

1 Inference of infectious disease transmission through a relaxed bottleneck using multiple genomes
2 per host

3 Jake Carson^{1,2,3}, Matt Keeling^{1,2,3}, David Wyllie⁴, Paolo Ribeca⁴, Xavier Didelot^{2,3,5}

4 ¹ Mathematics Institute, University of Warwick, Coventry CV4 7AL, United Kingdom

5 ² School of Life Sciences, University of Warwick, Coventry CV4 7AL, United Kingdom

6 ³ Zeeman Institute for Systems Biology and Infectious Disease Epidemiology Research
7 (SBIDER), University of Warwick, Coventry CV4 7AL, United Kingdom

8 ⁴ UK Health Security Agency, London NW9 5EQ, United Kingdom

9 ⁵ Department of Statistics, University of Warwick, Coventry CV4 7AL, United Kingdom

10 Keywords: genomic epidemiology, transmission analysis, infectious disease outbreak, within-
11 host diversity and evolution

12 **ABSTRACT**

13 In recent times, pathogen genome sequencing has become increasingly used to investigate
14 infectious disease outbreaks. When genomic data is sampled densely enough amongst infected
15 individuals, it can help resolve who infected whom. However, transmission analysis cannot
16 rely solely on a phylogeny of the genomes but must account for the within-host evolution of the
17 pathogen, which blurs the relationship between phylogenetic and transmission trees. When only
18 a single genome is sampled for each host, the uncertainty about who infected whom can be quite
19 high. Consequently, transmission analysis based on multiple genomes of the same pathogen per
20 host has a clear potential for delivering more precise results, even though it is more laborious to
21 achieve. Here we present a new methodology that can use any number of genomes sampled from
22 a set of individuals to reconstruct their transmission network. Furthermore, we remove the need
23 for the assumption of a complete transmission bottleneck. We use simulated data to show that
24 our method becomes more accurate as more genomes per host are provided, and that it can infer
25 key infectious disease parameters such as the size of the transmission bottleneck, within-host
26 growth rate, basic reproduction number and sampling fraction. We demonstrate the usefulness
27 of our method in applications to real datasets from an outbreak of *Pseudomonas aeruginosa*
28 amongst cystic fibrosis patients and a nosocomial outbreak of *Klebsiella pneumoniae*.

ACCEPTED MANUSCRIPT

29 INTRODUCTION

30 Pathogen genomic data has transformed our understanding of the epidemiology of infectious
31 diseases, whether they are caused by viruses (Grenfell et al., 2004; Pybus and Rambaut, 2009)
32 or bacteria (Didelot et al., 2012; Gardy and Loman, 2018). Most applications concern large-
33 scale pathogen populations, for example to estimate their demographic history (Pybus et al.,
34 2001; Ho and Shapiro, 2011) or the way that their ancestry relates to features of geography
35 (Lemey et al., 2009; De Maio et al., 2015), epidemiology (Volz et al., 2013; Rasmussen et al.,
36 2014) or host population (Mather et al., 2013; Dearlove et al., 2016). Genomic data can however
37 also be useful to perform much finer inference, down to the level of transmission analysis which
38 attempts to reconstruct who infected whom within an outbreak (Cottam et al., 2008; Jombart
39 et al., 2011). Phylogenetic methods have a long successful history and can reconstruct the
40 genealogy of a set of genomes given their sequences (Yang and Rannala, 2012; Kapli et al., 2020).
41 However, a phylogenetic tree is not identical to a transmission tree (Pybus and Rambaut, 2009;
42 Jombart et al., 2011; Romero-Severson et al., 2014). In particular, the nodes in a phylogenetic
43 tree do not correspond to transmission events, but rather to lineages diverging during the
44 evolutionary process that takes places within a host (Didelot et al., 2016). Several methods
45 have therefore been developed over the past few years specifically aimed at the reconstruction
46 of a transmission tree (Duault et al., 2022). Examples include SeqTrack (Jombart et al., 2011),
47 outbreaker (Jombart et al., 2014), beastlier (Hall et al., 2015), bitrugs (Worby et al., 2016),
48 SCOTTI (De Maio et al., 2016), phybreak (Klinkenberg et al., 2017), outbreaker2 (Campbell
49 et al., 2018) and TiTUS (Sashittal and El-Kebir, 2020).

50 Here we focus on one such method for transmission analysis called TransPhylo, which is based
51 on colouring the branches of a dated phylogeny to reveal the transmission tree (Didelot et al.,
52 2014). There are many software tools that can be used to construct such a dated phylogeny, for
53 example BEAST (Suchard et al., 2018), BEAST2 (Bouckaert et al., 2019), BactDating (Didelot
54 et al., 2018), treedater (Volz and Frost, 2017) and TreeTime (Sagulenko et al., 2018). An
55 advantage of the TransPhylo colouring approach is that it separates the initial phylogenetic
56 reconstruction from its epidemiological interpretation, which improves computational efficiency
57 and therefore scalability (Didelot and Parkhill, 2022). Furthermore, the original TransPhylo
58 model (Didelot et al., 2014) has been extended to deal with both partially sampled and ongoing
59 outbreaks (Didelot et al., 2017). Consequently, TransPhylo is a flexible and versatile software
60 to perform transmission analysis using pathogen genomic data (Didelot et al., 2021).

61 Following infection, many pathogens evolve within hosts on a time scale that is relevant to
62 transmission analysis (Lieberman et al., 2011; Bryant et al., 2013; Biek et al., 2015; Grote
63 and Earl, 2022). Consequently, when information is available about the within-host pathogen
64 diversity, this can help clarify who infected whom (Didelot et al., 2016; Leitner, 2019). This
65 information can come in two forms: either heterogeneities in the genomic sequencing of a
66 single clinical sample, or genomic sequencing of multiple separate clinical samples. Genetic
67 heterogeneities within a sample are relatively easy to survey, and a few methods have been
68 developed recently with the specific aim of exploiting this type of data to help infer transmission
69 (De Maio et al., 2018; Wymant et al., 2018; Torres Ortiz et al., 2023). However this approach
70 is based on the analysis of short sequencing reads individually which can be difficult and error-
71 prone; additionally the clinical sample may not represent the full within-host diversity of the
72 pathogen when it was collected, and it does not contain any information about evolution or
73 changes of diversity over time in the within-host pathogen population. The alternative approach

74 of sequencing several clinical samples can provide a more thorough and reliable overview of the
75 within-host diversity and evolution, especially if the samples are taken from multiple body sites
76 and/or at different points in time. Examples of such studies have been carried on infection
77 with *Staphylococcus aureus* (Young et al., 2012), *Helicobacter pylori* (Didelot et al., 2013) or
78 *Streptococcus pneumoniae* (Tonkin-Hill et al., 2022). Existing methods that can incorporate
79 such data include beastlier (Hall et al., 2015), bitrugs (Worby et al., 2016), SCOTTI (De Maio
80 et al., 2016), phyloscanner (Hall et al., 2019) and TiTUS (Sashittal and El-Kebir, 2020).

81 In principle, integrating multiple genomes into a joint model of phylogenetic and transmission
82 trees, such as TransPhylo, is possible by having as many leaves in the phylogenetic tree as
83 there are samples (Didelot et al., 2016; Leitner, 2019). However, this poses a significant number
84 of theoretical challenges to overcome, which is why TransPhylo was not previously able to
85 use more than one genome per host (Didelot et al., 2017; Xu et al., 2020). Furthermore,
86 TransPhylo previously assumed a complete transmission bottleneck to simplify the relationship
87 between transmission and phylogenetic trees (Didelot et al., 2014), but this assumption has
88 been disproved in some pathogens. Here we present a solution to these issues, which leads us to
89 formulate an extended version of the TransPhylo model, inference methodology and software,
90 so that any number of genomes per host can be used as input of a transmission analysis that
91 does not assume a complete transmission bottleneck.

92 NEW APPROACHES

93 We extend the latest TransPhylo framework (Didelot et al., 2017) to perform inference of
94 infectious disease transmission through a relaxed bottleneck using multiple genomes per host,
95 which may be sampled contemporaneously or longitudinally, or in any combination of both.
96 The model in TransPhylo has three basic ingredients which we detail below, before explaining
97 the changes needed to deal with multiple samples per host. Firstly, a coalescent model with
98 constant population size and temporally offset leaves (Drummond et al., 2002) to represent the
99 within-host evolution. Secondly, a branching process transmission model in which individuals
100 are sampled either once or not at all, so that unsampled individuals can be accounted for in the
101 transmission chains between sampled individuals. Thirdly, a complete transmission bottleneck
102 meaning that only a single lineage is ever transmitted between hosts. In other words the within-
103 host coalescent process is bounded so that the most recent common ancestor within a host occurs
104 after the date of infection (Carson et al., 2022).

105 The full bottleneck assumption can be problematic in settings where hosts are repeatedly
106 sampled, as the resulting phylogenetic trees may have no compatible transmission trees
107 (Romero-Severson et al., 2014, 2016). Therefore we remove this complete bottleneck assumption,
108 so that the phylogenetic trees are much more likely to have compatible transmission trees.
109 Removing this assumption was needed to allow for multiple samples per host, but it is also
110 important to note that a number of studies have found that the transmission bottleneck is only
111 partial for many pathogens including HIV (Boeras et al., 2011), FMDV (Cortey et al., 2019),
112 influenza (Ghafari et al., 2020) and *Staphylococcus aureus* (Hall et al., 2019). Relaxing the
113 transmission bottleneck assumption therefore leads to a more generally applicable model, in
114 which it is possible to additionally estimate the scale of the transmission bottleneck.

115 We also relax the assumption of a constant within-host population size by allowing linear growth,
 116 following previous work on HIV (Romero-Severson et al., 2014, 2016; Leitner, 2019). This linear
 117 growth model is a generalisation of the constant population size model which can be obtained if
 118 the linear growth rate parameter is set to zero. It is also a generalisation of a linear growth with
 119 complete transmission bottleneck model (Klinkenberg et al., 2017) since this can be obtained
 120 if the linear intersect is zero at the date of infection. The linear growth model therefore has
 121 several advantages, on top of being simple and statistically tractable, but other options such as
 122 an exponential or logistic growth model could also be used as will be discussed later.

123 Finally, in the transmission model we add the possibility that hosts are sampled multiple times,
 124 while also retaining the possibility that some hosts are sampled only once or not at all. We make
 125 the specific choice that the transmission model up to the first sample for each host is exactly the
 126 same as previously formulated (Didelot et al., 2017). The times of any further sampling depend
 127 only on the first observation times, and not the infection times. Since the infection times and
 128 secondary observation times are conditionally independent given the primary observation times,
 129 we can infer the infection times without the need to formally define this aspect of the model. In
 130 the Methods section we present a full mathematical description of this new extended model, and
 131 show how Bayesian inference can be performed using a Markov Chain Monte-Carlo (MCMC)
 132 scheme with reversible-jumps (Green, 1995) to accommodate the non-constant dimension of the
 133 parameter space.

134 RESULTS

135 Exemplary analysis of a single simulation

136 We simulate an outbreak with 100 observed hosts, each with five observations. The observation
 137 cut-off time T is determined by the simulation in order to return the correct number of observed
 138 hosts. The generation time and primary observation time are both Gamma distributed (see
 139 section “Epidemiological model” in the Materials and Methods) with shape and scale parameters
 140 equal to 2 and 1, respectively. Secondary observations are placed at intervals of 0.25 years
 141 following the primary observation. For the transmission model, the offspring distribution is
 142 negative binomial with mean equal to the basic reproduction number $R = 2$, and the sampling
 143 proportion is $\pi = 0.8$. The within-host pathogen population size is $\kappa + \lambda\tau$ at time τ after
 144 infection, with $\kappa = 0.1$ and $\lambda = 0.2$. The resulting simulation contains 124 hosts, four of which
 145 are infected with two lineages at the time of infection, one with three lineages, and the remaining
 146 119 with a single lineage.

