

1 **Title: Probabilistic reconstruction of measles transmission**
2 **clusters from routinely collected surveillance data.**
3

4

5 **Authors:**

6 Alexis Robert^{1,2}, Adam J. Kucharski^{1,2}, Paul A. Gastanaduy³, Prbasaj Paul⁴, Sebastian Funk^{1,2}

7 **Authors' affiliations:**

8 1. Centre for the Mathematical Modelling of Infectious Diseases, London School of Hygiene & Tropical
9 Medicine, Keppel Street, London, UK

10 2. Department of Infectious Disease Epidemiology, London School of Hygiene & Tropical Medicine,
11 Keppel Street, London, UK

12 3. Division of Viral Diseases, Centers for Disease Control and Prevention, Atlanta, Georgia

13 4. Division of Nutrition, Physical Activity, and Obesity, Centers for Disease Control and
14 Prevention, Atlanta, Georgia

15 **Abstract**

16 Pockets of susceptibility resulting from spatial or social heterogeneity in vaccine coverage can drive
17 measles outbreaks, as cases imported into such pockets are likely to cause further transmission and lead
18 to large transmission clusters. Characterising the dynamics of transmission is essential for identifying
19 which individuals and regions might be most at risk.

20 As data from detailed contact tracing investigations are not available in many settings, we combined
21 age, location, genotype, and onset date of cases in order to probabilistically reconstruct the importation
22 status and transmission clusters within a newly developed R package called *o2geosocial*.

23 We compared our inferred cluster size distributions to 737 transmission clusters identified through
24 detailed contact-tracing in the United States between 2001 and 2016. We were able to reconstruct the
25 importation status of the cases and found good agreement between the inferred and reference clusters.
26 The results were improved when the contact-tracing investigations were used to set the importation
27 status before running the model.

28 Spatial heterogeneity in vaccine coverage is difficult to measure directly. Our approach was able to
29 highlight areas with potential for local transmission using a minimal number of variables and could be
30 applied to assess the intensity of ongoing transmission in a region.

31 Introduction

32 Establishing who infected whom during an outbreak can help inform the design and evaluation of
33 control measures[1–5]. Transmission links can be reconstructed through contact tracing investigation,
34 whereby cases are asked their movements and contacts during their infectious period. Given that
35 contact-tracing investigations are not always carried out due to the logistical effort and cost involved,
36 inference methods have been developed to use epidemiological data to estimate the probability that a
37 transmission event occurred between any given pair of cases[6–12]. This makes it possible to establish
38 probabilistic transmission trees that link all observed cases.

39 Wallinga and Teunis first developed a likelihood-based estimation procedure to reconstruct
40 probabilistic transmission trees from a given distribution of generation times and observed symptom
41 onset dates of each case[2]. Since then, genomic, spatial or contact data have been used to supplement
42 the timing of symptoms, which helped identify determinants of transmission, mixing behaviour,
43 individual dispersion, evaluate control measures, anticipate future developments of outbreaks and study
44 viral evolutionary patterns[5,8,9,13–17].

45 As sequencing of pathogens has become more common, the use of such data to infer transmission trees
46 has increased. Methods developed to add genetic distance to a Wallinga-Teunis algorithm, where cases
47 with lower genetic distance are more likely to be grouped in the same transmission group, showed it
48 substantially increased the accuracy of the reconstructed transmission trees[8,18–21].

49 The utility of sequence data depends on the characteristics of the pathogen[22,23]. Based on N-450
50 sequence data, eight measles genotypes have been detected since 2009[24,25]; these genotype
51 designations are helpful in linking cases, as linked cases must be infected by virus of the same
52 genotype[25]; however, the diversity of measles genotypes is decreasing[26]. It has been suggested that
53 further sequencing the M-F non-coding region, or full genome sequencing, could help identify measles
54 virus transmission trees, but so far, extended sequencing during measles outbreaks has been
55 scarce[27,28]. In addition, the evolutionary rate of measles virus is very low[29], therefore, samples

56 from unrelated cases can be very close genetically and genetic sequences from measles cases are not
57 usually indicative of direct transmission links[27,28].

58 As measles is highly infectious, under-immunized communities (also called pockets of susceptibles)
59 resulting from local heterogeneity in vaccine coverage can lead to large, long-lasting outbreaks[30–34].
60 Detecting these pockets of susceptibles can be challenging, as historical local values of coverage
61 throughout a given country are rarely available. The size distribution of transmission trees resulting
62 from each importation during outbreaks (otherwise known as the cluster size distribution) will depend
63 both on individual factors (e.g. age of the imported case which might affect contact patterns) and
64 community factors (e.g. the history of coverage in the area)[35,36]. The size of a cluster can therefore
65 reflect the level of susceptibility of individuals directly and indirectly connected to the index case
66 [37,38].

67 Here we introduced a model combining age, location, genotype, and rash onset date of cases to
68 reconstruct probabilistic transmission trees. We chose these features to make the model applicable to a
69 wide range of settings as they are commonly reported and informative on transmission. We wrote the
70 R package *o2geosocial* to conduct inference on individual-level data using this model. It is based on
71 the package *outbreaker2* and is designed for outbreaks with partial sampling of cases, or uninformative
72 genetic sequences, such as measles outbreaks[9,39]. We used the likelihood of transmission links
73 between different cases to estimate their importation status. We compared the inferred importation
74 status and cluster size distribution to the transmission clusters identified via contact tracing during
75 measles outbreaks in the United States between 2001 and 2016.

76 **Methods**

77 **Presentation of the algorithm**

78 *Likelihood function and parameter definition*

79 We used a probabilistic model to infer the individual contribution to the log-likelihood L_i of every case
80 included in the list of cases.

$$81 \quad L_i(t_i, j, t_j, \theta) = \log(f(t_i - T_i)) + L_{ji}(t_i, t_j, \theta) \quad (1)$$

82 L_i was computed from $L_{ji}(\theta)$, the log-likelihood of case i being infected by case j as a function of
83 infection times t_i and t_j and model parameters θ , and the timing of t_i , the date of infection of i , relative
84 to the date of symptom onset T_i . We defined $f(t_i - T_i)$ as the probability density of observing T_i if
85 case i was infected at time t_i (i.e. f represents the distribution of incubation periods). The log-likelihood
86 of transmission L_{ji} was computed from five components reflecting the age group, genotype, location,
87 and inferred date of infection of cases i and j , and the generation time of the disease:

$$88 \quad L_{ji}(t_i, t_j, \theta) = \log\left(p(\kappa_{ji}|\rho) * w^{(\kappa_{ji})}(t_i - t_j) * a^{(\kappa_{ji})}(\alpha_i, \alpha_j) * G(g_i, g_{\tau_j}) * s^{(\kappa_{ji})}(r_i, r_j|a, b)\right) \quad (2)$$

89 We allowed for missing generations between cases due to an unreported individual, and κ_{ji} corresponds
90 to the number of generations between i and j . We calculated the temporal probability of transmission
91 between i and j from the number of days between the dates of infection of the two cases t_i and t_j and
92 the generation time of the disease $w(t)$. This probability of infection was quantified by $w^k(t_i - t_j, \kappa_{ji})$,
93 $w^{(\kappa)} = \prod_{\kappa} w$, where \prod is the convolution operator. We used an exponential distribution $p(\kappa_{ji}|\rho)$ to
94 quantify the probability of observing κ_{ji} missing generation between i and j from the conditional report
95 ratio ρ which quantifies the probability of missing generation between two connected cases in a cluster.
96 It does not correspond to the overall report ratio of an outbreak as entire missing clusters, or unreported
97 cases infected after the last case or before the ancestor of a cluster are not included in ρ . The “ancestor”
98 is the earliest identified case of a cluster.

