

Inferring the epidemiologic sources of infection from cross-sectionally sampled pathogen sequence data

Xiaoyue Xi

Department of Mathematics, Imperial College London, London SW72AZ, United Kingdom.

Simon EF Spencer

Department of Statistics, University of Warwick, Coventry CV47AL, United Kingdom.

M Kate Grabowski

Department of Medicine, Johns Hopkins School of Medicine, Baltimore, MD, USA; Rakai Health Sciences Program, Entebbe, Uganda.

Oliver Ratmann

Department of Mathematics, Imperial College London, London SW72AZ, United Kingdom.

E-mail: oliver.ratmann@imperial.ac.uk

on behalf of Rakai Health Sciences Program and PANGEA-HIV

Summary.

1. Introduction

1.1. The main burden of HIV transmission is in sub-Saharan Africa

HIV remains one of the largest public health threats, especially in sub-Saharan Africa where approximately 66% of all new cases worldwide occur (UNAIDS, 2018). In most countries, such as Uganda or Zambia, the epidemic spreads among sexually active men and women in general, reaching an HIV prevalence of 40% in particular areas such as fishing communities on Lake Victoria (Chang et al., 2016).

A large number of HIV interventions are now available to stop viral spread effectively, such as voluntary medical male circumcision (VMMC) to reduce the risk of HIV acquisition in men, or immediate provision of antiretroviral therapy (ART) to suppress the virus in infected individuals and thereby stop onward transmission (Grabowski et al.,

2 on behalf of Rakai Health Sciences Program and PANGEA-HIV

13 2017; Kagaayi et al., 2019). Since the widespread adoption of such prevention measures,
14 rates of incident cases have overall dropped considerably, although they remain too high
15 for HIV elimination (UNAIDS, 2014b).

16 1.2. *Inferring the sources, sinks and hubs of transmission flows to aid the design of* 17 *HIV prevention interventions*

18 In this context, there is increasing focus on identifying groups of individuals that are
19 at high risk of acquiring HIV and at high risk of spreading the virus, and then to
20 target tailored control interventions to these groups (UNAIDS, 2014a). Conceptually,
21 the first step in this strategy is to break down the epidemic into source, sink and hub
22 populations, according to the transmission flows that occur between them (figure 1).
23 Sources are population groups that disproportionately pass on infection, sinks are groups
24 that disproportionately acquire infection, and hubs are both sources and sinks (Abeler-
25 Dörner et al., 2019). The population groups can be defined in various ways.

26 For instance, Dwyer-Lindgren et al. (2019) provided sub-national estimates of HIV
27 prevalence across Africa, adding to data showing that the epidemic is highly heteroge-
28 neous across Africa, with small areas of very high prevalence (hotspots) that are sur-
29 rounded by neighbouring areas with substantially lower prevalence. The WHO and oth-
30 ers are recommending to target interventions to hotspots to maximise cost-effectiveness
31 of interventions and reduce infection burden in hotspots (UNAIDS, 2014a; The Office
32 of the US Global AIDS Coordinator, 2014). Although often assumed, it is not clear if
33 hotspots are also sources of epidemic spread to the neighbouring lower-prevalence com-
34 munities (Ratmann et al., 2020). If they are, spatial targeting could have a substantial
35 knock-on effect in reducing infection in the broader population (Aral et al., 2015). In
36 this application, interest centres on dividing the full population into individuals living
37 in high-prevalence areas (h) and low-prevalence areas (l), and estimate the transmission
38 flows within and between them while accounting for transmission flow from/to external
39 locations (e),

$$\boldsymbol{\pi} = \begin{pmatrix} \pi_{hh} & \pi_{hl} \\ \pi_{lh} & \pi_{ll} \end{pmatrix}, \quad (1)$$

40 where π_{ab} denotes the proportion of transmission flows from group a to group b subject

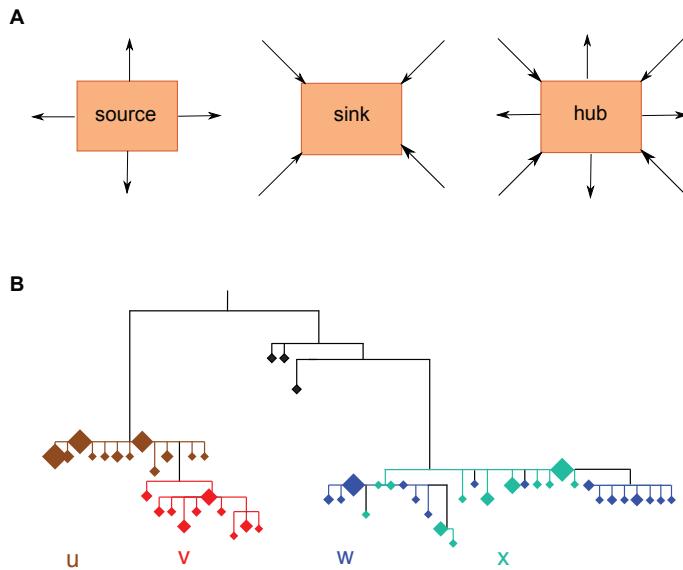


Fig. 1. Analysis aims and sketch of deep-sequence phylogenetic data to address these aims. (A) The overall aim of phylogenetic source attribution analyses is to infer how pathogens are passed on between population groups. Conceptually populations can be divided into source populations that predominantly transmit disease, sink populations that predominantly receive infection, and hub populations that both disproportionately transmit and receive infections. (B) Viral deep-sequencing generates many sequence samples per host, which can be used to establish phylogenetic orderings between individuals, and thereby estimate the direction of pathogen spread among sampled individuals. The figure sketches a deep-sequence phylogeny of pathogen sequences from individuals u , v , w , and x . Each tip (diamonds) represents a unique sequence, and the size of the tip copy number. Black tips correspond to out-of-sample reference sequences. Phylogenetic lineages are attributed to individuals (colours) using ancestral state reconstruction. Black lineages cannot be attributed to individuals. The subgraph of the tree associated with individual u is ancestral to that of v , suggesting that infection spread from u to v potentially via unsampled intermediates. Individual w has five subgraphs, some of which are ancestral to those of x and some of which are descended from those of x , indicating a complex ordering from which the direction of infection spread cannot be inferred. Subgraphs of v and w are not phylogenetically adjacent (disconnected), suggesting that one did not infect the other. With such information from a population-based sample of infected individuals, it is possible to quantify population-level transmission flows, sources, sinks, and hubs.

41 to $\sum_{ab} \pi_{ab} = 1$.

42 Another prominent application of this framework concerns the interruption of infec-
 43 tion cycles between men and women of different ages. De Oliveira et al. (2017) proposed
 44 the scenario that young women aged <25 years are predominantly infected by older men
 45 aged 25–40 years, and later spread the virus to similarly aged men in their late twenties
 46 and early thirties. If true, combination prevention interventions that seek to reduce
 47 HIV acquisition in young women from older men, and seek to reduce transmission from

4 on behalf of Rakai Health Sciences Program and PANGEA-HIV

48 older men and women could break the transmission cycle (De Oliveira et al., 2017). In
 49 this application, the full population is divided by sex-specific age bands (we consider
 50 1-year age bands between 15 to 49 years), and the main objective is to estimate the
 51 transmission flow vector between men and women across age groups,

$$\boldsymbol{\pi} = \begin{pmatrix} \boldsymbol{\pi}^{mf} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\pi}^{fm} \end{pmatrix}, \quad \boldsymbol{\pi}^{mf} = \begin{pmatrix} \pi_{11}^{mf} & \dots & \pi_{1K}^{mf} \\ \vdots & \ddots & \vdots \\ \pi_{K1}^{mf} & \dots & \pi_{KK}^{mf} \end{pmatrix} \quad \boldsymbol{\pi}^{fm} = \begin{pmatrix} \pi_{11}^{fm} & \dots & \pi_{1K}^{fm} \\ \vdots & \ddots & \vdots \\ \pi_{K1}^{fm} & \dots & \pi_{KK}^{fm} \end{pmatrix}, \quad (2)$$

where π_{ab}^{mf} denotes the proportion of transmissions from men in age band a to women in age band b , and similarly for π_{ab}^{fm} . The flow matrix (2) has $2K^2$ non-zero entries to estimate, which for 1-year age bands amounts to 2450 variables. Important summary statistics are the vector of sources of infection in group b individuals ($\boldsymbol{\delta}_b$), for example in young women aged <25 years; the vector of recipients of infection from group a individuals ($\boldsymbol{\omega}^a$), for example from men aged 25-40 years; and flow ratios from a to b (γ_{ab}), for example the ratio of transmissions from high-prevalence to low-prevalence areas compared to transmissions from low-prevalence to high-prevalence areas. Respectively these quantities are defined by

$$\boldsymbol{\delta}^b = (\delta_a^b)_{a=1,\dots,A}, \quad \delta_a^b = \pi_{ab} / \sum_c \pi_{cb}; \quad (3a)$$

$$\boldsymbol{\omega}^a = (\omega_b^a)_{b=1,\dots,A}, \quad \omega_b^a = \pi_{ab} / \sum_c \pi_{ac}; \quad (3b)$$

$$\gamma_{ab} = \pi_{ab} / \pi_{ba}. \quad (3c)$$

52 1.3. Inference from pathogen sequence data

53 Here, we are concerned in estimating the quantities (1– 3) from a population-based
 54 sample of pathogen sequences, obtained from a sample of infected individuals of a cross-
 55 sectionally surveyed study population. The statistical model that we describe is applic-
 56 able whenever the sampling frame is relatively short compared to the infectious disease
 57 dynamics, so that transmission flows can be considered constant in time.

58 Pathogen sequences are the result of an evolutionary process that occurs in infected
 59 individuals of an epidemic. Thus, for fast-evolving pathogens like HIV, the phylogenetic

60 relationship of pathogen sequences can be used in reverse to estimate transmission dy-
61 namics. These methods have, for example, been used to estimate the origins of HIV
62 (Faria et al., 2014), to quantify the contribution of different disease phases to onward
63 spread (Volz et al., 2013; Ratmann et al., 2016), to improve outbreak detection (Poon
64 et al., 2016), and they have been applied to a large range of other infectious diseases
65 (Uhlemann et al., 2014; Didelot et al., 2017; Dellicour et al., 2018).

66 Traditionally, HIV sequences are obtained through Sanger sequencing, which returns
67 for each sample one consensus nucleotide sequence that captures the entire viral diversity
68 in the sample from one individual. The genetic distance between two consensus sequences
69 can be used to estimate if the corresponding two individuals are epidemiologically closely
70 related (Hu   et al., 2005), however the data are insufficient to estimate the direction of
71 transmission between any two sampled individuals (Leitner and Romero-Severson, 2018).
72 For this reason most available methods propose to infer the flow parameters (1– 3) in-
73 directly from statistics of the entire phylogeny, usually the coalescent times (i.e. the
74 times when two lineages coalesce into one, backwards in time) and the disease states
75 of infected individuals at time of sampling (such as location or age in the two applica-
76 tions discussed above). In the migration model (Lemey et al., 2009), the states of viral
77 lineages at any time are described with a continuous-time Markov chain (CTMC) that
78 is independent of the evolutionary process. Flow estimates can be obtained from the
79 posterior distribution of the transition rates of the CTMC model via MCMC sampling,
80 as well as posterior estimates of the phylogeny and the states of its lineages, which are
81 latent variables in the model (Lemey et al., 2009). The MultiTypeTree model (Vaughan
82 et al., 2014) removes the independence assumption that the evolutionary history of the
83 genealogy is independent of population structure, however the induced model complexity
84 often renders sampling the latent phylogeny and state history computationally infeasible.
85 This limitation is addressed with the structured coalescent of Volz et al. (2009), which
86 integrates over the state histories of the phylogeny and describes the marginal probabili-
87 ties of each viral lineage to be in a particular state at a particular time. The changes
88 in the state probabilities along lineages and through coalescent events are derived under
89 an assumed ordinary differential equations (ODE) model of disease spread. The flow

6 *on behalf of Rakai Health Sciences Program and PANGEA-HIV*

90 parameters are obtained as by-products of the estimated latent states and parameters
91 of the compartment model, and in general vary in time over the phylogenetic history.
92 Müller et al. (2017) and others (Popinga et al., 2015) have clarified the assumptions that
93 underlie the marginalisation of state histories into lineage state probabilities, which lead
94 to more accurate maximum likelihood and Bayesian estimation routines (Müller et al.,
95 2018; Volz and Siveroni, 2018). With this marginalisation trick, a greater range of flow
96 models and data sets can be analysed, though computational run-times often remain on
97 the order of several weeks for data from hundreds of individuals and moderately com-
98 plex flow models (Stadler and Bonhoeffer, 2013; Ratmann et al., 2017; Rasmussen et al.,
99 2018).

100 An alternative strategy for estimating the flow parameters (1– 3) involves phyloge-
101 netic analysis of multiple distinct pathogen sequences per infected host, because such
102 data make possible to attribute viral lineages to individuals and infer their ancestral
103 relationships, which can provide direct evidence into the direction of transmission be-
104 tween two individuals (Romero-Severson et al., 2016; Leitner and Romero-Severson,
105 2018). This strategy is becoming broadly applicable to fast-evolving pathogens such
106 as HIV, because standard deep sequencing protocols (Gall et al., 2012; Bonsall et al.,
107 2018; Zhang et al., 2020) generate thousands to millions of distinct pathogen sequence
108 fragments. In particular the PANGEA consortium is providing access to deep sequences
109 from > 20,000 HIV-infected individuals from across sub-Saharan Africa, including dense
110 population-based samples from locations where the burden of HIV is among the high-
111 est in Africa (Abeler-Dörner et al., 2019). Prior work focused on software development
112 (Wymant et al., 2017; Skums et al., 2018), validation of the bioinformatics protocol for
113 inferring the direction of transmission (Ratmann et al., 2018; Rose et al., 2019; Todesco
114 et al., 2019; Zhang et al., 2020), and reconstruction of partially observed transmission
115 networks at the population level in the case of HIV (Ratmann et al., 2018).

