

1 Ancestral process for infectious disease outbreaks with superspreading

2 Xavier Didelot^{1,2,*}, David Helekal³, Ian Roberts²

3 ¹ School of Life Sciences, University of Warwick, Coventry, United Kingdom

4 ² Department of Statistics, University of Warwick, Coventry, United Kingdom

5 ³ Department of Immunology and Infectious Diseases, Harvard T. H. Chan School of Public Health,
6 Boston, Massachusetts, USA

7 * Corresponding author. Tel: 0044 (0)2476 572827. Email: `xavier.didelot@gmail.com`

8 Running title: Ancestry for outbreaks with superspreading

9 Keywords: infectious disease epidemiology modelling; offspring distribution; superspreading;
10 outbreaks; lambda-coalescent model; multiple mergers

Abstract

When an infectious disease outbreak is of relatively small size, describing the ancestry of a sample of infected individuals is difficult because most ancestral models assume large population sizes. Given a set of sampled individuals, we show that it is possible to express exactly the probability that a subset of them have the same infector, either inclusively (so that other sampled individuals may have the same infector too) or exclusively (so that they may not). To compute these probabilities requires knowledge of the offspring distribution, which determines how many infections each infected individual causes. We consider transmission both without and with superspreading, in the form of a Poisson and Negative-Binomial offspring distribution, respectively. We show how our results can be incorporated into a new lambda-coalescent model which allows multiple lineages to coalesce together. We compare this model with previously proposed alternatives, and advocate the use of this new model for future studies of small outbreaks.

1 Introduction

An outbreak of an infectious disease typically starts when a single or a small number of infected individuals appear within a susceptible population. Each infected individual may come in contact and infect each of the susceptible individuals, who will then become infected in their turn and spread the disease further. Most infectious disease modelling theory describes situations where the disease is at an equilibrium, when the number of infected individuals is high and/or with a significant part of the population already infected (Anderson and May 1991; Keeling and Rohani 2008). Here however we focus on the early stages of an epidemic, where the number of infected individuals is small and the number of susceptibles relatively high and unchanging. In this situation it is useful to think about the number of infections that each newly infected individual is likely to cause, and the probabilistic distribution for this number is often called the offspring distribution (Grassly and Fraser 2008). The mean of the offspring distribution is called the basic reproduction number R_0 and has been given much attention especially since it determines how likely the outbreak is to spread, and how much effort would be needed to bring it under control (Fraser et al. 2004; Ferguson et al. 2006).

If we consider that all individuals are infectious for the same duration and with the same infectiousness, the offspring distribution is Poisson distributed with mean R_0 , which means that the variance of the offspring distribution is also R_0 . We would then say that there is no transmission heterogeneity. However, in practice there are many reasons why this may not be the case, with some individuals being infectious for longer, or being more infectious than others, or having more contacts with susceptibles, or being less symptomatic and therefore less likely to reduce contact numbers, etc. All these factors cause the offspring distribution to be more dispersed than it would otherwise be, that is to have a variance greater than its mean R_0 . A frequent choice to capture this overdispersion is to model the offspring distribution using a Negative-Binomial distribution with mean R_0 and dispersion parameter r (Lloyd-Smith et al. 2005; Grassly and Fraser 2008). When r is close to zero the variance is high compared to the mean, whereas when r is high the variance becomes close to the mean. This transmission heterogeneity is often called superspreading, although this is perhaps misleading as it is the rule rather than the exception of how infectious diseases spread. Superspreading has indeed been described in many diseases (Woolhouse et al. 1997; Stein 2011; Kucharski and Althaus 2015; Wang et al. 2021), and most recently for SARS-CoV-2 (Wang et al. 2020; Lemieux et al. 2021; Gómez-Carballa et al. 2021; Du et al. 2022).

As an outbreak unfolds forward-in-time, a transmission tree is generated representing who-infected-whom, in which each node is an infected individual and points towards a number of nodes distributed

55 according to the offspring distribution. Here we consider the reverse problem of the transmission
56 ancestry, going backward-in-time, from a sample of infected individuals, until reaching the last common
57 transmission ancestor of the whole sample. Given a sample of n sampled individuals, we show how
58 to calculate the probability that a given subset of size k have the same infector, either inclusively (so
59 that the remaining $n - k$ may also have the same infector or not) or exclusively (so that none of the
60 remaining $n - k$ have the same infector). We start by considering the general case of an offspring
61 distribution with arbitrary form, and then the specific cases of offspring distributions that follow a
62 Poisson or a Negative-Binomial distribution. The main novelty of our approach is that we consider that
63 the overall population size is small, but we show that if the population size is large, our results agree
64 with several previous studies (Volz 2012; Koelle and Rasmussen 2012; Fraser and Li 2017). Finally, we
65 show how our results can be incorporated into a new lambda-coalescent model (Pitman 1999; Sagitov
66 1999; Donnelly and Kurtz 1999) and compare it with previously described models.

67 2 General case

68 Let time be measured in discrete units and denoted t . Each discrete value of t correspond to a unique
69 non-overlapping generations of infected individuals, so that individuals infected at t will have offspring
70 at $t + 1$, etc. Let N_t denote the number of infectious individuals at time t . Each of them creates a
71 number $s_{t,i}$ of secondary infections at time $t + 1$, following the offspring distribution $\alpha_t(s)$. The mean
72 of this distribution is the basic reproduction number R_t and the variance is V_t . We have:

$$N_{t+1} = \sum_{i=1}^{N_t} s_{t,i} \quad (1)$$

73 2.1 Inclusive coalescence probability

74 We define the inclusive coalescence probability $p_{k,t}(N_t, N_{t+1})$ as the probability that a specific set of
75 k individuals from generation $t + 1$ find a common ancestor in generation t , conditional on population
76 sizes N_t and N_{t+1} .

77 Given full information about offspring counts from individuals in generation t , $\mathbf{s}_t = (s_{t,1}, \dots, s_{t,N_t})$, we
78 have

$$\begin{aligned}
p_{k,t}(\mathbf{s}_t, N_t) &= \sum_{i=1}^{N_t} \frac{\binom{s_{t,i}}{k}}{\binom{N_{t+1}}{k}} \\
&= \sum_{i=1}^{N_t} \frac{\Gamma(s_{t,i} + 1) \Gamma(N_{t+1} - k + 1)}{\Gamma(s_{t,i} - k + 1) \Gamma(N_{t+1})}
\end{aligned} \tag{2}$$

79

80 Full information $\{s_{t,i}\}$ yields the population size N_{t+1} but is not feasible to observe in practice. We
81 can instead express the inclusive coalescence probability conditioning on the next population size N_{t+1}
82 by summing over possible offspring counts $\mathbf{s}_t = (s_{t,1}, \dots, s_{t,N_t})$ conditional on the total generation size.
83 Let $S_t^{-(1)} = (S_{t,2}, \dots, S_{t,N_t})$.

$$\begin{aligned}
p_{k,t}(N_t, N_{t+1}) &= \sum_{\mathbf{s}_t \in \mathbb{N}_0^{N_t}} \mathbb{P} \left[\mathbf{S}_t = \mathbf{s}_t \mid \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right] p_{k,t}(\mathbf{s}_t, N_t) \\
&= \sum_{\mathbf{s}_t \in \mathbb{N}_0^{N_t}} \mathbb{P} \left[\mathbf{S}_t = \mathbf{s}_t \mid \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right] \sum_{i=1}^{N_t} \frac{\binom{s_{t,i}}{k}}{\binom{N_{t+1}}{k}} \\
&= \sum_{i=1}^{N_t} \sum_{\mathbf{s}_t \in \mathbb{N}_0^{N_t}} \frac{\binom{s_{t,i}}{k}}{\binom{N_{t+1}}{k}} \mathbb{P} \left[S_{t,1} = s_{t,1}, \mathbf{S}_t^{-(1)} = \mathbf{s}_t^{-(1)} \mid \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right] \\
&= \frac{N_t}{\binom{N_{t+1}}{k}} \sum_{\mathbf{s}_t \in \mathbb{N}_0^{N_t}} \binom{s_{t,1}}{k} \mathbb{P} \left[S_{t,1} = s_{t,1} \mid \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right] \\
&\quad \times \mathbb{P} \left[\mathbf{S}_t^{-(1)} = \mathbf{s}_t^{-(1)} \mid S_{t,1} = s_{t,1}, \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right] \\
&= \frac{N_t}{\binom{N_{t+1}}{k}} \sum_{s_{t,1}=0}^{N_{t+1}} \binom{s_{t,1}}{k} \mathbb{P} \left[S_{t,1} = s_{t,1} \mid \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right] \\
&\quad \times \underbrace{\sum_{\mathbf{s}_t^{-(1)} \in \mathbb{N}_0^{N_t-1}} \mathbb{P} \left[\mathbf{S}_t^{-(1)} = \mathbf{s}_t^{-(1)} \mid \sum_{i=2}^{N_t} S_{t,i} = N_{t+1} - s_{t,1} \right]}_{=1} \\
&= \frac{N_t}{\binom{N_{t+1}}{k}} \mathbb{E} \left[\binom{S_{t,1}}{k} \mid \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right]
\end{aligned} \tag{3}$$

84

85 The k -th falling factorial moments $\mathbb{E} \left[\frac{S_{t,1}!}{(S_{t,1}-k)!} \mid \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right]$ in Equation 3 can be readily obtained

86 by differentiating the probability generating function of $S_{t,1} | (\sum_{i=1}^{N_t} S_{t,i} = N_{t+1})$.

