

1 The error in Bayesian phylogenetic reconstruction
2 when speciation co-occurs

3 Giovanni Laudanno¹, Richèl J.C. Bilderbeek¹, and Rampal S.
4 Etienne¹

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

7 October 5, 2018

8 **Abstract**

9
10 **Keywords:** computational biology, evolution, phylogenetics, Bayesian anal-
11 ysis, tree prior [GL: According to my fine graining approach we should
12 at each step deepen every small section. At a certain level I think
13 we can start to re-coarse-grain what we wrote to create the abstract.]
14 [RJCB: I enjoy this approach! Did some minor fine-graining] [RJCB:
15 Have you already looked up for a target journal? I know how a jour-
16 nal's constraints have helped me in writing an article, for example,
17 by having a maximum number of pictures]

1 Introduction

- There are many contemporary tools that provide the possibility to infer a phylogeny from genetic data (DNA, RNA, proteins). A popular Bayesian phylogenetic tool is called BEAST and its cousin BEAST2.
- BEAST is very flexible in setting up all possible phylogenetic priors (e.g. site/clock/speciation model).
- Current limits in current tools.
- BEAST2 gives us the possibility to introduce new tree priors to infer phylogenies based on different assumptions on how the speciation process takes place.
- One of such speciation processes is the multiple birth hypothesis, a new model (described below) and thus absent in BEAST.
- The Multiple birth hypothesis can be useful to explain a phenomenon that has always puzzled evolutionary biologists: what are the drivers of the diversification processes for those phylogenies that show an impressive amount of speciation events in relatively short times? The (constant-rate) birth-death (BD) model embodies the common assumption that only a single speciation event can occur at any given time. The multiple-birth-death (MBD) model [RJCB: I feel MBSD (Multiple Birth Single Death) may be a better name: extinctions are still one at a time] relaxes this assumption, allowing events in which large-scale environmental changes lead to a great number of species in relatively short time intervals. Such a hypothesis may be a better fit to describe the burst in cichlid fish diversification in systems like in the African Great Lakes: Malawi, Tanganyika and Victoria (Janzen *et al.* 2016, Janzen *et al.* 2017).

- However, it may be that current BD tree priors are good enough at detecting such events, with a (preferred) lower level of complexity. If this is the case one should always be more keen to adopt the simplest model.
- Here we present our study with the aim of exploring when using a more complex MBD tree prior is warranted.

2 Methods

2.1 Model

- Current phylogenetic tools assume that only a single speciation event can occur at any given time. While this assumption is useful to construct a wide variety of successful models (for example: Maddison *et al.* 2007, Valente *et al.* 2015, Etienne *et al.* 2012, Etienne *et al.* 2014), they disallow for environmental changes that trigger speciations in multiple clades at a same point in time.
- The (constant-rate) birth-death (BD) model embodies the common assumption that only a single speciation event can occur at any given time. The multiple-birth-death (MBD) model relaxes this assumption, allowing events in which large-scale environmental changes lead to a great number of species in relatively short time intervals. Such hypothesis can be useful to describe, for example, systems like cichlid fish diversification in the African Great Lakes: Malawi, Tanganyika and Victoria (Janzen *et al.* 2016, Janzen *et al.* 2017).
- In the MBD model, parameters λ and μ correspond, respectively, to the common per-species speciation and extinction rates. Additionally, ν is the rate at which an environmental change is triggered. When such event

- 67 is triggered, all species present in the phylogeny at that moment have a
68 probability q to speciate at that time, which is independent on λ .
- 69 • It is also possible to write down a likelihood function for such processes
70 as in Laudanno 2018.

