

1 **A biologist's guide to Bayesian phylogenetic analysis**

2 Fabrícia F. Nascimento^{1*}, Mario dos Reis² and Ziheng Yang^{3*}

3 1. Department of Zoology, University of Oxford, OX1 3PS, UK.

4 2. School of Biological and Chemical Sciences, Queen Mary University of London, E1 4NS,

5 UK

6 3. Department of Genetics, Evolution and Environment, University College London, WC1E

7 6BT, UK

8 *Correspondent author

9 FFN: nascimentoфф@yahoo.com.br

10 ZY: z.yang@ucl.ac.uk

11

12 **Abstract**

13 Bayesian methods have become very popular in molecular phylogenetics due to the
14 availability of user-friendly software implementing sophisticated models of evolution.
15 However, Bayesian phylogenetic models are complex, and analyses are often carried out
16 using default settings, which may not be appropriate. Here, we summarize the major features
17 of Bayesian phylogenetic inference and discuss Bayesian computation using Markov chain
18 Monte Carlo (MCMC), the diagnosis of an MCMC run, and ways of summarising the
19 MCMC sample. We discuss the specification of the prior, the choice of the substitution
20 model, and partitioning of the data. Finally, we provide a list of common Bayesian
21 phylogenetic software and provide recommendations as to their use.

22 **Introduction**

23 Bayesian phylogenetic methods were introduced in the 1990s^{1,2} and have since
24 revolutionised the way we analyse genomic sequence data³. Examples of such analyses
25 include phylogeographic analysis of virus spread in humans⁴⁻⁷, inference of phylogeographic
26 history and migration between species⁸⁻¹⁰, analysis of species diversification rates^{11,12},
27 divergence time estimation¹³⁻¹⁵, and inference of phylogenetic relationships among species or
28 populations^{13,16-20}. The popularity of Bayesian methods appears to be due to two factors: (1)
29 the development of powerful models of data analysis; and (2) the availability of user-friendly
30 computer programs implementing the models (Table 1).

31 Models implemented in Bayesian software programs are becoming increasingly
32 complicated, and the priors and model assumptions made in those programs are not always
33 clear to the user. Analyses are often conducted using default priors, which may not be
34 appropriate and may lead to biased or incorrect results. Likewise, over-simplified likelihood
35 models may produce biased results, while over-complicated models may lead to loss of
36 power as well as inefficient computation.

37 The workhorse underlying all modern Bayesian phylogenetic programs is the Markov
38 chain Monte Carlo (MCMC) or Metropolis-Hastings algorithm^{21,22}. However, MCMC is
39 both art and science, and a basic understanding of its workings is essential for the correct use
40 of those programs. In this review, we explain the basic concepts of Bayesian statistics and
41 discuss the major features of MCMC algorithms, such as the prior and the likelihood, MCMC
42 proposals, diagnosis of MCMC convergence and mixing, and summary of the posterior
43 sample. Our intended reader is the empirical biologist who needs to use Bayesian
44 phylogenetic programs to analyse their data. We lay out and answer a set of questions
45 important for setting up a Bayesian analysis. We focus on Bayesian estimation of
46 phylogenetic trees. However, the basic concepts discussed here apply to other phylogenetic
47 problems as well, such as divergence time estimation or species tree estimation under the
48 multi-species coalescent model. Extensive reviews of these are available elsewhere²³⁻²⁵.

49 **What is the Bayesian method?**

50 The Bayesian method is a statistical inference methodology. Its main feature is the use
51 of probability distributions to describe the uncertainty of all unknowns including the model
52 parameter(s). Let D be the observed data and θ the unknown parameter. We assign a
53 distribution $f(\theta)$, called the *prior distribution*, based on our knowledge about θ before
54 analysis of the data. After the data are observed, we use Bayes's theorem to calculate the
55 *posterior distribution* of θ given the data:

$$56 \quad f(\theta|D) = \frac{1}{z} f(\theta) f(D|\theta), \quad (1)$$

57 where the probability of the data given the parameter, $f(D|\theta)$, is called the *likelihood*. This
58 summarises the information about θ in the data. The *normalising constant*

59 $z = \int f(\theta) f(D|\theta) d\theta$ ensures that $f(\theta|D)$ integrates to 1 and is a proper statistical distribution.

60 Equation (1) indicates that the posterior is proportional to the prior times the likelihood, or
61 the posterior combines the information in the prior and in the data. An example of the prior,
62 likelihood and posterior for a two-parameter phylogenetic example is given in Figure 1.

63 In the above we assume that the model for generating the data is known. In the so-
64 called trans-model inference, we have several competing models, with each model m having
65 its own parameters θ_m . Then a prior, $f(m, \theta_m) = f(m) f(\theta_m|m)$, is assigned to both the model
66 (m) and its parameters (θ_m), and the posterior of the model and parameter is similarly given
67 by Bayes's theorem: $f(m, \theta_m|D) \propto f(m, \theta_m) f(D|m, \theta_m)$.

68 In phylogenetics, the tree topology and the substitution model together specify the
69 statistical model for the data. Different tree topologies thus correspond to different models,
70 while the branch lengths or divergence times as well as the substitution parameters (such as
71 the transition/transversion rate ratio) are parameters in the model. The data are usually a
72 molecular sequence alignment or an alignment of morphological characters (or a combination
73 of both).

74 An appealing property of Bayesian inference is that it makes direct probabilistic
75 statements about the model or unknown parameter. The posterior probability of a model,
76 $f(m|D)$, is the probability that the model is correct, given the data. The 95% *credibility*
77 *interval* (CI) of a parameter covers the true parameter with probability 0.95, given the data.
78 Such statements are impossible using confidence intervals and p -values in classical statistics,
79 which treat parameters as unknown constants²⁶.

80 **What type of data can I use?**

81 The most common type of data used in phylogenetic analyses is DNA and amino acid
82 sequence alignments. Morphological characters can also be used²⁷. Here, we focus on DNA
83 sequences. The sequences must be aligned before they are used as input data in phylogenetic
84 programs, and alignment accuracy is important in phylogenetic analysis. Much effort has

85 been made to develop models of insertions and deletions²⁸⁻³⁰. For species phylogeny
86 estimation, the sequences must be orthologs, as incorrect use of paralogs may lead to
87 incorrect phylogenies. Several methods are now available to infer paralogy/orthology^{31,32}.

88 **How do I select a substitution model for my data?**

89 A number of models have been developed to describe nucleotide or amino acid
90 substitutions^{26,33,34}. For nucleotide sequences, these range from the simple JC69 (for Jukes
91 and Cantor)³⁵ to the complex GTR (for General Time Reversible)³⁶⁻³⁸, and the unrestricted
92 model (UNREST)³⁷. In JC69 all nucleotide changes occur at the same rate, while in GTR or
93 UNREST substitutions occur at different rates depending on the source and target
94 nucleotides. It is also common to assume a gamma model of variable rates across sites, in
95 particular, in analysis of coding DNA or protein sequences³⁹⁻⁴¹.

96 Programs such as jModelTest⁴², Modelgenerator⁴³ or PartitionFinder⁴⁴ are commonly
97 used to choose a substitution model. Those programs examine the goodness of fit of the
98 model to the data but never consider the robustness of the analysis to model assumptions. For
99 example, it is well known that the transition/transversion bias typically has a greater impact on
100 the fit of the model to data (judged by the improvement in likelihood), but less effect on
101 estimation of the tree topology and branch lengths than rate variation among sites⁴¹.

102 Although there does not seem to be serious harm in mechanical use of those programs, it may
103 be unnecessary to do so in many cases. As a rule of thumb, different substitution models tend
104 to give very similar sequence distance estimates when sequence divergence is less than 10%,
105 so that a simple model can be used even though it may not fit the data. Complex models are
106 necessary in reconstruction of deep phylogenies. Two of the most complex nucleotide
107 substitution models, HKY+ Γ and GTR+ Γ , often produce similar estimates of phylogenetic
108 trees and branch lengths^{37,45}. When in doubt, note that it is more problematic to under-
109 specify than to over-specify the model in Bayesian phylogenetics⁴⁶.

