# Direct likelihood-based inference for discretely observed stochastic compartmental models of infectious disease

Lam Si Tung Ho

Department of Biostatistics

University of California, Los Angeles

Forrest W. Crawford

Departments of Biostatistics and Ecology & Evolutionary Biology

Yale University

Marc A. Suchard

Departments of Biomathematics, Biostatistics and Human Genetics

University of California, Los Angeles

**Abstract**

Stochastic compartmental models are important tools for understanding the course of infectious diseases epidemics in populations and in prospective evaluation of intervention policies. However, calculating the likelihood for discretely observed data from even simple models – such as the ubiquitous susceptible-infectious-removed (SIR) model – has been considered computationally intractable, since its formulation almost a century ago. Recently researchers have proposed methods to circumvent this limitation through data augmentation or approximation, but these approaches often suffer from high computational cost or loss of accuracy. We develop the mathematical foundation and an efficient algorithm to compute the likelihood for discretely observed data from a broad class of stochastic compartmental models. We also give expressions for the derivatives of the transition probabilities using the same technique, making possible inference via Hamiltonian Monte Carlo (HMC). We use the 17th century plague in Eyam, a classic example of the SIR model, to compare our recursion method to sequential Monte Carlo, analyze using HMC, and assess the model assumptions. We also apply our direct likelihood evaluation to perform Bayesian inference for the 2014-2015 Ebola outbreak in Guinea. The results suggest that the epidemic infectious rates have decreased since October 2014 in the Southeast region of Guinea, while rates remain the same in other regions, facilitating understanding of the outbreak and the effectiveness of Ebola control interventions.

**Keywords**   epidemic model, multivariate birth process, infectious disease, transition probabilities, Ebola

# 1   Introduction

Compartmental models have been used extensively in epidemiology to study the spread of infectious diseases such as plague (Raggett, 1982), measles (Cauchemez and Ferguson, 2008), influenza (Dukic et al., 2012), HIV (Blum and Tran, 2010), and Ebola (Althaus, 2014). These models stratify the population into separate groups according to differing health states. The famous susceptible-infectious-removed (SIR) model (McKendrick, 1926; Kermack and McKendrick, 1927) divides the population into three subpopulations: the susceptible (S) group including healthy persons who have no immunity to the disease, the infectious (I) group including infected persons who can transmit the disease to susceptible persons by contact, and the removed (R) group including recovered/dead persons who no longer affect disease dynamics. Important adaptions of the SIR model abound. For example, allowing for the loss of immunity in the removed group such that recovered persons can become susceptible again results in the susceptible-infectious-removed-susceptible (SIRS) model. As a simplification, the susceptible-infectious-susceptible (SIS) model assumes that individuals who recover from the disease have no immunity against reinfection, thus rejoin susceptible group immediately after recovery. The more complicated susceptible-exposed-infectious-removed (SEIR) model takes into account an incubation period by adding an exposed (E) group including individuals who are infected but not yet infectious.

Compartmental models have been studied in both deterministic and stochastic settings. One advantage of deterministic models is that they yield simpler statistical inference than their stochastic counterparts. However, "many infectious disease systems are fundamentally individual-based stochastic processes, and are more naturally described by stochastic models" (Roberts et al., 2015). Deterministic models are only appropriate when the populations of the compartments are sufficiently large (Brauer, 2008). Therefore, stochastic models remain preferable when their analysis is possible. If we are able to observe all transition events, likelihood-based inference for stochastic compartment models is straightforward. For example, Becker and Britton (1999) derive maximum likelihood estimates under complete

observation for the SIR model. Unfortunately, it is very unlikely that we know exactly when an individual contracts the disease. In general, surveillance data often include total counts of individuals in each compartment at several observation points. Calculation of the likelihood requires evaluating the transition probabilities of the underlying stochastic process between these time points and, thus, becomes intractable due to the requirement of integrating over all unobserved events (Cauchemez and Ferguson, 2008). Solving for the transition probabilities begins, as Renshaw (2011) reminds us, by innocuously writing out the Chapman-Kolmogorov equations for the compartmental model, but the "associated mathematical manipulations required to generate solutions can only be described as heroic."

One common solution considers stochastic compartmental models as finite, but very large, state-space Markov processes and approximates their transition probabilities using matrix exponentiation. Unfortunately, this method is extremely time consuming and numerically unstable in many instances (Schranz et al., 2008; Crawford and Suchard, 2012). Further, when the state-space is infinite, matrix exponentiation can suffer from truncation error (Crawford et al., 2016). Several alternative approaches have been developed to overcome the intractability of compartmental models, including data augmentation, diffusion approximation, sequential Monte Carlo (SMC) – namely, particle filters – and approximate Bayesian computation (ABC). However, these methods are limited and do not completely achieve tractability. In Section 2, we give a formal definition of stochastic compartmental models and discuss limitations of existing methods in more detail.

In this paper, we propose a method with polynomial complexity to compute the transition probabilities and their derivatives for stochastic compartmental models, making direct inference scalable to large epidemics. The main technique of our method is solving the Chapman-Kolmogorov equations in the Laplace domain and evaluating the inverse Laplace transform of these solutions numerically to get back the transition probabilities. Recently, this technique has been successfully applied to the SIS model (Crawford and Suchard, 2012) and the SIR model (Ho et al., 2017), where the solutions of the Chapman-Kolmogorov equa-

tions in the Laplace domain can be represented by continued fractions. Although these results make progress toward evaluating the likelihood function efficiently, applying the continued fraction representation for more complex models such as SEIR and SIRS remains an open problem. In this work, we bypass the need for an exotic continued fraction representation by constructing multivariate birth processes that are equivalent to epidemic processes of the compartmental models. Consequently, our method does not require evaluating continued fractions, and is therefore significantly faster and straightforward to apply to complex compartmental models. Section 3 explains the construction of multivariate birth process representations and the dynamic programming algorithm for computing the transition probabilities of compartmental models. In Section 4, we apply this new method to three prevailing infectious disease models (SIR, SEIR, and SIRS) and illustrate the computation gain for the SIR model compared to the method in Ho et al. (2017), the SMC method implemented in the increasingly popular R package pomp (King et al., 2016), and the matrix exponentiation method implemented in the state-of-the-art software Expokit (Sidje, 1998). We discuss two further statistical applications using our recursion which do not appear possible under previous approaches in Section 5. Specifically, we devise polynomial-time computable derivatives of the transition probabilities of the SIR model, enabling an analysis of the dynamics of an historical plague outbreak using Hamiltonian Monte Carlo (HMC). Further, the generality of our method equips us to explore the adequacy of the SIR model assumptions for this outbreak of plague. Finally, in Section 6, we turn to the 2014-2015 Ebola outbreak in Guinea and propose a time-inhomogeneous, hierarchical SIR extension that provides evidence for the slowing of this outbreak. Moreover, we find that the change in the trajectory only happened in the Southeast region of Guinea.

# 2 Stochastic compartmental models

In this section, we formally define stochastic compartmental models, discuss limitations of current inference methods when the data are observed discretely, and propose a new method of polynomial complexity for computing their transition probabilities.

## 2.1 Notation and definition

A stochastic $m$-compartmental model stratifies the population into $m$ homogeneous sub-populations called compartments. Let $\{\mathcal{C}_1, \mathcal{C}_2, \ldots, \mathcal{C}_m\}$ be the compartments and $\mathbf{Y}(t) = \{Y_1(t), Y_2(t), \ldots, Y_m(t)\}$ be their population at time $t \geq 0$, then the rate matrix $\mathbf{R}$ is an $m \times m$ matrix $[\mu_{ij}(\theta, \mathbf{Y})]_{1 \leq i,j \leq m}$ where $\mu_{ij}(\theta, \mathbf{Y}) \geq 0$ is a function of the parameter of interest $\theta$ and $\mathbf{Y}(t)$, representing an infinitesimal transition rate from $\mathcal{C}_i$ to $\mathcal{C}_j$. We set $\mu_{ii}(\theta, \mathbf{Y}) = 0$ for all $i = 1, \ldots, m$. Let $d$ count the number of positive elements of $\mathbf{R}$. Then, there are $d$ possible transitions of $\mathbf{Y}$ during a sufficient small time interval $(t, t + dt)$:

$$
\begin{aligned}
\Pr\{\mathbf{Y}(t+dt) = \mathbf{y} - \mathbf{e}_i + \mathbf{e}_j \mid \mathbf{Y}(t) = \mathbf{y}\} &= \mu_{ij}(\theta, \mathbf{y})dt + o(dt), \quad \mu_{ij} \neq 0 \\
\Pr\{\mathbf{Y}(t+dt) = \mathbf{y} \mid \mathbf{Y}(t) = \mathbf{y}\} &= 1 - \left(\sum_{i,j=1}^{m} \mu_{ij}(\theta, \mathbf{y})\right) dt + o(dt),
\end{aligned}
\tag{1}
$$

where $\mathbf{e}_i$ and $\mathbf{e}_j$ are the $i^{\text{th}}$ and $j^{\text{th}}$ coordinate vector of $\mathbb{R}^m$ respectively. We call $\mathbf{Y}(t)$ a compartmental process. We can visualize a compartmental model by a directed graph where nodes correspond to compartments and a directed edge from node $i$ to node $j$ means $\mu_{ij}$ is positive. Figure 1 gives an example of representing a 3-compartmental model by a directed graph.

## 2.2 Limitations of current approaches

The first approach for likelihood-based inference under discretely-observed stochastic compartmental models exploits data augmentation. This technique augments the observed data
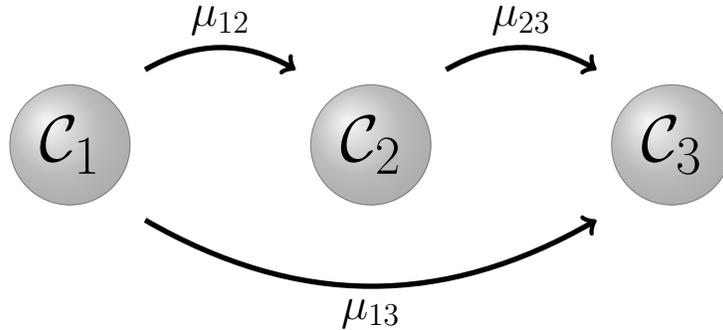
6

Figure 1: A directed graph representation of a 3-compartmental model. The rate matrix $\mathbf{R}$ of this model only has $d = 3$ positive elements: $\mu_{12}$, $\mu_{23}$, and $\mu_{13}$.

with the extensive unobserved information needed to evaluate the continuously-observed likelihood. This method often treats the times of all unobserved events as parameters and explores the joint posterior distribution by Markov chain Monte Carlo (MCMC) method (Gibson and Renshaw, 1998; O'Neill and Roberts, 1999; O'Neill, 2002). Although data augmentation works well for small epidemics, it has been criticized for being computationally prohibitive with large augmented data (Cauchemez and Ferguson, 2008; Blum and Tran, 2010).

An alternative approach to data augmentation entertains a diffusion approximation. This method approximates the discrete compartmental processes by continuous diffusion processes whose likelihood function is easy to calculate. For example, Cauchemez and Ferguson (2008) propose to mimic the SIR process by a Cox-Ingersoll-Ross process (Cox et al., 1985), and apply this approximation to study measles epidemics in London (1948-1964). However, a diffusion approximation is not applicable to epidemics in small communities because the approximation requires the state-space to be large enough to justify approximating a discrete process by a continuous one (Karev et al., 2005; Golightly and Wilkinson, 2005). Moreover, this method is often not sufficiently accurate for use even as a simulator (Golightly and Wilkinson, 2005).