116 The starting point of this paper is the output of a typical deep-sequence phylogenetic
117 analysis (Wymant et al., 2017), which includes for each ordered pair of sampled indi-
118 viduals a viral phylogenetic measure in [0, 1] that transmission likely occurred from the
119 first to the second individual, possibly via unsampled intermediate individuals (phylo-

120 genetic direction scores). Here we show that such deep-sequence data offer distinctive
 121 advantages for estimating population-level transmission flows, when compared to viral
 122 phylogenetic approaches based on standard Sanger sequencing. First, the data enable
 123 us to present the estimation problem in terms of a class of hierarchical Bayesian Pois-
 124 son models that can flexibly describe a range of epidemiological questions including
 125 transmission in space (1), or by age and sex (2). This class of models is in principle
 126 similar to previously developed, high-dimensional Bayesian models used for estimating
 127 contact rates from non-Gaussian response data (van de Kassteele et al., 2017). Sec-
 128 ond, the models can account for multi-level sampling heterogeneity, which are typically
 129 present in population-based disease occurrence data. We leverage Bayesian data aug-
 130 mentation to adjust for sampling heterogeneity (Givens et al., 1997), and exploit the
 131 fact that the additional latent variables can be integrated out in our framework, so that
 132 computational inference remains inexpensive. Third, while typical phylodynamic ap-
 133 proaches are limited to estimating transmission flows between coarse population strata,
 134 for example by age brackets 15 – 24 years and 25 – 40 years (De Oliveira et al., 2017;
 135 Le Vu et al., 2019), we can employ Gaussian-process-based regularisation techniques to
 136 capture fine detail in transmission flows by annual age increments. To illustrate, we
 137 present these innovations on combined epidemiologic and HIV deep sequence data from
 138 the Rakai Community Cohort Study (RCCS) of the Rakai Health Science program, sit-
 139 uated in south-eastern Uganda (Grabowski et al., 2017; Chang et al., 2016). Between
 140 August 10 2011 to January 30 2015, virus from 2652 HIV-infected individuals could
 141 be deep-sequenced, and 293 pairs of individuals with phylogenetically strong support
 142 for the direction of transmission were identified. We demonstrate that the new type
 143 of phylogenetic data and our statistical model enable characterisation of cross-sectional
 144 transmission flows at unprecedented detail while remaining computationally scalable.

145 2. Methodology

146 2.1. Notation and Definitions.

147 In this section we present the notation that is used to estimate transmission flows in a
 148 population \mathcal{P} during some observation period $\mathcal{T} = [t_1, t_2]$. We define by $i = 1, \dots, N$ the
 149 unique identifier of an individual in \mathcal{P} during \mathcal{T} . Individual-level characteristics such

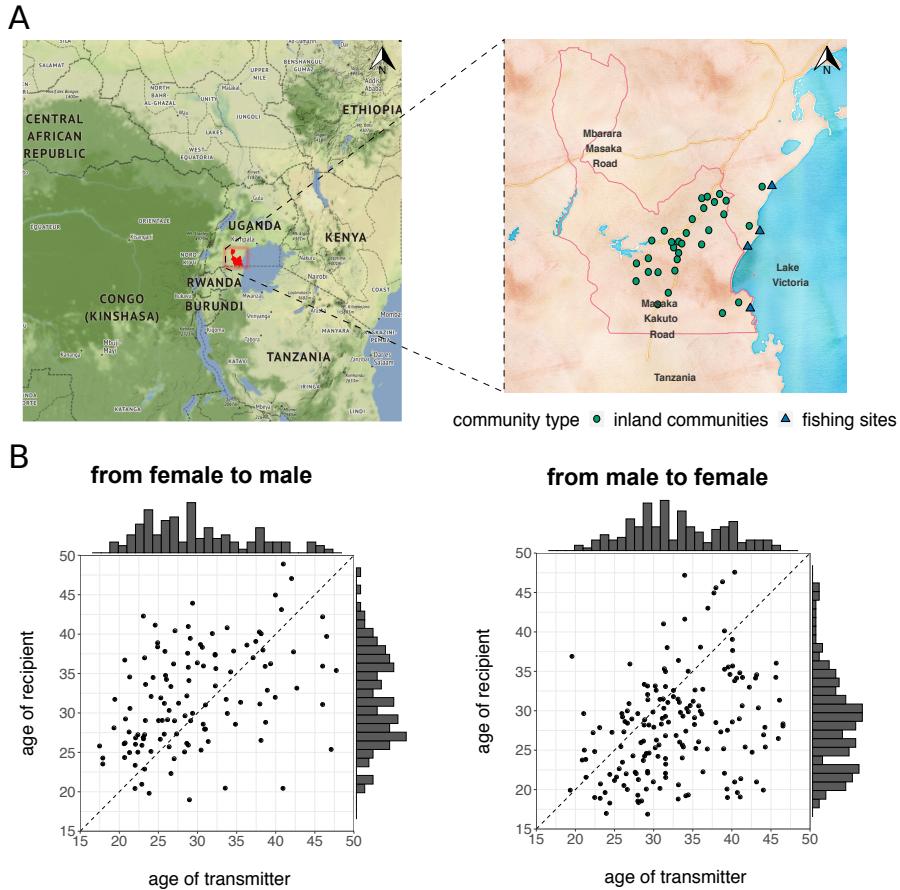


Fig. 2. Location of the Rakai Community Cohort Study, and data. (A) Location of Rakai District (red) in south-eastern Uganda at the shores of Lake Victoria. HIV surveillance data were obtained from 2 survey rounds in 36 inland communities of the Rakai Community Cohort Study (green circles) and three survey rounds in the main 4 fishing communities within 3km of Lake Victoria (green triangles) between August 10, 2011 and January 30, 2015. (B) The study did a household census, and all individuals aged 15-49 years capable to provide informed consent and resident for at least 1 month with the intention to stay were invited to participate. Viral deep-sequencing was performed on plasma blood samples from HIV infected participants who reported no ART use. Deep-sequence phylogenetic analysis returned phylogenetic transmission scores between individuals, and 293 pairs had strong support of phylogenetic linkage and transmission direction. 173 pairs were male-to-female and 120 were female-to-male. The figures show the phylogenetically likely source-recipient pairs by age of the source and recipient at the midpoint of the observation period.

as age or location of residence are described with p covariates, which we collect for all individuals in the $N \times p$ matrix \mathbf{X} . The sampling status of individual i is a binary variable, and individuals are either sampled ($s_i = 1$) or not ($s_i = 0$). We denote the sampling status vector for all individuals in \mathcal{P} by $\mathbf{s} = (s_i)$. The number of sampled

154 individuals is N^s , which corresponds in our exposition to the 2652 individuals for whom
 155 a viral deep sequence is available for analysis. The output of the phylogenetic deep
 156 sequence analysis can thus be described in a $N^s \times N^s$ direction score matrix \mathbf{W} that
 157 describes the evidence for transmission from i to j with the weight $w_{ij} \in [0, 1]$. The
 158 direction score matrix is in general not symmetric, and diagonal entries are zero.

159 We estimate transmission flows between population strata, defined for example by the
 160 cartesian product of gender and age (rounded to years), and denote the strata by a and
 161 the set of strata by \mathcal{A} , which is of dimension $A > 0$. Thus the primary object of interest
 162 is the $A \times A$ flow matrix $\boldsymbol{\pi}$, which is in general not symmetric, has positive diagonal
 163 entries, and is subject to $\sum_{a,b \in \mathcal{A}} \pi_{ab} = 1$. The matrix may contain structural zeros, for
 164 example in the case of HIV female to female transmission is typically not observed [],
 165 and we denote the number of structurally non-zero entries by L , which satisfies $L \leq A^2$.

166 In general the flow matrix is time-dependent due to changes in population composition
 167 and varying transmission rates [Anderson]. For instance in a compartment model of
 168 susceptible (S), infected (I) and treated (T) men and women of high (h) and low risk
 169 (l) of onward transmission, the ODE equations pertaining to the male (m) high risk
 170 population are

$$\begin{aligned}\dot{S}_{mh} &= -\lambda(t)S_{mh}(t) + \mu - \mu S_{mh}(t) \\ \dot{I}_{mh} &= \lambda(t)I_{mh}(t) - \gamma(t)I_{mh}(t) - \mu I_{mh}(t) \\ \dot{T}_{mh} &= \gamma(t)I_{mh}(t) - \mu T_{mh}(t),\end{aligned}\tag{4}$$

171 where the force of infection is $\lambda(t) = \beta_{fh}(t)I_{fh}(t)/N_{fh}(t) + \beta_{fl}(t)I_{fl}(t)/N_{fl}(t)$, the
 172 birth/death rate μ is constant, and the viral suppression rate γ and transmission rates
 173 β_{fh}, β_{fl} are time-dependent. The actual, unobserved number of transmissions from high
 174 risk women to high risk men in $\mathcal{T} = [t_1, t_2]$ are

$$z_{fh,mh}([t_1, t_2]) = \int_{t_1}^{t_2} \beta_{fh}(t)I_{fh}(t)S_{mh}(t)/N_{fh}(t)dt,\tag{5}$$

175 and the corresponding proportion of transmissions is

$$\pi_{fh,mh}([t_1, t_2]) = \frac{z_{fh,mh}([t_1, t_2])}{Z([t_1, t_2])},\tag{6}$$

176 where Z is the sum of transmission events in $\mathcal{T} = [t_1, t_2]$. Here, we focus on estimating
 177 transmission flows in a given observation window from data collected in the same period.

10 *on behalf of Rakai Health Sciences Program and PANGEA-HIV*

178 Such estimates provide a snapshot of transmission dynamics. For ease of notation we
179 now drop the dependence of our data and estimates on $\mathcal{T} = [t_1, t_2]$.

180 As a first approach to estimating the flow matrix, consider the observed flow counts

$$n_{ab} = \sum_{i \in a, j \in b} \mathbb{1}\{s_i = 1\} \mathbb{1}\{s_j = 1\} \mathbb{1}\{w_{ij} > \zeta\} \in \mathbb{N}_0^+, \quad (7)$$

181 where $\zeta \in (0, 1)$ is a threshold that can be used to select phylogenetically highly sup-
182 ported source-recipient pairs. The counts can be arranged into the $A \times A$ count matrix
183 \mathbf{n} , and sum to $N^f = \sum_{a,b} n_{ab}$. In previous studies, ζ was set to 0.5 or 0.6, and N^f was
184 between 100 to 500 (Hall et al., 2019; Ratmann et al., 2020). The naïve flow estimator
185 is defined by

$$\hat{\pi}_{ab} = \frac{n_{ab}}{\sum_{c,d} n_{cd}}. \quad (8)$$

186 If we suppose that each population group a is independently sampled at random with
187 probability ξ_a , and the actual flows from group a to group b are z_{ab} , then $\mathbb{E}(\hat{\pi}_{ab}) =$
188 $(z_{ab}\xi_a\xi_b)/(\sum_{c,d} z_{cd}\xi_c\xi_d)$. This clarifies that the naïve flow estimator (8) is only unbiased
189 when the population groups were homogeneously sampled, i. e. ξ_a is the same for all a ,
190 which is rarely the case (Ratmann et al., 2020).

191 2.2. *Inferring transmission flows from heterogeneously sampled deep-sequence data*

192 Considering the actual, unobserved number of transmissions between all population
193 groups, $\mathbf{z} = (z_{ab})$, and total $z^+ = \sum_{a,b} z_{ab}$, the complete data likelihood that arises under
194 mathematical models of the form (4) in a fixed observation period is the multinomial

$$p(\mathbf{z}|z^+, \boldsymbol{\pi}) = \prod_{a,b} (\pi_{ab})^{z_{ab}}. \quad (9)$$

195 This model ignores potential second-order correlations between transmission events, for
196 example that a female infected by an older male may be more likely to transmit to men
197 of older age. For computational efficiency, we consider the related Poisson model

$$p(\mathbf{z}|\boldsymbol{\lambda}) = \prod_{a,b} (\lambda_{ab})^{z_{ab}} \exp(-\lambda_{ab}) = \left[\prod_{a,b} (\pi_{ab})^{z_{ab}} \right] \left[\eta^{z^+} \exp(-\eta) \right], \quad (10)$$

198 where λ_{ab} can be interpreted as the transmission intensities from group a to group b ,
199 $\eta = \sum_{c,d} \lambda_{cd}$, and π_{ab} are recovered via $\pi_{ab} = \lambda_{ab}/\eta$.

200 The actual transmission flows \mathbf{z} are not observed. We assume that individuals are
 201 sampled at random within strata (SARWS) in a suitably chosen stratification of the
 202 population. SARWS implies in particular that sampling is independent of transmission
 203 status (i.e. being a transmitter or not), and the likelihood of the observed counts con-
 204 ditional on the complete data is $p(\mathbf{n}|\mathbf{z}, \boldsymbol{\xi}) = \prod_{a,b} \text{Binomial}(n_{ab}; z_{ab}, \xi_a \xi_b)$, where ξ_a is the
 205 sampling probability in group a . In this class of models the latent transmission counts
 206 z_{ab} can be conveniently integrated out, yielding for the observed flow counts the Poisson
 207 model

$$p(\mathbf{n}|\boldsymbol{\lambda}, \boldsymbol{\xi}) = \prod_{a,b} (\lambda_{ab} \xi_a \xi_b)^{n_{ab}} \exp(-\lambda_{ab} \xi_a \xi_b). \quad (11)$$

208 The sampling-adjusted maximum-likelihood estimates of λ_{ab} and π_{ab} under (11) can
 209 be derived under the SARWS assumption. The number of sampled individuals in a ,
 210 $N_a^s = \sum_{i \in a} \mathbb{1}\{s_i = 1\}$, is a Binomial sample of the number of all individuals in a , N_a ,
 211 which leads to

$$\hat{\pi}_{ab} = \frac{n_{ab}}{\hat{\xi}_a \hat{\xi}_b} / \left[\sum_{c,d} \frac{n_{cd}}{\hat{\xi}_c \hat{\xi}_d} \right] \quad (12)$$

212 where $\hat{\xi}_a = N_a^s / N_a$ (Supplementary Text Sxx).