71 2.2 Simulations

- 72 • To prove our hypothesis we simulate two twin datasets. All the simulations
73 are produced in continuous time, using the Doob-Gillespie algorithm.
- 74 • We start simulating 1000 **[RJC: I will measure the number of trees**
75 **we'll be able to simulate within a short enough time, when the**
76 **experiment is set up]** MBD trees. From each MBD tree, a a DNA se-
77 quence alignment is simulated, after which that alignment starts a Bayes-
78 ian analysis. We use the 'pirouette' package (Bilderbeek 2018) to call the
79 BEAST2 tool suite from R. We let the Bayesian analysis assume a BD
80 prior, to investigate the error this inference makes due to this.
- 81 • For each tree generated under the MBD model we aim to generate a "twin"
82 tree under the BD model in order to perform a fair comparison, using trees
83 with the same amount of information, which is the number of taxa and
84 the same (expected) number of DNA mutations. To obtain these twin
85 trees, **[RJC: I suggest to first start with the equal number of**
86 **taxa, and the calculation of the speciation rates first.]** we impose
87 that the expected number of mutations in an MBD tree, m_{MBD} equals
88 the expected number of mutations in a BD tree, m_{BD} :

$$m_{MBD} = m_{BD} \tag{1}$$

89 We first generate a set of MBD trees. For each of them we can measure
90 the amount of mutations m_{MBD} .

91 [GL: I think this should definitely go to the methods] [RJCB: I
92 put it there for use, hopefully at a spot you liked :-)]

93 The expected number of mutations m of a phylogeny [RJCB: I think
94 'expected number of mutations' would be more correct. Do you
95 agree?] with crown age T in fact is given by [RJCB: above stood
96 'crown age -T'. I feel that a crown age is a positive number, but
97 I know you have had a good reason. Perhaps better would be
98 to write something explicit like: $t_{\text{now}} - t_{\text{crown}} = T$. Looking
99 forward for a better suggestion than mine :-)]

$$m = L \cdot \rho \cdot \int_0^T n(t) dt \quad (2)$$

100 where L is the number of DNA nucleotides, ρ is the per-site per-species
101 mutation rate and $n(t)$ the number of species at each time.

102 [GL: This feels kind of a repetitions of what we wrote before the
103 formula. I comment it. We can think of reinsert it afterwards,
104 if needed (see comment above).] [RJCB: I suggest to remove
105 such commented-out lines. Although I sometimes get attached
106 to my sentences, they clog up the document by non-info and
107 I usually delete them anyways in the end. I cannot remember
108 ever regretting this (would I, I could find it in the git history).
109 Would you agree?]

110 Since we cannot know $n_{BD}(t)$ before running simulations we need to re-
111 place it with a proxy. For this reason we will use the average number of
112 species in time according to the BD model. It's well known that this is

113 equal to [GL: insert proper citation] [RJCB: I see you use angle
 114 bracket as a notation for the expected value. I usually see 'E(x)'
 115 as the expected value for 'x', and this is used at the beloved
 116 https://en.wikipedia.org/wiki/Expected_value. What are the
 117 reasons you prefer the notation with the angle brackets?]

$$< n_{BD} > (t) = n_0 \cdot e^{(\mu_{BD} - \lambda_{BD})t} \quad (3)$$

118 where $n_0 = n_{BD}(-T) = n_{MBD}(-T)$ is the initial number of species at
 119 the crown age. From 1, 2 and 3 follows:

$$m_{MBD} = L \cdot \rho \cdot \int_0^T < n_{BD} > (t) dt = L \cdot \rho \cdot n_0 \cdot \left[\frac{e^{(\mu_{BD} - \lambda_{BD})T} - 1}{\mu_{BD} - \lambda_{BD}} \right] \quad (4)$$

120 If we set $\mu_{BD} = \mu_{MBD}$ and reverse this relation we can extrapolate the
 121 value of λ_{BD} to use to generate BD trees.

