# SUPPLEMENTARY INFORMATION

## S1   Complete Data Likelihood for SIS-type Contagions

Here we consider an SIS-type infectious disease to illustrate how our methodology can be applied to similar model classes. The complete data likelihood can be derived following the same steps in Section 3.2. Alternatively, one can slightly modify (14) to arrive at the complete likelihood for an SIS-type contagion. Since an individual doesn't acquire immunity upon recovery, it is equivalent to setting $H(t) \equiv S(t)$ at any time $t$. Thus the complete data likelihood is

$$
\begin{aligned}
&\mathcal{L}(\beta, \gamma, \tilde{\alpha}, \tilde{\omega} | \mathcal{G}_0) \\
&= \gamma^{n_R} \beta^{n_E - 1} \alpha_{SS}^{C_{SS}} \alpha_{SI}^{C_{SI}} \alpha_{II}^{C_{II}} \omega_{SS}^{D_{SS}} \omega_{SI}^{D_{SI}} \omega_{II}^{D_{II}} \prod_{j=2}^{n} \left[ \tilde{M}(t_j) \left( I_{p_{j1}}(t_j) \right)^{F_j} \right] \\
&\quad \times \exp\left( - \int_0^{T_{\max}} \left[ \beta SI(t) + \gamma I(t) + \tilde{\alpha}^T \mathbf{M}_{\max}(t) + (\tilde{\omega} - \tilde{\alpha})^T \mathbf{M}(t) \right] dt \right).
\end{aligned}
\tag{20}
$$

## S2   Auxiliary Proofs and Derivations

**Proof for Theorem 3.1**   From (14), we can obtain the log-likelihood:

$$
\begin{aligned}
&\ell(\beta, \gamma, \tilde{\alpha}, \tilde{\omega} | \mathcal{G}_0) = \log \mathcal{L}(\beta, \gamma, \tilde{\alpha}, \tilde{\omega} | \mathcal{G}_0) \\
&= \sum_{j=2}^{n} \left[ \log \tilde{M}(t_j) + F_j \log \left( I_{p_{j1}}(t_j) \right) \right] + n_R \log \gamma + (n_E - 1) \log \beta \\
&\quad + C_{HH} \log \alpha_{SS} + C_{HI} \log \alpha_{SI} + C_{II} \log \alpha_{II} + D_{HH} \log \omega_{SS} + D_{HI} \log \omega_{SI} + D_{II} \log \omega_{II}
\end{aligned}
\tag{21}
$$

$$-\sum_{j=1}^{n} \left[\beta SI(t_j) + \gamma I(t_j) + \tilde{\alpha}^T(\mathbf{M}_{\max}(t_j) - \mathbf{M}(t_j)) + \tilde{\omega}^T\mathbf{M}(t_j)\right](t_j - t_{j-1}).$$

Taking partial derivatives of the right hand side of (21) with respect to the parameters and setting them to zero yield the results above.

## S3    Relaxing the Closed Population Assumption

Suppose the observed population is not fully closed, but is a subset of a larger yet unobserved population. Then it is possible for an individual to get infected by an outsider. Let $\xi$ be the "external infection" rate, the rate for any susceptible individual to be infected by any external infectious source, then the complete data likelihood is

$$\mathcal{L}(\beta, \xi, \gamma, \tilde{\alpha}, \tilde{\omega}|\mathcal{G}_0) = p(\text{epidemic events}, \text{network events}|\beta, \xi, \gamma, \alpha, \omega, \mathcal{G}_0)$$

$$= \gamma^{n_R} \alpha_{SS}^{C_{HH}} \alpha_{SI}^{C_{HI}} \alpha_{II}^{C_{II}} \omega_{SS}^{D_{HH}} \omega_{SI}^{D_{HI}} \omega_{II}^{D_{II}} \prod_{j=2}^{n} \left[\tilde{M}(t_j)\left(\beta I_{p_{j1}}(t_j) + \xi\right)^{F_j}\right]$$

$$\times \exp\left(-\int_0^{T_{\max}} \left[\beta SI(t) + \xi S(t) + \gamma I(t) + \tilde{\alpha}^T\mathbf{M}_{\max}(t) + (\tilde{\omega} - \tilde{\alpha})^T\mathbf{M}(t)\right]dt\right). \quad (22)$$

MLEs for $\{\gamma, \tilde{\alpha}, \tilde{\omega}\}$ remain unchanged, but estimating $\beta$ and $\xi$ is less straightforward. Let $\ell(\beta, \xi, \gamma, \tilde{\alpha}, \tilde{\omega}|\mathcal{G}_0)$ be the log-likelihood, then the partial derivatives of the log-likelihood w.r.t. $\beta$ and $\xi$ are

$$\frac{\partial \ell}{\partial \beta} = \sum_{j=2}^{n} \frac{F_j I_{p_{j1}}(t_j)}{\beta I_{p_{j1}}(t_j) + \xi} - \sum_{j=1}^{n} SI(t_j)(t_j - t_{j-1}),$$

$$\frac{\partial \ell}{\partial \xi} = \sum_{j=2}^{n} \frac{F_j}{\beta I_{p_{j1}}(t_j) + \xi} - \sum_{j=1}^{n} S(t_j)(t_j - t_{j-1}),$$

which do not directly lead to closed-form solutions.

Reparameterizing by $\xi = \kappa\beta$ leads to the following partially derivatives

$$\frac{\partial \ell}{\partial \beta} = \frac{n_E - 1}{\beta} - \sum_{j=1}^{n} [SI(t_j) + \kappa S(t_j)](t_j - t_{j-1}), \quad (23)$$

2

$$\frac{\partial \ell}{\partial \kappa} = \sum_{j=2}^{n} \frac{F_j}{I_{p_{j1}}(t_j) + \kappa} - \beta \sum_{j=1}^{n} S(t_j)(t_j - t_{j-1}), \tag{24}$$

which are slightly more straightforward in form, and can be solved numerically to obtain the MLEs.

