

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/301856822>

Lateral transfer in Stochastic Dollo models

Article in *The Annals of Applied Statistics* · January 2016

CITATIONS

7

READS

86

2 authors, including:



Geoff Nicholls

University of Oxford

70 PUBLICATIONS 1,789 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Phylogenetics and Languages - lexical trait data [View project](#)



Applied Bayesian Statistics [View project](#)

Lateral transfer in Stochastic Dollo models*

Luke J. Kelly[†] Geoff K. Nicholls

Department of Statistics, University of Oxford, United Kingdom

Abstract

Lateral transfer, a process whereby species exchange evolutionary traits through non-ancestral relationships, is a frequent source of model misspecification when inferring species ancestries. Lateral transfer obscures the phylogenetic signal in the data as the histories of affected traits are mosaics of the species phylogeny. We control for the effect of lateral transfer in a Stochastic Dollo model and a Bayesian setting. Our likelihood is highly intractable as parameters are given by the solution of a sequence of large systems of differential equations representing the expected evolution of traits along the a tree. We illustrate our method on a data set of lexical traits in Eastern Polynesian languages and obtain an improved fit over the corresponding model without lateral transfer.

1 Introduction

Complex evolutionary traits provide a basis for inferring the ancestry of a set of taxa. These traits may derive from sequences such as DNA or an unordered set of morphological characters, for example. When species evolve in isolation, it is reasonable to assume that traits are passed vertically from one generation to the next through ancestral relationships. The shared ancestry of the taxa may be described by a phylogenetic tree whose branches represent evolving species and nodes speciation events. We are concerned with inferring the phylogeny of taxa which also evolve through *lateral transfer*. Lateral transfer, such as *horizontal gene transfer* in biology or *borrowing* in linguistics, is one of the processes driving the evolution of populations whereby species acquire traits through non-vertical relationships.

In the absence of lateral transfer, individual trait histories are compatible with the species phylogeny and there are many statistical methods to infer phylogenies in this setting. When lateral transfer occurs, the histories of affected traits are a mosaic of the species phylogeny and while tree-like, may conflict with the overall phylogeny. This obscures the phylogenetic signal of the branching events and models based solely on vertical inheritance are misspecified in this setting. In our experience, this model error can result in overly high levels of confidence in poorly fitting trees. This article develops a fully model-based Bayesian method which explicitly accounts for lateral transfer in phylogeny reconstruction.

In this paper, we analyse a data set of lexical traits in Eastern Polynesian languages. There have been many phylogenetic studies of language families, some tracing the phylogenies of the languages themselves (Gray and Atkinson, 2003; Nicholls and Gray, 2008;

*Corresponding author: Luke J. Kelly, kelly@stats.ox.ac.uk.

[†]Supported in part by the St John's College and Engineering and Physical Sciences Research Council partnership award EP/J500495/1.

Ryder and Nicholls, 2011; Chang et al., 2015) and others the movements of the peoples who spoke them (Gray et al., 2009; Bouckaert et al., 2012). Lateral transfer is a frequent occurrence in language diversification (Greenhill et al., 2009), yet researchers typically discard known-transferred traits and fit a vertical transfer-based model to the remainder (Gray and Atkinson; Bouckaert et al.; and many others). This is problematic as recently transferred traits are more readily identified so the discarded traits may not represent a random thinning of the data.

Lathrop (1982) and Pickrell and Pritchard (2012) propose methods to infer the order of population splits and a finite number of instantaneous *hybridisation* events in allele frequency data. The authors describe how to test the significance of the inferred hybridisation events under their respective models. Patterson et al. (2012) review a number of tests for admixture in allele frequency data. Similarly, there are many methods which test for lateral transfer in sequence data by comparing gene trees to a species tree constructed *a priori* (Daubin et al., 2002; Beiko and Hamilton, 2006; Abby et al., 2010). These methods are not model based, however. This is also the case with *implicit* phylogenetic networks. Internal nodes in these networks accommodate incompatibilities in the data with the assumption of an underlying species tree, but do not necessarily represent the evolutionary history of the taxa (Huson and Bryant, 2006; Oldman et al., 2016). Kubatko (2009) uses the multispecies coalescent model (Rannala and Yang, 2003) to perform model selection on an *explicit* network with a fixed number of hybridisation events while Wen et al. (2016) perform Bayesian inference on a variable number of reticulation nodes. Again, both of these methods require the input gene trees to be inferred in advance.

Szöllősi et al. (2012) propose a model-based approach to inferring phylogenies which incorporates lateral transfer. In their model, they discretise time on the tree so as to limit the number of transfer events which may occur. From a set of input gene trees inferred *a priori*, they seek the species tree which maximises the likelihood under the model. They do not incorporate a molecular clock so their method returns a time ordering of the internal nodes. Sjöstrand et al. (2014) perform approximate Bayesian inference under a similar model, but in this case the species tree is fixed and they estimate the gene trees from sequence data.

Of particular interest to us is the *Stochastic Dollo* (SD) model for unordered sets of binary traits proposed by Nicholls and Gray and extended by Alekseyenko et al. (2008) for multiple character states and Ryder and Nicholls for missing data and rate heterogeneity. The SD model posits a birth-death process of traits along the branches of the tree. The set of traits present in a parent branch is copied into the offspring lineages at a speciation event. The basic process respects *Dollo parsimony*: each trait is born exactly once, and once extinct, remains so. Simulation studies of the SD model have shown that topology estimates are robust to *moderate* levels of lateral transfer when the underlying topology is balanced but the root time tends to be biased towards the present (Nicholls and Gray; Greenhill et al.; Ryder and Nicholls).

Nicholls and Gray describe how to simulate lateral transfer in the basic SD model. Each species randomly acquires copies of traits possessed by other contemporary species in a similar manner to the model considered by Roch and Snir (2013). We perform exact likelihood-based inference under this model. Our lateral transfer process is described in Section 3. We do not attempt to model processes such as *incomplete lineage sorting*, hybridisation or gene *introgression* directly. While our process can generate the trait histories which arise in these processes, it also generates many others and we recommend further case-specific modelling. We do not attempt to first infer trait trees then reconcile them to form a species tree; rather we integrate over all possible trait histories on a given

species tree under our model. By working in a continuous-time setting, we are able to infer the timing of speciation events in contrast to Szöllősi et al.. We use Markov chain Monte Carlo techniques to sample from the posterior distribution of trees and parameters under the model. Finally we assess the goodness-of-fit of our approach and perform model selection using Bayes factors.

To summarise, we build a detailed *ab initio* model of trait and tree dynamics which fully describe the data observation process. In building this model, we do not compromise the model to make it easier to fit nor do we make approximations at the inference stage. The price we pay is a massive integration over the unobserved trait histories, but our computational methods are up to this integration for moderately sized problems at least. In looking for competing methods, we focus on methods which infer dated trees, can quantify the uncertainty in their estimates and perform exact inference or use explicitly quantified approximations. There are no obvious benchmarks among the model-based inference schemes discussed above for the lexical trait data we are attempting to model. The SD model is a special case of our model and therefore a natural basis for assessing whether the effect of controlling for lateral transfer outweighs the increase in the computational cost of performing inference. The remainder of this article is arranged as follows: in Section 2, we describe the format of the data which motivates the model introduced in Section 3; the likelihood calculation is explained in Section 4 and some model extensions in Section 5; we discuss our inference method in Section 6 and tests to validate our implementation of it in Section 7; finally, in Section 8, we illustrate our model on a data set of lexical traits in Eastern Polynesian languages.

2 Homologous trait data

Homologous traits are derived from a common ancestral trait through a combination of vertical inheritance and lateral transfer events and are assigned a unique common label from the set of trait labels, \mathbb{Z} . We record the status of trait h in taxon i as

$$d_i^h = \begin{cases} 0, & \text{trait } h \text{ is absent in taxon } i, \\ 1, & \text{trait } h \text{ is present in taxon } i, \\ ?, & \text{the status of trait } h \text{ in taxon } i \text{ is unknown.} \end{cases}$$

We denote by \mathbf{D} the array recording the status of each trait across the observed taxa. A column \mathbf{d}^h of \mathbf{D} is a *site-pattern* recording the status of trait h across the taxa. These patterns shall form the basis of our model.

In the analysis in Section 8, each trait is a word in one of 210 meaning categories and each taxon is an Eastern Polynesian language. For example, the Maori and Hawaiian words for *woman* and *wife*, both *wahine*, are derived from a common ancestor h , say, so $d_{\text{Maori}}^h = d_{\text{Hawaiian}}^h = 1$. On the other hand, the Maori word for *mother*, *whaea*, is not related to its Hawaiian counterpart, *makuahine* so we record a 0 in the corresponding entry in the data array, and vice versa.

3 Generative model

A branching process on sets of traits determines the phylogeny of the observed taxa. Each set represents an evolving species. A branching event on the tree corresponds to a speciation event and a leaf represents an observed taxon. The set contents diversify according to a trait process, and finally the status of each trait is recorded in the taxa.

Figure 1 depicts a realisation of the model and the history of a single trait which, due to two lateral transfer events, bears little resemblance to the underlying phylogeny.

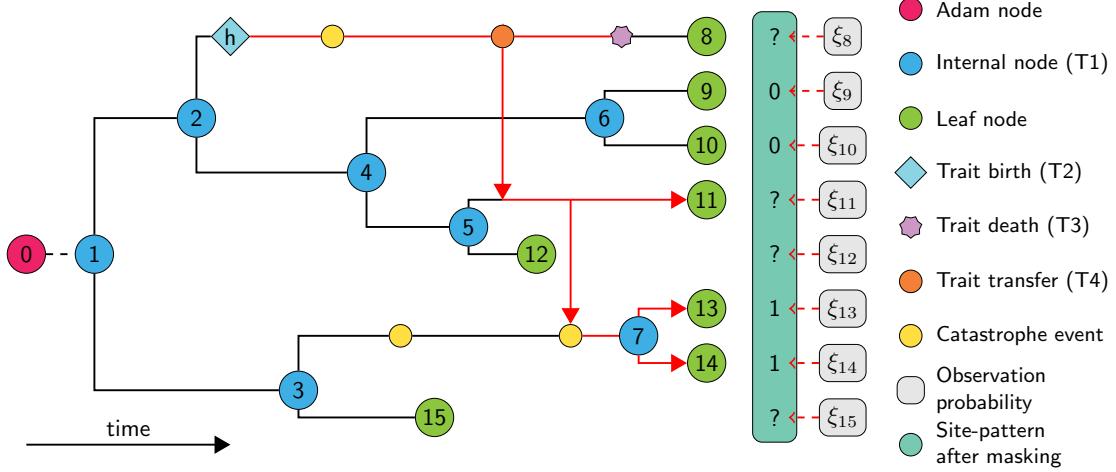


Figure 1: Illustration of the Stochastic Dollo with Lateral Transfer model. Catastrophes, missing data and offset leaves are introduced in Section 5.

We first define our model and inference method in terms of binary patterns of trait presence and absence in taxa which are observed simultaneously. It is then straightforward to extend the model to more complex scenarios.

A rooted phylogenetic tree $g = (V, E, T)$ on L leaves is a connected acyclic graph with node set $V = \{0, 1, \dots, 2L - 1\}$, directed edge set $E \subset V \times V$ and node times $T \in \{-\infty\} \times \mathbb{R}_{\leq 0}^{2L-1}$. The node set V comprises one *Adam* node labelled 0 of degree 1, the internal nodes $V_A = \{1, 2, \dots, L - 1\}$ of degree 3 and the leaf nodes $V_L = \{L, L + 1, \dots, 2L - 1\}$ of degree 1. Each node $i \in V$ is assigned a time $t_i \in T$ denoting when the corresponding event occurred relative to the current time, 0. For convenience, we label the internal nodes in such a way that t_1, \dots, t_{L-1} is a strictly increasing sequence of node times.

Edges represent evolving species and are directed forwards in time. We label each edge by its offspring node so if $\text{pa}(i)$ denotes the parent of node $i \in V \setminus \{0\}$, edge i runs from node $\text{pa}(i)$ at time $t_{\text{pa}(i)}$ to i at time t_i . We assume the Adam node arose at time $t_0 = -\infty$ so a branch of infinite length connects it to the *root* node 1 at time t_1 . The leaves are observed at time 0. At time t , there are $L^{(t)}$ species labelled $\mathbf{k}^{(t)} = (i \in E : t_{\text{pa}(i)} \leq t < t_i)$. For example, after the speciation event at time t_2 in Figure 1, there are $L^{(t_2)} = 3$ species labelled $\mathbf{k}^{(t_2)} = (8, 4, 3)$.

Letting $H_i(t) \subset \mathbb{Z}$ denote the set of traits possessed by species $i \in \mathbf{k}^{(t)}$ at time t , we now define four properties of the set-valued evolutionary process $H(t) = \{H_i(t) : i \in \mathbf{k}^{(t)}\}$.

Property T1 (Set branching event). At a speciation event, the traits present in the parent are copied into the offspring. Species $i \in \mathbf{k}^{(t_i^-)}$ branches at time t_i and is replaced by two identical offspring, j and $k \in \mathbf{k}^{(t_i)}$,

$$H_j(t_i) \leftarrow H_i(t_i^-), \\ H_k(t_i) \leftarrow H_i(t_i^-),$$

where t_i^- denotes the time just before the branching event.

Property T2 (Trait birth). New traits are born at rate λ in each extant species. If trait $h \in \mathbb{Z}$ is born in species i at time t ,

$$H_i(t) \leftarrow H_i(t^-) \cup \{h\}.$$

Property T3 (Trait death). A species kills off each trait it possesses independently at rate μ . If trait $h \in H_i(t^-)$ in species i dies at time t ,

$$H_i(t) \leftarrow H_i(t^-) \setminus \{h\}.$$

Property T4 (Lateral trait transfer). A species acquires a copy of a trait by lateral transfer at rate β scaled by the fraction of species which already possess it. If species i acquires a copy of trait $h \in \mathcal{H}^{(t^-)} = \bigcup_{i \in \mathbf{k}^{(t^-)}} H_i(t^-)$ at time t ,

$$H_i(t) \leftarrow H_i(t^-) \cup \{h\}.$$

Clearly if $h \in H_i(t^-)$ already then the transfer event has no effect.

Starting from a single set $H(-\infty) = \{\emptyset\}$, the process $H(t)$ evolves through a combination of branching (T1) and trait (T2–4) events to yield the diverse set of taxa $H(0) = \{H_i(0) : i \in V_L\}$ observed at time 0. When the lateral transfer rate $\beta = 0$, we recover the binary Stochastic Dollo process of [Nicholls and Gray](#).

4 Likelihood calculation

Traits displaying the same pattern of presence and absence across the leaves are exchangeable and their number is a Poisson random variable. We cannot compute the likelihood in closed form but now describe a representation of the process which allows us to evaluate it numerically.

4.1 Pattern evolution

If we cut through the tree at time t , each trait in $\mathcal{H}^{(t)}$ displays a *pattern* of presence and absence across the $L^{(t)}$ extant species $\mathbf{k}^{(t)} = (k_i^{(t)} : i \in [L^{(t)}])$, where $[L^{(t)}] = (1, \dots, L^{(t)})$. These patterns evolve over time as new branches arise and instances of h die and transfer. The pattern displayed by trait h at time t is $\mathbf{p}^h(t) = (p_i^h(t) : i \in [L^{(t)}])$, where

$$p_i^h(t) = \begin{cases} 1, & h \in H_{k_i^{(t)}}(t), \\ 0, & \text{otherwise,} \end{cases}$$

indicates the presence or absence of trait h on lineage $k_i^{(t)}$ at time t .

The space of binary patterns of trait presence and absence across $L^{(t)}$ lineages is $\mathcal{P}^{(t)} = \{0, 1\}^{L^{(t)}} \setminus \{\mathbf{0}\}$, where $\mathbf{0}$ is an $L^{(t)}$ -tuple of zeros notionally representing patterns displayed by traits in $\mathbb{Z} \setminus \mathcal{H}^{(t)}$. There are $N_{\mathbf{p}}(t) = |\{h \in \mathcal{H}^{(t)} : \mathbf{p}^h(t) = \mathbf{p}\}|$ traits displaying pattern $\mathbf{p} \in \mathcal{P}^{(t)}$ at time t . The dynamics of the pattern frequency process $\mathbf{N}(t) = (N_{\mathbf{p}}(t) : \mathbf{p} \in \mathcal{P}^{(t)})$ follow directly from Properties T1–4 of the trait process with the additional structure of $\mathbf{k}^{(t)}$.

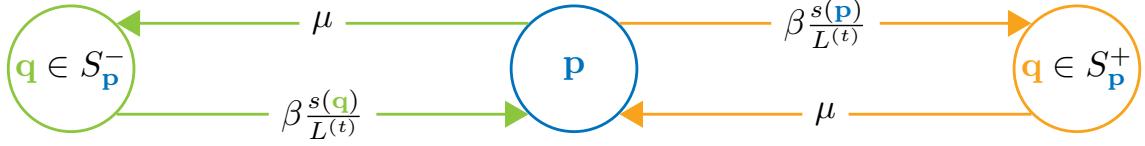


Figure 2: Transition rates between pattern states \mathbf{p} and $\mathbf{q} \in S_{\mathbf{p}}^- \cup S_{\mathbf{p}}^+$.

4.1.1 Patterns at branching events

At a branching event, patterns increase in length and the space of patterns expands to accommodate the new patterns which traits may display. The tuple $\mathbf{k}^{(t)}$ of branch labels is consistent across speciation events in the sense that when lineage $j = k_i^{(t_j^-)}$ branches at time t_j , the branch labels are

$$\mathbf{k}^{(t_j)} = \left(k_1^{(t_j^-)}, \dots, k_{i-1}^{(t_j^-)}, k_i^{(t_j)}, k_{i+1}^{(t_j)}, k_{i+1}^{(t_j^-)}, \dots, k_{L^{(t_j)}}^{(t_j^-)} \right),$$

where species $k_i^{(t_j)}$ and $k_{i+1}^{(t_j)}$ are the offspring of species $j = k_i^{(t_j^-)}$ (T1). It follows that each trait $h \in \mathcal{H}^{(t_j)}$ displays a pattern $\mathbf{p}^h(t_j)$ with entries $p_i^h(t_j) = p_{i+1}^h(t_j) \leftarrow p_i^h(t_j^-)$. For example, in Figure 1, at time t_4

$$\begin{aligned} \mathbf{k}^{(t_4^-)} &= (8, 4, 7, 15), & \mathbf{k}^{(t_4)} &= (8, 6, 5, 7, 15), \\ \mathbf{p}^h(t_4^-) &= (1, 0, 0, 0), & \mathbf{p}^h(t_4) &= (1, 0, 0, 0, 0). \end{aligned}$$

A pattern $\mathbf{p} \in \mathcal{P}^{(t_j)}$ with entries $p_i = p_{i+1}$ is consistent with the branching event on lineage $k_i^{(t_j^-)}$ as it may be formed by duplicating the i th entries of a pattern in $\mathcal{P}^{(t_j^-)}$. On the other hand, the trait process cannot generate a pattern $\mathbf{p} \in \mathcal{P}^{(t)}$ with $p_i \neq p_{i+1}$ at time t_j by definition (T1). We denote by $\mathbf{T}^{(j)}$ the operation which initialises the pattern frequencies $\mathbf{N}(t_j)$ with entries of $\mathbf{N}(t_j^-)$ for patterns consistent with the branching event and zeros otherwise. We return to this initialisation operation when computing expected pattern frequencies in Section 4.2.

4.1.2 Patterns between branching events

In order to formally describe the Markovian evolution of the pattern frequencies $\mathbf{N}(t)$ between branching events, we first define how patterns are related. The Hamming distance between patterns \mathbf{p} and $\mathbf{q} \in \mathcal{P}^{(t)}$ is $d(\mathbf{p}, \mathbf{q}) = |\{i \in [L^{(t)}] : p_i \neq q_i\}|$ and $s(\mathbf{p}) = d(\mathbf{p}, \mathbf{0})$ is the Hamming weight of \mathbf{p} . The pattern $\mathbf{p}^h(t) = \mathbf{p}$ displayed by trait h communicates with patterns in the sets

$$\begin{aligned} S_{\mathbf{p}}^- &= \{\mathbf{q} \in \mathcal{P}^{(t)} : s(\mathbf{q}) = s(\mathbf{p}) - 1, d(\mathbf{p}, \mathbf{q}) = 1\}, \\ S_{\mathbf{p}}^+ &= \{\mathbf{q} \in \mathcal{P}^{(t)} : s(\mathbf{q}) = s(\mathbf{p}) + 1, d(\mathbf{p}, \mathbf{q}) = 1\}, \end{aligned}$$

which may be formed from \mathbf{p} through a single death (T3) or transfer (T4) event amongst traits labelled h . Figure 3 describes the rates at which a trait displaying pattern $\mathbf{p}^h(t) = \mathbf{p}$ with $s(\mathbf{p}) > 1$ transition to state $\mathbf{q} \in S_{\mathbf{p}}^- \cup S_{\mathbf{p}}^+$ and back again. If $s(\mathbf{p}) = 1$, a death event would result in trait h becoming extinct and new traits displaying pattern \mathbf{p} arise at rate λ through birth events.

