

# Inferring Epidemiological Dynamics with Bayesian Coalescent Inference: The Merits of Deterministic and Stochastic Models

Alex Popinga,<sup>\*,†</sup> Tim Vaughan,<sup>\*,†,\*</sup> Tanja Stadler,<sup>§,\*\*</sup> and Alexei J. Drummond<sup>\*,†,1</sup>

<sup>\*</sup>Department of Computer Science, University of Auckland, Auckland, New Zealand 1010, <sup>†</sup>Allan Wilson Centre for Molecular Ecology and Evolution and <sup>‡</sup>Massey University, Palmerston North, New Zealand 4442, <sup>§</sup>Department of Biosystems Science and Engineering, ETH Zürich, Basel, Switzerland 4058, and <sup>\*\*</sup>Swiss Institute of Bioinformatics (SIB), Switzerland

ORCID ID: 0000-0003-4454-2576 (A.J.D.)

**ABSTRACT** Estimation of epidemiological and population parameters from molecular sequence data has become central to the understanding of infectious disease dynamics. Various models have been proposed to infer details of the dynamics that describe epidemic progression. These include inference approaches derived from Kingman's coalescent theory. Here, we use recently described coalescent theory for epidemic dynamics to develop stochastic and deterministic coalescent susceptible–infected–removed (SIR) tree priors. We implement these in a Bayesian phylogenetic inference framework to permit joint estimation of SIR epidemic parameters and the sample genealogy. We assess the performance of the two coalescent models and also juxtapose results obtained with a recently published birth–death–sampling model for epidemic inference. Comparisons are made by analyzing sets of genealogies simulated under precisely known epidemiological parameters. Additionally, we analyze influenza A (H1N1) sequence data sampled in the Canterbury region of New Zealand and HIV-1 sequence data obtained from known United Kingdom infection clusters. We show that both coalescent SIR models are effective at estimating epidemiological parameters from data with large fundamental reproductive number  $R_0$  and large population size  $S_0$ . Furthermore, we find that the stochastic variant generally outperforms its deterministic counterpart in terms of error, bias, and highest posterior density coverage, particularly for smaller  $R_0$  and  $S_0$ . However, each of these inference models is shown to have undesirable properties in certain circumstances, especially for epidemic outbreaks with  $R_0$  close to one or with small effective susceptible populations.

**KEYWORDS** Bayesian inference; phylodynamics; coalescent; epidemic; stochastic

## Phylodynamics and the Coalescent

**T**HE epidemiological and evolutionary processes that underpin rapidly evolving species occur on a shared spatiotemporal frame of reference. Unified analyses that include both the dynamics of an epidemic and the reconstruction of the pathogen phylogeny can therefore uncover otherwise inaccessible information to aid in outbreak prevention. Such information includes the rates of pathogen transmission and host recovery, effective population sizes, and the “time of origin” representing

the introduction of the first infected individual into a population of susceptible hosts.

The term *phylodynamics* was popularized by Grenfell *et al.* (2004) to describe the interlaced study of immunodynamics, epidemiology, and evolutionary mechanisms. Several phylodynamic models, both stochastic and deterministic in nature, have since been developed to characterize the phylogenetic history of the pathogen species and compartmentalizations of the host population throughout the epidemic. Such models grant the ability to infer key epidemiological parameters from genetic sequence data and include birth–death branching processes (Stadler *et al.* 2012, 2013; Gavryushkina *et al.* 2014; Kühnert *et al.* 2014), as well as coalescent approaches (Griffiths and Tavaré 1994; Pybus *et al.* 2001; Rasmussen *et al.* 2011, 2014; Koelle and Rasmussen 2012; Dearlove and Wilson 2013) derived from Kingman's coalescent theory (Kingman 1982).

Significant steps toward the unification of epidemiology and statistical phylogenetics were made by Pybus *et al.* (2001),

Copyright © 2015 by the Genetics Society of America

doi: 10.1534/genetics.114.172791

Manuscript received June 30, 2014; accepted for publication December 13, 2014; published Early Online December 19, 2014.

Available freely online through the author-supported open access option.

Supporting information is available online at <http://www.genetics.org/lookup/suppl/doi:10.1534/genetics.114.172791/-/DC1>.

<sup>1</sup>Corresponding author: Department of Computer Science, University of Auckland, 303-379, 38 Princes St., Auckland, New Zealand 1010.

E-mail: alexei@cs.auckland.ac.nz

Volz *et al.* (2009), and Dearlove and Wilson (2013), with the formalization and application of Kingman's  $n$ -coalescent to pathogen population dynamics. These methods involved numerical integration of a set of ordinary differential equations (ODEs) to find deterministic approximations to the variation in the number of sampled lineages through time. Volz (2012) extended the tree density calculation from previous work (Volz *et al.* 2009) to allow for serially sampled and spatially structured genetic sequence data. In this coalescent model, the birth and death rates can vary in time and by the state of the host, so that "the birth rate of a single gene copy is both time- and state-dependent" (Volz 2012, p. 7).

In this article, we assess the ability of coalescent-based phylodynamic models to infer, in a Bayesian setting, a range of epidemiological parameters from simulated data. While Dearlove and Wilson (2013) paved the way by implementing a coalescent approach for deterministic susceptible–infected (SI), susceptible–infected–susceptible (SIS), and susceptible–infected–removed (SIR) models for Bayesian inference, we implement and rigorously test both deterministic and stochastic coalescent SIR models of epidemic dynamics extended for heterochronously sampled data.

### Stochastic and Deterministic Models

Stochasticity and determinism in population sizes each maintain dominant roles in particular stages of an epidemic. Once the infected population has grown considerably large, on the order of 1000–10,000 lineages, the probability densities of stochastically expressed population size dynamics converge toward the deterministic interpretation (Rouzine *et al.* 2001). However, during the early stages the population size of infected individuals is small, and the dynamics of the epidemic are therefore governed by stochastic processes due to the relative significance of fluctuations in the demographic and rate parameters of the population model (Kühnert *et al.* 2014). Therefore, approximating the prevalence of infection by a deterministic function requires the number of infected hosts within the effective population to be assumed as very large throughout the duration of the described epidemic, *i.e.*, once the exponential growth phase has been reached (Rouzine *et al.* 2001).

Population size is critical to the epidemiological system and, as with any parameter in a Bayesian setting, yields the most accurate estimations when detailed prior information is available and incorporated into the inference (Drummond *et al.* 2006). In our extension and implementation of the coalescent model for epidemics, both stochastic and deterministic population size processes are used for the simulation of trees and/or trajectories for subsequent inference.

### Compartmental Population Models (SIR)

Host populations can be compartmentalized simply but effectively in mathematical models that describe epidemic progression. The specific division of the aggregate population depends on the contagion, spanning a range of scenarios

where hosts may or may not recover from infection, may or may not be reinfected, etc. Such examples include the SI, SIS, and SIR models (Anderson and May 1991; Keeling and Rohani 2008). Each of these compartments can be expressed either (a) by a set of ODEs that describe the deterministic time development of real-valued compartment occupancies or (b) in terms of integer-valued occupancies governed by continuous-time Markov chains (CTMC) that allow for a degree of uncertainty in the timing and number of events that occur over the course of the epidemic.

In this article, we concentrate on the SIR model, which describes epidemics that include infected individuals who are at some point in time removed from the effective population by way of immunity, death, behavioral changes, or some other termination of infectiousness. The deterministic variant of this model was introduced by Kermack and Mckendrick (1932) and is given by the trio of coupled ODEs,

$$\frac{d}{dt}S(t) = -\beta I(t)S(t), \quad (1)$$

$$\frac{d}{dt}I(t) = \beta I(t)S(t) - \gamma I(t), \quad (2)$$

$$\frac{d}{dt}R(t) = \gamma I(t), \quad (3)$$

where  $\beta$  and  $\gamma$  respectively represent the transition rates from susceptible  $S$  to infected  $I$  and infected  $I$  to removed  $R$ . The model fully defines the population dynamics with initial conditions  $S(z_0)$ ,  $I(z_0)$ , and  $R(z_0)$ . It is worth recognizing that, in the closed SIR model used here, there is no demographic change in the host population. Therefore,  $(d/dt)S(t) + (d/dt)I(t) + (d/dt)R(t) = 0$  and  $S(t) + I(t) + R(t) = N$ , where  $N$  is the constant total population size. Throughout this article we refer to the solutions to Equations 1–3 as *deterministic SIR trajectories*.

The comparable stochastic description is given in terms of the probability of the epidemic state at time  $t$  given its initial state and the rate parameters

$$\begin{aligned} \pi(s, i, r; t) \\ \equiv \Pr(S(t) = s, I(t) = i, R(t) = r | S(0), I(0), R(0), \beta, \gamma), \end{aligned} \quad (4)$$

which is governed by the following equation of motion:

$$\begin{aligned} \frac{d}{dt} \pi(s, i, r; t) \\ = \beta[(s+1)(i-1)\pi(s+1, i-1, r; t) - s i \pi(s, i, r; t)] \\ + \gamma[(i+1)\pi(s, i+1, r-1; t) - i \pi(s, i, r; t)]. \end{aligned} \quad (5)$$

An explicit sampling process is incorporated by allowing each removal event to coincide with a sampling event with a fixed probability  $\psi/(\psi + \mu)$ , where  $\psi$  and  $\mu$  are the overall rates of sampled and unsampled removals, respectively, such that  $\gamma = \psi + \mu$ . We refer to epidemic histories sampled from this model as *stochastic SIR trajectories*.

Both types of epidemic trajectories can be related to models of sampled transmission tree genealogies. In the deterministic case, this relationship is made via the coalescent distributions described in Volz (2012). We call this the *deterministic coalescent SIR model*. In the stochastic case, genealogies appear naturally from a branching process in which the branching events coincide with the transmission events in the CTMC and only those lineages ancestral to sampled removals are recorded. We call this the *stochastic SIR model*.

Another way of relating the stochastic SIR model to sampled transmission trees involves drawing a realization of a stochastic SIR epidemic and then using the coalescent distribution in Volz (2012) to produce a tree conditional on the particular piecewise constant infected compartment size corresponding to that realization. We call this approach the *stochastic coalescent SIR model*. Unlike BDSIR, the stochastic coalescent SIR model does not require the sampling process to be specified explicitly.

Both the transmission rate  $\beta$  and the removal rate  $\gamma$  can be estimated using each of the methods considered in this article from data ascribed to an SIR epidemic.

## Methods

### Inference framework

All phylodynamic inference discussed in this article is based on the joint posterior probability density

$$f(\mathcal{T}, \mathcal{V}, \eta, \theta | D) = \frac{\Pr(D | \mathcal{T}, \theta) f(\mathcal{T} | \mathcal{V}, \eta) f(\mathcal{V} | \eta) f(\eta) f(\theta)}{\Pr(D)}, \quad (6)$$

where the sampled transmission tree  $\mathcal{T}$ , the epidemic trajectory denoted  $\mathcal{V} = (\mathcal{S}, \mathcal{I}, \mathcal{R})$ , the substitution parameters  $\theta$ , and the epidemiological parameters  $\eta = \{\beta, \gamma, S_0, z_0\}$  are all estimated from the sequence data. The sampled transmission tree  $\mathcal{T}$  is assumed to be identical to the pathogen genealogy.

Here,  $\mathcal{S}$ ,  $\mathcal{I}$ , and  $\mathcal{R}$  represent the host compartment sizes from the present time  $\tau = 0$  back to the origin  $z_0$ , such that  $S(\tau) = S(z_0 - \tau)$ ,  $\mathcal{I}(\tau) = I(z_0 - \tau)$ , and  $\mathcal{R}(\tau) = R(z_0 - \tau)$ .

The various terms making up the right-hand side of Equation 6 are the tree likelihood  $\Pr(D | \mathcal{T}, \theta)$ , the tree prior  $f(\mathcal{T} | \mathcal{V}, \eta)$ , the epidemic trajectory density  $f(\mathcal{V} | \eta)$ , and the substitution and epidemiological parameter priors  $f(\eta)$  and  $f(\theta)$ . The probability  $\Pr(D)$  is merely a normalizing constant and can be ignored. It is the product of the tree prior and trajectory density  $f(\mathcal{T} | \mathcal{V}, \eta) f(\mathcal{V} | \eta)$  that distinguishes each of the models considered in this article.

For both the deterministic and stochastic coalescent SIR models, the tree prior  $f(\mathcal{T} | \mathcal{V}, \eta)$  is calculated in the following way. First, consider the time span of a tree divided into segments bracketed by both sampling and coalescent events. By considering intervals ending in sampling events as well as coalescent-ending intervals, we follow previous work that extended coalescent approaches to time-stamped, serially sampled data (Rodrigo and Felsenstein 1999; Drummond *et al.* 2002). Interval  $i$  is spanned by  $k_i$  lineages and is the  $i$ th interval

when ordered from the most recent tip to the root. The set of intervals  $A$  ending in sample events and the set of intervals  $Y$  ending in coalescent events together encompass all intervals,  $V = A \cup Y$ . Let the end time of an interval be  $\tau_i$  (going back in time), with  $\tau_0 = 0$  as the time of the most recent tip and with time increasing into the past. Then the probability density of a genealogy given an epidemic trajectory is

$$f(\mathcal{T} | \mathcal{V}, \eta) = \prod_{i \in Y} \lambda_{k_i}(\tau_i) \prod_{i \in V} \omega(\tau_i, k_i), \quad (7)$$

where  $\lambda_{k_i}(\tau)$  is the instantaneous coalescent rate at  $\tau$  prescribed by Volz (2012),

$$\lambda_{k_i}(\tau) = \binom{k_i}{2} \frac{2\beta S(\tau)}{\mathcal{I}(\tau)}, \quad (8)$$

and where  $\omega(\tau_i, k_i)$  is the survival probability

$$\omega(\tau_i, k_i) = \exp\left(-\int_{\tau_{i-1}}^{\tau_i} \lambda_{k_i}(\tau) d\tau\right). \quad (9)$$

The deterministic coalescent SIR model assumes that the SIR epidemic trajectories are found by integrating the ODEs in Equations 1–3. Therefore, under this model each epidemic trajectory is a deterministic function of its parameters  $\mathcal{V}(\eta)$ . This means that the trajectory density can be written as

$$f(\mathcal{V} | \eta) = \delta(\mathcal{V} - \mathcal{V}(\eta)), \quad (10)$$

where  $\delta(x)$  is the Dirac  $\delta$ -function and represents a point mass concentrated at  $x = 0$ .

