  
 132 character changes state along the path between  $e_1$  and  $e_2$  (see Lemma 1 in the  
 133 Supplementary Information).

134 Let  $\Phi_T^{(k)}$  be the probability that a character  $f$  that has evolved on  $T$  with 2 edge  
 135 changes generates a false split under MP, which means:

136 (C-i) it is a binary character,

137 (C-ii) the corresponding split is not a split of  $T$ , and

138 (C-iii) the split described by  $f$  is compatible with  $k$  characters that are perfectly evolved  
 139 on  $T$  (by the Markovian process described above).

140 In other words, we are interested in ‘false splits’ (i.e., splits in the reconstructed MP tree  
 141 that are not present in the—underlying and unknown—true tree  $T$ ). The split  
 142 corresponding to  $f$  (by condition C-i) should *not* be in  $T$  (condition C-ii); however,  
 143 condition C-iii would lead MP to add this false split into the reconstructed tree based on  
 144 the other ‘true splits’ since the false split is compatible with all of the latter.

145 Given a tree  $T \in B(n)$ , let  $d_T(e_1, e_2)$  denote the number of edges of  $T$  that lie  
 146 strictly within the path between  $e_1$  and  $e_2$  (i.e., excluding  $e_1$  and  $e_2$ ). Thus,  $e_1$  and  $e_2$  are  
 147 adjacent if and only if  $d_T(e_1, e_2) = 0$ . In addition, let  $\varphi_T = (\varphi_T(0), \varphi_T(1), \dots, \varphi_T(n-3))$ ,  
 148 where  $\varphi_T(i)$  is the number of (unordered) pairs of edges  $\{e, e'\}$  of  $T$  for which  $d_T(e, e') = i$ .  
 149 Finally, for  $i$  between 1 and  $n-3$ , let

$$\tilde{\varphi}_T(i) = \frac{\varphi_T(i)}{\binom{2n-3}{2}}. \quad (2)$$

150 The probability of a false split is then given by the following theorem (see SI for  
 151 proof).

**Theorem 1** For each  $T \in B(n)$ , and  $k \geq 1$  we have:

$$\Phi_T^{(k)} = \frac{1}{3} \cdot \sum_{i=1}^{n-3} \tilde{\varphi}_T(i) \left(1 - \frac{i}{(n-3)}\right)^k.$$

152 Theorem 1 shows that for fixed  $k$  and  $n$ , the shape of  $T$  plays a significant role in  
 153 determining  $\Phi_T^{(k)}$ ; in particular, unbalanced trees (such as caterpillars) will have a smaller  
 154 value of  $\Phi_T^{(k)}$  than more balanced trees. Indeed, it is possible to calculate the value of  $\Phi_T^{(k)}$   
 155 exactly for the two extreme cases of caterpillar trees and fully-balanced trees to determine  
 156 the extent of this dependence (see SI).



*Estimating the Expected False Positive Rate*

Given a binary phylogenetic tree  $T$ , and  $m$  characters evolved randomly on  $T$  by the model described earlier, the *false positive rate* ( $FP_T$ ) is the expected value of the ratio of false splits to all splits in the estimated tree (Eqn. 1; here we assume that if the reconstructed tree is a star, this proportion [which is technically  $0/0$ ] is zero). Recall that  $\xi$  is the expected number of state changes in the tree  $T$  per character, under the model described earlier.  $FP_T$  is a function of the three parameters  $T$  (specifically, its shape and number of leaves),  $m$ , and  $\lambda$  (equivalently,  $FP_T$  is a function of  $T$ ,  $m$ , and  $\xi$ ).

In general, it is mathematically complicated to describe  $FP_T$  in terms of these parameters. However, when the number of leaves in a tree grows faster than the number of perfectly compatible characters, it is possible to state a limit result to provide an approximation to  $FP_T$  for large trees.

In the following theorem, we consider the following setting:

(I)  $m\xi = \Theta(n^\beta)$  for some  $0 < \beta < \frac{1}{2}$ , and

(II)  $m\xi^2 = O(1)$ ,

where  $O(1)$  refers to dependence on  $n$  (thus  $m\xi^2$  is not growing with  $n$ ). Note that Condition (I) implies that the number of perfectly evolved characters grows with the number of leaves, but at a rate that is slower than linearly. Conditions (I) and (II) imply that  $\xi$  decreases as  $n$  increases.

In this setting, we show that the false positive rate is (asymptotically) of the form  $\frac{\xi}{3}$  times a function  $\Omega$  that involves  $T$  (via its shape),  $m$ , and  $\xi$ . If we now treat  $\xi$  as a variable, then for  $\xi = 0$ , the function  $\Omega$  is close to 1 (for large  $n$ ) and so  $FP_T$  initially grows like  $\xi/3$ . However, as  $\xi$  increases,  $\Omega$  begins to decline at an increasing rate, resulting in the false positive rate reaching a maximum value before starting to decrease.

To describe this result, we need to define this function  $\Omega$ . Let

$$\Omega(T_n, \xi, m) = \sum_{i=1}^{n-4} \tilde{\varphi}_{T_n}(i) \cdot \frac{e^{-i\mu/(n-3)} - e^{-\mu}}{1 - i/(n-3)},$$

where:

$$\mu = \frac{1}{2}m\xi$$

181 and where  $\tilde{\varphi}_{T_n}(i)$  is given in Eqn. (2). For example, for any caterpillar tree, we have  
 182  $\tilde{\varphi}_{T_n}(i) = 4(n-2-i)/\binom{2n-3}{2}$ .

183 Notice that  $\Omega(T_n, \xi, m)$  depends on  $T_n$  only via the coefficients  $\tilde{\varphi}_{T_n}(i)$ , and this  
 184 dependence is linear. Thus, if  $\mathcal{D}$  is a distribution on trees (e.g. the PDA or YH), then the  
 185 expected value of  $\Omega(T_n, \xi, m)$  is given by:

$$\mathbb{E}_{\mathcal{D}}[\Omega(T_n, \xi, m)] = \sum_{i=1}^{n-4} \mathbb{E}_{\mathcal{D}}[\tilde{\varphi}_{T_n}(i)] \cdot \frac{e^{-i\mu/(n-3)} - e^{-\mu}}{1 - i/(n-3)}. \quad (3)$$

186 For the PDA distribution, the term  $\mathbb{E}_{PDA}[\tilde{\varphi}_{T_n}(i)]$  has an explicit exact value,  
 187 namely,

$$\mathbb{E}_{PDA}[\tilde{\varphi}_{T_n}(i)] = \frac{(i+3)2^i(2n-i-4)!(n-2)!}{(2n-4)!(n-i-3)!\binom{2n-3}{2}}, \quad (4)$$

188 for all  $i$  between 1 and  $n-3$  (see SI for proof).

189 **Theorem 2** For each  $n \geq 1$ , let  $T_n$  be a binary phylogenetic tree with  $n$  leaves, and  
 190 suppose that Conditions (I) and (II) hold.

(i)

$$FP_{T_n} = \frac{\xi}{3} \cdot \Omega(T_n, \xi, m) \cdot (1 + o(1)),$$

191 where  $o(1)$  is a term that tends to 0 as  $n$  grows.

(ii) If  $T_n$  is sampled from a distribution  $\mathcal{D}$  (e.g. PDA, YH), then the expected value of  
 $FP_{T_n}$ , denoted  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ , satisfies

$$\mathbb{E}_{\mathcal{D}}[FP_{T_n}] = \frac{\xi}{3} \cdot \mathbb{E}_{\mathcal{D}}[\Omega(T_n, \xi, m)] \cdot (1 + o(1)).$$

192 **Remarks:** Note that  $FP_{T_n}$  depends only on the shape of the tree  $T_n$  (and not on  
 193 how its leaves are labelled), thus for a tree distribution  $\mathcal{D}$  on either the class of caterpillar  
 194 trees, or symmetric trees, we have  $FP_{T_n} = \mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ .

Notice also from Fig. 3 that as  $\xi$  increases from 0 the estimate of  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  given by  $\frac{\xi}{3} \cdot \Omega(T_n, \xi, m)$  for the YH, PDA distributions and for symmetric trees initially increases (approximately linearly) with  $\xi$  but then begins to decrease with increasing  $\xi$ . By contrast, when  $T_n$  has the caterpillar tree shape, the estimate of  $FP_{T_n}$  appears to be constant as  $\xi$  increases from 0 (see Fig. 3). Indeed, when  $T_n$  is a caterpillar tree, the expression for  $FP_{T_n}$  in Theorem 2(i) reduces to the following remarkably simple expression as  $n$  becomes large:

$$FP_{T_n} \sim 4/(3m),$$

195 which is independent of  $\xi$  (and  $n$ ). Details are provided in the SI.

## 196 MATERIALS AND METHODS II. SIMULATIONS, DATA, AND DATA ANALYSES

### 197 *Main Simulation Pipeline*

198 Simulations were run to assess goodness of fit and robustness of mathematical  
 199 predictions under various regimes of model parameters and tree inference criteria (MP or  
 200 ML), as well as to estimate expected accuracy in empirical datasets. Each of  $R$  simulation  
 201 replicates (with  $r$  sub-replicate tree searches in each) consisted of the following sequence of  
 202 steps: (i) generation of a random binary tree  $T$  with  $n$  leaves according to either a  
 203 “proportional-to-distinguishable-arrangement” (PDA) or Yule-Harding (YH) model  
 204 (Aldous, 2001) (as well as the two extreme cases of completely unbalanced caterpillar  
 205 trees, and completely balanced symmetric trees); (ii) assignment of edge lengths of  $T$   
 206 according to a gamma distribution with shape parameter  $\alpha_e$  and mean  $\bar{\lambda}$ ; (iii) generation  
 207 of a sequence alignment of  $m$  sites using Seq-Gen v. 1.3.4 (Rambaut and Grassly, 1997),  
 208 with either JC69, HKY, or GTR models (and base frequencies and rate matrix parameters  
 209 set or estimated from data), and with one of four across-site-rate (ASR) variation models:  
 210 no variation, invariant sites model, gamma model, or free-rate model ranging from 2-10  
 211 bins (Kalyaanamoorthy et al., 2017)—the free-rate model was implemented in Seq-Gen by  
 212 using 2-10 site partitions; (iv) reconstruction of estimated tree  $\hat{T}$  [using PAUP 4.0a, build

213 166: Swofford (2003)) for MP with options ‘hsearch add=simple swap=no nreps= $r$ ;contree  
214 all/strict’; and using IQ-TREE 2 (v. 2.0.6) (Minh et al., 2020) for ML with options ‘-m  
215 JC+FQ -nt 1 -redo -mredo -polytomy -blmin 1e-9’, replicated  $r$  times, followed by strict  
216 consensus]; (v) tallying  $N_{FP}$  and  $N_{FN}$  from  $T$  and  $\hat{T}$  and computing error rates. Mean  
217 rates across replicates were then tallied. All steps except (iii) and (iv) used custom PERL  
218 scripts (available in the Dryad repository).