99 $a(\alpha_i, \alpha_j, \kappa_{ji})$ was defined as the probability of transmission between age groups α_i and α_j . This
100 probability corresponds to the proportion of contacts to the age group α_i that originated from α_j and
101 can be deduced from studies such as Polymod[36]. We defined $G(g_i, g_{\tau_j})$ as the probability of
102 observing the pathogen genotype g_i in case i in the tree τ_j containing case j . There can only be one
103 measles virus genotype per transmission tree, or cases with unreported genotype.

$$G(g_i, g_{\tau_j}) = \begin{cases} 1 \text{ if } g_i \text{ unknown} \\ 1 \text{ if } g_{\tau_j} \text{ unknown} \\ 1 \text{ if } g_i \text{ and } g_{\tau_j} \text{ both known and } g_i = g_{\tau_j} \\ 0 \text{ otherwise} \end{cases}$$

105 $s(r_i, r_j, \kappa_{ij})$ was defined as the probability of connection from r_j to r_i , counties of residency of i and j .
106 We used a gravity model to quantify the connectivity of the different geographical units. In the simplest
107 form of the gravity approach, the number of connections between two counties k and l is proportional
108 to the product of the origin population m_k , the destination population m_l and a function of the distance
109 between k and l d_{kl} : $p_{kl} \propto f(d_{kl}, a) * m_k^b * m_l^c$, with a , b , and c parameters adjusting for the impact
110 of distance and population. From this definition, we deduced $s(k, l)$, the probability of transmission
111 from an individual from region k to another from region l :

$$s(k, l) = \frac{p_{kl}}{\sum_h p_{hl}} = \frac{f(d_{kl}, a) * m_k^b * m_l^c}{\sum_h f(d_{hl}, a) * m_h^b * m_l^c} = \frac{f(d_{kl}, a) * m_k^b}{\sum_h f(d_{hl}, a) * m_h^b}$$

113 Only the parameters a and b were required to compute the spatial probability of transmission. We used
114 the exponential gravity model ($f(d_{kl}, a) = e^{-a*d_{kl}}$)[40]. This approach showed good performance at
115 modelling short distance commuting, and was easy to parametrise[40–44].

116 In order to compute the log-posterior densities of the proposed trees, we summed the individual log-
117 likelihoods and added log-priors on the report ratio ρ , which quantified the percentage of cases in the
118 chains reported to the surveillance system; and the spatial parameters a and b (Table 1).

119 *Tree proposals*

120 We used a Metropolis Hastings algorithm with Markov chain Monte Carlo (MCMC) to sample from
121 the posterior distribution of parameters and the transmission trees. To do this, we developed a set of
122 proposal tree updates. These updates were accepted with acceptance probability as defined by the
123 Metropolis-Hastings algorithm[45]. We used eight types of tree proposal to ensure good mixing. Each
124 proposal conserved the overall number of trees, with a maximum of one unique genotype reported per
125 tree.

126 Five of the proposals had already been implemented in the *outbreaker2* package and were adapted to
127 this setting: i) change the number of generations between two cases; ii) change the conditional report
128 ratio ρ ; iii) change the time of infection; iv) change the infector of a case (if the case is not the ancestor
129 of a tree); v) swap infector-infectee (if none is the ancestor of a tree).

130 We added two proposals to change a and b , the spatial kernel parameters. For each proposal, the
131 probability of transmission between every geographical unit was re calculated with the new values.
132 Depending on the number of geographical units, this calculation considerably slowed down the
133 algorithm. Therefore, when a or b were estimated, we limited the maximal number of missing
134 generations to 1 ($\max(\kappa_{ji}) = 2$). Finally, the last proposal was designed to change the ancestor of the
135 tree whilst conserving the overall number of trees (Figure 1).

136 *Inference of importation status and cluster*

137 Unrelated measles cases stemming from different importations and different regions can be part of the
138 same dataset. Grouping cases and excluding unrealistic transmission links reduces the number of
139 possible trees and speeds up the MCMC runs. To do so, we listed each case's potential infectors using
140 three criteria: i) The potential infectors must be of the same genotype as the case, or have unreported
141 genotype, ii) The location of potential infectors must be less than γ km away from the case and iii) the
142 potential infectors must have been reported later than δ days before the case. This threshold should be
143 determined from the maximum plausible generation time of the disease. The spatial threshold γ should
144 be defined according to the relevance of long-distance transmissions. Cases with no potential infector
145 were considered as importations. Otherwise, they were grouped together with i) their potential infectors
146 and ii) cases with common potential infectors.

147 After grouping the cases, we estimated their importation status and the cluster size distribution using
148 two runs of MCMC (Figure 2). The first run was shorter and aimed at removing the most unlikely
149 connections among each group, as they can reflect unrealistic estimates for incubation periods or
150 generation times and corrupt the estimation of the date of infection. We defined a reference threshold λ ,
151 whereby if the individual value of log-likelihood L_i was worse than λ , then the connection between i

152 and their index was considered unlikely. In *Outbreaker2*, λ was a relative value, defined from a quantile
153 of the individual log-likelihoods. In *o2geosocial*, λ can be a relative value or an absolute value, chosen
154 from the number of components of the likelihood. For each sample saved from the short run, we
155 computed the number of unlikely connections n . If there was no iteration where all connections were
156 better than λ , $\min(n)$ new importations were added to the initial tree for the long run (Figure 2).
157 Finally, we ran a long MCMC chain and obtained samples from the posterior distribution. After
158 removing the burn-in period and thinning the chain, we deleted the unlikely transmission links in each
159 iteration and identified transmission clusters. Therefore, unlike the previous versions of *outbreaker2*,
160 the number of importations in each sample can vary and the individual probability of being an
161 importation can be computed (Figure 2).

162 **Validation case study: measles outbreaks in the United States between 2001 and 2016**

163 *Data*

164 To evaluate the performance of the model, we inferred the transmission clusters from a dataset that also
165 included information on whether measles cases were part of a cluster based on contact tracing
166 investigations. Measles cases in the United States are reported by healthcare providers and clinical
167 laboratories to their corresponding health department. Each case is investigated by local and state health
168 departments classified according to standard case definitions[46], and linked into clusters
169 epidemiologically (e.g., by establishing a direct contact or a shared location between cases, or when
170 cases are part of a specific community where an outbreak is occurring). Cases are considered
171 internationally imported if at least part of the exposure period (7–21 days before rash onset) occurred
172 outside the United States and rash occurred within 21 days of entry into the United States, with no
173 known exposure to measles in the United States during the exposure period.

174 Confirmed measles cases are routinely reported by state health departments to the CDC. 2,098 measles
175 cases were reported in the United States between January 2001 and December 2016. The number of
176 annual cases did not exceed 700 cases during this time period (Figure 3, Supplement Figure S1). The
177 importation status, 5-year age group, onset date, county, and state of residence were fully reported for

178 2,077 cases. The 21 cases with missing data were discarded. 25% of the cases were classified as
179 importations. 39% of the cases had their genotype reported. The dataset of 2,077 cases is referred to as
180 “reference dataset” in the results section, and was used to evaluate the performance of the inference
181 method.

182 Among cases with complete data, 737 independent clusters, containing 1 to 380 cases, were
183 reconstructed through contact tracing investigations. Not every identified case could be linked to an
184 importation, and some transmission clusters contained multiple imported cases (e.g. when related
185 individuals travel together to a foreign country and were infected there). Out of the 737 reference
186 clusters, 38 had several cases classified as importations, 256 had none identified.

187 *Model and parameters*

188 The distributions and priors used in the studies are listed in Table 1. As no studies quantifying the
189 probability of age-specific contacts have been carried out in the United States, we used the estimates
190 from the POLYMOD study in the UK[36]. The incubation period and the generation time of measles
191 were taken from previous studies [47–49]. We used the population centroid of each county to compute
192 the distance matrix[50]. We used a beta distribution as the prior of the conditional report ratio[8]. The
193 mean of the prior distribution was calculated using the number of clusters whose first case was not
194 classified as an imported case, meaning the investigations were not able to trace back to the first case
195 imported. As there was no prior information on the possible values of the spatial parameters, we used
196 uniform distributions as priors for a and b .

197 For pre-clustering of cases, we set the temporal threshold δ to 30 days, which is above the 97.5% upper
198 quantile of the generation time with a missing generation. We were interested in local transmission to
199 describe the impact of an imported case on a community. But we only had information on the county
200 of residency for each case. Counties are large geographical units: the average county land area is
201 2,911km² and the maximum values reach 50,000km². Therefore, we set the spatial threshold γ to 100km
202 to exclude long distance transmission, while still allowing for cross-county transmission.

203 Finally, we tested several relative and absolute importation thresholds λ . Absolute values were
204 calculated from a factor k , multiplied by the number of components in L_i , excluding the binary genetic
205 component. Tested values were $k = 0.05$ ($\lambda = -15$) and $k = 0.1$ ($\lambda = -11$). Connections were
206 considered unlikely if the log-likelihood was worse than λ . Relative values were quantiles of all
207 recorded log-likelihoods in the sampled trees (Table 1).

208 Using the contact tracing investigations, we considered three different initial distributions of the
209 importation status. In scenario 1, there was no inference of the importation status of cases, and the first
210 case of each epidemiological cluster was classified as importation (Ideal importation). In scenario 2:
211 there was no inference of the importation status of cases, and all cases identified as importation in the
212 contact tracing investigations were classified as importations (Epidemiological importation). Finally, in
213 Scenario 3, the importation status of cases was inferred, using different thresholds λ , and using no prior
214 information on the importation status of cases or the importation status from the contact tracing
215 investigations.