110 For discrete morphological data, the Mk model, an extension of the JC69 model to k
111 morphological character states, can be used²⁷. An extension that allows for unequal rates of
112 substitution is available in MrBayes⁴⁷. A correction for ascertainment bias is applied in
113 calculation of the likelihood function because only variable characters are used²⁷. For
114 continuous characters, diffusion process models (such as the Wiener or the Ornstein-
115 Uhlenbeck process) can be used⁴⁸. Definitions and detailed review of these models are given
116 elsewhere⁴⁹. There has been much interest in the joint analysis of morphological and
117 molecular data to estimate divergence times for extant and fossil species⁵⁰⁻⁵².

118 **What is over- and under-parameterisation?**

119 A model is non-identifiable if different values of parameters make the same predictions
120 about the data, so that such data can never be used to estimate those parameters; in other
121 words, the model is non-identifiable if $f(D|\theta_1) = f(D|\theta_2)$ for certain $\theta_1 \neq \theta_2$ and for all possible
122 data D [53]. A simple phylogenetic example is estimation of the geological time of
123 divergence between two species (t) and the molecular evolutionary rate (r) using data of a
124 pair of aligned sequences. The likelihood depends only on the molecular distance, $d = rt$, and
125 not on t and r separately, and is the same for, say, $t = 1$ and $r = 0.1$, or $t = 0.1$ and $r = 1$, or
126 any other combination of t and r such that $rt = d = 0.1$. In theory, non-identifiability (or over-
127 parameterisation) is not a serious problem for Bayesian analysis, especially if informative
128 priors are assigned on the parameters. In practice, over-parameterisation can cause both
129 inference difficulties (such as loss of power, strong correlations between parameters, large
130 variance in the posterior, and extreme sensitivity to the prior and model assumptions) and
131 computational problems (such as poor mixing of the MCMC). Sometimes, a model is
132 identifiable, but the data contain only weak information about the parameters with the
133 likelihood surface being nearly flat. Then similar symptoms will show up in the data
134 analysis.

135 An example is the popular I+G model of rate variation among sites, which assumes a
136 proportion of sites p_0 in the alignment are invariable with rate 0, while the other sites $(1 - p_0)$
137 evolve according to a discrete gamma distribution⁵⁴. Because the gamma distribution allows
138 for extremely conserved sites with rates close to 0, p_0 and the gamma shape parameter α are
139 strongly correlated⁵⁵. The MCMC algorithm may have to spend a long time exploring a ridge
140 on the posterior surface.

141 A similar case applies to the use of parameter-rich GTR+ Γ model in analysis of highly
142 similar sequences from closely related species as in Bayesian species delimitation or species
143 tree estimation under the multi-species coalescent model^{24,56}. The GTR model has eight
144 parameters that describe the exchangeabilities between nucleotides. If there are only a few
145 variable sites in the alignment, there will be little information about those parameters. Simple
146 models, such as JC69 and K80, may be adequate in such analysis.

147 On the other hand, the use of overly simplistic model or under-parametrisation can
148 cause systematically incorrect phylogenetic trees and seriously biased estimates of branch
149 lengths and substitution parameters, and over-confident assessment of uncertainties such as
150 spuriously high posterior probabilities for trees or clades⁴⁶. For example, ignoring variable
151 substitution rates among sites leads to underestimated branch lengths⁴¹. Systematic errors
152 tend to be greater when sequences are more divergent. In short, the substitution model is a
153 trade-off between bias on one hand and variance and computation expense on the other, and
154 should ideally be chosen by a careful consideration of its role on the analysis rather than
155 mechanistic use of a model selection procedure.

156 **How do I decide to concatenate or partition my data?**

157 The rationale for partitioned analysis is that sites in the same partition have similar
158 evolutionary characteristics while those in different partitions have different
159 characteristics^{40,44,57}. The characteristics here may be substitution rates, base composition,

160 branch lengths, or even the tree topology. The Bayesian program will estimate different
161 parameter values or even different gene tree topologies for the different partitions, thus
162 accounting for their heterogeneity in the evolutionary process.

163 For example, genes with different G+C compositions or evolutionary rates may be
164 analysed as separate partitions in phylogeny reconstruction. Vertebrate mitochondrial genes
165 coded on the same strand of the genome have similar G+C content and may be concatenated
166 and analysed as a single partition, although the three codon positions may be treated as
167 different partitions to account for their large differences in rate and in base compositions⁵⁸.
168 Non-coding mitochondrial genes (rRNAs and tRNAs) may be analysed as another partition.
169 Likewise, mitochondrial and nuclear sequences should also be analysed as different
170 partitions⁵⁹. For nuclear sequences, exons and introns should be analysed as different
171 partitions, and the three codon positions should be placed in their own partitions. Some
172 partitioning software may suggest the use of different substitution models for partitions⁴⁴
173 (e.g., HKY for one partition and GTR+G for another). This is unnecessary because with the
174 same model for all partitions, different parameter values will accommodate the heterogeneity
175 among partitions.

176 An important issue is whether partitions should share the same tree topology. In
177 traditional phylogenetic inference, topology is assumed to be the same across partitions.
178 However, a number of biological processes, such as gene duplication, horizontal gene
179 transfer, and incomplete lineage sorting can cause different genes to have different trees^{60,61}.
180 Recently, a number of methods for species tree estimation have been developed under the
181 multi-species coalescent (MSC) model^{24,62,63}, which account for the process of incomplete
182 lineage sorting (the so-called deep coalescent, due to polymorphism in ancestral species,
183 where coalescence may occur in ancient ancestors leading to gene trees that differ from the
184 species tree). Under the MSC different genomic regions (or exons) are placed into different

185 partitions and allowed to have their own gene-trees, which are embedded into the species
186 tree. The mitochondrial genome does not recombine and mitochondrial genes should be
187 treated as one partition within the MSC. In some viruses, such as influenza, different genome
188 segments can re-assort (i.e. be horizontally transferred) among related strains⁶⁴, and thus
189 different segments can have different topologies and should be treated as different partitions.

190 **How do I choose the prior for my Bayesian analysis?**

191 In theory the prior should summarize the biologist's best knowledge about the model or
192 parameters before the data are analysed^{26,65}. In practice, specification of the prior is often a
193 thorny issue, especially if there are multiple parameters with complex correlations or if little
194 is known about the parameters. While we are supposed to specify a joint prior distribution
195 for all parameters, the common practice is to ignore the correlation, and assign independent
196 priors for the parameters. When there are many parameters of the same kind, such
197 independent and identically distributed (i.i.d.) prior can sometimes cause problems because
198 they may make a strong statement about the mean or sum of those parameters. For example,
199 it is common to assign independent exponential or uniform priors for branch lengths in the
200 unrooted tree, but this i.i.d. prior can cause very long trees in analysis of highly similar
201 sequence data^{66,67}. In relaxed-clock dating analysis, the i.i.d prior for substitution rates
202 among different partitions makes a strong statement about the average rate over loci, leading
203 to biased but over-confident divergence time estimates⁶⁸, in particular as the number of
204 partitions increases. Such i.i.d. priors should be avoided.

205 Default priors in many Bayesian software packages may not be appropriate for the data
206 being analysed and should be used with caution. Specification of the prior is the biologist's
207 responsibility even though it may not be an easy task. Robustness analysis should also be an
208 important component of any Bayesian analysis. By evaluating the posteriors generated under
209 different priors, the biologist can evaluate whether the posterior is robust to the prior.

210 In Bayesian estimation of phylogenetic trees without the assumption of a molecular
211 clock, it is common to assign a uniform prior on the unrooted tree topologies. When
212 phylogenetic analysis is conducted on rooted trees under the clock or relaxed clock models⁶⁹,
213 rooted trees are commonly assigned a prior using a model of cladogenesis such as the Yule
214 process and the birth-death-sampling process⁷⁰. Note that all those models favour balanced
215 trees, and the impact of the prior on the posterior probabilities of the rooted trees can be
216 substantial if the tree is large. For coalescent-based species tree estimation, the MSC model
217 specifies a probability distribution for the rooted gene trees (topologies and node ages)⁷¹.
218 This is part of the model rather than a prior on gene trees to be specified. In molecular clock
219 dating analysis, fossils may be used to specify minimum and maximum bounds on clade age,
220 which are used to construct a so-called calibration density to calibrate the age of the clade, it
221 is also advisable to include a prior on the age of the root of the tree. For an overview on
222 calibration densities for use in divergence dating, see⁷². It is also necessary to specify a prior
223 on the evolutionary rates for the different loci or partitions. A gamma-Dirichlet prior can be
224 used instead of the i.i.d. prior mentioned above⁶⁸. In relaxed-clock models, the rates not only
225 vary among partitions, but also drift along branches on the tree. Current Bayesian
226 implementations assume that rates drift independently among partitions so that different
227 partitions are independent realizations of the rate-drift process^{73,74}. A discussion of the
228 different rate-drift models is given in⁶⁸.

