

Supplemental Text for: Mistreating birth-death models as priors in phylogenetic analysis compromises our ability to compare models

Contents

S§1Simulation Study	S2
S§1.1Molecular substitution models	S2
S§1.2Contemporaneous birth-death processes	S3
S§1.3Fossilized birth-death processes	S3
S§1.4Extended Results	S4
S§2Empirical Analysis	S5
S§2.1Analysis	S5
S§2.2Results	S6
S§3Samples As Data	S9
S§3.1The standard phylogenetic model	S9
S§3.2A paleontological model	S10
S§3.3Combining phylogenetic and paleontological models	S11
S§3.4When sample ages are uncertain	S12
S§4Sequential Bayesian Inference	S14
S§5Factorizing Bayes' Theorem	S16
S§6Posterior-Predictive Simulation with Samples as Data	S17
S§7Random Variables and Priors in Phylogenetic Inference	S19

1 S§1 Simulation Study

2 To demonstrate the issues that arise from not treating samples as data, we performed a series of
3 experiments in RevBayes (Höhna et al. 2016), which we selected for convenience; the problems we
4 demonstrate are inherent to the standard treatment of tree models as priors and are independent of
5 any particular implementation of these models. For each of the three experiments below, we per-
6 formed a series of simulations under a specific model, and then compared the fit of a pair of com-
7 peting models to these simulated data, using Bayes factors calculated with marginal-likelihood- and
8 posterior-probability-based approaches. For the marginal-likelihood-based approach, we computed
9 Bayes factors using a power-posterior algorithm, stepping-stone MCMC (SS MCMC; Xie et al. 2011),
10 and for the posterior-probability-based approach we used reversible-jump MCMC (RJ MCMC; Green
11 1995). In all cases, we fixed the tree topology to the true value for computational tractability, but es-
12 timated the node ages, the birth-death parameters, and the parameters that governed the process of
13 character evolution; for fossilized birth-death datasets, we also estimated whether each fossil was a
14 sampled ancestor (*i.e.*, whether it was sampled along a branch leading to another sample).

15 Initial experiments indicated that sufficiently precise marginal-likelihood estimates using SS
16 MCMC would be computationally prohibitive for large trees and sequence datasets. We therefore
17 simulated relatively small trees and character datasets (exact sizes described below). Further, we
18 implemented an adaptive power-posterior algorithm in RevBayes, similar to the one proposed by
19 Friel et al. (2014). Power-posterior algorithms work by running a set of Markov chains, each with
20 a “power”, β_i , ranging between 0 and 1. For any given β_i (a “stone”), the chain samples from the
21 distorted posterior distribution:

$$P_i(\theta | X) \propto P_i(X | \theta)^{\beta_i} P(\theta),$$

22 so that $\beta = 1$ corresponds to sampling from the posterior, and $\beta = 0$ corresponds to sampling from
23 the prior. The sampled likelihood values among the separate stones— $P_i(X | \theta)$ —can then be used to
24 estimate the marginal likelihood, *e.g.*, using the stepping-stone estimator (Xie et al. 2011). Usually, the
25 number of stones and the values of β are fixed in advance, but in our analyses we found that accurate
26 marginal-likelihood estimates demanded a large number of stones, so adopted an adaptive approach.
27 Briefly, our adaptive algorithm begins with two stones, $\beta_1 = 1$ and $\beta_2 = 0$, and then places additional
28 stones until the estimate of the marginal likelihood converges; as with the original algorithm, the
29 number of MCMC samples per stone is fixed in advance.

30 For each analysis described below, we performed two replicates to ensure stability of marginal-
31 likelihood and posterior-ratio estimates. Our simulated data and code (including specific param-
32 eter settings for simulations and analyses) are available at Zenodo ([http://doi.org/10.5281/
33 zenodo.5072533](http://doi.org/10.5281/zenodo.5072533)) and GitHub ([https://github.com/mikeryanmay/bd_bayes_factors/releases/
34 tag/initial_submission](https://github.com/mikeryanmay/bd_bayes_factors/releases/tag/initial_submission)).

35 S§1.1 Molecular substitution models

36 To demonstrate that SS and RJ MCMC compute the same BFs (and also to demonstrate that both
37 of these methods are implemented correctly, *i.e.*, that our results are not a consequence of program-
38 ming errors), we compared the fit of competing substitution models to simulated molecular datasets.
39 We simulated ten trees under a birth-death (BD) model for each of four numbers of extant samples,
40 $n = \{8, 16, 32, 64\}$. We assumed the tree began with two species at time $t = 1$, diversified at rates
41 $\lambda = 4$ and $\mu = 2$, and that all extant species were sampled. For each tree, we simulated a nucleotide
42 dataset with 100 sites under a Jukes-Cantor (JC; Jukes and Cantor 1969) model with rate parameter r

43 (scaled such that the expected number of substitutions per site was three). For each of the 40 simu-
44 lated datasets, we computed Bayes factors between JC69 and K80 (Kimura 1980) substitution models
45 using SS and RJ MCMC as described above, assuming the tree evolved under the true birth-death
46 model (Fig. 2A, main text). Positive values of $2 \ln \text{BF}$ indicate support for the variable-rate model.

47 S§1.2 Contemporaneous birth-death processes

48 Our second experiment considers the case of comparing two birth-death models for extant (contem-
49 poraneous) samples. We analyzed the same datasets simulated in the previous section, but in this
50 case compared two tree models. The first model, M_1 , is the Yule model (with speciation rate λ , and
51 no extinction rate parameter), and the second model, M_2 , is a BD model (with speciation rate λ and
52 extinction rate μ); M_2 is the same tree model used to simulate the data, as described above. We com-
53 puted Bayes factors between Yule (M_1) and BD (M_2) processes, again using both SS and RJ MCMC,
54 assuming the sequence data evolved under the true JC69 substitution model (Fig. 2B, main text). In
55 this case, positive values of $2 \ln \text{BF}$ indicate support for the true model.

56 S§1.3 Fossilized birth-death processes

57 Our third experiment considers the more complex case of comparing two birth-death models for non-
58 contemporaneous samples. We simulated four fossilized birth-death trees under a model that allowed
59 the fossilization rate to vary. Specifically, each tree began with one lineage at time $t = 1$ (in the past,
60 with $t = 0$ the present) and initially evolved under a fossilized birth-death model with $\lambda = 4$, $\mu = 2$,
61 $\phi = 3$; at time $t = 0.5$ in the past, the fossilization rate changed to $\phi = 0.5$ (*i.e.*, the fossilization rate
62 was high in the early part of the process, and decreased in the second half by a factor of six). We
63 then simulated stratigraphic uncertainty by dividing time into 20 equally sized bins and using the
64 boundaries of the bin that a given fossil sample fell into as the minimum and maximum ages of
65 the sample (we assigned extant samples minimum and maximum ages of zero). For each tree, we
66 simulated 100 binary characters under an Mk model (Lewis 2001) with rate r (scaled such that the
67 expected number of substitutions per site was three). We then compared the fit of two competing
68 fossilized birth-death models to the simulated data. The first model, M_1 , has constant speciation,
69 extinction, and fossilization rates (λ , μ , and ϕ , respectively). The second model, M_2 , is the same as
70 M_1 , but allows the fossilization rate to vary over time. Specifically, the initial fossilization rate (at
71 time $t = 1$ in the past) is ϕ_1 , and at time $t = 0.5$ units in the past, it changes to rate ϕ_2 , which persists
72 until the present ($t = 0$). M_2 is similar to the simulating process in that the fossilization rate is not
73 constant, and the time of the rate change is fixed; however, for both M_1 and M_2 , we assume that the
74 speciation, extinction, and fossilization rates are unknown. For both models, we also assume that
75 all extant species are sampled, $\rho = 1$. We computed the Bayes factors between M_1 and M_2 using SS
76 and RJ MCMC, assuming the morphological data evolves under the Mk model (Fig. 2C, main text).
77 Positive values of $2 \ln \text{BF}$ indicate support for the variable-rate model.

78 **S§1.4 Extended Results**

79 Our results indicate that SS and RJ MCMC provide essentially identical Bayes factors when comparing
 80 models of molecular evolution (simulation 1, Fig. 1A, main text), as expected based on the theoretical
 81 equivalence of these estimators. However, these method produce disparate estimates when compar-
 82 ing tree models, either for contemporaneous lineages (simulation 2, Fig. 1B, main text) or when
 83 including non-contemporaneous lineages (simulation 3, Fig. 1C, main text).

84 The discrepancy is not a consequence of a programming bug: the comparison of substitution mod-
 85 els demonstrates that both algorithms are correctly implemented. Likewise, it is not a consequence of
 86 numerical MCMC errors: we performed two replicates of each of the analysis to confirm that Bayes-
 87 factor estimates were sufficiently precise both for SS-based estimates (Fig. S1) and RJ-based estimates
 88 (Fig. S2). For the contemporaneous birth-death models, we can make precise quantitative predictions
 89 about the magnitude of the discrepancy (Fig. 1B, colored dashed lines, main text), as we explain in
 90 Section S§4; our simulation results match these predictions.