216 In order to compare the inferred and reference clusters, we calculated for each case i) the proportion of
217 the reference cluster correctly inferred (sensitivity) and ii) the proportion of the inferred cluster that was
218 part of the reference cluster (precision). These values were calculated at every iteration, and the median
219 values were used to evaluate the fit obtained with different values of λ . We also used the inferred cluster
220 size distribution to the reference data. The credibility intervals for each case are reported in the
221 Supplement (Supplement Figure S2).

222 **Results**

223 We clustered 2,077 measles cases reported in the United States between January 2001 and December
224 2016 using their onset date, age groups, location and genotype. Using the contact tracing investigations,
225 we considered three different initial importation status distribution: i) only the ancestors of each
226 epidemiological cluster (first case of each cluster) were importations (ideal importation), ii) all cases
227 classified as importation in the contact tracing investigations were importations (epidemiological
228 importation), iii) no prior information on importation status of cases. The importation status of the cases

229 was therefore not probabilistically inferred in scenario 1 and 2. The short preliminary run was 30,000
230 iterations and 70,000 iterations. For each run, the trace of the posterior distribution shows the
231 convergence of the algorithm (Supplement Figure S3).

232 In scenario 1, we did not infer the importation status of cases. The inferred cluster size distribution
233 matched the contact tracing investigations (Figure 4A); 98% of the reference singletons were also
234 isolated in the inferred cluster. For 94% (95% Credibility Interval: 91-98%) of cases, the inferred cluster
235 had a sensitivity and precision above 75%, meaning more than 75% of the cases in the inferred cluster
236 were in the reference cluster, and more than 75% of the cases in the reference cluster were in the inferred
237 cluster (Figure 4B). For 80% (78 – 93%) of cases, the inferred clusters were a perfect match with the
238 reference clusters. The cluster size distribution stratified by state was similar to the contact tracing
239 investigations (Supplement Figure S4). Therefore, when each ancestor was considered as an
240 importation, the inferred clusters were very close to the reference ones.

241 In scenario 2, we used the importation status distribution of cases reported in the contact tracing
242 investigations (539 importations). Pre-clustering highlighted 165 cases with no potential infector, which
243 were also classified as importations. We observed discrepancies between the inferred cluster size
244 distribution and the reference one: Among the 704 cases inferred as importation, 61 (9%) were not
245 importations in the reference cluster. Furthermore, 94 cases were the ancestor of a reference cluster and
246 were not classified as importations in the inferred clusters (13%). The overall cluster size distribution
247 matched the reference distribution, but 111 reference singletons were inferred as part of transmission
248 clusters (Figure 4A, Supplement Figure S5). Although the precision of the inferred cluster was above
249 75% for 93% (88-93%) of the cases, 31% (6-39%) had a sensitivity score below 0.5, meaning they were
250 classified with less than half of their reference clusters (Figure 4C). The discrepancies observed in this
251 scenario are due to inconsistencies between the importation status distribution and the clustering of
252 cases in the contact tracing investigations, as reference clusters that gathered several importations were
253 split into different inferred clusters in Scenario 2.

254 In scenario 3, the importation status of cases was inferred from a threshold λ . For each case i , if the log-
255 likelihood L_i was worse than λ , the connection between the case and its index was removed and the

256 case was considered imported. Firstly, using an absolute factor $k = 0.05$ ($\lambda = -15$), 586 (581-593)
257 cases were classified as importations, 361 (355-369) of them were singletons. These numbers are much
258 lower than the reference datasets that contains 737 clusters, and 539 singletons (Figure 5A, Supplement
259 Figure S6). We observed very few misclassifications of importation status and singletons (15 (10-22)
260 misclassified importations, 4 (0-14) misclassified singletons), and the cluster size distribution for
261 clusters including two cases and more was very similar to the reference one. The precision of the
262 reconstructed cluster was very high (above 75% for 88% (85-93%) of cases) (Figure 5B). Overall, the
263 algorithm was not able to accurately identify importations and singletons as the threshold was too low
264 to eliminate some unrealistic connections, but the inferred larger clusters matched their reference
265 counterparts.

266 We then observed the impact of increasing λ on the inferred cluster size distribution. Runs obtained
267 using an absolute threshold with $k = 0.10$ ($\lambda = -11$) and 95% relative threshold yielded very similar
268 results. The number of cases inferred as importations was higher than in previous runs, while all
269 remaining links showed good connection between cases. The number of importations was closer to the
270 reference dataset, and the number of singletons was greater than the reference. Nevertheless, the 11%
271 (10-12%) of the inferred importations was not classified as importation in the reference clusters.
272 Furthermore, the number of two-case chains was overestimated, and bigger clusters were likely to be
273 split because of the removal of weaker connections. Therefore, increasing λ did not improve the cluster
274 size distribution, as many importations in the reference clusters were not identified and the number of
275 mismatches increased (Supplement Figures S7).

276 Finally, we combined prior information and inference of importation status. Cases considered as
277 importations in the contact tracing investigations were set as importations, and we inferred the
278 importation status of the remaining cases. We used a low threshold, to remove the least likely
279 transmission links ($k = 0.05$). Including prior information led to some misclassification of importation
280 status due to the inconsistencies between the epidemiological importation status and the reference
281 clusters. As in scenario 2, some cases were classified with only part of their reference clusters because
282 clusters with several importations were split into different clusters. Indeed, the sensitivity score of 34%

283 (7-51%) of cases was below 0.5. Nevertheless, the cluster size distribution observed in the simulation
284 was the closest to the reference clusters. There were 725 (719-731) clusters, 89% of importations were
285 also ancestors of reference clusters and the number of singletons matched the reference clusters (Figure
286 5A-C). The inferred clusters of 88% (86-94%) of the cases had a precision score of 1, showing they
287 were clustered without any false positives. Despite discrepancies in several states (Massachusetts,
288 Ohio), the cluster size distribution stratified by state showed good agreement with the reference clusters
289 (Supplement Figures S8).

290 The conditional report ratio in the transmission chains ρ and the spatial parameters a and b was
291 estimated in each scenario. The parameter estimates did not depend on the prior importation status
292 distribution or the value of λ . ρ was consistently estimated above 90%, showing a low number of
293 missing generations between cases (Supplement Figure S9). This number is not representative of the
294 overall report ratio, which is usually much lower[51], and does not take into account missing
295 importations in singletons and chains. High values of ρ show that the reported cases can be connected
296 without missing generations.

297 There was little variation in the estimates of the spatial parameters between the different scenarios. The
298 population parameter a was estimated between 0.6 and 1 for every scenario, and the distance parameter
299 b was between 0.08 and 0.12. In every scenario, more than 80% of the inferred transmission were
300 between cases distant of less than 10km, and few long-distance transmissions were recorded (50-
301 100km), hence although most of the reconstructed connections were between cases from the same
302 county, the algorithm was able to identify clusters spreading over several counties or states (Supplement
303 Figure S10).

304 We highlighted the added value of including the spatial distance between cases in the likelihood by
305 comparing the cluster size distribution inferred by selecting certain components of L_i (Supplement
306 Figure S11). The credibility intervals were much wider when the distance between cases is not part of
307 the likelihood, and the number of chains containing 2 to 10 cases was over estimated. The important
308 impact of the spatial component of likelihood was also due to the widespread American territory, and
309 could be lower in a different setting.

310 We used the ratio of the number of importations over the number of subsequent cases per state to
311 evaluate the intensity of transmission in each state between 2001 and 2016 (Figure 6). The maps
312 obtained in the scenario 1 (ideal scenario) or in scenario 3 (estimation of importation, with
313 epidemiological importations and $k = 0.05$) were very similar. We only observed minor differences,
314 for example in South Dakota and in Massachusetts, where the ratios were higher in scenario 3. The
315 highest ratio (31.8 in scenario 1) was observed in Ohio, and is mostly due to a 383 case outbreak in
316 2014[32]. We observed major differences between the incidence map (Figure 3A) and the ratio per
317 state. Indeed, although 403 cases were reported in California (highest number in the US), importations
318 caused on average 1.32 subsequent cases in scenario 1 (1.60 in scenario 3), showing a high proportion
319 of reported cases were inferred as importations.

320 Similarly, we used the inferred transmission chain to compute the inferred reproduction number in each
321 state. According to the model, about 60% cases did not cause future transmission, and about 5% caused
322 more than 5 subsequent cases (Supplement Figure S12). These numbers were consistent in each run.
323 The geographical distribution of reproduction number was very similar to the importation - subsequent
324 cases ratio (Supplement Figure S13).