122 [RJCB: I suggest $n_{BD} = n_{MBD}$ and only change ρ_{BD} to reach
 123 $< m_{MBD} > = < m_{BD} >$] [GL: @Richel: Don't you think it might
 124 make more sense to set $\mu_{BD} = \mu_{MBD}$? What changes in the two
 125 model is the way we use to generate new species, not the way
 126 to remove them. Maybe one thing that's possible to do would
 127 be to make λ_{BD} a function of time, a bit like Tho is doing in
 128 his comparison between DD and TD4 (which, by the way, seem
 129 to yield very different results). In case you are wondering the
 130 theory can be found in Caesar's master thesis.] [RJCB: I fully
 131 agree to use the same extinction rates! The 'mu' used in the
 132 context of mutations (now 'rho') messed me up. I hope this is
 133 clear now. To recap: (1) calculate the speciation rate of the

134 twin tree as you wrote down excellently, (2) simulate a twin tree
135 with same number of taxa, (3) calculate the mutation rates of
136 the trees, so their alignments contain as much information]

137 [GL: My doubt is if we need to use m_{MBD} for the single tree or
138 the same quantity averaged on the full MBD dataset $\langle m_{MBD} \rangle$.
139 Do you think is better to use the individual m_{MBD} for each tree
140 or the average across the whole dataset?] [RJCB: I think a
141 per-tree calculation of the mutation rates is the best we can do.
142 As there was some noise between us above, I think the iteration
143 after the next will allow us to get a better idea about this]

144 • We explained how we set the parameters for each twin BD tree. Using
145 this rules we generate a BD dataset. We repeat the analysis, producing
146 alignments for each tree and subsequently using BEAST to produce a
147 posterior for each of them.

148 • Now we have two datasets of posteriors to compare, one for the BD model
149 and one for the MBD model.

150 • To compare the results for the two models we measure the inference error
151 using the nLTT statistic between known/true tree and posterior/inferred
152 trees. [RJCB: I would love to describe this more concrete. For
153 example, when do we say something has an effect? If we avoid
154 making such judgements, how will we visualize?]

155 [RJCB: I removed the Bayes factor text. It is useful when letting
156 BEAST2 pick more/overly complex models and see if that more
157 complex model fits the data better (penalized by its increased
158 complexity, similar to the AIC). It has its uses, but I am unsure
159 if we already want to discuss this now or first focus on the proper