213 2.3. Regularisation

214 Bayesian regularisation techniques play a central role in obtaining robust and suitably
 215 smoothed flow estimates. Considering population sampling, we exploit additional in-
 216 formation on the sampling vector \mathbf{s} . We assume that transmission is independent of
 217 sampling, allowing us to decompose the joint posterior distribution into

$$\begin{aligned} p(\boldsymbol{\lambda}, \boldsymbol{\xi} | \mathbf{n}, \mathbf{s}, \mathbf{X}) &\propto p(\mathbf{n} | \boldsymbol{\lambda}, \boldsymbol{\xi}, \mathbf{s}, \mathbf{X}) p(\boldsymbol{\lambda} | \boldsymbol{\xi}, \mathbf{s}, \mathbf{X}) p(\boldsymbol{\xi} | \mathbf{s}, \mathbf{X}) \\ &= p(\mathbf{n} | \boldsymbol{\lambda}, \boldsymbol{\xi}) p(\boldsymbol{\lambda} | \boldsymbol{\xi}) p(\boldsymbol{\xi} | \mathbf{s}, \mathbf{X}). \end{aligned} \quad (13)$$

218 A possible limitation of (12) is that the counts N_a, N_a^s can be small when the population
 219 is finely stratified. It is thus often advantageous to model individual-level sampling
 220 probabilities in terms of a linear combination of predictors. Using a logistic regression
 221 approach, we obtain

$$p(\xi_a | \mathbf{s}, \mathbf{X}) = \int \text{logit}^{-1}(\mathbf{x}_a \boldsymbol{\beta}) p(\boldsymbol{\beta} | \mathbf{s}, \mathbf{X}) d\boldsymbol{\beta}, \quad (14)$$

12 *on behalf of Rakai Health Sciences Program and PANGEA-HIV*

222 where \mathbf{x}_a is the row vector of population characteristics that specify group a , $\boldsymbol{\beta}$ are the
 223 regression coefficients, and $p(\boldsymbol{\beta}|\mathbf{s}, \mathbf{X})$ is the posterior density of the regression coefficients,
 224 estimated from the sampling status vector \mathbf{s} of all individuals in the study population.

225 Considering the prior density on the transmission intensities $\boldsymbol{\lambda}$, in some applications
 226 the population strata are unordered such as when estimating transmission flows within
 227 and between high- and low-prevalence areas (1). In this case we propose using

$$\lambda_{ab}|\xi_a, \xi_b \sim \text{Gamma}(\alpha_{ab}, \beta), \quad \alpha_{ab} = 0.8/L, \beta = 0.8/Z^p(\xi_a, \xi_b), \quad (15)$$

228 where L is the number of structurally non-zero entries in $\boldsymbol{\pi}$ and Z^p is the number of
 229 expected transmission events, $Z^p = \sum_{a,b:n_{ab} \neq 0} \frac{n_{ab}}{\xi_a \xi_b} + \sum_{a,b:n_{ab}=0} \frac{1-\xi_a \xi_b}{\xi_a \xi_b}$. This choice is
 230 motivated by the fact that (15) induces on $\boldsymbol{\pi}$ an objective Dirichlet prior density with
 231 parameters $\alpha_{ab} = 0.8/L$ (Berger et al., 2015). However when the population groups
 232 can be ordered, such as the 1-year age bands in (2), the structure of the flow model
 233 (13) enables using regularising prior densities that penalise against large deviations in
 234 transmission intensities between similar source and recipient populations. For (2), we
 235 opted for two-dimensional Gaussian-process priors on the non-zero entries of $\boldsymbol{\lambda}$,

$$\begin{aligned} \log \boldsymbol{\lambda} &= \mu + \mu^{mf} + \mathbf{f}, \quad \mathbf{f} \sim \mathcal{GP}(0, k), \\ k((a_1, b_1), (a_2, b_2)) &= \begin{cases} \sigma_{mf}^2 \exp \left(- \left[\frac{(a_2-a_1)^2}{2\ell_{mf,a}^2} + \frac{(b_2-b_1)^2}{2\ell_{mf,b}^2} \right] \right) & \text{if } (a_1, b_1), (a_2, b_2) \in mf \\ \sigma_{fm}^2 \exp \left(- \left[\frac{(a_2-a_1)^2}{2\ell_{fm,a}^2} + \frac{(b_2-b_1)^2}{2\ell_{fm,b}^2} \right] \right) & \text{if } (a_1, b_1), (a_2, b_2) \in fm \\ 0 & \text{otherwise} \end{cases} \\ \sigma_{mf}^2, \sigma_{fm}^2 &\sim \text{Half-Normal}(0, 10) \\ \ell_{d,i} &\sim \text{Inv-Gamma}(\alpha_{d,i}, \beta_{d,i}), \end{aligned} \quad (16)$$

236 where μ is the baseline log transmission intensity, μ^{mf} is a scalar on the elements of $\boldsymbol{\lambda}$
 237 in the male-female direction, and k is a gender-specific squared exponential kernel with
 238 variance parameters $\sigma_{mf}^2, \sigma_{fm}^2$ and length scales $\ell_{mf,a}, \ell_{mf,b}, \ell_{fm,a}, \ell_{fm,b}$ (Rasmussen,
 239 2003). The prior distributions on the length scale parameters were set so that their 95%
 240 credibility range matched the empirical 95% quantile range in Figure 2.

²⁴¹ **2.4. Numerical inference**

²⁴² The transmission flows and parameters for these highly structured Bayesian models can
²⁴³ be efficiently estimated with the Hamiltonian Monte Carlo sampler of the Stan com-
²⁴⁴ puting language (Carpenter et al., 2017). The implementation is available at <https://github.com/BDI-pathogens/phyloscanner/tree/master/phyloflows> and uses a spec-
²⁴⁵ tral basis function approximation for the Gaussian process prior (16). For comparison a
²⁴⁶ tailored Metropolis-within-Gibbs algorithm is also available. See Supplementary Text S2
²⁴⁷ for further details.

²⁴⁹ **3. Applications**

²⁵⁰ **3.1. Bias in source attribution when sampling differences are ignored**

²⁵¹ We first assessed the accuracy in estimating transmission flows from biased flow data in a
²⁵² series of simulation experiments, that are fully reported in Supplementary Text S3. The
²⁵³ first experiment was a minimal example involving flows between two population groups,
²⁵⁴ which for simplicity we refer to as individuals in rural areas (group *a*) and individuals in
²⁵⁵ large communities (group *b*). Transmission chains were simulated under the ODE model
²⁵⁶ (4) in, and the simulated flow matrix π_r^T was recorded in $r = 1, \dots, 100$ replicate simula-
²⁵⁷ tions. The sampling probability in rural areas was $\xi_a = 60\%$, and sampling probabilities
²⁵⁸ in large communities decreased from $\xi_b = 60\%$ to 35%, resulting in sampling differences
²⁵⁹ of (0%, 5%, 10%, 15%, 20%, 25%). First, we estimated the marginal median poste-
²⁶⁰ rior transmission flows $\hat{\pi}_{r,ab}$ from (13) under the assumption of no sampling differences,
²⁶¹ which we implemented by setting $\xi_a = \xi_b = 0.5$. Second, we estimated the marginal
²⁶² median posterior transmission flows with information on sampling differences included,
²⁶³ as in (??). Figure 3A compares the accuracy in transmission flow estimates when sam-
²⁶⁴ pling differences are ignored (light grey bars) compared to when they are not ignored
²⁶⁵ (dark grey bars) in terms of the worst case error (WCE) $\varepsilon_r = \max_{a,b} |\hat{\pi}_{r,ab} - \pi_{r,ab}^T|$. For
²⁶⁶ a 10% sampling difference between the two population groups, the median WCE was
²⁶⁷ 6.5% (3.5% – 10.2%) without adjusting for sampling differences, and 2.1% (0.5% – 5.1%)
²⁶⁸ after adjusting for sampling differences. More complex simulation experiments yielded
²⁶⁹ similar results (Supplementary Text S3).

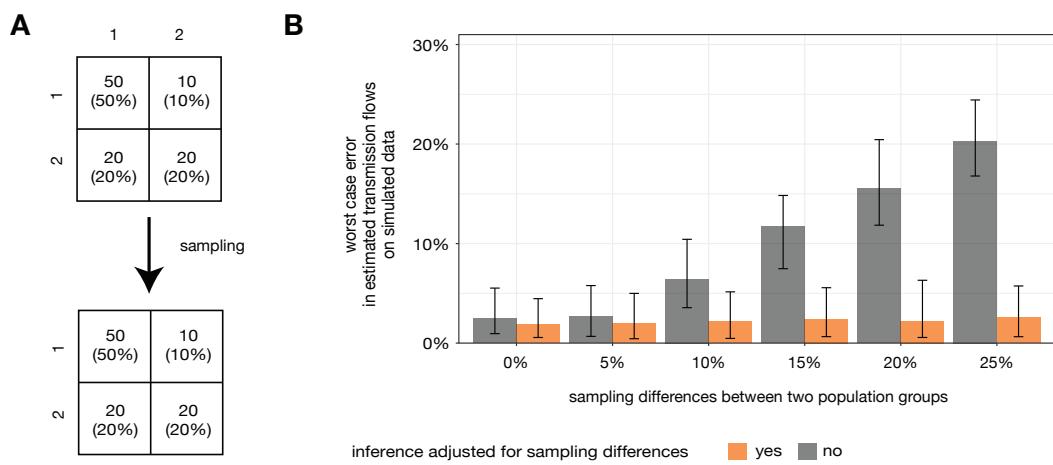


Fig. 3. Comparing errors in estimates of transmission flow. We assessed bias in estimating transmission flows with model (13) before and after accounting for sampling differences on simulations. **(A)** Design of simulation study. Transmission flows were simulated between two population groups according to an ODE model of disease spread. The number and proportion of actual transmissions was recorded (top 2 by 2 table, with source group in rows and recipient group in columns). The probability of sampling individuals in the first group was 60%, and decreased in the second group from 60% to 35%. Transmission events between two sampled events were recorded as observed (bottom 2 by 2 table). **(B)** The worst case error between the simulated transmission flows and median posterior estimates was calculated on 100 simulations in each sampling scenario, and is shown on the y-axis against increasing sampling differences on the x-axis. With increasing differences in sampling differences, we found increasing bias in flow estimates when sampling differences were not adjusted for. When sampling differences were adjusted for, the median worst case error in flow estimates remained on average below 5%.

270 3.2. Population-based deep-sequence data from Rakai, Uganda

271 We illustrate application of the Poisson flow model (13) on a population-based sam-
272 ple of HIV deep sequences from the RCCS in south-eastern Uganda at the shores of
273 Lake Victoria (Ratmann et al., 2018, 2020). Between August 10 2011 to January 30
274 2015, 29116 individuals aged 15-49 years were eligible to participate in one of the 36
275 inland communities surveyed, and 8526 individuals in 4 fishing communities in popula-
276 tion censuses immediately preceding the survey (Figure 2A). 19799 (68.0%) participated
277 in inland communities, and 6083 (71.3%) in fishing communities. 11404 (96.9%) of
278 non-participants were absent for school or work. Participation increased with age for
279 both men and women among residents, was higher and more uniform among recent
280 in-migrants by age, and was similar in fishing and inland communities (Supplementary
281 Text S1). 2703 (13.6%) individuals were HIV-infected in inland communities on their
282 last visit, and 2439 (40.1%) in fishing sites.

283 HIV sequencing was only feasible under the deep-sequencing protocol that we used
284 when individuals had detectable virus (Gall et al., 2012), but viral load measurements
285 were not available. We selected infected individuals who did not report ART use at their
286 first visit for sequencing. In inland communities, 1803 (66.7%) infected participants did
287 not report ART use, of whom 1138 (63.1%) could be deep-sequenced. In fishing com-
288 munities, 2059 (84.4%) infected participants did not report ART use, and 1514 (73.6%)
289 could be deep-sequenced (Ratmann et al., 2018). Deep-sequencing rates decreased with
290 age for both men and women, were higher among men than women, and higher in fishing
291 sites (Supplementary Text S1).