Particle filters, as a SMC approach, offer another popular tool for estimating the likelihood of stochastic models (Arulampalam et al., 2002). The R package pomp (King et al., 2016)

provides an increasingly popular SMC implementation for both frequentist and Bayesian inference settings. For example, Ionides et al. (2006) develop an iterated filtering method that uses a particle filter to approximate the maximum likelihood estimates of the parameters. In the Bayesian setting, Andrieu et al. (2010) construct a particle marginal Metropolis-Hastings sampler to explore the posterior distribution using estimates from a particle filter. The computational cost of these methods can be prohibitive when the convergence is slow because each iteration requires using a particle filter to estimate the likelihood (Owen et al., 2015).

Another alternative to data augmentation is ABC (Blum and Tran, 2010). This is a likelihood-free approach replacing the observations with summary statistics and approximating the posterior of the parameters given the summary statistics by a simulation-based method. Nonetheless, the ABC method can be biased because of non-zero tolerance and non-sufficient summary statistics (Sunnåker et al., 2013), especially in high dimensions (Blum and Tran, 2010). Therefore, credible interval estimates tend to be inflated (Csilléry et al., 2010), and model selection using the ABC method cannot be trusted (Robert et al., 2011).

Finally, Faddy (1977) proposes an approximation for the stochastic SIR model by assuming that each susceptible person becomes infected independently with the same rate $\beta \times i(t)$ where $i(t)$ is the number of infected individuals in the deterministic SIR model (Kermack and McKendrick, 1927). The transition probabilities of this approximated process have analytic formulae because of the independence assumption, but this approximation becomes less accurate as the epidemic progresses.

# 3 Evaluating transition probabilities

We present a new method for computing the transition probabilities of stochastic compartmental models. Our method achieves polynomial complexity, thus enabling direct likelihood-based inference for discretely observed data. The main idea is to recast a compartmental process whose rate matrix $\mathbf{R}$ has $d$ positive elements into a $d$-dimensional birth process by

keeping track of $d$ types of transition events between compartments. This idea has been used in chemical thermodynamics for almost 100 years, where the variable measuring the progress of all substances in a chemical reaction is called the *degree of advancement* or *extent of reaction* variable (de Donder et al., 1920). By doing this, we can evaluate the transition probabilities more efficiently because the resulting multivariate birth processes are monotonically non-decreasing, while the compartment populations may increase or decrease over time. This monotonicity affords us the opportunity to apply dynamic programming for building the transition probability matrix.

## 3.1  Multivariate birth process

**Definition 1.** *A $d$-dimensional birth process is a continuous-time Markov process counting the number of "birth" events for $d$ populations. Let $\mathbf{X}(t) = \{X_1(t), X_2(t), \ldots, X_d(t)\}$, $t \geq 0$ be a multivariate birth process, whose state-space is $\mathbb{N}^d$. Then, there are $d + 1$ possible transitions of $\mathbf{X}$ during a sufficiently small time interval $(t, t + dt)$:*

$$\Pr\left\{\mathbf{X}(t + dt) = \mathbf{x} + \mathbf{e}_k \mid \mathbf{X}(t) = \mathbf{x}\right\} = \lambda_{\mathbf{x}}^{(k)} dt + o(dt), \ \ k \in \{1, 2, \ldots, d\}$$

$$\Pr\left\{\mathbf{X}(t + dt) = \mathbf{x} \mid \mathbf{X}(t) = \mathbf{x}\right\} = 1 - \left(\sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)}\right) dt + o(dt), \tag{2}$$

*where $\lambda_{\mathbf{x}}^{(k)} \geq 0$ is the birth rate of the $k^{th}$ population given the current population is $\mathbf{x} = (x_1, x_2, \ldots, x_d)$.*

For two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{N}^d$, denote $P_{\mathbf{uv}}(t) = \Pr\{\mathbf{X}(t) = \mathbf{v} \mid \mathbf{X}(0) = \mathbf{u}\}$ be the transition probability of the multivariate birth process from $\mathbf{u}$ to $\mathbf{v}$ after $t$ units of time. We say $\mathbf{u} \leq \mathbf{v}$ if $u_k \leq v_k$ for every $k = 1, 2, \ldots, d$. Notice that $P_{\mathbf{uv}}(t) \neq 0$ if and only if $\mathbf{u} \leq \mathbf{v}$.

Let $\mathbf{B} \in \mathbb{N}^d$, and set $\lambda_{\mathbf{x}}^{(k)} = 0$ if $x_k = -1$. For $i \in \mathbb{N}$, we denote

$$D_i = \left\{\mathbf{x} : \sum_{k=1}^{d} x_k = i\right\}, \ \ \text{and} \ \ \lambda_i = \max_{\mathbf{x} \in D_i}\left\{\sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)}\right\}. \tag{3}$$

Throughout this section, we make the following assumption:

**Assumption 1** (Regularity condition).

$$\sum_{i=1}^{\infty} 1/\lambda_i = \infty.$$

This condition generalizes the classic regularity condition of a univariate birth process (Feller, 1968).

**Theorem 1.** *Under Assumption 1 (Regularity condition),*

(i) *the forward transition probabilities $\{P_{\mathbf{0}\mathbf{x}}(t)\}_{\mathbf{x} \leq \mathbf{B}}$ are the unique solution of the Chapman-Kolmogorov forward equations*

$$\frac{dP_{\mathbf{0}\mathbf{x}}(t)}{dt} = \sum_{k=1}^{d} \lambda_{\mathbf{x}-\mathbf{e}_k}^{(k)} P_{\mathbf{0},\mathbf{x}-\mathbf{e}_k}(t) - \left( \sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)} \right) P_{\mathbf{0}\mathbf{x}}(t), \text{ and} \tag{4}$$

(ii) *the backward transition probabilities $\{P_{\mathbf{x}\mathbf{B}}(t)\}_{\mathbf{x} \leq \mathbf{B}}$ are the unique solution of the Chapman-Kolmogorov backward equations*

$$\frac{dP_{\mathbf{x}\mathbf{B}}(t)}{dt} = \sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)} P_{\mathbf{x}+\mathbf{e}_k,\mathbf{B}}(t) - \left( \sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)} \right) P_{\mathbf{x}\mathbf{B}}(t). \tag{5}$$

*Proof.* It is sufficient to prove that the birth rates satisfying Assumption 1 uniquely determine the multivariate birth process. By Theorem 7 in Reuter (1957), we have to show that if for some $\zeta > 0$, $\{y_{\mathbf{x}}\} \in [0,1]$ satisfies the following equations

$$\left( \zeta + \sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)} \right) y_{\mathbf{x}} = \sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)} y_{\mathbf{x}+\mathbf{e}_k}, \tag{6}$$

then $y_{\mathbf{x}} = 0$. Let $y_i = \max_{\mathbf{x} \in D_i}\{y_{\mathbf{x}}\}$ and $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x} \in D_i}\{y_{\mathbf{x}}\}$, we have

$$\left(\zeta + \sum_{k=1}^{d} \lambda_{\mathbf{x}^*}^{(k)}\right) y_i = \sum_{k=1}^{d} \lambda_{\mathbf{x}^*}^{(k)} y_{\mathbf{x}^* + \mathbf{e}_k} \leq \left(\sum_{k=1}^{d} \lambda_{\mathbf{x}^*}^{(k)}\right) y_{i+1}. \tag{7}$$

Therefore,

$$\zeta y_i \leq \left(\sum_{k=1}^{d} \lambda_{\mathbf{x}^*}^{(k)}\right)(y_{i+1} - y_i) \leq \lambda_i(y_{i+1} - y_i). \tag{8}$$

Assume that there exists $i_0 > 0$ such that $y_{i_0} > 0$. From (8), we conclude that for every

$i > i_0$, $y_i > y_{i-1}$ and

$$y_i = \sum_{j=i_0}^{i-1} \frac{\zeta}{\lambda_j} + y_{i_0} \to \infty \quad \text{as} \quad i \to \infty, \tag{9}$$

which contradicts with $y_i \leq 1$. This contradiction completes the proof. $\qquad\square$

Theorem 1 shows that we can evaluate the forward and backward transition probabilities by solving the Chapman-Kolmogorov equations (4) and (5). However, traditional methods like matrix exponentiation and Euler's method are either computationally expensive or lack numerical accuracy. Instead, we first solve the Chapman-Kolmogorov equations in the Laplace domain and then apply an inverse Laplace transform to recover $P_{\mathbf{uv}}(t)$.

We define the Laplace transform of $P_{\mathbf{uv}}(t)$ as:

$$f_{\mathbf{uv}}(s) = \mathcal{L}[P_{\mathbf{uv}}(t)](s) = \int_0^\infty e^{-st} P_{\mathbf{uv}}(t)dt. \tag{10}$$

Note that $f_{\mathbf{uv}} \neq 0$ if and only if $\mathbf{u} \leq \mathbf{v}$.

**Corollary 1.** *For the multivariate birth process, we have the following recursive formulae:*

$$
\begin{aligned}
f_{\mathbf{00}}(s) &= \frac{1}{s + \sum_{j=1}^{d} \lambda_{\mathbf{0}}^{(j)}} \\
f_{\mathbf{BB}}(s) &= \frac{1}{s + \sum_{j=1}^{d} \lambda_{\mathbf{B}}^{(j)}} \\
f_{\mathbf{0x}}(s) &= \sum_{k=1}^{d} \frac{\lambda_{\mathbf{x}-\mathbf{e}_k}^{(k)}}{s + \sum_{j=1}^{d} \lambda_{\mathbf{x}}^{(j)}} f_{\mathbf{0},\mathbf{x}-\mathbf{e}_k}(s) \\
f_{\mathbf{xB}}(s) &= \sum_{k=1}^{d} \frac{\lambda_{\mathbf{x}}^{(k)}}{s + \sum_{j=1}^{d} \lambda_{\mathbf{x}}^{(j)}} f_{\mathbf{x}+\mathbf{e}_k,\mathbf{B}}(s),
\end{aligned}
\tag{11}
$$

*where* $\mathbf{0} \leq \mathbf{x} \leq \mathbf{B}$.