325 **Discussion**

326 We developed the R package *o2geosocial* to classify measles cases into transmission clusters and
327 estimate their importation status using routinely collected surveillance data (genotype, age, onset date
328 and location of the cases). As recently observed during the 2018-2019 measles outbreak in New York,
329 delays in childhood vaccination, local susceptibility, and increased contacts can lead to large outbreaks
330 following importations[52,53]. Therefore, we were interested in highlighting the effect of imported
331 cases on communities and we focused on short distance transmission to identify areas where they
332 repeatedly caused subsequent transmission chains. Although this is not predictive of future
333 transmission, it highlights communities with potential for large transmission clusters.

334 We compared the inferred transmission clusters to the contact tracing investigations of 2,077 confirmed
335 measles cases reported in the United States between 2001 and 2016. We were able to produce reliable

336 estimates of known transmission clusters using epidemiological features with only few
337 misclassifications. Estimating the importation status of cases without prior knowledge was challenging
338 and caused uncertainty on the results. We tested different threshold λ to eliminate unlikely
339 transmissions, and we were able to identify most of the imported cases. Nevertheless, if several cases
340 were imported in the same region at a similar time, we could not find all of them without discarding
341 valid transmission events, and increasing the number of false positives. When we used the importation
342 status as defined in the contact tracing investigations without probabilistic inference (scenario 1 and 2),
343 the reconstructed clusters were similar to the reference ones. Results were also conclusive when we
344 combined prior information and importation inference. The reconstruction of transmission greatly
345 depends on the epidemiological investigations to identify measles importations in a community.

346 We used the genotype to censor connections between cases when it was reported, as there can be only
347 one reported genotype per transmission cluster. Using a simulated dataset (*toy_outbreak_long* in
348 *o2geosocial*), we explored the impact of increasing the proportion of genotyped cases on clustering and
349 observed it could help identify the number of concurrent transmission trees when multiple genotypes
350 are co-circulating. Moreover, we introduced a spatial component to the likelihood of connection
351 between cases using an exponential gravity model. Previous studies showed this model was able to
352 capture short distance dynamics better than other gravity models, and was easy to parametrise.
353 Introducing the spatial component greatly improved the precision and the sensitivity of the
354 reconstructed clusters (Supplement Figure S11), and the parameter estimates were robust in the different
355 scenarios.

356 The final results on the clustering of the 2,077 cases using *o2geosocial* were obtained in 7 hours for
357 each run of 100,000 iterations on a standard desktop computer (Intel Core i7, 3.20 GHz 6 cores), which
358 is much faster than previous implementations of *outbreaker* and *outbreaker2*. With the addition of the
359 pre-clustering step, whereby we reduced the number of potential infectors for each case, the algorithm
360 ran faster. For smaller chains (50,000 iterations), 4 hours were needed to estimate the importation status
361 and cluster the cases. The code for the package and the analysis developed in this project is shared on
362 Github (<https://github.com/alxsrobert/o2geosocial> and [alxsrobert/datapaperMO](https://github.com/alxsrobert/datapaperMO)), with an illustrative

363 toy dataset, and can be used to analyse recent outbreaks where contact-tracing investigations were not
364 carried out.

365 Although the results obtained are promising, it should be noted that the dynamics of measles
366 transmission in the United States are likely to be very specific to this location. Indeed, there were less
367 than 700 annual cases between 2001 and 2016. These cases were scattered across a large area, which
368 made the pre-clustering of cases very efficient as we focused on short-distance transmission. In smaller
369 or more endemic settings, the number of potential infectors per cases after the pre-clustering step might
370 be higher, which would increase the running time.

371 Furthermore, as the location of each case was deduced from the population centroid of counties, we
372 assumed that the distance between cases from the same county was effectively zero. American counties
373 are large and widespread geographical units that can include more than 1 million individuals. For future
374 use of *o2geosocial*, more accurate information on the location of cases could improve cluster inference
375 by identifying multiple importations in a given county. Because cases are reported by the state of
376 residency, we had to ignore that cases may have been out of the reported county or state during their
377 incubation and infectious period, which has been seen during some outbreaks, such as the 2015 “Disney
378 outbreak” in California[54].

379 We did not include prior information on the local susceptibility of the different areas affected in
380 *o2geosocial*, and these could be estimated using historical values of local coverage. However, protocols
381 to estimate local vaccination coverage can differ in time and space and be difficult to compare, or
382 unavailable at the local level. Furthermore, these estimates are cross-sectional in nature, and might not
383 take into account catch-up vaccination campaigns, or immunity induced by previous outbreaks. Local
384 seroprevalence surveys could identify pockets of susceptibles, but they have not been carried out on a
385 subnational scale in most countries[55].

386 There has been no national quantitative analysis of age-specific contact patterns carried out in the United
387 States, so we relied on a contact matrix between age-groups available for Great Britain from the
388 POLYMOD study[36]. Nevertheless, little variation in the contact rates between age groups has been

389 observed between European countries, and a previous projection of the social contact matrix in the
390 United States yielded similar results[56]. POLYMOD data was probably the most reliable source of
391 information we could use to deduce an estimate of the contact matrix in the United States.

392 **Conclusions**

393 Heterogeneity in immunity can cause large outbreaks in countries with high national vaccine coverage,
394 and identifying potential foyers of transmission in post-elimination settings is key for outbreak
395 prevention and control. We have presented a method for estimating the cluster size distribution of past
396 measles outbreaks from routinely collected surveillance data. We found that adding prior knowledge
397 on the importation status of cases improved the inference of the transmission clusters. Although the
398 method was able to identify a proportion of importations, epidemiological investigations on the history
399 of travel and exposure reduced uncertainty on the clustering of cases. We believe these investigations
400 are needed to produce reliable estimates of past transmission clusters. In lieu of the importation status,
401 if multiple genotypes are co-circulating, increasing the proportion of genotyped cases could help discard
402 potential connections and find imported cases. Even with limited information, this method was able to
403 infer probabilistic transmission clusters in a fast and efficient way.

404 **Acknowledgements**

405 We acknowledge Thibaut Jombart for technical support and feedback on the analysis plan.

406 **Funding**

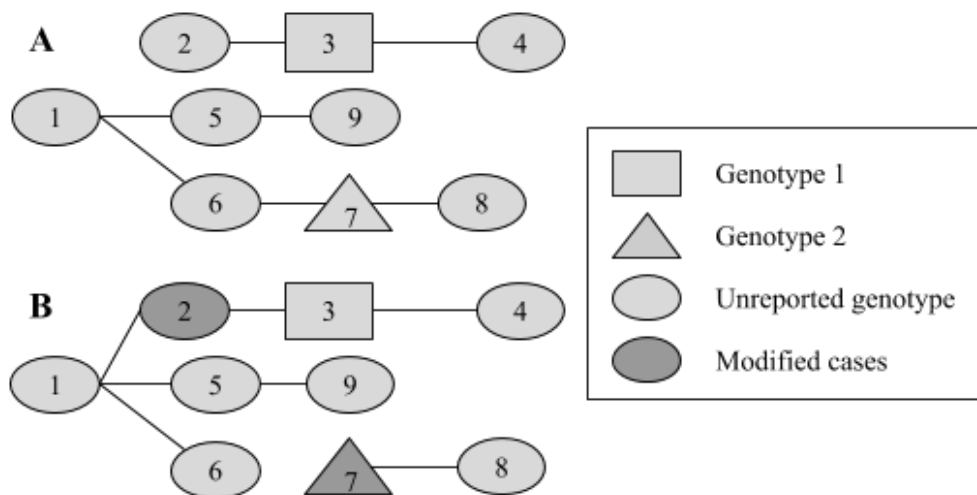
407 AR was supported by the Medical Research Council (MR/N013638/1). SF was supported by a
408 Wellcome Trust Senior Research Fellowship in Basic Biomedical Science (210758/Z/18/Z). AJK was
409 supported by a Sir Henry Dale Fellowship jointly funded by the Wellcome Trust and the Royal
410 Society (206250/Z/17/Z).

411 **Disclaimer**

412 The findings and conclusions in this report are those of the authors and do not necessarily represent
413 the official position of the Centers for Disease Control and Prevention, US Department of Health and
414 Human Services.