229 **What is Markov chain Monte Carlo (MCMC)?**

230 Once the biologist has decided on the data, model and prior, the next step is to obtain a
231 sample from the posterior. This is done by using MCMC, a simulation technique for
232 sampling from a probability distribution that is known up to a normalising constant^{21,22}. Note
233 that all terms on the right hand side of equation (1) are straightforward to calculate except the
234 normalizing constant z , which involves multidimensional integrals and may be too expensive

235 to compute. Thus, MCMC is particularly suitable for Bayesian computation. Instead of
 236 calculating the posterior distribution $f(\theta|D)$, the algorithm generates a sample from the
 237 posterior, which can be used to estimate the mean, the standard deviation of the posterior, or
 238 even the whole posterior distribution.

239 Here we illustrate the major features of MCMC by applying it to the problem of
 240 estimating the sequence distance d and the transition/transversion rate ratio κ under the K80
 241 model⁷⁵ using a pair of DNA sequences. The data (D) are an alignment of the human and
 242 orangutan mitochondrial 12S rRNA genes, summarized as $n_S = 84$ transitional differences
 243 and $n_V = 6$ transversional differences at $n = 948$ sites^{26, p.7}. We assign independent gamma
 244 priors, $d \sim G(2, 20)$ and $\kappa \sim G(2, 0.1)$, with densities (Fig. 1a):

$$245 \quad \begin{aligned} f(d) &= \frac{\beta^\alpha}{\Gamma(\alpha)} \times d^{\alpha-1} e^{-\beta d}, \quad \text{with } \alpha = 2, \beta = 20, \\ f(\kappa) &= \frac{\beta^\alpha}{\Gamma(\alpha)} \times \kappa^{\alpha-1} e^{-\beta \kappa}, \quad \text{with } \alpha = 2, \beta = 0.1. \end{aligned} \quad (2)$$

246 The likelihood (Fig. 1b) is given by the K80 model^{26,75} as

$$247 \quad f(D|d, \kappa) = \left(\frac{p_0}{4}\right)^{n-n_S-n_V} \left(\frac{p_1}{4}\right)^{n_S} \left(\frac{p_2}{4}\right)^{n_V}, \quad (3)$$

248 where

$$249 \quad \begin{aligned} p_0(t) &= \frac{1}{4} + \frac{1}{4} e^{-4\beta t} + \frac{1}{2} e^{-2(\alpha+\beta)t} = \frac{1}{4} + \frac{1}{4} e^{-4d/(\kappa+2)} + \frac{1}{2} e^{-2d(\kappa+1)/(\kappa+2)}, \\ p_1(t) &= \frac{1}{4} + \frac{1}{4} e^{-4\beta t} - \frac{1}{2} e^{-2(\alpha+\beta)t} = \frac{1}{4} + \frac{1}{4} e^{-4d/(\kappa+2)} - \frac{1}{2} e^{-2d(\kappa+1)/(\kappa+2)}, \\ p_2(t) &= \frac{1}{4} - \frac{1}{4} e^{-4\beta t} = \frac{1}{4} - \frac{1}{4} e^{-4d/(\kappa+2)}. \end{aligned} \quad (4)$$

250 Thus, the unnormalized posterior (Fig. 1c) is

$$251 \quad f(d, \kappa | D) \propto f(d) f(\kappa) f(D|d, \kappa). \quad (5)$$

252 We give a sketch of an MCMC algorithm in Box 1, and then discuss its main features.

253 We use two sliding windows (uniform distributions centred around the current parameter
 254 value) to update parameters d and κ . The sliding window (even with reflection) is a
 255 symmetrical proposal, in the sense that the probability density of proposing d^* from d is

256 equal to that of proposing d from d^* . If the proposal is asymmetrical, a correction term,
257 called the Hastings ratio²² needs to be applied.

258 Note that the parameter values (d and κ) visited in the next iteration depend on the
259 current values but not values visited in the past. The algorithm has no memory. This
260 memoryless property is called the *Markovian* property. As a result, the sequence of visited
261 parameter values form a Markov chain, and the algorithm is called Markov chain Monte
262 Carlo. An important feature of the algorithm is that it requires the calculation of the ratio of
263 posterior densities, but not the posterior density itself. The normalizing constant z of
264 equation (1) cancels in the calculation of the acceptance ratio α in steps 2a & 2b, and
265 algorithm thus avoids its calculation. It is easy to see that the algorithm visits parameter
266 values with high posterior more often than those with low posterior. Indeed, it visits the
267 parameter values exactly in proportion to their posterior. One runs the algorithm over many
268 iterations, and then uses the visited values of d and κ to construct a histogram to estimate the
269 posterior distribution or to calculate the mean and standard deviation of the posterior (Fig. 2).

270 The window size (or step length) in the sliding window proposal (w_d and w_κ) can affect
271 the mixing efficiency of the chain (Box 2). If the window is too large, most of the proposals
272 will fall in the tails of the posterior and be rejected. The chain then stays at the current value
273 and does not move (Fig 2a'). If the window is too small, the chain takes tiny baby steps,
274 almost all of which are accepted but the chain is ineffective in exploring the posterior surface
275 (Fig 2b'). Thus, both small steps (with high acceptance proportion) and large steps (with
276 very low acceptance proportion) lead to inefficient algorithms. The step lengths should be
277 adjusted to achieve a near optimal acceptance proportion, at about 30-40%. Fine-tuning a
278 phylogenetic MCMC chain to be efficient is important because MCMC runs may take weeks
279 or months. It is easy to monitor the acceptance proportion and use it to adjust the step length

280 automatically⁷⁶. Most current MCMC phylogenetic programs have automatic fine-tuning
281 algorithms and this is in most cases not a concern for the user.

282 In *trans-model MCMC* algorithms, both the model index m and the model parameters
283 θ_m change over the chain. The algorithm will involve both within-model proposals, which
284 change parameters of the current model, and trans-model proposals, which move from the
285 current model to another new model⁷⁷. In the long run, the frequency at which the MCMC
286 visits each model is an estimate of the posterior probability of that model. There are a
287 number of differences between within-model and trans-model algorithms²⁶, and here we note
288 a few concerning mixing efficiency and acceptance proportion. First, for a within-model
289 move (such as a sliding window changing the sequence distance or branch length), we can
290 make the window size small enough so that the acceptance proportion is arbitrarily close to
291 100%. However, in trans-model moves, the acceptance proportion is constrained by the
292 posterior model probabilities. If the maximum *a posteriori* (MAP) model (the model with the
293 highest posterior probability) has the posterior P_1 , then the acceptance proportion cannot
294 exceed $2(1 - P_1)$ [26]. Thus, if the MAP tree has posterior 99%, the highest acceptance
295 proportion for cross-tree moves is 2%. Second, while an acceptance proportion of near 0
296 indicates a poor proposal (e.g., the window size is too large) for a within-model move, this
297 may and may not indicate a mixing problem in cross-model moves because it may be caused
298 by the MAP model having posterior near 100%. Third, for a within-model move, the optimal
299 acceptance proportion is intermediate at 30-40%, but for a trans-model move, a mobile chain
300 is in general more efficient than a lazy chain, so that we should strive to achieve as high an
301 acceptance proportion as possible.

302 All those comments apply to Bayesian phylogenetic MCMC algorithms, which include
303 both within-tree moves that change the branch lengths and substitution parameters without
304 changing the tree topology and cross-tree moves that change the tree topology. The cross-

305 tree moves are typically constructed using tree-perturbation (branch-swapping) algorithms
306 such as nearest-neighbour interchange (NNI), subtree pruning and re-grafting (SPR) and tree
307 bisection and reconnection (TBR)^{26,78}. About a dozen MCMC phylogenetic programs are
308 now available (Table 1).