4.2 Expected pattern frequencies

Traits displaying a pattern $\mathbf{p} \in \mathcal{P}^{(t)}$ evolve independently of each other. By summing over the rates in Figure 2 for each trait displaying a given pattern, on a short interval of length dt between branching events, we have

$$\begin{aligned} \mathbb{P}[N_{\mathbf{p}}(t + dt) - N_{\mathbf{p}}(t) = k | g, \lambda, \mu, \beta] \\ = \begin{cases} s(\mathbf{p})N_{\mathbf{p}}(t) \left[\mu + \beta \left(1 - \frac{s(\mathbf{p})}{L^{(t)}} \right) \right] dt + o(dt), & k = -1, \\ \left[\lambda \mathbf{1}_{\{s(\mathbf{p})=1\}} + \beta \sum_{\mathbf{q} \in S_{\mathbf{p}}^-} \frac{s(\mathbf{q})}{L^{(t)}} N_{\mathbf{q}}(t) \right. \\ \left. + \mu \sum_{\mathbf{q} \in S_{\mathbf{p}}^+} N_{\mathbf{q}}(t) \right] dt + o(dt), & k = 1. \end{cases} \end{aligned} \quad (1)$$

We cannot compute these transition probabilities in practice as we only observe the terminal pattern frequencies, $\mathbf{N}(0)$. From Equation 1 we derive that $x_{\mathbf{p}}(t; g, \lambda, \mu, \beta) = \mathbb{E}[N_{\mathbf{p}}(t) | g, \lambda, \mu, \beta]$, the expected number of traits in $\mathcal{H}^{(t)}$ displaying pattern $\mathbf{p} \in \mathcal{P}^{(t)}$ at time t , evolves across the interval according to the following differential equation:

$$\begin{aligned} \dot{x}_{\mathbf{p}}(t) &= \lim_{dt \rightarrow 0} \frac{\mathbb{E}[N_{\mathbf{p}}(t + dt) - N_{\mathbf{p}}(t) | g, \lambda, \mu, \beta]}{dt} \\ &= -s(\mathbf{p})x_{\mathbf{p}}(t) \left[\mu + \beta \left(1 - \frac{s(\mathbf{p})}{L^{(t)}} \right) \right] + \lambda \mathbf{1}_{\{s(\mathbf{p})=1\}} \\ &\quad + \beta \sum_{\mathbf{q} \in S_{\mathbf{p}}^-} \frac{s(\mathbf{q})}{L^{(t)}} x_{\mathbf{q}}(t) + \mu \sum_{\mathbf{q} \in S_{\mathbf{p}}^+} x_{\mathbf{q}}(t). \end{aligned} \quad (2)$$

There are $|\mathcal{P}^{(t)}| = 2^{L^{(t)}} - 1$ coupled differential equations (2) describing the expected evolution of the pattern frequencies $\mathbf{N}(t)$. We may write these equations as $\dot{\mathbf{x}}(t) = \mathbf{A}^{(t)}\mathbf{x}(t) + \mathbf{b}^{(t)}$ where: $\mathbf{x}(t) = (x_{\mathbf{p}}(t) : \mathbf{p} \in \mathcal{P}^{(t)})$ is the vector of expected pattern frequencies; the sparse matrix $\mathbf{A}^{(t)}$ and vector $\mathbf{b}^{(t)}$ respectively describe the flow between patterns from trait death and transfer events and into patterns from trait births events.

The process $\mathbf{N}(t)$ is in equilibrium just before the first branching event by construction so $\mathbf{x}(t_1^-) = x_1(t_1^-) = \lambda/\mu$. With this initial condition at the root, we can write the expected pattern frequencies at the leaves $\mathbf{x}(0)$ recursively as a sequence of initial value problems between branching events:

$$\dot{\mathbf{x}}(t) = \mathbf{A}^{(t)}\mathbf{x}(t) + \mathbf{b}^{(t)} \quad \text{for } t \in [t_i, t_{i+1}) \quad \text{where } \mathbf{x}(t_i) = \mathbf{T}^{(i)}\mathbf{x}(t_i^-), \quad (3)$$

where we recall the pattern frequency initialisation operator $\mathbf{T}^{(i)}$ defined in Section 4.1.1 that propagates $\mathbf{N}(t)$ and $\mathbf{x}(t)$ across the i th branching event. We illustrate this procedure graphically in Figure 3.

4.3 Distribution of pattern frequencies

Theorem 1 describes the distribution of the pattern frequencies under the model. A proof is contained in Appendix 1.

Theorem 1 (Binary data distribution). *The pattern frequencies $\mathbf{N}(t)$ is a vector of independent Poisson variables with corresponding rate parameters $\mathbf{x}(t; g, \lambda, \mu, \beta)$ given by the solution of the initial value problems in Equation 3.*

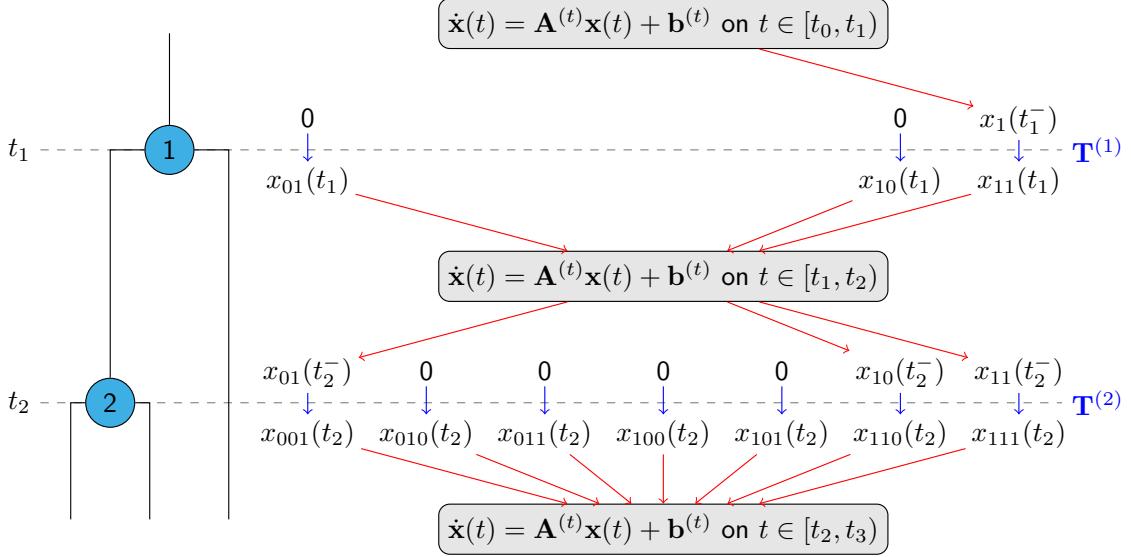


Figure 3: Computing the expected pattern frequencies $\mathbf{x}(t)$ as a sequence of initial value problems (3) on the tree. The initialisation operation $\mathbf{T}^{(i)}\mathbf{x}(t_i^-) = \mathbf{x}(t_i)$ from Section 4.1.1 provides the initial condition at the start of the i th interval between branching events.

5 Model extensions

It is straightforward to extend the model and likelihood calculation to allow for rate variation, missing data, offset leaves and the systematic removal of patterns from the data.

5.1 Rate heterogeneity

We introduce spikes of activity in the form of *catastrophes* (Ryder and Nicholls). Catastrophes, illustrated in Figure 1, occur at rate ρ along each branch of the tree. A catastrophe advances the trait process along a branch by $\delta = -\mu^{-1} \log(1 - \kappa)$ units of time, where the parameter $\kappa \in [0.25, 1]$ controls the severity of the catastrophe. The minimum catastrophe duration of $-\mu^{-1} \log(0.75)$ years is to avoid scenarios where the model attempts to over-explain the variation in the data with catastrophes. We could allow individual catastrophe durations to vary but do not pursue that approach here.

A branch may acquire traits through birth and transfer events and lose traits to death events during a catastrophe. We account for a catastrophe at time t on branch $k_i^{(t)}$ in the expected pattern frequency calculation (3) with

$$\begin{aligned}
 x_{\mathbf{p}}(t) &\leftarrow e^{-\mu\delta} x_{\mathbf{p}}(t^-) + (1 - e^{-\mu\delta}) \frac{\lambda}{\mu} & \mathbf{p} \in \mathcal{P}^{(t)}, \\
 & s(\mathbf{p}) = 1, p_i = 1, \\
 \begin{bmatrix} x_{\mathbf{q}}(t) \\ x_{\mathbf{r}}(t) \end{bmatrix} &\leftarrow \exp \left[\begin{pmatrix} -\beta \frac{s(\mathbf{q})}{L^{(t)}} & \mu \\ \beta \frac{s(\mathbf{q})}{L^{(t)}} & -\mu \end{pmatrix} \delta \right] \begin{bmatrix} x_{\mathbf{q}}(t^-) \\ x_{\mathbf{r}}(t^-) \end{bmatrix} & \mathbf{q}, \mathbf{r} \in \mathcal{P}^{(t)}, d(\mathbf{q}, \mathbf{r}) = 1 \\
 & q_i = 0, r_i = 1,
 \end{aligned}$$

where we exploit the property that each pattern communicates with at most one other at a catastrophe. The trait process at a catastrophe is equivalent to thinning the overall trait process to events on a single branch. It is straightforward then to adapt the proof of Theorem 1 to show that this update step correctly describes the effect of catastrophes on the pattern process.

5.2 Missing data

Following [Ryder and Nicholls](#), the true binary state of trait h at taxon $i \in V_L$ is recorded with probability $\xi_i = \mathbb{P}(d_i^h \in \{0, 1\})$ independently of the other taxa. Let $\Xi = (\xi_i : i \in V_L)$. The space of observable site-patterns with missing data across the L taxa at time 0 is $\mathcal{Q} = \{0, 1, ?\}^L \setminus \mathbf{0}$. The frequency of traits displaying pattern $\mathbf{q} \in \mathcal{Q}$ is a Poisson random variable with mean

$$\mathbf{x}_q(0; g, \lambda, \mu, \beta, \Xi) = \sum_{\mathbf{p} \in u(\mathbf{q})} x_{\mathbf{p}}(0; g, \lambda, \mu, \beta) \prod_{i=1}^L \xi_{k_i^{(0)}}^{\mathbf{1}_{\{q_i \in \{0, 1\}\}}} \left(1 - \xi_{k_i^{(0)}}\right)^{\mathbf{1}_{\{q_i = ?\}}}$$

where $u(\mathbf{q}) = \{\mathbf{p} \in \mathcal{P}^{(0)} : p_i = q_i \text{ if } q_i \neq ?, i \in [L]\}$ is the set of binary patterns consistent with \mathbf{q} before masking. This result, illustrated in [Appendix 2](#), follows from the restriction and superposition properties of Poisson processes.

5.3 Non-isochronous data

Data are non-isochronous when the taxa are not sampled simultaneously and appear as *offset* leaves in the phylogeny; nodes 12 and 15 in Figure 1, for example. Similar to catastrophes, the trait process is frozen on offset leaves. A pattern may now only communicate with the patterns identical to it on the extinct lineages and differing by one entry on the extant lineages.

The $L^{(t)} = |H(t)|$ extinct and evolving lineages at time t are labelled $\mathbf{k}^{(t)} = (i \in E : t_{\text{pa}(i)} \leq t < t_i \mathbf{1}_{\{i \in V_L\}})$. The Hamming distance between patterns \mathbf{p} and $\mathbf{q} \in \mathcal{P}^{(t)}$ across extant lineages only is $\hat{d}(\mathbf{p}, \mathbf{q}) = |\{i \in [L^{(t)}] : p_i \neq q_i, t < t_{k_i^{(t)}}\}|$ and the Hamming weight of \mathbf{p} is $\hat{s}(\mathbf{p}) = \hat{d}(\mathbf{p}, \mathbf{0})$. Recalling $S_{\mathbf{p}}^-$ and $S_{\mathbf{p}}^+$ from Section 4.1.2, pattern $\mathbf{p} \in \mathcal{P}^{(t)}$ communicates with patterns in the sets

$$\begin{aligned} \hat{S}_{\mathbf{p}}^- &= \{\mathbf{q} \in S_{\mathbf{p}}^- : \hat{s}(\mathbf{q}) = \hat{s}(\mathbf{p}) - 1, \hat{d}(\mathbf{p}, \mathbf{q}) = 1\}, \\ \hat{S}_{\mathbf{p}}^+ &= \{\mathbf{q} \in S_{\mathbf{p}}^+ : \hat{s}(\mathbf{q}) = \hat{s}(\mathbf{p}) + 1, \hat{d}(\mathbf{p}, \mathbf{q}) = 1\}, \end{aligned}$$

and its expected frequency evolves across the interval as

$$\begin{aligned} \dot{x}_{\mathbf{p}}(t) &= -\hat{s}(\mathbf{p})x_{\mathbf{p}}(t) \left[\mu + \beta \left(1 - \frac{\hat{s}(\mathbf{p})}{L^{(t)}} \right) \right] + \lambda \mathbf{1}_{\{s(\mathbf{p})=\hat{s}(\mathbf{p})=1\}} \\ &\quad + \beta \sum_{\mathbf{q} \in \hat{S}_{\mathbf{p}}^-} \frac{\hat{s}(\mathbf{q})}{L^{(t)}} x_{\mathbf{q}}(t) + \mu \sum_{\mathbf{q} \in \hat{S}_{\mathbf{p}}^+} x_{\mathbf{q}}(t). \end{aligned}$$

5.4 Data registration

Patterns which are uninformative or unreliable are typically removed from the data. For example, we discard all traits which are not marked present in at least one taxon. Given a registration rule R , which may be a composition of other simpler rules such as those in Table 1, we discard the columns in the data array \mathbf{D} not satisfying R , leaving $R(\mathbf{D})$, and restrict our analysis to patterns in $R(\mathcal{Q})$.

6 Bayesian inference

In order to estimate both node times and rate parameters, we calibrate the space Γ of rooted phylogenetic trees on L taxa with *clade constraints*. The constraint $\Gamma^{(0)} = (g \in \Gamma :$

Unregisterable traits	Unregisterable patterns $\mathcal{Q} \setminus R(\mathcal{Q})$
Traits observed in j taxa or fewer	$\{\mathbf{q} \in \mathcal{Q} : \{i \in [L] : q_i = 1\} \leq j\}$
Traits observed in j or more taxa	$\{\mathbf{q} \in \mathcal{Q} : \{i \in [L] : q_i = 1\} \geq j\}$
Traits potentially present in j taxa or greater	$\{\mathbf{q} \in \mathcal{Q} : \{i \in [L] : q_i \neq 0\} \geq j\}$
Traits absent in taxon $k_i^{(0)}$	$\{\mathbf{q} \in \mathcal{Q} : q_i = 0\}$

Table 1: Registration rules of [Alekseyenko et al.](#) and [Ryder and Nicholls](#).

Parameter	Prior	Reasoning
Trait birth rate	$\lambda \sim 1/\lambda$	Improper, scale invariant
Trait death rate	$\mu \sim \Gamma(10^{-3}, 10^{-3})$	Approximately $1/\mu$
Trait transfer rate	$\beta \sim \Gamma(10^{-3}, 10^{-3})$	Approximately $1/\beta$
Catastrophe rate	$\rho \sim \Gamma(1.5, 5 \times 10^3)$	$\mathbb{E}[\rho^{-1}] = 10^4$ years
Catastrophe severity	$\kappa \sim U[0.25, 1]$	$\mathbb{E}[\delta \mu] = \mu^{-1}[1 - \log(0.75)]$ years
Observation probabilities	$\Xi \sim U[0, 1]^L$	Independent, uniform

Table 2: Prior distributions on parameters in the Stochastic Dollo and Stochastic Dollo with Lateral Transfer models.

$\underline{t}_1 \leq t_1 < 0$) restricts the earliest admissible root time t_1 to \underline{t}_1 . Each additional constraint $\Gamma^{(c)}$ places either time or ancestry constraints on the remaining nodes. We denote by $\Gamma^C = \bigcap_c \Gamma^{(c)}$ the calibrated space of phylogenies satisfying the clade constraints.

[Nicholls and Ryder \(2011\)](#) describe a prior on trees with the property that the root time t_1 is marginally approximately uniform across a specified interval $[\underline{t}_1, \bar{t}_1] \subset \mathbb{R}_{\leq 0}$. For a given tree $g = (V, E, T, C)$, there are $Z(g)$ possible time orderings of the nodes amongst the admissible node times $T(g) = \{T' : (V, E, T', C) \in \Gamma^C\}$. For each node $i \in V$, $\underline{t}_i = \inf_{T \in T(g)} t_i$ and $\bar{t}_i = \sup_{T \in T(g)} t_i$ are the earliest and most recent times that i may achieve in an admissible tree with topology (V, E) . If $S(g) = \{i \in V_A : t_i = \underline{t}_1\}$ denotes the set of free internal nodes with times bounded below by \underline{t}_1 , the prior with density

$$f_G(g) \propto \frac{\mathbf{1}_{\{g \in \Gamma^C\}}}{Z(g)} \prod_{i \in S(g)} \frac{\underline{t}_1 - \bar{t}_i}{\underline{t}_1 - \bar{t}_i},$$

is approximately marginally uniform across topologies and root times provided $\underline{t}_1 \ll \min_{i \in V \setminus S} t_i$ ([Ryder and Nicholls](#)). Uniform priors on offset leaf times completes our prior specification on g . [Heled and Drummond \(2012\)](#) describe an exact method for computing uniform calibrated tree priors but we do not pursue that approach here. Table 2 lists the priors on the remaining parameters.

Inspecting the solution of the expected pattern frequency calculation (3) with initial condition $\mathbf{x}(\underline{t}_1^-) = \lambda/\mu$ at the root, we see that $\mathbf{x}(t; g, \lambda, \dots) = \lambda \mathbf{x}(t; g, 1, \dots)$. We can integrate λ out of the Poisson likelihood in Theorem 1 with respect to its prior in Table 2 above to obtain a multinomial likelihood whereby a pattern $\mathbf{p} \in R(\mathcal{Q})$ is observed with probability proportional to its expected frequency. Furthermore, we can integrate out the Gamma prior on the catastrophe rate ρ to obtain a Negative Binomial prior on the number of catastrophes $|C|$. We discuss these steps further in [Appendix 3](#). Let $n_{\mathbf{p}} = |\{h \in \mathcal{H}(0) : \mathbf{p} = \mathbf{d}^h \in R(\mathbf{D})\}|$ denote the frequency of traits in the registered data displaying pattern

$\mathbf{p} \in R(\mathcal{Q})$. Putting everything together, the posterior distribution is

$$\pi(g, \mu, \beta, \kappa, \Xi | R(\mathbf{D})) \propto f_G(g) f_M(\mu) f_B(\beta) \prod_{\mathbf{p} \in R(\mathcal{Q})} \left(\frac{x_{\mathbf{p}}}{\sum_{\mathbf{q} \in R(\mathcal{Q})} x_{\mathbf{q}}} \right)^{n_{\mathbf{p}}}, \quad (4)$$

where the expected pattern frequencies $\mathbf{x} \equiv \mathbf{x}(0; g, 1, \mu, \beta, \kappa, \Xi)$ (3) account for catastrophes, missing data and offset leaves, where necessary. This completes the specification of the Stochastic Dollo with Lateral Transfer (SDLT) model.

The posterior (4) is intractable but may be explored using standard Markov chain Monte Carlo (MCMC) sampling schemes for phylogenetic trees and Stochastic Dollo models ([Nicholls and Gray](#); [Ryder and Nicholls](#)). We describe the MCMC transition kernels for moves particular to the SDLT model in [Appendix 3](#). To assess convergence of the Markov chains, we monitor the sample autocorrelation functions of the parameters and log-likelihoods ([Geyer, 1992](#)).

Implementation

Code to implement the SDLT model in the software package [TraitLab](#) ([Nicholls et al., 2013](#)) is available from the authors.

7 Method testing

We compare the exact and empirical distributions of synthetic data to validate our implementation of the expected pattern frequency calculation (3). In addition, we test the identifiability of the SDLT model, its consistency with the SD model when the lateral transfer rate $\beta = 0$ and its robustness to a common form of model misspecification whereby recently transferred traits are discarded. In each case, we obtain a satisfactory fit to the data and recover the true parameters. Full details of these analyses are contained in [Appendix 4](#).

8 Applications

The order and timing of human settlement in Eastern Polynesia is a matter of debate. In the standard subgrouping of the Eastern Polynesian languages, Rapanui diverges first, followed by the split leading to the Marquesic (Hawaiian, Mangarevan, Marquesan) and Tahitic (Manihiki, Maori, Penrhyn, Rarotongan, Rurutuan, Tahitian, Tuamotuan) languages ([Marck, 2000](#)). This theory has recently been challenged in light of new linguistic and archaeological evidence. In an implicit phylogenetic network study of lexical traits, [Gray et al. \(2010\)](#) detect non-tree-like signals in the data and the Tahitic and Marquesic languages do not form clean clusters. From a meta-analysis of radiocarbon dates, [Wilmshurst et al. \(2011\)](#) claim that the islands of Eastern Polynesia were settled in two distinct phases: the Society Islands between 900 and 1000 years before the present (BP) and the remainder between 700 and 900 years BP. These dates are much later than previously thought ([Spriggs and Anderson, 1993](#)) and do not allow much time for the development of the Eastern Polynesian language subgroups. [Conte and Molle \(2014\)](#) present evidence of human settlement in the Marquesas Islands approximately 1100 years BP. On the basis of the above and further evidence of lateral transfer in primary source material, [Walworth \(2014\)](#) disputes Marquesic and Tahitic as distinct subgroups.

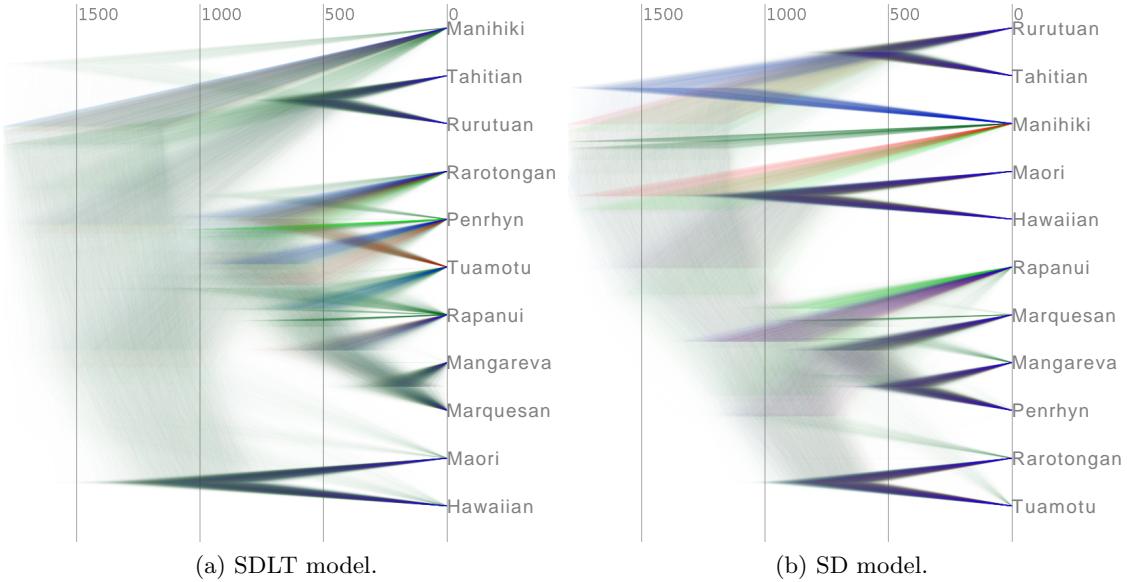


Figure 4: *DensiTree* (Bouckaert and Heled, 2014) plots of the marginal tree posterior under the SDLT and SD models fit to POLY-0. Time is in years before the present and the most frequently sampled topologies in each case are coloured blue, followed by red then green, with the remainder in dark green.

