



Published in final edited form as:

*J Comput Graph Stat.* 2017 ; 26(4): 918–929. doi:10.1080/10618600.2017.1328365.

## Efficient Data Augmentation for Fitting Stochastic Epidemic Models to Prevalence Data

Jonathan Fintzi<sup>1</sup>, Xiang Cui<sup>2</sup>, Jon Wakefield<sup>1,2</sup>, and Vladimir N. Minin<sup>2,3</sup>

<sup>1</sup>Department of Biostatistics, University of Washington, Seattle

<sup>2</sup>Department of Statistics, University of Washington, Seattle

<sup>3</sup>Department of Biology, University of Washington, Seattle

### Abstract

Stochastic epidemic models describe the dynamics of an epidemic as a disease spreads through a population. Typically, only a fraction of cases are observed at a set of discrete times. The absence of complete information about the time evolution of an epidemic gives rise to a complicated latent variable problem in which the state space size of the epidemic grows large as the population size increases. This makes analytically integrating over the missing data infeasible for populations of even moderate size. We present a data augmentation Markov chain Monte Carlo (MCMC) framework for Bayesian estimation of stochastic epidemic model parameters, in which measurements are augmented with subject-level disease histories. In our MCMC algorithm, we propose each new subject-level path, conditional on the data, using a time-inhomogeneous continuous-time Markov process with rates determined by the infection histories of other individuals. The method is general, and may be applied to a broad class of epidemic models with only minimal modifications to the model dynamics and/or emission distribution. We present our algorithm in the context of multiple stochastic epidemic models in which the data are binomially sampled prevalence counts, and apply our method to data from an outbreak of influenza in a British boarding school.

### Keywords

Bayesian data augmentation; continuous-time Markov chain; epidemic count data; hidden Markov model; stochastic epidemic model

## 1 Introduction

Stochastic epidemic models (SEMs) are classic tools for modeling the spread of infectious diseases. A SEM represents the time evolution of an epidemic in terms of the disease histories of individuals as they transition through disease states. Incorporating stochasticity into epidemic models is important when the disease prevalence is low or when the population size is small. In both cases, the stochastic variability in the evolution of an epidemic greatly influences the probability and severity of an outbreak, as well as the conclusions we draw about its dynamics (Keeling and Rohani, 2008, Allen, 2008). Moreover, many questions — e.g., what is the outbreak size distribution? What is the

probability that a disease has been eradicated? — cannot be answered using deterministic methods (Britton, 2010).

The task of fitting a SEM is typically complicated by the limited extent of epidemiological data, which are recorded at discrete observation times, commonly describe just one aspect of the disease process, e.g., infections, and usually capture only a fraction of cases. Complete subject–level data, which would consist of the exact times at which individuals transition through disease states, are of-ten unavailable (O’Neill, 2010). Fitting SEMs in the absence of complete subject–level data presents a complicated latent variable problem since it is usually impossible to analytically integrate over the missing data (O’Neill, 2002). This makes the observed data likelihood for a SEM intractable.

Existing approaches to fitting SEMs with intractable likelihoods have largely fallen into four groups: martingale methods, approximation methods, simulation based methods, and data augmentation (DA) methods (O’Neill, 2010). Martingale methods estimate the parameters of interest using estimating equations based on martingales for the counting processes within the SEM, e.g., infections and recoveries (Becker, 1977, Watson, 1981, Sudbury, 1985, Andersson and Britton, 2000, Linden-strand and Svensson, 2013). These methods are not easily implemented for SEMs with complex dynamics fit to partially observed count data. Approximation methods replace the SEM, typically represented as a Markov jump process, with a simpler model whose likelihood is more tractable. For example, Roberts and Stramer (2001) and Cauchemez and Ferguson (2008) use diffusion processes that approximate the SEM dynamics, while Jandarov et al. (2014) use a Gaussian process approximation of a related gravity model. Another typical simplification is to discretize time and to construct a transition model for the population flow between model compartments at successive times (Longini Jr. and Koopman, 1982, Held et al., 2005, Lekone and Finkenstaedt, 2006, Held and Paul, 2012). These methods are computationally efficient and in many cases yield sensible estimates. However, the simplifying assumptions used in the various approximations are not always appropriate. For instance, the diffusion approximation may not be valid in small populations where the system is far from its deterministic limit (Andersson and Britton, 2000), while the discretization of time makes it awkward to approximate systems in which the observation times are not evenly spaced or the rates of transition events span several orders of magnitude (Glass et al., 2003, Shelton and Ciardo, 2014). Simulation based methods use the underlying SEM to generate epidemic paths that serve as the basis for inference. This class of methods includes approximate Bayesian computation (ABC) methods (McKinley et al., 2009, Toni et al., 2009), pseudo–marginal methods (McKinley et al., 2014), and sequential Monte Carlo (or particle filter) methods (Toni et al., 2009, Andrieu et al., 2010, Ionides et al., 2011, Dukic et al., 2012, Koepke et al., 2016). Within this class of methods, the particle marginal Metropolis–Hastings algorithm of Andrieu et al. (2010) stands out in being a general method for Bayesian inference and is used as a benchmark method in this paper. Although simulation–based methods have been used to fit complex models, they are computationally intensive and suffer from well known pitfalls. ABC methods are sensitive to the choice of summary statistic, rejection threshold, and prior (Toni et al., 2009). Sequential Monte Carlo methods, on which pseudo–marginal methods often rely, are prone to “particle impoverishment” problems (Cappé et al., 2006, Dukic et al., 2012).

Traditional agent-based DA methods for fitting SEMs, first presented by O'Neill and Roberts (1999) and Gibson and Renshaw (1998), target the joint posterior distribution of the missing data and model parameters to obtain a tractable complete data likelihood. That the augmentation is agent-based refers to the fact that subject-level disease histories, rather than population-level epidemic paths, are introduced as latent variables in the model. The advantage of the agent-based approach is that household structure and subject-level covariates may be incorporated into the model (Auranen et al., 2000, Höhle and Jørgensen, 2002, Cauchemez et al., 2004, Neal and Roberts, 2004, O'Neill, 2009). Development of DA methods for SEMs is of continuing interest, and recent works by Pooley et al. (2015), Qin and Shelton (2015), and Shestopaloff and Neal (2016) have presented methods that could possibly be applied to epidemic count data. However, their algorithms forgo the flexibility of agent-based DA, and in the case of the latter two papers have not been applied to SEMs.

We present an agent-based DA Markov chain Monte Carlo (MCMC) framework for fitting SEMs to time series count data. We obtain a tractable complete data likelihood by augmenting the data with subject-level disease histories. Our MCMC targets the joint posterior distribution of the latent epidemic process and the model parameters as we alternate between updating subject-level paths and model parameters. We propose each new subject-path, conditionally on the data, using a time-inhomogeneous continuous-time Markov chain (CTMC) with rates determined by the disease histories of the other individuals. These data-driven path proposals result in highly efficient perturbations to the latent epidemic path, and enable us to analyze epidemic count data in the absence of any subject-level information. In contrast, traditional agent-based DA MCMC algorithms rely on data-agnostic trans-dimensional proposals and suffer from convergence issues as the fraction of missing information becomes large (Roberts and Stramer, 2001, McKinley et al., 2014, Pooley et al., 2015). The *de facto* need for some subject-level data has precluded the use of classical DA machinery in many settings. Thus, our MCMC algorithm enables exact Bayesian inference for SEMs fit to datasets that would have been impossible to study with existing agent-based DA methods. Finally, our algorithm is not specific to any particular SEM dynamics or measurement process, and may be applied, with minimal modifications, to a broad class of SEMs.

## 2 The Data Augmentation Algorithm for an SIR Model

For concreteness and clarity of exposition, we present our Bayesian DA algorithm (BDA) in the context of fitting a stochastic Susceptible-Infected-Recovered (SIR) model to binomially distributed prevalence counts. We also use our algorithm to fit Susceptible-Exposed-Infected-Recovered (SEIR) and Susceptible-Infected-Recovered-Susceptible (SIRS) models in Sections 3.1, 3.2, and 4, and outline the minimal adaptations required for these models in Section S6 of Supplementary Materials.

The SIR model describes the time evolution of an epidemic in terms of the disease histories of individuals as they transition through three states — susceptible (S), infected/infectious (I), and recovered (R). Under simple SIR dynamics, each individual becomes infectious immediately upon becoming infected, and acquires lifelong immunity upon recovery. For simplicity, we assume that the population is closed and mixes homogeneously, and that there

is no external force of infection. Therefore, the epidemic ceases once the pool of infectious individuals is depleted.

## 2.1 Measurement process and data

Our data,  $\mathbf{Y} = \{Y_1, \dots, Y_L\}$ , are disease prevalence counts recorded at times  $t_1, \dots, t_L \in [t_1, t_L]$ . It should not be a matter of belief that the data could be subject to measurement error, for example underreporting in settings where asymptomatic individuals escape detection. Let  $S_\tau$ ,  $I_\tau$ , and  $R_\tau$  denote the total susceptible, infected, and recovered people at time  $\tau$ . We model the observed prevalence as a binomial sample, with constant detection probability  $\rho$ , of the true prevalence at each observation time. Thus,

$$Y_\ell | I_{t_\ell}, \rho \sim \text{Binomial}(I_{t_\ell}, \rho). \quad (1)$$

## 2.2 Latent epidemic process

The data are sampled from a latent epidemic process,  $\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_N\}$ , that evolves continuously in time as individuals become infected and recover. The state space of this process is  $\mathcal{S} = \{S, I, R\}^N$ , the Cartesian product of  $N$  state labels taking values in  $\{S, I, R\}$ . The state space of a single subject,  $\mathbf{X}_j$ , is  $\mathcal{S}_j = \{S, I, R\}$ , and a realized subject-path is of the form

$$\mathbf{x}_j(\tau) = \begin{cases} S, & \tau < \tau_I^{(j)}, \\ I, & \tau_I^{(j)} \leq \tau < \tau_R^{(j)}, \\ R, & \tau_R^{(j)} \leq \tau, \end{cases} \quad (2)$$

where  $\tau_I^{(j)}$  and  $\tau_R^{(j)}$  are the infection and recovery times for subject  $j$  (though subject  $j$  may also never become infected or recover, or may become infected or recover outside of the observation period  $[t_1, t_L]$ ). We write the configuration of  $\mathbf{X}$  at time  $\tau$  as  $\mathbf{X}(\tau) = (\mathbf{X}_1(\tau), \dots, \mathbf{X}_N(\tau))$ , and adopt the convention that  $\mathbf{X}(\tau)$  and derived quantities, e.g.,  $I_\tau$ , depend on the configuration just before  $\tau$ . We use  $\tau^+$  for quantities evaluated just after a particular time. The waiting times between transition events are taken to be exponentially distributed, and we denote by  $\beta$  and  $\mu$  the per-contact infectivity and recovery rates. Thus, the latent epidemic process evolves according to a time-homogeneous CTMC, with transition rate from configuration  $\mathbf{X}$  to  $\mathbf{X}'$  given by

$$\lambda_{\mathbf{X}, \mathbf{X}'} = \begin{cases} \beta I, & \text{if } \mathbf{X} \text{ and } \mathbf{X}' \text{ differ only in subject } j, \text{ with } \mathbf{X}_j = S, \text{ and } \mathbf{X}'_j = I, \\ \mu, & \text{if } \mathbf{X} \text{ and } \mathbf{X}' \text{ differ only in subject } j, \text{ with } \mathbf{X}_j = I, \text{ and } \mathbf{X}'_j = R, \\ 0, & \text{for all other configurations } \mathbf{X} \text{ and } \mathbf{X}'. \end{cases} \quad (3)$$

