

Supplementary material

**Integrating epidemiological and genetic data with different sampling densities into a dynamic model of respiratory syncytial virus (RSV) transmission.**

Ivy K. Kombe<sup>1,2</sup>, Charles N. Agoti<sup>1</sup>, Patrick K. Munywoki<sup>1</sup>, Marc Baguelin<sup>3</sup>, D. James Nokes<sup>4</sup>,  
Graham F. Medley<sup>2</sup>

<sup>1</sup> KEMRI-Wellcome Trust Research Programme, KEMRI Centre for Geographical Medical Research-Coast. P.O. Box 230-80108, Kilifi, Kenya.

<sup>2</sup> Centre for Mathematical Modelling of Infectious Disease and Department of Global Health and Development, London School of Hygiene and Tropical Medicine. London, WC1H 9SH, United Kingdom.

<sup>3</sup> Centre for Mathematical Modelling of Infectious Disease and Department of Infectious Disease Epidemiology, London School of Hygiene and Tropical Medicine. London, WC1H 9SH, United Kingdom.

<sup>4</sup> School of Life Sciences and Zeeman Institute for Systems Biology & Infectious Disease Epidemiology Research, University of Warwick. Coventry, CV4 7AL, United Kingdom.

Correspondence: Ivy K. Kombe (ivkadzo@gmail.com)

22 **Table of contents**

23 A1. Data pre-processing.....3

24     Imputing complete shedding, ARI and presence durations. ....3

25     Imputing missing genetic information .....6

26 A2. Further details of the transmission model .....9

27     The rate of exposure.....9

28     The background community function.....12

29     The Likelihood.....17

30 A3. Further details of the adaptive MH-MCMC algorithm.....19

31     Choice of proposal distributions for the parameters .....20

32     Pseudo algorithm for our implementation of MH-MCMC.....21

33 A4. Further details of the HPTS .....24

34     Establishing the highest probability transmission source (HPTS).....24

35 A5. Extra results.....27

36     Parameter trace plots and convergence checks.....27

37 A6. Model validation .....37

38 A7. Model modification to fit pathogen data identified at group resolution .....46

39 A8. References.....49

40

41 **A1. Data pre-processing**

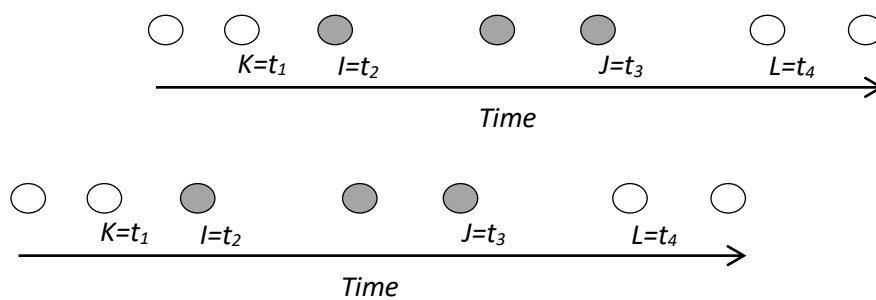
42 Imputing complete shedding, ARI and presence durations.

43 Given the 3-4 day sampling intervals, complete shedding and ARI durations had to be  
44 imputed, and missing viral loads linearly interpolated. For the model, we will assume that all  
45 the cases were observed, and ignore the possibility of short duration shedding episodes that  
46 could have been missed by the sampling intervals. During the sample-collection visits, if a  
47 household member was not present, they were recorded as being 'away' on that particular  
48 day. As with the shedding information, there was incomplete information on continuous  
49 periods of presence or absence from the household which was also imputed.

50

51 An RSV A/B shedding episode is defined as a period within which an individual provided PCR  
52 positive samples for RSV A/B that were no more than 14 days apart. Using the mid-point  
53 method, shedding was assumed to start mid-way between the last negative sample and the  
54 first positive sample, and it ended midway between the last positive sample and the first  
55 negative sample of an episode. This is illustrated below:

56



57

58

59 *Grey circles are positive samples in a single episode, empty circle are negative.  $t_1$ ,  $t_2$ ,  $t_3$  and  $t_4$*   
60 *are dates of sample collection.*

61

62 For  $(t_4 - t_3)$  and  $(t_2 - t_1) \leq 7$  days

63 
$$Duration = \left[ t_3 + \left( t_4 - t_3 / 2 \right) \right] - \left[ t_2 - \left( t_2 - t_1 / 2 \right) \right]$$

64 For  $(t_4 - t_3) > 7$

65 
$$Duration = \left[ t_3 + \left( x / 2 \right) \right] - \left[ t_2 - \left( t_2 - t_1 / 2 \right) \right] : \text{Right censoring}$$

66 For  $(t_2 - t_1) > 7$

67 
$$Duration = \left[ t_3 + \left( t_4 - t_3 / 2 \right) \right] - \left[ t_2 + \left( x / 2 \right) \right] : \text{Left censoring}$$

68 Where  $x$  = mean of sampling intervals for samples in an episode, which was found to be 3.45

69 days.

70 Any negative samples ( $C_t > 35$  or  $C_t = 0$ ) in between a shedding episode were ignored, i.e.

71 were not treated like true end of shedding

72

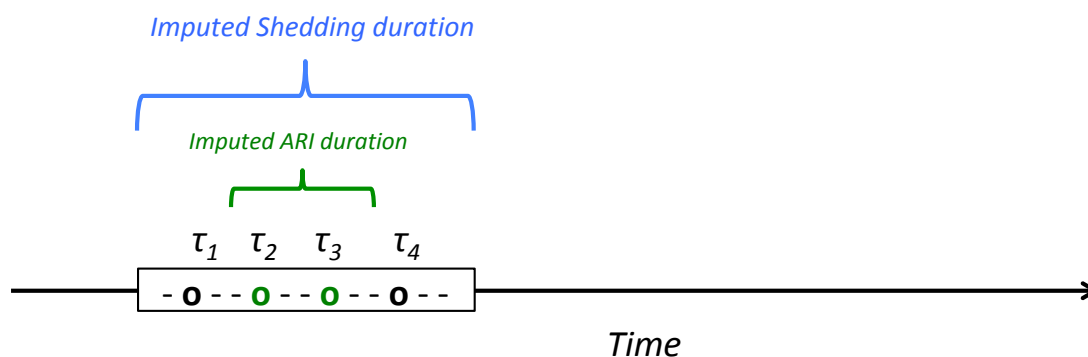
73 We imputed complete ARI episodes from intervals of recorded ARI. A virus shedding

74 episode that had no day where an ARI was reported was assumed to be asymptomatic. For

75 a virus shedding episode with at least one day of recorded ARI, the duration of symptoms

76 was imputed using the midpoint method described for shedding episodes. This is illustrated

77 below:



79 *Green open circles are reported ARI symptoms (ARI positive) within the shedding episode*  
 80 *and black open circles are confirmed absence of ARI (ARI negative).  $\tau_1, \tau_2, \tau_3$  and  $\tau_4$  are days*  
 81 *within the shedding episode where information on symptoms was collected.*

82

83 In this case, the mean sampling interval for ARI ‘samples’ within an episode was 3.78 days.

84 This was obtained from all ARI episodes not just the ones within shedding episodes

85

86 The imputation of continuous periods of presence or absence from the household was done

87 similar to the imputation of shedding durations, however, there was no left or right

88 censoring. Each participant had a set of days of recorded data, these days were either

89 marked as ‘away’ or ‘present’ in the household, e.g. a participant might have data on days

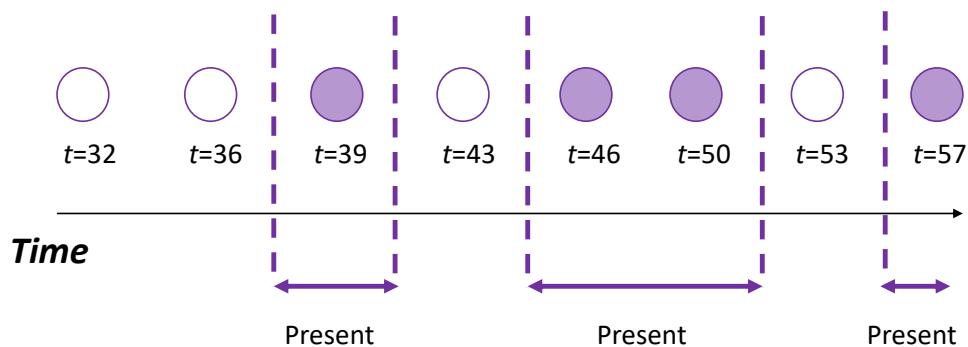
90 {32, 36, 39, 43, 46, 50, 53, 57} with status {away, away, present, away, present, present,

91 away, present}. Since no data is available for this individual before day 32 and after day 57,

92 no imputation is done outside this time window. For the days within the window,

93 imputation is done as illustrated below:

94



95

96 *Filled circles are days when the participant was recorded as being present while open ones is*

97 *when they were away. The present period starts halfway between the last ‘away’ and first*

98 *‘present’ and ends halfway between the last ‘present’ and first ‘away’.*

99 In order to include information of the amount of virus shed by an infected person into the  
100 transmission model, the Ct value need to be converted to  $\log_{10}$  RNA copy number which is a  
101 more direct measure of viral load. The formula used to convert Ct values to their  $\log_{10}$  RNA  
102 equivalent was  $y = -3.308x + 42.9$ , where  $y = \text{Ct values}$  and  $x = \log_{10}$  RNA copy number [1,2].  
103 Following conversion of the PCR Ct values to viral load, we proceeded to interpolate the  
104 viral loads for days in an episode that did not have data. Linear interpolation was used for all  
105 the shedding episodes. It was assumed that the starting and ending sample, if data was  
106 missing, had a viral load of 2.388  $\log_{10}$  RNA (baseline positive Ct value converted to viral  
107 load). For two samples of viral load  $V_a$  and  $V_b$  at times  $t_a$  and  $t_b$ ,  $t_b > t_a$ , the gap in between is  
108 filled out as follows:

109 For  $t_b - t_a = n$ , viral load  $V_j$  at time point  $t_j$  for  $j = 1 \dots (n-1)$  is given by

$$110 \quad V_j = V_a + \frac{j(V_b - V_a)}{n}$$

111 Viral loads lower than 2.388  $\log_{10}$  RNA in between an episode were not included in the  
112 interpolation

113

#### 114 Imputing missing genetic information

115 The WGS data was used to rule out transmission events where it was assumed that cases in  
116 different genetic clusters are not part of the same transmission cluster. It was also assumed  
117 that for cases within the same genetic cluster, the likelihood of a transmission event is  
118 weighed according to pairwise genetic distance  $d_{gen}(i,j)$ . As mentioned in the main text,  
119 genetic clusters were established based on a combination of criteria: nucleotide distance  
120 cut-off, clustering patterns on the global RSV phylogeny and the inferred date of sequence  
121 divergence. As a result of incomplete sequencing of all the positive samples, there are gaps

122 in the genetic data. To fill these in we classified missingness into 3 categories and exploited  
123 elements of the study design to fill in the gaps.

- 124 • Missing level 1: non-sequenced samples part of an infection episode with  $\geq 1$  other  
125 sequenced sample. The entire episode was assigned the cluster id of the sequenced  
126 sample(s), where there was more than 1 id, the episode was divided accordingly.
- 127 • Missing level 2: none of the samples in an episode were sequenced, but the episode  
128 is part of a spatial-temporal cluster with some genetic information. The entire  
129 episode was assigned the cluster id of the spatial-temporal cluster. This assumes that  
130 if an episode has a temporal overlap with other cases in the same household, they  
131 are likely part of the same infection cluster (household outbreak).
- 132 • Missing level 3: none of the samples in an episode were sequenced and there is no  
133 genetic information in the social-temporal cluster. The cluster id for the entire  
134 episode was treated as augmented data and inferred along with the model  
135 parameters.