To add to this debate, we illustrate the SDLT model on lexical traits in eleven Eastern Polynesian languages drawn from the approximately 1200 languages in the Austronesian Basic Vocabulary Database (Greenhill et al., 2008). We compare our results with the SD model to highlight the effect of the laterally transfer traits in the data. The data is a subset of the Polynesian language data set analysed by Gray et al.. We analyse the 968 traits marked present in at least one of the eleven languages, hereafter referred to as POLY-0. The data are assumed isochronous. Consistent with Gray et al., the sole clade constraint limits the root of the tree to lie between 1150 and 1800 years BP.

We plot samples from the marginal tree posterior under the SDLT and SD models in Figure 4. We summarise these distributions with *majority rule consensus trees* in the Appendix 5. In agreement with Gray et al. and Walworth, the traditional subgroupings do not appear as subtrees in either model, nor does Rapanui form an outgroup. There is little evidence in the tree posteriors to support the claims of Wilmshurst et al., however, as the posterior distribution of the root time t_1 is approximately uniform across its predefined range.

The majority of the uncertainty under the SDLT model is in the topology of the subtree containing Rarotongan, Penrhyn, Tuamotu, Rapanui, Mangareva and Marquesan. This subtree also has 100% posterior support under the SD model but here most of the uncertainty is in the relationships further up the tree. We obtain the 95% highest posterior probability sets for the tree topologies using BEAST (Drummond et al., 2012). This set comprises 135 topologies for the SDLT model and 19 for the SD model. This level of confidence in relatively few topologies is a likely result of the SD model's misspecification on the laterally transferred traits.

The effect of the laterally transferred traits is evident again in the histograms of samples from the marginal posterior distributions of the death rate μ and relative transfer rate β/μ in Figure 5. The death rate is approximately 50% higher under the SD model as traits must be born further up the tree and killed off at a higher rate to explain the variation in

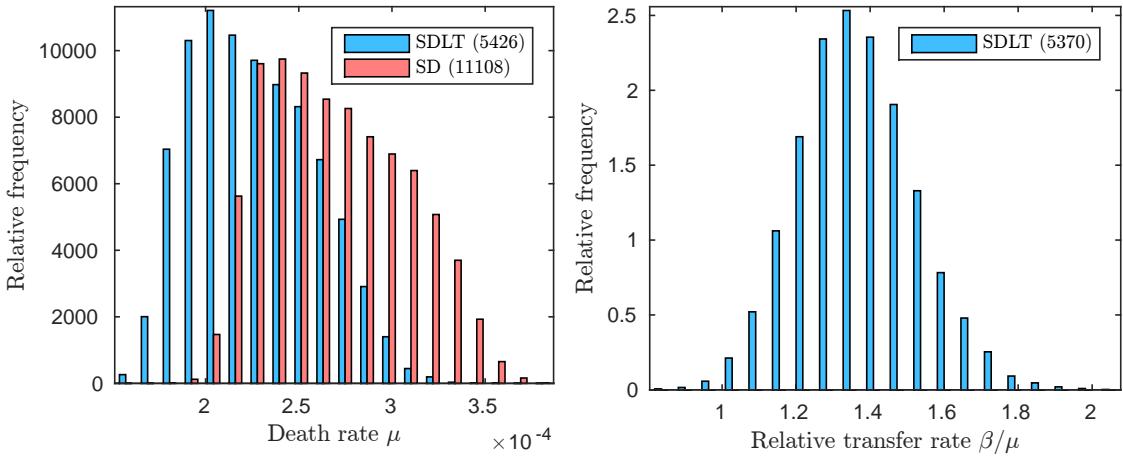


Figure 5: Marginal parameter posterior distributions under the SDLT and SD models fit to the Eastern Polynesian data set **POLY-0**. Effective sample sizes are in parentheses.

the data. The relative transfer rate is the expected number of times that an instance of a trait attempts to transfer before dying and has a Beta Prime prior with infinite mean and variance under the marginal priors on β and μ in Table 2. The posterior distribution in Figure 5 is well informed and centred on 1.35. In comparison, [Nicholls and Gray](#) and [Greenhill et al.](#) considered a relative transfer rate of 0.5 high. Histograms for the remaining parameters as well as the trace and autocorrelation plots used to diagnose convergence are contained in the [Appendix 5](#).

With the above concerns about the SD model in mind, we now assess the validity of our analyses with two tests of goodness-of-fit. For the first test, we relax each of the leaf constraints in turn and compute a Bayes factor comparing the relaxed and constrained models. The constraint $\Gamma^{(i)} = \{g \in \Gamma : t_i = 0\}$ fixes leaf $i \in V_L$ at time 0 and $\Gamma^{(i')} = \{g \in \Gamma : -10^4 \leq t_i \leq 10^3\}$ denotes its relaxation. We denote by $\Gamma^{C'}$ the calibrated space of phylogenies with $\Gamma^{(i)}$ replaced by $\Gamma^{(i')}$. Now, $\Gamma^C \subset \Gamma^{C'}$ so the Bayes factor

$$\begin{aligned} B_{i',i} &= \frac{\pi(R(\mathbf{D})|g \in \Gamma^{C'})}{\pi(R(\mathbf{D})|g \in \Gamma^C)} \\ &= \frac{\pi(R(\mathbf{D})|g \in \Gamma^{C'})}{\pi(R(\mathbf{D})|g \in \Gamma^C \cap \Gamma^{C'})} \\ &= \frac{\pi(g \in \Gamma^C | g \in \Gamma^{C'})}{\pi(g \in \Gamma^{C'} | R(\mathbf{D}), g \in \Gamma^{C'})}, \end{aligned} \quad (5)$$

a Savage–Dickey ratio of the marginal prior and posterior densities that the constraint is satisfied in the relaxed model. The marginal prior on the free leaf time is uniform across $[-10^4, 10^3]$. A large Bayes factor here therefore indicates a lack of support for the leaf constraint and is a sign of model misspecification. To compare the models on a given constraint, we simply compute the ratio of the corresponding Bayes factors.

We cannot compute the Savage–Dickey ratio (5) in closed form so in practice we estimate the densities by the corresponding proportion of sampled leaf times in the range $[-50, 50]$. We report log-Savage–Dickey ratios in Figure 6 and histograms of the marginal leaf ages in the [Appendix 5](#). The SD model does not support the constraints on Manihiki and Marquesan so we are unable to estimate the corresponding Bayes factors here. Comparing the two models on the remaining constraints, there are cases where the SDLT model outperforms the SD model and vice versa.

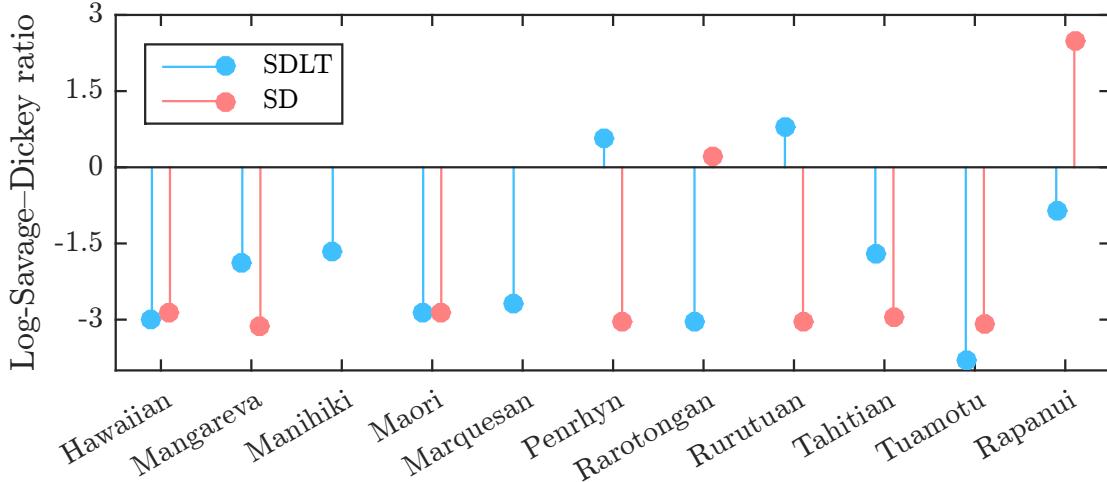


Figure 6: Bayes factors comparing the support for the leaf constraints used to fit the SDLT and SD models to POLY-0.

Training set	Test set	SDLT score	SD score	Log-Bayes factor
POLY-0/train	POLY-0/test	-3089.9	-3154.7	64.8
POLY-1/train	POLY-1/test	-1413.0	-1503.4	90.4

Table 3: Posterior predictive model assessment.

We assess the predictive performance of each model on a random splitting of the registered data $R(\mathbf{D})$ into evenly sized training and test sets labelled \mathbf{D}^{tr} and \mathbf{D}^{te} respectively. Following Madigan and Raftery (1994), we score each model by its log-posterior predictive probability, $\log \pi(\mathbf{D}^{\text{te}}|\mathbf{D}^{\text{tr}}) = \log \int \pi(\mathbf{D}^{\text{te}}|x)\pi(x|\mathbf{D}^{\text{tr}})dx$ where $x = (g, \mu, \beta, \kappa, \Xi)$. The difference in scores is a log-Bayes factor measuring the relative success of the models in predicting the test data (Kass and Raftery, 1995). The results in Table 3 support the superior fit of the SDLT model to POLY-0.

Traits marked present in a single language are often deemed unreliable and removed in the registration process. To address this concern, we repeat our analyses on the data set POLY-1 formed by discarding the *singleton* patterns from POLY-0. Singleton patterns play an important role in SDLT model inference as, although the outcome of the predictive model selection in Table 3 is unchanged, parameter credible intervals are affected.

9 Concluding remarks

Lateral transfer is an important problem but practitioners lack the statistical tools to perform fully likelihood-based inference for phylogenies in this setting. We address this issue with a novel model of species diversification which extends the Stochastic Dollo model for lateral transfer in trait presence/absence data. To our knowledge, our method is the first fully likelihood-based approach to control for lateral transfer in reconstructing a rooted phylogenetic tree. The second major contribution of this paper is the inference procedure whereby we integrate out the locations of the trait birth, death and transfer events using systems of differential equations. This marginal algorithm is more efficient than simultaneously inferring the locations of the trait events in a sampling scheme.

In the application we consider, accounting for lateral transfer results in an improved

fit over the regular Stochastic Dollo model but at a significant computational cost. The sequence of initial value problems to compute the likelihood parameters in the lateral transfer model is easy to state but difficult to solve in practice. On a tree with L leaves, we may exploit symmetry in the differential systems to compute the expected pattern frequencies exactly in $\mathcal{O}(2^{2L})$ operations. In practice, we use an ordinary differential equation solver to approximate their values within an error tolerance dominated by the Monte Carlo error. This approach requires $\mathcal{O}(L2^LC(L))$ operations, where the number of matrix-vector multiplications $C(L)$ taken by the **MATLAB ODE45** solver grows roughly linearly for the range of L we consider. This approach is feasible for up to approximately $L = 20$ leaves on readily available hardware. As we must evaluate the likelihood many times over the course of an MCMC analysis, this computational burden is a major stumbling block towards applying our model to data sets with more taxa or multiple character states and is the focus of our current research.

The model as described is not *projective* in the sense that we cannot marginalise out the effect of unobserved lineages. Consequently, as the number of unobserved lineages increases, the probability of a trait transferring to a sampled lineage decreases; and a trait which previously died out on the sampled lineages may transfer back into the system from an unobserved lineage. It would not be difficult to introduce *ghost* lineages in the style of [Szöllősi et al. \(2013\)](#) to our model to allow for lateral transfer between sampled and unsampled taxa. There are many other avenues for future work on the model. For example, one could partition the data across a mixture of models and trees; relax the global regime of lateral transfer; allow individual catastrophes to vary in their effect; model correlated evolution of traits in the style of [Cybis et al. \(2015\)](#), for example and allow for other types of missing data.

There are many open problems which have been ignored due to the expense of fitting models that account for lateral transfer. One such example occurs in the model proposed by [Chang et al.](#) whereby ancestral nodes may have data. Stochastic Dollo without lateral transfer cannot be used to model the observation process as traits absent in an ancestral state but present in both descendent and non-descendent leaves violate the Dollo parsimony assumption. Our method provides a model-based solution to this problem.

Acknowledgements

The authors wish to thank Robin Ryder for assistance with the implementation and Simon Greenhill for providing the Eastern Polynesian data set and feedback on the manuscript. We also wish to acknowledge the feedback of the associate editor and two anonymous reviewers.

References

- S.S. Abby, E. Tannier, M. Gouy, and V. Daubin. Detecting lateral gene transfers by statistical reconciliation of phylogenetic forests. *BMC Bioinform.*, 11(1):324, 2010.
- A.V. Alekseyenko, C.J. Lee, and M.A. Suchard. Wagner and Dollo: a stochastic duet by composing two parsimonious solos. *Syst. Biol.*, 57(5):772–784, 2008.
- R.G. Beiko and N. Hamilton. Phylogenetic identification of lateral genetic transfer events. *BMC Evol. Biol.*, 6(1), 2006.

- R. Bouckaert and J. Heled. DensiTree 2: Seeing Trees Through the Forest. *bioRxiv*, page 012401, 2014.
- R. Bouckaert, P. Lemey, M. Dunn, S.J. Greenhill, A.V. Alekseyenko, A.J. Drummond, R.D. Gray, M.A. Suchard, and Q.D. Atkinson. Mapping the origins and expansion of the Indo-European language family. *Science*, 337(6097):957–960, 2012.
- W. Chang, C. Cathcart, D. Hall, and A. Garrett. Ancestry-constrained phylogenetic analysis supports the indo-european steppe hypothesis. *Language*, 91(1):194–244, 2015.
- E. Conte and G. Molle. Reinvestigating a key site for Polynesian prehistory: new results from the Hane dune site, Ua Huka (Marquesas). *Archaeol. Oceania*, 49(3):121–136, 2014.
- G.B. Cybis, J.S. Sinsheimer, T. Bedford, A.E. Mather, P. Lemey, and M.A. Suchard. Assessing phenotypic correlation through the multivariate phylogenetic latent liability model. *Ann. Appl. Stat.*, 9(2):969–991, 2015.
- V. Daubin, M. Gouy, and G. Perrière. A phylogenomic approach to bacterial phylogeny: evidence of a core of genes sharing a common history. *Genome Res.*, 12(7):1080–1090, 2002.
- A.J. Drummond, M.A. Suchard, D. Xie, and A. Rambaut. Bayesian phylogenetics with BEAUTi and the BEAST 1.7. *Mol. Biol. Evol.*, 29(8):1969–1973, 2012.
- C.J. Geyer. Practical Markov chain Monte Carlo. *Statist. Sci.*, pages 473–483, 1992.
- D.T. Gillespie. Exact stochastic simulation of coupled chemical reactions. *J. Phys. Chem.*, 81(25):2340–2361, 1977.
- R.D. Gray and Q.D. Atkinson. Language-tree divergence times support the Anatolian theory of Indo-European origin. *Nature*, 426(6965):435–439, 2003.
- R.D. Gray, A.J. Drummond, and S.J. Greenhill. Language phylogenies reveal expansion pulses and pauses in Pacific settlement. *Science*, 323(5913):479–483, 2009.
- R.D. Gray, D. Bryant, and S.J. Greenhill. On the shape and fabric of human history. *Philos. T. R. Soc. B*, 365(1559):3923–3933, 2010.
- S.J. Greenhill, R. Blust, and R.D. Gray. The Austronesian basic vocabulary database: From bioinformatics to lexomics. *Evol. Bioinform. Online*, 4:271–283, 2008.
- S.J. Greenhill, T.E. Currie, and R.D. Gray. Does horizontal transmission invalidate cultural phylogenies? *Proc. Roy. Soc. B*, 276(1665):2299–2306, 2009.
- J. Heled and A.J. Drummond. Calibrated tree priors for relaxed phylogenetics and divergence time estimation. *Syst. Biol.*, 61(1):138–149, 2012.
- D.H. Huson and D. Bryant. Application of phylogenetic networks in evolutionary studies. *Mol. Biol. Evol.*, 23(2):254–267, 2006.
- R.E. Kass and A.E. Raftery. Bayes factors. *J. Am. Stat. Assoc.*, 90(430):773–795, 1995.
- J.F.C. Kingman. *Poisson Processes*. Oxford University Press, 1992.

- L.S. Kubatko. Identifying hybridization events in the presence of coalescence via model selection. *Syst. Biol.*, 58(5):478–488, 2009.
- G.M. Lathrop. Evolutionary trees and admixture: phylogenetic inference when some populations are hybridized. *Ann. Hum. Genet.*, 46(3):245–255, 1982.
- D. Madigan and A.E. Raftery. Model selection and accounting for model uncertainty in graphical models using Occam’s window. *J. Am. Stat. Assoc.*, 89(428):1535–1546, 1994.
- J.C. Marck. *Topics in Polynesian language and culture history*, volume 504. Canberra: Pacific Linguistics, 2000.
- MATLAB. *R2015a*. The MathWorks, Inc., Natick, Massachusetts, United States.
- G. Nicholls and R. Ryder. Phylogenetic models for Semitic vocabulary. In D. Conesa, A. Forte, A. López-Quílez, and F. Muñoz, editors, *Proceedings of the International Workshop on Statistical Modelling*, pages 431–436, 2011.
- G.K. Nicholls and R.D. Gray. Dated ancestral trees from binary trait data and their application to the diversification of languages. *J. Roy. Stat. Soc. B*, 70(3):545–566, 2008.
- G.K. Nicholls, R.J. Ryder, and D. Welch. *TraitLab: a MatLab Package for Fitting and Simulating Binary Trait-Like Data*, 2013.
- J. Oldman, T. Wu, L. van Iersel, and V. Moulton. TriLoNet: Piecing together small networks to reconstruct reticulate evolutionary histories. *Mol. Biol. Evol.*, 2016. doi: 10.1093/molbev/msw068.
- N. Patterson, P. Moorjani, Y. Luo, S. Mallick, N. Rohland, Y. Zhan, T. Genschoreck, T. Webster, and D. Reich. Ancient admixture in human history. *Genetics*, 192(3):1065–1093, 2012.
- J.K. Pickrell and J.K. Pritchard. Inference of population splits and mixtures from genome-wide allele frequency data. *PLoS Genet.*, 8(11):e1002967, 2012.
- B. Rannala and Z. Yang. Bayes estimation of species divergence times and ancestral population sizes using DNA sequences from multiple loci. *Genetics*, 164(4):1645–1656, 2003.
- S. Roch and S. Snir. Recovering the treelike trend of evolution despite extensive lateral genetic transfer: A probabilistic analysis. *J. Comput. Biol.*, 20(2):93–112, 2013.
- R.J. Ryder and G.K. Nicholls. Missing data in a stochastic Dollo model for binary trait data, and its application to the dating of Proto-Indo-European. *J. Roy. Stat. Soc. C*, 60(1):71–92, 2011.
- J. Sjöstrand, A. Tofigh, V. Daubin, L. Arvestad, B. Sennblad, and J. Lagergren. A Bayesian method for analyzing lateral gene transfer. *Syst. Biol.*, 63(3):409–420, 2014.
- M. Spriggs and A. Anderson. Late colonization of east polynesia. *Antiquity*, 67(255):200–217, 1993.
- G.J. Szöllősi, B. Boussau, S.S. Abby, E. Tannier, and V. Daubin. Phylogenetic modeling of lateral gene transfer reconstructs the pattern and relative timing of speciations. *Proc. Natl. Acad. Sci. USA*, 109(43):17513–17518, 2012.

- G.J. Szöllősi, E. Tannier, N. Lartillot, and V. Daubin. Lateral Gene Transfer from the Dead. *Syst. Biol.*, 62(3):386–397, 2013.
- M. Walworth. Eastern polynesian: The linguistic evidence revisited. *Ocean. Ling.*, 53(2):256–272, 2014.
- D. Wen, Y. Yu, and L. Nakhleh. Bayesian inference of reticulate phylogenies under the multispecies network coalescent. *PLoS Genet.*, 12(5):e1006006, 2016.
- J.M. Wilmshurst, T.L. Hunt, C.P. Lipo, and A.J. Anderson. High-precision radiocarbon dating shows recent and rapid initial human colonization of East Polynesia. *Proc. Natl. Acad. Sci. USA*, 108(5):1815–1820, 2011.

Appendix 1 Likelihood calculation — distribution of pattern frequencies

Theorem 1 states that under the SDLT model, the pattern frequencies $\mathbf{N}(t) = (N_{\mathbf{p}}(t) : \mathbf{p} \in \mathcal{P}^{(t)})$ at time t on a tree g is a vector of independent Poisson-distributed random variables with rate parameters $\mathbf{x}(t) = (x_{\mathbf{p}}(t) : \mathbf{p} \in \mathcal{P}^{(t)}) = \mathbb{E}[\mathbf{N}(t)|g, \lambda, \mu, \beta]$ given by the sequence of initial value problems in Equation 3.

Proof of Theorem 1. The trait process evolves forwards in time along the branches of the tree starting from the Adam node at time $t_0 = -\infty$. The pure birth-death trait process $N_1(t)$ is in equilibrium when it reaches the root at time t_1 so $N_1(t_1^-) \equiv N_{11}(t_1) \sim \text{Poisson}(x_{11}(t_1))$ where $x_{11}(t_1) = \lambda/\mu$. The patterns $(0, 1)$ and $(1, 0)$ are inconsistent with the branching event at the root so $N_{01}(t_1) \equiv N_{10}(t_1) \equiv 0$ and $x_{01}(t_1) \equiv x_{10}(t_1) \equiv 0$ by definition (T1). This provides us with the initial condition for the start of the next interval between branching events. By repeatedly solving the corresponding initial value problem (3) and applying the initialisation operators, $\mathbf{T}^{(1)}, \mathbf{T}^{(2)}, \dots, \mathbf{T}^{(L-1)}$, to obtain the initial conditions, we correctly calculate the expected pattern frequencies at any time t .