At the first observation time, we let  $\mathbf{X}(t_1) | \mathbf{p}_{t_1} \sim \text{Categorical}(\{S, I, R\}, \mathbf{p}_{t_1})$ , where  $\mathbf{p}_{t_1} = (p_S, p_I, p_R)$  are the probabilities that an individual is susceptible, infected, or recovered. Let  $\boldsymbol{\tau} = \{\tau_0, \dots, \tau_{K+1}\}$ , where  $t_1 \equiv \tau_0$  and  $t_L \equiv \tau_{K+1}$ , be the (ordered) set of  $K$  infection and recovery times of all individuals along with the endpoints of the observation period  $[t_1, t_L]$ . Let  $\mathbb{I}(\tau_k \hat{=} I)$  and  $\mathbb{R}(\tau_k \hat{=} R)$  indicate whether  $\tau_k$  is an infection or recovery time, and let  $\boldsymbol{\theta} = (\beta, \mu, \rho, \mathbf{p}_{t_1})$  denote the vector of unknown parameters. The complete data likelihood is

$$\begin{aligned} L(\mathbf{X}, \mathbf{Y} | \boldsymbol{\theta}) &= \Pr(\mathbf{Y} | \mathbf{X}, \rho) \times \Pr(\mathbf{X}(t_1) | \mathbf{p}_{t_1}) \times \pi(\mathbf{X} | \mathbf{X}(t_1), \beta, \mu) \\ &= \left[ \prod_{\ell=1}^L \binom{I_{t_\ell}}{Y_{t_\ell}} \rho^{Y_{t_\ell}} (1-\rho)^{I_{t_\ell} - Y_{t_\ell}} \right] \times \left[ p_S^{S_{t_1}} p_I^{I_{t_1}} p_R^{R_{t_1}} \right] \\ &\times \prod_{k=1}^K \left\{ \left[ \beta I_{\tau_k} \times \mathbb{I}(\tau_k \hat{=} I) + \mu \times \mathbb{I}(\tau_k \hat{=} R) \right] \exp \left[ -(\tau_k - \tau_{k-1}) \left( \beta I_{\tau_k} S_{\tau_k} + \mu I_{\tau_k} \right) \right] \right\} \\ &\times \exp \left[ -(t_L - \tau_K) \left( \beta I_{\tau_K} S_{\tau_K} + \mu I_{\tau_K} \right) \right]. \end{aligned} \quad (4)$$

We briefly reconcile what might seem like a discrepancy between the SIR model presented above and the lumped construction of the SIR model (see Andersson and Britton (2000)), which, for a number of computational and analytical reasons, is somewhat more common. Our model describes the time evolution of the subject-level collection of disease histories, and thus evolves on the state space of individual disease labels. The lumped SIR model describes the time evolution of the vector of compartment counts, the state space of which is defined as the partition of the original state space obtained by aggregating the individuals in each of the model compartments. The lumped construction would have been appropriate had we chosen to augment the data with the compartment counts (for example, as in Pooley et al. (2015)). Nonetheless, inference based on the full subject-level model will exactly match inference based on the lumped model. We discuss this further in Section S1 of Supplementary Materials.

### 2.3 Subject-path proposal framework

The observed data likelihood in the posterior  $\pi(\boldsymbol{\theta} | \mathbf{Y}) \propto \pi(\mathbf{Y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) = \int L(\mathbf{Y} | \mathbf{X}, \boldsymbol{\theta}) \pi(\mathbf{X} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\pi(\mathbf{X})$  is analytically intractable for even moderately sized  $N$  as it involves an extremely high dimensional integral over the collection of subject-paths,  $\mathbf{X}$ . The strategy employed in DA methods is to introduce the subject-paths,  $\mathbf{X}$ , as latent variables in the model. This enables us to work with the tractable complete data likelihood given by (4). The joint posterior distribution is

$$\pi(\boldsymbol{\theta}, \mathbf{X} | \mathbf{Y}) \propto \Pr(\mathbf{Y} | \mathbf{X}, \rho) \times \pi(\mathbf{X} | \mathbf{X}(t_1), \beta, \mu) \times \Pr(\mathbf{X}(t_1) | \mathbf{p}_{t_1}) \times \pi(\beta) \pi(\mu) \pi(\rho) \pi(\mathbf{p}_{t_1}), \quad (5)$$

where  $\pi(\beta)$ ,  $\pi(\mu)$ ,  $\pi(\rho)$ , and  $\pi(\mathbf{p}_{t_1})$  are prior densities. Our MCMC targets the joint posterior distribution, given by (5), as we alternate between updating  $\mathbf{X}|\theta, \mathbf{Y}$  and  $\theta|\mathbf{X}, \mathbf{Y}$ .

Given the current collection of subject–paths,  $\mathbf{x}^{\text{cur}}$ , we propose  $\mathbf{x}^{\text{new}}$  by sampling the path of a single subject  $\mathbf{X}_j$ , conditionally on the data, using a time–inhomogeneous CTMC with state space  $\mathcal{S}_j$  and rates conditioned on the collection of disease histories of the other individuals,  $\mathbf{x}_{(-j)} = \{\mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \mathbf{x}_{j+1}, \dots, \mathbf{x}_N\}$ . The proposed collection of paths is accepted or rejected in a Metropolis–Hastings step.

Let  $\tau^{(j)} = \{\tau_I^{(j)}, \tau_R^{(j)}\}$  be the (possibly empty) set of infection and recovery times for subject  $j$ , and define  $\tau^{(-j)} = \{\tau \mid \tau^{(j)}\} = \{\tau_0^{(-j)}, \tau_1^{(-j)}, \dots, \tau_M^{(-j)}, \tau_{M+1}^{(-j)}\}$ , where  $t_1 \equiv \tau_0^{(-j)}$  and  $t_L \equiv \tau_{M+1}^{(-j)}$ , to be the set of  $M \leq K$  (ordered) times at which other subjects become infected or recover, along with  $t_1$  and  $t_L$ . Let  $\mathcal{I} = \{\mathcal{I}_1, \dots, \mathcal{I}_{M+1}\}$  be the intervals that partition  $[t_1, t_L]$ , i.e.  $\mathcal{I}_1 = [\tau_0^{(-j)}, \tau_1^{(-j)})$ ,  $\mathcal{I}_2 = [\tau_1^{(-j)}, \tau_2^{(-j)})$ ,  $\dots$ ,  $\mathcal{I}_{M+1} = [\tau_M^{(-j)}, \tau_{M+1}^{(-j)})$ . Let  $I_\tau^{(-j)} = \sum_{i \neq j} \mathbb{1}(\mathbf{X}_i(\tau) = I)$  be the prevalence at time  $\tau$ , excluding subject  $j$ . Let  $\Lambda^{(-j)}(\theta) = \{\Lambda_1^{(-j)}(\theta), \dots, \Lambda_{M+1}^{(-j)}(\theta)\}$  be the sequence of rate matrices corresponding to each interval in  $\mathcal{I}$ , where for  $m = 1, \dots, M+1$ ,

$$\Lambda_m^{(-j)}(\theta) = \begin{matrix} & \begin{matrix} S & I & R \end{matrix} \\ \begin{matrix} S \\ I \\ R \end{matrix} & \begin{pmatrix} S & I & R \\ -\beta I_{\tau_m}^{(-j)} & \beta I_{\tau_m}^{(-j)} & 0 \\ 0 & -\mu & \mu \\ 0 & 0 & 0 \end{pmatrix} \end{matrix}. \quad (6)$$

We can construct the transition probability matrix for subject  $j$  over interval  $I_m$ ,

$$\mathbf{P}^{(j)}(\tau_{m-1}, \tau_m) = (p_{a,b}^{(j)}(\tau_{m-1}, \tau_m))_{a,b \in \mathcal{S}_j},$$

where  $p_{a,b}^{(j)}(\tau_{m-1}, \tau_m) = \Pr(\mathbf{X}_j(\tau_m) = b \mid \mathbf{X}_j(\tau_{m-1}) = a, \theta)$ , using the matrix exponential

$$\mathbf{P}^{(j)}(\tau_{m-1}, \tau_m) = \exp[(\tau_m - \tau_{m-1})\Lambda_m^{(-j)}(\theta)].$$

This computation requires an eigen–decomposition of each rate matrix. We may reduce the total computational burden by computing the eigen decompositions analytically, and by caching the decompositions to avoid duplicate computations. One additional point is that while the eigen– values of any SIR rate matrix are always real valued, this is not generally true, e.g., it is possible for the rate matrix of an SIRS model to have complex eigenvalues. In this case, we obtain a real valued transition probability matrix by first applying a rotation to

each rate matrix with complex eigenvalues to obtain its real canonical form (Hirsch et al., 2013). This is discussed in Section S2 of Supplementary Materials.

By the Markov property, the time–inhomogeneous CTMC density over the observation period  $[t_1, t_L]$ , denoted  $\pi(\mathbf{X}_j | \mathbf{x}_{(-j)}, \boldsymbol{\theta}) \equiv \pi(\mathbf{X}_j | \Lambda^{(-j)}(\boldsymbol{\theta}); \mathcal{J})$ , can be written as a product of time–homogeneous CTMC densities over the inter–event intervals  $\mathcal{J}_1, \dots, \mathcal{J}_M$ . Thus,

$$\pi(\mathbf{X}_j | \Lambda^{(-j)}; \mathcal{J}) = \Pr(\mathbf{X}_j(t_1) | \mathbf{p}_{t_1}) \prod_{m=1}^M \pi(\mathbf{X}_j | \mathbf{x}_j(\tau_{m-1}), \Lambda_m^{(-j)}(\boldsymbol{\theta}); \mathcal{J}_m). \quad (7)$$

Similarly, the transition probability matrix over an interval  $\mathcal{J}_\ell = [t_{\ell-1}, t_\ell]$  can be written as the product of transition probability matrices over the sub–intervals in  $\mathcal{J}_\ell$ , within which the subject–level CTMC is time–homogeneous. Thus, the transition probability matrix over an inter–observation interval,  $\mathcal{J}_\ell = [t_{\ell-1}, t_\ell]$ , partitioned by  $S$  transition events that define inter–event intervals with endpoints given by times

$t_{\ell-1} \equiv \tau_{\ell,0}^{(-j)} < \tau_{\ell,1}^{(-j)} < \dots < \tau_{\ell,S-1}^{(-j)} < \tau_{\ell,S}^{(-j)} \equiv t_\ell$ , is constructed as

$$\mathbf{P}^{(j)}(t_{\ell-1}, t_\ell) = \prod_{s=1}^S \mathbf{P}^{(j)}(\tau_{\ell,s-1}^{(-j)}, \tau_{\ell,s}^{(-j)}).$$

The MCMC algorithm for constructing a subject–path proposal proceeds in three steps (Figure 2):

1. *H4MM step*: sample the disease state of the subject under consideration at the observation times, conditional on the data and disease histories of other subjects.
2. *Discrete time skeleton step*: sample the state at times when the time–inhomogeneous CTMC rates change, conditional on the states sampled in the HMM step.
3. *Event time step*: sample the exact times of transition events conditional on the sequence of states sampled in the previous steps.