87 2.2 Exclusive coalescence probability

88 Generally, we observe a sample of individuals from each generation rather than the entire population.
 89 In this case, we are interested in the exclusive coalescence probability $p_{n,k,t}(N_t, N_{t+1})$ that exactly k
 90 individuals from a sample of n arose from a common ancestor one generation in the past given knowlege
 91 of the total population sizes N_t and N_{t+1} .

92 Given full information about offspring counts of the parents of sampled individuals at the present,
 93 $\mathbf{x}_t = (x_{t,1}, \dots, x_{t,N_t})$, we have

$$\begin{aligned} p_{n,k,t}(\mathbf{x}_t, N_t) &= \sum_{i=1}^{N_t} \frac{\binom{x_{t,i}}{k}}{\binom{n}{k}} \mathbb{I}\{x_{t,i} = k\} \\ &= \sum_{i=1}^{N_t} \frac{x_{t,i}!}{(x_{t,i} - k)!} \frac{(n - k)!}{n!} \mathbb{I}\{x_{t,i} = k\} \end{aligned} \quad (4)$$

94 Similarly to the exclusive coalescence probability, we can use this to evaluate the exclusive probability
 95 given N_t and N_{t+1} by summing over possible parent offspring configurations (for $k \leq n$),

$$\begin{aligned}
p_{n,k,t}(N_t, N_{t+1}) &= \sum_{\mathbf{x}_t \in \mathbb{N}_0^{N_t}} \mathbb{P} \left[\mathbf{X}_t = \mathbf{x}_t \middle| \sum_{i=1}^n X_{t,i} = n \right] p_{n,k,t}(\mathbf{x}_t, N_t) \\
&= \sum_{\mathbf{x}_t \in \mathbb{N}_0^{N_t}} \mathbb{P} \left[\mathbf{X}_t = \mathbf{x}_t \middle| \sum_{i=1}^n X_{t,i} = n \right] \sum_{i=1}^{N_t} \frac{\binom{x_{t,i}}{k}}{\binom{n}{k}} \mathbb{I}\{x_{t,i} = k\} \\
&= \frac{N_t}{\binom{n}{k}} \sum_{\mathbf{x}_t \in \mathbb{N}_0^{N_t}} \binom{x_{t,1}}{k} \mathbb{P} \left[\mathbf{X}_t = \mathbf{x}_t \middle| \sum_{i=1}^{N_t} X_{t,i} = n \right] \mathbb{I}\{x_{t,1} = k\} \\
&= \frac{N_t}{\binom{n}{k}} \sum_{\mathbf{x}_t^{-(1)} \in \mathbb{N}_0^{N_t-1}} \binom{k}{k} \mathbb{P} \left[X_{t,1} = k, \mathbf{X}_t^{-(1)} = \mathbf{x}_t^{-(1)} \middle| \sum_{i=1}^{N_t} X_{t,i} = n \right] \\
&= \frac{N_t}{\binom{n}{k}} \mathbb{P}[X_{t,1} = k \middle| \sum_{i=1}^{N_t} X_{t,i} = n] \underbrace{\sum_{\mathbf{x}_t^{-(1)} \in \mathbb{N}_0^{N_t-1}} \mathbb{P} \left[\mathbf{X}_t^{-(1)} = \mathbf{x}_t^{-(1)} \middle| \sum_{i=1}^{N_t} X_{t,i} = n, X_{t,1} = k \right]}_{=1} \\
&= \frac{N_t}{\binom{n}{k}} \mathbb{P} \left[X_{t,1} = k \middle| \sum_{i=1}^{N_t} X_{t,i} = n \right] \tag{5}
\end{aligned}$$

96 Note that $X_{t,i}$ does not follow the same offspring distribution as $S_{t,i}$. $(X_{t,1}, \dots, X_{t,N_t})$ consists of n
 97 individuals sampled from generation $t+1$ without replacement - there is no guarantee that all offspring
 98 from any given parent are included in the sample.

99 **2.3 Complementarity of exclusive coalescence probabilities**

100 If we consider one of the lines observed amongst a set of n , it can either remain uncoalesced (with
 101 probability $p_{n,1,t}$) or coalesce in an event of size k (with probability $p_{n,k,t}$) with any set of $k-1$ lines
 102 among the $n-1$ other lines, leading to the following complementarity equation:

$$\sum_{k=1}^n \binom{n-1}{k-1} p_{n,k,t} = 1 \tag{6}$$

103 We can show that it is indeed satisfied by the formula in Equation 5:

$$\begin{aligned}
\sum_{k=1}^n \binom{n-1}{k-1} p_{n,k,t} &= \sum_{k=1}^n \binom{n-1}{k-1} \frac{N_t}{\binom{n}{k}} \mathbb{P} \left[X_1 = k \middle| \sum_{i=1}^{N_t} X_i = n \right] \\
&= \sum_{k=1}^n N_t \frac{k}{n} \mathbb{P} \left[X_1 = k \middle| \sum_{i=1}^{N_t} X_i = n \right] \\
&= \frac{N_t}{n} \sum_{k=0}^n k \mathbb{P} \left[X_1 = k \middle| \sum_{i=1}^{N_t} X_i = n \right] \\
&= \frac{N_t}{n} \mathbb{E} \left[X_1 \middle| \sum_{i=1}^{N_t} X_i = n \right] \\
&= \frac{1}{n} \sum_{i=1}^{N_t} \mathbb{E} \left[X_i \middle| \sum_{i=1}^{N_t} X_i = n \right] \\
&= \frac{1}{n} \mathbb{E} \left[\sum_{i=1}^{N_t} X_i \middle| \sum_{i=1}^{N_t} X_i = n \right] \\
&= 1
\end{aligned} \tag{7}$$

104 3 Poisson case

105 In this section we consider that the offspring distribution is $\alpha_t = \text{Poisson}(R_t)$. In this case, we have:

$$\sum_{i=1}^{N_t} S_{t,i} \sim \text{Poisson}(N_t R_t) \tag{8}$$

106 and the conditional distribution:

$$\begin{aligned}
\mathbb{P} \left[S_{t,1} = s \middle| \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right] &= \frac{\mathbb{P} \left[S_{t,1} = s, \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right]}{\mathbb{P} \left[\sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right]} \\
&= \frac{\alpha_t(s) \mathbb{P} \left[\sum_{i=2}^{N_t} S_{t,i} = N_{t+1} - s \right]}{\mathbb{P} \left[\sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right]} \\
&= \frac{\frac{R_t^s e^{-R_t}}{s!} \cdot \frac{((N_t - 1)R_t)^{N_{t+1} - s}}{(N_{t+1} - s)!}}{\frac{(N_t R_t)^{N_{t+1}} e^{-N_t R_t}}{N_{t+1}!}}
\end{aligned}$$

$$= \binom{N_{t+1}}{s} \left(\frac{1}{N_t}\right)^s \left(1 - \frac{1}{N_t}\right)^{N_{t+1}-s} \quad (9)$$

107

108 This is the probability mass function of a Binomial distribution and therefore we deduce that:

$$S_{t,1} \left| \left(\sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right) \right. \sim \text{Binomial} \left(N_{t+1}, \frac{1}{N_t} \right) \quad (10)$$

109 The k -th falling factorial moments of $X \sim \text{Binomial}(n, p)$ are (Potts 1953):

$$\mathbb{E} \left[\frac{X!}{(X-k)!} \right] = \binom{n}{k} p^k k! \quad (11)$$

110 By applying this formula to the Binomial distribution in Equation 10 and injecting into Equation 3
111 we obtain the inclusive probability of coalescence for k lines:

$$\mathbb{E} \left[\binom{S_{t,1}}{k} \left| \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right. \right] = \frac{1}{k!} \mathbb{E} \left[\frac{S_{t,1}!}{(S_{t,1}-k)!} \left| \sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right. \right] = \frac{1}{k!} \frac{N_{t+1}!}{(N_{t+1}-k)!} \left(\frac{1}{N_t} \right)^k \quad (12)$$

112 Consequently, the inclusive probability of coalescence for k lines is

$$p_{k,t} = \frac{1}{N_t^{k-1}} \quad (13)$$

113 By injecting the probability mass function of the Binomial distribution in Equation 10 into Equation
114 5 we deduce that the exclusive probability of coalescence for k lines from a sample of n ($n \geq k$) is

$$p_{n,k,t} = \frac{(N_t - 1)^{n-k}}{N_t^{n-1}} \quad (14)$$

115 It is interesting to note that neither the inclusive nor the exclusive coalescence probability depend on
116 the mean R_t of the Poisson offspring distribution or the size N_{t+1} of the population at time $t+1$. The
117 inclusive coalescent probability in Equation 13 can also be obtained conceptually by considering that
118 among the k lines, the first one has an ancestor with probability one, and the remaining $k-1$ need to
119 have the same ancestor among a set of N_t from which they choose uniformly at random so that the

120 probability of picking the same ancestor is $1/N_t$. The exclusive coalescent probability in Equation 14
 121 can be derived likewise by considering that in addition to the above, each of the $n - k$ other lines need
 122 to choose a different ancestor, which happens with probability $(N_t - 1)/N_t$.