147 We investigate the ability of our methodology to recover the model parameters used in the
 148 simulation, and to recover transmission links between individuals. We also investigate what
 149 benefits are obtained by including multiple observations per host. To this end we construct
 150 additional phylogenetic trees by pruning the last observation for each host. Through repetition
 151 we obtain phylogenetic trees with four, three, two and one observations per host under the same
 152 transmission network. By comparing inference outcomes from these five trees we can establish
 153 the extent to which estimates are improved through the inclusion of secondary observations.

154 We perform 12,000 MCMC iterations for each phylogenetic tree, using the first 2,000 as a burn-

	Observations per host				
	1	2	3	4	5
π	0.85 [0.62, 0.99]	0.83 [0.62, 0.99]	0.85 [0.65, 0.99]	0.83 [0.63, 0.99]	0.84 [0.64, 0.99]
R	2.32 [1.84, 2.83]	2.32 [1.84, 2.86]	2.27 [1.78, 2.80]	2.25 [1.78, 2.77]	2.25 [1.79, 2.78]
κ	0.18 [0.01, 0.38]	0.15 [0.05, 0.29]	0.10 [0.03, 0.19]	0.10 [0.03, 0.17]	0.11 [0.05, 0.17]
λ	0.19 [0.01, 0.58]	0.18 [0.04, 0.30]	0.23 [0.14, 0.33]	0.20 [0.14, 0.27]	0.21 [0.15, 0.27]

Table 1: Posterior estimates of the simulation study given as the posterior mean and 95% credible interval. The model parameter is given in the left column, and the remaining columns indicate the number of observations per observed host. The values used in the simulation are $\pi = 0.8$, $R = 2$, $\kappa = 0.1$ and $\lambda = 0.2$.

155 in. The prior distribution for π is uniform between 0 and 1, and the prior distributions for R , κ
156 and λ are exponential with mean 1. The posterior means and 95% credible intervals are shown
157 in the Table 1. These results demonstrate that we are able to recover the model parameters used
158 in the simulation, even with no secondary observations. Comparing posterior estimates across
159 the different trees indicates that our estimates of the transmission model parameters R and π
160 are not considerably improved by the number of secondary observations. This makes sense, as
161 most of the relevant information for these parameters is contained in the primary observation.
162 However, the credible intervals for the coalescent model parameters κ and λ narrow as more
163 secondary observations are added. Secondary observations provide considerable information
164 about the within-host genomic diversity of infected hosts, leading to more precise estimates.

165 In order to evaluate our ability to reconstruct transmission links we look at transmissions
166 between observed hosts. Out of the 100 observed hosts, 67 are infected by another sampled
167 individual. From our estimated transmission trees we consider both directional transmission
168 links, where we must correctly establish the infector and infected host, and bidirectional
169 transmission links, where a transmission link is established but the roles of infector and
170 infected may swap. We define 0.5 as the posterior probability threshold for a transmission
171 being identified, and define the sensitivity as the proportion of correctly identified transmission
172 links (true positive rate). For the phylogenetic tree with one observation per host we obtain a
173 sensitivity of 0.51 for bidirectional transmission links, and 0.28 for directional transmission links
174 (Figure S1). For the phylogenetic tree with five observations per host the sensitivity increases
175 to 0.64 for bidirectional transmission links, and 0.55 for directional transmission links (Figure
176 1). The specificity (true negative rate) is greater than 0.996 in all cases. The full distributions
177 of posterior probability estimates in each setting are shown in Figure 2. Increasing the number
178 of secondary observations allows us to better reconstruct transmission links, and crucially, to
179 better distinguish the direction of transmission.

180 The within-host population model plays a key role in our ability to establish transmission
181 links. If the transmission of multiple lineages is more common, the posterior probabilities of
182 transmission links will tend to be lower. For example, repeating the simulation process above
183 with a full bottleneck (fixing $\kappa = 0$) results in a bidirectional (directional) sensitivity of 0.57
184 (0.43) with one observation per host, and 0.75 (0.63) with five observations per host, all higher
185 than in the previous results with a partial bottleneck. On the other hand, increasing to $\kappa = 0.4$
186 leads to a bidirectional (directional) sensitivity of 0.34 (0.25) with one observation per host,
187 and 0.54 (0.39) with five observations per host, all lower than the example with $\kappa = 0.1$.

188 When only a single genome per host is used, we are able to run the original TransPhylo algorithm

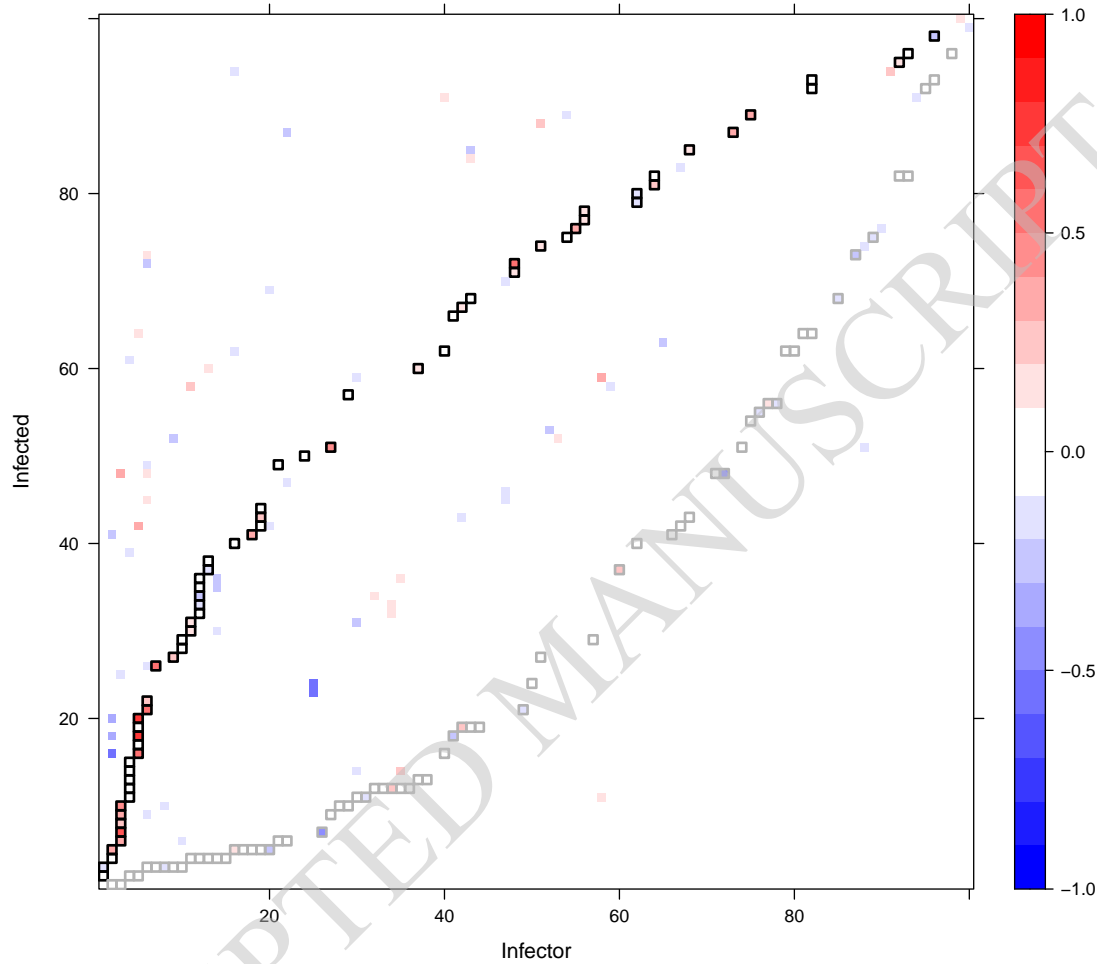


Figure 1: Difference in posterior probability estimates of transmission between a dataset with one observation per host and a dataset with five observations per host. The underlying transmission network remains the same; it is defined by the black squares, which show the true transmissions in the simulated dataset. The gray squares show the reverse relationship, switching the true infector and infected hosts. Black squares containing red demonstrate higher posterior probabilities being assigned to the true transmission links as a result of including more observations. Elsewhere, blue indicates lower posterior probabilities being assigned to incorrect transmission links.

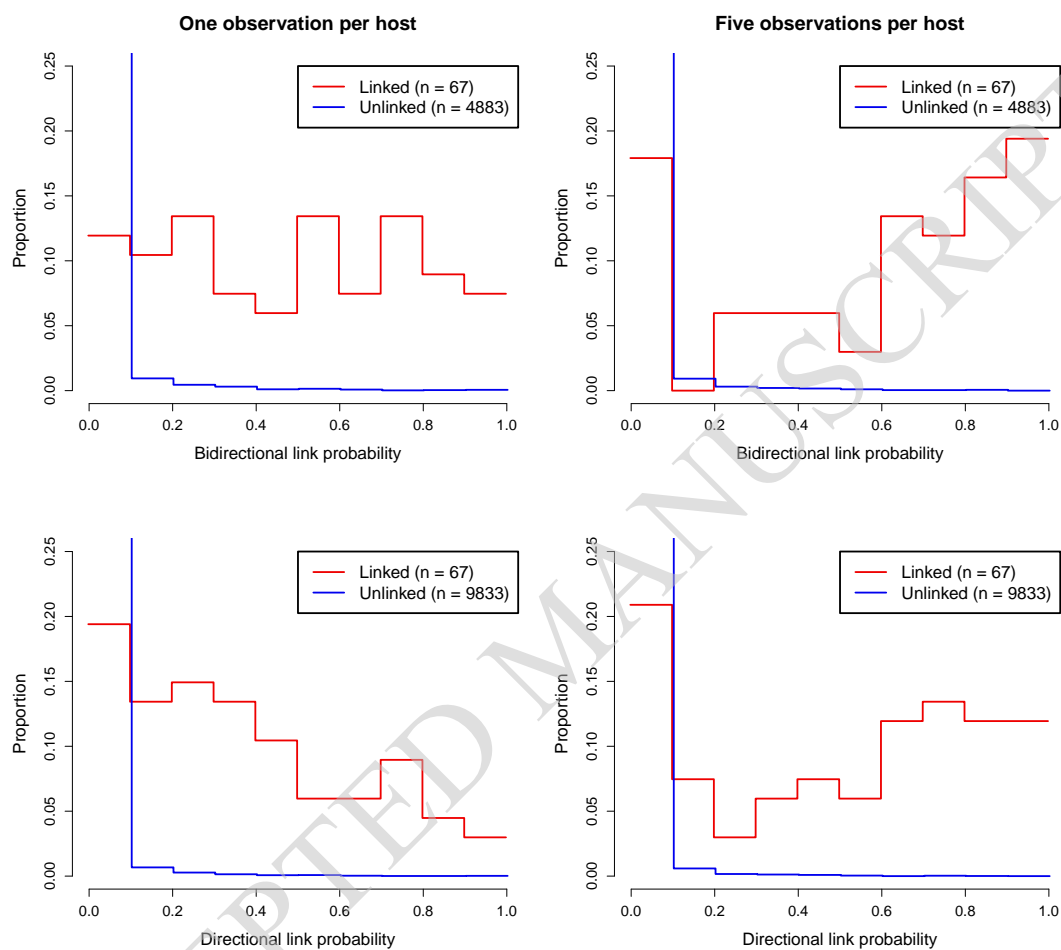


Figure 2: Distribution of posterior link probabilities inferred in the simulation studies with one (left) and five (right) observations per host. The top plots show bidirectional link probabilities in which the roles of infector and infected host may switch, the bottom plots show the directional link probabilities in which the infector and infected host must be correctly inferred. The red lines relate to pairs of individuals for which a transmission link exists, and the blue lines relate to pairs of individuals that are not linked.