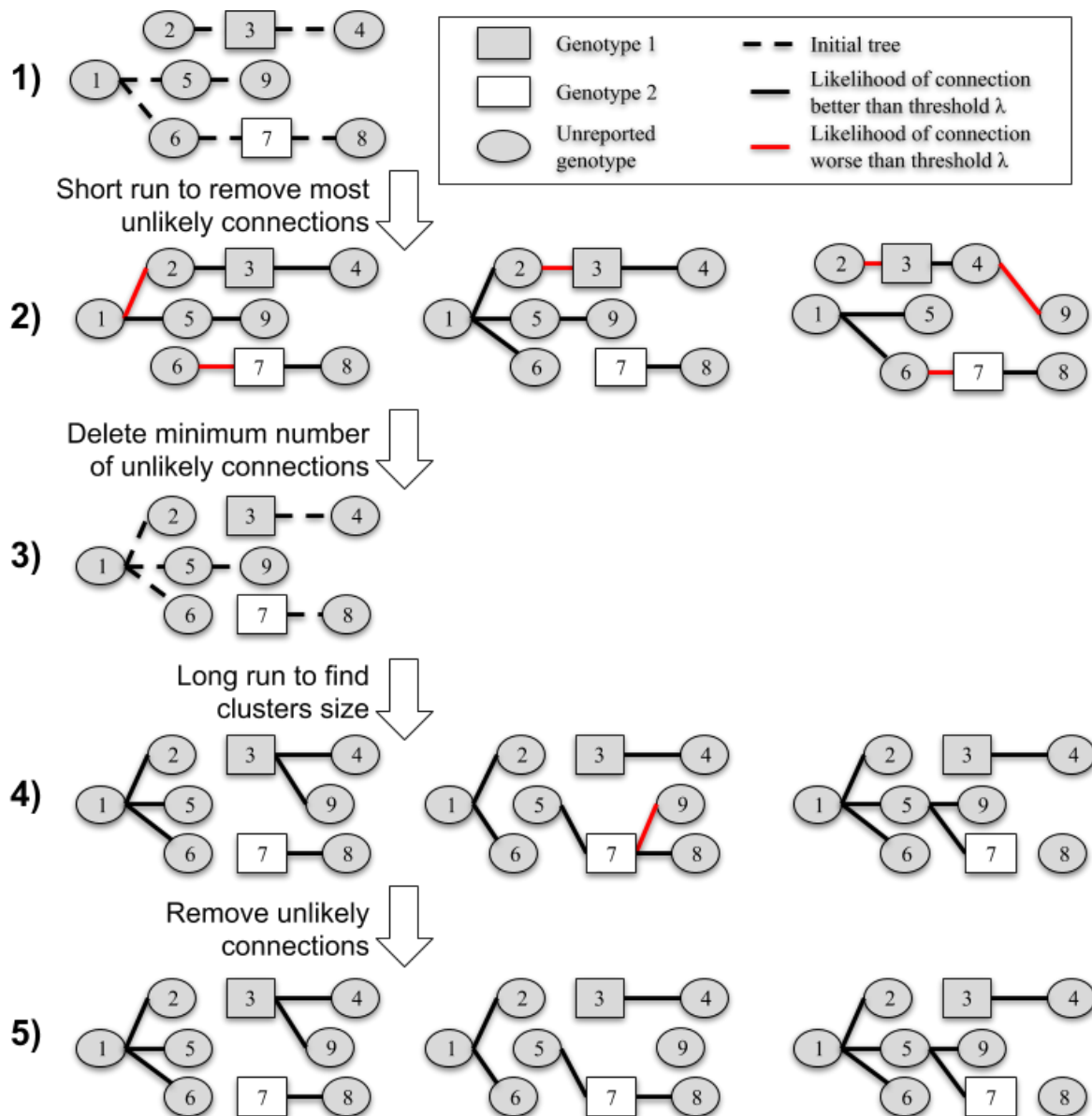
415

416 Figure legends



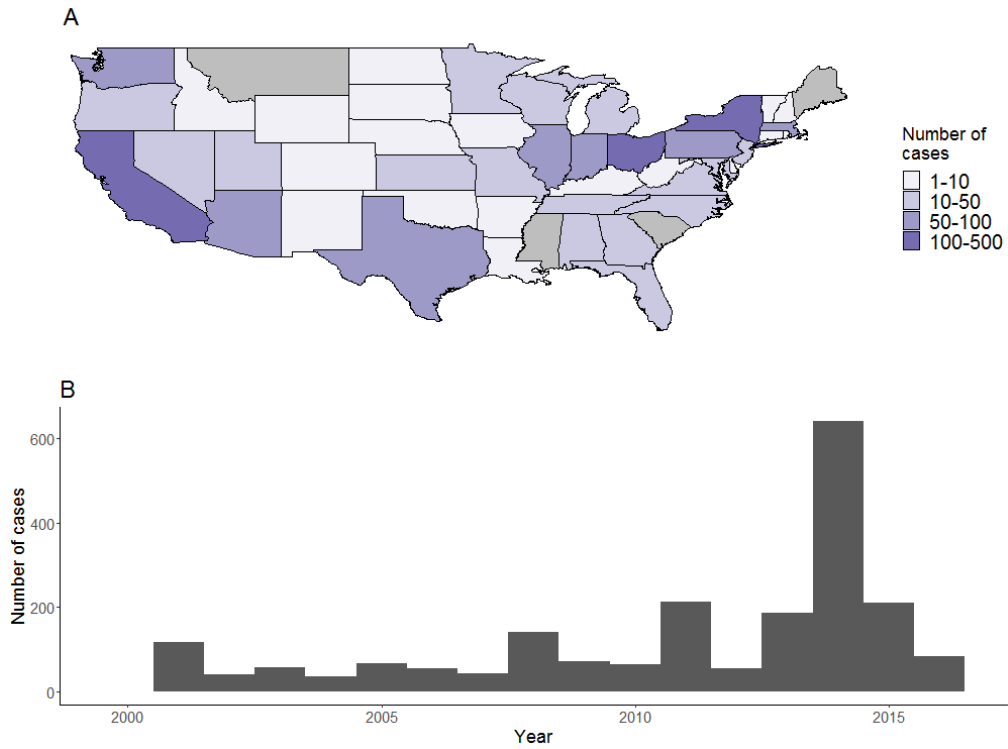
417

418 Figure 1: Example of the change of ancestors. Panel A represent the initial tree, B is the new tree proposed after
419 the movement. Initially, there are two ancestors (cases 1 and 2) in a group of 9 cases. 3 and 7 have different
420 genotypes and cannot be part of the same tree, the genotypes of the other cases are not reported. The date of
421 infection is in increasing order (1 is the first case, 9 is the last). Therefore, 1 is the only potential infector for 2.
422 One new ancestor was randomly drawn to conserve the number of trees. In this example, 7 is the new ancestor
423 (6 was the only other possibility). The ratio of the posterior densities of A and B were then used to determine
424 whether to accept or reject the proposal, according to the Metropolis-Hastings algorithm. This movement ensures
425 good mixing of the potential ancestors of the transmission clusters.



426

427 Figure 2: Estimating importation status and cluster size distributions in two MCMC runs. Step 1: Initial tree obtained
 428 after pre-clustering, with the minimum number of importations (here 2, as there are two reported genotypes). Step
 429 2: Samples from the first short run, with red line showing connection worse than the arbitrary threshold λ . Step 3:
 430 Initial tree for the final run, with 1 more importation than in step 1, which corresponds to the minimum number of
 431 unlikely transmissions at step 2. Step 4: Samples from the long run. Step 5: Final trees used to compute cluster
 432 size distribution and importation status of each case. Case 7 is an importation in one third of the final samples,
 433 whereas case 3 is an importation in all of them.