123 Figure 1 illustrates the inclusive and exclusive coalescence probabilities for the Poisson case for a set
 124 of size $k = 1$ to $k = 10$ amongst a total of $n = 10$ observed lines, in a population of size $N_t = 10$,
 125 $N_t = 20$ or $N_t = 30$.

126 4 Negative-Binomial case

127 In this section we consider that the offspring distribution is $\alpha_t = \text{Negative-Binomial}(r, p)$ with
 128 parameters (r, p) set by moment-matching mean R_t and variance V_t which are assumed constant
 129 over time. The resulting parameters for this distribution are $r = R_t^2/(V_t - R_t)$ and $p = R_t/V_t$. In this
 130 case, we have:

$$\sum_{i=1}^{N_t} S_{t,i} \sim \text{Negative-Binomial}(N_t r, p) \quad (15)$$

131 and similarly to the Poisson(λ) offspring distribution identify the conditional distribution of
 132 $S_{t,1} | \sum_{i=1}^{N_t} S_{t,i}$ is as follows:

$$\begin{aligned} \mathbb{P}\left[S_{t,1} = s \mid \sum_{i=1}^{N_t} S_{t,i} = N_{t+1}\right] &= \frac{\alpha_t(s) \cdot \mathbb{P}\left[\sum_{i=2}^{N_t} S_{t,i} = N_{t+1} - s\right]}{\mathbb{P}\left[\sum_{i=1}^{N_t} S_{t,i} = N_{t+1}\right]} \\ &= \frac{\frac{\Gamma(r+s)}{s!\Gamma(r)}(1-p)^s p^r \cdot \frac{\Gamma((N_t-1)r + (N_{t+1}-s))}{(N_{t+1}-s)!\Gamma((N_t-1)r)}(1-p)^{N_{t+1}-s} p^{(N_t-1)r}}{\frac{\Gamma(N_t r + N_{t+1})}{N_{t+1}!\Gamma(N_t r)}(1-p)^{N_{t+1}} p^{N_t r}} \\ &= \frac{N_{t+1}!}{s!(N_{t+1}-s)!} \frac{\Gamma(r+s)\Gamma((N_t-1)r + (N_{t+1}-s))}{\Gamma(N_t r + N_{t+1})} \frac{\Gamma(N_t r)}{\Gamma(r)\Gamma((N_t-1)r)} \\ &= \binom{N_{t+1}}{s} \frac{B(s+r, N_{t+1}-s+(N_t-1)r)}{B(r, (N_t-1)r)} \end{aligned} \quad (16)$$

133

134 where $B(x, y)$ denotes the Beta function defined as $B(x, y) = \Gamma(x)\Gamma(y)/\Gamma(x+y)$. This is the probability
 135 mass function of a Beta-Binomial distribution and therefore we deduce that:

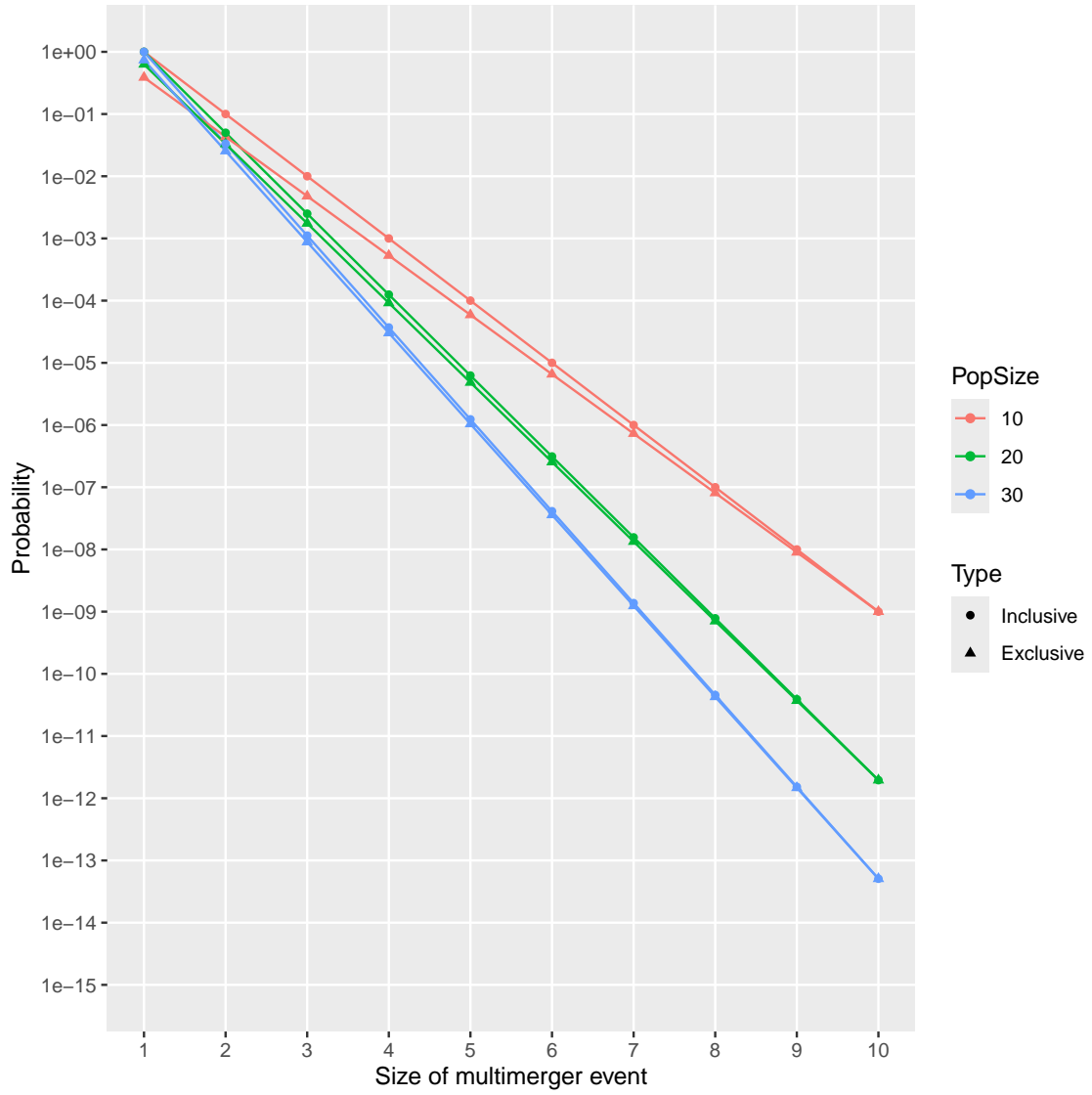


Figure 1: Inclusive and exclusive coalescence probabilities for the Poisson case.

$$S_{t,1} \left| \left(\sum_{i=1}^{N_t} S_{t,i} = N_{t+1} \right) \right. \sim \text{Beta-Binomial}(r, (N_t - 1)r) \quad (17)$$

136 The k -th falling factorial moments of $X \sim \text{Beta-Binomial}(\alpha, \beta)$ are (Tripathi et al. 1994):

$$\mathbb{E} \left[\frac{X!}{(X-k)!} \right] = \binom{n}{k} \frac{B(\alpha + k, \beta)k!}{B(\alpha, \beta)} \quad (18)$$

137 By applying this formula to the Beta-Binomial distribution in Equation 17 and injecting into Equation
138 3, we deduce that the inclusive probability of coalescence for k lines is:

$$p_{k,t} = \frac{B(N_t r + 1, r + k)}{B(r + 1, N_t r + k)} \quad (19)$$

139 By injecting the probability mass function of the Beta-Binomial distribution in Equation 17 into
140 Equation 5 we deduce that the exclusive probability of coalescence for k lines is:

$$p_{n,k,t} = \frac{N_t B(k + r, n - k + N_t r - r)}{B(r, N_t r - r)} \quad (20)$$

141 It is interesting to note that as for the Poisson case, the inclusive and exclusive coalescence probabilities
142 do not depend on the size N_{t+1} of the population at time $t + 1$. They both depend on the Negative-
143 Binomial offspring distribution only through the dispersion parameter r . If we consider that r is large
144 in Equations 19 and 20, we can derive that the asymptotic behaviour is the same as in the Poisson
145 case shown in Equations 13 and 14. For example this can be derived by rewriting the Beta functions
146 using Gamma functions, and using the following form of Stirling's approximation:

$$\lim_{a \rightarrow \infty} \frac{\Gamma(a + b)}{\Gamma(a)} = a^b e^{-b} \quad (21)$$

147 Figure 2 illustrates the inclusive and exclusive coalescence probabilities for the Negative-Binomial case
148 for a set of size $k = 1$ to $k = 10$ amongst a total of $n = 10$ observed lines, in a population of size
149 $N_t = 12$. Several Negative-Binomial offspring distributions are compared, all of which have the same
150 mean $R_t = 2$, and with the dispersion parameter equal to $r = 1$, $r = 2$, $r = 10$ and $r = 100$. When

151 $r = 1$ the Negative-Binomial reduces to a Geometric distribution. When r is high (for example $r = 100$
 152 as shown in Figure 2) the dispersion is low and the Negative-Binomial case behaves almost like the
 153 Poisson case. When r is lower the dispersion of the offspring distribution increases, so that both the
 154 inclusive and exclusive probabilities of larger multimerger events increase.