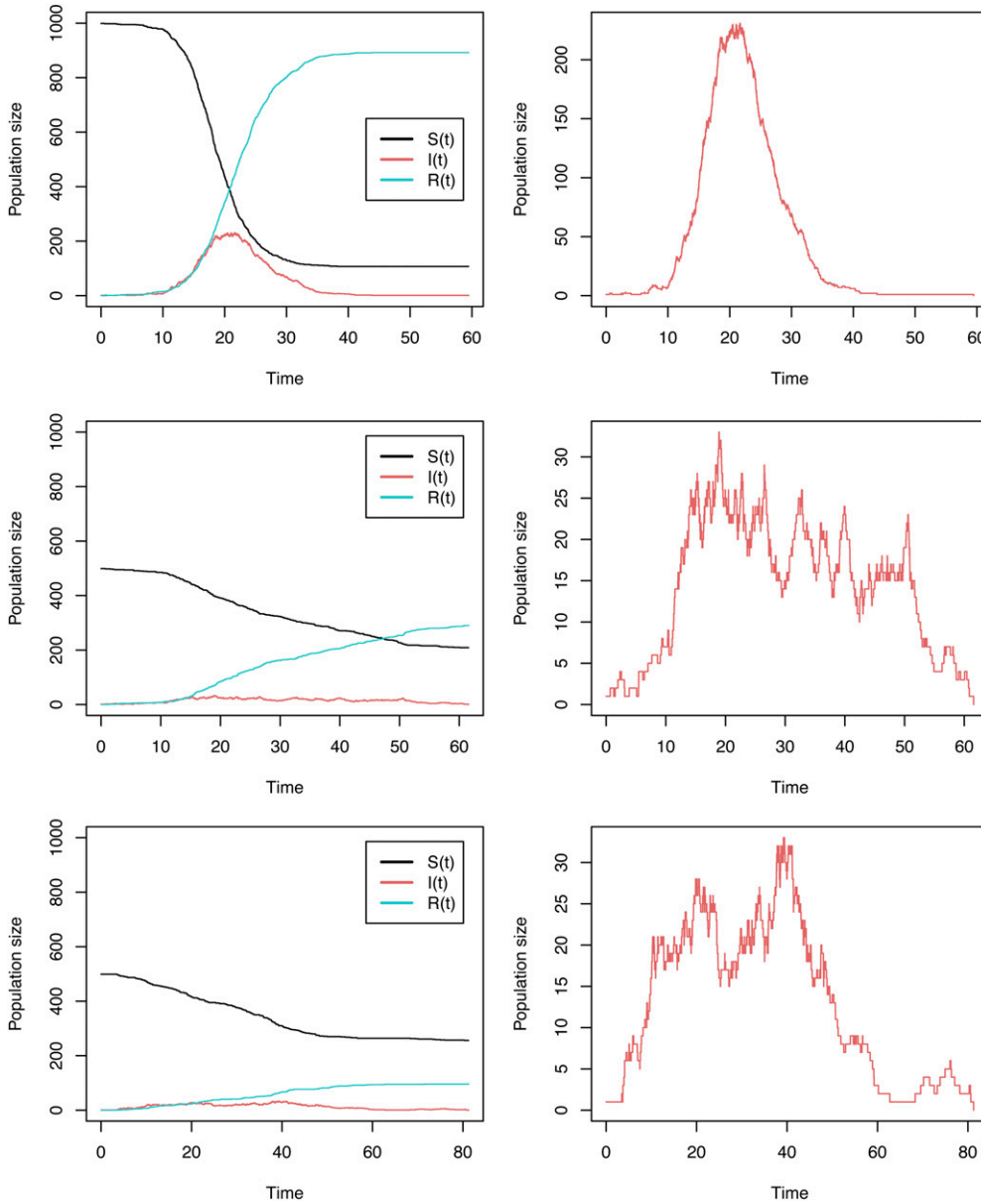
In contrast, the stochastic coalescent SIR model assumes that the epidemic is generated by a jump process corresponding to the master equation given in Equation 5. In this case, the probability  $f(\mathcal{V} | \eta)$  is nonsingular and thus contributes to the uncertainty in the final inference result.

In the BDSIR model introduced by Kuhnert *et al.* (2014),  $f(\mathcal{V} | \eta)$  is the same as for the stochastic coalescent SIR model, but  $f(\mathcal{T} | \mathcal{V}, \eta)$  is defined differently. See Kuhnert *et al.* (2014) for details.

### Markov chain Monte Carlo algorithm

We use Markov chain Monte Carlo (MCMC) to sample from the joint posterior density given in Equation 6. Many of the specifics of the algorithm used have been discussed previously, in particular the method for calculating the tree likelihood (Felsenstein 1981, 2004) and the mechanism for exploring tree space (Drummond *et al.* 2002). However, the model-specific product  $f(\mathcal{T} | \mathcal{V}, \eta) f(\mathcal{V} | \eta)$  requires special attention.

As we are primarily interested in parametric inference rather than the epidemic trajectory itself, we can regard  $\mathcal{V}$  as a nuisance parameter to be marginalized over. This marginalization can be achieved implicitly by sampling it using MCMC and then ignoring this component of the sampled state, which is the strategy we use when reporting the BDSIR results. It can also be made an explicit part of the likelihood calculation, which is the approach we take with the deterministic



**Figure 1** Stochastic SIR trajectories for susceptible  $S$ , infected  $I$ , and recovered  $R$  populations, with (top row)  $S_0 = 999$  and  $R_0 = 2.4975$ , (middle row)  $S_0 = 499$  and  $R_0 = 1.497$ , and (bottom row)  $S_0 = 499$  and  $R_0 = 1.0978$ . (The right column shows infected  $I$  only.)

and stochastic coalescent SIR models. This marginalization means that the product  $f(\mathcal{T}|\mathcal{V}, \eta)f(\mathcal{V}|\eta)$  becomes

$$f(\mathcal{T}|\eta) = \int f(\mathcal{T}|\mathcal{V}, \eta)f(\mathcal{V}|\eta)d\mathcal{V}, \quad (11)$$

the probability density of the tree given the epidemiological parameters.

In the case of the deterministic coalescent SIR model, this density reduces to  $f(\mathcal{T}|\mathcal{V}(\eta), \eta)$ , meaning that the density of the tree given epidemiological parameters  $\eta$  is obtained simply by substituting the numerical solution to Equations 1–3 for those parameters into Equation 7.

The stochastic coalescent SIR model is more complex, as in this case the trajectory density  $f(\mathcal{V}|\eta)$  is nonsingular, meaning that computing the integral in Equation 11 is non-trivial. We treat this here using the “pseudomarginal” ap-

proach (Beaumont 2003; Andrieu and Roberts 2009) in which, at each step in the MCMC chain, the marginalized tree density  $f(\mathcal{T}|\eta)$  is replaced by the Monte Carlo estimate

$$\hat{f}(\mathcal{T}|\eta) = \frac{1}{M} \sum_{r=1}^M f(\mathcal{T}|\mathcal{V}_r, \eta), \quad (12)$$

where each  $\mathcal{V}_r$  is a trajectory sampled independently from  $f(\mathcal{V}|\eta)$ , using a stochastic simulation algorithm (Sehl *et al.* 2009). Perhaps counterintuitively within an MCMC framework, this stochastic likelihood converges to the true marginal posterior distribution regardless of the number  $M$  of realizations used in the estimate. However, the magnitude of  $M$  can significantly affect the rate at which the chain produces effectively independent samples from the posterior and must be tuned carefully.

## Implementation and validation

We have implemented the schemes described above for performing inference under the deterministic and stochastic coalescent SIR models within the BEAST 2 phylogenetics package found at <http://github.com/CompEvol/phylogenetics>. This has a number of advantages over a stand-alone implementation. Foremost, we were able to avoid reimplementing components of the algorithm that are in common with other already-implemented phylogenetic and phylodynamic analyses, such as the MCMC proposal operators used to traverse the parameter space. Furthermore, this greatly increases the usefulness of the implementation, as it can be immediately used in conjunction with a wide variety of nucleotide and amino acid substitution models and parameter priors.

We have taken two steps to ensure our implementation is correct. First, we compared tree probability density  $f(T|\mathcal{V}, \eta)$  values calculated using the main implementation of each of the two models with those calculated using completely independent implementations in R (R Core Team 2014).

Second, we used the implemented MCMC algorithms to sample transmission trees from the tree density given in Equation 11 for each model. We then compared the distributions of tree height, total edge length, and binary clade count summary statistics from these sampled ensembles with sample distributions obtained directly via stochastic simulation. As shown in Supporting Information, File S1, Figure S1, Figure S2, and Figure S3 (*Sampling from the prior*) and in the associated figures, the resulting pairs of distributions agree, providing strong support for our claim that the implementations of the methods described above are correct.

Instructions for downloading and using this package are also available on the project website located at <http://github.com/CompEvol/phylogenetics>.

## Simulation study

To evaluate the implementation and extension of the coalescent models, we performed analyses on both sequence data and fixed trees simulated with known parameter values. The median estimated values produced by each model were then used to measure relative error and bias, along with the widths and coverage of 95% highest posterior density (HPD) intervals.

We used three methods for simulating the trees and trajectories, as shown below:

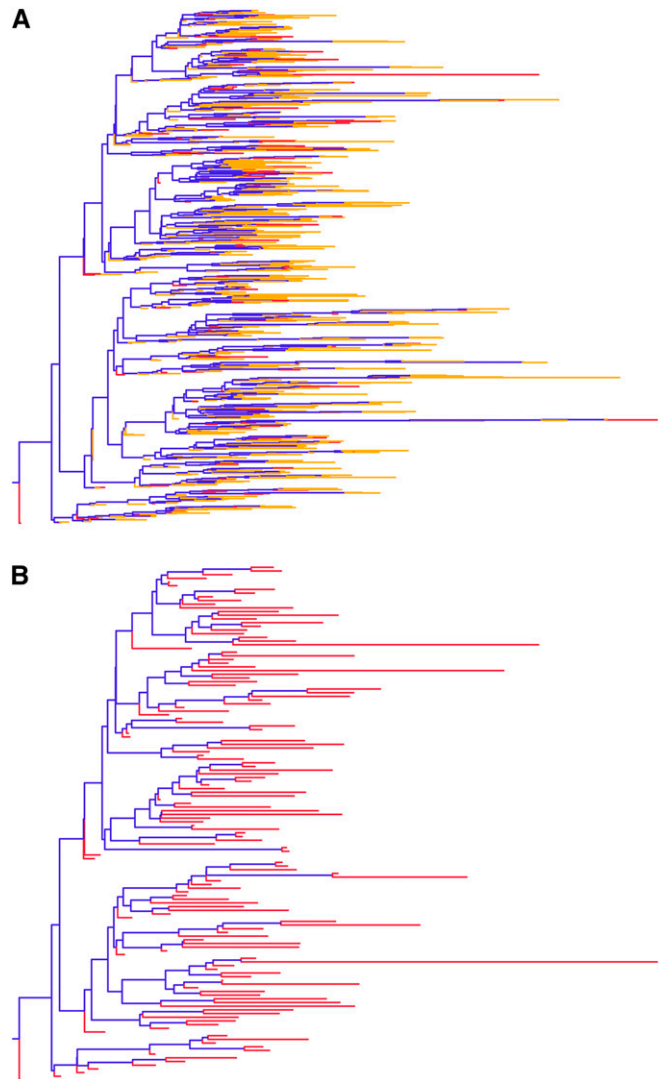
Inference model:

{ Stoch. Coal. SIR	Deter. Coal. SIR	BDSIR
--------------------	------------------	-------

Simulation scheme:

{	Stoch. Coal. SIR	Stoch. Coal. SIR	Stoch. Coal. SIR
	Deter. Coal. SIR	Deter. Coal. SIR	Deter. Coal. SIR
	Stochastic SIR	Stochastic SIR	Stochastic SIR.

The stochastic coalescent and deterministic coalescent simulation schemes were used to validate the coalescent



**Figure 2** (A) Full stochastic SIR transmission tree with both sampled  $\psi$ -tips, shown in red, and otherwise removed  $\mu$ -tips, shown in yellow. (B) The corresponding 140-tip sampled stochastic SIR tree. A and B were generated in FigTree (Rambaut 2007).

SIR inference models. The *stochastic SIR* scheme, contrarily, is emphasized for its realistic properties.

Stochastic SIR trees and trajectories were generated using master equations in the simulation package MASTER (Vaughan and Drummond 2013). Deterministic coalescent trajectories were generated using a Runge–Kutta integrator (Runge 1895; Kutta 1901) with adaptive step sizes to solve a system of first order ODEs. Stochastic coalescent trajectories were generated using Sehl *et al.*'s (2009) SAL  $\tau$ -leaping algorithm (Sehl *et al.* 2009).

To simulate the stochastic coalescent SIR trees, we used the stochastic SIR trajectories, which could be converted to effective population size with the mathematical expression used to obtain Volz's (2012) coalescent rate for the SIR model:  $N_e(\tau) = 1/\lambda_2(\tau) = \mathcal{I}(\tau)/(2\beta\mathcal{S}(\tau))$ . The sampling times, generated by a sampling rate  $\psi$ , for the stochastic

**Table 1 Results for simulated sequences:  $R_0 \approx 2.50$ ,  $S_0 = 999$**

$\eta$	Inference	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy (%)
$R_0$	Stoch.Coal.SIR	2.50	2.41	2.16	0.13	-0.11	0.97	97.00
	Deter.Coal.SIR	2.50	2.78	2.03	0.38	0.05	0.79	87.00
	BDSIR	2.50	3.21	2.84	0.15	0.14	1.86	100.00
$\gamma$	Stoch.Coal.SIR	0.30	0.16	0.13	0.52	-0.52	0.82	47.00
	Deter.Coal.SIR	0.30	0.25	0.16	0.56	-0.28	0.97	56.00
	BDSIR	0.30	0.17	0.14	0.52	-0.52	1.13	84.00
$S_{(0)}$	Stoch.Coal.SIR	999	1805	1148	0.32	0.21	5.12	99.00
	Deter.Coal.SIR	999	2384	1565	0.66	0.60	6.54	100.00
	BDSIR	999	4002	2611	1.70	1.70	10.38	99.00
$Z_{(0)}$	Stoch.Coal.SIR	Varies	51.67	48.89	0.26	0.23	0.61	37.00
	Deter.Coal.SIR	Varies	49.13	46.46	0.22	0.20	0.26	29.00
	BDSIR	Varies	31.16	29.52	0.51	0.51	0.79	18.00

coalescent SIR trees were also taken from the MASTER output to allow for direct comparison between the sets of trees. In other words, the underlying epidemic function was the same for both stochastic SIR and stochastic coalescent SIR trees, the latter of which were then simulated under a piecewise constant population function.

Likewise, for the simulation of deterministic coalescent trees we used deterministic SIR trajectories to construct a population function and the relation  $N_e = \mathcal{I}/(2\beta S)$  to convert infected and susceptible host population sizes to effective population size. The sampling times were randomly generated from a probability distribution so that the density of samples taken through time was proportional to the number of infected individuals through time, as with the stochastic SIR trees.

We simulated stochastic SIR trees, using multiple combinations of parameter values. We were particularly interested in varying the basic reproductive ratio  $R_0$  and the initial susceptible population size  $S_0$ , to observe the changes in relative error, bias, and uncertainty in stochastic and deterministic models. To alter the ratio  $R_0 = \beta S_0/\gamma$  and still generate sensible trees with a consistent number of tips, one or more of the other parameters (birth rate  $\beta$ , removal rate  $\gamma$ , or  $S_0$ ) must also change. Table 2, Table S6, Table S7, and Table S9 show the true values of the parameters for each set of simulations. (The birth rate  $\beta$  is not shown, as our implementation allows either  $\beta$  or  $R_0$  to serve as a parameter in the inference, and  $R_0$  is the parameter of interest. However,  $\beta$  can be calculated via the other three, using  $\beta = R_0\gamma/S_0$ . For example, when  $R_0 = 1.0978$ ,  $S_0 = 499$ , and  $\gamma = 0.25$ , then  $\beta = 5.50E-4$ .)

**Heterochronous trees:** We generated 100 trees under each of the three (stochastic SIR, stochastic coalescent SIR, and deterministic coalescent SIR) models with parameters  $S_0$ ,  $\beta$ , and  $\gamma$ . For heterochronously sampled trees, each removal generates a sample with probability  $\psi/(\psi + \mu)$ , where  $\psi$  is the overall rate of sampled removals and  $\mu$  is the rate of unsampled removals such that  $\gamma = \psi + \mu$ .

The simulations ended once the number of infected individuals reached zero, *i.e.*, when the last infected individ-

ual was removed. This ensured that the simulated trajectories spanned past the exponential growth phase of the epidemic and therefore included samples past the peak of infected individuals. This choice of procedure was motivated by (a) the suggestion of Stadler *et al.* (2014) that the behavior of the coalescent beyond the exponential phase could either inflate or reduce bias and (b) the observations of Dearlove and Wilson (2013) and Bošková *et al.* (2014) that deterministic coalescent SIR models might be properly fitted only once the epidemic has peaked. Figure 1 shows trajectories of susceptible, infected, and removed individuals underlying the simulation of stochastic SIR trees (Figure 2) generated in MASTER. An example XML for simulating these MASTER trees is provided in File S1.

We required that the trees had  $n \geq 100$  leaves, filtering out those in which the epidemic died out in the early stages, *i.e.*, when the initial infected individual was removed from the effective population too quickly to infect others. (Note that the inference procedures discussed in this article all implicitly condition on the number of leaves.) The probability that the first event in a given trajectory is the removal (by recovery, death, etc.) of patient zero is given by  $\delta/(\beta S_0 + \delta) = 1/(1 + R_0)$ . When  $R_0 \approx 2.50$ , this probability is  $\approx 30\%$ . In our case, 52/152 ( $\approx 34\%$ ) trees were “empty” or containing only one node. The filtering process left us with a mean of  $\approx 160$  leaves for the simulated trees.

