

The error in Bayesian phylogenetic inference
when speciation is protracted

Richèl J.C. Bilderbeek¹ and Rampal S. Etienne¹

⁴ ¹Groningen Institute for Evolutionary Life Sciences, University of
⁵ Groningen, Groningen, The Netherlands

6 December 16, 2019

Abstract

The tools for reconstructing phylogenetic relationships between taxonomic units (e.g. species) have become very advanced in the last three decades. Among the most popular tools are Bayesian approaches, such as BEAST, MrBayes and RevBayes, that use efficient tree sampling routines to create a posterior probability distribution of the phylogenetic tree. A feature of these approaches is the possibility to incorporate known or hypothesized structure of the phylogenetic tree through the tree prior. It has been shown that the effect of the prior on the posterior distribution of trees can be substantial.

Currently implemented tree priors assume that speciation is instantaneous, where we know that speciation can be a gradual process.

Here we explore the effects of ignoring the protractedness of the speciation process with an extensive simulation study.

We compare the inferred tree to the simulated tree, and find that

Keywords: computational biology, evolution, phylogenetics, Bayesian analysis, tree prior

1 Introduction

The computational tools that are currently available to the phylogeneticists go beyond the wildest imagination of those living four decades ago. Advances in computational power allowed the first cladograms to be inferred from DNA alignments in 1981 (Felsenstein 1981), and the first Bayesian tools emerged in 1996 (Rannala & Yang 1996), providing unprecedented flexibility in the setup of a phylogenetic model.

Currently, the most popular Bayesian phylogenetics tools are BEAST (Drummond & Rambaut 2007) and its offshoot BEAST2 (Bouckaert *et al.* 2014), MrBayes (Huelsenbeck & Ronquist 2001) and RevBayes (Höhna

34 *et al.* 2016). They allow to incorporate known or hypothesized structure of a
 35 phylogenetic tree-to-be-inferred through model priors. With these priors and
 36 an alignment of DNA, RNA or protein sequences, they create a sample of the
 37 posterior distribution of phylogenies and parameter estimates (of the models
 38 used as a prior), in which more probable combinations are represented more
 39 often. Each of these tools use efficient tree sampling routines to rapidly create
 40 an informative posterior.

41 The model priors in Bayesian phylogenetic reconstruction can be grouped
 42 into three categories: (1) site model, specifying nucleotide substitutions, (2)
 43 clock model, specifying the rate of mutation per lineage in time, and (3) tree
 44 model, constituting the speciation model underlying branching events (specia-
 45 tion) and branch termination (extinction). The choice of site model (Posada &
 46 Buckley 2004), clock model (Baele *et al.* 2012) or tree prior (Möller *et al.* 2018;
 47 Yang & Rannala 2005) is known to affect the posterior. There is evidence, how-
 48 ever, that the tree prior and molecular clock do not do so substantially affect
 49 the estimation of diversification rates (Sarver *et al.* 2019)).

50 Current phylogenetic tools use tree priors that assume speciation is instan-
 51 taneous, whilst we know that, speciation is often a gradual process (Schluter
 52 2009). The (constant-rate) birth-death (BD) model is a commonly used tree
 53 prior, but it ignores this temporal aspect of speciation. The protracted birth-
 54 death (PBD) model, an extension of the BD model, does incorporate the idea
 55 that speciation takes time. In this model, a branching event does not give rise
 56 to a new species, but to a new species-to-be, called an incipient species. Such an
 57 incipient species may go extinct, finish its speciation to become a good species,
 58 or give rise to new incipient species. Protracted speciation may explain observed
 59 declines in lineage accumulation (Etienne & Rosindell 2012).

60 Unfortunately, a tree prior according to this model, providing the probability

61 of a species tree under the PBD model, is unavailable in current Bayesian phy-
 62 logenetic tools. Whilst an approximate formula for this probability has been
 63 derived (Lambert *et al.* 2015) and the approximation is very good (Simonet
 64 *et al.* 2018), it has not been implemented as tree prior yet. There are various
 65 reasons for this. First, the computation of this probability involves solving a set
 66 of non-linear differential equations, and while this computation is quite fast, it
 67 still takes much more time than the corresponding probability of the BD model
 68 which is a simple analytical formula. In a Bayesian MCMC chain, the tree
 69 prior probability must be calculated many times, and hence the total compu-
 70 tation will take considerably longer with a PBD tree prior. Furthermore, the
 71 approximate probability is a probability for the species tree assuming an under-
 72 lying incipient species tree. It can be safely used as tree prior when only one
 73 individual per species is sampled, but if one has multiple samples per species
 74 - which is currently often the case - the methods to account for this such as
 75 the multi-species coalescent (Heled & Drummond 2009) may not be compatible
 76 with the underlying incipient species tree. More precisely, the phylogeny under
 77 the PBD model may contain paraphyly, while the multi-species coalescent was
 78 developed exactly to avoid these by explaining them as arising from incomplete
 79 lineage sorting. Because of these paraphyly there is no such thing as a true
 80 species tree in the PBD model. To get a species-level tree one must sample one
 81 incipient species per species. Which incipient species is sampled may therefore
 82 have an impact on the species tree.

83 Here we aim to explore the effect of using the BD prior on PBD simulated
 84 phylogenies, taking into account possible sampling effects. In brief, we simulate
 85 protracted phylogenies using the PBD process, from which we sample a species
 86 tree in two very different ways. Given this species tree, we simulate a DNA
 87 sequence alignment. Then, we use BEAST2 on these alignments to infer a pos-

terior of phylogenies, using a BD prior. We quantify the difference between the (BD) posterior phylogenies and the simulated (PBD) species tree. Furthermore, while we evidently know the clock and site models used in the simulation, using a different clock and/or site model prior in inference may compensate or increase this difference between inferred and simulated tree. To study this, we also explore the effect of a different clock and site model prior in inference.

