

1 Joint maximum-likelihood of phylogeny and 2 ancestral states is not consistent

3 David A. Shaw¹, Vu C. Dinh², and Frederick A. Matsen IV^{*1}

4 ¹Computational Biology Program, Fred Hutchinson Cancer
5 Research Center, Seattle, WA, USA

6 ²Department of Mathematical Sciences, University of Delaware,
7 Newark, DE, USA

8 Abstract

9 Maximum likelihood estimation in phylogenetics requires a means
10 of handling unknown ancestral states. Classical maximum likelihood
11 averages over these unknown intermediate states, leading to prov-
12 ably consistent estimation of the topology and continuous model pa-
13 rameters. Recently, a computationally-efficient approach has been
14 proposed to jointly maximize over these unknown states and phy-
15 logenetic parameters. Although this method of joint maximum like-
16 lihood estimation can obtain estimates more quickly, its properties as
17 an estimator are not yet clear. We show that this method of jointly
18 estimating phylogenetic parameters along with ancestral states is not
19 consistent in general. We find a sizeable region of parameter space
20 that generates data on a four-taxon tree for which this joint method
21 estimates a multifurcating topology in the limit of infinite-length se-
22 quences by estimating one or more branches to be zero length. More
23 generally, we show that this joint method only estimates branch lengths
24 correctly on a set of measure zero. We show empirically that branch
25 length estimates are biased even for short branches, with bias of the
26 same order as the branch lengths themselves.

*Corresponding author. Email: matsen@fredhutch.org

27 Introduction

28 Classical maximum likelihood (ML) estimation in phylogenetics operates
29 by integrating out latent ancestral states at the internal nodes of the tree,
30 obtaining an integrated likelihood [Goldman, 1990]. In a recent paper, Sag-
31 ulenko et al. [2018] suggest using an approximation to ML inference in
32 which the likelihood is maximized jointly across model parameters and
33 ancestral sequences on a fixed topology. This is attractive from a computa-
34 tional perspective: such joint inference can proceed according to an itera-
35 tive procedure in which ancestral sequences are first estimated and model
36 parameters are optimized conditional on these estimates. This latter con-
37 ditional optimization is simpler and more computationally efficient than
38 optimizing the integrated likelihood. But is it statistically consistent?

39 An estimator is said to be statistically consistent if it converges to the
40 generating model with probability one in the large-data limit; existing con-
41 sistency proofs for maximum likelihood phylogenetics [Allman et al., 2008,
42 Chai and Housworth, 2011, RoyChoudhury et al., 2015] apply only to es-
43 timating model parameters when the ancestral sequences have been inte-
44 grated out of the likelihood. These proofs do not readily extend to include
45 estimating ancestral states. Moreover, examples of inconsistency arising
46 from problems where the number of parameters increases with the amount
47 of data [Neyman and Scott, 1948] indicate that joint inference of trees and
48 ancestral states may not enjoy good statistical properties. In this case those
49 additional parameters are the states of ancestral sequences. Although Sag-
50 ulenko et al. [2018] explicitly warn that the approximation is for the case
51 where “branch lengths are short and only a minority of sites change on a
52 given branch,” their work motivates understanding the general properties
53 of such joint inference. In particular, one would like to know when this
54 approximate technique breaks down for both topology and branch length
55 inference, even when sequence data is “perfect,” i.e., is generated without
56 sampling error according to the exact model used for inference.

57 In this paper, we show that jointly inferring trees and ancestral sequences
58 is not consistent in general. To do so, we use a binary symmetric model

with data generated on a four-taxon tree: we compute closed form solutions to the joint objective function and demarcate a sizeable area of branch lengths in which joint inference is guaranteed to give a multifurcating tree in the case of perfect sequence data with an infinite number of sites by estimating one or more branch lengths to be exactly zero. We show that, when the topology is known and fixed, joint inference for branch length estimation cannot be consistent except on a set of measure zero (i.e. a set that occupies zero volume in parameter space). Empirically, we find areas where joint inference consistently underestimates interior branch lengths, including regions of short branch length where bias is on the same order as the branch length.

Phylogenetic maximum likelihood

Assume the binary symmetric model, namely with a character alphabet $\mathcal{A} = \{0, 1\}$ and a uniform stationary distribution [Semple and Steel, 2003]. Let m be the number of tips of the tree, and $p = m - 2$ be the number of internal nodes. We observe n independent and identically distributed samples of character data, i.e., an alignment with n columns, $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathcal{A}^{m \times n}$ distributed as the random variable Y . The corresponding unobserved ancestral states are $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_n] \in \mathcal{A}^{p \times n}$ and distributed as H with each $\mathbf{h}_i \in \mathcal{A}^p$.

We parameterize branches on the unique unrooted four-tip phylogenetic tree in ways known as the “inverse Felsenstein (InvFels)” tree (Figs. 1a and 1b) and the “Felsenstein” tree (Fig. 1c). The “inverse Felsenstein” terminology comes from Swofford et al. [2001], although it is also called the “Farris” tree [Siddall, 1998, Felsenstein, 2004]. In the standard configuration of this tree, the interior branch parameters are equal to the bottom two parameters as in Fig. 1a. We use this standard configuration as our data generating process, though we do not constrain our branch parameters to be equal when optimizing our objective function.

We parameterize the branches of these trees not with the standard notion of branch length in terms of number of substitutions per site, but with

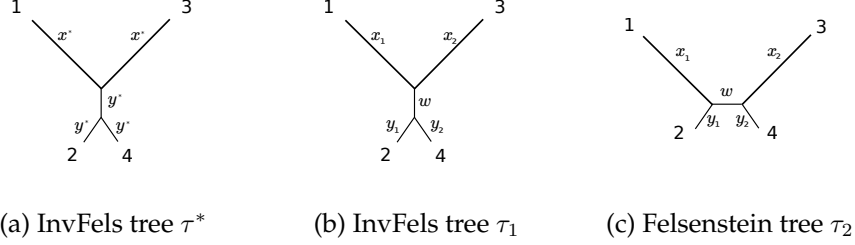


Figure 1: Three four-taxon trees with fidelities as labeled.

an alternate formulation called “fidelity.” The probability of a substitution on a branch with fidelity x is $(1 - x)/2$, while the probability of no substitution is $(1 + x)/2$ where $0 \leq x \leq 1$. This parameter quantifies the fidelity of transmission of the ancestral state across an edge [Matsen and Steel, 2007].

Fidelities have useful algebraic properties. As data becomes plentiful, we use the Hadamard transform (see (8) in the Appendix) to compute the exact probabilities that generate each particular configuration of taxa—we call these “generating probabilities”—and these have an especially simple form. For a four-taxon tree, define the general branch fidelity parameter $t = \{x_1, y_1, x_2, y_2, w\}$ where fidelities are ordered in the order of the taxa with the internal branch last (Figs. 1b and 1c). Although we use fidelities exclusively for our theoretical development, we have made our figures in terms of probabilities of substitution $p_x = (1 - x)/2$ as they are easier to interpret.

Two paths to maximum likelihood

The standard phylogenetic likelihood approach on unrooted trees under the usual assumption of independence between sites is as follows. For a topology τ and branch fidelities t the likelihood given observed ancestral states \mathbf{H} is

$$L_n(\tau, t; \mathbf{Y}, \mathbf{H}) = \prod_{i=1}^n \Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t). \quad (1)$$

109 The probability $\Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t)$ is a product of transition proba-
 110 bilities determined by \mathbf{Y} , \mathbf{H} , τ , and t [Felsenstein, 2004].

111 The classical approach is to maximize the likelihood marginalized across
 112 ancestral states

$$\tilde{L}_n(\tau, t; \mathbf{Y}) = \prod_{i=1}^n \sum_{\mathbf{h}_i \in \mathcal{A}^p} \Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t) \quad (2)$$

113 to estimate the tree τ and branch fidelities t .

114 The alternative approach [Sagulenko et al., 2018] does away with the
 115 marginalization and directly estimates the maximum likelihood paramete-
 116 ters of the fully-observed likelihood in (1). This is known in statistics as a
 117 profile likelihood [Murphy and van der Vaart, 2000] or a relative likelihood
 118 [Goldman, 1990], which exists here because \mathcal{A} is a finite set:

$$L'_n(\tau, t; \mathbf{Y}) = \prod_{i=1}^n \max_{\mathbf{h}_i \in \mathcal{A}^p} \Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t) = \max_{\mathbf{H} \in \mathcal{A}^{p \times n}} L_n(\tau, t; \mathbf{Y}, \mathbf{H}). \quad (3)$$

119 We use $\hat{\mathbf{H}}_n$ to denote an estimate for \mathbf{H} obtained by maximizing (3), and
 120 estimate a topology and branch fidelities using this profile likelihood as

$$(\hat{\tau}_n, \hat{t}_n) = \operatorname{argmax}_{\tau, t} L'_n(\tau, t; \mathbf{Y}). \quad (4)$$

121 In general, the functional form of (3) is determined by inequalities arising
 122 from taking maxima over ancestral states (Table S2) to obtain each condi-
 123 tional likelihood term, these terms depending on the unknown (τ, t) . For
 124 this reason, in practice, the joint inference strategy estimates $\hat{\mathbf{H}}_n$ for a fixed
 125 (τ, t) , then $(\hat{\tau}_n, \hat{t}_n)$ given $\hat{\mathbf{H}}_n$, maximizing each of these conditional objec-
 126 tives until convergence [Sagulenko et al., 2018].

127 Inconsistency of joint inference

128 We now state our results on the inconsistency of joint inference. All proofs
 129 are deferred to the Appendix.

130 Assume \mathbf{Y} is generated from the InvFels topology τ^* (Fig. 1a) and with
 131 true generating branch fidelities $t^* = \{x^*, y^*, x^*, y^*, y^*\}$. Let $\boldsymbol{\xi} = [\xi_j]_{j=1}^q$ be
 132 the vector of most likely ancestral state splits—the explicit definition for $\boldsymbol{\xi}$
 133 is given in the Appendix. Use $\ell_{\tau^*, t^*}(\tau, t; \boldsymbol{\xi})$ to denote the expected per-site
 134 log-likelihood, which can be thought of as the infinite-length sequence case
 135 because, as shown in the Appendix,

$$\frac{1}{n} \log L'_n(\tau, t; \mathbf{Y}) \rightarrow \ell_{\tau^*, t^*}(\tau, t; \boldsymbol{\xi}). \quad (5)$$

136 We give ℓ explicitly as (7) in the Appendix. For a fixed τ , let \hat{t}_n maximize
 137 the left-hand side of (5) and \hat{t} maximize the right-hand side. We show in
 138 the Appendix that $\hat{t}_n \rightarrow \hat{t}$, allowing us to focus on only the right-hand side
 139 above.

140 Inconsistent branch parameter estimation

141 When the topology is known and fixed and we estimate only the branch
 142 parameters, we show that for almost all generating parameter values, any
 143 branch parameter estimate is consistently biased. For the branch parameter,
 144 we use branch fidelities (discussed earlier) in all statements and proofs,
 145 though these trivially extend to branch lengths via transformation.

146 **Theorem 1.** *Let $\tau^* = \tau_1$, $t^* = \{x^*, y^*, x^*, y^*, y^*\}$, and $t = \{x_1, y_1, x_2, y_2, w\}$
 147 with $x_1, y_1, x_2, y_2, w > 0$. For $0 < x^*, y^* < 1$, the solution $\hat{t} := \{\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2, \hat{w}\}$
 148 given by*

$$\hat{t} = \arg \max_t \max_{\boldsymbol{\xi}} \ell_{\tau^*, t^*}(\tau_1, t; \boldsymbol{\xi})$$

149 *has the property that $\hat{t} \neq t^*$ everywhere except a set of measure zero.*

150 In words, the joint estimation procedure does not recover the true gen-
 151 erating t^* almost everywhere in the space of generating parameters. The
 152 proof can be seen intuitively through the fact that the estimator for a given
 153 branch fidelity cannot isolate the individual generating parameter as a lin-
 154 ear term since the estimator itself is a combination of nonlinear functions
 155 in (x^*, y^*) of generating fidelities. For example, in estimating x_1 , \hat{x}_1 is a lin-

ear combination of $p_{\tilde{y}_j}$ values. For the generating probabilities (top panel of Table S1), there is no linear combination that results in an isolated x^* term, as all terms are either quadratic, i.e., $(x^*)^2$ or $(x^*)^2(y^*)^2$, or have both x^* and y^* , i.e., $x^*(y^*)^2$. Thus, we cannot obtain a linear combination that in the general case yields x^* . We may have special cases where certain x^* and y^* values yield consistent estimates for \hat{t} , but consistency does not hold in general. We give an example of consistency in a degenerate case in the Appendix following the proof of Theorem 1.