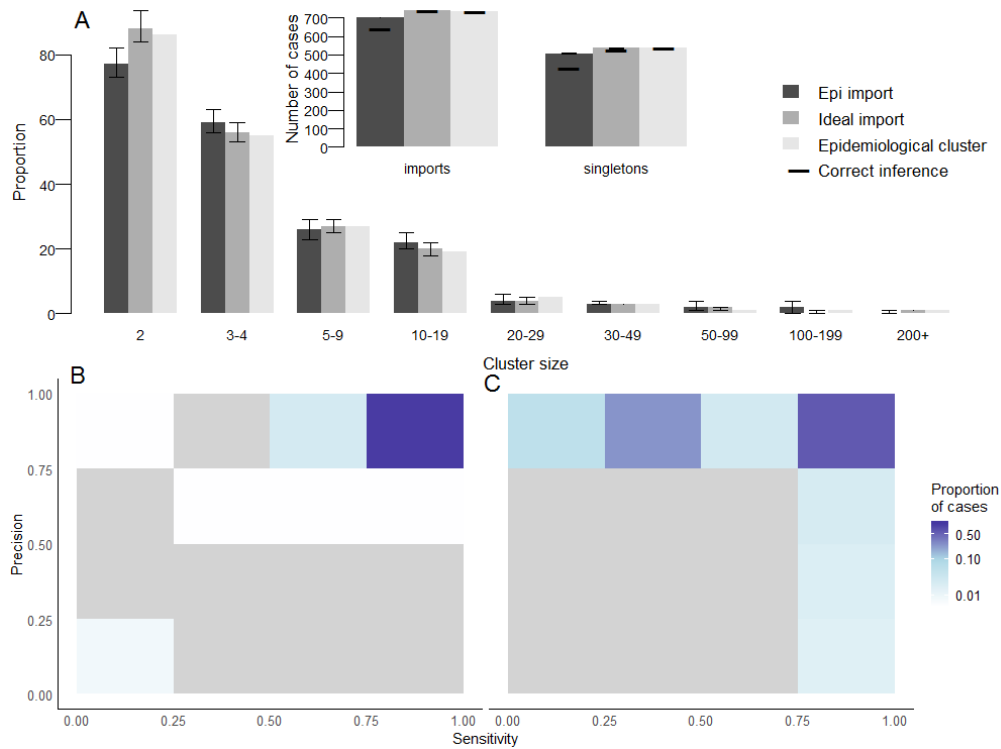


434

435 Figure 3: Panel A: number of cases per state and Panel B: Annual number of cases reported in the United States
436 between 2001 and 2016. Alaska and Hawaii are not shown on Panel A.

437

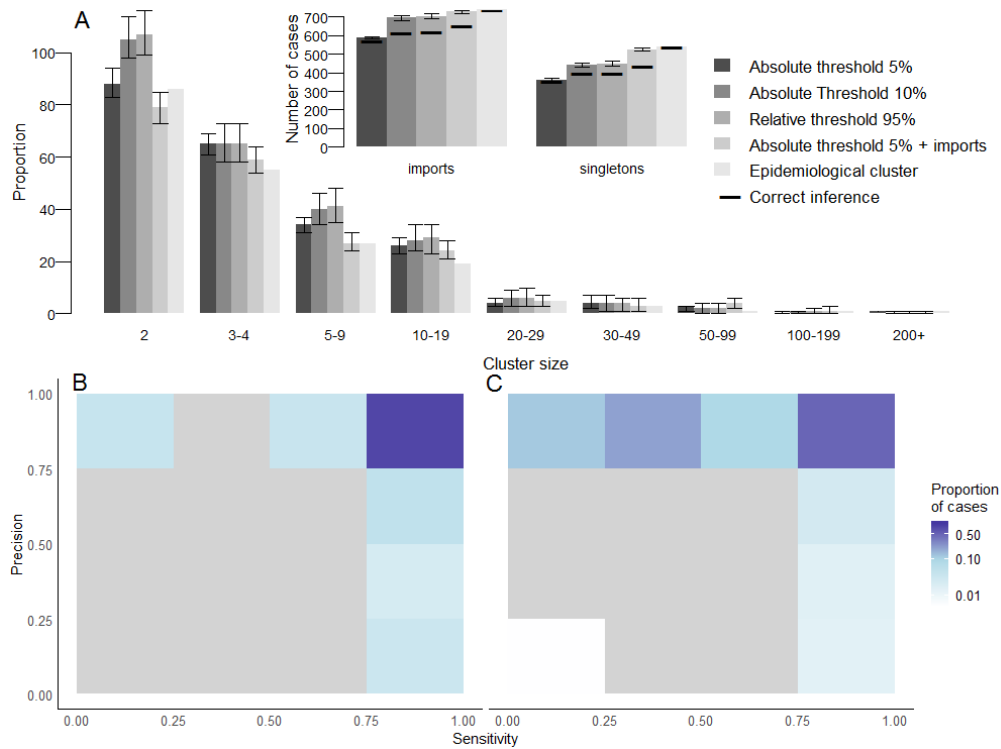
438



439

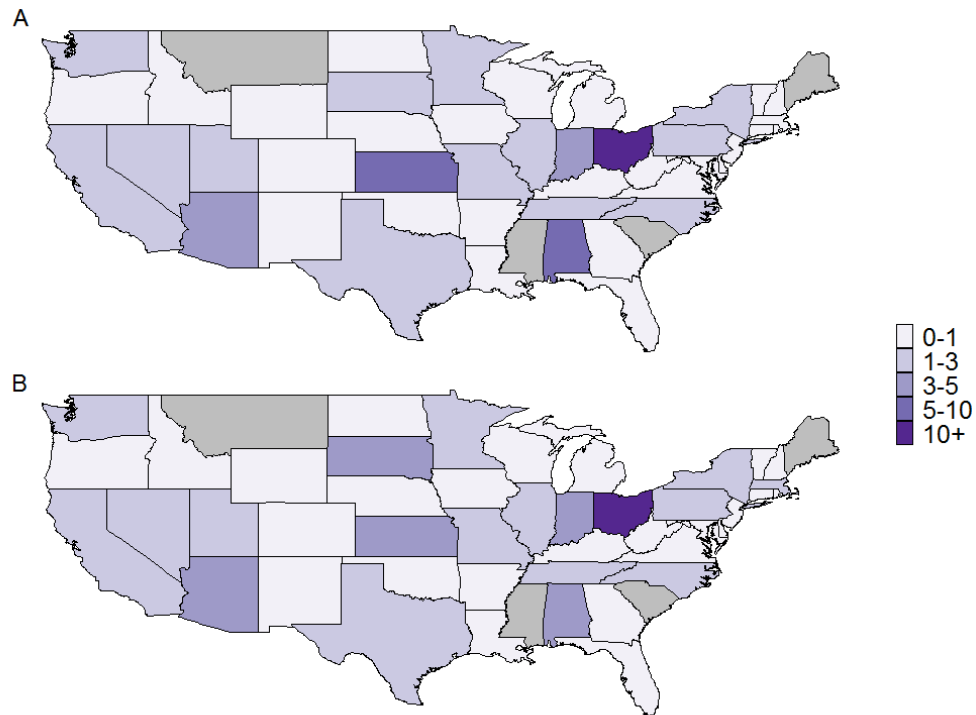
440 Figure 4: Description of transmission clusters inferred using prior knowledge on importation status of cases. Panel
 441 A: Cluster size distribution for the scenario 1 and 2 (grey and dark grey), compared to the reference clusters
 442 (lightgrey). Arrows represent the 95% credibility intervals of each estimate. Only clusters containing at least 2 cases
 443 are represented. Insert: Number of importations and number of isolated cases (singletons) in scenario 1 and 2, and
 444 in the reference clusters. For each scenario, the horizontal dark line represents the number of importations that are
 445 also importations in the reference clusters, same for singletons. Panel B: Heatmap representing the precision and
 446 sensitivity of the clusters for each case in scenario 1, cases are classified in a category depending on the proportion
 447 of their reference cluster that were inferred in the same cluster (x-axis) and the proportion of mismatches in the
 448 inferred cluster. Panel C: Same for scenario 2.

449



450

451 Figure 5: Description of transmission clusters generated with inferred importation status of cases. Panel A: Cluster
 452 size distribution for different value of threshold in the scenario 3 (sorted by shades of grey), compared to the
 453 reference clusters (lightgrey). Arrows represent the 95% credibility intervals of each estimate. Only clusters
 454 containing at least 2 cases are represented. Insert: Number of importations and number of isolated cases
 455 (singletons). For each scenario, the horizontal dark line represents the number of importation that are also
 456 importations in the reference clusters, same for singletons. Panel B: Heatmap representing the precision and
 457 sensitivity of the clusters for each case in scenario 3, with a 5% relative threshold, cases are classified in a category
 458 depending on the proportion of their reference cluster that were inferred in the same cluster. Panel C: Same when
 459 importation status is taken from the contact tracing investigations and inferred using a 5% relative threshold.



460

461 Figure 6: Ratio of the number of importations over the number of subsequent cases in each state in A/ Scenario 1
 462 (Ideal importations) and B/ Scenario 3 with epidemiological importations and $k = 0.05$. Grey states represent states
 463 that did not report any case.