To complete the proof, we derive the Kolmogorov forward equation describing the temporal evolution of $p_{\mathbf{n}}(t) = \mathbb{P}(\mathbf{N}(t) = \mathbf{n}|g, \lambda, \mu, \beta)$ for an integer vector $\mathbf{n} = (n_{\mathbf{p}} : \mathbf{p} \in \mathcal{P}^{(t)})$ and show that it is equivalent to the time-derivative of the hypothesised Poisson probability mass function

$$p_{\mathbf{n}}(t) = \prod_{\mathbf{p} \in \mathcal{P}^{(t)}} \frac{x_{\mathbf{p}}(t)^{n_{\mathbf{p}}} e^{-x_{\mathbf{p}}(t)}}{n_{\mathbf{p}}!}. \quad (6)$$

For patterns \mathbf{p} and $\mathbf{q} \in \mathcal{P}^{(t)}$, we require the operators $\mathbf{U}_{\mathbf{p}0}$, $\mathbf{U}_{\mathbf{pq}}$ and $\mathbf{U}_{0\mathbf{q}}$ which applied to \mathbf{n} yield

$$\begin{aligned} \mathbf{U}_{\mathbf{p}0}\mathbf{n} &= (\dots, n_{\mathbf{p}-1}, n_{\mathbf{p}} - 1, n_{\mathbf{p}+1}, \dots), \\ \mathbf{U}_{\mathbf{pq}}\mathbf{n} &= (\dots, n_{\mathbf{p}-1}, n_{\mathbf{p}} - 1, n_{\mathbf{p}+1}, \dots, n_{\mathbf{q}-1}, n_{\mathbf{q}} + 1, n_{\mathbf{q}+1}, \dots), \\ \mathbf{U}_{0\mathbf{q}}\mathbf{n} &= (\dots, n_{\mathbf{q}-1}, n_{\mathbf{q}} + 1, n_{\mathbf{q}+1}, \dots), \end{aligned}$$

where we have abused notation and used $\mathbf{p} - 1$ and $\mathbf{p} + 1$ to index the entries either side of $n_{\mathbf{p}}$ in \mathbf{n} . These operators respectively correspond to the change in \mathbf{n} observed if: a trait displaying pattern \mathbf{p} becomes extinct (T3); a trait which displayed pattern \mathbf{p} transitions to display pattern \mathbf{q} through either a death (T3) or transfer (T4) event; and a trait is born displaying pattern \mathbf{q} (T2). Of course, these transitions may only occur if the patterns communicate. If $\rho(\mathbf{n}, \mathbf{n}')$ denotes the transition rate from state $\mathbf{N}(t) = \mathbf{n}$ to \mathbf{n}' , then from Section 4.1.2,

$$\begin{aligned} \rho(\mathbf{n}, \mathbf{U}_{\mathbf{p}0}\mathbf{n}) &= n_{\mathbf{p}}\lambda_{\mathbf{p}0} \quad \text{where } \lambda_{\mathbf{p}0} = \begin{cases} \mu, & s(\mathbf{p}) = 1, \\ 0, & \text{otherwise,} \end{cases} \\ \rho(\mathbf{n}, \mathbf{U}_{\mathbf{pq}}\mathbf{n}) &= n_{\mathbf{p}}\lambda_{\mathbf{pq}} \quad \text{where } \lambda_{\mathbf{pq}} = \begin{cases} \mu, & \mathbf{q} \in S_{\mathbf{p}}^-, \\ \beta \frac{s(\mathbf{p})}{L}, & \mathbf{q} \in S_{\mathbf{p}}^+, \\ 0, & \text{otherwise,} \end{cases} \\ \rho(\mathbf{n}, \mathbf{U}_{0\mathbf{q}}\mathbf{n}) &= \lambda_{0\mathbf{q}} \quad \text{where } \lambda_{0\mathbf{q}} = \begin{cases} \lambda, & s(\mathbf{q}) = 1, \\ 0, & \text{otherwise.} \end{cases} \end{aligned} \quad (7)$$

The forward equation We derive the forward equation from first principles. For a short interval of length dt , we can use Equations 1 and 7 to obtain

$$\begin{aligned}
p_{\mathbf{n}}(t + dt) &= p_{\mathbf{n}}(t) \left[1 - \left(\sum_{\mathbf{p} \in \mathcal{P}(t)} \sum_{\mathbf{q} \in \mathcal{P}(t)} \rho(\mathbf{n}, \mathbf{U}_{\mathbf{pq}} \mathbf{n}) + \sum_{\mathbf{p} \in \mathcal{P}(t)} \rho(\mathbf{n}, \mathbf{U}_{\mathbf{p0}} \mathbf{n}) \right. \right. \\
&\quad \left. \left. + \sum_{\mathbf{q} \in \mathcal{P}(t)} \rho(\mathbf{n}, \mathbf{U}_{\mathbf{0q}} \mathbf{n}) \right) dt + o(dt) \right] \\
&+ \sum_{\mathbf{p} \in \mathcal{P}(t)} \sum_{\mathbf{q} \in \mathcal{P}(t)} p_{\mathbf{U}_{\mathbf{pq}} \mathbf{n}}(t) [\rho(\mathbf{U}_{\mathbf{pq}} \mathbf{n}, \mathbf{n}) dt + o(dt)] \\
&+ \sum_{\mathbf{p} \in \mathcal{P}(t)} p_{\mathbf{U}_{\mathbf{p0}} \mathbf{n}}(t) [\rho(\mathbf{U}_{\mathbf{p0}} \mathbf{n}, \mathbf{n}) dt + o(dt)] \\
&+ \sum_{\mathbf{q} \in \mathcal{P}(t)} p_{\mathbf{U}_{\mathbf{0q}} \mathbf{n}}(t) [\rho(\mathbf{U}_{\mathbf{0q}} \mathbf{n}, \mathbf{n}) dt + o(dt)].
\end{aligned} \tag{8}$$

Subtracting $p_{\mathbf{n}}(t)$ from both sides of Equation 8, dividing by dt and taking the limit as $dt \downarrow 0$, we obtain

$$\begin{aligned}
\dot{p}_{\mathbf{n}}(t) &= -p_{\mathbf{n}}(t) \left[\sum_{\mathbf{p} \in \mathcal{P}(t)} \sum_{\mathbf{q} \in \mathcal{P}(t)} \rho(\mathbf{n}, \mathbf{U}_{\mathbf{pq}} \mathbf{n}) + \sum_{\mathbf{p} \in \mathcal{P}(t)} \rho(\mathbf{n}, \mathbf{U}_{\mathbf{p0}} \mathbf{n}) \right. \\
&\quad \left. + \sum_{\mathbf{q} \in \mathcal{P}(t)} \rho(\mathbf{n}, \mathbf{U}_{\mathbf{0q}} \mathbf{n}) \right] + \sum_{\mathbf{p} \in \mathcal{P}(t)} \sum_{\mathbf{q} \in \mathcal{P}(t)} p_{\mathbf{U}_{\mathbf{pq}} \mathbf{n}}(t) \rho(\mathbf{U}_{\mathbf{pq}} \mathbf{n}, \mathbf{n}) \\
&+ \sum_{\mathbf{p} \in \mathcal{P}(t)} p_{\mathbf{U}_{\mathbf{p0}} \mathbf{n}}(t) \rho(\mathbf{U}_{\mathbf{p0}} \mathbf{n}, \mathbf{n}) + \sum_{\mathbf{q} \in \mathcal{P}(t)} p_{\mathbf{U}_{\mathbf{0q}} \mathbf{n}}(t) \rho(\mathbf{U}_{\mathbf{0q}} \mathbf{n}, \mathbf{n}).
\end{aligned} \tag{9}$$

We shall drop the dependence on t in our notation for the remainder of this section. From the hypothesised probability mass function (6), we see that

$$\begin{aligned}
p_{\mathbf{U}_{\mathbf{p0}} \mathbf{n}} &= p_{\mathbf{n}} \frac{n_{\mathbf{p}}}{x_{\mathbf{p}}}, \\
p_{\mathbf{U}_{\mathbf{pq}} \mathbf{n}} &= p_{\mathbf{n}} \frac{n_{\mathbf{p}}}{x_{\mathbf{p}}} \frac{x_{\mathbf{q}}}{n_{\mathbf{q}} + 1}, \\
p_{\mathbf{U}_{\mathbf{0q}} \mathbf{n}} &= p_{\mathbf{n}} \frac{x_{\mathbf{q}}}{n_{\mathbf{q}} + 1},
\end{aligned}$$

and substituting these identities into Equation 9, we obtain

$$\begin{aligned}
\dot{p}_{\mathbf{n}} &= -p_{\mathbf{n}} \left[\sum_{\mathbf{p} \in \mathcal{P}} \sum_{\mathbf{q} \in \mathcal{P}} n_{\mathbf{p}} \lambda_{\mathbf{pq}} + \sum_{\mathbf{p} \in \mathcal{P}} n_{\mathbf{p}} \lambda_{\mathbf{p0}} + \sum_{\mathbf{q} \in \mathcal{P}} \lambda_{\mathbf{0q}} \right] \\
&+ p_{\mathbf{n}} \sum_{\mathbf{p} \in \mathcal{P}} \sum_{\mathbf{q} \in \mathcal{P}} \frac{n_{\mathbf{p}}}{x_{\mathbf{p}}} \frac{x_{\mathbf{q}}}{n_{\mathbf{q}} + 1} (n_{\mathbf{q}} + 1) \lambda_{\mathbf{qp}} + p_{\mathbf{n}} \sum_{\mathbf{p} \in \mathcal{P}} \frac{n_{\mathbf{p}}}{x_{\mathbf{p}}} \lambda_{\mathbf{0p}} \\
&+ p_{\mathbf{n}} \sum_{\mathbf{q} \in \mathcal{P}} \frac{x_{\mathbf{q}}}{n_{\mathbf{q}} + 1} (n_{\mathbf{q}} + 1) \lambda_{\mathbf{q0}} \\
&= -p_{\mathbf{n}} \left[\sum_{\mathbf{p} \in \mathcal{P}} \sum_{\mathbf{q} \in \mathcal{P}} n_{\mathbf{p}} \lambda_{\mathbf{pq}} + \sum_{\mathbf{p} \in \mathcal{P}} n_{\mathbf{p}} \lambda_{\mathbf{p0}} + \sum_{\mathbf{q} \in \mathcal{P}} \lambda_{\mathbf{0q}} \right]
\end{aligned}$$

$$\begin{aligned}
& + p_{\mathbf{n}} \left[\sum_{\mathbf{p} \in \mathcal{P}} \sum_{\mathbf{q} \in \mathcal{P}} n_{\mathbf{p}} \frac{x_{\mathbf{q}}}{x_{\mathbf{p}}} \lambda_{\mathbf{q}\mathbf{p}} + \sum_{\mathbf{p} \in \mathcal{P}} \frac{n_{\mathbf{p}}}{x_{\mathbf{p}}} \lambda_{0\mathbf{p}} + \sum_{\mathbf{q} \in \mathcal{P}} x_{\mathbf{q}} \lambda_{\mathbf{q}0} \right] \\
& = p_{\mathbf{n}} \left[\sum_{\mathbf{p} \in \mathcal{P}} \sum_{\mathbf{q} \in \mathcal{P}} n_{\mathbf{p}} \left(-\lambda_{\mathbf{p}\mathbf{q}} + \frac{x_{\mathbf{q}}}{x_{\mathbf{p}}} \lambda_{\mathbf{q}\mathbf{p}} \right) + \sum_{\mathbf{p} \in \mathcal{P}} n_{\mathbf{p}} \left(-\lambda_{\mathbf{p}0} + \frac{1}{x_{\mathbf{p}}} \lambda_{0\mathbf{p}} \right) \right. \\
& \quad \left. + \sum_{\mathbf{q} \in \mathcal{P}} (-\lambda_{0\mathbf{q}} + x_{\mathbf{q}} \lambda_{\mathbf{q}0}) \right]. \tag{10}
\end{aligned}$$

Equation 10 is the forward equation for $\mathbf{N}(t)$.

Time derivative of the probability mass function Equation 2 describes the temporal evolution of $x_{\mathbf{p}}(t)$, the expected number of traits displaying pattern $\mathbf{p} \in \mathcal{P}^{(t)}$ at time t , which we can write

$$\dot{x}_{\mathbf{p}}(t) = -x_{\mathbf{p}}(t) \left(\lambda_{\mathbf{p}0} + \sum_{\mathbf{q} \in \mathcal{P}^{(t)}} \lambda_{\mathbf{p}\mathbf{q}} \right) + \lambda_{0\mathbf{p}} + \sum_{\mathbf{q} \in \mathcal{P}^{(t)}} x_{\mathbf{q}}(t) \lambda_{\mathbf{q}\mathbf{p}}, \tag{11}$$

using the identities in Equation 7 and the fact that $s(\mathbf{p}) = |S_{\mathbf{p}}^-| + \mathbf{1}_{\{s(\mathbf{p})=1\}}$ and $L^{(t)} - s(\mathbf{p}) = |S_{\mathbf{p}}^+|$. Differentiating the hypothesised probability mass function (6) with respect to time t , we obtain

$$p'_{\mathbf{n}}(t) = p_{\mathbf{n}}(t) \frac{d}{dt} \log(p_{\mathbf{n}}(t)) = p_{\mathbf{n}}(t) \sum_{\mathbf{p} \in \mathcal{P}^{(t)}} \frac{n_{\mathbf{p}}}{x_{\mathbf{p}}(t)} \dot{x}_{\mathbf{p}}(t) - p_{\mathbf{n}}(t) \sum_{\mathbf{p} \in \mathcal{P}^{(t)}} \dot{x}_{\mathbf{p}}(t). \tag{12}$$

Dropping the dependence on time t from our notation and substituting Equation 11 into Equation 12 yields

$$\begin{aligned}
p'_{\mathbf{n}} &= p_{\mathbf{n}} \sum_{\mathbf{p} \in \mathcal{P}} n_{\mathbf{p}} \left[- \left(\lambda_{\mathbf{p}0} + \sum_{\mathbf{q} \in \mathcal{P}} \lambda_{\mathbf{p}\mathbf{q}} \right) + \frac{1}{x_{\mathbf{p}}} \left(\lambda_{0\mathbf{p}} + \sum_{\mathbf{q} \in \mathcal{P}} x_{\mathbf{q}} \lambda_{\mathbf{q}\mathbf{p}} \right) \right] \\
&\quad + p_{\mathbf{n}} \sum_{\mathbf{p} \in \mathcal{P}} \left[x_{\mathbf{p}} \left(\lambda_{\mathbf{p}0} + \sum_{\mathbf{q} \in \mathcal{P}} \lambda_{\mathbf{p}\mathbf{q}} \right) - \lambda_{0\mathbf{p}} - \sum_{\mathbf{q} \in \mathcal{P}} x_{\mathbf{q}} \lambda_{\mathbf{q}\mathbf{p}} \right] \\
&= p_{\mathbf{n}} \sum_{\mathbf{p} \in \mathcal{P}} n_{\mathbf{p}} \left[\sum_{\mathbf{q} \in \mathcal{P}} \left(-\lambda_{\mathbf{p}\mathbf{q}} + \frac{x_{\mathbf{q}}}{x_{\mathbf{p}}} \lambda_{\mathbf{q}\mathbf{p}} \right) - \lambda_{\mathbf{p}0} + \frac{1}{x_{\mathbf{p}}} \lambda_{0\mathbf{p}} \right] \\
&\quad + p_{\mathbf{n}} \left[\sum_{\mathbf{p} \in \mathcal{P}} (-\lambda_{0\mathbf{p}} + x_{\mathbf{p}} \lambda_{\mathbf{p}0}) \right] + p_{\mathbf{n}} \sum_{\mathbf{p} \in \mathcal{P}} \left[x_{\mathbf{p}} \sum_{\mathbf{q} \in \mathcal{P}} \lambda_{\mathbf{p}\mathbf{q}} - \sum_{\mathbf{q} \in \mathcal{P}} x_{\mathbf{q}} \lambda_{\mathbf{q}\mathbf{p}} \right] \\
&= p_{\mathbf{n}} \left[\sum_{\mathbf{p} \in \mathcal{P}} \sum_{\mathbf{q} \in \mathcal{P}} n_{\mathbf{p}} \left(-\lambda_{\mathbf{p}\mathbf{q}} + \frac{x_{\mathbf{q}}}{x_{\mathbf{p}}} \lambda_{\mathbf{q}\mathbf{p}} \right) + \sum_{\mathbf{p} \in \mathcal{P}} n_{\mathbf{p}} \left(-\lambda_{\mathbf{p}0} + \frac{1}{x_{\mathbf{p}}} \lambda_{0\mathbf{p}} \right) \right. \\
&\quad \left. + \sum_{\mathbf{p} \in \mathcal{P}} (-\lambda_{0\mathbf{p}} + x_{\mathbf{p}} \lambda_{\mathbf{p}0}) \right] + p_{\mathbf{n}} \left[\sum_{\mathbf{p} \in \mathcal{P}} x_{\mathbf{p}} \sum_{\mathbf{q} \in \mathcal{P}} \lambda_{\mathbf{p}\mathbf{q}} - \sum_{\mathbf{p} \in \mathcal{P}} \sum_{\mathbf{q} \in \mathcal{P}} x_{\mathbf{q}} \lambda_{\mathbf{q}\mathbf{p}} \right]. \tag{13}
\end{aligned}$$

The final term in Equation 13 is 0 so it matches the forward equation (10). We therefore conclude that Equation 6 with parameters given by the solution of Equation 3 correctly describes the distribution of the pattern frequencies at any time t . \square

Appendix 2 Model extensions — missing data

We can drop all explicit dependence on time t from our notation as the data are recorded at time 0. In Section 5.2, we claimed that the frequency $N_{\mathbf{q}}$ of traits in \mathcal{H} displaying pattern $\mathbf{q} \in \mathcal{Q}$ across L leaves is Poisson distributed with mean

$$\mathbb{E}[N_{\mathbf{q}}|g, \lambda, \mu, \beta, \Xi] = \left(\sum_{\mathbf{p} \in u(\mathbf{q})} x_{\mathbf{p}} \right) \left[\prod_{i=1}^L \xi_{k_i}^{1_{\{q_i \neq ?\}}} (1 - \xi_{k_i})^{1_{\{q_i=?\}}} \right], \quad (14)$$

where $\Xi = (\xi_i : i \in V_L)$ denotes the state observation probabilities, and

$$u(\mathbf{q}) = \{\mathbf{p} \in \mathcal{P} : p_i = q_i \text{ if } q_i \neq ?, i \in [L]\}$$

is the set of patterns in \mathcal{P} consistent with \mathbf{q} .

The underlying pattern process is binary but when a trait is sampled, its true state at leaf $i \in V_L$ is recorded with probability $\xi_i \in \Xi$. Therefore, we can generate patterns in the set

$$\{\mathbf{q} \in \mathcal{Q} : q_i = p_i \text{ if } q_i \neq ?, i \in [L]\}$$

simply by first generating a binary pattern $\mathbf{p} \in u(\mathbf{q})$ from the trait process on the tree then masking its entries independently with probabilities given by Ξ . Let $\{N_{\mathbf{p}} = n_{\mathbf{p}}; N_{\mathbf{q}|\mathbf{p}} = n_{\mathbf{q}|\mathbf{p}}\}$ denote the event that $n_{\mathbf{p}}$ copies of pattern $\mathbf{p} \in \mathcal{P}$ are generated by the pattern process but $n_{\mathbf{q}|\mathbf{p}}$ have entries obscured so as to instead display $\mathbf{q} \in \mathcal{Q} \setminus \mathcal{P}$. By Theorem 1 and the restriction property of Poisson processes (Kingman, 1992),

$$\{N_{\mathbf{p}} = n_{\mathbf{p}}; N_{\mathbf{q}|\mathbf{p}} = n_{\mathbf{q}|\mathbf{p}}\} \sim \text{Poisson}(x_{\mathbf{p}} \xi(\mathbf{q})).$$

where $\xi(\mathbf{q}) = \prod_{i=1}^L \xi_{k_i}^{1_{\{q_i \neq ?\}}} (1 - \xi_{k_i})^{1_{\{q_i=?\}}}$. By the superposition property of independent Poisson processes, we add the rate parameters for each pattern $\mathbf{p} \in u(\mathbf{q})$ to conclude that $N_{\mathbf{q}}|g, \lambda, \mu, \beta, \Xi \sim \text{Poisson}(\sum_{\mathbf{p} \in u(\mathbf{q})} x_{\mathbf{p}} \xi(\mathbf{q}))$.

Appendix 3 Bayesian inference

Priors

We demonstrate how λ may be integrated out of the likelihood in Theorem 1. Let $\mathbf{y}(t) = \mathbf{x}(t; g, \lambda = 1, \mu, \beta, \dots)$ with entry $y_{\mathbf{p}}(t)$ for pattern $\mathbf{p} \in \mathcal{P}^{(t)}$. At the root,

$$\begin{aligned}\dot{\mathbf{x}}(t_1; g, \lambda = \lambda', \dots) &= \mathbf{A}^{(t_1)} \mathbf{x}(t_1; g, \lambda', \dots) + \mathbf{b}^{(t_1)} \\ &= \mathbf{A}^{(t_1)} \begin{bmatrix} 0 \\ 0 \\ \lambda'/\mu \end{bmatrix} + \begin{bmatrix} \lambda' \\ \lambda' \\ 0 \end{bmatrix} \\ &= \lambda' \dot{\mathbf{y}}(t_1),\end{aligned}$$

by construction, so $\mathbf{x}(t_2^-) = \lambda' \mathbf{y}(t_2^-)$. As the initial value problems (3) and initialisation operations $\mathbf{T}^{(i)}$ are linear, we can repeat the above argument to see that $\mathbf{x}(t; g, \lambda', \dots) = \lambda' \mathbf{y}(t)$ for any time t .