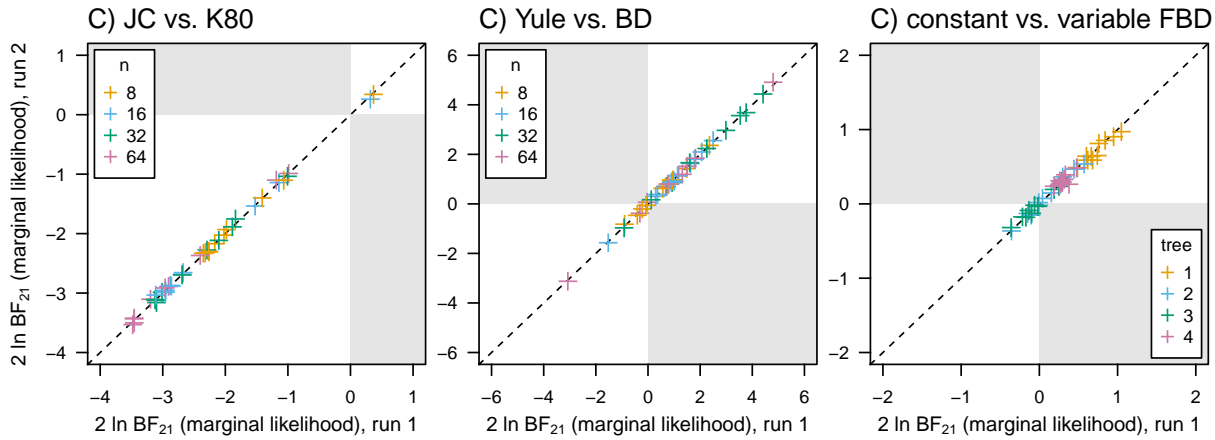


Figure S1: Precision of Bayes factor calculation using marginal likelihoods. Each analysis was performed twice, and the value from one run (x-axis) is plotted against the second run (y-axis). A) Bayes factors between JC and K80 models. B) Bayes factors between Yule and birth-death models. C) Bayes factors between a model with constant fossilization rates and one with variable fossilization rates. (see caption of Fig. 1)

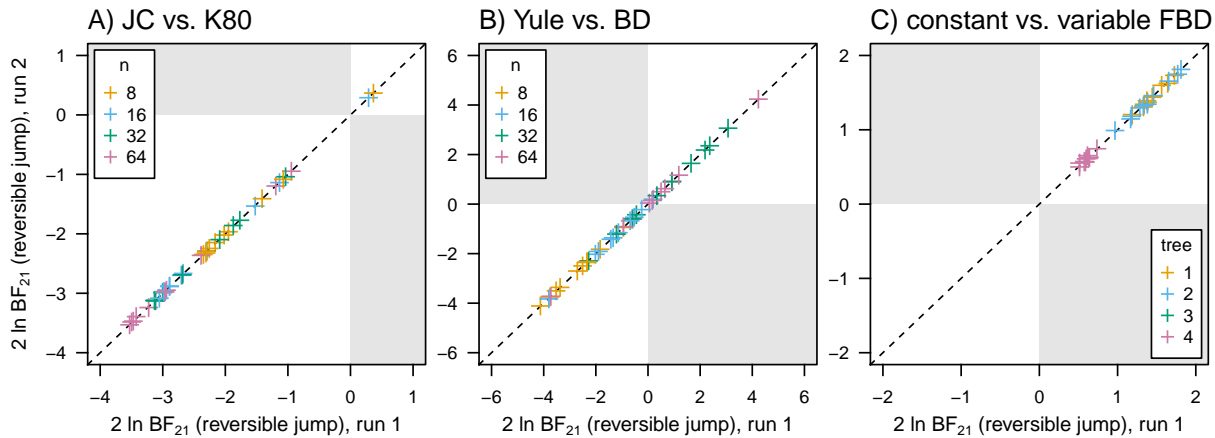


Figure S2: Precision of Bayes factor calculation using reversible-jump MCMC. Each analysis was performed twice, and the value from one run (x-axis) is plotted against the second run (y-axis). A) Bayes factors between JC and K80 models. B) Bayes factors between Yule and birth-death models. C) Bayes factors between a model with constant fossilization rates and one with variable fossilization rates. (see caption of Fig. 1)

91 S§2 Empirical Analysis

92 S§2.1 Analysis

93 We re-analyzed the empirical dataset of marattialean ferns from our previous study (May et al. 2021)
94 to demonstrate the impact of not treating samples as data in a realistic model-comparison scenario, as
95 well as to provide an example of the impact of the tree model on estimates of divergence times. For
96 the sake of computational tractability, we included only ingroup taxa—comprising 26 extant and 45
97 extinct samples—and only analyzed the binary morphological data from that study (*i.e.*, we excluded
98 the molecular data and multistate morphological characters).

99 We analyzed this dataset under an MkV model (Lewis 2001) with gamma-distributed rate varia-
100 tion among characters (Yang 1994), and assumed that rates of morphological evolution varied across
101 branches of the tree according to an uncorrelated-lognormal relaxed-clock model (UCLN Drummond
102 et al. 2006). We compared the fit of two competing birth-death models: 1) a fossilized birth-death
103 model where speciation and extinction rates varied over time, but the fossilization rate was constant
104 over time (the “constant” model), and; 2) a fossilized birth-death model where fossilization, specia-
105 tion, and extinction rates each varied over time (the “variable” model). Specifying arbitrary variable-
106 rate fossilized birth-death models that are amenable to efficient reversible-jump MCMC is non-trivial.
107 We therefore used the results from our previous study to constrain the rate variation so that it was
108 both appropriate for the dataset, and possible to specify in the existing reversible-jump machinery
109 available in RevBayes. For the speciation- and extinction-rate variation, we assumed that these rates
110 varied according to a piecewise-constant model defined by five time intervals intended to capture the
111 major patterns present in our prior analyses: $(\infty, 323.2]$, $(323.2, 298.9]$, $(298.9, 66.0]$, $(66.0, 5.3]$, $(5.3, 0.0]$.
112 Within each time interval, speciation and extinction rates were drawn independently from a shared
113 prior distribution. For the model that allowed the fossilization rate to vary, we assumed a piecewise
114 constant model with four time intervals: $(\infty, 323.2]$, $(323.2, 298.9]$, $(298.9, 66.0]$, $(66.0, 0.0]$, with corre-
115 sponding fossilization rates $\{\psi_1, \psi_2, \psi_3, \psi_4\}$. We assumed that the rates of the first and fourth interval
116 were the same ($\psi_1 = \psi_4$), but allowed the fossilization rates for the second and third intervals to be
117 different, reflecting an apparent peak in fossilization rates in the Pennsylvanian (the second interval),
118 followed by moderate fossilization rates from the Permian to the end of the Cretaceous (the third
119 interval).

120 For both tree models, we fixed the tree topology to the maximum-clade-credibility (MCC) tree
121 topology inferred in our previous study, but estimated the node ages, the fossilized birth-death param-
122 eters, the character-evolution parameters, and also whether each fossil was a sampled ancestor. We
123 then computed the Bayes factors between the constant and variable models using SS and RJ MCMC.
124 Our empirical data and code for these analyses (including specific parameter and prior settings) are
125 available at [XXXX](#).

126 **S§2.2 Results**

127 Under the variable model, we infer extreme variation in fossilization rates: rates are inferred to be
 128 substantially higher during the Pennsylvanian than in the other time intervals, and rates from the
 129 Permian to the Late Cretaceous are also elevated compared to the first and fourth intervals (Fig. S3).
 130 This result is unsurprising, given that 24 of the 45 fossil samples come from the 23 My window that
 131 constitutes the Pennsylvanian subperiod. Despite this evident rate variation, BFs based on marginal
 132 likelihoods favor the constant model ($2 \ln \text{BF} \approx 3$, Tables S.1 and S.2). By contrast, BFs based on
 133 posterior model probabilities decisively favor the variable model ($2 \ln \text{BF} \approx 18$, Table S.2). In other
 134 words, conventional marginal-likelihood-based BFs incorrectly indicate strong evidence for a deci-
 135 sively worse model.

Table S.1: Marginal likelihoods for the constant model computed with stepping-stone sampling. We performed four independent runs to ensure precise marginal likelihood estimates (runs 1 through 4); we report the mean and standard error of the mean (final column).

Model	Run 1	Run 2	Run 3	Run 4	Mean (\pm SD)
Constant	-883.0418	-882.8819	-882.9806	-882.9524	-882.9602 (± 0.03634)
Variable	-884.5965	-884.5084	-884.4403	-884.4645	-884.5025 (± 0.03136)

Table S.2: $2 \ln$ Bayes factors between constant and variable models, computed using two methods. We performed four independent runs to ensure precise Bayes factor estimates (runs 1 through 4); we report the mean and standard error of the mean (final column).