**Homochronous trees:** A major concern in the comparison between Kühnert *et al.* (2014)’s birth–death–sampling SIR inference model, which includes explicit sampling, and our implementations of Volz (2012)’s coalescent SIR models, which do not include explicit sampling, is that the former is given extra information via the sampling process. Volz and Frost (2014) addressed this issue by providing a coalescent SIR model that does incorporate sampling explicitly.

That being said, results from Bošková *et al.* (2014) indicate that the poor performance of the deterministic coalescent SIR model in comparison with birth–death models was due to the lack of handling stochastic population size changes through time rather than the lack of information about the sampling proportion. Their results showed that

**Table 2 Simulation study results for fixed trees:  $R_0 \approx 2.50$  and  $S_0 = 999$ ,  $R_0 \approx 1.50$  and  $S_0 = 499$ , and  $R_0 \approx 1.10$  and  $S_0 = 499$**

$\eta$	Inference	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy (%)
$R_0$	Stoch.Coal.SIR	2.50	2.84	2.68	0.12	0.09	0.98	100.00
	Deter.Coal.SIR	2.50	2.68	2.49	0.13	0.04	0.81	98.00
	BDSIR	2.50	2.73	2.67	0.12	0.08	0.55	94.00
$\gamma$	Stoch.Coal.SIR	0.30	0.27	0.25	0.19	-0.13	1.14	99.00
	Deter.Coal.SIR	0.30	0.32	0.29	0.16	3.14E-3	1.27	99.00
	BDSIR	0.30	0.28	0.27	0.13	-0.09	0.62	95.00
$S_{(0)}$	Stoch.Coal.SIR	999	1390	921	0.19	-0.03	3.85	100.00
	Deter.Coal.SIR	999	1807	1133	0.52	0.29	4.59	98.00
	BDSIR	999	1591	1142	0.39	0.24	3.42	99.00
$z_{(0)}$	Stoch.Coal.SIR	Varies	41.81	40.35	0.03	0.01	0.20	99.00
	Deter.Coal.SIR	Varies	41.17	39.99	0.03	0.01	0.07	76.00
	BDSIR	Varies	40.89	39.72	8.65E-4	-5.13E-4	3.43E-3	97.00
$R_0$	Stoch.Coal.SIR	1.50	1.48	1.37	0.09	-0.06	0.81	100.00
	Deter.Coal.SIR	1.50	1.80	1.49	0.24	0.15	0.52	85.00
	BDSIR	1.50	1.46	1.43	0.08	-0.03	0.47	99.00
$\gamma$	Stoch.Coal.SIR	0.30	0.19	0.17	0.40	-0.40	1.06	85.00
	Deter.Coal.SIR	0.30	0.26	0.23	0.27	-0.22	1.15	89.00
	BDSIR	0.30	0.26	0.25	0.18	-0.18	0.72	97.00
$S_{(0)}$	Stoch.Coal.SIR	499	599	390	0.25	-0.22	3.56	100.00
	Deter.Coal.SIR	499	562	361	0.44	-0.26	3.36	91.00
	BDSIR	499	996	714	0.51	0.49	4.63	100.00
$z_{(0)}$	Stoch.Coal.SIR	Varies	76.47	68.24	0.55	0.54	0.58	99.00
	Deter.Coal.SIR	Varies	91.03	72.51	0.39	0.38	0.42	88.00
	BDSIR	Varies	69.11	66.51	0.34	-0.31	0.20	94.00
$R_0$	Stoch.Coal.SIR	1.10	1.39	1.32	0.22	0.22	1.09	99.00
	Deter.Coal.SIR	1.10	1.68	1.44	0.46	0.46	0.59	25.00
	BDSIR	1.10	1.34	1.32	0.20	0.20	0.51	75.00
$\gamma$	Stoch.Coal.SIR	0.25	0.17	0.15	0.37	-0.36	1.11	84.00
	Deter.Coal.SIR	0.25	0.22	0.18	0.30	-0.22	1.16	86.00
	BDSIR	0.25	0.28	0.26	0.12	0.09	0.92	100.00
$S_{(0)}$	Stoch.Coal.SIR	499	608	398	0.24	-0.18	3.38	100.00
	Deter.Coal.SIR	499	553	337	0.42	-0.26	3.08	92.00
	BDSIR	499	1471	1040	1.21	1.21	6.52	99.00
$z_{(0)}$	Stoch.Coal.SIR	Varies	91.60	84.55	0.06	0.02	0.60	97.00
	Deter.Coal.SIR	Varies	112.79	90.37	0.26	0.26	0.94	85.00
	BDSIR	Varies	82.98	80.93	0.02	-0.01	0.08	88.00

the coalescent is “very robust to changes in sampling schemes” (Boskova *et al.* 2014, p. 8).

Regardless, to ensure a fair comparison of BDSIR and the coalescent SIR models, we simulated an SIR epidemic with homochronous, or contemporaneous, sampling. This type of simulation affords no additional information about the population size for explicit-sampling models, as there is only a single time of sampling.

We selected a simulation time of  $t = 20$  for the homochronously sampled trees, with the trajectories being sampled at high prevalence but also past the time of peak prevalence. This is important for distinguishing SIR from SI/SIS outbreaks, as it provides information about the removal parameter  $\gamma$ . In this set of simulations, each lineage was sampled at  $t = 20$  with probability 0.7 (the leaf count distribution for varied sampling probabilities is in File S1).

**Simulated sequences:** To assess the ability of each SIR model to infer epidemic parameters with the inclusion of phylogenetic uncertainty, we also simulated the evolution of

2000-bp sequences down each simulated tree. We time stamped the sequences with the tip dates of each corresponding tree and informed the inference with the true Hasegawa–Kishino–Yano (HKY) substitution model (Hasegawa *et al.* 1985), clock rate =  $5E-3$ , and  $\kappa = 5$ . These choices were made to reflect real data, specifically those of influenza (Vaughan *et al.* 2014).

Along with simulated sequence data, analyses were performed with the simulated trees fixed (results are in File S1), and the parameters  $R_0$ ,  $\gamma$ ,  $S_0$ , and the origin of the tree  $z_0$  were estimated with Bayesian prior distributions as listed in Table 4.

**Deterministic coalescent SIR on higher  $R_0$  and  $S_0$ :** Finally, we had particular interest in the effects of varying the population size parameter  $S_0$  on the deterministic coalescent SIR model, as comparisons from initial analyses with lower true  $R_0$  ( $\approx 1.5$  and  $\approx 1.1$ ) and  $S_0$  ( $= 499$ ) showed higher error and bias and lower 95% HPD coverage. Also, it is often assumed that deterministic descriptions will perform well

**Table 3 Epidemic parameter inference from H1N1 sequences in New Zealand**

Inference model	$R_0$	$\gamma$	$S_0$	Root of the tree (yr)	Origin $z_0$ of the epidemic (yr)
Stoch. Coal. SIR	1.46 (1.04–2.14)	27.08 (4.20–64.03)	6.90E4 (175–2.86E5)	0.53 (0.44–0.61)	0.69 (0.45–1.03)
Deter. Coal. SIR	1.35 (1.05–1.84)	34.50 (3.86–82.16)	1.20E5 (29–4.59E5)	0.54 (0.45–0.62)	0.73 (0.47–1.04)
BDSIR	1.61 (1.09–2.29)	27.72 (6.82–55.04)	2.22E4 (259–9.38E4)	0.49 (0.41–0.56)	0.53 (0.43–0.65)

Shown are mean estimates (and 95% HPD intervals) of each epidemic parameter inferred from seasonal influenza A (H1N1) sequence data collected in the Canterbury region of New Zealand throughout the 2001 flu season.

for higher  $R_0$  and larger population sizes. Table S7 and Table S9 detail the parameter values we used to explore the behavior of the deterministic coalescent on varied  $R_0$  and  $S_0$  combinations.

### Interpretation of results

We compared the coalescent SIR, as well as BDSIR, parameter estimations from the simulated data to the true values used to generate the SIR trajectories. Following Kühnert *et al.* (2014), the precision and accuracy of these methods were measured by relative error, bias, and HPD intervals. We used the posterior median value of the parameter value  $\hat{\eta}$  compared with the true parameter  $\bar{\eta} \in \{R_0, \gamma, S_0, z_0\}$ . Relative error and bias are then gauged by calculating the median value over medians from all 100 trees, such that

$$RE_{\hat{\eta}} = \frac{\sum_{\tau=1}^{100} |\hat{\eta} - \bar{\eta}| / \bar{\eta}}{100}$$

and

$$RB_{\hat{\eta}} = \frac{\sum_{\tau=1}^{100} |\hat{\eta} - \bar{\eta}| / \bar{\eta}}{100}.$$

Measures of HPD interval widths are given by

$$\frac{95\% \text{ HPD upper bound} - 95\% \text{ HPD lower bound}}{\bar{\eta}}.$$

Table 1, Table 2, and Table 3 show these results, along with the percentages of posterior estimates that produced 95% HPD intervals containing the true values (*i.e.*, 95% HPD coverage).

### H1N1 data analysis

To test the efficacy of the coalescent SIR models on real data, epidemic parameters  $R_0$ ,  $\gamma$ ,  $S_0$ , and time of origin  $z_0$  were estimated from 42 seasonal influenza A (H1N1) sequences sampled throughout the 2001 flu season in Canterbury, New Zealand.

Influenza infections are well known for their seasonal SIR behavior in nonequatorial populations, as each annual flu season begins with a supply of susceptible hosts and tapers off as the hosts recover with adaptive immunity (Iwasaki and Pillai 2014). Due partly to this seasonal pattern, the influenza virus is both a motivator for the development of specialized models and a prime subject for testing phylodynamic models (Koelle *et al.* 2006).

Sampling a particular region bypasses the necessity of specifying geographically structured populations, and New

Zealand is an area of particular interest due to its geographic location and relative isolation from other regions with potentially varying dynamics. It is also assumed to play a key role in the global circulation of influenza strains (Rambaut and Holmes 2009; Bedford *et al.* 2010).

We used an HKY nucleotide substitution model, with a substitution rate of  $5E-3$  as estimated in Vaughan *et al.* (2014), and informed the models with dated sequences. Priors used for the Bayesian inference are shown in Table 4.

### HIV-1 data analysis

In addition to our analysis of H1N1 sequence data, we selected HIV-1 subtype B nucleotide sequences collected from infected individuals located in the United Kingdom. The coalescent SIR results were collated with the results from the BDSIR data analysis performed by Kühnert *et al.* (2014), using the same sequences. More details of this analysis are provided in File S1.

## Results and Discussion

### Simulation study

Results for epidemic parameter inference from nucleotide sequences simulated from stochastic SIR trees are provided in Table 1 for  $R_0 \approx 2.50$ . Results for inference from fixed trees ( $R_0 \approx 2.50$ ,  $R_0 \approx 1.50$ ,  $R_0 \approx 1.10$ ) are shown in Table 2, with 95% HPD coverage shown for each analysis in Figure 3. Inference results for analyses with true  $R_0 = 1.0987$  and varying population size ( $S_0 = 499, 999, 1999$ ) are described in Tables S1 and S2 in the supporting information, along with results from trees simulated under the stochastic and deterministic coalescent models for validation.

**Heterochronous trees:** For  $R_0 \approx 2.50$ , all three inference methods performed similarly for parameters  $R_0$  and  $\gamma$ , with high 95% HPD coverage and low error and bias. The most weakly identifiable parameter  $S_0$  yielded the largest HPD intervals for all three inference models. The deterministic coalescent returned higher error (0.52) and bias (0.29) than the stochastic coalescent SIR (0.19,  $-0.03$ ) and BDSIR (0.39, 0.24) and recovered the origin parameter  $z_0$  for only 76 of 100 simulated trees, while the stochastic coalescent and BDSIR respectively recovered  $z_0$  for 99 and 97 of 100 simulations.

For  $R_0 \approx 1.50$ , the relative HPD widths (akin to variance) for three of the four estimated parameters ( $R_0$ ,  $\gamma$ , and  $z_0$ ) were smallest for BDSIR. For the parameter  $S_0$ , the relative HPD width is largest for BDSIR, although it also had slightly higher



**Table 4 Bayesian prior distributions**

Annalysis	$R_0$	$\gamma$	$S_{(0)}$	$z_{(0)}$	$\psi/(\psi + \mu)$
$R_0 \approx 2.5, S_0 = 999$	LogN(1, 1)	LogN(-1, 1)	LogN(7, 1)	Unif(0, 100)	Beta(1, 1)
$R_0 \approx 1.5, S_0 = 499$	LogN(0.5, 1)	LogN(-1, 1)	LogN(6, 1)	Unif(0, 500)	Beta(1, 1)
$R_0 \approx 1.1, S_0 = 499$	LogN(0.1, 1)	LogN(-1.5, 1)	LogN(6, 1)	Unif(0, 500)	Beta(1, 1)
$R_0 \approx 1.1, S_0 = 999^a$	LogN(0.1, 1)	LogN(-1.5, 1)	LogN(7, 1)	Unif(0, 500)	—
$R_0 \approx 1.1, S_0 = 1999^a$	LogN(0.1, 1)	LogN(-1.5, 1)	LogN(7.5, 1)	Unif(0, 500)	—
$R_0 \approx 1.2, S_0 = 499^a$	LogN(0.2, 1)	LogN(-1, 1)	LogN(6, 1)	Unif(0, 500)	—
H1N1	Unif(0, 10)	LogN(3, 0.75)	LogN(13, 2)	Unif(0, 10)	Beta(1, 1)
HIV-1	LogN(1, 1)	LogN(-1, 1)	LogN(7, 1)	Unif(0, 100)	Beta(1, 1)

Shown are prior distributions for the reestimation of SIR parameters—the reproductive ratio  $R_0$ , the rate of removal  $\gamma$ , the number of susceptible individuals at the start of the epidemic  $S_{(0)}$ , the time of origin  $z_{(0)}$ , and the sampling proportion  $\psi/(\psi + \mu)$  for BDSIR—from the simulated trees, seasonal influenza A (H1N1), and human immunodeficiency virus (HIV-1) data analyses. LogN( $M, S$ ) is a log-normal distribution with mean  $M$  and standard deviation  $S$  in log space.

<sup>a</sup> Only applies to deterministic coalescent SIR; see details in File S1.

95% HPD coverage than deterministic coalescent SIR and the same as stochastic coalescent SIR. The deterministic coalescent SIR method recovered the truth for 85, 89, 91, and 88 of 100 trees for parameters  $R_0$ ,  $\gamma$ ,  $S_0$ , and  $z_0$ , while its stochastic analog recovered the truth for 100, 85, 100, and 99 of 100 trees for the same parameters. Finally, for stochastic coalescent SIR and BDSIR, error and (absolute) bias were relatively low for  $R_0$ , arguably the parameter of most interest to epidemiologists since it represents the number of individuals each infected individual will infect in a naive population. Deterministic coalescent SIR has a higher error (0.24) and bias (0.15) and also has significantly lower coverage for  $R_0$  (85%).

For  $R_0 \approx 1.10$ , the two stochastic models again outperformed the deterministic coalescent in error, bias, and 95% HPD coverage. The stochastic coalescent most reliably recovered the truth for  $R_0$  (99 of 100 simulations), while the deterministic coalescent had more than double the error and bias and still recovered the truth for only 25 of the 100 simulations. BDSIR had the lowest error and bias for  $R_0$  under this scheme, although it recovered the truth for only 75 of 100 simulations. For removal parameter  $\gamma$ , BDSIR again yielded lower error and bias, in this case returning the truth for 100/100 trees (in contrast to 84 and 86 from the stochastic and deterministic coalescent, respectively).

In the stochastic models, there is a greater trade-off between parameters due to the impact the relationship between them has on the survival of trajectories at low  $R_0$ . A larger estimated removal rate tends to require a larger susceptible population for the epidemic to avoid dying out in the early stages. Likewise, a smaller susceptible population implies a smaller estimated  $\gamma$ .

**Deterministic coalescent SIR on higher  $R_0$  and  $S_0$ :** As mentioned in the preceding subsection, the deterministic coalescent model yielded higher error and bias than both the stochastic coalescent and BDSIR for most parameters with  $R_0 \approx 1.10$  and  $S_0 = 499$ .