155 5 Limit when the population size is large

156 If we consider that the population size N_t is fixed and large, we can show the connections between our
 157 results and several previous studies. In the Poisson case, from Equations 13 and 14 we can see that
 158 both inclusive and exclusive probabilities are of order $\mathcal{O}(N_t^{1-k})$. We can therefore ignore events with
 159 $k > 2$ and retain only the events with $k = 2$ which occur with probability:

$$p_{2,t} = p_{n,2,t} = \frac{1}{N_t} \quad (22)$$

160 For the Negative-Binomial case, from Equations 19 and 20 we can rewrite using Gamma functions
 161 and apply the form of Stirling's equation given in Equation 21 to show that once again both inclusive
 162 and exclusive probabilities are also of order $\mathcal{O}(N_t^{1-k})$. We can therefore ignore events with $k > 2$ and
 163 retain only the events with $k = 2$ which occur with probability:

$$p_{2,t} = p_{n,2,t} = \frac{r+1}{N_t r + 1} \approx \frac{r+1}{N_t r} \quad (23)$$

164 Fraser and Li (2017) calculated the effective population size $N_e(t)$ as a function of the actual population
 165 size $N(t)$ and the mean and variance of the offspring distribution R and σ^2 :

$$N_e(t) = \frac{N(t)}{\sigma^2/R + R - 1} \quad (24)$$

166 This formula was used to estimate the dispersion parameter from genetic data (Li et al. 2017). In our
 167 notation, this is equivalent to:

$$p_{2,t} = \frac{V_t/R_t + R_t - 1}{N_t R_t} \quad (25)$$

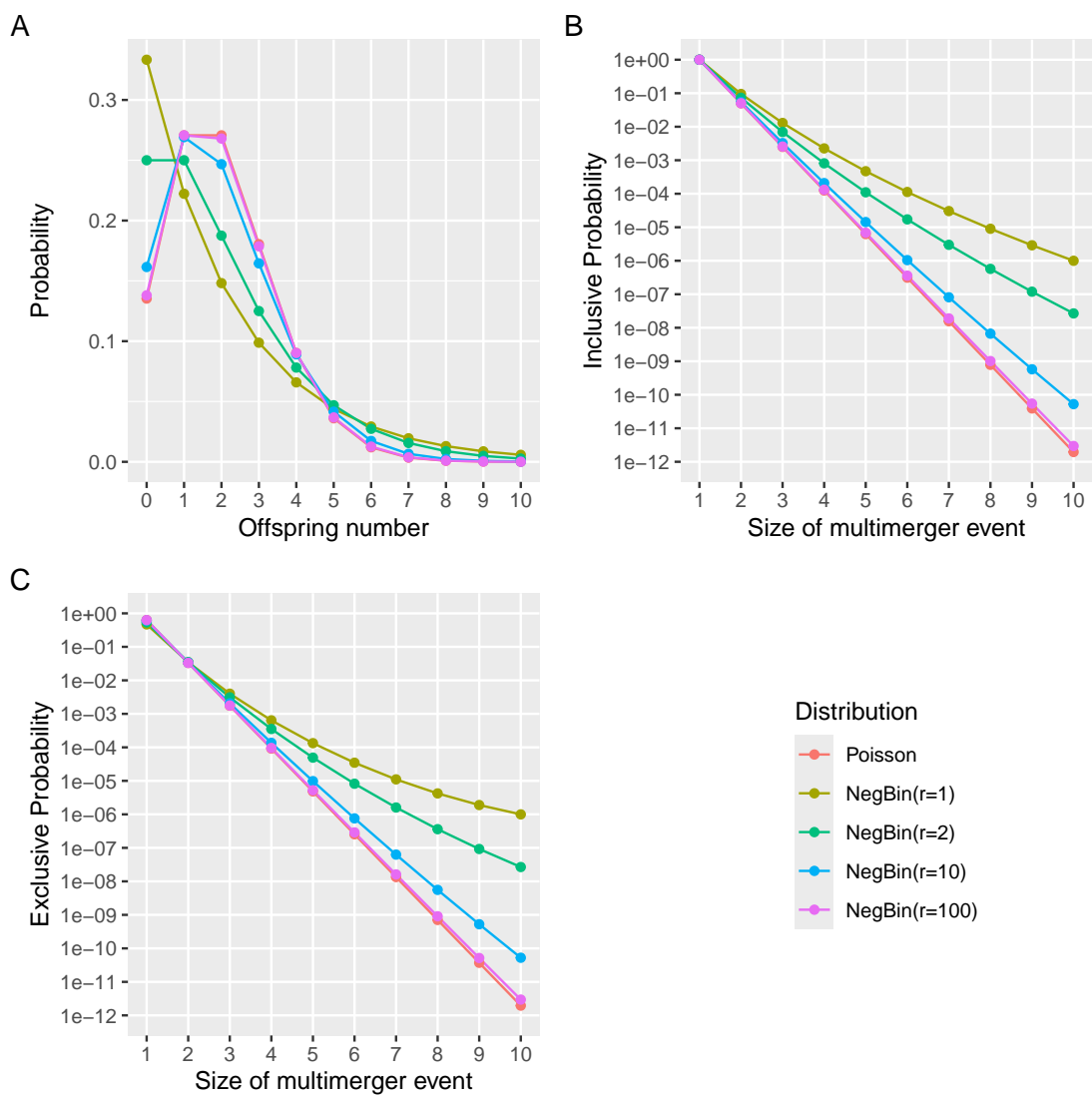


Figure 2: (A) Offspring distribution. (B) Inclusive probability of coalescence. (C) Exclusive probability of coalescence.

168 In the Poisson case we have $V_t = R_t$ so that Equation 25 simplifies to $p_{2,t} = 1/N_t$ which agrees with
169 Equation 22. In the Negative-Binomial case we have $V_t/R_t = 1/p = (r + R_t)/r$ so that Equation
170 25 simplifies to $(r + 1)/(rN_t)$ which agrees with our Equation 23. Conversely, if we substitute
171 $r = R_t^2/(V_t - R_t)$ in Equation 23 we obtain the formula Equation 25.

172 Koelle and Rasmussen (2012) derived the rates of coalescence of two lineages for several epidemiological
173 models, assuming a large population at equilibrium. For each model they use the equation $N_e = N/\sigma^2$
174 to relate the effective population size N_e to the actual population size N and the variance σ^2 in the
175 number of offspring. This relationship was first established by Kingman (1982a) to apply the coalescent
176 model to Cannings exchangeable models (Cannings 1974). From Equation 23 we can take $R_t = 1$ to
177 achieve equilibrium of the population size and $r = R_t^2/(V_t - R_t) = 1/(V_t - 1)$ to deduce the equivalent
178 $p_{2,t} = V_t/N_t$.

179 Volz (2012) showed that the rate of coalescence for two lineages under a continuous-time epidemic
180 coalescent model is $2f(t)/I(t)^2$ where $f(t)$ is the incidence and $I(t)$ the prevalence. Setting in this
181 formula the prevalence as $I(t) = N_{t+1} = N_t R_t$ and the incidence as $f(t) = R_t N_{t+1} = R_t^2 N_t$ we
182 get a coalescent rate of $2/N_t$. To apply the Equation 23 we need to set $r = 1$ so that the offspring
183 distribution is Geometric, which yields the same result.

184 6 Lambda-coalescent

185 The coalescent model (Kingman 1982a,b) describes the ancestry of a sample from a large population
186 evolving according to many forward-in-time models such as the Wright-Fisher model (Wright 1931;
187 Fisher 1930), the Moran model (Moran 1958) and the Cannings exchangeable model (Cannings 1974).
188 Since the coalescent considers a large population in which each individual only has a number of
189 offspring that is small compared to the population size, coalescent trees are always binary and do not
190 feature multimergers, making them unsuitable to represent the ancestry of outbreaks considered in
191 this study. However, the lambda-coalescent models are an extension of the coalescent model that do
192 allow multimergers (Pitman 1999; Sagitov 1999; Donnelly and Kurtz 1999).