**2.3.1 HMM step**—The key to sampling a sequence of disease states at the observation times is to rewrite the emission probability, given by (1), as

$$Y_\ell | X_j(t_\ell), I_{t_\ell}^{(-j)}, \rho \sim \text{Binomial}(\mathbb{1}(X_j(t_\ell) = I) + I_{t_\ell}^{(-j)}, \rho). \quad (8)$$

The emission probability in (8) only depends on whether subject  $j$  is infected at time  $t_\ell$  since we treat the paths of all other subjects, and the parameters, as fixed. Furthermore, the data are conditionally independent of one another, given  $\mathbf{x}$  and  $\boldsymbol{\theta}$  which induces a hidden Markov model (HMM) over the joint distribution  $\mathbf{X}$  and  $\mathbf{Y}$  (Figure 1b).



We sample the discrete path of  $\mathbf{X}_j$  at times  $t_1, \dots, t_L$  from the conditional distribution of  $\mathbf{X}_j$ , denoted  $\pi(\mathbf{X}_j | \mathbf{Y}, \mathbf{x}_{(-j)}, \theta; t_1, \dots, t_L)$ , using the stochastic forward–backward algorithm (Scott, 2002). The algorithm efficiently computes the conditional probabilities of the paths that  $\mathbf{X}_j$  can take through  $S_j$  in the forward recursion. A discrete path is then sampled in the backward recursion. We provide details about the HMM sampling step in Supplementary Material Section S3.

**2.3.2 Discrete-time skeleton step**—It would be straightforward to sample the exact infection and recovery times of subject  $j$ , conditional on the sequence of states at times  $t_1, \dots, t_L$ , if the subject–level CTMC rates did not possibly vary over each inter–observation interval. We may reduce our problem to the time–homogeneous case by first sampling the disease state at the intermediate event times when the CTMC rates change, and then sampling the full path within each inter–event interval. Consider an inter–observation interval,  $\mathcal{I}_\ell = [t_{\ell-1}, t_\ell]$ , containing inter–event intervals whose endpoints are given by times  $t_{\ell-1} \equiv \tau_{\ell,0}^{(-j)} < \tau_{\ell,1}^{(-j)} < \dots < \tau_{\ell,n-1}^{(-j)} < \tau_{\ell,n}^{(-j)} \equiv t_\ell$ . Let  $\tilde{\tau}_i = \tau_{\ell,i}^{(-j)}$  and  $x_i = \mathbf{x}_j(\tau_{\ell,i}^{(-j)})$ . We recursively sample  $\mathbf{X}_j$  at each intermediate event time, beginning at  $\tilde{\tau}_1$ , from the discrete distribution with masses

$$\begin{aligned} & \Pr(\mathbf{X}_j(\tilde{\tau}_i) = x_i | \mathbf{X}_j(\tilde{\tau}_{i-1}) = x_{i-1}, \mathbf{X}_j(\tilde{\tau}_n) = x_n) \\ &= \frac{\Pr(\mathbf{X}_j(\tilde{\tau}_i) = x_i, \mathbf{X}_j(\tilde{\tau}_{i-1}) = x_{i-1}, \mathbf{X}_j(\tilde{\tau}_n) = x_n)}{\Pr(\mathbf{X}_j(\tilde{\tau}_{i-1}) = x_{i-1}, \mathbf{X}_j(\tilde{\tau}_n) = x_n)} \\ &= \frac{\Pr(\mathbf{X}_j(\tilde{\tau}_i) = x_i | \mathbf{X}_j(\tilde{\tau}_{i-1}) = x_{i-1}) \Pr(\mathbf{X}_j(\tilde{\tau}_n) = x_n | \mathbf{X}_j(\tilde{\tau}_i) = x_i)}{\Pr(\mathbf{X}_j(\tilde{\tau}_n) = x_n | \mathbf{X}_j(\tilde{\tau}_{i-1}) = x_{i-1})} \\ &= \frac{[\mathbf{P}^{(j)}(\tilde{\tau}_{i-1}, \tilde{\tau}_i)]_{x_{i-1}, x_i} [\prod_{k=i}^{n-1} \mathbf{P}^{(j)}(\tilde{\tau}_k, \tilde{\tau}_{k-1})]_{x_i, x_n}}{[\prod_{k=i-1}^{n-1} \mathbf{P}^{(j)}(\tilde{\tau}_k, \tilde{\tau}_{k+1})]_{x_{i-1}, x_n}}. \end{aligned} \tag{9}$$

**2.3.3 Event time step**—The final step in constructing a subject–path is to sample the exact infection and recovery times given the discrete sequence of states obtained in the previous two steps. This amounts to simulating the path of an endpoint–conditioned time–homogeneous CTMC, a task for which there exist a variety of efficient methods (Hobolth and Stone, 2009). When fitting the SIR model, we chose to use modified rejection sampling, a modification of Gillespie’s direct algorithm (Gillespie, 1976) that explicitly avoids simulating constant paths. This method is known to be efficient when the states differ at the endpoints of small time intervals. We used uniformization–based sampling (Hobolth and Stone, 2009) when fitting SEIR and SIRS models, which was more robust when sampling paths in intervals with multiple transitions. Fast implementations of these methods are



available in the ECctmc package in R (Fintzi, 2016). We briefly summarize the algorithms in Section S4 of Supplementary Materials.

**2.3.4 Metropolis–Hastings step**—Having constructed a complete subject–path proposal, we decide whether to accept or reject the proposal via a Metropolis–Hastings step. It is important to understand that the true distribution of  $\mathbf{X}_j | \mathbf{x}_{(-j)}, \theta$  is neither Markovian nor analytically tractable, and therefore, does not match the time– inhomogeneous CTMC in our proposal. Suppressing the dependence on  $\theta$ , the target distribution of the subject–path proposal is  $\pi(\mathbf{X} | \mathbf{Y}) \propto \pi(\mathbf{Y} | \mathbf{X})\pi(\mathbf{X})$ . Thus, we accept a proposed subject–path with probability

$$a_{\mathbf{x}^{\text{cur}} \rightarrow \mathbf{x}^{\text{new}}} = \min \left\{ \frac{\pi(\mathbf{x}^{\text{new}} | \mathbf{y}) q(\mathbf{x}^{\text{cur}} | \mathbf{x}^{\text{new}}, \mathbf{y})}{\pi(\mathbf{x}^{\text{cur}} | \mathbf{y}) q(\mathbf{x}^{\text{new}} | \mathbf{x}^{\text{cur}}, \mathbf{y})}, 1 \right\} \quad (10)$$

$$= \min \left\{ \frac{\pi(\mathbf{x}^{\text{new}}) \pi(\mathbf{x}_j^{\text{cur}} | \Lambda^{(-j)}, \mathcal{F})}{\pi(\mathbf{x}^{\text{cur}}) \pi(\mathbf{x}_j^{\text{new}} | \Lambda^{(-j)}, \mathcal{F})}, 1 \right\}.$$

Hence, the Metropolis–Hastings ratio is equal to the ratio of population-level time–homogeneous CTMC densities, multiplied by the ratio of time–inhomogeneous CTMC proposal densities (see Supplementary Material Section S5 for the derivation).

**2.3.5 Initializing the collection of subject–paths**—We initialize the collection of subject paths at the start of our MCMC by simulating paths using Gillespie’s direct algorithm (Gillespie, 1976) until we have found one under which the data have non–zero probability. A sufficient condition for this under the binomial sampling model is that the number of infected individuals is greater than the observed prevalence at each observation time.

## 2.4 Parameter updates

One MCMC iteration includes a number of subject–path updates, followed by a set of parameter updates. The optimal number of subject–path updates per MCMC iteration is specific to the dynamics of the SEM and the epidemic setting (e.g., endemic vs. epidemic, high vs. low escape probability), but ultimately boils down to the cost of subject–path updates vis–a–vis parameter updates. We discuss this further in Section S7 of Supplementary Materials. In the case of the SIR model, as well as the other models we will fit in subsequent sections, conjugate priors are available for all our model parameters. Thus, we use Gibbs sampling to draw new parameter values from their univariate full conditional distributions (see Supplementary Material Section S8).

## 2.5 Implementation

We provide the R and C++ code base for this paper, along with examples and the code for reproducing the results we present in the following sections, in the form of an R package in a stable GitHub repository (<https://github.com/fintzij/BDAepimodel>). Future implementations, including ex-tensions to the algorithm presented in this paper, along with improvements to

the implementation, will be incorporated into the stemr package (<https://github.com/fintzij/stemr>).

### 3 Simulation Results

#### 3.1 Inference under various epidemic dynamics

We fit SIR, SEIR, and SIRS dynamics to binomially distributed prevalence counts sampled from epidemics simulated under corresponding dynamics in populations of 750, 500, and 200 individuals (details provided in Supplementary Material Section S9). Priors for the rate parameters and binomial sampling probability were chosen so that the priors spanned reasonable ranges of values (e.g. recovery durations ranging from days to weeks/months rather than seconds to eons under extremely diffuse priors), but were otherwise only mildly informative, while the initial distribution parameters were assigned informative priors (see Supplementary Tables S4, S6, and S8). The three datasets, depicted in Figure 3 along with the estimated pointwise posterior prevalence, presented a range of challenges. The SIR example was arguably the most “standard” example as the observation period captured the exponential growth and decline of the epidemic. Thus, much of the curvature in the latent path was reflected in the data. In contrast, data from the outbreak simulated under near-endemic SEIR dynamics contained very little information about the shape of the epidemic curve. The task of disentangling whether the data were sampled with low probability from a high-prevalence outbreak, or visa-versa, was further complicated by the inclusion of an additional disease state — the exposed state — that was not directly observed. Finally, the SIRS model was more computationally challenging for two reasons. First, the recurrent nature of the disease process demanded that the disease state at each event time, and the path within each inter-event interval, be sampled in the subject-path proposal. Second, it was possible for CTMC rate matrices to have complex eigen-decompositions, which made computing transition probability matrices more expensive. This affected the optimal number of subject-path updates per MCMC iteration (see Supplementary Material Section S7 for further discussion of this point). Simulation details, along with minor adaptations to our algorithm for fitting the SEIR and SIRS models, are presented in Supplementary Material Section S6.