Method	Run 1	Run 2	Run 3	Run 4	Mean (\pm SD)
SS	-3.0770	-3.3178	-2.9194	-3.0242	-3.0846 (± 0.0843)
RJ	18.1055	17.9870	18.0041	18.0804	18.0442 (± 0.0287)

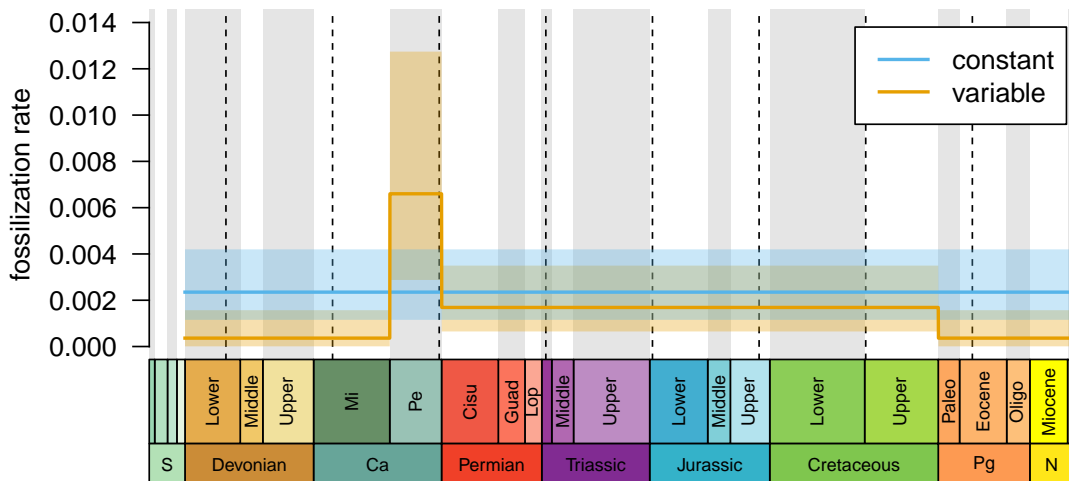


Figure S3: Estimated fossilization rates under the constant fossilization-rate model (blue) and the variable fossilization-rate model (orange). Dark lines correspond to the mean posterior rate at each time point, and colored regions correspond to the 95% credible interval. Dashed lines are placed at 50 My intervals.

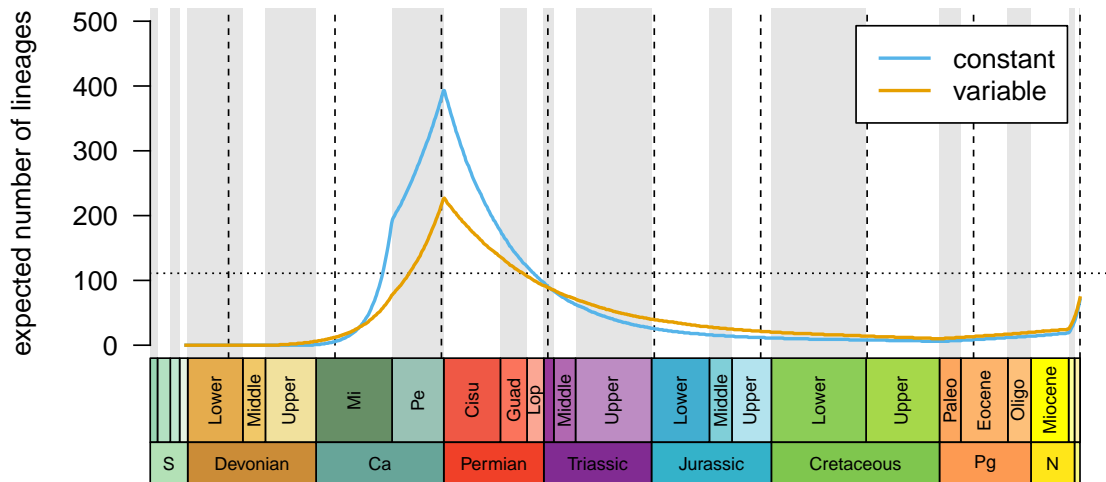


Figure S4: Median number of lineages over time, as predicted by the constant fossilization-rate model (blue) and the variable fossilization-rate model (orange). Vertical dashed lines are placed at 50 My intervals; horizontal dashed line is the number of extant lineages.

136 The different models produce significantly different inferences about the history of diversity for
 137 the group. We simulated 10,000 histories of lineage diversification under each model, and discretized
 138 time into many (10,000) small time intervals. We then computed the median number of lineages alive
 139 in each time interval over the group’s history (Fig. S4). Under M_1 , we predict an average peak diversity of ≈ 945 lineages in the Pennsylvanian; by contrast, M_2 predicts ≈ 546 lineages at that time. This
 140 discrepancy likely reflects the fact that M_1 requires a larger number of lineages during the Pennsylvanian
 141 in order to preserve the observed number of samples, given that fossilization rates are lower at
 142 that time relative to M_2 (Fig. S3).
 143

144 Beside providing qualitatively different inferences about the nature of the fossilization process
 145 and the underlying history of diversification, the two tree models strongly influence divergence-time
 146 estimates for this dataset. We computed the mean of the posterior distribution of the age of each node
 147 in the tree under both models as $\Delta = |a_v - a_c|$, where a_c and a_v represent the posterior-mean estimate
 148 under the constant and variable models, respectively. Posterior-mean estimates of divergence times
 149 under the constant model differ substantially from those under the variable model: some nodes ages
 150 are different by 15 million years between the two models (Fig. S5, left). Divergence-time estimates
 151 for young nodes are systematically more recent under the constant model, *i.e.*, the younger nodes are
 152 disproportionately pulled toward the present (Fig. S6). Presumably, this pattern reflects the fact that
 153 the fossilization rates in the Cenozoic are higher under the constant model than the variable model
 154 (Fig. S3): when fossilization rates are high, older clades in these time intervals imply more missing
 155 fossils, and are therefore “penalized” by the model (*i.e.*, younger clades are preferred).

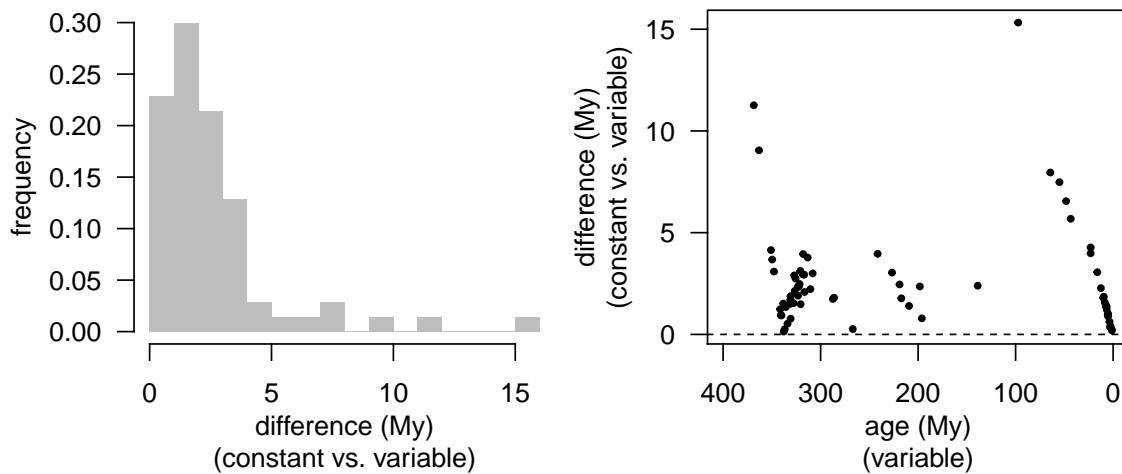


Figure S5: Left) Absolute differences in posterior mean node-age estimates under the two tree models. Right) Absolute differences in posterior mean node-age estimates under the two tree models as a function of age (the mean age estimated under the variable model).