Convergence to the degenerate topology

Given data generated on τ_1 there exist true nonzero branch lengths such that the estimator \hat{t} maximizing the right-hand side of (5) has an internal branch of length zero.

Theorem 2. *Let $\tau^* = \tau_1$, $t^* = \{x^*, y^*, x^*, y^*, y^*\}$, and $t = \{x_1, y_1, x_2, y_2, w\}$ with $x_1, y_1, x_2, y_2, w > 0$. There exists an open set of $0 < x^*, y^* < 1$ such that the solution $\hat{t} := \{\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2, \hat{w}\}$ given by*

$$\hat{t} = \arg \max_t \max_{\xi} \ell_{\tau^*, t^*}(\tau_1, t; \xi)$$

has the property $\hat{w} \equiv 1$.

This result implies an inconsistency because the joint estimate of the interior branch length is zero (i.e., interior branch fidelity is one) in an open set of values for x^* and y^* (Fig. 2). \hat{w} is available in closed form in the entire space of x^* and y^* (table in bottom panel of Fig. S2). As we consider different topologies τ_1 and τ_2 for \hat{t} , the incorrect topology τ_2 attains a likelihood value at its maximum equal to that of the true topology τ_1 in the limit. In other words, if $w = 1$ the objective functions $\ell_{\tau^*, t^*}(\tau_1, t; \xi)$ and $\ell_{\tau^*, t^*}(\tau_2, t; \xi)$ are equivalent. We elaborate on this point in the Appendix. The proof is through analytically reducing the general case to 81 separate cases (Table S3) to obtain a closed form maximal value for each.

We provide the following as an intuition for the theoretical development. For a particular site pattern, to obtain the joint maximum likelihood

function we maximize over ancestral states. For the internal branch—the branch between the two internal nodes—we have a choice of $(1 + w)$ or $(1 - w)$ in each of our likelihood terms depending on which ancestral state corresponds to the highest conditional log-likelihood. As $(1 + w) > (1 - w)$, a maximization procedure tends to prefer the $(1 + w)$ term, though this is not guaranteed because the maximum depends on the values of the unknown branch parameters t . Nevertheless, this tendency to include $(1 + w)$ terms in the likelihood results in a positive bias of branch fidelities, i.e., estimating branch lengths to be shorter than truth. This is apparent in the “long x^* , short y^* ” scenario as these are the cases in which the most likely ancestral states are the same for each internal node letting $x_1 = x_2 = x^*$ and $y_1 = y_2 = y^*$ ($\xi_j = \emptyset$ for all j in Table S3). If we allow multifurcating trees in our inference, then we can think of this as an instance of converging to the wrong topology, as the true $y^* \neq 1$.

Empirical validation

Direct numerical optimization confirms our theoretically-derived bounds and provides a more detailed picture compared to the analytically-derived region (Fig. S2). To verify the regions of inconsistency and obtain a clearer picture of the closed form parameter estimates, we plot the optimal \hat{w} via joint estimation (Fig. 2). As before, the region of inconsistency encompasses almost half of the branch fidelity space; given the correct topology, there are many situations where we estimate the interior branch length to be zero.

In our optimization procedure, we again consider the 81 separate cases (Table S3) and, for each function, we compute the closed form solution for \hat{t} . We compute these maxima over a lattice in steps of 10^{-2} for $x^*, y^* \in (0, 1)$. Our optimization code can be found at <https://github.com/matsengrp/joint-inf/>.

In estimating the interior branch length w , we find a systematic bias in the joint inference procedure even when the true branches are short (Fig. 3). As data are generated with parameters $\{x^*, y^*, x^*, y^*, y^*\}$, the true value for w is y^* . There are discontinuities in the fit (Fig. 2) due to the choice of

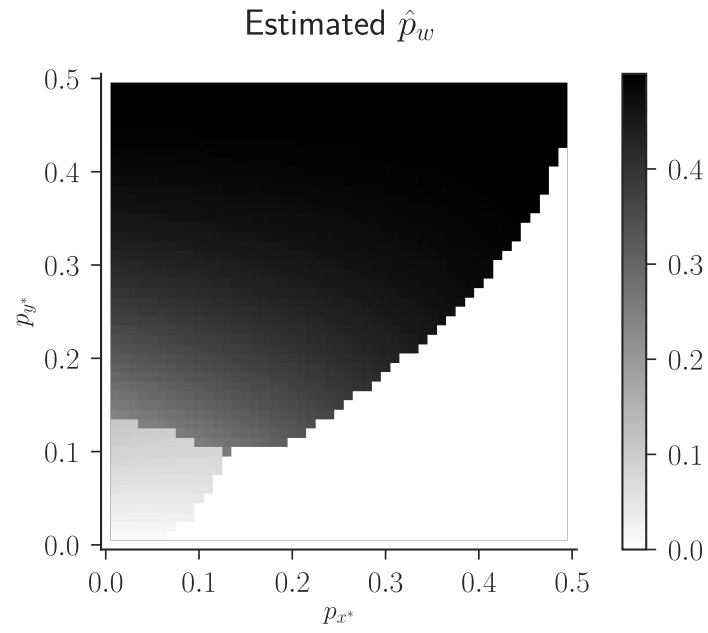


Figure 2: Estimates for $\hat{p}_w = (1 - \hat{w})/2$, the optimal probability of substitution along the inner edge for the profile likelihood (3), where the true value for p_w is p_{y^*} . Regions derived in terms of probabilities of a character change along a branch for “perfect” data generated on the InvFels topology (Fig. 1a). The white region in the lower right highlights which values of x^* and y^* result in an interior branch being estimated as length zero, resulting in an inconsistency.

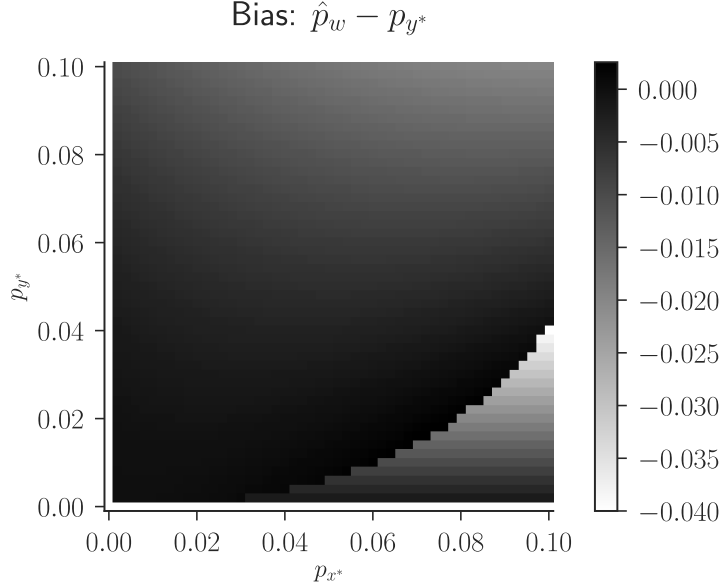


Figure 3: Bias in branch length estimation. Even in regions with short branch length ($p_{x^*}, p_{y^*} \leq .1$) where joint optimization should perform well, there is systematic bias toward shorter branch lengths.

215 which ancestral state splits are maximal, so we investigate the bias in the
 216 region where p_{x^*} and p_{y^*} are both small, i.e., $p_{x^*}, p_{y^*} \leq .1$, as these short-
 217 branch cases should be the best settings for joint optimization [Sagulenko
 218 et al., 2018]. Although the estimates for \hat{p}_w are better than the estimates
 219 when p_{y^*} is small and p_{x^*} is large (Fig. 2), joint inference still predictably
 220 underestimates the interior branch length. Additionally, the bias estimates
 221 $\hat{p}_w - p_{y^*}$ given $p_{x^*}, p_{y^*} \leq .1$ are on the same order as the branch lengths
 222 (Fig. 3), showing that even in cases where joint inference is supposed to do
 223 well, it still fails to achieve a low error from truth.

224 In contrast, inference on the integrated likelihood performs as expected,
 225 such that \hat{w} is equal to y^* regardless of the value of x^* (Fig. S3). The errors
 226 in this case (2) via optimization with L-BFGS-B are lower than machine
 227 tolerance.

228 Discussion

229 We have shown that jointly inferring ancestral states and phylogenetic pa-
230 rameters [Sagulenko et al., 2018] is not consistent in general. Specifically,
231 we have shown that the only generating parameters that yield consistent
232 branch length estimates given the correct topology lie in a set of measure
233 zero. In addition, in the case of four-taxon trees with infinite data, we
234 have obtained nontrivial regions of generating parameters that result in
235 topological ambiguity: the joint inference procedure estimates zero-length
236 branches, which can be considered as a multifurcating topology. Also, the
237 incorrect topology attains the same likelihood as the topology that gener-
238 ated the data by fixing this branch to have zero length. Since the parameters
239 with the highest likelihood given the generating topology include a zero-
240 length branch, we cannot exclude the possibility that the incorrect topol-
241 ogy with this branch having nonzero length is more likely to be observed,
242 though we have not found regions where this is the case. The regions of
243 inconsistency we found arise when one set of sister branches of the gen-
244 erating trees are “long,” that is, when the top branch fidelities tend to be
245 small, and when the lower branches are “short,” i.e., have large fidelities.
246 We see that this inconsistency occurs even if some branches are short. This
247 expands on the empirical findings of poor estimation given long branches
248 in Sagulenko et al. [2018] (their Figures 2 and 3). However, the problems
249 are not just for long branches as Sagulenko et al. [2018] imply: even when
250 all branches are short there is a consistent bias, and the bias is on the same
251 order as the magnitude of the parameters (Fig. 3).

252 Joint inference of tree parameters and ancestral sequences is a type of
253 profile likelihood, a well-studied subject in statistics [Murphy and van der
254 Vaart, 2000]. Many properties regarding the performance of maximum
255 likelihood estimates obtained using this approach are known, and many
256 methods exist to overcome their undesirable properties, e.g., the method of
257 sieves [Geman and Hwang, 1982]. A potential solution in this case using
258 the method of sieves could be to project the column-wise ancestral states
259 into a lower-dimensional space, allowing the degrees of freedom in the an-

260 central state columns to grow with n , albeit more slowly than $O(n)$. Else-
 261 where in statistics literature, the failure of maximum likelihood estimates
 262 to obtain consistent estimates as the number of parameters goes to infinity
 263 have been shown by the Neyman-Scott paradox [Neyman and Scott, 1948],
 264 though parameters tending to infinity is not a necessary condition for in-
 265 consistency [Le Cam, 1990]. Consistency proofs of standard maximum like-
 266 lihood estimates of phylogeny (2) are recent [Allman et al., 2008, Chai and
 267 Housworth, 2011, RoyChoudhury et al., 2015], and no results have been ob-
 268 tained for profile likelihood. We have furthered progress in understanding
 269 the limitations of this joint optimization procedure.