2 Hypotheses

[RJCB: but we are not allowed to use this header]

3 Hypotheses

\mathcal{H}_1 : RJCB expects that the inference error is lowest when the true tree is generated under PBD parameter settings without protractedness, i.e. $scr = \infty$. For such parameters settings, the true tree is in practice generated by a BD model, which matches the tree prior used in inference. Without a mismatch between true tree prior and assumed tree prior, this source of errors will be absent, where the other two sources of errors (stochasticity in the simulation of the alignment and the MCMC algorithm) will remain the same. [RJCB: figure 1a and b]

\mathcal{H}_2 : RJCB expects that, on average, the inference error is highest for parameter settings with a lower speciation rate, as for these settings, the mismatch between the the speciation model that generated the true tree (which is profoundly-P PBD) is biggest with the tree prior assumed to be generative (which is BD). [RJCB: figure 1a and b]

\mathcal{H}_3 : RJCB expects that the inference error is higher for true trees with an observable/actual higher percentage of extant incipient species, regardless of the

112 PBD parameters. These incipient species are one of the three (and the most
 113 interesting) sources of error, as a BD model -as a feature- will never infer the
 114 branch-length distribution of protracted species. **[RJCB: figure 2]**

115 \mathcal{H}_4 : RJCB expects that, on average, the inference error is equal for parame-
 116 ter settings with an equal SIRI, SIRG and SCR, i.e. extinction has no effect on
 117 the error made, because extinction affects all species equally. **[RJCB: figure**
 118 **1a and b, or a new figure]**

119 \mathcal{H}_5 : RJCB expects that an increased number of taxa has a negative effect
 120 on the variance of the errors made, as there is more information available to
 121 base inference on. **[RJCB: figure 3]**

122 \mathcal{H}_6 : RJCB expects that an increased extinction rate increases the variance
 123 of the errors made, due to the decrease of the number of taxa, reducing the
 124 amount of information to base inference on. **[RJCB: figure 4]**

125 \mathcal{H}_7 : RJCB expects that the nLTT statistic between a true and twin tree at
 126 the start of the pipeline, will correlate strongly with the difference between the
 127 highest posterior density (HPD) of the errors of the true and twin tree generated
 128 at the end of the pipeline. This hypothesis stems from the idea of assuming a
 129 setup without noise; that is, with a DNA alignment of infinite length. From
 130 such a DNA alignment with infinite information, the MCMC is able to infer
 131 the close-to-correct phylogeny. **[RJCB: if this is true, we can use this**
 132 **shortcut (instead of running multiple pirouette replicates) to draw**
 133 **stronger conclusions] [RJCB: figure 5]**

134 \mathcal{H}_8 : RJCB expects that for settings without extinctions, the candidate model
 135 with JC69, strict and Yule will come up as the best model most, as that model
 136 is the generative model, where the HKY, strict, BD will be selected least, as both
 137 its site and tree model are needlessly complex. For settings with extinctions,
 138 RJCB predict HKY, strict and Yule will come up as the worst model mostly, as

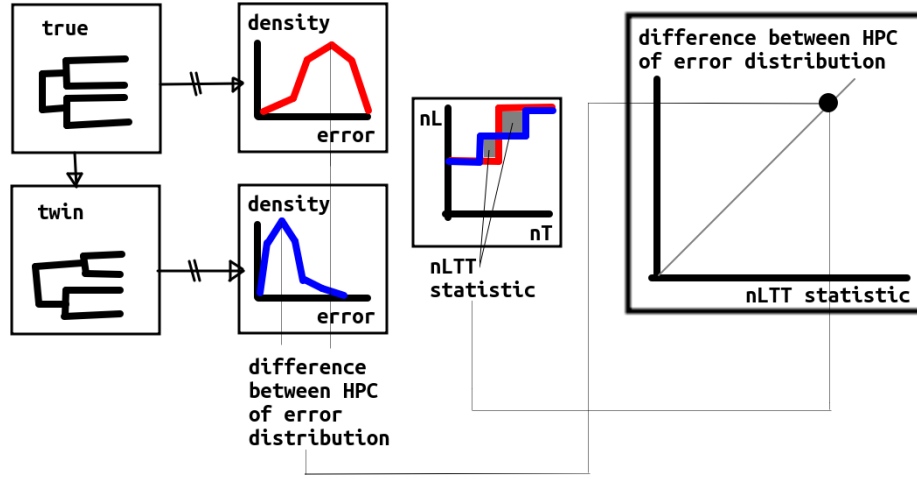


Figure 1: The nLTT statistic between true and twin tree correlates strongly with the difference between (the highest posterior density of) the errors of the true and twin tree.

139 it uses both an incorrect site model and incorrect tree model.

140 4 Methods

141 [RJC: but we are not allowed to use this header]

142 The PBD model has five parameters, depicted in table 2. The per species
 143 speciation completion rates λ we use are 0.1, 0.3, 1.0 and 10^9 . This means that
 144 there is a $0.15 dt$ ($0.35 dt$, $1.0 dt$, $10^9 dt$) probability of speciation completion
 145 occurring in an infinitesimal time dt . We use per species extinction rates of
 146 $\mu = \mu_g = \mu_i$ 0.0, 0.1 and 0.2. We use tree sizes of n of 50, 100 and 200 good
 147 taxa. For each combination of λ , μ and n , we use a speciation initiation rate
 148 $b = b_i = b_g$ so that the expected mean number of species $\mathbf{E}(\bar{n}; b, \lambda, \mu)$, given a b ,
 149 λ and μ , equals the desired number of species n . b is calculated using the PBD
 150 R package (Etienne 2015) for each parameter combination and shown in table
 151 3. We use $\lambda = 10^9 \approx \infty$ as our control for which the PBD model reduces to the
 152 BD model.