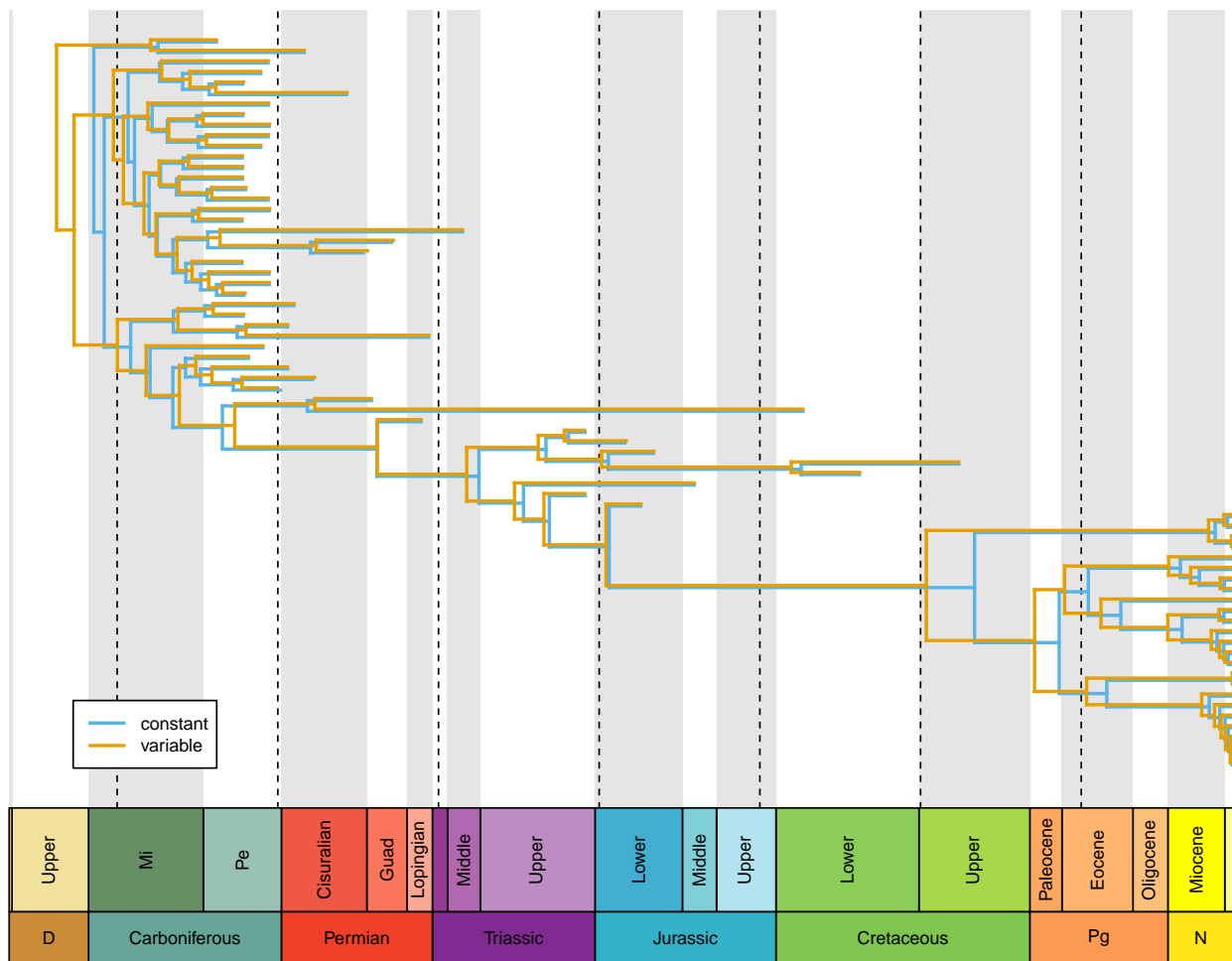


Figure S6: Phylogenies estimated under the constant (blue) and variable (orange) fossilized birth-death model. Node ages correspond to the the posterior mean age inferred under the model. Dashed lines are placed at 50 My intervals.

156 S§3 Samples As Data

157 S§3.1 The standard phylogenetic model

158 In a standard fossilized birth-death analysis (Heath et al. 2014; Zhang et al. 2016), we imagine that
159 a phylogeny, Ψ , evolves under a fossilized birth-death process with a set of parameters θ_Ψ , which
160 consists of speciation, extinction, and fossilization rates (λ , μ , and ϕ , respectively), and a sampling
161 probability for extant species, ρ . Additionally, we imagine that a set of characters (typically a com-
162 bination of molecular and discrete morphological characters) evolve along the branches of the tree
163 according to a defined model with a set of parameters θ_x . These processes give rise to a character
164 dataset, X .

165 Given a particular character dataset, X , our goal is to infer the phylogeny and model parameters
166 that gave rise to that dataset. To do so, we can apply Bayes' theorem:

$$P(\Psi, \theta_x, \theta_\Psi | X) = \frac{\overbrace{P(X | \Psi, \theta_x)}^{\text{likelihood}} \overbrace{P(\Psi | \theta_\Psi) P(\theta_x) P(\theta_\Psi)}^{\text{priors}}}{\underbrace{P(X)}_{\text{marginal likelihood}}}. \quad (\text{S.1})$$

167 The prior distributions $P(\theta_x)$ and $P(\theta_\Psi)$ represent our belief about these model components before
168 observing the data. $P(\Psi | \theta_\Psi)$ is generally considered a prior distribution on the phylogeny, and is
169 the probability density of a “ranked, labeled¹” tree (see equation [2] in Gavryushkina et al. 2014).
170 These prior distributions are updated by information in the character data by the likelihood function,
171 $P(X | \Psi, \theta_x)$, to produce our posterior belief about the tree and model parameters, $P(\Psi, \theta_x, \theta_\Psi | X)$.
172 The denominator— $P(X)$, *i.e.*, the marginal likelihood—is the likelihood function averaged over all of
173 the model parameters in proportion to their prior probability:

$$P(X) = \iiint P(X | \Psi, \theta_x) P(\Psi | \theta_\Psi) P(\theta_x) P(\theta_\Psi) d\theta_x d\theta_\Psi d\Psi \quad (\text{S.2})$$

174 (The integrals here represent multidimensional integration over θ_x and θ_Ψ , as well as summation
175 over all possible tree topologies and integration over all sets of branch lengths in Ψ ; we omit the
176 domains of integration throughout this document for the sake of simplicity.) Different models will
177 have different marginal likelihoods; we indicate the marginal likelihood for model i as $P_i(X)$ when
178 we need to distinguish among models.

179 Typically, we refer to probabilities of data (observations) as likelihoods (technically, in a frequen-
180 tist framework, the likelihood of the parameters is proportional to the probability of the data given
181 the parameters), and probabilities of parameters as prior probabilities. The labeling of terms affects
182 marginal likelihoods estimated using standard methods, which depend on treating likelihoods and
183 priors differently (*e.g.*, Lartillot and Philippe 2006; Xie et al. 2011). However, the number of samples
184 and their ages are observations that provide information about the lineage-diversification process,
185 independent from the character data; indeed, paleontologists regularly use this type of information,
186 by itself, to estimate parameters of lineage-diversification models. From this perspective, the “prior
187 probability” of the phylogeny, $P(\Psi | \theta_\Psi)$, actually represents the joint probability of the samples and
188 the phylogeny, and some of this probability—the portion related to the samples—belongs to the like-
189 lihood.

¹Here, “ranked” means the nodes in the tree have ages, and “labeled” means the tips and sampled ancestors in the tree are associated with specific named samples in our dataset (Murtagh 1984).

190 **S§3.2 A paleontological model**

191 There is a long history of studying lineage diversification from a purely paleontological perspective
 192 (e.g., Raup 1975, 1985; Foote 2000, 2001). Some of these methods, particularly PyRates (Silvestro et al.
 193 2014, 2019), use stochastic birth-death models that are effectively interchangeable with those used in
 194 phylogenetic methods. In this framework, the data are taken to be fossil occurrences of the clade(s)
 195 of interest, and the goal is to estimate speciation (origination) and extinction rates based on how the
 196 fossil occurrences are distributed over time. In contrast to phylogenetic methods, the relationships
 197 among the fossil occurrences are not of direct interest, so this method does not depend on character
 198 data or models of character evolution.

199 In a Bayesian framework, we can conceive of a generic “paleontological” model (similar to
 200 PyRates) that would be represented as:

$$P(\lambda, \mu, \phi, \rho | S) = \frac{\overbrace{P(S | \lambda, \mu, \phi, \rho)}^{\text{likelihood}} \overbrace{P(\lambda)P(\mu)P(\phi)P(\rho)}^{\text{priors}}}{\underbrace{P(S)}_{\text{marginal likelihood}}}, \quad (\text{S.3})$$

201 where S represents the set of samples, comprising individual fossil occurrences—the ages of the
 202 specimens, which for the time being we assume are known exactly, *i.e.*, there is no stratigraphic-age
 203 uncertainty—as well as any extant members of the group being analyzed, and the diversification pa-
 204 rameters are the same as those described above for the fossilized birth-death model. Here, it is clear
 205 that $P(S | \lambda, \mu, \phi, \rho)$ is the probability of the observations, and is thus (proportional to) the likelihood
 206 function. It may be difficult to compute the probability of the samples without knowing the complete
 207 tree (including unsampled lineages), $\tilde{\Psi}$. However, if we can compute the conditional probability of
 208 the samples given the complete tree, in principle we can write the unconditional probability of the
 209 samples as:

$$\begin{aligned} \underbrace{P(S | \lambda, \mu, \phi, \rho)}_{\text{probability of the samples}} &= \int \overbrace{P(S, \tilde{\Psi} | \lambda, \mu, \phi, \rho)}^{\text{joint probability of the tree and samples}} d\tilde{\Psi} \\ &= \int \underbrace{P(S | \tilde{\Psi}, \phi, \rho)}_{\substack{\text{conditional probability} \\ \text{of the samples,} \\ \text{given the tree}}} \underbrace{P(\tilde{\Psi} | \lambda, \mu)}_{\text{probability of the tree}} d\tilde{\Psi}, \end{aligned} \quad (\text{S.4})$$

210 where the integral represents integration over all possible complete phylogenies, $\tilde{\Psi}$, in proportion
 211 to their probability. Given the phylogeny, the probability of the samples is independent of the di-
 212 versification parameters, λ and μ ; likewise, the probability of the full tree does not depend on the
 213 sampling parameters, ϕ and ρ . (Note that this complete tree is different from the “reconstructed” tree
 214 in equation [S.1], Ψ , which hides all the unsampled lineages.)