309 **What are convergence, burn-in and mixing of the MCMC?**

310 An MCMC algorithm may suffer from two problems: slow convergence and poor
311 mixing. In the long run, the Markov chain should be spending most of the time visiting high-
312 probability regions of the posterior. The *convergence rate* is the rate at which a chain starting
313 from any initial position (which may be in the tails of the posterior) moves to the high-
314 posterior region of the parameter space⁷⁹. Parameter values sampled before reaching this
315 stationary phase are usually discarded as the *burn-in*. Thus, if convergence is slow, a long
316 burn-in will be necessary. Convergence rate is affected by the proposals used and by the
317 shape of the posterior in the tails⁶⁷. If the posterior is nearly flat in the tail, it will be difficult
318 for the chain to get out of the tail and move to the high-posterior region.

319 *Mixing efficiency* refers to how efficiently the chain traverses the posterior after it has
320 reached the stationary distribution. If the chain is more efficient, the estimate based on the
321 MCMC sample will have a smaller variance, and the results will show less variation among
322 independent runs (Box 2) and a relatively short chain will provide acceptable estimate. The
323 proposal (such as the uniform sliding window *vs.* the normal-distribution sliding window) as
324 well as the step length for the same proposal (such as the width of the sliding window) can
325 have a great effect on mixing efficiency⁷⁶.

326 Both convergence and mixing problems can be diagnosed by using a trace plot, in
327 which we plot the log likelihood or sampled parameter values against the MCMC iteration,
328 for example, using R⁸⁰ or Tracer⁸¹. It is also very important to run the same algorithm
329 multiple times to check consistency between runs. With fast convergence, different chains

330 that started from very different positions become indistinguishable very quickly. Efficient
331 mixing is indicated by different runs generated nearly identical means, standard deviations,
332 and histograms. If the runs are healthy, samples from different runs can be combined to
333 produce posterior summaries.

334 The trace plots of Figures 2a and 2b are from an efficient chain with good mixing,
335 while those of Figures 2a' and 2b' have poor mixing and low efficiency. The histograms from
336 the efficient algorithm match each other much better than those from the inefficient algorithm
337 (Fig. 2c and 2c'). In theory, the consistency among multiple runs could be because all runs
338 got stuck in a region of the parameter space, giving the false impression that convergence
339 was reached. This may happen when there are multiple peaks in the posterior. Thus, it is
340 important to initiate the runs from widely dispersed starting points.

341 **How many iterations should I run my chain for? How many samples should I take?**

342 Ideally one would like to run the MCMC long enough to obtain a reliable estimation of
343 the posterior distribution, but not overly too long as to waste computational resources.
344 However, currently reliable automatic stopping rules do not exist. As a result, the user has to
345 specify the number of iterations, and then decide whether the chain is long enough or
346 additional iterations are necessary using certain diagnosis tools. MCMC algorithms tend to
347 generate huge output files. To save disk space, one takes a sample only for every certain
348 number of iterations. For example, running an MCMC chain for 10^7 iterations and using a
349 sample frequency of 10^3 iterations will produce 10^4 samples.

350 Note that in some programs (such as MCMCtree and BPP), each MCMC iteration
351 consists of a fixed sequence of MCMC proposals, while in some others (such as MrBayes
352 and BEAST), it consists of one proposal, chosen at random from a collection of proposals.
353 Thus, if there are 1,000 parameters in the model and if each proposal changes one parameter,
354 each MCMC iteration in the former programs is worth about 1,000 iterations in the latter

355 programs. Thus, MCMC iterations from different programs are not comparable. The
356 biologist should instead aim to accumulate a reasonable (as large as practically possible)
357 effective sample size (ESS) for each parameter (Box 2).

358 **Why should an MCMC analysis be run with an “empty alignment”? Is the data**
359 **informative?**

360 It is useful to run the MCMC algorithm sampling from the prior. This is achieved by
361 setting the likelihood to 1 in equation (1). Some programs generate a dummy “empty”
362 alignment that can be used to achieve the same effect. Runs should also be assessed for good
363 convergence and mixing. Running the chain without data is a good way of checking the
364 correctness of the software, because the mean, variance, etc. of the prior are often analytically
365 available and can be checked against the MCMC sample. In molecular clock dating using
366 fossil calibrations, the prior on divergence times incorporates the calibration information and
367 is typically intractable. Running the program without using the sequences allows one to
368 generate the prior used by the program.

369 The sample from the prior can also be compared with the sample from the posterior
370 (which is generated by using the data) to assess how informative the data are, and whether
371 there are serious conflicts between the prior and the data. High similarity between the prior
372 and the posterior suggests that the data contain little information about the parameters.
373 Considerable overlap between the prior and posterior but with the posterior being much more
374 concentrated than the prior means that the data are informative and the prior is reasonable. In
375 contrast, if the prior and posterior do not overlap well, there may be a conflict between the
376 prior and the data, possibly caused by misspecified priors. One can also modify the prior to
377 assess the impact of the prior on the posterior. Note, however, that it is incorrect to specify
378 the prior by trying to match the posterior, since the prior is supposed to reflect our knowledge
379 before the analysis of the data.

380 **Conclusions**

381 Bayesian phylogenetics has undergone explosive growth during the past decade. The
382 implementation of sophisticated models in easy-to-use software programs has made the
383 method extremely appealing to biologists. The method is especially powerful in combining
384 different sources of information in an integrated data analysis. As a result, Bayesian MCMC
385 methods are the most commonly used framework for development of new models of data
386 analysis, especially in the areas of divergence time estimation integrating molecular,
387 morphological and fossil information⁸², species tree estimation using multi-locus genomic
388 sequence data²⁴, and species delimitation incorporating genetic and morphological/ecological
389 information⁸³. The potential of the Bayesian method to deal with these and future questions
390 has never been greater. For further reading on the Bayesian method and Bayesian
391 phylogenetics the reader may consult^{26,84,85}

392

393 A tutorial that helps the user to write a simple R program to conduct phylogenetic MCMC to
394 reproduce the figures of this paper is available at:

395 http://github.com/thednainus/Bayesian_tutorial

396

397 **Author contributions**

398 FFN conceived the idea, FFN, MdR and ZY wrote the paper.

399

400 **Correspondence** should be addressed to FFN or ZY

401

402

403 **Acknowledgements**

404 This work was supported by Biotechnology and Biological Sciences Research Council
405 (UK) grant BB/N000609/1. FFN was supported by a Royal Society and British Academy
406 Newton International Fellowship (UK) grant number NF140338.
407

408 **Table 1.** List of Bayesian programs

Program	Brief description	Refs
BEAST	Implements a vast number of models. Examples are simultaneous estimation of the tree topology and divergence times, phylodynamics, phylogeography, and species tree estimation under the multispecies coalescent model.	86
MrBayes	Implements a large number of models for analysis of nucleotide, amino acid, and morphological data. Estimates species phylogenies and species divergence times.	87
RevBayes	Similar to MrBayes, but with its own programming language to set up complex hierarchical Bayesian models.	88
MCMCTree	Estimates divergence times on a fixed phylogenetic tree.	89
Phycas	Estimates phylogenetic trees based on nucleotide data. This allows for multifurcating trees, helping to reduce spuriously high posterior probabilities for phylogenies.	90,91
PhyloBayes	Reconstructs phylogenetic trees using infinite mixture models to account for among-site and among-lineage heterogeneity in nucleotide or amino acid compositions, which may be important for inferring deep phylogenies.	92
BPP	Implements species tree estimation and species delimitation under the multi-species coalescent model using multi-loci genomic sequence data.	56
Migrate	Estimates population sizes and migration rates under the population-subdivision model based on molecular data.	93
IMa2	Estimates divergence times, population sizes and migration rates under the isolation-with-migration model using multi-loci DNA sequence data and a fixed phylogenetic tree for populations.	94
Structure	Estimates population structure from multi-locus genotype data.	95
BAMM	Estimates clade diversification rates on phylogenies.	96
Tracer	A program for MCMC diagnostics and summaries.	81
AWTY	A package for MCMC diagnostics for Bayesian phylogenetic inference.	97

409

410

411 **Box 1. MCMC algorithm to estimate d and κ under the K80**

412 • 1. (Initialization): Initialize window sizes w_d and w_κ . Choose random starting values
 413 (d, κ) .