153 We simulate protracted birth-death trees, using the PBD package (Etienne
 154 2015) in the R programming language (R Core Team 2013). The first tree
 155 has a random number generator seed of 1, which is incremented by 1 for each
 156 simulated tree. For each combination of λ, μ, b, n , we generate incipient species
 157 trees with a crown age of 15 million years. Only trees with the desired number
 158 of good taxa are kept.

159 We create two data sets: a general one, to explore parameter space, and one
 160 to investigate the effect of sampling incipient species (see below). For the general
 161 data set, all the trees with the correct number of good species are kept. There
 162 is an additional selection criterion, for the data set to investigate sampling: to
 163 generate that data set, only incipient species trees are kept, on which the two
 164 sampling methods (see below) result in different species trees. As sampling will
 165 never have an effect for $\lambda = \infty$, this parameter value is absent in that data set.

166 From each incipient species tree, we construct a species tree, by sampling one
 167 incipient/good species per good species. For example, when an incipient species
 168 branched off from its mother lineage, both of these subspecies are recognized
 169 as representing the species, and hence both can be picked as an (equally good)
 170 representative of the species. Here, we use three sampling scenarios, in which
 171 we pick the representative randomly or in such a way that this results in either
 172 the shortest or longest branch lengths.

173 See the supplementary information for a visualization of these sampling
 174 methods. Based on the sampled species tree, we simulate a DNA alignment that
 175 has the same history as this species tree, using the **phangorn** package (Schliep
 176 2011). We set the nucleotides of the DNA alignment to follow a Jukes-Cantor
 177 (Jukes *et al.* 1969) nucleotide substitution model, in which all nucleotide-to-
 178 nucleotide transitions are equally likely. The DNA sequence of the root ancestor
 179 consists of four equally sized single-nucleotide blocks of adenine, cytosine, gua-

180 nine and thymine respectively. For example, for a DNA sequence length of 12,
 181 this would be AAACCCGGGTTT. The order of nucleotides does not matter in
 182 this study, because we do not consider several partitions of the sequence with
 183 their own parameters. Only the frequency of occurrence matters. In our Bayes-
 184 ian inference (see below) we use the same site model as the (obviously correct)
 185 site model prior, but we also explore the effect of assuming a more complex site
 186 model prior. We predict with the more complex substitution model, that there
 187 will be more noise and hence our inference error will increase. On the other
 188 hand, we dare not rule out that the inference error will decrease, due to more
 189 flexibility in the more complex prior. We set the mutation rate in such a way
 190 to maximize the information contained in the alignment. To do so, we set the
 191 mutation rate such that we expect on average one (possibly silent) mutation per
 192 nucleotide between crown age and present, which equates to $\frac{1}{15}$ mutations per
 193 million years. The DNA sequence length is chosen to provide a resolution of 10^3
 194 years, that is, to have one expected nucleotide change per 10^3 years per lineage
 195 on average. As one nucleotide is expected to have on average one (possibly
 196 silent) mutation per 15 million years, $15 \cdot 10^3$ nucleotides result in 1 mutation
 197 per alignment per 10^3 years (which is coincidentally the same as Möller *et al.*
 198 2018). The simulation of these DNA alignments follows a strict clock model,
 199 which we will specify as one of the two clock models assumed in the Bayesian
 200 inference (see below).

201 From an alignment, we run a Bayesian analysis and create a posterior dis-
 202 tribution of trees and parameters using the **babette** (Bilderbeek & Etienne
 203 2018) package that sets the input parameters similar to BEAUti 2 and then
 204 runs BEAST2. For our site model, we assume either a Jukes-Cantor or GTR
 205 nucleotide substitution model. The Jukes-Cantor model is the correct one, as it
 206 is used for simulating that alignment, where the GTR model is the site model

207 that is picked as a default by most users. For our clock model, we assume either
 208 a strict or relaxed log-normal clock model. Also here, the strict clock model
 209 is the correct one, as it is used for simulating the alignment, but the relaxed
 210 log-normal clock model is the one most commonly used. We set the BD model
 211 as a tree prior, as gauging the effect of this incorrect assumption is the goal of
 212 this study. We assume an MRCA prior with a tight normal distribution around
 213 the crown age, by choosing the crown age as mean, and a standard deviation of
 214 $0.5 \cdot 10^{-3}$ time units, resulting in 95% of the crown ages inferred have the same
 215 resolution (of 10^{-3} time units) as the alignment. We ran the MCMC chain to
 216 generate 1111 states, of which we remove the first 10% (also called the 'burn-
 217 in'). Of the remaining 1000 MCMC states, the effective sample size (ESS) of
 218 the posterior must at least be 200 for a strong enough inference (Drummond &
 219 Bouckaert 2015). An ESS can be increased by increasing the number of samples
 220 or decreasing the autocorrelation between samples. If the ESS is less than 200,
 221 we decrease autocorrelation by doubling the MCMC sampling interval of that
 222 simulation, until the ESS exceeds 200.