To investigate the deterministic model’s sensitivity to population sizes, we also simulated a range of population sizes ( $S_0 = 499, 999, \text{ and } 1999$ ) for  $R_0 = 1.0987$ . Even with  $S_0 = 1999$ , the deterministic coalescent SIR model’s 95% HPD coverage was low. For parameters  $R_0, \gamma, S_0, \text{ and } z_0$ , this

coverage was respectively 40%, 64%, 66%, and 18%. Table S6 shows these results.

Additionally, we increased both  $R_0$  (to 3.5 and 5) and  $S_0$  (to 4999 and 9999). However, for parameters  $R_0, \gamma, \text{ and } S_0$ , the deterministic coalescent SIR showed increased error, bias, and HPD widths, and the HPD coverage for  $z_0$  did not improve. These results are shown in Table S9.

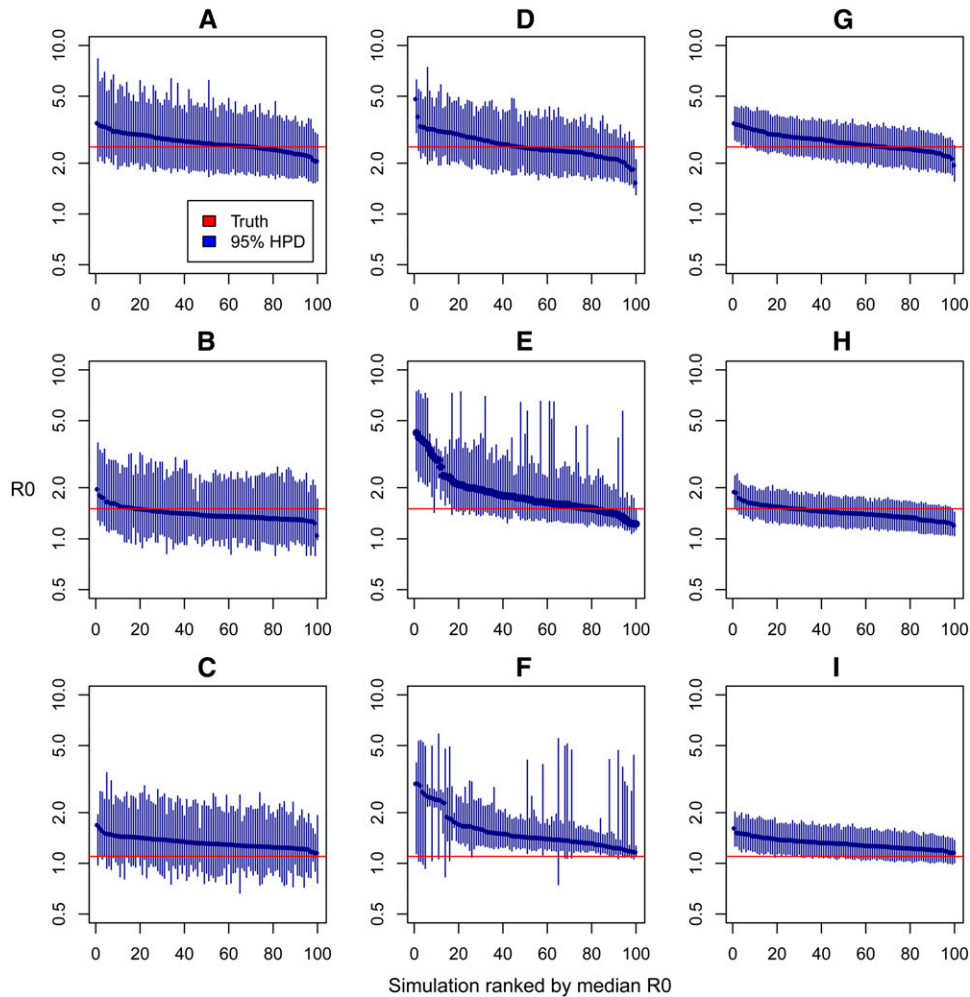
While each of these methods is an approximation, the deterministic coalescent particularly suffers from model misspecification since it does not account for the stochasticity that is always present in the early stages of epidemics, regardless of  $S_0$ .

**Homochronous trees:** Results for homochronously sampled trees are given in Table S3.

All three SIR inference models recover the truth for >95/100 trees within their respective 95% HPD widths for epidemic parameters  $R_0, \gamma, \text{ and } S_0$ . The time of origin  $z_0$  was recovered for 100/100 trees by BDSIR, 95/100 trees by stochastic coalescent SIR, and 73/100 trees by deterministic coalescent SIR. However, relative error and bias also increased consistently across all three models, along with the 95% HPD widths. The deterministic coalescent had the highest error, bias, and HPD width for  $R_0$  and highest error and HPD width for  $S_0$ , which is consistent with the heterochronously sampled data.

Further consideration of the effects of sampling rate changes and sampling model misspecification are warranted for BDSIR and coalescent SIR, the latter of which has been facilitated by Volz and Frost (2014).

**Simulated sequences:** Relative error and bias were inflated across all three inference models with the addition of phylogenetic uncertainty, and in certain cases the 95% HPD coverage was lower than with fixed trees. The deterministic coalescent model recovered the truth within its 95% HPD intervals only for  $\geq 90$  of the 100 trees in the case of  $S_0$ . The true values for the parameters  $R_0, \gamma, \text{ and } z_0$  were covered by 95% HPD intervals for 87, 56, and 29 of the 100 trees, respectively. This is contrasted with the performance of the stochastic coalescent (100, 97, 47, and 37 for parameters  $S_0, R_0, \gamma, \text{ and } z_0$ ) and BDSIR (99, 100, 84, and 18 for  $S_0, R_0, \gamma, \text{ and } z_0$ ), as shown in Table 1.



**Figure 3** Estimates of  $R_{(0)}$  from true stochastic SIR trees using inference methods by column, with stochastic coalescent SIR (A–C), deterministic coalescent SIR (D–F), and BDSIR (G–I). The truth varies by row, with  $R_0 = 2.4975$  (A, D, and G),  $R_0 = 1.4970$  (B, E, and H), and  $R_0 = 1.0978$  (C, F, and I).

Error, bias, and 95% HPD widths were higher with simulated sequences for all three inference models for parameters  $\gamma$ ,  $S_0$ , and  $z_0$  than with fixed trees. This indicates the importance of calibrating epidemic parameters of interest. In our case, we emphasize the basic reproductive number  $R_0$ , often the parameter of most interest to epidemiologists. For  $R_0$ , stochastic coalescent SIR and BDSIR recovered the truth within their 95% HPD intervals for 97 and 100 of the 100 simulations, respectively. They also showed only slight changes in error and bias compared to inference performed on the fixed trees used to generate the sequences. The deterministic coalescent SIR model recovered  $R_0$  for 87 of the 100 simulations (contrasted with 98/100 for the fixed trees) and with increased error.

**Priors and identifiability:** It is important to understand the impact of selected priors on inference results, as the priors are where the power of Bayesian inference lies. For example, we found relatively weak identifiability in the initial susceptible population parameter  $S_0$ , which must either be fixed or be estimated alongside the origin parameter  $z_0$ .

In addition to allowing each parameter to be either fixed or estimated, we have provided options for parameterization

of our models, with either the transmission rate  $\beta$  or  $R_0$  acting as operable parameters in MCMC analysis. For the deterministic coalescent, there is also an option to use the intrinsic growth parameter described by Dearlove and Wilson (2013).

The choice of parameterization necessarily affects the prior that will be used in the inference and should be considered carefully. However, we found that once a parameterization has been selected, our inference models are robust to different prior distributions placed on each parameter. We also used broader prior distributions on the deterministic coalescent to test whether this would increase its lower 95% HPD coverage relative to the stochastic models. We found that doing so increased the error and bias of the results without increasing the accuracy (shown in Table S4).

#### H1N1 data analysis

Epidemic parameter estimates from serially sampled influenza A (H1N1) virus sequence data are shown in Table 3.

The estimated means of the basic reproductive number were  $R_0 = 1.46$ , 1.35, and 1.61 for the stochastic coalescent, the deterministic coalescent, and BDSIR, respectively. Estimates of  $R_0$  from pandemic H1N1 in New Zealand range

from  $\sim 1.2$  to  $1.5$  (Paine *et al.* 2010; Opatowski *et al.* 2011; Roberts and Nishiura 2011; Roberts 2013; Biggerstaff *et al.* 2014), and estimates of  $R_0$  for seasonal H1N1 from other countries also range from  $\sim 1.2$  to  $1.5$  (Chowell *et al.* 2008). The 95% HPD intervals were very similar across each model, ranging from just over  $1.0$  to  $\sim 2.0$ .

The population of the Canterbury region in 2001 was reported to be  $\sim 481,431$  by the Environment Canterbury Regional Council (Ecan 2001) and  $521,832$  by Statistics New Zealand (StatsNZ 2001). The mean estimates of  $S_0$  were considerably lower using the stochastic coalescent ( $S_0 = 69,000$ ), the deterministic coalescent ( $S_0 = 120,000$ ), and BDSIR ( $S_0 = 22,200$ ). However, the *effective* population of susceptibles is assumed to be much smaller, as the total population contains individuals of various susceptibility, *e.g.*, those with partial immunity from vaccination and previous or secondary infections.

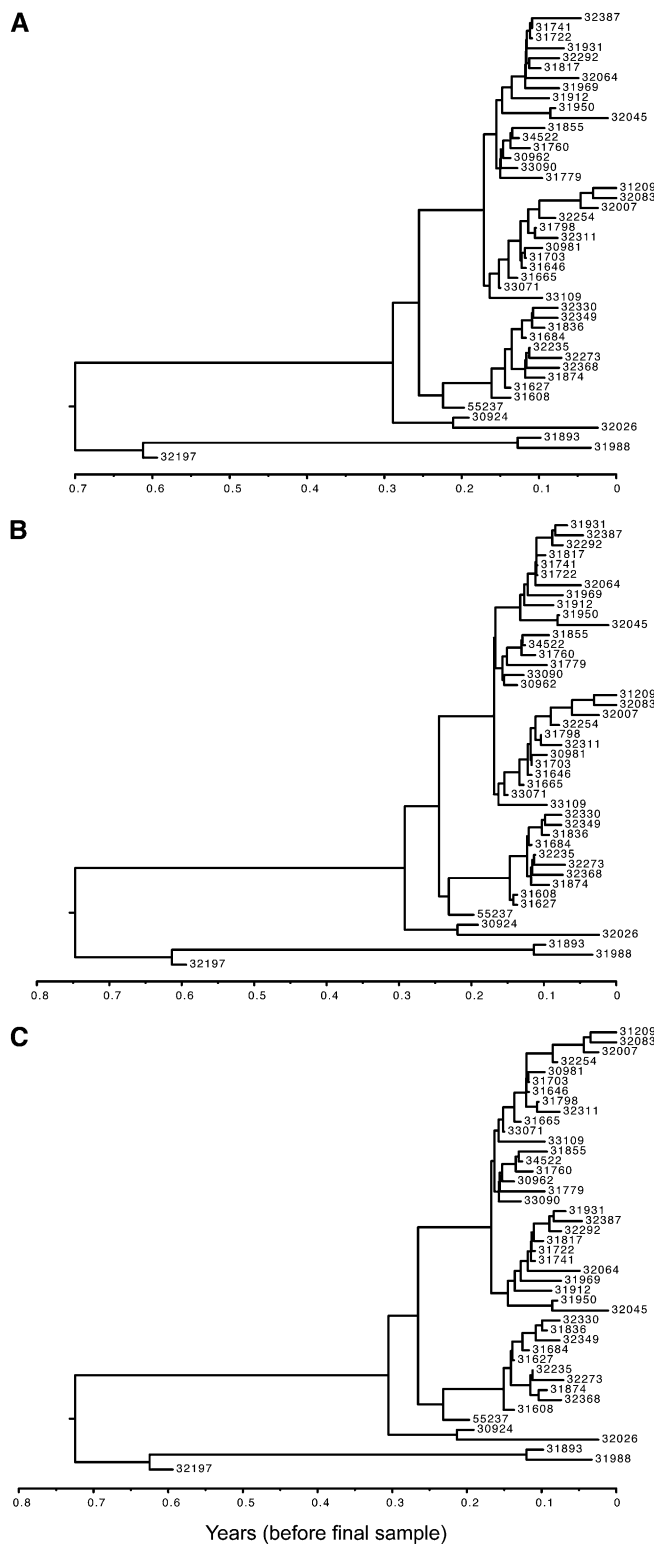
Most people recover from flu symptoms, the time they are likely to be most infectious, within a few days up to 2 weeks (CDC 2014; WHO 2014). This provides a range of probable true values for the removal parameter  $\gamma$ . The sequence data and molecular clock rate, and therefore the tree, are in units of years. Therefore, our  $\gamma$  range would be from  $365/14$  days to  $365/2$  days or from  $\gamma = 26.1$  to  $\gamma = 182.5$ . The stochastic coalescent, the deterministic coalescent, and BDSIR respectively inferred  $\gamma$  means of 27.08, 34.50, and 27.72. These estimates are on the low side compared to epidemiological models for influenza that include explicit spatial and household effects (Ferguson *et al.* 2005), but a moderate misfit of the model is not unexpected when fitting a simple closed SIR model with no population substructure.

The root of the tree was very similar across all inference models, respectively 0.53, 0.54, and 0.49 for stochastic coalescent SIR, deterministic coalescent SIR, and BDSIR. The same was true for the origin  $z_0$ , with: 0.69, 0.73, and 0.53 for the stochastic coalescent, the deterministic coalescent, and BDSIR. All three inference models returned tree root and origin estimates that are consistent with previous estimates from single flu seasons. That is, the tree age is young and the root coincides with the start of the (winter) influenza season in the Southern Hemisphere. The time of introduction of influenza into the region,  $z_0$ , was 1 or 2 months before the root. This supports the notion that the sequences selected represent a single introduction of the strain into the Canterbury population (see File S1 for details of data selection and Figure S4 for representative trees inferred from an alternate data selection.).

The trees estimated by each of the three models are typical for influenza (see Figure 4 for representative trees from each posterior), with branches that are quick to coalesce moving backward in time from the most recently sampled tip.

### HIV-1 data analysis

Results for inference from HIV-1 sequence data can be found in File S1. 95% HPD intervals are shown in Figure S5, Figure S6, Figure S7, and Table S8.



**Figure 4** Representative influenza A (H1N1) posterior trees from inference using the (A) BDSIR, (B) stochastic coalescent SIR, and (C) deterministic coalescent SIR inference models.

### Computational efficiency

Finally, Table S5 shows comparisons of computation times under each inference model for each type of data

analyzed. The deterministic coalescent SIR model is by far the fastest to sample and converge, with stochastic coalescent SIR and BDSIR varying, depending on the type of data.

### Closing remarks

A key reason for the success of coalescent theory in population genetics is its mathematical simplicity and the computational efficiency of calculating the probability density of a sample genealogy. Our results show that a stochastic variant of coalescent theory can be successfully adapted to estimate epidemiological parameters in a true Bayesian inference context. This stochastic coalescent SIR model performs better than the deterministic analog for estimating epidemic parameters in some circumstances. Unfortunately, the stochastic model relies on a computationally demanding Monte Carlo estimate of the coalescent density via simulation of an ensemble of epidemic trajectories, negating one of the main advantages of coalescent theory. In fact, the current implementation is less computationally efficient than the implementation of the BDSIR model. However, an advantage of the stochastic coalescent over the explicit sampling model in BDSIR is its robustness to biased sampling schemes, as has been shown for the case of pure exponential growth dynamics (Bošková *et al.* 2014).

A more computationally efficient approach to computing the coalescent probability of the sample genealogy in the stochastic setting would be to use particle filtering (Andrieu and Roberts 2009; Andrieu *et al.* 2010; Rasmussen *et al.* 2011, 2014), but there are no theoretical barriers to applying particle MCMC to the exact model (Stadler *et al.* 2014). Therefore, an obvious extension of this work would be to apply particle MCMC algorithms to the exact stochastic SIR model that was used in simulations in this work. We anticipate that the exact model would outperform all the methods tested here, especially when  $R_0$  is close to one.