If, somehow, we have information on which infection cases are caused by internal sources and which are caused by external sources, then we can directly obtain the MLEs and Bayesian posterior distributions for all the parameters. For an infection event $e_j$ (with $F_j = 1$), let $\mathrm{Int}_j = 1$ if it is "internal" and let $\mathrm{Int}_j = 0$ otherwise. Then the complete data likelihood can be re-written as

$$\mathcal{L}(\beta, \xi, \gamma, \tilde{\alpha}, \tilde{\omega} | \mathcal{G}_0)$$

$$= \beta^{\left(n_E^{\mathrm{int}} - \mathrm{Int}_1\right)} \xi^{\left(n_E^{\mathrm{ext}} - 1 + \mathrm{Int}_1\right)} \gamma^{n_R} \alpha_{SS}^{C_{HH}} \alpha_{SI}^{C_{HI}} \alpha_{II}^{C_{II}} \omega_{SS}^{D_{HH}} \omega_{SI}^{D_{HI}} \omega_{II}^{D_{II}} \prod_{j=2}^{n} \left[ \tilde{M}(t_j) I_{p_{j1}}(t_j)^{F_j \mathrm{Int}_j} \right]$$

$$\times \exp\left( -\int_0^{T_{\max}} \left[ \beta SI(t) + \xi S(t) + \gamma I(t) + \tilde{\alpha}^T \mathbf{M}_{\max}(t) + (\tilde{\omega} - \tilde{\alpha})^T \mathbf{M}(t) \right] dt \right), \tag{25}$$

where $n_E^{\mathrm{int}}$ and $n_E^{\mathrm{ext}}$ are the total numbers of internal and external infection events, respectively.

Estimation for all parameters remains unchanged except for $\beta$ and $\xi$. Their MLEs are

$$\hat{\beta} = \frac{n_E^{\mathrm{int}} - \mathrm{Int}_1}{\sum_{j=1}^{n} SI(t_j)(t_j - t_{j-1})}, \quad \hat{\xi} = \frac{n_E^{\mathrm{ext}} - 1 + \mathrm{Int}_1}{\sum_{j=1}^{n} S(t_j)(t_j - t_{j-1})}, \tag{26}$$

and with Gamma priors $\beta \sim Ga(a_\beta, b_\beta)$ and $\xi \sim Ga(a_\xi, b_\xi)$, their posterior distributions are

$$\beta | \{e_j\} \sim Ga\left( a_\beta + (n_E^{\mathrm{int}} - \mathrm{Int}_1), b_\beta + \sum_{j=1}^{n} SI(t_j)(t_j - t_{j-1}) \right), \tag{27}$$

$$\xi | \{e_j\} \sim Ga\left( a_\xi + (n_E^{\mathrm{ext}} - 1 + \mathrm{Int}_1), b_\xi + \sum_{j=1}^{n} S(t_j)(t_j - t_{j-1}) \right). \tag{28}$$

When there is missingness in recovery times, the Bayesian inference procedure described in Section 4 can still be carried out, with two slight modifications. First, in the data augmentation step, when drawing missing recovery times in an interval $(u, v]$, the DARCI algorithm ( Prop. 4.2) inspects $\mathcal{I}_p$ only for each $p \in \mathcal{P}^{\text{int}}$, where $\mathcal{P}^{\text{int}}$ is the group of individuals who get *internally* infected during $(u, v]$. Second, in each iteration, parameter values are drawn from the posterior distributions specified in (15) except for $\beta$ and $\xi$, for which the posteriors are stated in (27) and (28), respectively.

# S4 Flexible Network Dynamics under the Generative Model

As stated in the main text, the proposed generative model can accommodate **arbitrary** initial network structures. Moreover, due to the **coupled** nature of the epidemic process and the network process, different choices of parameters can lead to a wide variety of network behaviors.

Here we demonstrate via simulations that, even if the initial network follows a simple, idealized model, it can still evolve and adapt to exhibit characteristics that are no longer simple and idealistic throughout the process.

Setting the population size to $N = 500$, we assume two different models for the initial network: 1) a random Erdős–Rényi graph ("ER"), and 2) a random **scale-free** network (sampled using the Barabasi-Albert model ("BA") (Barabási and Albert, 1999)). The "ER" model leads to a degree distribution close to Poisson, whereas the "BA" model leads to power-law degree distribution.

Using the following parameter settings (same as those in most of our simulations),

$$\beta = 0.03, \gamma = 0.12;$$

$$\tilde{\alpha}^T = (\alpha_{SS}, \alpha_{SI}, \alpha_{II}) = (0.005, 0.001, 0.005),$$

$$\tilde{\omega}^T = (\omega_{SS}, \omega_{SI}, \omega_{II}) = (0.05, 0.1, 0.05).$$

we adopt a dynamic that 1) discourages new $S - I$ connections and thus effectively "isolates" super-spreaders probabilistically, and 2) mostly sustains the usual contact activities between $S-S$, $R-R$ and $I-I$ pairs. This, during the process, can lead to a **bimodal** degree distribution that obviously deviates from the initial degree distribution.

Below in Figure S1, we include the empirical degree distributions at the start and the end of a typical simulation, with a random Erdős–Rényi graph as the initial network. Here the red curve represents the actual degree distribution in the simulated population, and (in contrast) the lightblue curve represents the density of a degree sequence sampled from a Poisson with the same mean degree.

For comparison, in Figure S2 we include similar plots with the initial network as a **scale-free** network (sampled using the Barabasi-Albert model). Note that the initial network setting here is similar to that used in Volz and Meyers (2007), but our dynamic network model is fundamentally different from (and more flexible) than the neighbor-exchange framework proposed in that work.

