

1 Incorporating epidemiological data into the genomic analysis of partially sampled infectious  
2 disease outbreaks

3 Jake Carson<sup>1,2</sup>, Matt Keeling<sup>1,2</sup>, Paolo Ribeca<sup>3</sup>, Xavier Didelot<sup>2,4,\*</sup>

4 <sup>1</sup> Mathematics Institute, University of Warwick, Coventry CV4 7AL, United Kingdom

5 <sup>2</sup> School of Life Sciences, University of Warwick, Coventry CV4 7AL, United Kingdom

6 <sup>3</sup> UK Health Security Agency, London NW9 5EQ, United Kingdom

7 <sup>4</sup> Department of Statistics, University of Warwick, Coventry CV4 7AL, United Kingdom

8 \* Corresponding author: [xavier.didelot@gmail.com](mailto:xavier.didelot@gmail.com)

9 Keywords: genomic epidemiology, transmission analysis, infectious disease outbreak,  
10 epidemiological data

## 11 ABSTRACT

12 Pathogen genomic data is increasingly being used to investigate transmission dynamics in  
13 infectious disease outbreaks. Combining genomic data with epidemiological data should  
14 substantially increase our understanding of outbreaks, but this is highly challenging when the  
15 outbreak under study is only partially sampled, so that both genomic and epidemiological data  
16 are missing for intermediate links in the transmission chains. Here we present a new dynamic  
17 programming algorithm to perform this task efficiently. We implement this methodology into  
18 the well-established TransPhylo framework to reconstruct partially sampled outbreaks using a  
19 combination of genomic and epidemiological data. We use simulated datasets to show that  
20 including epidemiological data can improve the accuracy of the inferred transmission links  
21 compared to inference based on genomic data only. This also allows us to estimate parameters  
22 specific to the epidemiological data (such as transmission rates between particular groups) which  
23 would otherwise not be possible. We then apply these methods to two real-world examples.  
24 Firstly, we use genomic data from an outbreak of tuberculosis in Argentina, for which data was  
25 also available on the HIV status of sampled individuals, in order to investigate the role of HIV co-  
26 infection in the spread of this tuberculosis outbreak. Second, we use genomic and geographical  
27 data from the 2003 epidemic of avian influenza H7N7 in the Netherlands to reconstruct its  
28 spatial epidemiology. In both cases we show that incorporating epidemiological data into the  
29 genomic analysis allows us to investigate the role of epidemiological properties in the spread of  
30 infectious diseases.

## 31 INTRODUCTION

32 Over the past decade there has been considerable research interest and methodological  
33 development in the analysis of pathogen genomic sequences to reconstruct the transmission  
34 events that occurred during an infectious disease outbreak (Jombart et al., 2014; Croucher and  
35 Didelot, 2015; Campbell et al., 2018; Duault et al., 2022). Additional epidemiological data about  
36 the infected hosts is often available, and it can be useful to integrate such data into the genomic  
37 analysis for two complementary reasons. Firstly, it should allow the transmission trees to be  
38 reconstructed more precisely when using genomic and epidemiological data compared to using  
39 genomic data alone. An example of this was provided by the reconstruction of a tuberculosis  
40 outbreak in British Columbia, in which the matrix of who-infected-whom probabilities contained  
41 less uncertainty when geographic data and measures of infectiousness (provided by smear and  
42 skin tests) were used as additional input (Didelot et al., 2014; Biek et al., 2015; Hatherell  
43 et al., 2016). Secondly, using data on epidemiological properties can enable inference on the  
44 correlation between transmission and these epidemiological properties. For example, individuals  
45 with high-risk behaviours contribute disproportionately to the spread of sexual diseases such  
46 as gonorrhoea (Chan et al., 2012; Fingerhuth et al., 2016; Whittles et al., 2019). Integrating  
47 behavioural data into a transmission analysis could help quantify this effect, which could be  
48 used in predictive models for example to inform the design of control measures such as targeted  
49 vaccination (Craig et al., 2015; Whittles et al., 2020, 2022).

50 Previous attempts have been made to integrate epidemiological and genomic data into outbreak  
51 reconstructions. The simplest case occurs if we assume that all cases of the outbreak have been  
52 sampled and are therefore present in the transmission tree. In this case the likelihood of the  
53 genomic data can simply be multiplied by the likelihood of the epidemiological data (Ypma et al.,  
54 2012; Morelli et al., 2012; Didelot et al., 2014; Hall et al., 2015). The epidemiological component  
55 of the likelihood is easy to compute as a product over all links in the transmission tree (Ypma  
56 et al., 2012). Each link represents an infection from a sampled infector to a sampled infectee,  
57 and as epidemiological data are available for both hosts, the contribution to the likelihood  
58 is analytically tractable. However, the vast majority of infectious disease outbreaks are only  
59 partially observed, with the proportion of missing cases being typically unknown as well (O'Neill  
60 and Roberts, 1999; Jewell et al., 2009; Chis Ster et al., 2009). These missing intermediates in  
61 the transmission trees represent a challenge for the integration of epidemiological data, since by  
62 definition there is no epidemiological data available on unknown putative hosts. This difficulty  
63 was noted, for example, when a spatial-genetic framework assuming complete sampling (Morelli  
64 et al., 2012) was extended to handle incomplete sampling, with two extreme scenarios proposed  
65 as bounds on the probability of spatial dispersion for unknown cases (Mollentze et al., 2014).

66 A naive approach to incorporate epidemiological data into the transmission analysis of a  
67 partially sampled outbreak is to consider all possible combinations for the epidemiological  
68 data of the unsampled cases along with their associated probabilities. For each combination  
69 the epidemiological component of the likelihood can be calculated as previously described in  
70 the case of a fully sampled outbreak (Ypma et al., 2012; Morelli et al., 2012; Didelot et al.,  
71 2014; Hall et al., 2015). The unconditioned likelihood is then be obtained as the average of  
72 these conditioned likelihoods, weighted according to their probabilities, using the law of total  
73 probability. However, the number of combinations scales exponentially with the number of  
74 unsampled cases in the transmission tree, and is only be computationally feasible for very small  
75 outbreaks. Another approach is to rely on data augmentation techniques within a Markov Chain

76 Monte-Carlo (MCMC) framework, in order to treat the epidemiological data of unsampled cases  
77 as additional parameters ([van Dyk and Meng, 2001](#); [O’Neill, 2002](#)). Again this may not scale  
78 well to larger outbreaks, especially since the number of unsampled cases is unknown, so that  
79 efficient reversible jump proposals are required to deal with the transdimensional parameter  
80 space ([Green, 1995](#); [Sisson, 2005](#)). Instead, we present a computationally efficient approach to  
81 calculate the epidemiological component of the likelihood. We show that this computation can  
82 be used to incorporate epidemiological data into the transmission analysis, boosting the accuracy  
83 of the analysis and generating the type of who-acquires-infection-from-whom matrices that are  
84 the cornerstone of predictive modelling. We illustrate our method on simulated datasets, before  
85 considering real-world examples of tuberculosis and H7N7 outbreaks.

## 86 NEW APPROACHES

87 We take as our starting point the TransPhylo methodology ([Didelot et al., 2014](#)), which  
88 represents the transmission tree by colouring the branches of an input dated phylogeny ([Rieux  
89 and Balloux, 2016](#)). The first version of TransPhylo considered only fully sampled outbreaks,  
90 so that it was possible to incorporate epidemiological data ([Didelot et al., 2014](#)). With the  
91 extension of TransPhylo to the more generally useful situation of a partially sampled outbreak,  
92 this possibility to integrate epidemiological data was lost ([Didelot et al., 2017, 2021](#)). More  
93 recently, TransPhylo was further extended to allow some hosts to be sampled more than once  
94 and to remove the assumption of complete transmission bottleneck ([Carson et al., 2024](#)), and  
95 this is the version that we use as our starting point for the incorporation of epidemiological  
96 data.

97 We extend the TransPhylo framework to incorporate known discrete epidemiological data on  
98 the sampled hosts, or a subset of them. Note that we use the term ‘deme’ to represent data that  
99 could be any discrete property of the hosts, for example geographical location in different towns  
100 or hospital wards, age or gender categories, classification based on behavioural data, infectious  
101 status from other infectious diseases, etc. We let  $S$  denote the number of demes (number  
102 of discrete epidemiological states). The transmission model within TransPhylo is a continuous  
103 time branching process ([Farrington et al., 2003](#)), in which each infected host generates a number  
104 of offspring  $k$  from an offspring distribution function  $\alpha(k)$ , and their infection times  $\tau$  relative  
105 to the infection time of the infector are sampled from a generation time distribution  $\gamma(\tau)$ .  
106 The mean of the offspring distribution is the basic reproduction number  $R$ . We extend this  
107 branching process so that a deme is sampled for each offspring conditional on the deme of the  
108 infector. Specifically, the probability that a newly infected host belongs to deme  $j$  given that  
109 their infector belongs to deme  $i$  is denoted  $P_{ij}$ . This matrix may take any form, as long as each  
110 of the rows sums up to one, and may include some parameters that we wish to infer jointly with  
111 the transmission tree.

112 We consider two complementary cases. In the first case, hosts in every deme have the same  
113 offspring distribution function and probability of being sampled. In this case the likelihood  
114 can be decomposed as the product of the transmission tree and of the epidemiological data,  
115 with the latter being calculated efficiently using a dynamic programming algorithm similar to  
116 the Felsenstein pruning algorithm ([Felsenstein, 1973, 1981](#)). In the second case, the offspring  
117 distribution function and probability of being sampled depend on the deme. This dependency  
118 may once again involve some parameters that we wish to estimate, for example different values

119  $R_1, \dots, R_S$  for the basic reproduction number within each of the demes. In this case the  
120 likelihood can no longer be decomposed as previously, but we show that it can still be calculated  
121 analytically using a more complicated dynamic programming algorithm.

## 122 RESULTS

### 123 Exemplary analysis of a simulated dataset where all demes have the same 124 offspring distribution and sampling probabilities

125 We simulate an outbreak with 250 observed infected hosts across five demes, with each observed  
126 host being sampled once. The observation cut-off time  $T$  is determined by the simulation in  
127 order to return the correct number of observed infected hosts. The generation time and primary  
128 observation time are both Gamma-distributed with shape and scale parameters equal to 2 and 1,  
129 respectively. For the transmission model, the offspring distribution follows a Negative Binomial  
130 distribution with  $r = 2$  and  $p = 0.5$ , so the basic reproduction number  $R = r = 2$ , and the  
131 sampling proportion is  $\pi = 0.8$ . The within-host pathogen population size is  $\kappa + \lambda\tau$  at time  $\tau$   
132 after infection, with  $\kappa = 0.1$  and  $\lambda = 0.2$ . The probability of an offspring having the same deme  
133 as their infector is  $\rho = 0.8$ , otherwise one of the other four demes is sampled uniformly. The  
134 resulting simulation contains 302 infected hosts (of which 250 are sampled). The transmission  
135 and phylogenetic trees are shown in Figure 1, which is coloured according to the demes of the  
136 hosts. Note that only the deme data for the observed hosts is used in the analysis.



Figure 1: Combined transmission and phylogenetic tree coloured by host deme used in the first simulation study. The tree contains 302 infected hosts, of which 250 are sampled. The red stars correspond to transmission events.