In the meantime, the Bayesian coalescent inference methods developed here make it feasible to estimate epidemic parameters from time-stamped, serially sampled molecular sequence data, while accurately accounting for uncertainty in the topology and the divergence times of the phylogenetic tree.

### Acknowledgments

We thank Gabriel Leventhal and Louis Du Plessis (ETH Zürich) for constructive and valuable input and the New Zealand eScience Infrastructure for access to high-performance computing facilities (<http://www.nesi.org.nz/>). A.J.D. was funded by a Rutherford Discovery Fellowship from the Royal Society of New Zealand. A.P., T.V., T.S., and A.J.D. were also partially supported by Marsden grant UOA1324 from the Royal Society of New Zealand (<http://www.royalsociety.org.nz/programmes/funds/marsden/awards/2013-awards/>).

### Literature Cited

- Anderson, R. M., and R. M. May, 1991 *Infectious Diseases of Humans: Dynamics and Control*. Oxford University Press, Oxford.
- Andrieu, C., and G. O. Roberts, 2009 The pseudo-marginal approach for efficient Monte Carlo computations. *Ann. Stat.* 37: 697.
- Andrieu, C., A. Doucet, and R. Holenstein, 2010 Particle Markov chain Monte Carlo methods. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 72: 269–342.
- Beaumont, M. A., 2003 Estimation of population growth or decline in genetically monitored populations. *Genetics* 164: 1139–1160.
- Bedford, T., S. Cobey, P. Beerli, and M. Pascual, 2010 Global migration dynamics underlie evolution and persistence of human influenza A (H3N2). *PLoS Pathog.* 6: e1000918.
- Biggerstaff, M., S. Cauchemez, C. Reed, M. Gambhir, and L. Finelli, 2014 Estimates of the reproduction number for seasonal, pandemic, and zoonotic influenza: a systematic review of the literature. *BMC Infect. Dis.* 14: 480.
- Bošková, V., S. Bonhoeffer, and T. Stadler, 2014 Inference of epidemiological dynamics based on simulated phylogenies using birth-death and coalescent models. *PLoS Comput. Biol.* 10: e1003913.
- CDC, 2014 United States Centers for Disease Control and Prevention. Available at: <http://www.cdc.gov/flu/>. Accessed: November, 2014.
- Chowell, G., M. Miller, and C. Viboud, 2008 Seasonal influenza in the United States, France, and Australia: transmission and prospects for control. *Epidemiol. Infect.* 6: 852–864.
- Dearlove, B., and D. J. Wilson, 2013 Coalescent inference for infectious disease: meta-analysis of hepatitis C. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 368: 20120314.
- Drummond, A. J., G. K. Nicholls, A. G. Rodrigo, and W. Solomon, 2002 Estimating mutation parameters, population history and genealogy simultaneously from temporally spaced sequence data. *Genetics* 161: 1307–1320.
- Drummond, A. J., S. Y. W. Ho, M. J. Phillips, and A. Rambaut, 2006 Relaxed phylogenetics and dating with confidence. *PLoS Biol.* 4: e88.
- ECAN, 2001 Environment Canterbury Regional Council. Available at: <http://ecan.govt.nz/about-us/population/how-many/pages/census.aspx>. Accessed: November, 2014.
- Felsenstein, J., 1981 Evolutionary trees from DNA sequences: a maximum likelihood approach. *J. Mol. Evol.* 17: 368–376.
- Felsenstein, J., 2004 *Inferring Phylogenies*. Sinauer Associates, Sunderland, MA.
- Ferguson, N., D. Cummings, S. Cauchemez, C. Fraser, S. Riley *et al.*, 2005 Strategies for containing an emerging influenza pandemic in southeast Asia. *Nature* 437: 209–214.
- Gavryushkina, A., D. Welch, T. Stadler, and A. Drummond, 2014 Bayesian inference of sampled ancestor trees for epidemiology and fossil calibration. *arXiv:1406.4573*.
- Grenfell, B. T., O. G. Pybus, J. R. Gog, J. L. N. Wood, J. M. Daly *et al.*, 2004 Unifying the epidemiological and evolutionary dynamics of pathogens. *Science* 303: 327–332.
- Griffiths, R. C., and S. Tavaré, 1994 Ancestral inference in population genetics. *Stat. Sci.* 9: 307–319.
- Hasegawa, M., H. Kishino, and T. Yano, 1985 Dating of the human-ape splitting by a molecular clock of mitochondrial DNA. *J. Mol. Evol.* 22: 160–174.
- Iwasaki, A., and P. S. Pillai, 2014 Innate immunity to influenza virus infection. *Nat. Rev. Immunol.* 14: 315–328.
- Keeling, M. J., and P. Rohani, 2008 *Modeling Infectious Diseases in Humans and Animals*. Princeton University Press, Princeton.
- Kermack, W., and A. McKendrick, 1932 Contributions to the mathematical theory of epidemics. ii. The problem of endemicity. *Proc. R. Soc. A* 138: 55–83.

- Kingman, J. F. C., 1982 The coalescent. *Stoch. Proc. Appl.* 13: 235–248.
- Koelle, K., and D. A. Rasmussen, 2012 Rates of coalescence for common epidemiological models at equilibrium. *J. R. Soc. Interface* 9: 997–1007.
- Koelle, K., S. Cobey, B. Grenfell, and M. Pascual, 2006 Epochal evolution shapes the phylodynamics of interpandemic influenza a (h3n2) in humans. *Science* 314: 1898–1903.
- 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. *J. R. Soc. Interface* 11: 20131106.
- Kutta, M. W., 1901 Beitrag zur näherungsweise integration totaler differentialgleichungen. *Zeitschrift für Mathematik und Physik* 46: 435–453.
- Opatowski, L., C. Fraser, J. Griffin, E. de Silva, M. Van Kerkhove *et al.*, 2011 Transmission characteristics of the 2009 h1n1 influenza pandemic: comparison of 8 southern hemisphere countries. *PLoS Pathog.* 7: e1002225.
- Paine, S., G. Mercer, P. Kelly, D. Bandaranayake, M. Baker *et al.*, 2010 Transmissibility of 2009 pandemic influenza a(h1n1) in New Zealand: effective reproduction number and influence of age, ethnicity, and importations. *Eurosurveillance* 15(24). Available at: <http://www.eurosurveillance.org/ViewArticle.aspx?ArticleId=19591>.
- Pybus, O. G., M. A. Charleston, S. Gupta, A. Rambaut, E. C. Holmes *et al.*, 2001 The epidemic behavior of the hepatitis c virus. *Science* 292: 2323–2325.
- R Core Team, 2014 *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna.
- Rambaut, A., 2007 Figtree. Available at: <http://tree.bio.ed.ac.uk/software/figtree/>.
- Rambaut, A., and E. Holmes, 2009 The early molecular epidemiology of the swine-origin a/h1n1 human influenza pandemic. *PLoS Curr.* 1: RRR1003.
- Rasmussen, D. A., O. Ratmann, and K. Koelle, 2011 Inference for nonlinear epidemiological models using genealogies and time series. *PLoS Comput. Biol.* 7: e1002136.
- Rasmussen, D. A., E. M. Volz, and K. Koelle, 2014 Phylodynamic inference for structured epidemiological models. *PLoS Comput. Biol.* 10: e1003570.
- Roberts, M., 2013 Epidemic models with uncertainty in the reproduction number. *J. Math. Biol.* 66: 1463–1474.
- Roberts, M., and H. Nishiura, 2011 Early estimation of the reproduction number in the presence of imported cases: pandemic influenza h1n1–2009 in New Zealand. *PLoS ONE* 6: e17835.
- Rodrigo, A., and J. Felsenstein, 1999 Coalescent approaches to HIV population genetics, pp. 233–272 in *The evolution of HIV*, edited by K. A. Crandall. The Johns Hopkins University Press, Baltimore.
- Rouzine, I. M., A. Rodrigo, and J. M. Coffin, 2001 Transition between stochastic evolution and deterministic evolution in the presence of selection: general theory and application to virology. *Microbiol. Mol. Biol. Rev.* 65: 151–185.
- Runge, C., 1895 Ueber die numerische auflösung von differentialgleichungen. *Math. Ann.* 46: 167–178.
- Sehl, M., A. V. Alekseyenko, and K. L. Lange, 2009 Accurate stochastic simulation via the step anticipation tau-leaping (sal) algorithm. *J. Comput. Biol.* 16: 1195–1208.
- Stadler, T., R. Kouyos, V. von Wyl, S. Yerly, J. Böni *et al.*, 2012 Estimating the basic reproductive number from viral sequence data. *Mol. Biol. Evol.* 29: 347–357.
- 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). *Proc. Natl. Acad. Sci. USA* 110: 228–233.
- Stadler, T., T. G. Vaughan, A. Gavrushkin, S. Guindon, D. Kühnert, *et al.*, 2014 Population genetics vs. population dynamics: How well can coalescent-based models approximate population dynamic processes? *Genetics* 190: 187–201.
- StatsNZ, 2001 Statistics New Zealand. Available at: <http://stats.govt.nz/Census/>.
- Vaughan, T. G., and A. J. Drummond, 2013 A stochastic simulator of birth-death master equations with application to phylodynamics. *Mol. Biol. Evol.* 30: 1480–1493.
- Vaughan, T., D. Kühnert, A. Poppinga, D. Welch, and A. Drummond, 2014 Efficient Bayesian inference under the structured coalescent. *Bioinformatics* 30: 2272–2279.
- Volz, E., and S. D. Frost, 2014 Sampling through time and phylodynamic inference with coalescent and birth-death models. *J. R. Soc. Interface* 11: 20140945.
- Volz, E. M., 2012 Complex population dynamics and the coalescent under neutrality. *Genetics* 190: 187–201.
- 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.
- WHO, 2014 World Health Organization. Available at: <http://www.who.int/topics/influenza/en/>. Accessed: November, 2014.

Communicating editor: Y. S. Song

# GENETICS

Supporting Information

<http://www.genetics.org/lookup/suppl/doi:10.1534/genetics.114.172791/-/DC1>

## **Inferring Epidemiological Dynamics with Bayesian Coalescent Inference: The Merits of Deterministic and Stochastic Models**

Alex Poppinga, Tim Vaughan, Tanja Stadler, and Alexei J. Drummond

# Supporting material for “Inferring epidemiological dynamics with Bayesian coalescent inference: The merits of deterministic and stochastic models”

## File S1

### 1 Sampling from the prior

In order to assess the correctness of our implementation of the deterministic coalescent SIR and stochastic coalescent SIR models, for each model we used the MCMC algorithm to sample trees from the corresponding distribution  $f(\mathcal{T}|\eta)$ , and compared these samples with coalescent trees simulated directly under the model.

The chosen  $\eta$  included  $\beta = 7.5 \times 10^{-4}$ ,  $\gamma = 0.3$ ,  $S_0 = 999$  and  $z_0 = 30$ . The comparisons were performed for trees generated from 20 leaves, sampled at integer times 0 through 19, inclusive.

For the deterministic coalescent SIR model, the direct simulation involved numerically solving the Eqs. (1)–(3) in the main text for  $t \in [0, 30]$  and using this solution in combination with Eq. (10) in the main text to determine the instantaneous coalescent rate  $\lambda(\tau)$ . This rate was used to simulate each of the coalescent trees in the usual fashion for heterochronous leaf times. In the case that the MRCA was not reached before the origin time of the epidemic, the tree was discarded and the simulation repeated.

The direct simulation proceeded in a similar way for the stochastic coalescent SIR model, the major difference being that the stochasticity of this model required each coalescent tree to be simulated under a distinct realization of the stochastic trajectory.

Comparisons between the direct simulation and MCMC results are shown in Figures S1 and S2 for three different summary statistics and show very close agreement.

### 2 Validation through simulated data analysis

As part of the validation of our implementation of the two coalescent SIR models, trees were simulated by their own methods (using stochastically- and deterministically-generated SIR trajectories, as discussed in the Methods section of the main paper), and relevant epidemiological parameters were inferred

using the stochastic and deterministic coalescent SIR models. Tables 1 and 2 show the results of these analyses, indicative of correct implementations.

Analyses for varying  $R_0$  (and necessarily, slightly varied other parameters, such as the birth rate  $\beta$ ) are provided in Tables S3 and S4. Results from tests of the influence of broader priors (with larger standard deviations in log space) are shown in Table S4. It appears that allowance of broader priors reduces 95% HPD coverage in some cases (e.g., for parameter  $R_0$ ) when using the deterministic coalescent SIR inference model, as they increase error and bias.

Finally, it was noticed that even for the higher true parameter values of  $R_0 = 2.50$  and  $S_0 = 999$ , under which deterministic coalescent SIR is expected to perform relatively well, there was an inability to accurately estimate the origin parameter  $z_0$ . Figure S3 provides some insight into this conundrum by examining the trajectories used for tree simulation and subsequent analysis.

## 2.1 H1N1 data selection

Initially, the H1N1 dataset contained 45 sequences. The ages of the inferred trees (Figure S4) using the original 45 sequences extended more than 1.5 years into the past for each of the SIR models, which is contrary to what we expect for a single, current strain of seasonal influenza. Three taxa (labelled 32197, 31893, and 31988) were hypothesized to belong to a unique strain, e.g., an additional seeding from outside the Canterbury region or a low-lying previous strain. Removing these three taxa caused the inferred trees to behave as expected, i.e., tree heights and epidemic origin  $z_0$  less than a year old. It also raised the estimated  $R_0$  values for all three SIR models (initially 1.24, 1.10, and 1.55 for stochastic coalescent SIR, deterministic coalescent SIR, and BDSIR, respectively), as well as those for  $\gamma$  (initially 8.74, 12.65, and 11.33 for stochastic coalescent SIR, deterministic coalescent SIR, and BDSIR, respectively).

It will be interesting to further investigate the interplay between influenza strains and its contribution to the overall dynamics. For the closed SIR models discussed in this manuscript, however, this additional complexity leads to increased chance of model misspecification and misleading results. Therefore, we focused our attention on the analyses using 42 sequences.

## 2.2 HIV-1 data analysis

The original HIV-1 dataset (HUÉ *et al.* 2005) was agglomerated from both acute and chronic infections sampled in the United Kingdom (UK) and constitutes six phylogenetic clusters, from which the five used here (Clusters 1-4 and 6) were drawn. These particular clusters, with the omission of Cluster 5, were chosen simply for the purpose of direct comparison with KÜHNERT *et al.* (2014). Our extension to the models allowed us to imprint respective tip dates on the sequence data, sampled from 1999 to 2003, for inclusion in the likelihood computation.

For the selected five clusters, the nucleotide alignments contained 41, 62, 29, 26, and 35 sequences, respectively, each with 952 sites. The substitution



scheme chosen for phylogenetic analysis was the symmetric and independent general time reversible model (GTR), with gamma distributed rate variation and explicit proportion of invariable sites (GTR+G+I). Following HUÉ *et al.* (2005), the substitution rate was set to 2.55E-4 substitutions per site per year. All other parameters were estimated conjointly, and the Bayesian prior distributions are presented in Table 4: Bayesian prior distributions.

The pathophysiology of HIV is multifarious, and the patterns of its advancement within an infected host change throughout time. In addition to increased complexity potentially caused by recombination events, the transition between HIV’s acute and chronic phases alters the host’s infectivity (GUSS 1994). The SIR compartmental model used for this particular phylodynamic analysis on the UK cluster data does not allow for independent infection rates for the acute and chronic phases (but see VOLZ *et al.* (2012) and VOLZ *et al.* (2013)). However, in this study we did not attempt to estimate the infection rate  $\beta$  and thus did not expect such a difference to significantly impact the estimation of the parameters of interest: the basic reproductive number  $R_0$ , removal rate  $\gamma$ , size of the initial susceptible population  $S_0$ , and origin of the outbreak  $z_0$ .

### 2.2.1 HIV-1 inference results

In regard to parameter inference from the serially-sampled HIV-1 sequence data, the stochastic coalescent SIR, deterministic coalescent SIR, and BDSIR methods were most alike in light of the  $R_0$  results. The medians and HPD intervals for all clusters pertaining to this parameter, (especially Clusters 1, 2, 3, and 6), were very close, and those of Cluster 4 were still congruent across the three analyses (Figure S5).