We can observe that, once the network dynamics kicks in, the initial network structure doesn't matter that much, as changes in the network links are driven jointly by the parameters and the epidemic dynamics in time.
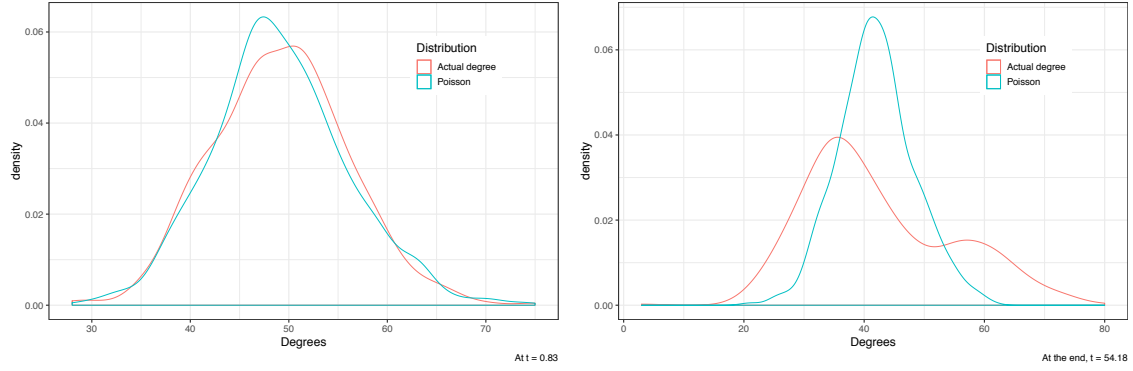
At t = 0.83

At the end, t = 54.18

Figure S1: Actual empirical degree distribution in simulation (red) versus empirical degree distribution drawn from Poisson (blue/green). Here the initial network is a random Erdős–Rényi graph. **Left: beginning of process; right: end of process.**
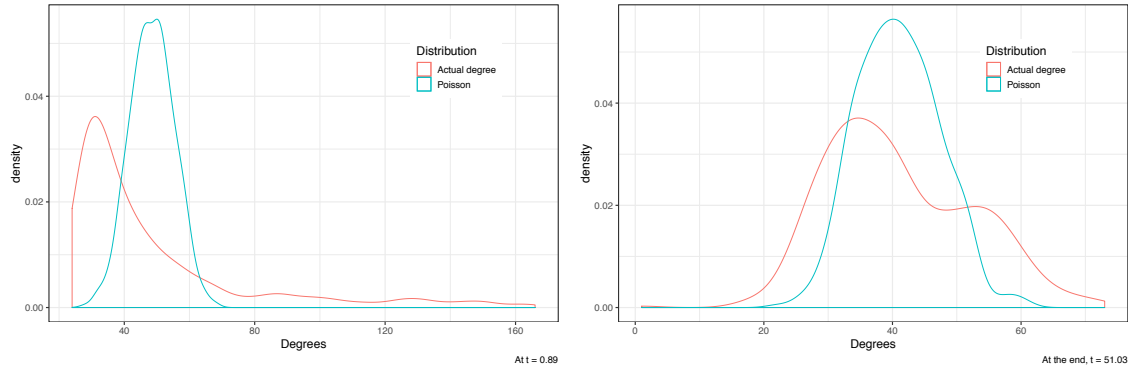


At t = 0.89

At the end, t = 51.03

Figure S2: Actual empirical degree distribution in simulation (red) versus empirical degree distribution drawn from Poisson (blue/green). Here the initial network is a random scale-free network (sampled using the Barabasi-Albert model (Barabási and Albert, 1999)). **Left: beginning of process; right: end of process.**

# S5  More Results on Simulation Experiments

**Supplement for "inference from complete event data"**   Figure S3 and S4 complement Figure 2 and 3 in the main text, showing inference results for all the parameters in the corresponding experiments.

**Experiments on larger networks**   Figure S5 shows MLEs and 95% confidence bands for parameters with complete data generated on a network with $N = 500$ individuals. Other experimental settings are the same as those in Section 5.1. With a larger population, there tends to be more events available for inference, so the accuracy is in fact improved.

**Experiments on different initial network configurations**   Still set population size $N = 100$, but instead of a random Erdős–Rényi graph as $\mathcal{G}_0$, the initial network is a "hubnet": one individual (the "hub") is connected to everyone else in the population while the others form an $ER(N - 1, p)$ random graph, with edge probability $p = 0.1$. Figure S6 summarizes results of Bayesian inference carried out on complete event data generated in this setting.

**Supplement for "Assessing model flexibility"**   Estimate parameters $\Theta$ of the full model on datasets generated from 1) the decoupled temporal network epidemic process with type-independent edge rates, and 2) the static network epidemic process where the network remains unchanged. For both simpler models, fix $\beta = 0.03$ and $\gamma = 0.12$, and for the former model, let link activation rate $\alpha = 0.005$ and termination rate $\omega = 0.05$. Still, set population size $N = 100$ and let the initial network be a random Erdős–Rényi graph with edge probability $p = 0.1$.