136 Within a given RSV group, infection by a particular cluster is assumed to be a mutually  
137 exclusive process, an individual can only shed one cluster type at a time. The genetic data  
138 available is consensus whole genome sequences as such, only one cluster can be identified  
139 from a single sample.

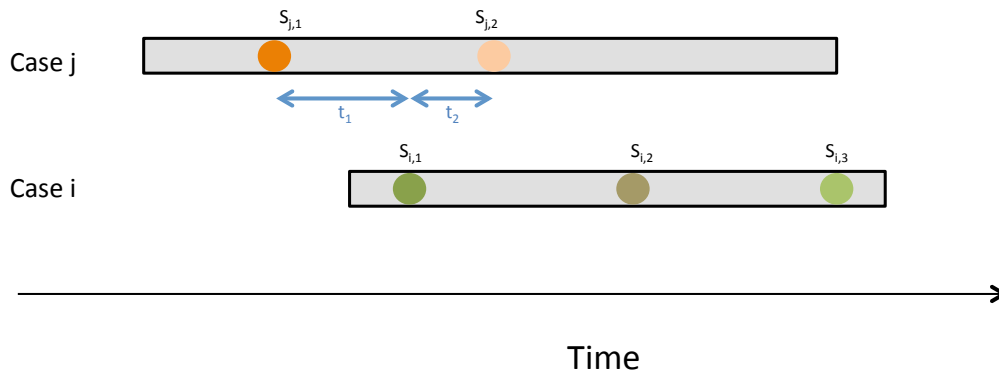
140

141 Consider a *case i* who had an onset after *case j*, both of whom have sequences. The genetic  
142 distance between case *i* and *j* is obtained by comparing the first sequence available from  
143 case *i* and any sequence from *j* whose sampling time is closest to the first sequence from *i*.

144 In the illustration below, this would mean comparing sequence  $S_{i,1}$  to  $S_{j,2}$  to obtain genetic  
145 distance  $d_{gen}(i,j)$ . The phylogenetic analysis of *Agoti et al* [3] found that long shedding

146 episodes do not have drastically differing genetic sequences (<6 SNPs) as such it should not  
 147 make a significant difference whether we compare sequences forward ( $S_{i,1}$  to  $S_{j,2}$ ) or  
 148 backward ( $S_{i,1}$  to  $S_{j,1}$ ) in time.

149



150

151

152 If either one or both of the cases do not have sequences, then  $d_{gen}(i,j)$  is randomly selected  
 153 from the set of all pair-wise genetic distances from the specific genetic cluster. For cases  
 154 with sequence data,  $d_{gen}(i,j)$  is fixed, but for cases where one or both is missing sequences,  
 155  $d_{gen}(i,j)$  changes every time the likelihood is calculated to reflect uncertainty. In this way  
 156 only pairs of cases with sequence data contribute definitive genetic information to the  
 157 parameter inference algorithm while the rest will not. We use nucleotide differences as the  
 158 distance  $d_{gen}(i,j)$ . Once we have  $d_{gen}(i,j)$ , we then use this to obtain a genetic weight for the  
 159 probability of a transmission event given by  $P_{j \rightarrow i} = \exp^{-d_{gen}(i,j)*\vartheta}$  where  $\vartheta$  is the rate of  
 160 exponential decay and is estimated along with other model parameters. This function form  
 161 results in a negative exponential relationship between the genetic weight and the genetic  
 162 distance between a pair of cases.



163 **A2. Further details of the transmission model**

164 The rate of exposure

165 In our model, an individual can get infected by someone they share a household with or  
166 from a source outside the household, resulting in a two-component rate of exposure: a  
167 within household exposure component and a community exposure component.

$$\lambda_{i,c}(t) = [\textit{baseline household rate of exposure} * \textit{number of infectious household contacts}(t)]$$

+

168  $[\textit{baseline community rate of exposure} * \textit{number of infectious community contacts}(t)]$

169

$$\lambda_{i,c}(t) = \left( \eta * \sum_{\substack{j \in \textit{infectious} \\ \textit{household} \\ \textit{contact}}} I_{j,c}(t) \right) + \left( \varepsilon * \sum_{\substack{j \in \textit{infectious} \\ \textit{community} \\ \textit{contact}}} I_{j,c}(t) \right)$$

170

171 The number of infectious household contacts is observed in the data. Though there are  
172 cases from different households in the data, the sample in the study is small relative to the  
173 number of households in the community, as such the true number of infectious community  
174 contacts is unknown. We further split the community rate of exposure into two  
175 components: exposure from sampled neighbours and exposure from unknown sources  
176 represented by a time varying function  $f_d(\mathbf{t})$ . We assumed that the rate of exposure from a  
177 sampled neighbour is dependent on the spatial distance between individuals. The per capita  
178 rate of exposure now takes the form:

$$\lambda_{i,c}(t) = [\textit{baseline household rate of exposure} * \textit{number of infectious household contacts}(t)]$$

$$+$$

$$\{\textit{baseline community rate of exposure} * [\textit{number of infectious neighbour contacts}(t) +$$

$$\textit{background community function}(t)]\}$$

179

180

$$\lambda_{i,c}(t) = \left( \eta * \sum_{\substack{j \in \textit{infectious} \\ \textit{household} \\ \textit{contact}}} I_{j,c}(t) \right) + \varepsilon * \left( \left( \sum_{\substack{j \in \textit{infectious} \\ \textit{sampled} \\ \textit{neighbour} \\ \textit{contact}}} I_{j,c}(t) \right) + f_c(t) \right)$$

182

183 We extended this basic formulation to include our assumptions on RSV natural history and  
 184 explore factors that could influence the rate of exposure such as household size. In detail,  
 185 we present the model by specifying the rate of exposure to a particular RSV cluster  $c$  acting  
 186 on a susceptible person  $i$  from household  $h$  at time  $t$ , denoted  $\lambda_{i,h,c}(t)$  as:

187

$$\lambda_{i,h,c}(t) = S_{i,g}(t) \left[ M_{i,h}(t) \sum_{j \neq i} HH\_Rate_{h,c,j \rightarrow i}(t) + Comm\_Rate_{i,c}(t) \right] \quad \dots \quad (Eq \ A2.1)$$

189 Where:

190  $S_{i,g}(t)$  is the factor modifying exposure by recent group specific infection history, age and  
 191 group specific shedding status at time  $t$  given by:

192

$$S_{i,g}(t) = \exp \left( \phi_{Y,hist}(Infection\_History_i(t)) + \phi_{X,age}(Age\_group_{S,i}) \right)$$

$$+ \phi_{W,curr}(Shedding\_status_i(t))$$

194

195

196  $HH\_Rate_{h,c,j \rightarrow i}(t)$  is the cluster specific within household exposure rate from infectious  
197 individual  $j$  present in the household at time  $t$ , and is given by:

198

$$\begin{aligned} 199 \quad & HH\_Rate_{h,c,j \rightarrow i}(t) \\ 200 \quad & = \eta_g \times \psi_H(Household\_size_i) \times \psi_{I,inf}(Infectivity_{j,h,c}(t)) \times M_{j,h}(t) \end{aligned}$$

201

202  $Comm\_Rate_{i,c}(t)$  is the cluster specific community (external to the household) exposure  
203 rate given by:

204

$$205 \quad Comm\_Risk_{i,c}(t)$$

$$206 \quad = \varepsilon_g$$

$$\begin{aligned} 207 \quad & \times \psi_{E,age}(Age\_group_{E,i}) \left( \left( M_{i,h}(t) \sum_{\substack{j \neq i, j \text{ not in} \\ i's \text{ house}}} Sampled\_Neighbour\_Risk_{h,c,j \rightarrow i}(t) \right) \right. \\ 208 \quad & \left. + f_c(t) \right) \end{aligned}$$

209 Where:

210  $Sampled\_Neighbour\_Risk_{h,c,j \rightarrow i}(t)$  is the cluster specific exposure rate from sampled  
211 infectious individual  $j$  present in a neighbouring household at time  $t$ , and is given by:

$$212 \quad Sampled\_Neighbour\_Risk_{h,c,j \rightarrow i}(t) = \psi_{I,inf}(Infectivity_{j,h,c}(t)) \times K(d_{i,j}, \kappa) \times M_{j,h}(t)$$

213 The parameter  $\kappa$  is the rate of exponential decay for the spatial distance kernel given by

$$214 \quad K(d_{i,j}, \kappa) = e^{-\kappa * d_{i,j}}.$$

215 The background community function

216 We defined a background cluster-specific rate of exposure,  $f_c(t)$ , which affects susceptible  
217 individuals outside their household. This background function allows for introduction of new  
218 transmission clusters. The function form for a cluster  $c$  at time  $t$  is given as

$$219 \quad f_c(t) = \delta + \sum_{\substack{i \text{ shedding} \\ \text{RSV cluster } c}} e^{(t-\tau_{i,c})\beta}$$

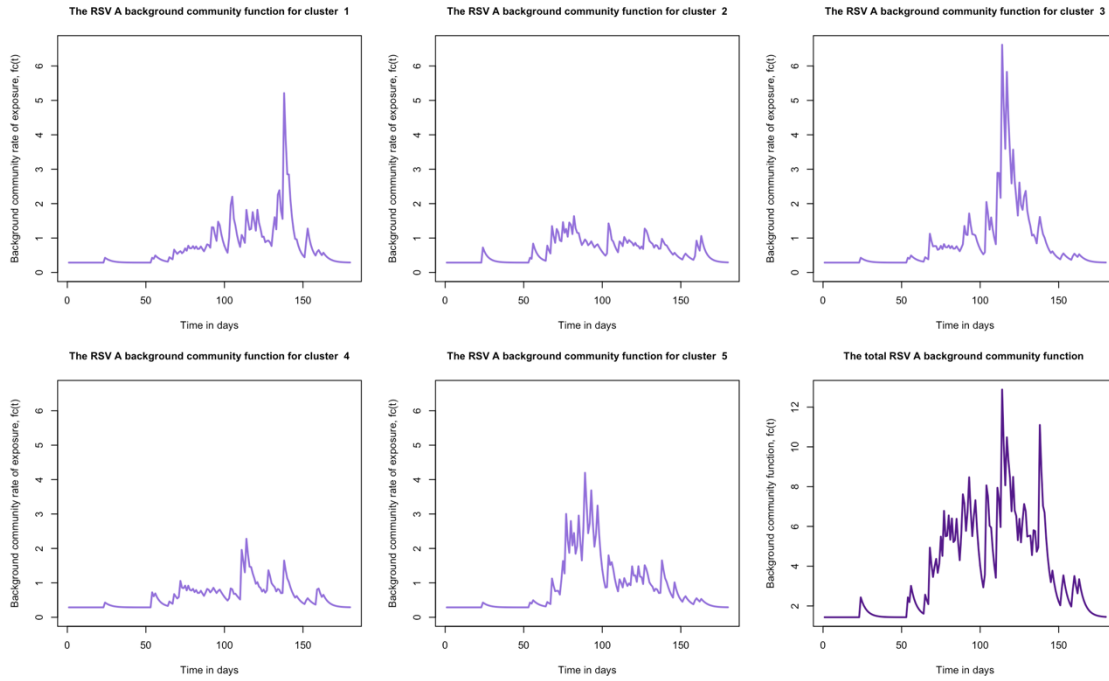
220 Where  $\delta$  is the basic risk prior to any observed onsets and  $\beta$  is the rate of exponential decay  
221 related to the time since onset of a case shedding cluster type  $c$ ,  $\beta$  is a measure of the rate  
222 at which the cluster might disappear from the community and  $\tau_{i,c}$  is the onset time of RSV  
223 cluster type  $c$  by person  $i$ . The parameters  $\delta$  and  $\beta$  are not cluster or group specific. The sum  
224 of the cluster specific curves has to add up to the group specific curve, otherwise using  
225 clusters could lead to an over or under representation of the background community  
226 exposure rate. To ensure that  $\sum f_c(t) = f_g(t)$  we need to normalize the cluster level curves  
227 such that their sum adds up to the group level curve. The equation for the normalized  
228 function  $\hat{f}_c(t)$  is given as:

$$229 \quad \hat{f}_c(t) = \left( \delta + \sum_{\substack{i \text{ shedding} \\ \text{RSV cluster } c}} e^{(t-\tau_{i,c})\beta} \right) \times \left( \sum_{c \in C'} \left( \delta + \sum_{\substack{i \text{ shedding} \\ \text{RSV cluster } c}} e^{(t-\tau_{i,c})\beta} \right) \right)$$

230 Where  $C'$  is the set of all clusters in a given RSV group.

231

232 An example of the shapes of the background community rate of exposure curves is shown in  
233 Figure A2. 1 for the 5 clusters in RSV A and Figure A2. 2 for the 7 clusters in RSV B.