292 Using (7), there were 293 heterosexual pairs with phylogenetic support for linkage
293 and direction of transmission above the threshold $\zeta = 0.6$ (source-recipient pairs), of
294 whom 173 (59.0%) were male-to-female, and 120 (41.0%) were female-to-male. The
295 estimated infection times of the recipients were between October 2009 and January
296 2015, which defined the observation period \mathcal{T} . Figure 2B illustrates the reconstructed
297 source-recipient pairs by age of both individuals at the midpoint of the study period. We
298 estimated that our population-based sample contained approximately 800 transmission
299 pairs (Supplementary Text Sxx), suggesting that not all transmission pairs could be

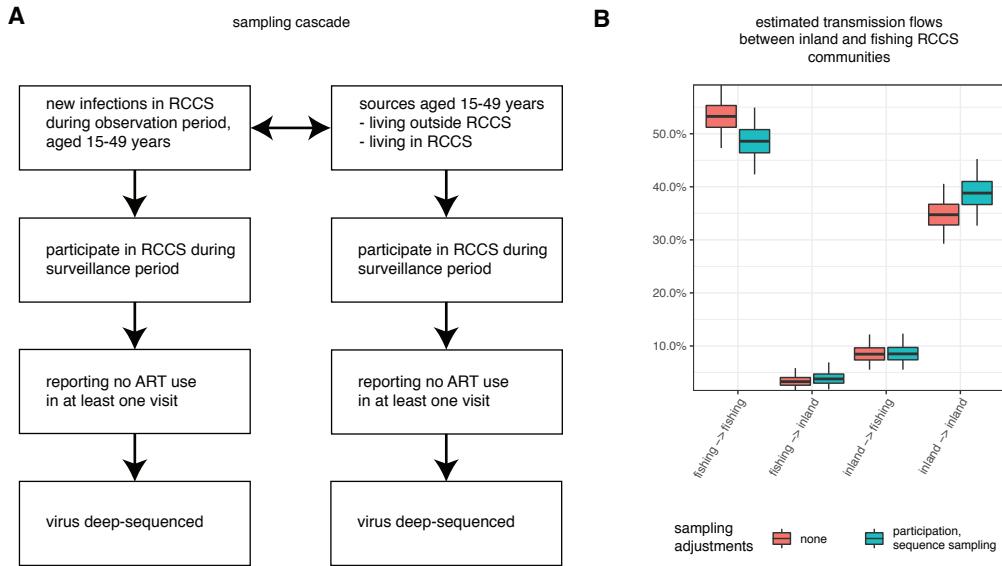


Fig. 4. Sampling cascade of transmission events, and impact on flow estimates. (A) The sampling cascade formalises the steps involved in sampling sources and recipients of transmission events that occurred in the observation period October 2009 to January 2015. Recipients were defined as individuals aged 15-49 years who acquired infection in one of the RCCS communities during the observation period. Sources were defined as individuals aged 15-49 years who transmitted to one of the recipients. Each arrow corresponds to a sampling step of the source and recipient populations, which are described in terms of individual-level predictive variables; see further Supplementary Text S1. (B) Estimated transmission flow estimates between inland and fishing communities of the RCCS. We first estimated transmission flows without sampling adjustments (red), using model (13) and setting the sampling probabilities ζ_a to the overall empirical average. We next iteratively drew sampling probabilities associated with each step of the sampling cascade for both sources and recipients, and used the resulting overall sampling probabilities to adjust transmission flow estimates with model (13). The resulting flow estimates (blue) indicated higher transmission flows within the inland communities of the RCCS and lower transmission flows within the fishing communities, which primarily reflected higher participation rates and higher sequence sampling rates in fishing communities.

300 identified through phylogenetic analysis. Using the thresholds $\zeta = 0.55$ and $0.5, 356$ and
 301 417 source-recipient pairs were found.

302 To interpret these observations, we formalised the individual steps in the sampling
 303 cascade of infection events to individuals aged 15-49 years between October 2009 and
 304 January 2015 in the RCCS and their sources (Figure 4A). The sources and recipients
 305 had to participate in at least one survey round between August 2011 and January 2015,
 306 report no ART use, and have virus sequenced successfully. For each step, we estimated
 307 based on detailed cohort data population-level sampling probabilities independently for
 308 sources and recipients under the assumption that sampling was independent of infection

309 and transmission status (i.e. being a source or recipient, Supplementary Text S1).

310 ***3.3. Transmission flows between areas with high and low disease prevalence***

311 We used the source-recipient data of Figure 2B to address problem (1) and estimate
 312 transmission HIV flows within and between high- and low-prevalence RCCS commu-
 313 nities. The high-prevalence communities comprised the four fishing communities, and
 314 the low-prevalence communities included the remaining 36 inland communities. Migra-
 315 tion is common in the study area (Grabowski et al., 2017). To account for sampling
 316 heterogeneity, the high- and low-prevalence communities were sub-stratified by gender,
 317 migration status, and three age bands (15 – 24, 25 – 34, 35 – 49), which resulted in 24
 318 sampling groups and 576 flow variables. For each step of the sampling cascade, Monte
 319 Carlo draws from the posterior distribution of the conditional sampling distributions
 320 were obtained for each of the 24 population groups, multiplied, and used to adjust the
 321 flow estimates. Numerical inference of the joint posterior density (13) took 3.7 hours on
 322 a 2.4 Ghz processor (Supplementary Text S4). The flow estimates by the 24 sampling
 323 groups were then aggregated into 4 flow estimates between and within the RCCS fish-
 324 ing and inland communities. Figure 4B compares the marginal posterior distributions
 325 for each of the 4 transmission flows obtained with and without adjusting for sampling
 326 differences of the population. Participation and sequencing rates were overall higher in
 327 fishing communities compared to inland communities, which led to significantly differ-
 328 ent flow estimates after adjusting for sampling heterogeneity. The estimated flow ratio
 329 (inland→fishing / fishing→inland) was 2.58 (1.24 - 5.91) without sampling adjustments
 330 and 2.25 (1.04-5.23) with sampling adjustments. Therefore both analyses supported,
 331 contrary to what has commonly been assumed (Uganda AIDS Commission, 2014), the
 332 finding that the high-prevalence fishing communities were net sinks of local infection
 333 flows (Ratmann et al., 2020).

334 ***3.4. Transmission flows between age groups***

335 We next turned to estimating transmission flows by age from the source-recipient data
 336 shown in Figure 5, our initial problem (2). Here, we divided the data by gender and
 337 1-year age bands (15, . . . , 49), and accounted for sampling heterogeneity using the same

18 *on behalf of Rakai Health Sciences Program and PANGEA-HIV*

338 strata, which resulted in 70 sampling groups and 2450 flow variables. To obtain smooth
339 estimates, we specified a GP prior on the log transmission intensities as in (15), and
340 determined the prior distributions on the GP hyper-parameters from the data shown in
341 Figure 2C. Specifically, the possible ranges of length-scale are $1 - 25$ (ρ_1), $1 - 15$ (ρ_2
342 when the source of infection is women) and $1 - 20$ (ρ_2 when the source of infection is
343 men). Numerical inference of the joint posterior density (13) took 51 hours on a 2.4 Ghz
344 processor (Supplementary Text S4).

345 Figure 5ACE shows the inferred age-distribution of male sources to female infections.
346 For example, the red “wave” shows the mean posterior density estimate that women of
347 age 16 are infected by men of particular ages (x-axis), i.e. $\mathbb{E}(\gamma_{ab}|\mathbf{n}, \mathbf{s}, \mathbf{X})$ for $b = (f, 16)$
348 in our notation of (3a). Corresponding coefficients of variation are shown in Figure
349 5E . For young women aged < 25 years, the mean posterior age-distribution of their
350 sources plateaued at age 35. Our data thus contain evidence that some women aged
351 < 25 years were infected by men who are at least 10 years older, often referred to as
352 sugar-daddies. However our non-parametric approach also indicates that transmission
353 from sugar-daddies was uncommon. Considering women aged 20 years, an estimated
354 XX% were infected from similarly aged men < 24 years, XX% from slightly older men
355 24 – 29 years, and XX% from considerably older men 30 – 49 years. Further, as infected
356 women age, our data contradict the idea that these women would then continue to infect
357 their peers. Figure 5BDF shows the mean posterior age-distribution of female sources
358 to male infections. Corresponding coefficients of variation are shown in Figure 5F. For
359 men aged > 35 years, the mean posterior age-distribution of their female source is not
360 centred among their peers. In fact, it is very broad, and much broader than for men
361 aged < 25 years.

362 **4. Discussion**

363 In this study we introduce a hierarchical Bayesian Poisson model for estimating time
364 homogeneous disease flows in human populations from pathogen deep-sequence data.
365 The model (13) allows inference of transmission flows between arbitrarily defined strata
366 of human populations, which enables addressing a range of epidemiological questions on

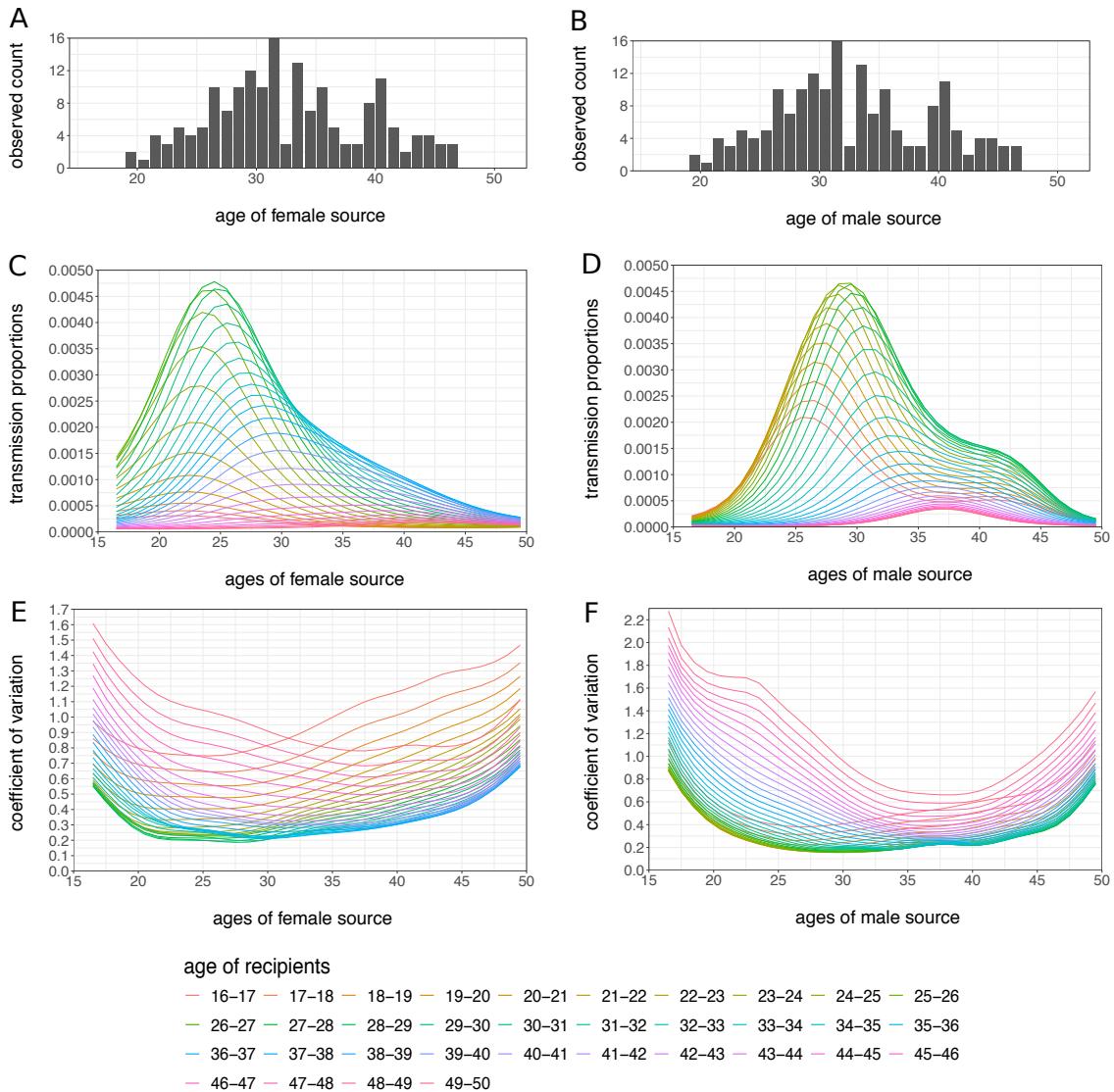


Fig. 5. transmission flows between age groups. Panel A/B shows the observed transmission counts from female/male one-year increment age groups. Gaussian process regression (16) was applied to estimate transmission flows between age groups. Panel C/D shows the proportion of transmissions from women/men in an age group (in horizontal axis), given the recipients' age groups (in color). Panel E/F plots the associated coefficients of variation of the estimated flows.

20 *on behalf of Rakai Health Sciences Program and PANGEA-HIV*

367 pathogen spread between geographic areas (CITE), by age categories (CITE), or other
368 discretely-valued sociodemographic characteristics (CITE). The modelling framework is
369 flexible and easy to extend. Here, we showed how Gaussian-process based regularisation
370 techniques (CITE) can be incorporated to obtain smoothed flow estimates at high res-
371 olution. We also showed how data sampling bias can be adjusted for in this framework
372 under the assumption that data are missing at random within strata (CITE). To fit the
373 hierarchical Bayesian Poisson flow model, we provide an R interface to the NUTS Hamil-
374 tonian Monte Carlo algorithm of the Stan computing language (CITE). Means, medians,
375 standard deviations and credibility intervals of the flow estimates and epidemiologically
376 important, related quantities (3) are straightforward to obtain from numerical samples
377 of the posterior distribution of the flow parameters.