414 • 2. (Main loop)

415 ○ 2a (Proposal to change distance d): Propose a new value d^* by sampling from a
 416 uniform sliding window (with reflection) around the current value: $d^* = U(d -$
 417 $w_d/2, d + w_d/2)$, where w_d is the width of the window. If $d^* < 0$, set $d^* = -d^*$
 418 (reflection). If the unnormalised posterior is higher at the new value, accept the
 419 proposal. Otherwise accept with probability equal to the ratio of the posteriors:

$$420 \quad \alpha = \frac{f(d^*, \kappa | D)}{f(d, \kappa | D)} = \frac{f(d^*)f(\kappa)f(D | d^*, \kappa)}{f(d)f(\kappa)f(D | d, \kappa)} \quad (6)$$

421 If the proposal is accepted, set $d = d^*$. If it is rejected, stay where it is ($d = d$).

422 2b (Proposal to change κ): Use a similar sliding window of width w_κ to propose a
 423 new value $\kappa^* = U(\kappa - w_\kappa/2, \kappa + w_\kappa/2)$. If $\kappa^* < 0$, reflect by setting $\kappa^* = -\kappa^*$.

424 Accept the proposal with probability $\min\{1, \alpha\}$, where

$$425 \quad \alpha = \frac{f(d, \kappa^* | D)}{f(d, \kappa | D)} = \frac{f(d)f(\kappa^*)f(D | d, \kappa^*)}{f(d)f(\kappa)f(D | d, \kappa)} \quad (7)$$

426 If the proposal is accepted, set $\kappa = \kappa^*$. Otherwise stay where it is ($\kappa = \kappa$).

427 ○ 2c (Save the state of the chain): Print out d and κ . Go back to 2a and iterate to
 428 obtain as many samples as desired.

429

430 **Box 2. Efficiency of the MCMC and the effective sample size (ESS)**

431 Parameter values sampled during the MCMC are autocorrelated because the current value is
432 either the same as the previous value (if the proposed value is rejected) or a modification of it
433 (e.g., a value sampled from the sliding window around the current value). Stronger
434 autocorrelations mean that the Markov chain is less efficient in traversing the posterior space.
435 More formally, we use the mean of the MCMC sample (\tilde{x}) to estimate the posterior mean of
436 any parameter. This has the variance

437
$$v_{\text{MCMC}} = v_{\text{IND}} \times [1 + 2(\rho_1 + \rho_2 + \dots)], \quad (8)$$

438 where v_{IND} is the variance for an independent sample of the same size from the posterior
439 distribution, and where $\rho_k = \text{corr}(x_t, x_{t+k})$ is the correlation between the values of the
440 parameter in the MCMC sample that are k iterations apart, known as the lag k autocorrelation.
441 Both the independent-sample variance v_{IND} and the MCMC-sample variance v_{MCMC} are
442 typically proportional to $1/n$, with n to be the sample size. The efficiency of an MCMC chain
443 is defined as the variance ratio

444
$$\text{Eff} = \frac{v_{\text{IND}}}{v_{\text{MCMC}}} = \frac{1}{1 + 2(\rho_1 + \rho_2 + \dots)}. \quad (9)$$

445 For example, an $\text{Eff} = 0.25$ means that an MCMC sample of size n is as efficient as an
446 independent sample of size $n/4$, so that we need to generate an MCMC sample four times as
447 large as the independent sample to have the same variance. The effective sample size, ESS, is
448 simply

449
$$\text{ESS} = n \times \text{Eff}.$$

450 As a rule of thumb, one should aim for $\text{ESS} = 1,000$ or $10,000$ [98]. Bayesian phylogenetic
451 algorithms are computationally intensive, so that $\text{ESS} = 200$ is commonly recommended, but
452 this may be too small for calculation of the 95% or 99% credibility intervals. A good
453 strategy may be to conduct multiple runs of the same analysis, and then combine the samples

454 before producing the posterior summary. If $ESS = 200$ for each sample, 10 replicate runs

455 will give a combined sample of $ESS = 2000$.

456

457 **Figure 1 | Prior, likelihood and posterior distribution for a two-parameter phylogenetic**
458 **example.** The data of the 12s RNA mitochondrial genes from human and orang-utan are
459 used to estimate of the evolutionary distance (d) and the transition/transversion ratio (κ)
460 model⁷⁵.

461

462 **Figure 2 | Trace plots and histograms for parameters d and κ sampling the posterior**
463 **distribution of Figure 1c using efficient and inefficient MCMC chains.** Parts **a** and **b**
464 show the trace plots of d and κ for an efficient chain with good mixing. The window sizes
465 are $w_d = 0.12$ and $w_\kappa = 180$, with acceptance proportions $P_{jump} = 30.4\%$ for d and 29.8% for κ ,
466 achieving efficiency $\text{Eff} = 23\%$ for d and 20% for κ . Parts **a'** and **b'** show the trace plots for
467 an inefficient chain with poor mixing, with $w_d = 5$ and $w_\kappa = 1$. In **a'**, the window for d is too
468 wide, and most proposals are rejected ($P_{jump} = 1.5\%$), so that the chain is often stuck at the
469 same value for many iterations, leading to poor mixing with $\text{Eff} = 1.79\%$. In **b'**, the window
470 for κ is too small, so that most of the proposals are accepted (with $P_{jump} = 98.6\%$), but the
471 chain makes small baby steps and is very slow in traversing the posterior parameter space,
472 with $\text{Eff} = 1.28\%$. Parts **c** and **c'** show histograms of κ for two runs of the efficient and
473 inefficient chains (sample size $n = 10,000$). The posterior mean (and standard deviation)
474 calculated using a very long run of the efficient chain is 0.104 (0.0114) for d , and 29.2 (10.0)
475 for κ .

476

477 **References**

- 478
479 1 Rannala, B. & Yang, Z. Probability distribution of molecular evolutionary trees: a
480 new method of phylogenetic inference. *J. Mol. Evol.* **43**, 304-311 (1996).
- 481 2 Mau, B. & Newton, M. A. Phylogenetic inference for binary data on dendograms
482 using Markov chain Monte Carlo. *J. Comp. Graph. Stat.* **6**, 122-131 (1997).
- 483 3 Huelsenbeck, J. P., Ronquist, F., Nielsen, R. & Bollback, J. P. Bayesian inference of
484 phylogeny and its impact on evolutionary biology. *Science* **294**, 2310-2314 (2001).
- 485 4 Wilfert, L. *et al.* Deformed wing virus is a recent global epidemic in honeybees driven
486 by *Varroa* mites. *Science* **351**, 594-597 (2016).
- 487 5 Pybus, O. G. *et al.* Unifying the spatial epidemiology and molecular evolution of
488 emerging epidemics. *Proc. Natl. Acad. Sci. USA* **109**, 15066-15071 (2012).
- 489 6 Faria, N. R. *et al.* HIV epidemiology. The early spread and epidemic ignition of HIV-
490 1 in human populations. *Science* **346**, 56-61 (2014).
- 491 7 Lemey, P., Rambaut, A., Welch, J. J. & Suchard, M. A. Phylogeography takes a
492 relaxed random walk in continuous space and time. *Mol. Biol. Evol.* **27**, 1877-1885
493 (2010).
- 494 8 Bloomquist, E. W., Lemey, P. & Suchard, M. A. Three roads diverged? Routes to
495 phylogeographic inference. *Trends Ecol. Evol.* **25**, 626-632 (2010).
- 496 9 Nascimento, F. F. *et al.* The role of historical barriers in the diversification processes
497 in open vegetation formations during the Miocene/Pliocene using an ancient rodent
498 lineage as a model. *PLoS One* **8**, e61924 (2013).
- 499 10 Werneck, F. P., Leite, R. N., Geurgas, S. R. & Rodrigues, M. T. Biogeographic
500 history and cryptic diversity of saxicolous Tropicoduridae lizards endemic to the
501 semiarid Caatinga. *BMC Evol. Biol.* **15**, 94 (2015).