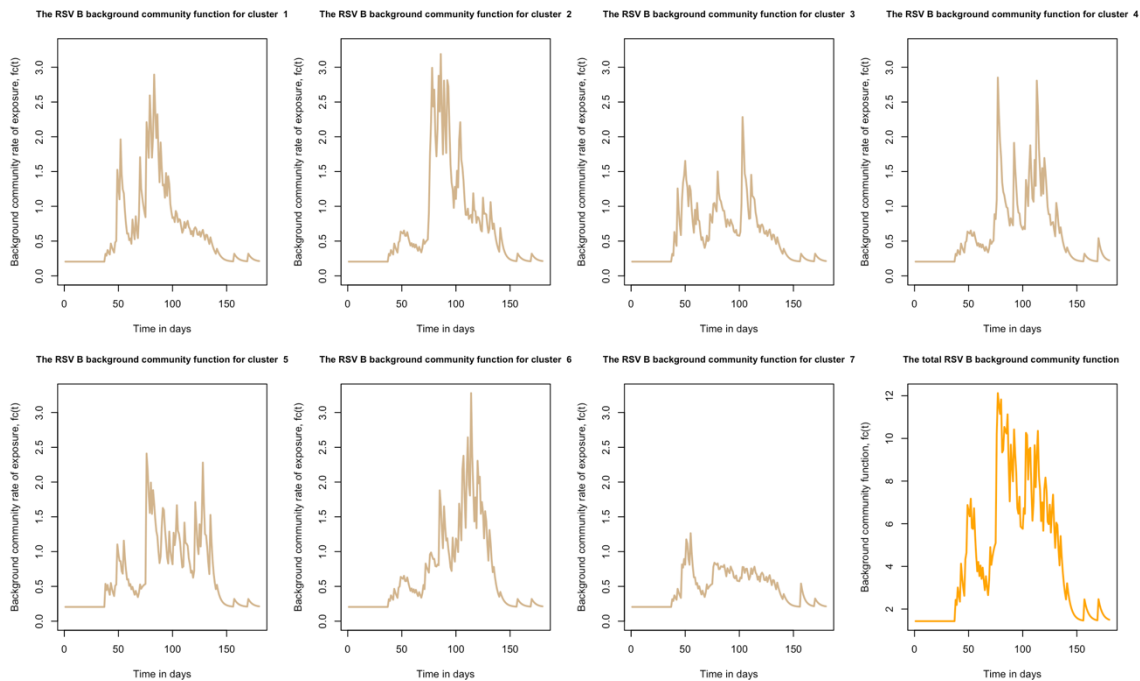


234

235 **Figure A2. 1: The background cluster-specific rate of exposure curves for RSV A. The**

236 normalized  $f_c(t)$  curves are shown for the 5 different clusters and the group.

237



238

239 **Figure A2. 2: The background cluster-specific rate of exposure curves for RSV B. The**

240 normalized  $f_c(t)$  curves are shown for the 7 different clusters and the group.

241 Table A2. 1 lists all the parameters in the model and gives a brief description. Despite  
 242 identifying the infection pathogen at the cluster level, we do not have any cluster-specific  
 243 parameters in the model.

244

245 **Table A2. 1: Model parameters and their descriptions**

Parameter (symbol)	Parameter (name)	Description
$\phi_Y$	<i>Prev.hom</i> , <i>Prev.het</i>	Coefficients modifying susceptibility to infection by a particular RSV group depending on infection history. <i>Prev.hom</i> estimates the effect of a previous homologous group infection, while <i>Prev.het</i> estimates the effect of a previous heterologous group infection
$\phi_X$	<i>Sus.age.2</i> , <i>Sus.age.3</i> , <i>Sus.age.4</i>	Coefficients modifying susceptibility to RSV depending on age. <i>Sus.age.2</i> estimates the effect being in age group 1-4 years, <i>Sus.age.3</i> the effect of group 5-14 and <i>Sus.age.4</i> of group $\geq 15$ relative to group $< 1$ year.
$\phi_W$	<i>Curr.het</i>	Coefficient modifying susceptibility to a particular RSV group based on shedding status of the heterologous group type
$\eta_g$	<i>HH.rsv.a</i> , <i>HH.rsv.b</i>	Baseline rate of within household exposure by RSV group, per person per day.
$\psi_H$	<i>HH.size</i>	Coefficient modifying the amount of within household exposure by household size. <i>HH.size</i> estimates the effect

		of being in a large household(>8 inhabitants) relative to a small one
$\vartheta$	<i>Gen.rate</i>	For $P_{j \rightarrow i} = \exp^{-d_{gen}(i,j)*\vartheta}$ the genetic distance kernel giving the genetic weight on probability of transmission, <i>Gen.rate</i> is the rate of exponential decay.
$\psi_I$	<i>Low.Sym</i> <i>High.Sym</i>	Coefficients modifying infectiousness by viral load and symptom status. Relative to being asymptomatic, <i>Low.Sym</i> estimates the effect of shedding low viral load and being symptomatic and <i>High.Sym</i> the effect of shedding high viral load and being symptomatic
$\varepsilon_g$	<i>Comm.rsv.a</i> <i>Comm.rsv.b</i>	Baseline rate of community exposure by RSV group, per person per day.
$\psi_E$	<i>Exp.age.2</i> <i>Exp.age.3</i>	Coefficients modifying the rate of community exposure by age group. <i>Exp.age.2</i> estimates the effect being in age group 1-4 years and <i>Exp.age.3</i> the effect of group $\geq 5$ , relative to the <1-year age group
$\kappa$	<i>Dist.rate</i>	The rate of exponential decay for the spatial distance kernel given by $K(d_{i,j}, \kappa) = e^{-\kappa*d_{i,j}}$
$\delta, \beta$	<i>Delta</i> , <i>Beta</i>	For the cluster specific background community function given by $f_c(t) = \delta + \sum_{\substack{i \text{ shedding} \\ \text{RSV cluster } c}} e^{(t-\tau_{i,c})\beta}$

*Delta*( $\delta$ ) is the basic risk and *Beta*( $\beta$ ) is the rate of exponential decay related to the time since onset of a case shedding cluster type  $c$ .

246

247 Following from the rate of exposure is the probability of exposure to cluster  $c$  given an  
 248 exposure event has occurred, expressed as:

249 Probability of exposure = *prob(any exposure event)* \* *prob(exposure to cluster  $c$ )*

$$250 \quad \alpha_{i,h,c}(t) = (1 - \exp^{-\sum_{c'} \lambda_{i,h,c}(t)}) * \left( \frac{\lambda_{i,h,c}(t)}{\sum_{c'} \lambda_{i,h,c}(t)} \right) \quad \dots \text{ (Eq A2.2)}$$

251

252 Where  $C'$  is the set of all clusters in a given RSV group.

253 This formulation factors in the fact that on any given day, an individual can only be shedding  
 254 virus from a single cluster, in the respective group. The clusters are therefore competing for  
 255 susceptible hosts. Exposure events are mutually exclusive and distributed according to a  
 256 multinomial distribution. We thus have

$$257 \quad \left[ \text{prob(no exposure)} + \sum_{\text{All clusters}} \text{prob(exposure to cluster } c) \right] = 1$$

258

259

260 Assuming that the duration of latency can range from 0 to 5 days with probabilities [0,  
 261 0,0.33,0.33,0.25,0.083] [4], we then have the following probability of onset at time  $t$  given  
 262 no onsets or shedding until  $t$ :

$$263 \quad p_{i,h,c}(t) = \sum_{l=0}^L \theta_l \alpha_{i,h,c}(t-l) \quad \dots \text{ (Eq A2.3)}$$



264 Where  $L$  is the maximum latency period and  $\theta_l$  is the probability that the latency period is  
 265 exactly  $l$  days. In this way, the genetic clusters are used together with the spatial/social  
 266 clusters (households) and the latency distribution (which implicitly works based on temporal  
 267 clusters) to make joint inference on transmission parameters.

268

269 The Likelihood

270 Since the model is focused on the determinants of infection onset process, the data whose  
 271 likelihood we are interested in is the onset data. Given the model described, the likelihood  
 272 of an individual's observed cluster  $c$  data is the probability of all the onsets, and days of no  
 273 onsets where the individual was at risk of infection, i.e. not shedding RSV cluster  $c$ . For a  
 274 particular cluster, this follows a Bernoulli distribution with probability  $p_{i,h,c}(u)$ .

275

276 For  $i$  with no onset of type  $c$ :

$$277 \quad L_{i,c} = \prod_{t=1}^T [1 - p_{i,h,c}(t)]$$

278 Where  $T$  is the end of the observation period.

279 For  $i$  with an onset of type  $c$ , the likelihood is give as:

$$280 \quad L_{i,c} = \left[ \left( \prod_{u \in \text{Onsets}_{i,h,c}} p_{i,h,c}(u) \right) * \left( \prod_{a \in \text{AtRisk}_{i,h,c}} (1 - p_{i,h,c}(a)) \right) \right]$$

281 In this instance, to factor in the genetic data we modify the rate of exposure given in (Eq

282 4.1) such that:

$$283 \quad HH\_Risk_{h,c,j \rightarrow i}(t)$$

$$284 \quad = \eta_g \times \psi_H(\text{Household\_size}_i) \times P_{j \rightarrow i} \times \psi_{I,inf}(\text{Infectivity}_{j,h,c}(t))$$

$$285 \quad \times M_{j,h}(t)$$

286

$$287 \quad \text{Sampled\_Neighbour\_Risk}_{h,c,j \rightarrow i}(t) = P_{j \rightarrow i} \times \psi_{l,c,j}(t) \times K(d_{i,j}, \kappa) \times M_{j,h}(t)$$

288 With this formulation, the genetic components of the model are dependent on the  
289 epidemiological in that they are not expressed independently in the likelihood function as is  
290 the case with modular approaches such as the kind implemented in the Outbreaker  
291 package[5,6]. We introduce  $P_{j \rightarrow i}$  into the rate of exposure equation as opposed to directly  
292 into the likelihood because for a given case, we are not making direct inference on the  
293 source of infection or the exact date of exposure: we consider all likely dates and sources  
294 given the latency distribution.

295

296 The total likelihood is thus given by the product of  $L_{i,c}$  over all the genetic clusters and

297 individuals in the data

298

$$299 \quad L = \prod_i \left[ \prod_c \left[ \left( \prod_{u \in \text{Onsets}_{i,h,c}} p_{i,h,c}(u) \right) * \left( \prod_{a \in \text{AtRisk}_{i,h,c}} (1 - p_{i,h,c}(a)) \right) \right] \right]$$

300

### 301 **A3. Further details of the adaptive MH-MCMC algorithm**

302 We used Bayesian inference to obtain estimates of the model parameters  $\varphi = \{Prev.hom,$   
303 *Prev.het, Sus.age.2, Sus.age.3, Sus.age.4, Curr.het, HH.rsv.a, HH.rsv.b, HH.size, Gen.rate,*  
304 *Low.Sym, High.Sym, Comm.rsv.a, Comm.rsv.b, Exp.age.2, Exp.age.3, Dist.rate, Delta, Beta\}  
305 and the augmented data  $D_A$  given the observed data  $D$ . We assume that all the cases were  
306 observed but that for some of the cases, there is no information on the cluster id of the  
307 shedding episode, as such, the augmented data is the set of all shedding episodes whose  
308 cluster id was left unassigned by the imputation process previously described. These include  
309 cases that are part of household outbreaks with no genetic information and cases that are  
310 part of household outbreaks with more than one possible genetic cluster id. For cases that  
311 are part of an outbreak with no genetic information, a single cluster id is inferred for all the  
312 cases in the household outbreak.*

313

314 Bayesian inference results in an updated distribution of the parameter of interest (posterior  
315 distribution) given prior assumptions/knowledge of the parameter (prior distribution) and  
316 an expression giving the probability of a parameter value given data (likelihood) i.e.