378 Existing phylodynamic estimation approaches (CITE) are tailored for pathogen con-
379 sensus sequences. The approach described here is tailored for pathogen deep-sequence
380 data that captures in contrast to consensus sequences the genetic diversity of pathogens
381 within infected hosts. The higher resolution of these data make possible to estimate the
382 direction of pathogen spread between pairs of individuals (CITE). Although it is not pos-
383 sible to determine transmission direction with certainty between two individuals (CITE),
384 deep-sequence data provide an alternative starting point for estimating population-level
385 patterns of pathogen spread. The main advantages are first, that population-level spread
386 can be directly estimated from individual source-recipient relationships and associated
387 individual covariates. This allows modelling of and accounting for individual-level vari-
388 ables that shape disease transmission in remarkable detail, such as the sources of HIV
389 infection by 1-year age bands. Second, little computational effort is needed to fit the Pois-
390 son flow model XX to deep-sequence data, because it falls within the class of Bayesian
391 hierarchical models for binary data, for which efficient fitting procedures exist (CITE).
392 This makes it computationally feasible to investigate complex models of disease spread
393 such as the gender-and-age specific flow model XX that we considered here. It is also
394 possible to automate numerical inference using for instance Stan (CITE) as we did here,
395 which makes available well-tested and versatile numerical techniques for phylodynamic
396 inference. In comparison, when using pathogen consensus sequences, existing inference

397 techniques typically require solving complex systems of ordinary differential equations
398 (CITE BEAST PHYDYN), which limits modelling of disease transmission computa-
399 tionally to a relatively small number of covariates (CITE). Alternative phylogenetic
400 clustering techniques are computationally faster, but it remains difficult to interpret the
401 inferred associations between cluster membership and covariates associated with sampled
402 pathogen consensus sequences (CITE Poon, Le Vu).

403 The method we propose has limitations. First, deep-sequencing protocols generate
404 short sequence fragments, usually of 200 to 300 base pairs in length after trimming
405 adaptors and low quality ends, and merging paired end fragments. This implies that
406 pathogens need to evolve at a rate above 1/200 mutations per site within the time scales
407 of interest in human hosts, because otherwise reconstructed deep-sequence phylogenies
408 do not contain the pattern of ancestral subgraphs that is characteristic of pathogen
409 spread in one direction. Such high evolutionary rates are typical for viral pathogens that
410 infect and evolve in humans over long periods of time, such as hepatitis C or HIV (CITE).
411 We expect that the methods developed here will become applicable to a broad range of
412 viral and bacterial infectious diseases as existing deep-sequencing methods that generate
413 substantially longer pathogen sequence fragments become cheaper (CITE), or alternative
414 approaches are being developed. Second, our method requires that deep-sequence data
415 from a large population-based sample of infected individuals (> 30%) are available.
416 Such data are emerging, for example from the HPTN072 Popart HIV prevention trial
417 (CITE), several HIV surveillance cohorts in Africa (CITE), the molecular epidemiologic
418 HIV surveillance program in Germany (CITE), or TODO HPC EXAMPLE. Third, our
419 inferences are based on source-recipient pairs with strong evidence for the direction of
420 transmission, which is a subset of all the data available. We cannot exclude that this
421 selection step introduces bias into the estimates obtained with the flow model XX, and
422 recommend performing sensitivity analysis on the selection threshold in XX. Fourth, no
423 continuous variables. Fifth, time-homogeneous.

424 **Acknowledgements**

425 **References**

- 426 Abeler-Dörner, L., Grabowski, M. K., Rambaut, A., Pillay, D., Fraser, C. et al. (2019)
 427 Pangea-hiv 2: Phylogenetics and networks for generalised epidemics in africa. *Current*
 428 *Opinion in HIV and AIDS*, **14**, 173–180.
- 429 Aral, S. O., Torrone, E. and Bernstein, K. (2015) Geographical targeting to improve pro-
 430 gression through the sexually transmitted infection/hiv treatment continua in different
 431 populations. *Current Opinion in HIV and AIDS*, **10**, 477.
- 432 Berger, J. O., Bernardo, J. M. and Sun, D. (2015) Overall objective priors. *Bayesian*
 433 *Analysis*, **10**, 189–221.
- 434 Bonsall, D., Golubchik, T., de Cesare, M., Limbada, M., Kosloff, B., MacIntyre-Cockett,
 435 G., Hall, M., Wymant, C., Ansari, M. A., Abeler-Dörner, L. et al. (2018) A compre-
 436 hensive genomics solution for hiv surveillance and clinical monitoring in a global health
 437 setting. *bioRxiv*, 397083.
- 438 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M.,
 439 Brubaker, M., Guo, J., Li, P. and Riddell, A. (2017) Stan: A probabilistic program-
 440 ming language. *Journal of statistical software*, **76**.
- 441 Chang, L. W., Grabowski, M. K., Ssekubugu, R., Nalugoda, F., Kigozi, G., Nantume,
 442 B., Lessler, J., Moore, S. M., Quinn, T. C., Reynolds, S. J. et al. (2016) Heterogeneity
 443 of the HIV epidemic: an observational epidemiologic study of agrarian, trading, and
 444 fishing communities in Rakai, Uganda. *The lancet. HIV*, **3**, e388.
- 445 De Oliveira, T., Kharsany, A. B., Gräf, T., Cawood, C., Khanyile, D., Grobler, A.,
 446 Puren, A., Madurai, S., Baxter, C., Karim, Q. A. et al. (2017) Transmission networks
 447 and risk of hiv infection in kwazulu-natal, south africa: a community-wide phylogenetic
 448 study. *The lancet HIV*, **4**, e41–e50.
- 449 Dellicour, S., Baele, G., Dudas, G., Faria, N. R., Pybus, O. G., Suchard, M. A., Rambaut,
 450 A. and Lemey, P. (2018) Phylodynamic assessment of intervention strategies for the
 451 west african ebola virus outbreak. *Nature communications*, **9**, 1–9.

- 452 Didelot, X., Fraser, C., Gardy, J. and Colijn, C. (2017) Genomic infectious disease
453 epidemiology in partially sampled and ongoing outbreaks. *Molecular biology and evo-*
454 *lution*, **34**, 997–1007.
- 455 Dwyer-Lindgren, L., Cork, M. A., Sligar, A., Steuben, K. M., Wilson, K. F., Provost,
456 N. R., Mayala, B. K., VanderHeide, J. D., Collison, M. L., Hall, J. B. et al. (2019)
457 Mapping hiv prevalence in sub-saharan africa between 2000 and 2017. *Nature*, **570**,
458 189.
- 459 Faria, N. R., Rambaut, A., Suchard, M. A., Baele, G., Bedford, T., Ward, M. J., Tatem,
460 A. J., Sousa, J. D., Arinaminpathy, N., Pépin, J. et al. (2014) The early spread and
461 epidemic ignition of hiv-1 in human populations. *science*, **346**, 56–61.
- 462 Gall, A., Ferns, B., Morris, C., Watson, S., Cotten, M., Robinson, M., Berry, N., Pillay,
463 D. and Kellam, P. (2012) Universal amplification, next-generation sequencing, and
464 assembly of hiv-1 genomes. *Journal of clinical microbiology*, **50**, 3838–3844.
- 465 Givens, G. H., Smith, D. and Tweedie, R. (1997) Publication bias in meta-analysis: a
466 bayesian data-augmentation approach to account for issues exemplified in the passive
467 smoking debate. *Statistical Science*, 221–240.
- 468 Grabowski, M. K., Serwadda, D. M., Gray, R. H., Nakigozi, G., Kigozi, G., Kagaayi,
469 J., Ssekubugu, R., Nalugoda, F., Lessler, J., Lutalo, T. et al. (2017) HIV prevention
470 efforts and incidence of HIV in uganda. *New England Journal of Medicine*, **377**,
471 2154–2166.
- 472 Hall, M. D., Holden, M. T., Srisomang, P., Mahavanakul, W., Wuthiekanun, V., Lim-
473 mathurotsakul, D., Fountain, K., Parkhill, J., Nickerson, E. K., Peacock, S. J. et al.
474 (2019) Improved characterisation of mrsa transmission using within-host bacterial se-
475 quence diversity. *eLife*, **8**.
- 476 Hué, S., Pillay, D., Clewley, J. P. and Pybus, O. G. (2005) Genetic analysis reveals the
477 complex structure of hiv-1 transmission within defined risk groups. *Proceedings of the*
478 *National Academy of Sciences*, **102**, 4425–4429.

24 *on behalf of Rakai Health Sciences Program and PANGEA-HIV*

- 479 Kagaayi, J., Chang, L. W., Ssempijja, V., Grabowski, M. K., Ssekubugu, R., Nakigozi,
480 G., Kigozi, G., Serwadda, D. M., Gray, R. H., Nalugoda, F. et al. (2019) Impact of
481 combination hiv interventions on hiv incidence in hyperendemic fishing communities
482 in uganda: a prospective cohort study. *The Lancet HIV*, **6**, e680–e687.
- 483 van de Kassteele, J., van Eijkelen, J., Wallinga, J. et al. (2017) Efficient estimation
484 of age-specific social contact rates between men and women. *The Annals of Applied
485 Statistics*, **11**, 320–339.
- 486 Le Vu, S., Ratmann, O., Delpech, V., Brown, A. E., Gill, O. N., Tostevin, A., Dunn, D.,
487 Fraser, C., Volz, E. M. and Database, U. H. D. R. (2019) Hiv-1 transmission patterns
488 in men who have sex with men: Insights from genetic source attribution analysis.
489 *AIDS research and human retroviruses*, **35**, 805–813.
- 490 Leitner, T. and Romero-Severson, E. (2018) Phylogenetic patterns recover known hiv
491 epidemiological relationships and reveal common transmission of multiple variants.
492 *Nature microbiology*, **3**, 983–988.
- 493 Lemey, P., Rambaut, A., Drummond, A. J. and Suchard, M. A. (2009) Bayesian phylo-
494 geography finds its roots. *PLoS computational biology*, **5**.
- 495 Müller, N. F., Rasmussen, D. and Stadler, T. (2018) Mascot: parameter and state
496 inference under the marginal structured coalescent approximation. *Bioinformatics*,
497 **34**, 3843–3848.
- 498 Müller, N. F., Rasmussen, D. A. and Stadler, T. (2017) The structured coalescent and
499 its approximations. *Molecular biology and evolution*, **34**, 2970–2981.
- 500 Poon, A. F., Gustafson, R., Daly, P., Zerr, L., Demlow, S. E., Wong, J., Woods, C. K.,
501 Hogg, R. S., Krajden, M., Moore, D. et al. (2016) Near real-time monitoring of hiv
502 transmission hotspots from routine hiv genotyping: an implementation case study.
503 *The lancet HIV*, **3**, e231–e238.
- 504 Popinga, A., Vaughan, T., Stadler, T. and Drummond, A. J. (2015) Inferring epidemi-
505 ological dynamics with bayesian coalescent inference: the merits of deterministic and
506 stochastic models. *Genetics*, **199**, 595–607.

- 507 Rasmussen, C. E. (2003) Gaussian processes in machine learning. In *Summer School on*
508 *Machine Learning*, 63–71. Springer.
- 509 Rasmussen, D. A., Wilkinson, E., Vandormael, A., Tanser, F., Pillay, D., Stadler, T.
510 and De Oliveira, T. (2018) Tracking external introductions of hiv using phylodynamics
511 reveals a major source of infections in rural kwazulu-natal, south africa. *Virus*
512 *evolution*, **4**, vey037.
- 513 Ratmann, O., Grabowski, M. K., Hall, M., Golubchik, T., Wymant, C., Hoppe, A.,
514 Brown, A. L., de Oliveira, T., Gall, A., Kellam, P., Pillay, D., Quinn, T., Wawer,
515 M., Laeyendecker, O., Serwadda, D., Gray, R. and Fraser, C. (2018) Inferring HIV-
516 1 transmission networks and sources of ongoing viral spread in Africa with next-
517 generation sequencing. In preparation.
- 518 Ratmann, O., Hodcroft, E. B., Pickles, M., Cori, A., Hall, M., Lycett, S., Colijn, C.,
519 Dearlove, B., Didelot, X., Frost, S. et al. (2017) Phylogenetic tools for generalized
520 hiv-1 epidemics: findings from the pangea-hiv methods comparison. *Molecular biology*
521 *and evolution*, **34**, 185–203.
- 522 Ratmann, O., Kagaayi, J., Hall, M., Golubchick, T., Kigozi, G., Xi, X., Wymant, C.,
523 Nakigozi, G., Abeler-Dörner, L., Bonsall, D. et al. (2020) Quantifying hiv transmission
524 flow between high-prevalence hotspots and surrounding communities: a population-
525 based study in rakai, uganda. *The Lancet HIV*.
- 526 Ratmann, O., Van Sighem, A., Bezemer, D., Gavryushkina, A., Jurriaans, S., Wensing,
527 A., De Wolf, F., Reiss, P., Fraser, C. et al. (2016) Sources of hiv infection among men
528 having sex with men and implications for prevention. *Science translational medicine*,
529 **8**, 320ra2–320ra2.
- 530 Riutort-Mayol, G., Vehtari, A., Bürkner, P.-C. and Riis, A. M. (2020) .
- 531 Romero-Severson, E. O., Bulla, I. and Leitner, T. (2016) Phylogenetically resolving
532 epidemiologic linkage. *Proceedings of the National Academy of Sciences*, **113**, 2690–
533 2695.

