# Unifying Phylogenetic Birth-Death Models in Epidemiology and

## Macroevolution

Ailene MacPherson $^{1,2,*}$ , Stilianos Louca $^{3,4}$ , Angela McLaughlin $^{5,6}$ , Jeffrey B. Joy,  $^{5,6,7}$ , and Matthew W. Pennell, $^1$ 

<sup>1</sup> Zoology, University of British Columbia, Vancouver, V6T 1Z4, Canada

<sup>2</sup> Ecology and Evolutionary Biology, University of Toronto, Toronto, M5S 3B2, Canada

<sup>3</sup> Biology, University of Oregon, Eugene, 97403, USA

<sup>4</sup> Institute of Ecology and Evolution, University of Oregon

<sup>5</sup> British Columbia Centre for Excellence in HIV/AIDS, Vancouver, V6H 3X8, Canada

<sup>6</sup> Bioinformatics, University of British Columbia

<sup>7</sup> Medicine, University of British Columbia \*Allene.macpherson@utoronto.ca

#### Abstract

- Birth-death stochastic processes are the foundation of many phylogenetic models and are widely
- used to make inferences about epidemiological and macroevolutionary dynamics. There are a
- large number of birth-death model variants that have been developed; these impose different
- assumptions about the temporal dynamics of the parameters and about the sampling process. As
- each of these variants was individually derived, it has been difficult to understand the
- relationships between them as well as their precise biological and mathematical assumptions.
- Without a common mathematical foundation, deriving new models is non-trivial. Here we unify
- these models into a single framework, prove that many previously developed epidemiological and
- macroevolutionary models are all special cases of a more general model, and illustrate the
- connections between these variants. This unification includes both models where the process is
- the same for all lineages and those in which it varies across types. We also outline a
- straightforward procedure for deriving likelihood functions for arbitrarily complex
- birth-death(-sampling) models that will hopefully allow researchers to explore a wider array of
- scenarios than was previously possible. By re-deriving existing single-type birth-death sampling
- models we clarify and synthesize the range of explicit and implicit assumptions made by these
- 16 models.
- Key words: epidemiology; macroevolution; phylogenetics; birth-death processes; statistical
- 18 inference

21

25

26

27

29

31

40

43

44

46

47

48

49

51

55

57

59

Evolutionary, demographic, and epidemiological processes leave a footprint in the branch length distribution and topology of reconstructed phylogenetic trees. This insight has inspired a huge effort to extract information about these processes by fitting stochastic models. For example, in molecular epidemiology, researchers have leveraged the fact that for many viral pathogens, such as HIV and SARS-CoV-2, accumulate genetic diversity on the timescale of transmission (Drummond et al., 2003; Duffy et al., 2008). This genetic diversity can be used to reconstruct the evolutionary relationships between viral variants sampled from different hosts, which in turn can help elucidate the epidemiological dynamics of a pathogen over time (Grenfell et al., 2004; Volz, 2012). Similarly, phylogenetic trees can provide unique insights into the variation in speciation and extinction rates (Morlon, 2014).

Phylogenetic branching models can be broadly grouped into two classes. The first, based on Kingman's coalescent process (Kingman, 1982), has been widely used to examine changes in the historical population size of pathogens (Pybus et al., 2000; Strimmer and Pybus, 2001; Drummond et al., 2005; Volz et al., 2009). These coalescent methods have also been applied to reconstruct macroevolutionary dynamics (Morlon et al., 2010). Coalescent models are well suited for estimating deterministic population dynamics; however, fitting highly stochastic processes, such as the dynamics of an emerging pathogen, is computationally intensive and in some cases the assumptions of the coalescent may not be appropriate (Stadler et al., 2015; Boskova et al., 2014; Volz and Frost, 2014). The second class of models, which are collectively referred to as birth-death-sampling (BDS) models (Kendall, 1948; Maddison et al., 2007; Stadler, 2009, 2010), is well suited for stochastic scenarios, and are thus becoming an increasingly favorable and popular alternative to coalescent models in epidemiology (Stadler et al., 2012) and have long been the foundation of most macroevolutionary studies — both for inferring speciation and extinction dynamics (Raup, 1985; Nee et al., 1994) and for estimating divergence times (Gernhard, 2008; Heath et al., 2014). As the name implies, the BDS process includes three types of events: birth (pathogen transmission between hosts, or speciation in a macroevolutionary context), death (host death or recovery, or extinction in macroevolution), and sampling (including fossil collection in macroevolution).

In the context of epidemiology, BDS models have the additional property that the model parameters, which can be estimated from viral sequence data, explicitly correspond to parameters in classic structured epidemiological models that are often fit to case surveillance data. If we re-parameterize these models, we can describe the dynamics of the basic and effective reproductive ratios ( $R_0$  and  $R_e$ , respectively) over time (Stadler et al., 2012, 2013) (see Box 1). A common research aim is to describe how the frequency of birth, death, and sampling events, and other derived variables such as  $R_e$ , change throughout the course of an epidemic. There has been less work in macroevolution linking the parameters of a BDS model to those of an underlying more mechanistic model (but see Ezard et al., 2016) but this seems like a promising avenue for future development.

As we detail below and in the Supplementary Material, there has been an astounding rise in the variety and complexity of BDS model variants. A key assumption in the specification of BDS sub-models is whether all lineages alive at some time point are exchangeable (Stadler, 2013) (such models are hereafter "single-type" models), meaning they diversify according to the same process, or if rather the diversification process is variable ("multi-type" models; e.g., Maddison et al., 2007; FitzJohn, 2012; Stadler and Bonhoeffer, 2013; Rasmussen and Stadler, 2019; Barido-Sottani et al., 2018), with lineages belonging to one of multiple possible states each characterized by a unique process. Each of these diversification processes can then be characterized by different dynamical assumptions. In the epidemiological case, these assumptions specify, for example, the nature of viral transmission and the sampling procedure (Stadler et al., 2013; Kühnert et al., 2014; Gavryushkina et al., 2014). While typically not explicitly tied to mechanistic evolutionary processes, there are a similar abundance of dynamical assumptions employed in the macroevolutionary context specifying the nature of biodiversity change through time (Nee, 2006; Gernhard, 2008; Morlon et al., 2011; Stadler, 2011; Morlon, 2014; Heath et al.,

2014; Louca, 2020).

71

81

93

100

101

102

103

104

105

106

108

110

112

114

This flourishing of methods and models has facilitated critical insights into epidemics (du Plessis and Stadler, 2015; Joy et al., 2016) and the origins of contemporary biodiversity (Morlon, 2014; Schluter and Pennell, 2017). However, this diversity of models has made it difficult to trace the connections between variants and to understand the precise epidemiological, evolutionary, and sampling processes that differ between them. Furthermore, despite their apparent similarities, these models have been derived on a case-by-case basis using different notation and techniques; this creates a substantial barrier for researchers working to develop novel models for new situations. And critically, it is imperative that we understand the general properties of BDS phylogenetic models and the limits of inferences from them (Louca and Pennell, 2020a; Louca et al., 2021) and this is difficult to do without considering the full breadth of possible scenarios.

Here we address all of these challenges by unifying the whole class of phylogenetic BDS models. We do so by first deriving a likelihood for general single- and multi-type BDS models; in the general case, we do not assume anything about the functional forms (i.e., temporal dynamics) of the various parameters including the sampling rate through time, the possibility of sampling ancestors (or not), or how the process was conditioned. While such general models may be useful for studying the mathematical properties of BDS models as a whole (Lambert and Stadler, 2013; Louca and Pennell, 2020a; Louca et al., 2021), statistical inference from these models requires researchers to make further constraints on the process. We prove that existing BDS model variants are indeed sub-models of the more general case — and thereby clarify the specific assumptions made by different models — and provide a standardized notation and technique for deriving these and other sub-models that have not previously been considered in the literature.

#### The single-type birth-death-sampling model

*Model Specification:* The BDS stochastic process begins with a single lineage at time T before the present day. We note that this may be considerably older than the age of the most recent common ancestor of an observed sample which is given by  $t_{\rm MRCA}$ . While we focus primarily on applications to epidemiology, our approach is agnostic to whether the rates are interpreted as describing pathogen transmission or macroevolutionary diversification.

In the model, transmission/speciation results in the birth of a lineage and occurs at rate  $\lambda(\tau)$ , where  $\tau$  ( $0 \le \tau \le T$ ) is measured in units of time before the present day ( $\tau = 0$ ), such that  $\lambda$  can be time-dependent. We make the common assumption that lineages in the viral phylogeny coalesce exactly at transmission events, thus ignoring the within-host coalescent processes in the donor (Romero-Severson et al., 2016). Throughout, we will use  $\tau$  as a general time variable and  $t_{\times}$  to denote the time at which a specific event  $\times$  occurs as measured in units of time before the present day (see Table S1). Lineage extinction, resulting from host recovery or death in the epidemiological case or the death of all individuals in a population in the macroevolutionary case, occurs at time-dependent rate  $\mu(\tau)$ . We allow for two distinct types of sampling: lineages are either sampled according to a Poisson process through time  $\psi(\tau)$  or binomially at very short intervals, which we term "concerted sampling attempts" (CSAs), where lineages at some specified time  $t_l$  are sampled with probability  $\rho_l$  ( $\vec{\rho}$  denotes a vector of concerted sampling events at different time points). In macroevolutionary studies based only on extant lineages, there is no Poissonian sampling, but a CSA at the present ( $\rho_0 > 0$ ). In epidemiology, CSAs correspond to large-scale testing efforts (relative to the background rate of testing) in a short amount of time (relative to the rates of viral sequence divergence); for full explanation, see Appendix. We call these attempts rather than events because if  $\rho$  is small or the infection is rare in the population, few or no samples may be obtained. CSAs can also be incorporated into the model by including infinitesimally short spikes in the sampling rate  $\psi$  (more precisely, appropriately scaled Dirac distributions). Hence, for simplicity, in the main text we focus on the seemingly simpler case of pure Poissonian sampling through time except at present-day, where we allow for a CSA to

124

128

130

131

132

133

135

137

143

144

150

152

154

156

160

facilitate comparisons with macroevolutionary models; the resulting formulas can then be used to derive a likelihood formula for the case where past CSAs are included (see Appendix).

In the epidemiological case, sampling may be concurrent (or not) with host treatment or behavioural changes resulting in the effective extinction of the viral lineage. Hence, we assume that sampling results in the immediate extinction of the lineage with probability  $r(\tau)$ . As with the CSAs, this arbitrary time dependence allows for the incorporation of Dirac spikes in any of these variables, for example with mass extinctions ( $\mu$ ) and lagerstätten in the fossil record ( $\psi(1-r)$ ) (Magee and Höhna, 2021). Similarly, in the case of past CSAs we must include the probability,  $r_l$ , that sampled hosts are removed from the infectious pool during the CSA at time  $t_l$ . Poissonian sampling without the removal of lineages ( $r(\tau) < 1$ ) can be employed in the macroevolutionary case to explicitly model the collection of samples from the fossil record (such as the fossilized-birth-death process; Heath et al., 2014).

For our derivation, we make no assumption about the temporal dynamics of  $\lambda$ ,  $\mu$ ,  $\psi$ , or r; each may be constant, or vary according to any arbitrary function of time given that it is biologically valid (non-negative and between 0 and 1 in the case of r). Specifically, the time-varying functions may be any piecewise-continuous functions of time with at most finite number of discontinuities (see 1). Note that these functions need not be differentiable. We make the standard assumption that at any given time any given lineage experiences a birth, death or sampling event independently of (and with the same probabilities as) all other lineages. We revisit this assumption in Box 1 where we discus how the implicit assumptions of the single-type BDS process are well summarized by the diversification model's relationship to the SIR epidemiological model. Our resulting general time-variable BDS process can be fully defined by the parameter set  $\Theta_{\rm BDS} = \{\lambda(\tau), \mu(\tau), \psi(\tau), r(\tau), \bar{\rho}\}$ .

In order to make inference about the model parameters, we need to calculate the likelihood,  $\mathcal{L}$ , that an observed phylogeny,  $\mathcal{T}$ , is the result of a given BDS process. With respect to the BDS process there are two ways to represent the information contained in the phylogeny T, both of which have been used in the literature, which we call the "edge" and "critical time" representations, respectively. We begin by deriving the likelihood in terms of the edge representation and later demonstrate how to reformulate the likelihood in terms of critical times. In the edge representation, the phylogeny is summarized as a set of edges in the mathematical graph that makes up the phylogeny, numbered 1-11 in Figure B1C, and the types of events that occurred at each node. We define  $g_e(\tau)$  as the probability that the edge e which begins at time  $s_e$ and ends at time  $t_e$  gives rise to the subsequently observed phylogeny between time  $\tau$ ,  $(s_e < \tau < t_e)$  and the present day. The likelihood of the model for the observed tree is then, is by definition  $g_{\text{stem}}(T)$ : the probability density that the stem lineage (stem = 1 in Figure B1c) gives rise to the observed phylogeny from the origin, T, to the present day. We find that it is more intuitive to derive the likelihood in terms of the edge representation, as we show below; from this it is straightforward to derive the critical times formulation which results in mathematical simplification. Below we present our five-step technique for the derivation of the tree likelihood. Step 1. Deriving the Initial Value Problem (IVP) for  $g_e(\tau)$ : We derive the IVP for the likelihood density  $g_e(\tau)$  using an approach first developed by Maddison et al. (2007). We begin by deriving the recursion equation for  $q_e$  by considering all the possible events that could occur along edge e between time  $\tau$  and  $\tau + \Delta \tau$  assuming that that  $\Delta \tau$  is small enough such that at most one event is likely to occur.

$$g_{e}(\tau + \Delta \tau) \approx \underbrace{(1 - \lambda(\tau)\Delta\tau)(1 - \mu(\tau)\Delta\tau)(1 - \psi(\tau)\Delta\tau) \times g_{e}(\tau)}_{\text{nothing happens}} + \underbrace{\lambda(\tau)\Delta\tau(1 - \mu(\tau)\Delta\tau)(1 - \psi(\tau)\Delta\tau) \times 2g_{e}(\tau)E(\tau)}_{\text{birth event}} + \underbrace{\mu(\tau)\Delta\tau(1 - \lambda(\tau)\Delta\tau)(1 - \psi(\tau)\Delta\tau) \times 0}_{\text{death event}} + \underbrace{\psi(\tau)\Delta\tau(1 - \lambda(\tau)\Delta\tau)(1 - \mu(\tau)\Delta\tau) \times 0}_{\text{sampling event}} + 0(\Delta\tau^{2}).$$

Here,  $E(\tau)$  is the probability that a lineage alive at time  $\tau$  leaves no sampled descendants at the present day. We will examine this probability in more detail below. Assuming  $\Delta \tau$  is small, we can approximate the above recursion equation as the following difference equation.

$$\Delta g_e(\tau) \approx -(\lambda(\tau) + \mu(\tau) + \psi(\tau))\Delta \tau g_e(\tau) + 2\lambda(\tau)g_e(\tau)E(\tau)\Delta \tau + \mathcal{O}(\Delta \tau^2). \tag{2}$$

By the definition of the derivative we have:

164

168

170

178

179

$$\frac{dg_e(\tau)}{d\tau} = -(\lambda(\tau) + \mu(\tau) + \psi(\tau))g_e(\tau) + 2\lambda(\tau)g_e(\tau)E(\tau). \tag{3}$$

Equation (3) is known as the Kolmogorov backward equation of the BDS process (Feller, 1949; Louca and Pennell, 2020b). Beginning at time  $s_e$ , the initial condition of  $g_e$  depends on which event occurred at the beginning of edge e.

$$g_e(s_e) = \begin{cases} \lambda(s_e)g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ (1 - r(s_e))\psi(s_e)g_{e1}(s_e) & \text{ancestral sampling event} \\ \psi(s_e)r(s_e) + \psi(s_e)(1 - r(s_e))E(s_e) & \text{terminal sampling event} \\ \rho_0 & s_e = 0, \text{ extant sample} \end{cases}$$

$$(4)$$

Together Equations (3) and (4) define the initial value problem for  $g_e(\tau)$  as a function of the probability  $E(\tau)$ .

Because the likelihood density  $g_e$  is the solution to a linear differential equation with initial condition at time  $s_e$ , we can express its solution as follows:

$$g_e(\tau) = \Psi(s_e, \tau)g_e(s_e),\tag{5}$$

186

188

196

where the auxiliary function,  $\Psi$ , is given by:

$$\Psi(s_e, \tau) = \exp\left[\int_{s_e}^{\tau} 2\lambda(x)E(x) - (\lambda(x) + \mu(x) + \psi(x)) dx\right]. \tag{6}$$

This function,  $\Psi(s,t)$ , maps the value of  $g_e$  at time s to its value at t, and hence is known as the probability "flow" of the Kolmogorov backward equation (Louca and Pennell, 2020b).

Step 2. Deriving the IVP for  $E(\tau)$ : We derive the IVP for  $E(\tau)$  in a similar manner as above, beginning with a difference equation.

$$E(\tau + \Delta \tau) = \underbrace{(1 - \lambda(\tau)\Delta\tau)(1 - \mu(\tau)\Delta\tau)(1 - \psi(\tau)\Delta\tau) \times E(\tau)}_{\text{nothing happens}} + \underbrace{\lambda(\tau)\Delta\tau(1 - \mu(\tau)\Delta\tau)(1 - \psi(\tau)\Delta\tau) \times E(\tau)^2}_{\text{birth event}} + \underbrace{\mu(\tau)\Delta\tau(1 - \lambda(\tau)\Delta\tau)(1 - \psi(\tau)\Delta\tau) \times 1}_{\text{death event}} + \underbrace{\psi(\tau)\Delta\tau(1 - \lambda(\tau)\Delta\tau)(1 - \mu(\tau)\Delta\tau) \times 0}_{\text{sampling event}}.$$

$$(7)$$

By the definition of a derivative we have:

$$\frac{dE(\tau)}{d\tau} = -\left(\lambda(\tau) + \mu(\tau) + \psi(\tau)\right)E(\tau) + \lambda(\tau)E(\tau)^2 + \mu(\tau),$$

$$E(0) = 1 - \rho_0,$$
(8)

where  $\rho_0$  is the probability a lineage is sampled at the present day. The initial condition at time 0 is therefore the probability that a lineage alive at the present day is not sampled. Given an analytical or numerical general solution to  $E(\tau)$ , we can find the likelihood by evaluating  $g_{stem}(T)$ , as follows.

Step 3. Deriving the expression for  $g_{stem}(T)$ : Given the linear nature of the differential equation for  $g_e(\tau)$  and hence the representation in Equation (5)), the likelihood  $g_{stem}(\tau)$  is given by the product over all the initial conditions times the product over the probability flow for each edge.

$$g_{stem}(T) = \underbrace{\rho_0^{N_0}}_{\text{extant}} \underbrace{\prod_{i=1}^{I} \lambda(x_i)}_{\text{births}} \underbrace{\prod_{j=1}^{n} \left[ \psi(y_j)(1 - r(y_j)) E(y_j) + \psi(y_j) r(y_j) \right]}_{\text{extinct tips}} \times \underbrace{\prod_{k=1}^{m} \psi(z_k)(1 - r(z_k))}_{\text{ancestral samples}} \underbrace{\prod_{e \in \mathcal{T}} \Psi(s_e, t_e)}_{\text{edges}}. \tag{9}$$

where  $x_i$ ,  $y_j$  and  $z_k$  are the times at which individual birth, terminal sampling and ancestral sampling events occur as we elaborated below.

Step 4. Representing  $g_{stem}(T)$  in terms of critical times: Equation (9) can be further simplified by removing the need to enumerate over all the edges of the phylogeny (the last term of Equation (9)) and writing  $\mathcal{L}$  in terms of the tree's critical times (horizontal lines in Figure B1). The critical times of the tree are made up of three vectors,  $\vec{x}$ ,  $\vec{y}$ , and  $\vec{z}$ , as well as the time of origin T. The vector  $\vec{x}$  gives the time of each birth event in the phylogeny and has length  $I = N_0 + n - 1$  where  $N_0$  is the number of lineages sampled at the present day and n is the number of terminal samples. Unless noted otherwise the elements of vector  $\vec{x}$  are listed in decreasing order, such that  $x_1 > x_2 > ...x_I$  and hence  $x_1$  is the time of the most recent common ancestor  $t_{\text{MRCA}}$ . The vector  $\vec{y}$  gives the timing of each terminal sample and hence has length n whereas vector  $\vec{z}$  gives the timing of each ancestral sample and has length m. With respect to the BDS likelihood then the sampled tree is summarized by  $\mathfrak{T} = \{\vec{x}, \vec{y}, \vec{z}, T\}$ . We note that the critical times only contain the same information as the edges as a result of the assumptions of the BDS process but are not generally equivalent representations of  $\mathfrak{T}$ .

