

1 Supplementary material for “Genomic infectious disease
2 epidemiology in partially sampled and ongoing outbreaks”

3 Xavier Didelot¹, Christophe Fraser^{1,2}, Jennifer Gardy^{3,4}, Caroline Colijn⁵

4 **1** Department of Infectious Disease Epidemiology, Imperial College London, Norfolk Place,
5 London, W2 1PG, United Kingdom

6
7 **2** Oxford Big Data Institute, Li Ka Shing Centre for Health Information and Discovery,
8 Nuffield Department of Medicine, University of Oxford, Oxford OX3 7BN, United Kingdom

9
10 **3** Communicable Disease Prevention and Control Services, British Columbia Centre for
11 Disease Control, Vancouver, British Columbia, Canada

12
13 **4** School of Population and Public Health, University of British Columbia, Vancouver, British
14 Columbia, Canada

15
16 **5** Department of Mathematics, Imperial College, London SW7 2AZ, UK

17 **Numerical calculation of ω_***

18 In the finished outbreak scenario, the probability of being excluded is given in Equation 2. We
 19 calculate ω_* numerically by solving Equation 2 using the R command uniroot over the interval
 20 $[0,1]$.

21 **Numerical approach to $\bar{\omega}_t$**

22 In the ongoing outbreak scenario, let ω_t be the probability of being unsampled and having all
 23 descendants unsampled, conditional on being infected at time t . As $t \rightarrow -\infty$, $\omega_t \rightarrow \omega_*$ which is
 24 the solution to Equation 2 described above. However, as the study ends at a finite time T , we
 25 know that $\omega_T = 1$ because a case infected at T will be excluded.

26 We do not know the times of an individual's descendants, but we still condition on the total
 27 number of descendants. Integrating those out, we have

$$\omega_t = (1 - \pi_t) \sum_{k=0}^{\infty} \alpha(k) \prod_{j=1}^k \left[\int_t^{\infty} \gamma(\tau_j - t) \omega_{\tau_j} d\tau_j \right] \quad (\text{S1})$$

28 Let the term in square brackets be $\bar{\omega}_t$. We need to determine what this function is because
 29 ultimately, we need it to compute the probability of having k sampled descendants.

30 We have (with g the probability generating function of the negative binomial offspring
 31 distribution, and p and r its parameters)

$$\omega_t = (1 - \pi_t) g(\bar{\omega}_t) = (1 - \pi_t) \left(\frac{1 - p}{1 - p\bar{\omega}_t} \right)^r. \quad (\text{S2})$$

32 We have $\omega_t = 1$ for $t \geq T$, so

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

33 We substitute this into Equation S2. We use the trapezoid method for the first term, and the
 34 second term we can compute explicitly: $\int_{T-t}^{\infty} \gamma(u) du \equiv F(t)$.

35 Let $t_i = T - i\Delta t$.

36 The trapezoid method gives:

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

37 where $c_i = 1$ unless $i = 0$ or $i = k$, where $c_i = 1/2$. The k 'th term drops out because $\gamma(0) = 0$
 38 by assumption, so:

$$\omega(t_k) \approx (1 - \pi_t) \left(\frac{1 - p}{1 - pF(t) - p \sum_{i=0}^{k-1} c_i \gamma((k - i)\Delta t) \omega(t_i) \Delta t} \right)^r \quad (\text{S3})$$

39 This is straightforward to compute with iteration. We should find that $\omega_t \rightarrow \omega_*$ as $t \rightarrow -\infty$.

40 This gives the probability of being excluded and having no sampled descendants, conditional
 41 on having been infected at time t . We now follow the same approach as above (ie as in the
 42 asymptotic case), and condition on the total number of descendants. We need to compute the
 43 "modified offspring function" of Equation 10:

$$\alpha_t(d) = \sum_{k=d}^{\infty} \binom{k}{d} \alpha(k) \bar{\omega}_t^{k-d}$$

44 Since k follows a negative binomial distribution, we have

$$\alpha_t(d) = \sum_{k=d}^{\infty} \binom{k+r-1}{r-1} \binom{k}{d} p^k (1-p)^r \bar{\omega}_t^{k-d} p_s(d)$$

45 which we can compute without difficulty (typically d is small).