270 Previous work in phylogenetics has developed consistency counterex-
 271 amples using similar four-taxon topologies to the one used here [Felsen-
 272 stein, 1978]. In this previous work, when simulating data under the Felsen-
 273 stein topology τ_2 , as the number of observations increases, the InvFels topol-
 274 ogy τ_1 becomes more likely when performing a particular estimation pro-
 275 cedure. We have shown cases in which, when generating from the In-
 276 vFels topology, we converge to a multifurcating topology, with one or more
 277 branch lengths estimated to be zero. Moreover, the inconsistency demon-
 278 strated by Felsenstein [1978] is attributed to long branch attraction, i.e., the
 279 fact that there may be multiple long branches where parallel changes are
 280 more likely than a single change along a short branch. This is not the case
 281 here; while analytically the inconsistency occurs on a four-taxon tree when
 282 one pair of sister branches are long and the other three are short, we see em-
 283 pirically that this inconsistency is present in roughly half of the entire pa-
 284 rameter space, and occurs when the true branches generate data that more
 285 likely has no change along the interior branch. Additionally, we generate
 286 data on the InvFels tree τ_1 while Felsenstein [1978] generates data on the
 287 Felsenstein tree τ_2 . Difficulties in phylogenetic estimation when generating
 288 data on the InvFels tree have been found by Siddall [1998], though Swof-
 289 ford et al. [2001] show that the difficulties come from insufficient sequence
 290 length, which is not the case here.

291 The case of joint inference of a phylogenetic likelihood is discussed in
 292 Goldman [1990]. There, Goldman provides a worked example in which es-

293 timating a topology with fixed branch lengths is equivalent to parsimony
294 and thus not guaranteed to be consistent, though he does not discuss the
295 inconsistency of joint inference in general. We show cases where the in-
296 correct topology attains an equal likelihood value at the maximum as the
297 correct topology, and, moreover, if we know the correct topology, we show
298 cases where branch lengths are severely biased and cannot be consistent.
299 Finally, just prior to his conclusion, he discusses when parsimony gives
300 the same answer as maximum likelihood, concluding that the question is
301 ill-posed since parsimony estimates different parameters than maximum
302 likelihood, i.e., it assumes equal branch lengths. Our question, in contrast,
303 is well-posed: the joint inference procedure outlined here estimates the
304 same parameters as classical maximum likelihood—topology and branch
305 lengths—albeit implicitly estimating ancestral states as well. We are able to
306 provide much more detail on how large branch lengths must be for general
307 joint inference to fail to be consistent.

308 We have shown an inconsistency when performing joint inference on
309 branch lengths given an InvFels topology and investigated the performance
310 of branch parameter estimation. There is substantial scope for future work
311 to make these results more precise and more general. All of these results
312 hold only for a simple binary symmetric model on four-taxon trees, and
313 extensive simulation is necessary to understand how these results extend
314 to more complicated general cases, such as applied examples with larger
315 trees or more realistic mutation models that are of interest to practition-
316 ers. Also, given that many of the bounds presented here are in the form of
317 level sets of multivariate polynomials, a more formal approach using alge-
318 braic geometric techniques may reveal more stable or interesting patterns
319 of inconsistency; see [Sturmfels \[2002\]](#) for a thorough treatment of solving
320 systems of polynomial equations. Finally, all of the material presented here
321 concerns joint estimation under maximum likelihood, and does not pose
322 any problem for other settings, such as joint sampling of trees and ances-
323 tral sequences in a Bayesian framework.

324 Acknowledgements

325 We thank Richard Neher, Vladimir Minin, and Joe Felsenstein for helpful
326 discussions. Insight from the reviewers and the editors greatly improved
327 this manuscript.

328 This work was supported by National Institutes of Health grants R01-
329 GM113246, R01-AI120961, U19-AI117891, and U54-GM111274 as well as
330 National Science Foundation grants CISE-1561334 and CISE-1564137. The
331 research of Frederick Matsen was supported in part by a Faculty Scholar
332 grant from the Howard Hughes Medical Institute and the Simons Founda-
333 tion.

334 References

335 Elizabeth S Allman, Cécile Ané, and John A Rhodes. Identifiability of a
336 markovian model of molecular evolution with Gamma-Distributed rates.
337 *Adv. Appl. Probab.*, 40(1):229–249, 2008. ISSN 0001-8678. URL [http:
338 //www.jstor.org/stable/20443578](http://www.jstor.org/stable/20443578).

339 Juanjuan Chai and Elizabeth A Housworth. On rogers’ proof of identi-
340 fiability for the GTR + Γ + I model. *Syst. Biol.*, 60(5):713–718, October
341 2011. ISSN 1063-5157, 1076-836X. doi: 10.1093/sysbio/syr023. URL
342 <http://dx.doi.org/10.1093/sysbio/syr023>.

343 Joseph Felsenstein. Cases in which parsimony or compatibility methods
344 will be positively misleading. *Systematic Zoology*, 27(4):401–410, 1 De-
345 cember 1978. ISSN 0039-7989. doi: 10.2307/2412923. URL [http:
346 //www.jstor.org/stable/2412923](http://www.jstor.org/stable/2412923).

347 Joseph Felsenstein. *Inferring phylogenies*. Sinauer Associates, Inc., Sunder-
348 land, Massachusetts, 2004.

349 Stuart Geman and Chii-Ruey Hwang. Nonparametric maximum likelihood
350 estimation by the method of sieves. *The Annals of Statistics*, 10(2):401–414,
351 1982.

352 Nick Goldman. Maximum likelihood inference of phylogenetic
 353 trees, with special reference to a poisson process model of DNA
 354 substitution and to parsimony analyses. *Syst. Biol.*, 39(4):345–
 355 361, December 1990. ISSN 1063-5157. doi: 10.2307/2992355.
 356 URL [https://academic.oup.com/sysbio/article-abstract/
 357 39/4/345/1646997?redirectedFrom=fulltext](https://academic.oup.com/sysbio/article-abstract/39/4/345/1646997?redirectedFrom=fulltext).

358 Lucien Le Cam. Maximum likelihood: An introduction. *International Sta-*
 359 *tistical Review*, 58(2):153–171, Aug 1990.

360 Frederick A. Matsen and Mike Steel. Phylogenetic mixtures on a sin-
 361 gle tree can mimic a tree of another topology. *Systematic Biology*,
 362 56(5):767–775, 1 October 2007. ISSN 1063-5157. doi: 10.1080/
 363 10635150701627304. URL [http://sysbio.oxfordjournals.org/
 364 content/56/5/767.abstract](http://sysbio.oxfordjournals.org/content/56/5/767.abstract).

365 Susan A. Murphy and Aad W. van der Vaart. On profile likelihood. *Journal*
 366 *of the American Statistical Association*, 95(450):449–465, 2000. ISSN 0162-
 367 1459. doi: 10.2307/2669386. URL [http://www.jstor.org/stable/
 368 2669386](http://www.jstor.org/stable/2669386).

369 Jerzy Neyman and Elizabeth L. Scott. Consistent estimates based on par-
 370 tially consistent observations. *Econometrica*, 16(1):1–32, 1948. ISSN 0012-
 371 9682, 1468-0262. doi: 10.2307/1914288. URL [http://www.jstor.
 372 org/stable/1914288](http://www.jstor.org/stable/1914288).

373 Arindam RoyChoudhury, Amy Willis, and John Bunge. Consistency of
 374 a phylogenetic tree maximum likelihood estimator. *Journal of Statisti-*
 375 *cal Planning and Inference*, 161:73–80, June 2015. ISSN 0378-3758. doi:
 376 10.1016/j.jspi.2015.01.001. URL [http://www.sciencedirect.com/
 377 science/article/pii/S0378375815000038](http://www.sciencedirect.com/science/article/pii/S0378375815000038).

378 Pavel Sagulenko, Vadim Puller, and Richard A Neher. TreeTime:
 379 Maximum-likelihood phylodynamic analysis. *Virus Evol*, 4(1):vex042,
 380 January 2018. ISSN 2057-1577. doi: 10.1093/ve/vex042. URL [http:
 381 //dx.doi.org/10.1093/ve/vex042](http://dx.doi.org/10.1093/ve/vex042).

- 382 Charles Semple and Mike Steel. *Phylogenetics*. Oxford University Press,
383 New York, NY, 2003.
- 384 Mark E Siddall. Success of parsimony in the four-taxon case: Long-
385 branch repulsion by likelihood in the Farris zone. *Cladistics*, 14(3):209–
386 220, 1 September 1998. ISSN 0748-3007, 1096-0031. doi: 10.1111/
387 j.1096-0031.1998.tb00334.x. URL [http://dx.doi.org/10.1111/j.](http://dx.doi.org/10.1111/j.1096-0031.1998.tb00334.x)
388 [1096-0031.1998.tb00334.x](http://dx.doi.org/10.1111/j.1096-0031.1998.tb00334.x).
- 389 Bernd Sturmfels. Solving systems of polynomial equations. In *American*
390 *Mathematical Society, CBMS Regional Conferences Series, No. 97*, 2002.
- 391 David L. Swofford, Peter J. Waddell, John P. Huelsenbeck, Peter G. Foster,
392 Paul O. Lewis, and James S. Rogers. Bias in phylogenetic estimation and
393 its relevance to the choice between parsimony and likelihood methods.
394 *Systematic Biology*, 50(4):525–539, August 2001. ISSN 1063-5157. URL
395 <http://www.ncbi.nlm.nih.gov/pubmed/12116651>.
- 396 A.W. van Der Vaart. *Asymptotic statistics*. Cambridge Series in Statisti-
397 cal and Probabilistic Mathematics, 3. Cambridge University Press, 1998.
398 ISBN 9780521496032. URL [https://books.google.com/books?](https://books.google.com/books?id=udhfQgAACAAJ)
399 [id=udhfQgAACAAJ](https://books.google.com/books?id=udhfQgAACAAJ).
- 400 Ziheng Yang. Statistical properties of the maximum likelihood method of
401 phylogenetic estimation and comparison with distance matrix methods.
402 *Systematic Biology*, 43(3):329–342, 1994.

403 Appendix

404 Site split formulation

405 We begin by introducing “site splits.” We use site splits to formalize the
406 notion that a given site pattern is equally probable to its complement under
407 the binary symmetric model. This is a standard step in the description of
408 the Hadamard transform (Section 8.6 of [Semple and Steel \[2003\]](#)), although
409 our approach is complicated slightly by the inclusion of ancestral states.

410 Since we have a finite character alphabet, for a given column i there are
411 a finite number of possible assignments of characters to tips \mathbf{y}_i or internal
412 nodes \mathbf{h}_i . For the binary symmetric model, the alphabet \mathcal{A} is $\{0, 1\}$. Take
413 the tip labels of τ to be $\{1, \dots, m\}$. For likelihood calculation under the
414 binary symmetric model, we describe a given \mathbf{y}_i as a subset of indices $\tilde{y} \subseteq$
415 $\mathcal{Y} := \{1, \dots, m-1\}$, commonly called a “site split.” Define the complement
416 of \mathbf{y} as $\bar{\mathbf{y}}$, and let $\mathbf{y}_{i,k}$ be the label of the k th tip in the i th alignment column.
417 We define the site split \tilde{y} for a \mathbf{y}_i as the set of tips labeled with 1 in \mathbf{y}_i if the
418 m th tip is not labeled with 1, and as the set of tips labeled with 1 in $\bar{\mathbf{y}}_i$ if the
419 m th tip is labeled with 1. Taking such a complement simplifies but does
420 not change the result of likelihood computation because the probability of
421 observing a particular collection of binary characters is equivalent to the
422 probability of its complement under the binary symmetric model.

423 For a fixed topology τ , we define an ordered set of internal node labels
424 $\{1, \dots, p\}$ for \mathbf{h}_i and similarly use a subset of characters $\tilde{h} \subseteq \mathcal{H} := \{1, \dots, p\}$
425 to describe a realization \mathbf{h}_i . In this case we cannot use the same complement
426 trick as before: the probability of observing an ancestral state split condi-
427 tional on a site split is not invariant to taking its complement. We thus
428 define an “ancestral state split” \tilde{h} for an internal node \mathbf{h}_i to be the set of
429 internal nodes labeled with 1 if the m th tip is not labeled with 1, and as the
430 set of internal nodes labeled with 1 in $\bar{\mathbf{h}}_i$ if the m th tip is labeled with 1. We
431 emphasize that the ancestral state split complementing procedure depends
432 on tip states, not ancestral states: both site splits and ancestral state splits
433 are defined by whether the m th element of \mathbf{y}_i is labeled as 1.