The true epidemic paths and parameter values fell well within the 95% Bayesian credible intervals in all three simulations (Figure 3 presents the estimated latent posterior prevalence; Figure 4 presents posterior estimates of model parameters; Supplementary Material Figure S12 presents estimated latent posterior distributions and true epidemic paths for all model compartments). The acceptance rates for subject-path proposals were roughly 92% for the SIR model, 91% for the SEIR model, and 77% for the SIRS model. Our posterior estimates of the model parameters also closely match estimates obtained using the particle marginal Metropolis-Hastings (PMMH) algorithm of Andrieu et al. (2010), implemented using the pomp package in R (King et al., 2016). We simulated particle paths in the PMMH algorithm in two ways; exactly using Gillespie’s direct algorithm (Gillespie, 1976), and approximately using a multinomial modification of  $\tau$ -leaping (Bretó and Ionides, 2011). In these small population examples, the exact algorithm is arguably more appropriate, as the leap conditions for  $\tau$ -leaping may not be met in small populations, but it is also substantially

slower. In these simple settings, PMMH tended to outperform our algorithm in terms of log-posterior effective sample size (ESS) per CPU time. When PMMH particle paths were simulated by  $\tau$ -leaping, the average ESS per CPU compared to BDA was roughly 350 $\times$  greater for the SIR model, 4.4 $\times$  greater for the SEIR model, and 13 $\times$  greater for the SIRS model. Exact simulation of PMMH particle paths reduced the computational advantage of PMMH substantially. In this case, the average log-posterior ESS per CPU time was 10.5 $\times$  greater for PMMH in fitting the SIR model, 2 $\times$  for the SEIR model, and 0.7 $\times$  for the SIRS model. These comparisons did not include the time required to tune the MCMC for PMMH, which was nontrivial. In contrast, our algorithm required no tuning beyond selecting the number of subject-paths to update per MCMC iteration. We also note that in fitting the models using PMMH, we were required to make several implementation decisions to prevent particle degeneracy and to balance speed with precision. These included selecting the number of particles and the time-step in the approximate  $\tau$ -leaping algorithm. For example, when using  $\tau$ -leaping to simulate particle paths, the number of particles required to obtain good mixing for the SIRS model fit with PMMH was much higher than for the other two models. Details of the PMMH implementations and further results are discussed in Supplementary Material Section S9.

### 3.2 Inference under model misspecification

In practice, every stochastic epidemic model is misspecified with respect to the real world epidemic process from which the data arise, and the malignancy of the model misspecification is often impossible to diagnose a priori. We can build up an understanding of an epidemic's dynamics by fitting SEMS under a range of dynamics, beginning with simple, easily interpretable models. The results of each model are interpreted counterfactually — e.g. “If the true epidemic followed SIR dynamics, our best guess of the dynamics that gave rise to the data would be...”. The iterative nature of epidemic modeling suggests that some minimal criteria for the usefulness of any computational algorithm would be that, for a reasonable model, the MCMC should converge to the posterior of the model parameters, and that the estimated latent posterior distribution under the hypothetical dynamics should reflect the true epidemic.

However, it is precisely the inherent misspecification of SEMs that leads simulation-based methods to struggle in many instances, and it is here that we highlight a critical advantage of our DA algorithm. Our subject-path proposals are driven, not just by the SEM dynamics, but also by the data. This enables us to overcome model misspecification in situations in which simulation-based methods degenerate due to their reliance on an adequately accurate model for simulating epidemic paths. We demonstrate this in a simple example in which we fit SIR and SEIR models to four years of weekly prevalence data sampled from an epidemic simulated under time-varying SEIR dynamics, where the latent period, infectious period, and per-contact infectivity rate were modulated over four discrete epochs (depicted in Figure 5, details presented in Supplementary Material Section S10).

We fit SIR and SEIR models to the data using our DA algorithm, and using PMMH with 2,500 particles, the paths for which were simulated approximately via  $\tau$ -leaping with a time-step of 1 day. We assigned weakly informative priors for the rate parameters governing

the epidemic dynamics in both models, and informative priors for the binomial sampling probability and the initial state probabilities (Supplementary Material Table S11). The MCMC chains for models fit with PMMH suffered from severe particle degeneracy and did not converge (see Supplementary Material Figures S13 and S15).

Both models fit via DA yield reasonable estimates for the within–subject disease dynamics (i.e. the infectious period, as well as the latent period in the case of the SEIR model). The posterior median average infectious period duration was estimated to be 292 days (95% BCI: 263 days, 323 days) under SIR dynamics, and 287 days (95% BCI: 260 days, 318 days) under SEIR dynamics. The posterior median average latent period under SEIR dynamics was 211 days (95% BCI: 165 days, 260 days). The posterior median estimate of  $R_0$  under SIR dynamics was 4.05 (95% BCI: 3.40, 4.81), while under SEIR dynamics, the posterior median estimate of  $R_0$  was 23.8 (95% BCI: 15.1, 37.0). While the true prevalence fell well within the pointwise 95% credible interval for both models (Figure 6), we notice that the degree of model misspecification drastically affected our ability to estimate the history of the numbers of noninfectious people over the course of the epidemic. Under SIR dynamics, we drastically overestimate the number of susceptible individuals. The SEIR model much more closely resembles the time–varying SEIR model used to simulate the epidemic. Although the true path for the number of susceptible still falls outside the 95% credible interval at times, we are still able to reconstruct a reasonable range of paths for the number of exposed individuals. This contrasts with the models fit in Section 3.1, which were not misspecified with respect to the true epidemic dynamics. In that case, the complete path of the epidemic fell well within the estimated credible intervals for all disease states for all three models (Supplementary Material Figure S12). Therefore, we advise caution in reconstructing the epidemic history for disease states that were not measured, particularly when severe model misspecification is suspected.

### 3.3 Inference under population size misspecification

Model misspecification often extends not only to the SEM dynamics, but also to the assumed population size. This is often the case in settings where subject–level data is unavailable, for example, in resource limited settings or surveillance settings, and may result in biased estimates of the SEM dynamics. This bias is the result of a mismatch between the intensive dynamics of the epidemic process, which are a function of the fractions of people in the population in each disease state, and the extensive scale of prevalence counts, which are not normalized by the population size. Without knowing the true population size, it is difficult to know whether the scale of the observed counts reflects a high prevalence/low detection rate setting, or visa–versa. Moreover, wrongly assuming too large, or too small, of a population size could bias posterior inference of the epidemic dynamics.

We simulated weekly prevalence counts under a binomial measurement process with detection probability  $\rho = 0.3$  from an epidemic with SIR dynamics in a population of  $N = 1,250$  individuals. We then fit SIR models using a series of assumed population sizes under a flat prior for the binomial sampling probability and diffuse priors for the epidemic dynamics (see Supplementary Material Section S11 for complete simulation details and prior specifications), and compared the resulting scaled parameter estimates. The per–contact

infectivity rate,  $\beta$ , was rescaled by the population size,  $N$ , so that it could be interpreted as the rate of disease transmission. We computed  $R_0$  using the assumed population size. Finally, we scaled the binomial sampling probability by the assumed population size to give the expected number of observed infections in a completely infected population.

We are able to obtain approximately valid inference under moderate misspecification of the population size. However, estimates of the epidemic dynamics and the case detection probability become severely biased as the magnitude of the misspecification increases. Furthermore, the widths of the credible intervals for the model parameters shrink as misspecification of the population size becomes more severe. The constrained ranges of model dynamics also manifest in a narrowing of the widths of the pointwise credible intervals for disease prevalence (Figure 8). Under severe misspecification of the population size ( $N = 150$ ), the latent posterior distribution has 95% of its mass within only a narrow band of epidemic paths. In contrast, under moderate misspecification of the population size, the widths of the latent posterior credible intervals are quite similar to the estimated range using the true population size.

There are two final points that we wish to make based on this simulation. The first is that it might be possible to deliberately misspecify the true population size in order to speed up computation time and still obtain approximately valid inference. The average run time using the true population size of 1250 individuals was roughly  $2\times$  and  $7\times$  longer than the average run times in populations of 900 and 500 individuals. Yet, posterior inferences about the epidemic dynamics were not substantially affected. Longer run times in large populations result from having to sample more subject–paths per MCMC iteration at a relatively higher cost per subject–path. The second point is that in situations where the true population size is unknown, SEM likelihood–based inference has some robustness to misspecification of the population size, at least in a neighborhood of population sizes around the true number of individuals. Thus, comparing posterior inferences under a range of population sizes could be a useful heuristic diagnostic for population size misspecification.

### 3.4 Effect of prior specification on posterior inference

Given the relatively limited extent of aggregated prevalence counts compared to a setting in which subject–level data are available, we must consider how our choices of prior distributions influence our posterior inferences. We simulated an outbreak with SIR dynamics in a population of 750 individuals for which  $R_0 = \beta \times 763 / \mu \approx 1.84$  and the mean infectious period was  $1/\mu = 7$  days. We fit SIR models to binomially distributed weekly prevalence data, sampled with detection probability  $\rho = 0.2$ , under the following four prior regimes: Regime 1 — informative priors for all model parameters; Regime 2 — vague priors for the rate parameters and an informative prior for the sampling probability; Regime 3 — informative priors for the rate parameters and a flat prior for the sampling probability; Regime 4 — vague priors for the rate parameters and a flat prior for the sampling probability. The same prior for the initial state probabilities was used in all four regimes. Complete simulation details and convergence diagnostics are supplied in Section S12.

The true values for all model parameters fell within the 95% credible intervals under all four prior regimes. Unsurprisingly, informative priors tended to result in narrower credible

intervals for the parameters (Figure 9) as well as for the latent process (Figure 10). The strength of prior information about the sampling probability affected the widths of credible intervals to a much greater extent than the priors for the rate parameters. Strong prior information about the sampling probability also resulted in substantially narrower credible intervals for disease prevalence under each of the prior regimes for the rate parameters. In contrast, informative priors for the rate parameters yielded only slightly narrower credible intervals for disease prevalence when holding constant the strength of the sampling probability prior. The effects on the initial state probability parameters seem to reverse this pattern, although we caution against overinterpretation given the paucity of data available for estimating those parameters. MCMC chains with strong priors for the binomial sampling probability also appeared to mix somewhat better than chains with diffuse priors for the sampling probability (see traceplots in Supplementary Material Section S12).

#### 4 Influenza in a British boarding school

As an application, we analyze data from an outbreak of influenza in a British boarding school (Anon., 1978, Davies et al., 1982). This outbreak took place shortly after the Easter term began in January 1978, and was estimated to eventually infect roughly 90% of the 763 boys aged 10–18. Daily counts of the boys who were confined to the infirmary from January 22<sup>nd</sup> through February 4<sup>th</sup> were accessed via the `pomp` package in R (King et al., 2016), and are displayed in Figure 11.