193 A lambda-coalescent model is defined by a probability measure $\Lambda(dx)$ on the interval $[0, 1]$, from which
194 we can deduce the rate $\lambda_{n,k}$ at which any subset of k lineages within a set of n observed lineages
195 coalesce:

$$\lambda_{n,k} = \int_0^1 x^{k-2} (1-x)^{n-k} \Lambda(dx) \quad (26)$$

196 The beta-coalescent (Schweinsberg 2003) is a specific type of lambda-coalescent that has been used
 197 recently in several studies analysing genetic data from infectious disease agents (Hoscheit and Pybus
 198 2019; Menardo et al. 2021; Helekal et al. 2024; Zhang and Palacios 2024). The Beta-coalescent model
 199 has a single parameter $\alpha \in [0, 2]$ and is defined as:

$$\Lambda(dx) = \frac{x^{1-\alpha} (1-x)^{\alpha-1}}{B(2-\alpha, \alpha)} dx \quad (27)$$

200 By combining Equations 26 and 27 we can deduce that:

$$\lambda_{n,k} = \frac{B(k-\alpha, n-k+\alpha)}{B(2-\alpha, \alpha)} \quad (28)$$

201 Special cases of the beta-coalescent include $\alpha = 2$ corresponding to the Kingman coalescent, $\alpha = 1$
 202 which is known as the Bolthausen-Sznitman coalescent and $\alpha = 0$ for which the phylogeny is always
 203 star-shaped.

204 We now define our own lambda-coalescent based on the Negative-Binomial case described previously.
 205 For ease of comparison with other coalescent models, we consider that time is continuous and that
 206 the population size remains constant equal to N_t . The exclusive coalescent probability $p_{n,k,t}$ in the
 207 Negative-Binomial case given by Equation 20 can be used to determine the corresponding rate of our
 208 lambda-coalescent, if we consider that the probability of each event in discrete time is the result of the
 209 event happening at a constant rate in continuous time:

$$\lambda_{n,k} = -\log(1 - p_{n,k,t}) \quad (29)$$

210 In order to compare our lambda-coalescent with other models, we consider the distribution of the size
 211 k of the next event among a set of n lineages. For any lambda-coalescent this can be computed as:

$$p(k|n) = \frac{\binom{n}{k} \lambda_{n,k}}{\sum_{i=2}^n \binom{n}{i} \lambda_{n,i}} \quad (30)$$

Figure 3 compares this distribution for $n = 10$ in the Beta-coalescent with parameter $\alpha \in \{0.5, 1, 1.5\}$ and for our lambda-coalescent with parameters $N_t \in \{15, 25, 50\}$ and $r \in \{0.1, 0.5, 1\}$. In the Beta-coalescent, the distribution shifts towards more larger multimerger events as the parameter α decreases. In our new model a wider range of behaviours is obtained when varying the two parameters N_t and r . For a given value of N_t , decreasing the value of r results in more larger events. Conversely, for a given value of r we can see that increasing the value of N_t reduces the probability of larger events.

Genealogies can be simulated from our lambda-coalescent model defined in Equation 29 using the same algorithm as for other lambda-coalescent models (Pitman 1999). Figure 4 shows examples of trees simulated for a sample of size $n = 20$, constant population size $N_t = 40$ and dispersion parameter $r \in \{0.1, 1, 5, 10\}$. It is already clear from these single realisations that the lower values of r result in trees with more larger multimerger events and lower time to the most recent common ancestor, but to quantify these properties we need to consider many trees.

Figure 5 shows summary statistics for 10,000 trees simulated in the same conditions as the individual trees shown in Figure 4. As the dispersion parameter increases from $r = 0.1$ to $r = 10$ multimerger events become less and less likely and large. Simultaneously, the time to the most recent common ancestor increases, as well as the stemminess of the tree (ie the proportion of branch lengths in non-terminal branches).

7 Parameter inference

Consider a genealogy T with n leaves and c coalescent nodes, with $t_0 = 0$ the sampling time, t_1, \dots, t_c the times of the coalescent nodes in increasing order and k_i the number of lineages coalescing at time t_i . The number of lineages existing between time t_{i-1} and t_i is then $n_i = n - \sum_{j=1}^{i-1} k_j$. Under a lambda-coalescent model, the genealogy T has likelihood:

$$p(T|\Lambda) = \prod_{i=1}^c \binom{n_i}{k_i} \lambda_{n_i, k_i} \exp \left(- \sum_{j=2}^{n_i} \binom{n_i}{j} \lambda_{n_i, j} (t_i - t_{i-1}) \right) \quad (31)$$

Estimating the lambda measure in general is a difficult problem (Koskela 2018; Miró Pina et al. 2023). Here however we focus on estimation under our lambda-coalescent model, where the $\lambda_{n, k}$ terms are given by Equation 29. There are therefore two parameters to estimate which have direct

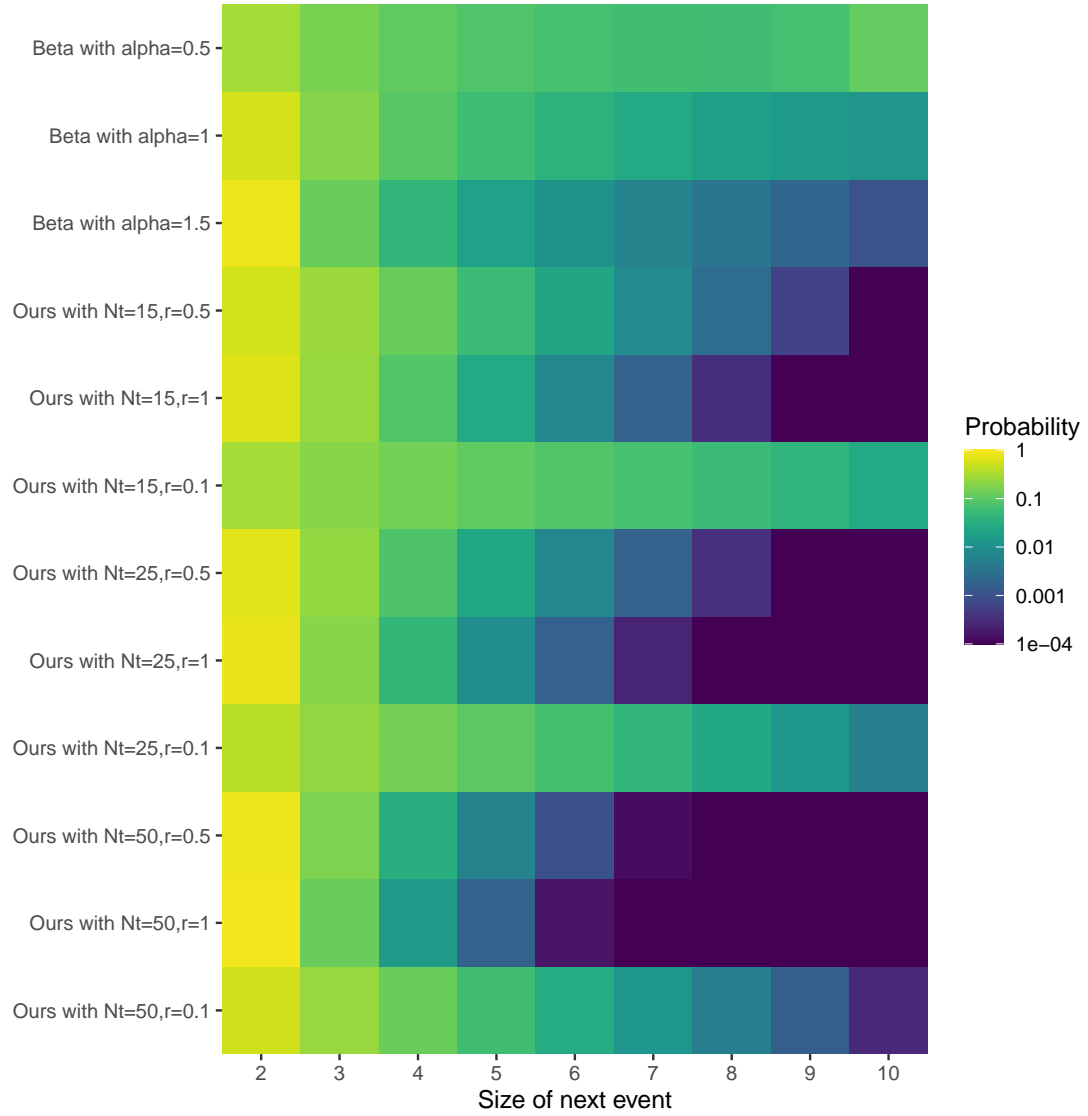


Figure 3: Distribution of the size of the next event among a set of $n = 10$ lineages, compared between the Beta-coalescent and our lambda-coalescent model with various parameters.