Dropping the dependence on time from our notation, we now integrate λ out of the likelihood in Theorem 1 with respect to its prior in Table 2 to obtain

$$\begin{aligned}\pi(\mathbf{D}|g, \mu, \beta, \dots) &= \int_0^\infty \pi(\mathbf{D}|g, \lambda, \mu, \beta, \dots) \pi(\lambda) d\lambda \\ &\propto \int_0^\infty \frac{1}{\lambda} \prod_{\mathbf{p}} (\lambda y_{\mathbf{p}})^{n_{\mathbf{p}}} e^{-\lambda y_{\mathbf{p}}} d\lambda \\ &= \left(\prod_{\mathbf{p}} y_{\mathbf{p}}^{n_{\mathbf{p}}} \right) \int_0^\infty \lambda^{(\sum_{\mathbf{p}} n_{\mathbf{p}})-1} e^{-\lambda \sum_{\mathbf{p}} y_{\mathbf{p}}} d\lambda \\ &= \left(\prod_{\mathbf{p}} y_{\mathbf{p}}^{n_{\mathbf{p}}} \right) \left(\sum_{\mathbf{p}} y_{\mathbf{p}} \right)^{-\sum_{\mathbf{p}} n_{\mathbf{p}}} \Gamma\left(\sum_{\mathbf{p}} n_{\mathbf{p}}\right) \\ &\propto \prod_{\mathbf{p}} \left(\frac{y_{\mathbf{p}}}{\sum_{\mathbf{q}} y_{\mathbf{q}}} \right)^{n_{\mathbf{p}}},\end{aligned}\tag{15}$$

which we recognise as the likelihood term in the marginal posterior (4).

Catastrophes occur according to a Poisson(ρ) process. Let $n^{(i)}$ denote the number of catastrophes on branch $i \in V \setminus \{0, 1\}$ of length $\Delta_i = t_{\text{pa}(i)} - t_i$. Let n denote the total number of the tree and Δ its length below the root. Then conditional on the rate ρ , the prior distribution on the number of catastrophes and their locations is

$$\pi(C|\rho) = \prod_{i \in V \setminus \{0, 1\}} \frac{(\rho \Delta_i)^{n^{(i)}} e^{-\rho \Delta_i}}{n^{(i)}!} \frac{n^{(i)}!}{\Delta^{n^{(i)}}} = \rho^n e^{-\rho \Delta}.\tag{16}$$

where we recall that conditional on their number, catastrophes are uniformly distributed across a branch. The factorial terms in the numerator in Equation 16 account for the fact that the set of catastrophes on each branch is invariant to relabelling.

The prior on ρ in Table 2 is $\Gamma(a, b)$ where $a = 1.5$ and $b = 5 \times 10^3$. We now integrate ρ out of the prior on the set catastrophes (16) with respect to its prior,

$$\pi(C) = \int_0^\infty \pi(C|\rho) \pi(\rho) d\rho$$

$$\begin{aligned}
&= \frac{b^a}{\Gamma(a)} \int_0^\infty \rho^{n+a-1} e^{-\rho(\Delta+b)} d\rho \\
&= \frac{b^a}{\Gamma(a)} \frac{\Gamma(n+a)}{(\Delta+b)^{n+a}} \\
&= \frac{\Gamma(n+a)}{\Gamma(a)n!} \left(\frac{\Delta}{\Delta+b} \right)^n \left(\frac{b}{\Delta+b} \right)^a \frac{n!}{\Delta^n},
\end{aligned}$$

to obtain a Negative Binomial distribution on the number of catastrophes with catastrophe locations distributed uniformly across the tree.

MCMC transition kernels

We extend the sampling algorithms described by [Nicholls and Gray](#) and [Ryder and Nicholls](#) for the Stochastic Dollo model to construct a Markov chain whose invariant distribution is the posterior $\pi(g, \mu, \beta, \kappa, \Xi | R(\mathbf{D}))$ in Equation 4. In the following, $x = [(V, E, T, C), \mu, \beta, \kappa, \Xi]$ is the current state of the chain and a move to a new state x^* drawn from the proposal distribution $Q(x, x^*)$ is accepted with probability $\min(1, r(x, x^*))$ where

$$r(x, x^*) = \frac{\pi(x^* | \mathbf{D})}{\pi(x | \mathbf{D})} \frac{Q(x^*, x)}{Q(x, x^*)}.$$

We apply the same scaling update to the lateral transfer rate β and the death rate μ . If $x^* = [(V, E, T, C), \mu, \beta^*, \kappa, \Xi]$ where $\beta^* \sim U[\varrho^{-1}\beta, \varrho\beta]$ for some constant $\varrho > 1$, the Hastings ratio for the move is

$$\frac{Q(x^*, x)}{Q(x, x^*)} = \frac{\beta}{\beta^*}.$$

A catastrophe $c = (b, u) \in C$ in state x occurs on branch $b \in E$ at time $t_b + u(t_{\text{pa}(b)} - t_b)$ where $u \in (0, 1)$. The location for a new catastrophe $c^* = (b^*, u^*)$ is chosen uniformly at random across the branches of the tree to form the proposed state x^* with catastrophe set $C \cup \{c^*\}$. Catastrophes are chosen uniformly at random for deletion in the reverse move so

$$\frac{Q(x^*, x)}{Q(x, x^*)} = \frac{p_{DC}}{p_{AC}} \frac{1}{|C| + 1} \sum_{i \in V \setminus \{0, 1\}} (t_{\text{pa}(i)} - t_i),$$

where p_{AC} and p_{DC} denote the probabilities of proposing to add and delete a catastrophe respectively.

We chose catastrophe $c = (b, u)$ uniformly from the catastrophe set C to move to branch b^* chosen uniformly from the $\deg(b) + \deg(\text{pa}(b)) - 2$ branches neighbouring branch b . This is equivalent to deleting a randomly chosen catastrophe and adding it to a neighbouring branch, although we do not resample the position factor u . If $c^* = (b^*, u)$ replaces c in the proposed state x^* ,

$$\frac{Q(x^*, x)}{Q(x, x^*)} = \frac{\deg(b) + \deg(\text{pa}(b)) - 2}{\deg(b^*) + \deg(\text{pa}(b^*)) - 2} \frac{t_{\text{pa}(b^*)} - t_{b^*}}{t_{\text{pa}(b)} - t_b},$$

where the scale factor $(t_{\text{pa}(b^*)} - t_{b^*}) / (t_{\text{pa}(b)} - t_b)$ accounts for the change in the catastrophe branch's length.

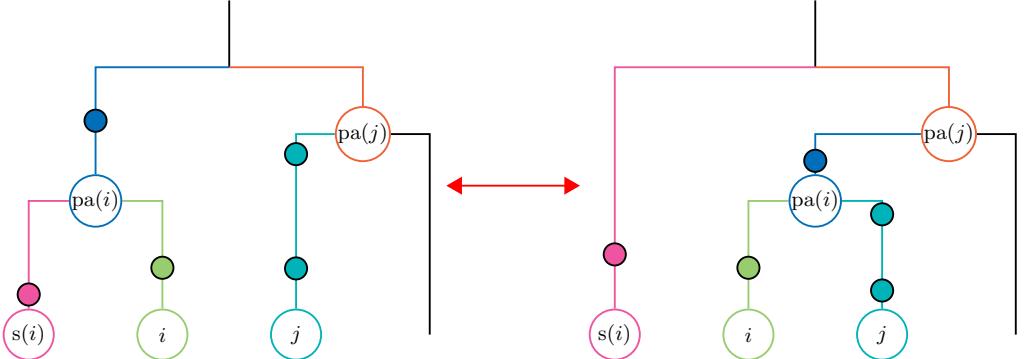
For every proposed move x to $x^* = [(V', E', T', C), \dots]$ which affects the tree, the Hastings ratio includes a term

$$\prod_{i \in V \setminus \{0, 1\}} \frac{|C^{(i)}|!}{(t_{\text{pa}(i)} - t_i)^{|C^{(i)}|}} \frac{(t_{\text{pa}(i)}^* - t_i^*)^{|C^{*(i)}|}}{|C^{*(i)}|!}$$

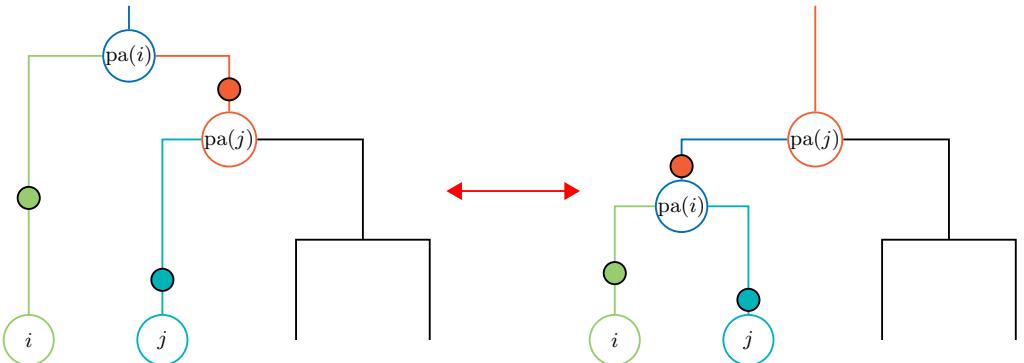
to account for the relative probabilities of sampling the catastrophe sets in each state, where $C^{(i)}$ and $C^{*(i)}$ denote the sets of catastrophes on branch i in the current and proposed states.

We construct subtree-prune-and-regraft moves on the tree in such a way that the total number of catastrophes on the tree remains constant and does not affect the ratio of proposal distributions outside the scaling term above. Let $\langle \text{pa}(i), i \rangle \in E$ denote a time-directed branch. From the current state x , we choose a node $i \in V \setminus \{0, 1\}$ below the root, prune the subtree beneath its parent $\text{pa}(i)$ and reattach it at a location chosen uniformly along a randomly chosen branch $\langle \text{pa}(j), j \rangle \in E$ to create state x^* . Now, recall that catastrophes are indexed by the offspring node of the branch they lie on and suppose that vertices retain their labels in the move from state x to state x^* . There are three possible outcomes of the proposed move.

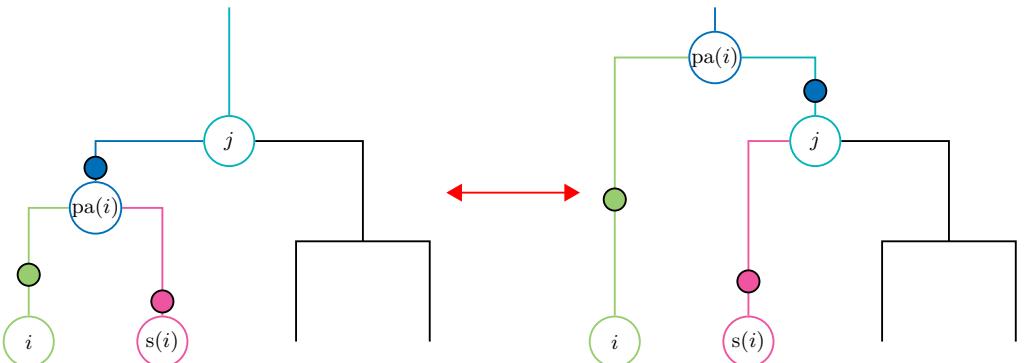
- If neither of nodes i or $\text{pa}(j)$ is the root in x then catastrophes remain on their assigned branches in state x^* . In more detail, if a catastrophe $c = (b, u)$ is on branch $\langle \text{pa}(b), b \rangle$ in state x then $c = (b, u)$ in state x^* also, with the possible difference that node b 's parent may change thereby affecting when the catastrophe occurs. We illustrate this in Figure 7a.
- If $\text{pa}(i)$ is the root in state x then i 's sibling $s(i)$ is the root in state x^* . There are no catastrophes on $\langle 0, \text{pa}(i) \rangle$ in x by definition so we move the catastrophes currently on $\langle \text{pa}(i), s(i) \rangle$ in x to $\langle \text{pa}(j), \text{pa}(i) \rangle$ in x^* , and the catastrophes on $\langle \text{pa}(j), j \rangle$ in x to $\langle \text{pa}(i), j \rangle$ in x^* . We illustrate this move in Figure 7b.
- Finally, if j is the root in x then $\text{pa}(i)$ becomes the root in x^* so we move the catastrophes on $\langle \text{pa}(\text{pa}(i)), \text{pa}(i) \rangle$ in x to $\langle \text{pa}(i), j \rangle$ in x^* . This move, the reverse of the above, is illustrated in Figure 7c.



(a) Neither node $pa(i)$ nor $pa(j)$ is the root in either state. Catastrophes remain on their assigned branches.



(b) Node $pa(i)$ is the root in the left-hand state and node $pa(j)$ is i 's sibling. The catastrophes on branch $\langle pa(i), pa(j) \rangle$ in the left-hand state are moved to edge $\langle pa(j), pa(i) \rangle$ in the right-hand state when node j becomes the root.



(c) Node j is the root in the left-hand state. The catastrophes on branch $\langle j, pa(i) \rangle$ in the left-hand state are moved to edge $\langle pa(i), i \rangle$ in the right-hand state.

Figure 7: Subtree-prune-and-regraft moves do not affect the number of catastrophes on the tree. Only the catastrophe branch index b changes in a move; the location u along the branch remains fixed.

Appendix 4 Method testing

We validate our model and computer implementation using three coupled synthetic data sets. We simulated the data set **SIM-B** using a [Gillespie](#)-type algorithm on the tree in Figure 9a with the trait process parameters in Table 4 ([Gillespie, 1977](#)). The data set **SIM-B** is an exact draw from the **SDLT** process and we use it to test the identifiability of the **SDLT** model. The relative transfer rate β/μ is high and although the shortcomings of the **SD** model in this setting have already been established ([Nicholls and Gray, 2008](#); [Greenhill et al., 2009](#)), we also fit it here to highlight the differences between the models.

Parameter	Value	Parameter	Value
Trait birth rate	$\lambda = 10^{-1}$	Root time	$t_1 = -10^3$
Trait death rate	$\mu = 5 \times 10^{-4}$	Catastrophe severity	$\kappa = 0.2212$
Trait transfer rate	$\beta = 5 \times 10^{-4}$	Observation probabilities	$\Xi \sim \beta(1, 1/3)^L$

Table 4: The parameter settings we used to simulate data set **SIM-B**. The catastrophe severity was chosen so that the effective duration of a catastrophe was 500 years.

From **SIM-B**, we create two additional data sets: **SIM-N** and **SIM-T**. To form **SIM-N**, we remove all the laterally transferred traits from **SIM-N**; that is, we do not discard all instances of a given trait, only the copies which transferred. This is equivalent to ignoring all the lateral transfer events when simulating **SIM-B** so **SIM-N** is a draw from the **SD** process coupled to **SIM-B**. As the **SD** model is nested within the **SDLT** model, we use **SIM-N** to test the consistency of the two models when the lateral transfer rate $\beta = 0$.

Recently transferred traits are more readily identified and discarded in practice. This potential bias is a common source of model misspecification. To this end, we only discarded instances of traits in **SIM-B** which transferred in the final 250 years to create **SIM-T**. We fit the **SDLT** and **SD** models here to test their robustness to this common form of model misspecification. We summarise the synthetic data sets in Table 5.

Data set	Traits	True model	Purpose
SIM-B	678	SDLT	Identifiability
SIM-N	672	SD	Consistency
SIM-T	675	SDLT before time -250 , SD thereafter	Robustness

Table 5: Summary of synthetic data sets for model testing.

In Figure 8, we compare the exact Poisson cumulative distribution function of each pattern in Ω under the **SDLT** model (Theorem 1) with empirical estimates based on 10^3 replicates of **SIM-B**. On the basis of these tests, we are satisfied that our simulation model and expected pattern frequency calculation (3) are correct.

For the MCMC analyses, we discard traits not marked present in at least one taxon (Section 5.4). In addition to the clade constraints depicted in Figure 9a, we enforce a minimum root time $t_1 = -2000$ years (maximum root age $-t_1$ is 2000 years). The results of our MCMC analyses are displayed in Figures 9 to 21. Figures in parentheses denote the effective sample size unless stated otherwise otherwise.

Of particular interest among the marginal tree posteriors in Figure 9 is the contrasting supports for the true topology in each case, particularly **SIM-B** where the **SD** model focuses on the wrong topology entirely. We return to this point in the goodness-of-fit analyses at

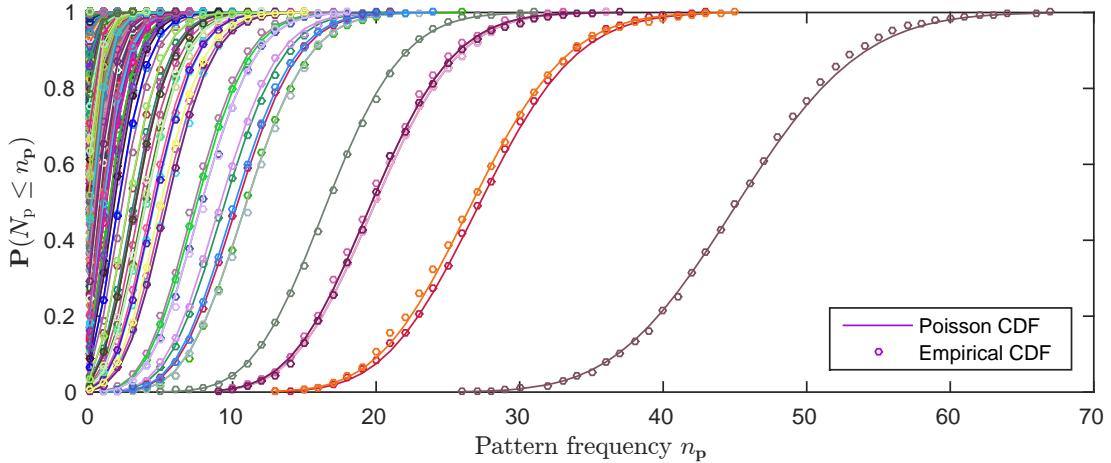


Figure 8: Exact and empirical cumulative distribution functions (CDFs) for 10^3 replicates of **SIM-B**.

the end of this section. We also note that the SD model applied to **SIM-B** underestimates the root age $-t_1$ in Figure 10.

Posterior estimates of the relative transfer rate β/μ in Figures 10–12 are consistent with their true values for **SIM-B** and **SIM-N**. On **SIM-T**, the posterior resembles a mixture distribution, which is unsurprising given its nature. Of the remaining parameters, we detect slight differences in the posterior distributions of the root age $-t_1$ and death rate μ on **SIM-B** and **SIM-T**. There is no cause for concern among the trace and autocorrelation plots displayed in Figures 13–18, so it remains to assess the quality of our analyses.

We repeat the Bayes factor tests on leaf constraints in Section 8 on the clade constraints in the models applied to the synthetic data sets. Each of the constraints in Figure 9a restricts the time t_i of an internal node $i \in V_A$ to lie on an interval $[\underline{t}_i, \bar{t}_i] \subset \mathbb{R}$. The clades are labelled 1, 2 and 3, where

- Clade $\Gamma^{(1)}$ fixes the time of the leaf labelled ‘1’ to lie on the interval $[-550, -450]$. Clade $\Gamma^{(1')}$ relaxes this to $[-800, -200]$.
- Clade $\Gamma^{(2)}$ fixes the time of the leaf labelled ‘8’ to lie on the interval $[-150, -50]$. Clade $\Gamma^{(2')}$ relaxes this to $[-400, 200]$.
- Clade $\Gamma^{(3)}$ fixes the time of the most recent ancestor of the leaves labelled ‘6’ to ‘10’ to lie on the interval $[-500, -200]$. Clade $\Gamma^{(3')}$ retains the ancestral constraint in $\Gamma^{(3)}$ but removes the time constraint. The node must therefore be greater than the root time and less than that of leaf ‘8’.

We plot histograms of the node ages under each model in Figure 19 and report Bayes factors computed according to Equation 5 in Figure 20. Unsurprisingly, the SD model performs poorly compared to the SDLT model in predicting the leaf constraints when fit to **SIM-B**. There is little to distinguish between the models applied to the other combinations of data sets and constraints.

We repeat the predictive performance checks described in Section 8 when the registered data $R(\mathbf{D})$ is randomly split into evenly-sized training and test portions, \mathbf{D}^{tr} and \mathbf{D}^{te} respectively. For the plots in Figure 21, the training sets are **SIM-B**, **SIM-N** and **SIM-T** analysed above, with the corresponding test sets drawn from the same distributions. The

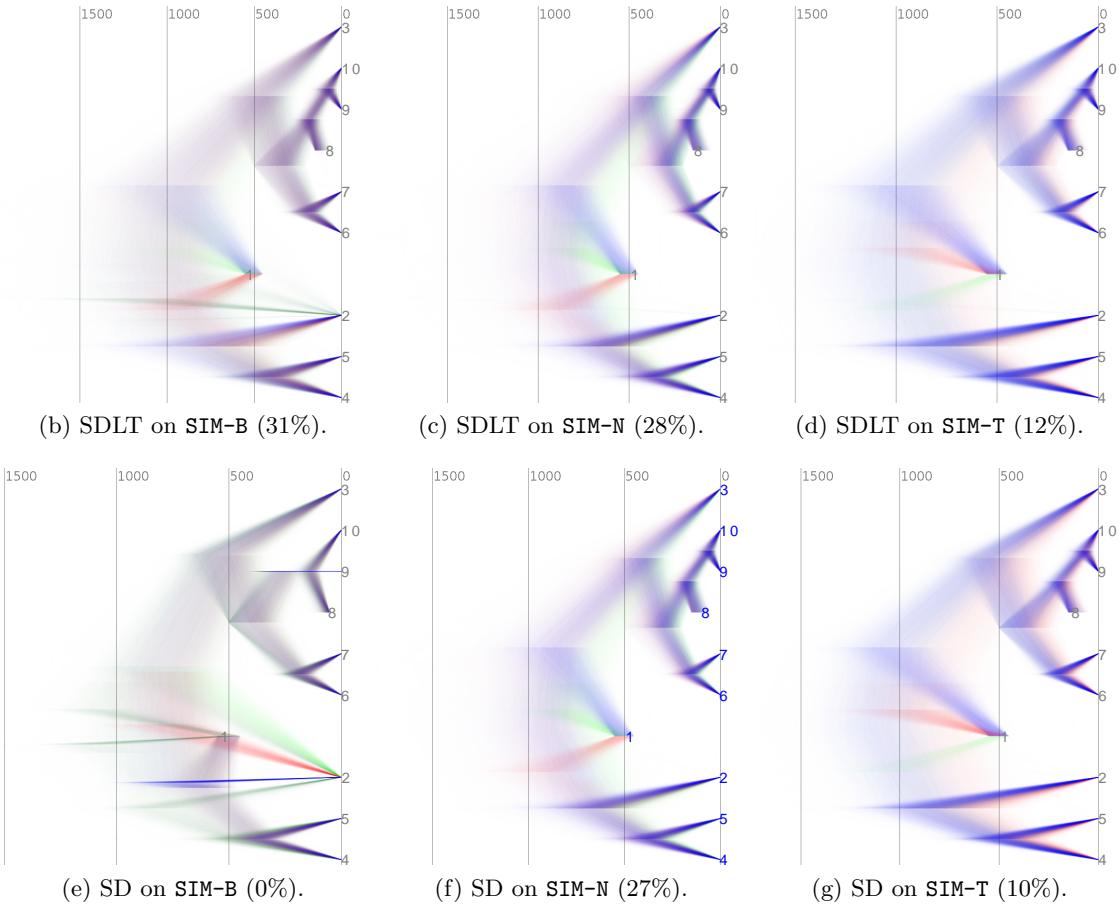
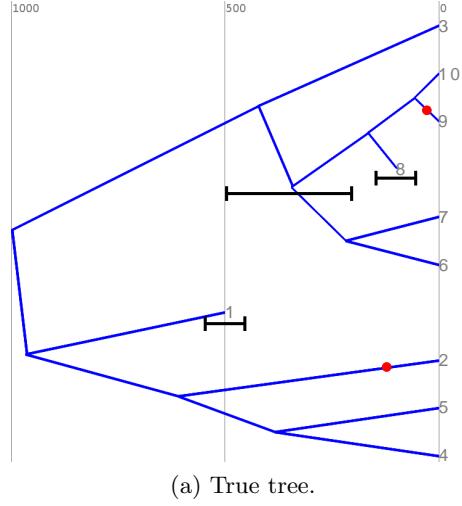


Figure 9: The tree (a) used to create the synthetic datasets. Catastrophe locations are marked in red and clade constraints for the MCMC analyses in black. **DensiTree** (Bouckaert and Heled, 2014) plots of the marginal tree posteriors for each synthetic dataset under the SDLT (b)–(d) and SD (e)–(g) models. Figures in parentheses denote posterior support for the true topology.