317  $P(\varphi|D, D_A) \propto P(\varphi) \times L(\varphi|D, D_A)$ . Where there is no exact expression for the posterior  
318 distribution, numerical methods are used to find an approximation of the target  
319 distribution, adaptive MH-MCMC is a popular first step.

320

321 We specified the target distribution as  $p(\varphi|D, D_A) = P(D|D_A)L(\varphi|D, D_A)P(\varphi)$ ;  $P(D|D_A) =$   
322 probability of the observed data given the augmented data;  $L(\varphi|D, D_A) =$  the likelihood of  
323 the parameters given the observed and augmented data;  $P(\varphi) =$  the prior probability of the  
324 parameters. The augmented and observed data are independent and we have no

325 information to inform what the missing cluster ids could be, making every combination of  $D$   
326 and  $D_A$  equally likely. Consequently, we did not include  $P(D|A)$  when calculating the  
327 posterior probability. We used weakly informative priors in the form of a normal  
328 distribution with mean 0 and a standard deviation of  $\sim 3$  for the log of parameters. We  
329 initiated 3 chains and set the algorithm to start adapting the proposal distribution based on  
330 accepted parameters after 10000, 15000 and 10000 iterations respectively. Burn-in was  
331 assessed visually after which the results of the three concurrent chains were combined to  
332 infer the posterior distribution. The three chains were run for 250,000 iterations each.

333

### 334 Choice of proposal distributions for the parameters

335 For the parameter set  $\varphi$  we used a multivariate normal distribution as the proposal  
336 distribution. For iteration  $n$  in the chain a new set  $\varphi^*$  will be proposed such that  
337  $\varphi^* \sim Normal(\varphi^{n-1} | \Sigma)$ . The choice of the variance-covariance matrix  $\Sigma$  will determine the  
338 size of the space that is explored and how fast the MCMC chain converges. After a certain  
339 number of iterations,  $\Sigma$  was modified to ensure proper mixing. The modification was  
340 automated through an adaptive random walk MH-MCMC algorithm. There are several  
341 adaptation algorithms [7], we chose one that learns from the empirical distribution of  
342 values up to the  $(n-1)^{th}$  iteration to modify the  $\Sigma$  at iteration  $n$ . For samples  
343  $\{\varphi_1, \varphi_2, \varphi_3, \dots, \varphi_{n-1}\}$  in the MCMC chain so far, at iteration  $n$  the proposal density  $g(\cdot)$  is  
344 given by

345

$$346 \quad g_n(\cdot) = (1 - \varepsilon)N(\varphi^{n-1} | 2.38^2 \Sigma_{n-1} / d) + \varepsilon N(\varphi^{n-1} | 0.1^2 \Sigma_0 / d)$$

347

348 Where:

349  $\varepsilon =$  A small positive constant, chosen to be 0.05 as in [7].  
 350  $\Sigma_{n-1} =$  The empirical variance-covariance matrix derived from samples  
 351  $\{\varphi_1, \varphi_2, \varphi_3, \dots, \varphi_{n-1}\}$   
 352  $d =$  The dimension of the parameter set  
 353  $\Sigma_0 =$  The initial guess of the parameter variance-covariance matrix. This is usually a  
 354 diagonal matrix of variances.

355

356 This notation means for a fraction of the time  $(1 - \varepsilon)$ , the proposal distribution will be  
 357  $N(\varphi^{n-1} | 2.38^2 \Sigma_{n-1} / d)$  and the rest of the time it will be  $N(\varphi^{n-1} | 0.1^2 \Sigma_0 / d)$ . Prior to  
 358 adaptation beginning at iteration  $n$ , the proposal distribution at iteration  $k$  is given by

359  $g_k(\cdot) = N(\varphi^{k-1} | 0.1^2 \Sigma_0 / d)$

360

361 Pseudo algorithm for our implementation of MH-MCMC

362 For each MCMC chain

363 1. Set initial values for the parameters and assign cluster ids at random for the  
 364 outbreaks with no sequence information (uninformed outbreaks).

365 2. For every iteration  $n$

366 a. Update parameter values

367 i. Propose a new set of parameters by sampling from the proposal  
 368 distribution:  $\varphi^* \sim Normal(\varphi^{n-1} | \Sigma)$

369 ii. Calculate the acceptance probability  $\rho(\varphi^{n-1}, \varphi^*) =$

370 
$$\min \left\{ 1, \frac{p(\varphi^* | D, D_A^{n-1})}{p(\varphi^{n-1} | D, D_A^{n-1})} \right\}$$

371                   iii. If  $\rho(\varphi^{n-1}, \varphi^*) > r \sim \text{Uniform}(0,1)$  update  $\varphi^n = \varphi^*$  otherwise  $\varphi^n =$   
 372    $\varphi^{n-1}$

373           b. Update cluster id for a single uniformed outbreak

374                   i. Randomly select an uniformed outbreak from the set of uninformed  
 375   outbreaks, all with the same probability of being selected.

376                   ii. Given the present cluster id for the chosen outbreak  $C_r$ , randomly  
 377   select a new cluster id from the set of all possible clusters excluding  
 378    $C_r$ .

379                   iii. With  $C_s$  as the proposed cluster id, the proposed change to the  
 380   augmented data is accepted with probability

381   
$$\rho'(D_A^{n-1}, D_A^*) = \min \left\{ 1, \frac{p(\varphi^n | D, D_A^*)}{p(\varphi^n | D, D_A^{n-1})} \frac{|C_r|}{|C_s| + 1} \right\}$$

382                   Where  $|C_r|$  is the number of household outbreaks in  $C_r$  in the present  
 383   permutation of the augmented data  $D_A^{n-1}$  and  $|C_s|$  is the number of  
 384   household outbreaks in  $C_s$ .

385  
 386                   iv. If  $\rho'(D_A^{n-1}, D_A^*) > r' \sim \text{Uniform}(0,1)$  update  $D_A^n, D_A^*$  otherwise

387    $D_A^n, D_A^{n-1}$

388

389   The correction factor  $\frac{|C_r|}{|C_s|+1}$  is introduced into the acceptance ratio for a proposed change in  
 390   cluster id because the proposal distributions are not symmetric. For an update of cluster id  
 391   from  $C_s$  to  $C_r$ , the proposed change is uniformly distributed over the set of all household  
 392   outbreaks/cases in cluster  $C_s$  that are part of the augmented dataset. Conversely the  
 393   reverse move of a change of cluster id from  $C_r$  to  $C_s$  is uniformly distributed over the set of

394 all household outbreaks/cases in cluster  $C_r$  that are part of the augmented dataset. As such,  
395 the proposal distributions are dependent on the number of uniformed household outbreaks  
396 in each cluster.

#### 397 **A4. Further details of the HPTS**

398 Establishing the highest probability transmission source (HPTS)

399 Per case, we identified the transmission source that had the highest likelihood given the  
400 data and a parameter set  $\varphi^*$  sampled from the joint parameter posterior distribution  
401 (highest probability transmission source: HPTS). Consider a case  $i$ , with onset date  $T_i^O$ . Given  
402 our assumption of a maximum latency duration of 5 days, we define a time window where  
403 potential infection could have occurred. For each day in the time window, potential sources  
404 of infection are  $\{\Omega_i^1, \Omega_i^2 \dots \Omega_i^n\}$ . An infection source is assigned if it gives the highest value of  
405  $i$ 's likelihood defined as " the likelihood of  $i$ 's onset date, infection date and infection source  
406 given sample parameter set  $\varphi^*$ .

407

408 We modified the likelihood to establish the most likely infection source (HPTS) for every  
409 case. For a given case  $i$  infected with RSV cluster  $c$  within group  $g$ , there are three possible  
410 sources of infection ( $\Omega_i$ ), either a sampled housemate, a sampled neighbour or an unknown  
411 community source. The total rate of exposure is given as:

$$412 \quad \lambda_{i,h,c}(t) = S_{i,g}(t) \left[ M_{i,h}(t) \sum_{j \neq i} HH_{Rate_{h,c,j \rightarrow i}}(t) + Comm\_Rate_{i,c}(t) \right] \quad \dots (Eq A4.1)$$

413 Where (as in the main text):

414  $S_{i,g}(t)$  is the factor modifying exposure by recent group specific infection history, age and  
415 group specific shedding status at time  $t$

416  $Comm\_Rate_{i,c}(t)$  is the cluster specific community (external to the household) exposure  
417 rate.

418



419 The probability of exposure is =  $prob(\text{any exposure event}) * prob(\text{exposure to cluster } c)$

$$420 \quad \alpha_{i,h,c}(t) = (1 - \exp^{-\sum_{c'} \lambda_{i,h,c}(t)}) * \left( \frac{\lambda_{i,h,c}(t)}{\sum_{c'} \lambda_{i,h,c}(t)} \right) \quad \dots \text{ (Eq A4.2)}$$

421

422 For a given source of infection  $\Omega_i$  in the same household as  $i$ , the rate of exposure is given  
423 by:

424

$$425 \quad \lambda_{\Omega_i \rightarrow i,h,c}(t) = S_{i,g}(t) [M_{i,h}(t) \times P_{\Omega_i \rightarrow i} \times \eta_g \times \psi_H(\text{Household\_size}_i) \\ 426 \quad \times \psi_{I,inf}(\text{Infectivity}_{\Omega_i,h,c}(t)) \times M_{\Omega_i,h}(t)]$$

427

428 For  $\Omega_i$  not in the same household as  $i$  but among the sampled individuals, the rate of  
429 exposure is given by:

430

$$431 \quad \lambda_{\Omega_i \rightarrow i,h,c}(t) = S_{i,g}(t) \left[ \varepsilon_g \times \psi_{E,age}(\text{Age}_{group_{E,i}}) \times M_{i,h}(t) \times P_{\Omega_i \rightarrow i} \right. \\ 432 \quad \left. \times \psi_{I,inf}(\text{Infectivity}_{\Omega_i,h,c}(t)) \times K(d_{i,\Omega_i}, \kappa) \times M_{\Omega_i,h}(t) \right]$$

433

434 For  $\Omega_i$  an unknown source external to the household, the rate of exposure is given by:

435

$$436 \quad \lambda_{\Omega_i \rightarrow i,h,c}(t) = S_{i,g}(t) \left[ \varepsilon_g \times \psi_{E,age}(\text{Age}_{group_{E,i}}) \times f_c(t) \right]$$

437

438 The probability of transmission from a single source  $\Omega_i$  at time  $t$  thus becomes:

$$439 \quad Pr_{\Omega_i \rightarrow i,h,c}(t) = \frac{\lambda_{\Omega_i \rightarrow i,h,c}(t)}{\lambda_{i,h,c}(t)} \quad \dots \text{ (Eq A4.3)}$$

440

441 *The likelihood function*

442 The probability given in (Eq A4.1) is calculated for a time point  $t$  = exposure time of  
443 individual  $i$ ,  $t_i^E$ . This is not observed in the data, however, given our assumption on the  
444 latency duration, we can define a 6-day window of possibility. If case  $i$  had a shedding onset  
445 at time  $T_i^O$ , then the window for transmission is from day  $(T_i^O - 5)$  to  $(T_i^O - 0)$ . For each  
446 day in the window, potential sources are identified based on shedding status and for each  
447 combination of infection source  $\Omega_i$  and exposure date  $t_i^E$ , the likelihood is calculated using  
448 the formula below:

449

$$450 \quad L(\varphi|\{T_i^O, t_i^E, \Omega_i\}) = \alpha_{i,h,c}(t) * \left( \prod_{t_i \neq t_i^E} (1 - \alpha_{i,h,c}(t)) \right) * (\theta_l(T_i^O - t_i^E)) * \left( \frac{\lambda_{\Omega_i \rightarrow i,h,c}(t_i^E)}{\lambda_{i,h,c}(t_i^E)} \right)$$

451

452 The first part of the product is the probability of infection with cluster  $c$  at time  $t_i^E$ , the  
453 second part is the probability of escaping infection at any time  $t_i \neq t_i^E$ , the third is the  
454 probability of a latency duration of length  $(T_i^O - t_i^E)$  and the last term is the probability of  
455 transmission from source  $\Omega_i$  to  $i$ .