26 *on behalf of Rakai Health Sciences Program and PANGEA-HIV*

- 534 Rose, R., Hall, M., Redd, A. D., Lamers, S., Barbier, A. E., Porcella, S. F., Hudelson,
535 S. E., Piwowar-Manning, E., McCauley, M., Gamble, T. et al. (2019) Phylogenetic
536 methods inconsistently predict the direction of hiv transmission among heterosexual
537 pairs in the hptn 052 cohort. *The Journal of infectious diseases*, **220**, 1406–1413.
- 538 Skums, P., Zelikovsky, A., Singh, R., Gussler, W., Dimitrova, Z., Knyazev, S., Mandric,
539 I., Ramachandran, S., Campo, D., Jha, D. et al. (2018) Quentin: reconstruction of
540 disease transmissions from viral quasispecies genomic data. *Bioinformatics*, **34**, 163–
541 170.
- 542 Solin, A. and Särkkä, S. (2014) Hilbert space methods for reduced-rank gaussian process
543 regression. *Statistics and Computing*, 1–28.
- 544 Stadler, T. and Bonhoeffer, S. (2013) Uncovering epidemiological dynamics in heteroge-
545 neous host populations using phylogenetic methods. *Philosophical Transactions of the
546 Royal Society B: Biological Sciences*, **368**, 20120198.
- 547 The Office of the US Global AIDS Coordinator (2014) Pepfar 3.0 con-
548 trolling the epidemic: delivering on the promise of an aids-free gen-
549 eration. URL: <https://aidsfree.usaid.gov/resources/prevention-update/editions/december-2014/pepfar-30-controlling-epidemic-delivering>.
- 551 Todesco, E., Wirden, M., Calin, R., Simon, A., Sayon, S., Barin, F., Katlama, C.,
552 Calvez, V., Marcellin, A.-G. and Hué, S. (2019) Caution is needed in interpreting hiv
553 transmission chains by ultradeep sequencing. *Aids*, **33**, 691–699.
- 554 Uganda AIDS Commission (2014) National strategic plan for hiv aids, 2014.
555 URL: http://www.nationalplanningcycles.org/sites/default/files/country_docs/Uganda/national_strategic_plan_for_hiv_aids_2011_2015.pdf.
- 557 Uhlemann, A.-C., Dordel, J., Knox, J. R., Raven, K. E., Parkhill, J., Holden, M. T., Pea-
558 cock, S. J. and Lowy, F. D. (2014) Molecular tracing of the emergence, diversification,
559 and transmission of s. aureus sequence type 8 in a new york community. *Proceedings
560 of the National Academy of Sciences*, **111**, 6738–6743.

- 561 UNAIDS (2014a) 90-90-90: an ambitious treatment target to help end the aids epidemic.
562 geneva: Unaids; 2014.
- 563 — (2014b) Fast-track: ending the aids epidemic by 2030. URL: https://www.unaids.org/en/resources/documents/2014/JC2686_WAD2014report.
- 565 — (2018) Unaids data 2018. URL: https://www.unaids.org/sites/default/files/media_asset/unaids-data-2018_en.pdf.
- 567 Vaughan, T. G., Kühnert, D., Popinga, A., Welch, D. and Drummond, A. J. (2014)
568 Efficient bayesian inference under the structured coalescent. *Bioinformatics*, **30**, 2272–
569 2279.
- 570 Volz, E. M., Ionides, E., Romero-Severson, E. O., Brandt, M.-G., Mokotoff, E. and
571 Koopman, J. S. (2013) Hiv-1 transmission during early infection in men who have sex
572 with men: a phylodynamic analysis. *PLoS medicine*, **10**.
- 573 Volz, E. M., Pond, S. L. K., Ward, M. J., Brown, A. J. L. and Frost, S. D. (2009)
574 Phylodynamics of infectious disease epidemics. *Genetics*, **183**, 1421–1430.
- 575 Volz, E. M. and Siveroni, I. (2018) Bayesian phylodynamic inference with complex mod-
576 els. *PLoS computational biology*, **14**, e1006546.
- 577 Wymant, C., Hall, M., Ratmann, O., Bonsall, D., Golubchik, T., de Cesare, M., Gall, A.,
578 Cornelissen, M., Fraser, C., STOP-HCV Consortium, T. M. P. C. and Collaboration,
579 T. B. (2017) PHYLOSCANNER: inferring transmission from within-and between-host
580 pathogen genetic diversity. *Molecular biology and evolution*, **35**, 719–733.
- 581 Zhang, Y., Wymant, C., Laeyendecker, O., Grabowski, M. K., Hall, M., Hudelson, S.,
582 Piwowar-Manning, E., McCauley, M., Gamble, T., Hosseinpour, M. C. et al. (2020)
583 Evaluation of phylogenetic methods for inferring the direction of hiv transmission:
584 Hptn 052. *Clinical Infectious Diseases*.

585 **Supplementary material**

586 **S1. Regression models for sampling**

587 Transmission networks in our study are based on three survey rounds in Rakai Com-
 588 munity Cohort Study between August 10, 2011 and January 30, 2015. Despite the
 589 enormous efforts devoted to enumerating all the household, unsampled individuals still
 590 remain in the survey. Specifically, potential sampling bias could be reduced by modelling
 591 approach. Informative priors of sampling proportions for each sampling category were
 592 built through Bayesian logistic regression, and prior samples were then propagated into
 593 the transmission flow estimation. This supplementary text explains how we modelled
 594 the sampling procedure.

595 All the reconstructed transmission events occurred between October 1, 2009 and
 596 January 30, 2015 where the median time of infected was estimated by [CITE]. Thus, we
 597 considered transmission events between eligible individuals in the observational period
 598 during the period between October 1, 2009 and January 30. However, 293 reconstructed
 599 pairs were between successfully sequenced participants who are not on the treatment.
 600 We start by providing an overview of the sampling cascade for eligible recipients.

601 The newly infected eligible individuals between October 1, 2009 and January 30,
 602 2015 are recipients of interests. A fraction of eligible recipients participated. Each
 603 RCCS survey round is preceded by a population census, which enabled us to obtain
 604 posterior estimates of participation probabilities of the form (14) with a Bayesian lo-
 605 gistic regression model implemented in Stan (Supplementary Text S1.2). Ideally, the
 606 probability of an eligible newly infected individuals participating is required. However,
 607 there is no knowledge of the infection status of individuals who were absent. Thus, we
 608 assumed participation was independent of being newly HIV infected for simplicity, and
 609 the probability of eligible individuals participating is an approximation to the targeted
 610 probability. Next, participating sink cases were sequencing-eligible only if they reported
 611 to be ART-naive at the first visit. Finally, a proportion of ART-naive newly infected par-
 612 ticipants were deep sequenced successfully. The final two steps were combined together,
 613 and the aim is to estimate the proportion of newly infected participants being sequenced

614 (Supplementary Text S1.3). Though the time of infection were estimated for sequenced
615 individuals only instead of the newly infected participants, whether an infection is in the
616 infection window could be predicted in Supplementary Text S1.1.

617 We then considered the source cases of new infections. Of all source cases infecting
618 individuals in the cohort, approximately 30% were previously estimated to reside outside
619 the surveyed area (Grabowski PLoS Med) and thus ineligible in RCCS. Of those in
620 the surveyed area, we assumed that all sources were eligible to participate. Next, a
621 fraction of eligible sources participated. Similarly, we assumed infectious individuals
622 participated in the same manner as other eligible individuals, and estimated the posterior
623 distribution of the probability of eligible individuals participating (Supplementary Text
624 S1.2). Next, participating source cases were able to be sequenced if they were ART-
625 naive at the first visit. Assume ART-naive patients are infectious individuals, because
626 all the infected people on the treatment at the first visit were predicted to be infected
627 before 2009.10, which indicates all the participating source cases were sequencing-eligible.
628 Then, a fraction of ART naive individuals who participated were sequenced for virus
629 samples satisfying minimum quality criteria. These prompt us to model the probability
630 of possible infectious participants being sequenced (Supplementary Text S1.4).

631 In applying the sequence-success probabilities to the sampling cascade, we assumed
632 that sequencing success was independent of individuals being a source or not. We further
633 modelled the sampling probabilities of each transmission pair as the product of the
634 sampling probabilities of the source and those of the recipient. Subject to the sampling
635 probability of transmission events, a fraction of transmission events between eligible
636 individuals in the observational window occurring during the infection window were
637 observed.

638 *S1.1. Estimation of infection in 2009.10-2015.1*

639 To obtain whether an infected participant acquired infection between 2009.10 and 2015.1,
640 sequenced individuals were used as the training set due to their known infection time
641 from the estimation by XXX. In addition, the training set included individuals who
642 were tested to be HIV-positive before 2009.10. 948 individuals who were infected before
643 2009.10 and 1691 newly infected individuals were used to predict newly infection status

644 of all other 2386 infected participants through Bayesian logistic regression,

$$\begin{aligned}
 N_i^n &\sim \text{Bernoulli}(1, \xi_i^n), \forall a \\
 \text{logit}(\xi^n) &= \boldsymbol{\beta} \mathbf{X} \\
 \beta_0 &\sim \mathcal{N}(0, 100) \\
 \beta_j &\sim \mathcal{N}(0, 10), j \neq 0
 \end{aligned} \tag{S1}$$

645 where $\boldsymbol{\beta}$ and \mathbf{X} refers to coefficients and design matrix, N_i^n is the indicator of the i th
 646 person were newly infected, and ξ_i^n is the probability of the i th person being newly
 647 infected. Specifically,

$$\begin{aligned}
 \boldsymbol{\beta} &= (\beta_0, \beta_m, \beta_i, \beta_a, \beta_f, \beta_d, \beta_{fp}, \beta_{art}) \in \mathbb{R}^8 \\
 \mathbf{X}_{i \cdot} &= (1, x_{i,m}, x_{i,i}, x_{i,a}, x_{i,f}, x_{i,d}, x_{i,fp}, x_{i,art}) \in \mathbb{R}^8
 \end{aligned}$$

648 where $x_{i,m}$ is the indicator of individual i being male, $x_{i,i}$ is the indicator of individual
 649 i being immigrant, $x_{i,f}$ is the indicator of individual i being in fishing sites, $x_{i,a}$ is the
 650 age of individual i , $x_{i,d}$ is the visit date of individual i in the study, $x_{i,fp}$ is the date of
 651 individual i first found to be HIV-positive and $x_{i,art}$ is the indicator of individual i being
 652 on ART. β_0 is the intercept, $\beta_m, \beta_i, \beta_d, \beta_a, \beta_f, \beta_{fp}$ and β_{art} are regression coefficients for
 653 corresponding input variables. 10-fold cross validation were implemented to assess the
 654 model performance: accuracy of prediction on the test set is 96.96% (95.61%-98.47%),
 655 F1 score for the new infection class is 97.71% (96.59% - 98.90%) and 95.79% (93.82 % -
 656 97.71%). Details of other competing models are in Table XXX.

657 *S1.2. Probability of an eligible individual participating*

658 To obtain the probability of an eligible recently infected (or infectious) individuals par-
 659 ticipating, the probability of an eligible individuals participating in the survey was used
 660 instead due to unknown infection status for absent individuals. However, the approxi-
 661 mation assumes participation is independent of newly infection or infectious status. The
 662 following figures show the differences of participation rates between XXXXXX. Other
 663 competing models are in Table XXX.

664 **INSERT A FIGURE SHOWING THE PARTICIPATION DIFF**

Table S1. Model comparison

Quantities	Model	Accuracy	F1 score (new infection)	F1 score (old infection)
infection between 2009,10 and 2015,1	model 1 (no interaction)			
	model 2			
	(age \times sex interaction)			
	model 3 (age \times sex \times immigrant interaction)			
participation fractions	model 1 (no interaction)			
	model 2			
	(age \times sex interaction)			
	model 3 (age \times sex interaction, dispersion)			
sequencing fractions (source)	model 1 (no interaction)			
	model 2			
	(age \times sex interaction)			
sequencing fractions (recipient)	model 1 (no interaction)			
	model 2			
	(age \times sex interaction)			

665 The number of individuals in a who were eligible N_a^e and participated N_a^p , and pre-
 666 dictors \mathbf{X} enables us to model the proportion of participation $\xi^p = \{\xi_a^p\}_a$ through a
 667 Bayesian logistic model. Specifically,

$$\begin{aligned} N_a^p &\sim \text{BB}(N_a^n, \xi_a^p, \gamma), \forall a \\ \text{logit}(\xi^p) &= \boldsymbol{\beta} \mathbf{X} \\ \beta_0 &\sim \mathcal{N}(0, 100) \\ \beta_j &\sim \mathcal{N}(0, 10), j \neq 0 \\ \gamma &\sim \exp(1) \end{aligned} \tag{S2}$$

668 where $\boldsymbol{\beta}$ and \mathbf{X} refers to coefficients and design matrix, and BB is beta-binomial
 669 distribution with dispersion parameter γ . Specifically,

$$\boldsymbol{\beta} = (\beta_0, \beta_f, \beta_{m1i}, \beta_{m2i}, \beta_{m3i}, \beta_{f1i}, \beta_{f2i}, \beta_{f3i}, \beta_{m1r}, \beta_{m2r}, \beta_{m3r}, \beta_{f1r}, \beta_{f2r}) \in \mathbb{R}^{11}$$

$$\mathbf{X}_{a \cdot} = (1, x_{a,f}, x_{a,m1i}, x_{a,m2i}, x_{a,m3i}, x_{a,f1i}, x_{a,f2i}, x_{a,f3i}, x_{a,m1r}, x_{a,m2r}, x_{a,m3r}, x_{a,f1r}, x_{a,f2r}) \in \mathbb{R}^{11}$$

670 where $x_{a,f}$ is the indicator of group a being in fishing sites, and $x_{a,xyz}$ is 1 if individuals in
 671 group a are gender x , in age group y and with migration status z where x takes m for male
 672 and f for female, y takes 1(15 – 24), 2(25 – 34), 3(35+), and z is i for immigrants and r for
 673 residents. β_0 is the intercept, $\beta_{a,f}$ and $\beta_{a,xyz}$ are regression coefficients for corresponding

674 input variables. The number of participants in 97.73% (90.91% - 100%) groups were
 675 correctly predicted in the 10-fold cross validation. Effective sample sizes were between
 676 4192.857 and 20089.852, and Rhat were between 0.9999312 and 1.0007630. Figure XXXX
 677 shows the density plots of the posterior distribution of regression coefficients.