189 (Didelot et al., 2017) for comparison. The estimate of π is 0.93 with credible interval [0.76, 1.00],
 190 and the estimate of R is 2.38 with credible interval [1.88, 2.95], which are similar to the estimates
 191 obtained previously with one observation per host (Table 1). The probabilities for who infected
 192 whom are shown in Figure S2. The bidirectional (directional) sensitivity is 0.61 (0.37), as
 193 illustrated in Figure S3. Since a small value of $\kappa = 0.1$ is used in the simulation, the strict
 194 bottleneck assumption in TransPhylo is advantageous here, whereas using a relaxed bottleneck
 195 leads to additional uncertainty on who infected whom. TransPhylo would perform comparatively
 196 less well if the true bottleneck was more relaxed.

197 **Benchmarking using multiple simulations**

198 We now repeat this process, again using a simulated dataset with 100 hosts and five observations
 199 per host; but performing the inference on simulations generated from a range of key parameters
 200 (π , R , λ , and κ), totalling 43 datasets. As previously, both the generation time distribution
 201 and primary observation time distribution follow a Gamma distribution with shape parameter
 202 2 and scale parameter 1, and secondary observations occur 0.25 years later than the previous
 203 sample.

204 For the MCMC chains we obtain 12,000 samples, and discard the first 2,000 as a burn-in.
 205 Figure 3 shows the posterior parameter estimates. The vertical lines show central 95% credible
 206 intervals for each parameter, and the posterior mean is shown with a solid circle. The horizontal
 207 and diagonal lines indicate the true parameter values used to generate the data. These results
 208 demonstrate strong performance of the algorithm across very different simulation settings.

209 The linear growth assumption of the within-host population size model is unlikely to resemble a
 210 real-world population, and so we also test for robustness to the mis-specification of the within-
 211 host population model. We repeat the inference, but fix the within-host population growth
 212 rate λ at either half or double the true value. The posterior estimates are shown in Figure
 213 S4. Most notably, the mis-specification biases our estimates of the initial pathogen population
 214 size κ . There is a strong negative correlation between λ and κ , so that when λ is set lower
 215 (higher) κ is overestimated (underestimated). There are smaller changes in the transmission
 216 model parameters, with a lower λ resulting in higher estimates of π and lower estimates of r , but
 217 the true values for these parameters usually remain within the 95% credible intervals. These
 218 results suggest that estimates of the transmission model parameters are reasonably robust to
 219 the mis-specification of the within-host population model. However, caution is warranted when
 220 interpreting the estimates of the within-host model parameters. We can reasonably conclude,
 221 for instance, that different estimates of the initial population size κ may be obtained under
 222 different growth models.

223 **Application to *Pseudomonas aeruginosa* transmission between cystic fibrosis** 224 **patients**

225 We reanalysed previously published genomic data from Danish cystic fibrosis (CF) patients
 226 infected with *Pseudomonas aeruginosa* (Marvig et al., 2013). This dataset included 42 genomes
 227 from 14 patients, sampled over almost 40 years between 1972 and 2008, after exclusion of

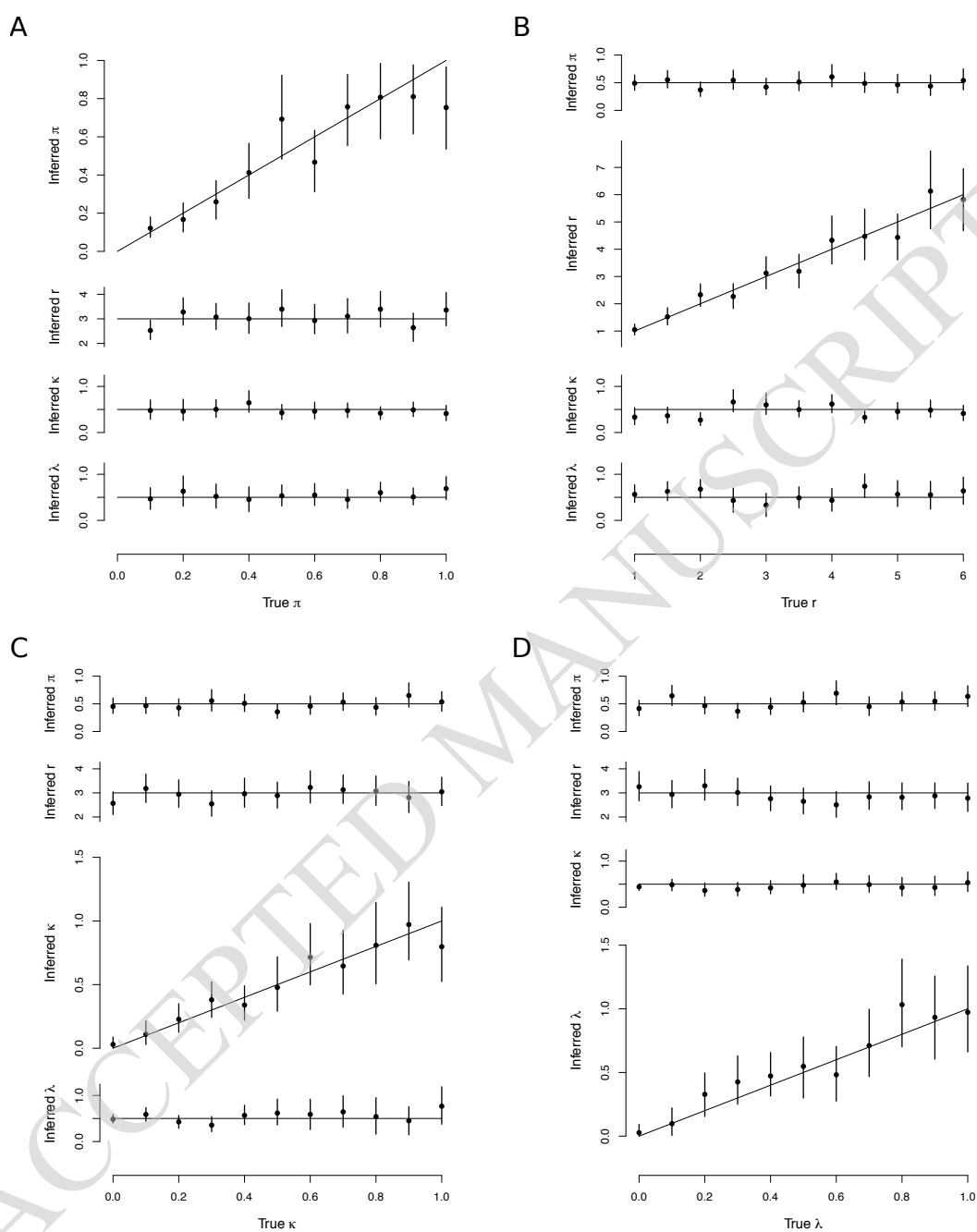


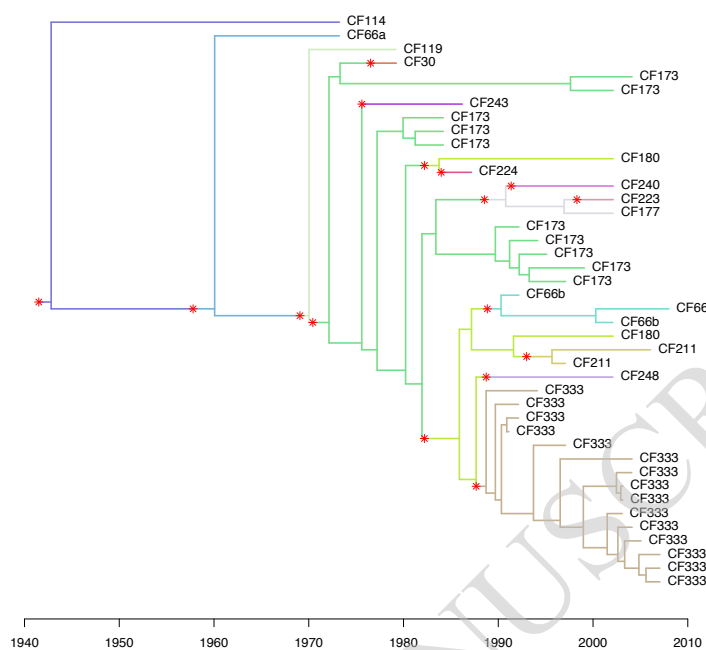
Figure 3: Varying the four key simulation parameters. Vertical bars show 95% central credible intervals, while solid circles show posterior means. Horizontal or diagonal lines show true values for simulations. (A) Varying π . (B) Varying R . (C) Varying κ . (D) Varying λ .

228 hypermutator and recombinant isolates (Marvig et al., 2013). Previous studies explored within-
 229 host evolutionary dynamics (Yang et al., 2011), variations in gene content (Rau et al., 2012)
 230 and comparative adaptation in CF human hosts (Marvig et al., 2013). The hosts are designated
 231 CFXXX as in these previous studies. We use as our starting point the dated phylogeny
 232 previously computed (Marvig et al., 2013) using BEAST (Suchard et al., 2018) and shown
 233 in Figure S5. It was previously noted (Yang et al., 2011) that one of the individual (CF66) had
 234 been infected twice in the 1970s and the 1990s, and so we modelled this as two separate hosts
 235 (labeled CF66a and CF66b). Infection with *P. aeruginosa* can be stable over long periods of
 236 time in CF patients (Rossi et al., 2021) and indeed some of the patients had been sampled, and
 237 found positive, over a period of more than 20 years (Marvig et al., 2013). We therefore set the
 238 generation time distribution to be Gamma with shape 2 and scale 5, resulting in a mean of 10
 239 years, standard deviation of 7 years, and 95% range of 1.2 to 27.9 years. The last samples were
 240 from 2008 and the exact end of the sampling period was unclear from previous publications but
 241 we set it to the end of 2009.

242 We performed four separate runs of 100,000 iterations, which took approximately 3 hours on a
 243 standard laptop computer. For each of the four parameters π , R , κ and λ we checked that the
 244 effective sample size in each run was over 1,000 and the multivariate Gelman-Rubin statistic
 245 comparing runs was less than 1.1 (Brooks and Gelman, 1998). Figure 4A shows the dated
 246 tree, coloured by host according to the MCMC iteration with the highest posterior probability.
 247 Changes in colours along the branches of the tree correspond to transmission events and are
 248 highlighted with red stars. Note that there are two simultaneous stars leading to the two
 249 genomes from patient CF180. These both correspond to infection from CF173, with the two
 250 lineages being transmitted through the relaxed transmission bottleneck. Figure 4A is useful
 251 to illustrate the colouring process which relates the phylogenetic tree to the transmission tree.
 252 However, this only represents a single transmission configuration explored by the MCMC, and
 253 other iterations of the MCMC would look different, maybe with some of the same transmission
 254 events and others being different. It is therefore important to consider the probability of the
 255 transmission events. Figure 4B shows the matrix of probabilities of infection from each host to
 256 another, computed as the frequency of each transmission event across all MCMC iterations.

257 Figure S6 shows the trace and density of the parameters estimated in a single MCMC run.
 258 The sampling proportion was estimated to be $\pi = 0.65$, with a wide 95% credible interval
 259 $[0.30 - 0.96]$. The reproduction number was $R = 1.20$ $[0.58 - 1.99]$; as the credible interval
 260 includes one, it is not clear if the outbreak has the potential to cause a self-sustained epidemic.
 261 The within-host linear growth rate was $\lambda = 0.56$ $[0.16 - 1.09]$ per year, which is lower than
 262 the prior exponential with mean one. On the other hand, the within-host starting population
 263 size was $\kappa = 2.16$ $[0.41 - 5.05]$ which is higher than the prior exponential with mean one. This
 264 suggest that the bottleneck was not complete, and indeed attempting to fit the model with $\kappa = 0$
 265 is impossible as it leads to a likelihood of zero. This is caused by the two samples from CF180
 266 and the ten samples from CF173 being “inconsistent” as previously designated for samples from
 267 two hosts that cannot be explained by transmission of a single lineage (Romero-Severson et al.,
 268 2014, 2016). The individual CF173 was found to have infected at least three other hosts (CF30,
 269 CF224 and CF243) with probability higher than 50% (Figure 4B). These transmission events
 270 and their directionality are made clear by the paraphyletic relationship of the ten samples from
 271 CF173 as shown in Figure 4A (Leitner, 2019). In contrast, the 15 samples from CF333 formed
 272 a single monophyletic clade (Figure 4A) so that they are unlikely to have infected many others
 273 except maybe CF248 (Figure 4B).