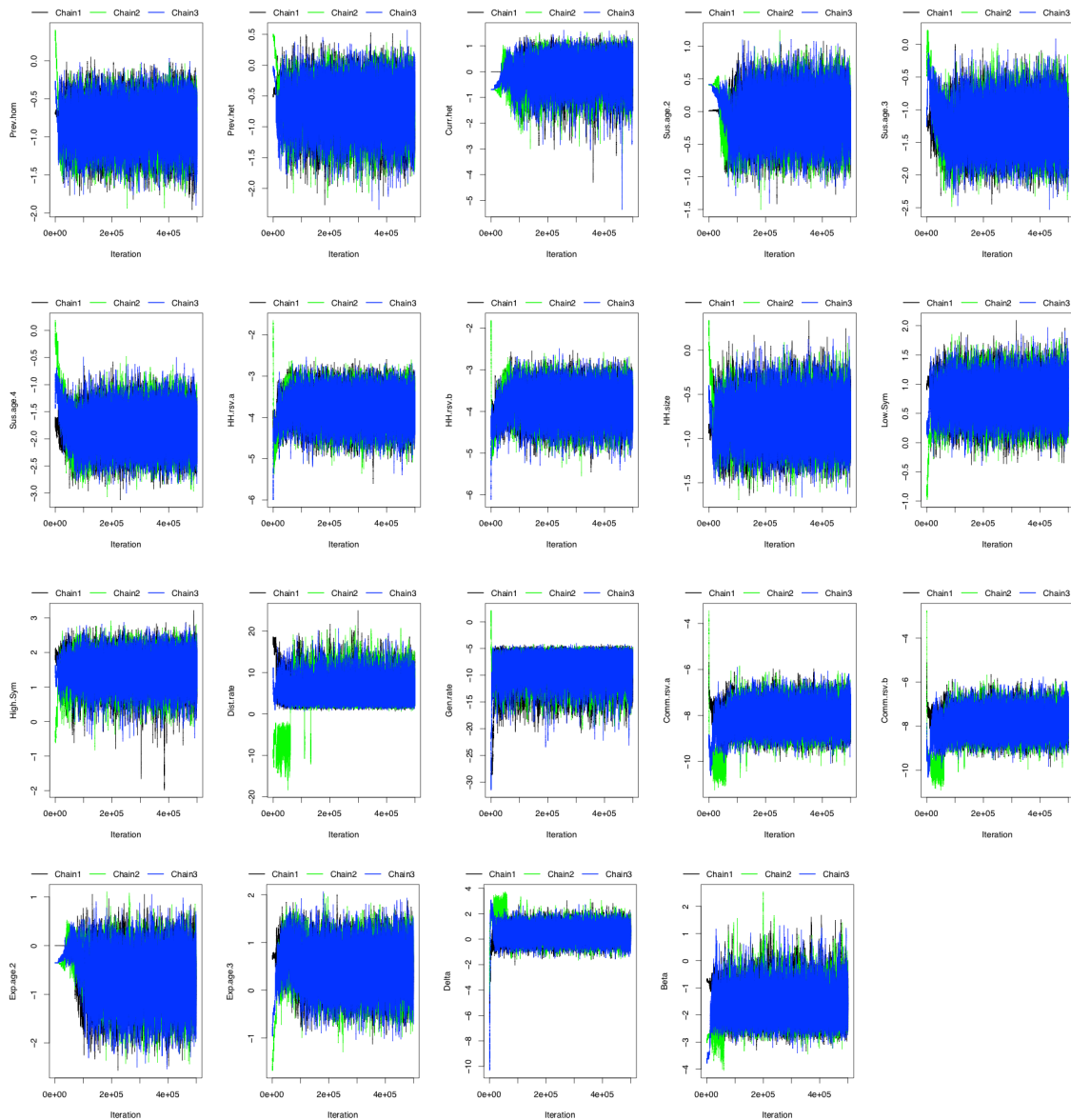
456

457 Given the likelihood, the highest-probability-source is chosen as the infection source that  
458 gives the highest value of the likelihood.

459 **A5. Extra results**

460 Parameter trace plots and convergence checks

461 Three MCMC chains were run, and the burn-in point assessed for each, after which, the  
462 remainder of the three chains were combined to give the posterior estimates for the  
463 parameters presented as median and 95% credible intervals. The figures below show the  
464 evolution of the parameter value with increasing number of iterations for the model with  
465 pathogen identification at the genetic cluster level (cluster model) and at the group level  
466 (group model).

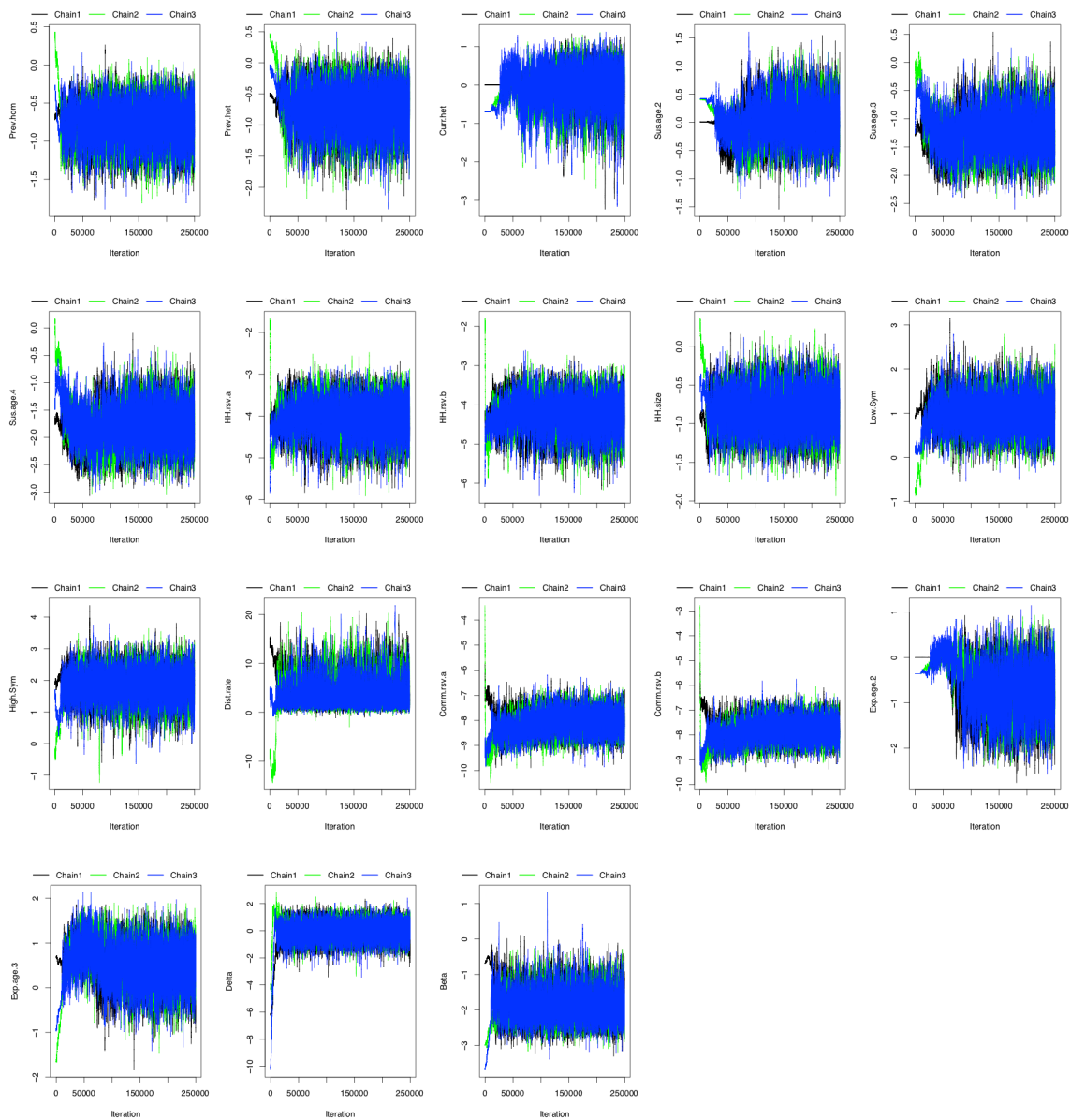


467

468 **Figure A5. 1: Trace plots of parameters in the cluster model.**

469 Three chains were initiated at different parameter values and these are shown in black  
 470 (Chain 1), green (Chain 2) and blue (Chain 3) lines. The x-axis shows the iteration number,  
 471 while the y-axis shows the log parameter value.

472



473

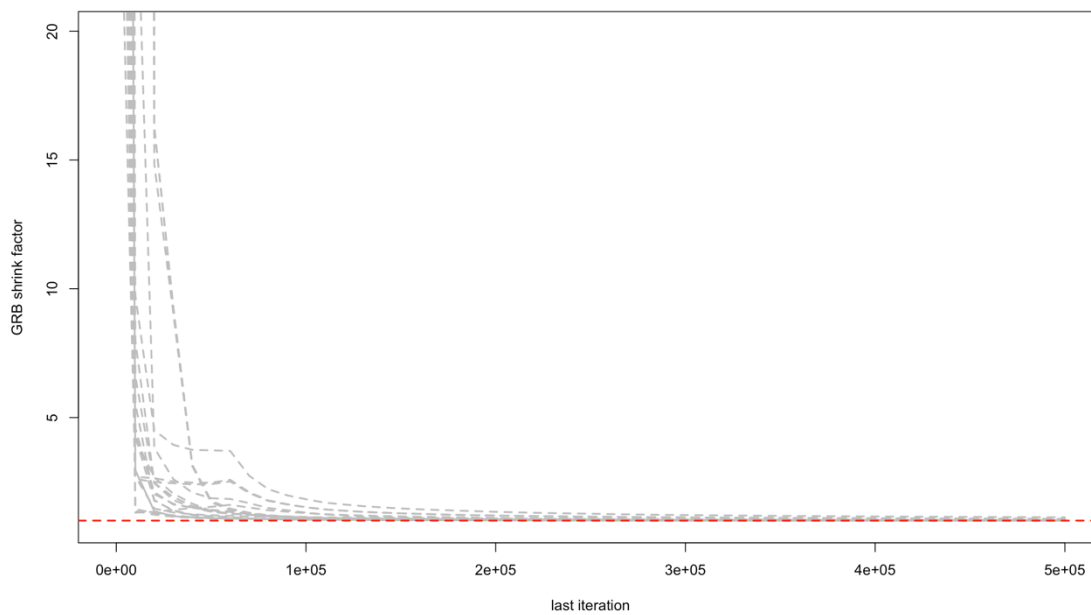
474 **Figure A5. 2: Trace plots of parameters in the group level data model.**

475 Three chains were initiated at different parameter values and these are shown in black  
 476 (Chain 1), green (Chain 2) and blue (Chain 3) lines. The x-axis shows the iteration number,  
 477 while the y-axis shows the log parameter value.