434 We enumerate the site splits \tilde{y}_j of which there are $q = |\mathcal{P}(\mathcal{Y})|$ in total
 435 where \mathcal{P} denotes the power set. Similarly we enumerate ancestral state
 436 splits \tilde{h}_k of which there are $r = |\mathcal{P}(\mathcal{H})|$ in total.
 437 We first fix notation.

438 **Definition.** *Let the mapping from site patterns to site splits*

$$\psi : \mathcal{A}^m \rightarrow \mathcal{P}(\mathcal{Y})$$

439 *be*

$$\psi(\mathbf{y}) = \begin{cases} \{i' \in \{1, \dots, m-1\} : \mathbf{y}_{i,i'} = 1\} & \text{if } \mathbf{y}_{i,m} = 0, \\ \{i' \in \{1, \dots, m-1\} : \bar{\mathbf{y}}_{i,i'} = 1\} & \text{if } \mathbf{y}_{i,m} = 1, \end{cases}$$

440 *and the mapping from ancestral states and tip states to ancestral state splits*

$$\xi : \mathcal{A}^m \times \mathcal{A}^p \rightarrow \mathcal{P}(\mathcal{H})$$

441 *be*

$$\xi(\mathbf{y}, \mathbf{h}) = \begin{cases} \{i' \in \{1, \dots, p\} : \mathbf{h}_{i,i'} = 1\} & \text{if } \mathbf{y}_{i,m} = 0, \\ \{i' \in \{1, \dots, p\} : \bar{\mathbf{h}}_{i,i'} = 1\} & \text{if } \mathbf{y}_{i,m} = 1. \end{cases}$$

442 *Then, given a site pattern-valued random variable Y and an ancestral state-valued*
 443 *random variable H , define the random variables*

$$\Psi := \psi(Y)$$

444 *and*

$$\Xi := \xi(Y, H).$$

445 The mapping ψ operates by returning the tips labeled as 1 in a site pat-
 446 tern to obtain a site split in $\mathcal{P}(\mathcal{Y})$ if the set of tips labeled 1 is not in $\mathcal{P}(\mathcal{Y})$.
 447 The mapping ξ is defined by whether the tip states have their complements
 448 taken or not: if the set of tips labeled 1 in \mathbf{y} is in $\mathcal{P}(\mathcal{Y})$, $\xi(\mathbf{y}, \mathbf{h})$ is the set of
 449 tips labeled 1 in \mathbf{h} ; otherwise, the set of tips labeled 1 in $\bar{\mathbf{y}}$ necessarily is in
 450 $\mathcal{P}(\mathcal{Y})$ and so $\xi(\mathbf{y}, \mathbf{h})$ is $\bar{\mathbf{h}}$.

451 We now consider the i th factor of (1). As a consequence of assuming a

452 binary symmetric model, for some $\tilde{y}_j \in \mathcal{P}(\mathcal{Y})$ the mapping $\psi(\mathbf{y}_i)$ has the
 453 property

$$\begin{aligned} \Pr(\Psi = \tilde{y}_j, \Xi = \tilde{h}_k \mid \tau, t) &= \Pr(\Psi = \psi(\mathbf{y}_i), \Xi = \xi(\mathbf{y}_i, \mathbf{h}_i) \mid \tau, t) \\ &= \Pr((Y = \mathbf{y}_i, H = \mathbf{h}_i) \cup (\bar{Y} = \mathbf{y}_i, \bar{H} = \mathbf{h}_i) \mid \tau, t) \\ &= \Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t) + \Pr(\bar{Y} = \mathbf{y}_i, \bar{H} = \mathbf{h}_i \mid \tau, t) \\ &= 2 \cdot \Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t) \end{aligned}$$

454 where \bar{Y} is the complement of the site pattern-valued random variable Y
 455 and has the same distribution as Y (similarly for H). Since

$$2 \cdot \Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t) = \Pr(\Psi = \psi(\mathbf{y}_i), \Xi = \xi(\mathbf{y}_i, \mathbf{h}_i) \mid \tau, t),$$

456 given (τ, t) , there exist sets $\eta_1(\tau, t), \dots, \eta_q(\tau, t)$ such that $\xi_{\tilde{y}_j} \in \eta_j(\tau, t)$ satis-
 457 fies

$$\max_{\tilde{h}_k \in \mathcal{P}(\mathcal{H})} \Pr(\Psi = \tilde{y}_j, \Xi = \tilde{h}_k \mid \tau, t) = \Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau, t).$$

458 In other words, for the j th site split, $\eta_j(\tau, t) \subseteq \mathcal{P}(\mathcal{H})$ is the set of most likely
 459 ancestral state splits for that particular site split, topology and set of branch
 460 lengths, i.e., $\eta_j(\tau, t)$ is a set of sets of most likely internal node states. Here,
 461 $\xi_{\tilde{y}_j}$ is one of possibly many equiprobable ancestral state splits in $\eta_j(\tau, t)$. For
 462 each \mathbf{y}_i , $\xi(\mathbf{y}_i, \cdot)$ is surjective as it can map values from \mathcal{A}^p to all elements in
 463 $\mathcal{P}(\mathcal{H})$. This can be seen by using the definition of $\xi(\mathbf{y}_i, \cdot)$ and assuming
 464 $\mathbf{y}_{i,m} = 0$, where in this case each of the 2^p values of \mathbf{h} correspond to each
 465 of the 2^p elements of $\mathcal{P}(\{1, \dots, p\})$. The same can be done for the case of
 466 $\mathbf{y}_{i,m} = 1$, implying $\xi(\mathbf{y}_i, \cdot)$ is surjective. From this we have

$$\begin{aligned} \max_{\mathbf{h}_i} 2 \cdot \Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t) &= \max_{\mathbf{h}_i} \Pr(\Psi = \psi(\mathbf{y}_i), \Xi = \xi(\mathbf{y}_i, \mathbf{h}_i) \mid \tau, t) \\ &= \max_{\tilde{h}_k \in \mathcal{P}(\mathcal{H})} \Pr(\Psi = \tilde{y}_j, \Xi = \tilde{h}_k \mid \tau, t) \\ &= \Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau, t) \end{aligned}$$

467 for some j . Thus, each term in the likelihood can be collapsed into terms re-
 468 lating only to site splits and ancestral state splits, indexed by j , as opposed
 469 to individual observations, indexed by i .

470 Example

471 We follow with an example computing these probabilities and likelihoods.
 472 Consider the fixed, binary four-taxon tree τ_1 in Fig. 1a. The set of all possi-
 473 ble character assignments is

$$\mathcal{P}(\{1, 2, 3, 4\}) = \{\emptyset, \{1, 2, 3, 4\}, \{1\}, \{2, 3, 4\}, \{2\}, \{1, 3, 4\}, \{3\}, \{1, 2, 4\}, \\ \{1, 2\}, \{3, 4\}, \{1, 3\}, \{2, 4\}, \{2, 3\}, \{1, 4\}, \{1, 2, 3\}, \{1, 4\}\}$$

474 where each set indicates the tips assigned the character 1. For example,
 475 \emptyset is the labeling 0000 and $\{1, 3, 4\}$ is the labeling 1011. Symmetry allows
 476 us to group adjacent pairs in $\mathcal{P}(\{1, 2, 3, 4\})$ into equiprobable splits, letting
 477 $\mathcal{Y} = \{1, 2, 3\}$. The unique site splits, collapsing complements, are

$$\mathcal{P}(\mathcal{Y}) = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\} \\ =: \{\tilde{y}_1, \dots, \tilde{y}_8\}.$$

478 Since we identify character complements, we do not consider the addi-
 479 tional splits

$$\mathcal{P}(\{1, 2, 3, 4\}) \setminus \mathcal{P}(\mathcal{Y}) = \\ \{\{1, 2, 3, 4\}, \{2, 3, 4\}, \{1, 3, 4\}, \{1, 2, 4\}, \{3, 4\}, \{2, 4\}, \{1, 4\}, \{4\}\},$$

480 the symmetry of the binary character model allowing us to focus only on
 481 the elements of $\mathcal{P}(\mathcal{Y})$. This tree has two internal nodes with $\mathcal{H} = \{1, 2\}$ and
 482 unique ancestral state splits

$$\mathcal{P}(\mathcal{H}) = \{\emptyset, \{1\}, \{2\}, \{1, 2\}\}.$$

Internal node 1 is the node connected to leaves 1 and 3 while internal node 2 is connected to leaves 2 and 4. The mapping from characters to splits in this case depends on the characters at the tips and the ancestral states. For example, we take both $\psi(0000) = \emptyset$ and $\psi(1111) = \emptyset$. Similarly, we have $\xi(0000, 00) = \emptyset$ and $\xi(1111, 11) = \emptyset$, needing to take the complement of all the characters present on the tree to identify splits. We cannot identify complements for ancestral states in the same way as tip states since, for $\tilde{y} \in \mathcal{P}(\mathcal{Y})$,

$$\Pr(\Psi = \tilde{y}, \Xi = \emptyset \mid \tau, t) \neq \Pr(\Psi = \tilde{y}, \Xi = \{1, 2\} \mid \tau, t)$$

in general.

For each site split $\tilde{y} \in \mathcal{P}(\mathcal{Y})$, we maximize the likelihood over all $\tilde{h} \in \mathcal{P}(\mathcal{H})$. A maximum occurs at one of possibly several ancestral state splits in $\mathcal{P}(\mathcal{H})$, defined via $\eta_j(\tau, t)$ for the j th site split. As a simple example, say all branch lengths correspond to a probability $p (< 1/2)$ of changing character along that branch, with $t = \{p, p, p, p, p\}$. The probabilities of observing ancestral state splits for $\tilde{y}_1 = \emptyset$ are

$$\Pr(\Psi = \emptyset, \Xi = \emptyset \mid \tau, t) = (1 - p)^5,$$

$$\Pr(\Psi = \emptyset, \Xi = \{1\} \mid \tau, t) = \Pr(\Psi = \emptyset, \Xi = \{2\} \mid \tau, t) = p^3(1 - p)^2,$$

$$\Pr(\Psi = \emptyset, \Xi = \{1, 2\} \mid \tau, t) = p^4(1 - p).$$

The set of most likely ancestral states contains a single element, here $\eta_1(\tau, t) = \{\emptyset\}$. Then, taking $\xi_\emptyset \in \eta_1(\tau, t)$ we have

$$\Pr(\Psi = \emptyset, \Xi = \xi_\emptyset \mid \tau, t) = \Pr(\Psi = \emptyset, \Xi = \emptyset \mid \tau, t) = (1 - p)^5.$$

For $\tilde{y}_5 = \{1, 2\}$ we have

$$\Pr(\Psi = \{1, 2\}, \Xi = \emptyset \mid \tau, t) = \Pr(\Psi = \{1, 2\}, \Xi = \{1, 2\} \mid \tau, t) = p^2(1 - p)^3,$$

$$\Pr(\Psi = \{1, 2\}, \Xi = \{1\} \mid \tau, t) = \Pr(\Psi = \{1, 2\}, \Xi = \{2\} \mid \tau, t) = p^3(1 - p)^2.$$

504 Here, the set of most likely ancestral states is $\eta_5(\tau, t) = \{\emptyset, \{1, 2\}\}$, and, for
 505 $\xi_{12} \in \eta_5(\tau, t)$,

$$\Pr(\Psi = \{1, 2\}, \Xi = \xi_{12} \mid \tau, t) = p^2(1 - p)^3.$$

506 Site split likelihood

507 The likelihood (3) can be written as

$$\begin{aligned} L'_n(\tau, t; \mathbf{Y}) &= \max_{\mathbf{H}} L_n(\tau, t; \mathbf{Y}, \mathbf{H}) \\ &= \prod_{i=1}^n \max_{\mathbf{h}_i} \Pr(Y = \mathbf{y}_i, H = \mathbf{h}_i \mid \tau, t) \\ &\propto \prod_{i=1}^n \max_{\mathbf{h}_i} \Pr(\Psi = \psi(\mathbf{y}_i), \Xi = \xi(\mathbf{y}_i, \mathbf{h}_i) \mid \tau, t) \\ &= \prod_{i=1}^n \Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau, t) \\ &= \prod_{j=1}^q [\Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau, t)]^{n_j(\mathbf{Y})} \end{aligned} \quad (6)$$

508 for $\tilde{y}_j \in \mathcal{P}(\mathcal{Y})$ and some $\xi_{\tilde{y}_j} \in \eta_j(\tau, t)$ with $1 \leq j \leq q$ where $n_j(\mathbf{Y})$ is the
 509 number of columns in \mathbf{Y} that project to site split \tilde{y}_j .