223 We compare each posterior phylogeny to the (sampled) species tree using the
 224 nLTT statistic (Janzen *et al.* 2015), from the nLTT package (Janzen 2015). The
 225 nLTT statistic equals the area between the normalized lineages-through-time-
 226 plots of two phylogenies, which has a range from zero (for identical phylogenies)
 227 to one. We use inference error and nLTT statistic interchangeably. Comparing
 228 the simulated species tree with each of the posterior species trees yields a dis-
 229 tribution of nLTT statistics. The input trees generated with a $\lambda = 10^9$ allow
 230 us to measure the noise of the experiment. For $\lambda = \infty$, the PBD model that
 231 generates the starting trees reduces to a BD model. In the subsequent steps,
 232 sampling will have no effect in this case, because BEAST2 will assume the cor-
 233 rect speciation model, so the difference between inferred tree and true species

234 tree are explained purely due to this experimental noise.

235 TWINNING

236 As described above, per alignment, we do four different Bayesian analyses, as
237 we use two different site models and two different clock models. We know the
238 generative site model (which is JC69) and clock model (a strict clock), but want
239 to explore how often the correct model is indeed preferred. Per alignment, we
240 measure the estimated marginal likelihood (which is the probability of the data
241 given the model) for each of the four models and measure their relative proper-
242 tions. To estimate the marginal likelihood, we use the novel Nested Sampling
243 approach (Russel *et al.* 2018), which we configure to have a relative error ϵ of
244 1^{-12} .

245 We produce two data sets as a comma-separated file. The general data set
246 has 144 [RJC*B*: recal*c*] different combinations of biological parameter combi-
247 nations, site and clock models. The data set to investigate sampling has 552
248 [RJC*B*: recal*c*] different combinations of biological parameter combinations,
249 site models, clock models and sampling methods. The experiment is compu-
250 tationally intensive: pilot experiments show that the experiment takes roughly
251 100 days of CPU time and 20 days of wall clock time (which includes the queued
252 waiting for computational resources) per replicate. Due to this, we choose to
253 perform ten replicates, so that the complete experiment will take an acceptable
254 time of roughly seven months.

255 In this article, we showcase the effect of sampling and the certainty when
256 selecting a (site and clock) model. We show the effect sampling has on the
257 inference error, for the nLTT distributions generated from assuming the (cor-
258 rect) Jukes-Cantor site model and a strict clock model. We do so per parameter
259 combination, displaying the distribution as a violin plot. In the supplementary
260 information, we display these error distributions of the general data set all com-

261 bined, or separated per site or clock model. The certainty in the model selection
 262 is again shown per parameter setting, as a stacked bar.

263 HIERO: DESCRIBE PLOTTING.

Plotting a stacked bar with uncertainty

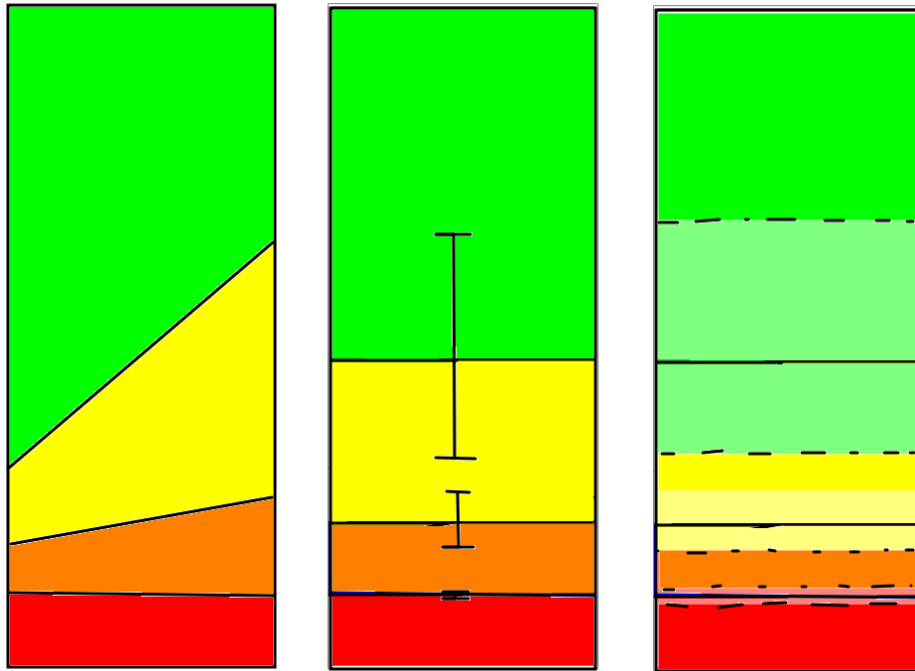


Figure 2: Showing a stacked bar with uncertainty. These three different ways display the same underlying data.

264 5 Results

265 6 Acknowledgements

266 We would like to thank the Center for Information Technology of the University
 267 of Groningen for their support and for providing access to the Peregrine high