478

479 To confirm convergence observed in the trace plots, we calculated the Gelman-Rubin-  
 480 Brooks statistic and the effective sample size. When using the GRB statistic, convergence is  
 481 said to have occurred if the ratio of pooled/within chain variance is close to 1. The GRB

482 statistic assumes that the target distribution is Normal. The plot below shows the value of  
483 the GRB statistic as the number of iterations increases for each parameter. This is to check  
484 whether a value close to one was reached by chance or if the trend line had truly stabilized  
485 close to 1.  
486



487  
488 **Figure A5. 3: The evolution of the Gelman-Rubin-Brooks (GRB) statistic (shrink factor) as**  
489 **the number of iterations increases.**

490 Each grey line represents a model parameter in the cluster level data model and the dashed  
491 red line shows the value 1.

492  
493 The point estimated of the GRB and the values of the ESS after burn in are given in the table  
494 below.

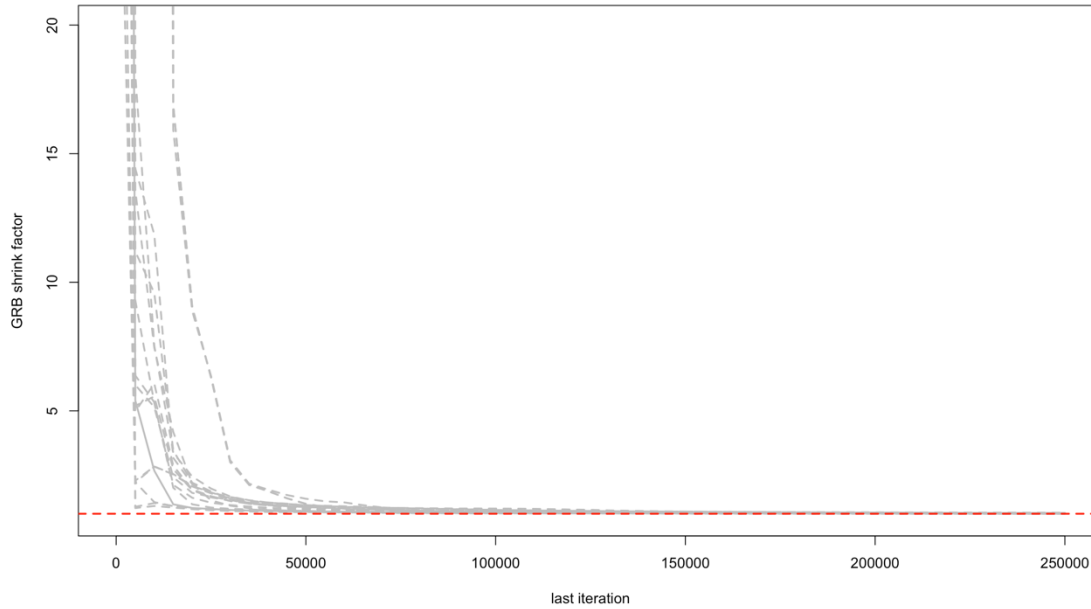
495

496 **Table A5. 1: The value of the GRB statistic (to 3 significant figures) and the ESS after burn-**  
 497 **in are shown for the parameters in the cluster level data model.**

Parameter	Point estimate	ESS
	GRB statistic	
Prev.hom	1	10607
Prev.het	1	10073
Curr.het	1.01	7131
Sus.age.2	1.01	9154
Sus.age.3	1.02	9771
Sus.age.4	1.02	10384
HH.rsv.a	1	9476
HH.rsv.b	1.01	9765
HH.size	1	10147
Low.Sym	1.02	9987
High.Sym	1.01	9774
Dist.rate	1.16	10455
Gen.rate	1.04	10436
Comm.rsv.a	1.09	7847
Comm.rsv.b	1.09	7823
Exp.age.2	1	8432
Exp.age.3	1.01	9863
Delta	1.04	7908
Beta	1.03	6678

498 The mGRB is 1.07 and the mESS is 10008.

499



500

501 **Figure A5. 4: The evolution of the Gelman-Rubin-Brooks (GRB) statistic (shrink factor) as**  
502 **the number of iterations increases.**

503 Each grey line represents a model parameter in the group level data model and the dashed  
504 red line shows the value 1.

505

506 **Table A5. 2: The value of the GRB statistic (to 3 significant figures) and the ESS after burn-**  
507 **in are shown for the parameters in the group level data model.**

Parameter	Point estimate	ESS
	GRB statistic	
Prev.hom	1.01	3713
Prev.het	1.02	3978
Curr.het	1.07	2309



Sus.age.2	1.02	2998
Sus.age.3	1.03	3617
Sus.age.4	1.04	3694
HH.rsv.a	1.01	3426
HH.rsv.b	1.01	3361
HH.size	1.02	3673
Low.Sym	1.04	3957
High.Sym	1.03	3744
Dist.rate	1.07	3374
Comm.rsv.a	1.05	4069
Comm.rsv.b	1.05	4093
Exp.age.2	1.02	2858
Exp.age.3	1.02	3476
Delta	1.04	5331
Beta	1.04	3873

508 *The mGRB is 1.09 and the mESS is 4146.*

509

510 As a rule of thumb, a GRB of <1.1 is generally considered good, as such, it is safe to conclude

511 that there was convergence.

512

513

514 **Table A5. 3: Median and 95% credible intervals for parameters estimated using the model**  
 515 **with sequence data.**

Symbol	Description	Name	Median (95% Credible interval)
$\phi_Y$	Coefficients modifying susceptibility to infection by a particular RSV group depending on infection history. <i>Prev.hom</i> estimates the effect of a previous homologous group infection, and <i>Prev.het</i> the effect of a previous heterologous infection	<i>Prev.hom</i>	0.4328 (0.2665,
		<i>Prev.het</i>	0.6727) 0.5126 (0.2601, 0.8985)
$\phi_W$	Coefficient modifying susceptibility to a particular RSV group based on shedding status of the heterologous group type	<i>Curr.het</i>	0.9520 (0.2494, 2.262)
$\phi_X$	Coefficients modifying susceptibility to RSV by age. <i>Sus.age.2</i> estimates modification to group 1-4 years, <i>Sus.age.3</i> 5-15 years and <i>Sus.age.4</i> $\geq 15$ years relative to group $< 1$ year.	<i>Sus.age.2</i>	0.8804 (0.4997,
		<i>Sus.age.3</i>	1.616)
		<i>Sus.age.4</i>	0.2741 (0.1591, 0.4946)
			0.1562 (0.08867, 0.2852)

$\eta_g$	Baseline rate of within household exposure by RSV group, per person per day.	<i>HH.rsv.a</i> <i>HH.rsv.b</i>	0.02360 (0.0119, 0.04361) 0.02272 (0.01120, 0.04196)
$\psi_H$	Coefficient modifying the amount of within household exposure by household size for households of 8 or more relative to <8.	<i>HH.size</i>	0.4457 (0.2892, 0.6843)
$\psi_I$	Coefficients modifying infectiousness by viral load and symptom status. Relative to being asymptomatic, <i>Low.Sym</i> estimates the effect of shedding low viral load and being symptomatic and <i>High.Sym</i> the effect of shedding high viral load and being symptomatic	<i>Low.Sym</i> <i>High.Sym</i>	2.1 (1.214, 3.67) 4.437 (1.8, 8.959)
$\kappa$	The rate of exponential decay on the spatial distance kernel	<i>Dist.rate</i>	207.7 (7.819, 169100)
$\vartheta$	The rate of exponential decay on the genetic weight function.	<i>Gen.rate*</i>	0.0002631 (0.000001027, 0.003817)

$\varepsilon_g$	Baseline rate of community exposure by RSV group, per person per day.	<i>Comm.rsv.a</i> <i>Comm.rsv.b</i>	0.0003091 (0.0001198, 0.0008682) 0.0003849 (0.0001525, 0.001072)
$\psi_E$	Coefficients modifying the rate of community exposure by age group. <i>Exp.age.2</i> for 1-4 years and <i>Exp.age.3</i> for $\geq 5$ years, relative <1 year	<i>Exp.age.2</i> <i>Exp.age.3</i>	0.5311 (0.2179, 1.221) 1.64 (0.7705, 3.386)
$\delta, \beta$	Parameters for the cluster specific background community function.	<i>Delta</i> <i>Beta</i>	1.58 (0.5466, 4.693) 0.1929 (0.08315, 0.7321)

517 **A6. Model validation**

518 To validate the model, we simulated multiple epidemics and checked to see if the observed  
519 epidemic was captured by the range of simulated dynamics. In addition to comparing the  
520 time course of cases, we also looked at the total number of cases in an epidemic, the  
521 proportion of individuals with multiple onsets and the number of cases in the first and last  
522 week of the time period. These values from the data were compared to the range of  
523 simulated values to check that key aspects of the epidemic were being reproduced by the  
524 simulations.

525

526 The results of the model fitting are the posterior parameter distribution and corresponding  
527 augmented data for the cluster ids of cases with no genetic information. A simulation based  
528 on a set of parameter values will also be based on the corresponding augmented data which  
529 will be used to derive a complete set of shedding profiles from the observed data. A single  
530 shedding profile is a combination of duration of shedding, viral loads and symptom status,  
531 and genetic cluster. The simulation pseudo code per simulation is as follows:

532

- 533 1. Initiate system such that everyone one is susceptible to RSV.
- 534 2. At every time step keep track of the following variables:
  - 535 a. Exposure status (by RSV cluster)
  - 536 b. Shedding status by group
  - 537 c. Shedding status by genetic cluster
  - 538 d. Infectiousness status (combination of viral load and symptom status)
  - 539 e. Infection history (by RSV group)
  - 540 f. The background rate of exposure from the community

- 541 3. At every time step:
- 542 a. Update the background community function to reflect any new shedding
- 543 onsets
- 544 b. Calculate the cluster specific rate of exposure,  $\lambda_{i,h,c}(t)$ , as defined in the
- 545 main text.
- 546 c. Determine the number of group specific transmission events  $E_g$  where

547 
$$E_g = \text{Poisson} \left( \sum_{i \in S_{E_g}} P_{E_g,i} \right)$$

548  $S_{E_g}$  = set of all individuals susceptible to infection event  $E_g$ .

549  $P_{E_g,i}$  = probability of person  $i$  experiencing event  $E_g$

550 
$$P_{E_g,i} = \sum_{\substack{c = \text{clusters} \\ \text{in } g}} \left( \left( 1 - \exp^{-\sum_{c'} \lambda_{i,h,c}(t)} \right) * \left( \frac{\lambda_{i,h,c}(t)}{\sum_{c'} \lambda_{i,h,c}(t)} \right) \right)$$

551 Where  $\lambda_{i,h,c}(t)$  = rate at which person  $i$  is exposed to infection of

552 cluster type  $C$ .

- 553 d. Given the number of group specific transmission events, determine the
- 554 cluster id of each through weighted sampling. E.g. if  $E_g = 4$  and  $c = \{1,2,3\}$  are
- 555 the cluster ids in the group, the probability of a case being any one if the
- 556 three clusters is:

557 
$$\left\{ \frac{\lambda_{h,1}(t)}{\sum_{c'} \lambda_{h,c}(t)}, \frac{\lambda_{h,2}(t)}{\sum_{c'} \lambda_{h,c}(t)}, \frac{\lambda_{h,3}(t)}{\sum_{c'} \lambda_{h,c}(t)} \right\}, \text{ for } \lambda_{h,1}(t) = \sum_i \lambda_{i,h,c}(t)$$

- 558 e. Determine who experiences each cluster specific transmission event. For a
- 559 given event, order individuals capable of experiencing the event. For a given
- 560 person  $p$  to experience the event, the following inequality has to be satisfied.