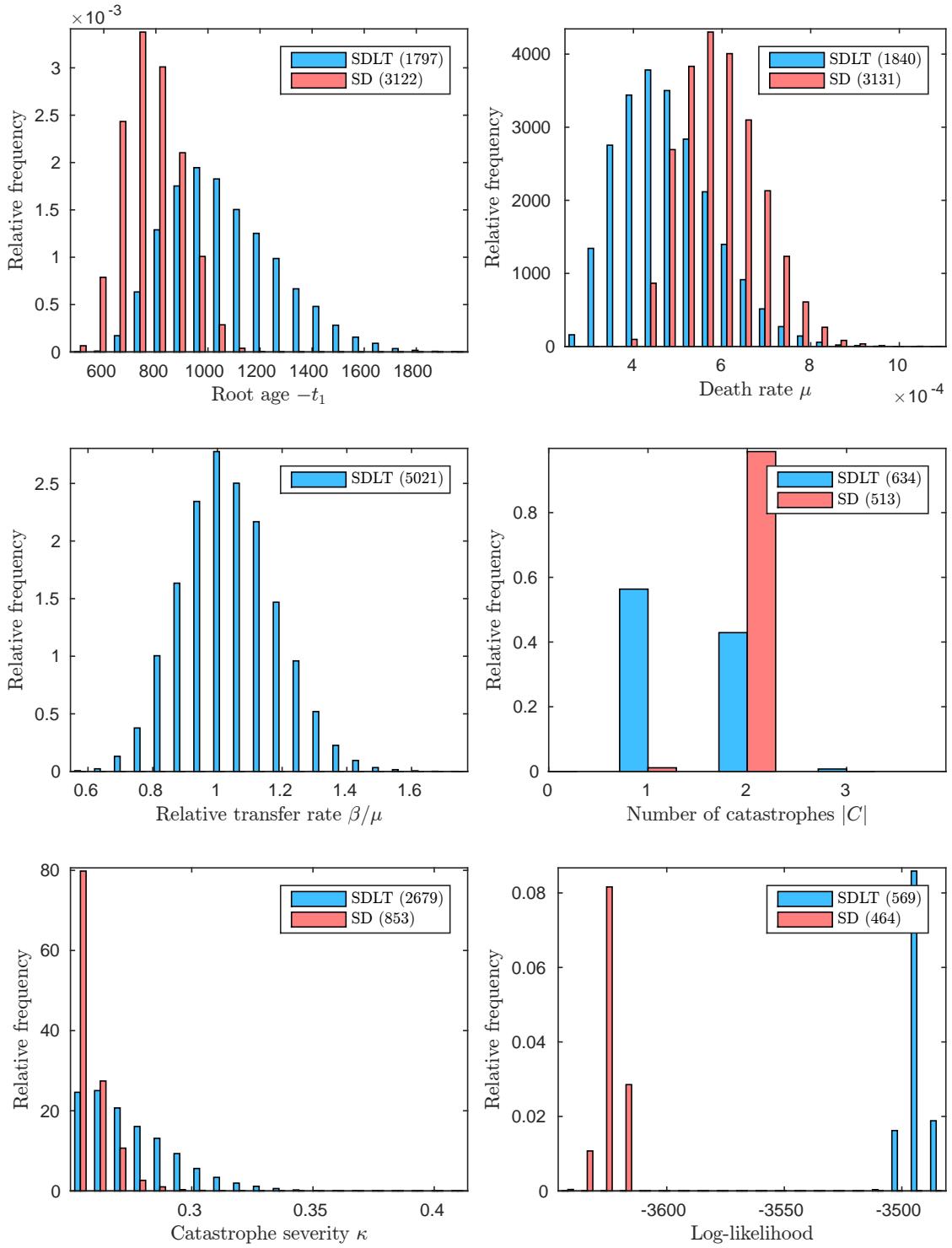


Figure 10: Histograms of samples in our analyses of SIM-B.

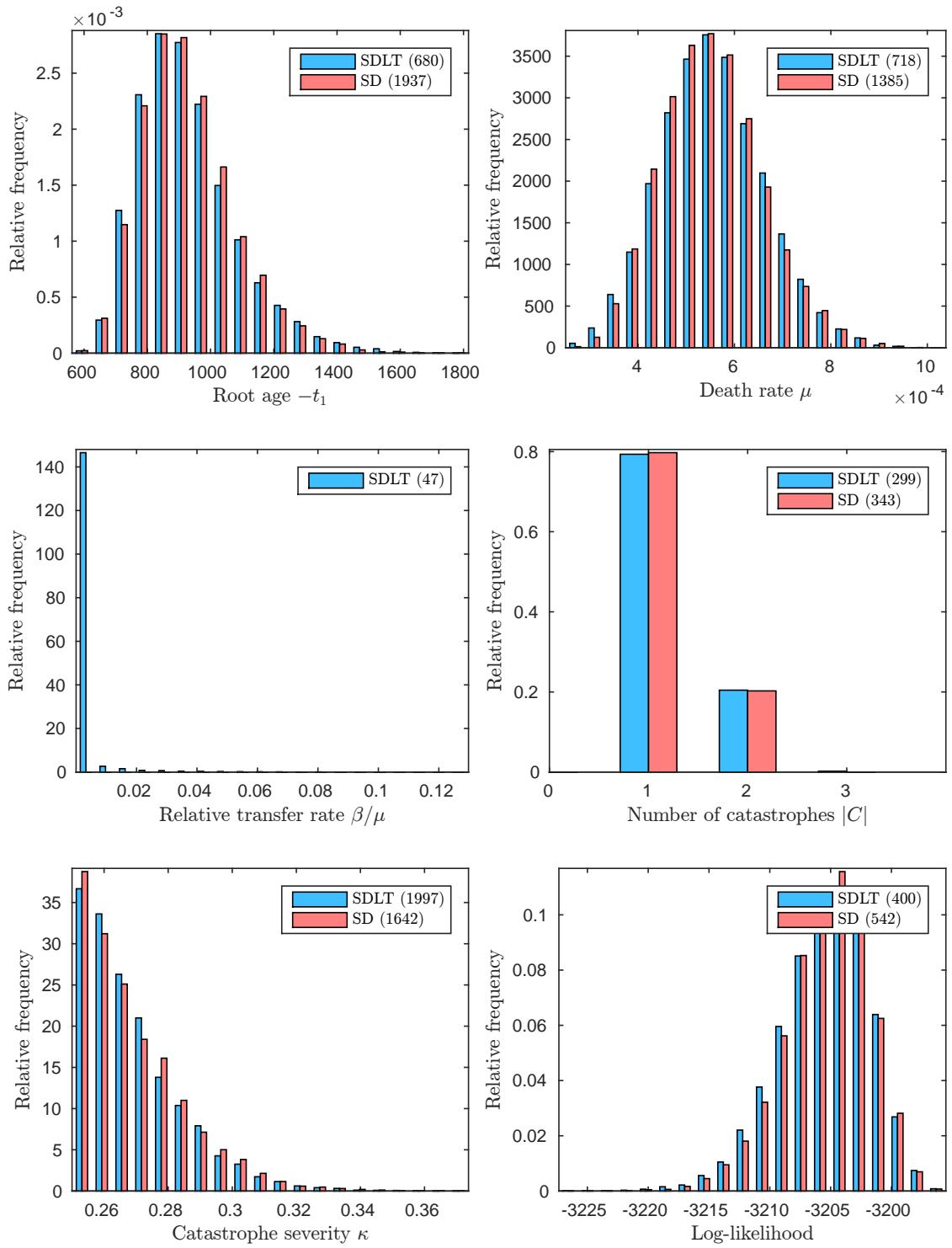


Figure 11: Histograms of samples in our analyses of SIM-N.

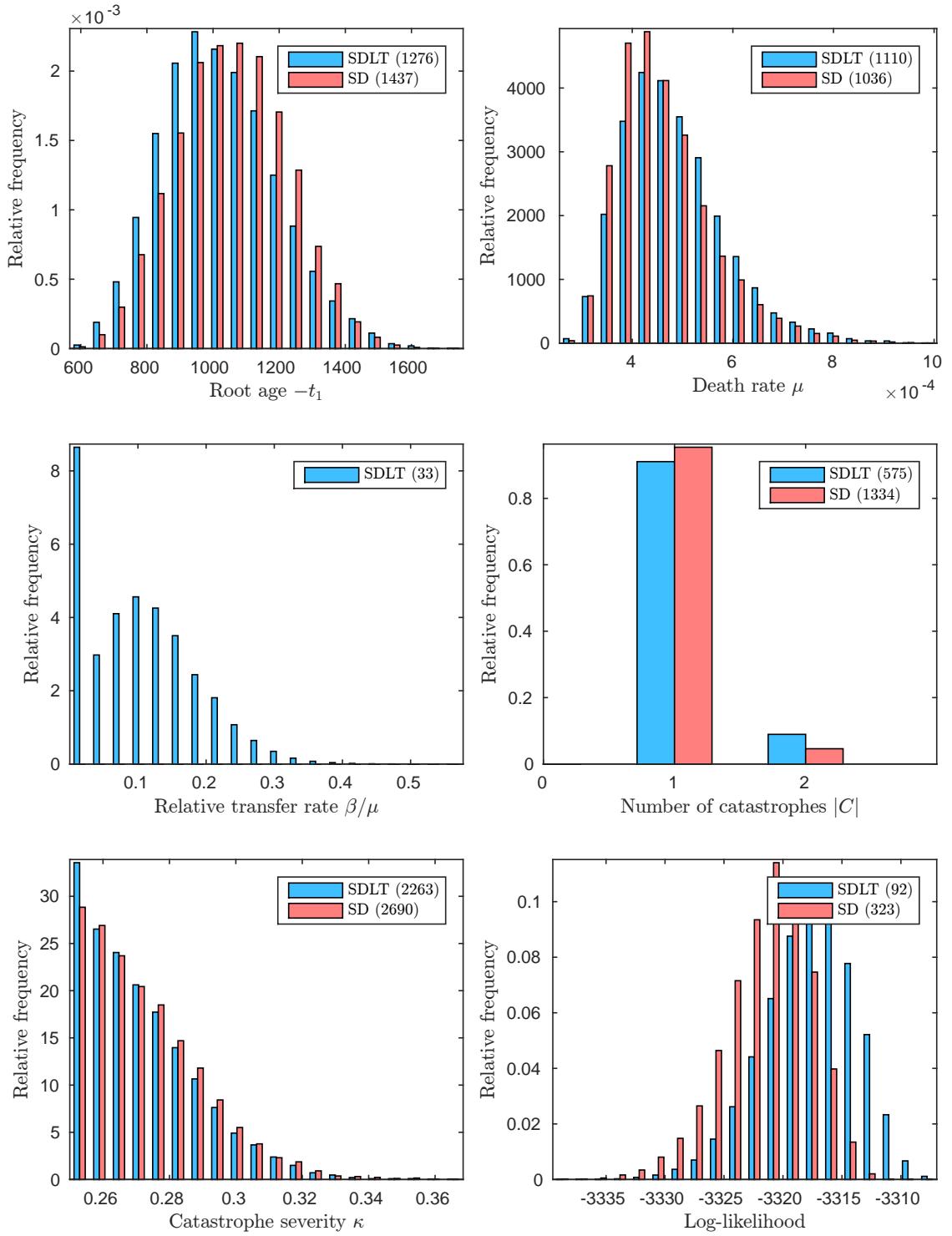


Figure 12: Histograms of samples in our analyses of SIM-T.

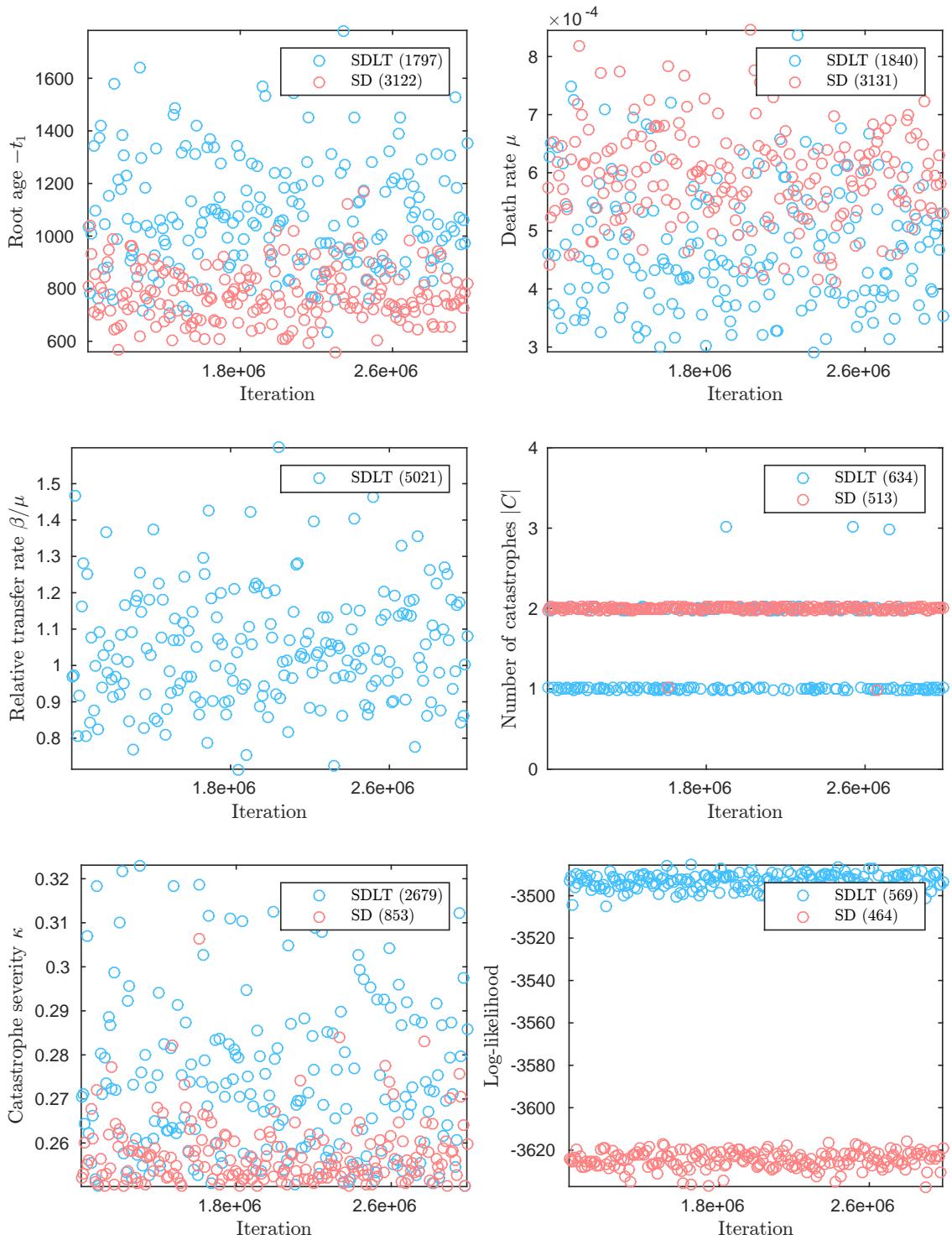


Figure 13: Trace plots of samples in our analyses of **SIM-B**.

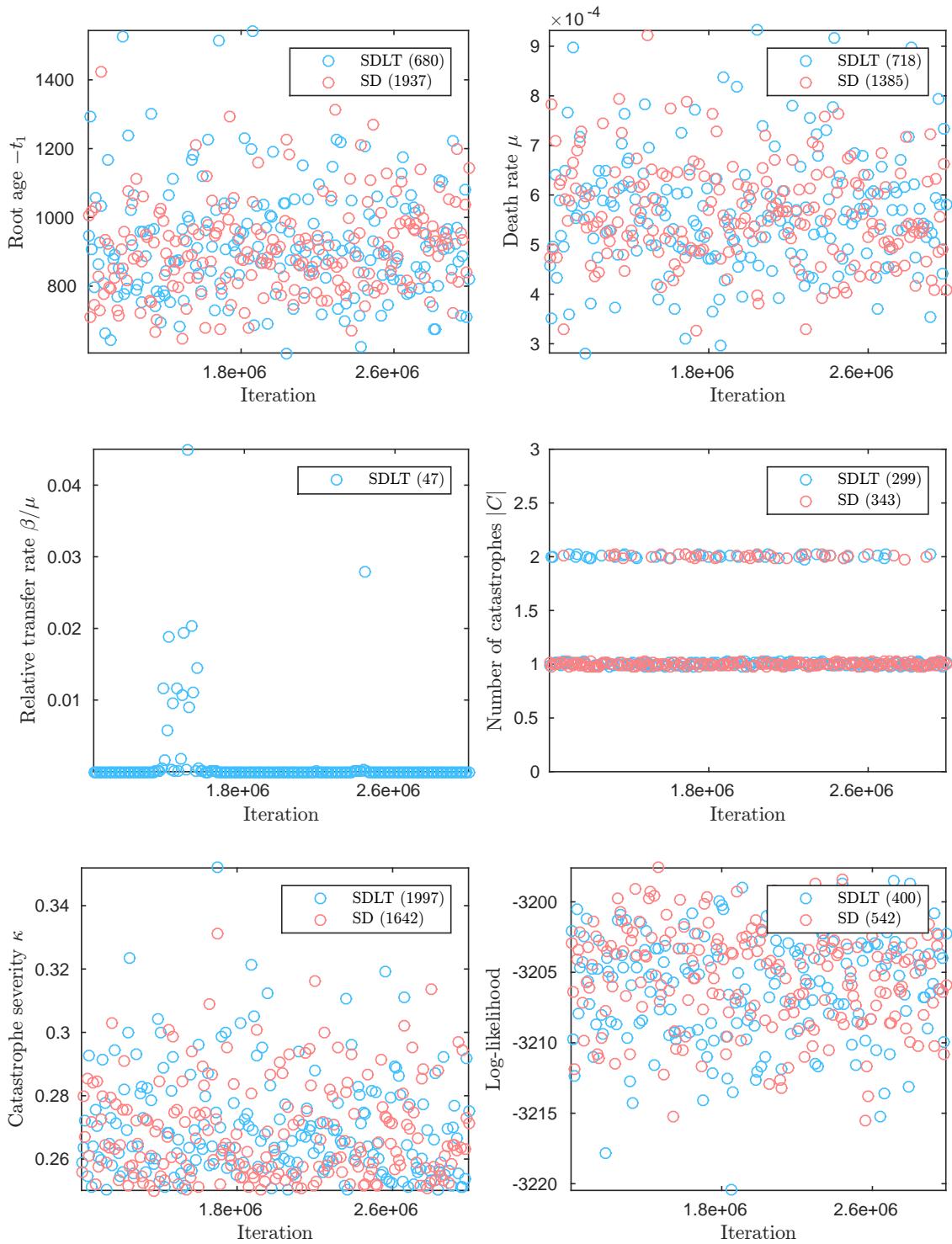


Figure 14: Trace plots of samples in our analyses of SIM-N.