678 **S1.3. Probability of a newly infected individual being sequenced**

679 A proportion of newly infected participants were not on treatment and eligible to be
 680 sequenced, of whom a proportion were taken sequence samples to the acceptable stan-
 681 dard. These two proportions could be combined into a single quantity, the probability of
 682 a newly infected participant being sequenced successfully ξ^{sr} , which may vary between
 683 subpopulations. Figure XXXX shows **INSERT A FIGURE SHOWING THE DIFF**

684 The number of newly infected participants N^{pr} predicted in Supplementary Text
 685 S1.1, those who were sequenced N^{sr} and attributes of strata were used to estimate ξ^{sr} .
 686 Specifically,

$$\begin{aligned} N_a^{sr} &\sim \text{Bin}(N_a^{pr}, \xi_a^{sr}), \forall a \\ \text{logit}(\xi^{sr}) &= \boldsymbol{\beta} \mathbf{X} \\ \beta_0 &\sim \mathcal{N}(0, 100) \\ \beta_j &\sim \mathcal{N}(0, 10), j \neq 0 \end{aligned} \tag{S3}$$

687 where $\boldsymbol{\beta}$ and \mathbf{X} refers to coefficients and design matrix. Specifically,

$$\boldsymbol{\beta} = (\beta_0, \beta_f, \beta_i, \beta_m, \beta_1, \beta_2) \in \mathbb{R}^6$$

$$\mathbf{X}_{a \cdot} = (1, x_{a,f}, x_{a,i}, x_{a,m}, x_{a,1}, x_{a,2}) \in \mathbb{R}^6$$

688 where $x_{a,f}$, $x_{a,i}$, $x_{a,m}$, $x_{a,1}$, $x_{a,2}$ are taking value one if all the people in group a
 689 are in fishing sites, immigrants, males, aged 15-24, aged 25-34. β_0 is the intercept,
 690 $\beta_f, \beta_i, \beta_m, \beta_1, \beta_2$ are regression coefficients for the input variables. The number of se-
 691 quenced in 95.24% (85.16% - 100%) groups were correctly predicted in the 10-fold cross
 692 validation. Effective sample sizes were between 11833.35 and 31646.49, and Rhat were
 693 between 0.9999662 and 1.0000076. Figure XXXX shows the density plots of the posterior
 694 distribution of regression coefficients.

695 *S1.4. Probability of an infectious individual being sequenced*

696 Assume all the infectious participants were not on treatment at their first visit and el-
 697 igible to be sequenced, and a proportion of those were sequenced. Taking both into
 698 account, we modelled the probability of an ART-naive participant being sequenced suc-
 699 cessfully ξ^{st} . Figure XXXX shows the varying ξ^{st} between strata. **INSERT A FIGURE**
 700 **SHOWING THE DIFF**

701 The number of ART-naive participants N^{pt} , those who were sequenced N^{st} and
 702 attributes of strata were used to obtain the posterior distribution of ξ^{st} . Specifically,

$$\begin{aligned} N_a^{st} &\sim \text{Bin}(N_a^{pt}, \xi_a^{st}), \forall a \\ \text{logit}(\xi^{st}) &= \boldsymbol{\beta} \mathbf{X} \\ \beta_0 &\sim \mathcal{N}(0, 100) \\ \beta_j &\sim \mathcal{N}(0, 10), j \neq 0 \end{aligned} \tag{S4}$$

703 where $\boldsymbol{\beta}$ and \mathbf{X} refers to coefficients and design matrix. Specifically,

$$\begin{aligned} \boldsymbol{\beta} &= (\beta_0, \beta_f, \beta_i, \beta_m, \beta_1, \beta_2) \in \mathbb{R}^6 \\ \mathbf{X}_{a\cdot} &= (1, x_{a,f}, x_{a,i}, x_{a,m}, x_{a,1}, x_{a,2}) \in \mathbb{R}^6 \end{aligned}$$

704 where $x_{a,f}$, $x_{a,i}$, $x_{a,m}$, $x_{a,1}$, $x_{a,2}$ are taking value one if all the people in group a
 705 are in fishing sites, immigrants, males, aged 15-24, aged 25-34. β_0 is the intercept,
 706 $\beta_f, \beta_i, \beta_m, \beta_1, \beta_2$ are regression coefficients for the input variables. The number of se-
 707 quenced in 100% (95.05% - 100%) groups were correctly predicted in the 10-fold cross
 708 validation. Effective sample sizes were between 13289.51 and 32010.51, and Rhat were
 709 between 0.9999662 and 1.0000500. Figure XXXX shows the density plots of the posterior
 710 distribution of regression coefficients.

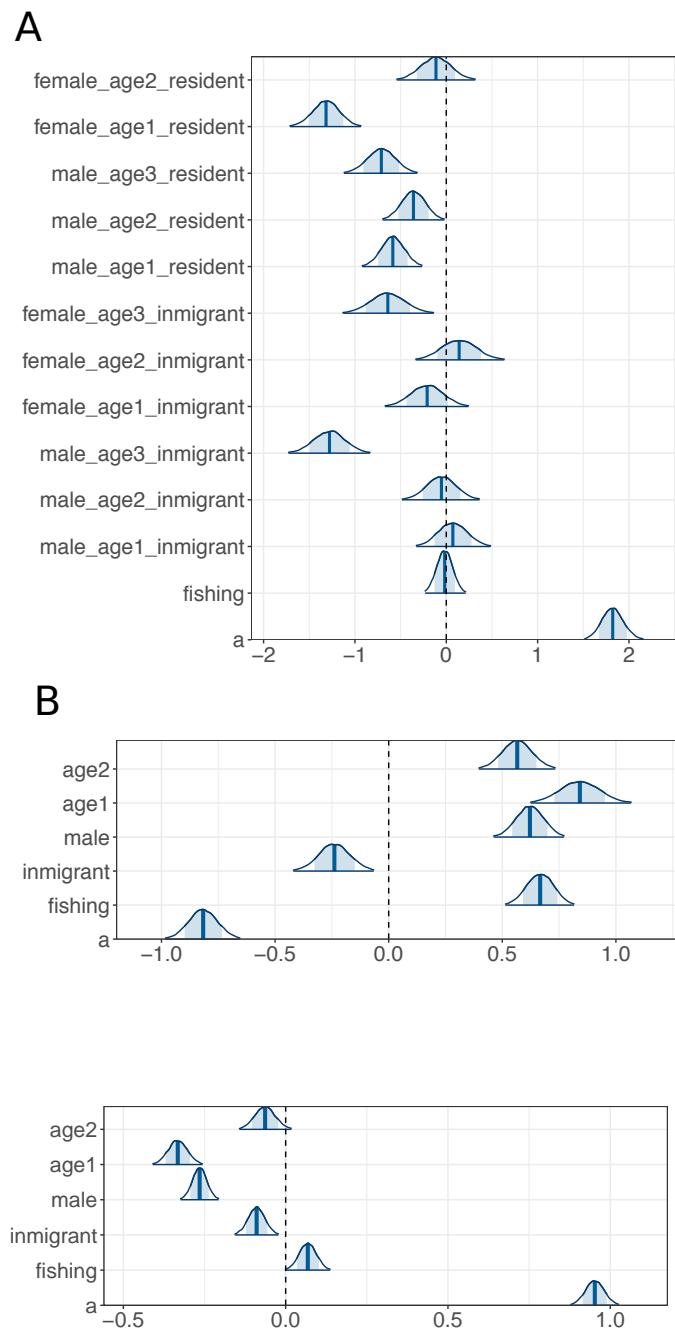


Fig. S1. Densities of the posterior distribution of regression coefficients. Panel A, B, C show the density plots of posterior distributions of parameters in (S2) and (??), (??)

711 **S2. Algorithm**

712 *S2.1. Markov Chain Monte Carlo algorithm*

To estimate the posterior density (13) through Markov Chain Monte Carlo (MCMC), we considered the full conditional distributions

$$p(\xi_a | \boldsymbol{\xi}_{-a}, \boldsymbol{\lambda}, \mathbf{s}, \mathbf{n}, \mathbf{X}) \text{ for all } a, \quad (\text{S5})$$

$$p(\lambda_{ab} | \boldsymbol{\xi}, \boldsymbol{\lambda}_{-ab}, \mathbf{s}, \mathbf{n}, \mathbf{X}), \quad (\text{S6})$$

713 where ξ_a is the sampling probability of the a th population group, λ_{ab} denotes the average
 714 transmission counts from population group a to population group b , $\boldsymbol{\xi}_{-a} = \boldsymbol{\xi} \setminus \xi_a$ and
 715 $\boldsymbol{\lambda}_{-ab} = \boldsymbol{\lambda} \setminus \lambda_{ab}$. The density (S5) is not available in closed form, and we performed
 716 Metropolis-within-Gibbs updates for each ξ_a . The density (S6) is of Gamma form, and so
 717 a Gibbs step can be used to update $\boldsymbol{\pi}$. The resulting MCMC algorithm iterates through
 718 Metropolis-within-Gibbs updates for each ξ_a , followed by a series of Gibbs updates for
 719 each λ_{ab} ; see Supplementary Text S2.1.1 for full details.

This Supplementary Text describes implementation details of the Metropolis-within-Gibbs algorithm to estimate the posterior distribution (13),

$$p(\boldsymbol{\lambda}, \boldsymbol{\xi} | \mathbf{n}, \mathbf{s}, \mathbf{X}) \propto \Pi_{a,b} p(n_{ab} | \lambda_{ab}, \xi_a \xi_b) \Pi_a p(\xi_a | \mathbf{X}_a, \mathbf{s}_a) \Pi_{a,b} p(\lambda_{ab} | \boldsymbol{\xi}).$$

720 with unknown parameters

$\boldsymbol{\xi}$ vector of sampling proportions for each population group a , $a = 1, \dots, A$

$\boldsymbol{\lambda}$ vector of average actual transmission counts from group a to group b , with length L
 721 corresponding to the number of observed counts among all
 722 possible pairwise combinations of counts a, b between population strata

$\boldsymbol{\pi}$ vector of transmission flows from group a to group b , of length L .

722 The algorithm is implemented in the `phyloflows` R package, [https://github.com/](https://github.com/BDI-pathogens/phyloscanner/tree/master/phyloflows)
 723 `BDI-pathogens/phyloscanner/tree/master/phyloflows`. Calculations in this paper
 724 were performed with version 1.2.0.

725 *S2.1.1. Overall structure of algorithm*

726 Owing to the large number of model parameters in the presence of deep sampling heterogeneity, the algorithm exploits the factorisation of the posterior density (13) into the
 727 728 full conditionals (S5-S6) to update parameters in blocks.

729 *S2.1.2. Gibbs step*730 The full conditional distribution $p(\lambda_{ab}|\boldsymbol{\xi}, \boldsymbol{\lambda}_{-ab}, \mathbf{s}, \mathbf{n}, \mathbf{X})$ is

$$\begin{aligned}
p(\lambda_{ab}|\boldsymbol{\xi}, \boldsymbol{\lambda}_{-ab}, \mathbf{s}, \mathbf{n}, \mathbf{X}) &\propto p(\lambda_{ab})p(n_{ab}|\lambda_{ab}, \xi_a, \xi_b) \\
&\propto \exp(-\lambda_{cd}\xi_c\xi_d)(\lambda_{cd}\xi_c\xi_d)^{n_{cd}}\lambda_{cd}^{\alpha_{cd}-1}\exp(-\beta\lambda_{cd}) \\
&\propto \exp(-\lambda_{cd}(\xi_c\xi_d + \beta))\lambda_{cd}^{n_{cd}+\alpha_{cd}-1} \\
&\sim \text{Gamma}(n_{cd} + \alpha_{cd}, \xi_c\xi_d + \beta)
\end{aligned} \tag{S7}$$

731 so that each λ_{ab} can be updated in a single Gibbs step by sampling from the right hand
732 side.733 *S2.1.3. Metropolis-Hastings within Gibbs steps*

734 Next, the algorithm updates in turn the sampling proportions of the source populations.

735 The full conditional distribution of ξ_a follows from (13),

$$\begin{aligned}
p(\xi_a|\boldsymbol{\xi}_{-a}, \boldsymbol{\lambda}, \mathbf{n}, \mathbf{s}, \mathbf{X}) & \\
&\propto \Pi_{c=a \text{ or } d=a} p(n_{cd}|\lambda_{cd}, \xi_c\xi_d)p(\xi_a|\mathbf{X}) \\
&\propto \Pi_{c=a \text{ or } d=a} \text{Poi}(n_{cd}; \lambda_{cd}\xi_c\xi_d)\text{Gamma}(\lambda_{cd}; \xi_c\xi_d)p(\xi_a|\mathbf{X}).
\end{aligned} \tag{S8}$$

736 The full conditional (S8) does not have closed form and so the algorithm performs a
737 Metropolis-Hastings update. We sought to avoid tuning parameters as much as possible,
738 and for this reason adopted the independence proposal density

$$q(\xi'_a|\xi_a) = q(\xi'_a) \tag{S9}$$

739 where ξ'_a is proposed from the sampling prior $p(\xi'_a)$ and

$$\boldsymbol{\xi}' = \begin{cases} \xi_c & c \neq a \\ \xi'_c & c = a \end{cases}$$

740 The resulting Metropolis Hasting ratio is

$$\begin{aligned}
 r &= \frac{p(\xi'_a | \boldsymbol{\xi}_{-a}, \boldsymbol{\lambda}, \mathbf{n}, \mathbf{s}, \mathbf{X})}{p(\xi_a | \boldsymbol{\xi}_{-a}, \boldsymbol{\lambda}, \mathbf{n}, \mathbf{s}, \mathbf{X})} \\
 &\quad \times \frac{q(\xi_a)}{q(\xi'_a)} \\
 &= \frac{\Pi_b p(n_{ab} | \lambda_{ab}, \xi'_a \xi_b) \Pi_{b \neq a} p(n_{ba} | \lambda_{ba}, \xi_b \xi'_a) \Pi_b p(\lambda_{ab} | \boldsymbol{\xi}') \Pi_{b \neq a} p(\lambda_{ba} | \boldsymbol{\xi}') p(\xi'_a)}{\Pi_b p(n_{ab} | \lambda_{ab}, \xi_a \xi_b) \Pi_{b \neq a} p(n_{ba} | \lambda_{ba}, \xi_b \xi_a) \Pi_b p(\lambda_{ab} | \boldsymbol{\xi}) \Pi_{b \neq a} p(\lambda_{ba} | \boldsymbol{\xi}) p(\xi_a)} \\
 &\quad \times \frac{p(\xi_a)}{p(\xi'_a)} \\
 &= \frac{\Pi_b p(n_{ab} | \lambda_{ab}, \xi'_a \xi_b) \Pi_{b \neq a} p(n_{ba} | \lambda_{ba}, \xi_b \xi'_a) \Pi_b p(\lambda_{ab} | \boldsymbol{\xi}') \Pi_{b \neq a} p(\lambda_{ba} | \boldsymbol{\xi}')}{\Pi_b p(n_{ab} | \lambda_{ab}, \xi_a \xi_b) \Pi_{b \neq a} p(n_{ba} | \lambda_{ba}, \xi_b \xi_a) \Pi_b p(\lambda_{ab} | \boldsymbol{\xi}) \Pi_{b \neq a} p(\lambda_{ba} | \boldsymbol{\xi})}
 \end{aligned} \tag{S10}$$

741 Proposed moves are accepted with probability $\min(1, r)$.

742 The sampling fraction in subpopulation a , ξ_a could be separated into that for trans-
 743 mitters and recipients, ξ_a^t and ξ_a^r if different sampling cascades were applied to source
 744 and recipients. Gibbs steps remain same and updates for ξ_a will be divided into updating
 745 ξ_a^t and ξ_a^r sequentially.