561

562 
$$\sum_{i=1}^{i \leq p-1} P_{E_c,i} < \left( RAND \times \sum_{i \in S_{E_c}} P_{E_c,i} \right) \leq \sum_{i=1}^{i \leq p} P_{E_c,i}$$

563 Where:

564

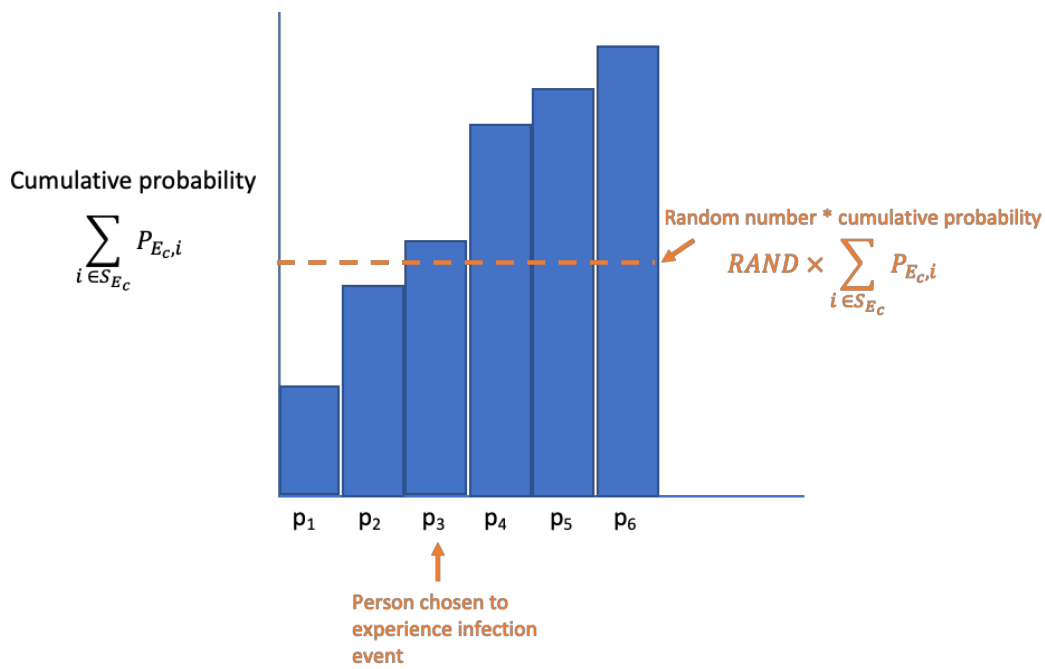
565 
$$P_{E_c,i} = \left( 1 - \exp^{-\sum_{c'} \lambda_{i,h,c}(t)} \right) * \left( \frac{\lambda_{i,h,c}(t)}{\sum_{c'} \lambda_{i,h,c}(t)} \right)$$

566  $S_{E_c}$  = all individuals susceptible to infection of cluster type c.

567  $RAND$  = a random number between (but not including) 0 and 1.

568

569 This is illustrated in the figure below.



570

571

572 Repeat this until the required number of events

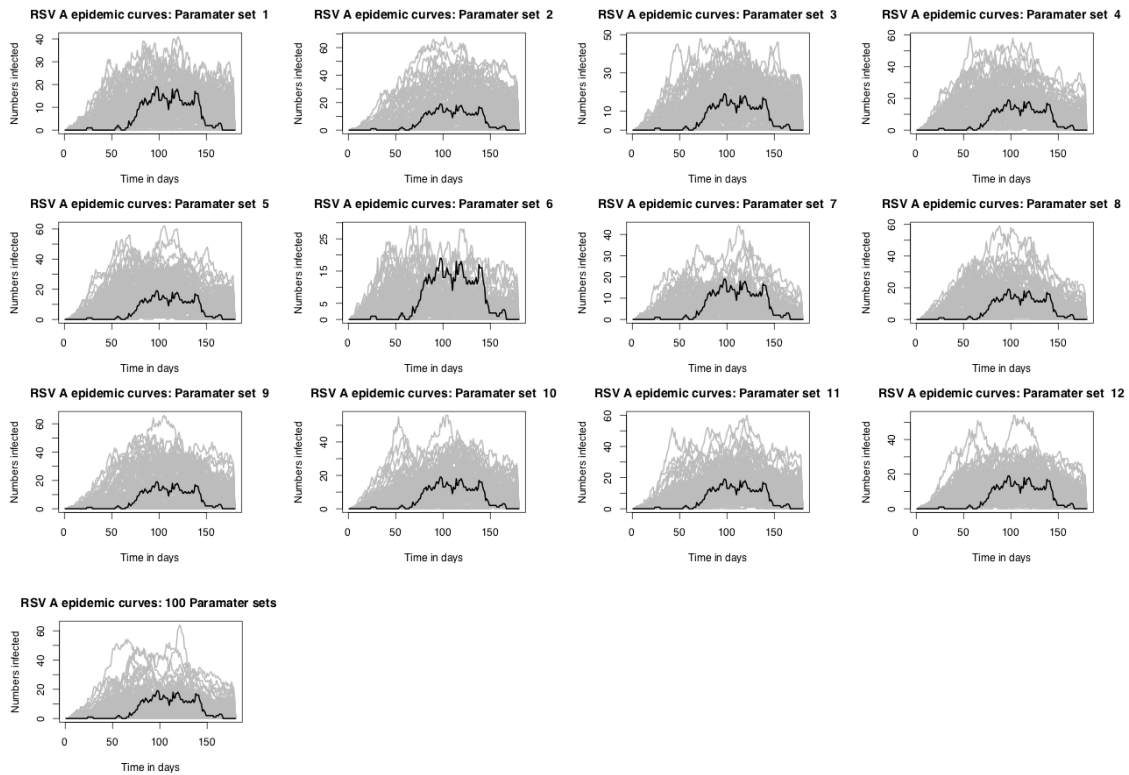
573

574 f. For each individual experiencing a transmission event, assign a latency

575 duration and shedding profile by sampling from the relevant empirical

576 distributions. The empirical latency distribution is the same as was used in  
577 estimating the parameters and is homogeneous for every individual.  
578 Shedding profiles are derived from the observed data and a combination of  
579 duration of shedding, viral loads and symptom status, and genetic cluster.  
580 The shedding profiles are grouped by age in the following 4 groups <1,1-5, 5-  
581 15 and  $\geq 15$  years. Once latency durations and shedding profiles have been  
582 assigned, the state variables for each individual are updated accordingly.  
583  
584 To explore how much variation there can be in the simulations from a single parameter set,  
585 a set of 12 parameter set samples were used, and for each set, 100 simulations were run,  
586 giving a total of 1200 simulations. We then sampled 100 parameter sets and run single  
587 simulations from each to explore between-parameter-set variation. The results of the  
588 simulations are presented in the form of epidemic curves and summary measures that are  
589 used to compare the main features of the outbreak.  
590

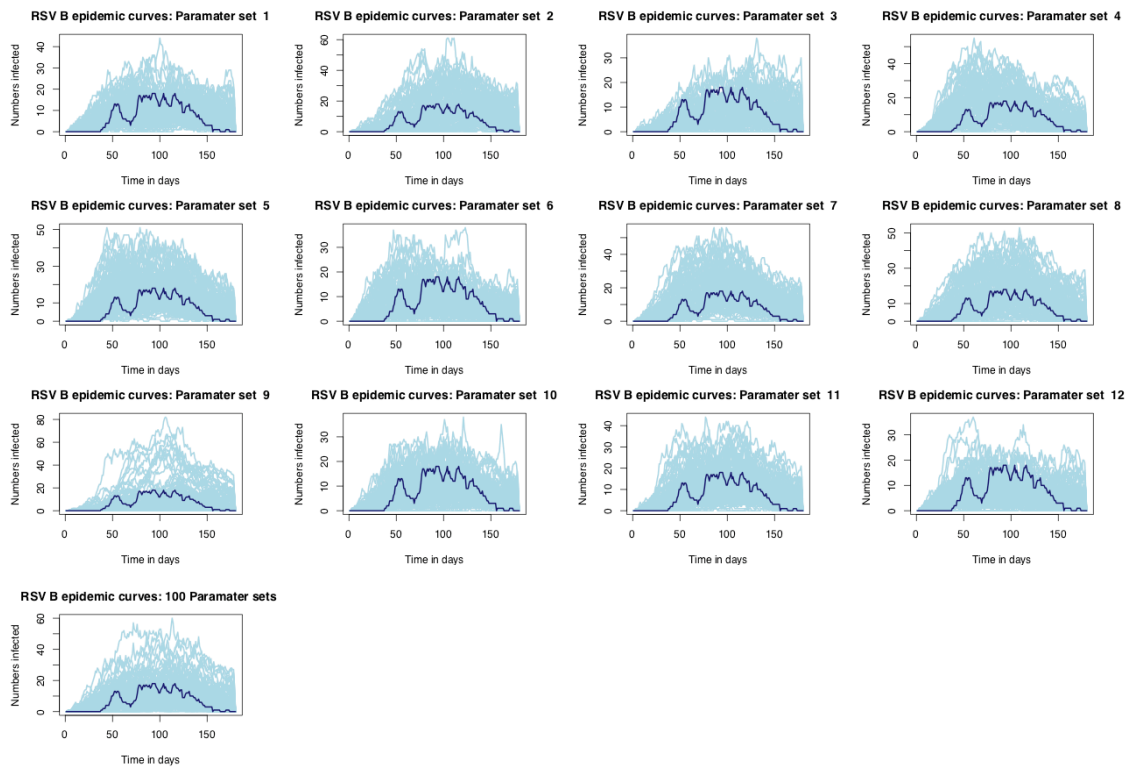




591

592 **Figure A6. 1: A comparison of simulated and observed data for RSV A.**

593 Each panel shows the results of 100 simulations from a single parameter set. The grey lines  
 594 show the simulated data while the black lines show the observed data. Time is shown on the  
 595 x-axis while the y-axis shows the total number of people who are shedding at a given point  
 596 in time.

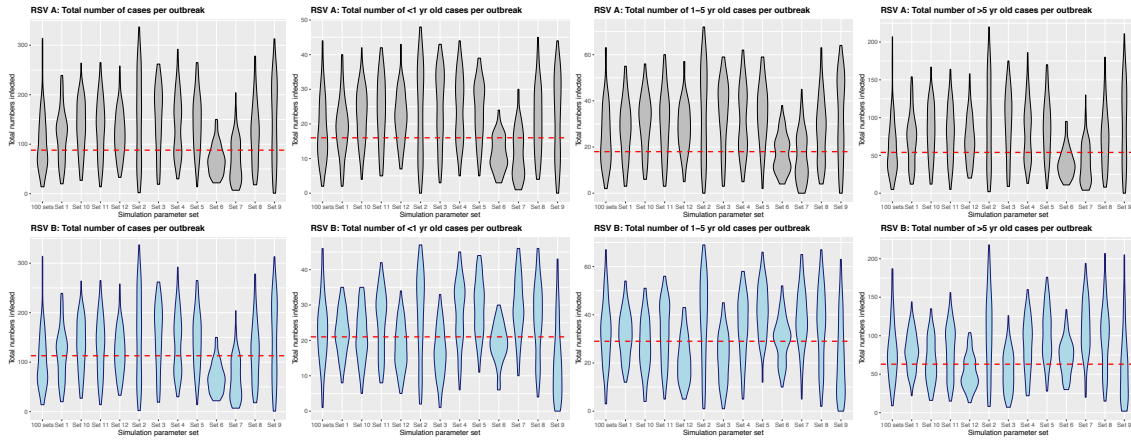


597

598 **Figure A6. 2: A comparison of simulated and observed data for RSV B.**

599 Each panel shows the results of 100 simulations from a single parameter set. The light blue  
 600 lines show the simulated data while the dark blue lines show the observed data. Time is  
 601 shown on the x-axis while the y-axis shows the total number of people who are shedding at  
 602 a given point in time.