A



B

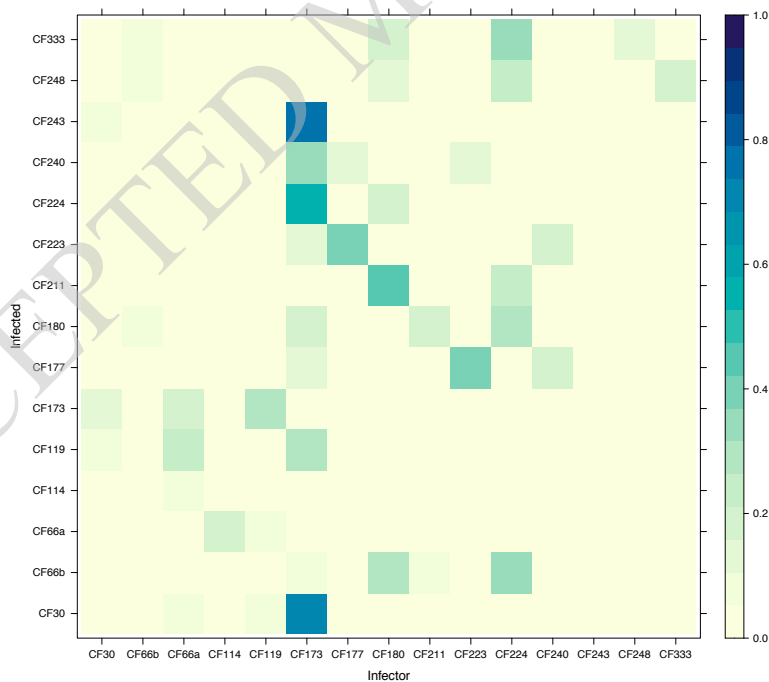


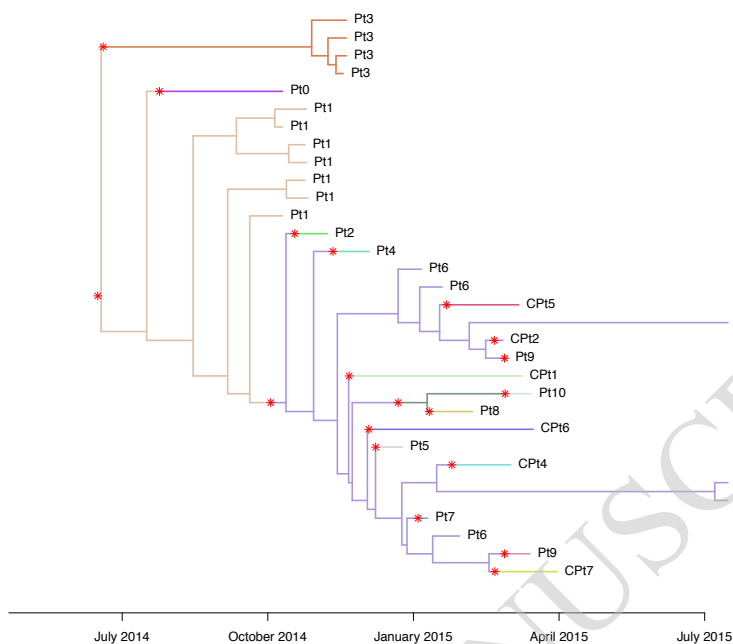
Figure 4: Transmission analysis of *P. aeruginosa*. (A) Dated phylogeny coloured by host according to the iteration with highest posterior probability. (B) Matrix of transmission probabilities from each host (row) to any other (column).

275 An outbreak of carbapenem-resistant *Klebsiella pneumoniae* expressing the *bla*_{OXA-232} gene
276 was identified over the course of 40 weeks at a single healthcare institution in California (Yang
277 et al., 2017). A total of 17 infected patients were identified, from which 32 isolates were taken
278 between 12th October 2014 and 17th July 2015. Case finding was performed using all samples
279 in the 2014 and 2015 calendar years (Yang et al., 2017) and so we set the date for the end
280 of the sampling period to the end of 2015. Whole-genome sequencing was applied to these
281 *K. pneumoniae* isolates and a dated phylogeny was computed previously (Yang et al., 2017)
282 using BEAST (Suchard et al., 2018) which is shown in Figure S7. The hosts are labeled either
283 PtXXX if they were symptomatic or CPtXXX if they were colonized, as in the previous study
284 (Yang et al., 2017). We set the generation time distribution to be exponential with mean 0.5
285 year, following a previous study of another *K. pneumoniae* hospital outbreak (van Dorp et al.,
286 2019). This diffuse distribution is well suited to capture transmission via hospital equipment
287 contamination as was previously suggested (Yang et al., 2017). We used the same number of
288 MCMC runs, length of runs, and convergence diagnostics as in the previous application.

289 Figure S8 shows the trace and density of the parameters estimated in a single MCMC run.
290 The sampling proportion was estimated to be high, with $\pi = 0.88$ [0.60 – 0.99], suggesting that
291 there were only few missing transmission links between the 17 sampled patients. The basic
292 reproduction number was $R = 0.97$ [0.37 – 1.74], with the credible interval including the value
293 of one needed for an outbreak to spread beyond a few cases. The within-host linear growth rate
294 was $\lambda = 0.49$ [0.03 – 1.28] per year and the within-host population size at time of infection was
295 $\kappa = 0.066$ [0.009 – 0.158]. This is lower than the prior exponential with mean one and suggests
296 that the transmission bottleneck was almost complete during this small outbreak. However, the
297 transmission bottleneck was not absolutely complete, as indicated by the fact that fitting our
298 model with $\kappa = 0$ would result in a likelihood equal to zero. This is because the six samples
299 from Pt6 and the two samples from Pt9 are inconsistent, as can be seen in the dated phylogeny
300 on Figure S7.

301 Figure 5A shows the dated tree coloured by host according to the MCMC iteration with highest
302 posterior probability, while Figure 5B shows the posterior probabilities of infection from any
303 host to any other. For example, a high probability of transmission was found from Pt8 to Pt10,
304 which is consistent with the fact that these two patients were staying in neighboring rooms for
305 two weeks (Yang et al., 2017). Strikingly, according to our analysis patient Pt6 had a greater
306 than 50% posterior probability of having infected seven other patients (CPt2, CPt4, CPt5,
307 CPt6, Pt5, Pt7 and Pt9). There were six genomes isolated from Pt6, with dates ranging from
308 7th January 2015 to 17th July 2015 which is more than half of the overall sampling period. The
309 specimen types for these isolates were quite diverse: three from blood, one rectal and two from
310 bile (Yang et al., 2017), suggesting that the patient was infected long enough for the pathogen to
311 spread throughout their body. While other patients in the study do present a similar number of
312 samples, a comparable variety of originating tissues, and a similarly long infection duration —
313 for instance patient Pt1, with seven genomes from respiratory, abdominal and blood specimen
314 over a period of several months — that does not translate in a similar amount of infection
315 events estimated by our method. In fact, the genetic diversity of isolates from Pt6 appears to
316 be very high (Figure 5A), thus backing our inference that Pt6 is a superspreading individual
317 (Lloyd-Smith et al., 2005). This could not have been detected without the use of multiple
318 genomes.

A



B

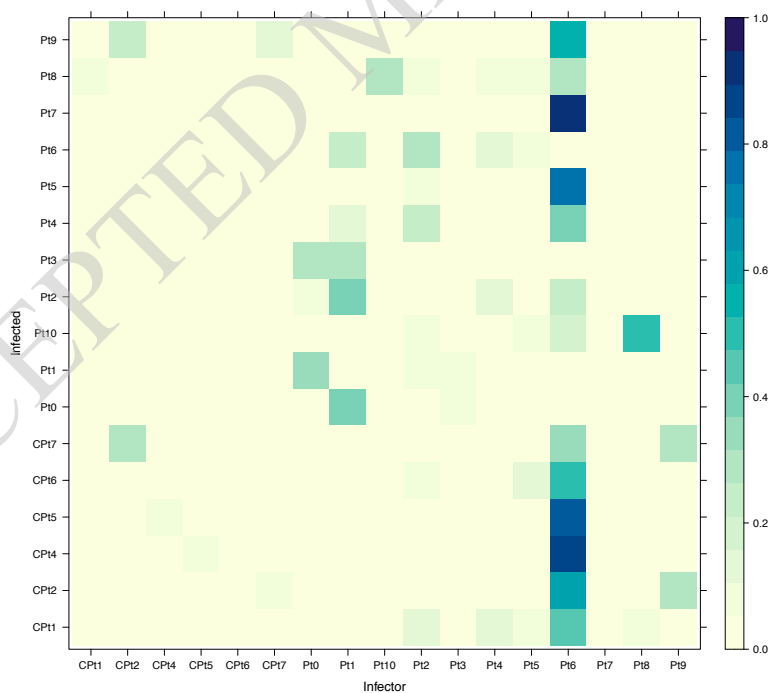


Figure 5: Transmission analysis of *K. pneumoniae*. (A) Dated phylogeny coloured by host according to the iteration with highest posterior probability. (B) Matrix of transmission probabilities from each host (row) to any other (column).

319 DISCUSSION

320 We have described new methodology for inferring who infected whom from a dated phylogenetic
321 tree in which hosts have potentially been sampled multiple times. A key change compared to
322 previous work (Didelot et al., 2014, 2017) is the removal of the full transmission bottleneck,
323 meaning that hosts may be infected with multiple lineages from the transmission donor.
324 Without this change many phylogenetic trees with multiple samples per host would not support
325 compatible transmission trees (Romero-Severson et al., 2014, 2016; Leitner, 2019). Indeed the
326 two real datasets we analysed, corresponding to outbreaks of *Pseudomonas aeruginosa* and
327 *Klebsiella pneumoniae*, could not be explained without relaxing the transmission bottleneck.
328 Most previous transmission analysis methods could not accommodate more than a single genome
329 per host, so that leaves would need to be pruned from the phylogenetic tree in order to undertake
330 transmission inference (Xu et al., 2020), leading to less informative outcomes. Under our new
331 methodology we are able to incorporate multiple samples per host, resulting in the stronger
332 identification of transmission links and their direction, as was showed when analysing simulated
333 datasets.

334 We build upon previous work (Didelot et al., 2014, 2016) that performs transmission analysis
335 by colouring the branches of a pre-established dated phylogeny. This allows us to model
336 the relationship between transmission tree and phylogeny through an explicit within-host
337 evolutionary model, to develop an explicit transmission model in which sampled and unsampled
338 individuals are featured, and to achieve better scalability by separating phylogenetic inference
339 from its epidemiological interpretation. On the other hand, relying on a fixed dated tree could be
340 problematic as this does not account for the uncertainty in the phylogeny or the dates of common
341 ancestors. When this uncertainty is captured using a Bayesian phylogenetic method (Suchard
342 et al., 2018; Didelot et al., 2018; Bouckaert et al., 2019), this effect can be tested by applying
343 analysis to multiple samples instead of a single fixed tree (Nylander et al., 2008). However, this
344 was found in practice to make little difference to the inferred transmission probabilities and
345 parameters (Didelot and Parkhill, 2022).

346 Our method implements a general pathogen population growth model rather than using the
347 constant bounded coalescent model, in which the population size is constant and the most
348 recent common ancestor is forced to occur after the infection time (Carson et al., 2022). By
349 removing this restriction, we were able to model transmission through a relaxed bottleneck. The
350 main restriction on the choice of model is that we must be able to calculate the likelihood of the
351 phylogenetic tree, which in turn means that the coalescence rate must be integrable. However,
352 this is not a strong requirement, as many widely used models satisfy it — among them the
353 exponential growth model, the logistic growth model, or any piecewise models with separate
354 growth and decay phases. For the work presented here we used a linear growth model, which
355 has been used before in HIV work (Romero-Severson et al., 2014, 2016; Leitner, 2019), but for
356 most other pathogens there is little information about which within-host population size model
357 is most realistic (Didelot et al., 2016). We demonstrated that using phylogenetic trees with
358 multiple samples per host improves the estimation of the population model parameters. With
359 sufficient samples per host it should be possible to determine which within-host population size
360 models are more strongly supported by the data, for example and comparing the evidence of
361 each model (Friel and Wyse, 2012).