137 Our goals are to estimate the five parameters  $R$ ,  $\pi$ ,  $\rho$ ,  $\kappa$  and  $\lambda$ , and to correctly identify  
138 transmission links between sampled hosts. We performed four separate MCMC runs of 100,000  
139 iterations, which each took approximately 36 hours on a 3 GHz processor core. Mixing was  
140 relatively slow for the coalescent parameters  $\kappa$  and  $\lambda$ , with effective sample sizes of 300-600 in  
141 each run. This is in part due to these parameters being highly correlated to the transmission  
142 tree, which is updated separately within the MCMC algorithm, and in part due to having  
143 only one sample per host, leading to a wide posterior to explore for these parameters (Carson  
144 et al., 2024). The effective sample size was between 1300-1900 in each chain for  $\pi$ , 6000-  
145 7000 for  $R$ , and 6000-10000 for  $\rho$ . The multivariate Gelman-Rubin statistic comparing runs  
146 was 1.01 (Brooks and Gelman, 1998). The inferred means (95% credible intervals) for each  
147 parameter are  $R : 1.94 (1.68, 2.22)$ ,  $\pi : 0.78 (0.63, 0.93)$ ,  $\rho : 0.78 (0.73, 0.83)$ ,  $\kappa : 0.09 (0.00, 0.21)$ ,  
148  $\lambda : 0.20 (0.02, 0.41)$ . This shows that we are able to recover the simulated parameter values  
149 effectively, since the posterior means are close to the correct values and the credible intervals  
150 cover the correct values.

151 In order to evaluate our ability to reconstruct transmission links, we focus on transmissions  
152 between observed hosts. Out of the 250 observed hosts, 184 are infected by another observed  
153 host. The posterior probability estimates for the transmission links are summarised in Figure 2.  
154 If we define 0.5 as the posterior probability threshold for a transmission event being identified,  
155 we correctly identify 71 transmission links including the direction of transmission, giving a  
156 directional sensitivity of 39%. With only one observation per host it is common to identify a  
157 transmission link between two hosts, but be unsure of the direction of transmission (Didelot

158 et al., 2014, 2017; Carson et al., 2024). If we ignore the direction of transmission we identify  
 159 104 transmission links, giving a bidirectional sensitivity of 57%. We incorrectly establish 30  
 160 directional transmission links, and 37 bidirectional transmission links. However, as there are  
 161 62,250 possible host combinations, specificity is high ( $> 99.9\%$ ) in both cases. The resulting  
 162 precision is 70% when including the direction of transmission, and 74% when ignoring the  
 163 direction of transmission.

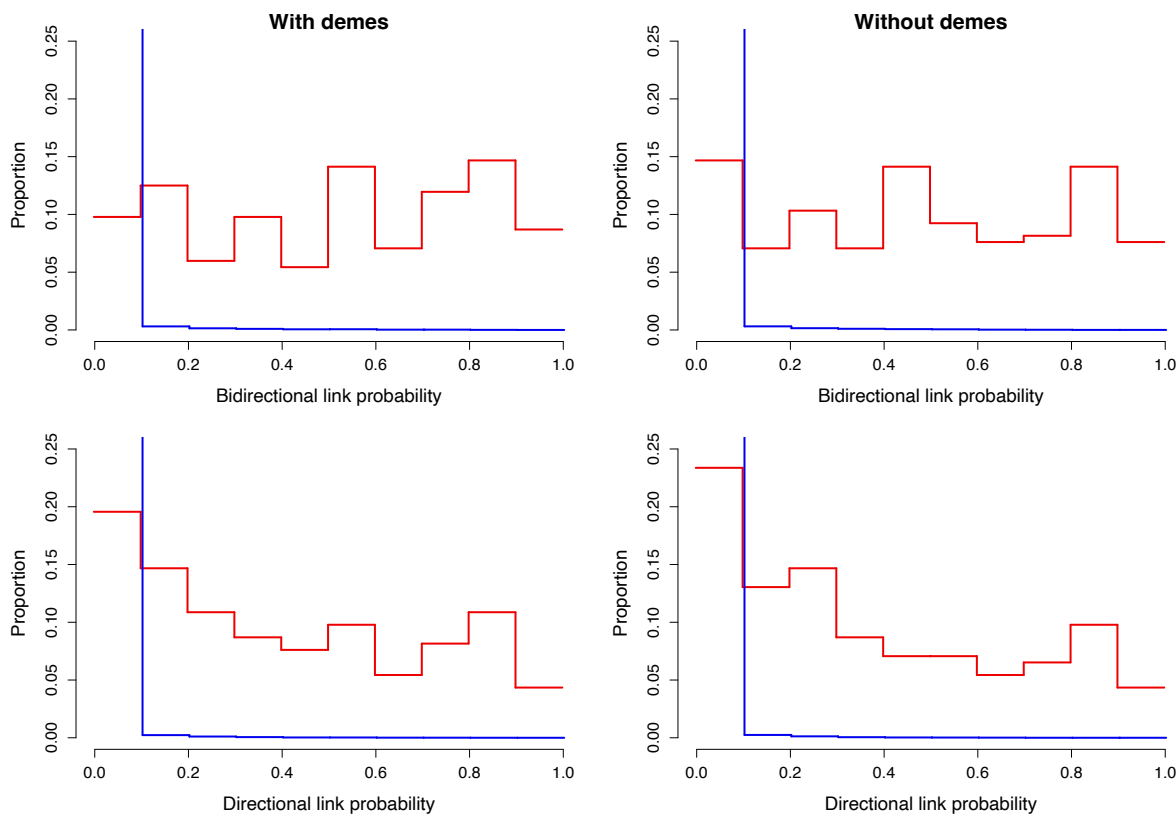


Figure 2: Posterior probability estimates for the transmission links in the first simulation study. The red lines correspond to correct transmission links, and the blue lines correspond to incorrect transmission links. In each case the posterior probability estimates are binned in 0.1 width intervals, and the vertical height indicates the proportion of transmission links contained within each bin. The left plots show estimates with deme information, and the right plots show estimates without deme information. The top plots show bidirectional transmission link estimates, and the bottom plots show directional transmission link estimates.

164 If inference is undertaken without the deme data we obtain similar parameter estimates for  $R$ ,  
 165  $\pi$ ,  $\kappa$ , and  $\lambda$ , but no longer obtain an estimate for  $\rho$  as we assume that all hosts belong to a single  
 166 deme. We correctly identify fewer transmission links, with 61 correct links being established if we  
 167 consider direction (directional sensitivity of 33%) and 86 if we do not (bidirectional sensitivity of  
 168 47%). However, we also obtain a slightly smaller number of false positives: 25 in the directional  
 169 case and 32 in the bidirectional case. The resulting precision changes very little, with a precision  
 170 of 71% when including the direction of transmission, and 73% otherwise.

171 **Exemplary analysis of a simulated dataset where demes have different**  
172 **offspring distributions and sampling probabilities**

173 We simulate a second outbreak with 250 observed hosts. The hosts belong to two demes, with  
174 each deme having its own  $R$ ,  $\pi$ , and  $\rho$  parameters. Specifically we set  $R_1 = 1.2$ ,  $R_2 = 2.2$ ,  
175  $\pi_1 = 0.4$ ,  $\pi_2 = 0.9$ ,  $\rho_1 = 0.9$ ,  $\rho_2 = 0.7$ , whilst maintaining  $\kappa = 0.1$  and  $\lambda = 0.2$ . Consequently,  
176 deme 1 has a low transmission rate and is poorly surveyed, whilst deme 2 has a high transmission  
177 rate and is well surveyed. Offspring are more likely than not to be in the same deme as the  
178 infecting host, but transmissions from deme 2 to deme 1 are more likely than transmissions  
179 from deme 1 to deme 2. The resulting simulation contains 325 hosts, and the transmission and  
180 phylogenetic trees are shown in Figure 3, which is coloured according to the demes.

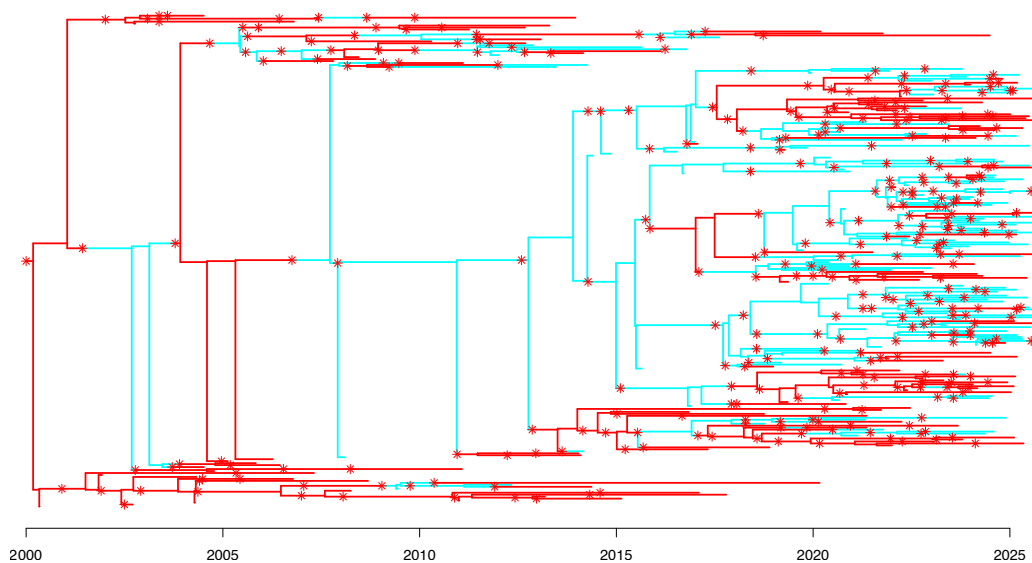


Figure 3: Combined transmission and phylogenetic tree coloured by host deme used in the second simulation study. The tree contains 325 infected hosts, of which 250 are sampled. The red stars correspond to transmission events.

181 We performed four separate MCMC runs of 100,000 iterations, which took approximately 46  
182 hours on a 3 GHz processor core. As with the previous simulation study, mixing was slowest  
183 for  $\kappa$  and  $\lambda$ , with effective sample sizes between 300-1100 for  $\kappa$  and 300-400 for  $\lambda$  in each  
184 chain. For the remaining parameters the effective sample sizes were 4000-5500 for  $R_1$ , 500-  
185 4200 for  $R_2$ , 2000-2300 for  $\pi_1$ , 1700-2300 for  $\pi_2$ , 900-3700 for  $\rho_1$ , and 1500-3300 for  $\rho_2$ . The  
186 multivariate Gelman-Rubin statistic comparing runs was 1.01 (Brooks and Gelman, 1998). The  
187 inferred means (and 95% credible intervals) for each parameter are  $R_1 : 1.15$  (0.92, 1.39),  
188  $R_2 : 2.21$  (1.78, 2.68),  $\pi_1 : 0.50$  (0.35, 0.68),  $\pi_2 : 0.88$  (0.68, 0.99),  $\rho_1 : 0.93$  (0.87, 0.96),  
189  $\rho_2 : 0.72$  (0.62, 0.81),  $\kappa : 0.10$  (0.00, 0.28),  $\lambda : 0.23$  (0.03, 0.48). Once again this shows that  
190 we are able to recover the simulated parameter values effectively since the inferred values are  
191 close to the correct values used in the simulation. Furthermore, since the credible intervals for  
192  $R_1$  and  $R_2$ , and for  $\pi_1$  and  $\pi_2$  do not overlap, we can deduce that we have correctly inferred



193 that the reproduction number and sampling probability are both higher in the second deme  
194 than in the first deme.