215 While the probability of the samples given the complete tree— $P(S | \tilde{\Psi}, \phi, \rho)$ —may be relatively
 216 easy to compute, the above integral may be quite difficult to compute. However, a simple solution
 217 would be to include the complete phylogeny in the posterior distribution as an additional parameter:

$$P(\tilde{\Psi}, \lambda, \mu, \phi, \rho | S) = \frac{P(S | \tilde{\Psi}, \phi, \rho)P(\tilde{\Psi} | \lambda, \mu)P(\lambda)P(\mu)P(\phi)P(\rho)}{P(S)}. \quad (\text{S.5})$$

218 Importantly, this will produce the same posterior estimates of λ , μ , ϕ , and ρ as produced by equation
 219 (S.3) when we integrate over the phylogeny:

$$\int P(\tilde{\Psi}, \lambda, \mu, \phi, \rho | S) d\tilde{\Psi} = \int \frac{P(S | \tilde{\Psi}, \phi, \rho) P(\tilde{\Psi} | \lambda, \mu) P(\lambda) P(\mu) P(\phi) P(\rho)}{P(S)} d\tilde{\Psi}$$

$$P(\lambda, \mu, \phi, \rho | S) = \frac{\left[\int P(S | \tilde{\Psi}, \phi, \rho) P(\tilde{\Psi} | \lambda, \mu) d\tilde{\Psi} \right] P(\lambda) P(\mu) P(\phi) P(\rho)}{P(S)}$$

$$P(\lambda, \mu, \phi, \rho | S) = \frac{P(S | \lambda, \mu, \phi, \rho) P(\lambda) P(\mu) P(\phi) P(\rho)}{P(S)}.$$

220 Additionally, both representations have the same marginal likelihood, $P(S)$. The equation (S.5) im-
 221 plies an approach that is very similar to the approach used by the Bayesian program PyRates (Silve-
 222 stro et al. 2014), which estimates the times of origin and extinction of every sampled lineages (much
 223 like branch lengths in the complete phylogeny), but not the relationships between lineages (*i.e.*, the
 224 phylogenetic topology).

225 S§3.3 Combining phylogenetic and paleontological models

226 While the phylogenetic and paleontological models use the same underlying models of lineage diver-
 227 sification and sampling, they do not appear to use the same data (the phylogenetic model uses the
 228 character data, X , whereas the paleontological model uses the sample data, S). We can resolve this
 229 apparent discrepancy by simply adding character data, X , and a model of character evolution to the
 230 paleontological model that includes an implicit tree (equation [S.5]). In this case, the data are both S
 231 and X , and the corresponding posterior distribution is:

$$P(\tilde{\Psi}, \theta_x, \lambda, \mu, \phi, \rho | S, X) = \frac{\overbrace{P(X | \tilde{\Psi}, \theta_x) P(S | \tilde{\Psi}, \phi, \rho)}^{\text{likelihood}} \overbrace{P(\tilde{\Psi} | \lambda, \mu) P(\theta_x) P(\lambda) P(\mu) P(\phi) P(\rho)}^{\text{priors}}}{\underbrace{P(S, X)}_{\text{marginal likelihood}}}.$$

232 We can attempt to derive something like the standard phylogenetic model—equation (S.1)—from
 233 this combined model by first recognizing $P(S | \tilde{\Psi}, \phi, \rho) P(\tilde{\Psi} | \lambda, \mu)$ as the joint probability of the sam-
 234 ples and full tree, $P(S, \tilde{\Psi} | \lambda, \mu, \phi, \rho)$ (see equation [S.4]):

$$P(\tilde{\Psi}, \theta_x, \lambda, \mu, \phi, \rho | S, X) = \frac{\overbrace{P(X | \tilde{\Psi}, \theta_x)}^{\text{likelihood}} \overbrace{P(S, \tilde{\Psi} | \lambda, \mu, \phi, \rho)}^{\text{some likelihood, some prior}} \overbrace{P(\theta_x) P(\lambda) P(\mu) P(\phi) P(\rho)}^{\text{prior}}}{P(S, X)}.$$

235 Next, we must reduce $\tilde{\Psi}$ to Ψ . For any given Ψ , there are an infinite number of unobserved histories
 236 consistent with Ψ , each of which produces a unique $\tilde{\Psi}$. We label the unobserved history Ψ^c , and
 237 say that $\tilde{\Psi} = \{\Psi, \Psi^c\}$. The posterior distribution of the phylogenetic model should integrate over all

238 possible unobserved histories in proportion to their probability:

$$\begin{aligned}
P(\Psi, \theta_x, \lambda, \mu, \phi, \rho \mid S, X) &= \int P(\Psi, \Psi^c, \theta_x, \lambda, \mu, \phi, \rho \mid S, X) d\Psi^c \\
&= \frac{\left[\int P(X \mid S, \Psi, \Psi^c, \theta_x) P(S, \Psi, \Psi^c \mid \lambda, \mu, \phi, \rho) d\Psi^c \right] P(\theta_x) P(\lambda) P(\mu) P(\phi) P(\rho)}{P(S, X)} \\
&= \frac{P(X \mid S, \Psi, \theta_x) P(S, \Psi \mid \lambda, \mu, \phi, \rho) P(\theta_x) P(\lambda) P(\mu) P(\phi) P(\rho)}{P(S, X)} \\
P(\Psi, \theta_x, \theta_\Psi \mid S, X) &= \frac{\underbrace{P(X \mid S, \Psi, \theta_x)}_{\text{likelihood}} \underbrace{P(S, \Psi \mid \theta_\Psi)}_{\substack{\text{some likelihood,} \\ \text{some prior}}} \underbrace{P(\theta_x) P(\theta_\Psi)}_{\text{prior}}}{\underbrace{P(S, X)}_{\text{marginal likelihood}}}. \tag{S.6}
\end{aligned}$$

239 (We include S as a dependency in $P(X \mid S, \Psi, \theta_x)$ because Ψ is a function of S and $\tilde{\Psi}$.)

240 Equation (S.6) is very similar to the standard phylogenetic representation, but with critical differ-
241 ences. First, the data include both X and S , rather than just X ; consequently, the marginal likelihoods
242 must be different. Second, the prior probability of the sampled tree— $P(\Psi \mid \theta_\Psi)$ in equation (S.1)—has
243 been replaced with the joint probability of the samples and the sampled tree, $P(S, \Psi \mid \theta_\Psi)$. Because
244 this joint probability includes observations, some part of it should be regarded as part of the likelihood
245 of the model.

246 These equations demonstrate that probabilities we are used to thinking of as prior probabilities—
247 specifically, the probabilities of trees under a birth-death model—are actually an ambiguous mixture
248 of likelihood-like and prior-like quantities. That is, the likelihood and prior functions in the standard
249 Bayesian model are mislabeled. We explore the consequences of this mislabeling in Section S§4.

250 S§3.4 When sample ages are uncertain

251 So far, we have assumed that the ages of the fossil occurrences are known without error, which helps to
252 clarify our main argument that samples should be treated as data. However, in real datasets, the ages
253 of fossil specimens are often uncertain, because the age of the sediments in which the specimens are
254 found can only be known within a certain interval. This phenomenon—referred to as stratigraphic-
255 age uncertainty—is somewhat orthogonal to our argument, but we mention it here because previous
256 work has argued that stratigraphic-age uncertainty should be treated a part of the likelihood function
257 [Drummond and Stadler \(2016\)](#). We agree with this perspective, and show how it fits in to the frame-
258 work we outlined above. Unfortunately, stratigraphic-age uncertainty leads to additional challenges
259 when computing marginal likelihoods.