46 MCMC results for the tuberculosis outbreak

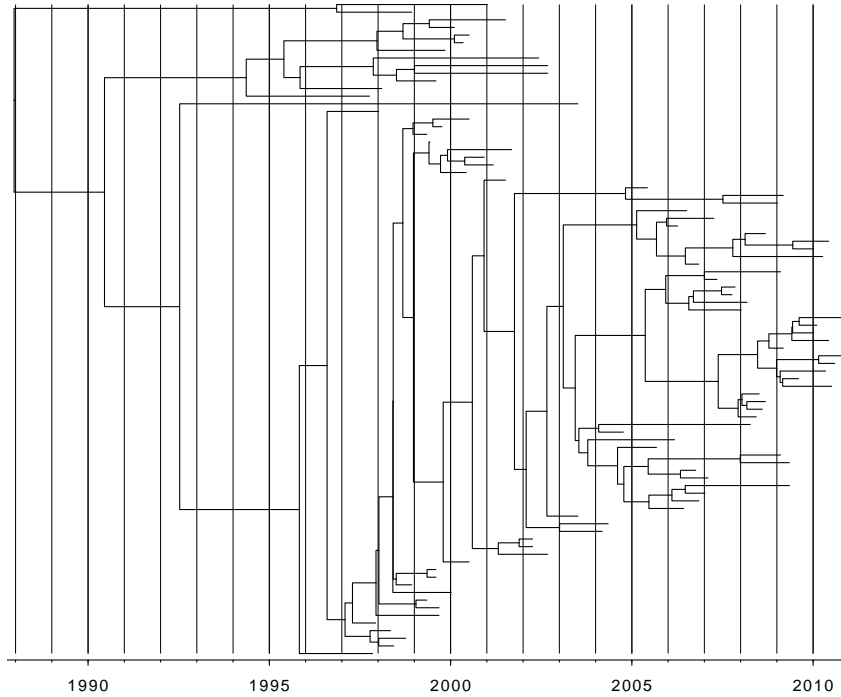


Figure S1. BEAST tree

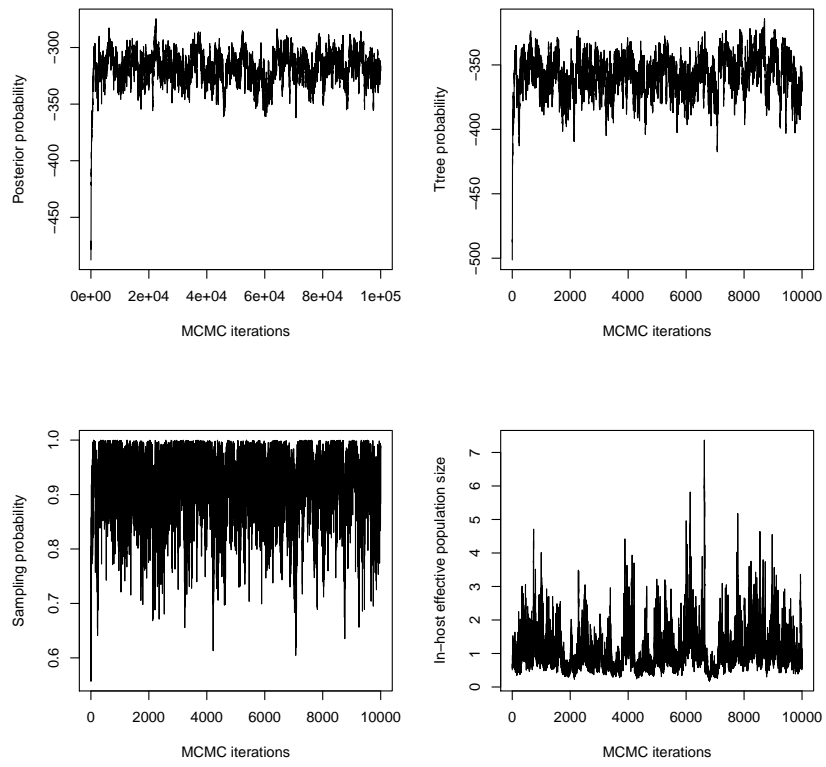


Figure S2. MCMC traces

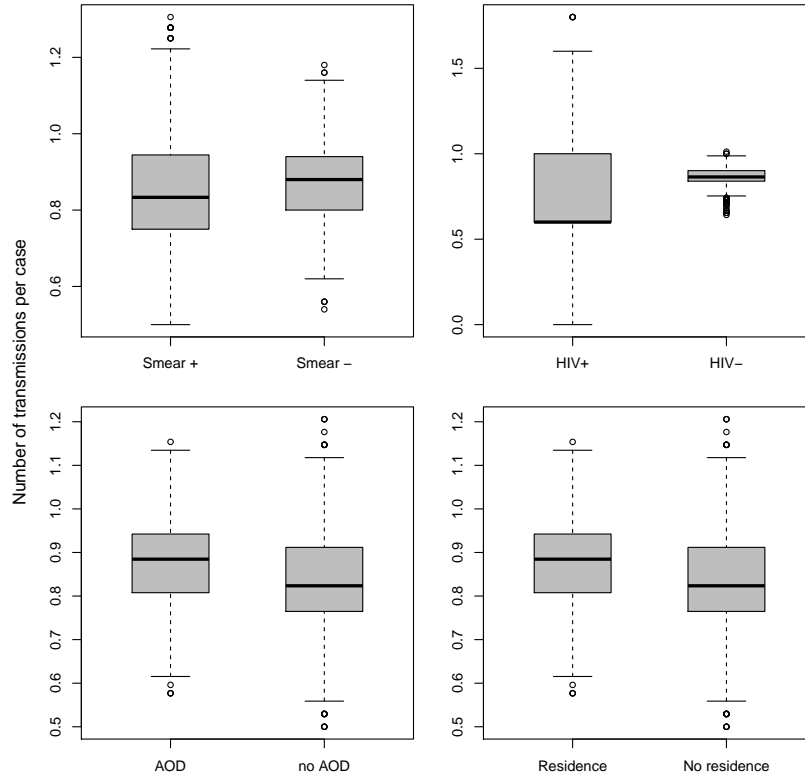


Figure S3. How direct transmissions among sampled cases are affected by potential risk factors. Among the inferred transmission events between sampled individuals, we computed the number of transmissions by (eg) smear-positive and smear-negative individuals divided by the number of smear-positive and -negative individuals to obtain a per-individual average number of transmissions. We computed this for 1000 samples from the posterior MCMC chain (chosen uniformly at random from the latter half of the posterior).