We used our DA algorithm and PMMH to fit SIR and SEIR models with a binomial emission distribution to the data (see Supplementary Material Section S13 of the supplement for complete details). All of the parameters were assigned diffuse priors, which are plotted over the posterior ranges in Figure 12. The PMMH algorithm failed to converge for both models, which we suspect was due to a combination of model misspecification and the constrained state space of the binomial measurement process. We also fit a set of supplementary SIR and SEIR models in Section S13.2, in which we assumed a negative–binomial emission distribution. This was done in order to facilitate comparison with PMMH, although we feel that a negative binomial emission distribution is not appropriate in such a closely monitored outbreak setting since it does not rule out over–reporting of cases.

Together, the SIR and SEIR models suggest that cases were detected with high probability and that the outbreak, though aggressive, was not atypical given the closed environment in which it occurred. The posterior median estimates of the detection probability, roughly 0.98 for both models (SIR 95% BCI: 0.92, 1.00; SEIR 95% BCI: 0.91, 1.00), suggested that while almost all of the infectious boys were detected, a handful of cases went unnoticed. The posterior median recovery rate under SIR dynamics corresponds to an average period of 2.16 days (95% BCI: 1.99, 2.37) during which an infectious boy could transmit an infection to other boys before being confined to the infirmary. Under SEIR dynamics, the posterior median average infectious period was 2.12 days (95% BCI: 1.95, 2.33), and the posterior median average latent period was 1.19 days (95% BCI: 0.84, 1.51). These results are consistent with the typical progression of influenza, in which individuals typically incubate for between one to four days before symptoms manifest, and are typically infectious for one day before, and up to a week after, symptom onset (Centers for Disease Control and



Prevention, 2014). The posterior median estimates of  $R_0$  were 3.89 (95% BCI: 3.40, 4.47) under SIR dynamics, and 10.38 (95% BCI: 7.40, 14.11) under SEIR dynamics. Previous analyses of this dataset with trajectory matching estimate  $R_0$  to be roughly 3.7 for the SIR model and 35.9 for the SEIR model (Wearing et al., 2005, Keeling and Rohani, 2008), though we note that these estimates are based on deterministic models that do not properly account for distributional properties of the data. Our results for both models are also in agreement with estimates of SIR and SEIR model dynamics under a negative binomial emission distribution (see Section S13.2).

## 5 Conclusion

We have presented an agent-based Bayesian DA algorithm for fitting SEMs to disease prevalence time series counts. This was previously difficult, if not computationally infeasible, to carry out using traditional agent-based DA methods in the absence of subject-level data. Although we outlined the algorithm in the context of fitting an SIR model to binomially distributed prevalence data, our algorithm represents a general solution for fitting SEMs to prevalence counts. In simulations and the applied example, we fit SEIR and SIRS models to prevalence data, and in the supplement also fit SIR and SEIR models with a negative binomial emission distribution to the British boarding school data. We have demonstrated that our algorithm yields approximately valid inference when the population size is misspecified. Moreover, our algorithm is usable in settings in which simulation-based methods, such as PMMH, break down due to misspecification of the SEM. Finally, our DA algorithm is carried out entirely at the subject level, making it possible to also incorporate subject-level covariates and household structure, or to fit models to subject-level data.

There are two fundamental limitations of agent-based DA methods from which our algorithm is not excepted. First, the bookkeeping required to track the collection of subject-paths increases in size and complexity as the number of events grows large. Attempts to fit stochastic epidemic models in large populations using agent-based DA may be thwarted by prohibitive computational overhead. MCMC run times using our implementation, which was coded for reliability rather than speed, substantially degraded once the assumed population size was greater than a few thousand people. Second, we suspect that MCMC mixing in large populations could eventually become too slow for agent-based DA to be of practical use, even if solutions could be found for the computational bottlenecks. As the population size gets large, perturbations to the likelihood from re-sampling one subject at a time become relatively less significant. For this reason, we view extensions for jointly sampling multiple subject-paths as a critical step in mitigating slow MCMC mixing in large populations.

Finally, we would like to comment on directions for future work that we intend to pursue. The DA algorithm in this paper addresses the problem of fitting SEMs to prevalence data. This type of data summarizes total number of infections in the population at a particular time. However, outbreak data often consist of incidence counts, which are the number of new cases accumulated in each inter-observation interval. Extending our DA algorithm to accommodate incidence data is an important next step and should be straightforward in



situations where the state space for the subject level process is finite — for instance, if a subject cannot become reinfected more than once or twice in a given inter-observation interval. We also believe it is important to investigate whether there is a way to make our DA algorithm more efficient by selecting the subjects whose paths are resampled in each iteration in a way that maximizes the perturbation to the population-level path and does not invalidate the MCMC. Designing an optimal schedule of subject-path updates could be critical to being able to use our algorithm in fitting more complex models to data from epidemics in large, structured populations.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## 6 Acknowledgements

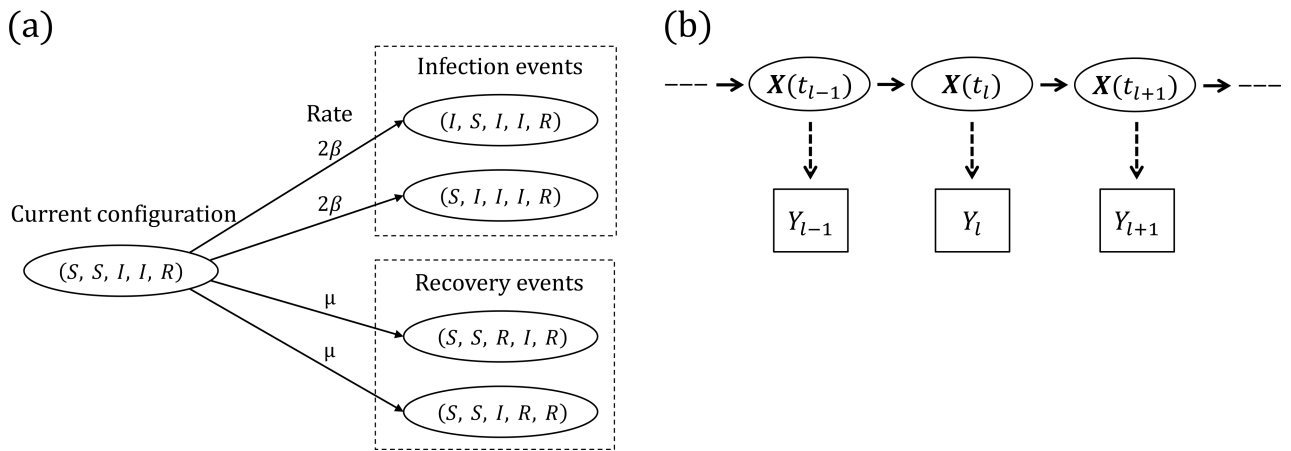
J.F., J.W., and V.N.M. were supported by the NIH grant U54 GM111274. J.W. was supported by the NIH grant R01 CA095994. V.N.M. was supported by the NIH grant R01 AI107034. We would also like to thank Aaron King and the rest of the authors of the pomp package for their help with the PMMH algorithm that served as a benchmark for the methods presented in this paper.

## References

- Allen LJS. An introduction to stochastic epidemic models In *Mathematical Epidemiology*, pages 81–130. Springer, New York, 2008.
- Andersson H and Britton T. *Stochastic Epidemic Models and Their Statistical Analysis Lecture Notes in Statistics*. Springer, New York, 2000.
- Andrieu C, Doucet A, and Holenstein R. Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72:269–342, 2010.
- Anon. Influenza in a boarding school. *The British Medical Journal*, 1:587, 1978.
- Auranen K, Arjas E, Leino T, and Takala AK. Transmission of pneumococcal carriage in families: a latent Markov process model for binary longitudinal data. *Journal of the American Statistical Association*, 95:1044–1053, 2000.
- Becker NG. On a general stochastic epidemic model. *Theoretical Population Biology*, 11:23–36, 1977. [PubMed: 854859]
- Bretó C and Ionides EL. Compound Markov counting processes and their applications to modeling infinitesimally over-dispersed systems. *Stochastic Processes and their Applications*, 121:2571–2591, 2011.
- Britton T. Stochastic epidemic models: a survey. *Mathematical Biosciences*, 225:24–35, 2010. [PubMed: 20102724]
- Cappé O, Moulines E, and Ryden T. *Inference in Hidden Markov Models Springer Series in Statistics*. Springer, New York, 2006.
- Cauchemez S and Ferguson NM. Likelihood-based estimation of continuous-time epidemic models from time-series data: application to measles transmission in London. *Journal of the Royal Society Interface*, 5:885–897, 2008.
- Cauchemez S, Carrat F, Viboud C, and Valleron AJ. A Bayesian MCMC approach to study transmission of influenza: application to household longitudinal data. *Statistics in Medicine*, 23: 3469–3487, 2004. [PubMed: 15505892]
- Centers for Disease Control and Prevention. How flu spreads, 2014 URL <http://www.cdc.gov/flu/about/disease/spread.htm>. Accessed on January 3, 2016.
- Davies JR, Smith AJ, Grilli EA, and Hoskins TW. Christ’s Hospital 1978–79: An account of two outbreaks of influenza A H1N1. *Journal of Infection*, 5:151–156, 1982.