464 Tables

465

Table 1: Values of parameters used to cluster cases declared in the United States

Parameter	Symbol	Distribution
Incubation period	$f(t)$	Gamma, mean = 11.5, sd = 2.24
Generation time	$w(t)$	Normal, Mean = 11.7, sd = 2.0
Conditional report ratio	ρ	Prior: Beta distribution, Mean = 0.65, sd = 0.15
Spatial parameter 1	a	Prior: Uniform distribution
Spatial parameter 2	b	Prior: Uniform distribution
Spatial pre clustering	γ	Fixed: 100 km
Temporal pre clustering	δ	Fixed: 30 days
Importation threshold	λ	Absolute: <ul style="list-style-type: none"> • $5 * \log 0.05$ • $5 * \log 0.1$ Relative: <ul style="list-style-type: none"> • 5%

466

467

468 Reference

- 469 [1] Ferguson NM, Donnelly CA, Anderson RM. Transmission intensity and impact of control
470 policies on the foot and mouth epidemic in Great Britain. *Nature* 2001.
471 <https://doi.org/10.1038/35097116>.
- 472 [2] Wallinga J, Teunis P. Different Epidemic Curves for Severe Acute Respiratory Syndrome
473 Reveal. *Am J Epidemiol* 2004;160:509–16.
- 474 [3] Lloyd-Smith JO, Schreiber SJ, Kopp PE, Getz WM. Superspreading and the effect of individual
475 variation on disease emergence. *Nature* 2005;438:355–9. <https://doi.org/10.1038/nature04153>.
- 476 [4] Faye O, Boëlle P-Y, Heleze E, Faye O, Loucoubar C, Magassouba N, et al. Chains of
477 transmission and control of Ebola virus disease in Conakry, Guinea, in 2014: an observational
478 study. *Lancet Infect Dis* 2015;15:320–6. [https://doi.org/10.1016/S1473-3099\(14\)71075-8](https://doi.org/10.1016/S1473-3099(14)71075-8).
- 479 [5] Ypma RJF, van Ballegooijen WM, Wallinga J. Relating phylogenetic trees to transmission trees
480 of infectious disease outbreaks. *Genetics* 2013;195:1055–62.
481 <https://doi.org/10.1534/genetics.113.154856>.
- 482 [6] Wallinga J, Lipsitch M. How generation intervals shape the relationship between growth rates
483 and reproductive numbers. *Proc R Soc B Biol Sci* 2007;274:599–604.
484 <https://doi.org/10.1098/rspb.2006.3754>.
- 485 [7] Cauchemez S, Ferguson NM. Methods to infer transmission risk factors in complex outbreak
486 data. *J R Soc Interface* 2012;9:456–69. <https://doi.org/10.1098/rsif.2011.0379>.
- 487 [8] Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N. Bayesian Reconstruction of
488 Disease Outbreaks by Combining Epidemiologic and Genomic Data. *PLoS Comput Biol*
489 2014;10. <https://doi.org/10.1371/journal.pcbi.1003457>.
- 490 [9] Campbell F, Cori A, Ferguson N, Jombart T. Bayesian inference of transmission chains using
491 timing of symptoms, pathogen genomes and contact data. *PLoS Comput Biol* 2019.

- 492 <https://doi.org/10.1371/journal.pcbi.1006930>.
- 493 [10] Haydon DT, Chase-Topping M, Shaw DJ, Matthews L, Friar JK, Wilesmith J, et al. The
494 construction and analysis of epidemic trees with reference to the 2001 UK foot-and-mouth
495 outbreak. *Proc R Soc B Biol Sci* 2003. <https://doi.org/10.1098/rspb.2002.2191>.
- 496 [11] Cauchemez S, Boëlle PY, Donnelly CA, Ferguson NM, Thomas G, Leung GM, et al. Real-time
497 estimates in early detection of SARS. *Emerg Infect Dis* 2006.
- 498 [12] Heijne JCM, Rondy M, Verhoef L, Wallinga J, Kretzschmar M, Low N, et al. Quantifying
499 transmission of norovirus during an outbreak. *Epidemiology* 2012.
500 <https://doi.org/10.1097/EDE.0b013e3182456ee6>.
- 501 [13] Kendall M, Ayabina D, Colijn C. Estimating transmission from genetic and epidemiological
502 data: a metric to compare transmission trees 2016:1–22. <https://doi.org/10.1214/17-STS637>.
- 503 [14] Worby CJ, O’Neill PD, Kypraios T, Robotham J V., De Angelis D, Cartwright EJP, et al.
504 Reconstructing transmission trees for communicable diseases using densely sampled genetic
505 data. *Ann Appl Stat* 2016. <https://doi.org/10.1214/15-AOAS898>.
- 506 [15] Lau MSY, Marion G, Streftaris G, Gibson G. A Systematic Bayesian Integration of
507 Epidemiological and Genetic Data. *PLoS Comput Biol* 2015.
508 <https://doi.org/10.1371/journal.pcbi.1004633>.
- 509 [16] Spada E, Saggiocca L, Sourdis J, Garbuglia AR, Poggi V, De Fusco C, et al. Use of the minimum
510 spanning tree model for molecular epidemiological investigation of a nosocomial outbreak of
511 hepatitis C virus infection. *J Clin Microbiol* 2004. <https://doi.org/10.1128/JCM.42.9.4230-4236.2004>.
- 513 [17] Mollentze N, Nel LH, Townsend S, le Roux K, Hampson K, Haydon DT, et al. A bayesian
514 approach for inferring the dynamics of partially observed endemic infectious diseases from
515 space-time-genetic data. *Proc R Soc B Biol Sci* 2014. <https://doi.org/10.1098/rspb.2013.3251>.
- 516 [18] Gire SK, Goba A, Andersen KG, Sealfon RSG, Park DJ, Kanneh L, et al. Genomic surveillance

- 517 elucidates Ebola virus origin and transmission during the 2014 outbreak. *Science* (80-)
518 2014;345:1369--1372. <https://doi.org/10.1126/science.1259657>.
- 519 [19] Carroll MW, Matthews DA, Hiscox JA, Elmore MJ, Pollakis G, Rambaut A, et al. Temporal
520 and spatial analysis of the 2014-2015 Ebola virus outbreak in West Africa. *Nature* 2015;524:97–
521 101. <https://doi.org/10.1038/nature14594>.
- 522 [20] Ruan YJ, Wei CL, Ee LA, Vega VB, Thoreau H, Yun STS, et al. Comparative full-length
523 genome sequence analysis of 14 SARS coronavirus isolates and common mutations associated
524 with putative origins of infection. *Lancet* 2003;361:1779–85. [https://doi.org/10.1016/S0140-](https://doi.org/10.1016/S0140-6736(03)13414-9)
525 [6736\(03\)13414-9](https://doi.org/10.1016/S0140-6736(03)13414-9).
- 526 [21] Pybus OG, Rambaut A. Evolutionary analysis of the dynamics of viral infectious disease. *Nat*
527 *Rev Genet* 2009;10:540–50. <https://doi.org/10.1038/nrg2583>.
- 528 [22] Grenfell BT, Pybus OG, Gog JR, Wood JLN, Daly JM, Mumford JA, et al. Unifying the
529 Epidemiological and Evolutionary Dynamics of Pathogens. *Science* (80-) 2004;303.
- 530 [23] Campbell F, Strang C, Ferguson N, Cori A, Jombart T. When are pathogen genome sequences
531 informative of transmission events? *PLoS Pathog* 2018.
532 <https://doi.org/10.1371/journal.ppat.1006885>.
- 533 [24] Rota PA, Brown K, Mankertz A, Santibanez S, Shulga S, Muller CP, et al. Global distribution
534 of measles genotypes and measles molecular epidemiology. *J Infect Dis* 2011;204.
535 <https://doi.org/10.1093/infdis/jir118>.
- 536 [25] Hiebert J, Severini A. Measles molecular epidemiology : What does it tell us and why is it
537 important? *Canada Commun Dis Rep CCDR* 2014;40.
- 538 [26] Brown KE, Rota PA, Goodson JL, Williams D, Abernathy E, Takeda M, et al. Genetic
539 characterization of measles and rubella viruses detected through global measles and rubella
540 elimination surveillance, 2016-2018. *Morb Mortal Wkly Rep* 2019;68:587–91.
541 <https://doi.org/10.15585/mmwr.mm6826a3>.