As a result of the linear nature of  $g_e(\tau)$  it is straightforward to rewrite the likelihood in Equation (9) in terms of the critical-time representation of the sampled tree. Defining

$$\Phi(\tau) = \Psi(0, \tau) = \exp\left[\int_0^\tau 2\lambda(x)E(x) - (\lambda(x) + \mu(x) + \psi(x)) dx\right],\tag{10}$$

the probability flow  $\Psi$  can be rewritten as the following ratio:

$$\Psi(s,\tau) = \frac{\Psi(0,\tau)}{\Psi(0,s)} = \frac{\Phi(\tau)}{\Phi(s)}.$$
(11)

This relationship allows us to rewrite the likelihood by expressing the product over the edges as two separate products, one over the start of each edge and the other over the end of each edge which in turn allows us to rearrange and cancel terms to obtain an alternative likelihood expression. Edges begin (value of  $t_e$ ) at either: 1) the tree origin, 2) a birth event resulting to two lineages, or 3) an ancestral sampling event. Edges end (values of  $s_e$ ) at either: 1) a birth event, 2) an ancestral sampling event, 3) a terminal sampling event, or 4) the present day. Hence we have:

$$g_{stem}(T) = \underbrace{\Phi(T)}_{\text{root}} \times \underbrace{\left(\frac{\rho_0}{\Phi(\emptyset)}\right)^{N_0}}_{\text{extant tips}} \times \underbrace{\prod_{i=1}^{I} \lambda(x_i) \frac{\Phi(x_i)^2}{\Phi(x_i)}}_{\text{births}} \times \underbrace{\prod_{j=1}^{n} \frac{\psi(y_j)}{\Phi(y_j)} \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right]}_{\text{extinct tips}} \times \underbrace{\prod_{k=1}^{m} \underbrace{\Phi(z_k)}_{\Phi(z_k)} \psi(z_k) (1 - r(z_k))}_{\text{ancestral samples}}.$$
(12)

Note  $\Phi(0) = 1$ . While Equations (9) and (12) are numerically identical, the critical time expression is more convenient for application as it requires numerically evaluating only a single function  $\Phi(\tau)$  as given by Equation (10).

Step 5. Conditioning the likelihood: While Equation (12) is equal to the basic model likelihood for the phylogeny  $\mathcal{T}$ , it is often appropriate to condition the tree likelihood on the tree exhibiting some property, for example the condition there being at least sampled lineage. Imposing a

235

236

237

239

240

242

244

246

248

250

251

252

253

254

255

257

258

259

260

261

263

265

267

269

270

271

condition on the likelihood is done by multiplying by a factor S. Various conditioning schemes are considered in section 1 and listed in Table S3 with the value of S ranging in complexity from a constant to a general function of the model parameters. The resulting likelihood expression for the general BDS model is:

$$\mathcal{L}(\Theta_{\text{BDS}}, \mathcal{S}|\vec{x}, \vec{y}, \vec{z}, N_0) = \mathcal{S}\rho_0^{N_0} \Phi(T) \prod_{i=1}^{I} \lambda(x_i) \Phi(x_i) \times \prod_{j=1}^{n} \frac{\psi(y_j)}{\Phi(y_j)} \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right] \prod_{k=1}^{m} \psi(z_k) \left( 1 - r(z_k) \right)$$
(13)

Many existing models are special cases of this general BDS model

A large variety of previously published BDS models in epidemiology and macroevolution are special cases of the general model presented here (for a summary of the models we investigated see Table S2; proofs in Supplemental Material). Indeed, we can obtain the likelihood of these models by adding mathematical constraints (i.e., simplifying assumptions) to the terms in Equation (13). Our work thus not only provides a consistent notation for unifying a multitude of seemingly disparate models, it also provides a concrete and numerically straightforward recipe for computing their likelihood functions. We recognize that there are many valid approaches for deriving tree likelihoods for BDS models with share many similarities with our own (e.g., Nee et al., 1994; Maddison et al., 2007; Gernhard, 2008; Morlon et al., 2011; Lambert and Stadler, 2013; Lambert, 2018; Laudanno et al., 2020; Louca and Pennell, 2020b) and do not claim ours is superior to these; however, we have found our technique to be intuitive and flexible. We have implemented the single-type BDS likelihood in the R package castor (Louca and Doebeli, 2018), including routines for maximum-likelihood fitting of BDS models with arbitrary functional forms of the parameters given a phylogeny and routines for simulating phylogenies under the general BDS models (functions fit\_hbds\_model\_on\_grid, fit\_hbds\_model\_parametric and generate\_tree\_hbds).

Figure 1 summarizes the simplifying assumptions that underlie common previously published BDS models; these assumptions generally fall into four categories: 1) assumptions about the functional form of birth, death, and sampling rates over time, 2) assumptions pertaining to the sampling of lineages, 3) the presence of mass-extinction events, and 4) the nature of the tree-conditioning as given by  $\mathcal{S}$ . Here we provide a brief overview of the type of previously-invoked constraints which are consistent (or not) with our unified framework; for full details on each specific case, we refer readers to the Supplementary Material. While we illustrate these constraints within the single-type context, analogous assumptions can be made within the multi-type context examined in the following section.

In regards to rate assumptions, many early BDS models (Stadler, 2009, 2010; Stadler et al., 2012) assumed that the birth, death, and sampling rates remained constant over time. This is mathematically and computationally convenient since an analytical solution can easily be obtained for  $E(\tau)$ . In the epidemiological case, holding  $\lambda$  constant, however, implies that the number of susceptible hosts is effectively constant throughout the epidemic and/or that the population does not change its behavior over time; this is an unrealistic assumption given seasonal changes or changes in response to the disease itself. As such, this assumption is only really valid for small time periods or the early stages of an epidemic. This is useful for estimating the basic reproductive number,  $R_0$ , of the SIR model (Box 1) but not for the effective reproductive number  $R_e$  at later time points (Stadler et al., 2012).

A similarly tractable, but more epidemiologically relevant, model is known as the "birth-death-skyline" variant (Stadler and Bonhoeffer, 2013; Gavryushkina et al., 2014), in which

rates are piecewise-constant functions through time (like the constant rate model, there is also an analytical way to calculate the likelihood of this model; see Appendix). The BDS skyline model has been implemented under a variety of additional assumptions in the Bayesian phylogenetics software BEAST2 Bouckaert et al. (2019). The BDS skyline model has also been extended by Kuhnert et al. (Kühnert et al., 2014) to infer the the parameters of an underlying stochastic SIR model. In this case the diversification model parameters  $\Theta_{BDS}$  are random variables that emerge from stochastic realizations of the epidemiological model given by  $\Theta_{SIR}$ , see Equation (B1). Finally, the birth-death skyline model with piecewise constant rates can also be applied in the macroevolutionary case when no sampling occurs through time,  $\psi(\tau)=0$  (Stadler, 2011).

275

277

280

281

282

283

286

288

290

292

295

297

301

303

305

307

309

310

312

314

316

318

In addition to imposing constraints on the temporal variation in the rates, previously derived sub-models have considered a variety of different assumptions about the nature of the sampling process. Most notably, in macroevolutionary studies, sampling of molecular data typically occurs only at the present day (Stadler, 2009, 2011; Morlon et al., 2011) whereas past Poissonian sampling can be introduced to include the sampling of fossil data (Heath et al., 2014). In epidemiology, concerted sampling at the present day is likely biologically unrealistic (Stadler et al., 2012), though in some implementations of the models, such a sampling scheme has been imposed. These concerted sampling attempts prior to the present day as well as mass extinction events can be incorporated via the inclusion of Dirac distributions in the sampling and death rates, respectively. Finally, previous models often multiply the likelihood by a factor S in order to condition on a particular observation (e.g., observing at least one lineage or exactly  $N_0$  lineages), enumerate indistinguishable trees (e.g., accounting for possible orientations or unlabeled trees) (Gavryushkina et al., 2013, 2014; Stadler, 2009), or to reflect known uncertainties. The "fossilized-birth-death" likelihood derived by Heath et al. (2014) for example, includes a factor that reflects the uncertainty in the attachment and placement of fossils on the macroevolutionary tree. This fossilized-birth-death process has been used to estimate divergence times and to model lineage diversification (Gavryushkina et al., 2017; Landis et al., 2021). Variants of the fossilized-birth-death process, for example including mass extinction events, are feasible and can be derived using our approach. We also note that models similar to the time-variable fossilized-birth-death process have been developed for cases when phylogenetic data is not available (i.e., when only including fossil occurrence data; see Silvestro et al., 2014; Lehtonen et al., 2017); we have not investigated how these models relate to our generalized BDS model but we speculate that it would be possible to also bring these models into a common framework with those that we have discussed. The Supplementary Material demonstrates how these sub-models can be re-derived by either imposing the necessary constraints on the general likelihood formula given in Equation (13) or, alternatively, by starting from the combinations of assumptions and using the five-step procedure outlined above.

#### The multi-type birth-death-sampling model

A common extension of the single-type diversification models explored above is to consider cases where the diversification rates  $(\lambda,\mu,\psi)$  and probabilities  $(r,\rho)$  vary among lineages as a function of a categorical "lineage type". This lineage type can be defined in terms of specific (Maddison et al., 2007; Rasmussen and Stadler, 2019) or unspecified traits (Beaulieu and O'Meara, 2016) or trait combinations (FitzJohn, 2012) (for reviews of these models see Morlon, 2014; Ng and Smith, 2014). Representing these lineage types as colours at nodes and along branches of the tree, we first extend the the single-type model above by deriving the likelihood of a fully coloured tree with topology T where the states along all edges of the phylogeny are known as given by  $\mathcal C$ . The resulting likelihood is an extension of the likelihood first developed by Barido-Sottani et al. (2018), where the diversification rates and probabilities are allowed to vary arbitrarily through time. To illustrate that our derivation is indeed quite general, we follow the model developed (independently) by Magnuson-Ford and Otto (2012) and Goldberg and Igić

(2012), where the state of lineages can change either anagenetically, with a lineage of type a mutating to a type b at rate  $\gamma_{a,b}(\tau)$  or cladogenetically, with a lineage of type a giving rise to a daughter lineage of type b at rate  $\lambda_{a,b}(\tau)$ . Lineages go extinct at a state-dependent rate  $\mu_a(\tau)$  and are sampled at rate  $\psi_a(\tau)$ . As in the single-type model, upon sampling lineages are removed from the population with probability  $r_a(\tau)$  whereas all lineages alive at the present day are sampled with a probability  $\rho_a(\tau)$ . As discussed in depth by Goldberg and Igić (2012), the other discrete variations of state-dependent diversification models (FitzJohn et al., 2009; Goldberg et al., 2011; FitzJohn, 2012) fall out as special cases of this model. (See Ng and Smith, 2014, Caetano et al., 2018, and Louca and Pennell, 2020b for further discussion of the connection between multi-type models.)

We use the five-step technique specified above for the single-type case to derive the probability of observing a given coloured tree under a general multi-type model (see supplementary material). We first derive the initial value problem for the probability  $g_{e,a}(\tau)$  that an edge e of type a in the tree at time  $\tau$  gives rise to the subsequently observed phylogeny. The edge e here refers not to an edge in the topological tree, but to a segment of the tree all of one state between birth, sampling, or mutation events.

$$\frac{dg_{e,a}(\tau)}{d\tau} = -\left(\sum_b \lambda_{a,b}(\tau) + \mu_a(\tau) + \psi_a(\tau) + \sum_b \gamma_{a,b}(\tau)\right)g_{e,a}(\tau) + \sum_b \iota_{a,b}\lambda_{a,b}(\tau)g_{e,a}(\tau)E_b(\tau)$$
 birth event  $a \to a + b$  
$$(1 - r_a(s_e))\psi_a(s_e)g_{e_1,a}(s_e)$$
 ancestral sampling event 
$$r_a(s_e)\psi_a(s_e) + (1 - r_a(s_e))\psi_a(s_e)E_a(s_e)$$
 terminal sampling event 
$$(\gamma_{a,b}(s_e) + \lambda_{a,b}E_a)g_{e_1,b}(s_e)$$
 mutation/hidden birth event  $a \to b$  sampled at present day

Equation (16) distinguishes between multiple types of birth events as pictured in Figure S1. Birth events may be symmetric, with both daughter lineages inheriting the parental type. The exchangeability of the resulting daughter lineages is reflected in the indicator variable  $\iota_{a,b}$  which takes on value of 2 if a=b and 1 otherwise. In contrast asymmetric birth events the resulting daughter lineages differ in type due to caldogenetic change. Importantly the differential equation for  $g_{e,a}$  is linear and hence has a known general solution  $g_{e,a}=g_{e,a}(s_e)\Psi(s_e,\tau)$ . As in the single-type model  $\Psi(s_e,\tau)$  is the probability flow (Louca and Pennell, 2020b) mapping the probability  $g_{e,a}$  from the initial state at time  $s_e$  to the probability at time  $\tau$ .

An analogous initial value problem can be derived for the probability  $E_a(\tau)$ , that a lineage

of type a alive at time au leaves no observed descendants in the sampled tree.

$$\frac{dE_a(\tau)}{d\tau} = -\left(\sum_b \lambda_{a,b}(\tau) + \mu_a(\tau) + \psi_a(\tau) + \sum_b \gamma_{a,b}(\tau)\right) E_a(\tau) 
+ \sum_b \lambda_{a,b}(\tau) E_a(\tau) E_b(\tau) + \mu_a(\tau) + \sum_b \gamma_{a,b}(\tau) E_b(\tau) 
E_a(0) = 1 - \rho_a$$
(15)

This is a non-linear differential equation and must be solved numerically. Given the solution of  $g_{e,a}$  and  $E_a$  the likelihood for the fully coloured tree is characterized by a series of critical times: first,  $\vec{x}_{a,b}$  the times at which a lineage of type a gives birth to a lineage of type a,  $\vec{y}_a$  the ages of tip samples of type a,  $\vec{z}_a$  the ages of ancestral samples of type a, and  $\vec{w}_{a,b}$  the times at which lineages are observed to transition events from type a to type b. The resulting likelihood is given by:

$$\mathcal{L}(\Theta_{\text{MBDS}}|\mathcal{T},\mathcal{C}) = S \times \Phi_{c^*}(T) \times \left[\prod_{a=1}^{A} \rho_a^{N_a}\right] \times \left[\prod_{a=1}^{A} \prod_{b=1}^{A} \prod_{i=1}^{I} \lambda_{a,b}(x_{a,b,i}) \Phi_b(x_{a,b,i})\right] \\
\times \left[\prod_{a=1}^{A} \prod_{j=1}^{J_a} \left[\psi_a(y_{a,j})(1 - r_a(y_{a,j})) E_a(y_{a,j}) + \psi_a(y_{a,j}) r_a(y_{a,j})\right] \frac{1}{\Phi_a(y_{a,j})}\right] \\
\times \left[\prod_{a=1}^{A} \prod_{k=1}^{K_a} \psi_a(z_{a,k})(1 - r_a(z_{a,k}))\right] \\
\times \left[\prod_{a=1}^{A} \prod_{b\neq a} \prod_{l=1}^{L_{a,b}} \left[\gamma_{a,b}(w_{a,b,l}) + \lambda_{a,b}(w_{a,b,l}) E_a(w_{a,b,l})\right] \frac{\Phi_b(w_{a,b,l})}{\Phi_a(w_{a,b,l})}\right]$$
(16)

Here S is an arbitrary form of conditioning as in Equation (13) and  $\Phi_a(\tau) = \Psi_a(\tau, 0)$ , a complete list of notation is given in Table S4.

Equation (16) gives the likelihood of a fully coloured tree, the tree topology plus the state along each branch and at each node in the tree. This likelihood is a generalization of that presented by Barido-Sottani et al. (2018; 2020). Maximizing Equation (16) while incrementally adding and removing changes in state along the branches of the tree can be used to identify clades with distinct diversification parameters. This method can be used, for example, to identify transmission clusters within a disease outbreak (Barido-Sottani et al., 2018). This likelihood is distinct from but related to post-traversal likelihood methods developed to infer state-dependent diversification rates given the known state of sampled lineages (e.g., Maddison et al., 2007; Magnuson-Ford and Otto, 2012; Stadler and Bonhoeffer, 2013). Specifically, these methods give the likelihood  $\mathcal{L}$  ( $\Theta_{\text{MBDS}}|\mathcal{T}, \mathcal{C}_{\bullet}$ ) where  $\mathcal{C}_{\bullet} = \{\mathcal{C}_{\rho}, \mathcal{C}_{y}, \mathcal{C}_{z}\}$  is the state of present-day,  $\mathcal{C}_{\rho}$ , past  $\mathcal{C}_{y}$ , and ancestral,  $\mathcal{C}_{z}$ , sampled lineages. The relationship between the numerically obtained post-traversal likelihood and the closed-form fully coloured likelihood (Equation (16)) is given by:

$$\mathcal{L}\left(\Theta_{\text{MBDS}}|\mathcal{T}, \mathcal{C}_{\bullet}\right) = \frac{\mathcal{L}\left(\Theta_{\text{MBDS}}|\mathcal{T}, \mathcal{C}^{*}\right)}{Pr(\mathcal{C}^{*}|\mathcal{T}, \mathcal{C}_{\bullet}, \Theta_{\text{MBDS}})}.$$
(17)

Here  $\mathbb{C}^*$  is one specific colouring of the tree  $\mathbb{T}$  (e.g., a maximum parsimony ancestral state

 MACPHERSON ET AL.

reconstruction) that is consistent with the observed states. We include Equation (17) as it clarifies the relationship between these two different approaches that have been used to calculate multi-type likelihoods in phylogenetics. Whether or not this is useful for inference is an open question as  $Pr(\mathcal{C}^*|\mathcal{T},\mathcal{C}_{\bullet},\Theta_{\text{MBDS}})$  is challenging to compute (the details of which are beyond the scope of the present paper).

## Concluding remarks

In this paper we have unified a broad class of BDS models that have been widely used both in epidemiology and macroevolution. And in doing so, we have also presented a standardized notation and approach that can be used both for deriving the various sub-models that have previously been studied as well as novel combinations of assumptions about the model parameters. The unification of these models clarifies the connections between BDS variants, facilitates the development of new variants tailored to specific scenarios, and provides a structure for understanding how results depend on model assumptions (Kirkpatrick et al., 2002; Lafferty et al., 2015; Louca and Pennell, 2020a). And importantly, given the recent discovery of widespread non-identifiability in birth-death processes fit to extant-only (Louca and Pennell, 2020a) and serially-sampled (Louca et al., 2021) phylogenetic data, there is a critical need to explore a much broader range of BDS models than were previously considered and the mathematical generalization presented here will be enable this.

## Box 1: The connection between BDS and SIR models

The single-type BDS model is intimately related to the SIR compartmental model used in classic theoretical epidemiology. This connection illustrates the explicit and implicit assumptions of the general BDS model and its sub models. Here we define the SIR epidemiological model, discuss how it can inform and be informed by these diversification models, and examine the shared assumptions of the two frameworks.

#### The SIR model:

The SIR model partitions the host population via infection status into susceptible (S), infected (I), and recovered (R) hosts. Infection of susceptible hosts occurs at a per-capita rate  $\beta I$ . Infected hosts may recover (at rate  $\gamma$ ), die of virulent cases (at rate  $\alpha$ ), or be sampled (at rate  $\psi$ ). The cumulative number of sampled hosts is represented in the SIR model (Figure B1 top) by  $I^*$ . Upon sampling, infected hosts may be treated and hence effectively recover with probability r. Hosts that have recovered from infection exhibit temporary immunity to future infection which wanes at rate  $\sigma$ . The special case of the SIR model with no immunity (the SIS model) is obtained in the limit as  $\sigma \to \infty$ . In addition to these epidemiological processes, the SIR model includes demographic processes, such as host birth (rate B) and death from natural causes (rate  $\delta$ ). While not shown explicitly in the figure, these epidemiological and demographic rates may change over time as a result of host behavioural change, pharmaceutical and non-pharmaceutical interventions, or host/pathogen evolution.