195 Looking once more at inferred transmission links, out of the 250 observed hosts 155 are infected  
196 by another observed host. The full transmission link probabilities are shown in Figure 4. Using  
197 the same posterior probability threshold as the first simulation study we correctly identify 65  
198 transmission links including the direction of transmission, giving a directional sensitivity of  
199 42%. If we ignore the direction of transmission then we identify 89 transmission links, giving a  
200 bidirectional sensitivity of 57%. We incorrectly establish 39 directional transmission links, and  
201 50 bidirectional transmission links. These values indicate a precision of 63% when including the  
202 direction of transmission, and 64% otherwise.

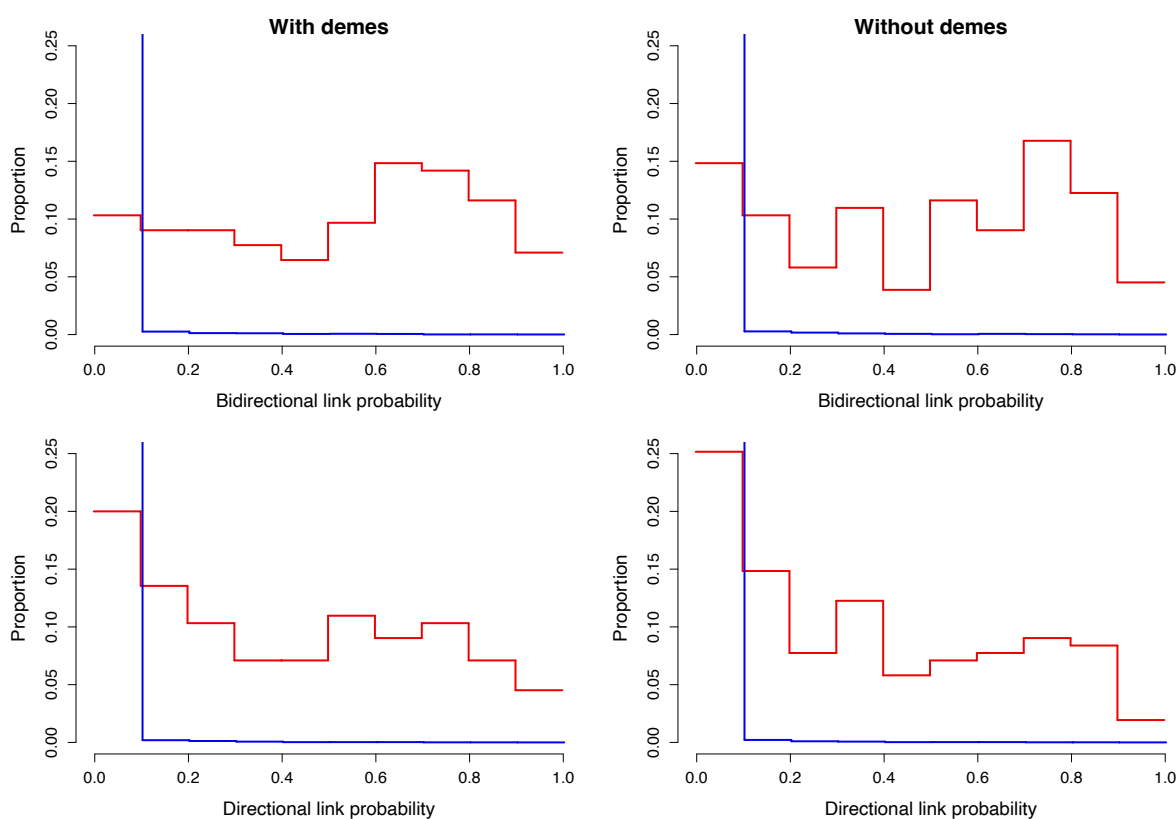


Figure 4: Posterior probability estimates for the transmission links in the second simulation study. The red lines correspond to correct transmission links, and the blue lines correspond to incorrect transmission links. In each case the posterior probability estimates are binned in 0.1 width intervals, and the vertical height indicates the proportion of transmission links contained within each bin. The left plots show estimates with deme information, and the right plots show estimates without deme information. The top plots show bidirectional transmission link estimates, and the bottom plots show directional transmission link estimates..

203 When undertaking inference without deme data we assume that all hosts belong to the same  
204 deme, giving single  $R$  and  $\pi$  parameters, and removing all  $\rho$  parameters. For the posterior  
205 means and credible intervals we find  $R : 1.56 (1.35, 1.79)$ ,  $\pi : 0.73 (0.58, 0.9)$ ,  $\kappa : 0.9 (0.00, 0.25)$ ,  
206 and  $\lambda : 0.28 (0.04, 0.58)$ . That is, we estimate  $R$  and  $\pi$  as somewhere between the pairs of values  
207 used in the simulation. The true values of  $\kappa$  and  $\lambda$  remain in the 95% credible intervals, but

208 the estimate of  $\lambda$  has slightly increased. In this instance we identify fewer correct transmission  
209 links, and more incorrect transmission links. Specifically, 53 links are established if we consider  
210 direction (directional sensitivity of 34%) and 84 if we do not (bidirectional sensitivity of 54%).  
211 We find 47 false positives in the directional case and 50 in the bidirectional case. This means  
212 that the precision is significantly smaller when including the direction of transmission (53%),  
213 but similar when ignoring the direction of transmission (63%).

## 214 Benchmarking using multiple simulations

215 We undertake 50 simulation studies across a range of parameter values, similar in design to  
216 the second simulation study. The parameter sets are sampled from an orthogonal array Latin  
217 hypercube using the lhs R package. In each case we simulate an outbreak with 250 observed  
218 hosts across two demes, with each observed host being sampled once. Each deme has separate  
219  $R$  values sampled between 1 and 6, separate  $\pi$  values sampled between 0.1 and 1, and separate  
220  $\rho$  values sampled between 0.5 and 0.9. Smaller values of  $\rho$  are not considered, as if  $\rho$  is high  
221 in one deme and low in the other this will lead to samples being dominated by one deme. In  
222 such cases we would expect to obtain poor parameter estimates for the less sampled deme. The  
223 remaining parameters are the same for both demes, namely both  $\kappa$  and  $\lambda$  are sampled between  
224 0 and 1. For each simulated dataset we estimate the eight parameters used in the simulation.

225 The marginal posterior results credible intervals are shown in Figure 5 and compared to the  
226 correct values of the parameters used in the simulations. In general we are able to recover  
227 the parameter values used in each simulation, but there is considerable uncertainty on some  
228 of the parameters, as can be seen by the wide credible intervals in Figure 5. The uncertainty  
229 is particular high for parameters of a deme with a small number of representative samples,  
230 as shown by the more lightly shaded bars being longer than the more darkly shaded bars in  
231 Figures 5A-F. The parameter  $\rho$  tends to be inferred more precisely when its correct value is  
232 high (Figure 5E-F). The parameters  $\lambda$  and  $\kappa$  of the within-host population size are only weakly  
233 informed (Figures 5G-H) as previously noted in the two exemplary analyses.

234 Averaging across the 50 simulated data sets, 106.8 out of 250 observed hosts are infected by another  
235 sampled host. Using a posterior probability threshold of 0.5, on average we correctly identify  
236 15.2 transmission links including the direction of transmission, and 26.6 transmission links  
237 ignoring the direction of transmission. These correspond to an average directional sensitivity  
238 of 14%, and an average bidirectional sensitivity of 24%. These relatively low sensitivities are  
239 typical of data sets with one sample per observed host, and data sets with a relaxed bottleneck  
240 (Carson et al., 2024). On average we also find 9.3 false positives including the direction  
241 of transmission, and 14.1 false positives when ignoring the direction of transmission. The  
242 corresponding directional and bidirectional specificities are  $> 99.9\%$ . The average directional  
243 precision is 61%, and the average bidirectional precision is 64%.

244 We obtain the sensitivity (recall), specificity, and precision for each simulation under a series  
245 of different posterior probability thresholds between 0 and 1. Averaging over the 50 simulated  
246 datasets, we present the resulting Receiver Operating Characteristic (ROC) and Precision-  
247 Recall (PR) curves in Figure S1 for both directional and bidirectional transmission links. The  
248 Area Under the Curve (AUC) is 0.99 for both ROC curves, 0.33 for the PR curve of directional  
249 transmission links, and 0.44 for the PR curve of bidirectional transmission links.

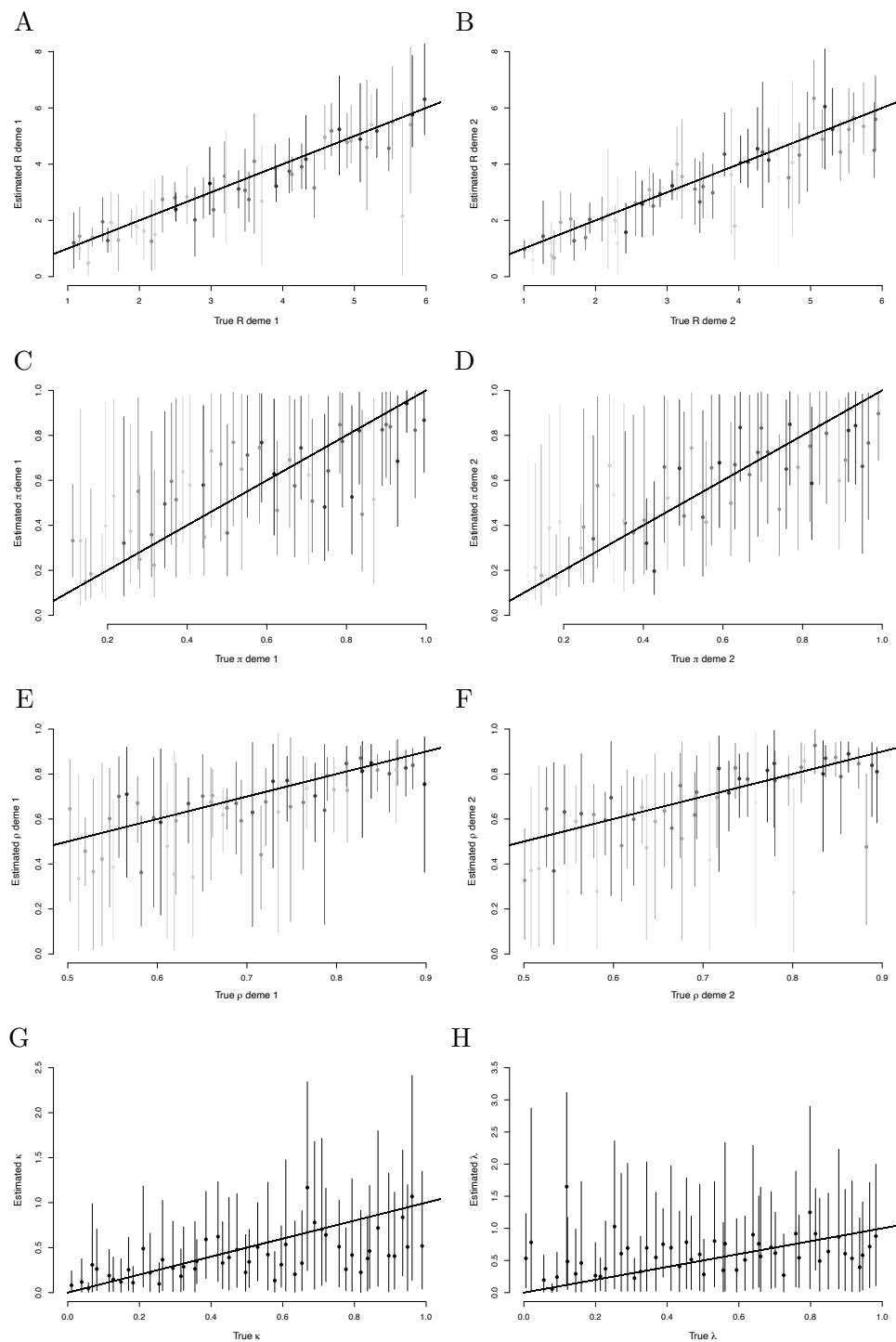


Figure 5: Benchmarking results for the parameters  $R$  of deme 1 (A),  $R$  of deme 2 (B),  $\pi$  of deme 1 (C),  $\pi$  of deme 2 (D),  $\rho$  of deme 1 (E),  $\rho$  of deme 2 (F),  $\kappa$  (G) and  $\lambda$  (H). Mean values are shown by dots, 95% credible intervals are shown by vertical lines. In A-F the shade indicates the proportion of sampled hosts in the associated deme (darker implying a greater proportion in the deme).