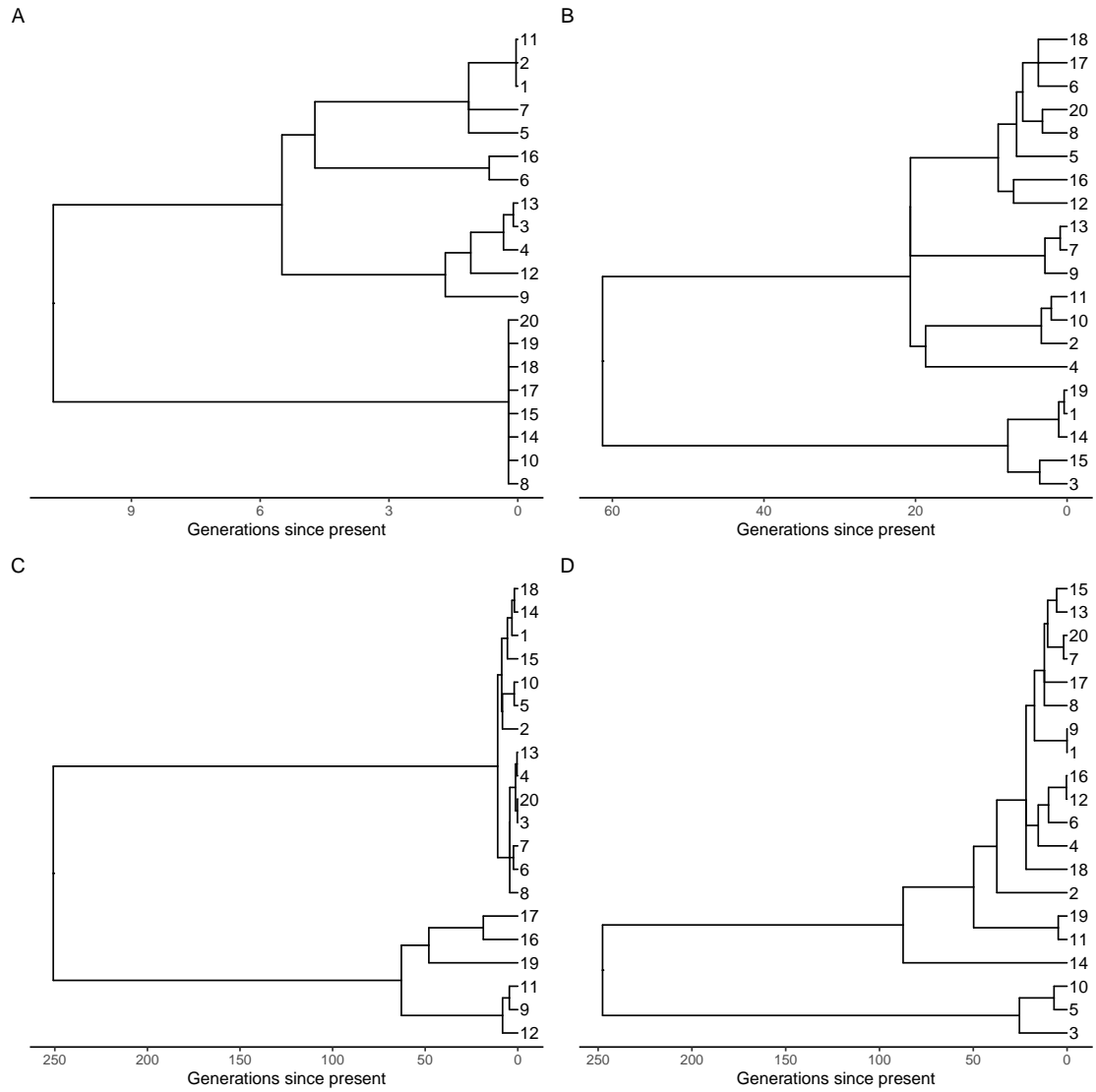


Figure 4: Example of trees simulated under our lambda-coalescent with $r = 0.1$ (A), $r = 1$ (B), $r = 5$ (C) and $r = 10$ (D).

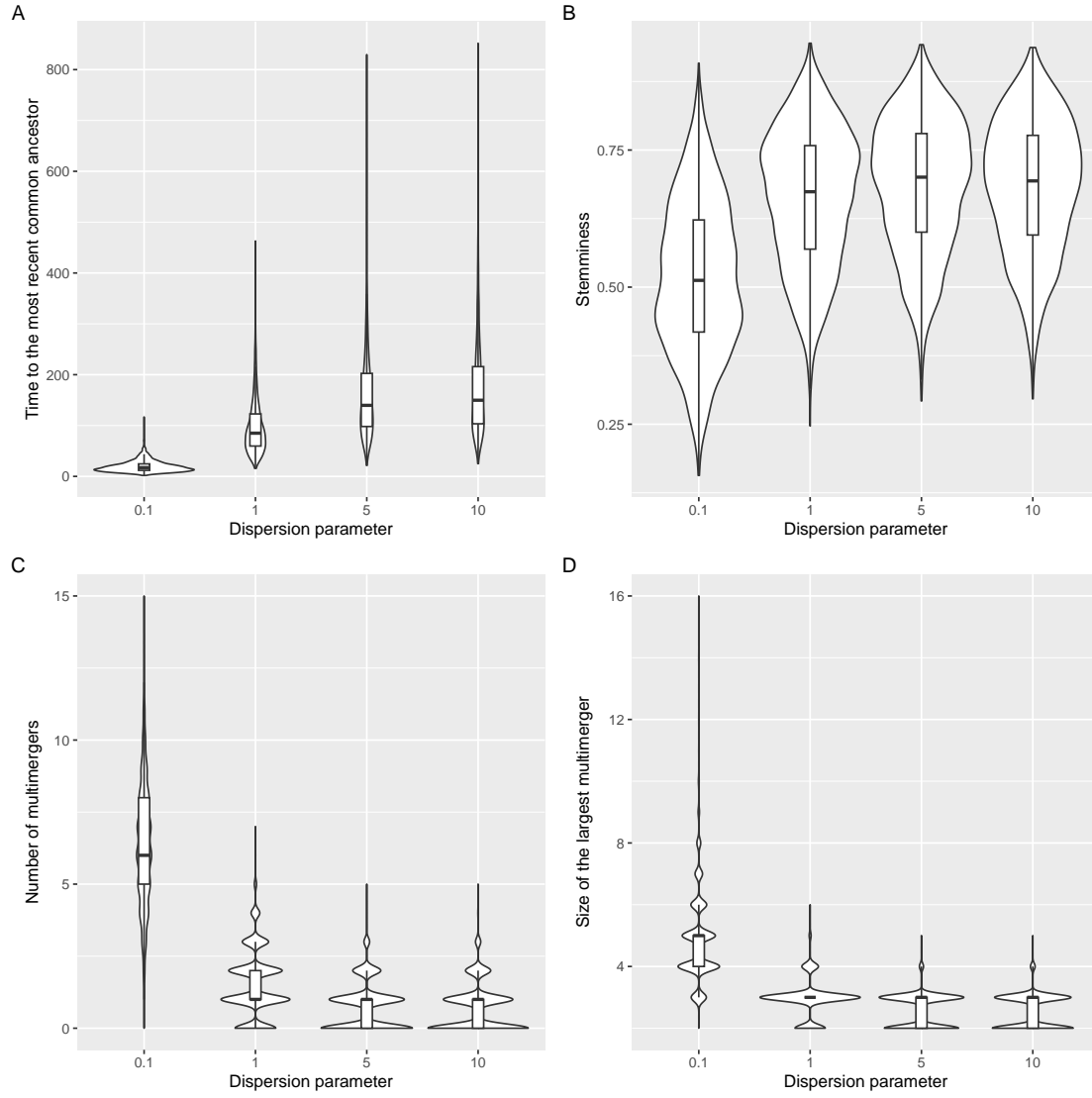


Figure 5: Summary statistics for trees simulated under our lambda-coalescent with $r = 0.1, r = 1, r = 5$ and $r = 10$, namely the time to the most recent common ancestor (A), stemminess (B), number of multimerers (C) and the size of the largest multimerger (D).

and important biological meaning: the effective population size N_t (which remains constant) and the dispersion parameter r of the Negative-Binomial offspring distribution. We perform estimation simply by maximising the likelihood in Equation 31, using the Brent algorithm (Brent 1971) when estimating a single parameter and the L-BFGS-B algorithm when (Byrd et al. 1995) estimating both parameters.

We simulated 100 genealogies from our lambda-coalescent model each of which had $n = 100$ leaves, with parameter N_e drawn uniformly at random between 100 and 500 and parameter r drawn uniformly at random between 0.01 and 2. If we assume knowledge of the dispersion parameter, then estimating the population size works really well (Figure 6A). Conversely we obtain good result when estimating the dispersion parameter given a known population size (Figure 6B). However, attempting to estimate both parameters at the same time performed significantly less well (Figures 6C and D). To illustrate the cause of this, we consider a simulation for which the true N_t was 200 and the true r was 0.5, and we construct the likelihood surface (Figure 6E). This shows a strong inverse tradeoff between the two parameters, which explains why one can be estimated given the other, but not jointly.