362 Our methodology maintains some of the assumptions from previous work (Didelot et al., 2017),

363 for example the sampling proportion and reproduction number are assumed to remain constant
 364 through time. In many settings, users would have knowledge about whether and how the
 365 sampling proportion varied over time, for example by looking at the number cases for which
 366 genomic sequences are available divided by the number of confirmed cases (Jelley et al., 2022).
 367 This information could be integrated relatively easily into an analysis, by having users supply
 368 a function $\pi(t)$ instead of the constant π . On the other hand, it would often be interesting
 369 to infer variations in the reproduction number $R(t)$, since this would provide an additional
 370 genomic-based estimate compared to existing methods based on incidence data (Wallinga and
 371 Teunis, 2004; Cori et al., 2013). A simple approach would be to use a step-wise constant
 372 function. The dates of these steps may be fixed based on real-world policy changes, such as
 373 intensifying monitoring in response to an outbreak, or potentially inferred via change point
 374 detection (Tartakovsky and Moustakides, 2010).

375 In conclusion, we presented a new Bayesian inference method for the reconstruction of
 376 transmission trees from dated phylogenetic trees in which hosts are sampled multiple times.
 377 This method is implemented in a R package that extends TransPhylo and is available at
 378 <https://github.com/DrJCarson/TransPhyloMulti>. When applied to multiple sampled genomes
 379 from several infected individuals, our method has the potential to improve our understanding of
 380 both the within-host and between-host dynamics of many pathogens causing infectious disease.

381 MATERIALS AND METHODS

382 Notation

383 Let us denote \mathcal{P} as the dated phylogenetic tree, \mathcal{T} as a transmission tree, θ_P as the coalescent
 384 model parameters, and θ_T as the transmission model parameters. We want to sample from the
 385 posterior distribution

$$p(\theta_P, \theta_T, \mathcal{T} \mid \mathcal{P}) \propto p(\mathcal{P} \mid \mathcal{T}, \theta_P) p(\mathcal{T} \mid \theta_T) p(\theta_T) p(\theta_P), \quad (1)$$

386 where the term $p(\mathcal{P} \mid \mathcal{T}, \theta_P)$ is the likelihood of the coalescent model conditional on a given
 387 transmission tree, the term $p(\mathcal{T} \mid \theta_T)$ is the likelihood of the transmission model, and the terms
 388 $p(\theta_P)$ and $p(\theta_T)$ are prior distributions.

389 We parameterise the transmission tree \mathcal{T} as follows. Let x be a vector of infection times such
 390 that element x^j gives the infection time of host j . Likewise let A be a vector of infectors, so
 391 that if $A^j = i$ then host j was infected by host i . We indicate the root host by setting $A^j = 0$.
 392 Primary observation times are denoted by vector y , with the corresponding host denoted by
 393 vector H_y . Secondary observation times are denoted by vector z , with host H_z .

394 For the phylogenetic tree \mathcal{P} we need to consider the leaf and coalescent times. The leaves
 395 correspond to observations under the transmission tree. We denote the vector of leaf times s
 396 and corresponding hosts H_s , noting that $s = (y, z)$ and that $H_s = (H_y, H_z)$. We indicate the
 397 parent node of each sample using vector C_s . The coalescent node times are denoted by vector
 398 u , and their parent nodes C_u . We again denote the root node with $C_u^j = 0$.

399 Figure 6A demonstrates a transmission tree with

$$x = \begin{pmatrix} 0.0 \\ 0.8 \\ 1.5 \\ 2.6 \\ 2.5 \\ 0.6 \end{pmatrix}, \quad A = \begin{pmatrix} 0 \\ 1 \\ 6 \\ 3 \\ 3 \\ 1 \end{pmatrix}.$$

400 That is, host 1 infects hosts 2 and 6, host 6 infects host 3, and host 3 infects hosts 4 and 5. In
401 addition we have primary and secondary observations (not shown), for example

$$y = \begin{pmatrix} 1.9 \\ 2.6 \\ 3.2 \\ 3.1 \\ 3.0 \end{pmatrix}, \quad H_y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, \quad z = \begin{pmatrix} 3.5 \\ 3.4 \end{pmatrix}, \quad H_z = \begin{pmatrix} 3 \\ 4 \end{pmatrix},$$

402 indicates that hosts 1, 2 and 5 are observed once, hosts 3 and 4 are observed twice, and host 6
403 is unobserved.

404 Figure 6B shows an example phylogenetic tree obtained by combining the primary and secondary
405 observations from the transmission tree. Here

$$s = \begin{pmatrix} 1.9 \\ 2.6 \\ 3.2 \\ 3.1 \\ 3.0 \\ 3.5 \\ 3.4 \end{pmatrix}, \quad u = \begin{pmatrix} 0.2 \\ 0.4 \\ 0.9 \\ 2.3 \\ 2.9 \\ 3.1 \end{pmatrix}, \quad H_s = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 3 \\ 4 \end{pmatrix}, \quad C_s = \begin{pmatrix} 2 \\ 1 \\ 6 \\ 5 \\ 4 \\ 6 \\ 5 \end{pmatrix}, \quad C_u = \begin{pmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 3 \\ 4 \end{pmatrix}.$$

406 We can represent both the transmission and phylogenetic trees as a coloured phylogenetic tree,
407 as shown in Figure 6C. Doing so highlights that each coalescent event is now assigned to a host.

408 Epidemiological model

409 The epidemiological model is a stochastic branching process in which infected individuals
410 transmit to secondary cases (offspring). The number of offspring k is sampled from the offspring
411 distribution $\alpha(k)$, assumed to be a negative binomial distribution with parameters (r, p) , i.e.

$$\alpha(k) = \binom{k+r-1}{k} p^k (1-p)^r. \quad (2)$$

412 The time between the primary and any secondary infection is sampled from the generation time
413 distribution $\gamma(\tau)$, which typically follows a Gamma distribution with known parameters.

414 Under a *finished outbreak* scenario, each host is assumed to be observed with probability π . The
415 time between the host being infected and first being observed is sampled from the observation

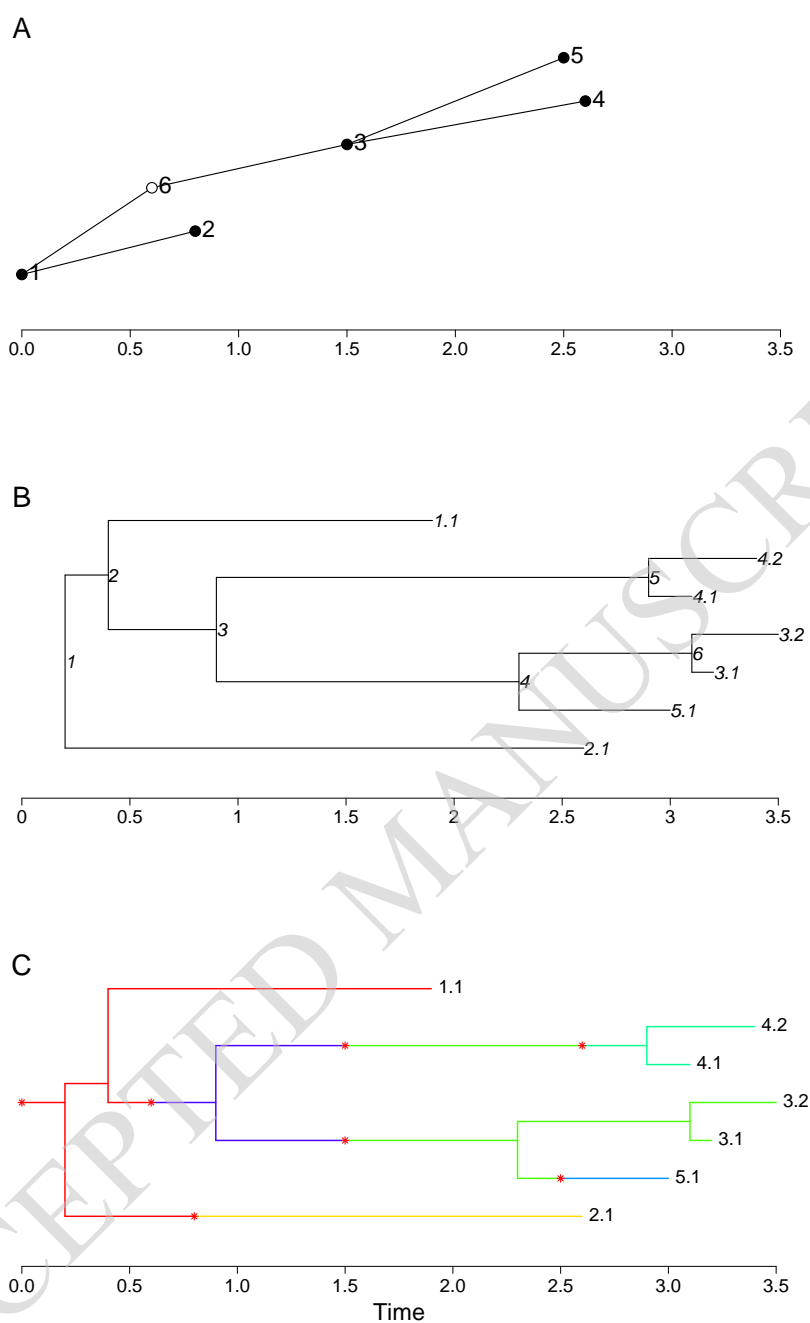


Figure 6: (A) Example transmission tree with six hosts. Points indicate the infected times of each host. Filled circles show observed hosts, and empty circles show unobserved hosts. (B) Example phylogenetic tree with seven leaves from five observed hosts. Leaf labels indicate the host, followed by the sample number for that host. Each coalescence node is given a label. (C) Example coloured phylogenetic host with seven leaves from five observed hosts, and six hosts overall. The branch colour indicates the host, and the asterisks indicate transmissions. Here host 3 is infected with two lineages.

416 time distribution $\sigma(\tau)$. As with the generation time distribution this is typically a Gamma
 417 distribution with known parameters.

418 In some applications observations occur over a restricted time interval, or possibly set of time
 419 intervals. In such applications the probability of a host being observed depends on their
 420 infection time. An example we will look at is the *ongoing outbreak* scenario, in which there
 421 is an observation cut-off time T . In this scenario a host infected at time t is observed with
 422 probability

$$\zeta(t) = \pi \int_0^{T-t} \sigma(\tau) d\tau.$$

423 In other words, we use the same observation distribution as the finished outbreak scenario, but
 424 treat observations later than T as censored.

425 Finally, hosts may be observed multiple times. We assume that any host can only be infected
 426 once, and that any subsequent observations relate to the same infected period. We define $\beta(b)$ as
 427 the distribution for the number of secondary observations $b \geq 0$, and $\rho(\tau_{1:b})$ as the distribution
 428 for the times between the secondary observations and the primary observation assuming that
 429 $b \geq 1$. Note that it is possible for the time between observations to be zero, meaning that
 430 multiple observations occur at the primary observation time.