219 A typical dataset size of  $n = 513$  (chosen to allow perfectly symmetrical trees plus  
220 one outgroup, when such were needed), and  $m = 1000$  was used to model trees large  
221 enough to potentially satisfy the near-perfect assumptions, and to have a sufficient number  
222 of sites to infer a range of accuracy when combined with  $\bar{\lambda}$  values ranging from  $10^{-5}$  to  
223 0.316 substitutions per site. Gamma shape parameters were set at 0.1, 1.0, and 10.0, which  
224 encompasses distributions ranging from highly variable to nearly constant. For edge length  
225 variation this range encompasses what we observed in the empirical virus datasets. For  
226 ASR variation, it captures much of the range of inferred values we have seen in the  
227 literature. Finally,  $R$  was generally set to 1000 and  $r$  to 100.

### 228 *Support Simulations*

229 Phylogenetic support measures were estimated in trees simulated via the main  
230 pipeline described above with  $n = 513$ ,  $m = 1000$ , a JC69 model with no rate variation,  
231 and PDA random trees. Ten values of  $\lambda$  in the interval  $[10^{-5}, 0.31622]$  were analyzed.  
232 PAUP (Swofford, 2003) was used for MP bootstrapping (same heuristic search as above  
233 but with 100 replicates  $\times$  10 subreplicates); IQ-TREE 2 (Minh et al., 2020) was used (50  
234 random tree replicates) for SH-aLRT (‘-alrt 1000’), aBayes, and ultrafast bootstrapping  
235 (‘-B 1000’), with additional options enforcing minimum branch lengths of  $10^{-9}$  and  
236 collapsed polytomies. Mean support across replicates was computed.

237 Perfect four-taxon alignments were generated in which each of the five branches had  
238 a single, non-homoplastic nucleotide substitution in the alignment and all other sites were

239 constant. Alignment lengths ranged between 40 nt and 30,000 nt. ML trees were inferred in  
240 IQTree2 with a JC69 model, minimum branch lengths of  $10^{-9}$ , and collapsed polytomies.  
241 Clade support was determined using Felsenstein's bootstrapping (1,000 replicates),  
242 ultrafast bootstrapping (10,000 replicates), transfer bootstrap exchange (TBE; 1,000  
243 replicates), SH-aLRT (10,000 replicates), and aBayes. Full Bayesian inference was also  
244 performed in MrBayes v3.2.7 (Ronquist and Huelsenbeck, 2003) with a single run per  
245 replicate of 2.5 million generations, with the first 10% of generations discarded as burnin.

246 Alignments for larger perfect symmetrical and asymmetrical (caterpillar) trees were  
247 generated with 8, 16, 32, 64, and 128 taxa. Each branch, including terminal branches, had  
248 a single nonhomoplastic nucleotide substitution in the alignment with all other sites  
249 constant. Alignment lengths ranged from 236 to 32,768 nt. ML trees were inferred as  
250 described above for the four-taxon alignments, and support was assessed by Felsenstein's  
251 bootstrap, ultrafast bootstrapping, TBE, SH-aLRT, and aBayes.

252 All Python scripts related to perfect tree simulations are available in the Dryad  
253 repository.

### 254 *Virus Datasets*

255 Viral phylogenies were obtained from the NextStrain (Hadfield et al., 2018) website  
256 (accessed 05 May 2020) (Table 1). Phylograms were downloaded for dengue virus, dengue  
257 virus serotype 1, Ebolavirus (Dudas et al., 2017), Enterovirus 68 (Dyrdak et al., 2019),  
258 measles morbillivirus, mumps virus, respiratory syncytial virus, West Nile virus (Hadfield  
259 et al., 2019), and Zika virus. In addition, we also analyzed an iatrogenic HIV-1 outbreak in  
260 Cambodia (Rouet et al., 2018) and the first wave of the SARS-CoV-2 epidemic in China  
261 (Pekar et al., 2021). The SARS-CoV-2 phylogeny is the ML tree used in Pekar et al. (2021)  
262 (see Data S1 on Dryad [<https://doi.org/10.6076/D12S3M>] for list of GISAID Accession  
263 IDs). Publicly available genomic sequences (or genetic sequences for HIV-1) were  
264 downloaded from GenBank and aligned with mafft v7.407 (Kato and Standley, 2013)

(accession numbers can be found in Data S2 on Dryad).

False positive rates for the virus phylogenies were estimated with our simulation pipeline, setting parameters to values estimated from published trees and publicly available sequences used to construct them (Table 1, Table S1). For each virus, we used IQ-TREE 2 to infer the six rate parameters of a GTR substitution model with empirical base frequencies. The optimal site-to-site rate variation model, including free-rate models, was determined using the Bayesian Information Criterion (BIC) in IQ-TREE 2 (Kalyaanamoorthy et al., 2017). These models were used to parameterize sequence simulation in Seq-Gen, as described above.

Edge length (per site) variation was assumed to follow a gamma distribution:  $\lambda \sim \Gamma(\alpha_e, \alpha_e/\bar{\lambda})$  having mean  $\bar{\lambda}$  and variance  $\bar{\lambda}^2/\alpha_e$ . The distribution of substitutions is a mixture of Poisson and gamma distributions, which is a negative binomial with a variance to mean ratio of

$$1 + \frac{m\bar{\lambda}}{\alpha_e} \quad (5)$$

which was shown by Bedford and Hartl (2008) for an equivalent parameterization. Virus trees were preprocessed, setting any edge lengths  $< 1.1 \times 10^{-6}$  to zero, assuming these reflected ML numeric artifacts. Then,  $\bar{\lambda}$  was estimated from the observed sum of per site edge lengths divided by  $2n - 3$ , and Eqn. 5 was then used to estimate  $\alpha_e$ .

Ideally, we would fit the data to the random tree model, but standard methods either assume binary trees or model polytomies with an a priori assumption about the tree model itself (e.g., Bortolussi et al., 2006). Therefore, we repeated simulations using both PDA and YH models.

## RESULTS

*Overview of Results on Accuracy*

Simulations of tree inference with MP over a large range of tree lengths,  $\xi$ , and other parameters, illustrate several known results (Fig. 2) and perhaps a few less well known ones. First, resolution of the inferred tree increases with tree length. Second, “overall” accuracy, as measured by the RF distance, is optimal at an intermediate tree length,  $\xi^*$  (Yang, 1998; Bininda-Emonds et al., 2001; Steel and Leuenberger, 2017). Moreover, when  $\xi \gg \xi^*$ , the false positive error rate,  $FP_T$ , is similar to the false negative rate,  $FN_T$ , as might be expected because the true and estimated trees are nearly binary; therefore  $N_{FP} \cong N_{FN}$ .

However, when  $\xi \ll \xi^*$ , then  $FP_T \ll FN_T$ , and the false positive error rate can remain quite good ( $< 0.05$ ) over a large range of  $\xi$  even when the false negative error rate is very high. However, the range of tree lengths for which this result holds depends critically on rate variation across edges and sites. When  $\xi \leq 1$ , the false positive rate is low and insensitive to the presence of rate variation; but, when  $\xi > 1$ , the false positive rate is much more sensitive to rate variation—high when variation is present and low when absent (contrast Fig. 2A and Fig. 2B). In real-world data, as  $\xi$  increases, we expect that evidence of rate variation will become more apparent.

Key elements of these findings can be shown analytically in a “near-perfect” zone described by a simple evolutionary model.

*Overview of the Mathematical Theory*

First we define “near-perfect” more formally. Assume the data consist of an alignment of  $m$  independent and identically distributed nucleotide sites that have evolved according to a Jukes-Cantor model (JC69; Felsenstein, 2004) on an unrooted binary tree  $T$ , with  $n$  leaves. Each of the  $2n - 3$  edges of  $T$  have length  $\lambda$ , and thus the total tree length is  $\xi = \lambda(2n - 3)$ . When  $n$  is large and  $\xi \leq 1$ , the expected number of substitutions