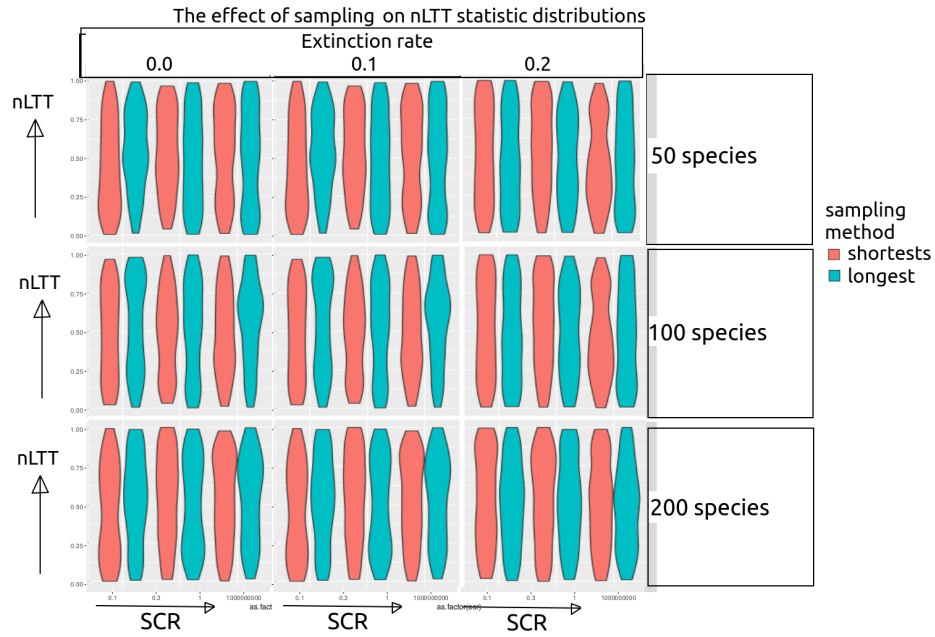


Figure 3: nLTT statistic distribution per biological parameter set per sampling regime, using the data set conditioned on sampling regime having an effect, under the (correct) assumptions of a strict clock and Jukes-Cantor site model.

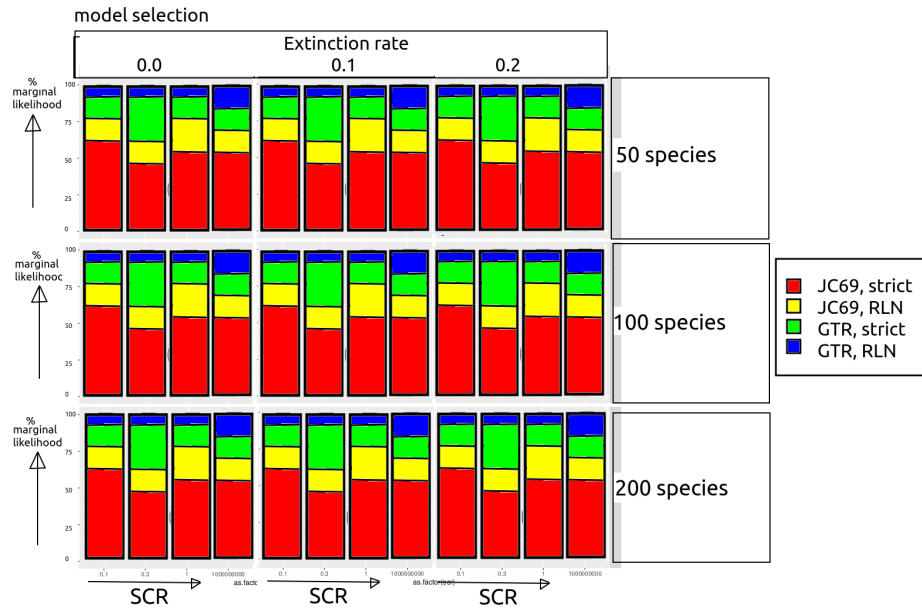


Figure 4: Model preference on the general data set.

Term	Definition
Phylogenetics	The inference of evolutionary relationships of groups of organisms using genetics
Model prior	Knowledge or assumptions about the ontogeny of evolutionary histories
Posterior	A collection of phylogenies and parameter estimates, in which more probable combinations (determined by the data and the model prior) are presented more frequently
Protracted speciation	The process in which speciation takes two events: a speciation-initiation event and a speciation-completion event
Speciation initiation	The start of a speciation event creating an incipient species
Speciation completion	The end of a speciation event, in which an incipient species becomes or is recognized as a good species

Table 1: Glossary

performance computing cluster.

7 Authors' contributions

RSE conceived the idea for this experiment. RJCB created and tested the experiment, and wrote the first draft of the manuscript. RSE contributed substantially to revisions.

8 Glossary

References

Baele, G., Li, W.L.S., Drummond, A.J., Suchard, M.A. & Lemey, P. (2012) Accurate model selection of relaxed molecular clocks in bayesian phylogenetics. *Molecular biology and evolution*, **30**, 239–243.

278 Bilderbeek, R.J. & Etienne, R.S. (2018) babette: Beau ti 2, beast 2 and tracer
 279 for r. *Methods in Ecology and Evolution*.

280 Bouckaert, R., Heled, J., Kühnert, D., Vaughan, T., Wu, C.H., Xie, D., Suchard,
 281 M.A., Rambaut, A. & Drummond, A.J. (2014) Beast 2: a software platform
 282 for bayesian evolutionary analysis. *PLoS computational biology*, **10**, e1003537.

283 Drummond, A.J. & Bouckaert, R.R. (2015) *Bayesian evolutionary analysis with*
 284 *BEAST*. Cambridge University Press.

285 Drummond, A.J. & Rambaut, A. (2007) Beast: Bayesian evolutionary analysis
 286 by sampling trees. *BMC evolutionary biology*, **7**, 214.

287 Etienne, R.S. (2015) *PBD: Protracted Birth-Death Model of Diversification*. R
 288 package version 1.1.

289 Etienne, R.S. & Rosindell, J. (2012) Prolonging the past counteracts the pull of
 290 the present: protracted speciation can explain observed slowdowns in diver-
 291 sification. *Systematic Biology*, **61**, 204–213.

292 Felsenstein, J. (1981) Evolutionary trees from dna sequences: a maximum like-
 293 lihood approach. *Journal of molecular evolution*, **17**, 368–376.