#### The BDS Model:

The BDS model follows the number of sampled and unsampled viral lineages over time, analogous to the I and  $I^*$  classes of the SIR model. A key element of general BDS model is that birth and death rates may vary over time. This time dependence may be either continuous (Morlon et al., 2011; Rabosky and Lovette, 2008b) or discrete (Stadler, 2011; Stadler and Bonhoeffer, 2013; Gavryushkina et al., 2014; Kühnert et al., 2014) Although arbitrarily time-dependent, the birth, death, and sampling rates in the general BDS model are assumed to be diversity-independent, analogous to the assumption of density-dependent transmission (pseudo mass action) in the SIR model (Keeling and Rohani, 2008). Incorporating such diversity dependence into macroevolutionary models has been shown to increase the accuracy of extinction rate estimates and are necessary to accurately capture the saturation of diversity (Etienne et al., 2012). While some forms of diversity-dependence in diversification rates may be incorporated implicitly capturing deterministic diversity dependence as time dependence (Rabosky and Lovette, 2008a), stochastic diversity-dependence (Etienne and Rosindell, 2012) goes beyond the scope of the BDS models considered here.

The single-type BDS model assumes all viral lineages are exchangeable - this has several implications. First, all viral lineages are epidemiologically identical hence all mutations between them are neutral. Incorporating non-neutral genetic variation requires a multi-type approach as in Equation (16). Second, transmission is independent of lineage age. In the macroevolutionary case, such age-dependence has been suggested to reflect niche differentiation in novel species (Hagen et al., 2015) and in the epidemiological case may reflect adaptation towards increased transmissibility following a host species-jumping event. Third, lineage exchangeablity is reflected in the absence of an exposed (E) class in the SIR model in which hosts can, for example, transmit infections but not be sampled or vice versa. Finally, the single-type BDS model assumes all lineages are sampled at random and does not include sub-models with non-random representation of lineages (Stadler et al., 2012).

#### Model Connections

Given their shared model assumptions, the single-type BDS model can be constrained explicitly to reflect an underlying SIR epidemic by setting the viral birth rate equal to the per-capita transmission rate of the infectious class  $\lambda(\tau) = \beta S(\tau)$  and the viral death rate to the infectious recovery or removal rate  $\mu(\tau) = \gamma + \delta + \alpha$ , whereas the sampling rate  $\psi(\tau)$  is identical across models (Figure B1a). While constraining the birth, death, and sampling rates in this manner can be used to parameterize compartmental models (Kühnert et al., 2014) doing so is an approximation assuming independence between the exact timing of transmission, recovery or removal from population, and sampling events in the SIR model and birth, death, and sampling events in the diversification model. The resulting tree likelihood in terms of the compartmental model is given by:

$$Pr(\mathcal{T}|\Theta_{SIR}) = \underbrace{Pr(\mathcal{T}|\Theta_{BDS})}_{\text{BDS likelihood}} \underbrace{P(\Theta_{BDS}|\Theta_{SIR})}_{\text{SIR process}}.$$
 (B1)

While they are not sub-models of the general BDS process, likelihood models have been developed that capture the full non-independence of viral diversification and epidemiological dynamics for the SIR model specifically (Leventhal et al., 2012) and in compartmental models in general (Vaughan et al., 2019). The connection between the BDS process and SIR epidemiological models can also be used after the diversification rates are inferred to estimate the basic and effective reproductive rates (Stadler et al., 2012; Stadler and Bonhoeffer, 2013). Specifically, the effective reproductive rate at time  $\tau$  before the present day is given by  $R_e(\tau) = \frac{\lambda(\tau)}{\mu(\tau) + r(\tau)\psi(\tau)}$ . Although the SIR model is a useful epidemiological model for is simplicity, realistically modelling epidemic dynamics requires far more complex compartmental models. As reflected by their shared structure, the application of the single-type BDS model is restricted, however, to the assumptions of the SIR model alone and further methodological advances in multi-type modelling are necessary for direct inference for the larger class of epidemiological models.

#### ACKNOWLEDGMENTS

We would like to thank Sally Otto for her thoughtful comments on this work. This work was supported by a *Grant for Catalyzing Research Clusters* awarded to the UBC Biodiversity Research Centre, NSF DEB Grant #2028986 award to SL and MWP. MWP was supported by a NSERC Discovery Grant and AM was supported in part by the EEB department Postdoctoral Fellowship from the University of Toronto. AMcL was supported in part by a Canadian Institutes of Health Research (CIHR) doctoral award (#6557). JBJ is supported by a Genome Canada Bioinformatics and Computational Biology grant (287PHY), CIHR coronavirus rapid response program grant (440371) and is grateful to the British Columbia Centre for Excellence in HIV/AIDS for additional funding support.

#### APPENDIX: ADDING ASSUMPTIONS TO THE GENERAL MODEL

In this appendix, we demonstrate how one can obtain the likelihood of sub-models with different sets of assumptions by applying constraints to the general likelihood. There are four classes of assumptions that are commonly applied in epidemiological and macroevolutionary studies. First, researchers can make assumptions about the functional form of the birth, death, and sampling rates. Here we address two such unique assumptions: i) Birth, death, and sampling rates are constant (see 1, , , ); and ii) birth, death, and sampling rates are piecewise-constant functions of time (see 1 and ). The cases where birth, death, and sampling rates are defined by a stochastic or deterministic SIR model are mathematically analogous to the cases of the piecewise-constant and general time-variable models respectively. All additional constraints imposed will depend on the exact compartmental model used and hence we will not discuss them in detail in this section. The second major class of assumptions pertains to sampling. There are four such sampling assumptions: i) sampling happens only at the present day as in a birth-death model (see 1 and , , ) or as implemented in the "Birth Death Skyline Contemporary" prior in the BDSKY package in BEAST2; ii) the absence of concerted present-day sampling (see 1 and ); iii) the inclusion of ancestral samples with sampled descendants (, and ); and iv) concerted sampling attempts (CSA) during which all lineages are sampled with a given probability (see 1 and ). The third assumption class considers the presence of mass extinction events (see 1 and ). The fourth and final major class of assumptions deal with the conditioning of the likelihood. The various conditioning schemes are explored in below and summarized in Table S3.

#### RATE ASSUMPTIONS

#### Constant rates

- Model Assumptions: Constant diversification rates:  $\lambda(t) = \lambda$ ,  $\mu(t) = \mu$ ,  $\psi(t) = \psi$ , and constant removal probability r(t) = r.
- The IVP for  $g_e(\tau)$ :

470

471

472

473

474

475

477

479

481

483

487

488

490

491

492

493

498

$$\frac{dg_e(\tau)}{d\tau} = -\left(\lambda + \mu + \psi\right)g_e(\tau) + 2\lambda g_e(\tau)E(\tau)$$
 birth event giving rise to edges e1 and e2 
$$g_e(s_e) = \begin{cases} \lambda g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ \psi(1-r)g_{e1}(s_e) & \text{ancestral sampling event} \\ \psi r + \psi(1-r)E(s_e) & \text{terminal sampling event} \\ \rho_0 & s_e = 0, \text{edge sampled at present day} \end{cases}$$

• The IVP for  $E(\tau)$ :

$$\frac{dE(\tau)}{d\tau} = -(\lambda + \mu + \psi)E(\tau) + \lambda E(\tau)^2 + \mu \quad E(0) = 1 - \rho_0.$$

In this case the IVP for  $E(\tau)$  is a Bernoulli differential equation and has a known analytical

solution. As given by Equation 1 in Stadler (2010) this solution is given by:

$$E(\tau) = \frac{\lambda + \mu + \psi}{2\lambda} + \frac{c_1}{2\lambda} \frac{e^{-c_1 t} (1 - c_2) - (1 + c_2)}{e^{-c_1 t} (1 - c_2) + (1 + c_2)}$$

$$c_1 = \left| \sqrt{(\lambda - \mu - \psi)^2 + 4\lambda \psi} \right| \quad c_2 = -\frac{\lambda - \mu - 2\lambda \rho_0 - \psi}{c_1}.$$
(A1)

• *The Probability Flow:* 

$$\Phi(\tau) = \exp\left[\int_0^{\tau} 2\lambda E(x) - (\lambda + \mu + \psi) dx\right].$$

• The Likelihood:

$$\mathcal{L}_{C} = \mathcal{S}\rho_{0}^{N_{0}}\Phi(T)\lambda^{I}\psi^{n+m}(1-r)^{m}\prod_{i=1}^{I}\Phi(x_{i})\prod_{j=1}^{n}\frac{1}{\Phi(y_{j})}\left[(1-r)E(y_{j})+r\right]$$
(A2)

#### Piecewise-constant rates

• Model assumptions: Divide time into L+1 intervals defined by transition times  $0=t_0 < t_1 < t_2 < ... < t_L < t_{L+1} = T$ . Define rates and removal probabilities constant within a given interval.

$$\lambda(\tau) = \lambda_l \quad t_l < \tau \leqslant t_{l+1}$$

$$\mu(\tau) = \mu_l \quad t_l < \tau \leqslant t_{l+1}$$

$$\psi(\tau) = \psi_l \quad t_l < \tau \leqslant t_{l+1}$$

$$r(\tau) = r_l \quad t_l < \tau < t_{l+1}$$

- The IVP and Solution for  $g_e(\tau)$ : Given the definitions of  $\lambda(\tau)$ ,  $\mu(\tau)$ ,  $\psi(\tau)$ , and  $r(\tau)$  within each time interval the IVP for  $g_e(\tau)$  is identical to that given in Equations (3) and (4). If  $g_{l,e}(\tau)$  is the probability density within time interval l than  $g_{l,e}(t_l) = g_{l-1,e}(t_l)$ .
- The IVP and Solution for  $E(\tau)$ : As with  $g_e(\tau)$ , the IVP for  $E(\tau)$  is given by Equation (8). With the piecewise-constant rate assumptions, however, the general solution for  $E(\tau)$  between  $t_l < \tau \leqslant t_{l+1}$  is known (similar to Equation (A1)). Defining  $E_l(\tau) = E(\tau)$  where  $t_l < \tau \leqslant t_{l+1}$  and  $E_l(t_l) = E_{l-1}(t_l)$  we have:

$$E_{l}(\tau) = \frac{\lambda_{l} + \mu_{l} + \psi_{l}}{2\lambda_{l}} + \frac{c_{1}}{2\lambda_{l}} \frac{e^{-c_{1}t}(1 - c_{2}) - (1 + c_{2})}{e^{-c_{1}t}(1 - c_{2}) + (1 + c_{2})}$$

$$c_{1} = \left| \sqrt{(\lambda_{l} - \mu_{l} - \psi_{l})^{2} + 4\lambda_{l}\psi_{l}} \right| \quad c_{2} = -\frac{\lambda_{l} - \mu_{l} - 2\lambda_{l}(1 - E_{l}(t_{l})) - \psi_{l}}{c_{1}},$$

where  $E_l(t_l) = E_{l-1}(t_l)$  for l > 0 and  $E_0(t_0) = 1 - \rho_0$ .

• The Probability Flow: We define a probability sub-flow within each time interval. Specifically, in the  $l^{th}$  time interval.

$$\Phi_l(\tau) = \exp\left[\int_{t_l}^{\tau} 2\lambda_l E_l(x) - (\lambda_l + \mu_l + \psi_l) dx\right].$$

The complete flow can be expressed as a function of the sub-flows in the following manner:

$$\Phi(\tau) = \Phi_{L_{\tau}}(\tau) \prod_{l=1}^{L_{\tau}} \Phi_{l-1}(t_l) 
\Phi(t_k) = \underbrace{\Phi_k(t_k)}_{1} \prod_{l=1}^{k} \Phi_{l-1}(t_l) = \prod_{l=1}^{k} \Phi_{l-1}(t_l),$$
(A3)

where  $L_t$  is the index of the time  $t_l$  at or after time t, i.e. the largest index such that  $t_l \leq \tau$ .

• *The Likelihood:* Given these piecewise definitions we substitute them into the general BDS likelihood (13).

$$\begin{split} \mathcal{L}_{PC} = & \underbrace{8 \underbrace{\rho_0^{N_0}}_{\text{extant}} \underbrace{\Phi_L(T) \prod_{l=1}^L \Phi_{l-1}(t_l)}_{\text{root}} \times \underbrace{\prod_{i=1}^I \left[ \lambda_{L_{x_i}} \Phi_{L_{x_i}}(x_i) \prod_{l=1}^{L_{x_i}} \Phi_{l-1}(t_l) \right]}_{\text{births}} \\ & \times \underbrace{\prod_{j=1}^n \frac{\psi_{L_{y_j}} \left( (1-r_{L_{y_j}}) E_{L_{y_j}}(y_j) + r_{L_{y_j}} \right)}_{\Phi_{L_{y_j}} \prod_{l=1}^{L_{y_j}} \Phi_{l-1}(t_l)}_{\text{extinct}} \times \underbrace{\prod_{i=1}^m \psi_{L_{z_k}} \left( 1-r_{L_{z_k}} \right),}_{\text{ancestral samples}}, \end{split}$$

where we use PC to denote the piecewise-constant assumption.

We can simplify several of these products. Let  $\alpha_l$  be the number of birth events  $\geqslant t_l$  and  $\sigma_l$  the number of sampling events  $\geqslant t_l$ .

$$\prod_{i=1}^{I} \prod_{l=1}^{L_{x_i}} \Phi_{l-1}(t_l) = \prod_{l=1}^{L} \left[ \Phi_{l-1}(t_l) \right]^{\alpha_l} \\
\prod_{j=1}^{n} \prod_{l=1}^{L_{y_j}} \frac{1}{\Phi_{l-1}(t_l)} = \prod_{l=1}^{L} \left[ \Phi_{l-1}(t_l) \right]^{-\sigma_l}.$$
(A4)

Let  $n_l$  be the number of observed lineages alive at time  $t_l$ . Because the number of observed lineages increases with each birth and decreases with each sampled tip, counting the root we have  $n_l = \alpha_l - \sigma_l + 1$ . Substituting the expressions for the into the likelihood and using

the definition of  $n_l$  we have:

$$\mathcal{L}_{PC} = \mathcal{S}\rho_0^{N_0} \Phi_L(T) \prod_{i=1}^{I} \lambda_{L_{x_i}} \Phi_{L_{x_i}}(x_i)$$

$$\times \prod_{j=1}^{n} \frac{\psi_{L_{y_j}}}{\Phi_{L_{y_j}}} \left( (1 - r_{L_{y_j}}) E_{L_{y_j}}(y_j) + r_{L_{y_j}} \right) \prod_{k=1}^{m} \psi_{L_{z_k}} \left( 1 - r_{L_{z_k}} \right)$$

$$\times \prod_{l=1}^{L} \Phi_{l-1}(t_l)^{n_l}.$$
(A5)

#### SAMPLING ASSUMPTIONS

#### Birth-death models

- Model Assumptions: The birth-death model assumes that  $\psi(\tau) = 0$ . Note that the probability of sampling a lineage given it is alive at the present day remains as  $\rho_0$  (incomplete sampling).
- IVP for  $q_e(\tau)$ :

$$\begin{split} \frac{dg_e(\tau)}{d\tau} &= -\left(\lambda(\tau) + \mu(\tau)\right)g_e(\tau) + 2\lambda(\tau)g_e(\tau)E(\tau) \\ g_e(s_e) &= \begin{cases} \lambda(s_e)g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ \rho_0 & s_e = 0, \text{edge sampled at present day} \end{cases} \end{split}$$

• IVP for  $E(\tau)$ :

$$\frac{dE(\tau)}{d\tau} = -(\lambda(\tau) + \mu(\tau))E(\tau) + \lambda(\tau)E(\tau)^2 + \mu(\tau) \quad E(0) = 1 - \rho_0.$$

Note in this case  $E(\tau)$  equals  $E(\tau)$ , the probability a lineage leaves no sampled extant descendants. As demonstrated by Morlon et al. (2011) there exists a general solution to this initial value problem, see section for more details. This general solution is given by:

$$E(\tau) = 1 - \frac{\rho_0 \exp\left[\int_0^{\tau} (\lambda(u) - \mu(u)) \, du\right]}{1 + \int_0^{\tau} \rho_0 \exp\left[\int_0^{x} (\lambda(u) - \mu(u)) \, du\right] \, dx}.$$

• *The Probability Flow:* From (Morlon et al., 2011), the probability flow can be written as the following:

$$\Phi(\tau) = \exp\left[\int_0^{\tau} (\lambda(\sigma) - \mu(\sigma)) d\sigma\right] \left[1 + \frac{\int_0^{\tau} \rho_0 \lambda(u) \exp\left[\int_0^u (\lambda(\sigma) - \mu(\sigma)) d\sigma\right] du}{1 + \rho_0}\right]^{-2}.$$

• The Likelihood:

$$\mathcal{L}_{BD} = \mathcal{S}\rho_0^{N_0}\Phi(T)\prod_{i=1}^I \lambda(x_i)\Phi(x_i)$$
(A6)

No sampling at the present

Here we consider the case when  $\rho_0 = 0$ . The likelihood follows exactly as in the general model case. The resulting likelihood expression is given by:

$$\mathcal{L}_{\rho_0=0} = \$\Phi(T) \prod_{i=1}^{I} \lambda(x_i) \Phi(x_i) \times \prod_{j=1}^{n} \frac{\psi(y_j)}{\Phi(y_j)} \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right] \prod_{k=1}^{m} \psi(z_k) \left( 1 - r(z_k) \right).$$
(A7)

Note that in this case I = n - 1.

## Concerted sampling attempts

• Model Assumptions: Here we introduce L concerted sampling attempts (CSA) at known points in time,  $t_l$   $l \in \{1, 2, ...L\}$ . Like the CSA at the present day, and in contrast to the background Poissonian sampling rate, during the CSA at time  $t_l$  every lineages is sampled with a fixed probability  $\rho_l$ . In the derivation of the likelihood below, we must distinguish between three different sampling event types. First, past Poissonian sampling events are those that do not occur during CSAs. Second, past concerted sampling events are those that occur during a CSA at time  $t_l$   $l \in \{1, 2, ..., L\}$ . Finally, present concerted sampling events are those that occur at the present day  $\tau = 0$ . Past concerted sampling attempts can be included in the general model above by adding L Dirac distributions to the Poisson sampling rate function. Namely,

$$\psi(\tau) = \bar{\psi}(\tau) + \sum_{l=1}^{L} [w_l * \delta(\tau - t_l)],$$
 (A8)

where  $\bar{\psi}(\tau)$  is the background Poissonian sampling rate and  $w_l = -ln(1 - \rho_l)$ . The definition of  $w_l$  comes from solving the CDF of the exponentially distribution for the 'effective sampling rate' such that the probability of a lineage being sampled is  $\rho_l$ .

• IVP for  $g_e(\tau)$ :

$$\frac{dg_e(\tau)}{d\tau} = -\left(\lambda(\tau) + \mu(\tau) + \psi(\tau)\right)g_e(\tau) + 2\lambda(\tau)g_e(\tau)E(\tau)$$

$$g_e(s_e) = \begin{cases} \lambda(s_e)g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ (1 - r(s_e))\bar{\psi}(s_e)g_{e1}(s_e) & \text{Poisson ancestral sampling event} \\ \bar{\psi}(s_e)r(s_e) + \bar{\psi}(s_e)(1 - r(s_e))E(s_e) & \text{Poisson terminal sampling event} \\ (1 - r(t_l))\rho_lg_{e1}(t_l) & \text{ancestral sample at } t_l \\ \rho_lr(t_l) + \rho_l(1 - r(t_l))E(t_l) & \text{terminal sample at } t_l \\ \rho_0 & s_e = 0, \text{edge sampled at present day} \end{cases}$$

The solution to  $g_e(\tau)$  is given by Equations (5) and (6).

• IVP for  $E(\tau)$ : As with  $g_e(\tau)$ , the IVP for  $E(\tau)$  is identical to that given for the general model in Equation (8). Except in rare cases the IVP must be solved numerically hence requiring numerical integration over Dirac distributions which can prove to be problematic.

Note however, that when examining the integrals over the CSAs, a priori, it is a matter of convention whether the Dirac distribution should be considered as "integrated over" when located at the upper integration bound  $\int_a^b \delta(s-b)ds=1$  or at the lower integration bound  $\int_a^b \delta(s-a)ds=1$ . Whichever convention we chose, we must rigorously obey it so that the ratio  $\Phi(t)/\Phi(s)$  correctly evaluates to  $\Psi(s,t)$  whenever  $s\leqslant t$ . Using the former convention, we can rewrite the probability  $E(t_l)$  at each concerted sampling time  $t_l$  as:

$$E(t_l) = E(t_l^-)e^{w_l} = E(t_l^-)(1 - \rho_l),$$

where  $t_l^-$  denotes the limit as time approaches  $t_l$  from below. Hence the probability  $E(\tau)$  at any time  $\tau$  can be evaluated numerically by considering the dynamics between successive CSAs and at each CSA separately.