*Proof.* Applying a Laplace transform to both sides of (4) and (5), we arrive at

$$
\begin{aligned}
\mathcal{L}\left[\frac{dP_{\mathbf{0x}}(t)}{dt}\right](s) &= \sum_{k=1}^{d} \lambda_{\mathbf{x}-\mathbf{e}_k}^{(k)} \mathcal{L}[P_{\mathbf{0},\mathbf{x}-\mathbf{e}_k}(t)](s) - \left(\sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)}\right) \mathcal{L}[P_{\mathbf{0x}}(t)](s) \text{ and} \\
\mathcal{L}\left[\frac{dP_{\mathbf{xB}}(t)}{dt}\right](s) &= \sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)} \mathcal{L}[P_{\mathbf{x}+\mathbf{e}_k,\mathbf{B}}(t)](s) - \left(\sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)}\right) \mathcal{L}[P_{\mathbf{xB}}(t)](s).
\end{aligned}
\tag{12}
$$

Noting that

$$
\begin{aligned}
\mathcal{L}\left[\frac{dP_{\mathbf{0x}}(t)}{dt}\right](s) &= s\mathcal{L}[P_{\mathbf{0x}}(t)](s) - P_{\mathbf{0x}}(0) \text{ and} \\
\mathcal{L}\left[\frac{dP_{\mathbf{xB}}(t)}{dt}\right](s) &= s\mathcal{L}[P_{\mathbf{xB}}(t)](s) - P_{\mathbf{xB}}(0)
\end{aligned}
\tag{13}
$$

enables us to write

$$
\begin{aligned}
s f_{\mathbf{0x}}(s) - P_{\mathbf{0x}}(0) &= \sum_{k=1}^{d} \lambda_{\mathbf{x}-\mathbf{e}_k}^{(k)} f_{\mathbf{0},\mathbf{x}-\mathbf{e}_k}(s) - \left(\sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)}\right) f_{\mathbf{0x}}(s) \text{ and} \\
s f_{\mathbf{xB}}(s) - P_{\mathbf{xB}}(0) &= \sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)} f_{\mathbf{x}+\mathbf{e}_k,\mathbf{B}}(s) - \left(\sum_{k=1}^{d} \lambda_{\mathbf{x}}^{(k)}\right) f_{\mathbf{xB}}(s).
\end{aligned}
\tag{14}
$$

From (14), we have

$$sf_{\mathbf{00}}(s) - P_{\mathbf{00}}(0) = \sum_{k=1}^{d} \lambda_{-\mathbf{e}_k}^{(k)} f_{\mathbf{0},-\mathbf{e}_k}(s) - \left( \sum_{k=1}^{d} \lambda_{\mathbf{0}}^{(k)} \right) f_{\mathbf{00}}(s) \text{ and}$$

$$sf_{\mathbf{BB}}(s) - P_{\mathbf{BB}}(0) = \sum_{k=1}^{d} \lambda_{\mathbf{B}}^{(k)} f_{\mathbf{B}+\mathbf{e}_k,\mathbf{B}}(s) - \left( \sum_{k=1}^{d} \lambda_{\mathbf{B}}^{(k)} \right) f_{\mathbf{BB}}(s). \tag{15}$$

Since $P_{\mathbf{00}}(0) = P_{\mathbf{BB}}(0) = 1$ and $P_{\mathbf{0},-\mathbf{e}_k}(t) = P_{\mathbf{B}+\mathbf{e}_k,\mathbf{B}}(t) = 0$, we deduce

$$f_{\mathbf{00}}(s) = \frac{1}{s + \sum_{j=1}^{d} \lambda_{\mathbf{0}}^{(j)}} \text{ and}$$

$$f_{\mathbf{BB}}(s) = \frac{1}{s + \sum_{j=1}^{d} \lambda_{\mathbf{B}}^{(j)}}. \tag{16}$$

Moreover, $P_{\mathbf{0x}}(0) = 0$ for $\mathbf{x} \neq \mathbf{0}$ and $P_{\mathbf{xB}}(0) = 0$ for $\mathbf{x} \neq \mathbf{B}$. Hence, from (14), we obtain

$$f_{\mathbf{0x}}(s) = \sum_{k=1}^{d} \frac{\lambda_{\mathbf{x}-\mathbf{e}_k}^{(k)}}{s + \sum_{j=1}^{d} \lambda_{\mathbf{x}}^{(j)}} f_{\mathbf{0},\mathbf{x}-\mathbf{e}_k}(s) \text{ and}$$

$$f_{\mathbf{xB}}(s) = \sum_{k=1}^{d} \frac{\lambda_{\mathbf{x}}^{(k)}}{s + \sum_{j=1}^{d} \lambda_{\mathbf{x}}^{(j)}} f_{\mathbf{x}+\mathbf{e}_k,\mathbf{B}}(s). \tag{17}$$

Thus, the proof is completed. $\qquad\square$

From Corollary 1, we can derive analytic formulae for all $\{f_{\mathbf{0x}}(s)\}_{\mathbf{x}\leq\mathbf{B}}$ and $\{f_{\mathbf{xB}}(s)\}_{\mathbf{x}\leq\mathbf{B}}$. For $\mathbf{u} \leq \mathbf{v}$, let a path from $\mathbf{u}$ to $\mathbf{v}$ be an increasing sequence $\mathbf{p} = \{\mathbf{p}_i\}_{i=1}^{n}$ such that

$$\mathbf{p}_1 = \mathbf{u}, \quad \mathbf{p}_n = \mathbf{v}, \quad \mathbf{p}_i \leq \mathbf{p}_{i+1}, \quad \text{and} \quad \mathbf{p}_{i+1} - \mathbf{p}_i \in \{\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_d\}.$$

Denote $\mathcal{P}_{\mathbf{uv}}$ and $\mathcal{I}_i$ to be the set of all paths from $\mathbf{u}$ to $\mathbf{v}$ and the index of the only non-zero

13

coordinate of $\mathbf{p}_{i+1} - \mathbf{p}_i$ respectively. We have

$$
\begin{aligned}
f_{\mathbf{0x}}(s) &= \frac{1}{s + \sum_{j=1}^{d} \lambda_{\mathbf{0}}^{(j)}} \left( \sum_{\mathbf{p} \in \mathcal{P}_{\mathbf{0x}}} \prod_{i=2}^{n} \frac{\lambda_{\mathbf{p}_{i-1}}^{(\mathcal{I}_{i-1})}}{s + \sum_{j=1}^{d} \lambda_{\mathbf{p}_i}^{(j)}} \right) \\
f_{\mathbf{xB}}(s) &= \frac{1}{s + \sum_{j=1}^{d} \lambda_{\mathbf{B}}^{(j)}} \left( \sum_{\mathbf{p} \in \mathcal{P}_{\mathbf{xB}}} \prod_{i=1}^{n-1} \frac{\lambda_{\mathbf{p}_i}^{(\mathcal{I}_i)}}{s + \sum_{j=1}^{d} \lambda_{\mathbf{p}_i}^{(j)}} \right).
\end{aligned} \tag{18}
$$

However, evaluating $\{f_{\mathbf{0x}}(s)\}_{\mathbf{x} \leq \mathbf{B}}$ and $\{f_{\mathbf{xB}}(s)\}_{\mathbf{x} \leq \mathbf{B}}$ using (18) is infeasible because the number of paths from $\mathbf{0}$ to $\mathbf{B}$ is extremely large. For example, when all the birth rates are positive, the number of paths is

$$
\prod_{i=1}^{d} \frac{\left( \sum_{j=i}^{d} B_j \right)!}{B_i! \left( \sum_{j=i+1}^{d} B_j \right)!}. \tag{19}
$$

For example, when $d = 2$ and $B_1 = B_2 = B$, the number of paths (19) becomes $(B+1)(B+2)\cdots(2B) > B^B$.

The sum-product structure in (18) suggests that dynamic programming may lead to efficient computation of $\{f_{\mathbf{0x}}(s)\}_{\mathbf{x} \leq \mathbf{B}}$ and $\{f_{\mathbf{xB}}(s)\}_{\mathbf{x} \leq \mathbf{B}}$ that we achieve through the recursive formulae (11). The computation cost of the recursion is only $\mathcal{O}(\prod_{k=1}^{d} B_k)$ because we need one loop for each coordinate. Algorithm 1 presents pseudo-code for computing $\{f_{\mathbf{0x}}(s)\}_{\mathbf{x} \leq \mathbf{B}}$ via dynamic programming. The algorithm for evaluating $\{f_{\mathbf{xB}}(s)\}_{\mathbf{x} \leq \mathbf{B}}$ is similar.

Then, we approximate the inverse Laplace transform of $f_{\mathbf{uv}}(s)$ by the method proposed in Abate and Whitt (1992, equation (4.6)):

$$
P_{\mathbf{uv}}(t) = \mathcal{L}^{-1}(f_{\mathbf{uv}})(t) \approx \frac{e^{M/2}}{2t} \mathcal{R}\left[ f_{\mathbf{uv}}\left( \frac{M}{2t} \right) \right] + \frac{e^{M/2}}{t} \sum_{k=1}^{\infty} (-1)^k \mathcal{R}\left[ f_{\mathbf{uv}}\left( \frac{M + 2k\pi i}{2t} \right) \right], \tag{20}
$$

where $\mathcal{R}[z]$ is the real part of $z$. Here, the positive number $M$ is used to control the dis-

**Algorithm 1** Dynamic programming algorithm for computing $\{f_{\mathbf{0x}}(s)\}_{\mathbf{x} \leq \mathbf{B}}$.

**Require:** $s > 0$, $\{\lambda_{\mathbf{x}}^{(j)}\}_{j=1}^d$
1: $f_{\mathbf{00}} \leftarrow 1$
2: **for** $i_1 = 0$ to $B_1$ **do**
3:     **for** $i_2 = 0$ to $B_2$ **do**
4:       $\vdots$
5:         **for** $i_d = 0$ to $B_d$ **do**
6:           $\mathbf{x} \leftarrow (i_1, i_2, \ldots, i_d)$
7:           $m \leftarrow s + \sum_{j=1}^d \lambda_{\mathbf{x}}^{(j)}$
8:           $f_{\mathbf{0x}} \leftarrow f_{\mathbf{0x}}/m$
9:           **for** $k = 1$ to $d$ **do**
10:             **if** $i_k < B_k$ **then**
11:               $f_{\mathbf{0},\mathbf{x}+\mathbf{e}_k} \leftarrow f_{\mathbf{0},\mathbf{x}+\mathbf{e}_k} + \lambda_{\mathbf{x}}^{(k)} \times f_{\mathbf{0x}}$
12:             **end if**
13:           **end for**
14:         **end for**
15:       $\vdots$
16:     **end for**
17: **end for**

cretization error. Specifically, the discretization error is

$$\sum_{k=1}^{\infty} e^{-kM} P_{\mathbf{uv}}((2k+1)t),$$

which can be bounded by $1/(e^M - 1)$. However, Abate and Whitt (1992) warn that we should not choose $M$ too large because it makes the infinite sum (20) harder to evaluate. They suggest to aim for $10^{-7}$ to $10^{-8}$ accuracy on a machine with 14-digit precision. Follow this instruction, we choose $M = 20$ throughout this paper. We opt to use a Levin acceleration method (Levin, 1973) to improve the convergence rate of (20). Let $L$ be the number of iterations required from Levin acceleration to achieve a certain error bound for the approximation (20), then we have the following corollary:

**Corollary 2.** *The total complexity of our algorithm to compute $\{P_{\mathbf{0x}}(t)\}_{\mathbf{x} \leq \mathbf{B}}$ and $\{P_{\mathbf{xB}}(t)\}_{\mathbf{x} \leq \mathbf{B}}$ is $\mathcal{O}(L \prod_{k=1}^d B_k)$.*

Note that when we aim for $10^{-8}$ accuracy, $L$ usually ranges from 100 to 1000.

## 3.2 Re-parameterization

Given an $m$-compartmental process $\mathbf{Y}(t)$ with $d$ possible types of transition between compartments, computing the transition probability $\Pr\{\mathbf{Y}(t) = \mathbf{v} \mid \mathbf{Y}(0) = \mathbf{u}\}$ by solving the compartmental Chapman-Kolmogorov equations is generally intractable because, unlike multivariate birth processes, individual compartment population $Y_i(t)$ may increase or decrease over time. Here, we recast $\mathbf{Y}(t)$ into a $d$-dimensional birth process $\mathbf{X}(t)$ and aim to compute the transition probabilities of $\mathbf{Y}(t)$ from the transition probabilities of $\mathbf{X}(t)$.