- Dukic V, Lopes HF, and Polson NG. Tracking epidemics with Google flu trends data and a state-space SEIR model. *Journal of the American Statistical Association*, 107:1410–1426, 2012.
- Fintzi J. ECctmc: Simulation from Endpoint-Conditioned Continuous Time Markov Chains, 2016 URL <https://github.com/fintzij/ECctmc>. R package, version 0.2.2.
- Gibson GJ and Renshaw E. Estimating parameters in stochastic compartmental models using Markov chain methods. *Mathematical Medicine and Biology*, 15:19–40, 1998.
- Gillespie DT. A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of Computational Physics*, 22:403–434, 1976.
- Glass K, Xia Y, and Grenfell B. Interpreting time-series analyses for continuous-time biological models — measles as a case study. *Journal of Theoretical Biology*, 223:19–25, 2003. [PubMed: 12782113]
- Held L and Paul M. Modeling seasonality in space-time infectious disease surveillance data. *Biometrical Journal*, 54:824–843, 2012. [PubMed: 23034894]
- Held L, Höhle M, and Hofmann M. A statistical framework for the analysis of multivariate infectious disease surveillance counts. *Statistical modelling*, 5:187–199, 2005.
- Hirsch MW, Smale S, and Devaney RL. *Differential Equations, Dynamical Systems, and an Introduction to Chaos*. Academic Press. Academic Press, Waltham, 2013.
- Hobolth A and Stone EA. Simulation from endpoint-conditioned, continuous-time Markov chains on a finite state space, with applications to molecular evolution. *The Annals of Applied Statistics*, 3:1204–1231, 2009. [PubMed: 20148133]
- Höhle M and Jørgensen E. Estimating parameters for stochastic epidemics. Technical Report 102, The Royal Veterinary and Agricultural University, 11 2002.
- Ionides EL, Bhadra A, Atchadé Y, King AA, et al. Iterated filtering. *The Annals of Statistics*, 39:1776–1802, 2011.
- Jandarov R, Haran M, Bjørnstad O, and Grenfell B. Emulating a gravity model to infer the spatiotemporal dynamics of an infectious disease. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 63:423–444, 2014.
- Keeling MJ and Rohani P. *Modeling Infectious Diseases in Humans and Animals*. Princeton University Press, Princeton, 2008.
- King AA, Nguyen D, and Ionides EL. Statistical inference for partially observed Markov processes via the R package pomp. *Journal of Statistical Software*, 69:1–43, 2016.
- Koepke AA, Longini IM, Jr, Halloran ME, Wakefield J, and Minin VN. Predictive modeling of cholera outbreaks in Bangladesh. *The Annals of Applied Statistics*, 10:575–595, 2016. [PubMed: 27746850]
- Lekone PE and Finkenstädt BF. Statistical inference in a stochastic epidemic SEIR model with control intervention: Ebola as a case study. *Biometrics*, 62:1170–1177, 2006. [PubMed: 17156292]
- Lindenstrand D and Svensson Å. Estimation of the Malthusian parameter in an stochastic epidemic model using martingale methods. *Mathematical Biosciences*, 246:272–279, 2013. [PubMed: 24427788]
- Longini IM, Jr and Koopman JS. Household and community transmission parameters from final distributions of infections in households. *Biometrics*, 38:115–126, 1982. [PubMed: 7082755]
- McKinley T, Cook AR, and Deardon R. Inference in epidemic models without likelihoods. *The International Journal of Biostatistics*, 5:1–40, 2009.
- McKinley TJ, Ross JV, Deardon R, and Cook AR. Simulation-based Bayesian inference for epidemic models. *Computational Statistics & Data Analysis*, 71:434–447, 2014.
- Neal PJ and Roberts GO. Statistical inference and model selection for the 1861 Hagelloch measles epidemic. *Biostatistics*, 5:249–261, 2004. [PubMed: 15054029]
- O’Neill PD. A tutorial introduction to Bayesian inference for stochastic epidemic models using Markov chain Monte Carlo methods. *Mathematical Biosciences*, 180:103–114, 2002. [PubMed: 12387918]
- O’Neill PD. Bayesian inference for stochastic multitype epidemics in structured populations using sample data. *Biostatistics*, 10:779–791, 2009. [PubMed: 19648227]
- O’Neill PD. Introduction and snapshot review: relating infectious disease transmission models to data. *Statistics in Medicine*, 29:2069–2077, 2010. [PubMed: 20809536]

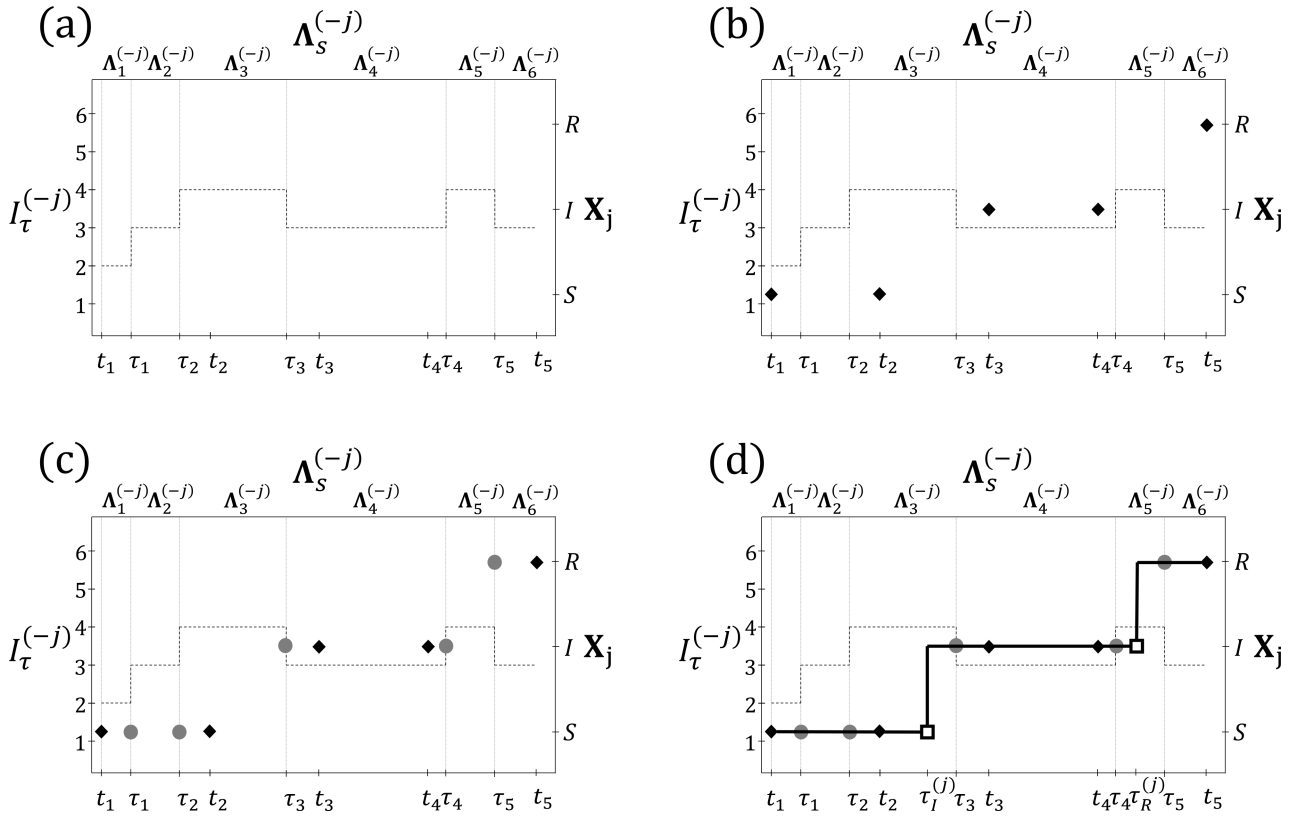
- O'Neill PD and Roberts GO. Bayesian inference for partially observed stochastic epidemics. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 162:121–129, 1999.
- Pooley CM, Bishop SC, and Marion G. Using model-based proposals for fast parameter inference on discrete state space, continuous-time Markov processes. *Journal of The Royal Society Interface*, 12:20150225, 2015.
- Qin Z and Shelton CR. Auxiliary Gibbs sampling for inference in piecewise-constant conditional intensity models. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, 2015.
- Roberts GO and Stramer O. On inference for partially observed nonlinear diffusion models using the Metropolis-Hastings algorithm. *Biometrika*, 88:603–621, 2001.
- Scott SL. Bayesian methods for hidden Markov models: Recursive computing in the 21st century. *Journal of the American Statistical Association*, 97:337–351, 2002.
- Shelton CR and Ciardo G. Tutorial on structured continuous-time Markov processes. *Journal of Artificial Intelligence Research*, 51:725–778, 2014.
- Shestopaloff AY and Neal RM. Sampling latent states for high-dimensional non-linear state space models with the embedded HMM method. arXiv preprint arXiv:1602.06030v2, 2016.
- Sudbury A. The proportion of the population never hearing a rumour. *Journal of Applied Probability*, 22:443–446, 1985.
- Toni T, Welch D, Strelkowa N, Ipsen A, and Stumpf MPH. Approximate Bayesian computation scheme for parameter inference and model selection in dynamical systems. *Journal of the Royal Society Interface*, 6:187–202, 2009.
- Watson R. An application of a martingale central limit theorem to the standard epidemic model. *Stochastic Processes and Their Applications*, 11:79–89, 1981.
- Wearing HJ, Rohani P, and Keeling MJ. Appropriate models for the management of infectious diseases. *PLOS Medicine*, 2:e174, 2005. [PubMed: 16013892]



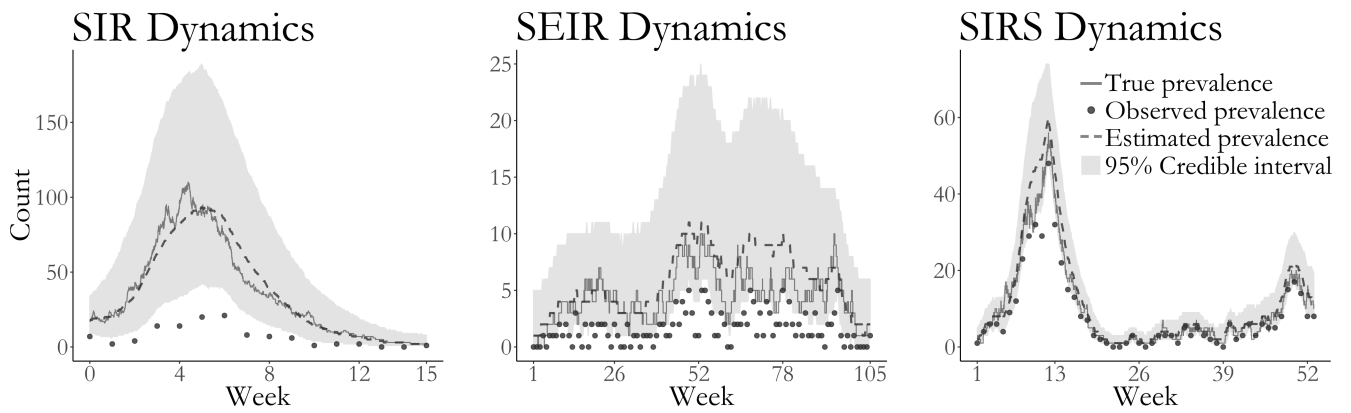
**Figure 1:**

(a) SIR dynamics in a population of five subjects. The number of infecteds can increase from two to three via an infection of the first or second subject, reaching each of those configurations at rate  $2\beta$ . The number of recovered individuals can increase from one to two via a recovery of the third or fourth subject, reaching each of those configurations at rate  $\mu$ .