8 Implementation

We implemented the analytical methods described in this paper in a new R package entitled *EpiLambda* which is available at <https://github.com/xavierdidelot/EpiLambda> for R version 3.5 or later. All code and data needed to replicate the results are included in the “run” directory of the *EpiLambda* repository. The R package **ape** was used to store, manipulate and visualise phylogenetic trees (Paradis and Schliep 2019).

9 Discussion

Our lambda-coalescent could be defined in a varying population size following the same approach as previously described for the coalescent (Griffiths and Tavaré 1994) and the beta-coalescent (Hoscheit and Pybus 2019; Zhang and Palacios 2024). Could also extend to temporally offset leaves following work on the coalescent (Drummond et al. 2003) and the beta-coalescent (Hoscheit and Pybus 2019).

The Xi-coalescent models admit multiple simultaneous mergers (Schweinsberg 2000).

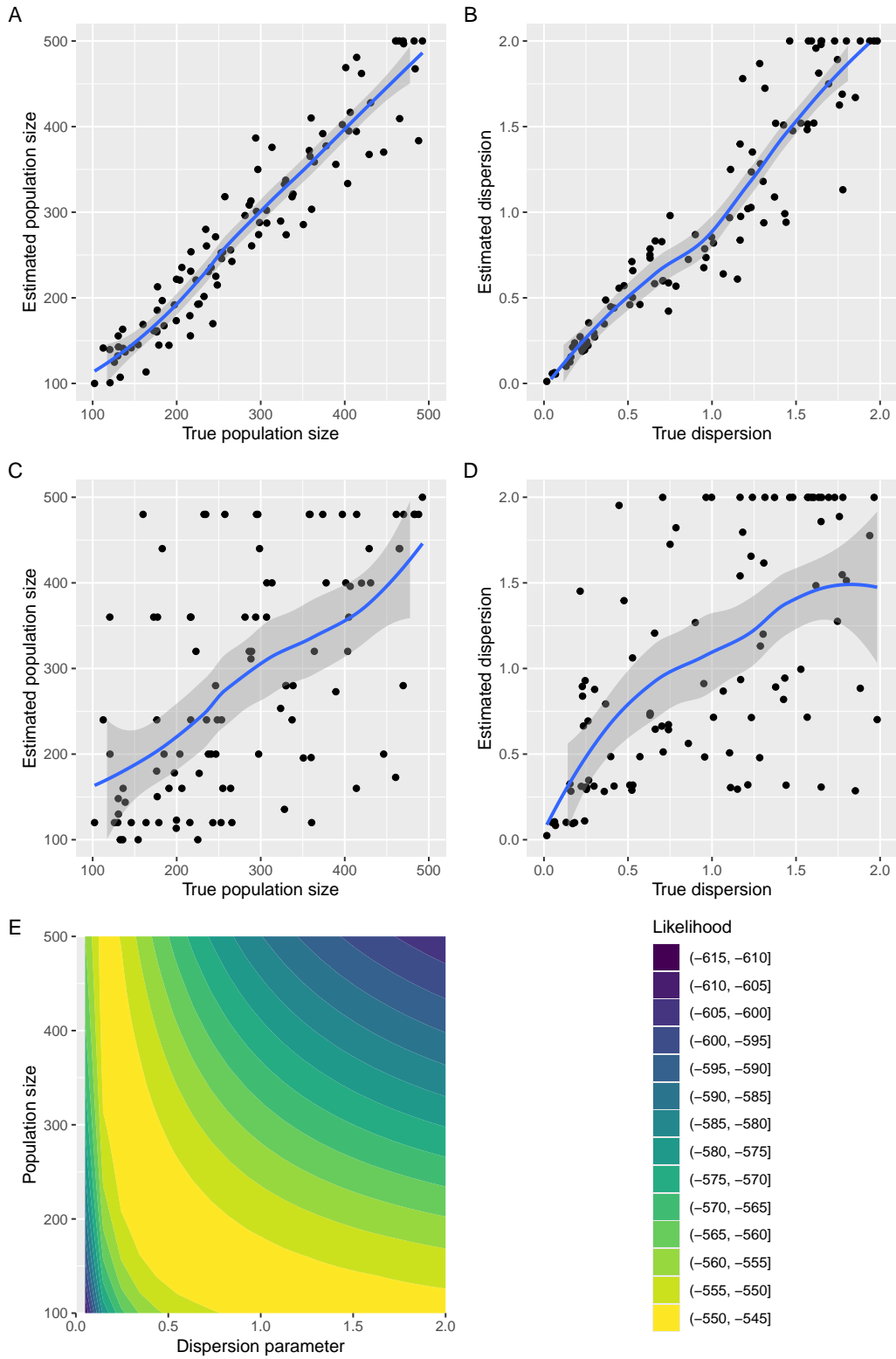


Figure 6: Maximum likelihood estimation of parameters. (A) Estimation of the population size given the dispersion parameter. (B) Estimation of the dispersion parameter given the population size. (C and D) Joint estimation of both the population size and dispersion parameters. (E) Example of likelihood surface as a function of both parameters.

262 Difference between transmission tree and phylogenetic tree (Jombart et al. 2011). Modelling within-
263 host evolution to bridge the gap (Didelot et al. 2014; Hall et al. 2015; Didelot et al. 2017).
264 Superspreading individuals vs superspreading events (Riley et al. 2003; Wallinga and Teunis 2004;
265 Ho et al. 2023).

266 **Acknowledgements**

267 We acknowledge funding from the National Institute for Health Research (NIHR) Health Protection
268 Research Unit in Genomics and Enabling Data.

References

- Anderson, R.M., May, R.M., 1991. Infectious Diseases of Humans: Dynamics and Control. Oxford University Press, USA.
- Brent, R.P., 1971. An algorithm with guaranteed convergence for finding a zero of a function. The computer journal 14, 422–425.
- Byrd, R.H., Lu, P., Nocedal, J., Zhu, C., 1995. A limited memory algorithm for bound constrained optimization. SIAM Journal on scientific computing 16, 1190–1208.
- Cannings, C., 1974. The latent roots of certain Markov chains arising in genetics: a new approach, I. Haploid models. Adv. Appl. Probab. 6, 260–290. doi:10.2307/1426293.
- Didelot, X., Fraser, C., Gardy, J., Colijn, C., 2017. Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks. Molecular Biology and Evolution 34, 997–1007. doi:10.1093/molbev/msw275.
- Didelot, X., Gardy, J., Colijn, C., 2014. Bayesian inference of infectious disease transmission from whole genome sequence data. Molecular Biology and Evolution 31, 1869–1879. doi:10.1093/molbev/msu121.
- Donnelly, P., Kurtz, T.G., 1999. Particle Representations for Measure-Valued Population Models. The Annals of Probability 27. doi:10.1214/aop/1022677258.
- Drummond, A.J., Pybus, O.G., Rambaut, A., Forsberg, R., Rodrigo, A.G., 2003. Measurably evolving populations. Trends in Ecology and Evolution 18, 481–488. doi:10.1016/S0169-5347(03)00216-7.
- Du, Z., Wang, C., Liu, C., Bai, Y., Pei, S., Adam, D.C., Wang, L., Wu, P., Lau, E.H.Y., Cowling, B.J., 2022. Systematic review and meta-analyses of superspreading of SARS-CoV-2 infections. Transboundary and Emerging Diseases 69. doi:10.1111/tbed.14655.
- Ferguson, N.M., Cummings, D.A.T., Fraser, C., Cajka, J.C., Cooley, P.C., Burke, D.S., 2006. Strategies for mitigating an influenza pandemic. Nature 442, 448–452. doi:10.1038/nature04795.
- Fisher, R.A., 1930. The genetical theory of natural selection. Clarendon Press. doi:10.5962/bhl.title.27468.
- Fraser, C., Li, L.M., 2017. Coalescent models for populations with time-varying population sizes and arbitrary offspring distributions. bioRxiv , 10.1101/131730doi:10.1101/131730.

Fraser, C., Riley, S., Anderson, R.M., Ferguson, N.M., 2004. Factors that make an infectious disease outbreak controllable. *Proceedings of the National Academy of Sciences* 101, 6146–6151. doi:10.1073/pnas.0307506101.

Gómez-Carballa, A., Pardo-Seco, J., Bello, X., Martínón-Torres, F., Salas, A., 2021. Superspreading in the emergence of COVID-19 variants. *Trends in Genetics* 37, 1069–1080. doi:10.1016/j.tig.2021.09.003.

Grassly, N.C., Fraser, C., 2008. Mathematical models of infectious disease transmission. *Nature Reviews Microbiology* 6, 477–87. doi:10.1038/nrmicro1845.

Griffiths, R.C., Tavaré, S., 1994. Sampling theory for neutral alleles in a varying environment. *Philosophical Transactions of the Royal Society B* 344, 403–410.

Hall, M., Woolhouse, M., Rambaut, A., 2015. Epidemic Reconstruction in a Phylogenetics Framework: Transmission Trees as Partitions of the Node Set. *PLOS Computational Biology* 11, e1004613. doi:10.1371/journal.pcbi.1004613.