- 502 11 Merckx, V. S. F. T. *et al.* Evolution of endemism on a young tropical mountain.
503 *Nature* **524**, 347-350 (2015).
- 504 12 Hoorn, C. *et al.* Amazonia through time: Andean uplift, climate change, landscape
505 evolution, and biodiversity. *Science* **330**, 927-931 (2010).
- 506 13 Prum, R. O. *et al.* A comprehensive phylogeny of birds (Aves) using targeted next-
507 generation DNA sequencing. *Nature* **526**, 569-573 (2015).
- 508 14 dos Reis, M. *et al.* Uncertainty in the timing of origin of animals and the limits of
509 precision in molecular timescales. *Curr. Biol.* **25**, 2939-2950 (2015).
- 510 15 Meredith, R. W. *et al.* Impacts of the Cretaceous terrestrial revolution and KPg
511 extinction on mammal diversification. *Science* **334**, 521-524 (2011).
- 512 16 Nascimento, F. F. *et al.* Evolution of endogenous retroviruses in the Suidae: evidence
513 for different viral subpopulations in African and Eurasian host species. *BMC Evol.*
514 *Biol.* **11**, 139 (2011).
- 515 17 Jarvis, E. D. *et al.* Whole-genome analyses resolve early branches in the tree of life of
516 modern birds. *Science* **346**, 1320-1331 (2014).
- 517 18 Misof, B. *et al.* Phylogenomics resolves the timing and pattern of insect evolution.
518 *Science* **346**, 763-767 (2014).
- 519 19 Raymann, K., Brochier-Armanet, C. & Gribaldo, S. The two-domain tree of life is
520 linked to a new root for the Archaea. *Proc. Natl. Acad. Sci. USA* **112**, 6670-6675
521 (2015).
- 522 20 Foley, N. M., Springer, M. S. & Teeling, E. C. Mammal madness: is the mammal tree
523 of life not yet resolved? *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **371** (2016).
- 524 21 Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H. & Teller, E.
525 Equation of state calculations by fast computing machines. *J. Chem. Phys.* **21**, 1087-
526 1092 (1953).

- 527 22 Hastings, W. K. Monte Carlo sampling methods using Markov chains and their
528 applications. *Biometrika* **57**, 97-109 (1970).
- 529 23 Liu, L., Xi, Z., Wu, S., Davis, C. C. & Edwards, S. V. Estimating phylogenetic trees
530 from genome-scale data. *Ann. NY. Acad. Sci.* **1360**, 36-53 (2015).
- 531 24 Xu, B. & Yang, Z. Challenges in species tree estimation under the multispecies
532 coalescent model. *Genetics* **204**, 1353-1368 (2016).
- 533 25 Szöllösi, G. J., Tannier, E., Daubin, V. & Boussau, B. The inference of gene trees
534 with species trees. *Syst. Biol.* **64**, e42-e62 (2015).
- 535 26 Yang, Z. *Molecular Evolution: A statistical Approach*. (Oxford Univ. Press, 2014).
- 536 27 Lewis, P. O. A likelihood approach to estimating phylogeny from discrete
537 morphological character data. *Syst. Biol.* **50**, 913-925 (2001).
- 538 28 Redelings, B. D. & Suchard, M. A. Joint Bayesian estimation of alignment and
539 phylogeny. *Syst. Biol.* **54**, 401-418 (2005).
- 540 29 Löytynoja, A. & Goldman, N. Uniting alignments and trees. *Science* **324**, 1528-1529
541 (2009).
- 542 30 Chatzou, M. *et al.* Multiple sequence alignment modeling: methods and applications.
543 *Brief. Bioinform.* **17**, 1009-1023 (2016).
- 544 31 Altenhoff, A. M. & Dessimoz, C. Inferring orthology and paralogy. *Methods Mol.*
545 *Biol.* **855**, 259-279 (2012).
- 546 32 Altenhoff, A. M. *et al.* The OMA orthology database in 2015: function predictions,
547 better plant support, synteny view and other improvements. *Nucleic Acids Res.* **43**,
548 D240-D249 (2015).
- 549 33 Dimmic, M. in *Statistical Methods in Molecular Evolution* (ed R. Nielsen)
550 (Springer-Verlag, 2005).

- 551 34 Liò, P. & Goldman, N. Models of molecular evolution and phylogeny. *Genome Res.*
552 **8**, 1233-1244 (1998).
- 553 35 Jukes, T. H. & Cantor, C. R. in *Mammalian Protein Metabolism* (ed H.N. Munro)
554 21-132 (Academic Press, 1969).
- 555 36 Tavaré, S. Some probabilistic and statistical problems in the analysis of DNA
556 sequences. *Lect. Math. Life Sci.* **17**, 57-86 (1986).
- 557 37 Yang, Z. Estimating the pattern of nucleotide substitution. *J. Mol. Evol.* **39**, 105-111
558 (1994).
- 559 38 Zharkikh, A. Estimation of evolutionary distances between nucleotide sequences. *J.*
560 *Mol. Evol.* **39**, 315-329 (1994).
- 561 39 Mayrose, I., Graur, D., Ben-Tal, N. & Pupko, T. Comparison of site-specific rate-
562 inference methods for protein sequences: empirical Bayesian methods are superior.
563 *Mol. Biol. Evol.* **21**, 1781-1791 (2004).
- 564 40 Yang, Z., Lauder, I. J. & Lin, H. J. Molecular evolution of the hepatitis B virus
565 genome. *J. Mol. Evol.* **41**, 587-596 (1995).
- 566 41 Yang, Z. Among-site rate variation and its impact on phylogenetic analyses. *Trends*
567 *Ecol. Evol.* **11**, 367-372 (1996).
- 568 42 Darriba, D., Taboada, G. L., Doallo, R. & Posada, D. jModelTest 2: more models,
569 new heuristics and parallel computing. *Nat. Methods* **9**, 772 (2012).
- 570 43 Keane, T. M., Creevey, C. J., Pentony, M. M., Naughton, T. J. & McLnerney, J. O.
571 Assessment of methods for amino acid matrix selection and their use on empirical
572 data shows that ad hoc assumptions for choice of matrix are not justified. *BMC Evol.*
573 *Biol.* **6**, 29 (2006).

574 44 Lanfear, R., Calcott, B., Ho, S. Y. & Guindon, S. Partitionfinder: combined selection
575 of partitioning schemes and substitution models for phylogenetic analyses. *Mol. Biol.*
576 *Evol.* **29**, 1695-1701 (2012).

577 45 Hoff, M., Orf, S., Riehm, B., Darriba, D. & Stamatakis, A. Does the choice of
578 nucleotide substitution models matter topologically? *BMC Bioinformatics* **17**, 143
579 (2016).

580 46 Huelsenbeck, J. & Rannala, B. Frequentist properties of Bayesian posterior
581 probabilities of phylogenetic trees under simple and complex substitution models.
582 *Syst. Biol.* **53**, 904-913 (2004).

583 47 Wright, A. M., Lloyd, G. T. & Hillis, D. M. Modeling character change heterogeneity
584 in phylogenetic analyses of morphology through the use of priors. *Syst. Biol.* **65**, 602-
585 611 (2016).

586 48 Felsenstein, J. Maximum-likelihood estimation of evolutionary trees from continuous
587 characters. *Am. J. Hum. Genet.* **25**, 471-492 (1973).

588 49 Felsenstein, J. *Inferring Phylogenies*. (Sinauer Associates Sunderland, 2004).

589 50 Ronquist, F. *et al.* A total-evidence approach to dating with fossils, applied to the
590 early radiation of the Hymenoptera. *Syst. Biol.* **61**, 973-999 (2012).

591 51 Heath, T. A., Huelsenbeck, J. P. & Stadler, T. The fossilized birth-death process for
592 coherent calibration of divergence-time estimates. *Proc. Natl. Acad. Sci. USA* **111**,
593 E2957-2966 (2014).

594 52 O'Reilly, J. E., dos Reis, M. & Donoghue, P. C. Dating tips for divergence-time
595 estimation. *Trends Genet.* **31**, 637-650 (2015).

596 53 Rannala, B. Identifiability of parameters in MCMC Bayesian inference of phylogeny.
597 *Syst. Biol.* **51**, 754-760 (2002).

598 54 Gu, X., Fu, Y. X. & Li, W. H. Maximum likelihood estimation of the heterogeneity of
599 substitution rate among nucleotide sites. *Mol. Biol. Evol.* **12**, 546-557 (1995).

600 55 Sullivan, J., Swofford, D. L. & Naylor, G. J. The effect of taxon sampling on
601 estimating rate heterogeneity parameters of maximum-likelihood models. *Mol. Biol.*
602 *Evol.* **16**, 1347-1356 (1999).