The coalescent SIR models and BDSIR disagreed with respect to the age of the most recent common ancestor and the origin  $z_0$  (Figure S6). The coalescent SIR models also exhibited much larger 95% HPD intervals for  $z_0$  in each of the clusters; while BDSIR encompassed an average of 16 years, the stochastic coalescent SIR and deterministic coalescent SIR models had averages of 49 and 37 years, respectively. Furthermore, the estimated age of the common ancestor of the tree was older under the coalescent SIR models than the estimates reported by either BDSIR or the original data analysis (HUÉ *et al.* 2005) for each cluster. This was also true for the time of origin for the epidemic, although for certain clusters the differences between the coalescent estimates of the origin  $z_0$  and the birth-death estimates were much greater than others (e.g., Cluster 3).

The estimates of removal rate  $\gamma$  from Clusters 1 and 6 were very similar across the three methods (Figure S7). However, both coalescent SIR models estimated considerably higher  $\gamma$  values for Clusters 2-4 than BDSIR. This is reflective of the simulation study results, where the two coalescent models did not perform as well as BDSIR for the removal parameter.

Median estimates for the initial susceptible population  $S_0$  were quite similar in all methods for Clusters 1-4, although BDSIR displayed much wider HPD intervals than stochastic coalescent SIR and deterministic coalescent SIR (Figure S8). In Cluster 6, the coalescent SIR models showed the smallest HPD intervals

for their individual analyses on each cluster, while the opposite was true for BDSIR. There was also a disparity between the median estimates for the two coalescent approaches and that of BDSIR for Cluster 6. To this effect, it should be noted that the number of infections accrued throughout the duration of the epidemic was reported as  $N_e = 1,350$  by Hué *et al.* This casts some suspicion on the low susceptible population estimates obtained by the stochastic coalescent SIR and deterministic coalescent SIR methods (median estimates of  $S_0 = 727$  and  $S_0 = 693$ , respectively), since they appear lower than the estimated number of infected individuals from the original study.

There is disagreement in the literature in regard to the modelling of HIV-1 evolutionary dynamics under stochastic or deterministic processes (NIJHUIS *et al.* 1998; ROUZINE and COFFIN 1999; ACHAZ *et al.* 2004; SHRINER *et al.* 2004). The predicament dwells in the observation that the actual effective population size  $N_e$  for HIV-1 is often smaller than the total population size (KOUYOS *et al.* 2006). While most of this debate has focused on within-host population dynamics, many of the arguments hold when considering the broader epidemic dynamics of host-to-host transmission. As previously mentioned, the appropriateness of these descriptions is hinged on the magnitude of the infected population, precisely, the effective infected population size. Consequently, even when the total infected population is quite large there may yet be significant stochastic effects in play.

Finally, as mentioned in the main article, the existence of two distinct infectious stages and the possibility of large effects due to recombination are reasons for any discrepancy produced by these SIR inference models.

## 2.2.2 Example XML

Below is an example XML for simulating 100 trees and trajectories in MASTER (VAUGHAN and DRUMMOND 2013). This example is for  $R_0 = 2.4975$  and  $S_0 = 999$ . The simulation ends when the infected  $I$  population returns to zero, i.e., when the last infected individual is removed.

```
<beast version='2.0'
namespace='master.beast:beast.core.parameter:beast.evolution.tree.TreeHeightLogger'>

  <run spec='InheritanceEnsemble'
    nTraj='100'
    samplePopulationSizes='true'
    verbosity='1'>

    <model spec='InheritanceModel' id='model'>
      <population spec='Population' id='S' populationName='S'/>
      <population spec='Population' id='I' populationName='I'/>
      <population spec='Population' id='R' populationName='R'/>
      <population spec='Population' id='Rh' populationName='Rh'/>

      <!-- infection reaction -->
      <reaction spec='InheritanceReaction' reactionName='Infection' rate='0.00075'>
        S + I -> 2I
      </reaction>
    </model>
  </run>
</beast>
```

```

</reaction>

<!-- recovery reaction -->
<reaction spec='InheritanceReaction' reactionName='Recovery' rate='0.25'>
  I -> R
</reaction>

<!-- sampling reaction -->
<reaction spec='InheritanceReaction' reactionName='Sampling' rate='0.05'>
  I -> Rh
</reaction>
</model>

<initialState spec='InitState'>
  <populationSize spec='PopulationSize' population='@S' size='999' />
  <lineageSeed spec='Individual' population='@I' />
</initialState>

<populationEndCondition spec='PopulationEndCondition'
  population='@I'
  threshold='0'
  exceedCondition='false' />

<inheritancePostProcessor spec='LineageFilter'
  reactionName='Sampling'
  discard='false' />

<output spec='NewickOutput' fileName='SIR.newick' />
<output spec='NexusOutput' fileName='SIR.nexus' />
<output spec='JsonOutput' fileName='SIR.json' />

</run>
</beast>

```

[Figure 1 about here.]

[Figure 2 about here.]

[Figure 3 about here.]

[Figure 4 about here.]

[Figure 5 about here.]

[Figure 6 about here.]

[Figure 7 about here.]

[Figure 8 about here.]

[Table 1 about here.]

[Table 2 about here.]

[Table 3 about here.]

[Table 4 about here.]

[Table 5 about here.]

[Table 6 about here.]

## References

- ACHAZ, G., S. PALMER, M. KEARNEY, F. MALDARELLI, J. W. MELLORS, *et al.*, 2004 A robust measure of HIV-1 population turnover within chronically infected individuals. *Mol Biol Evol* **21**: 1902–12.
- GUSS, D. A., 1994 The acquired immune deficiency syndrome: An overview for the emergency physician, part 1. *The Journal of Emergency Medicine* **12**: 375–384.
- HUÉ, S., D. PILLAY, J. P. CLEWLEY, and O. G. PYBUS, 2005 Genetic analysis reveals the complex structure of hiv-1 transmission within defined risk groups. *PNAS* **102**: 4425–4429.
- KOUYOS, R. D., C. L. ALTHAUS, and S. BONHOEFFER, 2006 Stochastic or deterministic: what is the effective population size of HIV-1? *Trends Microbiol* **14**: 507–11.
- 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. *J R Soc Interface* **11**: 20131106.
- NIJHUIS, M., C. A. BOUCHER, P. SCHIPPER, T. LEITNER, R. SCHURMAN, *et al.*, 1998 Stochastic processes strongly influence hiv-1 evolution during suboptimal protease-inhibitor therapy. *Proc Natl Acad Sci U S A* **95**: 14441–6.
- ROUZINE, I., and J. COFFIN, 1999 Linkage disequilibrium test implies a large effective population number for hiv in vivo. *PNAS* **96**: 10758–10763.
- SHRINER, D., R. SHANKARAPPA, M. A. JENSEN, D. C. NICKLE, J. E. MITTLER, *et al.*, 2004 Influence of random genetic drift on human immunodeficiency virus type 1 env evolution during chronic infection. *Genetics* **166**: 1155–64.
- VAUGHAN, T. G., and A. J. DRUMMOND, 2013 A stochastic simulator of birth-death master equations with application to phylodynamics. *Molecular Biology and Evolution* .
- VOLZ, E. M., 2012 Complex population dynamics and the coalescent under neutrality. *Genetics* **190**: 187–201.
- VOLZ, E. M., E. IONIDES, E. O. ROMERO-SEVERSON, M.-G. BRANDT, E. MOKOTOFF, *et al.*, 2013 Hiv-1 transmission during early infection in men who have sex with men: a phylogenetic analysis. *PLoS Med* **10**: e1001568; discussion e1001568.
- VOLZ, E. M., J. S. KOOPMAN, M. J. WARD, A. L. BROWN, and S. D. W. FROST, 2012 Simple epidemiological dynamics explain phylogenetic clustering of hiv from patients with recent infection. *PLoS Comput Biol* **8**: e1002552.

## List of Figures

S1	Comparison between distributions of summary statistics of trees sampled using MCMC employing our implementation of the <i>deterministic coalescent SIR model</i> likelihood and those calculated, and those of trees sampled using direct simulation. Summary statistics shown are (a) the age of the MRCA of the transmission tree, (b) the sum of all edge lengths in the tree, and (c) the total number of two-leaf clades in the tree. . . . .	8
S2	Comparison between distributions of summary statistics of trees sampled using MCMC employing our implementation of the <i>stochastic coalescent SIR model</i> likelihood and those calculated, and those of trees sampled using direct simulation. Summary statistics shown are (a) the age of the MRCA of the transmission tree, (b) the sum of all edge lengths in the tree, and (c) the total number of two-leaf clades in the tree. . . . .	9
S3	(a) True stochastic SIR trajectories simulated jointly alongside phylogenies, with the corresponding trajectories used by deterministic coalescent SIR. Adjusting deterministic coalescent SIR to fit the underlying stochastic trajectories causes major shifts to the origin $z_0$ . (b) Deterministic residuals with $z_0$ either fitted or not. . . . .	10
S4	Representative influenza A (H1N1) posterior trees from inference using the stochastic coalescent SIR (left), deterministic coalescent SIR (right), BDSIR (bottom) models. . . . .	11
S5	95% HPD intervals of $R_0$ for the HIV-1 subtype B UK cluster analyses using coalescent and birth-death methods. . . . .	12
S6	95% HPD intervals for coalescent [VOLZ (2012)] and birth-death [KÜHNERT <i>et al.</i> (2014)] estimations of the time into the past at which the root of the HIV-1 tree and introduction of the first infection occurred. . . . .	13
S7	95% HPD intervals of $\gamma$ for the HIV-1 subtype B UK cluster analyses using coalescent and birth-death methods. . . . .	14
S8	95% HPD intervals of $S_0$ for the HIV-1 subtype B UK cluster analyses using coalescent and birth-death methods. . . . .	15

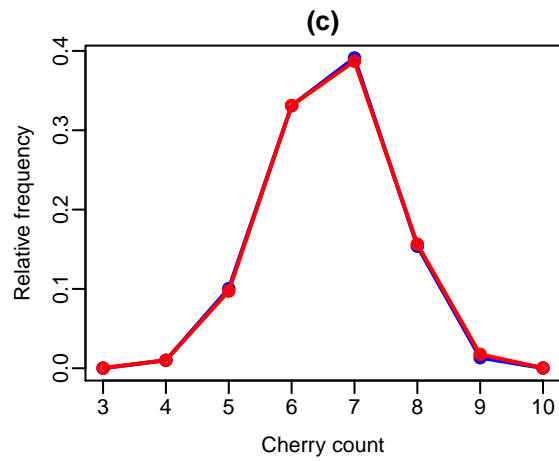
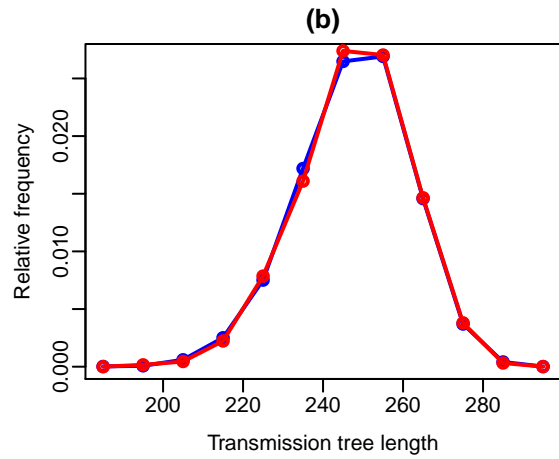
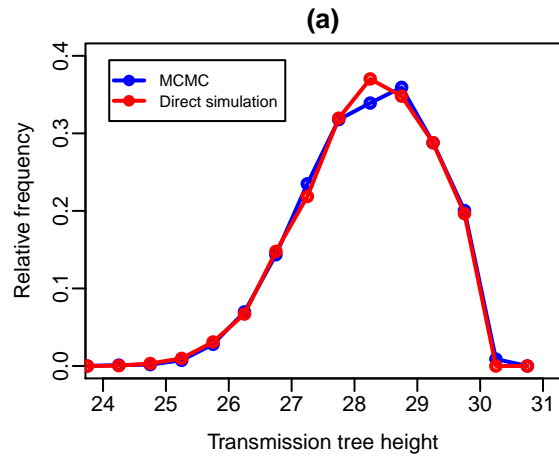


Figure S1: Comparison between distributions of summary statistics of trees sampled using MCMC employing our implementation of the *deterministic coalescent SIR model* likelihood and those calculated, and those of trees sampled using direct simulation. Summary statistics shown are (a) the age of the MRCA of the transmission tree, (b) the sum of all edge lengths in the tree, and (c) the total number of two-leaf clades in the tree.

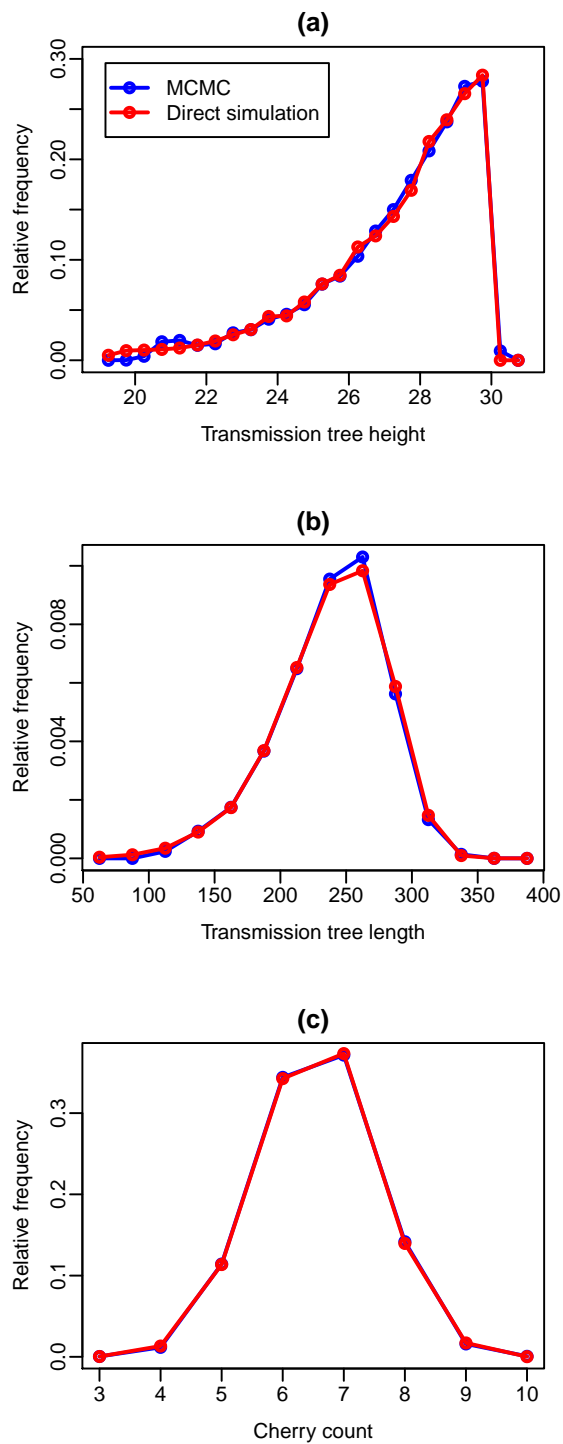


Figure S2: Comparison between distributions of summary statistics of trees sampled using MCMC employing our implementation of the *stochastic coalescent SIR model* likelihood and those calculated, and those of trees sampled using direct simulation. Summary statistics shown are (a) the age of the MRCA of the transmission tree, (b) the sum of all edge lengths in the tree, and (c) the total number of two-leaf clades in the tree.