260 Following [Drummond and Stadler \(2016\)](#), we represent stratigraphic-age data as
261 $A = \{a_1, a_2, \dots, a_n\}$ for the n samples, where $a_i = \{\hat{a}_i, \hat{a}_i^\vee\}$ are the minimum and maximum ages
262 of the i^{th} sample, respectively. The probability of the data (the stratigraphic ranges) would then be:

$$\begin{aligned}
P(A \mid \lambda, \mu, \phi, \rho) &= \int P(A, S \mid \lambda, \mu, \phi, \rho) dS \\
&= \int P(A \mid S) P(S \mid \lambda, \mu, \phi, \rho) dS, \tag{S.7}
\end{aligned}$$

263 where the integration is over the exact ages of all the samples, S , and $P(A \mid S)$ is a product of indicator

264 functions:

$$P(A | S) = \prod_i^n P(a_i | S_i),$$

265 with

$$P(a_i | S_i) = \begin{cases} 1 & \text{if } \check{a}_i \leq s_i \leq \hat{a}_i \\ 0 & \text{otherwise.} \end{cases}$$

266 Including the full tree, the likelihood with stratigraphic age uncertainty becomes:

$$P(A | \lambda, \mu, \phi, \rho) = \iint P(A | S) P(S | \tilde{\Psi}, \phi, \rho) P(\tilde{\Psi} | \lambda, \mu) d\tilde{\Psi} dS.$$

267 Even if this integral were analytically tractable, we could not use it with character data, because the
268 probability of the character data will generally depend on the exact ages of the samples. However, we
269 can use data augmentation (Tanner and Wong 1987) to include the exact ages in the model, and write
270 the full posterior distribution:

$$P(S, \Psi, \theta_x, \theta_\Psi | A, X) = \frac{P(X | S, \Psi, \theta_x) P(A | S) P(S, \Psi | \theta_\Psi) P(\theta_x) P(\theta_\Psi)}{P(A, X)},$$

271 where now the data are A and X . This approach amounts to a data augmentation because the likeli-
272 hood should average over the exact sample ages, S , as implied by equation (S.7).

273 Without special machinery, generic methods for computing the marginal likelihood that depend
274 on raising the likelihood function to a power cannot effectively deal with data-augmented models
275 (Rodrigue and Aris-Brosou 2011). As a consequence, correct solutions for marginal-likelihood estima-
276 tors with stratigraphic uncertainty are currently unavailable.

277 **S§4 Sequential Bayesian Inference**

278 Above, we showed that the traditional phylogenetic model mistreats the probability of the samples
 279 as part of the prior rather than the likelihood function. Here, we use the principle of sequential
 280 Bayesian inference to understand the quantitative consequences of this error. The parameters of the
 281 prior distributions (hyperparameters) we choose for a Bayesian model represent our prior belief about
 282 plausible parameter values, and in principle reflect our previous experiences with analyzing relevant
 283 data (or ignorance, if we have no previous experience). In a sense, when informed by previous analy-
 284 sis, these hyperparameters encapsulate the information in the previous datasets about the parameters,
 285 *i.e.*, they can be viewed as “old” data. When we analyze a “new” dataset, we update our prior beliefs
 286 accordingly. We can repeat this process indefinitely, as we collect additional datasets. This sequential
 287 Bayesian updating process is the basis of Lindley’s aphorism that “today’s posterior is tomorrow’s
 288 prior” (Lindley 1972).

289 When we perform a Bayesian phylogenetic analysis under a birth-death model, we can imagine
 290 collecting two datasets. We first collect samples, represented by their ages S . We may infer the tree
 291 model (birth-death and sampling) parameters, θ_Ψ , directly from this dataset. We can write the poste-
 292 rior distribution of this model as:

$$P(\theta_\Psi | S) = \frac{\overbrace{P(S | \theta_\Psi)}^{\text{likelihood}} \overbrace{P(\theta_\Psi)}^{\text{prior}}}{\underbrace{P(S)}_{\text{marginal likelihood of the samples}}}, \tag{S.8}$$

293 which corresponds to the posterior distribution of a “paleontological” model (equation [S.3]).

294 However, if we then become additionally interested in the phylogenetic relationships themselves,
 295 we can assemble a character dataset, X . Rather than re-doing the initial analysis, we may apply the
 296 principle of sequential Bayesian updating and use the first posterior as a prior in our second analysis:

$$P(\Psi, \theta_x, \theta_\Psi | X, S) = \frac{\overbrace{P(X | S, \Psi, \theta_x)}^{\text{likelihood}} \overbrace{P(\Psi | \theta_\Psi) P(\theta_\Psi | S) P(\theta_x)}^{\text{priors}}}{\underbrace{P(X | S)}_{\text{marginal likelihood of the characters given the samples}}}. \tag{S.9}$$

297 In this equation, S is effectively treated as a hyperparameter, *i.e.*, a fixed parameter of the prior dis-
 298 tribution on θ_Ψ . (We note that *all* prior distributions have hyperparameters, but we usually exclude
 299 them from our notation for simplicity.)

300 Alternatively, we could start again and do both analyses simultaneously (jointly). The posterior of
 301 such a joint analysis would be

$$P(\Psi, \theta_x, \theta_\Psi | X, S) = \frac{\overbrace{P(X | S, \Psi, \theta_x)}^{\text{likelihood}} \overbrace{P(S, \Psi | \theta_\Psi)}^{\text{some likelihood, some prior}} \overbrace{P(\theta_\Psi) P(\theta_x)}^{\text{priors}}}{\underbrace{P(X, S)}_{\text{marginal likelihood of samples and characters}}}, \tag{S.10}$$

302 where we view the likelihood function as $P(X | \Psi, \theta_x)$, as well as some contribution from $P(S, \Psi | \theta_\Psi)$,
 303 as we explain in Section S§3. We can verify that the posterior distribution from the joint analysis, equa-
 304 tion (S.10), is equivalent to the posterior distribution after the second step of the sequential analysis

305 by substituting equation (S.8) into equation (S.9), and recognizing that

$$P(X, S) = P(X | S)P(S),$$

306 *i.e.*, that the marginal likelihood of the joint analysis is the product of the marginal likelihoods of each
307 step in the sequential analysis.

308 The remaining task is to explain why SS and RJ MCMC estimate different Bayes factors. Methods
309 for calculating the marginal likelihood, such as SS, require that we clearly distinguish the probability
310 terms that are “likelihood” from those that are “prior”. In the standard phylogenetic notation (equa-
311 tion [S.1]), the probability of the character data is labeled the likelihood, while the joint probability of
312 the tree and samples is labeled the prior. This corresponds to the labeling in the second step of the
313 sequential analysis, equation (S.9), in which case the marginal likelihood is $P(X | S)$; this marginal
314 likelihood perceives the samples as “old” data, and only computes the marginal likelihood of the
315 “new” data, X . When we compare two models in this way, we are essentially imagining that for each
316 model, we first update the priors according to S , and then compute the marginal likelihood of X given
317 the corresponding posteriors from the first step. The resulting Bayes factor between the two models i
318 and j will be

$$\text{BF}_{ij} = \frac{P_i(X | S)}{P_j(X | S)}, \quad (\text{S.11})$$

319 where $P_k(X | S)$ is the marginal likelihood of model k . By contrast, RJ MCMC does not depend on the
320 labeling of probability terms. In this case, the dataset implicitly includes both samples and character
321 data, and the Bayes factors will be

$$\text{BF}_{ij} = \frac{P_i(X, S)}{P_j(X, S)} = \frac{P_i(X | S) P_i(S)}{P_j(X | S) P_j(S)}, \quad (\text{S.12})$$

322 Equations (S.11) and (S.12) predict that the discrepancy in Bayes factors between SS and RJ MCMC
323 should be equal to the ratio of the marginal likelihoods of the samples, $P_i(S) \div P_j(S)$.

324 For birth-death processes that generate contemporaneous samples, the probability of the samples
325 for a given set of parameters is straightforward to compute (Kendall 1948; Höhna 2015), and we can
326 relatively easily compute the marginal likelihood. Indeed, when we calculate the marginal likelihood
327 of the samples, it corresponds exactly to the discrepancy we observe between SS and RJ MCMC esti-
328 mates of the Bayes factor (Fig. 1B, main text, middle, dashed lines). For birth-death processes generat-
329 ing non-contemporaneous samples, the marginal probability of the samples is not possible to compute
330 analytically, and therefore we cannot make precise numerical predictions about the discrepancy.

331 **S§5 Factorizing Bayes’ Theorem**

332 The problem with equation (S.6) is that $P(S, \Psi | \theta_\Psi)$ combines likelihood and prior quantities. This
 333 could be resolved by factoring this quantity as:

$$P(S, \Psi | \theta_\Psi) = P(\Psi | S, \theta_\Psi)P(S | \theta_\Psi),$$

334 in which case each term on the right is unambiguously likelihood (the marginal probability of the
 335 samples) or not (the conditional probability of the tree, given the samples). This can be directly sub-
 336 stituted into the posterior:

$$P(\Psi, \theta_x, \theta_\Psi | S, X) = \frac{P(X | S, \Psi, \theta_x)P(\Psi | S, \theta_\Psi)P(S | \theta_\Psi)P(\theta_x)P(\theta_\Psi)}{P(S, X)}, \quad (\text{S.13})$$

337 which is compatible with standard numerical methods for computing the marginal likelihood that
 338 rely on raising the likelihood function to a power (“power-posterior” methods, for example, path-
 339 sampling and stepping-stone-sampling algorithms; [Lartillot and Philippe 2006](#); [Xie et al. 2011](#)). How-
 340 ever, the marginal probability of the samples and the conditional probability of the tree and the sam-
 341 ples are not generally easy to compute; in particular, analytical solutions are only currently available
 342 for simple models of contemporaneous samples, and may be impossible for more complex models.