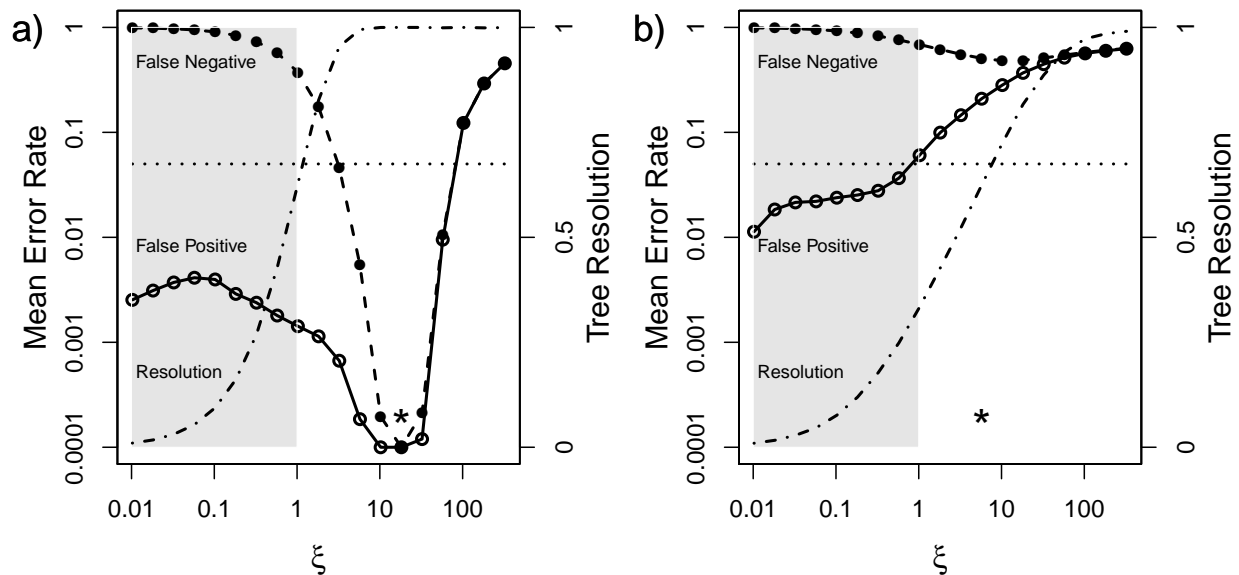


Fig. 2. Accuracy of maximum parsimony phylogeny reconstruction in simulations over a wide range of per site tree length,  $\xi$ , and other parameters. Solid and dashed curves are mean false positive and negative error rates, respectively (log scale left); dashed sigmoidal curve is fractional resolution of estimated tree (linear scale right). Trees are generated by a random proportional-to-distinguishable-arrangement (PDA) model for 513 taxa, from which a sequence alignment length of 1000 sites is generated. The dotted horizontal line is placed at an error rate of 0.05. Asterisk marks the location of the optimal tree length with best overall Robinson–Foulds accuracy,  $\xi^*$ . Each point is mean of 1000 replicates  $\times$  100 sub-replicates (see Methods). “Near-perfect” values ( $\xi \leq 1.0$ ) are shaded. a) JC69 model with no edge length or across-site-rate variation [because of  $y$ -axis log scaling, two  $y$  values of zero were set to 0.0001]. b) JC69 model with substantial edge length and across-site-rate variation, both modeled as a gamma distribution with shape parameters  $\alpha_e = \alpha_{ASR} = 0.25$ ).

312 per site is  $\leq 1$ ; the number of edges on which a site changes state is approximately Poisson  
 313 distributed with mean  $\xi$ ; and the probability of more than one change on an edge is low,  
 314 meaning multiple changes at a site occur on distinct edges. Though these conditions will  
 315 generate alignments dominated by “perfect” sites exhibiting no homoplasy, a few sites may  
 316 exhibit homoplasy even with  $\xi \leq 1$ , which motivates the term “near- perfect”. Under these  
 317 conditions, tree reconstruction methods will tend to infer relatively unresolved trees unless  
 318 the number of sites is very large.

Rare sites that exhibit homoplasy can introduce false positive splits on the inferred tree (Fig. 1). A naïve argument using Equation 1 might suggest that  $FP_T$  would depend on  $\xi$  roughly as  $O(\xi^2)/O(\xi) = O(\xi)$ , namely the ratio of the expected numbers of sites having changes on two edges (i.e., those that are potentially homoplastic and misleading) to those



sites having only a single change (those that are reliable), for sufficiently small  $\xi$ . But because only one-third of those two-edge sites are actually homoplastic in a JC69 model,

$$FP_T \cong \xi/3,$$

319 which implies  $FP_T$  is small when  $\xi$  is small enough (e.g.,  $FP_T < 0.05$  whenever  $\xi < 0.15$ ).

320 This approximation can be improved further by recognizing that not all two-edge  
321 homoplastic sites induce false positives, depending on their position in the true tree (Fig.  
322 1). Given the evolutionary model, the probability that  $k$  perfect sites, and another site  $f$   
323 that has evolved with two edge changes will produce a “false positive” under MP is  
324 denoted  $\Phi_T^{(k)}$  (Theorem 1 above). Because this probability is often less than one,  $FP_T$  can  
325 remain below 0.05 at higher values of  $\xi$  than the naïve argument suggests.

326 If the true tree were known with some precision, the first part of Theorem 2 could  
327 be used directly to calculate false positive rates. However, in the “near-perfect” parameter  
328 space of large  $n$  and  $\xi \leq 1$ , estimates of the true tree are likely to be only partially resolved  
329 (Fig. 2). We therefore derive the expected false positive rate for a distribution,  $\mathcal{D}$ , of  
330 randomly generated trees of size  $n$ ,  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ , generated from parameters based on the  
331 inferred tree. In the remainder of this paper, the “expected false positive rate” will  
332 generally refer to  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ . We assume that  $\mathcal{D}$  is usually either a  
333 “proportional-to-distinguishable-arrangement” (PDA) or Yule-Harding (YH) distribution  
334 (Aldous, 2001), but also consider the two extreme cases of completely unbalanced  
335 (caterpillar) trees, and completely balanced (symmetric) trees. Unlike PDA and YH trees,  
336 these last two have a constant tree shape (with random leaf labels). From the second part  
337 of Theorem 2, we see that, for a JC69 model and trees inferred with MP, the following  
338 approximation holds increasingly well as  $n$  increases:

$$\mathbb{E}_{\mathcal{D}}[FP_{T_n}] \cong \frac{\xi}{3} \cdot \mathbb{E}_{\mathcal{D}}[\Omega(T_n, \xi, m)] \quad (6)$$

339 given the assumption that  $\xi$  is sufficiently small and the number of sites does not grow too  
340 quickly with the size of the tree. The function  $\Omega(T_n, \xi, m)$ , defined in Materials and

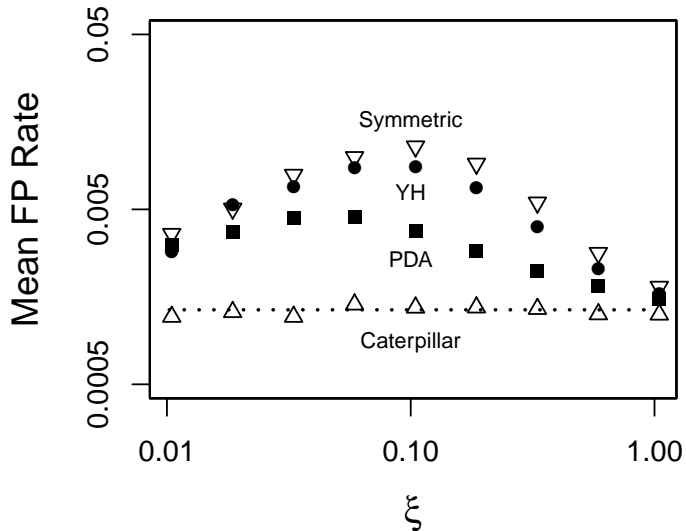


Fig. 3. Mean false positive rate in four tree models. Fit to theoretical predictions from Equation 6 (or the limit expression of  $4/3m$  for caterpillar trees: see Methods) are shown by dashed lines. Each point is mean of 1000 replicates  $\times$  100 sub-replicates. Simulation conditions were  $n = 513$ ,  $m = 1000$ , with a JC69 model. Predicted values are not known for YH model.

341 Methods I, is monotonically decreasing in  $\xi$  and  $m$ , and depends on the shape of  $T$ .  
 342 Simulations indicate that the approximation is close for  $\xi \leq 1$  (Fig. 3), but if many equally  
 343 parsimonious trees are present, the search algorithm should take a strict consensus of a  
 344 broad sample of those solutions (Fig. S3).  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  is better on average for PDA than YH  
 345 trees, and both are bounded between a theoretical worst case error rate for symmetric and  
 346 best case error rate for caterpillar trees. In fact, the expected false positive rate for the  
 347 latter is just  $4/(3m)$  in the limit of large  $n$ , which is independent of  $\xi$ .

### 348 *Robustness to Violation of Assumptions*

349 Violations of assumptions tend to increase the expected false positive rate above the  
 350 predictions of Equation 6. For example, adding edge length (EL) variation or  
 351 across-site-rate (ASR) variation increases  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  (Figs. 2, 4 and Fig. S4). The difference

352 between predicted  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  based on Eqn. (6), with no edge length variation, and  
 353 simulation-based estimates with edge length variation included is small when  $\xi \ll 1$  but  
 354 increases substantially as  $\xi$  increases. When edge length variation is large (gamma shape  
 355 parameter  $\alpha_e = 0.1$ ), there is no longer a local maximum value of  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  around  
 356  $\xi = 0.1$ ; instead,  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  increases monotonically with  $\xi$  and eventually exceeds 5% for  
 357 the simulated dataset sizes. The impact of ASR variation is deleterious at all values of  $\xi$ ,  
 358 but even when ASR variation is large (gamma shape parameter  $\alpha_{ASR} = 0.1$ ), the false  
 359 positive rate remains slightly below 5% for simulated dataset sizes in the absence of EL  
 360 variation (Fig. S4).

361 Departure of the substitution model from the JC69 model assumed in the  
 362 “near-perfect” zone can also increase the expected false positive rate. For example, a  
 363 strong transition–transversion bias increases  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  substantially, though it still remains  
 364 well below 5% under our typical simulation conditions when  $\xi \leq 1$  (Fig. S5).

365 Thus, the near-perfect tree length of  $\xi \leq 1$  is a region in which rate variation  
 366 appears to have less of an impact on false positive rates than when tree lengths are longer.  
 367 This suggests that the definition of near-perfect zone in practice can include substantial  
 368 rate variation.