We denote $i \to j$ be a transition from compartment $\mathcal{C}_i$ to compartment $\mathcal{C}_j$. For $k = 1, 2, \ldots, d$, let $i_k \to j_k$ be the $k$-th type of transition. We construct $\mathbf{X}(t)$ by letting $X_k(t)$ be the number of $k$-type transition events happening from time $0$ to $t$. Define an $m \times d$ matrix $\mathcal{A} = [a_{lk}]$ as follows:

$$
a_{lk} = \begin{cases} -1, & \text{if } l = i_k \\ 1, & \text{if } l = j_k \\ 0, & \text{otherwise,} \end{cases} \tag{21}
$$

then we have the following lemma:

**Lemma 1.** $\mathbf{Y}(t) = \mathbf{Y}(0) + [\mathcal{A}\mathbf{X}(t)]^T$ where $T$ denotes the matrix transpose. Moreover, the birth rates for $\mathbf{X}(t)$ are $\lambda_{\mathbf{x}}^{(k)} = \mu_{i_k j_k}(\theta, \mathbf{Y}(0) + [\mathcal{A}\mathbf{x}]^T)$.

Define $W = \{\mathbf{w} \in \mathbb{N}^d : \mathcal{A}\mathbf{w} = (\mathbf{u} - \mathbf{v})^T\}$. By Lemma 1, we deduce that

$$
Pr\{\mathbf{Y}(t) = \mathbf{v} \mid \mathbf{Y}(0) = \mathbf{u}\} = \sum_{\mathbf{w} \in W} Pr\{\mathbf{X}(t) = \mathbf{w} \mid \mathbf{X}(0) = \mathbf{0}\}. \tag{22}
$$

We want to employ Equation (22) for computing the transition probabilities of $\mathbf{Y}(t)$. However, evaluating the summation in (22) is infeasible when the set $W$ has infinitely many elements. To limit the cardinality of $W$, we proffer a small restriction on the class of compartmental models for which we can compute their transition probabilities in polynomial complexity.

**Assumption 2** (Finite loops)**.** *Each individual visits each compartment at most $\mathcal{U}$ times between two consecutive observations.*

Assumption 2 is rarely restrictive for many compartmental models for infectious diseases. Infected individuals usually develop at least partial immunity to re-inflection that wanes at a rate commensurate with or slower than the observation process. Further, it is notable that if a compartmental model can be represented by a directed acyclic graph, then an individual never returns to a compartment after leaving. In this case, this assumption is satisfied with $\mathcal{U} = 1$.

**Theorem 2.** *For a compartmental model satisfying Assumption 2 (Finite loops), the complexity for computing its transition probabilities via Equation (22) is $\mathcal{O}(L\mathcal{U}^d N^d)$, where $N$ is the total population of all compartments.*

*Proof.* By Assumption 2, $\mathbf{X} \leq \mathcal{U}\mathbf{N}$ where $\mathbf{N}$ is a $d$-dimensional vector $(N_1, N_2, \ldots, N_d)$. Hence $\lambda_{\mathbf{x}}^{(k)} = 0$ when $\mathbf{x} \geq \mathcal{U}\mathbf{N}$. By Theorem 1 and Corollary 2, we can compute the transition probabilities $(Pr\{\mathbf{X}(t) = \mathbf{x} \mid \mathbf{X}(0) = \mathbf{0}\})_{\mathbf{x} \leq \mathcal{U}\mathbf{N}}$ at a cost of $\mathcal{O}(L\mathcal{U}^d N^d)$. Then, we can compute the transition probabilities of $\mathbf{Y}$ through Equation (22) at the same cost. Therefore, the total complexity is $\mathcal{O}(L\mathcal{U}^d N^d)$. $\qquad\square$

# 4 Compartmental models of infectious diseases

We apply our recursion method to three prevailing compartmental models of infectious diseases including the SIR, SEIR and SIRS models.

## 4.1 Susceptible-infectious-removed model

Proposed by McKendrick (1926), the stochastic SIR model is probably the most famous compartmental model in epidemiology. This model divides the population into three different compartments: susceptible ($S$), infectious ($I$), and removed ($R$), and allows two possible

transitions: infection $(S \to I)$ with rate $\beta SI$ and removal $(I \to R)$ with rate $\gamma I$. Here, $\beta > 0$ is the infection rate and $\gamma > 0$ is the removal rate of the disease. Figure 2 visualizes the directed graph representing this model.
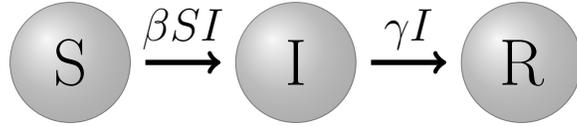
$$S \xrightarrow{\beta SI} I \xrightarrow{\gamma I} R$$

Figure 2: A directed graph representation of the SIR model.

Because the total population $S(t) + I(t) + R(t)$ is constant, Ho et al. (2017) consider $\{S(t), I(t)\}$ as a death/birth-death process and propose an algorithm to compute its transition probabilities using a continued fraction representation. The computational cost of this algorithm for evaluating the full transition probability matrix is $\mathcal{O}(LN^3)$. Our present method re-parameterizes the SIR model using number of infection events $N_{SI}(t)$ and removal events $N_{IR}(t)$. Note that there is a one-to-one correspondence between $\{S(t), I(t)\}$ and $\{N_{SI}(t), N_{IR}(t)\}$:

$$\begin{pmatrix} S(t) \\ I(t) \end{pmatrix} = \begin{pmatrix} s_0 \\ i_0 \end{pmatrix} + \begin{pmatrix} -1 & 0 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} N_{SI}(t) \\ N_{IR}(t) \end{pmatrix}, \tag{23}$$

where $(s_0, i_0)$ is the realized value of $\{S(0), I(0)\}$. It follows that $\{N_{SI}(t), N_{IR}(t)\}$ is a bivariate birth process with birth rates $\beta(s_0 - N_{SI})^+(i_0 + N_{SI} - N_{IR})^+$ and $\gamma(i_0 + N_{SI} - N_{IR})^+$ where $a^+ = \max\{a, 0\}$. Since the directed graph representing the SIR model is acyclic, Assumption 2 is satisfied with $\mathcal{U} = 1$. By Theorem 2, we have the following Corollary:

**Corollary 3.** *The complexity for evaluating the full transition probability matrix of the SIR model using our method is $\mathcal{O}(LN^2)$.*

We remark that our present method is an order of magnitude in $N$ faster than that of Ho et al. (2017) for computing the entire transition probability matrix. In practice, however, we often only need to compute the transition probabilities between observations. In this

case, the computational cost of our present method decreases further to $\mathcal{O}(L\Delta_S\Delta_R)$ where $\Delta_S$ and $\Delta_R$ are the changes in susceptible and removed populations between observations. In many situations, $\Delta_S$ and $\Delta_R$ are significantly smaller than the total population $N$, for example, when tracing the dynamics of a rare disease across an entire nation. We implement our method in the R function SIR_prob (MultiBD package) https://github.com/msuchard/MultiBD.
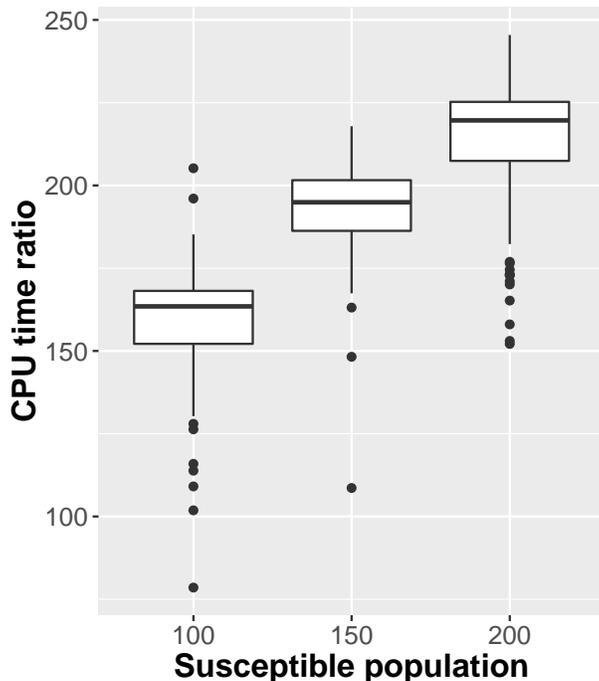


Figure 3: CPU time ratios of the continued fraction method (dbd_prob) to the proposed recursion method (SIR_prob) for computing the full transition probabilities matrix of the SIR model with $\gamma = 2.73$ and $\beta = 0.0178$. We set $I(0) = 1$ and $S(0) = 100, 150, 200$.

To illustrate the computation gain of our recursion method compared to the continued fraction representation of Ho et al. (2017), we evaluate the full forward transition probability matrix of the SIR model with $\gamma = 2.73$ and $\beta = 0.0178$ (estimated values from the Eyam plague data by Ho et al., 2017) using both methods. The death/birth-death method in Ho et al. (2017) is implemented in the R function dbd_prob (MultiBD package). We set the starting infectious population $i_0$ to be 1 and consider 3 different starting susceptible

populations $s_0 = 100, 150, 200$. For each scenario, we repeat the evaluation a hundred times and compare the computing times and the results from both methods. Figure 3 summarizes this comparison, and we see that `SIR_prob` is more than 150 times faster than `dbd_prob`. On the other hand, the two methods return similar transition probability matrices whose $L_1$ distance is less than $10^{-12}$. Here, the $L_1$ distance between two matrices $\mathbf{A} = (a_{ij})$ and $\mathbf{B} = (b_{ij})$ is $\sum_{ij} |a_{ij} - b_{ij}|$.

## 4.2 Susceptible-exposed-infectious-removed model

The SEIR model extends the SIR model by adding an exposed (E) compartment. We visualize the SEIR model by the directed acyclic graph in Figure 4.
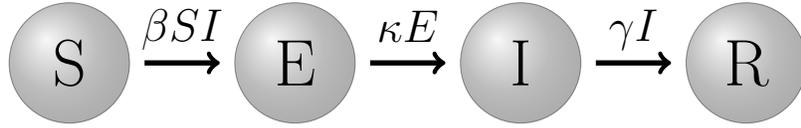


$$S \xrightarrow{\beta SI} E \xrightarrow{\kappa E} I \xrightarrow{\gamma I} R$$

Figure 4: A directed graph representation of the SEIR model.