• *The Probability Flow:* The probability flow is given by:

$$\Phi(\tau) = \exp\left[\int_0^\tau 2\lambda(x)E(x) - \left(\lambda(x) + \mu(x) + \bar{\psi}(x) + \sum_{l=1}^L w_l \delta(x - t_l)\right) dx\right].$$

As with  $E(\tau)$  integration over the dirac distributions can be problematic and hence we rewrite this expression separating out these terms. Let  $L_{\tau}$  be the oldest CSA occurring at or

after time  $\tau$ , i.e. the largest index for which  $t_l \leqslant \tau$ .

$$\Phi(\tau) = \exp\left[\int_0^{\tau} 2\lambda(x)E(x) - \left(\lambda(x) + \mu(x) + \bar{\psi}(x) + \sum_{l=1}^{L_{\tau}} w_l \delta(x - t_l)\right) dx\right] 
= \exp\left[\int_0^{\tau} 2\lambda(x)E(x) - \left(\lambda(x) + \mu(x) + \bar{\psi}(x)\right) dx\right] \prod_{l=1}^{L_{\tau}} e^{-w_l} 
= \exp\left[\int_0^{\tau} 2\lambda(x)E(x) - \left(\lambda(x) + \mu(x) + \bar{\psi}(x)\right) dx\right] \prod_{l=1}^{L_{\tau}} (1 - \rho_l).$$
(A9)

We define:

595

596

597

599

600

602

603

605

$$\bar{\Phi}(\tau) = \exp\left[\int_0^\tau 2\lambda(x)E(x) - \left(\lambda(x) + \mu(x) + \bar{\psi}(x)\right)dx\right],\tag{A10}$$

which means that we can rewrite Equation (A9) as:

$$\Phi(\tau) = \bar{\Phi}(\tau) \prod_{l=1}^{L_{\tau}} (1 - \rho_l). \tag{A11}$$

• The Likelihood: The edge representation of  $g_{stem}$  is given by:

$$g_{stem}(T) = \rho_0^{N_0} \prod_{i=1}^{L} \lambda(x_i) \prod_{j=1}^{n} \psi(y_j) \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right] \prod_{k=1}^{m} \psi(z_j) (1 - r(y_j))$$

$$\times \prod_{l=1}^{L} \rho_l \left[ (1 - r_l) E(t_l) + r_l \right]^{N_l} \prod_{l=1}^{L} \left[ \rho_l (1 - r_l) \right]^{M_l} \prod_{edges} \Psi(s_e, t_e).$$

The critical time representation of  $g_{stem}$  is given by:

$$g_{stem}(T) = \underbrace{\rho_0^{N_0}}_{\text{extant tips}} \underbrace{\bar{\Phi}(T) \prod_{l=1}^{L} (1-\rho_l) \prod_{i=1}^{I} \lambda(x_i) \bar{\Phi}(x_i) \left[ \prod_{l=1}^{L_{x_i}} (1-\rho_l) \right]}_{\text{births}}$$

$$\times \underbrace{\prod_{j=1}^{n} \frac{\psi(y_j)}{\bar{\Phi}(y_j)} \left[ (1-r(y_j)) E(y_j) + r(y_j) \right] \left[ \prod_{l=1}^{L_{y_j}} \frac{1}{1-\rho_l} \right] \prod_{k=1}^{m} \psi(z_k) (1-r(z_k))}_{\text{Pois. extinct tips}}$$

$$\times \underbrace{\prod_{l=1}^{L} \left( \frac{\rho_l}{\bar{\Phi}(t_l)} \left[ (1-r_l) E(t_l) + r_l \right] \right)^{N_l} \left[ \prod_{j=1}^{l} \frac{1}{(1-\rho_j)^{N_l}} \right] \prod_{l=1}^{L} \left[ \rho_l (1-r_l) \right]^{M_l},}_{\text{CSA extinct tips}}$$

$$\underbrace{\text{CSA extinct tips}}_{\text{CSA ances. samples}}$$

where  $N_l$  is the number of tip samples (samples without descendants) obtained during the

 $l^{th}$  CSA and  $M_l$  is the number of ancestral samples (sequences with descendants). By changing how we enumerate birth, death, and sampling events we can greatly simplify this likelihood. First, let  $\alpha_l$  be the number of branching events at or before the the  $l^{th}$  CSA. In other words,  $\alpha_l$  is the number of branching events if the tree were trimmed at the  $l^{th}$  CSA. Then:

$$\prod_{i} \left[ \prod_{l=1}^{L_{x_i}} (1 - \rho_l) \right] = \prod_{l=1}^{L} (1 - \rho_l)^{\alpha_l}.$$
 (A12)

Second, let  $\sigma_l$  be the number of past Poissonian sampling events before time  $t_l$ . Then:

$$\prod_{j}^{n} \left[ \prod_{l=1}^{L_{y_{j}}} \frac{1}{(1-\rho_{l})} \right] = \prod_{l=1}^{L} \frac{1}{(1-\rho_{l})^{\sigma_{l}}}.$$
 (A13)

Finally, let  $\beta_l$  be the number of past lineages sampled during a CSA at or before the CSA at time  $t_l$ . Hence,  $\beta_l = N_l + N_{l+1} + ... + N_L$ . Then:

$$\prod_{l=1}^{L} \left[ \prod_{j=1}^{l} \frac{1}{(1-\rho_j)^{N_l}} \right] = \prod_{l=1}^{L} \frac{1}{(1-\rho_l)^{\beta_l}}.$$
 (A14)

The likelihood hence simplifies to:

$$g_{stem}(T) = \rho_0^{N_0} \bar{\Phi}(T) \prod_{l=1}^{L} (1 - \rho_l)^{\alpha_l - \beta_l - \sigma_l + 1} \prod_{i=1}^{I} \lambda(x_i) \bar{\Phi}(x_i)$$

$$\times \prod_{j=1}^{n} \frac{\psi(y_j)}{\bar{\Phi}(y_j)} \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right] \prod_{k=1}^{m} \psi(z_k) (1 - r(z_k))$$

$$\times \prod_{l=1}^{L} \left( \frac{\rho_l}{\bar{\Phi}(t_l)} \left[ (1 - r_l) E(t_l) + r_l \right] \right)^{N_l} \prod_{l=1}^{L} \left[ \rho_l (1 - r_l) \right]^{M_l}.$$

Let  $n_l$  be the number of lineages that cross  $t_l$ , i.e., the number of lineages alive at time  $t_l$  with sampled descendants at some younger age. Note that by this definition  $n_0 = 0$ . Then  $b_l + \sigma_l + n_l$  is the number of tips in the tree had it been trimmed at age  $t_l$  whereas  $\alpha_l$  is the number of branching events. Therefore we must have  $\alpha_l = b_l + \sigma_l + n_l - 1$ . This allows us to simplify the conditioned likelihood given below.

$$\mathcal{L}_{CSA} = \mathcal{S}\rho_0^{N_0} \bar{\Phi}(T) \prod_{l=1}^{L} (1 - \rho_l)^{n_l} \prod_{i=1}^{I} \lambda(x_i) \bar{\Phi}(x_i)$$

$$\times \prod_{j=1}^{n} \frac{\psi(y_j)}{\bar{\Phi}(y_j)} \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right] \prod_{k=1}^{m} \psi(z_k) (1 - r(z_k))$$

$$\times \prod_{l=1}^{L} \left( \frac{\rho_l}{\bar{\Phi}(t_l)} \left[ (1 - r_l) E(t_l) + r_l \right] \right)^{N_l} \prod_{l=1}^{L} \left[ \rho_l (1 - r_l) \right]^{M_l}$$
(A15)

#### MASS EXTINCTION

• Model Assumptions: In addition to the Poisson birth death and sampling events considered in the general model, there are L mass extinctions events occurring at times  $t_1 > t_2 > ..t_L$ . During the  $l^{th}$  mass extinction event each lineage goes extinct with probability  $\nu_l$ . As with concerted sampling such mass extinction events can be introduced into the model by adding a set of dirac-delta functions to the Poisson death rate,  $\bar{\mu}(\tau)$ .

$$\mu(\tau) = \bar{\mu}(\tau) + \sum_{l=1}^{L} m_l \delta(\tau - t_l), \tag{A16}$$

where  $m_l = -ln(1 - \nu_l)$ .

- *IVP for*  $g_e(\tau)$ : The initial value problem for  $g_e(\tau)$  is identical to that given in equation by Equations (3) and (4) except that  $\mu$  is now includes the mass extinction events.
- IVP for  $E(\tau)$ : The IVP for  $E(\tau)$  to that given by Equation (8) except where the extinction rate is given by Equation (A16). The solution to  $E(\tau)$  is obtained by numerical integration. Given the dirac-delta functions this numerical integration can be carried out in a piecewise manner integrating separately between and over each mass extinction event. Defining  $E(t_l^-)$  as the solution up to but not including the mass extinction event at time  $t_l$  we have:

$$E(t_l) = (1 - \nu_l)E(t_l^-) + \nu_l.$$

The first term reflects the probability that a lineage that does not go extinct during the  $l^{th}$  mass extinction event leaves no observable offspring (with probability  $E(t_l^-)$ ) whereas the second term reflects the fact that all lineages that go extinct during the  $l^{th}$  mass extinction leave no observed descendants with probability 1.

• The Probability Flow: The solution to the IVP is once again given by  $g_e(\tau) = g_e(s_e) \Psi(s_e, \tau) = g_e(s_e) \frac{\Phi(\tau)}{s_e}$  where:

$$\Phi(\tau) = \exp\left[\int_0^{\tau} 2\lambda(x)E(x) - \left(\lambda(x) + \bar{\mu}(x) + \sum_{l=1}^{L} m_l \delta(x - t_l) + \psi(x)\right) dx\right].$$

As with the CSAs, let  $L_{\tau}$  be the last index l such that  $t_l < \tau$ . We can separate out the mass

651

652

654

655

656

657

658

659

665

extinction terms in the following way.

$$\Phi(\tau) = \exp\left[\int_0^{\tau} 2\lambda(x)E(x) - (\lambda(x) + \bar{\mu}(x) + \psi(x)) dx\right] \prod_{l=1}^{L_{\tau}} e^{-m_l} 
= \exp\left[\int_0^{\tau} 2\lambda(x)E(x) - (\lambda(x) + \bar{\mu}(x) + \psi(x)) dx\right] \prod_{l=1}^{L_{\tau}} (1 - \nu_l) 
= \bar{\Phi}(\tau) \prod_{l=1}^{L_{\tau}} (1 - \nu_l) ,$$

where  $\bar{\Phi}(\tau)$  is defined as in Equation (A10).

• *The Likelihood:* Given these initial value problems the likelihood follows as in the general model.

$$\mathcal{L}_{ME} = \mathcal{S}\rho_0^{N_0} \bar{\Phi}(T) \prod_{l=1}^{L} (1 - \nu_l) \prod_{i=1}^{I} \left[ \lambda(x_i) \bar{\Phi}(x_i) \prod_{l=1}^{L_{x_i}} (1 - \nu_l) \right]$$

$$\times \prod_{j=1}^{n} \left[ \frac{\psi(y_j) \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right]}{\bar{\Phi}(y_j) \prod_{l=1}^{L_{y_j}} (1 - \nu_l)} \right] \prod_{k=1}^{m} \psi(z_k) (1 - r(z_k)).$$

As with the CSAs we can use relations analogous to Equations A12-A14 to rewrite the likelihood:

$$\mathcal{L}_{ME} = \mathcal{S}\rho_0^{N_0}\bar{\Phi}(T) \prod_{l=1}^{L} (1 - \nu_l)^{n_l} \times \prod_{i=1}^{I} \lambda(x_i)\bar{\Phi}(x_i) \prod_{j=1}^{n} \frac{\psi(y_j)}{\bar{\Phi}(y_j)} \left[ (1 - r(y_j))E(y_j) + r(y_j) \right] \prod_{k=1}^{m} \psi(z_k)(1 - r(z_k)),$$
(A17)

where  $n_l$  is defined as before as the number of lineages present at time  $t_l$ .

#### ALTERNATIVE CONDITIONING

Table S3 lists a number of possible conditionings,  $\mathcal{S}$  that can be applied to the tree likelihood. First, is the trivial case of no conditioning  $\mathcal{S}_0=1$  which gives the probability of the observed tree including the stem edge between time T and  $t_{\mathrm{MRCA}}$ . To obtain the model likelihood excluding the stem edge, i.e., conditioning of the  $t_{\mathrm{MRCA}}$ , can be obtained by setting  $\mathcal{S}=\mathcal{S}_1=\frac{\Phi(x_1)}{\Phi(T)}$ . Recall that the elements of  $\vec{x}$  are ordered such that  $x_1=t_{\mathrm{MRCA}}$  is the first (oldest) birth event.

Acknowledging that one would not reconstruct a phylogeny without any sampled lineages, we can condition the likelihood on observing at least one sampled lineage (either at or before the present day) given the time of origin,  $S_2 = \frac{1}{1-E(T)}$ . Or as with  $S_1$ , conditioning on at least one sampled lineage given the  $t_{\rm MRCA}$ . In order to have at least one sampled lineage *and* a most recent common ancestor, however, each daughter lineage of the common ancestor must have at least one

descendent. Hence we have  $S_3=rac{\Phi(x_1)}{\Phi(T)}rac{1}{(1-E(x_1))^2}$ . The general birth-death-sampling model assumes that all lineages alive at the present day are sampled with probability  $\rho_0$ . As with concerted sampling attempts (CSAs) prior to the present day, this present day CSA may include 672 the sampling of multiple lineages as well as possibly resulting in no sampled lineages. As with  $S_2$ and  $S_3$  we can condition the tree likelihood on observing at least one extant lineage at the present 674 day. To do so, we define  $\hat{E}(\tau) = E(\tau | \psi = 0)$ , the probability that a lineage alive at time  $\tau$  has no 675 extant descendants. Conditioning on the time of origin we have:  $S_4 = \frac{1}{1 - \hat{E}(T)}$ . Conditioning on the time of the most recent common ancestor we have:  $S_5 = \frac{\Phi(x_1)}{\Phi(T)} \frac{1}{2(1-\hat{E}(x_1))(1-E(x_1))}$ , where now at least one of the two daughter lineages of the common ancestor has a present day sample. In 678 many cases  $S_5$  is modified, however, to condition on both daughter lineages having an extant 679 sampled descendent:  $S_5' = \frac{\Phi(x_1)}{\Phi(T)} \frac{1}{\left(1 - \hat{E}(x_1)\right)^2}$ .

As an alternative to conditioning on at least one extant sampled descent, tree likelihoods 680

can be conditioned on having exactly  $N_0$  sampled (extant) descendants. Let  $\hat{E}_{N_0}(\tau)$  be the probability a lineage alive at time  $\tau$  has exactly  $N_0$  descendants. Although a general expression for  $\hat{E}_{N_0}(\tau)$  is unknown, in the case of constant birth, death, and sampling rates (the case in which this form of conditioning has been applied), the expression for  $\hat{E}_{N_0}(\tau)$  is given by (Gernhard, 2008; Kendall, 1948) and **Theorem 3.3** by Stadler Stadler (2010):

681

683

687

693

695

696

698

$$\begin{split} \hat{E}_{N_0}(\tau) = & \rho_0 \hat{\Phi}(\tau) \left( \frac{\lambda}{\mu} \hat{E}(\tau) \right)^{N_0 - 1} \\ \hat{E}_{N_0}(\tau) = & \rho_0 \hat{\Phi}(\tau) \left( \frac{\rho_0 \lambda \left( 1 - e^{-(\lambda - \mu)t} \right)}{\lambda \rho_0 + \left( \lambda (1 - \rho_0) - \mu \right) e^{-(\lambda - \mu)t}} \right)^{N_0 - 1}, \end{split}$$

where, like  $\hat{E}$ ,  $\hat{\Phi}$  is given by Equation (10) evaluated with where  $\psi=0$ . Given the time of origin we can condition on observing exactly  $N_0$  lineages by setting  $S = S_6 = \frac{1}{\hat{E}_{N_0}(T)}$ . When  $t_{\text{MRCA}}$  is given instead, then the number of descendants of the two daughter lineages must add up to  $N_0$ while both daughter lineages must still have at least one descendant(see Stadler (2010) Corollary 3.9). 692

$$S = S_7 = \frac{\Phi(x_1)}{\Phi(t_{\text{or}})} \left( \sum_{i=1}^{N_0 - 1} \hat{E}_i(x_1) \hat{E}_{N_0 - i}(x_1) \right)^{-1}$$

$$= \frac{\Phi(x_1)}{\Phi(t_{\text{or}})} \left[ (N_0 - 1)(\rho_0 \hat{\Phi}(x_1))^2 \left( \frac{\rho_0 \lambda \left( 1 - e^{-(\lambda - \mu)t} \right)}{\lambda \rho_0 + (\lambda (1 - \rho_0) - \mu) e^{-(\lambda - \mu)t}} \right)^{N_0 - 2} \right]^{-1}.$$

While early BDS models often employed such conditioning (Stadler, 2009, 2010), this form of conditioning has not been employed in many later models perhaps because the biological justification for such conditioning is vague.

The final form of conditioning used in the literature, which we will represent simply as  $S_8$ , is the multiplication of the BDS likelihood by a constant to account for the enumeration over the possible indistinguishable representations of a given tree. The value of this constant depends on

705

whether the tree considered is "labeled" and "oriented" (Gavryushkina et al., 2013) and whether, as in the derivation here, the vector of birth events,  $\vec{x}$ , is (un)ordered. Inclusion of such a constant should have no effect on the maximum likelihood inference of the model parameters given a specified tree. In cases where the constant is a function of the critical times (Heath et al., 2014), it can influence the inference when the parameters and the tree are jointly estimated.

- Barido-Sottani, J., T. G. Vaughan, and T. Stadler. 2018. Detection of HIV transmission clusters from phylogenetic trees using a multi-state birth–death model. Journal of The Royal Society Interface 15:20180512.
- Barido-Sottani, J., T. G. Vaughan, and T. Stadler. 2020. A Multitype Birth–Death Model for Bayesian Inference of Lineage-Specific Birth and Death Rates. Systematic Biology.
- Beaulieu, J. M. and B. C. O'Meara. 2016. Detecting Hidden Diversification Shifts in Models of Trait-Dependent Speciation and Extinction. Systematic Biology 65:583–601.
- Boskova, V., S. Bonhoeffer, and T. Stadler. 2014. Inference of Epidemiological Dynamics Based on Simulated Phylogenies Using Birth-Death and Coalescent Models. PLoS Computational Biology 10.
- Bouckaert, R., T. G. Vaughan, J. Barido-Sottani, S. Duchêne, M. Fourment, A. Gavryushkina,
   J. Heled, G. Jones, D. Kühnert, N. D. Maio, M. Matschiner, F. K. Mendes, N. F. Müller, H. A.
   Ogilvie, L. du Plessis, A. Popinga, A. Rambaut, D. Rasmussen, I. Siveroni, M. A. Suchard,
   C.-H. Wu, D. Xie, C. Zhang, T. Stadler, and A. J. Drummond. 2019. BEAST 2.5: An advanced software platform for Bayesian evolutionary analysis. PLOS Computational Biology
   15:e1006650.
- Caetano, D. S., B. C. O'Meara, and J. M. Beaulieu. 2018. Hidden state models improve state-dependent diversification approaches, including biogeographical models. Evolution 72:2308–2324.
- Drummond, A. J., O. G. Pybus, A. Rambaut, R. Forsberg, and A. G. Rodrigo. 2003. Measurably evolving populations. Trends in Ecology & Evolution 18:481–488.
- Drummond, A. J., A. Rambaut, B. Shapiro, and O. G. Pybus. 2005. Bayesian Coalescent Inference of Past Population Dynamics from Molecular Sequences. Molecular Biology and Evolution 22:1185–1192.
- du Plessis, L. and T. Stadler. 2015. Getting to the root of epidemic spread with phylodynamic analysis of genomic data. Trends in Microbiology 23:383–386.
- Duffy, S., L. A. Shackelton, and E. C. Holmes. 2008. Rates of evolutionary change in viruses: Patterns and determinants. Nature Reviews Genetics 9:267–276.
- Etienne, R. S., B. Haegeman, T. Stadler, T. Aze, P. N. Pearson, A. Purvis, and A. B. Phillimore. 2012. Diversity-dependence brings molecular phylogenies closer to agreement with the fossil record. Proceedings of the Royal Society B: Biological Sciences 279:1300–1309.
- Etienne, R. S. and J. Rosindell. 2012. Prolonging the Past Counteracts the Pull of the Present:
  Protracted Speciation Can Explain Observed Slowdowns in Diversification. Systematic
  Biology 61:204–204.
- Ezard, T. H. G., T. B. Quental, and M. J. Benton. 2016. The challenges to inferring the regulators of biodiversity in deep time. Philosophical Transactions of the Royal Society B: Biological Sciences 371:20150216.