## 250 Application to an outbreak of tuberculosis in Argentina

251 A multidrug resistant *Mycobacterium tuberculosis* outbreak in Argentina has been described and  
252 studied in detail using genomic epidemiology (Eldholm et al., 2015). A total of 252 genomes  
253 were sequenced, with collection dates ranging between 1996 and 2010. 153 of the genomes  
254 originated from HIV positive individuals, whereas the remaining 99 genomes were sampled  
255 from HIV negative individuals. A dated phylogeny was previously reconstructed using BEAST  
256 (Drummond et al., 2012), with the root of this tree being estimated to have existed around  
257 1970 (Eldholm et al., 2015). This dated phylogeny is the starting point of our analysis and  
258 reproduced in Figure S2, with leaves colored according to the HIV status of the hosts. The  
259 role of HIV co-infection in the transmission of this tuberculosis outbreak has been previously  
260 investigated and found to be not statistically significant (Eldholm et al., 2016). However, this  
261 analysis was based on a rough reconstruction of transmission events, with a-posteriori testing  
262 of the effect of HIV status, limiting its statistical power (Eldholm et al., 2016). It is therefore  
263 interesting to reanalyse this dataset with the new methodology presented here, considering two  
264 demes for the HIV positive and negative individuals. The sampling window was set from 1st  
265 October 1996 to 1st December 2009 to include all samples and reflect the original sampling  
266 collection methodology (Eldholm et al., 2015). We used the same generation time and sampling  
267 time distributions as was used in previous analyses of tuberculosis outbreaks (Didelot et al.,  
268 2017; Séraphin et al., 2018; Sobkowiak et al., 2023; Chitwood et al., 2024).

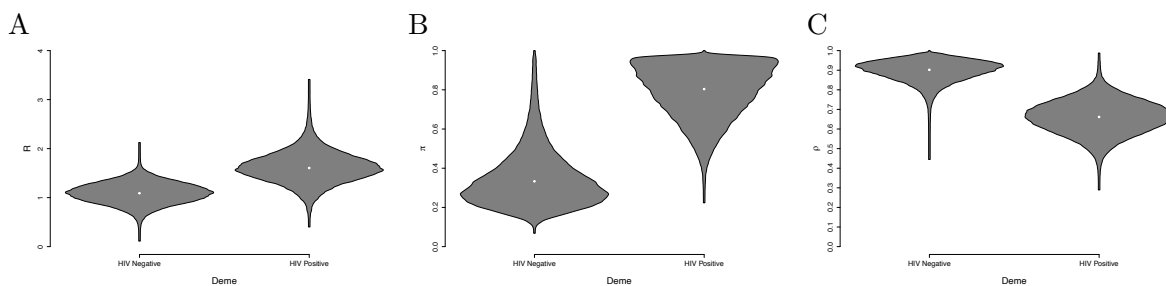


Figure 6: Parameter estimates in the tuberculosis analysis: the reproduction number  $R$  (A), sampling fraction  $\pi$  (B) and probability to remain in a deme  $\rho$  (C) are shown for both the HIV negative (left) and HIV positive (right) demes.

269 The means (95% credible intervals) of the parameters of the within-host population size function  
270 are  $\kappa$  : 5.38 (2.82, 8.79) and  $\lambda$  : 0.63 (0.04, 1.69). The initial pathogen population size  $\kappa$  is large  
271 compared to the per-year linear growth rate  $\lambda$ , suggesting a relaxed transmission bottleneck  
272 (Carson et al., 2024). Figure 6 shows the posterior distribution for the parameters specific to  
273 both demes. The reproduction number  $R$  for the HIV negative and HIV positive demes are  
274 1.10 (0.65, 1.52) and 1.63 (0.86, 2.32), respectively (Figure 6A). The probability that the HIV  
275 positive deme has a greater reproduction number than the HIV negative deme is 0.86. The  
276 sampling probability  $\pi$  for the HIV negative and HIV positive demes are 0.32 (0.13, 0.70) and  
277 0.79 (0.46, 0.99), respectively (Figure 6B). The probability that the HIV positive deme has a  
278 greater sampling probability than the HIV negative deme is 0.98. These results suggest that  
279 HIV positive individuals have a greater reproduction number and are more likely to be observed,  
280 as would be expected from the fact that HIV co-infection accelerates the transition from latent  
281 to active tuberculosis (Bruchfeld et al., 2015).

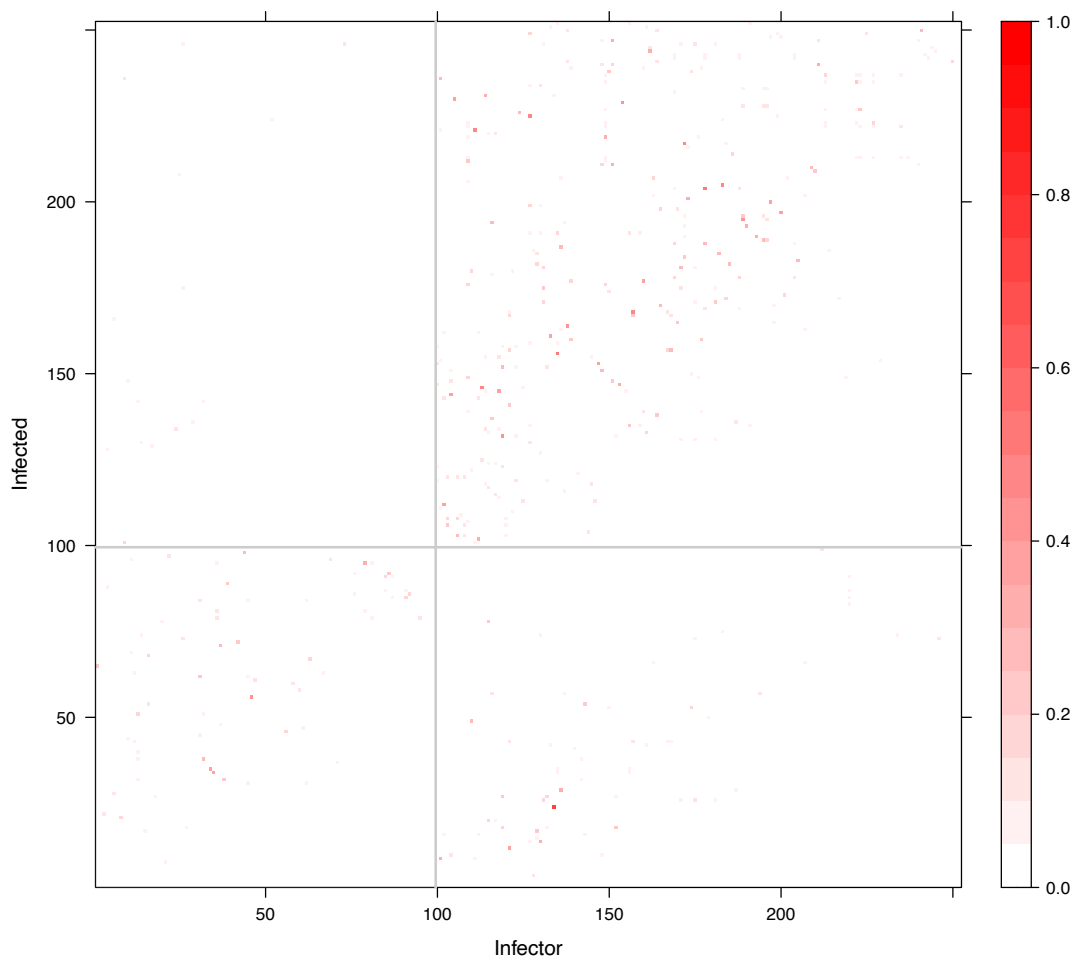


Figure 7: Posterior direct transmission probabilities in the tuberculosis analysis. The grey lines separate the two demes, with HIV negative being left/bottom and HIV positive being right/top.

282 Figure 7 shows the posterior probabilities of direct transmission from any individual to any  
283 other. It is visually clear that infector/infected pairs tend to have the same HIV status. This is  
284 confirmed by the estimate of the probabilities that the pathogen remains in the same deme at  
285 transmission which are 0.90 (0.77, 0.98) and 0.66 (0.48, 0.82), respectively for the HIV negative  
286 and HIV positive demes (Figure 6C). Therefore HIV negative hosts are highly likely to transmit  
287 to other HIV negative hosts. HIV positive hosts are also more likely to transmit to other HIV  
288 positive hosts, but to a lesser extent. This may be in part due to the HIV negative population  
289 being larger than the HIV positive population.

## 290 Application to an outbreak of avian flu H7N7 outbreak in the Netherlands

291 An outbreak of avian influenza H7N7 occurred in the Netherlands during 2003, infecting 255  
292 Dutch farms in less than 3 months, and leading to drastic control measures including the  
293 culling of 30 million birds (Stegeman et al., 2004). Genetic data is available from GISAID (Shu

294 and McCauley, 2017) from 227 farms, for genes HA, NA and PB2 which were concatenated.  
295 Most sequences are from the Gelderland (G) area ( $n = 186$ ) with smaller numbers from the  
296 Limburg (L) area ( $n = 33$ ), Central (C) area ( $n = 7$ ) and Southwest (S) area ( $n = 1$ ). The  
297 phylogeography of this outbreak has been described before in a number of studies (Bataille  
298 et al., 2011; Ypma et al., 2012, 2013; Hall et al., 2015; Klinkenberg et al., 2017). We built  
299 a dated tree using BEAST2 (Bouckaert et al., 2019) which is shown in Figure S3 with leaves  
300 colored by location. We used the same generation time and sampling time distributions as in  
301 a recent study of this outbreak (Klinkenberg et al., 2017). The sampling window was set from  
302 the 50th to the 125th day from the root of the dated tree, which included all samples (Figure  
303 S3).