### 369 *Expected False Positive Rates in Virus Phylogenies*

370 We estimated key parameters from the trees and underlying data for 11 empirical  
 371 virus phylogenies (Table 1, Table S1) and used simulation to estimate expected false  
 372 positive rates (Figs. 5). The studies span a wide range of tree size and resolution and  
 373 alignment length, and their tree lengths span three orders of magnitude. Seven of these  
 374 viruses fell within the “near-perfect” tree length zone of  $\xi \leq 1.0$ , and six of those had  
 375  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}] \leq 0.05$  irrespective of random tree model.  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  was generally lower for PDA  
 376 vs. YH models. As expected,  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  increased roughly with  $\xi$ , despite the large  
 377 differences in these datasets.

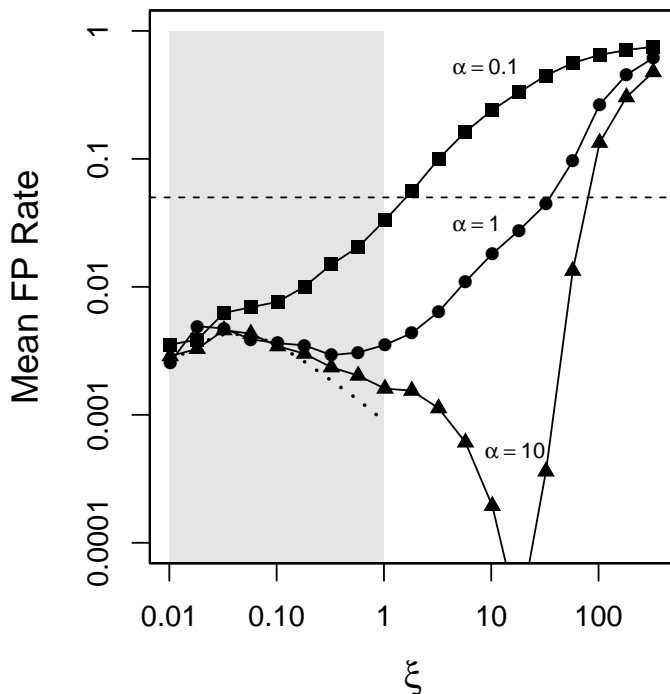


Fig. 4. Effect of edge length variation on expected false positive rate for different values of the shape parameter of the edge length gamma distribution,  $\alpha_e$ . Smaller values of  $\alpha_e$  correspond to higher rate variation. ASR variation is assumed absent. The dashed curve is the prediction from Eqn. (6), in which both sources of variation are absent. Simulation conditions assumed PDA trees with  $n = 513$ ,  $m = 1000$ , 1000 replicates, 100 subreplicates. Gray rectangle shows “near-perfect” values of  $\xi \leq 1$ .

378 Epidemics with young crown group ages on the order of years or decades (e.g., Zika  
 379 virus, West Nile virus, and mumps virus) had expected false positive rates below 5%, even  
 380 though West Nile virus had a  $\xi$  slightly above 1. Viruses encompassing single epidemics  
 381 (e.g., SARS-CoV-2 in China, EBOV in West Africa, and HIV-1 in Cambodia) also had  
 382 expected false positive rates below 5%. Remarkably, HIV-1 had a low expected false  
 383 positive rate even though the tree was constructed using the fewest number of sites in our  
 384 sample (from only a single partial gene). Number of site affects accuracy through the  
 385  $\Omega(T_n, \xi, m)$  term in Eqn. 6.

386 Trees with lowest levels of resolution (Table 1) had the highest expected false

Table 1. Parameters of 11 empirical virus phylogenies

| Abbreviation | Study   | Leaves | Sites | Resolution |
|--------------|---|--------|-------|------------|
| DENV         | Dengue virus                                    | 1197   | 10264 | 0.8795     |
| DENV-1       | Dengue virus serotype 1                         | 1067   | 10264 | 0.8160     |
| EBOV         | Ebolavirus                                      | 1610   | 18164 | 0.3632     |
| EV-D68       | Enterovirus 68                                  | 824    | 7293  | 0.8029     |
| HIV-1        | Human immunodeficiency virus type 1             | 189    | 1038  | 0.2193     |
| MeV          | Measles morbillivirus                           | 109    | 15782 | 0.7009     |
| MuV          | Mumps virus                                     | 458    | 15154 | 0.2961     |
| RSV          | Respiratory syncytial virus                     | 997    | 14986 | 0.6121     |
| SARS-CoV-2   | Severe acute respiratory syndrome coronavirus 2 | 583    | 29668 | 0.2324     |
| WNV          | West Nile virus                                 | 2512   | 10395 | 0.5960     |
| ZIKV         | Zika virus                                      | 543    | 10320 | 0.5453     |

387 positive rates. For example, dengue virus serotype 1, which does not represent a single  
 388 epidemic, had low phylogenetic resolution, a  $\xi > 1$ , and a correspondingly high expected  
 389 false positive rate. The phylogenetically more diverse dengue virus tree representing all  
 390 four DENV serotypes had an even higher tree length and expected false positive rate.

391 The measles virus tree was an outlier with  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  above 5%, even though its tree  
 392 lengths was below one. Notable, MeV had the fewest taxa of any virus analyzed (Table 1)  
 393 and subsequently lower phylogenetic resolution. This combination of factors implies  
 394 sensitivity to the assumption of large  $n$  in our results.

### 395 *Extension to Maximum Likelihood (ML) Inference*

396 Theoretical results hint that ML and MP should reconstruct the same tree under  
 397 “near-perfect” assumptions. For example, ML provably converges to MP when there are  
 398 enough constant characters in an alignment, a condition similar to  $\xi \ll 1$  (Tuffley and  
 399 Steel, 1997, Thm. 3). Further arguments presented in the SI support this conjecture.

400 We used simulation to check how well Equation 6, derived for MP, predicted the  
 401 expected false positive rate under ML inference in the near-perfect zone. Simulations with

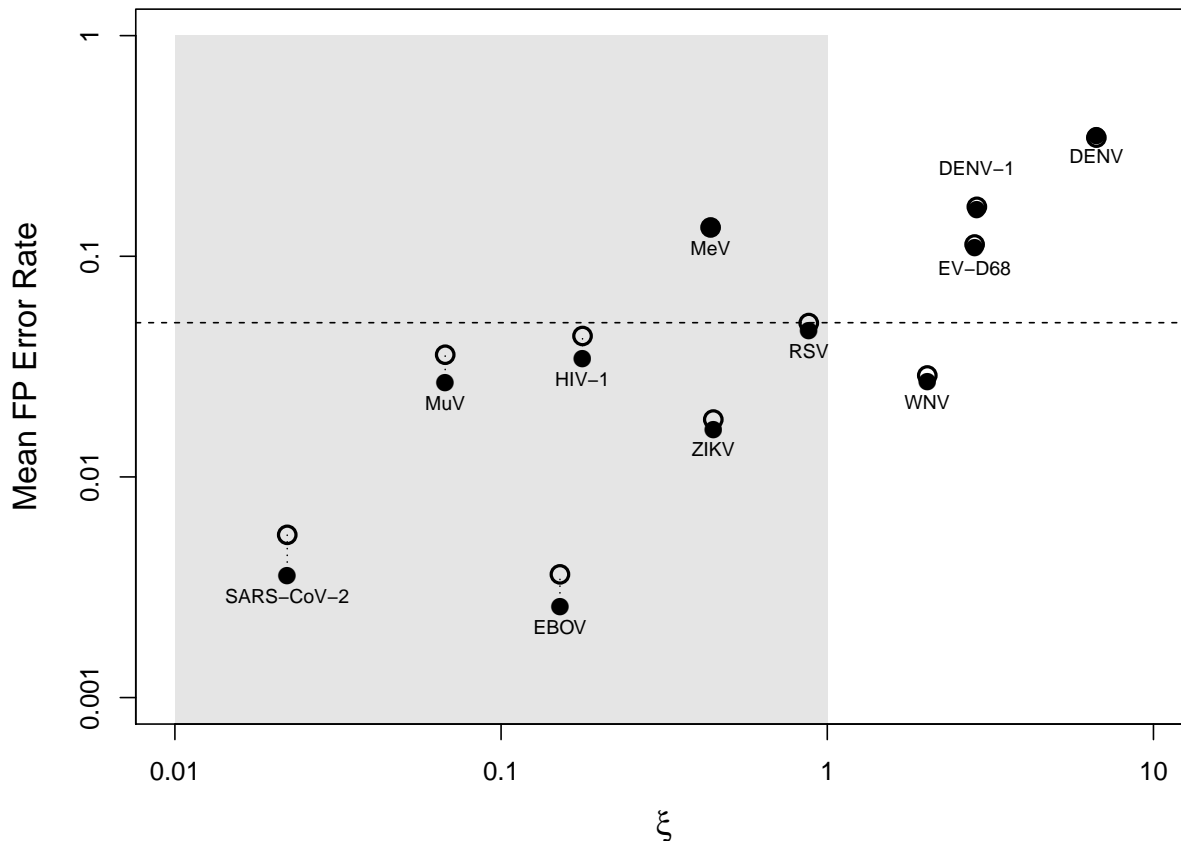


Fig. 5. Expected false positive rates,  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ , for 11 empirical virus phylogenetic datasets (Table 1) for maximum parsimony (MP) inference, estimated by simulation using parameters estimated from the data (Table S1). Abbreviations given in Table 1. Simulation experiments used either a Yule–Harding random tree distribution (open circles) or PDA distribution (closed circles: some data points have indistinguishable differences between random tree models). Each point is mean of 500 replicates  $\times$  100 sub-replicates. The near-perfect zone of  $\xi \leq 1.0$  is shaded. Horizontal dashed line indicates a 0.05 expected false positive rate.