- Feller, W. 1949. Proceedings of the [First] Berkeley Symposium on Mathematical Statistics and Probability. University of California Press, Berkeley California.
- FitzJohn, R. G. 2012. Diversitree: Comparative phylogenetic analyses of diversification in R. Methods in Ecology and Evolution 3:1084–1092.
- FitzJohn, R. G., W. P. Maddison, and S. P. Otto. 2009. Estimating Trait-Dependent Speciation and Extinction Rates from Incompletely Resolved Phylogenies. Systematic Biology 58:595–611.
- Gavryushkina, A., T. A. Heath, D. T. Ksepka, T. Stadler, D. Welch, and A. J. Drummond. 2017.
  Bayesian Total-Evidence Dating Reveals the Recent Crown Radiation of Penguins. Systematic Biology 66:57–73.
- Gavryushkina, A., D. Welch, and A. J. Drummond. 2013. Recursive algorithms for phylogenetic tree counting. Algorithms for Molecular Biology 8:26.
- Gavryushkina, A., D. Welch, T. Stadler, and A. J. Drummond. 2014. Bayesian Inference of
   Sampled Ancestor Trees for Epidemiology and Fossil Calibration. PLoS Computational
   Biology 10.
- Gernhard, T. 2008. The conditioned reconstructed process. Journal of Theoretical Biology 253:769–778.
- Goldberg, E. E. and B. Igić. 2012. Tempo and Mode in Plant Breeding System Evolution. Evolution 66:3701–3709.
- Goldberg, E. E., L. T. Lancaster, and R. H. Ree. 2011. Phylogenetic Inference of Reciprocal Effects between Geographic Range Evolution and Diversification. Systematic Biology 60:451–465.
- Grenfell, B. T., O. G. Pybus, J. R. Gog, J. L. N. Wood, J. M. Daly, J. A. Mumford, and E. C. Holmes. 2004. Unifying the epidemiological and evolutionary dynamics of pathogens. Science 303:327–332.
- Hagen, O., K. Hartmann, M. Steel, and T. Stadler. 2015. Age-Dependent Speciation Can Explain the Shape of Empirical Phylogenies. Systematic Biology 64:432–440.
- Heath, T. A., J. P. Huelsenbeck, and T. Stadler. 2014. The fossilized birth–death process for coherent calibration of divergence-time estimates. Proceedings of the National Academy of Sciences 111:E2957–E2966.
- Joy, J. B., R. M. McCloskey, T. Nguyen, R. H. Liang, Y. Khudyakov, A. Olmstead, M. Krajden,
  J. W. Ward, P. R. Harrigan, J. S. G. Montaner, and A. F. Y. Poon. 2016. The spread of hepatitis
  C virus genotype 1a in North America: A retrospective phylogenetic study. The Lancet.
  Infectious Diseases 16:698–702.
- Keeling, M. J. and P. Rohani. 2008. Modeling Infectious Diseases: In Humans and Animals. Princeton University Press.
- Kendall, D. G. 1948. On the Generalized "Birth-and-Death" Process. Annals of Mathematical Statistics 19:1–15.
- Kingman, J. F. C. 1982. On the Genealogy of Large Populations. Journal of Applied Probability 19:27–43.
- Kirkpatrick, M., T. Johnson, and N. Barton. 2002. General models of multilocus evolution. Genetics 161:1727–1750.

- Kühnert, D., T. Stadler, T. G. Vaughan, and A. J. Drummond. 2014. Simultaneous reconstruction of evolutionary history and epidemiological dynamics from viral sequences with the birth–death SIR model. Journal of The Royal Society Interface 11:20131106.
- Lafferty, K. D., G. DeLeo, C. J. Briggs, A. P. Dobson, T. Gross, and A. M. Kuris. 2015. A general consumer-resource population model. Science 349:854–857.
- Lambert, A. 2018. The coalescent of a sample from a binary branching process. Theoretical Population Biology 122:30–35.
- Lambert, A. and T. Stadler. 2013. Birth–death models and coalescent point processes: The shape and probability of reconstructed phylogenies. Theoretical Population Biology 90:113–128.
- Landis, M. J., D. A. R. Eaton, W. L. Clement, B. Park, E. L. Spriggs, P. W. Sweeney, E. J.
  Edwards, and M. J. Donoghue. 2021. Joint Phylogenetic Estimation of Geographic Movements
  and Biome Shifts during the Global Diversification of Viburnum. Systematic Biology
  70:67–85.
- Laudanno, G., B. Haegeman, D. L. Rabosky, and R. S. Etienne. 2020. Detecting Lineage-Specific Shifts in Diversification: A Proper Likelihood Approach. Systematic Biology.
- Lehtonen, S., D. Silvestro, D. N. Karger, C. Scotese, H. Tuomisto, M. Kessler, C. Peña,
  N. Wahlberg, and A. Antonelli. 2017. Environmentally driven extinction and opportunistic origination explain fern diversification patterns. Scientific Reports 7:4831.
- Leventhal, G. E., R. Kouyos, T. Stadler, V. von Wyl, S. Yerly, J. Böni, C. Cellerai, T. Klimkait, H. F. Günthard, and S. Bonhoeffer. 2012. Inferring Epidemic Contact Structure from Phylogenetic Trees. PLOS Computational Biology 8:e1002413.
- Louca, S. 2020. Simulating trees with millions of species. Bioinformatics.
- Louca, S. and M. Doebeli. 2018. Efficient comparative phylogenetics on large trees. Bioinformatics 34:1053–1055.
- Louca, S., A. McLaughlin, A. MacPherson, J. B. Joy, and M. W. Pennell. 2021. Fundamental identifiability limits in molecular epidemiology. Molecular Biology and Evolution.
- Louca, S. and M. W. Pennell. 2020a. Extant timetrees are consistent with a myriad of diversification histories. Nature Pages 1–4.
- Louca, S. and M. W. Pennell. 2020b. A General and Efficient Algorithm for the Likelihood of Diversification and Discrete-Trait Evolutionary Models. Systematic Biology 69:545–556.
- Maddison, W. P., P. E. Midford, and S. P. Otto. 2007. Estimating a Binary Character's Effect on Speciation and Extinction. Systematic Biology 56:701–710.
- Magee, A. F. and S. Höhna. 2021. Impact of K-Pg Mass Extinction Event on Crocodylomorpha Inferred from Phylogeny of Extinct and Extant Taxa. bioRxiv Page 2021.01.14.426715.
- Magnuson-Ford, K. and S. P. Otto. 2012. Linking the investigations of character evolution and species diversification. The American Naturalist 180:225–245.
- Morlon, H. 2014. Phylogenetic approaches for studying diversification. Ecology Letters 17:508–525.
- Morlon, H., T. L. Parsons, and J. B. Plotkin. 2011. Reconciling molecular phylogenies with the fossil record. Proceedings of the National Academy of Sciences of the United States of America 108:16327–16332.

- Morlon, H., M. D. Potts, and J. B. Plotkin. 2010. Inferring the Dynamics of Diversification: A Coalescent Approach. PLOS Biology 8:e1000493.
- Nee, S. 2006. Birth-Death Models in Macroevolution. Annual Review of Ecology, Evolution, and Systematics 37:1–17.
- Nee, S., R. M. May, and P. H. Harvey. 1994. The reconstructed evolutionary process. Phil. Trans. R. Soc. B Page 7.
- Ng, J. and S. D. Smith. 2014. How traits shape trees: New approaches for detecting character state-dependent lineage diversification. Journal of Evolutionary Biology 27:2035–2045.
- Pybus, O. G., A. Rambaut, and P. H. Harvey. 2000. An Integrated Framework for the Inference of Viral Population History From Reconstructed Genealogies. Genetics 155:1429–1437.
- Rabosky, D. L. and I. J. Lovette. 2008a. Density-dependent diversification in North American wood warblers. Proceedings of the Royal Society B: Biological Sciences 275:2363–2371.
- Rabosky, D. L. and I. J. Lovette. 2008b. Explosive evoltuionary radiation: Decreasing speciation or increasing extinction through time? Evolution 62:1866–1875.
- Rasmussen, D. A. and T. Stadler. 2019. Coupling adaptive molecular evolution to phylodynamics using fitness-dependent birth-death models. eLife 8.
- Raup, D. M. 1985. Mathematical models of cladogenesis. Paleobiology 11:42–52.
- Romero-Severson, E. O., I. Bulla, and T. Leitner. 2016. Phylogenetically resolving epidemiologic linkage. Proceedings of the National Academy of Sciences 113:2690–2695.
- Schluter, D. and M. W. Pennell. 2017. Speciation gradients and the distribution of biodiversity.

  Nature 546:48–55.
- Silvestro, D., J. Schnitzler, L. H. Liow, A. Antonelli, and N. Salamin. 2014. Bayesian Estimation of Speciation and Extinction from Incomplete Fossil Occurrence Data. Systematic Biology 63:349–367.
- Stadler, T. 2009. On incomplete sampling under birth-death models and connections to the sampling-based coalescent. Journal of Theoretical Biology 261:58–66.
- Stadler, T. 2010. Sampling-through-time in birth-death trees. Journal of Theoretical Biology 267:396–404.
- Stadler, T. 2011. Mammalian phylogeny reveals recent diversification rate shifts. Proceedings of the National Academy of Sciences of the United States of America 108:6187–6192.
- Stadler, T. 2013. Recovering speciation and extinction dynamics based on phylogenies. Journal of Evolutionary Biology 26:1203–1219.
- Stadler, T. and S. Bonhoeffer. 2013. Uncovering epidemiological dynamics in heterogeneous host populations using phylogenetic methods. Phil. Trans. R. Soc. B 368:20120198.
- Stadler, T., R. D. Kouyos, V. von Wyl, S. Yearly, and J. Böni. 2012. Estimating the Basic
   Reproductive Number from Viral Sequence Data. Mol. Biol. Evol. .
- Stadler, T., D. Kühnert, S. Bonhoeffer, and A. J. Drummond. 2013. Birth-death skyline plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). PNAS 110:228–233.

- Stadler, T., T. G. Vaughan, A. Gavryushkin, S. Guindon, D. Kühnert, G. E. Leventhal, and A. J. Drummond. 2015. How well can the exponential-growth coalescent approximate constant-rate birth-death population dynamics? Proceedings. Biological Sciences 282:20150420.
- Strimmer, K. and O. G. Pybus. 2001. Exploring the Demographic History of DNA Sequences
  Using the Generalized Skyline Plot. Molecular Biology and Evolution 18:2298–2305.
- Vaughan, T. G., G. E. Leventhal, D. A. Rasmussen, A. J. Drummond, D. Welch, and T. Stadler. 2019. Estimating Epidemic Incidence and Prevalence from Genomic Data. Molecular Biology and Evolution 36:1804–1816.
- Volz, E. M. 2012. Complex Population Dynamics and the Coalescent Under Neutrality. Genetics 190:187–201.
- Volz, E. M. and S. D. W. Frost. 2014. Sampling through time and phylodynamic inference with coalescent and birth—death models. Journal of The Royal Society Interface 11:20140945.
- Volz, E. M., S. L. Kosakovsky Pond, M. J. Ward, A. J. Leigh Brown, and S. D. W. Frost. 2009. Phylodynamics of Infectious Disease Epidemics. Genetics 183:1421–1430.

### FIGURE CAPTIONS

878

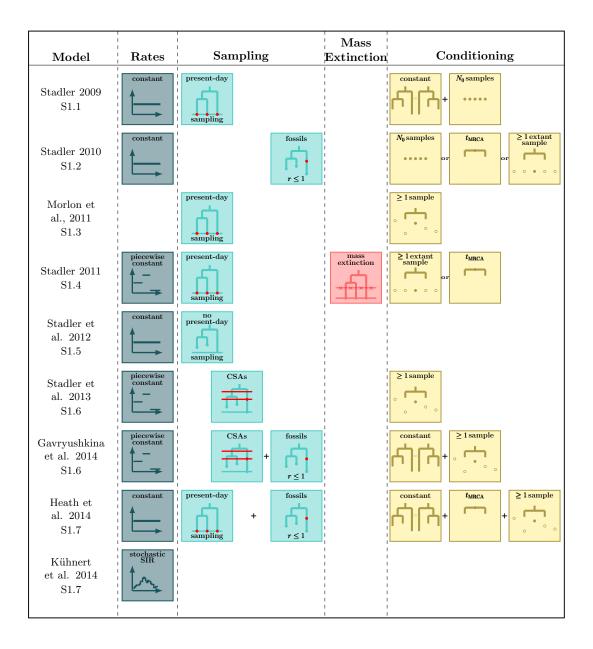


Figure 1: **Sub-model assumptions.** Rate, sampling, mass extinction, and conditioning assumptions of existing sub-models of the general time-variable BDS process. The key points are that i) each of the previously developed models we considered can be obtained by adding specific combinations of constraints to the various parameters of the general BDS model; and ii) that there are many plausible, and potentially biologically informative combinations of constraints that have not been considered by researchers in epidemiology or macroevolution.

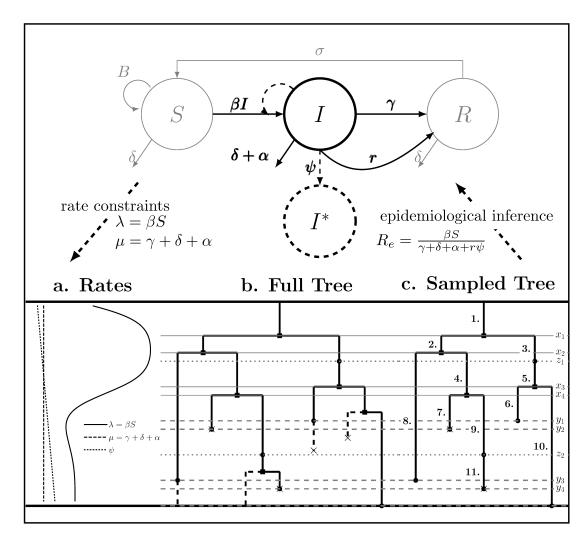


Figure B1: **BDS-SIR model connection.** Top: The SIR epidemiological model. Black (gray) lines and classes represent rates and variables followed (in)directly by the BDS model. The SIR model can be used to constrain the rates of the BDS model (panel a). Simulated forward in time, the result of the BDS stochastic processes is a *full tree* (panel b) giving the complete genealogy of the viral population. Pruning away extinct and unsampled lineages produces the *sampled tree* (panel c). Arising from a BDS process, this sampled tree can be summarized in two ways. First by the set of edges (labeled 1-11) or as a set of critical times (horizontal lines) including: 1) the time of birth events (solid,  $x_i$ ) 2) terminal sampling times (dashed,  $y_j$ ), and 3) ancestral sampling times (dotted,  $z_k$ ). Given the inferred rates from a reconstructed sampled tree, these rates can be used to estimate characteristic parameters of the SIR model, for example the basic or effective reproductive number.

## **Single-Type Model**

#### RELATIONSHIPS BETWEEN EXISTING MODELS

In the appendix, we proved that one could go from the most general case to specific sub-models by incorporating additional constraints to the parameters. In this section, we illustrate how to work in the other direction — that is, to start with the most assumptions of a particular sub-model and derive its likelihood function using the same five-step procedure used to derive the general BDS model in the Main Text. In addition to illustrating the utility of our mathematical technique, by deriving the likelihoods of previously developed models, we are able to unify a diverse and, occasionally opaque, literature using a common terminology, notation, and formulation.

Stadler 2009

Here we re-derive the likelihood given by **Equation 2** in (Stadler, 2009). Note throughout all equation, corollary, and theorem references in other publications will be placed in bold face type.

• Step 1: Specify the model.

11

12

15

16

17

18

20

21

22

23

25

- Constant rates:  $\lambda(\tau) = \lambda$ ,  $\mu(\tau) = \mu$ .
- Birth-death model with incomplete sampling at present day:  $\psi(\tau) = 0$  and  $\rho_0 \leqslant 1$ .
- Conditioning on there being exactly  $N_0$  lineages at the present day given the time of origin,  $S_6$  and un-ordered birth events  $S_8 = (N_0 1)!$ .
- Step 2: IVP for  $g_e(\tau)$ .

$$\begin{split} \frac{dg_e(\tau)}{d\tau} &= -\; (\lambda - \mu)g_e(\tau) + 2\lambda g_e(\tau)E(\tau) \\ g_e(\tau) &= \begin{cases} \lambda g_{e1}(s_e)g_{e2}(s_e) & \text{birth} \\ \rho_0 & \text{present day} \end{cases} \end{split}$$

• Step 3: IVP for  $E(\tau)$ .

$$\frac{dE(\tau)}{d\tau} = -(\lambda + \mu)E(\tau) + \lambda E(\tau)^2 + \mu \quad \text{where } E(0) = 1 - \rho_0.$$

Given the constant rate assumption there exists a general solution for  $E(\tau)$ .

$$E(\tau) = \frac{\lambda + \mu + c_1 \frac{\exp[-c_1 t](1 - c_2) - (1 + c_2)}{\exp[-c_1 t](1 - c_2) + (1 + c_2)}}{2\lambda}$$
$$c_1 = \lambda - \mu \quad c_2 = \frac{\lambda - \mu - 2\lambda \rho_0}{\lambda - \mu}$$

• Step 4: Derive  $g_{stem}(T)$ .

$$g_{stem}(T) = \rho_0^{N_0} \lambda^{N_0 - 1} \prod_{edges} \Psi(s_e, t_e).$$

• Step 5: Calculate  $g_{stem}(T)$  wrt the critical time representation. Given the assumption of constant rates and no Poisson sampling the expression for  $\Phi(\tau)$  simplifies to:

$$\Phi(\tau) = \frac{e^{-x(\lambda-\mu)}(\lambda-\mu)^2}{((\lambda(1-\rho_0)-\mu)e^{-x(\lambda-\mu)} + \lambda\rho_0)^2}.$$

Hence we have:

28

30

31

32

33

34

35

36

37

40

42

43

44

45

47

48

49

51

$$g_{stem}(T) = \rho_0^{N_0} \lambda^{N_0 - 1} \prod_{i=1}^{N_0 - 1} \frac{e^{-x_i(\lambda - \mu)} (\lambda - \mu)^2}{\left( (\lambda(1 - \rho_0) - \mu)e^{-x_i(\lambda - \mu)} + \lambda \rho_0 \right)^2}.$$

• Step 6: Impose conditioning. The likelihood is conditioned on having exactly  $N_0$  lineages at the present day ( $S_6$ ) and a constant  $S_8 = (N_0 - 1)!$  as the birth events are left un-ordered.

$$\mathcal{L}_{S09} = \frac{(N_0 - 1)!}{\hat{E}_{N_0}(T)} \rho_0^{N_0} \lambda^{N_0 - 1} \prod_{i=1}^{N_0 - 1} \frac{e^{-x_i(\lambda - \mu)} (\lambda - \mu)^2}{((\lambda(1 - \rho_0) - \mu)e^{-x_i(\lambda - \mu)} + \lambda \rho_0)^2}.$$
 (S1)

Stadler 2010

- Step 1: Specify the model.
  - Constant rates:  $\lambda(\tau) = \lambda$ ,  $\mu(\tau) = \mu$ ,  $\psi(\tau) = \psi$ .
  - No removal upon sampling, r = 0.
  - Multiple presented.
    - \* **Equation 3** (Stadler, 2010): No conditioning.
    - \* **Equation 4** (Stadler, 2010): Exactly  $N_0$  extant sampled tips.
    - \* Corollary 3.7 (Stadler, 2010): At least one extant tip conditioning on the time of origin.
    - \* Equation 5 (Stadler, 2010): At least one extant tip conditioning on the  $t_{MRCA}$ .
    - \* **Equation 6** (Stadler, 2010):  $N_0$  extant tips conditioning on the time the  $t_{MRCA}$ .
- Step 2: IVP for  $q_e(\tau)$ .

$$\frac{dg_e(\tau)}{d\tau} = -\left(\lambda + \mu + \psi\right)g_e(\tau) + 2\lambda g_e(\tau)E(\tau)$$
 
$$g_e(s_e) = \begin{cases} \lambda g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ \psi g_e(s_e) & \text{ancestral sampling event} \\ \psi E(s_e) & \text{terminal sampling event} \\ \rho_0 & s_e = 0, \text{edge sampled at present day} \end{cases}$$

• Step 3: IVP for  $E(\tau)$ . Given the constant rate assumption there exists a general solution for  $E(\tau)$ .

$$\frac{dE(\tau)}{d\tau} = -(\lambda + \mu + \psi) E(\tau) + \lambda E(\tau)^2 + \mu \quad \text{where } E(0) = 1 - \rho_0.$$

This is a Bernoulli differential equation and has a known solution as given by Stadler (Stadler, 2010).

$$E(t) = \frac{\lambda + \mu + c_1 \frac{\exp[-c_1 t](1 - c_2) - (1 + c_2)}{\exp[-c_1 t](1 - c_2) + (1 + c_2)}}{2\lambda}$$

$$c_1 = \left| \sqrt{(\lambda - \mu - \psi)^2 + 4\lambda\psi} \right| \quad c_2 = \frac{\lambda - \mu - 2\lambda\rho_0 - \psi}{c_1}.$$

• Step 4: Derive  $g_{stem}(T)$ .