603 56 Yang, Z. The BPP program for species tree estimation and species delimitation. *Curr.*
604 *Zoo.* **61**, 854-865 (2015).

605 57 Shapiro, B., Rambaut, A. & Drummond, A. J. Choosing appropriate substitution
606 models for the phylogenetic analysis of protein-coding sequences. *Mol Biol Evol* **23**,
607 7-9 (2006).

608 58 Yang, Z. & Rannala, B. Bayesian estimation of species divergence times under a
609 molecular clock using multiple fossil calibrations with soft bounds. *Mol. Biol. Evol.*
610 **23**, 212-226 (2006).

611 59 Nylander, J. A., Ronquist, F., Huelsenbeck, J. P. & Nieves-Aldrey, J. L. Bayesian
612 phylogenetic analysis of combined data. *Syst. Biol.* **53**, 47-67 (2004).

613 60 Maddison, W. P. Gene trees in species trees. *Syst. Biol.* **46**, 523-536 (1997).

614 61 Nichols, R. Gene trees and species tree are not the same. *Trends Ecol. Evol.* **16**, 358-
615 364 (2001).

616 62 Liu, L. & Pearl, D. K. Species trees from gene trees: reconstructing Bayesian
617 posterior distributions of a species phylogeny using estimated gene tree distributions.
618 *Syst. Biol.* **56**, 504-514 (2007).

619 63 Edwards, S. V. *et al.* Implementing and testing the multispecies coalescent model: A
620 valuable paradigm for phylogenomics. *Mol. Phylogenet. Evol.* **94**, 447-462 (2016).

621 64 Vijaykrishna, D., Mukerji, R. & Smith, G. J. D. RNA virus reassortment: an
622 evolutionary mechanism for host jumps and immune evasion. *PLoS Pathog.* **11**,
623 e1004902 (2015).

624 65 Ronquist, F., van der Mark, P. & Huelsenbeck, J. P. in *The Phylogenetic Handbook:
625 A Practical Approach to Phylogenetic Analysis and Hypothesis Testing* (eds P.
626 Lemey, M. Salemi, & A.M. Vandamme) (Cambridge Univ. Press, 2009).

627 66 Brown, J. M., Hedtke, S. M., Lemmon, A. R. & Lemmon, E. M. When trees grow too
628 long: investigating the causes of highly inaccurate bayesian branch-length estimates.
629 *Syst. Biol.* **59**, 145-161 (2010).

630 67 Rannala, B., Zhu, T. & Yang, Z. Tail paradox, partial identifiability, and influential
631 priors in Bayesian branch length inference. *Mol. Biol. Evol.* **29**, 325-335 (2012).

632 68 dos Reis, M., Zhu, T. & Yang, Z. The impact of the rate prior on Bayesian estimation
633 of divergence times with multiple loci. *Syst. Biol.* **63**, 555-565 (2014).

634 69 Drummond, A. J., Ho, S. Y., Phillips, M. J. & Rambaut, A. Relaxed phylogenetics
635 and dating with confidence. *PLoS Biol.* **4**, e88 (2006).

636 70 Yang, Z. & Rannala, B. Bayesian phylogenetic inference using DNA sequences: a
637 Markov chain Monte Carlo method. *Mol. Biol. Evol.* **14**, 717-724 (1997).

638 71 Rannala, B. & Yang, Z. Bayes estimation of species divergence times and ancestral
639 population sizes using DNA sequences from multiple loci. *Genetics* **164**, 1645-1656
640 (2003).

641 72 Ho, S. Y. & Phillips, M. J. Accounting for calibration uncertainty in phylogenetic
642 estimation of evolutionary divergence times. *Syst. Biol.* **58**, 367-380 (2009).

643 73 Thorne, J. L., Kishino, H. & Painter, I. S. Estimating the rate of evolution of the rate
644 of molecular evolution. *Mol. Biol. Evol.* **15**, 1647-1657 (1998).

- 645 74 Rannala, B. & Yang, Z. Inferring speciation times under an episodic molecular clock.
646 *Syst. Biol.* **56**, 453-466 (2007).
- 647 75 Kimura, M. A simple method for estimating evolutionary rates of base substitutions
648 through comparative studies of nucleotide sequences. *J. Mol. Evol.* **16**, 111-120
649 (1980).
- 650 76 Yang, Z. & Rodriguez, C. E. Searching for efficient Markov chain Monte Carlo
651 proposal kernels. *Proc. Natl. Acad. Sci. USA* **110**, 19307-19312 (2013).
- 652 77 Green, P. J. Reversible jump Markov chain Monte Carlo computation and Bayesian
653 model determination. *Biometrika* **82**, 711-732 (1995).
- 654 78 Lakner, C., van der Mark, P., Huelsenbeck, J. P., Larget, B. & Ronquist, F. Efficiency
655 of Markov chain Monte Carlo tree proposals in Bayesian phylogenetics. *Syst. Biol.*
656 **57**, 86-103 (2008).
- 657 79 Green, P. J. & Han, X. L. in *Stochastic Models, Statistical Methods, and Algorithms*
658 *in Image Analysis* (eds P. Barone, A. Frigessi, & M. Piccioni) (Springer, 1992).
- 659 80 R Core Team. *R: A language and environment for statistical computing*, <
660 <http://www.r-project.org/>> (2015).
- 661 81 Rambaut, A., Suchard, M. A., Xie, D. & Drummond, A. J. *Tracer v1.6*,
662 <<http://beast.bio.ed.ac.uk/Tracer>> (2014).
- 663 82 dos Reis, M., Donoghue, P. C. & Yang, Z. Bayesian molecular clock dating of species
664 divergences in the genomics era. *Nat. Rev. Genet.* **17**, 71-80 (2016).
- 665 83 Solís-Lemus, C., Knowles, L. L. & Ané, C. Bayesian species delimitation combining
666 multiple genes and traits in a unified framework. *Evolution* **69**, 492-507 (2015).
- 667 84 Chen, M.-H., Kuo, L. & Lewis, P. *Bayesian Phylogenetics: Methods, Algorithms, and*
668 *Applications*. (Chapman & Hall/CRC, 2014).
- 669 85 Gelman, A. *et al. Bayesian Data Analysis*. (Chapman and Hall/CRC, 2013).

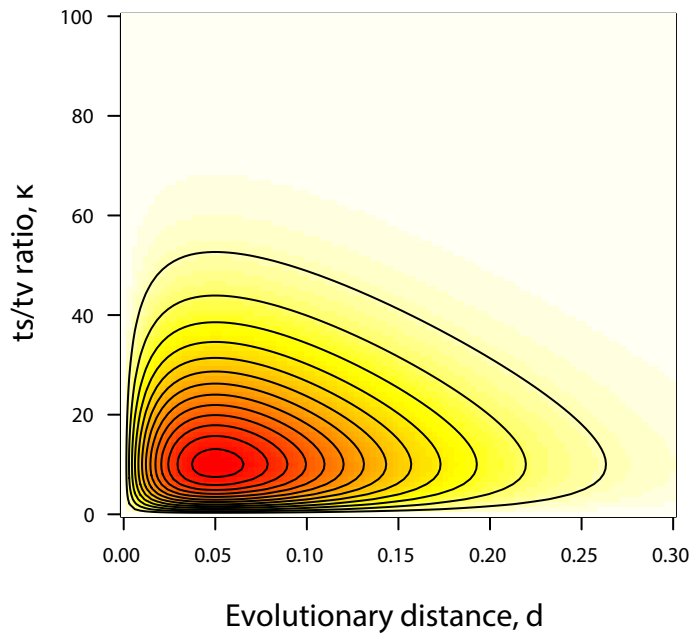
- 670 86 Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary
671 analysis. *PLoS Comput. Biol.* **10**, e1003537 (2014).
- 672 87 Ronquist, F. *et al.* MrBayes 3.2: efficient Bayesian phylogenetic inference and model
673 choice across a large model space. *Syst. Biol.* **61**, 539-542 (2012).
- 674 88 Höhna, S. *et al.* RevBayes: Bayesian phylogenetic inference using graphical models
675 and an interactive model-specification language. *Syst. Biol.* **65**, 726-736 (2016).
- 676 89 Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.*
677 **24**, 1586-1591 (2007).
- 678 90 Lewis, P. O., Holder, M. T. & Swofford, D. L. Phycas: software for Bayesian
679 phylogenetic analysis. *Syst. Biol.* **64**, 525-531 (2015).
- 680 91 Lewis, P. O., Holder, M. T. & Holsinger, K. E. Polytomies and Bayesian phylogenetic
681 inference. *Syst. Biol.* **54**, 241-253 (2005).
- 682 92 Lartillot, N., Lepage, T. & Blanquart, S. PhyloBayes 3: a Bayesian software package
683 for phylogenetic reconstruction and molecular dating. *Bioinformatics* **25**, 2286-2288
684 (2009).
- 685 93 Beerli, P. Comparison of Bayesian and maximum-likelihood inference of population
686 genetic parameters. *Bioinformatics* **22**, 341-345 (2006).
- 687 94 Hey, J. & Nielsen, R. Integration within the Felsenstein equation for improved
688 Markov chain Monte Carlo methods in population genetics. *Proc. Natl. Acad. Sci.*
689 *USA* **104**, 2785-2790 (2007).
- 690 95 Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of population structure using
691 multilocus genotype data. *Genetics* **155**, 945-959 (2000).
- 692 96 Rabosky, D. L. Automatic detection of key innovations, rate shifts, and diversity-
693 dependence on phylogenetic trees. *PLoS One* **9**, e89543 (2014).