510 Let

$$L''_n(\tau, t; \mathbf{Y}) = \prod_{j=1}^q [\Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau, t)]^{n_j(\mathbf{Y})}$$

511 be the final product in (6). Assume n observations are generated from a
 512 model with parameters (τ^*, t^*) . We have

$$\frac{1}{n} \log L''_n(\tau, t; \mathbf{Y}) = \sum_{j=1}^q \frac{n_j(\mathbf{Y})}{n} \cdot \log \Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau, t)$$

513 so that, in the $n \rightarrow \infty$ limit,

$$\begin{aligned} & \frac{1}{n} \log L_n''(\tau, t; \mathbf{Y}) \\ & \rightarrow \sum_{j=1}^q \Pr(\Psi = \tilde{y}_j \mid \tau^*, t^*) \cdot \log \Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau, t). \end{aligned} \quad (7)$$

514 Hadamard representation

515 We state the Hadamard representation of site split generating probabilities—
516 that is, probabilities of obtaining particular site splits given a tree—following
517 Section 8.6 of [Semple and Steel \[2003\]](#). For each edge e define the edge “fi-
518 delity” for that edge as

$$\theta(e) = 1 - 2p(e)$$

519 where $p(e)$ is the probability of a character change along edge e . For an
520 even-sized subset of $Y \subseteq \mathcal{S}$, let the path set $P(Y)$ be the set of edges in the
521 path connecting both elements of Y . For n taxa, the probability of observing
522 site split $A \in \mathcal{P}(\mathcal{Y})$ is

$$p_A = \frac{1}{2^{n-1}} \sum_{Y \subseteq \mathcal{S}: |Y| \equiv 0 \pmod{2}} \left[(-1)^{|Y \cap A|} \prod_{e \in P(Y)} \theta(e) \right]. \quad (8)$$

523 By convention, we set $P(\emptyset) = \emptyset$ and $\prod_{e \in \emptyset} \theta(e) = 1$. For notational convenience, let
524

$$p_{\tilde{y}_j} := \Pr(\Psi = \tilde{y}_j \mid \tau_1, t),$$

525 for any site split \tilde{y}_j . Table [S1](#) contains calculations of site split probabilities
526 for the trees in Fig. [1](#).

527 Likelihood computations

528 To compute the likelihood of observing a set of data, we need $\Pr(\Psi =$
529 $\tilde{y}_j, \Xi = \tilde{h}_k \mid \tau, t)$ for each \tilde{h}_k and \tilde{y}_j . Using branch fidelities, the probability
530 of a character change along a branch with fidelity parameter x is $(1 - x)/2$,
531 while the probability of a character remaining the same is $(1 + x)/2$. See

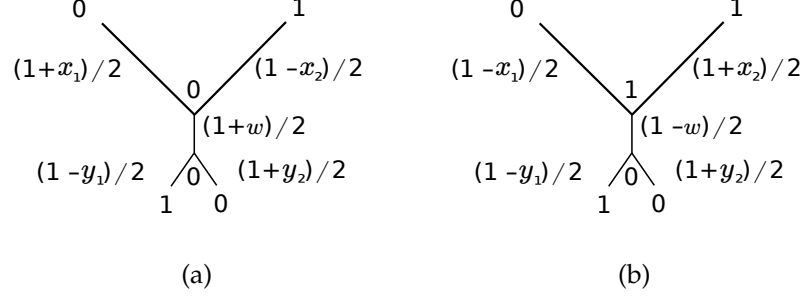


Figure S1: Example likelihood computations on the InvFels tree τ_1 for fidelities $t = \{x_1, y_1, x_2, y_2, w\}$. Edges labeled by the probability of substitution along that edge. In (a), we compute the product to obtain $\Pr(\Psi = \{2, 3\}, \Xi = \emptyset \mid \tau_1, t) = (1 + x_1)(1 - x_2)(1 + y_1)(1 - y_2)(1 + w)/32$. In (b), the same process yields $\Pr(\Psi = \{2, 3\}, \Xi = \{1\} \mid \tau_1, t) = (1 + x_1)(1 - x_2)(1 + y_1)(1 - y_2)(1 - w)/32$.

Fig. S1 for the parameters on an example site pattern on the InvFels tree.
Likelihood computations for all site splits and ancestral state splits are in
Table S2 for the InvFels tree.

Convergence of branch parameters

For a fixed τ , we show that $\hat{t}_n \rightarrow \hat{t}$ for

$$\hat{t}_n = \arg \max_{t \in \mathcal{T}} \frac{1}{n} \log L'_n(\tau, t; \mathbf{Y})$$

and

$$\hat{t} = \arg \max_{t \in \mathcal{T}} \ell_{\tau^*, t^*}(\tau, t; \boldsymbol{\xi}).$$

Using the notation in Section 5.2.1 in van Der Vaart [1998], we let

$$m_t(\mathbf{y}) = \sum_{j=1}^q 1\{\psi(\mathbf{y}) = \tilde{y}_j\} \cdot \log \Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau, t)$$

so that

$$\frac{1}{n} \log L'_n(\tau, t; \mathbf{Y}) = \frac{1}{n} \sum_{i=1}^n m_t(\mathbf{y}_i)$$

540 and

$$\ell_{\tau^*, t^*}(\tau, t; \boldsymbol{\xi}) = E[m_t].$$

541 To show $\hat{t}_n \rightarrow \hat{t}$, we use Wald's consistency proof [p. 48, Theorem 5.14 of
 542 [van Der Vaart, 1998](#)], which requires four conditions. The first is that \mathcal{T} is
 543 compact, which is obviously true. The second is that

$$E \left[\sup_{t \in \mathcal{T}} m_t \right] < \infty,$$

544 and, since $m_t(\mathbf{y})$ is nonpositive for all t and \mathbf{y} , this property holds. The
 545 remaining conditions are on the maps

$$\mathbf{y} \mapsto \sup_t m_t(\mathbf{y})$$

546 and

$$t \mapsto m_t(\mathbf{y}).$$

547 We need the first map to be measurable, which is evident since the do-
 548 main \mathcal{A}^m of the mapping is a finite set, and so all subsets of the domain
 549 are also finite and thus measurable. Finally, we must have the the second
 550 mapping be upper-semicontinuous for almost all \mathbf{y} . For a fixed ancestral
 551 state split $t \mapsto m_t(\mathbf{y})$ is continuous for all \mathbf{y} . If we move about in \mathcal{T} , a
 552 different ancestral state split becomes more likely, though when we maxi-
 553 mize over ancestral state splits we obtain a continuous function since the
 554 maximum over continuous functions is also continuous. This ensures the
 555 upper-semicontinuous property of this mapping, and shows $\hat{t}_n \rightarrow \hat{t}$, allow-
 556 ing our consistency results to be proved using $\ell_{\tau^*, t^*}(\tau, t; \boldsymbol{\xi})$.

557 This Wald-type consistency does not extend readily to convergence of
 558 topology, i.e., we cannot use the above arguments to say that, for

$$(\hat{\tau}_n, \hat{t}_n) = \arg \max_{\tau, t} \frac{1}{n} \log L'_n(\tau, t; \mathbf{Y})$$

559 and

$$(\hat{\tau}, \hat{t}) = \arg \max_{\tau, t} \ell_{\tau^*, t^*}(\tau, t; \boldsymbol{\xi})$$

560 then

$$(\hat{\tau}_n, \hat{t}_n) \rightarrow (\hat{\tau}, \hat{t})$$

561 as we did for \hat{t}_n with τ fixed due to the properties of tree space [Yang, 1994].
 562 For this reason we do not claim joint inference is inconsistent in estimating
 563 topology, only that there exist parameters that yield degenerate maxima
 564 when performing joint optimization.

565 Properties of the joint objective function

566 Consider the InvFels tree τ_1 with arbitrary fidelities, i.e., $t = \{x_1, y_1, x_2, y_2, w\}$.
 567 Next we show that the likelihood $\ell_{\tau_1, t}(\tau_1, t; \xi)$ remains unchanged if x_1 and
 568 x_2 are exchanged or if y_1 and y_2 are. Although this property should not be
 569 surprising due to symmetry, we write it out for completeness. This holds
 570 for a general t , and thus holds setting $t = t^*$. Using the Hadamard trans-
 571 form, we calculate the generating probabilities on the InvFels tree. For site
 572 split \emptyset ,

$$\begin{aligned} \Pr(\Psi = \emptyset \mid \tau_1, t) &= \frac{1}{8}(1 + x_1x_2 + y_1y_2 + x_1y_1w + x_1y_2w + y_1x_2w + x_2y_2w + x_1y_1x_2y_2) \\ &= \frac{1}{8}(1 + x_1x_2 + y_1y_2 + w[x_1y_1 + x_1y_2 + y_1x_2 + x_2y_2] + x_1y_1x_2y_2) \\ &= \frac{1}{8}(1 + x_1x_2 + y_1y_2 + w[x_1 + x_2][y_1 + y_2] + x_1y_1x_2y_2), \end{aligned}$$

573 and this probability is unchanged when x_1 is exchanged with x_2 and y_1 is
 574 exchanged with y_2 . Similarly, for site split $\{1, 3\}$,

$$\Pr(\Psi = \{1, 3\} \mid \tau_1, t) = \frac{1}{8}(1 + x_1x_2 + y_1y_2 - w[x_1 + x_2][y_1 + y_2] + x_1y_1x_2y_2),$$

575 which also is invariant to exchanging x_1 with x_2 and y_1 with y_2 .

576 All other generating probabilities differ only in the signs of each term
 577 (see Table S1). For example, for site split $\{1\}$ we have

$$\Pr(\Psi = \{1\} \mid \tau_1, t) = \frac{1}{8}(1 - x_1x_2 + y_1y_2 + w[-x_1 + x_2][y_1 + y_2] - x_1y_1x_2y_2)$$

578 and for site split $\{3\}$ we have

$$\Pr(\Psi = \{3\} \mid \tau_1, t) = \frac{1}{8}(1 - x_1x_2 + y_1y_2 + w[x_1 - x_2][y_1 + y_2] - x_1y_1x_2y_2)$$

579 meaning if we exchange the values of x_1 and x_2 then these probabilities
 580 swap values, regardless of what we do with y_1 and y_2 . We show that for site
 581 splits $\{1\}$ and $\{3\}$, exchanging x_1 and x_2 also swaps the values of the like-
 582 lihood terms, again independent of what happens to y_1 and y_2 (Table S2).
 583 Indeed, the corresponding possibilities for the likelihood values are

$$\begin{aligned}\Pr(\Psi = \{1\}, \Xi = \emptyset \mid \tau_1, t) &= \frac{1}{32}(1 - x_1)(1 + x_2)(1 + w)(1 + y_1)(1 + y_2); \\ \Pr(\Psi = \{1\}, \Xi = \{1\} \mid \tau_1, t) &= \frac{1}{32}(1 + x_1)(1 - x_2)(1 - w)(1 + y_1)(1 + y_2); \\ \Pr(\Psi = \{1\}, \Xi = \{2\} \mid \tau_1, t) &= \frac{1}{32}(1 - x_1)(1 + x_2)(1 - w)(1 - y_1)(1 - y_2); \\ \Pr(\Psi = \{1\}, \Xi = \{1, 2\} \mid \tau_1, t) &= \frac{1}{32}(1 + x_1)(1 - x_2)(1 + w)(1 - y_1)(1 - y_2);\end{aligned}$$

584 for site split $\{1\}$ and

$$\begin{aligned}\Pr(\Psi = \{3\}, \Xi = \emptyset \mid \tau_1, t) &= \frac{1}{32}(1 + x_1)(1 - x_2)(1 + w)(1 + y_1)(1 + y_2); \\ \Pr(\Psi = \{3\}, \Xi = \{1\} \mid \tau_1, t) &= \frac{1}{32}(1 - x_1)(1 + x_2)(1 - w)(1 + y_1)(1 + y_2); \\ \Pr(\Psi = \{3\}, \Xi = \{2\} \mid \tau_1, t) &= \frac{1}{32}(1 + x_1)(1 - x_2)(1 - w)(1 - y_1)(1 - y_2); \\ \Pr(\Psi = \{3\}, \Xi = \{1, 2\} \mid \tau_1, t) &= \frac{1}{32}(1 - x_1)(1 + x_2)(1 + w)(1 - y_1)(1 - y_2);\end{aligned}$$

585 for site split $\{3\}$, which shows the likelihood remains unchanged if x_1 and
 586 x_2 are swapped.