343 We implemented this solution for Yule and birth-death models producing contemporaneous sam-
 344 ples. For these models, $P(S | \theta_\Psi)$ is the probability of realizing n samples of a given age, for which
 345 there are available analytical solutions (*e.g.*, equation [8] from [Höhna 2015](#)). Likewise, $P(\Psi | S, \theta_\Psi)$ —
 346 the probability of the tree conditional on n samples—also has an available analytical solution (*e.g.*,
 347 equation [3] from [Yang and Rannala 1997](#)). We re-analyzed our simulated data from Section S§1.2
 348 using this formulation, demonstrating that it provides correct marginal-likelihood estimates (Fig. S7).

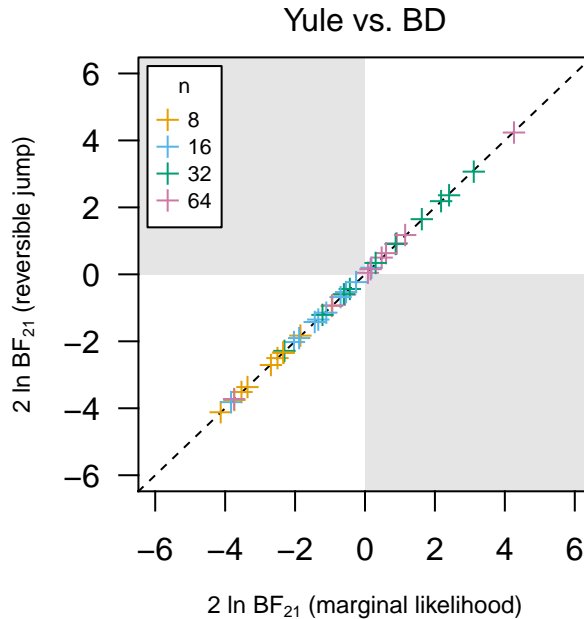


Figure S7: Bayes factor discrepancies are resolved by refactoring Bayes’ theorem. We compared the fit of two birth-death processes—the Yule model (with no extinction rate parameter) and the standard birth-death (BD) model—to datasets simulated under the BD model (as described in Section S§1.2). We corrected the likelihood function according to equation (S.13). As expected, there is no discrepancy between the BFs calculated using marginal likelihoods and reversible-jump MCMC.

349 S§6 Posterior-Predictive Simulation with Samples as Data

350 While Bayes factors are useful for comparing the *relative* fit of competing models, they provide no
351 guarantee that the best model adequately describes the process that gave rise to the observed data.
352 Posterior-predictive simulation (PPS; Gelman et al. 1996) is a Bayesian tool that fills this gap by as-
353 sessing model adequacy—whether our inference model provides an adequate description of the true
354 process that produced our observed dataset—and is therefore useful for assessing *absolute* model fit.
355 Generally, the procedure works by drawing parameters of the model from their joint posterior dis-
356 tribution (*e.g.*, as produced by an MCMC analysis), simulating new datasets under these parameters,
357 and checking whether the simulated data resembles the observed dataset: are the values of a partic-
358 ular summary statistic computed from the simulated datasets reasonably close to the value of that
359 statistic computed from the empirical data?

360 In phylogenetics, the PPS has been largely limited to morphological or molecular character
361 datasets (*e.g.*, Brown 2014; Höhna et al. 2018; Slater and Pennell 2014; May et al. 2021). This lim-
362 ited application of PPS is understandable, given that the character datasets are the only component of
363 the study that is considered to be data under the standard phylogenetic model. However, for studies
364 that rely on the tree model, such as diversification-rate analyses or divergence-time estimation, a more
365 natural summary statistic would be one that relates to characteristics of the sample, rather than the
366 morphological or molecular data. For example, if we wanted to assess the adequacy of a diversifica-
367 tion model, we might use the number of samples at a particular time as a test statistic. The availability
368 of this application of sample-based test statistics is one of the primary benefits—to theoreticians and
369 empiricists alike—of recognizing the samples themselves as data that inform the tree model.

370 To demonstrate the utility of posterior-predictive distributions for samples under birth-death
371 models, we applied this technique to the Marattiales analyses described above (S§2). We simulated
372 datasets by simulating trees under the sampled fossilized birth-death model parameters, and keeping
373 track of the number of fossils recorded in each geological epoch. The posterior-predictive distributions
374 of the number of samples shows that the model with constant fossilization rates does a poor job of pre-
375 dicting the observed number of fossils in the Mississippian, Pennsylvanian, and Cisuralian (Fig. S8).
376 By contrast, the model with variable fossilization rates does a much better job at predicting the number
377 of fossils in these (and subsequent) intervals (Fig. S8, right). This result is concordant with our relative
378 measures of model fit (using Bayes factors), which we report in the main text, and demonstrates that
379 the variable-rate model is not only better-fitting than its constant-rate counterpart, but moreover that
380 it is an adequate representation of the process that generated our data. (We present these results as a
381 proof-of-concept rather than as a method: developing appropriate posterior-predictive methods is a
382 significant task that requires validation and evaluation of statistic properties.)

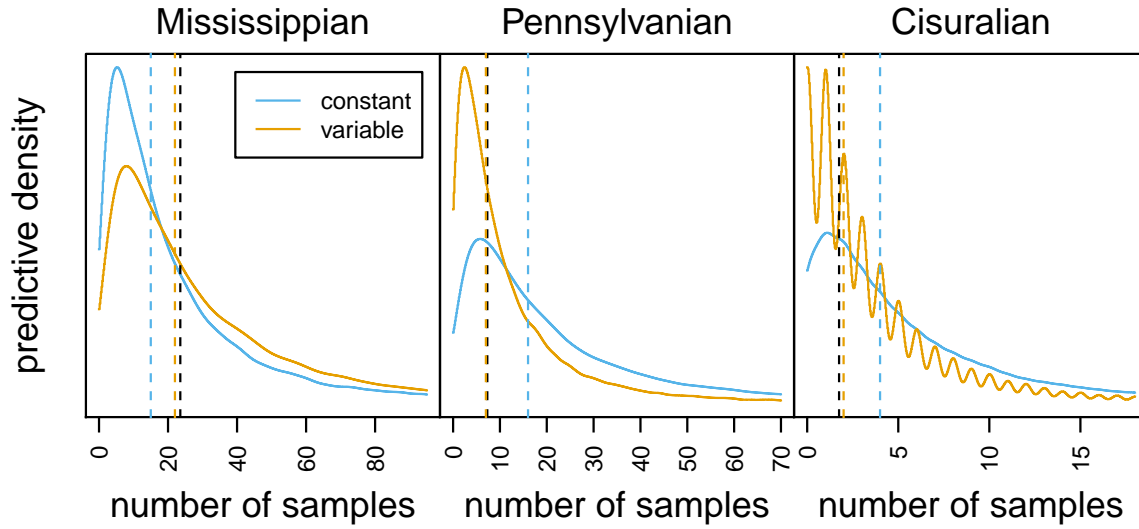


Figure S8: Posterior-predictive simulation for the Marattiales dataset under two models. We simulated fossil and extant marattialean datasets under a model with constant fossilization rates (blue) and variable fossilization rates (orange). Each density represents the posterior-predictive distribution for the number of samples in the given epoch. The black vertical line represent the observed number of samples in the epoch, and the colored vertical lines represent the median number of samples under the corresponding model.

383 S§7 Random Variables and Priors in Phylogenetic Inference

384 The problem that we detail in this manuscript—that the standard phylogenetic model mislabels sam-
385 ples as belonging to the prior rather than to the likelihood, which prevents accurate marginal likeli-
386 hood calculation for tree models—is a specific example of a more general inconsistency in likelihood-
387 based phylogenetics of distinguishing random variables (which have probabilities) from paramete-
388 rs (which have likelihoods). For example, there is a history of mistreating data as parameters in
389 maximum-likelihood inference of ancestral states. Ancestral states are an outcome of the model: they
390 are random variables, just as the character states at the tips are, an equivalence that is apparent when
391 one recognizes that today’s tip data are tomorrow’s ancestral states. It is true, of course, that ancestral
392 states are not observed, but that does not make them any less data-like; they are perhaps best con-
393 ceived of as missing data, just as there can be missing data at the tips. There is nonetheless a strong
394 history of treating ancestral states as parameters and, for example, comparing among different an-
395 cestral states with likelihood-ratio tests (*e.g.*, [Pagel 1999](#)). We contend that this approach is incorrect:
396 ancestral states are not parameters and thus do not have likelihoods; rather, ancestral states are ran-
397 dom variables to which we can assign different probabilities, given the tip data and model of character
398 evolution ([Yang et al. 1995](#); [Yang 2014](#)).