Let $\{N_{SE}(t), N_{EI}(t), N_{IR}(t)\}$ be the number of transition events $S \to E$, $E \to I$, and $I \to R$ respectively. Then, we have an one-to-one correspondence with $\{S(t), E(t), I(t)\}$ as follows:

$$\begin{pmatrix} S(t) \\ E(t) \\ I(t) \end{pmatrix} = \begin{pmatrix} s_0 \\ e_0 \\ i_0 \end{pmatrix} + \begin{pmatrix} -1 & 0 & 0 \\ 1 & -1 & 0 \\ 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} N_{SE}(t) \\ N_{EI}(t) \\ N_{IR}(t) \end{pmatrix}, \tag{24}$$

where $(s_0, e_0, i_0)$ is the realized value of $\{S(0), E(0), I(0)\}$. Again, $\{N_{SE}, N_{EI}, N_{IR}\}$ is a trivariate birth process with birth rates $\beta(s_0 - N_{SE})^+ (i_0 + N_{EI} - N_{IR})^+$, $\kappa(e_0 + N_{SE} - N_{EI})^+$, and $\gamma(i_0 + N_{EI} - N_{IR})^+$. By Theorem 2, we have:

**Corollary 4.** *The complexity for evaluating the full transition probability matrix of the SEIR model using our method is $\mathcal{O}(LN^3)$.*

## 4.3 Susceptible-infectious-removed-susceptible model

For some diseases, removed persons can lose immunity, making possible transition from the "recovered" (R) to "susceptible" (S) compartments. The SIRS model takes into account these scenarios by allowing the transition $R \to S$. Figure 5 visualizes the directed graph representing the SIRS model.
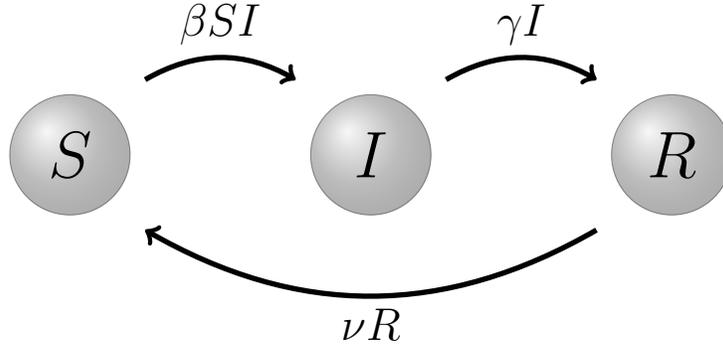


Figure 5: A directed graph representation of the SIRS model.

Denote $\{N_{SI}(t), N_{IR}(t), N_{RS}(t)\}$ as the number of transition events $S \to I$, $I \to R$, and $R \to S$ respectively. We have

$$
\begin{pmatrix} S(t) \\ I(t) \\ R(t) \end{pmatrix} = \begin{pmatrix} s_0 \\ i_0 \\ r_0 \end{pmatrix} + \begin{pmatrix} -1 & 0 & 1 \\ 1 & -1 & 0 \\ 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} N_{SI}(t) \\ N_{IR}(t) \\ N_{RS}(t) \end{pmatrix},
\tag{25}
$$

where $(s_0, i_0, r_0)$ is the realized value of $\{S(0), I(0), R(0)\}$. In this situation, $\{N_{SI}, N_{IR}, N_{RS}\}$ is a trivariate birth process with birth rates $\beta(s_0 + N_{RS} - N_{SI})^+(i_0 + N_{SI} - n_{IR})^+$, $\gamma(i_0 + n_{SI} - n_{IR})^+$, and $\nu(r_0 + n_{IR} - n_{RS})^+$. In practice, $\nu$ is much smaller than $\beta \times I(t)$ and $\gamma$. Hence, we can assume that during $(0, t)$ each individual can only be infected at most $\mathcal{U}$ times. By Theorem 2, we arrive at

**Corollary 5.** *The complexity for evaluating the full transition probability matrix of the SIRS model using our method is $\mathcal{O}(L\mathcal{U}^3 N^3)$.*

## 4.4 Comparisons

We use prevalence counts from the plague in Eyam from June 18th to October 20th, 1666 (Raggett, 1982) to compare our recursion method with the SMC algorithm implemented in the R function `pfilter` (`pomp` package King et al., 2016) and the matrix exponentiation method implemented in the state-of-the-art software `Expokit` (Sidje, 1998). Plague is a deadly infectious disease caused by the bacterium *Yersinia pestis*. It is mainly spread by infected fleas from small animals, particularly rodents, and has killed 100s of millions of people through human history. In Eyam, only 83 of the original 350 villagers survived at the end of the plague. The data contain the susceptible and infectious populations $\{(s_m, i_m)\}_{m=1}^n$ in Eyam at time $\{t_m\}_{m=1}^n$. The log likelihood function is

$$\log l(\beta, \gamma | \{(s_m, i_m)\}_{m=1}^n) = \sum_{m=1}^{n-1} \log \Pr \left\{ \begin{array}{c|c} S(t_{m+1}) = s_{m+1} & S(t_m) = s_m \\ I(t_{m+1}) = i_{m+1} & I(t_m) = i_m \end{array} \right\}. \qquad (26)$$

We compute the log likelihood (26) under the stochastic SIR model with $\beta = 0.0178$ and $\gamma = 2.73$ (estimated values from the Eyam plague data by Ho et al., 2017).

### 4.4.1 Comparing to sequential Monte Carlo

The likelihood calculation is repeated a thousand times and the number of attempted simulant particles for each estimation for `pfilter` is set as $1000, 2000, 3000,$ and $4000$. For these data and parameter estimates, `pfilter` fails to achieve a 100% success rate for approximating the likelihood. The success rate is low with 1000 particles (only 20.1%), and increases as the number of particles increases (see Table 1). Filtering failure occurs when all particles become incompatible with the data counts; this can happen frequently when the counts are observed without error. When filtering succeeds, the approximation is fairly similar to our method, and the standard deviation of these approximations, while sizable, decreases from 1.28 to 0.97 as the number of particles increases. When filtering fails, the approximation is

off target by a large margin. The computation time of `pfilter` is about 10 times slower compared to our algorithm for every 1000 particles (Table 1). This comparison shows that our recursion method is faster than the SMC method. Moreover, our method is stable while approximations using SMC are very unstable due to a high failure rate. It is worth mentioning that SMC is known to be an inefficient algorithm for computing the likelihood when the observations have no error.

| Number of particles | 1000 | 2000 | 3000 | 4000 |
|---|---|---|---|---|
| Success rate | 20.1% | 53.2% | 71.1% | 78.8% |
| Average time ratio | 10.14 | 20.11 | 30.09 | 40.1 |
| Standard deviation | 1.28 | 1.11 | 1.06 | 0.97 |

Table 1: Success rates of sequential Monte Carlo method (`pfilter`) and its average computing time ratios compared to our algorithm.

### 4.4.2 Comparing to matrix exponentiation method

To evaluate the log likelihood (26) via matrix exponentiation, we use the function `expv` in `expoRkit`, an `R`-interface to the Fortran package `Expokit`, to compute the transition probabilities. Again, the likelihood calculation is repeated a thousand times. Our method and matrix exponentiation method produce similar results: the difference is less than $1.53 \times 10^{-7}$. In term of speed, the average CPU computation time ratio of matrix exponentiation method to our method is 15 and the standard deviation is 1. Therefore, our method is more efficient in computing the likelihood function of the stochastic SIR model than matrix exponentiation method.

## 5 Further statistical applications

The ability to efficiently compute the likelihood function makes it straightforward to use maximum likelihood estimators and Metropolis-Hasting algorithms for Bayesian inference. In this section, we provide two additional extensions that the recursion opens up to us that

were unavailable with previous methods. The first application is inference via HMC, which requires evaluating the derivative of the posterior distribution with respect to the unknown model parameters. The second application is accessing model adequacy for the classic SIR model using Bayes factors.

## 5.1 Inference via Hamiltonian Monte Carlo

HMC is a MCMC method using Hamiltonian dynamics to produce proposals for sampling from a continuous distribution on $\mathbb{R}^d$. Hamiltonian dynamics contain "location" variables $q$, that are the parameters of interest, and nuisance "momentum" variables $p$ (see Neal et al., 2011, for an excellent review). In a Bayesian setting, we may treat the negative log of the posterior distribution as the potential energy function:

$$U(q) = -\log\left[l(q|\mathbf{D})\pi(q)\right], \tag{27}$$

where $l(q|\mathbf{D})$ is the likelihood given data $\mathbf{D}$ and $\pi(q)$ is the prior distribution. On the other hand, researchers often place a multivariate Normal distribution $\mathcal{N}(0, \Sigma)$ on $p$ and let $p$ be independent of $q$. Typically, $\Sigma$ is the identity matrix and the corresponding kinetic energy function is

$$K(p) = \sum_{i=1}^{d} \frac{p_i^2}{2}.$$

The Hamiltonian is defined as $H(q, p) = U(q) + K(p)$, and the Hamiltonian dynamics follow the following system of partial differential equations:

$$\begin{aligned} \frac{dq_i}{dt} &= \frac{\partial H}{\partial p_i} = p_i \\ \frac{dp_i}{dt} &= -\frac{\partial H}{\partial q_i} = -\frac{\partial U}{\partial q_i}. \end{aligned} \tag{28}$$

The HMC algorithm consists two steps. In the first step, a proposal for $p$ is sampled from $\mathcal{N}(0, \Sigma)$. In the second step, $(q_t, p_t)$ is obtained from the Hamiltonian dynamics (28) starting

at the current value $(q_c, p_c)$. In practice, we may use a leapfrog integration scheme to approximate the solution of (28). The proposal $(q^*, p^*)$ is set as $(q_t, -p_t)$ and is accepted with probability

$$\min \left[ 1, e^{H(q_c, p_c) - H(q^*, p^*)} \right]. \tag{29}$$

The ability to efficiently compute the derivatives of the transition probabilities with respect to $q$ opens the possibility of using HMC for studying infectious disease epidemics. To illustrate, we employ HMC to analyze the 17ᵗʰ century plague in Eyam. Denote

$$P_m = \Pr \left\{ \begin{array}{c} S(t_{m+1}) = s_{m+1} \;\middle|\; S(t_m) = s_m \\[2mm] I(t_{m+1}) = i_{m+1} \;\middle|\; I(t_m) = i_m \end{array} \right\}. \tag{30}$$

Then, the log likelihood function (26) can be written as

$$\log l(\beta, \gamma | \{(s_m, i_m)\}_{m=1}^n) = \sum_{m=1}^{n-1} \log P_m. \tag{31}$$

To satisfy positivity constraints, we opt to use $(u, v) := (\log \beta, \log \gamma)$ as our parameters instead of $(\beta, \gamma)$. To apply HMC, we derive the derivatives of $\log l$ with respect to $u$ and $v$:

$$\begin{aligned} \frac{\partial \log l}{\partial u} &= \frac{\partial \log l}{\partial \beta} \frac{\partial \beta}{\partial u} = \sum_{m=1}^{n-1} \frac{P_m^{(\beta)}}{P_m} \beta \\ \frac{\partial \log l}{\partial v} &= \frac{\partial \log l}{\partial \gamma} \frac{\partial \gamma}{\partial v} = \sum_{m=1}^{n-1} \frac{P_m^{(\gamma)}}{P_m} \gamma. \end{aligned} \tag{32}$$

We assume *a priori* that $u \sim \mathcal{N}(0, 100^2)$ and $v \sim \mathcal{N}(0, 100^2)$. We explore the posterior distribution of $(u, v)$ using HMC with 10000 iterations and discard the first 2000 iterations. Figure 6 visualizes the posterior density of $(u, v)$. This result is similar to the density estimation using a Metropolis-Hasting algorithm performed in Ho et al. (2017), but at a significant time cost savings. The average effective sample size per unit-time of HMC is 60-fold larger, mostly owing to substantial computational order reduction in the likelihood

evaluation under our multivariate-birth process formulation. The posterior means of $\beta$ and $\gamma$ are 0.0197 and 3.22. The 95% Bayesian credible intervals are $(0.0164, 0.0234)$ and $(2.69, 3.83)$ respectively.
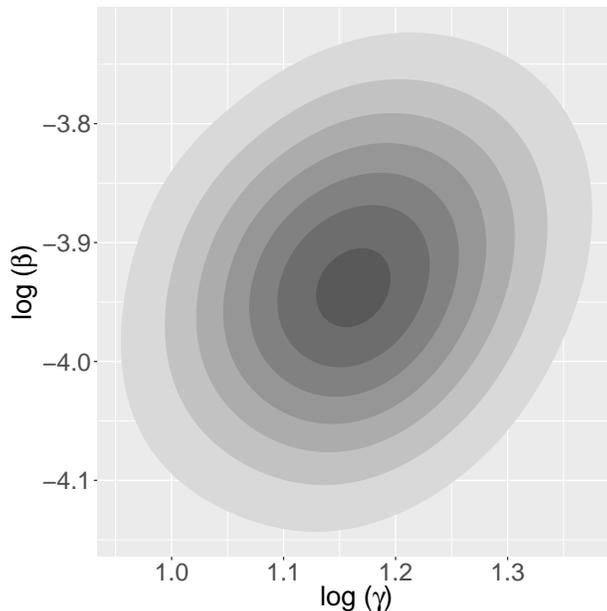
Figure 6: Posterior density of infection $\beta$ and removal $\gamma$ rates of the Eyam plague.