52

53

54

55

56

57

60

61

62

63

71

73

$$g_{stem}(T) = \underbrace{\rho_0^{N_0}}_{\text{extant tips}} \underbrace{\prod_{i=1}^{I} \lambda \prod_{j=1}^{n} \psi E(y_j)}_{\text{births}} \underbrace{\prod_{e \in \mathcal{T}}_{\text{extinct tips}} \Psi(s_e, t_e)}_{\text{ancestral samples}} \underbrace{\prod_{e \in \mathcal{T}} \Psi(s_e, t_e)}_{\text{edges}}.$$

• Step 5: Calculate  $g_{stem}(T)$  wrt the critical time representation. Given the assumption of constant rates and no Poisson sampling the expression for  $\Phi(\tau)$  simplifies to:

$$\Phi(\tau) = \exp\left[\int_0^{\tau} 2\lambda E(x) - (\lambda + \mu + \psi) dx\right].$$

Hence we have:

$$g_{stem}(T) = \underbrace{\Phi(T)}_{\text{root}} \underbrace{\rho_0^{N_0}}_{\text{extant}} \underbrace{\prod_{i=1}^{I} \lambda \Phi(x_i)}_{\text{births}} \underbrace{\prod_{j=1}^{n} \frac{\psi}{\Phi(y_j)} E(y_j)}_{\text{extinct}} \underbrace{\prod_{\substack{k=1 \\ \text{ancestral samples}}}^{m} \psi \; .$$

• Step 6: Impose conditioning. Imposing an arbitrary conditioning we have the following likelihood. Note that  $I = N_0 + n - 1$ .

$$\mathcal{L}_{S10} = \mathcal{S}\Phi(T)\rho_0^{N_0}\lambda^{N_0+n-1}\psi^{n+m} \prod_{i=1}^{N_0+n-1} \Phi(x_i) \prod_{j=1}^n \frac{E(y_j)}{\Phi(y_j)}$$
(S2)

- **Equation 3** (Stadler, 2010): No conditioning.  $S = S_0 = 1$ .
- Equation 4 (Stadler, 2010): Exactly  $N_0$  extant sampled tips.  $S = S_6 = \frac{1}{\hat{E}_{N_0}(T)}$  where in the case of constant rates we have:

$$\hat{E}_{N_0}(T) = \hat{\phi}(T) \left(\frac{\mu}{\lambda} E(T)\right)^{N_0 - 1}.$$

where  $\hat{\phi}(\tau)$  equals  $\phi(\tau)$  where  $\psi = 0$ .

- Corollary 3.7 (Stadler, 2010): At least one extant tip conditioning on the time of origin.  $S = S_4 = \frac{1}{1-\hat{E}(T)}$  where  $\hat{E}(T)$  is given by the solution for  $E(\tau)$  above given that  $\psi = 0$ .
- Equation 5 (Stadler, 2010): At least one extant tip conditioning on the  $t_{MRCA}$ ,

$$S = S_{5'} = \frac{\phi(x_1)}{\phi(T)(1-\hat{E}(x_1)^2}.$$

75

76

77

78

79

80

85

86

88

90

91

92

93

94

95

- **Equation 6** (Stadler, 2010):  $N_0$  extant tips conditioning on the time the  $t_{\rm MRCA}$ ,

$$S = S_7 = \frac{\Phi(x_1)}{\Phi(t_{or})} \left( \sum_{i=1}^{N_0 - 1} \hat{E}_i(x_1) \hat{E}_{N_0 - i}(x_1) \right)^{-1},$$

#### Morlon et al. 2011

Here we derive **Equation 1** from (Morlon et al., 2011).

- Step 1: Specify the model.
  - Time variable rates.
  - Birth-death only  $\psi(\tau) = 0$ .
  - At least one extant sample  $S_2 = S_4$ .
  - Step 2: IVP for  $g_e(\tau)$ .

$$\begin{split} \frac{dg_e(\tau)}{d\tau} &= -\left(\lambda(\tau) + \mu(\tau)\right)g_e(\tau) + 2\lambda(\tau)g_e(\tau)E(\tau) \\ g_e(s_e) &= \begin{cases} \lambda g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ \rho_0 & s_e = 0, \text{edge sampled at present day} \end{cases} \end{split}$$

• Step 3: IVP for  $E(\tau)$ .

$$\frac{dE(\tau)}{d\tau} = -\left(\lambda(\tau) + \mu(\tau)\right)E(\tau) + \lambda(\tau)E(\tau)^2 + \mu(\tau) \quad \text{where } E(0) = 1 - \rho_0.$$

The general solution of this differential equation is given by:

$$E(\tau) = 1 - \frac{\rho_0 \exp\left[\int_0^{\tau} (\lambda(u) - \mu(u)) \, du\right]}{1 + \int_0^{\tau} \rho_0 \exp\left[\int_0^{x} (\lambda(u) - \mu(u)) \, du\right] dx},$$

see **Equation 2** in (Morlon et al., 2011).

• Step 4: Derive  $g_{stem}(T)$ .

$$g_{stem}(T) = \underbrace{\rho_0^{N_0}}_{\text{extant}} \underbrace{\prod_{i=1}^{I} \lambda(x_i)}_{\text{births}} \underbrace{\prod_{e \in \mathcal{T}} \Psi(s_e, t_e)}_{\text{edges}},$$

where the expression for  $\Psi(s_e, t_e)$  is given by **Equation 3** in (Morlon et al., 2011).

• Step 5: Calculate  $g_{stem}(T)$  wrt the critical time representation. Given  $\Phi(\tau) = \Psi(0, \tau)$  from **Equation 3** for in (Morlon et al., 2011) we have:

$$\Phi(\tau) = \exp\left[\int_0^{\tau} \left(\lambda(\sigma) - \mu(\sigma)\right) d\sigma\right] \left[1 + \frac{\int_0^{\tau} \rho_0 \lambda(u) \exp\left[\int_0^u \left(\lambda(\sigma) - \mu(\sigma)\right) d\sigma\right] du}{1 + \rho_0}\right]^{-2}.$$

Hence we have:

97

100

101

102

103

104

105 106

107

109

110

111

112

113

114

115

116

118

$$g_{stem}(T) = \underbrace{\Phi(T)}_{\text{root}} \underbrace{\rho_0^{N_0}}_{\text{extant tips}} \underbrace{\prod_{i=1}^{N_0-1} \lambda(x_i) \Phi(x_i)}_{\text{births}}.$$

Note  $I = N_0 - 1$ .

• Step 6: Impose conditioning. The likelihood given by **Equation 1** in (Morlon et al., 2011) is conditioned on the existence of at least one sampled lineage,  $S = S_3 = \frac{1}{1 - E(t_{or})}$ .

$$\mathcal{L}_{M11} = \frac{\rho_0^{N_0}}{1 - E(t_{\text{or}})} \Phi(T) \prod_{i=1}^{N_0 - 1} \lambda(x_i) \Phi(x_i)$$
 (S3)

Stadler et al. 2011

Here we derive the likelihoods given by **Theorem 2.6** and **2.7** in (Stadler, 2011).

- Step 1: Specify the model.
  - piecewise-constant Poissonian rates.  $\lambda(\tau)=\lambda_l$  and  $\bar{\mu}(\tau)=\bar{\mu}_l$  if  $t_l\leqslant \tau < t_{l+1} \quad l=0,2,...L+1$  where  $t_0=0$  and  $t_{L+1}=T$ .
  - No Poisson sampling,  $\psi(\tau) = 0$ .
  - Mass extinction events at times  $t_l$  l = 1, 2, ...L as specified above.

$$\mu(\tau) = \bar{\mu} + \sum_{l}^{L} m_l \delta(\tau - t_l)$$
$$m_l = -\ln(1 - \nu_l).$$

- Theorem 2.6 (Stadler, 2011) imposes no additional conditioning  $S = S_0$  whereas Theorem 2.7 (Stadler, 2011) conditions on observing at least one descendent given the time of the most recent common ancestor  $S = S_2 = S_4$ .
- Step 2: IVP for  $g_e(\tau)$ .

$$\begin{split} \frac{dg_e(\tau)}{d\tau} &= -\left(\lambda(\tau) + \mu(\tau)\right)g_e(\tau) + 2\lambda(\tau)g_e(\tau)E(\tau) \\ g_e(s_e) &= \begin{cases} \lambda(s_e)g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ \rho_0 & s_e = 0, \text{edge sampled at present day}, \end{cases} \end{split}$$

where  $\mu(\tau)$  is given above.

In terms of the probability flow  $g_e(\tau) = g_e(s_e)\Psi(s_e,\tau)$ , where:

$$\Psi(s_e, \tau) = \exp\left[\int_{s_e}^{\tau} 2\lambda(x)E(x) - (\lambda(x) + \mu(x)) dx\right].$$

Here  $\mu(\tau)$  includes the mass extinction events.

• Step 3: IVP for  $E(\tau)$ .

119

121

122

123

124

125

127

128

129

131

132

133

134

135

136

$$\frac{dE(\tau)}{d\tau} = -(\lambda(\tau) + \mu(\tau))E(\tau) + \lambda(\tau)E(\tau)^2 + \mu(\tau)$$
  

$$E(0) = 1 - \rho_0.$$

Given the piecewise constant nature there is a known general solution. Let  $E(\tau) = E_l(\tau)$  where  $t_l < \tau \le t_{l+1}$ . Then define  $E_{l-1}(t_l^-)$  as the solution up to but not including the mass extinction event at time  $t_l$  we have:

$$E_l(t_l) = E_{l-1}(t_l) = (1 - \nu_l)E_{l-1}(t_l^-) + \nu_l,$$

where  $E_l(\tau)$  is given by a solution similar to that in Equation (A1).

$$E_{l}(\tau) = \frac{\lambda_{l} + \bar{\mu}_{l}}{2\lambda_{l}} + \frac{c_{1}}{2\lambda_{l}} \frac{e^{-c_{1}t}(1 - c_{2}) - (1 + c_{2})}{e^{-c_{1}t}(1 - c_{2}) + (1 + c_{2})}$$

$$c_{1} = \left| \sqrt{(\lambda_{l} - \bar{\mu}_{l})^{2}} \right| \quad c_{2} = -\frac{\lambda_{l} - \bar{\mu}_{l} - 2\lambda_{l}(1 - E_{l}(t_{l}))}{c_{1}},$$

where  $E_0(t_0) = 1 - \rho_0$ .

• Step 4: Derive  $g_{stem}(T)$ . The expression for  $g_{stem}(T)$  is given by:

$$g_{stem}(T) = \underbrace{\rho_0^{N_0}}_{\text{extant tips}} \underbrace{\prod_{i=1}^{I} \lambda(x_i)}_{\text{births}} \underbrace{\prod_{e \in \mathcal{T}} \Psi(s_e, t_e)}_{\text{edges}},$$

where  $\lambda(\tau)$  and  $\Psi(s_e, \tau)$  are specified above.

• Step 5: The critical time representation. We once again define the sub-flow  $\Phi_l(\tau)$  where  $t_l < \tau \leqslant t_{l+1}$ :

$$\Phi_l(\tau) = \exp\left[\int_{t_l}^{\tau} 2\lambda(x)E(x) - (\lambda(x) + \bar{\mu}(x) + m_l\delta(x - t_l)) dx\right]$$

$$= \exp\left[\int_{t_l}^{\tau} 2\lambda(x)E(x) - (\lambda(x) + \bar{\mu}(x)) dx\right] (1 - \nu_l)$$

$$= \bar{\Phi}_l(\tau) (1 - \nu_l).$$

Note  $\nu_0 = 0$ . The complete flow is, as given by Equation (A3).

$$\Phi(\tau) = \bar{\Phi}_{L_{\tau}}(\tau) \prod_{l=1}^{L_{\tau}} \bar{\Phi}_{l-1}(t_l) (1 - \nu_l)$$

$$\Phi(t_k) = \underbrace{\bar{\Phi}_k(t_k)}_{l} \prod_{l=1}^{k} \bar{\Phi}_{l-1}(t_l) (1 - \nu_l) = \prod_{l=1}^{k} \bar{\Phi}_{l-1}(t_l) (1 - \nu_l),$$

where  $L_{\tau}$  is once again the largest index l such that  $t_l < \tau$ .

The critical time representation of  $g_{stem}$  then is:

$$g_{stem}(T) = \rho_0^{N_0} \bar{\Phi}_L(T) \prod_{l=1}^L \bar{\Phi}_{l-1}(t_l) (1 - \nu_l) \prod_{i=1}^I \left[ \lambda(x_i) \bar{\Phi}_{L_{x_i}}(x_i) \prod_{l=1}^{L_{x_i}} \bar{\Phi}_{l-1}(t_l) (1 - \nu_l) \right].$$

Defining  $\alpha_l$  as the number of observed birth events  $\geqslant t_l$  we can rewrite the product:

$$\prod_{i=1}^{I} \prod_{l=1}^{L_{x_i}} \bar{\Phi}_{l-1}(t_l) = \prod_{l=1}^{L} \left(\bar{\Phi}_{l-1}(t_l)\right)^{\alpha_l}$$

$$\prod_{i=1}^{I} \prod_{l=1}^{L_{x_i}} (1 - \nu_l) = \prod_{l=1}^{L} (1 - \nu_l)^{\alpha_l}$$

Hence we have:

137

138

139

140

141

142

143

145

146

148

149

150

151

153

$$g_{stem}(T) = \rho_0^{N_0} \bar{\Phi}_L(T) \prod_{l=1}^L \left( \bar{\Phi}_{l-1}(t_l) (1 - \nu_l) \right)^{\alpha_l + 1} \prod_{i=1}^I \left[ \lambda(x_i) \bar{\Phi}_{L_{x_i}}(x_i) \right]$$
$$= \rho_0^{N_0} \bar{\Phi}_L(T) \prod_{l=1}^L \left( \bar{\Phi}_{l-1}(t_l) (1 - \nu_l) \right)^{n_l} \prod_{i=1}^I \left[ \lambda(x_i) \bar{\Phi}_{L_{x_i}}(x_i) \right],$$

where  $n_l$  is the number of lineages in the observed phylogeny at time  $t_l$  then  $n_l = \alpha_l + 1$ .

• Step 6: Likelihood conditioning.

For **Theorem 2.6** (Stadler, 2011) the likelihood is given by:

$$\mathcal{L}_{S11} = \rho_0^{N_0} \bar{\Phi}_L(T) \prod_{l=1}^L \left( \bar{\Phi}_{l-1}(t_l) (1 - \nu_l) \right)^{n_l} \prod_{i=1}^I \left[ \lambda(x_i) \bar{\Phi}_{L_{x_i}}(x_i) \right]$$

For **Theorem 2.7** (Stadler, 2011) the likelihood is given by:

$$\mathcal{L}_{S11} = \frac{\rho_0^{N_0}}{(1 - E(x_1))^2} \bar{\Phi}_{L_{x_1}}(x_1) \prod_{l=1}^L \left( \bar{\Phi}_{l-1}(t_l)(1 - \nu_l) \right)^{n_l} \prod_{i=1}^I \left[ \lambda(x_i) \bar{\Phi}_{L_{x_i}}(x_i) \right]$$

Stadler et al. 2012

Here we derive **Equation 1** in (Stadler et al., 2012).

- Step 1: Specify the model.
  - Constant birth, death, and sampling rates.
- No present day sampling. All lineages removed upon sampling r=1.
- No conditioning.

• Step 2: IVP for  $g_e(\tau)$ .

155

156

157

159

160

161

162

165

166

168

169

170

$$\begin{split} \frac{dg_e(\tau)}{d\tau} &= -\left(\lambda + \mu + \psi\right)g_e(\tau) + 2\lambda g_e(\tau)E(\tau) \\ g_e(s_e) &= \begin{cases} \lambda(s_e)g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ \psi(s_e) & \text{terminal sampling event} \end{cases} \end{split}$$

• Step 3: IVP for  $E(\tau)$ .

$$\frac{dE(\tau)}{d\tau} = -(\lambda + \mu + \psi)E(\tau) + \lambda E(\tau)^2 + \mu \quad \text{where } E(0) = 0.$$

The solution to this differential equation is given by:

$$E(t) = \frac{\lambda + \mu + c_1 \frac{\exp[-c_1 t](1 - c_2) - (1 + c_2)}{\exp[-c_1 t](1 - c_2) + (1 + c_2)}}{2\lambda}$$

$$c_1 = \left| \sqrt{(\lambda - \mu - \psi)^2 + 4\lambda\psi} \right| \quad c_2 = \frac{\lambda - \mu - \psi}{c_1}.$$

• Step 4: Derive  $g_{stem}(T)$ .

$$g_{stem}(T) = \underbrace{\prod_{i=1}^{I} \lambda \prod_{j=1}^{n} \psi \prod_{e \in \mathcal{T}} \Psi(s_e, t_e)}_{\text{births}} \underbrace{\prod_{e \in \mathcal{T}} \Psi(s_e, t_e)}_{\text{edges}}$$
$$= \lambda^{n-1} \psi^n \prod_{e \in \mathcal{T}} \Psi(s_e, t_e),$$

where I = n - 1.

• Step 5: Critical time representation. Given the assumption of constant rates and no Poisson sampling the expression for  $\Phi(t)$  simplifies to:

$$\Phi(\tau) = \exp\left[\int_0^{\tau} 2\lambda E(x) - (\lambda + \mu + \psi) dx\right].$$

Hence we have:

$$g_{stem}(T) = \underbrace{\Phi(T)}_{\text{root}} \lambda^{n-1} \psi^n \underbrace{\prod_{i=1}^{n-1} \Phi(x_i)}_{\text{births}} \underbrace{\prod_{j=1}^{n} \frac{1}{\Phi(y_j)}}_{\text{extinct tips}}$$

• Step 6: Impose conditioning  $S = S_0$ .

$$\mathcal{L}_{S12} = \lambda^{n-1} \psi^n \Phi(T) \prod_{i=1}^{n-1} \Phi(x_i) \prod_{j=1}^n \frac{1}{\Phi(y_j)}$$
 (S4)

#### Stadler et al. 2013 and Gavryushkina et al. 2014

Here we derive the tree likelihood given by **Theorem 1** in (Stadler and Bonhoeffer, 2013) and due to their shared the likelihood given by **Equation 4** in (Gavryushkina et al., 2014).

• Step 1: Specify the model.

171

172

174

175

176

177

178

180

181

183

184

185

186

187

188

189

190

191

194

195

Piecewise constant Poissonian birth, death, and sampling rates.

\* define transition times 
$$0 = t_0 < t_1 < ... < t_{L+1} = T$$
.

\* 
$$\lambda(\tau) = \lambda_l t_l < \tau \leqslant t_{l+1}$$
.

$$* \mu(\tau) = \mu_l t_l < \tau \leqslant t_{l+1}.$$

\* 
$$\bar{\psi}(\tau) = \bar{\psi}_l t_l < \tau \leqslant t_{l+1}$$
.