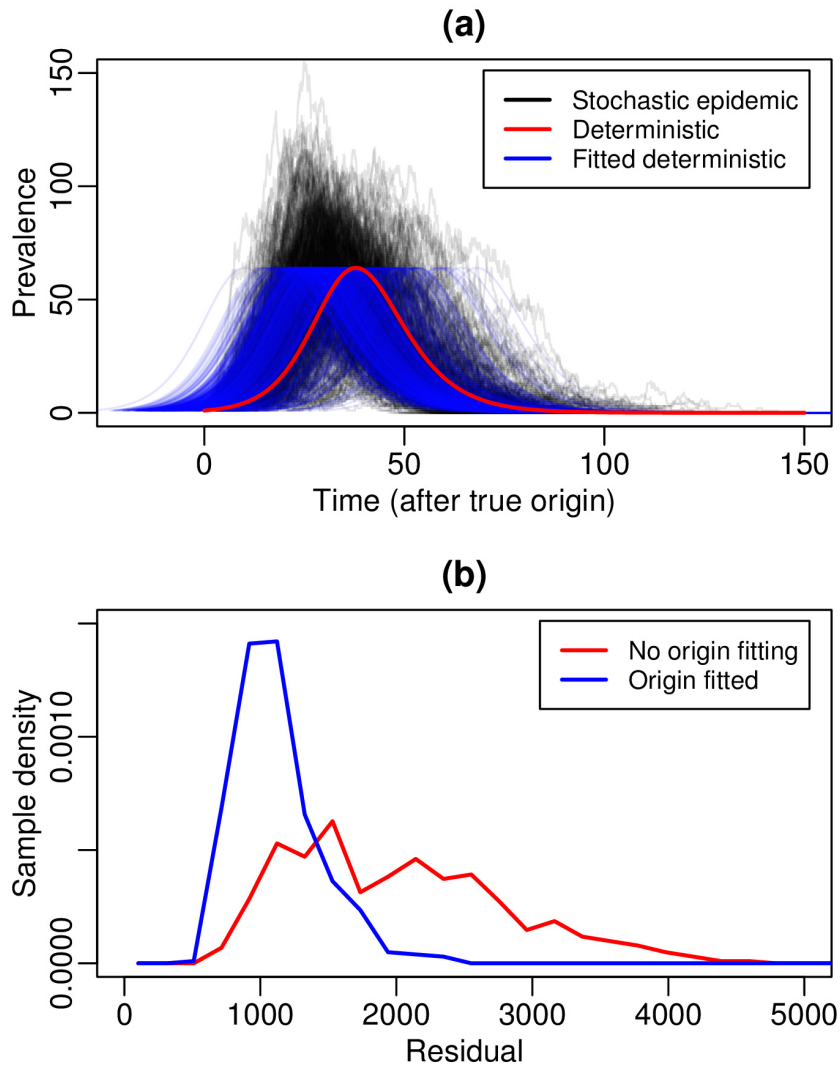


Figure S3: (a) True stochastic SIR trajectories simulated jointly alongside phylogenies, with the corresponding trajectories used by deterministic coalescent SIR. Adjusting deterministic coalescent SIR to fit the underlying stochastic trajectories causes major shifts to the origin  $z_0$ . (b) Deterministic residuals with  $z_0$  either fitted or not.



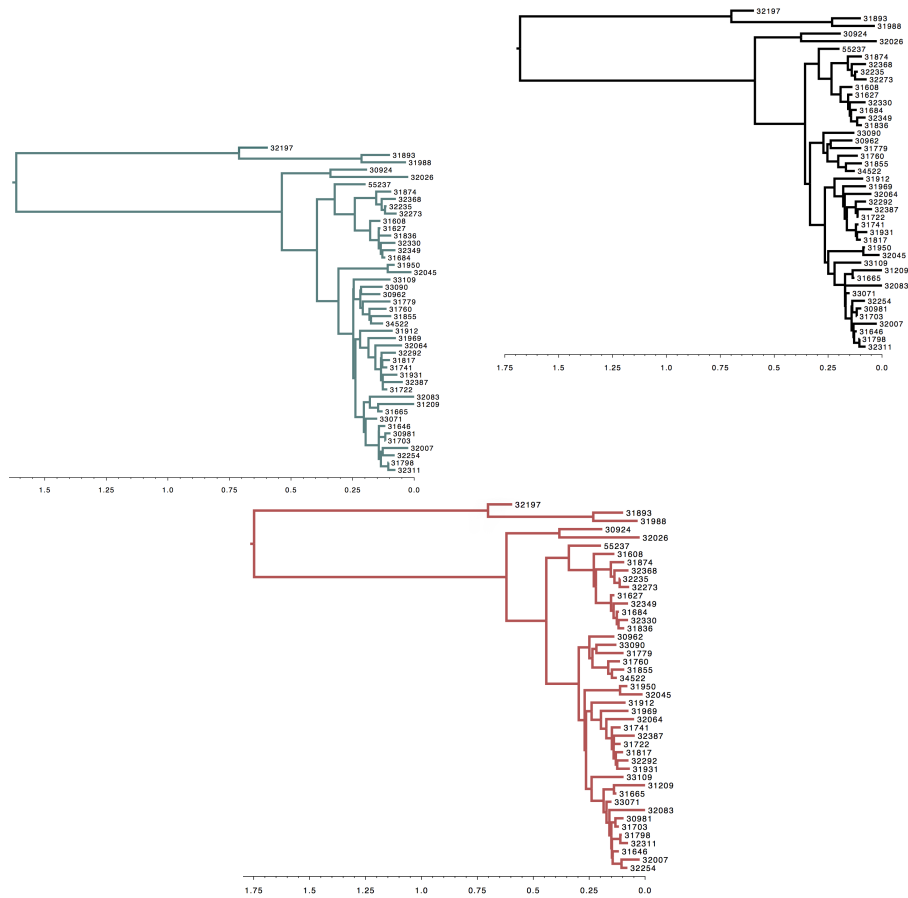


Figure S4: Representative influenza A (H1N1) posterior trees from inference using the stochastic coalescent SIR (left), deterministic coalescent SIR (right), BDSIR (bottom) models.

### HIV-1 subtype B: Estimates of $R_0$

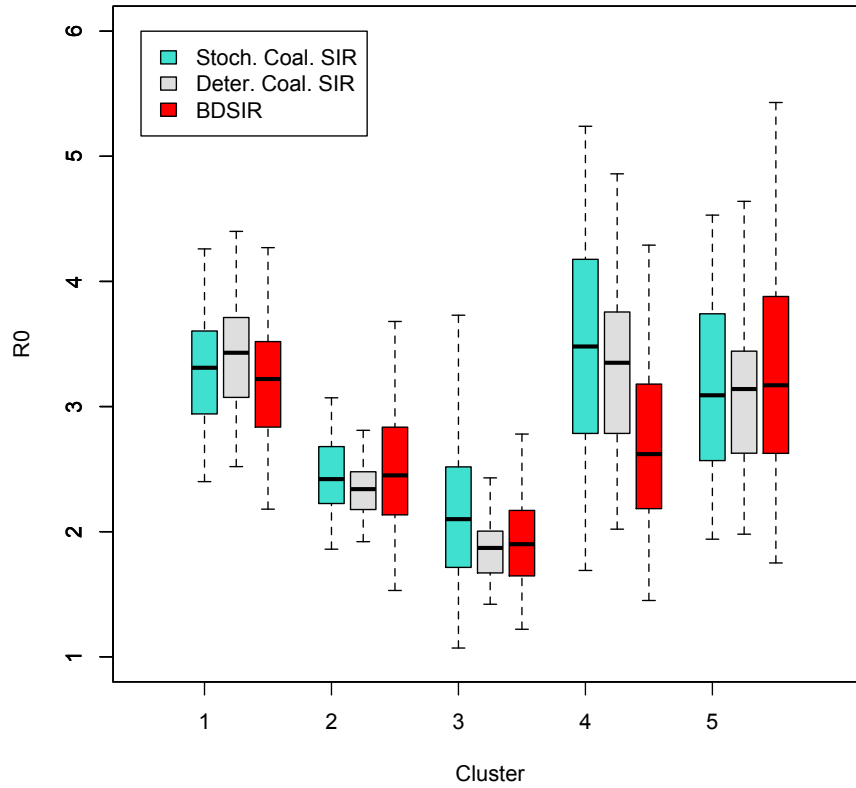


Figure S5: 95% HPD intervals of  $R_0$  for the HIV-1 subtype B UK cluster analyses using coalescent and birth-death methods.

### HIV-1 subtype B: Estimates for Tree Height and the Origin of the Epidemic

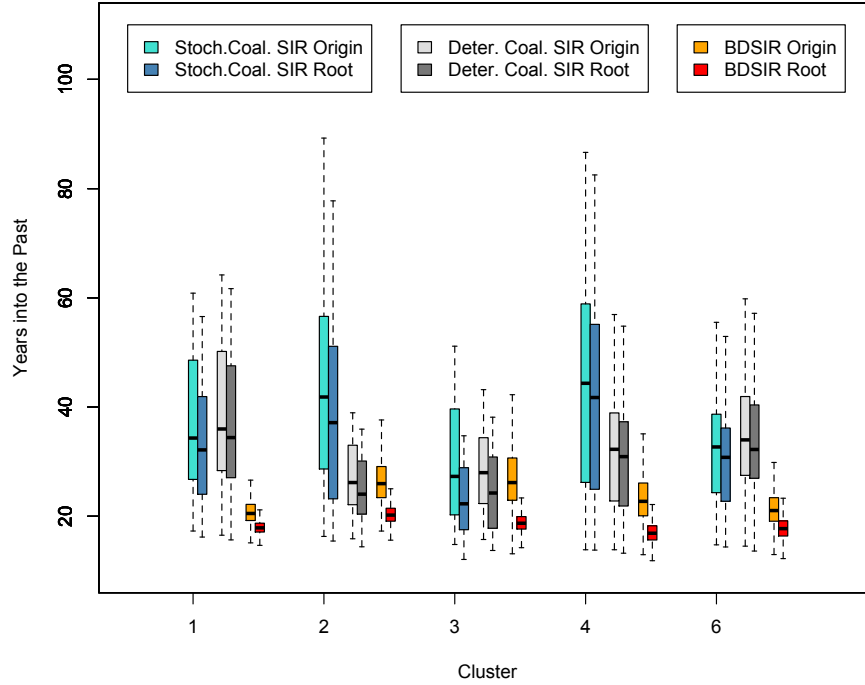


Figure S6: 95% HPD intervals for coalescent [VOLZ (2012)] and birth-death [KÜHNERT *et al.* (2014)] estimations of the time into the past at which the root of the HIV-1 tree and introduction of the first infection occurred.

### HIV-1 subtype B: Estimates of Gamma

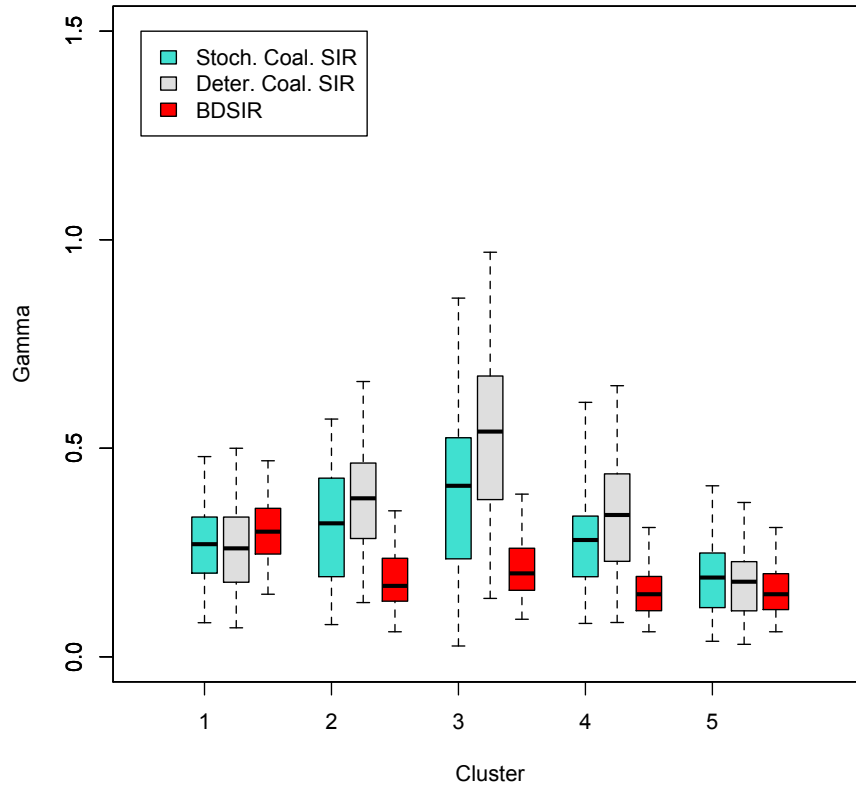


Figure S7: 95% HPD intervals of  $\gamma$  for the HIV-1 subtype B UK cluster analyses using coalescent and birth-death methods.

### HIV-1 subtype B: Estimates of $S_0$

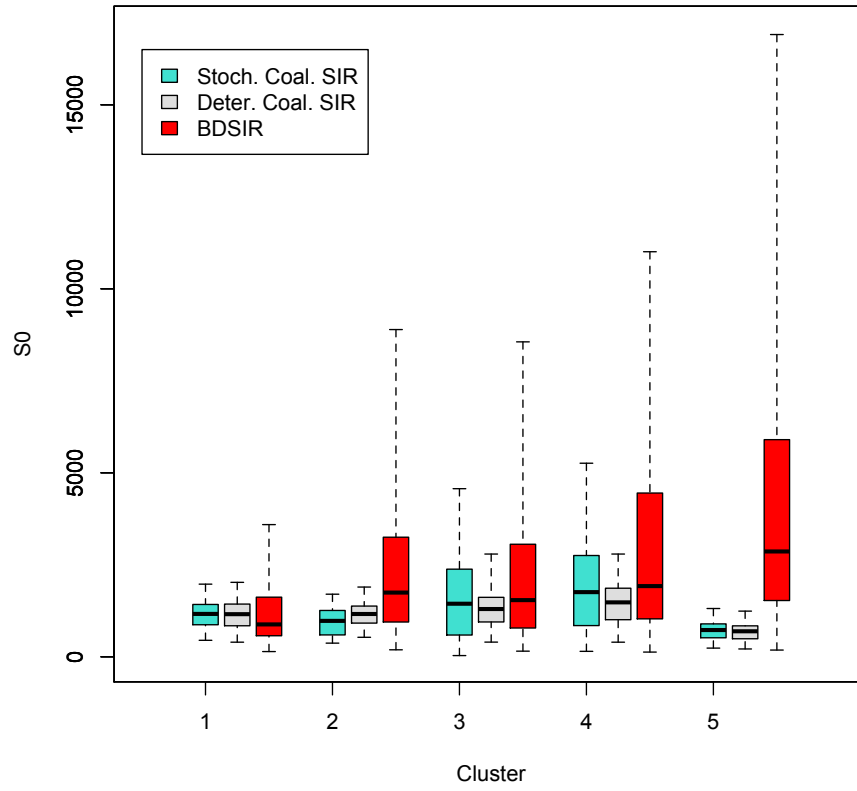


Figure S8: 95% HPD intervals of  $S_0$  for the HIV-1 subtype B UK cluster analyses using coalescent and birth-death methods.