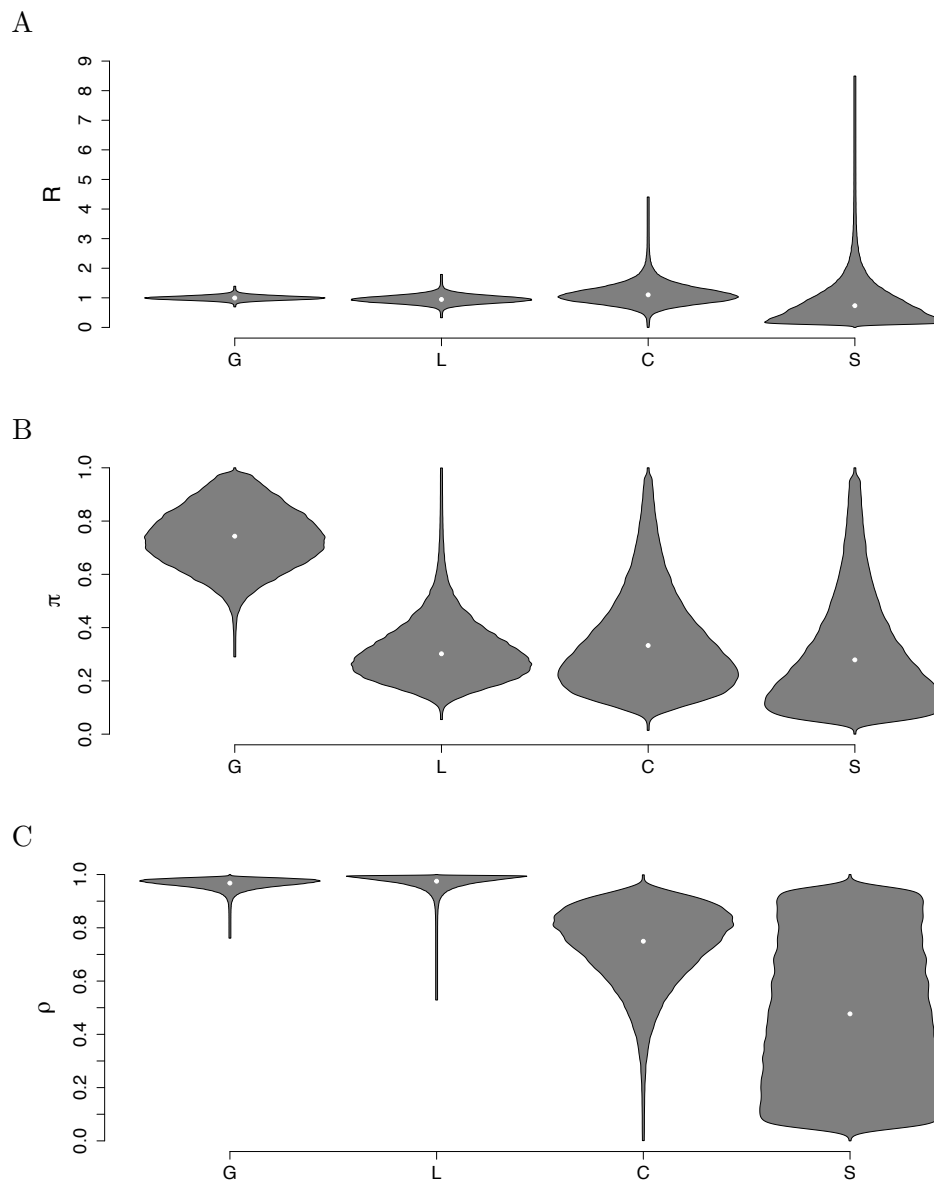


Figure 8: Parameter estimates in the H7N7 analysis: the reproduction number  $R$  (A), sampling proportion  $\pi$  (B) and probability to remain in a deme  $\rho$  (C) are shown for each of the four locations.

304 The means (95% credible intervals) of the parameters of the within-host population size  
305 function are  $\kappa$  : 5.38 (2.23, 9.60) and  $\lambda$  : 3.55 (1.24, 6.28). Figure 8 shows the posterior  
306 distribution for the parameters specific to the four locations. The per-location reproduction  
307 numbers are 1.00 (0.86, 1.16), 0.98 (0.71, 1.25), 1.13 (0.54, 1.93) and 0.74 (0.02, 2.37) for regions  
308 G, L, C and S, respectively. These reproduction numbers are close to 1, with uncertainty  
309 increasing as the sample numbers decrease (Figure 8A). In location S we approximately  
310 recover the prior exponential with mean 1, as would be expected given that there was only  
311 a single representative of this location. The sampling probabilities  $\pi$  are 0.71 (0.36, 0.97),  
312 0.28 (0.10, 0.55), 0.33 (0.07, 0.82), 0.27 (0.02, 0.83) for regions G, L, C and S, respectively.  
313 Location G is best sampled, which makes sense given that it has the largest number of sampled  
314 cases (Figure 8B). The remaining demes are likely less well sampled, although there was high  
315 uncertainty due to the small sample numbers. Our analysis was conducted on 227 farms for  
316 which genetic data was available, whereas during the outbreak 255 farms were confirmed to  
317 be infected [Stegeman et al. \(2004\)](#); [Bataille et al. \(2011\)](#), suggesting an upper bound of the  
318 sampling fraction of 0.89. Our estimates are compatible with this, and further suggest that  
319 some infection went undetected as would be expected from the large scale culling that took  
320 place at farms even in absence of detection [Stegeman et al. \(2004\)](#).

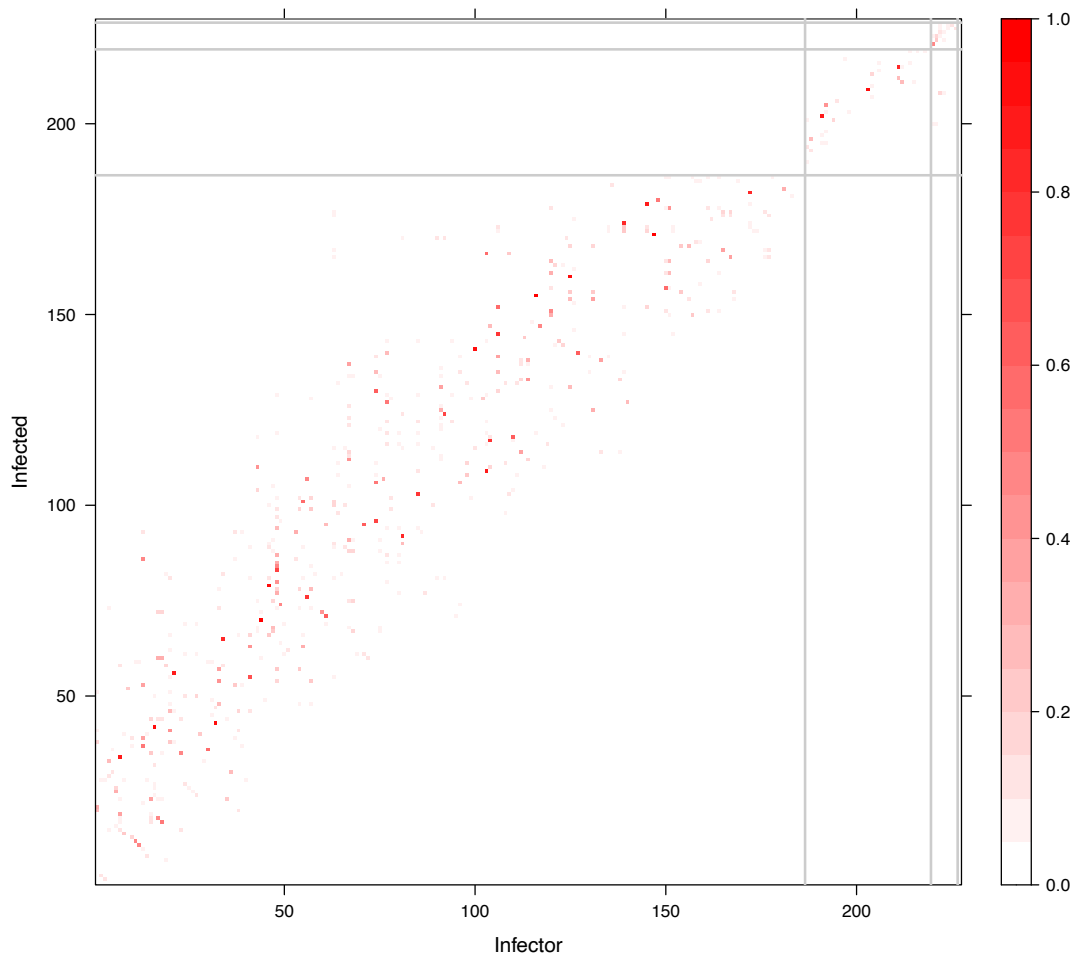


Figure 9: Posterior direct transmission probabilities in the H7N7 analysis. The grey lines separate the four demes.

321 Figure 9 shows the posterior probabilities of direct transmission from any farm to any other,  
322 with a clear tendency for farms to infect other farms from the same region. This is confirmed  
323 by the parameters  $\rho$  representing the probability that the pathogen infects in the same location  
324 which are estimated to be 0.97 (0.92, 0.99), 0.98 (0.92, 1.00), 0.75 (0.4, 0.95) and 0.48 (0.02, 0.97)  
325 for regions G, L, C and S, respectively (Figure 8C). Note that we assume that transmissions to  
326 different locations are evenly distributed across the other locations. Farms in locations G and L  
327 almost always transmit to offspring in the same deme. The probability that a farm in location  
328 C infects another farms in location C also seems high, but is more uncertain due to the small  
329 number of samples. For location S represented by only a single farm we approximately recover  
330 the prior uniform between 0 and 1.

## 331 DISCUSSION

332 Combining genomic and epidemiological data to reconstruct the transmission events within  
333 an infectious disease epidemic is an idea that was formulated over a decade ago, when the  
334 first methods to use genomic data for outbreak reconstruction were proposed (Ypma et al.,  
335 2012; Jombart et al., 2014; Didelot et al., 2014; Hall et al., 2015). It is however difficult to use  
336 epidemiological data when considering partially sampled or ongoing outbreaks (Mollentze et al.,  
337 2014; Didelot et al., 2017), since unsampled cases do not have any associated epidemiological  
338 data. A naive approach to this issue quickly becomes intractable as larger numbers of unsampled  
339 cases need to be considered. Instead we presented a new dynamic programming algorithm that  
340 can efficiently resolve this problem. We implemented this new methodology by extending the  
341 TransPhylo framework which can be applied to partially sampled outbreaks (Didelot et al., 2017,  
342 2021), even when multiple genomes per host are provided or when the transmission bottleneck  
343 is not complete (Carson et al., 2024).

344 We used simulations to show that the combined approach has improved statistical power to  
345 infer the correct transmission events compared to the previous approach based on genomic data  
346 only. We also showed that the parameters governing the epidemiological data can be inferred  
347 with accuracy that increases with the amount of data available for analysis. We applied our  
348 new algorithm to two widely different real datasets to showcase the range of scenarios in which  
349 it can be useful. First we analysed data from a tuberculosis outbreak in Argentina (Eldholm  
350 et al., 2015) to investigate the role of HIV co-infection on the spread of the bacterial causative  
351 agent *M. tuberculosis*. Second we analysed data from the avian influenza H7N7 epidemic that  
352 hit the Netherlands in 2003 (Stegeman et al., 2004), to infer the parameters involved in the  
353 spatial spread of this virus from farm to farm.