We present, in Figure S7, the results of Bayesian inference on datasets generated from the decoupled process model. Across four different realizations, it can be observed that, the posterior samples of link activation rates $(\alpha_{SS}, \alpha_{SI}, \alpha_{II})$ concentrate around the same mean, and uncertainty is reduced with more events available for inference. Same can be said about the link termination rates, $\omega_{SS}, \omega_{SI}, \omega_{II}$. This verifies that the proposed model is indeed a generalization of the aforementioned two simpler processes, and the inference method is able to recover the truth under mild model misspecification.
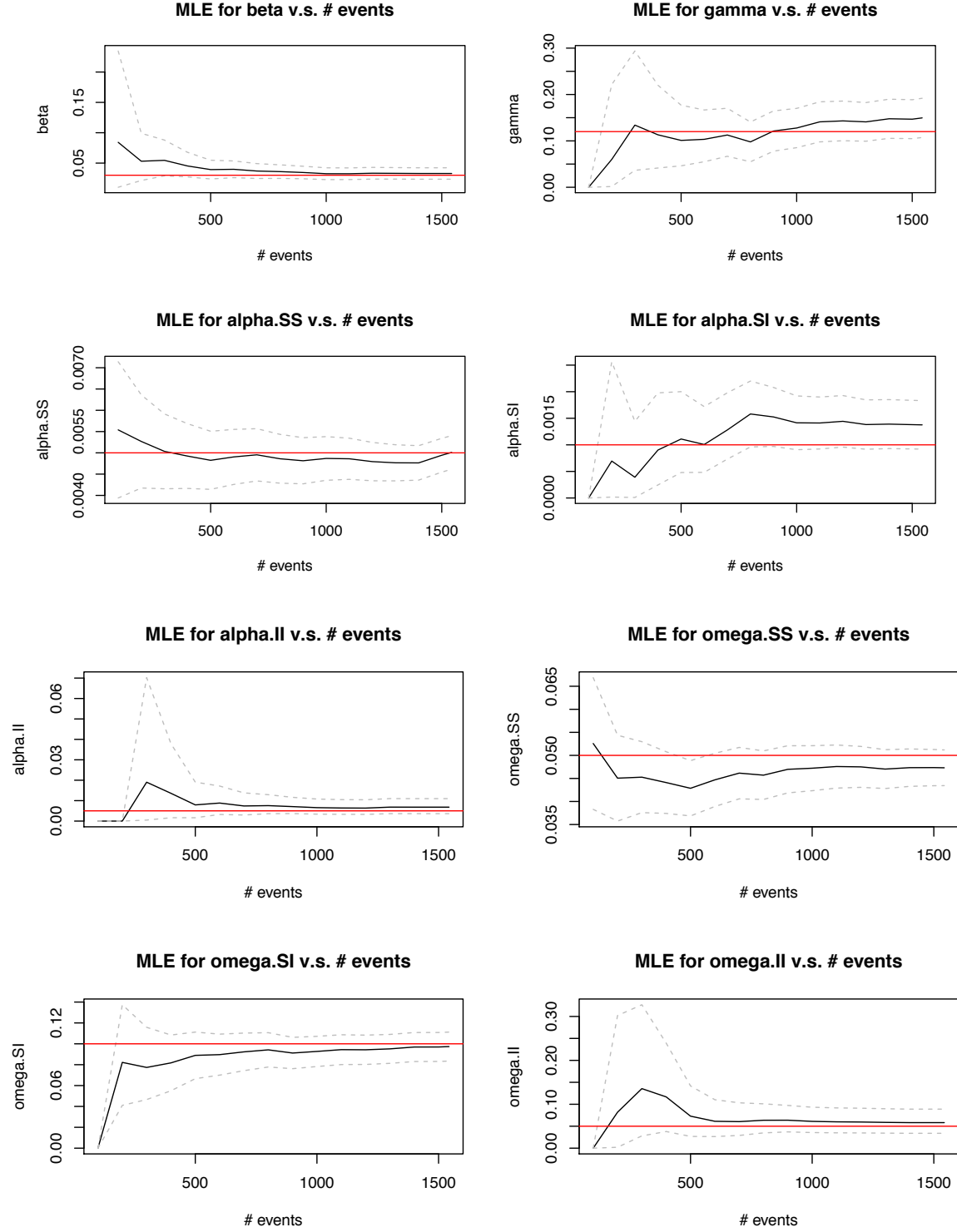
Figure S3: MLEs versus number of events used for inference. Dashed gray lines show the lower and upper bounds for 95% frequentist confidence intervals, and red lines mark the true parameter values.
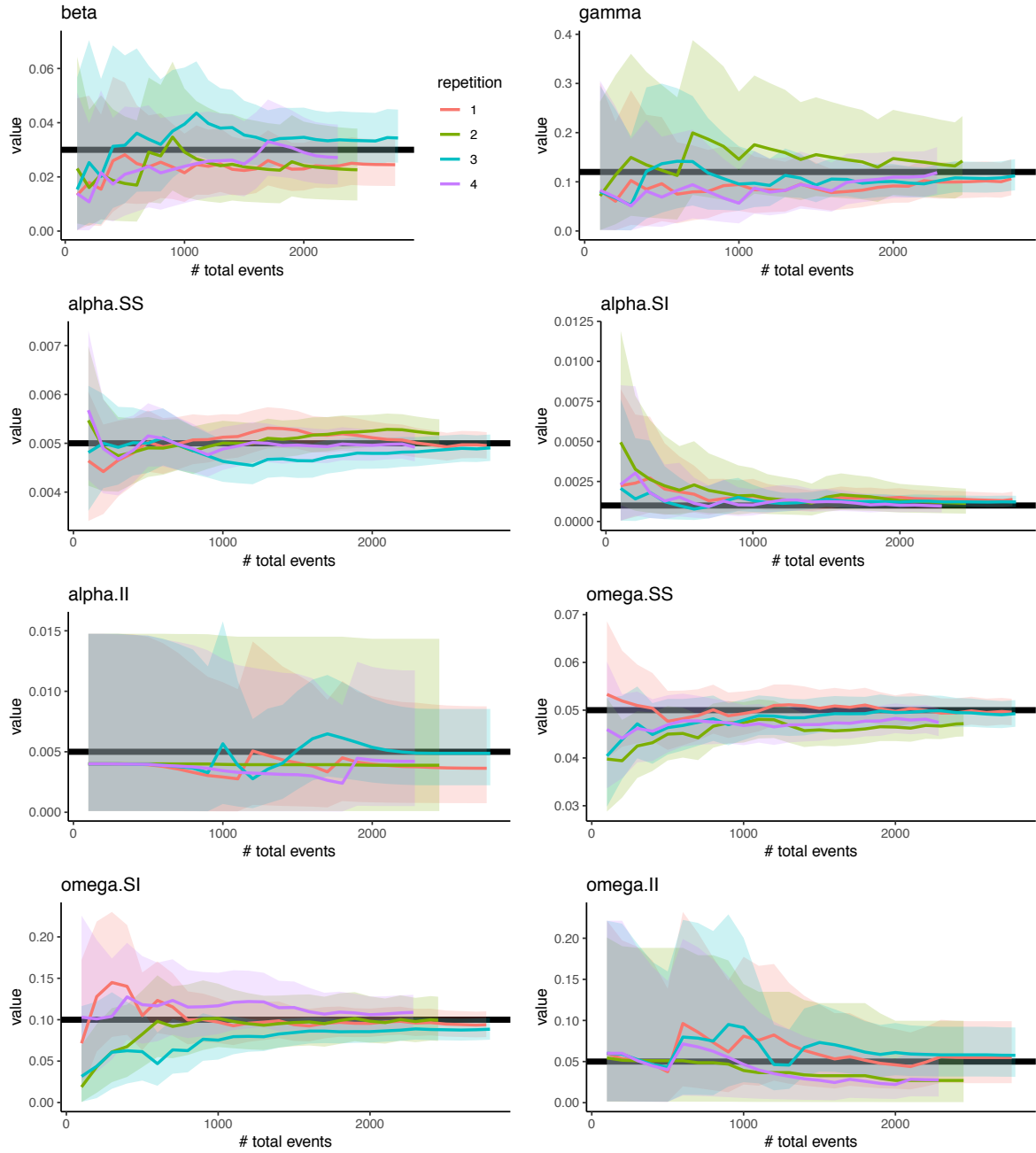
Figure S4: Posterior sample means v.s. number of total events used for inference. True parameter values are marked by **bold dark** horizontal lines, along with 95% credible bands. Results are presented for 4 different complete datasets.