694 97 Nylander, J. A., Wilgenbusch, J. C., Warren, D. L. & Swofford, D. L. AWTY (are we
695 there yet?): a system for graphical exploration of MCMC convergence in Bayesian
696 phylogenetics. *Bioinformatics* **24**, 581-583 (2008).

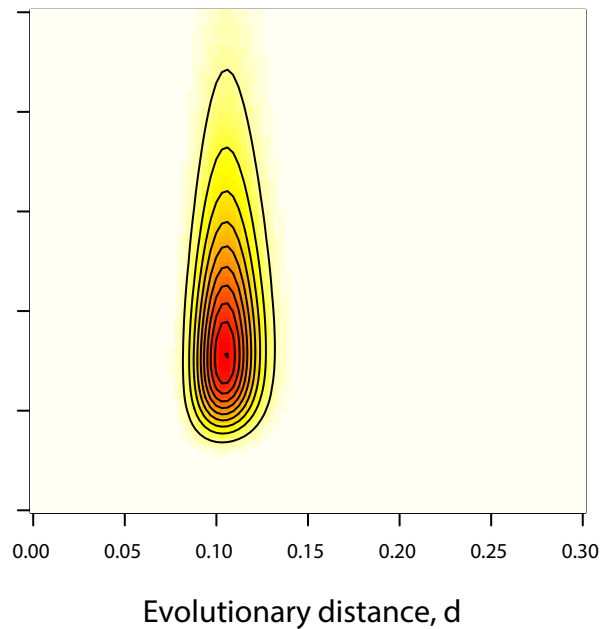
697 98 Efron, B. & Tibshirani, R. J. *An Introduction to the Bootstrap*. (Chapman &
698 Hall/CRC, 1994).

699

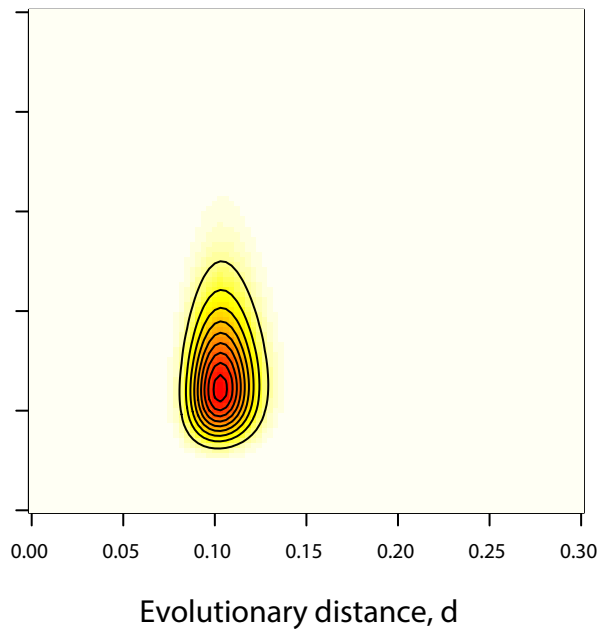
a Prior, $f(d) \times f(\kappa)$

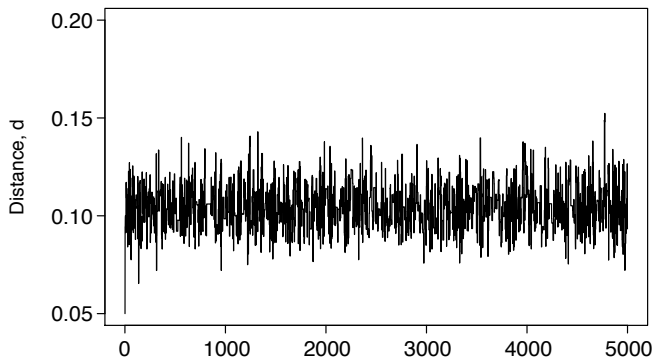
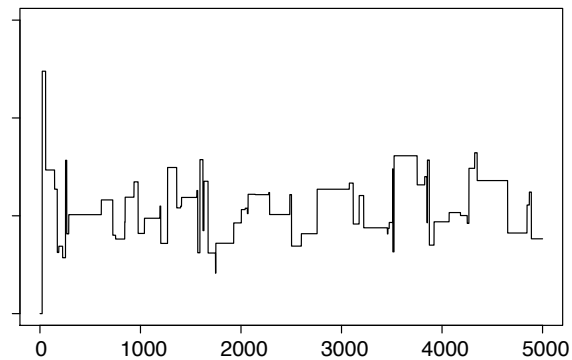
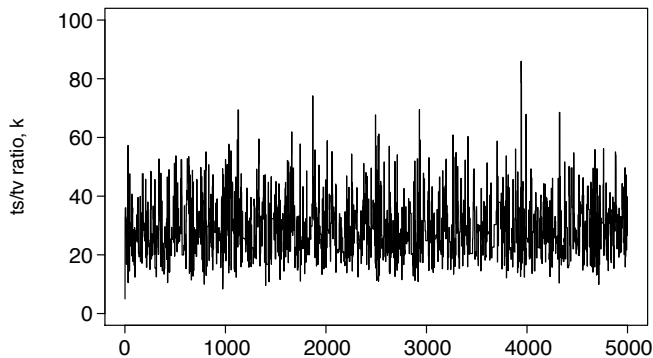
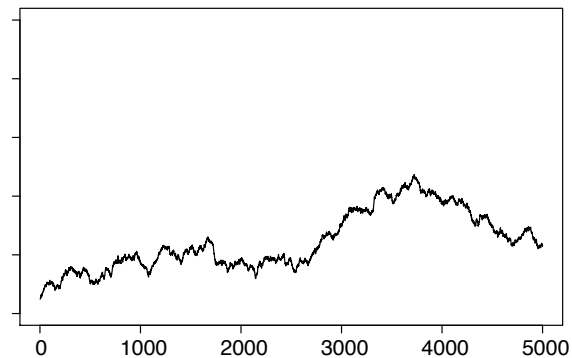
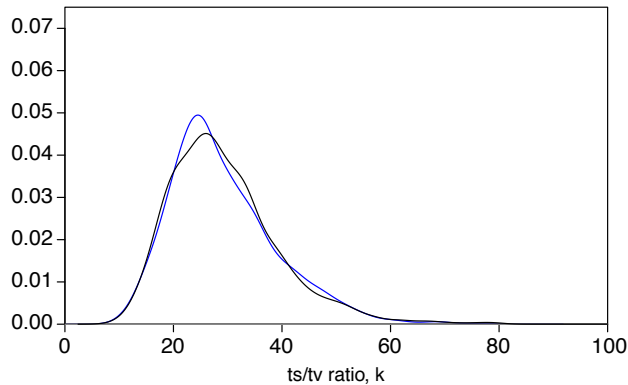


b Likelihood, $f(D | d, \kappa)$



c Posterior, $f(d, \kappa | D)$



a Trace of d , efficient chain**a' Trace of d , inefficient chain****b Trace of k , efficient chain****b' Trace of k , inefficient chain****c Histograms of k , efficient chain****c' Histograms of k , inefficient chain**