431 Secondary observations are an additional modelling component to the previous version of
 432 TransPhylo (Didelot et al., 2017). However, by assuming that the secondary observation
 433 times depend only on the primary observation times, we can undertake inference in a similar
 434 manner without formally specifying these distributions. Under our modelling assumptions we
 435 can express the likelihood of the transmission tree as

$$\begin{aligned} p(\mathcal{T} | \theta_T) &= p(x, y, z, A, H_y, H_z | \theta_T) \\ &= p(z, H_z | y, H_y) p(y, H_y | x, A, \theta_T) p(x, A | \theta_T), \end{aligned} \quad (3)$$

436 where x , A , and θ_T are parameters we are trying to estimate, and y , z , H_y , and H_z are fixed
 437 by the dated phylogenetic tree. Within a Metropolis-Hastings algorithm, when we propose new
 438 values x' and A' (giving a new transmission tree \mathcal{T}') or θ'_T , the term $p(z, H_z | y, H_y)$ will cancel
 439 in the likelihood ratio, i.e.

$$\frac{p(\mathcal{T} | \theta_T)}{p(\mathcal{T}' | \theta'_T)} = \frac{p(y, H_y | x, A, \theta_T) p(x, A | \theta_T)}{p(y, H_y | x', A', \theta'_T) p(x', A' | \theta'_T)}. \quad (4)$$

440 Consequently, $p(z, H_z | y, H_y)$ does not need to be explicitly calculated to determine if proposals
 441 are accepted or rejected, and practically can be excluded from the transmission tree likelihood
 442 altogether.

443 Host inclusion and exclusion

444 Our goal is to infer a transmission tree from a dated phylogenetic tree. This can be visualised as
 445 *colouring* the branches of the phylogenetic tree, where each colour represents a distinct host. For
 446 a host to appear on the phylogenetic tree they must either be observed directly or be an ancestor
 447 to a different observed host. We refer to such hosts as *included* hosts. In many applications the
 448 number included hosts is dwarfed by the number of hosts implied by the epidemiological model

449 to not appear on the phylogenetic tree (*excluded* hosts). Examples include when π is small, or
 450 when r is large in an ongoing outbreak scenario. In the latter case, a large number of hosts
 451 will be infected shortly before the observation cut-off time, and so will be excluded with high
 452 probability. For this reason we instead formalise a transmission model for only the included
 453 hosts.

454 Define $\omega(t)$ as the exclusion probability of a host infected at time t . Assuming that T is the
 455 cut-off time for observations $\omega(t) = 1$ for $t \geq T$. We can then define the following recursive
 456 relationships.

457 The exclusion probability of an offspring from a host infected at time t is

$$\bar{\omega}(t) = \int_0^\infty \omega(t + \tau)\gamma(\tau)d\tau. \quad (5)$$

458 The probability that all offspring from an individual infected at time t are excluded is

$$\phi(t) = \sum_{k=0}^{\infty} \alpha(k)\bar{\omega}(t)^k. \quad (6)$$

459 The exclusion probability of an individual infected at time t is

$$\begin{aligned} \omega(t) &= (1 - \zeta(t))\phi(t) \\ &= (1 - \zeta(t)) \sum_{k=0}^{\infty} \alpha(k) \left(\int_0^\infty \omega(t + \tau)\gamma(\tau)d\tau \right)^k. \end{aligned} \quad (7)$$

460 That is, the probability of the host being unobserved and having no included offspring. In the
 461 finished outbreak scenario the recursive relationship is simply

$$\omega_* = (1 - \pi) \sum_{k=0}^{\infty} \alpha(k)\omega_*^k, \quad (8)$$

462 with ω_* being the exclusion probability for every host. Note that these calculations do not
 463 depend on the secondary observation times or their distribution.

464 Numerical approximations

465 The exclusion probabilities are intractable, and so we use numerical approximations. For
 466 example, consider the ongoing outbreak scenario with observation cut-off time T . For $t \geq T$,
 467 $\omega_t = 1$, and so

$$\bar{\omega}(t) = \int_t^T \gamma(\tau - t)\omega(\tau)d\tau + \int_T^\infty \gamma(\tau - t)d\tau. \quad (9)$$

468 The second term can be computed explicitly, and the first term can be approximated using the
 469 trapezoid method:

$$\int_t^T \gamma(\tau - t)\omega(\tau)d\tau \approx \sum_{i=0}^k c_i \gamma((k - i)\Delta t)\omega(t_i)\Delta t, \quad (10)$$

470 where $c_i = 1$ for $0 < i < k$ and $c_i = 0.5$ otherwise, and $t_i = T - i\Delta t$. Assuming $\gamma(0) = 0$:

$$\bar{\omega}(t) \approx F(t) + \sum_{i=0}^{k-1} c_i \gamma((k-i)\Delta t) \omega(t_i) \Delta t. \quad (11)$$

471 where $F(t) = \int_T^\infty \gamma(\tau - t) d\tau$.

472 Using the probability generating function of a negative binomial distribution with parameters
473 r and p , we can evaluate

$$\phi(t) = \left(\frac{p}{1 - (1-p)\bar{\omega}(t)} \right)^r, \quad (12)$$

474 and finally

$$\omega(t) = (1 - \zeta(t))\phi(t). \quad (13)$$

475 Both will be approximate owing to the approximation of $\bar{\omega}(t)$. All three exclusion probabilities
476 are therefore approximated by iterating backwards through time from T in discrete steps of size
477 Δt .

478 **Transmission tree likelihood**

479 We can now define a likelihood for the transmission tree for only included individuals.
480 Throughout we will set T as the cut-off time for observations. Consider first the root host
481 (the first infected individual in our transmission chain) with infection time x^1 , and let $I^1 = 1$
482 denote that the root host is included. The probability that the root host is unobserved (denoted
483 by $S^1 = 0$) given that they are included is

$$\begin{aligned} p(S^1 = 0 \mid I^1 = 1, x^1) &= \frac{p(I^1 = 1 \mid S^1 = 0, x^1)p(S^1 = 0 \mid x^1)}{p(I^1 = 1 \mid x^1)} \\ &= \frac{(1 - \phi(x^1))(1 - \zeta(x^1))}{1 - \omega(x^1)}, \end{aligned} \quad (14)$$

484 and the probability that the root host is observed ($S^1 = 1$) is

$$\begin{aligned} p(S^1 = 1 \mid I^1 = 1, x^1) &= \frac{p(I^1 = 1 \mid S^1 = 1, x^1)p(S^1 = 1 \mid x^1)}{p(I^1 = 1 \mid x^1)} \\ &= \frac{\zeta(x^1)}{1 - \omega(x^1)}. \end{aligned} \quad (15)$$

485 In the event the root host is observed we also need to calculate the density of the primary
486 observation time y^1 ,

$$p(y^1 \mid S^1 = 1, x^1) = \frac{\sigma(y^1 - x^1)}{\int_0^{T-x^1} \sigma(\tau) d\tau}, \quad x^1 < y^1 < T. \quad (16)$$

487 Additionally the full transmission tree likelihood incorporates the density of the secondary
488 observation times. However, when it comes to undertaking inference these terms will cancel
489 out, and so we skip this step.

490 Second, we calculate the probability that the root host has d^1 included offspring. The probability
 491 of a host infected at time t producing d included offspring is

$$\begin{aligned} p(d | t) &= \sum_{k=d}^{\infty} \alpha(k) p(d | k, t) \\ &= \sum_{k=d}^{\infty} \alpha(k) \binom{k}{d} \bar{\omega}(t)^{k-d} (1 - \bar{\omega}(t))^d. \end{aligned} \quad (17)$$

492 We then need to condition on whether or not the root host was sampled. If the root host was
 493 not sampled, they must produce at least one included offspring to be included, and so

$$\begin{aligned} p(d^1 | I^1 = 1, S^1 = 0, x^1) &= \frac{p(I^1 = 1 | d^1, S^1 = 0, x^1) p(d^1 | S^1 = 0, x^1)}{p(I^1 = 1 | S^1 = 0, x^1)} \\ &= \frac{p(d^1 | x^1)}{1 - \phi(x^1)}, \quad d^1 > 0. \end{aligned} \quad (18)$$

494 If the root host was sampled, then it is included for any value of d^1 , and so

$$\begin{aligned} p(d^1 | I^1 = 1, S^1 = 1, x^1) &= \frac{p(I^1 = 1 | d^1, S^1 = 1, x^1) p(d^1 | S^1 = 1, x^1)}{p(I^1 = 1 | S^1 = 1, x^1)} \\ &= p(d^1 | x^1), \quad d^1 \geq 0. \end{aligned} \quad (19)$$

495 In the event $d^1 > 0$, we also calculate the density of the transmission times for any included
 496 offspring. Denoting \mathcal{H}^1 as the offspring labels, $\bar{x}^1 = \{x^j | j \in \mathcal{H}^1\}$ as the set of offspring
 497 infection times, and $\bar{I}^1 = 1$ that the set of offspring are included, the likelihood contribution is

$$\begin{aligned} p(\bar{x}^1 | \bar{I}^1 = 1, x^1) &= d^1! \prod_{j \in \mathcal{H}^1} \frac{p(I^j = 1 | x^j) p(x^j | x^1)}{p(I^j = 1 | x^1)} \\ &= d^1! \prod_{j \in \mathcal{H}^1} \frac{(1 - \omega(x^j)) \gamma(x^j - x^1)}{1 - \bar{\omega}(x^1)}. \end{aligned} \quad (20)$$

498 The $d^1!$ term arises from the fact that the infection times are labelled according to host, and the
 499 host labels are arbitrary. If we imagine simulating a transmission tree, the offspring infection
 500 times can be generated in any order (of which there are $d^1!$ possible orderings) to produce the
 501 same transmission tree.

502 In summation, the likelihood contribution (sans secondary observations) for the root host in
 503 the unobserved case is

$$\begin{aligned} \mathcal{L}_T^1(\theta_T) &= \frac{(1 - \phi(x^1))(1 - \zeta(x^1))}{1 - \omega(x^1)} \times \\ &\quad \frac{1}{1 - \phi(x^1)} \sum_{k=d^1}^{\infty} \alpha(k) \binom{k}{d^1} \bar{\omega}(x^1)^{k-d^1} (1 - \bar{\omega}(x^1))^{d^1} \times \\ &\quad d^1! \prod_{j \in \mathcal{H}^1} \frac{(1 - \omega(x^j)) \gamma(x^j - x^1)}{1 - \bar{\omega}(x^1)} \\ &= \frac{(1 - \zeta(x^1))}{1 - \omega(x^1)} \sum_{k=d^1}^{\infty} \alpha(k) \binom{k}{d^1} \bar{\omega}(x^1)^{k-d^1} d^1! \prod_{j \in \mathcal{H}^1} (1 - \omega(x^j)) \gamma(x^j - x^1), \end{aligned} \quad (21)$$

504 and for the observed case is

$$\begin{aligned}
\mathcal{L}_T^1(\theta_T) &= \frac{\zeta(x^1)}{1 - \omega(x^1)} \frac{\sigma(y^1 - x^1)}{\int_0^{T-x^1} \sigma(\tau) d\tau} \times \\
&\quad \sum_{k=d^1}^{\infty} \alpha(k) \binom{k}{d^1} \bar{\omega}(x^1)^{k-d^1} (1 - \bar{\omega}(x^1))^{d^1} \times \\
&\quad d^1! \prod_{j \in \mathcal{H}^1} \frac{(1 - \omega(x^j)) \gamma(x^j - x^1)}{1 - \bar{\omega}(x^1)} \\
&= \frac{\pi \sigma(y^1 - x^1)}{1 - \omega(x^1)} \sum_{k=d^1}^{\infty} \alpha(k) \binom{k}{d^1} \bar{\omega}(x^1)^{k-d^1} d^1! \prod_{j \in \mathcal{H}^1} (1 - \omega(x^j)) \gamma(x^j - x^1).
\end{aligned} \tag{22}$$

505 The full likelihood is calculated by recursion, applying the same density calculations to each
506 included host, i.e.

$$p(\mathcal{T} | \theta_T) = \prod_{j=1}^N \mathcal{L}_T^j(\theta_T), \tag{23}$$

507 with N being the total number of included hosts. Note that in doing so, with the exception of
508 the root host, the terms $1 - \omega(x^j)$ will cancel in the likelihood.