**MLE for beta v.s. # events**

**MLE for gamma v.s. # events**

**MLE for alpha.SS v.s. # events**

**MLE for alpha.SI v.s. # events**

**MLE for alpha.II v.s. # events**

**MLE for omega.SS v.s. # events**

**MLE for omega.SI v.s. # events**
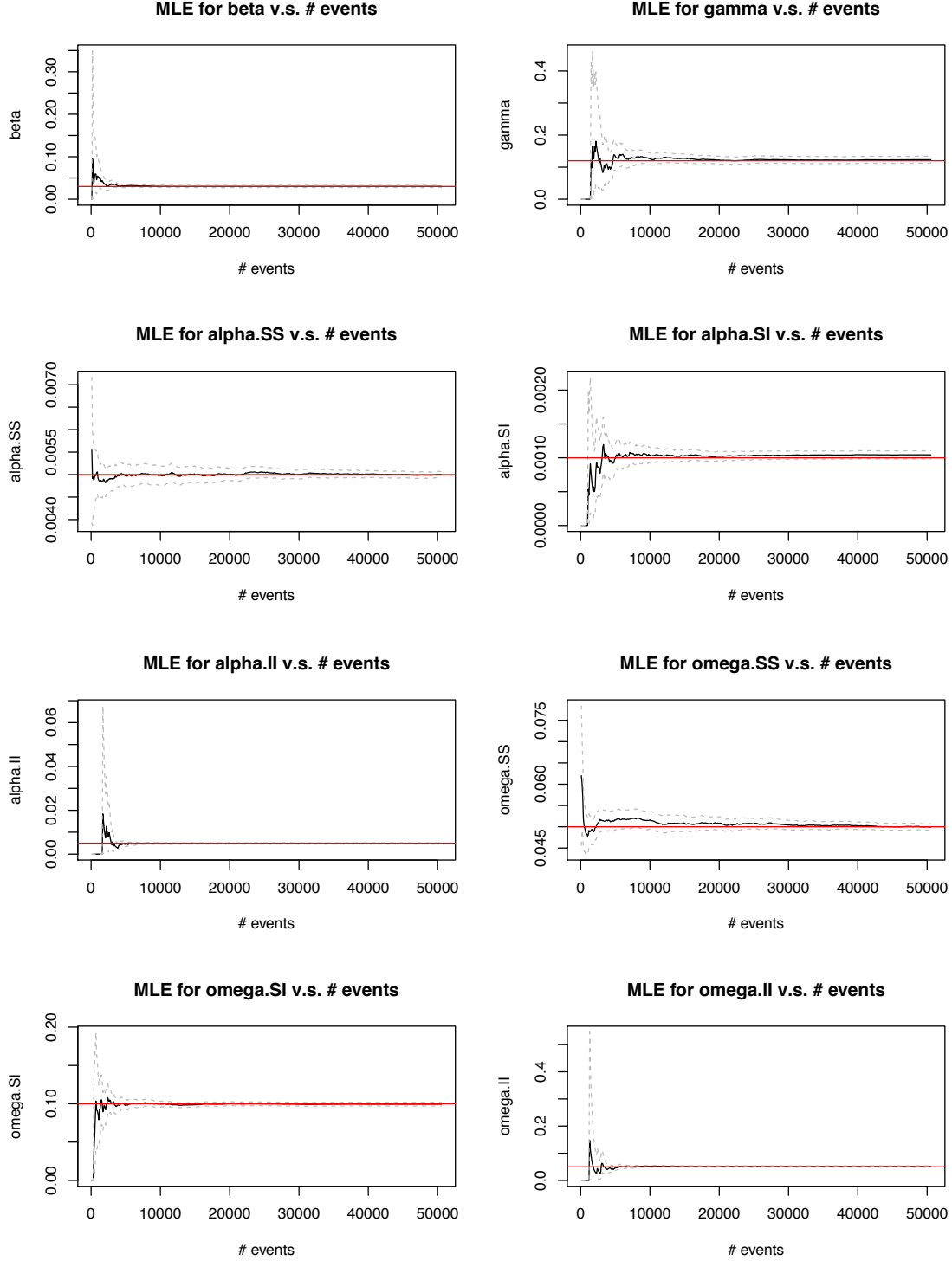
**MLE for omega.II v.s. # events**

Figure S5: MLEs versus number of total events, on a larger population with $N = 500$. Dashed gray lines show the lower and upper bounds for 95% confidence intervals, and red