402  $\xi \leq 1$ , a JC69 model, and no edge length or ASR variation, with trees inferred by  
 403 IQ-TREE 2 (Minh et al., 2020) under the same model, are close to the equation’s  
 404 predictions (Fig. S6). Nonetheless, some differences were observed, which tended to imply  
 405 better accuracy for MP. These differences could largely be attributed to technical or  
 406 implementation issues in ML software. First, the computational expense of ML searches  
 407 makes it tempting to undertake fewer replicate searches for local optima, but this was as  
 408 critical to improve the fit to Equation 6 for ML as it was for MP (Fig. S6). Second, ML  
 409 programs set hard numerical lower bounds strictly greater than zero on edge lengths, often

410 (by default) on the same order as  $\bar{\lambda}$  for the virus datasets, so these must be reset downward  
 411 to obtain correct tree likelihoods (Morel et al., 2021). Finally, inferred edge lengths that  
 412 are larger than these programs’ lower bounds but still smaller than about  $1/m$  tend to be  
 413 included in the ML tree despite weak evidence (IQ-TREE 2 issues a warning about this).  
 414 We saw this in ML searches roughly when  $\xi \geq 0.1$ , when three-state sites become more  
 415 common in alignments than they were at lower values of  $\xi$ . Even without homoplasy, ML  
 416 tends to over-resolve trees in a way that elevates  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ . By collapsing short edge  
 417 lengths inferred by ML to be less than  $1/m$ , this behavior can be mitigated (Fig. S6).

418 In general, ML is expected to be more accurate than MP under more realistic  
 419 model conditions and higher rates, something we observed commonly in simulations in  
 420 which  $\xi > \xi^*$ . However, simulations also suggest that in the near-perfect zone, MP can  
 421 achieve an  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  comparable with ML but with much faster running times.

### 422 *Accuracy and Support in Near-perfect and Perfect Trees*

423 False positive “accuracy”, defined as  $1 - \mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ , is very high in the near-perfect  
 424 zone of small tree lengths, whereas conventional support values are quite variable in this  
 425 zone under the same simulation conditions (Fig. 6). At very low  $\xi$ , the average bootstrap  
 426 support for MP is about the theoretically expected 64% for a single nonhomoplastic  
 427 substitution supporting an edge (Felsenstein, 1985). Model-based support measures had  
 428 higher values, with aBayes (Anisimova et al., 2011) being greater than ultrafast bootstrap  
 429 (Hoang et al., 2018), which, in turn, was greater than SH-aLRT (Guindon et al., 2010),  
 430 but only aBayes was close to our  $1 - \mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  false positive accuracy across the range of  
 431 tree lengths in the near-perfect zone. Notably, aBayes is the only one of the metrics that is  
 432 not based on resampling.

433 We explored other factors impacting support in the boundary case of perfect trees.  
 434 For sequence length, we computed standard support metrics in an ML framework in  
 435 perfect four-taxon datasets, in which each branch was defined by a single change, and

436 alignments range between 40 nt and 30,000 nt (Fig. S7). As observed for MP, Felsenstein’s  
437 ML bootstrap support is approximately 63%, regardless of sequence length, in accordance  
438 with theoretical predictions (Felsenstein, 1985). Transfer bootstrap exchange (TBE)  
439 (Lemoine et al., 2018; Lutteropp et al., 2020) values were indistinguishable from  
440 Felsenstein’s bootstrap. Of the other ML model-based support metrics, aBayes provided  
441 higher values than ultrafast bootstrap and SH-aLRT, both of which rely on bootstrap  
442 resampling. The aBayes support reached  $\geq 95\%$  for alignments as short as 100 nt, which  
443 tracked the full Bayesian posterior support estimates that had support  $\geq 95\%$  in  
444 alignments as short as 60 nt. The discrepancy between the Bayesian estimates and those  
445 that use bootstrap resampling, in light of our other results, suggests that resampling  
446 methods used in the presence of splits defined by only a single informative site may fail to  
447 integrate relevant information about low tree lengths.

448 On the other hand, in perfect trees from 8–128 taxa, in which the mean edge length  
449 remained the same (but therefore  $\xi$  grew with  $n$ ), mean SH-aLRT and aBayes support was  
450 unchanged, but mean ultrafast bootstrap support increased (Fig. S8). The TBE method  
451 was developed to correct for a downward bias of bootstrap values often seen in large trees.  
452 As expected, TBE exceeds conventional bootstrap support as taxon number increases.  
453 However, this increase is modest in perfectly symmetrical trees compared with perfectly  
454 asymmetric trees and only surpasses 95% in the largest asymmetric trees (Fig. S8).

## 455 DISCUSSION

456 In this paper, we study a “near-perfect” parameter space for phylogenetic inference  
457 on large trees with small tree lengths and no rate variation within or between sites or  
458 edges. The “near-perfect” tree length of  $\xi \leq 1$  means that few sites exhibit homoplasy and,  
459 for MP inference, the false positive rate can be much better than the false negative rate  
460 and well under 5% for typical datasets with thousands of sites. The near-perfect conditions  
461 defined here to allow mathematical derivations appear to be sufficient but not necessary.



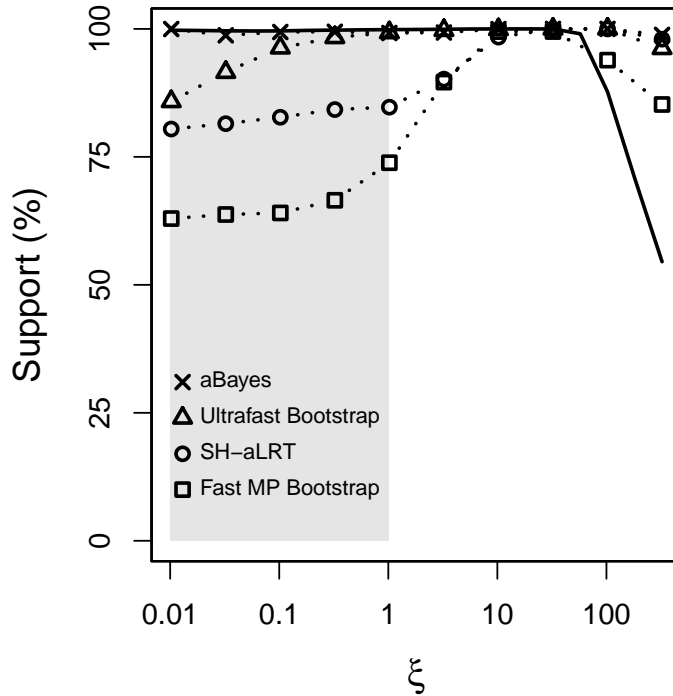


Fig. 6. Statistical support measures compared to expected false positive accuracy, as a function of tree length. The solid curve is the mean value of  $(1 - \mathbb{E}_{\mathcal{D}}[FP_{T_n}]) \times 100$  in simulations. The near-perfect parameter space is shaded.

462 For example, with no rate variation, the false positive rate can be very good even when  
 463  $\xi > 1$  (Fig. 2A, S5), and, if  $\xi < 1$ , a substantial level of rate variation can be present  
 464 without elevating the false positive rate by nearly as much as when  $\xi > 1$  (Fig. 2,4, S4).

465 The second case is clearly more relevant in real-world data. The 11 empirical virus  
 466 datasets all had substantial rate variation and showed a general increase in false positive  
 467 rate with  $\xi$ , with almost all rates below 5% occurring when  $\xi \leq 1$ , much like the predicted  
 468 patterns seen in Fig. 2B and Fig. 4. This observation accords with our simulation results  
 469 suggesting that the good “near-perfect” false positive rates may emerge even when relaxing  
 470 the strict near-perfect assumption of no rate variation—as long as  $\xi \leq 1$ .

471 These and many other empirical findings about RNA virus phylogenies sampled  
 472 intensively in epidemics postdate much of the extensive body of other work on accuracy

473 and support in phylogenetics. Not surprisingly, little note has been made about the stark  
474 contrast between false positive and false negative rates in phylogenies in which tree length  
475 is well below the optimal tree length for “overall accuracy”, since published examples have  
476 been relatively rare. The goal of much of the field of phylogenetics is, after all, to maximize  
477 tree resolution, even if this effort requires adding (or switching to) sequence data with  
478 more variation and thus longer tree lengths.

479 Because “near-perfect” datasets reflect a combination of the number of taxa and  
480 sites, evolutionary rate and time parameters, and assumptions about the substitution  
481 model, they also implicitly reflect sampling of the true tree, which is particularly relevant  
482 in epidemic trees in which sampling is far below disease incidence. Sampling can continue  
483 over time, increasing  $n$ , and the viruses continue to evolve over time, increasing the depth  
484 of the tree. Both of these increase  $\xi$  but in different ways; therefore, it is possible for the  
485 same RNA virus to have near-perfect and not near-perfect datasets depending on the  
486 study. For example, the SARS-Cov2 dataset we included had  $n = 583$  and  $\xi = 0.02$ , well  
487 within the “near-perfect” zone, but a much more intensively sampled tree over a longer  
488 period of time (Lanfear, 2020) with  $n = 147156$  has a tree length of  $\xi = 3.89$  (after  
489 collapsing any edges with  $\lambda \leq 1.1 \times 10^{-6}$ ), which is remarkably small for such a large tree  
490 but lies just outside our definition of near-perfect. This finding suggests that large-scale  
491 phylogenetic approaches for SARS-CoV-2 surveillance are appropriate (Ferreira et al.,  
492 2021; Turakhia et al., 2021) and that such approaches are unlikely to falsely suggest close  
493 relatedness (i.e., transmission clusters) where none exists.

494 Other mathematical results on phylogenetic accuracy have largely focused on either  
495 the limiting case of infinite sequence length (“consistency”), or the number of sites needed  
496 for accurate inference (the “sequence length requirement”). For MP, for example, the  
497 shortest edge length is critical and  $\lim_{m \rightarrow \infty} \text{Prob}(\hat{T}_{MP} = T) = 1$  as long as  
498  $\lambda_{\min} > \xi^2/(1 - \xi)$  (Steel, 2000, Thm. 1(A)). More generally, let  $m'$  be the number of sites  
499 needed for  $\text{Prob}(\hat{T}_{MP} = T)$  to exceed some fixed required accuracy. For the

500 neighbor-joining method  $m'$  grows exponentially with  $n$  (Lacey and Chang, 2006); for ML,  
 501  $m'$  is polynomial or better in  $n$ , depending on edge lengths (Roch and Sly, 2017). Moreover,  
 502  $m'$  also grows as  $O(1/\lambda_{\min}^2)$  for ML and some more ad hoc estimators (Erdős et al., 1999;  
 503 Roch, 2019), implying again that short edges tend to degrade accuracy when accuracy is  
 504 defined in terms of total agreement between  $T$  and  $\hat{T}$ , in contrast to our findings here.