354 The methodological framework we developed makes few assumptions, and we therefore envisage  
355 that it can be useful in a wide range of situations. The epidemiological model at the heart of  
356 TransPhylo is a flexible branching process (Didelot et al., 2017) based on offspring distribution  
357 and generation time distribution whose parameters can be set to appropriately model many  
358 infectious diseases transmitted directly from host to host (Wallinga and Teunis, 2004; Grassly  
359 and Fraser, 2008; Cori et al., 2013). An example of application concerns the inference of  
360 the different reproduction numbers for different components of the population. This can help  
361 determine their relative contribution to the overall disease burden, and therefore inform how  
362 to target public health policies for maximum effect (Fraser et al., 2004; Grassly and Fraser,  
363 2008; Hollingsworth, 2009). Another application likely to be useful is to estimate the sampling



364 proportions for different components of the population, which can reveal if sampling is currently  
365 biased and how it could be improved (Magnani et al., 2005; Brooks-Pollock et al., 2021; Layan  
366 et al., 2023).

367 There are however some limitations to the methodology we presented. First, the analysis is based  
368 on a dated phylogeny that needs to be correctly precomputed. This requires consideration of  
369 how such a tree is computed, under which prior model if a Bayesian method is used, and to  
370 what extent a single point estimate can be used without quantification of uncertainty. These  
371 questions arise for all of the many recently developed phylodynamic methods that take a dated  
372 tree as input (Didelot and Parkhill, 2022). However, this step-by-step approach is necessary  
373 to be able to analyse state-of-the-art large genomic datasets. Second, we only considered  
374 discrete epidemiological data. Continuous variables can always be discretised to circumvent  
375 this limitation, but doing so may lose some information and requires to define potentially  
376 arbitrary discrete classes. This situation is analogous to the almost ubiquitous use of discrete  
377 locations in phylogeography (Lemey et al., 2009; De Maio et al., 2015; Baele et al., 2017). In our  
378 applications, we also assume that transmission outside of a deme is evenly distributed across  
379 the other demes. We make this choice in order to restrict the number of parameters to linear in  
380 the number of demes instead of quadratic, but recognise that it will not be appropriate for every  
381 situation. The dynamic programming algorithm allows for this assumption to be relaxed, but a  
382 more complex MCMC algorithm would be required to efficiently estimate the additional model  
383 parameters, and these estimates would likely exhibit greater uncertainty. Finally it should be  
384 noted that our methodology has a non-negligible computational cost. For example, the largest  
385 analysis we performed, on the H7N7 dataset, took several days to achieve acceptable MCMC  
386 convergence and mixing properties. However, our algorithm currently runs only on a single  
387 CPU core, whereas most standard desktop and laptop computers have 8 to 16 cores, with many  
388 more cores available on servers dedicated to computer-intensive tasks. Future work should  
389 therefore seek to exploit multiple cores to reduce the overall runtime, for example by following  
390 recent progress in parallel MCMC algorithms (Schwedes and Calderhead, 2021; Syed et al.,  
391 2022; Glatt-Holtz et al., 2024).

## 392 METHODS

### 393 Case where all demes have the same offspring distribution and sampling 394 probabilities

395 Let us start with the simpler case where all demes are assumed to have the same offspring  
396 distribution and sampling probabilities. In this case most of the calculations in TransPhylo  
397 (Didelot et al., 2017; Carson et al., 2024) remain unchanged, we simply include an additional  
398 likelihood term obtained by using an efficient dynamic programming algorithm similar to  
399 Felsenstein’s tree-pruning algorithm (Felsenstein, 1973, 1981). This is necessary to integrate  
400 over demes for unsampled individuals that form part of the transmission tree and for sampled  
401 individuals with missing deme data.

402 Let hosts be labelled  $1, \dots, N$ , and define the deme of host  $n$  by  $s^n \in \{1, \dots, S\}$ . Let  $P_{ij}$  be the  
403 probability that an offspring of a host in deme  $i$  is in deme  $j$ . Finally, let  $\mathcal{L}_s^n$  be the likelihood  
404 from the deme data of host  $n$  and their descendants, conditional on host  $n$  being in deme  $s$ .

405 The algorithm is initialised at the leaf nodes, which in this case are hosts with no offspring in  
 406 the transmission tree, noting that all such hosts must have been sampled in TransPhylo. If the  
 407 deme of such a host is known, then  $\mathcal{L}_{s^n}^n = 1$  at the hosts deme  $s^n$ , and  $\mathcal{L}_s^n = 0$  for  $s \neq s^n$  (all  
 408 other demes). If the deme of the host is unknown, then  $\mathcal{L}_s^n = 1$  for all possible demes  $s$ .

409 The algorithm then proceeds backwards in time to evaluate the conditional likelihoods of the  
 410 internal nodes (hosts with offspring). Let  $\mathcal{H}^n$  denote the set of offspring of Host  $n$ . If Host  $n$   
 411 has a known deme  $s^n$  then

$$\mathcal{L}_{s^n}^n = \prod_{j \in \mathcal{H}^n} \sum_{s^j=1}^S P_{s^n s^j} \mathcal{L}_{s^j}^j, \quad (1)$$

412 and  $\mathcal{L}_s^n = 0$  for  $s \neq s^n$ . If on the other hand Host  $n$  has no known deme then

$$\mathcal{L}_s^n = \prod_{j \in \mathcal{H}^n} \sum_{s^j=1}^S P_{s s^j} \mathcal{L}_{s^j}^j, \quad (2)$$

413 for all possible values of  $s$ .

414 The algorithm terminates at the root host, assumed here to be  $n = 1$ . The overall likelihood of  
 415 the demes on the transmission tree is given by

$$\mathcal{L} = \sum_{s=1}^S \zeta_s \mathcal{L}_s^1, \quad (3)$$

416 where  $\zeta_s$  is the prior probability of the root host being in deme  $s$ .

## 417 Illustrative example

418 An example of the dynamic programming algorithm is shown in Figure [S4A](#). The target  
 419 transmission tree contains 10 hosts and we assume that there are three possible demes. The  
 420 transition probability between the three demes is given by

$$P = \begin{pmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{pmatrix}, \quad (4)$$

421 meaning that a host has probability 0.8 of an offspring having the same deme as its infector,  
 422 and a probability 0.1 of an offspring being in either of the other two demes. We know that Host  
 423 5 is in deme 1, Hosts 3 and 7 are in deme 2, and Host 10 is in deme 3. Hosts 2, 4, and 9 are  
 424 sampled hosts, but their deme is missing. Hosts 1, 6 and 8 are unsampled hosts.

425 Each host has an associated vector for the conditional likelihood at the three demes. Working  
 426 backwards in time, Host 10 is known to be in deme 3 and is a leaf, and so:

$$\mathcal{L}^{10} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}. \quad (5)$$

427 Host 9 is a leaf, but does not have a known deme, so that:

$$\mathcal{L}^9 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}. \quad (6)$$

428 Host 8 is an unsampled individual, whose only offspring is Host 10:

$$\mathcal{L}^8 = \begin{pmatrix} 0.8 \cdot 0 + 0.1 \cdot 0 + 0.1 \cdot 1 \\ 0.1 \cdot 0 + 0.8 \cdot 0 + 0.1 \cdot 1 \\ 0.1 \cdot 0 + 0.1 \cdot 0 + 0.8 \cdot 1 \end{pmatrix} = \begin{pmatrix} 0.1 \\ 0.1 \\ 0.8 \end{pmatrix}. \quad (7)$$

429 Host 7 is another leaf with deme 2:

$$\mathcal{L}^7 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}. \quad (8)$$

430 Host 6 is an unsampled individual, whose only offspring is Host 9:

$$\mathcal{L}^6 = \begin{pmatrix} 0.8 \cdot 1 + 0.1 \cdot 1 + 0.1 \cdot 1 \\ 0.1 \cdot 1 + 0.8 \cdot 1 + 0.1 \cdot 1 \\ 0.1 \cdot 1 + 0.1 \cdot 1 + 0.8 \cdot 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}. \quad (9)$$

431 Host 5 is a leaf with deme 1:

$$\mathcal{L}^5 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}. \quad (10)$$

432 Host 4 is a leaf with no deme data:

$$\mathcal{L}^4 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}. \quad (11)$$

433 Host 3 has three offspring: Hosts 4, 8, and 7. Additionally, Host 3 is in deme 2, and so we only  
434 calculate the second element of the vector:

$$\mathcal{L}^3 = \begin{pmatrix} 0.0 \\ (0.1 \cdot 1 + 0.8 \cdot 1 + 0.1 \cdot 1)(0.1 \cdot 0 + 0.8 \cdot 1 + 0.1 \cdot 0)(0.1 \cdot 0.1 + 0.8 \cdot 0.1 + 0.1 \cdot 0.8) \\ 0.0 \end{pmatrix} = \begin{pmatrix} 0.000 \\ 0.136 \\ 0.000 \end{pmatrix}. \quad (12)$$

435 Host 2 has two offspring: Hosts 5 and 6. Host 2 is sampled, but has no deme data, and so all  
436 elements of the vector are evaluated:

$$\mathcal{L}^2 = \begin{pmatrix} (0.8 \cdot 1 + 0.1 \cdot 0 + 0.1 \cdot 0)(0.8 \cdot 1 + 0.1 \cdot 1 + 0.1 \cdot 1) \\ (0.1 \cdot 1 + 0.8 \cdot 0 + 0.1 \cdot 0)(0.1 \cdot 1 + 0.8 \cdot 1 + 0.1 \cdot 1) \\ (0.1 \cdot 1 + 0.1 \cdot 0 + 0.8 \cdot 0)(0.1 \cdot 1 + 0.1 \cdot 1 + 0.8 \cdot 1) \end{pmatrix} = \begin{pmatrix} 0.8 \\ 0.1 \\ 0.1 \end{pmatrix}. \quad (13)$$

437 Finally, Host 1 has two offspring: Hosts 2 and 3:

$$\mathcal{L}^1 = \begin{pmatrix} (0.8 \cdot 0.8 + 0.1 \cdot 0.1 + 0.1 \cdot 0.1)(0.8 \cdot 0.0 + 0.1 \cdot 0.136 + 0.1 \cdot 0.0) \\ (0.1 \cdot 0.8 + 0.8 \cdot 0.1 + 0.1 \cdot 0.1)(0.1 \cdot 0.0 + 0.8 \cdot 0.136 + 0.1 \cdot 0.0) \\ (0.1 \cdot 0.8 + 0.1 \cdot 0.1 + 0.8 \cdot 0.1)(0.1 \cdot 0.0 + 0.1 \cdot 0.136 + 0.8 \cdot 0.0) \end{pmatrix} = \begin{pmatrix} 0.008976 \\ 0.018496 \\ 0.002312 \end{pmatrix}. \quad (14)$$

438 Figure [S4B](#) shows the transmission tree annotated with the conditional likelihoods calculated  
439 in the dynamic programming algorithm. If we assume that the prior for the deme of the root

440 host is 1/3 for the three demes, then the likelihood of the deme is the mean of the values for  
441 Host 1. In this case  $\mathcal{L} = 0.009928$ . This is verified by brute force by calculating

$$\mathcal{L} = \sum_{s^1=1}^S \sum_{s^2=1}^S \sum_{s^3=1}^S \sum_{s^4=1}^S \sum_{s^5=1}^S \sum_{s^6=1}^S \sum_{s^7=1}^S \sum_{s^8=1}^S \sum_{s^9=1}^S \pi_{s^1} P_{s^1 s^2} P_{s^1 s^3} P_{s^2 s^5} P_{s^2 s^6} P_{s^3 s^4} P_{s^3 s^7} P_{s^3 s^8} P_{s^6 s^9} P_{s^8 s^{10}}. \quad (15)$$

442 Note that leaves with no deme data do not ultimately contribute to the likelihood, and can  
443 therefore be excluded.