lines mark the true parameter values. With a larger population size, there tends to be more events, which in fact facilitates estimation.
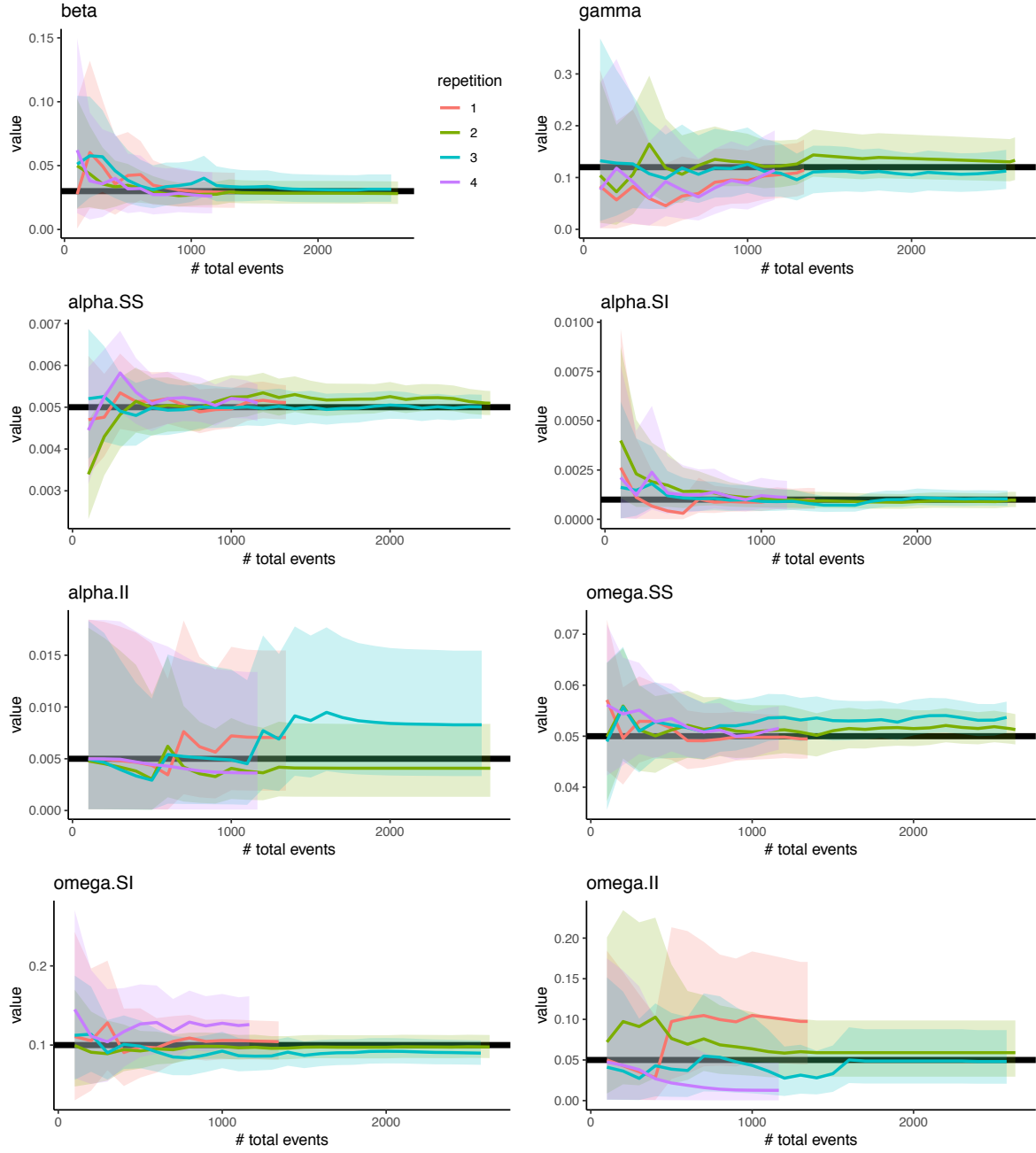
Figure S6: Posterior sample means v.s. number of total events, with $\mathcal{G}_0$ as a $N = 100$-node "hubnet". True parameter values are marked by **bold dark** horizontal lines, along with 95% credible bands. Results are presented for 4 different complete datasets.
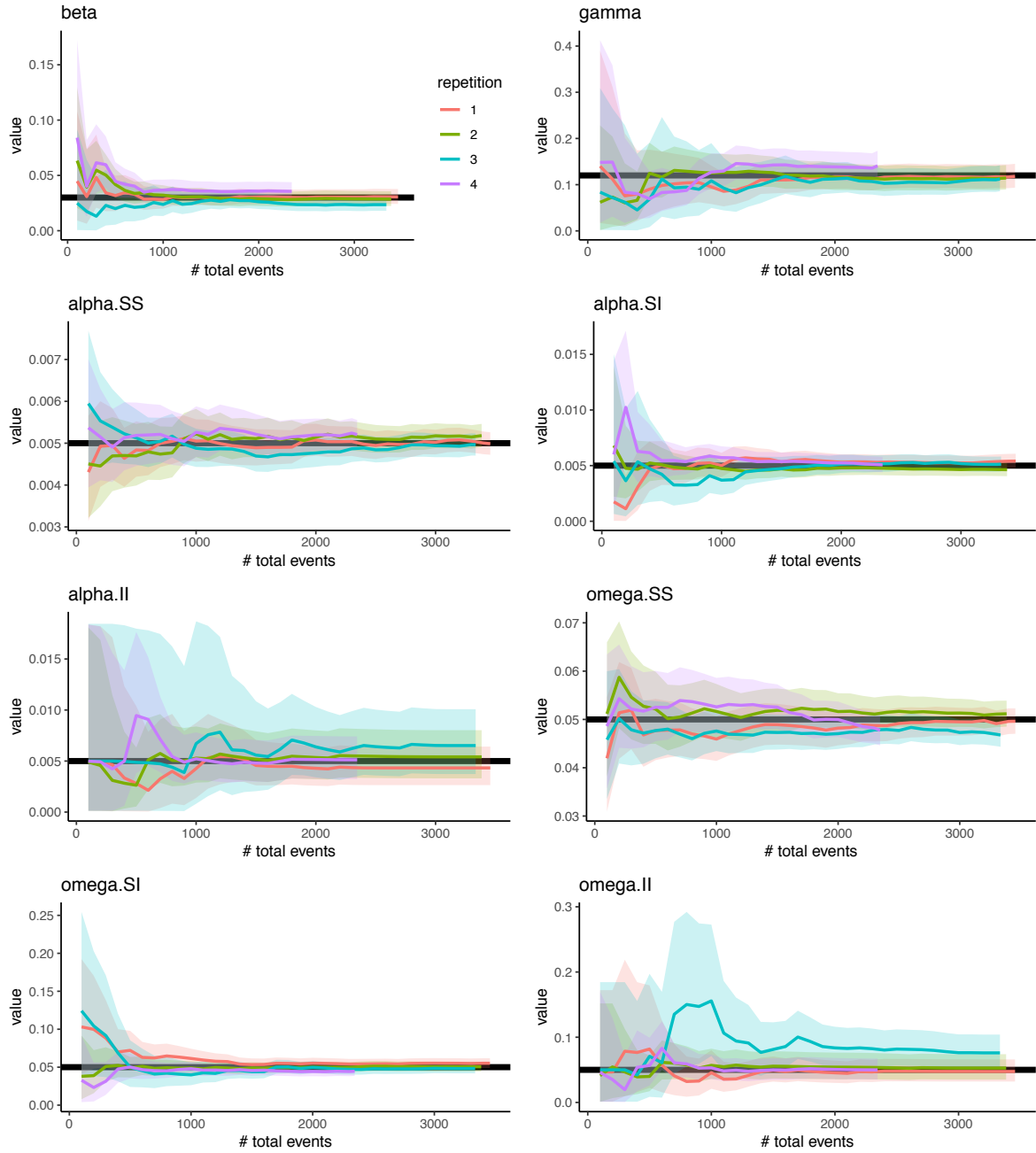
12