- 542 [27] Gardy JL, Naus M, Amlani A, Chung W, Kim H, Tan M, et al. Whole-genome sequencing of
543 measles virus genotypes H1 and D8 during outbreaks of infection following the 2010 Olympic
544 Winter Games reveals viral transmission routes. *J Infect Dis* 2015;212:1574–8.
545 <https://doi.org/10.1093/infdis/jiv271>.
- 546 [28] Penedos AR, Myers R, Hadeef B, Aladin F, Brown KE. Assessment of the Utility of Whole
547 Genome Sequencing of Measles Virus in the Characterisation of Outbreaks 2015:1–16.
548 <https://doi.org/10.1371/journal.pone.0143081>.
- 549 [29] World Health Organisation. Measles virus nomenclature Update: 2012. *Wkly Epidemiol Rec*
550 2012;87:73–80. <https://doi.org/10.1016/j.actatropica.2012.04.013>.
- 551 [30] Hagemann C, Streng A, Kraemer A, Liese JG. Heterogeneity in coverage for measles and
552 varicella vaccination in toddlers - Analysis of factors influencing parental acceptance. *BMC*
553 *Public Health* 2017;17. <https://doi.org/10.1186/s12889-017-4725-6>.
- 554 [31] Glasser JW, Feng Z, Omer SB, Smith PJ, Rodewald LE. The effect of heterogeneity in uptake
555 of the measles, mumps, and rubella vaccine on the potential for outbreaks of measles: A
556 modelling study. *Lancet Infect Dis* 2016;16:599–605. [https://doi.org/10.1016/S1473-](https://doi.org/10.1016/S1473-3099(16)00004-9)
557 [3099\(16\)00004-9](https://doi.org/10.1016/S1473-3099(16)00004-9).
- 558 [32] Gastañaduy PA, Budd J, Fisher N, Redd SB, Fletcher J, Miller J, et al. A Measles Outbreak in
559 an Underimmunized Amish Community in Ohio. *N Engl J Med* 2016;375:1343–54.
560 <https://doi.org/10.1056/NEJMoa1602295>.
- 561 [33] Woudenberg T, Van Binnendijk RS, Sanders EAM, Wallinga J, De Melker HE, Ruijs WLM, et
562 al. Large measles epidemic in the Netherlands, May 2013 to March 2014: Changing
563 epidemiology. *Eurosurveillance* 2017;22:1–9. [https://doi.org/10.2807/1560-](https://doi.org/10.2807/1560-7917.ES.2017.22.3.30443)
564 [7917.ES.2017.22.3.30443](https://doi.org/10.2807/1560-7917.ES.2017.22.3.30443).
- 565 [34] Keenan A, Ghebrehewet S, Vivancos R, Seddon D, MacPherson P, Hungerford D. Measles
566 outbreaks in the UK, is it when and where, rather than if? A database cohort study of childhood

- 567 population susceptibility in Liverpool, UK. *BMJ Open* 2017;7.
568 <https://doi.org/10.1136/bmjopen-2016-014106>.
- 569 [35] Kucharski AJ, Edmunds WJ. Characterizing the Transmission Potential of Zoonotic Infections
570 from Minor Outbreaks. *PLoS Comput Biol* 2015;11:1–17.
571 <https://doi.org/10.1371/journal.pcbi.1004154>.
- 572 [36] Mossong J, Hens N, Jit M, Beutels P, Auranen K, Mikolajczyk R, et al. Social contacts and
573 mixing patterns relevant to the spread of infectious diseases. *PLoS Med* 2008;5:0381–91.
574 <https://doi.org/10.1371/journal.pmed.0050074>.
- 575 [37] Blumberg S, Lloyd-Smith JO. Inference of R0 and Transmission Heterogeneity from the Size
576 Distribution of Stuttering Chains. *PLoS Comput Biol* 2013;9:1–17.
577 <https://doi.org/10.1371/journal.pcbi.1002993>.
- 578 [38] Blumberg S, Enanoria WTA, Lloyd-Smith JO, Lietman TM, Porco TC. Identifying
579 postelimination trends for the introduction and transmissibility of measles in the United States.
580 *Am J Epidemiol* 2014;179:1375–82. <https://doi.org/10.1093/aje/kwu068>.
- 581 [39] Campbell F, Didelot X, Fitzjohn R, Ferguson N, Cori A, Jombart T. outbreaker2: A modular
582 platform for outbreak reconstruction. *BMC Bioinformatics* 2018;19.
583 <https://doi.org/10.1186/s12859-018-2330-z>.
- 584 [40] Lenormand M, Bassolas A, Ramasco JJ. Systematic comparison of trip distribution laws and
585 models. *J Transp Geogr* 2016;51:158–69. <https://doi.org/10.1016/j.jtrangeo.2015.12.008>.
- 586 [41] Zipf GK. The P 1 P 2/D hypothesis: On the intercity movement of persons. *Am Sociol Rev*
587 1946;11:677–86. <https://doi.org/10.2307/2657358>.
- 588 [42] Barthélemy M. Spatial networks. *Phys Rep* 2011;499:1–79.
589 <https://doi.org/10.1016/j.physrep.2010.11.002>.
- 590 [43] Xia Y, Bjørnstad ON, Grenfell BT. Measles Metapopulation Dynamics: A Gravity Model for
591 Epidemiological Coupling and Dynamics. *Am Nat* 2004;164:267–81.

- 592 <https://doi.org/10.1086/422341>.
- 593 [44] Lenormand M, Huet S, Gargiulo F, Deffuant G. A Universal Model of Commuting Networks.
594 PLoS One 2012;7. <https://doi.org/10.1371/journal.pone.0045985>.
- 595 [45] Andrieu C, De Freitas N, Doucet A, Jordan MI. An introduction to MCMC for machine learning.
596 Mach Learn 2003;50:5–43. <https://doi.org/10.1023/A:1020281327116>.
- 597 [46] Centers for Disease Control and Prevention (CDC). National Notifiable Disease Surveillance
598 System: measles/rubeola 2013. [https://www.cdc.gov/nndss/conditions/measles/case-](https://www.cdc.gov/nndss/conditions/measles/case-definition/2013/)
599 [definition/2013/](https://www.cdc.gov/nndss/conditions/measles/case-definition/2013/) (accessed October 23, 2019).
- 600 [47] Lessler J, Reich NG, Brookmeyer R, Perl TM, Nelson KE. Incubation periods of acute
601 respiratory viral infections: a systematic review 2015;9:291–300.
602 [https://doi.org/10.1016/S1473-3099\(09\)70069-6](https://doi.org/10.1016/S1473-3099(09)70069-6).Incubation.
- 603 [48] Klinkenberg D, Nishiura H. The correlation between infectivity and incubation period of
604 measles, estimated from households with two cases. J Theor Biol 2011;284:52–60.
605 <https://doi.org/10.1016/j.jtbi.2011.06.015>.
- 606 [49] Fine PEM. The Interval between Successive Cases of an Infectious Disease. Am J Epidemiol
607 2003;158:1039–47. <https://doi.org/10.1093/aje/kwg251>.
- 608 [50] US Census Bureau. Centers of Population for the 2010 Census 2010.
609 <https://www.census.gov/geographies/reference-files/2010/geo/2010-centers-population.html>
610 (accessed August 22, 2019).
- 611 [51] Woudenberg T, Woonink F, Kerkhof J, Cox K, Ruijs WLM. The tip of the iceberg :
612 incompleteness of measles reporting during a large outbreak in The Netherlands in 2013 – 2014.
613 Epidemiol Infect 2018;146:716–22. [https://doi.org/https://doi.org/10.1017/](https://doi.org/https://doi.org/10.1017/S0950268818002698)
614 [S0950268818002698](https://doi.org/https://doi.org/10.1017/S0950268818002698).
- 615 [52] Gastañaduy PA, Funk S, Paul P, Tatham L, Fisher N, Budd J, et al. Impact of public health
616 responses during ameasles outbreak in an amish community in Ohio: Modeling the dynamics of

- 617 transmission. *Am J Epidemiol* 2018. <https://doi.org/10.1093/aje/kwy082>.
- 618 [53] Patel M, Lee AD, Clemmons NS, Redd SB, Poser S, Blog D, et al. National Update on Measles
619 Cases and Outbreaks - United States, January 1-October 1, 2019. *MMWR Morb Mortal Wkly*
620 *Rep* 2019;68:893–6. <https://doi.org/10.15585/mmwr.mm6840e2>.
- 621 [54] Zipprich J, Winter K, Hacker J, Xia D, Watt J, Harriman K. Measles outbreak--California,
622 December 2014-February 2015. *vol.* 64. 2015.
623 <https://doi.org/10.1016/j.annemergmed.2015.04.002>.
- 624 [55] Durrheim D. Measles elimination, immunity, serosurveys, and other immunity gap diagnostic
625 tools. *J Infect Dis* 2018;218:341–3. <https://doi.org/10.1093/infdis/jiy138>.
- 626 [56] Prem K, Cook AR, Jit M. Projecting social contact matrices in 152 countries using contact
627 surveys and demographic data. *PLoS Comput Biol* 2017.
628 <https://doi.org/10.1371/journal.pcbi.1005697>.
- 629