509 Methods for simulating transmission trees are provided in Supplementary Text S1.

510 Coalescent model

511 In the original version of TransPhylo the coalescent model used was the bounded coalescent
512 (Carson et al., 2022). This model follows the standard coalescent model with heterochronous
513 sampling (Drummond et al., 2002), but conditions all lineages to coalesce before the infection
514 time of each host. Here we need to choose a coalescent model that allows for the transmission of
515 multiple lineages between hosts. With a bottleneck assumption many dated phylogenetic trees
516 would not permit the overlaying of a transmission tree under our stochastic branching model.

517 Here we assume that the within-host pathogen population size $q(\tau)$ grows linearly:

$$q(\tau) = \kappa + \lambda\tau, \tag{24}$$

518 where τ is the time since the host was infected. Should $\kappa = 0$ all lineages will coalesce by
519 the host's infection time. We could adopt alternative population models, so long as they are
520 integrable.

521 The likelihood of the phylogenetic tree conditional on the set of transmissions is calculated
522 by taking the product of the likelihood of each *subtree* for each host. The subtree of any
523 host j is formed by taking the parts of the phylogenetic tree assigned (coloured) by host j .
524 Each subtree is rooted at the host's infection time x^j , with the number of roots being the
525 number of lineages transmitted to the host. Leaves correspond to observations of the host and
526 transmissions to the hosts included offspring, noting that each transmission may contribute
527 multiple leaves (transmitting multiple lineages).

528 Let v_j^m , $m = 1, \dots, M_j$ be the times leaves are added within the subtree of host j , and let u_j^n ,
 529 $n = 1, \dots, N_j$ be the coalescence times, supposing $N_j > 0$. Then we define the number of extant
 530 lineages at time t as

$$L_j(t) = \sum_{m=1}^{M_j} \mathbb{I}(v_j^m \geq t) - \sum_{n=1}^{N_j} \mathbb{I}(u_j^n > t), \quad (25)$$

531 so that if t is the time of a coalescence, $L_j(t)$ is the number of lineages that could have coalesced.
 532 Denoting $\tau_j = t - x^j$, the phylogenetic likelihood contribution from each host is then

$$\mathcal{L}_{P|T}^j(\theta_P) = \exp\left(-\int_0^\infty \binom{L_j(x^j + \tau_j)}{2} \frac{1}{q(\tau_j)} d\tau_j\right) \prod_{n=1}^{N_j} \frac{1}{q(u_j^n - x^j)}, \quad (26)$$

533 and the full phylogenetic likelihood conditional on transmission tree \mathcal{T} is given by the product

$$p(\mathcal{P} | \mathcal{T}, \theta_P) = \prod_{j=1}^N \mathcal{L}_{P|T}^j(\theta_P). \quad (27)$$

534 Let w_j^k , $k = 0, \dots, K$ be the ordered set of root, leaf, and coalescence times, with $w_j^0 = x^j$. Let
 535 L_j^k be the number of lineages in the interval (w_j^{k-1}, w_j^k) . The integral in the exponent can then
 536 be partitioned accordingly

$$\int_0^\infty \binom{L_j(x^j + \tau_j)}{2} \frac{1}{q(\tau_j)} d\tau_j = \sum_{k=1}^n \int_{w_j^{k-1} - x^j}^{w_j^k - x^j} \binom{L_j^k}{2} \frac{1}{q(\tau_j)} d\tau_j. \quad (28)$$

537 For the linear growth model, these terms are then

$$\int_{w_j^{k-1} - x^j}^{w_j^k - x^j} \binom{L_j^k}{2} \frac{1}{q(\tau_j)} d\tau_j = \frac{\binom{L_j^k}{2}}{\lambda} \left(\log(\kappa + \lambda(w_j^k - x^j)) - \log(\kappa + \lambda(w_j^{k-1} - x^j)) \right) \quad (29)$$

538 Phylogenetic tree simulation is described in Supplementary Text S2.

539 Inference

540 Inference is undertaken using reversible-jump Markov chain Monte Carlo (Green, 1995). We
 541 iterate through the following update steps:

- 542 1. Update the transmission model parameters according to $p(\theta_T | \mathcal{T})$.
- 543 2. Update the coalescent model parameters according to $p(\theta_P | \mathcal{P}, \mathcal{T})$.
- 544 3. Update the transmission tree according to $p(\mathcal{T} | \mathcal{P}, \theta_T, \theta_P)$.

545 Steps 1 and 2 are performed using multivariate Gaussian random walks, conditional on the
 546 current transmission and phylogenetic trees. The scale and covariance in each case is determined

547 using the accelerated shaping and scaling algorithm of Spencer (2021) with target acceptance
548 $a = 0.234$ and forgetting sequence $f(n) = \lfloor 0.5n \rfloor$.

549 In Step 3 we randomly select from three proposals that update the transmission tree conditional
550 on the current model parameters: an add proposal for adding a new transmission to the current
551 transmission tree, a remove proposal for removing a transmission, and a local move proposal for
552 moving a transmission within the bounds set by its upstream and downstream transmissions.
553 The add and remove proposals form a reversible pair that change the dimension of the model,
554 whereas the local move proposal is its own reverse and maintains the dimension of the model.
555 Each proposal ensures that the new transmission tree is compatible with the phylogenetic tree.
556 For instance, observations from a single host cannot be split among multiple hosts when adding
557 a transmission. Likewise, observations from different hosts cannot be assigned to the same
558 host when removing a transmission. Full details including the acceptance probabilities for each
559 proposal are provided in Supplementary Text S3.

560 Step 3 makes relatively small changes to the transmission tree with each update. Additionally,
561 the computational cost is relatively cheap as we only need to evaluate the likelihood
562 contributions from the one or two affected hosts. Consequently it is beneficial to perform
563 Step 3 multiple times in each scan, in order to improve the mixing of the MCMC. In general,
564 we find that performing $\mathcal{O}(N)$ Step 3 updates in each scan works well, where N is the number
565 of primary observations.

566 **Implementation**

567 We implemented the methods above into a new R package called TransPhyloMulti which extends
568 TransPhylo. TransPhyloMulti is available at <https://github.com/DrJCarson/TransPhyloMulti>.
569 This repository also contains all the code and data needed to reproduce all results shown in
570 this paper. The R package `ape` was used to store, manipulate and visualise phylogenetic trees
571 (Paradis and Schliep, 2019).

572 **ACKNOWLEDGEMENTS**

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

References

- 575
- 576 Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. 2015. Measurably evolving pathogens in the
577 genomic era. *Trends in Ecology & Evolution*. 30:306–313.
- 578 Boeras DI, Hraber PT, Hurlston M, Evans-Strickfaden T, Bhattacharya T, Giorgi EE, Mulenga
579 J, Karita E, Korber BT, Allen S, et al. (11 co-authors). 2011. Role of donor genital tract
580 HIV-1 diversity in the transmission bottleneck. *Proceedings of the National Academy of
581 Sciences*. 108:E1156–E1163.
- 582 Bouckaert R, Vaughan TG, Barido-Sottani J, Duchêne S, Fourment M, Gavryushkina A,
583 Heled J, Jones G, Kühnert D, De Maio N, et al. (25 co-authors). 2019. BEAST 2.5 :
584 An Advanced Software Platform for Bayesian Evolutionary Analysis. *PLoS Computational
585 Biology*. 15:e1006650.
- 586 Brooks SP, Gelman A. 1998. General methods for monitoring convergence of iterative
587 simulations. *Journal of Computational and Graphical Statistics*. 7:434–455.
- 588 Bryant JM, Grogono DM, Greaves D, Foweraker J, Roddick I, Inns T, Reacher M, Haworth
589 CS, Curran MD, Harris SR, et al. (13 co-authors). 2013. Whole-genome sequencing to
590 identify transmission of *Mycobacterium abscessus* between patients with cystic fibrosis: A
591 retrospective cohort study. *The Lancet*. 381:1551–1560.
- 592 Campbell F, Didelot X, Fitzjohn R, Ferguson N, Cori A, Jombart T. 2018. Outbreaker2: A
593 Modular Platform for Outbreak Reconstruction. *BMC Bioinformatics*. 19:363.
- 594 Carson J, Ledda A, Ferretti L, Keeling M, Didelot X. 2022. The bounded coalescent
595 model: Conditioning a genealogy on a minimum root date. *Journal of Theoretical Biology*.
596 548:111186.
- 597 Cori A, Ferguson NM, Fraser C, Cauchemez S. 2013. A new framework and software to estimate
598 time-varying reproduction numbers during epidemics. *American journal of epidemiology*.
599 178:1505–12.
- 600 Cortey M, Ferretti L, Pérez-Martín E, Zhang F, de Klerk-Lorist LM, Scott K, Freimanis G,
601 Seago J, Ribeca P, van Schalkwyk L, et al. (11 co-authors). 2019. Persistent infection of
602 African buffalo (*Syncerus caffer*) with foot-and-mouth disease virus: limited viral evolution
603 and no evidence of antibody neutralization escape. *Journal of virology*. 93:10–1128.
- 604 Cottam EM, Wadsworth J, Shaw AE, Rowlands RJ, Goatley L, Maan S, Maan NS, Mertens
605 PPC, Ebert K, Li Y, et al. (18 co-authors). 2008. Transmission pathways of foot-and-mouth
606 disease virus in the United Kingdom in 2007. *PLoS Pathogens*. 4:e1000050.
- 607 De Maio N, Worby CJ, Wilson DJ, Stoesser N. 2018. Bayesian reconstruction of transmission
608 within outbreaks using genomic variants. *PLOS Computational Biology*. 14:e1006117.
- 609 De Maio N, Wu CH, O'Reilly KM, Wilson D. 2015. New Routes to Phylogeography: A Bayesian
610 Structured Coalescent Approximation. *PLoS Genetics*. 11:e1005421.
- 611 De Maio N, Wu CH, Wilson DJ. 2016. SCOTTI: Efficient Reconstruction of Transmission within
612 Outbreaks with the Structured Coalescent. *PLoS Computational Biology*. 12:e1005130.
- 613 Dearlove BL, Cody AJ, Pascoe B, Méric G, Wilson DJ, Sheppard SK, Daniel J, Sheppard SK.
614 2016. Rapid host switching in generalist *Campylobacter* strains erodes the signal for tracing
615 human infections. *The ISME journal*. 10:721–729.

- 616 Didelot X, Bowden R, Wilson DJ, Peto TEA, Crook DW. 2012. Transforming clinical
617 microbiology with bacterial genome sequencing. *Nature Reviews Genetics*. 13:601–612.
- 618 Didelot X, Croucher NJ, Bentley SD, Harris SR, Wilson DJ. 2018. Bayesian inference of
619 ancestral dates on bacterial phylogenetic trees. *Nucleic Acids Research*. 46:e134.
- 620 Didelot X, Fraser C, Gardy J, Colijn C. 2017. Genomic infectious disease epidemiology in
621 partially sampled and ongoing outbreaks. *Molecular Biology and Evolution*. 34:997–1007.
- 622 Didelot X, Gardy J, Colijn C. 2014. Bayesian inference of infectious disease transmission from
623 whole-genome sequence data. *Molecular Biology and Evolution*. 31:1869–1879.
- 624 Didelot X, Kendall M, Xu Y, White PJ, McCarthy N. 2021. Genomic Epidemiology Analysis
625 of Infectious Disease Outbreaks Using TransPhylo. *Current Protocols*. 1:e60.
- 626 Didelot X, Nell S, Yang I, Woltemate S, van der Merwe S, Suerbaum S. 2013. Genomic evolution
627 and transmission of *Helicobacter pylori* in two South African families. *Proceedings of the
628 National Academy of Sciences*. 110:13880–13885.
- 629 Didelot X, Parkhill J. 2022. A scalable analytical approach from bacterial genomes to
630 epidemiology. *Philosophical Transactions of the Royal Society B: Biological Sciences*.
631 377:20210246.
- 632 Didelot X, Walker AS, Peto TE, Crook DW, Wilson DJ. 2016. Within-host evolution of bacterial
633 pathogens. *Nature Reviews Microbiology*. 14:150–162.
- 634 Drummond AJ, Nicholls GK, Rodrigo AG, Solomon W. 2002. Estimating mutation parameters,
635 population history and genealogy simultaneously from temporally spaced sequence data.
636 *Genetics*. 161:1307–1320.
- 637 Duault H, Durand B, Canini L. 2022. Methods Combining Genomic and Epidemiological Data
638 in the Reconstruction of Transmission Trees: A Systematic Review. *Pathogens*. 11:252.
- 639 Friel N, Wyse J. 2012. Estimating the evidence—a review. *Statistica Neerlandica*. 66:288–308.
- 640 Gardy JL, Loman NJ. 2018. Towards a genomics-informed, real-time, global pathogen
641 surveillance system. *Nature Reviews Genetics*. 19:9–20.
- 642 Ghafari M, Lumby CK, Weissman DB, Illingworth CJR. 2020. Inferring Transmission Bottleneck
643 Size from Viral Sequence Data Using a Novel Haplotype Reconstruction Method. *Journal of
644 Virology*. 94:e00014–20.
- 645 Green PJ. 1995. Reversible Jump Markov Chain Monte Carlo Computation and Bayesian Model
646 Determination. *Biometrika*. 82:711–732.
- 647 Grenfell BT, Pybus OG, Gog JR, Wood JLN, Daly JM, Mumford JA, Holmes EC. 2004.
648 Unifying the epidemiological and evolutionary dynamics of pathogens. *Science*. 303:327–332.
- 649 Grote A, Earl AM. 2022. Within-host evolution of bacterial pathogens during persistent
650 infection of humans. *Current Opinion in Microbiology*. 70:102197.
- 651 Hall M, Woolhouse M, Rambaut A. 2015. Epidemic Reconstruction in a Phylogenetics
652 Framework: Transmission Trees as Partitions of the Node Set. *PLOS Computational Biology*.
653 11:e1004613.