505 A cryptic factor affecting the false positive rate is tree shape. Highly asymmetric  
 506 trees have better expected false positive rates than highly symmetric trees, because  
 507 expected path lengths are longer and it is harder to induce false positive splits by chance  
 508 (Fig. 1). Thus, a random sample of PDA trees will have a better  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$  than more  
 509 symmetrical YH trees. Differences in tree shape among RNA virus phylogenies have long  
 510 been noted (Grenfell et al., 2004), such as the typically more asymmetric influenza trees.

511 Perfect and near-perfect phylogenies have been studied as discrete optimization  
 512 problems (Gusfield, 1997; Fernandez-Baca and Lagergren, 2003) in which the goal is to  
 513 find an optimal tree when, at most, some small number of sites exhibit homoplasy. Little of  
 514 this work has considered accuracy per se, but Gronau et al. (Gronau et al., 2012)  
 515 highlighted the connection between short edge lengths and false positives, and developed a  
 516 “fast converging” algorithm (i.e., having an  $O(\text{poly}(n))$  sequence length requirement) that  
 517 returns a tree with short edges collapsed when they do not meet a threshold probability of  
 518 being correct, thus minimizing false positives. The connection between this tree and those  
 519 built by more conventional methods is unclear, but it may be a promising approach for  
 520 building trees in the near-perfect zone.

521 Model-based phylogenetic inference methods such as ML and Bayesian inference are  
 522 generally regarded as theoretically superior to MP, especially for datasets that fit  
 523 substitution models much more complex than our “near-perfect” JC69 model with no rate  
 524 variation. Though our mathematical results for expected false positive rates were derived  
 525 for MP, there is both relevant theory and considerable simulation evidence to suggest that  
 526 in the near-perfect zone, the ML expected false positive rate is approximated by the MP

theory, both in terms of its absolute value and its shape as a function of tree length. As  $\xi$  increases, especially above  $\xi^*$ , ML consistently has better accuracy than MP, but we conjecture that the false positive rates of MP and ML differ much less as  $\xi$  gets very small. Further work is needed to test this conjecture.

The connection between the false positive rate as a measure of accuracy and conventional measures of phylogenetic support appears to be sensitive to the choice of support method when  $\xi \ll 1$  (Fig. 6). The aBayes method corresponds well to what is implied by  $1 - \mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ , but resampling methods using either likelihood or parsimony correspond less well. The connection between phylogenetic accuracy and support in frequentist and Bayesian settings has been studied in detail (Felsenstein, 1985; Hillis and Bull, 1993; Felsenstein and Kishino, 1993; Efron et al., 1996; Susko, 2008, 2009; Alfaro and Holder, 2006; Simmons and Norton, 2014), but remains somewhat fraught. We hesitate to draw firm conclusions without a formal analysis of support in the “near-perfect” parameter space, but we do note the variability in support estimates we found and suspect that Bayesian measures may be better reflections of false positive accuracy in practice (Fig. 6). If individual clade support needs to be invoked in near-perfect viral phylogenies, we recommend Bayesian approaches that do not rely on bootstrap resampling of sparse substitutions. In near-perfect trees, Bayesian approaches can make use of the limited amount of genetic diversity to draw strongly supported inference, as opposed to bootstrapping approaches which require multiple sites supporting a clade before inferring similarly strong support. When phylogenetically informative data are limited, as in near-perfect trees, the consistency of the data supporting a clade appears more relevant than their prevalence.

The low false positive rate in near-perfect trees suggests that phylogenies describing viral epidemics in this zone can be interpreted directly without defaulting to identifying clades with strong support values. This finding supports the current practice in SARS-CoV-2 nomenclature, whereby clades (e.g., denoting variants or migration events)

554 are defined with reference to specific synapomorphies (Rambaut et al., 2020; O’Toole  
555 et al., 2021; Worobey et al., 2020). We acknowledge that frequent convergent evolution,  
556 and recombination in positive-strand RNA viruses, can complicate phylogenetic inference  
557 and may increase the false positive rate in real-world trees (Morel et al., 2021).

558 The benefit of real-time viral genomic sequencing for public health action became  
559 apparent during the 2014–2015 West African Ebola epidemic (Gire et al., 2014), and is a  
560 critical component of tracking the COVID-19 pandemic (Oude Munnink et al., 2020;  
561 Grubaugh et al., 2021). Consequently, the viruses responsible for these diseases, Ebolavirus  
562 and SARS-CoV-2, epitomize near-perfect phylogenetic trees in our analysis. We can expect  
563 a greater intensity of genomic sequencing accompanying future viral outbreaks, increasing  
564 the importance and relevance of near-perfect phylogenies.

565 In conclusion, we have shown that many RNA virus datasets satisfy assumptions  
566 used to derive results on near-perfect phylogenetic accuracy. These criteria include  
567 sufficiently low substitution rates across a large enough tree and no recombination. Any set  
568 of genomes sampled in a clade on a short enough time scale, or highly conserved regions of  
569 genomes sampled across a deeper clade, can also satisfy the first assumption, but  
570 recombination would remain problematic in many taxa. Springer et al. (2020), illustrate a  
571 potential path forward in their study of “low-homoplasy” retroelement characters in  
572 mammal genomes. They pursue a species tree inference approach to such data, which  
573 would likely be “near-perfect” were it not for recombination. It may be possible to derive  
574 additional results on accuracy when local near-perfect trees (or sub-alignments) are  
575 combined under the multi-species coalescent (Liu et al., 2019).

## 576 ACKNOWLEDGEMENTS

577 We gratefully acknowledge the authors from the originating laboratories and the  
578 submitting laboratories, who generated and shared via GISAID the SARS-CoV-2 genomic  
579 sequence data on which this research is based. A complete list acknowledging the authors

580 who submitted data analyzed in this study can be found in Data S1. MJS thanks the  
581 University of Arizona's HPC facility, Bio5 Institute, and Rod Wing's lab for computing  
582 support. JOW was supported by an NIH-NIAID R01 AI135992.

#### 583 DISCLOSURE STATEMENT

584 The authors have no conflicts of interest related to this work.

#### 585 SUPPLEMENTARY MATERIAL

586 Supplementary material, including data files, scripts, and online-only appendices  
587 containing mathematical proofs and additional figures, can be found in the Dryad data  
588 repository: <https://doi.org/10.6076/D12S3M>

## REFERENCES

589

590 Aldous, D. J. 2001. Stochastic models and descriptive statistics for phylogenetic trees, from  
591 Yule to today. *Stat. Sci.* 16:23–34.

592 Alfaro, M. E. and M. T. Holder. 2006. The posterior and the prior in Bayesian  
593 phylogenetics. *Annu. Rev. Ecol. Evol. Syst.* 37:19–42.

594 Anisimova, M., M. Gil, J. F. Dufayard, C. Dessimoz, and O. Gascuel. 2011. Survey of  
595 branch support methods demonstrates accuracy, power, and robustness of fast  
596 likelihood-based approximation schemes. *Syst. Biol.* 60:685–99.

597 Awasthi, P., A. Blum, J. Morgenstern, and O. Sheffet. 2012. Additive approximation for  
598 near-perfect phylogeny construction. Pages 25–36 *in* Approximation, randomization, and  
599 combinatorial optimization. Algorithms and techniques (M. Goemans, K. Jansen,  
600 J. Rolim, and L. Trevisan, eds.). Springer, Berlin.

601 Bedford, T. and D. L. Hartl. 2008. Overdispersion of the molecular clock: Temporal  
602 variation of gene-specific substitution rates in *Drosophila*. *Mol. Biol. Evol.* 25:1631–1638.

603 Berry, V. and O. Gascuel. 1996. On the interpretation of bootstrap trees: Appropriate  
604 threshold of clade selection and induced gain. *Mol. Bio. Evol.* 13:999–1011.

605 Bininda-Emonds, O. R. P., S. G. Brady, J. Kim, and M. J. Sanderson. 2001. Scaling of  
606 accuracy in extremely large phylogenetic trees. *Pacific Symposium on Biocomputing*  
607 6:547–558.

608 Bortolussi, N., E. Durand, M. Blum, and O. François. 2006. apTreeshape: Statistical  
609 analysis of phylogenetic tree shape. *Bioinformatics* 22:363–364.

610 Campbell, F., C. Strang, N. Ferguson, A. Cori, and T. Jombart. 2018. When are pathogen  
611 genome sequences informative of transmission events? *PLoS Pathog.* 14:e1006885.

612 Dudas, G. and T. Bedford. 2019. The ability of single genes vs full genomes to resolve time  
613 and space in outbreak analysis. *BMC Evol. Biol.* 19:232.