## 5.2 Adequacy of the classic SIR model

Although the classic SIR model has been used extensively in practice, it makes a strong assumption that each infected person can independently transmit the disease to one susceptible person with rate $\beta$. ONeill and Wen (2012) argue that this assumption may not be realistic in settings where a saturation effect occurs; that is, a newly infected person contributes less to the overall infection pressure. Therefore, the authors propose to consider a general SIR model with infection rate $\beta S I^\omega$. This model is a special case of a more general SIR model where the infection rate is $\beta S^\alpha I^\omega$ and the removal rate is $\gamma I^\eta$ (Severo, 1969).

Our computational method does not require any special structure for the infection and removal rates, thus can also be applied to evaluate the likelihood function under these general SIR models. In particular, $\{N_{SI}(t), N_{IR}(t)\}$ in the general SIR model from Severo (1969) is
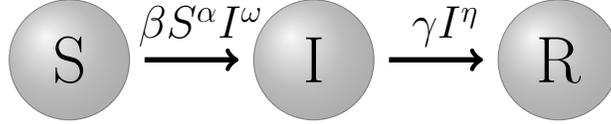
Figure 7: A directed graph representation of the general SIR model (Severo, 1969).

a bivariate birth process with birth rates $\beta[(s_0 - N_{SI})^+]^\alpha[(i_0 + N_{SI} - N_{IR})^+]^\omega$ and removal rates $\gamma[(i_0 + N_{SI} - N_{IR})^+]^\eta$. So, we can address some questions about model adequacy of the classic SIR model. To illustrate, we use Bayes factors to assess if the classic SIR model is appropriate for the Eyam plague dynamics (Raggett, 1982). In particular, we test between the general SIR model against its nested sub-models. Since the tests are between nested models, we apply the Savage-Dickey density ratio to evaluate the Bayes factors (Verdinelli and Wasserman, 1995). To be specific, if model $\mathcal{M}_0$ with parameter $(\theta = 0, \phi)$ is nested within model $\mathcal{M}_1$ with parameter $(\theta, \phi)$ and the prior $p_0(\phi)$ under $\mathcal{M}_0$ is proportional to the prior $p_1(\theta = 0, \phi)$ under $\mathcal{M}_1$, then the Bayes factor $B_{01}$ in favor of $\mathcal{M}_0$ over $\mathcal{M}_1$ can be estimated via the marginal posterior distribution under $\mathcal{M}_1$ as follows

$$B_{01} = \frac{p(\theta = 0|\mathbf{Y}, \mathcal{M}_1)}{p(\theta = 0|\mathcal{M}_1)} \tag{33}$$

where $p(\theta = 0|\mathbf{Y}, \mathcal{M}_1)$ and $p(\theta = 0|\mathcal{M}_1)$ are marginal posterior and prior densities of $\theta$ evaluated at 0 under model $\mathcal{M}_1$. Here, we posit independent log-normal priors $\ln \mathcal{N}(0, 100^2)$ for each parameter. Therefore, the condition for applying the Savage-Dickey density ratio is satisfied. To estimate the posterior distribution under the general SIR model, we use our MCMC tools. The marginal posterior densities are estimated using kernel density estimation implemented in the R package `ks` (Duong et al., 2007) and these estimates are then used to compute the Bayes factors via the Savage-Dickey density ratio. Table 5.2 lists these Bayes factors, and we can see that they strongly support the classic SIR model over the general SIR model. Although Savage-Dickey density ratio is not the best approximation method for Bayes factors, we can safely ignore its drawback because the evidence supporting the classic SIR model is overwhelming.

| Model $\mathcal{M}_0$ | $\log_{10} B_{01}$ |
|---|---|
| $\alpha = \omega = \eta = 1$ | 6.9 |
| $\alpha = \omega = 1$ | 4.4 |
| $\omega = \eta = 1$ | 4.7 |
| $\alpha = \eta = 1$ | 4.6 |
| $\alpha = 1$ | 2.2 |
| $\omega = 1$ | 2.2 |
| $\eta = 1$ | 2.6 |

Table 2: Bayes factors $B_{01}$ in favor of nested models $\mathcal{M}_0$ over the general SIR model $\mathcal{M}_1$ estimated using the Savage-Dickey density ratio.

# 6   Ebola outbreak in Guinea

Ebola is a contagious viral hemorrhagic fever caused by *Zaire ebolavirus*. The fatality rate of Ebola is very high, up to 70.8% (WHO Ebola Response Team, 2014). The 2014-2015 Ebola outbreak in West Africa is the largest Ebola epidemic in history. In this section, we focus on the outbreak in Guinea from January 2014 to May 2015 (73 weeks). During this period, the World Health Organization (WHO) has convened 5 meetings of the IHR Emergency Committee regarding the Ebola outbreak in West Africa. The first three meetings happened in three consecutive month August, September and October 2014. During the fourth meeting in January 2015, World Health Organization (2015) noted that the number of Ebola cases in Guinea had decreased since the third meeting. WHO Ebola Response Team (2015) also confirmed that the Ebola outbreak has slowed down since October 2014.

We study this change in the trajectory of the outbreak using the number of reported Ebola cases reported weekly in 19 prefectures across Guinea. To be specific, we are interested in finding evidence that the outbreak in Guinea became less severe after the third WHO meeting and in what regions this happened. These 19 prefectures are the only places in Guinea where Ebola cases were reported both before and after the third WHO meeting.

We employ a hierarchical and time-inhomogeneous, but still Markovian, SIR model to analyze these data. We re-parameterize the SIR model by replacing the infection rate $\beta$ with the basic reproduction number $R_0 := \beta N/\gamma$. Basic reproduction number is an important concept

in epidemiology and can be interpreted as the average number of secondary infections caused by a new infectious individual in a susceptible population. When $R_0 < 1$ the disease will die out, and when $R_0 > 1$ the disease will be able to spread in the population. Researchers often use the value of $R_0$ to measure the severity of an epidemic. To simplify the analysis, we assume that the population of each prefecture is closed. In other words, we naïvely assume that the movement between prefectures and the movement in and out of Guinea are negligible. This assumption is violated if large number of healthy persons or small number of infected persons enter (or leave) a prefecture. We obtain the total populations of these prefectures from the 2014 census https://en.wikipedia.org/wiki/Prefectures_of_Guinea.

We use a "week" as the unit for time in this analysis. Letting $t_0$ be the week when the third WHO meeting happened, our model proceeds as follows: the Ebola cases of each prefecture follow a conditionally independent SIR process with parameters $R_{0p}(t)$ and $\gamma_p$ for prefecture $p = 1, \ldots, 19$. Further, $R_{0p}(t)$ is a time-inhomogeneous function that satisfies

$$\log R_{0p}(t) = \log r_{0p} + 1_{\{t \geq t_0\}} \log \delta_p, \tag{34}$$

where $r_{0p}$ quantifies the basic reproduction number before $t_0$ and $\delta_p$ is the scale factor by which the basic reproduction number changes after $t_0$ in prefecture $p$. Moreover, we assume a simple hierarchical prior distribution

$$(\log r_{0p}, \log \delta_p, \log \gamma_p)^t \sim \mathcal{N}\left(\mathbf{M}, \mathrm{diag}(\mathbf{\Sigma})\right), \tag{35}$$

where $\mathbf{M} = (\mu_r, \mu_\delta, \mu_\gamma)$ is the grand-mean on the log-scale across prefectures and $\mathbf{\Sigma} = (\sigma_r^2, \sigma_\delta^2, \sigma_\gamma^2)$ is the variance, with relatively uninformative conjugate hyperpriors

$$\mu_\phi \sim \mathcal{N}\left(\mathbf{0}, 10^2\right), \text{ and } \sigma_\phi^2 \sim \mathrm{InverseGamma}\left(10^{-3}, 10^{-3}\right), \quad \phi \in \{r, \delta, \gamma\}. \tag{36}$$

Of primary scientific interest, $\delta_p < 1$ corresponds to a reduction in the basic reproduction

number in prefecture $p$, suggesting that the Ebola outbreak slowed down in that prefecture. However, an important limitation of the data arises, in that field epidemiologists were only able to record the number of new cases between time points. The number of removals is unknown. To overcome this limitation, we use a Metropolis-within-Gibbs scheme to sample the posterior distribution of the rate parameters and the number of removals (see Appendix B for more details). Because we can compute the joint transition probability matrix between time points, we can draw directly from the full conditional distribution of the removal number, leading to substantially more efficient numerical integration than previous data augmentation approaches that require all sufficient statistics of the completely observed likelihood. Further, we can speed up this sampling scheme by updating the unknown parameters in each prefecture in parallel. The result is summarized in Figure 8, where we plot estimates of the basic reproduction number for each prefecture on the map of Guinea. Yellow circles represent $r_{0p}$ and blue circles represent $r_{0p} \times \delta_p$ when the posterior probability that $\delta_p < 1$ is greater than 97.5% Note that there is no posterior evidence supporting $\delta_p > 1$ for any $p$ because the posterior probability that $\delta_p > 1$ is less than 0.5 for all $p$. The radius of each circle reports a posterior mean estimate. We present the posterior means and 95% Bayesian credible interval of $\mathbf{M}$ and of $\boldsymbol{\Sigma}$ in Table 3.

| Parameter | Posterior mean | 95% Bayesian credible interval |
|---|---|---|
| $\mu_r$ | $7.47 \times 10^{-2}$ | $(-0.425,\ 15.1\ ) \times 10^{-2}$ |
| $\mu_\delta$ | $-1.25 \times 10^{-1}$ | $(\ -2.36,\ -0.0844) \times 10^{-1}$ |
| $\mu_\gamma$ | $-6.76 \times 10^{-1}$ | $(\ -10.4,\ -3.19\ ) \times 10^{-1}$ |
| $\sigma_r^2$ | $4.67 \times 10^{-3}$ | $(\ 0.399,\ 24.1\ ) \times 10^{-3}$ |
| $\sigma_\delta^2$ | $2.24 \times 10^{-2}$ | $(\ 0.216,\ 8.18\ ) \times 10^{-2}$ |
| $\sigma_\gamma^2$ | $5.98 \times 10^{-1}$ | $(\ 2.84,\ 12.0\ ) \times 10^{-1}$ |

Table 3: Posterior mean and 95% Bayesian credible interval of hierarchical parameters $\mathbf{M}$ and $\boldsymbol{\Sigma}$.