294 Heled, J. & Drummond, A.J. (2009) Bayesian inference of species trees from
 295 multilocus data. *Molecular biology and evolution*, **27**, 570–580.

296 Höhna, S., Landis, M.J., Heath, T.A., Boussau, B., Lartillot, N., Moore, B.R.,
 297 Huelsenbeck, J.P. & Ronquist, F. (2016) Revbayes: Bayesian phylogenetic
 298 inference using graphical models and an interactive model-specification lan-
 299 guage. *Systematic biology*, **65**, 726–736.

300 Huelsenbeck, J.P. & Ronquist, F. (2001) Mrbayes: Bayesian inference of phylo-
 301 genetic trees. *Bioinformatics*, **17**, 754–755.

302 Janzen, T. (2015) *nLTT: Calculate the NLTT Statistic*. R package version 1.1.

303 Janzen, T., Höhna, S. & Etienne, R.S. (2015) Approximate bayesian compu-
 304 tation of diversification rates from molecular phylogenies: introducing a new
 305 efficient summary statistic, the nltt. *Methods in Ecology and Evolution*, **6**,
 306 566–575.

307 Jukes, T.H., Cantor, C.R. *et al.* (1969) Evolution of protein molecules. *Mam-*
 308 *malian protein metabolism*, **3**, 132.

309 Lambert, A., Morlon, H. & Etienne, R.S. (2015) The reconstructed tree in
 310 the lineage-based model of protracted speciation. *Journal of mathematical*
 311 *biology*, **70**, 367–397.

312 Möller, S., du Plessis, L. & Stadler, T. (2018) Impact of the tree prior on
 313 estimating clock rates during epidemic outbreaks. *Proceedings of the National*
 314 *Academy of Sciences*, p. 201713314.

315 Posada, D. & Buckley, T.R. (2004) Model selection and model averaging in
 316 phylogenetics: advantages of akaike information criterion and bayesian ap-
 317 proaches over likelihood ratio tests. *Systematic biology*, **53**, 793–808.

318 R Core Team (2013) *R: A Language and Environment for Statistical Computing*.
 319 R Foundation for Statistical Computing, Vienna, Austria.

320 Rannala, B. & Yang, Z. (1996) Probability distribution of molecular evolution-
 321 ary trees: a new method of phylogenetic inference. *Journal of molecular*
 322 *evolution*, **43**, 304–311.

323 Russel, P.M., Brewer, B.J., Klaere, S. & Bouckaert, R.R. (2018) Model selection
 324 and parameter inference in phylogenetics using nested sampling. *Systematic*
 325 *Biology*, p. syy050.

- 326 Sarver, B.A., Pennell, M.W., Brown, J.W., Keeble, S., Hardwick, K.M., Sulli-
 327 van, J. & Harmon, L.J. (2019) The choice of tree prior and molecular clock
 328 does not substantially affect phylogenetic inferences of diversification rates.
 329 *PeerJ*, **7**, e6334.
- 330 Schliep, K. (2011) phangorn: phylogenetic analysis in r. *Bioinformatics*, **27**,
 331 592–593.
- 332 Schluter, D. (2009) Evidence for ecological speciation and its alternative. *Sci-*
 333 *ence*, **323**, 737–741.
- 334 Simonet, C., Scherrer, R., Rego-Costa, A. & Etienne, R. (2018) Robustness of
 335 the approximate likelihood of the protracted speciation model. *Journal of*
 336 *evolutionary biology*, **31**, 469–479.
- 337 Tavaré, S. (1986) Some probabilistic and statistical problems in the analysis of
 338 dna sequences. *Lectures on mathematics in the life sciences*, **17**, 57–86.
- 339 Yang, Z. & Rannala, B. (2005) Branch-length prior influences bayesian posterior
 340 probability of phylogeny. *Systematic Biology*, **54**, 455–470.