Figure S7: Posterior sample means versus number of total events, estimated using datasets generated by the decoupled process model. True parameter values are marked by **bold dark** horizontal lines, along with 95% credible bands. Results are presented for 4 different complete datasets.

13

**Scalability of DARCI and the data-augmented inference scheme**   In the main text, most simulations are conducted on a population of size $N = 100$, in order to mimic the population size of the real data, but experiments have been carried out on larger networks (for example, with $N = 200$ and $N = 500$). Here we include some inference results for $N = 500$ (see Figure S8).
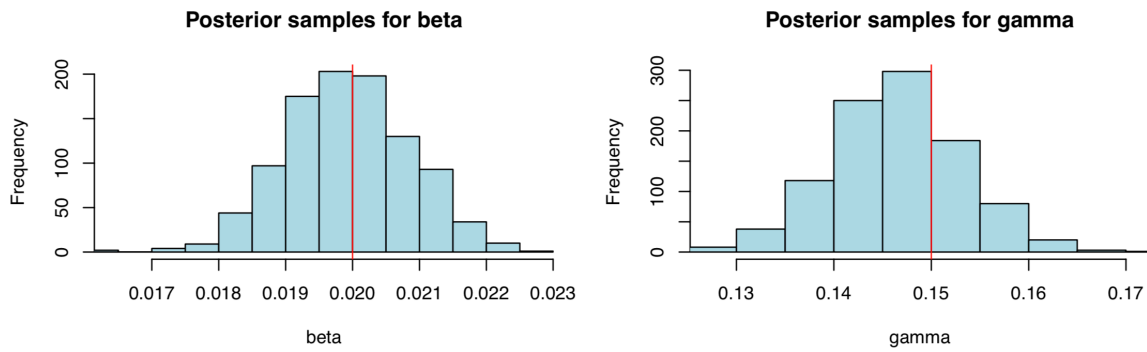


Figure S8: Inference performance with population size $N = 500$. Here we show histograms of posterior samples for parameter $\beta$ (infection rate, left panel) and $\gamma$ (recovery rate, right panel); ground truth is represented by red vertical lines. The proposed inference scheme can certainly handle a relatively large population with a lot of events, and can recover parameter values accurately.