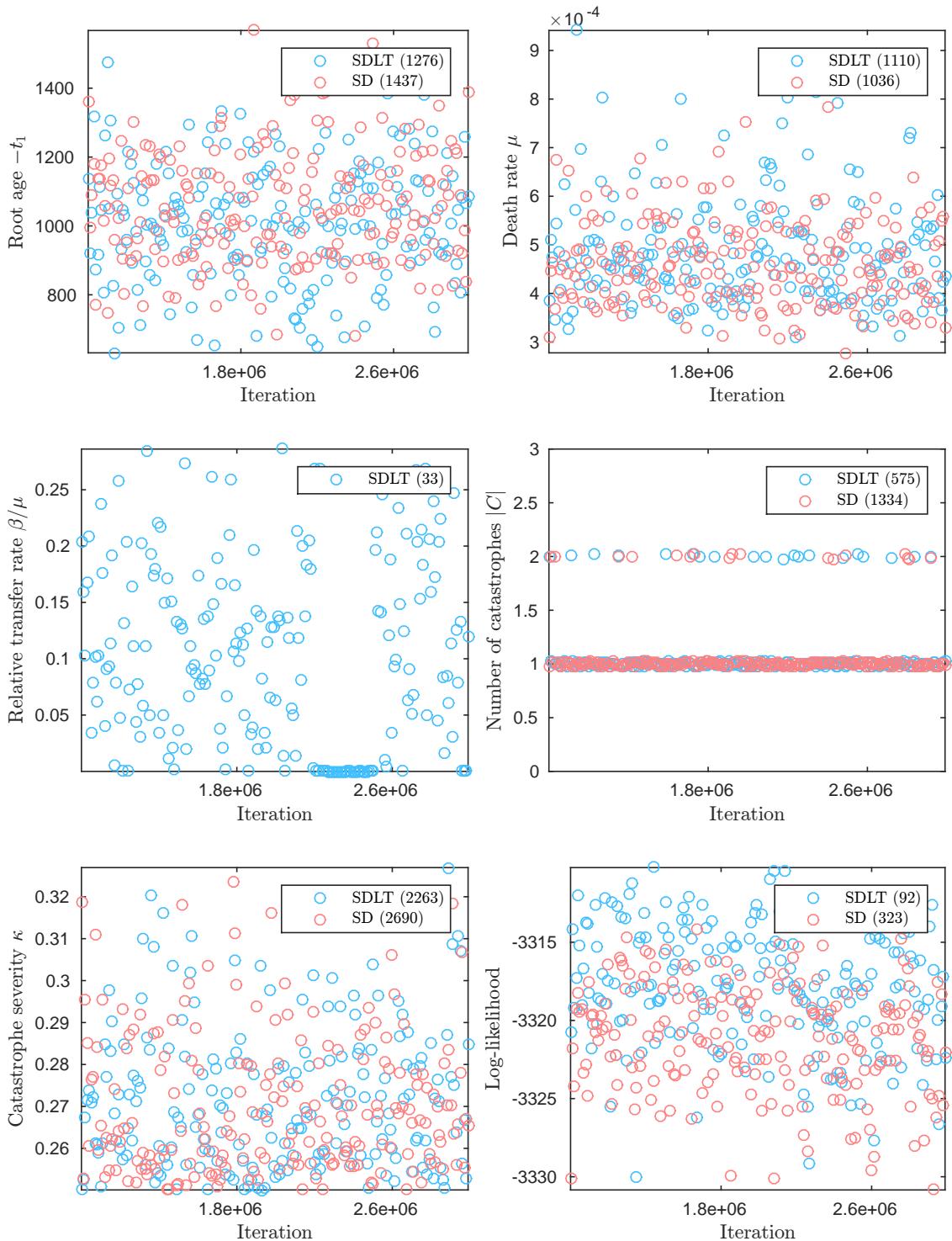


Figure 15: Trace plots of samples in our analyses of **SIM-T**.

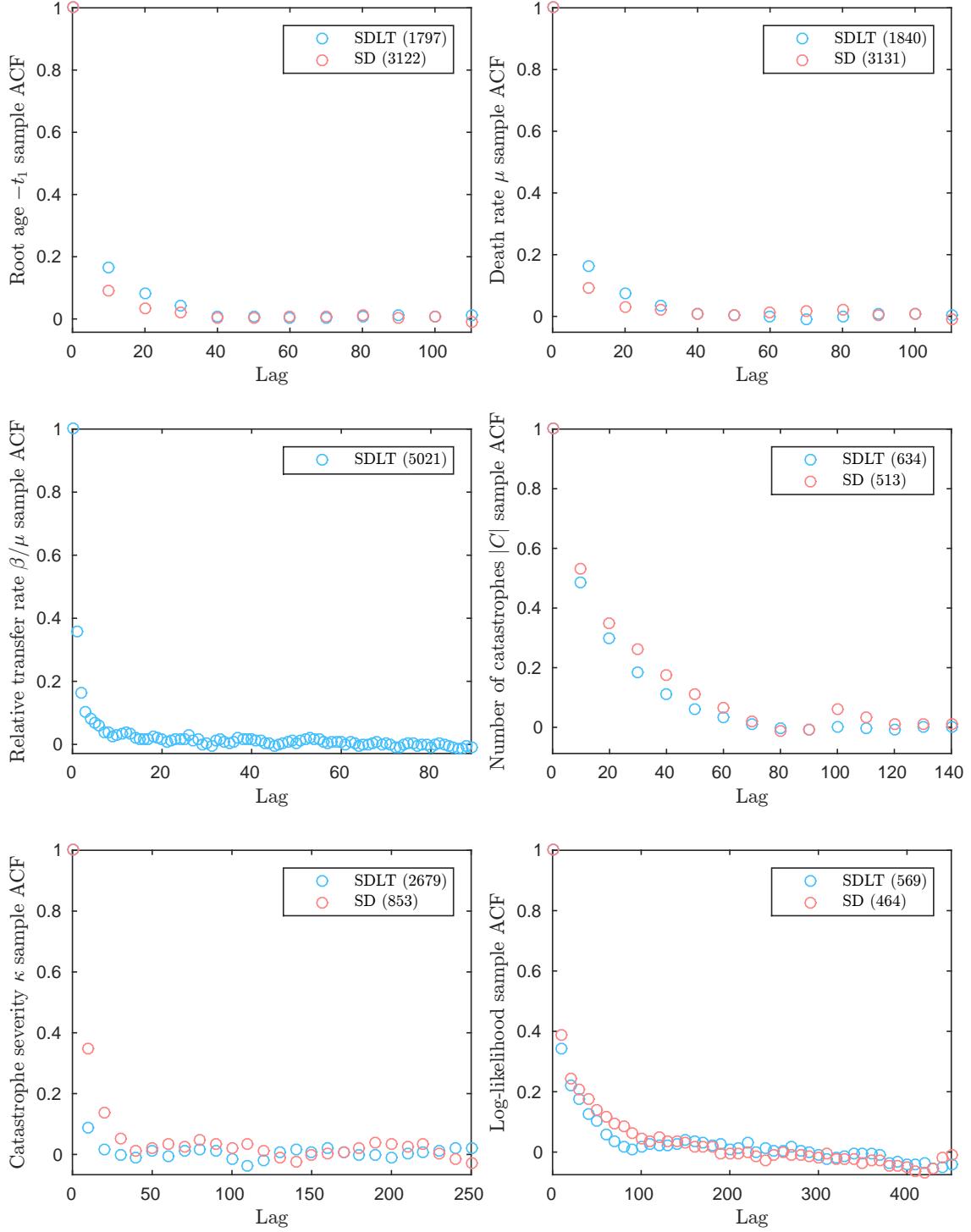


Figure 16: Autocorrelation plots of samples in our analyses of **SIM-B**.

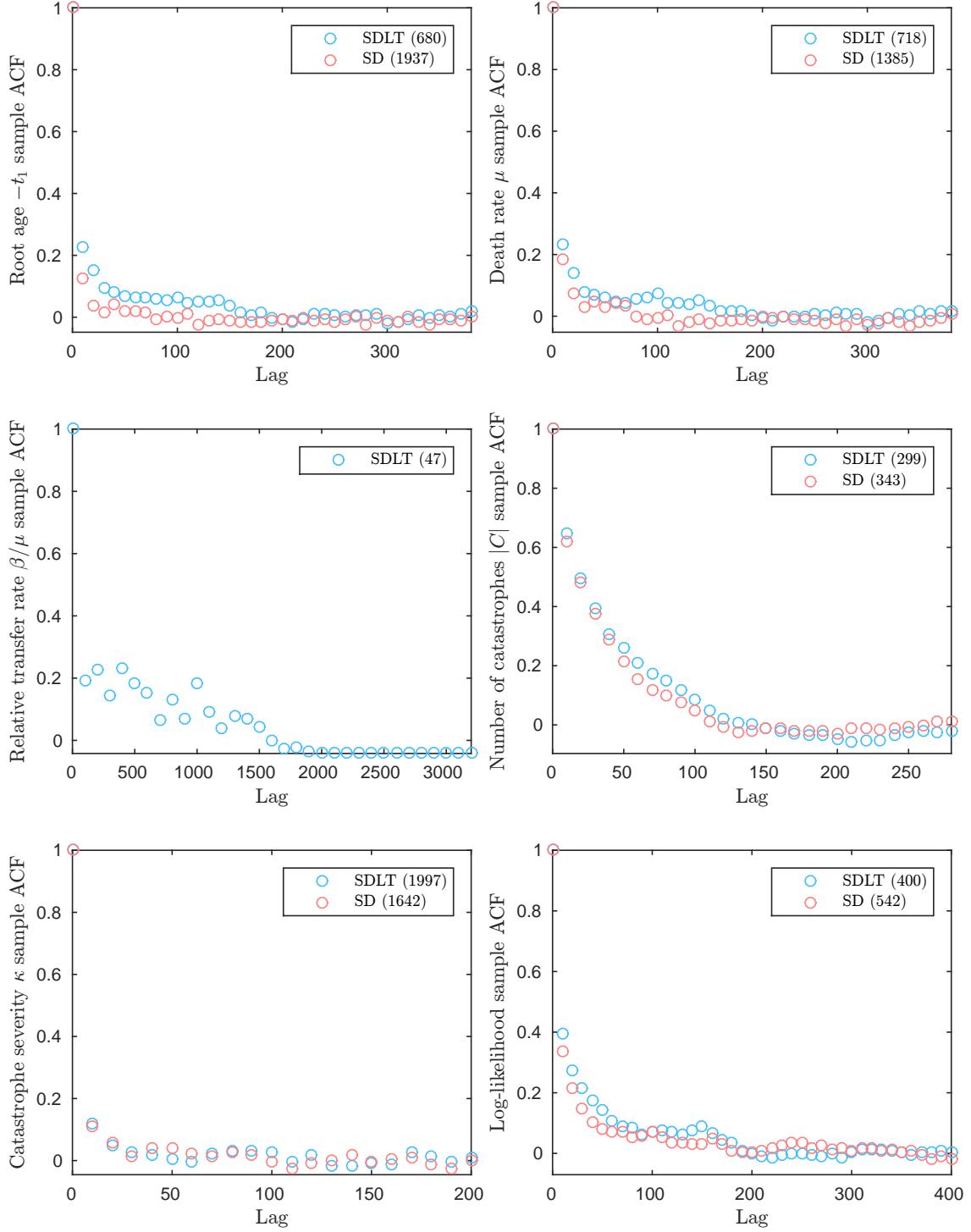


Figure 17: Autocorrelation plots of samples in our analyses of **SIM-N**.

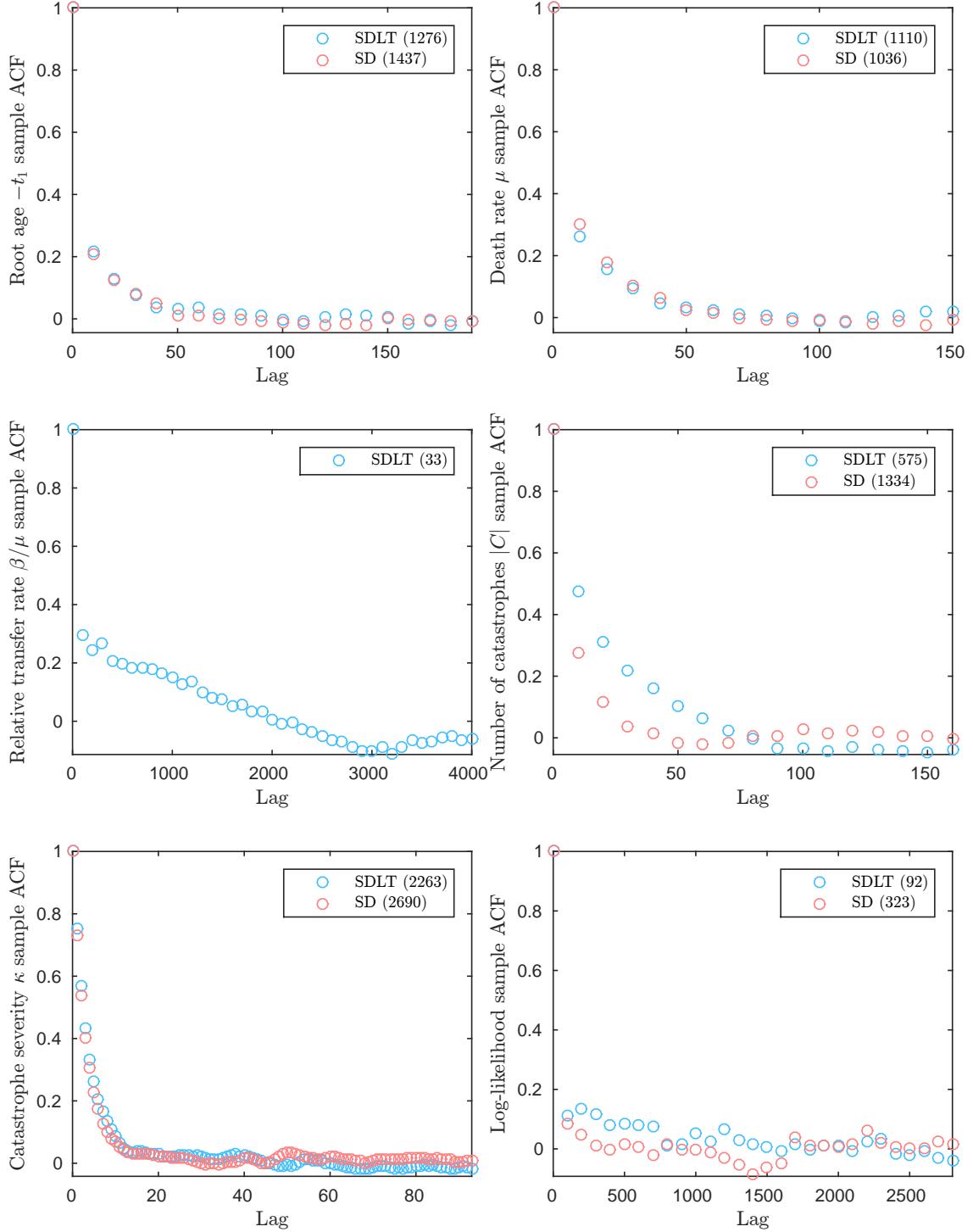


Figure 18: Autocorrelation plots of samples in our analyses of **SIM-T**.

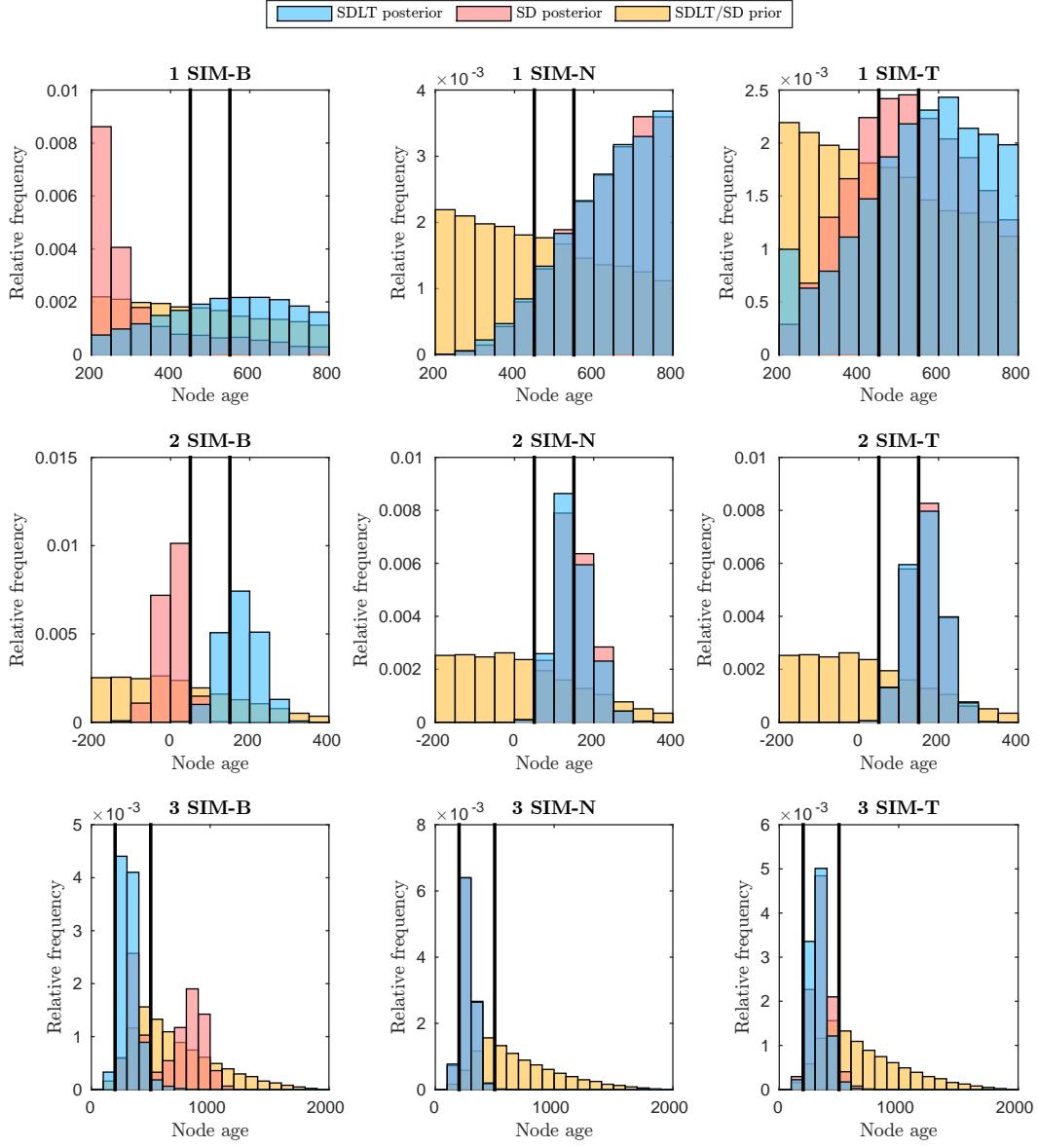


Figure 19: We relax each clade constraint in turn and compute the histogram of the corresponding node time under the prior and posterior for each model fit to each of the synthetic data sets.

more flexible SDLT model outperforms the SD model in each case, although marginally in the case of **SIM-N**.

An alternative approach to model comparison is to use reversible jump MCMC to sample from a posterior on models and parameters. The SD model is nested within the SDLT model so it is not difficult to implement such an algorithm. However, with our current set of proposal distributions we are unable to bridge the gap between the SDLT and SD models fit to **SIM-B** to obtain an accurate estimate of the corresponding Bayes factor.

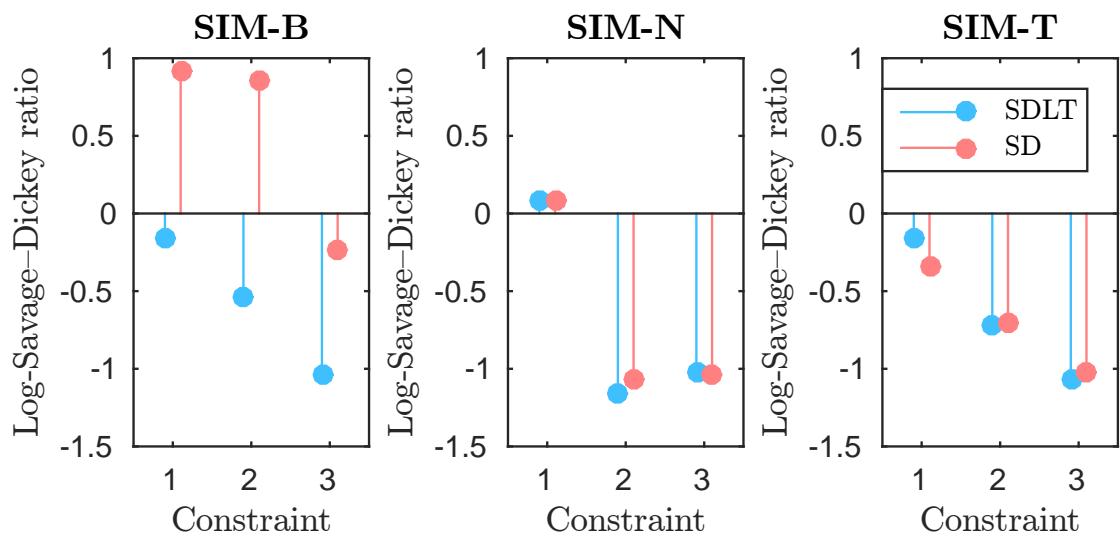


Figure 20: Bayes factors describing the lack of support for the clade constraints used to fit the SDLT and SD models to the synthetic data sets.

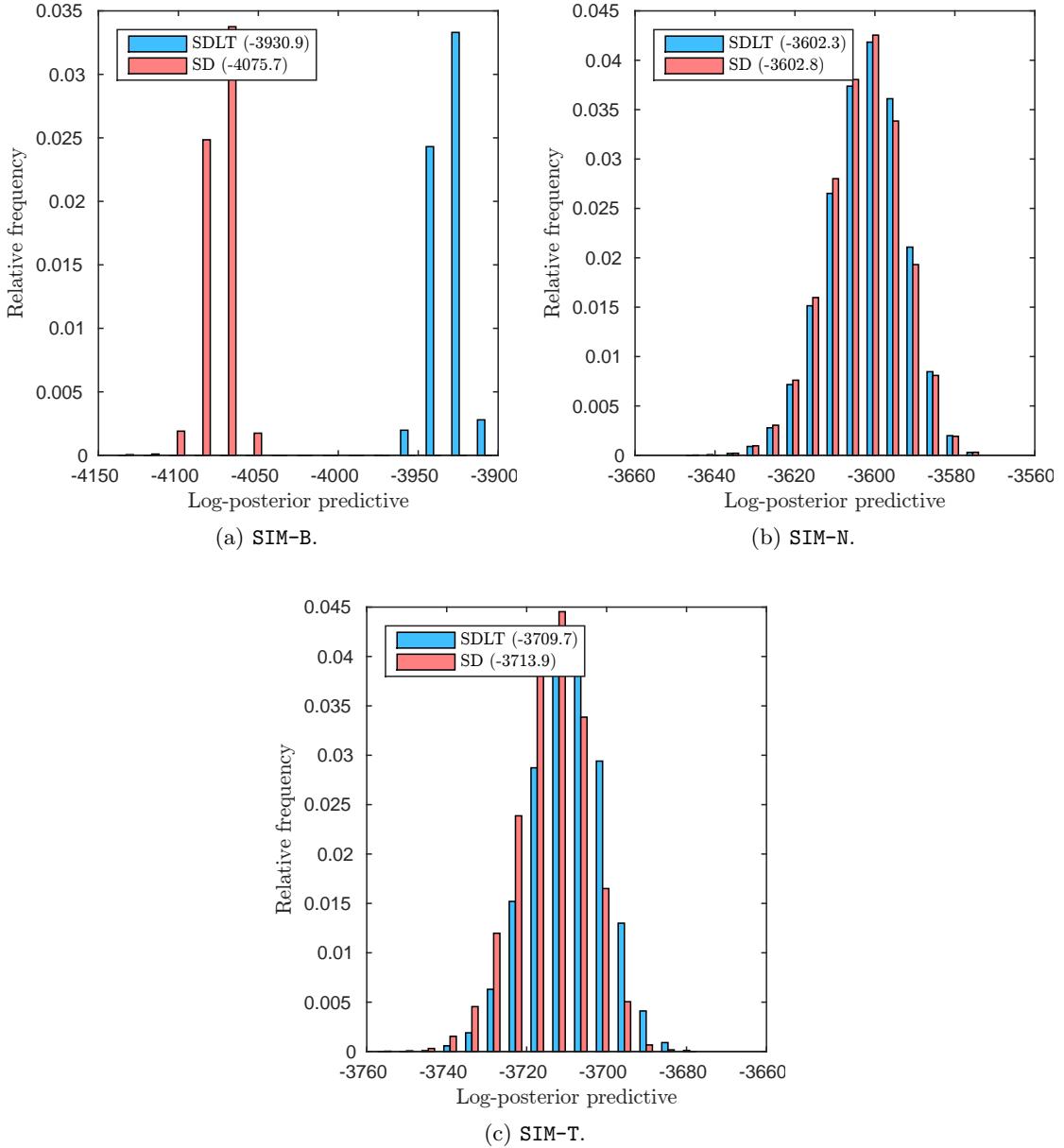


Figure 21: Posterior predictive model evaluation on the models fit to the synthetic data sets. We plot the joint posterior distribution of the test data and parameters given the training data, with the predictive score displayed in parentheses in each case.