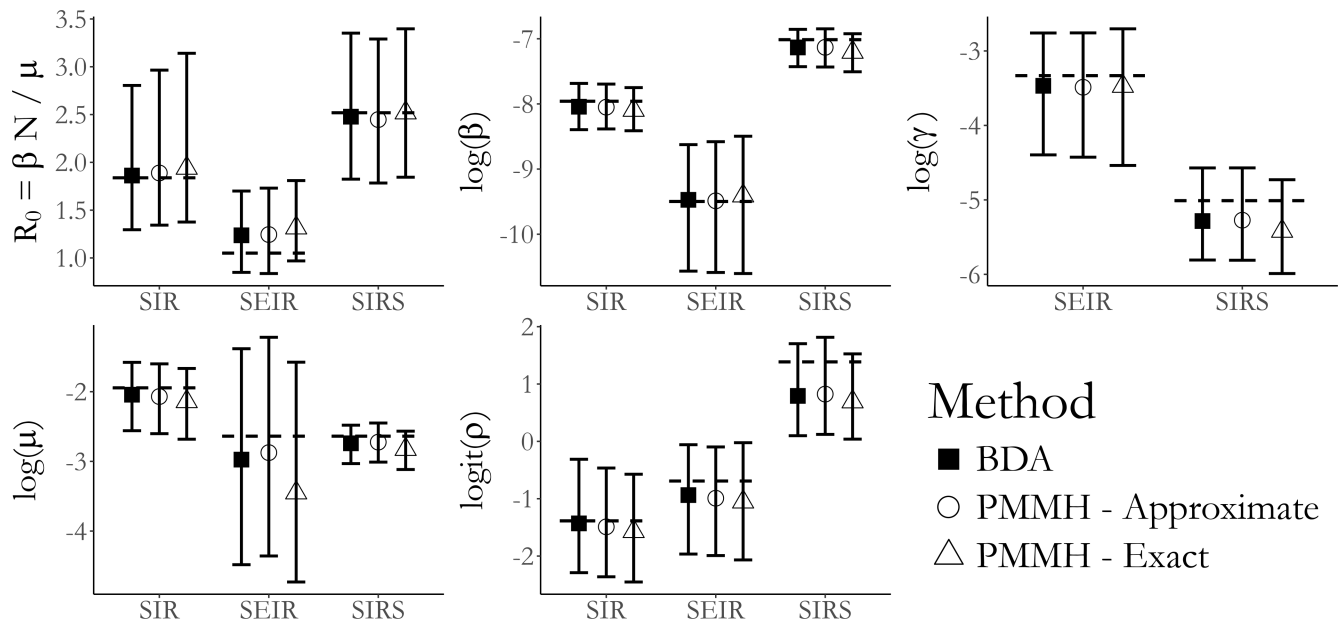
(b) Hidden Markov model for the joint distribution of the latent epidemic process and the data. The observations,  $\mathbf{Y}_\ell \ell=1, \dots, L$ , are conditionally independent given  $\mathbf{X}(t)$ , and  $Y_\ell | I_{t_\ell}, \rho \sim \text{Binomial}(I_{t_\ell}, \rho)$ .



**Figure 2:** Procedure for constructing a subject–path proposal with SIR dynamics. (a) The dashed line depicts the number of infected individuals, excluding  $\mathbf{X}_j$ , the subject whose path is being sampled. The observation times,  $t_1, \dots, t_5$ , and times at which other subjects change disease states,  $\tau_1, \dots, \tau_5$ , are shown on the bottom axis. Rate matrices of the time–inhomogeneous CTMC (top axis) are constant within inter–event intervals (vertical lines). The state space of the subject–level process,  $\mathbf{X}_j$ , is shown on the right axis. (b) *HMM step*: Sample the state of  $\mathbf{X}_j$  at  $t_1, \dots, t_5$ , conditional on the data and on the disease histories of other subjects. (c) *Discrete time skeleton step*: Sample the infection status at  $\tau_1, \dots, \tau_5$ , conditional on the sequence of states sampled in the HMM step. (d) *Event time step*: Sample the infection and recovery times from endpoint–conditioned time–homogeneous CTMC distributions, conditional on the sequence of disease states sampled in the HMM and discrete time skeleton steps.

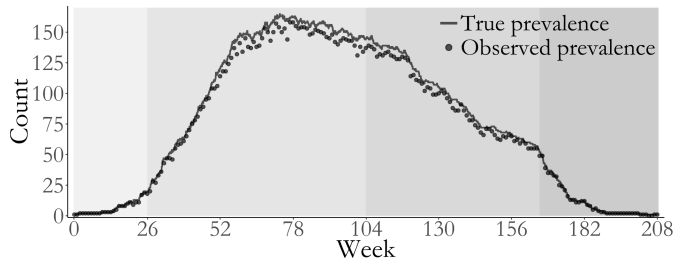


**Figure 3:** Estimated latent posterior distributions of disease prevalence in outbreaks simulated under SIR (left), SEIR (middle), and SIRS (right) dynamics. Depicted are the true unobserved prevalence (solid line), observed data (dots), pointwise posterior median prevalence (dashed line), and pointwise 95% credible intervals (shaded region). Latent posterior estimates are based on a thinned sample, with every 250<sup>th</sup> sample retained.



**Figure 4:** Posterior medians and 95% credible intervals of parameters in the SIR, SEIR, and SIRS models fit with Bayesian data augmentation (BDA) and particle marginal Metropolis–Hastings (PMMH) with particle paths simulated approximately (using  $\tau$ -leaping) and exactly (using Gillespie’s direct algorithm). Displayed are estimates of the basic reproductive number,  $R_0$ , the rate parameters, and the binomial sampling probability. In all models,  $\beta$  is the per-contact infectivity rate,  $\mu$  is the recovery rate, and  $\rho$  is the binomial sampling probability. In the SEIR model,  $\gamma$  denotes the rate at which an exposed individual becomes infectious, while in the SIRS model  $\gamma$  denotes the rate at which immunity is lost.





Parameter	Epoch			
	1	2	3	4
$R_0^{\text{Eff}}$	14.9	9.2	0.1	0
$1/\gamma$ (days)	210	210	90	180
$1/\mu$ (days)	150	330	300	70

**Figure 5 & Table 1:**

Simulated outbreak with SEIR dynamics that varied over four epochs (shaded regions).

Weekly prevalence counts (points) were binomially sampled with sampling probability  $\rho$

“ 0.95 from the true unobserved prevalence (solid line). The table presents the effective reproductive number computed based on the number of susceptibles at the beginning of each epoch,

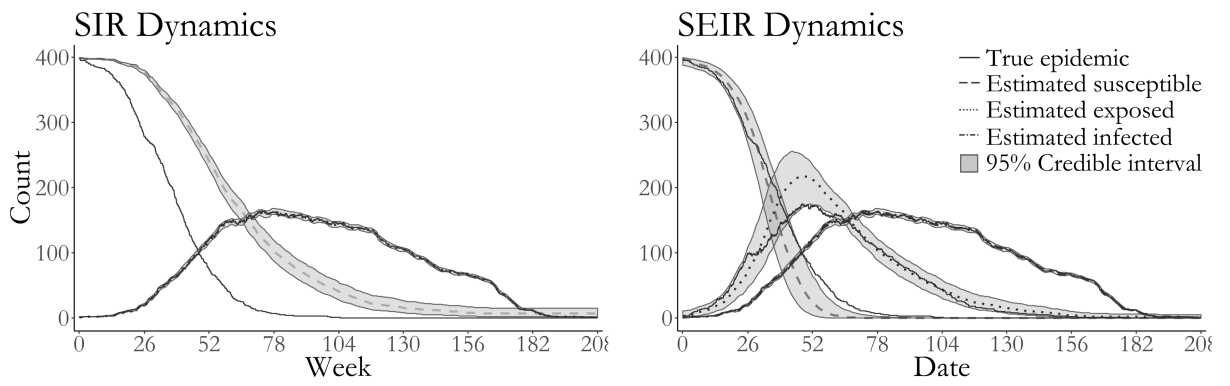
$R_0^{\text{Eff}} = \beta(\tau)S(\tau)/\mu(\tau)$ , the mean latent period,  $1/\gamma$ , and the mean infectious period,  $1/\mu$ .

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript



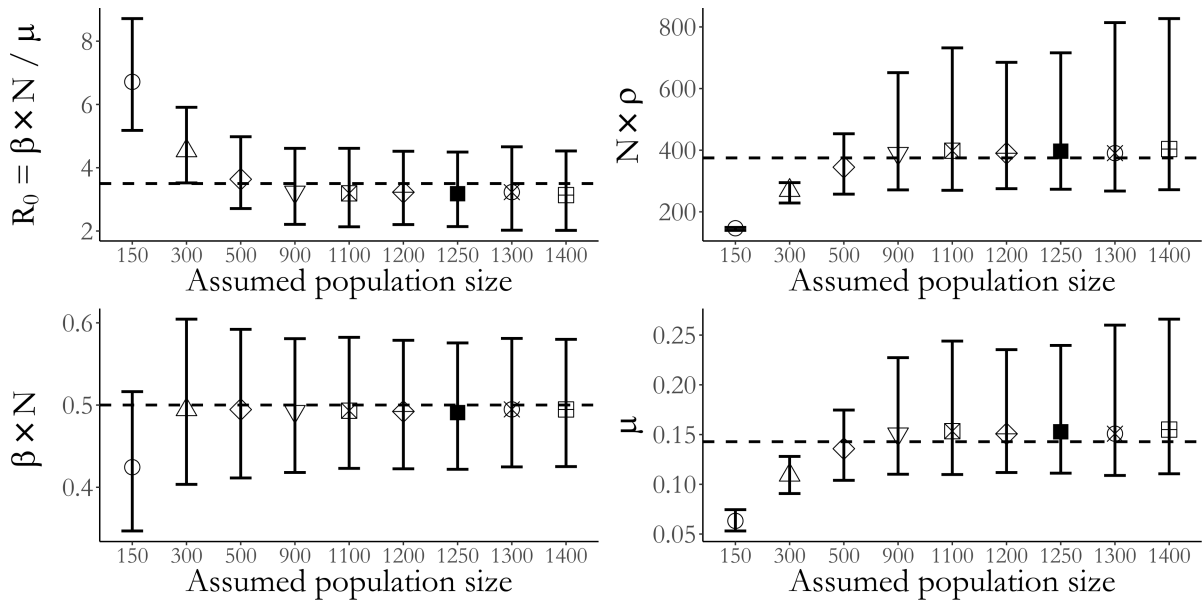
**Figure 6:** True epidemic path (solid lines), pointwise posterior median estimate of the numbers of susceptibles (dashed line), exposed (dotted line), and infected individuals (dash–dotted line) and pointwise 95% credible intervals (shaded regions) under SIR and SEIR dynamics.