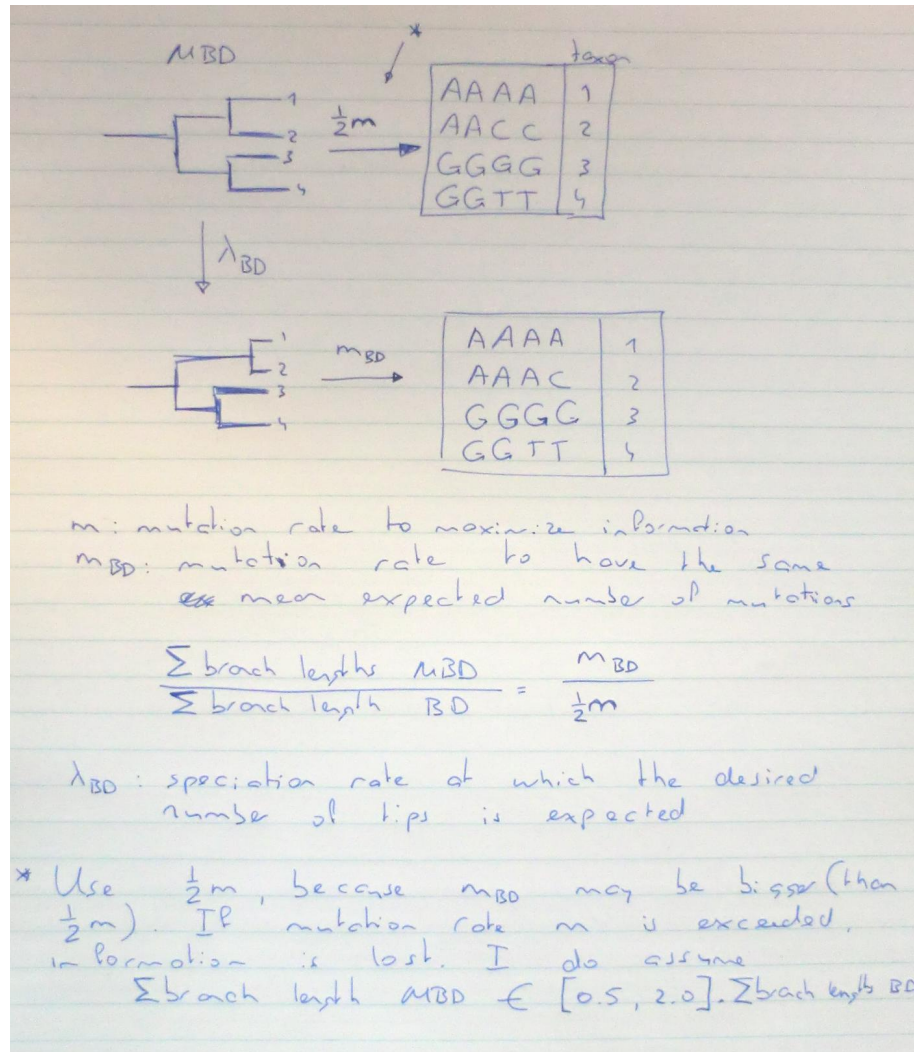


Figure 1: How to create twin trees and alignments. From a focal MBD tree, a twin tree is produced as such: (1) estimate the λ_{BD} to get the same expected number of tips, (2) simulate a BD tree with that amount of tips (discard trees with different number of tips), (3) estimate a mutation rate to get an alignment with the same expected number of mutations, (4) simulate alignments with that amount of mutations (discard those that don't, the picture shows an alignment that should be discarded)

160 tree twinning]

161 3 Results

- 162 • [RJC: I guess you know I am a fan of the Open Science Frame-
163 work, in which you first register your work before you do the
164 experiment (note: I will do some small pilots to estimate the
165 complete time of the experiment). I think it is the proper and
166 superior science, which helps us against writing down bullshit
167 stories after having obtained the results (e.g. 'We expected A
168 and indeed found it!'). It also helps me structure my work: first
169 think deeply about the experiment, then do it (instead of the
170 mixing up the two phases). What are your thoughts on that?]

171 •

172 References

- 173 Bilderbeek, R.J. (2018) *pirouette: create a posterior from a phylogeny*.
- 174 Etienne, R.S., Haegeman, B., Stadler, T., Aze, T., Pearson, P.N., Purvis, A.
175 & Phillimore, A.B. (2012) Diversity-dependence brings molecular phylogenies
176 closer to agreement with the fossil record. *Proc R Soc Lond B: Biol Sci*, **279**,
177 1300–1309.
- 178 Etienne, R.S., Morlon, H. & Lambert, A. (2014) Estimating the duration of
179 speciation from phylogenies. *Evolution*, **68**, 2430–2440.
- 180 Janzen, T., Alzate, A., Muschick, M., Maan, M.E., van der Plas, F. & Etienne,
181 R.S. (2017) Community assembly in lake tanganyika cichlid fish: quantifying

182 the contributions of both niche-based and neutral processes. *Ecology and*
 183 *Evolution*, **7**, 1057–1067.

184 Janzen, T., Alzate, A., Muschick, M., van der Plas, F. & Etienne, R.S. (2016)
 185 Stochastic processes dominate community assembly in cichlid communities in
 186 lake tanganyika.

187 Laudanno, G. (2018) *MBD: Multiple Birth Death Diversification*. R package
 188 version 0.1.

189 Maddison, W.P., Midford, P.E. & Otto, S.P. (2007) Estimating a binary char-
 190 acter’s effect on speciation and extinction. *Systematic Biology*, **56**, 701–710.

191 Valente, L.M., Phillimore, A.B. & Etienne, R.S. (2015) Equilibrium and non-
 192 equilibrium dynamics simultaneously operate in the galápagos islands. *Ecol-*
 193 *ogy Letters*, **18**, 844–852.