341 9 Supplement

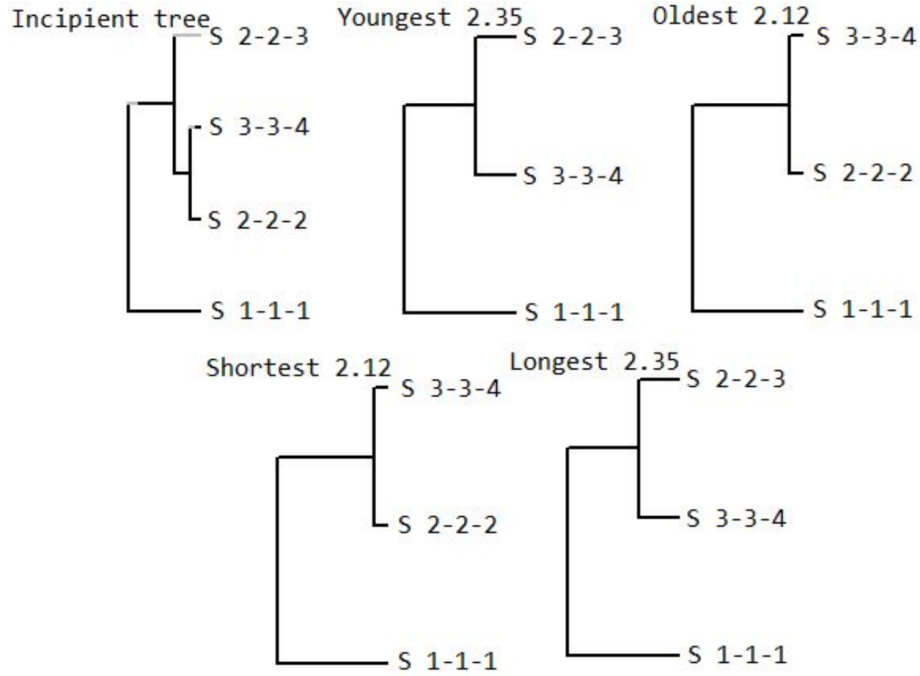


Figure 5: Sampling a species tree from an incipient species tree. At the top left, an incipient species tree is shown, of three different good species (the first and second number in the taxon label) and four different subspecies (the third number in the taxon label). The other four trees are species trees, that use a different sampling method to determine which sub-species is picked to represent a good species. These are: 'Youngest', 'Oldest', 'Shortest' and 'Longest'. With 'Youngest' the youngest sub-species is picked to represent the good species. With 'Oldest' the oldest sub-species is picked to represent the good species. 'Shortest' is the sampling method in which the sub-species are picked to assure the shortest branch lengths. 'Longest' is the sampling method in which the sub-species are picked to assure the longest branch lengths.

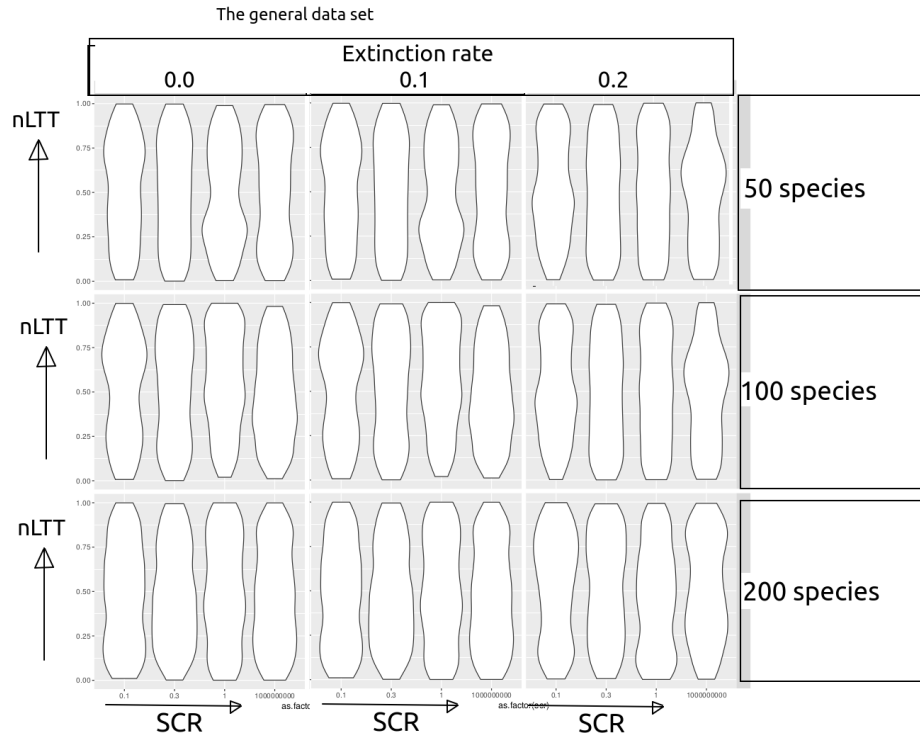


Figure 6: nLTT statistic distribution per biological parameter set, using the general data set, under the (correct) assumptions of a strict clock and Jukes-Cantor site model.

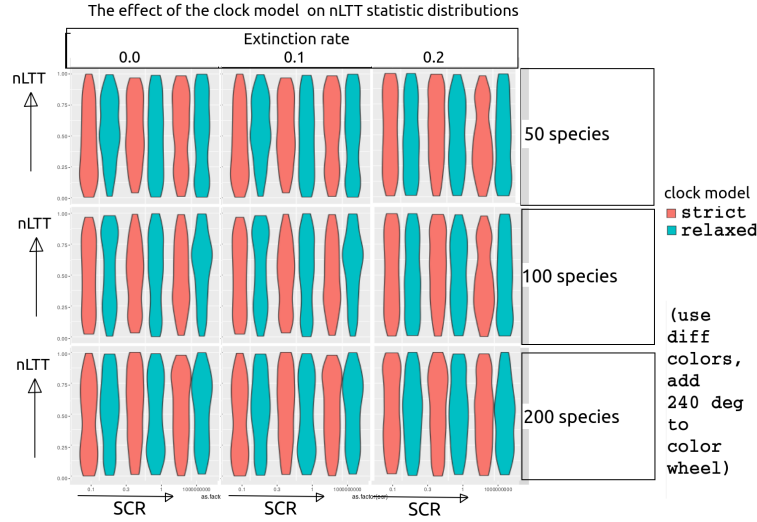


Figure 7: nLTT statistic distribution per biological parameter set per clock model, using the general data set, under the (correct) assumption of a Jukes-Cantor site model.