We note that *a posteriori* $\sigma_\gamma^2$ is larger than $\sigma_r^2$ or $\sigma_\delta^2$, with probability approaching 1. Therefore, the removal rate $\gamma$ varies across the country more than the reproduction number
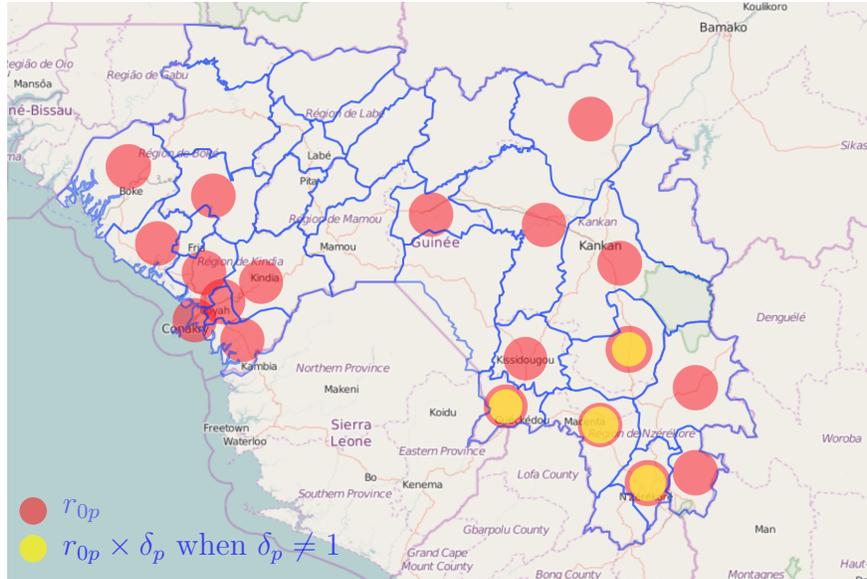
Figure 8: Basic reproduction numbers for 19 prefectures in Guinea before and after the third WHO meeting.

$R_0$. The posterior of $(\mu_\delta, \sigma_\delta^2)$ provides evidence for the slowing down of the Ebola outbreak in Guinea after the third WHO meeting. However, Figure 8 suggests that the epidemic only slowed down in the Southeast region of Guinea while the epidemic in other regions seems to stay the same. This finding gives a clearer picture of the change in the trajectory of the Ebola epidemic in Guinea. It raises a very practical question: what made the outbreak in the Southeast region of Guinea to slow down? Answering this question could help in efforts to find a more effective method for controlling Ebola epidemic.

# 7 Discussion

In this paper, we develop an algorithm to compute the transition probabilities of stochastic compartmental models for inference from surveillance data. We introduce a new representation for compartmental processes using multivariate birth processes and, through this representation, avoid the need for continued fraction evaluation to solve the Chapman-Kolmogorov equations. With quadratic complexity in number of transitions between observations, our approach emerges as computationally more efficient than previous methods for the ubiqui-

tous SIR model and applicable to a larger class of compartmental models, such as the SEIR and SIRS models. Further performance gains through embarrassingly parallel evaluation of the series in Equation (20) remain open.

Since the formulation of the SIR model over 90 years ago, many have viewed its transition probabilities as beyond reach. We provide some brief intuition on why the Laplace transform of the transition probabilities carries mere quadratic complexity $\mathcal{O}(\Delta_S \Delta_R)$. Viewed as a multivariate birth process that conveniently only increases, the transition probabilities we seek are related to the waiting time until the $\Delta_S$ and $\Delta_R$ births have occurred. Inter-birth times are independent exponential random variables with potentially unique rates, and we can arrive at the distribution of the total waiting time through taking a convolution of $\Delta_S + \Delta_R$ of these exponentials. However, the rates depend on the order of births and there are $(\Delta_S + \Delta_R)!/\Delta_S!\Delta_R!$ possible orderings. Putting these pieces together, the transition probabilities are then exponential sums of multiple convolutions. We recall several properties of Laplace transformations. First, they are linear operators, so sums in probability-space remain sums in the transformed space. Second, convolutions metamorphose into simple multiplication in the transformed space. These properties leave us with a sum-product expression, suggesting a distribution of the sums within the products. To gain insight into this dynamic programming, consider the $\Delta_S \times \Delta_R$ lattice graph. Each lattice path from $(0,0)$ to $(\Delta_S, \Delta_R)$ represents one possible ordering of the birth events. If we want the transformed probability of ending at $(\Delta_S, \Delta_R)$, there are only two possible one-shorter paths that could have gotten us there, specifically $(0,0)$ to $(\Delta_{S-1}, \Delta_R)$ or $(\Delta_S, \Delta_{R-1})$. So, the resulting transformed probability becomes the sum of the two shorter-path transformed probabilities, each multiplied by the Laplace transform of an exponential random variable that has a simple, closed-form expression. Consequentially, in filling out the whole lattice graph, we need to visit each point once in increasing order and there are only $\Delta_S \Delta_R$ points.

Because differentiation is also a linear operator, our recursion method remains pertinent for computing the derivatives of the transition probabilities with respect to the unknown pa-

32

rameters of the compartmental model. This feature makes HMC-based Bayesian inference feasible. As the number of unknown parameters in the compartmental models grows, we suspect HMC to generally outperform Metropolis-Hastings algorithms using standard transition kernels. Equally noteworthy, our algorithm does not require any specific structure in the birth rates $\lambda_x^{(\cdot)}$ of the multivariate birth processes. Therefore, we can apply our method to other general stochastic epidemic models such as one proposed by Severo (1969). This opens the possibility to access the model adequacy of traditional epidemic models. It is worth noticing that our method only works for time-homogeneous rates between observations. When the rates depend on time, the Chapman-Kolmogorov equations in the Laplace domain do not have analytic formulae making the current tool inapplicable. Therefore, an important subject for future direction of this work is extending to time-inhomogeneous processes.

Finally, we examine the 2014-2015 Ebola outbreak in Guinea using a marginalized, hierarchical and time-inhomogeneous Markovian SIR model. By applying our recursion method, we can effectively explore the posterior distribution of the basic reproductive number and removal rate across the country, while simultaneously integrating out the unobserved removed population sizes using a Metropolis-within-Gibbs scheme. This example highlights the flexibility of a Bayesian framework for direct likelihood-based inference for a compartmental model when one or more of the compartments are missing or immeasurable, as is common in infectious disease surveillance. Our results provide evidence for the slowing down of this epidemic in the Southeast region of Guinea. Several important extensions are immediately obvious. For example, we assume no error in the reported Ebola case counts, but a simple modification similar to that we accomplished for missing compartments can relax this assumption.

# Acknowledgments

# A  Derivatives of the transition probabilities of SIR model

We propose an efficient method to evaluate the derivatives of the transition probabilities of the SIR model. Again, we use the bivariate birth presentation for this model. Denote $X = N_{SI}$ and $Y = N_{IR}$, and consider the forward transition probability $P_{xy}(t) = \Pr\{X(t) = x, Y(t) = y \mid X(0) = 0, Y(0) = 0\}$. The forward Chapman-Kolmogorov equations are:

$$
\begin{aligned}
\frac{dP_{xy}(t)}{dt} =& \beta(s_0 - x + 1)^+(i_0 + x - 1 - y)^+ P_{x-1,y}(t) \\
&+ \gamma(i_0 + x - y + 1)^+ P_{x,y-1}(t) \\
&- [\beta(s_0 - x)^+(i_0 + x - y)^+ + \gamma(i_0 + x - y)^+]P_{xy}(t).
\end{aligned}
\tag{37}
$$

Let $P_{xy}^{(\beta)}$ be the derivative of $P_{xy}$ with respect to $\beta$. From (37), we have

$$
\begin{aligned}
\frac{dP_{xy}^{(\beta)}(t)}{dt} =& \beta(s_0 - x + 1)^+(i_0 + x - 1 - y)^+ P_{x-1,y}^{(\beta)}(t) \\
&+ \gamma(i_0 + x - y + 1)^+ P_{x,y-1}^{(\beta)}(t) \\
&- [\beta(s_0 - x)^+(i_0 + x - y)^+ + \gamma(i_0 + x - y)^+]P_{xy}^{(\beta)}(t) \\
&+ (s_0 - x + 1)^+(i_0 + x - 1 - y)^+ P_{x-1,y}(t) \\
&- (s_0 - x)^+(i_0 + x - y)^+ P_{xy}(t)
\end{aligned}
\tag{38}
$$

Denote $f_{xy}$ and $f_{xy}^{(\beta)}$ be the Laplace transform of $P_{xy}$ and $P_{xy}^{(\beta)}$ respectively. Taking Laplace transform to both sides of (38), we have

$$
\begin{aligned}
s f_{xy}^{(\beta)}(s) - P_{xy}^{(\beta)}(0) =& \beta(s_0 - x + 1)^+(i_0 + x - 1 - y)^+ f_{x-1,y}^{(\beta)}(s) \\
&+ \gamma(i_0 + x - y + 1)^+ f_{x,y-1}^{(\beta)}(s) \\
&- [\beta(s_0 - x)^+(i_0 + x - y)^+ + \gamma(i_0 + x - y)^+] f_{xy}^{(\beta)}(s) \\
&+ (s_0 - x + 1)^+(i_0 + x - 1 - y)^+ f_{x-1,y}(s) \\
&- (s_0 - x)^+(i_0 + x - y)^+ f_{xy}(s).
\end{aligned}
\tag{39}
$$

Since $P_{xy}(0) = 1_{\{x=0,y=0\}}$ for all $\beta$, we deduce that $P_{xy}^{(\beta)}(0) = 0$. Therefore, we can compute $f_{xy}^{(\beta)}$ using the following recursion

$$
\begin{aligned}
f_{00}^{(\beta)}(s) =& -\frac{s_0 i_0 f_{00}(s)}{s + \beta s_0 i_0 + \gamma i_0} \\
f_{xy}^{(\beta)}(s) =& \frac{\beta(s_0 - x + 1)^+(i_0 + x - 1 - y)^+ f_{x-1,y}^{(\beta)}(s)}{s + \beta(s_0 - x)^+(i_0 + x - y)^+ + \gamma(i_0 + x - y)^+} \\
&+ \frac{\gamma(i_0 + x - y + 1)^+ f_{x,y-1}^{(\beta)}(s)}{s + \beta(s_0 - x)^+(i_0 + x - y)^+ + \gamma(i_0 + x - y)^+} \\
&+ \frac{(s_0 - x + 1)^+(i_0 + x - 1 - y)^+ f_{x-1,y}(s)}{s + \beta(s_0 - x)^+(i_0 + x - y)^+ + \gamma(i_0 + x - y)^+} \\
&- \frac{(s_0 - x)^+(i_0 + x - y)^+ f_{xy}(s)}{s + \beta(s_0 - x)^+(i_0 + x - y)^+ + \gamma(i_0 + x - y)^+}.
\end{aligned}
\tag{40}
$$

Then, we can compute $P_{xy}^{(\beta)}$ by approximating the inverse Laplace transform using (20). Similarly, we can derive the recursive formulae for $f_{xy}^{(\gamma)}$:

$$
\begin{aligned}
f_{00}^{(\gamma)}(s) &= -\frac{i_0 f_{00}(s)}{s + \beta s_0 i_0 + \gamma i_0} \\
f_{xy}^{(\beta)}(s) &= \frac{\beta(s_0 - x + 1)^+ (i_0 + x - 1 - y)^+ f_{x-1,y}^{(\beta)}(s)}{s + \beta(s_0 - x)^+ (i_0 + x - y)^+ + \gamma(i_0 + x - y)^+} \\
&\quad + \frac{\gamma(i_0 + x - y + 1)^+ f_{x,y-1}^{(\beta)}(s)}{s + \beta(s_0 - x)^+ (i_0 + x - y)^+ + \gamma(i_0 + x - y)^+} \\
&\quad + \frac{(i_0 + x - y + 1)^+ f_{x,y-1}(s)}{s + \beta(s_0 - x)^+ (i_0 + x - y)^+ + \gamma(i_0 + x - y)^+} \\
&\quad - \frac{(i_0 + x - y)^+ f_{xy}(s)}{s + \beta(s_0 - x)^+ (i_0 + x - y)^+ + \gamma(i_0 + x - y)^+},
\end{aligned}
\tag{41}
$$

and evaluate $P_{xy}^{(\gamma)}$ using (20).