- 654 Hall MD, Holden MT, Srisomang P, Mahavanakul W, Wuthiekanun V, Limmathurotsakul D,
655 Fountain K, Parkhill J, Nickerson EK, Peacock SJ, et al. (11 co-authors). 2019. Improved
656 characterisation of MRSA transmission using within-host bacterial sequence diversity. *eLife*.
657 8:e46402.
- 658 Ho SYW, Shapiro B. 2011. Skyline-plot methods for estimating demographic history from
659 nucleotide sequences. *Molecular Ecology Resources*. 11:423–434.
- 660 Jelley L, Douglas J, Ren X, Winter D, McNeill A, Huang S, French N, Welch D, Hadfield J,
661 de Ligt J, et al. (11 co-authors). 2022. Genomic epidemiology of delta sars-cov-2 during
662 transition from elimination to suppression in aotearoa new zealand. *Nature Communications*.
663 13:4035.
- 664 Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N. 2014. Bayesian
665 Reconstruction of Disease Outbreaks by Combining Epidemiologic and Genomic Data. *PLoS*
666 *Computational Biology*. 10:e1003457.
- 667 Jombart T, Eggo RM, Dodd PJ, Balloux F. 2011. Reconstructing disease outbreaks from genetic
668 data: A graph approach. *Heredity*. 106:383–390.
- 669 Kapli P, Yang Z, Telford MJ. 2020. Phylogenetic tree building in the genomic age. *Nature*
670 *Reviews Genetics*. 21:428–444.
- 671 Klinkenberg D, Backer JA, Didelot X, Colijn C, Wallinga J. 2017. Simultaneous inference of
672 phylogenetic and transmission trees in infectious disease outbreaks. *PLoS Computational*
673 *Biology*. 13:e1005495.
- 674 Leitner T. 2019. Phylogenetics in HIV transmission: Taking within-host diversity into account.
675 *Current Opinion in HIV and AIDS*. 14:181–187.
- 676 Lemey P, Rambaut A, Drummond AJ, Suchard MA. 2009. Bayesian phylogeography finds its
677 roots. *PLoS computational biology*. 5:e1000520.
- 678 Lieberman TD, Michel JB, Aingaran M, Potter-Bynoe G, Roux D, Davis MR Jr, Skurnik D,
679 Leiby N, LiPuma JJ, Goldberg JB, et al. (13 co-authors). 2011. Parallel bacterial evolution
680 within multiple patients identifies candidate pathogenicity genes. *Nature genetics*. 43:1275–
681 1280.
- 682 Lloyd-Smith JO, Schreiber SJ, Kopp PE, Getz WM. 2005. Superspreading and the effect of
683 individual variation on disease emergence. *Nature*. 438:355–359.
- 684 Marvig RL, Johansen HK, Molin S, Jelsbak L. 2013. Genome analysis of a transmissible lineage
685 of *pseudomonas aeruginosa* reveals pathoadaptive mutations and distinct evolutionary paths
686 of hypermutators. *PLoS genetics*. 9:e1003741.
- 687 Mather AE, Reid SWJ, Maskell DJ, Parkhill J, Fookes MC, Harris SR, Brown DJ, Coia JE,
688 Mulvey MR, Gilmour MW. 2013. Distinguishable epidemics of multidrug-resistant *Salmonella*
689 *Typhimurium* DT104 in different hosts. *Science*. 341:1514–1517.
- 690 Nylander JAA, Olsson U, Alström P, Sanmartín I. 2008. Accounting for phylogenetic
691 uncertainty in biogeography: A Bayesian approach to dispersal-vicariance analysis of the
692 thrushes (Aves: *Turdus*). *Systematic Biology*. 57:257–68.
- 693 Paradis E, Schliep K. 2019. Ape 5.0: An environment for modern phylogenetics and evolutionary
694 analyses in R. *Bioinformatics*. 35:526–528.

- 695 Pybus OG, Charleston MA, Gupta S, Rambaut A, Holmes EC, Harvey PH. 2001. The Epidemic
696 Behavior of the Hepatitis C Virus. *Science*. 292:2323–2325.
- 697 Pybus OG, Rambaut A. 2009. Evolutionary analysis of the dynamics of viral infectious disease.
698 *Nature Reviews Genetics*. 10:540–550.
- 699 Rasmussen DA, Volz EM, Koelle K. 2014. Phylodynamic Inference for Structured
700 Epidemiological Models. *PLoS Computational Biology*. 10:e1003570.
- 701 Rau MH, Marvig RL, Ehrlich GD, Molin S, Jelsbak L. 2012. Deletion and acquisition of
702 genomic content during early stage adaptation of *Pseudomonas aeruginosa* to a human host
703 environment. *Environmental microbiology*. 14:2200–2211.
- 704 Romero-Severson E, Skar H, Bulla I, Albert J, Leitner T. 2014. Timing and order of transmission
705 events is not directly reflected in a pathogen phylogeny. *Molecular Biology and Evolution*.
706 31:2472–2482.
- 707 Romero-Severson EO, Bulla I, Leitner T. 2016. Phylogenetically resolving epidemiologic linkage.
708 *Proceedings of the National Academy of Sciences*. 113:2690–2695.
- 709 Rossi E, La Rosa R, Bartell JA, Marvig RL, Haagensen JAJ, Sommer LM, Molin S, Johansen
710 HK. 2021. *Pseudomonas aeruginosa* adaptation and evolution in patients with cystic fibrosis.
711 *Nature Reviews Microbiology*. 19:331–342.
- 712 Sagulenko P, Puller V, Neher RA. 2018. TreeTime: Maximum-likelihood phylodynamic analysis.
713 *Virus Evolution*. 4:vex042.
- 714 Sashittal P, El-Kebir M. 2020. Sampling and summarizing transmission trees with multi-strain
715 infections. *Bioinformatics (Oxford, England)*. 36:i362–i370.
- 716 Spencer SEF. 2021. Accelerating adaptation in the adaptive Metropolis-Hastings random walk
717 algorithm. *Australian & New Zealand Journal of Statistics*. 63:468–484.
- 718 Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. 2018. Bayesian
719 phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evolution*.
720 4:vey016.
- 721 Tartakovsky AG, Moustakides GV. 2010. State-of-the-art in Bayesian changepoint detection.
722 *Sequential Analysis*. 29:125–145.
- 723 Tonkin-Hill G, Ling C, Chaguza C, Salter SJ, Hinfonthong P, Nikolaou E, Tate N, Pastusiak A,
724 Turner C, Chewapreecha C, et al. (15 co-authors). 2022. Pneumococcal within-host diversity
725 during colonization, transmission and treatment. *Nature Microbiology*. 7:1791–1804.
- 726 Torres Ortiz A, Kendall M, Storey N, Hatcher J, Dunn H, Roy S, Williams R, Williams
727 C, Goldstein RA, Didelot X, et al. (13 co-authors). 2023. Within-host diversity improves
728 phylogenetic and transmission reconstruction of SARS-CoV-2 outbreaks. *eLife*. 12:e84384.
- 729 van Dorp L, Wang Q, Shaw LP, Acman M, Brynildsrud OB, Eldholm V, Wang R, Gao H, Yin
730 Y, Chen H, et al. (15 co-authors). 2019. Rapid phenotypic evolution in multidrug-resistant
731 *Klebsiella pneumoniae* hospital outbreak strains. *Microbial Genomics*. 5:e000263.
- 732 Volz EM, Frost SDW. 2017. Scalable relaxed clock phylogenetic dating. *Virus Evolution*.
733 3:vex025.

- 734 Volz EM, Koelle K, Bedford T. 2013. Viral Phylodynamics. *PLoS Computational Biology*.
735 9:e1002947.
- 736 Wallinga J, Teunis P. 2004. Different Epidemic Curves for Severe Acute Respiratory Syndrome
737 Reveal Similar Impacts of Control Measures. *American Journal of Epidemiology*. 160:509–516.
- 738 Worby CJ, O’Neill PD, Kypraios T, Robotham JV, De Angelis D, Cartwright EJP, Peacock SJ,
739 Cooper BS. 2016. Reconstructing transmission trees for communicable diseases using densely
740 sampled genetic data. *Annals of Applied Statistics*. 10:395–417.
- 741 Wymant C, Hall M, Ratmann O, Bonsall D, Golubchik T, de Cesare M, Gall A, Cornelissen M,
742 Fraser C, STOP-HCV Consortium, et al. (12 co-authors). 2018. PHYLOSCANNER: Inferring
743 Transmission from Within- and Between-Host Pathogen Genetic Diversity. *Molecular Biology
744 and Evolution*. 35:719–733.
- 745 Xu Y, Stockdale JE, Naidu V, Hatherell H, Stimson J, Stagg HR, Abubakar I, Colijn C. 2020.
746 Transmission analysis of a large tuberculosis outbreak in London: A mathematical modelling
747 study using genomic data. *Microbial Genomics*. 6:e000450.
- 748 Yang L, Jelsbak L, Marvig RL, Damkiær S, Workman CT, Rau MH, Hansen SK, Folkesson A,
749 Johansen HK, Ciofu O, et al. (13 co-authors). 2011. Evolutionary dynamics of bacteria in a
750 human host environment. *Proceedings of the National Academy of Sciences*. 108:7481–7486.
- 751 Yang S, Hemarajata P, Hindler J, Li F, Adisetiyo H, Aldrovandi G, Sebra R, Kasarskis A,
752 MacCannell D, Didelot X, et al. (13 co-authors). 2017. Evolution and Transmission of
753 Carbapenem-Resistant *Klebsiella pneumoniae* Expressing the bla_{OXA-232} Gene During an
754 Institutional Outbreak Associated With Endoscopic Retrograde Cholangiopancreatography.
755 *Clinical Infectious Diseases*. 64:894–901.
- 756 Yang Z, Rannala B. 2012. Molecular phylogenetics: Principles and practice. *Nature Reviews
757 Genetics*. 13:303–314.
- 758 Young BC, Golubchik T, Batty EM, Fung R, Larner-Svensson H, Votintseva AA, Miller RR,
759 Godwin H, Knox K, Everitt RG, et al. (21 co-authors). 2012. Evolutionary dynamics
760 of *Staphylococcus aureus* during progression from carriage to disease. *Proceedings of the
761 National Academy of Sciences*. 109:4550–4555.