399 Similarly, there has been some historical disagreement about whether the phylogeny itself is a
400 parameter or a random variable. The current dominant perspective, which derives from Felsenstein
401 ([Felsenstein 1973a,b](#)), is that the tree is a parameter. In a maximum-likelihood framework, this per-
402 spective implies that the tree has a maximum-likelihood estimate; in a Bayesian framework, this per-
403 spective suggests that the tree should have a prior distribution. However, a minority perspective is
404 that the tree should be viewed as a random variable just as the ancestral states are viewed as a random
405 variable ([Edwards 1970](#); [Rannala and Yang 1996](#)), and therefore should be associated with a proba-
406 bility distribution even in a maximum-likelihood framework. This latter perspective is even more
407 germane today, considering the significant development of character-state-dependent diversification
408 models (*i.e.*, the binary-state-specific speciation-and-extinction model [BiSSE; [Maddison et al. 2007](#)]
409 and its derivatives). These models assume that rates of speciation and extinction are a function of an
410 evolving character; since the evolution of the character is a random process, the resulting tree must be
411 a random variable.

412 Our arguments are in alignment with the latter perspective: the tree should be viewed as the out-
413 come of a random process, even in a maximum likelihood framework. However, while both [Edwards](#)
414 ([1970](#)) and [Rannala and Yang \(1996\)](#) used birth-death models for the tree, they condition the model on
415 achieving exactly the observed number of extant species. This is equivalent to performing the second
416 step of a sequential Bayesian analysis, which we describe above. Our view is therefore an extension
417 of Edwards’, to include the samples themselves as part of the outcome.

418 References

- 419 Brown, J. M. (2014). Predictive approaches to assessing the fit of evolutionary models. *Systematic*
420 *Biology*, 63(3):289–292.
- 421 Drummond, A. J., Ho, S. Y. W., Phillips, M. J., and Rambaut, A. (2006). Relaxed phylogenetics and
422 dating with confidence. *PLoS Biology*, 4(5):699–710.
- 423 Drummond, A. J. and Stadler, T. (2016). Bayesian phylogenetic estimation of fossil ages. *Philosophical*
424 *Transactions of the Royal Society B: Biological Sciences*, 371(1699):20150129.
- 425 Edwards, A. W. (1970). Estimation of the branch points of a branching diffusion process. *Journal of the*
426 *Royal Statistical Society: Series B (Methodological)*, 32(2):155–164.
- 427 Felsenstein, J. (1973a). Maximum likelihood and minimum-steps methods for estimating evolutionary
428 trees from data on discrete characters. *Systematic Biology*, 22(3):240–249.
- 429 Felsenstein, J. (1973b). Maximum-likelihood estimation of evolutionary trees from continuous char-
430 acters. *American Journal of Human Genetics*, 25(5):471.
- 431 Foote, M. (2000). Origination and extinction components of taxonomic diversity: general problems.
432 *Paleobiology*, 26(sp4):74–102.
- 433 Foote, M. (2001). Inferring temporal patterns of preservation, origination, and extinction from taxo-
434 nomic survivorship analysis. *Paleobiology*, 27(4):602–630.
- 435 Friel, N., Hurn, M., and Wyse, J. (2014). Improving power posterior estimation of statistical evidence.
436 *Statistics and Computing*, 24(5):709–723.
- 437 Gavryushkina, A., Welch, D., Stadler, T., and Drummond, A. J. (2014). Bayesian inference of sampled
438 ancestor trees for epidemiology and fossil calibration. *PLoS Computational Biology*, 10(12):e1003919.
- 439 Gelman, A., Meng, X.-L., and Stern, H. (1996). Posterior predictive assessment of model fitness via
440 realized discrepancies. *Statistica Sinica*, pages 733–760.
- 441 Green, P. J. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model
442 determination. *Biometrika*, 82(4):711–732.
- 443 Heath, T. A., Huelsenbeck, J. P., and Stadler, T. (2014). The fossilized birth-death process for coher-
444 ent calibration of divergence-time estimates. *Proceedings of the National Academy of Sciences, USA*,
445 111(29):E2957–E2966.
- 446 Höhna, S. (2015). The time-dependent reconstructed evolutionary process with a key-role for mass-
447 extinction events. *Journal of Theoretical Biology*, 380:321–331.
- 448 Höhna, S., Coghill, L. M., Mount, G. G., Thomson, R. C., and Brown, J. M. (2018). P3: Phylogenetic
449 posterior prediction in RevBayes. *Molecular Biology and Evolution*, 35(4):1028–1034.
- 450 Höhna, S., Landis, M. J., Heath, T. A., Boussau, B., Lartillot, N., Moore, B. R., Huelsenbeck, J. P., and
451 Ronquist, F. (2016). RevBayes: Bayesian phylogenetic inference using graphical models and an
452 interactive model-specification language. *Systematic Biology*, 65(4):726–736.
- 453 Jukes, T. H. and Cantor, C. R. (1969). Evolution of protein molecules. *Mammalian Protein Metabolism*,
454 3:21–132.

- 455 Kendall, D. G. (1948). On the generalized “birth-and-death” process. *The Annals of Mathematical*
456 *Statistics*, 19(1):1–15.
- 457 Kimura, M. (1980). A simple method for estimating evolutionary rates of base substitutions through
458 comparative studies of nucleotide sequences. *Journal of Molecular Evolution*, 16(2):111–120.
- 459 Lartillot, N. and Philippe, H. (2006). Computing Bayes factors using thermodynamic integration.
460 *Systematic Biology*, 55(2):195–207.
- 461 Lewis, P. O. (2001). A likelihood approach to estimating phylogeny from discrete morphological
462 character data. *Systematic Biology*, 50(6):913–925.
- 463 Lindley, D. V. (1972). *Bayesian statistics: A review*. SIAM.
- 464 Maddison, W. P., Midford, P. E., and Otto, S. P. (2007). Estimating a binary character’s effect on
465 speciation and extinction. *Systematic Biology*, 56(5):701–710.
- 466 May, M. R., Contreras, D. L., Sundue, M. A., Nagalingum, N. S., Looy, C. V., and Rothfels, C. J. (2021).
467 Inferring the total-evidence timescale of marattialean fern evolution in the face of model sensitivity.
468 *Systematic Biology*.
- 469 Murtagh, F. (1984). Counting dendrograms: a survey. *Discrete Applied Mathematics*, 7(2):191–199.
- 470 Pagel, M. (1999). The maximum likelihood approach to reconstructing ancestral character states of
471 discrete characters on phylogenies. *Systematic Biology*, 48(3):612–622.
- 472 Rannala, B. and Yang, Z. (1996). Probability distribution of molecular evolutionary trees: a new
473 method of phylogenetic inference. *Journal of Molecular Evolution*, 43(3):304–311.
- 474 Raup, D. M. (1975). Taxonomic survivorship curves and Van Valen’s Law. *Paleobiology*, 1(1):82–96.
- 475 Raup, D. M. (1985). Mathematical models of cladogenesis. *Paleobiology*, 11(1):42–52.
- 476 Rodrigue, N. and Aris-Brosou, S. (2011). Fast Bayesian choice of phylogenetic models: Prospecting
477 data augmentation–based thermodynamic integration. *Systematic Biology*, 60(6):881–887.
- 478 Silvestro, D., Salamin, N., Antonelli, A., and Meyer, X. (2019). Improved estimation of macroevolu-
479 tionary rates from fossil data using a Bayesian framework. *Paleobiology*, 45(4):546–570.
- 480 Silvestro, D., Schnitzler, J., Liow, L. H., Antonelli, A., and Salamin, N. (2014). Bayesian estimation of
481 speciation and extinction from incomplete fossil occurrence data. *Systematic Biology*, 63(3):349–367.
- 482 Slater, G. J. and Pennell, M. W. (2014). Robust regression and posterior predictive simulation increase
483 power to detect early bursts of trait evolution. *Systematic Biology*, 63(3):293–308.
- 484 Tanner, M. A. and Wong, W. H. (1987). The calculation of posterior distributions by data augmenta-
485 tion. *Journal of the American statistical Association*, 82(398):528–540.
- 486 Xie, W., Lewis, P. O., Fan, Y., Kuo, L., and Chen, M.-H. (2011). Improving marginal likelihood estima-
487 tion for Bayesian phylogenetic model selection. *Systematic Biology*, 60(2):150–160.
- 488 Yang, Z. (1994). Maximum likelihood phylogenetic estimation from DNA sequences with variable
489 rates over sites: approximate methods. *Journal of Molecular evolution*, 39(3):306–314.
- 490 Yang, Z. (2014). *Molecular Evolution: A Statistical Approach*. Oxford University Press.

- 491 Yang, Z., Kumar, S., and Nei, M. (1995). A new method of inference of ancestral nucleotide and amino
492 acid sequences. *Genetics*, 141(4):1641–1650.
- 493 Yang, Z. and Rannala, B. (1997). Bayesian phylogenetic inference using DNA sequences: a Markov
494 Chain Monte Carlo method. *Molecular Biology and Evolution*, 14(7):717–724.
- 495 Zhang, C., Stadler, T., Klopstein, S., Heath, T. A., and Ronquist, F. (2016). Total-evidence dating under
496 the fossilized birth–death process. *Systematic Biology*, 65(2):228–249.