- Fossils/Ancestral sampling:
  - \* (Stadler and Bonhoeffer, 2013): No fossils r = 1.
  - \* (Gavryushkina et al., 2014): piecewise-constant rate  $r(\tau) = r_l t_l < \tau \leqslant t_{l+1}$ .
- Concerted sampling attempts at each internal transition time  $t_l$  where l = 1, 2...L.
  - \* The probability of a lineage being sampled during the CSA at time  $t_l$  is  $\rho_l$ .
  - \* The resulting total sampling rate is given by:

$$\psi(\tau) = \bar{\psi}(\tau) + \sum_{l=1}^{L} w_l \delta(\tau - t_l)$$
$$w_l = -\ln(1 - \rho_l).$$

- Conditioning:
  - \* Stadler et al. (Stadler and Bonhoeffer, 2013)—At least one sampled lineage,  $S = S_2$ .
  - \* Gavryushkina et al. (Gavryushkina et al., 2014)— At least one sample ( $S_2$ ) and a constant giving the number of un-oriented phylogenies  $S = S_8$ .
- Step 2: Derive IVP for  $g_e(\tau)$ . As in section 1 we have:

$$\frac{dg_e(\tau)}{d\tau} = -\left(\lambda(\tau) + \mu(\tau) + \psi(\tau)\right)g_e(\tau) + 2\lambda(\tau)g_e(\tau)E(\tau)$$
 birth event giving rise to edges e1 and e2 
$$g_e(s_e) = \begin{cases} \lambda(s_e)g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ (1 - r(s_e))\bar{\psi}(s_e)g_{e1}(s_e) & \text{Poisson ancestral sampling event} \\ \bar{\psi}(s_e)r(s_e) + \bar{\psi}(s_e)(1 - r(s_e))E(s_e) & \text{Poisson terminal sampling event} \\ (1 - r(t_l))\rho_lg_{e1}(t_l) & \text{ancestral sample at } t_l \\ \rho_lr(t_l) + \rho_l(1 - r(t_l))E(t_l) & \text{terminal sample at } t_l \\ \rho_0 & s_e = 0, \text{edge sampled at present day} \end{cases}$$

• Step 3: Derive IVP for  $E(\tau)$ . The IVP for  $E(\tau)$  follows from Equation (8).

$$\frac{dE(\tau)}{d\tau} = -\left(\lambda(\tau) + \mu(\tau) + \psi(\tau)\right)E(\tau) + \lambda(\tau)E(\tau)^2 + \mu(\tau)$$
  
$$E(0) = 1 - \rho_0.$$

Given the piecewise constant nature there is a known general solution. Let  $E(\tau) = E_l(\tau)$  where  $t_l < \tau \leqslant t_{l+1}$ . Then define  $E_{l-1}(t_l^-)$  as the solution up to but not including the CSA at time  $t_l$  we have:

$$E_l(t_l) = E_{l-1}(t_l) = (1 - \rho_l)E_{l-1}(t_l^-),$$

where  $E_l(\tau)$  is given by a solution similar to that in Equation (A1).

$$E_{l}(\tau) = \frac{\lambda_{l} + \mu_{l} + \bar{\psi}_{l}}{2\lambda_{l}} + \frac{c_{1}}{2\lambda_{l}} \frac{e^{-c_{1}t}(1 - c_{2}) - (1 + c_{2})}{e^{-c_{1}t}(1 - c_{2}) + (1 + c_{2})}$$

$$c_{1} = \left| \sqrt{(\lambda_{l} - \mu_{l} - \bar{\psi}_{l})^{2} + 4\lambda_{l} + \bar{\psi}_{l}} \right| \quad c_{2} = -\frac{\lambda_{l} - \mu_{l} - 2\lambda_{l}(1 - E_{l}(t_{l})) - \bar{\psi}_{l}}{c_{1}},$$

where  $E_0(t_0) = 1 - \rho_0$ .

196

197

199

201

202

204

205

206

207

208

209

210

211

• Step 4: Derive  $g_{stem}(T)$ . As in section 1 the edge representation of  $g_{stem}$  is given by:

$$g_{stem}(T) = \rho_0^{N_0} \prod_{i=1}^{L} \lambda(x_i) \prod_{j=1}^{n} \psi(y_j) \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right] \prod_{k=1}^{m} \psi(z_j) (1 - r(y_j))$$

$$\times \prod_{l=1}^{L} \rho_l \left[ (1 - r_l) E(t_l) + r_l \right]^{N_l} \prod_{l=1}^{L} \left[ \rho_l (1 - r_l) \right]^{M_l} \prod_{edges} \Psi(s_e, t_e),$$

where  $\Psi(s_e, t_e)$  is given by Equation (6).

• Step 5: Critical time representation. We once again define the sub-flow  $\Phi_l(\tau)$  where  $t_l < \tau \leqslant t_{l+1}$ :

$$\Phi_{l}(\tau) = \exp\left[\int_{t_{l}}^{\tau} 2\lambda(x)E(x) - \left(\lambda(x) + \mu(x) + \bar{\psi}(x) + w_{l}\delta(x - t_{l})\right)dx\right]$$

$$= \exp\left[\int_{t_{l}}^{\tau} 2\lambda(x)E(x) - \left(\lambda(x) + \mu(x) + \bar{\psi}(s)\right)dx\right](1 - \rho_{l})$$

$$= \bar{\Phi}_{l}(\tau)(1 - \rho_{l}).$$

The complete flow is, as given by Equation (A3).

$$\Phi(\tau) = \bar{\Phi}_{L_{\tau}}(\tau) \prod_{l=1}^{L_{\tau}} \bar{\Phi}_{l-1}(t_l) (1 - \rho_l)$$

$$\Phi(t_k) = \underbrace{\bar{\Phi}_k(t_k)}_{1} \prod_{l=1}^{k} \bar{\Phi}_{l-1}(t_l) (1 - \rho_l) = \prod_{l=1}^{k} \bar{\Phi}_{l-1}(t_l) (1 - \rho_l),$$

where  $L_{\tau}$  is once again the largest index l such that  $t_l < \tau$ . The critical time representation

of  $g_{stem}$  from the previous stem gives the following.

$$g_{stem}(T) = \rho_0^{N_0} \bar{\Phi}_L(T) \prod_{l=1}^L \bar{\Phi}_{l-1}(t_l) (1-\rho_l) \prod_{i=1}^I \lambda(x_i) \bar{\Phi}_{L_{x_i}}(x_i) \prod_{i=1}^I \left[ \prod_{l=1}^{L_{x_i}} \bar{\Phi}_{l-1}(t_l) (1-\rho_l) \right]$$

$$\times \prod_{j=1}^n \frac{\psi(y_j)}{\bar{\Phi}_{L_{y_j}}(y_j)} \left[ (1-r(y_j)) E(y_j) + r(y_j) \right] \prod_{j=1}^n \left[ \prod_{l=1}^{L_{y_j}} \frac{1}{\bar{\Phi}_{l-1}(t_l) (1-\rho_l)} \right] \prod_{k=1}^m \psi(z_k) (1-r(z_k))$$

$$\times \prod_{l=1}^L \left( \frac{\rho_l}{\bar{\Phi}_l(t_l)} \left[ (1-r_l) E(t_l) + r_l \right] \right)^{N_l} \prod_{l=1}^L \left[ \prod_{j=1}^l \frac{1}{(1-\rho_j) \bar{\Phi}_{j-1}(t_j)} \right] \prod_{l=1}^L \left[ \rho_l (1-r_l) \right]^{M_l}.$$

We can then simplify the likelihood using the relationships similar to those given by Equation (A4).

$$\prod_{i=1}^{I} \prod_{l=1}^{L_{x_i}} \bar{\Phi}_{l-1}(t_l) = \prod_{l=1}^{L} \left[ \bar{\Phi}_{l-1}(t_l) \right]^{\alpha_l}$$

$$\prod_{j=1}^{n} \prod_{l=1}^{L_{y_j}} \frac{1}{\bar{\Phi}_{l-1}(t_l)} = \prod_{l=1}^{L} \left[ \bar{\Phi}_{l-1}(t_l) \right]^{-\sigma_l}.$$

$$\prod_{i=1}^{I} \prod_{l=1}^{L_{x_i}} (1 - \rho_l) = \prod_{l=1}^{L} (1 - \rho_l)^{\alpha_l}$$

$$\prod_{j=1}^{n} \prod_{l=1}^{L_{y_j}} \frac{1}{(1 - \rho_l)} = \prod_{l=1}^{L} (1 - \rho_l)^{-\sigma_l}$$

$$\prod_{l=1}^{L} \left[ \prod_{j=1}^{l} \frac{1}{(1 - \rho_j)^{N_l}} \right] = \prod_{l=1}^{L} \frac{1}{(1 - \rho_l)^{\beta_l}},$$

where  $\alpha_l$  is the number of birth events before time  $t_l$  and  $\sigma_l$  is the number of Poisson sampling events before time  $t_l$  and  $\beta_l$  is the number of lineages sampled during the CSAs up to and including at time  $t_l$ .

The resulting simplified likelihood is given by:

$$g_{stem}(T) = \rho_0^{N_0} \bar{\Phi}_L(T) \prod_{l=1}^{L} \left[ \bar{\Phi}_{l-1}(t_l) (1 - \rho_l) \right]^{\alpha_l + 1 - \sigma_l - \beta_l}$$

$$\times \prod_{j=1}^{n} \frac{\psi(y_j)}{\bar{\Phi}_{L_{y_j}}(y_j)} \left[ (1 - r(y_j)) E(y_j) + r(y_j) \right] \prod_{k=1}^{m} \psi(z_k) (1 - r(z_k))$$

$$\times \prod_{l=1}^{L} \left( \frac{\rho_l}{\bar{\Phi}_l(t_l)} \left[ (1 - r_l) E(t_l) + r_l \right] \right)^{N_l} \prod_{l=1}^{L} \left[ \rho_l (1 - r_l) \right]^{M_l}.$$

- Step 6: Conditioned likelihood.
  - Stadler et al. (Stadler and Bonhoeffer, 2013) conditions on observing at least one sample. In addition the likelihood given by Stadler et al. assumes that all lineages are removed upon sampling  $r_l = 1$  an assumption we will apply now. The likelihoods can

226

227

228

229

230

231

232

233

235

236

238

239

240

242

243

be simplified by letting  $n_l = \alpha_l + 1 - \sigma_l - \beta_l$  be the number of lineages that are alive immediately following the concerted sampling attempt at time  $t_l$ 

$$\mathcal{L}_{S13} = \frac{\rho_0^{N_0}}{1 - E(T)} \bar{\Phi}_L(T) \prod_{l=1}^L \left[ \bar{\Phi}_{l-1}(t_l)(1 - \rho_l) \right]^{n_l} \prod_{j=1}^n \frac{\psi(y_j)}{\bar{\Phi}_{L_{y_j}}(y_j)} \prod_{l=1}^L \left( \frac{\rho_l}{\bar{\Phi}_l(t_l)} \right)^{N_l}.$$
(S5)

- Gavryushkina et al. (Gavryushkina et al., 2014) also condition on observing at least one lineage as well as multiply by a constant giving the number of un-oriented trees.

$$\mathcal{L}_{G14} = 8_8 \frac{\rho_0^{N_0}}{1 - E(T)} \bar{\Phi}_L(T) \prod_{l=1}^L \left[ \bar{\Phi}_{l-1}(t_l)(1 - \rho_l) \right]^{n_l} \\
\times \prod_{j=1}^n \frac{\psi(y_j)}{\bar{\Phi}_{L_{y_j}}(y_j)} \left[ (1 - r(y_j))E(y_j) + r(y_j) \right] \prod_{k=1}^m \psi(z_k)(1 - r(z_k)) \\
\times \prod_{l=1}^L \left( \frac{\rho_l}{\bar{\Phi}_l(t_l)} \left[ (1 - r_l)E(t_l) + r_l \right] \right)^{N_l} \prod_{l=1}^L \left[ \rho_l(1 - r_l) \right]^{M_l}$$
(S6)

Heath et al. 2014

- Step 1: Specify the model.
  - Constant rates:  $\lambda(\tau) = \lambda$ ,  $\mu(\tau) = \mu$ ,  $\psi(\tau) = \psi$ .
    - \* Here  $\psi$  denotes the sampling rate of fossils.
  - Birth-death model with fossils, r = 0.
  - Conditioned on observing  $\geqslant 1$  sample given  $t_{MRCA}$  ( $S_3$ ) and enumerated over all possible attachments of fossils (terminal sampling events before the present day ( $S_8$ )).
- Step 2: Derive IVP for  $q_e(\tau)$ .

$$\frac{dg_e(\tau)}{d\tau} = -\left(\lambda + \mu + \psi\right)g_e(\tau) + 2\lambda g_e(\tau)E(\tau)$$
 
$$g_e(s_e) = \begin{cases} \lambda g_{e1}(s_e)g_{e2}(s_e) & \text{birth event giving rise to edges e1 and e2} \\ \psi g_{e1}(s_e) & \text{ancestral sampling event} \\ \psi E(s_e) & \text{terminal sampling event} \\ \rho_0 & s_e = 0, \text{edge sampled at present day} \end{cases}$$

• Step 3: Derive IVP for  $E(\tau)$ .

$$\frac{dE(\tau)}{d\tau} = -(\lambda + \mu + \psi)E(\tau) + \lambda E(\tau)^2 + \mu$$
  
 
$$E(0) = 1 - \rho_0.$$

This has the same general solution as given above in section.

• Step 4: Derive  $g_{stem}(T)$ .

$$g_{stem}(T) = (\rho_0)^{N_0} \prod_{i=1}^{I} \lambda \prod_{j=1}^{n} \psi E(y_j) \prod_{k=1}^{m} \psi \prod_{edges} \Psi(s_e, t_e).$$

• Step 5: Critical time representation.

$$g_{stem}(T) = \Phi(T) (\rho_0)^{N_0} \prod_{i=1}^{I} \lambda \Phi(x_i) \prod_{j=1}^{n} \frac{\psi}{\Phi(y_j)} E(y_j) \prod_{k=1}^{m} \psi.$$

• Step 6: Conditioned likelihood. Conditioning on the probability that both daughter clades of the MRCA have at least one sampled extant descendant.  $S_5'$ . The general solution to  $\hat{E}(\tau)$  is known as given by the solution to  $E(\tau)$  in section . Conditioning on all the possible attachments of the fossils (terminal samples before the present day). This requires multiplication by a constant  $S_8$ . Let  $\gamma(\tau)$  be the number of lineages alive at time  $\tau$ . The the resulting number of fossil trees is.

$$S = \prod_{i=1}^{n} 2\gamma(y_j) \prod_{k=1}^{m} \gamma(z_k).$$
 (S7)

$$\mathcal{L}_{H14} = \frac{\Phi(x_1)}{\left(1 - \hat{E}(x_1)\right)^2} \left(\rho_0\right)^{N_0} \prod_{i=1}^{I} \lambda \Phi(x_i) \prod_{j=1}^{n} \frac{2\gamma(y_j)\psi}{\Phi(y_j)} E(y_j) \prod_{k=1}^{m} \gamma(z_k)\psi$$
 (S8)

256

Tables

# **General Birth-Death-Sampling Model Notation**

Variable	Definition
$\lambda(\tau)$	The rate at which lineages speciate at time $\tau^a$ . $\lambda: \mathbb{R} \to [0, \infty)$
$\mu(\tau)$	The rate at which lineages go extinct at time $\tau$ . $\mu: \mathbb{R} \to [0, \infty)$
$\psi(\tau)$	The rate at which lineages are sampled at time $\tau$ . $\psi: \mathbb{R} \to [0, \infty)$
$r(\tau)$	The probability that sampling is associated with host recovery/viral extinction at
	time $\tau$ . $r: \mathbb{R} \to [0,1]$
$ ho_0$	The probability of sampling a lineage alive at the present day. $\rho_0: \mathbb{R} \to (0,1]$
T	The time of origin of the phylogeny/epidemic.
$\vec{x}$	The vector of $I$ branching times. $\vec{x} = \{x_1, x_2,, x_i,, x_I\}^d$
$\vec{y}$	The vector of $n$ internal (Poisson) sampling times without sampled descendants (ter-
	minal nodes). $\vec{y} = \{y_1, y_2,, y_j,, y_n\}$
$\vec{z}$	The vector of $m$ internal (Poisson) sampling times of lineages with sampled descen-
	dants (sampled ancestors). $\vec{z} = \{z_1, z_2,, z_k,, z_m\}$
$N_0$	The number of lineage sampled at the present day.
$g_e( au)$	The likelihood density of observing a given phylogeny by the present day descending
	from a single edge $e$ alive at age $\tau$ $^{\rm e}$ .
$E(\tau)$	The probability that a lineage alive at time $\tau$ leaves no sampled descendants in the
	phylogeny <sup>f</sup> .
S	The conditioning of the phylogeny.
	$\{\lambda(\tau), \mu(\tau), \psi(\tau), r(\tau), T, \rho_0\}$
CSA <sup>b</sup> , m	ass extinctions, and piecewise constant models
L	The total number of past concerted sampling events.
$L_t$ $\vec{t}$	The index of the time $t_l$ at or after time $t$ , i.e. the largest index such that $t_l \leqslant \tau$ .
$ \vec{t} $	Vector of times of CSAs in order of most recent to oldest $t_1 < t_2 < < t_L$ .
	$\vec{t} = \{t_1, t_2,, t_l,, t_L\}$ c
$ec{ ho}$	Vector of sampling probabilities of lineages sampled during each CSA. $\rho_l : \mathbb{R} \to \mathbb{R}$
	$(0,1]. \vec{\rho} = \{\rho_0, \rho_1,, \rho_l,, \rho_L\}$
$\vec{r}$	The probability that a lineage that is sampled during each CSA is removed upon
	sampling. $r_l: \mathbb{R} \to []0,1]$ . $\vec{r} = \{r_1, r_2,, r_l,, r_L\}$
$N_l$	The number of tips sampled during the CSA at time $t_l$ .
$M_l$	The number of ancestral samples during the CSA at time $t_l$ .
$n_l$	The number of lineages that cross time $t_l$ .
$\alpha_l$	The number of observed birth events prior to time $t_l$ ( $x_i > t_l$ ).
$\sigma_l (\beta_l)$	The number of Poissonian (concerted) sampled lineages prior to or at time $\tau \geqslant t_l$ .
$\Theta_{\mathrm{CSA}} = \{$	$\{\lambda(\tau), \mu(\tau), \psi(\tau), r(\tau), T, \vec{\tau}, \vec{\rho}, \vec{r}\}$

#### Table S1: model notation.

- <sup>a</sup>Throughout  $\tau$  is measured in units of time before the present day ( $\tau=0$ ).
- <sup>b</sup> CSA: Concerted Sampling Attempt.
- $^{\mathrm{c}}$  When used  $t_0$  indicates the present day and  $t_{L+1}=T$  .
- <sup>d</sup> When used  $x_1$  indicates the time of the most recent common ancestor  $\tau_{\mathrm{MRCA}}$
- $^{\mathrm{e}}$  The edge e spans time  $\tau,\,s_{e}\geqslant\tau\geqslant t_{e}$
- $^{\mathrm{f}}$   $\hat{E}(\tau)$  is the special case where  $\psi(\tau)=0$

Model	$\lambda(\tau), \mu(\tau), \psi(\tau)$	r( au)	$\rho_0$	$\rho_l$	$r_l$	$\nu_l$	S	Section
Stadler 2009	constant, $\psi = 0$	NA	> 0	0	NA	0	$S_6S_8$	
Stadler 2010	constant	0	> 0	0	NA	0	multiple	
Morlon et al. 2011	$\psi = 0$	NA	> 0	0	NA	0	$S_1 = S_2$	
Stadler 2011	piecewise, $\psi = 0$	NA	> 0	0	NA	$\geqslant 0$	$S_3 = S_5$	
Stadler et al. 2012	constant	1	0	0	NA	0	$S_1$	
Stadler and Bonhoeffer 2013	piecewise	1	> 0	$\geqslant 0$	1	0	$S_1$	
Gavryushkina et al. 2014	piecewise	piecewise	> 0	$\geqslant 0$	$\geqslant 0$	0	$S_1 \times S_7$	
Kühnert et al. 2014	stochastic SIR	constant	0	0	NA	0	?a	
Heath et al. 2014	constant	0	> 0	0	NA	0	$S_6$	

Table S2: Relationship of single-type sub-models to the general BDS model. <sup>a</sup> No specific conditioning mentioned.