Author Manuscript

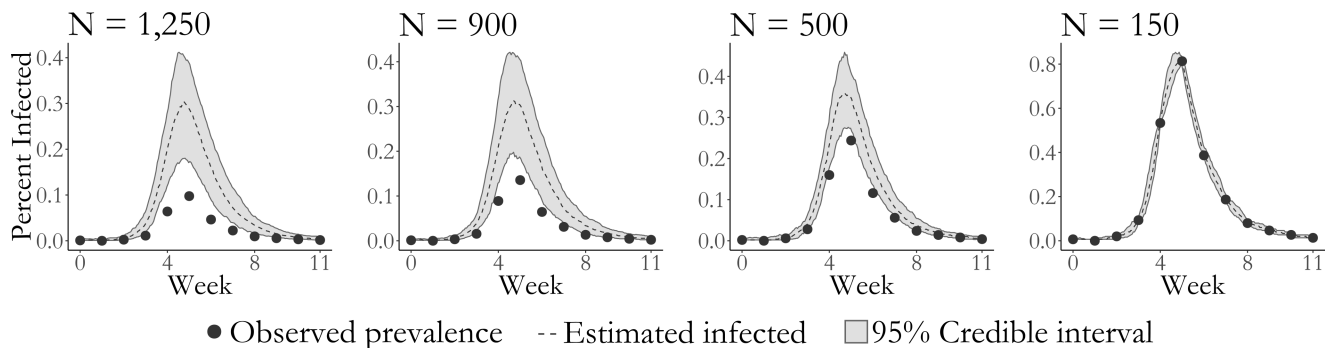
Author Manuscript

Author Manuscript

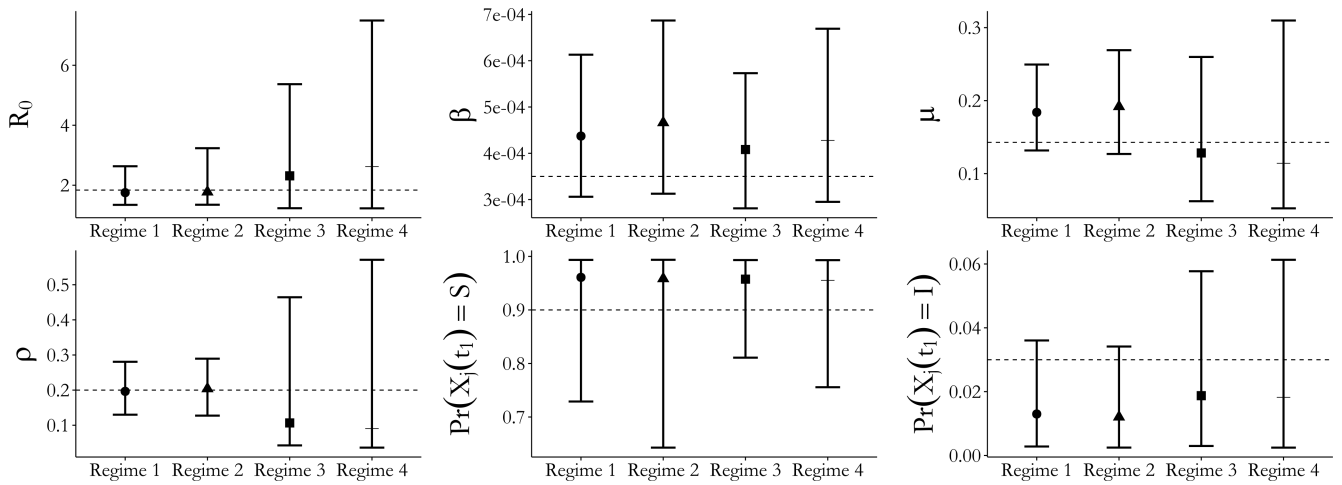
Author Manuscript



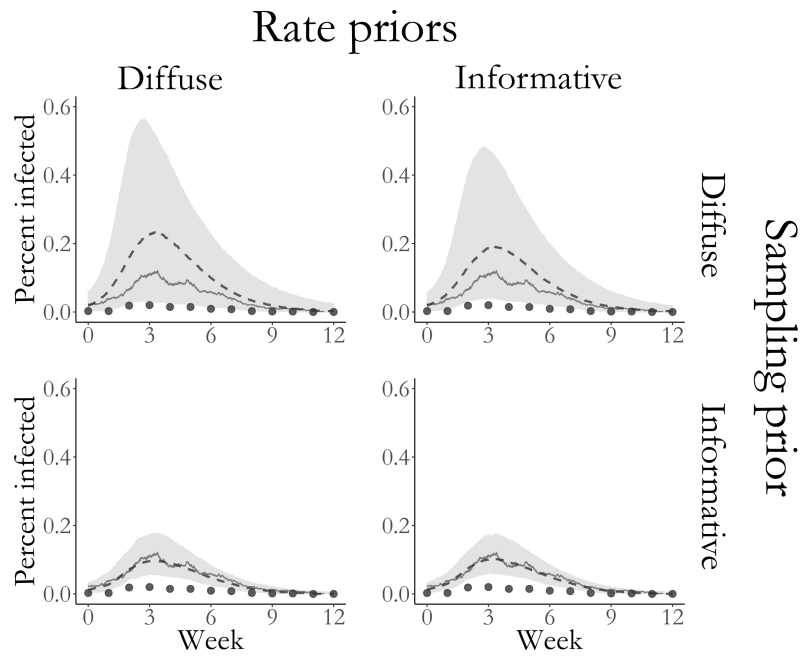
**Figure 7:** Posterior medians and 95% credible intervals for the basic reproductive number,  $R_0$ , infectivity rate, recovery rate, and binomial sampling probability scaled by the assumed population size. The dashed lines indicate the true values in the population of size 1,250. The population size,  $N$ , indicates the assumed population size used in fitting the model.

**Figure 8:**

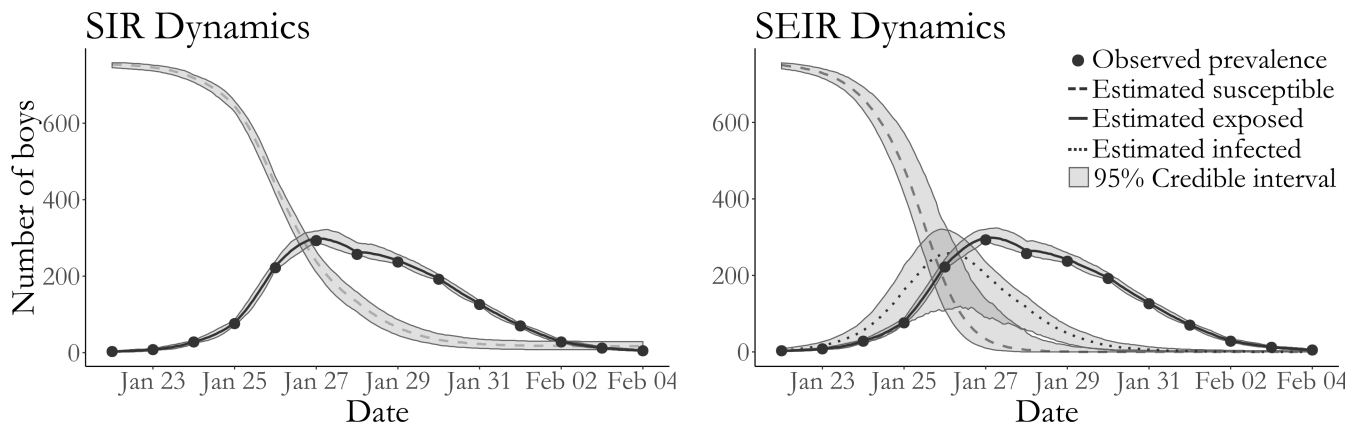
Estimated latent posterior distributions of disease prevalence under SIR dynamics. The true population size is 1,250. Depicted are the observed prevalence (dots), pointwise posterior median prevalence (dashed line), and pointwise 95% credible intervals (shaded region) all scaled by the assumed population size. Latent posterior estimates are based on a thinned sample, with every 250<sup>th</sup> sample retained.



**Figure 9:** Posterior median estimates and 95% credible intervals for all SIR model parameters under four different prior regimes (Table S14). Regimes 1 and 3 set informative priors for the per-contact infectivity and recovery rates. Regimes 1 and 2 set informative priors for the binomial sampling probability. The same mildly informative prior for the initial state probabilities was used in all four regimes.



**Figure 10:** Estimated latent posterior distributions of disease prevalence in outbreaks simulated under four prior regimes for SIR model rate parameters and the binomial sampling probability. Depicted are the true unobserved prevalence (solid line), observed data (dots), pointwise posterior median prevalence (dashed line), and pointwise 95% credible intervals (shaded region). Latent posterior estimates are based on a thinned sample, with every 250<sup>th</sup> sample retained.



**Figure 11:** Boarding school data, pointwise posterior median estimates and pointwise 95% credible intervals (grey shaded areas) under SIR and SEIR dynamics of the numbers of susceptible boys (dashed line), exposed boys (dotted line), and infected boys (solid line). Posterior estimates based on a thinned sample, with every 250<sup>th</sup> configuration retained.

Author Manuscript

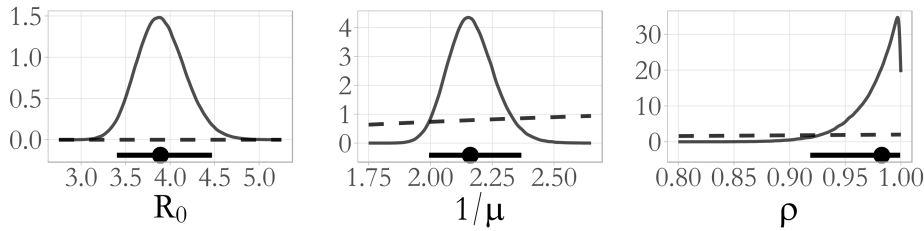
Author Manuscript

Author Manuscript

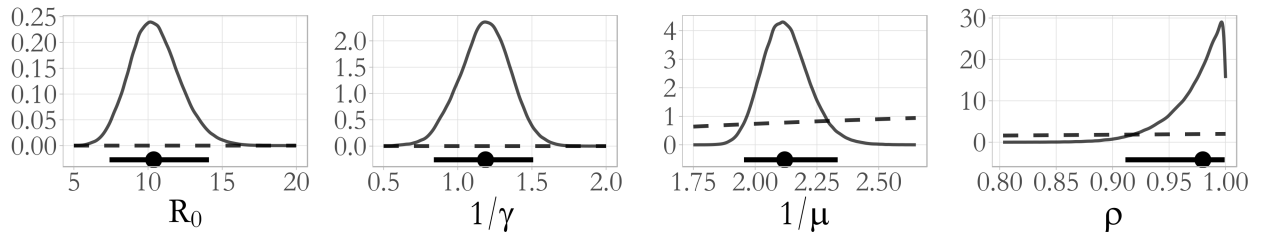
Author Manuscript



### SIR model



### SEIR Model



**Figure 12:**

Posterior density estimates for  $R_0 = \beta N/\mu$ , the mean latent and infectious periods,  $1/\gamma$  and  $1/\mu$ , and the binomial sampling probability,  $\rho$ , from SIR and SEIR model parameters fit to the British boarding school data (solid lines). The posterior median and 95% Bayesian credible intervals are drawn below the density plots (solid lines with circles). The implied prior densities (dashed lines) for  $R_0$  and the latent and infectious periods, and the prior density for the binomial sampling probability, are plotted over the posterior ranges.