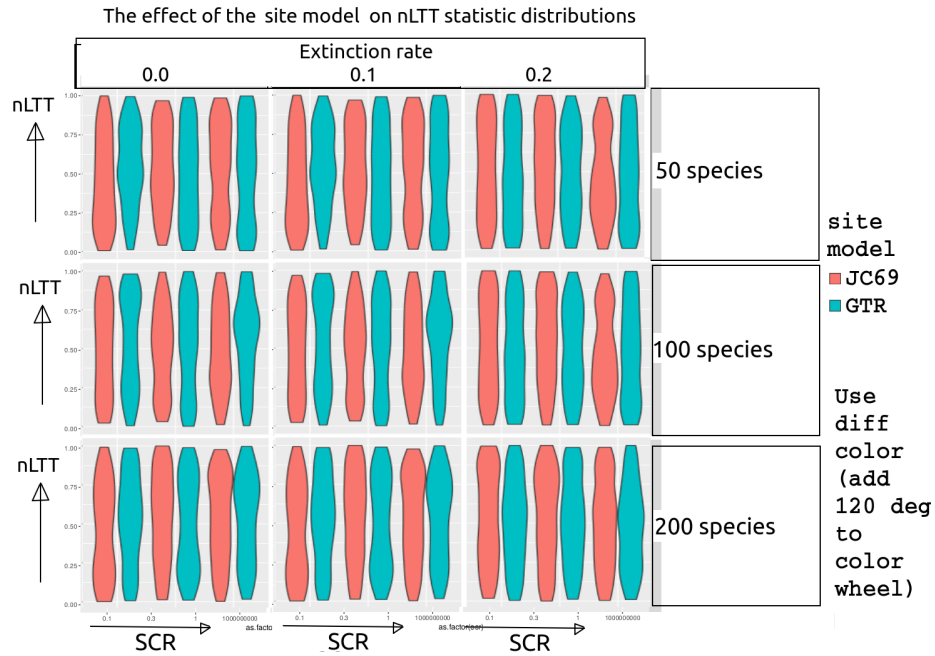


Figure 8: nLTT statistic distribution per biological parameter set per site model, using the general data set, under the (correct) assumption of a strict clock model.

	Description	Values
b_g	Speciation initiation rate of a good species	derived, see 3
b_i	Speciation initiation rate of an incipient species	derived, see 3
λ	Speciation completion rate	0.1, 0.3, 1.0, 10^9
μ_g	Extinction rate of a good species	0.0, 0.1, 0.2
μ_i	Extinction rate of an incipient species	0.0, 0.1, 0.2
n	Number of good taxa	50, 100, 200
t_c	Crown age	15
σ_c	Standard deviation around crown age	0.001
M_s	Sampling method	S, L, R
M_c	Clock model	S, RLN
M_t	Site model	JC69, GTR
r	Mutation rate	$\frac{1}{15}$
l_a	DNA alignment length	15K
f_i	MCMC sampling interval	1K or more
R_i	RNG seed incipient tree and randomly sampled species tree	1, 2, etc.
R_a	RNG seed alignment simulation	R_i
R_b	RNG seed BEAST2	R_i

Table 2: Overview of the simulation parameters. Above the horizontal line is the biological parameter set. The RNG seed R_i is 1 for the first simulation of the general data set, 2 for the next, and so on. The RNG seeds for the data set investigating the effect of sampling continue from there, but only those RNG seeds are used in which sampling has an effect. The sampling methods are abbreviated as such: 'R' denotes random sampling, 'S' is 'shortest' and 'L' is 'longest'. Sampling method M_s is random for the general data set. For the data set exploring the effect of sampling, we use 'shortest' and 'longest' for each value of R_i (which are random seeds in which sampling has an effect). The clock models are abbreviated as 'S' for a strict and 'RLN' for a relaxed log-normal model. The site models are abbreviated as 'JC69' for Jukes-Cantor (Jukes *et al.* 1969) and 'GTR' for the generalized time-reversible model (Tavaré 1986).

	μ	n	λ	b
1	0	50	0.1	0.30944
2	0.1	50	0.1	0.39674
3	0.2	50	0.1	0.48667
4	0	100	0.1	0.36344
5	0.1	100	0.1	0.45283
6	0.2	100	0.1	0.54425
7	0	200	0.1	0.41669
8	0.1	200	0.1	0.50759
9	0.2	200	0.1	0.6001
10	0	50	0.3	0.25717
11	0.1	50	0.3	0.34003
12	0.2	50	0.3	0.42648
13	0	100	0.3	0.30862
14	0.1	100	0.3	0.39455
15	0.2	100	0.3	0.48328
16	0	200	0.3	0.35991
17	0.1	200	0.3	0.44804
18	0.2	200	0.3	0.53841
19	0	50	1	0.2297
20	0.1	50	1	0.30759
21	0.2	50	1	0.38984
22	0	100	1	0.2778
23	0.1	100	1	0.35961
24	0.2	100	1	0.44481
25	0	200	1	0.32617
26	0.1	200	1	0.41078
27	0.2	200	1	0.49818
28	0	50	10^9	0.21589
29	0.1	50	10^9	0.28896
30	0.2	50	10^9	0.36635
31	0	100	10^9	0.26146
32	0.1	100	10^9	0.33872
33	0.2	100	10^9	0.41945
34	0	200	10^9	0.30733
35	0.1	200	10^9	0.38768
36	0.2	200	10^9	0.47099

Table 3: The speciation parameters used. Starting from extinction rate μ ($\mu = \mu_g = \mu_i$), the expected mean number of good species n , speciation completion rate λ , the speciation initiation rate b ($b = b_g = b_i$) follows.

n	Description
12	simulation parameters, see table 2
1000	nLTT statistic values
11	ESSes of all parameters estimated by BEAST2 (see specs below)
1	Marginal likelihood estimate
1	Marginal likelihood estimation uncertainty
1	Marginal likelihood ESS

Table 4: Specification of the data sets. Each row will contain one experiment, where the columns contain parameters, measurements and diagnostics. This table displays the content of the columns. n denotes the number of columns a certain item will occupy, resulting in a table of 1023 **[RJC]: recal** columns and 20K rows.

#	Description
1	posterior
2	likelihood
3	prior
4	treeLikelihood
5	TreeHeight
6	BirthDeath
7	BDBirthRate
8	BDDeathRate
9	logP.mrca
10	mrcatime
11	clockRate

Table 5: Overview of the 11 parameters estimated by BEAST2