It seems that the inference scheme is effective in recovering parameters and is capable of handling a large population and many events. Our finding is that computing time scales linearly with the total number of events observed. Figure S9 shows the computing time on a single processor for each iteration with population size $N = 500$, using a naive implementation in R. Since DARCI is parallelizable across time intervals, a more efficient implementation can further reduce computing time, and the algorithm bottleneck will be the maximum number of recovery events that occur in a time interval; moreover, the

computation in one iteration only involves vector operations and random number sampling, which can be significantly sped up in other programming languages if necessary.
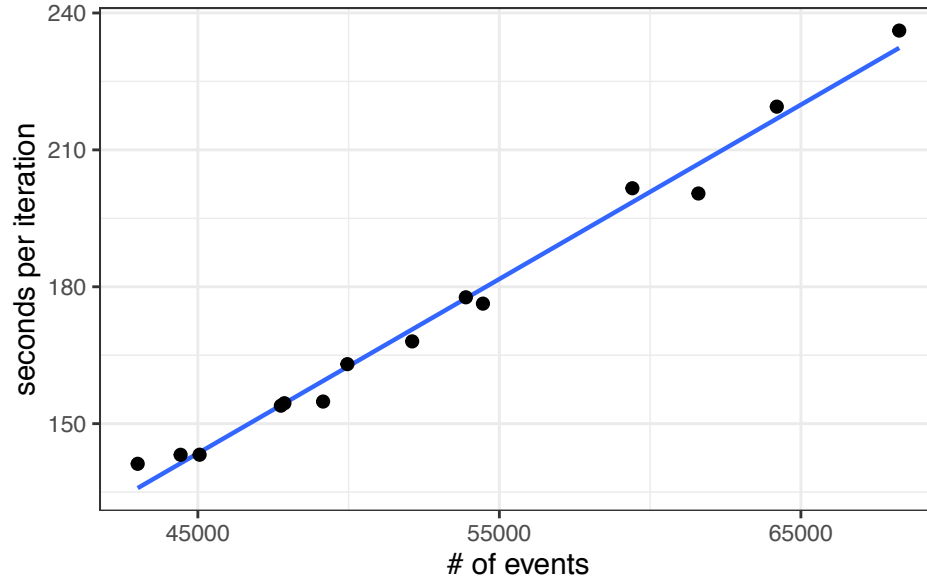


Figure S9: Run time per iteration (in seconds) on **one processor** versus number of observed events in simulations, for population size $N = 500$. Computing time scales linearly with the number of events, and the algorithm can handle data at least at the scale of $10^4$ events.

# S6   Real Data Experiments

## S6.1   Data Pre-processing

All infection events and weekly health statuses of all $N = 103$ individuals are extracted from the weekly surveys. In every survey, study participants were asked if they ever felt ill at all in the past week, if they ever experienced certain symptoms, and, if there were symptoms, when the approximate illness onset time was. We take an "infection"

15

as a positive ILI (influenza-like illness) case, which, following the protocol in Aiello et al. (2016), is defined as a cough plus at least one of the following symptoms: fever or feverishness, chills, or body aches. We further examine each ILI case and only accept one as a positive infection if the individual also indicated that they "felt ill" in the past week, thus eliminating a small number of reoccurring ILI cases for the same participants [7]. Moreover, since an individual may start exhibiting symptoms at most 3 days *after* getting infected and becoming infectious, for each infection event, we set the "real" infection time as the reported onset time minus a random "delay time" uniformly sampled between 0 and 3 days.

Social link activation and termination events are obtained from the iEpi Bluetooth contact records. Each time two study devices were paired, the iEpi application recorded the unique identifiers of the devices, a timestamp, and a received signal strength indicator (RSSI). Since Bluetooth detection can be activated whenever two devices are within a few meters of each other while the two users may not actually be in contact, we only keep those Bluetooth records with relatively strong signals (high values of RSSIs) [8]. If two consecutive Bluetooth records for one pair of devices are no more than 7.5 minutes apart in time [9], then the two records are considered to belong to one single continuous contact; a social link between two individuals is activated at the time of the first Bluetooth detection record in a series of consecutive records that belong to a single contact, and the link is terminated at a random time point between 1 and 6 minutes

---

[7]One particular individual had positive ILI cases and felt ill in week 2, 3, and 5, but not in week 4. We therefore treat his/her illness as an extended one, starting in week 2 and lasting till week 5.

[8]The RSSIs range from -109 to 6, and we set the threshold as -90, so only those records with RSSIs larger than -90 are kept.

[9]We choose 7.5 minutes as a threshold instead of 5 minutes to accommodate potential lapses in Bluetooth detection.

after the last Bluetooth detection of a continuous contact.

The resulting processed data contain 24 infection events in total, with 14 before the spring break week and 10 after, as well as 45,760 social link activation and termination events. The weekly disease status (healthy or ill) of every participant can be acquired from the weekly surveys, so we know, for example, if an individual recovered sometime after day 7 and before day 14, but the exact times of all recoveries are unknown.

## S6.2   Maximum Likelihood Estimation

Instead of assuming the knowledge of which infection cases are internal and which are external, we directly estimate all the parameters based on the likelihood function in (22), solving (23) and (24) for the MLEs of $\beta$ and $\xi$.

However, the real data are incomplete, with the exact times of all the recoveries unobserved. We resolve this issue using a naive imputation method—for each recovery, an event time is randomly sampled from a uniform distribution between the time of infection and the earliest time point the individual no longer felt ill (in response to the weekly surveys). Such imputation, of course, is subject to a considerable level of uncertainty, so we randomly generate 10 differently imputed datasets , obtain the MLEs from every dataset, and then report the averages over the 10 runs (see Table S1).

We can see that the MLEs acquired in this manner generally agree with the Bayesian estimates in Section 6.2.

Table S1: MLEs for model parameters using imputed data with all recovery times randomly sampled. The table presents average estimates as well as the standard deviations of estimates over 10 different, randomly imputed datasets. Results generally agree with those acquired using the proposed Bayesian data augmentation inference method.

| Parameter | Avg. estimate | Std. deviation |
|---|---|---|
| $\beta$ (internal infection) | 0.0676 | 0.0092 |
| $\xi$ (external infection) | 0.00320 | $1.11 \times 10^{-6}$ |
| $\gamma$ (recovery) | 0.236 | 0.012 |
| $\alpha_{SS}$ (S-S link activation) | 0.0530 | 0.0001 |
| $\omega_{SS}$ (S-S link termination) | 42.15 | 0.105 |
| $\alpha_{SI}$ (S-I link activation) | 0.0704 | 0.0028 |
| $\omega_{SI}$ (S-I link termination) | 52.21 | 3.83 |