- 614 Dudas, G., L. M. Carvalho, T. Bedford, A. J. Tatem, G. Baele, N. R. Faria, D. J. Park,  
615 J. T. Ladner, A. Arias, D. Asogun, F. Bielejec, S. L. Caddy, M. Cotten, J. D'Ambrozio,  
616 S. Dellicour, A. Di Caro, J. W. Diclaro, S. Duraffour, M. J. Elmore, L. S. Fakoli,  
617 O. Faye, M. L. Gilbert, S. M. Gevao, S. Gire, A. Gladden-Young, A. Gnirke, A. Goba,  
618 D. S. Grant, B. L. Haagmans, J. A. Hiscox, U. Jah, J. R. Kugelman, D. Liu, J. Lu,  
619 C. M. Malboeuf, S. Mate, D. A. Matthews, C. B. Matranga, L. W. Meredith, J. Qu,  
620 J. Quick, S. D. Pas, M. V. T. Phan, G. Pollakis, C. B. Reusken, M. Sanchez-Lockhart,  
621 S. F. Schaffner, J. S. Schieffelin, R. S. Sealfon, E. Simon-Loriere, S. L. Smits,  
622 K. Stoecker, L. Thorne, E. A. Tobin, M. A. Vandi, S. J. Watson, K. West, S. Whitmer,  
623 M. R. Wiley, S. M. Winnicki, S. Wohl, R. Wolfel, N. L. Yozwiak, K. G. Andersen, S. O.  
624 Blyden, F. Bolay, M. W. Carroll, B. Dahn, B. Diallo, P. Formenty, C. Fraser, G. F. Gao,  
625 R. F. Garry, I. Goodfellow, S. Gunther, C. T. Happi, E. C. Holmes, B. Kargbo, S. Keita,  
626 P. Kellam, M. P. G. Koopmans, J. H. Kuhn, N. J. Loman, N. Magassouba, D. Naidoo,  
627 S. T. Nichol, T. Nyenswah, G. Palacios, O. G. Pybus, P. C. Sabeti, A. Sall, U. Stroher,  
628 I. Wurie, M. A. Suchard, P. Lemey, and A. Rambaut. 2017. Virus genomes reveal factors  
629 that spread and sustained the Ebola epidemic. *Nature* 544:309–315.
- 630 Dyrdak, R., M. Mastafa, E. B. Hodcroft, R. A. Neher, and J. Albert. 2019. Intra- and  
631 interpatient evolution of enterovirus D68 analyzed by whole-genome deep sequencing.  
632 *Virus Evol.* 5:vez007.
- 633 Efron, B., E. Halloran, and S. Holmes. 1996. Bootstrap confidence levels for phylogenetic  
634 trees. *Proc. Natl. Acad. Sci. USA* 93:13429–13434.
- 635 Erdős, P. L., M. A. Steel, L. A. Szekely, and T. J. Warnow. 1999. A few logs suffice to  
636 build (almost) all trees (I). *Random Structures and Algorithms* 14:153–184.
- 637 Felsenstein, J. 1973. Maximum likelihood and minimum steps methods for estimating  
638 evolutionary trees from data on discrete characters. *Syst. Zool.* 22:240–249.



- 639 Felsenstein, J. 1985. Confidence limits on phylogenies: An approach using the bootstrap.  
640 *Evolution* 39:783–791.
- 641 Felsenstein, J. 2004. *Inferring phylogenies*. Sinauer Press, Sunderland, MA.
- 642 Felsenstein, J. and H. Kishino. 1993. Is there something wrong with the bootstrap on  
643 phylogenies? A reply to Hillis and Bull. *Syst. Biol.* 42:182–192.
- 644 Fernandez-Baca, D. and J. Lagergren. 2003. A polynomial-time algorithm for near-perfect  
645 phylogeny. *SIAM J. Comput.* 32:1115–1127.
- 646 Ferreira, R.-C., E. Wong, G. Guban, K. Wade, M. Liu, L. M. Baena, C. Chato, B. Lu,  
647 A. S. Olabode, and A. F. Y. Poon. 2021. CoVizu: Rapid analysis and visualization of the  
648 global diversity of SARS-CoV-2 genomes. *bioRxiv* .
- 649 Gire, S. K., A. Goba, K. G. Andersen, R. S. Sealfon, D. J. Park, L. Kanneh, S. Jalloh,  
650 M. Momoh, M. Fullah, G. Dudas, S. Wohl, L. M. Moses, N. L. Yozwiak, S. Winnicki,  
651 C. B. Matranga, C. M. Malboeuf, J. Qu, A. D. Gladden, S. F. Schaffner, X. Yang, P. P.  
652 Jiang, M. Nekoui, A. Colubri, M. R. Coomber, M. Fonnies, A. Moigboi, M. Gbakie, F. K.  
653 Kamara, V. Tucker, E. Konuwa, S. Saffa, J. Sellu, A. A. Jalloh, A. Kovoma, J. Koninga,  
654 I. Mustapha, K. Kargbo, M. Foday, M. Yillah, F. Kanneh, W. Robert, J. L. Massally,  
655 S. B. Chapman, J. Bochicchio, C. Murphy, C. Nusbaum, S. Young, B. W. Birren, D. S.  
656 Grant, J. S. Scheffelin, E. S. Lander, C. Happi, S. M. Gevao, A. Gnirke, A. Rambaut,  
657 R. F. Garry, S. H. Khan, and P. C. Sabeti. 2014. Genomic surveillance elucidates Ebola  
658 virus origin and transmission during the 2014 outbreak. *Science* 345:1369–72.
- 659 Grenfell, B. T., O. G. Pybus, J. R. Gog, J. L. Wood, J. M. Daly, J. A. Mumford, and E. C.  
660 Holmes. 2004. Unifying the epidemiological and evolutionary dynamics of pathogens.  
661 *Science* 303:327–32.
- 662 Gronau, I., S. Moran, and S. Snir. 2012. Fast and reliable reconstruction of phylogenetic  
663 trees with indistinguishable edges. *Random Structures and Algorithms* 40:350–384.

- 664 Grubaugh, N. D., E. B. Hodcroft, J. R. Fauver, A. L. Phelan, and M. Cevik. 2021. Public  
665 health actions to control new SARS-CoV-2 variants. *Cell* 184:1127–1132.
- 666 Grubaugh, N. D., J. T. Ladner, P. Lemey, O. G. Pybus, A. Rambaut, E. C. Holmes, and  
667 K. G. Andersen. 2019. Tracking virus outbreaks in the twenty-first century. *Nat.*  
668 *Microbiol.* 4:10–19.
- 669 Guindon, S., J. F. Dufayard, V. Lefort, M. Anisimova, W. Hordijk, and O. Gascuel. 2010.  
670 New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the  
671 performance of PhyML 3.0. *Syst. Biol.* 59:307–21.
- 672 Gusfield, D. 1997. Algorithms on strings, trees and sequences. Cambridge University Press,  
673 New York.
- 674 Hadfield, J., A. F. Brito, D. M. Swetnam, C. B. F. Vogels, R. E. Tokarz, K. G. Andersen,  
675 R. C. Smith, T. Bedford, and N. D. Grubaugh. 2019. Twenty years of West Nile virus  
676 spread and evolution in the Americas visualized by Nextstrain. *PLoS Pathog.*  
677 15:e1008042.
- 678 Hadfield, J., C. Megill, S. M. Bell, J. Huddleston, B. Potter, C. Callender, P. Sagulenko,  
679 T. Bedford, and R. A. Neher. 2018. Nextstrain: real-time tracking of pathogen evolution.  
680 *Bioinformatics* 34:4121–4123.
- 681 Hillis, D. M. and J. J. Bull. 1993. An empirical test of bootstrapping as a method for  
682 assessing confidence in phylogenetic analysis. *Syst. Biol.* 42:182–192.
- 683 Hoang, D. T., O. Chernomor, A. von Haeseler, B. Q. Minh, and L. S. Vinh. 2018.  
684 UFBoot2: improving the ultrafast bootstrap approximation. *Mol. Bio. Evol.* 35:518–522.
- 685 Huelsenbeck, J. P. and D. M. Hillis. 1993. Success of phylogenetic methods in the 4-taxon  
686 case. *Syst. Biol.* 42:247–264.
- 687 Kalyaanamoorthy, S., B. Q. Minh, T. K. F. Wong, A. von Haeseler, and L. S. Jermin.

- 688 2017. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat.*  
689 *Methods* 14:587–589.
- 690 Katoh, K. and D. M. Standley. 2013. MAFFT multiple sequence alignment software  
691 version 7: improvements in performance and usability. *Mol. Biol. Evol.* 30:772–80.
- 692 Lacey, M. R. and J. T. Chang. 2006. A signal-to-noise analysis of phylogeny estimation by  
693 neighbor-joining: insufficiency of polynomial length sequences. *Math. Biosci.*  
694 199:188–215.
- 695 Lanfear, R. 2020. A global phylogeny of SARS-CoV-2 sequences from GISAID.
- 696 Lemoine, F., J. B. Domelevo Entfellner, E. Wilkinson, D. Correia, M. Davila Felipe,  
697 T. De Oliveira, and O. Gascuel. 2018. Renewing Felsenstein’s phylogenetic bootstrap in  
698 the era of big data. *Nature* 556:452–456.
- 699 Liu, L., C. Anderson, D. Pearl, and S. Edwards. 2019. Modern phylogenomics: building  
700 phylogenetic trees using the multispecies coalescent model. *Methods Mol. Biol.*  
701 1910:211–239.
- 702 Lutteropp, S., A. M. Kozlov, and A. Stamatakis. 2020. A fast and memory-efficient  
703 implementation of the transfer bootstrap. *Bioinformatics* 36:2280–2281.
- 704 Minh, B. Q., H. A. Schmidt, O. Chernomor, D. Schrempf, M. D. Woodhams, A. von  
705 Haeseler, and R. Lanfear. 2020. IQ-TREE 2: new models and efficient methods for  
706 phylogenetic inference in the genomic era. *Mol. Biol. Evol.* 37:1530–1534.
- 707 Morel, B., P. Barbera, L. Czech, B. Bettisworth, L. Hubner, S. Lutteropp, D. Serdari, E. G.  
708 Kostaki, I. Mamais, A. M. Kozlov, P. Pavlidis, D. Paraskevis, and A. Stamatakis. 2021.  
709 Phylogenetic analysis of SARS-CoV-2 data is difficult. *Mol. Biol. Evol.* 38:1777–1791.
- 710 Oude Munnink, B. B., D. F. Nieuwenhuijse, M. Stein, A. O’Toole, M. Haverkate,  
711 M. Mollers, S. K. Kamga, C. Schapendonk, M. Pronk, P. Lexmond, A. van der Linden,