#### 444 Case where the demes may have different offspring distributions and sampling 445 probabilities

446 A useful extension would be to allow  $R$  and/or  $\pi$  to change based on deme. For now, let us  
447 assume that there are  $S = 2$  demes with offspring distribution  $\alpha_1(k)$  and  $\alpha_2(k)$ . Host infected  
448 at time  $t$  are observed with probability  $\pi_1$  and  $\pi_2$  (leading to time-dependent probabilities  $\zeta_1(t)$   
449 and  $\zeta_2(t)$ ). Unlike the previous case, the transmission tree likelihood and deme likelihood can  
450 not be calculated separately. To evaluate the combined likelihood, we start by calculating the  
451 exclusion probabilities as follows.

452 Define  $\omega_1(t)$  as the exclusion probability of a host infected at time  $t$  in deme 1, and  $\omega_2(t)$  as  
453 the exclusion probability of a host infected at time  $t$  in deme 2. Assuming that  $T$  is the cut-off  
454 time for observations  $\omega_1(t) = \omega_2(t) = 1$  for  $t \geq T$ . We can then define the following recursive  
455 relationships.

456 The exclusion probability of an offspring from a host in deme 1 infected at time  $t$  is

$$\bar{\omega}_1(t) = \int_0^\infty (P_{11}\omega_1(t+\tau) + P_{12}\omega_2(t+\tau))\gamma(\tau)d\tau, \quad (16)$$

457 and for deme 2,

$$\bar{\omega}_2(t) = \int_0^\infty (P_{21}\omega_1(t+\tau) + P_{22}\omega_2(t+\tau))\gamma(\tau)d\tau. \quad (17)$$

458 The probability that all offspring from an individual in deme  $i$  infected at time  $t$  are excluded  
459 is

$$\phi_i(t) = \sum_{k=0}^\infty \alpha(k)\bar{\omega}_i(t)^k. \quad (18)$$

460 The exclusion probability of an individual in deme  $i$  infected at time  $t$  is then

$$\omega_i(t) = (1 - \zeta_i(t))\phi_i(t). \quad (19)$$

461 That is, the probability of the host being unobserved and having no included offspring.

462 As established in [Carson et al. \(2024\)](#), the transmission tree likelihood contribution from an  
463 unsampled Host  $n$  is

$$\frac{(1 - \zeta(x^n))}{1 - \omega(x^n)} \sum_{k=d^n}^\infty \alpha(k) \binom{k}{d^n} \bar{\omega}(x^n)^{k-d^n} d^n! \prod_{j \in \mathcal{H}^n} (1 - \omega(x^j))\gamma(x^j - x^n), \quad (20)$$

464 where  $x^n$  is the host's infection time, and  $d^n$  is the number of included offspring. If Host  $n$  is  
 465 sampled, the likelihood contribution is

$$\frac{\pi\sigma(y^n - x^n)}{1 - \omega(x^n)} \sum_{k=d^n}^{\infty} \alpha(k) \binom{k}{d^n} \bar{\omega}(x^n)^{k-d^n} d^n! \prod_{j \in \mathcal{H}^n} (1 - \omega(x^j)) \gamma(x^j - x^n), \quad (21)$$

466 where  $y^n$  is the host's primary observation time, and  $\sigma(\tau)$  is the observation time distribution.  
 467 Here, we define

$$\mathcal{T}_s^n = \frac{(1 - \zeta_s(x^n))}{1 - \omega_s(x^n)} \sum_{k=d^n}^{\infty} \alpha_s(k) \binom{k}{d^n} \bar{\omega}_s(x^n)^{k-d^n} d^n! \quad (22)$$

468 for an unobserved Host  $n$  in deme  $s$ , and

$$\mathcal{T}_s^n = \frac{\pi_s\sigma(y^n - x^n)}{1 - \omega_s(x^n)} \sum_{k=d^n}^{\infty} \alpha_s(k) \binom{k}{d^n} \bar{\omega}_s(x^n)^{k-d^n} d^n! \quad (23)$$

469 for an observed Host  $n$  in deme  $s$ . In addition we define

$$\mathcal{U}_{ss^j}^{n,j} = (1 - \omega_{s^j}(x^j)) \gamma(x^j - x^n) P_{ss^j} \quad (24)$$

470 for  $j \in \mathcal{H}^n$  being the offspring of Host  $n$ , and  $s^j$  being the deme of the offspring. Finally, define  
 471

$$\mathcal{L}_s^n = \mathcal{T}_s^n \quad (25)$$

472 for leaf hosts, and

$$\mathcal{L}_s^n = \mathcal{T}_s^n \prod_{j \in \mathcal{H}^n} \sum_{s^j=1}^S \mathcal{U}_{ss^j}^{n,j} \mathcal{L}_{s^j}^j \quad (26)$$

473 for hosts with offspring. The combined transmission tree and deme likelihood is then calculated  
 474 using dynamic programming with this replacement definition of  $\mathcal{L}_s^n$ .

## 475 Illustrative example

476 We return to the transmission tree presented in Figure [S4A](#). We include observation times as  
 477 follows:

Host	Deme	Infection time	Observation time	Offspring
1	-	0.0	-	2, 3
2	-	2.6	3.5	5, 6
3	2	3.2	5.3	4, 7, 8
4	-	5.1	6.9	-
5	1	5.2	8.5	-
6	-	5.5	-	9
7	2	6.3	7.1	-
8	-	7.1	-	10
9	-	8.9	11.0	-
10	3	9.8	11.6	-

479 We assume that the demes have basic reproduction number  $R$  equal to 2, 1.5 and 1, respectively,  
480 and sampling proportion  $\pi$  equal to 0.5, 0.7 and 0.9, respectively. We again set  $P$  as in  
481 Equation (4). Both the generation time distribution and observation time distribution are  
482 Gamma distributed with shape 2 and scale 1. The resulting exclusion probabilities are shown  
483 in Figure S5.

484 The resulting conditional likelihoods are as follows:

Host	Deme 1	Deme 2	Deme 3
1	$1.35 \times 10^{-19}$	$1.24 \times 10^{-19}$	$2.60 \times 10^{-21}$
2	$9.71 \times 10^{-8}$	$1.07 \times 10^{-8}$	$4.86 \times 10^{-9}$
3	0	$3.08 \times 10^{-9}$	0
4	$5.67 \times 10^{-2}$	$9.84 \times 10^{-2}$	$1.53 \times 10^{-1}$
485 5	$2.34 \times 10^{-2}$	0	0
6	$2.41 \times 10^{-3}$	$1.72 \times 10^{-3}$	$6.20 \times 10^{-4}$
7	0	$1.30 \times 10^{-1}$	0
8	$1.40 \times 10^{-3}$	$7.30 \times 10^{-4}$	$1.80 \times 10^{-3}$
9	$1.25 \times 10^{-1}$	$1.68 \times 10^{-1}$	$2.14 \times 10^{-1}$
10	0	0	$3.61 \times 10^{-1}$

486 The overall likelihood is  $\mathcal{L} = 8.69 \times 10^{-20}$ , which again is confirmed by using a brute force  
487 calculation.

488 As a further check we recalculate the likelihood under  $R = (2, 2, 2)$  and  $\pi = (0.8, 0.8, 0.8)$ .  
489 As the demes now have the same offspring distribution and sampling probabilities, we should  
490 obtain the same likelihood by taking the product of the transmission tree and deme likelihoods,  
491 as in the case where all demes have the same offspring distribution and sampling probabilities.  
492 We find that both approaches do indeed return the same likelihood.

## 493 Implementation

494 We implemented the methods above into a new R package called TransPhylo2, which  
495 extends TransPhyloMulti (Carson et al., 2024) and therefore inherits the same advantages  
496 over the previous implementation of TransPhylo (Didelot et al., 2017) in terms of allowing  
497 multiple samples per host and relaxing the assumption of a complete transmission bottleneck.  
498 TransPhylo2 is available at <https://github.com/DrJCarson/TransPhylo2>. This repository  
499 also contains all the code and data needed to reproduce all results shown in this paper. The  
500 R package ape was used to store, manipulate and visualise phylogenetic trees (Paradis and  
501 Schliep, 2019).

## 502 ACKNOWLEDGEMENTS

503 We acknowledge funding from the National Institute for Health Research (NIHR) Health  
504 Protection Research Unit in Genomics and Enabling Data (grant number NIHR200892).

## 505 References

- 506 Baele G, Suchard MA, Rambaut A, Lemey P. 2017. Emerging concepts of data integration in  
507 pathogen phylodynamics. *Systematic biology*. 66:e47–e65.
- 508 Bataille A, Van Der Meer F, Stegeman A, Koch G. 2011. Evolutionary Analysis of Inter-Farm  
509 Transmission Dynamics in a Highly Pathogenic Avian Influenza Epidemic. *PLoS Pathogens*.  
510 7:e1002094.
- 511 Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. 2015. Measurably evolving pathogens in the  
512 genomic era. *Trends in Ecology & Evolution*. 30:306–313.
- 513 Bouckaert R, Vaughan TG, Fourment M, Gavryushkina A, Heled J, Denise K, Maio ND,  
514 Matschiner M, Ogilvie H, Plessis L, et al. (11 co-authors). 2019. BEAST 2.5 : An  
515 Advanced Software Platform for Bayesian Evolutionary Analysis. *PLoS computational  
516 biology*. 15:e1006650.
- 517 Brooks SPB, Gelman AG. 1998. General methods for monitoring convergence of iterative  
518 simulations. *Journal of computational and graphical statistics*. 7:434–455.
- 519 Brooks-Pollock E, Danon L, Jombart T, Pellis L. 2021. Modelling that shaped the early COVID-  
520 19 pandemic response in the UK. *Philosophical Transactions of the Royal Society B: Biological  
521 Sciences*. 376:20210001.
- 522 Bruchfeld J, Correia-Neves M, Källenius G. 2015. Tuberculosis and HIV Coinfection. *Cold  
523 Spring Harbor Perspectives in Medicine*. 5:a017871.
- 524 Campbell F, Strang C, Ferguson N, Cori A, Jombart T. 2018. When are pathogen genome  
525 sequences informative of transmission events? *PLOS Pathogens*. 14:e1006885.
- 526 Carson J, Keeling M, Wyllie D, Ribeca P, Didelot X. 2024. Inference of infectious disease  
527 transmission through a relaxed bottleneck using multiple genomes per host. *Molecular Biology  
528 and Evolution*. 41:msad288.
- 529 Chan CH, McCabe CJ, Fisman DN. 2012. Core groups, antimicrobial resistance and rebound  
530 in gonorrhoea in North America. *Sexually Transmitted Infections*. 88:200–204.
- 531 Chis Ster I, Singh BK, Ferguson NM. 2009. Epidemiological inference for partially observed  
532 epidemics: The example of the 2001 foot and mouth epidemic in Great Britain. *Epidemics*.  
533 1:21–34.
- 534 Chitwood MH, Corbett EL, Ndhlovu V, Sobkowiak B, Colijn C, Andrews JR, Burke RM,  
535 Cudahy PG, Dodd PJ, Imai-Eaton JW, et al. (23 co-authors). 2024. Distribution and  
536 transmission of *M. tuberculosis* in a high-HIV prevalence city in Malawi: A genomic and  
537 spatial analysis. medRxiv. medRxiv:2024.05.17.24307525.
- 538 Cori A, Ferguson NM, Fraser C, Cauchemez S. 2013. A new framework and software to estimate  
539 time-varying reproduction numbers during epidemics. *American journal of epidemiology*.  
540 178:1505–12.
- 541 Craig AP, Gray RT, Edwards JL, Apicella MA, Jennings MP, Wilson DP, Seib KL. 2015. The  
542 potential impact of vaccination on the prevalence of gonorrhoea. *Vaccine*. 33:4520–4525.
- 543 Croucher NJ, Didelot X. 2015. The application of genomics to tracing bacterial pathogen  
544 transmission. *Current Opinion in Microbiology*. 23:62–67.