Helekal, D., Koskela, J., Didelot, X., 2024. Inference of multiple mergers while dating a pathogen phylogeny. *bioRxiv* , 2023.09.12.557403doi:10.1101/2023.09.12.557403.

Ho, F., Parag, K.V., Adam, D.C., Lau, E.H.Y., Cowling, B.J., Tsang, T.K., 2023. Accounting for the Potential of Overdispersion in Estimation of the Time-varying Reproduction Number. *Epidemiology* 34, 201–205. doi:10.1097/EDE.0000000000001563.

Hoscheit, P., Pybus, O.G., 2019. The multifurcating skyline plot. *Virus Evolution* 5, 1–10. doi:10.1093/ve/vez031.

Jombart, T., Eggo, R.M., Dodd, P.J., Balloux, F., 2011. Reconstructing disease outbreaks from genetic data: A graph approach. *Heredity* 106, 383–90. doi:10.1038/hdy.2010.78.

Keeling, M.J., Rohani, P., 2008. Modeling infectious diseases in humans and animals. Princeton university press.

Kingman, J., 1982a. The coalescent. *Stochastic Processes and their Applications* 13, 235–248. doi:10.1016/0304-4149(82)90011-4.

Kingman, J.F.C., 1982b. On the genealogy of large populations. *Journal of Applied Probability* 19, 27–43. doi:10.2307/3213548.

Koelle, K., Rasmussen, D.A., 2012. Rates of coalescence for common epidemiological models at equilibrium. *Journal of The Royal Society Interface* 9, 997–1007. doi:10.1098/rsif.2011.0495.

327 Koskela, J., 2018. Multi-locus data distinguishes between population growth and multiple merger
328 coalescents. *Statistical Applications in Genetics and Molecular Biology* 17, 1–24. doi:10.1515/
329 **sagmb-2017-0011**.

330 Kucharski, A.J., Althaus, C.L., 2015. The role of superspreading in Middle East respiratory syndrome
331 coronavirus (MERS-CoV) transmission. *Eurosurveillance* 20. doi:10.2807/1560-7917.ES2015.20.
332 **25.21167**.

333 Lemieux, J.E., Siddle, K.J., Shaw, B.M., Loreth, C., Schaffner, S.F., Gladden-Young, A., Adams,
334 G., Fink, T., Tomkins-Tinch, C.H., Krasilnikova, L.A., DeRuff, K.C., Rudy, M., Bauer, M.R.,
335 Lagerborg, K.A., Normandin, E., Chapman, S.B., Reilly, S.K., Anahtar, M.N., Lin, A.E., Carter,
336 A., Myhrvold, C., Kembal, M.E., Chaluvadi, S., Cusick, C., Flowers, K., Neumann, A., Cerrato,
337 F., Farhat, M., Slater, D., Harris, J.B., Branda, J.A., Hooper, D., Gaeta, J.M., Baggett, T.P.,
338 O’Connell, J., Gnirke, A., Lieberman, T.D., Philippakis, A., Burns, M., Brown, C.M., Luban, J.,
339 Ryan, E.T., Turbett, S.E., LaRocque, R.C., Hanage, W.P., Gallagher, G.R., Madoff, L.C., Smole, S.,
340 Pierce, V.M., Rosenberg, E., Sabeti, P.C., Park, D.J., MacInnis, B.L., 2021. Phylogenetic analysis
341 of SARS-CoV-2 in Boston highlights the impact of superspreading events. *Science* 371, eabe3261.
342 doi:10.1126/science.abe3261.

343 Li, L.M., Grassly, N.C., Fraser, C., 2017. Quantifying Transmission Heterogeneity Using Both
344 Pathogen Phylogenies and Incidence Time Series. *Molecular Biology and Evolution* 34, 2982–2995.
345 doi:10.1093/molbev/msx195.

346 Lloyd-Smith, J., Schreiber, S., Kopp, P., Getz, W., 2005. Superspreading and the effect of individual
347 variation on disease emergence. *Nature* 438, 355–9. doi:10.1038/nature04153.

348 Menardo, F., Gagneux, S., Freund, F., 2021. Multiple Merger Genealogies in Outbreaks of
349 *Mycobacterium tuberculosis*. *Molecular Biology and Evolution* 38, 290–306. doi:10.1093/molbev/
350 **msaa179**.

351 Miró Pina, V., Joly, É., Siri-Jégousse, A., 2023. Estimating the Lambda measure in multiple-merger
352 coalescents. *Theoretical Population Biology* 154, 94–101. doi:10.1016/j.tpb.2023.09.002.

353 Moran, P., 1958. Random Processes in Genetics. *Mathematical Proceedings of the Cambridge*
354 *Philosophical Society* 54, 60–71.

355 Paradis, E., Schliep, K., 2019. Ape 5.0: An environment for modern phylogenetics and evolutionary
356 analyses in R. *Bioinformatics* 35, 526–528. doi:10.1093/bioinformatics/bty633.

357 Pitman, J., 1999. Coalescents with multiple collisions. *The Annals of Probability* 27, 1870–1902.

358 Potts, R.B., 1953. Note on the Factorial Moments of Standard Distributions. Australian Journal of
359 Physics 6, 498–499. URL: <https://www.publish.csiro.au/ph/ph530498>, doi:10.1071/ph530498.
360 publisher: CSIRO PUBLISHING.

361 Riley, S., Fraser, C., a Donnelly, C., Ghani, A.C., Abu-Raddad, L.J., Hedley, A.J., Leung, G.M.,
362 Ho, L.M., Lam, T.H., Thach, T.Q., Chau, P., Chan, K.P., Lo, S.V., Leung, P.Y., Tsang, T., Ho,
363 W., Lee, K.H., Lau, E.M.C., Ferguson, N.M., Anderson, R.M., 2003. Transmission dynamics of the
364 etiological agent of SARS in Hong Kong: Impact of public health interventions. Science 300, 1961–6.
365 doi:10.1126/science.1086478.

366 Sagitov, S., 1999. The general coalescent with asynchronous mergers of ancestral lines. Journal of
367 Applied Probability 36, 1116–1125. doi:10.1239/jap/1032374759.

368 Schweinsberg, J., 2000. Coalescents with Simultaneous Multiple Collisions. Electronic Journal of
369 Probability 5. doi:10.1214/EJP.v5-68.

370 Schweinsberg, J., 2003. Coalescent processes obtained from supercritical Galton–Watson processes.
371 Stochastic Processes and their Applications 106, 107–139. doi:10.1016/S0304-4149(03)00028-0.

372 Stein, R.A., 2011. Super-spreaders in infectious diseases. International Journal of Infectious Diseases
373 15, e510–e513. doi:10.1016/j.ijid.2010.06.020.

374 Tripathi, R.C., Gupta, R.C., Gurland, J., 1994. Estimation of parameters in the beta binomial model.
375 Annals of the Institute of Statistical Mathematics 46, 317–331. URL: [https://doi.org/10.1007/](https://doi.org/10.1007/BF01720588)
376 BF01720588, doi:10.1007/BF01720588.

377 Volz, E.M., 2012. Complex population dynamics and the coalescent under neutrality. Genetics 190,
378 187–201. doi:10.1534/genetics.111.134627.

379 Wallinga, J., Teunis, P., 2004. Different Epidemic Curves for Severe Acute Respiratory Syndrome
380 Reveal Similar Impacts of Control Measures. American Journal of Epidemiology 160, 509–516.

381 Wang, J., Chen, X., Guo, Z., Zhao, S., Huang, Z., Zhuang, Z., Wong, E.L.y., Zee, B.C.Y., Chong,
382 M.K.C., Wang, M.H., Yeoh, E.K., 2021. Superspreading and heterogeneity in transmission of SARS,
383 MERS, and COVID-19: A systematic review. Computational and Structural Biotechnology Journal
384 19, 5039–5046. doi:10.1016/j.csbj.2021.08.045.

385 Wang, L., Didelot, X., Yang, J., Wong, G., Shi, Y., Liu, W., Gao, G.F., Bi, Y., 2020. Inference of
386 person-to-person transmission of COVID-19 reveals hidden super-spreading events during the early
387 outbreak phase. Nature Communications 11, 5006. doi:10.1038/s41467-020-18836-4.

388 Woolhouse, M.E.J., Dye, C., Etard, J.F., Smith, T., Charlwood, J.D., Garnett, G.P., Hagan, P., Hii,
389 J.L.K., Ndhlovu, P.D., Quinnell, R.J., Watts, C.H., Chandiwana, S.K., Anderson, R.M., 1997.
390 Heterogeneities in the transmission of infectious agents: Implications for the design of control
391 programs. *Proceedings of the National Academy of Sciences* 94, 338–342. doi:10.1073/pnas.94.1.
392 338.

393 Wright, S., 1931. Evolution in Mendelian populations. *Genetics* 16, 97–159. doi:10.1093/genetics/
394 16.2.97.

395 Zhang, J., Palacios, J.A., 2024. Multiple merger coalescent inference of effective population size. *arXiv*
396 , 2407.14976doi:10.48550/arXiv.2407.14976, arXiv:2407.14976.