- 712 T. Bestebroer, I. Chestakova, R. J. Overmars, S. van Nieuwkoop, R. Molenkamp, A. A.  
713 van der Eijk, C. GeurtsvanKessel, H. Vennema, A. Meijer, A. Rambaut, J. van Dissel,  
714 R. S. Sikkema, A. Timen, M. Koopmans, and Dutch-Covid-19 response team. 2020.  
715 Rapid SARS-CoV-2 whole-genome sequencing and analysis for informed public health  
716 decision-making in the Netherlands. *Nat. Med.* 26:1405–1410.
- 717 O’Toole, Á., E. Scher, A. Underwood, B. Jackson, V. Hill, J. T. McCrone, R. Colquhoun,  
718 C. Ruis, K. Abu-Dahab, B. Taylor, C. Yeats, L. Du Plessis, D. Maloney, N. Medd, S. W.  
719 Attwood, D. M. Aanensen, E. C. Holmes, O. G. Pybus, and A. Rambaut. 2021.  
720 Assignment of epidemiological lineages in an emerging pandemic using the Pangolin  
721 tool. *Virus Evol.* veab064.
- 722 Pekar, J., M. Worobey, N. Moshiri, K. Scheffler, and J. O. Wertheim. 2021. Timing the  
723 SARS-CoV-2 index case in Hubei province. *Science* 372:412–417.
- 724 Rambaut, A. and N. C. Grassly. 1997. Seq-Gen: an application for the Monte Carlo  
725 simulation of dna sequence evolution along phylogenetic trees. *Comput. Appl. Biosci.*  
726 13:235–238.
- 727 Rambaut, A., E. C. Holmes, Á. O’Toole, V. Hill, J. T. McCrone, C. Ruis, L. du Plessis,  
728 and O. G. Pybus. 2020. A dynamic nomenclature proposal for SARS-CoV-2 lineages to  
729 assist genomic epidemiology. *Nat. Microbiol.* 5:1403–1407.
- 730 Robinson, D. F. and L. R. Foulds. 1981. Comparison of phylogenetic trees. *Math. Biosci.*  
731 53:131–147.
- 732 Roch, S. 2019. Hands-on introduction to sequence-length requirements in phylogenetics.  
733 Pages 47–86 *in* *Bioinformatics and Phylogenetics: Seminal Contributions of Bernard*  
734 *Moret* (T. Warnow, ed.). Springer International Publishing.
- 735 Roch, S. and A. Sly. 2017. Phase transition in the sample complexity of likelihood-based  
736 phylogeny inference. *Probability Theory and Related Fields* 169:3–62.

- 737 Ronquist, F. and J. P. Huelsenbeck. 2003. MrBayes 3: Bayesian phylogenetic inference  
738 under mixed models. *Bioinformatics* 19:1572–4.
- 739 Rouet, F., J. Nouhin, D. P. Zheng, B. Roche, A. Black, S. Prak, M. Leoz,  
740 C. Gaudy-Graffin, L. Ferradini, C. Mom, S. Mam, C. Gautier, G. Lesage, S. Ken,  
741 K. Phon, A. Kerleguer, C. Yang, W. Killam, M. Fujita, C. Mean, D. Fontenille,  
742 F. Barin, J. C. Plantier, T. Bedford, A. Ramos, and V. Saphonn. 2018. Massive  
743 iatrogenic outbreak of Human Immunodeficiency Virus Type 1 in rural Cambodia,  
744 2014–2015. *Clin. Infect. Dis.* 66:1733–1741.
- 745 Simmons, M. P. and A. P. Norton. 2014. Divergent maximum-likelihood-branch-support  
746 values for polytomies. *Mol. Phylogenetics Evol.* 73:87–96.
- 747 Smirnov, D. and T. Warnow. 2021. Phylogeny estimation given sequence length  
748 heterogeneity. *Syst. Biol.* 70:268–282.
- 749 Springer, M. S., E. K. Molloy, D. B. Sloan, M. P. Simmons, and J. Gatesy. 2020. ILS-aware  
750 analysis of low-homoplasy retroelement insertions: inference of species trees and  
751 introgression using quartets. *J. Hered.* 111:147–168.
- 752 Steel, M. 2000. Sufficient conditions for two tree reconstruction techniques to succeed on  
753 sufficiently long sequences. *SIAM J. Discrete Math.* 14:36–48.
- 754 Steel, M. and C. Leuenberger. 2017. The optimal rate for resolving a near-polytomy in a  
755 phylogeny. *J. Theor. Biol.* 420:174–179.
- 756 Susko, E. 2008. On the distributions of bootstrap support and posterior distributions for a  
757 star tree. *Syst. Biol.* 57:602–612.
- 758 Susko, E. 2009. Bootstrap support is not first-order correct. *Syst. Biol.* 58:211–223.
- 759 Swofford, D. L. 2003. PAUP\*. Phylogenetic analysis using parsimony (\*and other  
760 methods).

- 761 Tuffley, C. and M. Steel. 1997. Links between maximum likelihood and maximum  
762 parsimony under a simple model of site substitution. *Bull. Math. Biol.* 59:581–607.
- 763 Turakhia, Y., B. Thornlow, A. S. Hinrichs, N. De Maio, L. Gozashti, R. Lanfear,  
764 D. Haussler, and R. Corbett-Detig. 2021. Ultrafast sample placement on existing trees  
765 (USHER) enables real-time phylogenetics for the SARS-CoV-2 pandemic. *Nat. Genet.*  
766 53:809–816.
- 767 Wake, D. B., M. H. Wake, and C. D. Specht. 2011. Homoplasy: from detecting pattern to  
768 determining process and mechanism of evolution. *Science* 331:1032–1035.
- 769 Warnow, T. 2013. Large-scale multiple sequence alignment and phylogeny estimation.  
770 Pages 85–146 *in* *Models and algorithms for genome evolution* (C. Chauve,  
771 N. El-Mabrouk, and E. Tannier, eds.). Springer, London.
- 772 Worobey, M., J. Pekar, B. B. Larsen, M. I. Nelson, V. Hill, J. B. Joy, A. Rambaut, M. A.  
773 Suchard, J. O. Wertheim, and P. Lemey. 2020. The emergence of SARS-CoV-2 in  
774 Europe and North America. *Science* 370:564–570.
- 775 Yang, Z. 1998. On the best evolutionary rate for phylogenetic analysis. *Syst. Biol.*  
776 47:125–133.

777 FIGURE LEGENDS (COPIED FROM INLINE FIGURE LEGENDS)

778 Fig. 1. How a false positive split is inferred by maximum parsimony (MP). On true  
 779 tree (top) sites 1–3 are binary and “perfect”; that is, they have only a single change  
 780 (locations marked by black circles), but site 4 is binary and homoplastic, changing twice  
 781 (open circles), on edges  $e_1$  and  $e_2$ . The dotted line is the path between the two homoplastic  
 782 changes in site 4. If no perfect sites change along the dotted line path on the true tree, a  
 783 false positive split is inferred on the MP tree (bottom).

784 Fig. 2. Accuracy of maximum parsimony phylogeny reconstruction in simulations  
 785 over a wide range of per site tree length,  $\xi$ , and other parameters. Solid and dashed curves  
 786 are mean false positive and negative error rates, respectively (log scale left); dashed  
 787 sigmoidal curve is fractional resolution of estimated tree (linear scale right). Trees are  
 788 generated by a random proportional-to-distinguishable-arrangement (PDA) model for 513  
 789 taxa, from which a sequence alignment length of 1000 sites is generated. The dotted  
 790 horizontal line is placed at an error rate of 0.05. Asterisk marks the location of the optimal  
 791 tree length with best overall Robinson–Foulds accuracy,  $\xi^*$ . Each point is mean of 1000  
 792 replicates  $\times$  100 sub-replicates (see Methods). “Near-perfect” values ( $\xi \leq 1.0$ ) are shaded.  
 793 a) JC69 model with no edge length or across-site-rate variation [because of  $y$ -axis log  
 794 scaling, two  $y$  values of zero were set to 0.0001]. b) JC69 model with substantial edge  
 795 length and across-site-rate variation, both modeled as a gamma distribution with shape  
 796 parameters  $\alpha_e = \alpha_{ASR} = 0.25$ ).

797 Fig. 3. Mean false positive rate in four tree models. Fit to theoretical predictions  
 798 from Equation 6 (or the limit expression of  $4/3m$  for caterpillar trees: see Methods) are  
 799 shown by dashed lines. Each point is mean of 1000 replicates  $\times$  100 sub-replicates.  
 800 Simulation conditions were  $n = 513$ ,  $m = 1000$ , with a JC69 model. Predicted values are  
 801 not known for YH model.

802 Fig. 4. Effect of edge length variation on expected false positive rate for different  
 803 values of the shape parameter of the edge length gamma distribution,  $\alpha_e$ . Smaller values of  
 804  $\alpha_e$  correspond to higher rate variation. ASR variation is assumed absent. The dashed curve  
 805 is the prediction from Eqn. (6), in which both sources of variation are absent. Simulation  
 806 conditions assumed PDA trees with  $n = 513$ ,  $m = 1000$ , 1000 replicates, 100 subreplicates.  
 807 Gray rectangle shows “near-perfect” values of  $\xi \leq 1$ .

808 Fig. 5. Expected false positive rates,  $\mathbb{E}_{\mathcal{D}}[FP_{T_n}]$ , for 11 empirical virus phylogenetic  
 809 datasets (Table 1) for maximum parsimony (MP) inference, estimated by simulation using  
 810 parameters estimated from the data (Table S1). Abbreviations given in Table 1.  
 811 Simulation experiments used either a Yule–Harding random tree distribution (open circles)  
 812 or PDA distribution (closed circles: some data points have indistinguishable differences  
 813 between random tree models). Each point is mean of 500 replicates  $\times$  100 sub-replicates.  
 814 The near-perfect zone of  $\xi \leq 1.0$  is shaded. Horizontal dashed line indicates a 0.05  
 815 expected false positive rate.

816 Fig. 6. Statistical support measures compared to expected false positive accuracy,  
 817 as a function of tree length. The solid curve is the mean value of  $(1 - \mathbb{E}_{\mathcal{D}}[FP_{T_n}]) \times 100$  in  
 818 simulations. The near-perfect parameter space is shaded.