- 545 De Maio N, Wu CH, O'Reilly KM, Wilson D. 2015. New Routes to Phylogeography: A Bayesian  
546 Structured Coalescent Approximation. *PLoS Genetics*. 11:e1005421.
- 547 Didelot X, Fraser C, Gardy J, Colijn C. 2017. Genomic infectious disease epidemiology in  
548 partially sampled and ongoing outbreaks. *Molecular Biology and Evolution*. 34:997–1007.
- 549 Didelot X, Gardy J, Colijn C. 2014. Bayesian inference of infectious disease transmission from  
550 whole genome sequence data. *Molecular Biology and Evolution*. 31:1869–1879.
- 551 Didelot X, Kendall M, Xu Y, White PJ, McCarthy N. 2021. Genomic Epidemiology Analysis  
552 of Infectious Disease Outbreaks Using TransPhylo. *Current Protocols*. 1:e60.
- 553 Didelot X, Parkhill J. 2022. A scalable analytical approach from bacterial genomes to  
554 epidemiology. *Philosophical Transactions of the Royal Society B: Biological Sciences*.  
555 377:20210246.
- 556 Drummond AJ, Suchard MA, Xie D, Rambaut A. 2012. Bayesian phylogenetics with BEAUti  
557 and the BEAST 1.7. *Molecular Biology and Evolution*. 29:1969–1973.
- 558 Duault H, Durand B, Canini L. 2022. Methods Combining Genomic and Epidemiological Data  
559 in the Reconstruction of Transmission Trees: A Systematic Review. *Pathogens*. 11:252.
- 560 Eldholm V, Monteserin J, Rieux A, Lopez B, Sobkowiak B, Ritacco V, Balloux F. 2015. Four  
561 decades of transmission of a multidrug-resistant *Mycobacterium tuberculosis* outbreak strain.  
562 *Nature communications*. 6:7119.
- 563 Eldholm V, Rieux A, Monteserin J, Lopez JM, Palmero D, Lopez B, Ritacco V, Didelot X,  
564 Balloux F. 2016. Impact of HIV co-infection on the evolution and transmission of multidrug-  
565 resistant tuberculosis. *eLife*. 5:e16644.
- 566 Farrington CP, Kanaan MN, Gay NJ. 2003. Branching process models for surveillance of  
567 infectious diseases controlled by mass vaccination. *Biostatistics (Oxford, England)*. 4:279–95.
- 568 Felsenstein J. 1973. Maximum Likelihood and Minimum-Steps Methods for Estimating  
569 Evolutionary Trees from Data on Discrete Characters. *Systematic Zoology*. 22:240–49.
- 570 Felsenstein J. 1981. Evolutionary trees from DNA sequences: A maximum likelihood approach.  
571 *Journal of Molecular Evolution*. 17:368–376.
- 572 Fingerhuth SM, Bonhoeffer S, Low N, Althaus CL. 2016. Antibiotic-Resistant *Neisseria*  
573 *gonorrhoeae* Spread Faster with More Treatment, Not More Sexual Partners. *PLOS*  
574 *Pathogens*. 12:e1005611.
- 575 Fraser C, Riley S, Anderson RM, Ferguson NM. 2004. Factors that make an infectious disease  
576 outbreak controllable. *Proceedings of the National Academy of Sciences*. 101:6146–6151.
- 577 Glatt-Holtz NE, Holbrook AJ, Krometis JA, Mondaini CF. 2024. Parallel MCMC algorithms:  
578 Theoretical foundations, algorithm design, case studies. *Transactions of Mathematics and Its*  
579 *Applications*. 8:tnae004.
- 580 Grassly NC, Fraser C. 2008. Mathematical models of infectious disease transmission. *Nature*  
581 *Reviews Microbiology*. 6:477–87.
- 582 Green PJ. 1995. Reversible Jump Markov Chain Monte Carlo Computation and Bayesian Model  
583 Determination. *Biometrika*. 82:711–732.



- 584 Hall M, Woolhouse M, Rambaut A. 2015. Epidemic Reconstruction in a Phylogenetics  
585 Framework: Transmission Trees as Partitions of the Node Set. *PLOS Computational Biology*.  
586 11:e1004613.
- 587 Hatherell HA, Didelot X, Pollock SL, Tang P, Crisan A, Johnston JC, Colijn C, Gardy J. 2016.  
588 Declaring a tuberculosis outbreak over with genomic epidemiology. *Microbial Genomics*.  
589 1:10.1099/mgen.0.000060.
- 590 Hollingsworth TD. 2009. Controlling infectious disease outbreaks: Lessons from mathematical  
591 modelling. *Journal of Public Health Policy*. 30:328–341.
- 592 Jewell CP, Keeling MJ, Roberts GO. 2009. Predicting undetected infections during the 2007  
593 foot-and-mouth disease outbreak. *Journal of The Royal Society Interface*. 6:1145–1151.
- 594 Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N. 2014. Bayesian  
595 Reconstruction of Disease Outbreaks by Combining Epidemiologic and Genomic Data. *PLoS*  
596 *Computational Biology*. 10:e1003457.
- 597 Klinkenberg D, Backer JA, Didelot X, Colijn C, Wallinga J. 2017. Simultaneous inference of  
598 phylogenetic and transmission trees in infectious disease outbreaks. *PLoS Computational*  
599 *Biology*. 13:e1005495.
- 600 Layan M, Müller NF, Dellicour S, De Maio N, Bourhy H, Cauchemez S, Baele G. 2023. Impact  
601 and mitigation of sampling bias to determine viral spread: Evaluating discrete phylogeography  
602 through CTMC modeling and structured coalescent model approximations. *Virus Evolution*.  
603 9:vead010.
- 604 Lemey P, Rambaut A, Drummond AJ, Suchard MA. 2009. Bayesian phylogeography finds its  
605 roots. *PLoS computational biology*. 5:e1000520.
- 606 Magnani R, Sabin K, Saidel T, Heckathorn D. 2005. Review of sampling hard-to-reach and  
607 hidden populations for HIV surveillance. *AIDS*. 19:S67–S72.
- 608 Mollentze N, Nel LH, Townsend S, le Roux K, Hampson K, Haydon DT, Soubeyrand S. 2014.  
609 A Bayesian approach for inferring the dynamics of partially observed endemic infectious  
610 diseases from space-time-genetic data. *Proceedings of the Royal Society B: Biological Sciences*.  
611 281:20133251.
- 612 Morelli MJ, Thébaud G, Chadœuf J, King DP, Haydon DT, Soubeyrand S. 2012. A Bayesian  
613 Inference Framework to Reconstruct Transmission Trees Using Epidemiological and Genetic  
614 Data. *PLoS Computational Biology*. 8:e1002768.
- 615 O’Neill PD. 2002. A tutorial introduction to Bayesian inference for stochastic epidemic models  
616 using Markov chain Monte Carlo methods. *Mathematical biosciences*. 180:103–14.
- 617 O’Neill PD, Roberts GO. 1999. Bayesian inference for partially observed stochastic epidemics.  
618 *Journal of the Royal Statistical Society: Series A*. 162:121–129.
- 619 Paradis E, Schliep K. 2019. Ape 5.0: An environment for modern phylogenetics and evolutionary  
620 analyses in R. *Bioinformatics*. 35:526–528.
- 621 Rieux A, Balloux F. 2016. Inferences from tip-calibrated phylogenies: A review and a practical  
622 guide. *Molecular Ecology*. 25:1911–1924.
- 623 Schwedes T, Calderhead B. 2021. Rao-Blackwellised parallel MCMC. *Aistats*. 130.

- 624 Séraphin MN, Didelot X, Nolan DJ, May JR, Khan MSR, Murray ER, Salemi M, Morris  
625 JG, Lauzardo M. 2018. Genomic Investigation of a *Mycobacterium tuberculosis* Outbreak  
626 Involving Prison and Community Cases in Florida, United States. *American Journal of*  
627 *Tropical Medicine and Hygiene*. 99:867–874.
- 628 Shu Y, McCauley J. 2017. GISAID: Global initiative on sharing all influenza data – from vision  
629 to reality. *Eurosurveillance*. 22.
- 630 Sisson SA. 2005. Transdimensional Markov chains: A decade of progress and future perspectives.  
631 *Journal of the American Statistical Association*. 100:1077–1089.
- 632 Sobkowiak B, Romanowski K, Sekirov I, Gardy JL, Johnston JC. 2023. Comparing  
633 *Mycobacterium tuberculosis* transmission reconstruction models from whole genome sequence  
634 data. *Epidemiology and Infection*. 151:e105.
- 635 Stegeman A, Bouma A, Elbers ARW, de Jong MCM, Nodelijk G, de Klerk F, Koch G, van Boven  
636 M. 2004. Avian Influenza A Virus (H7N7) Epidemic in The Netherlands in 2003: Course of  
637 the Epidemic and Effectiveness of Control Measures. *The Journal of Infectious Diseases*.  
638 190:2088–2095.
- 639 Syed S, Bouchard-Côté A, Deligiannidis G, Doucet A. 2022. Non-Reversible Parallel Tempering:  
640 A Scalable Highly Parallel MCMC Scheme. *Journal of the Royal Statistical Society Series B:*  
641 *Statistical Methodology*. 84:321–350.
- 642 van Dyk DA, Meng XL. 2001. The Art of Data Augmentation. *Journal of Computational and*  
643 *Graphical Statistics*. 10:1–50.
- 644 Wallinga J, Teunis P. 2004. Different Epidemic Curves for Severe Acute Respiratory Syndrome  
645 Reveal Similar Impacts of Control Measures. *American Journal of Epidemiology*. 160:509–516.
- 646 Whittles LK, Didelot X, White PJ. 2022. Public health impact and cost-effectiveness of  
647 gonorrhoea vaccination: An integrated transmission-dynamic health-economic modelling  
648 analysis. *Lancet Infectious Diseases*. 22:1030–1041.
- 649 Whittles LK, White PJ, Didelot X. 2019. A dynamic power-law sexual network model of  
650 gonorrhoea outbreaks. *PLoS Computational Biology*. 15:e1006748.
- 651 Whittles LK, White PJ, Didelot X. 2020. Assessment of the Potential of Vaccination to Combat  
652 Antibiotic Resistance in Gonorrhoea: A Modeling Analysis to Determine Preferred Product  
653 Characteristics. *Clinical Infectious Diseases*. 71:1912–1919.
- 654 Ypma RJF, Bataille AMA, Stegeman A, Koch G, Wallinga J, van Ballegooijen WM. 2012.  
655 Unravelling transmission trees of infectious diseases by combining genetic and epidemiological  
656 data. *Proceedings of the Royal Society B*. 279:444–450.
- 657 Ypma RJF, van Ballegooijen WM, Wallinga J. 2013. Relating Phylogenetic Trees to  
658 Transmission Trees of Infectious Disease Outbreaks. *Genetics*. 195:1055–1062.