## List of Tables

S1	Simulation Study Results for Stochastic Coalescent Trees	17
S2	Simulation Study Results for Deterministic Coalescent Trees . . . . .	18
S3	Results for Homochronous Sampling . . . . .	19
S4	Simulation Study Results: The Effect of Broader Priors on Deterministic Coalescent SIR . . . . .	20
S5	Comparison of Computation Times for Bayesian Inference of Epidemic Parameters from Genetic Sequence Data using SIR Models . . . . .	21
S6	Deterministic Coalescent SIR Results for Simulated Sequences: $R_0 = 1.0987$ and $S_0 = 499$ , $R_0 = 1.0989$ and $S_0 = 999$ , $R_0 = 1.09945$ and $S_0 = 1999$ . . . . .	22
S7	Simulation Study Details . . . . .	23
S8	Epidemic Parameter Estimations from HIV-1 Subtype B Sequence Data . . . . .	24
S9	Deterministic Coalescent SIR Results from Trees Simulated with Higher $S_0$ (with Fixed $R_0$ ) and Higher $R_0$ (with Fixed $S_0$ ) . . . . .	25

Table S1: Simulation Study Results for Stochastic Coalescent Trees

$\eta$	Inference	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy
$\mathcal{R}_0$	Stoch.Coal.SIR	2.50	2.81	2.64	0.11	0.08	0.95	100.00%
	Deter.Coal.SIR	2.50	2.73	2.65	0.14	0.06	0.85	96.00%
$\gamma$	Stoch.Coal.SIR	0.30	0.28	0.26	0.16	-0.11	1.17	99.00%
	Deter.Coal.SIR	0.30	0.30	0.28	0.18	-0.03	1.20	99.00%
$S_{(0)}$	Stoch.Coal.SIR	999	1456	986	0.21	0.02	3.93	100.00%
	Deter.Coal.SIR	999	1720	1057	0.48	0.24	4.28	99.00%
$z_{(0)}$	Stoch.Coal.SIR	(varies)	42.36	40.43	0.03	0.02	0.20	98.00%
	Deter.Coal.SIR	(varies)	41.25	39.77	0.03	0.01	0.07	64.00%

Table S2: Simulation Study Results for Deterministic Coalescent Trees

$\eta$	Inference	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy
$\mathcal{R}_0$	Stoch.Coal.SIR	2.50	2.44	2.37	0.06	-0.05	0.67	100.00%
	Deter.Coal.SIR	2.50	2.51	2.46	0.08	-0.01	0.59	99.00%
$\gamma$	Stoch.Coal.SIR	0.30	0.33	0.31	0.07	0.05	1.00	100.00%
	Deter.Coal.SIR	0.30	0.32	0.30	0.10	0.02	0.79	100.00%
$\tilde{S}_{(0)}$	Stoch.Coal.SIR	999	1586	1142	0.26	0.20	3.83	100.00%
	Deter.Coal.SIR	999	1426	1030	0.36	0.13	3.03	100.00%
$z_{(0)}$	Stoch.Coal.SIR	44.12	45.52	44.74	0.02	0.01	0.19	93.00%
	Deter.Coal.SIR	44.12	44.34	44.11	0.02	1.93E-3	0.08	92.00%



Table S3: Results for Homochronous Sampling

$\eta$	Inference	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy
$\mathcal{R}_0$	Stoch.Coal.SIR	2.50	3.04	2.74	0.13	0.11	1.32	100.00%
	Deter.Coal.SIR	2.50	4.05	3.29	0.34	0.32	2.38	100.00%
	BDSIR	2.50	2.84	2.49	0.16	0.03	1.45	97.00%
$\gamma$	Stoch.Coal.SIR	0.30	0.26	0.23	0.25	-0.21	1.43	100.00%
	Deter.Coal.SIR	0.30	0.26	0.19	0.36	-0.31	2.03	100.00%
	BDSIR	0.30	0.23	0.17	0.42	-0.42	2.04	100.00%
$S_{(0)}$	Stoch.Coal.SIR	999	1660	1065	0.18	0.09	4.75	100.00%
	Deter.Coal.SIR	999	4127	679	0.78	0.09	10.24	100.00%
	BDSIR	999	1907	1320	0.41	0.41	4.86	100.00%
$z_{(0)}$	Stoch.Coal.SIR	20.0	20.17	19.82	0.09	-0.03	0.43	95.00%
	Deter.Coal.SIR	20.0	19.09	19.21	0.09	-0.05	0.19	73.00%
	BDSIR	20.0	36.56	29.38	0.55	0.54	4.24	100.00%

Table S4: Simulation Study Results: The Effect of Broader Priors on Deterministic Coalescent SIR

$\eta$	St. Dev.	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy
$\mathcal{R}_0$	2	1.50	2.06	1.75	0.40	0.35	0.86	79.00%
$\mathcal{R}_0$	1	1.50	1.80	1.49	0.24	0.15	0.52	85.00%
$\mathcal{R}_0$	2	2.50	3.31	2.85	0.34	0.24	1.43	95.00%
$\mathcal{R}_0$	1	2.50	2.68	2.49	0.13	0.04	0.80	99.00%
$\gamma$	2	0.30	0.31	0.23	0.37	-0.12	1.59	96.00%
$\gamma$	1	0.30	0.26	0.23	0.27	-0.22	1.15	89.00%
$\gamma$	2	0.30	0.31	0.25	0.33	-0.09	1.59	95.00%
$\gamma$	1	0.30	0.32	0.29	0.16	3.14E-3	1.27	99.00%
$S_{(0)}$	2	499	2041	249	1.40	0.49	7.75	85.00%
$S_{(0)}$	1	499	562	361	0.44	-0.26	3.36	91.00%
$S_{(0)}$	2	999	3028	717	1.05	0.33	6.60	94.00%
$S_{(0)}$	1	499	553.38	337	0.42	-0.26	3.08	92.00%
$z_{(0)}$	2	(varies)	65.10	62.01	0.04	0.03	0.25	86.00%
$z_{(0)}$	1	(varies)	91.03	72.51	0.39	0.38	0.42	88.00%
$z_{(0)}$	2	(varies)	40.97	39.85	0.03	-6.78E-4	0.08	81.00%
$z_{(0)}$	1	(varies)	112.79	90.37	0.26	0.26	0.94	85.00%

Table S5: Comparison of Computation Times for Bayesian Inference of Epidemic Parameters from Genetic Sequence Data using SIR Models

Data Type	Inference Model	Mean time per million samples (MCMC)
Sim. Study ( $R_0 \approx 2.50$ )	Stoch.Coal.SIR	20m 41s
	Deter.Coal.SIR	3m 27s
	BDSIR	56m 27s
Sim. Study ( $R_0 \approx 1.50$ )	Stoch.Coal.SIR	1h 43m 30s
	Deter.Coal.SIR	3m 47s
	BDSIR	41m 35s
Sim. Study ( $R_0 \approx 1.10$ )	Stoch.Coal.SIR	1h 50m 41s
	Deter.Coal.SIR	6m 45s
	BDSIR	41m 21s
H1N1	Stoch.Coal.SIR	1h 20m 55s
	Deter.Coal.SIR	9m 44s
	BDSIR	47m 33s
HIV-1	Stoch.Coal.SIR	14h 37m 45s
	Deter.Coal.SIR	7m 56s
	BDSIR	1h 38m 54s

Table S6: Deterministic Coalescent SIR Results for Simulated Sequences:  $R_0 = 1.0987$  and  $S_0 = 499$ ,  $R_0 = 1.0989$  and  $S_0 = 999$ ,  $R_0 = 1.09945$  and  $S_0 = 1999$

$\eta$	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy
$\mathcal{R}_0$	$\approx 1.10$	1.89	1.28	0.62	0.63	0.40	52.00%
$\gamma$	0.30	0.57	0.44	0.59	0.52	2.14	95.00%
$S_{(0)}$	499	1830	1222	1.50	1.31	11.49	96.00%
$z_{(0)}$	(varies)	109.55	76.21	0.61	0.54	0.35	37.00%
$\mathcal{R}_0$	$\approx 1.10$	1.55	1.35	0.25	0.25	0.45	16.00%
$\gamma$	0.30	0.27	0.24	0.20	-0.12	1.23	61.00%
$S_{(0)}$	999	1293	804	0.27	-0.10	3.58	64.00%
$z_{(0)}$	(varies)	117.39	99.75	0.25	0.18	0.23	23.00%
$\mathcal{R}_0$	$\approx 1.10$	1.37	1.22	0.16	0.16	0.28	40.00%
$\gamma$	0.30	0.26	0.24	0.18	-0.14	1.13	64.00%
$S_{(0)}$	1999	2292	1531	0.23	-0.18	3.45	66.00%
$z_{(0)}$	(varies)	150.39	138.69	0.23	0.20	0.32	18.00%

Table S7: Simulation Study Details

Type of simulated data	Inference models used
1. Varying $R_0$ and $S_0$ (orig.) (a) $R_0 \approx 1.1$ , $S_0 = 499$ , $\gamma = 0.25$ , $\psi = 0.15$ (b) $R_0 \approx 1.2$ , $S_0 = 499$ , $\gamma = 0.30$ , $\psi = 0.15$ (c) $R_0 \approx 1.5$ , $S_0 = 499$ , $\gamma = 0.30$ , $\psi = 0.15$ (d) $R_0 \approx 1.5$ , $S_0 = 999$ , $\gamma = 0.30$ , $\psi = 0.20$ (e) $R_0 \approx 2.5$ , $S_0 = 999$ , $\gamma = 0.30$ , $\psi = 0.05$	Deter.Coal.SIR, Stoch.Coal.SIR, BDSIR Deter.Coal.SIR Deter.Coal.SIR, Stoch.Coal.SIR, BDSIR Deter.Coal.SIR, Stoch.Coal.SIR, BDSIR Deter.Coal.SIR, Stoch.Coal.SIR, BDSIR
2. Varying $S_0$ for fixed $R_0$ (a) $R_0 \approx 1.1$ , $S_0 = 499$ , $\gamma = 0.25$ , $\psi = 0.15$ (f) $R_0 \approx 1.1$ , $S_0 = 999$ , $\gamma = 0.30$ , $\psi = 0.20$ (g) $R_0 \approx 1.1$ , $S_0 = 1999$ , $\gamma = 0.30$ , $\psi = 0.09$	Deter.Coal.SIR, Stoch.Coal.SIR, BDSIR Deter.Coal.SIR Deter.Coal.SIR
3. Contemporaneous sampling (d) $R_0 \approx 1.5$ , $S_0 = 999$ , $\gamma = 0.30$ , $\psi = 0.20$ (e) $R_0 \approx 2.5$ , $S_0 = 999$ , $\gamma = 0.30$ , $\psi = 0.05$	Deter.Coal.SIR, Stoch.Coal.SIR, BDSIR Deter.Coal.SIR, Stoch.Coal.SIR, BDSIR
4. Phylogenetic uncertainty (e) $R_0 \approx 2.5$ , $S_0 = 999$ , $\gamma = 0.30$ , $\psi = 0.05$	Deter.Coal.SIR, Stoch.Coal.SIR, BDSIR
5. Reparameterization (growth rate) (e) $R_0 \approx 2.5$ , $S_0 = 999$ , $\gamma = 0.30$ , $\psi = 0.05$	Deter.Coal.SIR

Table S8: Epidemic Parameter Estimations from HIV-1 Subtype B Sequence Data

<u>Inference Model</u> HIV cluster	$R_0$	$\gamma$	$S_0$	Root of the tree (yr)	Origin $z_0$ of the epidemic (yr)
<b>Stoch. Coal. SIR</b>					
Cluster 1	3.31 (2.40 - 4.26)	0.27 (8.17E-2 - 0.48)	1165 (448 - 1974)	1971 (1946-1987)	1969 (1942-1986)
Cluster 2	2.42 (1.86 - 3.07)	0.32 (7.72E-2 - 0.57)	976 (371 - 1701)	1975 (1953 - 1988)	1972 (1947 - 1988)
Cluster 3	2.10 (1.07 - 3.73)	0.41 (2.59E-2 - 0.86)	1442 (33 - 4568)	1979 (1959 - 1990)	1973 (1943 - 1989)
Cluster 4	3.48 (1.69 - 5.24)	0.28 (0.08 - 0.61)	1757 (148 - 5260)	1964 (1922 - 1990)	1961 (1918 - 1991)
Cluster 6	3.09 (1.94 - 4.53)	0.19 (3.72E-2 - 0.41)	727 (236 - 1312)	1972 (1950 - 1989)	1970 (1947 - 1988)
<b>Deter. Coal. SIR</b>					
Cluster 1	3.43 (2.52 - 4.40)	0.26 (6.95E-2 - 0.50)	1158 (397 - 2023)	1969 (1941-1987)	1967 (1939-1986)
Cluster 2	2.34 (1.92 - 2.81)	0.38 (0.13 - 0.66)	1163 (530 - 1895)	1979 (1967 - 1989)	1977 (1964 - 1987)
Cluster 3	1.87 (1.42 - 2.43)	0.54 (0.14 - 0.97)	1298 (399 - 2267)	1979 (1965 - 1989)	1975 (1960 - 1987)
Cluster 4	3.35 (2.02 - 4.86)	0.34 (8.22E-2 - 0.65)	1479 (397 - 2792)	1972 (1948 - 1990)	1971 (1946 - 1989)
Cluster 6	3.14 (1.98 - 4.64)	0.18 (2.99E-2 - 0.37)	693 (213 - 1241)	1971 (1949 - 1989)	1969 (1943 - 1988)
<b>BDSIR</b>					
Cluster 1	3.22 (2.18-4.27)	0.30 (0.15-0.47)	880 (142-3592)	1986 (1983-1988)	1983 (1978-1987)
Cluster 2	2.45 (1.53-3.68)	0.17 (0.06-0.35)	1745 (190-8892)	1983 (1979-1986)	1978 (1968-1984)
Cluster 3	1.90 (1.22-2.78)	0.20 (0.09-0.39)	1540 (153-8558)	1985 (1981-1988)	1978 (1962-1986)
Cluster 4	2.62 (1.45-4.29)	0.15 (0.06-0.31)	1921 (128-11007)	1987 (1983-1990)	1981 (1970-1988)
Cluster 6	3.17 (1.73-5.43)	0.15 (0.06-0.31)	2862 (183-16909)	1986 (1981-1989)	1983 (1975-1989)

Table S9: Deterministic Coalescent SIR Results from Trees Simulated with Higher  $S_0$  (with Fixed  $R_0$ ) and Higher  $R_0$  (with Fixed  $S_0$ )

$\eta$	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy
$\mathcal{R}_0$	2.50	2.68	2.49	0.13	0.04	0.81	98.00%
$\gamma$	0.30	0.32	0.29	0.16	3.14E-3	1.27	99.00%
$S_{(0)}$	999	1807	1133	0.52	0.29	4.59	98.00%
$z_{(0)}$	(varies)	41.17	39.99	0.03	0.01	0.07	76.00%
$\mathcal{R}_0$	2.50	3.28	2.97	0.23	0.20	1.42	100.00%
$\gamma$	0.35	0.30	0.28	0.24	-0.20	1.28	99.00%
$S_{(0)}$	4999	7733	4838	0.34	0.03	4.18	100.00%
$z_{(0)}$	(varies)	37.45	36.15	0.03	1.48e-3	0.06	56.00%
$\mathcal{R}_0$	2.50	3.76	3.05	0.26	0.23	1.50	100.00%
$\gamma$	0.40	0.33	0.31	0.26	-0.22	1.22	100.00%
$S_{(0)}$	9999	12,609	7405	0.35	-0.15	3.31	100.00%
$z_{(0)}$	(varies)	34.99	34.28	0.04	-1.94e-3	0.05	43.00%
$\eta$	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy
$\mathcal{R}_0$	2.50	2.68	2.49	0.13	0.04	0.81	98.00%
$\gamma$	0.30	0.32	0.29	0.16	3.14E-3	1.27	99.00%
$S_{(0)}$	999	1807	1133	0.52	0.29	4.59	98.00%
$z_{(0)}$	(varies)	41.17	39.99	0.03	0.01	0.07	76.00%
$\mathcal{R}_0$	3.50	3.92	3.76	0.18	0.06	0.95	95.00%
$\gamma$	0.30	0.31	0.29	0.21	-0.01	1.16	99.00%
$S_{(0)}$	999	1909	1060	0.64	0.36	4.26	100.00%
$z_{(0)}$	(varies)	30.65	30.35	0.04	-6.35E-3	0.05	45.00%
$\mathcal{R}_0$	5.00	6.13	5.53	0.20	0.12	1.39	100.00%
$\gamma$	0.30	0.28	0.27	0.29	-0.09	1.18	100.00%
$S_{(0)}$	999	2144	1220	0.68	0.49	4.94	99.00%
$z_{(0)}$	(varies)	26.28	25.26	0.03	-0.01	0.03	52.00%
$\eta$	Truth	Mean	Median	Error	Bias	Relative HPD width	95% HPD accuracy
$\mathcal{R}_0$	5.00	7.20	6.37	0.27	0.23	2.53	100.00%
$\gamma$	0.30	0.26	0.22	0.26	-0.19	1.41	100.00%
$S_{(0)}$	9999	17,339	10,518	0.31	0.12	4.91	100.00%
$z_{(0)}$	(varies)	28.20	26.93	0.03	-0.01	1.05	36.00%