# B  Metropolis-within-Gibbs algorithm for inference of Ebola dynamics in West Africa

Let $\mathbf{t}^{(p)} = (t_1^{(p)}, t_2^{(p)}, \ldots, t_{m_p}^{(p)})$ be the times when the counts of Ebola cases in prefecture $p$ are reported. We define $\mathbf{N}_{SI}^{(p)}$ and $\mathbf{N}_{IR}^{(p)}$ be the total numbers of new infection and removal events at $\mathbf{t}_{-1}^{(p)}$ respectively. Here, $\mathbf{t}_{-j}^{(p)}$ denotes the vector $\mathbf{t}^{(p)}$ without the $j^{\text{th}}$ coordinate. Notice that we only observe the total of Ebola cases at $\mathbf{t}^{(p)}$, thus we only know $\mathbf{N}_{SI}^{(p)}$. So, our unknown parameters are $\{\mathbf{N}_{IR}^{(p)}, r_{0p}, \delta_p, \gamma_p\}$ for all $p$ and $(\mathbf{M}, \boldsymbol{\Sigma})$. We update our parameters using a Metropolis-within-Gibbs algorithm as follows:

1. For every $p = 1, \ldots, 19$ in parallel,

   (i) For every $j = 2, 3, \ldots, m_p$, we can compute $\mathbf{P}(\mathbf{N}_{IR}^{(p)}(t_j) = n \mid \mathbf{N}_{SI}^{(p)}, \mathbf{N}_{IR}^{(p)}(\mathbf{t}_{-j}^{(p)}), r_{0p}, \delta_p, \gamma_p)$ using the forward and backward transition probabilities of the SIR model. Therefore, we sample from $\mathbf{N}_{IR}^{(p)}(t_j) \mid \mathbf{N}_{SI}^{(p)}, \mathbf{N}_{IR}^{(p)}(\mathbf{t}_{-j}^{(p)}), r_{0p}, \delta_p, \gamma_p$ directly to update the value of $\mathbf{N}_{IR}^{(p)}(t_j)$.

(ii) Then, we update $r_{0p}, \delta_p, \gamma_p \mid \mathbf{N}_{SI}^{(p)}, \mathbf{N}_{IR}^{(p)}$ on the log-scale using a random-walk Metropolis-Hasting algorithm with Gaussian proposals or HMC. This step is straight forward because we can evaluate the density $l(r_{0p}, \delta_p, \gamma_p \mid \mathbf{N}_{SI}^{(p)}, \mathbf{N}_{IR}^{(p)})$ efficiently.

2. Finally, since we choose conjugate priors for the hierarchical parameters, we Gibbs sample $\mathbf{M}$ and $\mathbf{\Sigma}$ .

Note that we update $\mathbf{N}_{IR}^{(p)}(t_j)$ sequentially instead of sampling from the joint distribution of $\mathbf{N}_{IR}^{(p)}$ because sampling sequentially only requires transition probability matrices between counts of Ebola cases, which is much smaller compared to the full transition probability matrix of size $N^2 \times N^2$, where $N$ is the total population, required for sampling from the joint distribution.

# References

Abate, J. and W. Whitt (1992). The Fourier-series method for inverting transforms of probability distributions. *Queueing Systems 10*(1-2), 5–87.

Althaus, C. L. (2014). Estimating the reproduction number of Ebola virus (EBOV) during the 2014 outbreak in West Africa. *PLOS Currents Outbreaks 6*.

Andrieu, C., A. Doucet, and R. Holenstein (2010). Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology) 72*(3), 269–342.

Arulampalam, M. S., S. Maskell, N. Gordon, and T. Clapp (2002). A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking. *Signal Processing, IEEE Transactions on 50*(2), 174–188.

Becker, N. G. and T. Britton (1999). Statistical studies of infectious disease incidence. *Journal of the Royal Statistical Society: Series B (Statistical Methodology) 61*(2), 287–307.

Blum, M. G. and V. C. Tran (2010). HIV with contact tracing: a case study in approximate Bayesian computation. *Biostatistics 11*(4), 644–660.

Brauer, F. (2008). Compartmental models in epidemiology. In *Mathematical Epidemiology*, pp. 19–79. Springer.

Cauchemez, S. and N. M. Ferguson (2008). Likelihood-based estimation of continuous-time epidemic models from time-series data: application to measles transmission in London. *Journal of the Royal Society Interface 5*(25), 885–897.

Cox, J. C., J. E. Ingersoll, and S. A. Ross (1985). A theory of the term structure of interest rates. *Econometrica 53*(2), 385–407.

Crawford, F. W., T. C. Stutz, and K. Lange (2016). Coupling bounds for approximating birth-death processes by truncation. *Statistics & probability letters 109*, 30–38.

Crawford, F. W. and M. A. Suchard (2012). Transition probabilities for general birth–death processes with applications in ecology, genetics, and evolution. *Journal of Mathematical Biology 65*(3), 553–580.

Csilléry, K., M. G. Blum, O. E. Gaggiotti, and O. François (2010). Approximate Bayesian computation (ABC) in practice. *Trends in Ecology & Evolution 25*(7), 410–418.

de Donder, T., F. van den Dungen, and G. van Lerberghe (1920). *Leçons de thermodynamique et de chimie physique*. Number v. 1 in Leçons de thermodynamique et de chimie physique. Gauthier-Villars et cie.

Dukic, V., H. F. Lopes, and N. G. Polson (2012). Tracking epidemics with Google flu

trends data and a state-space SEIR model. *Journal of the American Statistical Association 107*(500), 1410–1426.

Duong, T. et al. (2007). ks: Kernel density estimation and kernel discriminant analysis for multivariate data in R. *Journal of Statistical Software 21*(7), 1–16.

Faddy, M. (1977). Stochastic compartmental models as approximations to more general stochastic systems with the general stochastic epidemic as an example. *Advances in Applied Probability 9*(3), 448–461.

Feller, W. (1968). *An Introduction to Probability Theory and its Applications*, Volume 1. John Wiley & Sons.

Gibson, G. J. and E. Renshaw (1998). Estimating parameters in stochastic compartmental models using Markov chain methods. *Mathematical Medicine and Biology 15*(1), 19–40.

Golightly, A. and D. J. Wilkinson (2005). Bayesian inference for stochastic kinetic models using a diffusion approximation. *Biometrics 61*(3), 781–788.

Ho, L. S. T., J. Xu, F. W. Crawford, V. V. Minin, and M. A. Suchard (2017). Birth/birth-death processes and their computable transition probabilities with biological applications. *Journal of Mathematical Biology*, in press.

Ionides, E., C. Bretó, and A. King (2006). Inference for nonlinear dynamical systems. *Proceedings of the National Academy of Sciences 103*(49), 18438–18443.

Karev, G. P., F. S. Berezovskaya, and E. V. Koonin (2005). Modeling genome evolution with a diffusion approximation of a birth-and-death process. *Bioinformatics 21*(Suppl 3), iii12–iii19.

Kermack, W. and A. McKendrick (1927). A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London. Series A 115*(772), 700–721.

King, A. A., D. Nguyen, and E. L. Ionides (2016). Statistical inference for partially observed Markov processes via the R package pomp. *Journal of Statistical Software 69*(12), 1–43.

Levin, D. (1973). Development of non-linear transformations for improving convergence of sequences. *International Journal of Computer Mathematics 3*(1-4), 371–388.

McKendrick, A. (1926). Applications of mathematics to medical problems. *Proceedings of the Edinburgh Mathematics Society 44*, 98–130.

Neal, R. M. et al. (2011). MCMC using Hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo 2*, 113–162.

O'Neill, P. D. (2002). A tutorial introduction to Bayesian inference for stochastic epidemic models using Markov chain Monte Carlo methods. *Mathematical Biosciences 180*(1), 103–114.

O'Neill, P. D. and G. O. Roberts (1999). Bayesian inference for partially observed stochastic epidemics. *Journal of the Royal Statistical Society: Series A (Statistics in Society) 162*(1), 121–129.

Owen, J., D. J. Wilkinson, and C. S. Gillespie (2015). Scalable inference for Markov processes with intractable likelihoods. *Statistics and Computing 25*(1), 145–156.

ONeill, P. and C. Wen (2012). Modelling and inference for epidemic models featuring non-linear infection pressure. *Mathematical biosciences 238*(1), 38–48.

Raggett, G. (1982). A stochastic model of the Eyam plague. *Journal of Applied Statistics 9*(2), 212–225.

Renshaw, E. (2011). *Stochastic Population Processes: Analysis, Approximations, Simulations.* Oxford University Press Oxford, UK.

Reuter, G. E. H. (1957). Denumerable Markov processes and the associated contraction semigroups on l. *Acta Mathematica 97*(1), 1–46.

Robert, C. P., J.-M. Cornuet, J.-M. Marin, and N. S. Pillai (2011). Lack of confidence in approximate Bayesian computation model choice. *Proceedings of the National Academy of Sciences 108*(37), 15112–15117.

Roberts, M., V. Andreasen, A. Lloyd, and L. Pellis (2015). Nine challenges for deterministic epidemic models. *Epidemics 10*, 49–53.

Schranz, H. W., V. B. Yap, S. Easteal, R. Knight, and G. A. Huttley (2008). Pathological rate matrices: from primates to pathogens. *BMC Bioinformatics 9*(1), 550.

Severo, N. C. (1969). Generalizations of some stochastic epidemic models. *Mathematical Biosciences 4*(3-4), 395–402.

Sidje, R. B. (1998). Expokit: a software package for computing matrix exponentials. *ACM Transactions on Mathematical Software (TOMS) 24*(1), 130–156.

Sunnåker, M., A. G. Busetto, E. Numminen, J. Corander, M. Foll, and C. Dessimoz (2013). Approximate Bayesian computation. *PLoS Comput Biol 9*(1), e1002803.

Verdinelli, I. and L. Wasserman (1995). Computing Bayes factors using a generalization of the Savage-Dickey density ratio. *Journal of the American Statistical Association 90*, 614–618.

WHO Ebola Response Team (2014). Ebola virus disease in West Africa-the first 9 months of the epidemic and forward projections. *N Engl J Med 371*(16), 1481–95.

WHO Ebola Response Team (2015). West African Ebola epidemic after one year-slowing but not yet under control. *N Engl J Med 372*(6), 584–7.

World Health Organization (2015). Statement on the 4th meeting of the IHR Emergency Committee on the 2014 Ebola outbreak in West Africa. *World Health Organization, IHR Emergency Committee regarding Ebola*.