587 For site splits $\{2\}$ and $\{1, 2, 3\}$, exchanging y_1 and y_2 swaps the values of
 588 the generating probabilities, independent of what happens to x_1 and x_2 . In
 589 the case of the likelihood values, we see that the values for these site splits
 590 swap as well, though, we look at the complement of the most likely ances-
 591 tral state split. In other words, the function value for the likelihood also
 592 swaps between site splits $\{2\}$ and $\{1, 2, 3\}$, though the most likely ancestral

593 state split is different. Indeed,

$$\begin{aligned}\Pr(\Psi = \{2\}, \Xi = \emptyset \mid \tau_1, t) &= \frac{1}{32}(1+x_1)(1-y_1)(1+x_2)(1+y_2)(1+w); \\ \Pr(\Psi = \{2\}, \Xi = \{1\} \mid \tau_1, t) &= \frac{1}{32}(1-x_1)(1-y_1)(1-x_2)(1+y_2)(1-w); \\ \Pr(\Psi = \{2\}, \Xi = \{2\} \mid \tau_1, t) &= \frac{1}{32}(1+x_1)(1+y_1)(1+x_2)(1-y_2)(1-w); \\ \Pr(\Psi = \{2\}, \Xi = \{1, 2\} \mid \tau_1, t) &= \frac{1}{32}(1-x_1)(1+y_1)(1-x_2)(1-y_2)(1+w);\end{aligned}$$

594 for site split $\{2\}$ and

$$\begin{aligned}\Pr(\Psi = \{1, 2, 3\}, \Xi = \emptyset \mid \tau_1, t) &= \frac{1}{32}(1-x_1)(1-y_1)(1-x_2)(1+y_2)(1+w); \\ \Pr(\Psi = \{1, 2, 3\}, \Xi = \{1\} \mid \tau_1, t) &= \frac{1}{32}(1+x_1)(1-y_1)(1+x_2)(1+y_2)(1-w); \\ \Pr(\Psi = \{1, 2, 3\}, \Xi = \{2\} \mid \tau_1, t) &= \frac{1}{32}(1-x_1)(1+y_1)(1-x_2)(1-y_2)(1-w); \\ \Pr(\Psi = \{1, 2, 3\}, \Xi = \{1, 2\} \mid \tau_1, t) &= \frac{1}{32}(1+x_1)(1+y_1)(1+x_2)(1-y_2)(1+w);\end{aligned}$$

595 for site split $\{1, 2, 3\}$, which shows the likelihood remains unchanged if y_1
596 and y_2 are swapped.

597 For site splits $\{1, 2\}$ and $\{2, 3\}$ we see the following. By exchanging
598 only x_1 with x_2 , the generating probabilities and likelihood values swap
599 between these two site splits. The same is true of the generating probabili-
600 ties if we exchange only y_1 and y_2 , except, for the case of the likelihood val-
601 ues, we again look at the complement of the most likely ancestral state split
602 as in the case of splits $\{2\}$ and $\{1, 2, 3\}$. Now, if we exchange both x_1 with
603 x_2 and y_1 with y_2 , we see these generating probabilities remain unchanged,
604 and, for the likelihood values, we look at the complement of the most likely
605 ancestral state split and see these values also remain unchanged.

606 Thus exchanging x_1 with x_2 and y_1 with y_2 does not change the value
607 of the log-likelihood $\ell_{\tau_1, t}(\tau_1, t; \xi)$. Therefore we can reduce the number of
608 candidate likelihoods we need to search by, without loss of generality, as-
609 suming $x_2 \geq x_1$ and $y_2 \geq y_1$, with these likelihoods given in Table S3 after
610 maximizing over ancestral state splits.

611 Theorems and proofs

612 We begin by showing an inconsistency in branch length estimation on the
613 InvFels tree.

614 **Theorem 1.** Let $\tau^* = \tau_1$, $t^* = \{x^*, y^*, x^*, y^*, y^*\}$, and $t = \{x_1, y_1, x_2, y_2, w\}$
615 with $x_1, y_1, x_2, y_2, w > 0$. For $0 < x^*, y^* < 1$, the solution $\hat{t} := \{\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2, \hat{w}\}$
616 given by

$$\hat{t} = \arg \max_t \max_{\xi} \ell_{\tau^*, t^*}(\tau_1, t; \xi)$$

617 has the property that $\hat{t} \neq t^*$ everywhere except a set of measure zero.

618 *Proof.* For a fixed, known ξ , there exists a closed form solution to $\hat{t} :=$
619 $\{\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2, \hat{w}\}$ solving

$$\hat{t}_{\xi} = \arg \max_t \ell_{\tau^*, t^*}(\tau_1, t; \xi).$$

620 We show in this case that the log-likelihood ℓ attains a unique maximum at
621 \hat{t}_{ξ} . For fixed ξ , the log-likelihood can be decomposed into a sum of func-
622 tions of each variable, i.e.,

$$\begin{aligned} \ell_{\tau^*, t^*}(\tau^*, t, \xi) &= \sum_{j=1}^q p_{\tilde{y}_j} \cdot \log h_{j, x_1}(x_1) + \sum_{j=1}^q p_{\tilde{y}_j} \cdot \log h_{j, y_1}(y_1) + \sum_{j=1}^q p_{\tilde{y}_j} \cdot \log h_{j, x_2}(x_2) \\ &\quad + \sum_{j=1}^q p_{\tilde{y}_j} \cdot \log h_{j, y_2}(y_2) + \sum_{j=1}^q p_{\tilde{y}_j} \cdot \log h_{j, w}(w). \end{aligned}$$

623 Due to this additive form, all off-diagonal terms of the Hessian for this
624 function are zero, so we show that the diagonal terms are nonpositive. For
625 example, if we focus on the variable x_1 , holding other variables constant,
626 then

$$\ell(x_1) \sim \sum_{j=1}^q p_{\tilde{y}_j} \cdot \log h_{j, x_1}(x_1).$$

627 where \sim denotes equality up to an additive constant. Doing calculation as
628 in Figure S1, each functional form, suppressing constants with respect to

629 the focal variable (here x_1) and the initial $1/32$ constant, takes the form

$$h_{j,x_1}(x_1) = (1 + x_1)^{e_j} (1 - x_1)^{1-e_j}$$

630 for $e_j \in \{0, 1\}$, which, simplifying, results in

$$\ell(x_1) \sim \left(\sum_{j=1}^q p_{\tilde{y}_j} e_j \right) \log(1 + x_1) + \left(\sum_{j=1}^q p_{\tilde{y}_j} (1 - e_j) \right) \log(1 - x_1) \quad (9)$$

$$= \left(\sum_{j=1}^q p_{\tilde{y}_j} e_j \right) \log(1 + x_1) + \left(1 - \sum_{j=1}^q p_{\tilde{y}_j} e_j \right) \log(1 - x_1), \quad (10)$$

631 which has second derivative

$$\ell''(x_1) = - \left(\frac{\sum_j p_{\tilde{y}_j} e_j}{(1 + x_1)^2} + \frac{1 - \sum_j p_{\tilde{y}_j} e_j}{(1 - x_1)^2} \right).$$

632 As $x_1 \in (0, 1]$, we need only $0 \leq \sum_j p_{\tilde{y}_j} e_j \leq 1$ to imply the diagonal terms
 633 of the Hessian are nonpositive. Since $p_{\tilde{y}_j}$ are probabilities, then $\sum_j p_{\tilde{y}_j} = 1$.
 634 As $e_j \in \{0, 1\}$ by definition, this implies $0 \leq \sum_j p_{\tilde{y}_j} e_j \leq 1$ and $\ell''(x_1) \leq 0$.
 635 Applying similar arguments to the other variables, the Hessian for the log-
 636 likelihood has nonpositive diagonal terms and off-diagonal terms equal to
 637 zero, and \hat{t} uniquely maximizes ℓ .

638 Now, by straightforward calculus, we solve for the unique maximum
 639 \hat{x}_1 by setting the first derivative of (10) to zero to obtain

$$\hat{x}_1 = 2 \cdot \left(\sum_{j=1}^q p_{\tilde{y}_j} e_j \right) - 1$$

640 where

$$\sum_{j=1}^q p_{\tilde{y}_j} e_j = \sum_{j=1}^q \mathbf{1}\{\text{site split } j \text{ has term } (1 + x_1)\} \cdot p_{\tilde{y}_j}.$$

641 Next we show that solutions of this form never obtain $\hat{t} = t^*$ except on

642 a set of measure zero. Given Table S1, all solutions \hat{x}_1 have the form

$$\hat{x}_1(x^*, y^*) = a_{x_1,0} + a_{x_1,1}(x^*)^2 + a_{x_1,2}(y^*)^2 + a_{x_1,3}x^*(y^*)^2 + a_{x_1,4}(x^*)^2(y^*)^2,$$

643 where $a_{x_1,k}$ are constants independent of x^* and y^* —in fact, $a_{x_1,k}$ takes
 644 values in the set $\{i/8 : i = -4, -3, \dots, 7, 8\}$. The true branch fidelity for x_1
 645 is x^* , and so we have consistency when

$$f_{y^*}(x^*) = \hat{x}_1(x^*, y^*) - x^*$$

646 is zero. As the number of zeros of $f_{y^*}(x^*)$ is finite for $0 < x^*, y^* < 1$, con-
 647 sistency holds only on a set of measure zero for a fixed ξ . Since we have
 648 inconsistency in \hat{t}_ξ for each ξ , this implies an inconsistency when maximiz-
 649 ing over ξ as it takes values on a finite set.

650 The same is true for x_2 , and a similar argument for y_1, y_2 , and w shows
 651 that estimates can only be consistent on a set of measure zero.