603



604

605 **Figure A6. 3: Violin plots showing the distribution of the total number of people infected**

606 **in the simulations by RSV group and age.**

607 Each panel shows the distribution of the total numbers infected in the simulations run using

608 12 different parameter sets (violin plots) compared to the total number from the observed

609 data (dashed red line). The y-axis shows the total number and the x-axis is labeled by

610 parameter set used. Top row: RSV A results for all the cases (1<sup>st</sup> column), cases < 1 year old

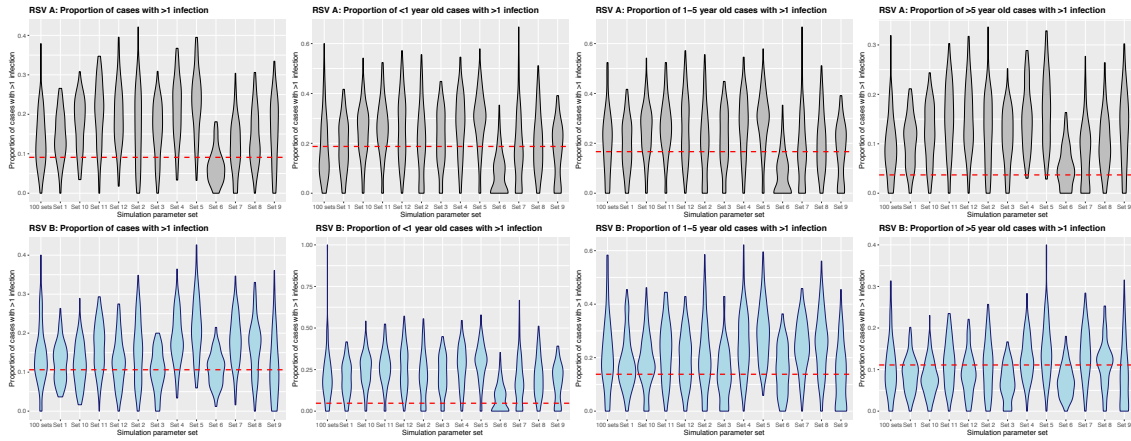
611 (2<sup>nd</sup> column), cases between 1-5 years old (3<sup>rd</sup> column) and cases > 5 years old (4<sup>th</sup> column).

612 Bottom row: RSV B results. Violin plots are a combination of box plots and density

613 distributions, the shapes should therefore be interpreted as density plots would while the

614 ranges should be interpreted as the tips of whiskers in a box and whisker plots.

615

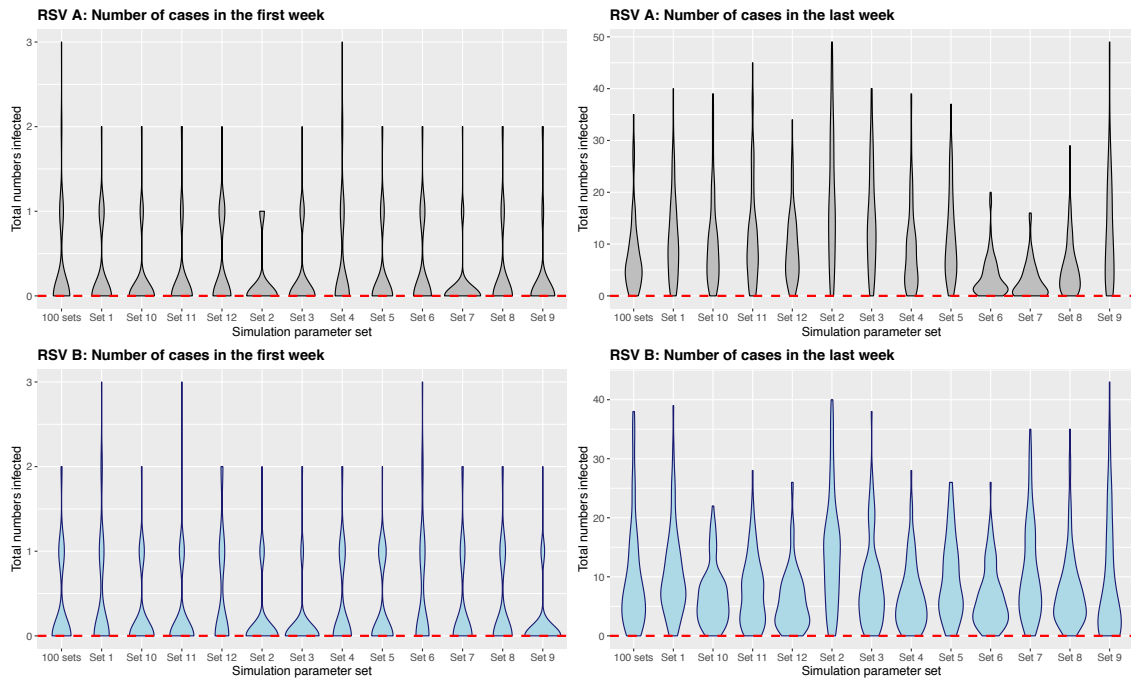


616

617 **Figure A6. 4: Violin plots showing the distribution of the proportion of cases that had**  
 618 **multiple onsets in the simulations by RSV group and age.**

619 Each panel shows the distribution of the proportion of cases that had multiple onsets in the  
 620 simulations run using 12 different parameter sets (violin plots) compared to the proportion  
 621 from the observed data (dashed red line). The y-axis shows the proportion and the x-axis is  
 622 labeled by parameter set used. Top row: RSV A results for all the cases (1<sup>st</sup> column), cases <  
 623 1 year old (2<sup>nd</sup> column), cases between 1-5 years old (3<sup>rd</sup> column) and cases > 5 years old (4<sup>th</sup>  
 624 column). Bottom row: RSV B results.

625



626

627 **Figure A6. 5: Violin plots showing the distribution of the number of cases in the first (1<sup>st</sup>**

628 **column) and last (2<sup>nd</sup> column) week of the observation/simulation period in the**

629 **simulations by RSV group.**

630 The y-axis shows the total number of people infected and the x-axis is labeled by parameter

631 set used. The dashed red line shows what was observed in the data, i.e. there were no cases

632 observed in the first and last week of the 180-day observation period.

633 **A7. Model modification to fit pathogen data identified at group resolution**

634 The null model is similar in structure to the model of sequence data presented in the main  
 635 text, however, there is no identification of the infecting pathogen at the cluster level, only at  
 636 the group level. The rate of exposure to a particular RSV cluster  $g$  acting on a susceptible  
 637 person  $i$  from household  $h$  at time  $t$ :

638

$$639 \quad \lambda_{i,h,g}(t) = S_{i,g}(t) \left[ M_{i,h}(t) \sum_{j \neq i} HH_{Rate}_{h,g,j \rightarrow i}(t) + Comm_{Rate}_{i,g}(t) \right] \quad \dots (Eq A7.1)$$

640 Where:

641  $S_{i,g}(t)$  is the factor modifying exposure by recent group specific infection history, age and  
 642 group specific shedding status at time  $t$  given by:

643

$$644 \quad S_{i,g}(t) = \exp \left( \phi_{Y,hist}(Infection\_History_i(t)) + \phi_{X,age}(Age\_group_{S,i}) \right. \\
 645 \quad \left. + \phi_{W,curr}(Shedding\_status_i(t)) \right)$$

646

647  $HH_{Rate}_{h,g,j \rightarrow i}(t)$  is the group specific within household exposure rate given by:

648

$$649 \quad HH_{Rate}_{h,g,j \rightarrow i}(t) \\
 650 \quad = \eta_g \times \psi_H(Household\_size_i) \times \psi_{I,inf}(Infectivity_{j,h,g}(t)) \times M_{j,h}(t)$$

651

652  $Comm_{Rate}_{i,g}(t)$  is the cluster specific community (external to the household) exposure  
 653 rate given by:

$$\begin{aligned}
654 \quad & Comm\_Rate_{i,g}(t) \\
655 \quad & = \varepsilon_g \\
656 \quad & \times \psi_{E,age}(Age\_group_{E,i}) \left( \left( M_{i,h}(t) \sum_{\substack{j \neq i, j \text{ not in} \\ i's \text{ house}}} Sampled\_Neighbour\_Rate_{h,g,j \rightarrow i}(t) \right) \right. \\
657 \quad & \left. + f_g(t) \right)
\end{aligned}$$

658 Where:

659

$$660 \quad Sampled\_Neighbour\_Rate_{h,g,j \rightarrow i}(t) = \psi_{I,g,j}(t) \times K(d_{i,j}, \kappa) \times M_{j,h}(t)$$

661

662 The background function  $f_g(t)$  is derived the same way  $f_c(t)$  is, as described in the main text.

663 Since we do not use genetic distances in this version of the model, we do not estimate  $\vartheta$  for

664  $P_{j \rightarrow i} = \exp^{-d_{gen}(i,j)*\vartheta}$  or  $P_{j \rightarrow i} = 1$  if  $d_{gen}(i,j) \leq \vartheta$ , 0 otherwise, making the total

665 number of parameters 17.

666

667 Following from the rate of exposure is the probability of exposure give by:

$$668 \quad \alpha_{i,h,g}(t) = (1 - \exp^{-\lambda_{i,h,g}(t)}) \quad \dots (Eq A7.2)$$

669

670

671 The probability of onset is given as:

$$672 \quad p_{i,h,g}(t) = \sum_{l=0}^L \theta_l \alpha_{i,h,g}(t-l) \quad \dots (Eq A7.3)$$

673 Where  $L$  is the maximum latency period and  $\theta_l$  is the probability that the latency period is  
674 exactly  $l$  days.

675

676 The likelihood for individual  $i$ 's data is given as:

$$677 \quad L_i = \prod_g \left[ \prod_{u \in U_{i,h,g}} p_{i,h,g}(u) \prod_{a \in A_{i,h,g}} (1 - p_{i,h,g}(a)) \right]$$

678 The total likelihood is thus given by the product of  $L_i$  over all the individuals in the data

679

$$680 \quad L = \prod_i \left[ \prod_g \left[ \prod_{u \in U_{i,h,g}} p_{i,h,g}(u) \prod_{a \in A_{i,h,g}} (1 - p_{i,h,g}(a)) \right] \right]$$

681



682 **A8. References**

- 683 1. Nolan T, Hands RE, Bustin SA. Quantification of mRNA using real-time RT-PCR. Nat  
684 Protoc. Nature Research; **2006**; 1(3):1559–1582.
- 685 2. Wathuo M, Medley GF, Nokes DJ, Munywoki PK. Quantification and determinants of  
686 the amount of respiratory syncytial virus (RSV) shed using real time PCR data from a  
687 longitudinal household study. Wellcome Open Res [Internet]. **2017**; 1(0):27. Available  
688 from: <https://wellcomeopenresearch.org/articles/1-27/v2>
- 689 3. Agoti CN. Genomic analysis of respiratory syncytial virus infections in households and  
690 utility in inferring who infects the infant. Sci Rep. **2019**; .
- 691 4. Lee FE, Walsh EE, Falsey AR, Betts RF, Treanor JJ. Experimental infection of humans  
692 with A2 respiratory syncytial virus. Antivir Res [Internet]. 2004/09/29. **2004**;  
693 63(3):191–196. Available from:  
694 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=C](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=15451187)  
695 [itation&list\\_uids=15451187](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=15451187)
- 696 5. Campbell F, Didelot X, Fitzjohn R, Ferguson N, Cori A, Jombart T. outbreaker2 : a  
697 modular platform for outbreak reconstruction. **2018**; 19(Suppl 11).
- 698 6. Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N. Bayesian  
699 Reconstruction of Disease Outbreaks by Combining Epidemiologic and Genomic Data.  
700 PLoS Comput Biol [Internet]. **2014**; 10(1):e1003457. Available from:  
701 <http://dx.plos.org/10.1371/journal.pcbi.1003457>
- 702 7. Roberts GO, Rosenthal JS. Examples of Adaptive MCMC. J Comput Graph Stat  
703 [Internet]. **2009**; 18(2):349–367. Available from:  
704 <http://www.tandfonline.com/doi/abs/10.1198/jcgs.2009.06134>

705