Appendix 5 Applications

This section contains figures to support our analyses of the Polynesian data set POLY-0 in Section 8. A majority rule consensus tree depicts the relationships which appear in the majority of the sampled trees. The tree is multifurcating when an edge does not appear in the majority of the samples. In Figure 22, we report the consensus trees for the samples in Figure 4. We plot histograms of parameter samples in Figure 23, trace plots in Figure 24, autocorrelation plots in Figure 25, marginal leaf times when constraints are relaxed in Figure 26 and posterior predictive model checks in Figure 27. Figures in parentheses denote effective sample size unless stated otherwise otherwise.

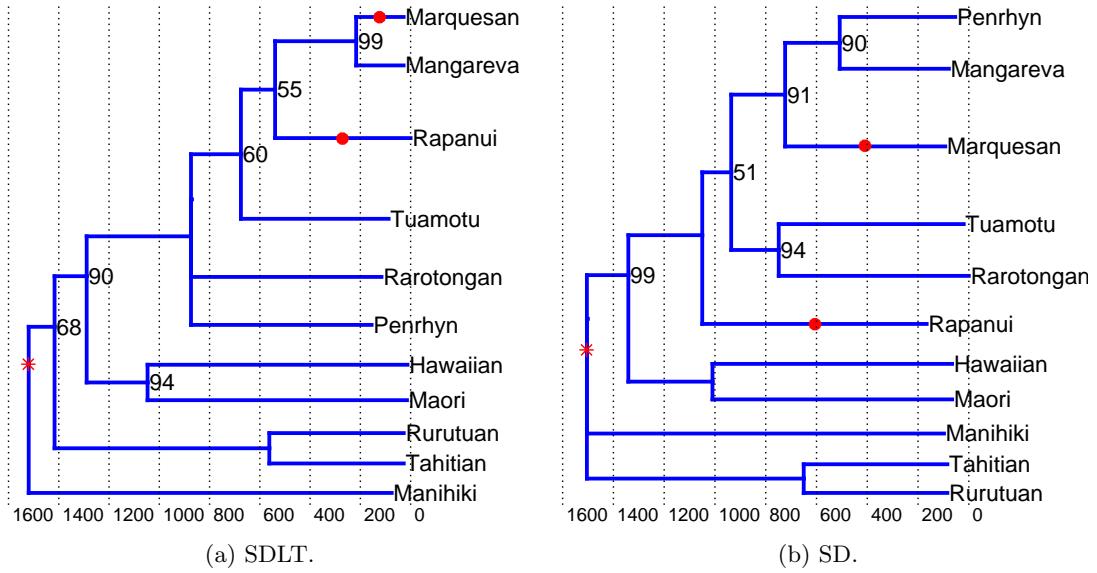


Figure 22: Majority rule consensus trees for our analyses of POLY-0.

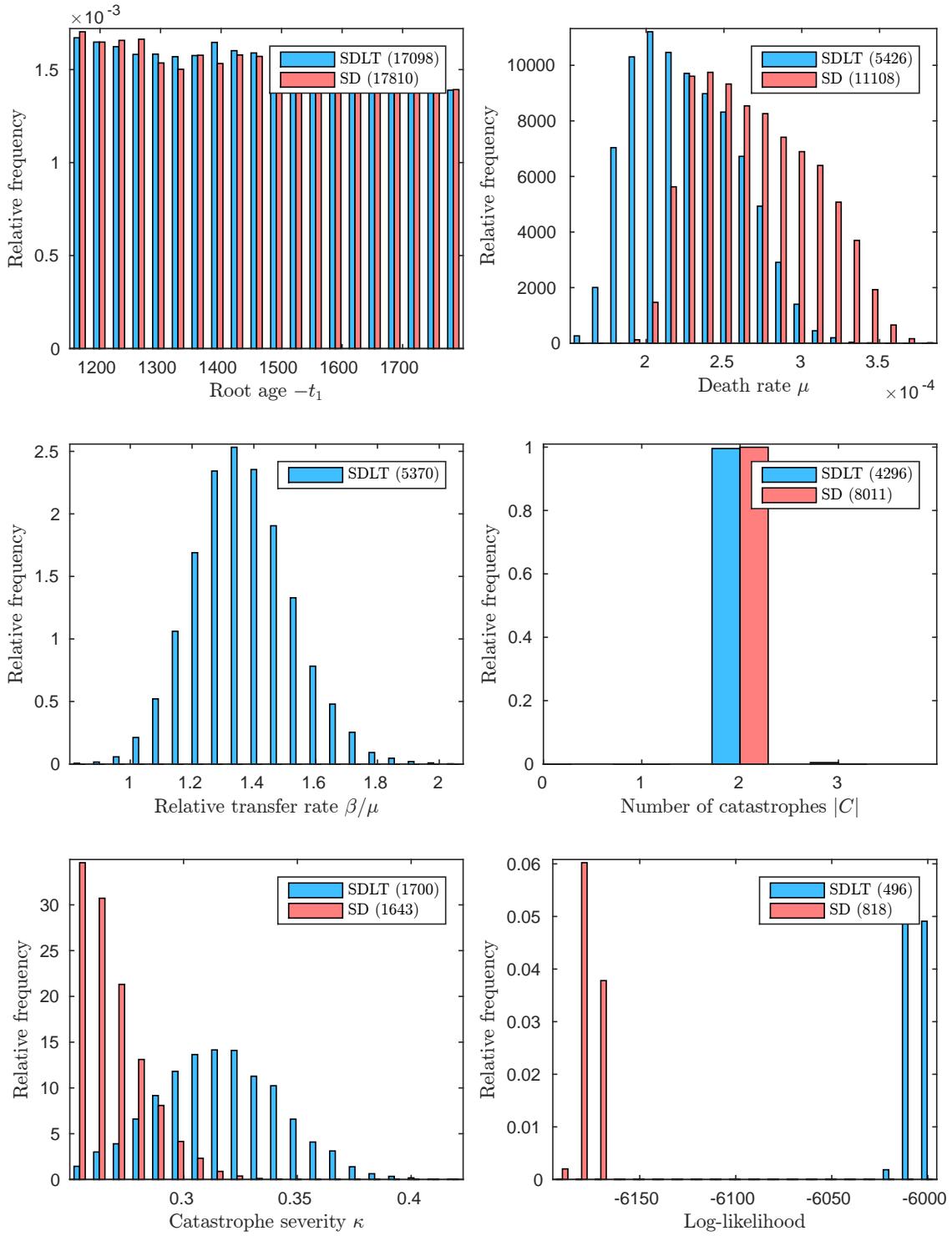


Figure 23: Histograms of samples in our analyses of POLY-0.

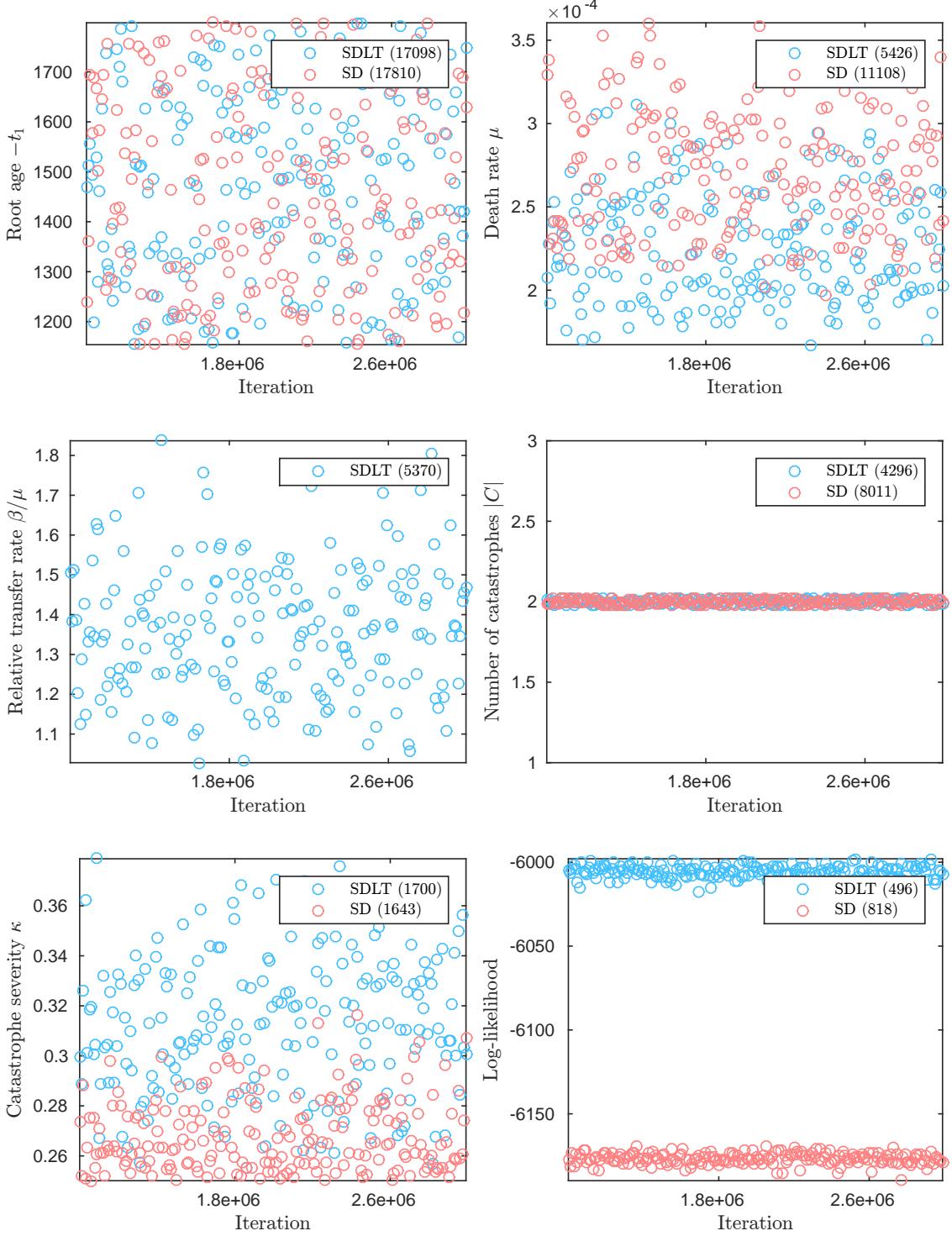


Figure 24: Trace plots of samples in our analyses of POLY-0. For clarity, the trace plots only depict a subset of the samples and this explains the discrepancy between the effective sample size for the number of catastrophes and the corresponding plot.

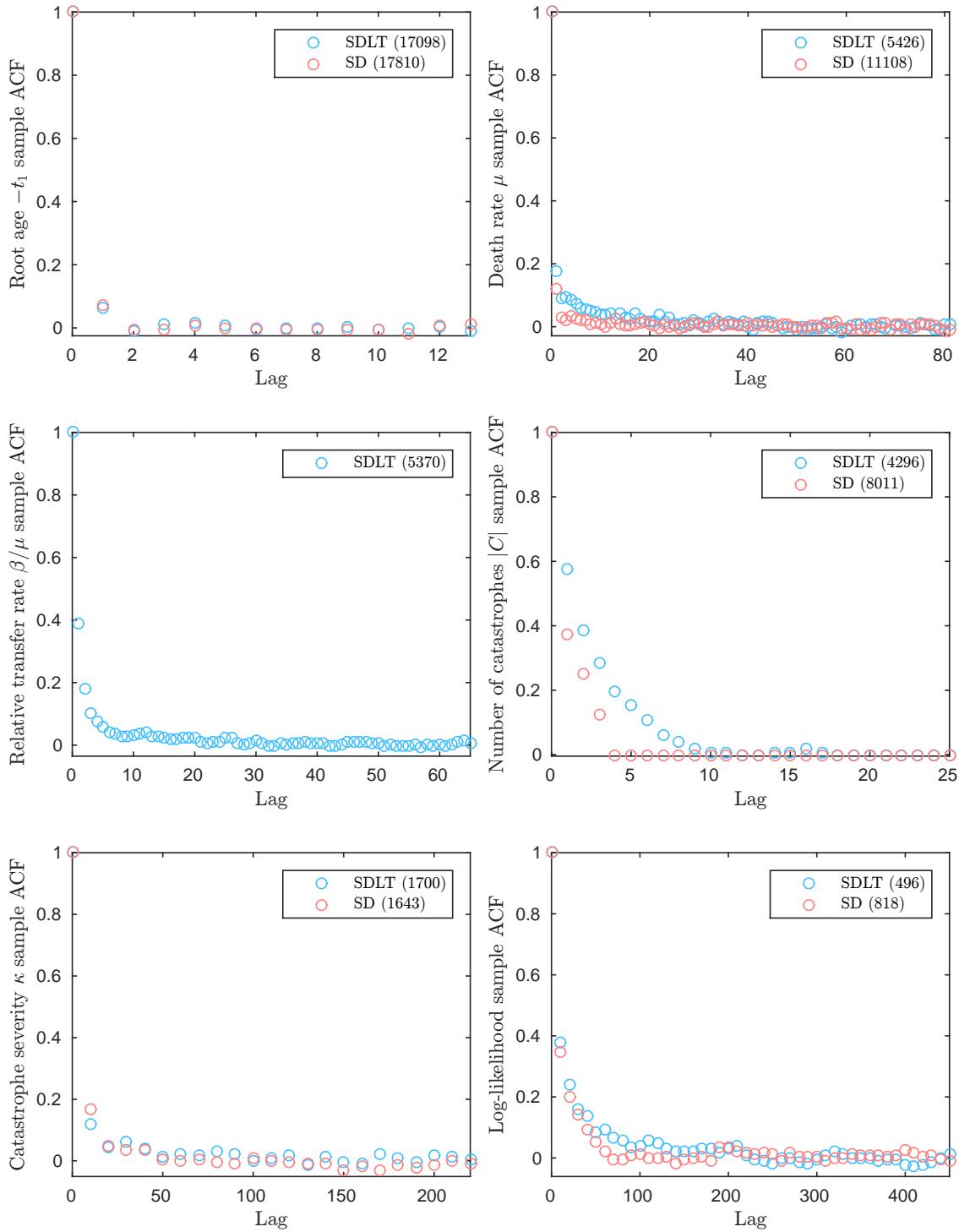


Figure 25: Sample autocorrelation plots in our analyses of POLY-0.

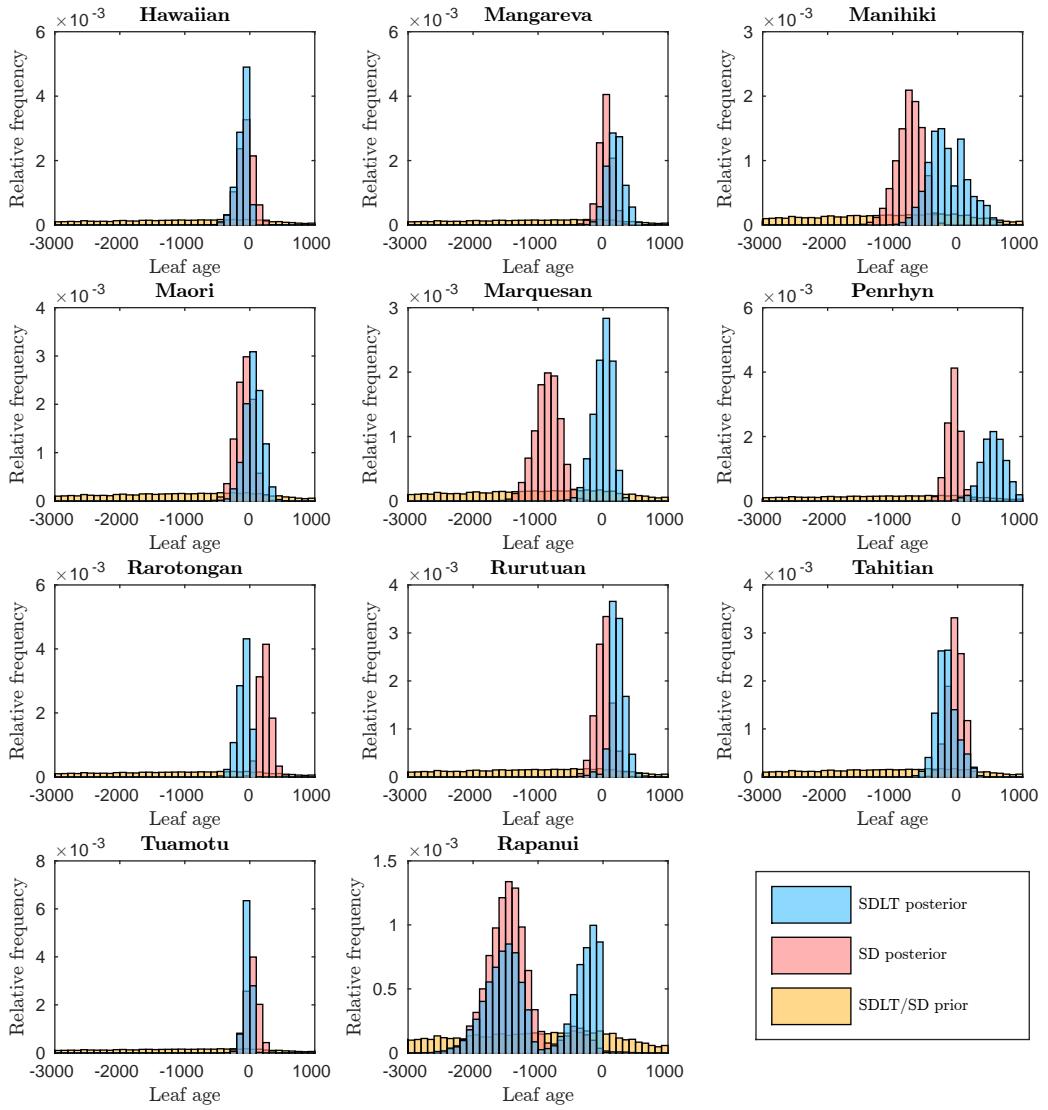


Figure 26: We relax the constraint on each leaf in turn and compute the histogram of its time under the prior and posterior for each model fit to POLY-0. Time in years is on the horizontal axis and relative frequency on the vertical axis.

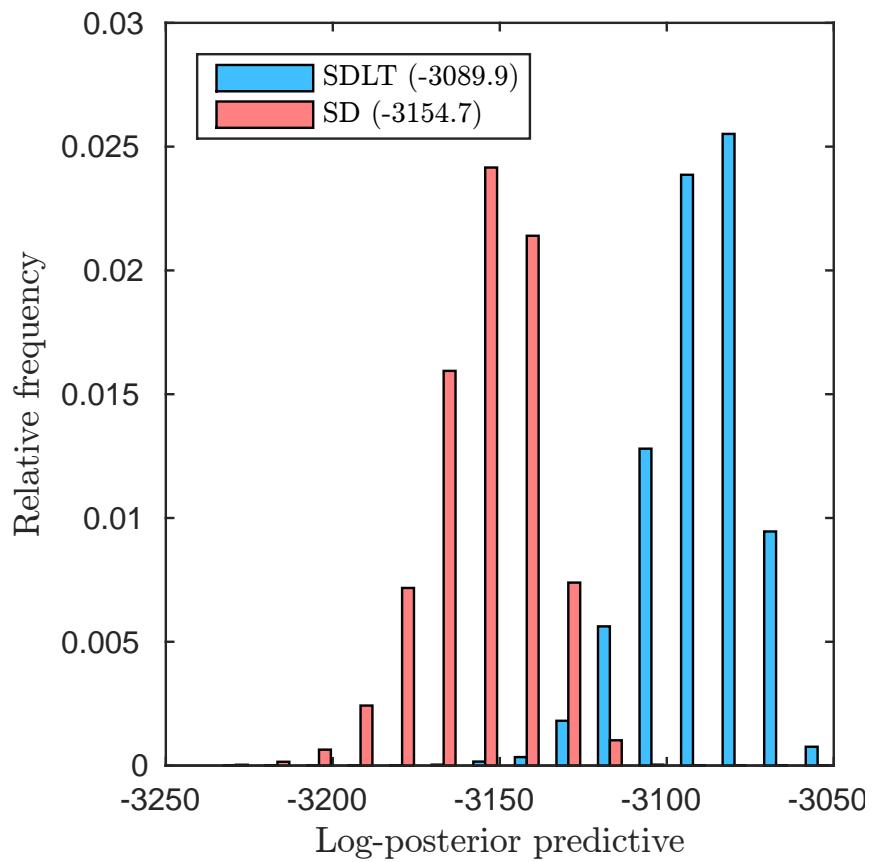


Figure 27: Assessing posterior predictive performance of the SDLT and SD models fit to POLY-0/train and tested on POLY-0/test. Predictive scores appear in parentheses.