652 □

653 As an example solution for maximizing x_1 , the maximal ancestral state
 654 splits for $\hat{\xi}_1 = \{\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset\}$ (Table S3) yields the log-likelihood

$$\begin{aligned} \ell_{\tau^*, t^*}(x_1; \hat{\xi}_1) &= \left(\frac{1}{2} + \frac{1}{2}x^*(y^*)^2\right) \log(1 + x_1) + \left(\frac{1}{2} - \frac{1}{2}x^*(y^*)^2\right) \log(1 - x_1) \\ &\quad + \left(\frac{1}{2} + \frac{1}{2}(y^*)^2\right) \log(1 + y_1) + \left(\frac{1}{2} - \frac{1}{2}(y^*)^2\right) \log(1 - y_1) \\ &\quad + \left(\frac{1}{2} + \frac{1}{2}x^*(y^*)^2\right) \log(1 + x_2) + \left(\frac{1}{2} - \frac{1}{2}x^*(y^*)^2\right) \log(1 - x_2) \\ &\quad + \log(1 + y_2) + \log(1 + w) - \log 32 \\ &:= \left(\frac{1}{2} + \frac{1}{2}x^*(y^*)^2\right) \log(1 + x_1) + \left(\frac{1}{2} - \frac{1}{2}x^*(y^*)^2\right) \log(1 - x_1) \\ &\quad + C_{x_1}(x_2, y_1, y_2, w). \end{aligned}$$

655 where C_{x_1} is a function of y_1, x_2, y_2 , and w and

$$\sum_{j=1}^q p_{\tilde{y}_j} e_j = p_{\emptyset} + p_2 + p_3 + p_{23} = \frac{1}{2} + \frac{1}{2} x^* (y^*)^2.$$

656 Thus, $\hat{x}_1(x^*, y^*) = x^* (y^*)^2$, and it is possible to obtain $\hat{x}_1(x^*, y^*) = x^*$,
 657 such as when $y^* = 1$, but we cannot have $\hat{x}_1(x^*, y^*) = x^*$ for all x^* where
 658 $0 < x^*, y^* < 1$.

659 We now proceed to show there exist x^* and y^* such that the interior
 660 branch parameter w is estimated as exactly one, indicating convergence to
 661 a multifurcating topology.

662 **Theorem 2.** Let $\tau^* = \tau_1$, $t^* = \{x^*, y^*, x^*, y^*, y^*\}$, and $t = \{x_1, y_1, x_2, y_2, w\}$
 663 with $x_1, y_1, x_2, y_2, w > 0$. There exists an open set of $0 < x^*, y^* < 1$ such that
 664 the solution $\hat{t} := \{\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2, \hat{w}\}$ given by

$$\hat{t} = \arg \max_t \max_{\xi} \ell_{\tau^*, t^*}(\tau_1, t; \xi)$$

665 has the property $\hat{w} \equiv 1$.

666 *Proof.* As we have a closed form solution to our likelihood problem, we
 667 compute the optimal solution given Table S2. Let

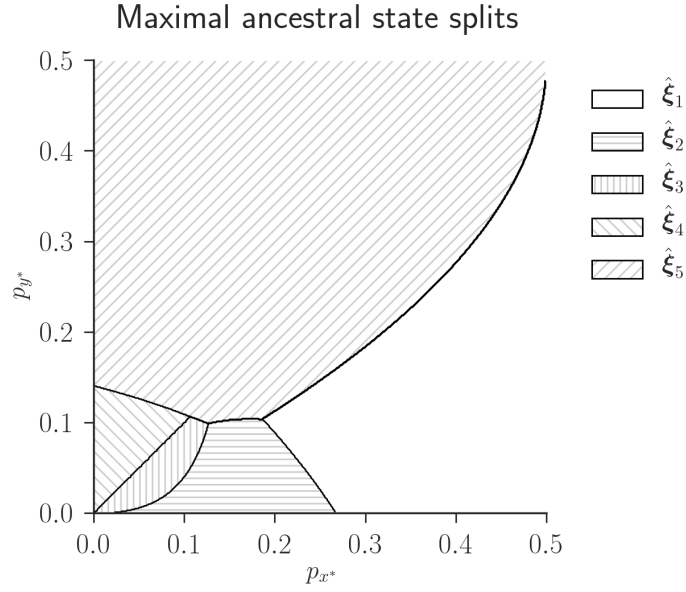
$$\hat{t}_{\xi} = \arg \max_t \ell_{\tau^*, t^*}(\tau, t; \xi).$$

668 be the closed form solution for t for a fixed maximal ancestral state split ξ .
 669 We need only consider the possibilities for choices of ancestral state splits
 670 in Table S3 as opposed to Table S2. Upon excluding cases of infinite branch
 671 lengths (i.e., any of x_1, y_1, x_2, y_2, w equal to zero) and the redundant cases
 672 of $x_1 > x_2$ and $y_1 > y_2$, we obtain

$$\hat{\xi} = \arg \max_{\xi} \ell_{\tau^*, t^*}(\tau_1, \hat{t}_{\xi}; \xi).$$

673 We show the maximal ancestral states in Fig. S2.

674 Mapping each maximal ancestral state split to each likelihood value,



Maximal ancestral state split definitions and \hat{w}

$\hat{\xi} = \{\xi_\emptyset, \xi_1, \xi_2, \xi_3, \xi_{123}, \xi_{12}, \xi_{23}, \xi_{13}\}$	\hat{w}
$\hat{\xi}_1 = \{\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset\}$	1
$\hat{\xi}_2 = \{\emptyset, \emptyset, \emptyset, \emptyset, \{1, 2\}, \emptyset, \emptyset, \emptyset\}$	1
$\hat{\xi}_3 = \{\emptyset, \emptyset, \emptyset, \emptyset, \{1, 2\}, \emptyset, \emptyset, \{1\}\}$	$\frac{3}{4} + x^*(y^*)^2 - \frac{1}{4}((x^*)^2 + (y^*)^2 + (x^*)^2(y^*)^2)$
$\hat{\xi}_4 = \{\emptyset, \emptyset, \emptyset, \emptyset, \{1, 2\}, \emptyset, \{1, 2\}, \{1\}\}$	$\frac{3}{4} + x^*(y^*)^2 - \frac{1}{4}((x^*)^2 + (y^*)^2 + (x^*)^2(y^*)^2)$
$\hat{\xi}_5 = \{\emptyset, \emptyset, \emptyset, \{1\}, \{1\}, \emptyset, \{1\}, \{1\}\}$	$x^*(y^*)^2$

Figure S2: Regions of maximal ancestral state splits on the InvFels tree τ_1 .
Data generated as in Fig. 2.

we see that $\hat{w} \equiv 1$ if $\hat{\xi} = \hat{\xi}_1$ or $\hat{\xi} = \hat{\xi}_2$, which encompasses the bottom-right region of Figure S2. \square

The regions in Fig. S2 are analytically-derived regions of inconsistency in terms of probabilities of a character change along a branch for “perfect” data generated on the InvFels topology (Fig. 1) with $p_{w^*} = p_{y^*}$ (in terms of fidelities, $w^* = y^*$). As the region of degeneracy in Fig. S2 gives the values of x^* and y^* where \hat{w} is guaranteed to be one, we converge on a multifurcating topology in these cases. It is easy to see that when \emptyset is the maximal ancestral state split, we have the same log-likelihood for τ_1 and τ_2 . Moreover, if $w = 1$, the internal branch becomes zero-length and the two topologies are indistinguishable.

The boundaries determining maximal ancestral state splits (Fig. S2) are obtained through maximizing over 81 separate functional values (Table S3). Referring to the proof of Theorem 1, we see that

$$\ell(x_1) \sim \left(\sum_{j=1}^q p_{\tilde{y}_j} e_j \right) \log(1 + x_1) + \left(1 - \sum_{j=1}^q p_{\tilde{y}_j} e_j \right) \log(1 - x_1)$$

and

$$\hat{x}_1 = 2 \cdot \left(\sum_{j=1}^q p_{\tilde{y}_j} e_j \right) - 1$$

so that the maximal value for $\ell(x_1)$ is

$$\frac{1 + \hat{x}_1}{2} \log(1 + \hat{x}_1) + \frac{1 - \hat{x}_1}{2} \log(1 - \hat{x}_1)$$

with similar forms for the remaining variables. Our likelihood (7) has an exact, closed form maximum of

$$\begin{aligned} \ell_{\tau^*, t^*}(\tau_1, \hat{t}; \hat{\xi}) &= \frac{1 + \hat{x}_1}{2} \log(1 + \hat{x}_1) + \frac{1 - \hat{x}_1}{2} \log(1 - \hat{x}_1) \\ &\quad + \frac{1 + \hat{y}_1}{2} \log(1 + \hat{y}_1) + \frac{1 - \hat{y}_1}{2} \log(1 - \hat{y}_1) \\ &\quad + \frac{1 + \hat{x}_2}{2} \log(1 + \hat{x}_2) + \frac{1 - \hat{x}_2}{2} \log(1 - \hat{x}_2) \end{aligned}$$

$$\begin{aligned}
& + \frac{1 + \hat{y}_2}{2} \log(1 + \hat{y}_2) + \frac{1 - \hat{y}_2}{2} \log(1 - \hat{y}_2) \\
& + \frac{1 + \hat{w}}{2} \log(1 + \hat{w}) + \frac{1 - \hat{w}}{2} \log(1 - \hat{w}). \quad (11)
\end{aligned}$$

693 For example, in the lower right region of Fig. S2, the maximal ancestral
694 state split is $\hat{\xi}_1 = \{\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset\}$ so that

$$\{\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2, \hat{w}\} = \{x^*(y^*)^2, (y^*)^2, x^*(y^*)^2, 1, 1\}$$

695 and the exact, closed form value of the likelihood is, as a function of x^* and
696 y^* ,

$$\begin{aligned}
g_{\xi_1}(x^*, y^*) &:= (1 + x^*(y^*)^2) \log(1 + x^*(y^*)^2) + (1 - x^*(y^*)^2) \log(1 - x^*(y^*)^2) \\
&+ \frac{1 + (y^*)^2}{2} \log(1 + (y^*)^2) + \frac{1 - (y^*)^2}{2} \log(1 - (y^*)^2) \\
&+ 2 \log(2). \quad (12)
\end{aligned}$$

697 We obtain similar maximal values as functions of x^* and y^* for each ances-
698 tral state split to get g_{ξ} for all relevant ξ . The curves delineating maximal
699 ancestral state splits (Fig. S2) are determined by these 81 functions.

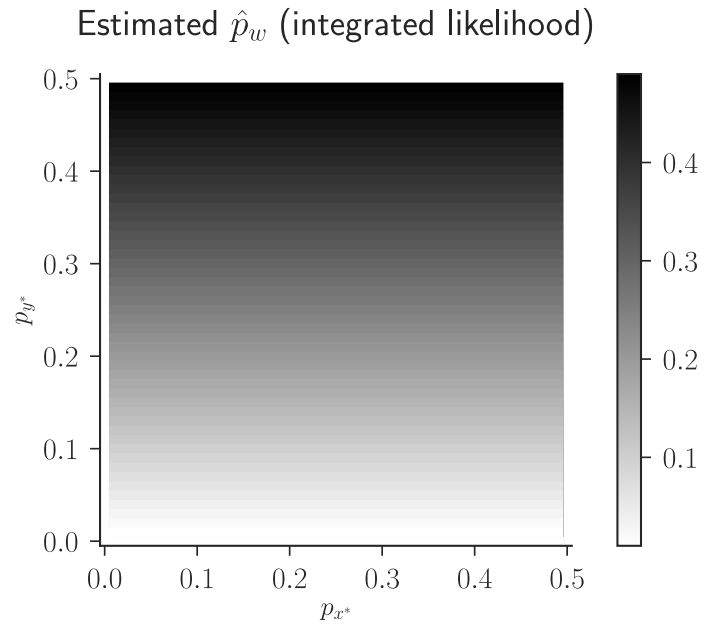


Figure S3: Estimates for \hat{p}_w when computing $(\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2, \hat{w})$ using L-BFGS-B optimizing the classical integrated likelihood (2) rather than a joint optimization procedure.