Condition	Description	Examples
$S_0 = 1$	No conditioning	Eq.3(Stadler, 2010), Eq.1(Stadler et al., 2012) & Thrm. 2.6(Stadler, 2011)
$S_1 = \frac{\Phi(x_1)}{\Phi(T)}$	Likelihood given the $t_{\mathrm{MRCA}}$ rather than on the time of origin.	
$\mathcal{S}_2 = rac{1}{1 - E(T)}$	At least one sampled descendent either at or before the present day (given $T$ ).	Eq.2Morlon et al. (2011), Thrm.1Stadler and Bonhoeffer (2013)
$S_3 = \frac{\Phi(x_1)}{\Phi(T)(1 - E(x_1))^2}$	At least one sampled descendent at or before the present day given $\tau_{\text{MRCA}} = x_1$	(Heath et al., 2014)
$S_4 = \frac{1}{1 - \hat{E}(T)}$	At least one <i>extant</i> sampled lineage (given $T$ ).	<b>Cor.3.7</b> (Stadler, 2010) & <b>Thrm.2.7</b> (Stadler, 2011)
$S_5 = \frac{\Phi(x_1)}{\Phi(T)2(1-\hat{E}(x_1))(1-E(x_1))}$	At least one <i>extant</i> sampled lineage given $\tau_{\text{MRCA}} = x_1$	
$S_5' = \frac{\Phi(x_1)}{\Phi(T)(1-\hat{E}(x_1))^2}$ $S_6 = \frac{1}{\hat{E}_{N_0}(T)}$	Both daughters of MRCA have at least one <i>extant</i> sampled lineage Exactly $N_0$ <i>extant</i> sampled lineages (given $T$ ).	Eq. 5 Stadler (2010)  Eq.2Stadler (2009),  Eq.4(Stadler, 2010) & Cor.3.6  Stadler (2010)
$\delta_7 = \frac{\Phi(x_1)}{\Phi(t_{\text{or}})} \left( \sum_{i=1}^{N_0 - 1} \hat{E}_i(x_1) \hat{E}_{N_0 - i}(x_1) \right)^{-1}$	Exactly $N_0$ extant sampled lineages given $t_{\text{MRCA}}$ .	<b>Eq.6</b> Stadler (2010)
$S_8 = constant$	Multiply by a constant	Stadler (2009) & Gavryushkina et al. (2014) <sup>a</sup> & Heath et al. (2014)

Table S3: **Alternative conditioning of tree likelihood**.

<sup>a</sup> See Gavryushkina et al. (2013) for algorithms for enumerating trees.

### The multi-type model

#### THE MULTI-TYPE DIVERSIFICATION MODEL

Here we consider the diversification of lineages of A discrete types. Lineages of type  $a \in \{1,2,...A\}$  speciate/give birth to lineages of type  $b \in \{1,2,...A\}$  at the time-variable rate  $\lambda_{a,b}(\tau)$  at time  $\tau$  before the present day. As above we will use  $\tau$  to denote time moving backward from the present day ( $\tau=0$ ) to the origin of the phylogeny  $\tau=T$ . When a=b speciation occurs without cladogenetic change, when  $a\neq b$  speciation is coincident with state change. In addition to the lineages changing state at birth events, lineages can mutate anagenetically from state a to state b at rate  $\gamma_{a,b}(\tau)$ . Lineages of type a alive at time  $\tau$ , go extinct/die at rate  $\mu_a(\tau)$ . For  $\tau>0$  lineages are sampled at rate  $\psi_a(\tau)$ . Upon sampling a lineage may be removed from the population (e.g., sampling is coincident with treatment) with the state-dependent probability  $r_a(\tau)$ . Finally, all lineages alive at the present day are sampled with a state-dependent probability  $\rho_a$ . Model notation is summarized in Table S4.

As with the single-type model, the result of the mutli-type diversification process is a *full* and a *sampled* tree. Now however the tree is characterized by its topology  $\mathcal T$  and the *colouring* of the tree  $\mathcal C$  denoting the states of each lineage through time. Below we first derive an expression for the likelihood of a given coloured tree,  $\mathcal L$   $(\mathcal T, \mathcal C|\Theta_{MBDS})$ . However, as plausibly attainable data consists of knowledge at only some or all the sampled ancestral nodes and/or tips, we must then integrate the likelihood  $\mathcal L$   $(\mathcal T, \mathcal C|\Theta_{MBDS})$  over all possible tree colourings consistent with the observed data.

## Derivation of $\mathcal{L}$ ( $\mathfrak{I}, \mathfrak{C}|\Theta_{MRDS}$ )L(T,C)

We derive the likelihood using the steps used above for the single-type birth death sampling model. As with the single-type model for the likelihood calculation be begin by representing the phylogeny as a series of edges. However, as we are now referring to the coloured phylogeny we consider the set of all coloured edges, with each edge being a a segment of the phylogeny that is all of one colour beginning and ending at birth, sampling, or mutation events. Specifically, moving backward in time towards the tree origin let edge e start at time  $s_e$  before the present at a birth event, sampling (ancestral or tip), or mutation event, and continue toward the tree origin until time  $\tau_e$  ending at either a birth event, ancestral sampling event, mutation event or at the tree origin.

As with the single-type model the tree likelihood depends on two different functions. First,  $g_{e,a}(\tau)$  is the probability that an edge e with state a alive at time  $\tau$  (hence  $t_e > \tau > s_e$ ) gives rise to the subsequently observed phylogeny between  $\tau$  and the present day. Second,  $E_a(\tau)$  is the probability a lineage of type a alive at time  $\tau$  has no sampled descendants between  $\tau$  and the present day. We begin below by first deriving the initial value problems for these two functions. Step 1: Derive the Initial Value Problem for  $g_{e,a}(\tau)$ .

To simplify the notation we define the total birth and mutation rates of a given type  $\Lambda_a = \sum_b \lambda_{a,b}$  and  $\Gamma_a = \sum_b \gamma_{a,b}$  and the relative probability a birth or mutation event was of a given type,  $p_{\lambda_{a,b}} = \frac{\lambda_{a,b}}{\Lambda_a}$  and  $p_{\gamma_{a,b}} = \frac{\gamma_{a,b}}{\Gamma_a}$ . For some small amount of time  $\Delta \tau$  the recursion equation for  $g_{e,a}(\tau)$  is given by:

$$g_{e,a}(\tau + \Delta \tau) \approx \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \Gamma_a \Delta \tau)) g_{e,a}(\tau)}_{\text{No Event}} + \underbrace{((\Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \Gamma_a \Delta \tau)) \times \sum_{b}^{A} p_{\lambda_{a,b}} \iota_{a,b} g_{e,a}(\tau) E_b(\tau)}_{\text{Birth Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(\mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \Gamma_a \Delta \tau)) \times 0}_{\text{Death Event}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(\psi_a \Delta \tau)(1 - \Gamma_a \Delta \tau)) \times 0}_{\text{Sampling Event}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \mu_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)) \times 0}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)(1 - \psi_a \Delta \tau)(\Gamma_a \Delta \tau)}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)}_{\text{Mutation Events}} + \underbrace{((1 - \Lambda_a \Delta \tau)(1 - \psi_a \Delta \tau)(1 - \psi_a \Delta \tau)}_{\text{Mutati$$

where  $\iota_{a,b}$  is an indicator variable that has value 2 when a=b and 1 otherwise. This recursion equation uses the fact that edges are all in state with mutation events resulting in the end of an edge and the origin of a daughter edge with a different state. This is reflected by the initial conditions of  $g_{e,a}(\tau)$  as given below in equation (S11).

As above we can then use the definition of a derivative to obtain the corresponding differential equation

$$\frac{dg_{e,a}(\tau)}{d\tau} = -\left(\Lambda_a + \mu_a + \psi_a + \Gamma_a\right)g_{e,a}(\tau) + \sum_b \iota_{a,b}\lambda_{a,b}g_{e,a}(\tau)E_b(\tau) \tag{S10}$$

The initial conditions for  $g_{e,a}$  at the start of the edge, time  $s_e$ , depend on the event that occurs in the sampled tree at that time. Specifically there are five types of events in the observed tree. 1) Observed birth events where a lineage of type a speciates producing a new lineage of type b. 2) Ancestral sampling events where a lineage of type a is sampled, is not removed from the population, and then has subsequently observed descendants. 3) Terminal sampling events where a lineage of type a is sampled and then either is immediately removed from the population or remains in the population but has no sampled descendants. 4) Transition events where a lineage of type a transitions to a lineage of type a. There are two possible ways a transition event can occur. First, a lineage can switch states due to a anagenetic mutation event or second an apparent transition can arise due to a birth event with cladogenetic state change followed by subsequent extinction of the parental lineage. We refer to this second form of transition event as a hidden birth event.

$$g_{e,a}(s_e) = \begin{cases} \lambda_{a,b}(s_e)g_{e_1,a}(s_e)g_{e_2,b}(s_e) & \text{birth event } a \to a+b \\ (1-r_a(s_e))\psi_a(s_e)g_{e_1,a}(s_e) & \text{ancestral sampling event} \\ r_a(s_e)\psi_a(s_e) + (1-r_a(s_e))\psi_a(s_e)E_a(s_e) & \text{terminal sampling event} \\ (\gamma_{a,b}(s_e) + \lambda_{a,b}E_a)\,g_{e_1,b}(s_e) & \text{mutation/hidden birth event } a \to b \\ \rho_a & \text{sampled at present day} \end{cases}$$
(S11)

Importantly, as the ODES in equation (S10) are linear the corresponding IVP has the

following general solution given the initial conditions above.

$$g_{e,a}(\tau) = g_{e,a}(s_e)\Psi_a(s_e, \tau), \quad \text{where}$$

$$\Psi_a(s_e, \tau) = exp \left[ \int_{s_e}^{\tau} -(\Lambda_a(x) + \mu_a(x) + \psi_a(x) + \Gamma_a(x)) + \sum_b \iota_{a,b} \lambda_{a,b}(x) E_b(x) dx \right]. \tag{S12}$$

The IVP for the function  $E_a(\tau)$  can be derived in an analogous manner.

Step 2: Derive the Initial Value Problem for  $E_a(\tau)$  Once again we begin by deriving a recursion equation for the change in  $E_a(\tau)$  over from time  $\tau$  to  $\tau + \Delta \tau$ .

$$E_{a}(\tau + \Delta \tau) \approx \underbrace{((1 - \Lambda_{a} \Delta \tau)(1 - \mu_{a} \Delta \tau)(1 - \psi_{a} \Delta \tau)(1 - \Gamma_{a} \Delta \tau)) E_{a}(\tau)}_{\text{No Event}}$$

$$+ \underbrace{((\Lambda_{a} \Delta \tau)(1 - \mu_{a} \Delta \tau)(1 - \psi_{a} \Delta \tau)(1 - \Gamma_{a} \Delta \tau)) \times \sum_{b}^{A} p_{\lambda_{a,b}} E_{a}(\tau) E_{b}(\tau)}_{\text{Birth Events}}$$

$$+ \underbrace{((1 - \Lambda_{a} \Delta \tau)(\mu_{a} \Delta \tau)(1 - \psi_{a} \Delta \tau)(1 - \Gamma_{a} \Delta \tau)) \times 1}_{\text{Death Event}}$$

$$+ \underbrace{((1 - \Lambda_{a} \Delta \tau)(1 - \mu_{a} \Delta \tau)(\psi_{a} \Delta \tau)(1 - \Gamma_{a} \Delta \tau)) \times 0}_{\text{Sampling Event}}$$

$$+ \underbrace{((1 - \Lambda_{a} \Delta \tau)(1 - \mu_{a} \Delta \tau)(\psi_{a} \Delta \tau)(1 - \Gamma_{a} \Delta \tau)) \times 0}_{\text{Mutation Events}}$$

$$(S13)$$

Using the definition of a derivative we have:

$$\frac{dE_a(\tau)}{d\tau} = -\left(\Lambda_a + \mu_a + \psi_a + \Gamma_a\right)E_a(\tau) + \sum_b \lambda_{a,b}E_a(\tau)E_b(\tau) + \mu_a + \sum_b \gamma_{a,b}E_b(\tau), \quad (S14)$$

with initial conditions given at the present day by:

$$E_a(0) = 1 - \rho_a \tag{S15}$$

Step 3: Derive Expression for  $g_{stem}(T)$  As in equation (9) for the single-type model, the tree likelihood is given by the value of the function g of the stem edge evaluated at the time of the origin of the phylogeny, T. From the solution to the initial value problem given by equation (S12) and the initial conditions (S11) we can write  $g_{stem}(T)$  as a product over the events in the tree and all the probability flow  $\Psi$  of the edges in-between. To do so let  $\vec{x}_{a,b}$  be a vector of length  $I_{a,b}$ , giving the time before the present day of all the observed birth events where lineages of type a give rise to lineages of type a. Let  $\vec{y}_a$  be a vector of length a giving the sampling time of tips of type a and a and a a vector of length a the sampling type of sampled ancestors of type a. Finally let a with a be the time of observed transition events where a edge of type a transitions becomes an edge of type a. Once again, transition events can arise due to both direct mutation and hidden birth events. The resulting expression for the likelihood a signer by:

$$g_{\textit{stem}}(T) = \prod_{a=1}^{A} \left[ \underbrace{\rho_{a}^{N_{a}}}_{\text{present day sampling}} \times \underbrace{\prod_{b=1}^{A} \prod_{i=1}^{I_{a,b}} \lambda_{a,b}(x_{a,b,i})}_{\text{birth events}} \right] \times \underbrace{\prod_{j=1}^{J_{a}} \psi_{a}(y_{a,j}) \left(1 - r_{a}(y_{a,j})\right) E_{a}(y_{a,j}) + \psi_{a}(y_{a,j}) r_{a}(y_{a,j})}_{\text{terminal samples}} \times \underbrace{\prod_{k=1}^{K_{a}} \psi_{a}(z_{a,k}) \left(1 - r_{a}(z_{a,k})\right)}_{\text{ancestral samples}} \times \underbrace{\prod_{l=1}^{L_{a,b}} \left[\gamma_{a,b}(w_{a,b,l}) + \lambda_{a,b}(w_{a,b,l}) E_{a}(w_{a,b,l})\right]}_{\text{transitions}} \underbrace{\prod_{edges \text{ of type } a} \Psi_{a}(s_{e}, \tau_{e})}_{\text{edges of type } a}$$

$$(S16)$$

Step 4: Rewrite  $g_{stem}(T)$  in Terms of Critical Times

Rather than enumerate  $\Psi$  over the edges of the phylogeny we can rewrite equation (S16) in terms of only the critical times  $\vec{x}, \vec{y}, \vec{z}, \vec{w}$  Written in this form the likelihood also depends on  $c^*$  the colour of the phylogeny at the origin. To do so we define  $\Phi_a(\tau) = \Psi_a(0,\tau)$  and rewrite the probability flow  $\Psi_a(s_e,\tau)$  as a ratio:

$$\Psi_a(s_e, \tau) = \frac{\Phi_a(\tau)}{\Phi_a(s_e)} \tag{S17}$$

Substitution into equation (S16).

333

334

335

336

337

340

$$g_{\textit{stem},c^*}(T) = \underbrace{\left[\prod_{a=1}^{A} \rho_a^{N_a}\right]}_{\text{present day sampling}} \times \underbrace{\left[\Phi_{c^*}(T)\right]}_{\text{stem}} \times \underbrace{\left[\prod_{a=1}^{A} \prod_{b=1}^{I} \lambda_{a,b}(x_{a,b,i}) \underbrace{\Phi_a(x_{a,b,i})}_{\Phi_a(x_{a,b,i})}\right]}_{\text{birth events}}$$

$$\times \underbrace{\left[\prod_{a=1}^{A} \prod_{j=1}^{J_a} \left[\psi_a(y_{a,j})(1-r_a(y_{a,j}))E_a(y_{a,j})+\psi_a(y_{a,j})r_a(y_{a,j})\right]}_{\text{terminal sampling events}}\right]}_{\text{terminal sampling events}}$$

$$\times \underbrace{\left[\prod_{a=1}^{A} \prod_{k=1}^{K_a} \psi_a(z_{a,k})(1-_a(z_{a,k})) \underbrace{\Phi_a(z_{a,k})}_{\Phi_a(z_{a,k})}\right]}_{\text{ancestral sampling events}}$$

$$\times \underbrace{\left[\prod_{a=1}^{A} \prod_{l=1}^{L_{a,b}} \left[\gamma_{a,b}(w_{a,b,l}) + \lambda_{a,b}(w_{a,b,l})E_a(w_{a,b,l})\right] \underbrace{\Phi_b(w_{a,b,l})}_{\Phi_a(w_{a,b,l})}\right]}_{\text{transition sampling events}}$$

#### **Step 5: Condition the Likelihood**

We conclude, as in the single type model by multiplying the likelihood by including a general form of conditioning S.

$$\mathcal{L}(\Theta_{\text{MBDS}}|\mathcal{T}, \mathcal{C}) = \mathcal{S} \times \left[ \prod_{a=1}^{A} \rho_{a}^{N_{a}} \right] \times \left[ \Phi_{c^{*}}(T) \right] \times \left[ \prod_{a=1}^{A} \prod_{b=1}^{A} \prod_{i=1}^{I_{a,b}} \lambda_{a,b}(x_{a,b,i}) \Phi_{b}(x_{a,b,i}) \right]$$

$$\times \left[ \prod_{a=1}^{A} \prod_{j=1}^{J_{a}} \left[ \psi_{a}(y_{a,j})(1 - r_{a}(y_{a,j})) E_{a}(y_{a,j}) + \psi_{a}(y_{a,j}) r_{a}(y_{a,j}) \right] \frac{1}{\Phi_{a}(y_{a,j})} \right]$$

$$\times \left[ \prod_{a=1}^{A} \prod_{k=1}^{K_{a}} \psi_{a}(z_{a,k})(1 - r_{a}(z_{a,k})) \right]$$

$$\times \left[ \prod_{a=1}^{A} \prod_{b\neq a} \prod_{l=1}^{L_{a,b}} \left[ \gamma_{a,b}(w_{a,b,l}) + \lambda_{a,b}(w_{a,b,l}) E_{a}(w_{a,b,l}) \right] \frac{\Phi_{b}(w_{a,b,l})}{\Phi_{a}(w_{a,b,l})} \right]$$
(S19)

Due to the combination of cladogenetic change and different possibility of hidden birth events due to extinction and sampling, there are six different types of birth events included in the likelihood of the sampled tree. The following diagram summarizes these different birth events and how each is included in the equation (S19).

343

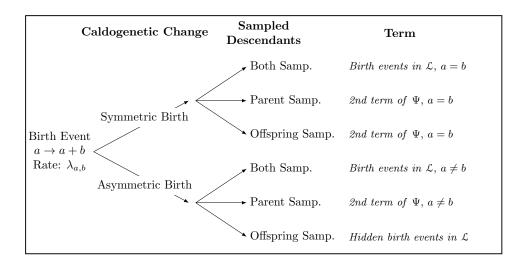


Figure S1: The six possible birth events in the sampled tree and how each in included in the likelihood given by equation (S19).

Tables

346

Vari- able	Definition
A	The number of discrete lineage types.
	The time-varying rate at which lineage of type $a$ gives rise to daughter lineages of
$\lambda_{a,\{b,c\}}(\tau)$	type $b$ and $c$ at time $\tau$ .
$\mu_a(\tau)$	The time varying rate at which lineages of type $a$ go extinct at time $\tau$
$\psi_a(\tau)$	The time varying rate at which lineages of type $a$ are sampled.
$r_a(\tau)$	The time varying probability that upon sampling lineages are removed.
$\gamma_{a,b}(\tau)$	The time varying rate at which lineages of type $a$ transition to type $b$ .
$\rho_a$	The probability a lineage of type $a$ alive at the present day is sampled.
$N_a$	The number of lineages of type $a$ sampled at the present day.
$\vec{x}_{a,b,c}$	A vector of length $I_{a,b}$ giving the timing of speciation events where a lineage of
	type $a$ gives rise to a daughter of type $b$ .
$ec{y}_a$	A vector of length $J_a$ giving the times at which lineages of type $a$ are sampled for
	which there are no sampled descendants (sampled tips).
$ec{z}_a$	A vector of length $K_a$ giving the times at which lineages of type $a$ are sampled for
$\sim a$	which there <i>are</i> sampled descendants (sampled ancestors).
$ec{w}_{a,b}$	A vector of length $L_{a,b}$ giving the times at which lineages of type $a$ mutate to
$\bigcup_{a,b} \omega_{a,b}$	lineages of type $b$ .
c	The state at the origin (stem node) of the phylogeny.
T	The tree topology.
C	The colouring of the tree.
$g_e(a,  au)$	The probability an edge $e$ of type $a$ alive at time $\tau$ gives rise to the subsequently
$g_e(a, r)$	observed phylogeny
$E_a(\tau)$	The probability a lineage of type $a$ alive at time $\tau$ has no sampled descendants
	between time $\tau$ and the present day.
$\iota_{a,b}$	An indicator variable with value 2 if $a = b$ and 1 otherwise.
$\mathfrak{C}_y$	The known tip states
$\mathcal{C}_z$	The known states of sampled ancestors
$\Theta_{\mathrm{MBDS}}$	The parameters of the multi-type BDS model
SSE Algor	
$\pi_c$	The probability the stem node has state $c$
$D_{N,a}( au)$	The probability a "topological" edge $N$ of type $a$ at time $\tau$ gives rise to the
	subsequently observed phylogeny.

Table S4: Notation for the general multi-type model. Throughout t represents time moving forward from the origin (t=0) of the tree to the tips (t=T) and  $\tau$  represents time moving backward from the tips  $(\tau=0)$  to the origin  $(\tau=T)$ . The indices a,b denote lineage types and can take on integer values between 1 and A.