InvFels tree $\tau = \tau^*, t^* = \{x^*, y^*, x^*, y^*, y^*\}$

\tilde{y}_j	$p_{\tilde{y}_j}$	$8 \cdot \Pr(\Psi = \tilde{y}_j \mid \tau, t)$
\emptyset	p_{\emptyset}	$1 + (x^*)^2 + (y^*)^2 + 4x^*(y^*)^2 + (x^*)^2(y^*)^2$
$\{1\}$	p_1	$1 - (x^*)^2 + (y^*)^2 - (x^*)^2(y^*)^2$
$\{2\}$	p_2	$1 + (x^*)^2 - (y^*)^2 - (x^*)^2(y^*)^2$
$\{3\}$	p_3	$1 - (x^*)^2 + (y^*)^2 - (x^*)^2(y^*)^2$
$\{1, 2, 3\}$	p_{123}	$1 + (x^*)^2 - (y^*)^2 - (x^*)^2(y^*)^2$
$\{1, 2\}$	p_{12}	$1 - (x^*)^2 - (y^*)^2 + (x^*)^2(y^*)^2$
$\{2, 3\}$	p_{23}	$1 - (x^*)^2 - (y^*)^2 + (x^*)^2(y^*)^2$
$\{1, 3\}$	p_{13}	$1 + (x^*)^2 + (y^*)^2 - 4x^*(y^*)^2 + (x^*)^2(y^*)^2$

InvFels tree $\tau = \tau_1, t = \{x_1, y_1, x_2, y_2, w\}$

\tilde{y}_j	$p_{\tilde{y}_j}$	$8 \cdot \Pr(\Psi = \tilde{y}_j \mid \tau, t)$
\emptyset	p_{\emptyset}	$1 + x_1x_2 + y_1y_2 + w[x_1 + x_2][y_1 + y_2] + x_1y_1x_2y_2$
$\{1\}$	p_1	$1 - x_1x_2 + y_1y_2 + w[-x_1 + x_2][y_1 + y_2] - x_1y_1x_2y_2$
$\{2\}$	p_2	$1 + x_1x_2 - y_1y_2 + w[x_1 + x_2][-y_1 + y_2] - x_1y_1x_2y_2$
$\{3\}$	p_3	$1 - x_1x_2 + y_1y_2 + w[x_1 - x_2][y_1 + y_2] - x_1y_1x_2y_2$
$\{1, 2, 3\}$	p_{123}	$1 + x_1x_2 - y_1y_2 + w[x_1 + x_2][y_1 - y_2] - x_1y_1x_2y_2$
$\{1, 2\}$	p_{12}	$1 - x_1x_2 - y_1y_2 + w[-x_1 + x_2][-y_1 + y_2] + x_1y_1x_2y_2$
$\{2, 3\}$	p_{23}	$1 - x_1x_2 - y_1y_2 + w[x_1 - x_2][-y_1 + y_2] + x_1y_1x_2y_2$
$\{1, 3\}$	p_{13}	$1 + x_1x_2 + y_1y_2 + w[-x_1 - x_2][y_1 + y_2] + x_1y_1x_2y_2$

Felsenstein tree $\tau = \tau_2, t = \{x_1, y_1, x_2, y_2, w\}$

\tilde{y}_j	$p_{\tilde{y}_j}$	$8 \cdot \Pr(\Psi = \tilde{y}_j \mid \tau, t)$
\emptyset	p_{\emptyset}	$1 + x_1y_1 + x_2y_2 + w[x_1 + y_1][x_2 + y_2] + x_1y_1x_2y_2$
$\{1\}$	p_1	$1 - x_1y_1 + x_2y_2 + w[-x_1 + y_1][x_2 + y_2] - x_1y_1x_2y_2$
$\{2\}$	p_2	$1 - x_1y_1 + x_2y_2 + w[x_1 - y_1][x_2 + y_2] - x_1y_1x_2y_2$
$\{3\}$	p_3	$1 + x_1y_1 - x_2y_2 + w[x_1 + y_1][-x_2 + y_2] - x_1y_1x_2y_2$
$\{1, 2, 3\}$	p_{123}	$1 + x_1y_1 - x_2y_2 + w[-x_1 - y_1][-x_2 + y_2] - x_1y_1x_2y_2$
$\{1, 2\}$	p_{12}	$1 + x_1y_1 + x_2y_2 + w[-x_1 - y_1][x_2 + y_2] + x_1y_1x_2y_2$
$\{2, 3\}$	p_{23}	$1 - x_1y_1 - x_2y_2 + w[x_1 - y_1][-x_2 + y_2] + x_1y_1x_2y_2$
$\{1, 3\}$	p_{13}	$1 - x_1y_1 - x_2y_2 + w[-x_1 + y_1][-x_2 + y_2] + x_1y_1x_2y_2$

Table S1: 8 times the site split probabilities $p_{\tilde{y}_j}$ on the true InvFels tree τ^* with $t^* = \{x^*, y^*, x^*, y^*, y^*\}$, and on the InvFels tree τ_1 and Felsenstein tree τ_2 with $t = \{x_1, y_1, x_2, y_2, w\}$ obtained using the Hadamard transform.

\tilde{y}_j	\tilde{h}_k	$32 \cdot \Pr(\Psi = \tilde{y}_j, \Xi = \tilde{h}_k \mid \tau_1, t)$
\emptyset	\emptyset	$(1+x_1)(1+y_1)(1+x_2)(1+y_2)(1+w)$
	$\{1\}^*$	$(1-x_1)(1+y_1)(1-x_2)(1+y_2)(1-w)$
	$\{2\}^*$	$(1+x_1)(1-y_1)(1+x_2)(1-y_2)(1-w)$
	$\{1,2\}^*$	$(1-x_1)(1-y_1)(1-x_2)(1-y_2)(1+w)$
$\{1\}$	\emptyset	$(1-x_1)(1+y_1)(1+x_2)(1+y_2)(1+w)$
	$\{1\}$	$(1+x_1)(1+y_1)(1-x_2)(1+y_2)(1-w)$
	$\{2\}^*$	$(1-x_1)(1-y_1)(1+x_2)(1-y_2)(1-w)$
	$\{1,2\}$	$(1+x_1)(1-y_1)(1-x_2)(1-y_2)(1+w)$
$\{2\}$	\emptyset	$(1+x_1)(1-y_1)(1+x_2)(1+y_2)(1+w)$
	$\{1\}^*$	$(1-x_1)(1-y_1)(1-x_2)(1+y_2)(1-w)$
	$\{2\}$	$(1+x_1)(1+y_1)(1+x_2)(1-y_2)(1-w)$
	$\{1,2\}$	$(1-x_1)(1+y_1)(1-x_2)(1-y_2)(1+w)$
$\{3\}$	\emptyset	$(1+x_1)(1+y_1)(1-x_2)(1+y_2)(1+w)$
	$\{1\}$	$(1-x_1)(1+y_1)(1+x_2)(1+y_2)(1-w)$
	$\{2\}^*$	$(1+x_1)(1-y_1)(1-x_2)(1-y_2)(1-w)$
	$\{1,2\}$	$(1-x_1)(1-y_1)(1+x_2)(1-y_2)(1+w)$
$\{1,2,3\}$	\emptyset	$(1-x_1)(1-y_1)(1-x_2)(1+y_2)(1+w)$
	$\{1\}$	$(1+x_1)(1-y_1)(1+x_2)(1+y_2)(1-w)$
	$\{2\}^*$	$(1-x_1)(1+y_1)(1-x_2)(1-y_2)(1-w)$
	$\{1,2\}$	$(1+x_1)(1+y_1)(1+x_2)(1-y_2)(1+w)$
$\{1,2\}$	\emptyset	$(1-x_1)(1-y_1)(1+x_2)(1+y_2)(1+w)$
	$\{1\}$	$(1+x_1)(1-y_1)(1-x_2)(1+y_2)(1-w)$
	$\{2\}$	$(1-x_1)(1+y_1)(1+x_2)(1-y_2)(1-w)$
	$\{1,2\}$	$(1+x_1)(1+y_1)(1-x_2)(1-y_2)(1+w)$
$\{2,3\}$	\emptyset	$(1+x_1)(1-y_1)(1-x_2)(1+y_2)(1+w)$
	$\{1\}$	$(1-x_1)(1-y_1)(1+x_2)(1+y_2)(1-w)$
	$\{2\}$	$(1+x_1)(1+y_1)(1-x_2)(1-y_2)(1-w)$
	$\{1,2\}$	$(1-x_1)(1+y_1)(1+x_2)(1-y_2)(1+w)$
$\{1,3\}$	\emptyset	$(1-x_1)(1+y_1)(1-x_2)(1+y_2)(1+w)$
	$\{1\}$	$(1+x_1)(1+y_1)(1+x_2)(1+y_2)(1-w)$
	$\{2\}^*$	$(1-x_1)(1-y_1)(1-x_2)(1-y_2)(1-w)$
	$\{1,2\}$	$(1+x_1)(1-y_1)(1+x_2)(1-y_2)(1+w)$

Table S2: 32 times likelihood values for all site splits \tilde{y}_j and ancestral state splits \tilde{h}_k of the InvFels tree τ_1 . Ancestral states with * are never maximal provided parameters are in $(0, 1]$. By combinations of \tilde{h}_k , there are $3^5 \cdot 4^2 = 3,888$ possible forms for the likelihood.

\tilde{y}_j	$\eta_j(\tau_1, t)$	$\xi_{\tilde{y}_j}$	$32 \cdot \Pr(\Psi = \tilde{y}_j, \Xi = \xi_{\tilde{y}_j} \mid \tau_1, t)$
\emptyset	$\{\emptyset\}$	\emptyset	$(1 + x_1)(1 + y_1)(1 + x_2)(1 + y_2)(1 + w)$
$\{1\}$	$\{\emptyset\}$	\emptyset	$(1 - x_1)(1 + y_1)(1 + x_2)(1 + y_2)(1 + w)$
$\{2\}$	$\{\emptyset\}$	\emptyset	$(1 + x_1)(1 - y_1)(1 + x_2)(1 + y_2)(1 + w)$
$\{3\}$	$\{\emptyset, \{1\}, \{1, 2\}\}$	\emptyset	$(1 + x_1)(1 + y_1)(1 - x_2)(1 + y_2)(1 + w)$
		$\{1\}$	$(1 - x_1)(1 + y_1)(1 + x_2)(1 + y_2)(1 - w)$
		$\{1, 2\}$	$(1 - x_1)(1 - y_1)(1 + x_2)(1 - y_2)(1 + w)$
$\{1, 2, 3\}$	$\{\emptyset, \{1\}, \{1, 2\}\}$	\emptyset	$(1 - x_1)(1 - y_1)(1 - x_2)(1 + y_2)(1 + w)$
		$\{1\}$	$(1 + x_1)(1 - y_1)(1 + x_2)(1 + y_2)(1 - w)$
		$\{1, 2\}$	$(1 + x_1)(1 + y_1)(1 + x_2)(1 - y_2)(1 + w)$
$\{1, 2\}$	$\{\emptyset\}$	\emptyset	$(1 - x_1)(1 - y_1)(1 + x_2)(1 + y_2)(1 + w)$
$\{2, 3\}$	$\{\emptyset, \{1\}, \{1, 2\}\}$	\emptyset	$(1 + x_1)(1 - y_1)(1 - x_2)(1 + y_2)(1 + w)$
		$\{1\}$	$(1 - x_1)(1 - y_1)(1 + x_2)(1 + y_2)(1 - w)$
		$\{1, 2\}$	$(1 - x_1)(1 + y_1)(1 + x_2)(1 - y_2)(1 + w)$
$\{1, 3\}$	$\{\emptyset, \{1\}, \{1, 2\}\}$	\emptyset	$(1 - x_1)(1 + y_1)(1 - x_2)(1 + y_2)(1 + w)$
		$\{1\}$	$(1 + x_1)(1 + y_1)(1 + x_2)(1 + y_2)(1 - w)$
		$\{1, 2\}$	$(1 + x_1)(1 - y_1)(1 + x_2)(1 - y_2)(1 + w)$

Table S3: 32 times likelihood values on the InvFels tree τ_1 . Due to the symmetry of the likelihood, WLOG we let $x_2 \geq x_1$ and $y_2 \geq y_1$ and maximize over ancestral state splits to reduce the number of possible functional forms to consider. Likelihoods with multiple entries have maxima determined by unknown branch length parameters. Because in 4 cases there are 3 possibilities for $\xi_{\tilde{y}_j}$, there are $3^4 = 81$ possible forms for the likelihood.