746 *S2.1.4. Discussion*

747 (a) The observed phylogenetic data do not provide information on the sampling proba-
 748 bilities and consequently the posterior distributions of the group sampling proba-
 749 bilities will be close to the corresponding prior distributions. For this reason we do not
 750 think that a random walk proposal would confer advantages over the independence
 751 proposal (S9).

752 *S2.2. Hamiltonian Monte Carlo*

753 *S2.2.1. Independent prior*

754 The posterior distribution (13) under the prior (15) could also be alternatively inferred
 755 through Hamiltonian Monte Carlo in Stan. However, it is impossible to incorporate
 756 prior samples in Stan. Instead, we fitted beta distributions to sampling probabilities,
 757 and provided parameters of beta distributions.

```

758 data {
759   int<lower=1> N;  \\ number of observations
760   int Y[N];  \\ response variable
761   // sampling fractions
  
```

```

762     int<lower=1> N_xi; \\ number of sampling categories
763     matrix[N_xi,2] shape; \\ parameters of sampling rates
764     int xi_id_src[N]; \\ sampling category indices for source
765     int xi_id_rec[N]; \\ sampling category indices for recipients
766     real alpha; \\ the parameter of gamma prior
767 }
768 parameters {
769     vector<lower=0>[N] lambda; \\ poisson rate
770     vector<lower=0,upper=1>[N_xi] xi; \\ sampling fractions
771 }
772 transformed parameters {
773     vector[N] lxi_pair; \\ log sampling fractions
774     vector[N] betav; \\ the parameter of gamma prior
775     real beta;
776     for (i in 1:N){
777         lxi_pair[i] = log(xi[xi_id_src[i]]) + log(xi[xi_id_rec[i]]);
778         if (Y[i]==0){
779             betav[i] = 1/(xi[xi_id_src[i]]*xi[xi_id_rec[i]])-1;
780         } else{
781             betav[i] = Y[i]/(xi[xi_id_src[i]]*xi[xi_id_rec[i]]);
782         }
783         beta = sum(betav);
784     }
785 model {
786     for (i in 1:N_xi){
787         target += beta_lpdf(xi[i]|shape[i,1],shape[i,2]);
788     }
789     for (i in 1:N){
790         target += gamma_lpdf(lambda[i]|alpha,beta);
791     }
792     // likelihood including all constants
793     target += poisson_log_lpmf(Y | lxi_pair + log(lambda));
794 }
```

796 S2.2.2. Gaussian process prior

797 The aim to estimate the flow vector under the prior (16),

$$p(\boldsymbol{\lambda}, \boldsymbol{\xi}, \mathbf{f}, \mu, \alpha, \boldsymbol{\rho} | \mathbf{n}, \mathbf{s}, \mathbf{X}) \quad (\text{S11})$$

798 with unknown parameters

- ξ vector of sampling proportions for each population group a , $a = 1, \dots, A$
- λ vector of average actual transmission counts from group a to group b , with length L corresponding to the number of observed counts among all possible pairwise combinations of counts a, b between population strata
- π vector of transmission flows from group a to group b , of length L .
- ⁷⁹⁹ f vector of Gaussian process realisations over the deviation of log average transmission counts from the baseline, of length L .
- μ a scalar representing the baseline of log average transmission counts for all groups.
- α a scalar representing the marginal standard deviation of Gaussian process.
- ρ a vector of length-scales of Gaussian process.
- θ set of Gaussian process hyperparameters, $\theta = \{\alpha, \rho\}$

⁸⁰⁰ The Gaussian process regression is computationally expensive to implement in Stan for
⁸⁰¹ a considerable number of data points while hyperparameters were estimated at the same
⁸⁰² time. Therefore, we used basis function approximation approach proposed by Solin and
⁸⁰³ Särkkä (2014) and implemented by Riutort-Mayol et al. (2020).

⁸⁰⁴ Specifically, we considered the Gaussian process regression under squared exponential
⁸⁰⁵ kernel in Equation (16). Bochner's theorem demonstrates a stationary covariance
⁸⁰⁶ function is a Fourier transform of a positive finite measure which is called spectral den-
⁸⁰⁷ sity if the density exist. In addition, a covariance operator defined for each stationary
⁸⁰⁸ covariance function is translation invariant, and thus expressed as a series of Laplace
⁸⁰⁹ operators, deduced from the polynomial expansion of spectral density. Laplace opera-
⁸¹⁰ tors could represented as an eigenfunction expension with Dirichlet boundary conditions.
⁸¹¹ These leads to the theorem that a stationary covariance function can be represented as

$$k(\mathbf{x}, \mathbf{x}') \approx \sum_{j=1}^{\infty} S_{\theta} \left(\sqrt{\lambda_j} \right) \phi_j(\mathbf{x}) \phi_j(\mathbf{x}'), \quad (\text{S12})$$

⁸¹² where S_{θ} is the spectral density of the stationary covariance function k , θ is the set of
⁸¹³ hyperparameters of k , $\{\lambda_j\}_{j=1}^{\infty}$ and $\{\phi_j(\mathbf{x})\}_{j=1}^{\infty}$ are respectively the j th eigenvalues and
⁸¹⁴ eigenfunctions of the Laplacian operator (Solin and Särkkä, 2014).

⁸¹⁵ The spectral density S_{θ} for a squared exponential kernel k with hyperparameters
⁸¹⁶ $\theta = (\alpha, \rho)$ where α represents marginal standard deviation and ρ denotes length-scales

817 has the form

$$S_\theta(\boldsymbol{\omega}) = \alpha^2 \sqrt{2\pi}^2 \prod_{i=1}^2 \rho_i \exp \left(-\frac{1}{2} \sum_{i=1}^2 \rho_i^2 \omega_i^2 \right), \quad (\text{S13})$$

818 By contrast, the form of the j th eigenvalues and eigenfunctions of the Laplace operator is independent of covariance functions,

$$\begin{aligned} \boldsymbol{\lambda}_j &= \{\lambda_{\mathbf{I}_{jd}}\}_{d=1}^2 = \left\{ \left(\frac{\pi \mathbf{I}_{jd}}{2L_d} \right)^2 \right\}_{d=1}^2 \\ \phi_j(\mathbf{x}) &= \prod_{d=1}^2 \phi_{\mathbf{I}_{jd}}(x_d) = \prod_{d=1}^2 \sqrt{\frac{1}{L_d}} \sin \left(\sqrt{\lambda_{\mathbf{I}_{jd}}} (x_d + L_d) \right) \end{aligned} \quad (\text{S14})$$

820 where $[-L_d, L_d]$ is the domain for the d th dimension and \mathbf{I} denotes the matrix containing 821 indices.

$$\mathbf{I} = \begin{pmatrix} 1 & 1 & \cdots & 2 & 2 & \cdots & 3 & 3 & \cdots & \cdots \\ 1 & 2 & \cdots & 1 & 2 & \cdots & 1 & 2 & \cdots & \cdots \end{pmatrix}^T$$

822 The fact that eigenvalues of the Laplace operator rise with j and spectral densities 823 decrease faster to zero with j for a bounded covariance function suggests Equation (S12) 824 could be truncated to the first m terms,

$$k(\mathbf{x}, \mathbf{x}') \approx \sum_{j=1}^m S_\theta(\sqrt{\boldsymbol{\lambda}_j}) \phi_j(\mathbf{x}) \phi_j(\mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^T \Delta \boldsymbol{\phi}(\mathbf{x}'), \quad (\text{S15})$$

825 where $\boldsymbol{\phi}(\mathbf{x}) = (\phi_j(\mathbf{x}))_{j=1}^m \in \mathbb{R}^m$ is the column vector of eigenfunctions and $\Delta \in \mathbb{R}^{m \times m}$ 826 is the diagonal matrix whose (j, j) th entry is $S_\theta(\sqrt{\boldsymbol{\lambda}_j})$.

827 Assume m_1 and m_2 are the number of basis functions for the first and second dimensions respectively, the number of eigenvalues and eigenfunctions is $m = m_1 \times m_2$. We 829 define a matrix $\mathbf{I} \in \mathbb{R}^{m \times 2}$ containing indices of univariate basis functions. For example, 830 if $m_1 = 3$ and $m_2 = 4$, then

$$\mathbf{I} = \begin{pmatrix} 1 & 1 & 1 & 1 & 2 & 2 & 2 & 2 & 3 & 3 & 3 & 3 \\ 1 & 2 & 3 & 4 & 1 & 2 & 3 & 4 & 1 & 2 & 3 & 4 \end{pmatrix}^T$$

831 That is, the j th eigenvalues $\boldsymbol{\lambda}_j$ is a vector of length 2 whose elements are the \mathbf{I}_{j1} th and 832 \mathbf{I}_{j2} th univariate eigenvalues. The eigenfunction $\{\phi_j(\mathbf{x})\}_j$ is the product of the \mathbf{I}_{j1} th and 833 \mathbf{I}_{j2} th univariate eigenfunctions.

834 Once approximating the kernel function, Gaussian process could be approximated by
 835 m basis functions of Laplacian (Solin and Särkkä, 2014)

$$f(\mathbf{x}) \approx \sum_{j=1}^m (S_\theta(\sqrt{\lambda_j}))^{\frac{1}{2}} \phi_j(\mathbf{x}) \beta_j = \boldsymbol{\phi}(\mathbf{x})^T \Delta^{\frac{1}{2}} \boldsymbol{\beta} \quad (\text{S16})$$

836 where $\beta_j \sim \mathcal{N}(0, 1)$ is the j th element of a column vector $\boldsymbol{\beta}$. Therefore,

$$f(\mathbf{X}) \approx \boldsymbol{\Phi} \Delta^{\frac{1}{2}} \boldsymbol{\beta} = \boldsymbol{\Phi}(\Delta^{\frac{1}{2}} \boldsymbol{\beta}) \quad (\text{S17})$$

where

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \cdots & \phi_m(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_n) & \cdots & \phi_m(\mathbf{x}_n) \end{pmatrix} \text{ and } \Delta = \begin{pmatrix} S_\theta(\sqrt{\lambda_1}) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & S_\theta(\sqrt{\lambda_m}) \end{pmatrix}$$

837 We provided an example code to demonstrate the basis function approximation. Con-
 838 sider seven age groups called "15-19", "20-24", "25-29", "30-34", "35-39", "40-44", "45-49".
 839 Assume $\alpha = 2.5$, $\rho = (12, 9)$, $\mu = -1$ and sampling probabilities are 0.35, 0.45, 0.5,
 840 0.55, 0.5, 0.55, 0.4 for these age groups. We simulated transmission counts between age
 841 groups through (11) and (16), and then re-estimate model parameters. Figure S2AB
 842 shows the posterior distributions of α , ρ , μ and Poisson rates (black) and assumed
 843 parameter values (red).

```
844 \\ eigenvalues
845 functions {
846   vector lambda_nD(real[] L, int[] m, int D) {
847     vector[D] lam;
848     for(i in 1:D){
849       lam[i] = ((m[i]*pi())/(2*L[i]))^2; }
850     return lam;
851   }
852   \\ spectral functions
853   real spd_nD(real alpha, row_vector rho, vector w, int D) {
854     real S;
855     S = alpha^2 * sqrt(2*pi())^D * prod(rho) * exp(-0.5*((rho .* rho) * (w .* w)));
856     return S;
857   }
858   \\ eigenfunctions
859   vector phi_nD(real[] L, int[] m, matrix x) {
860     int c = cols(x);
861     int r = rows(x);
862     matrix[r,c] fi;
863     vector[r] fi1;
864     for (i in 1:c){
865       fi[,i] = 1/sqrt(L[i])*sin(m[i]*pi()*(x[,i]+L[i])/(2*L[i]));
```

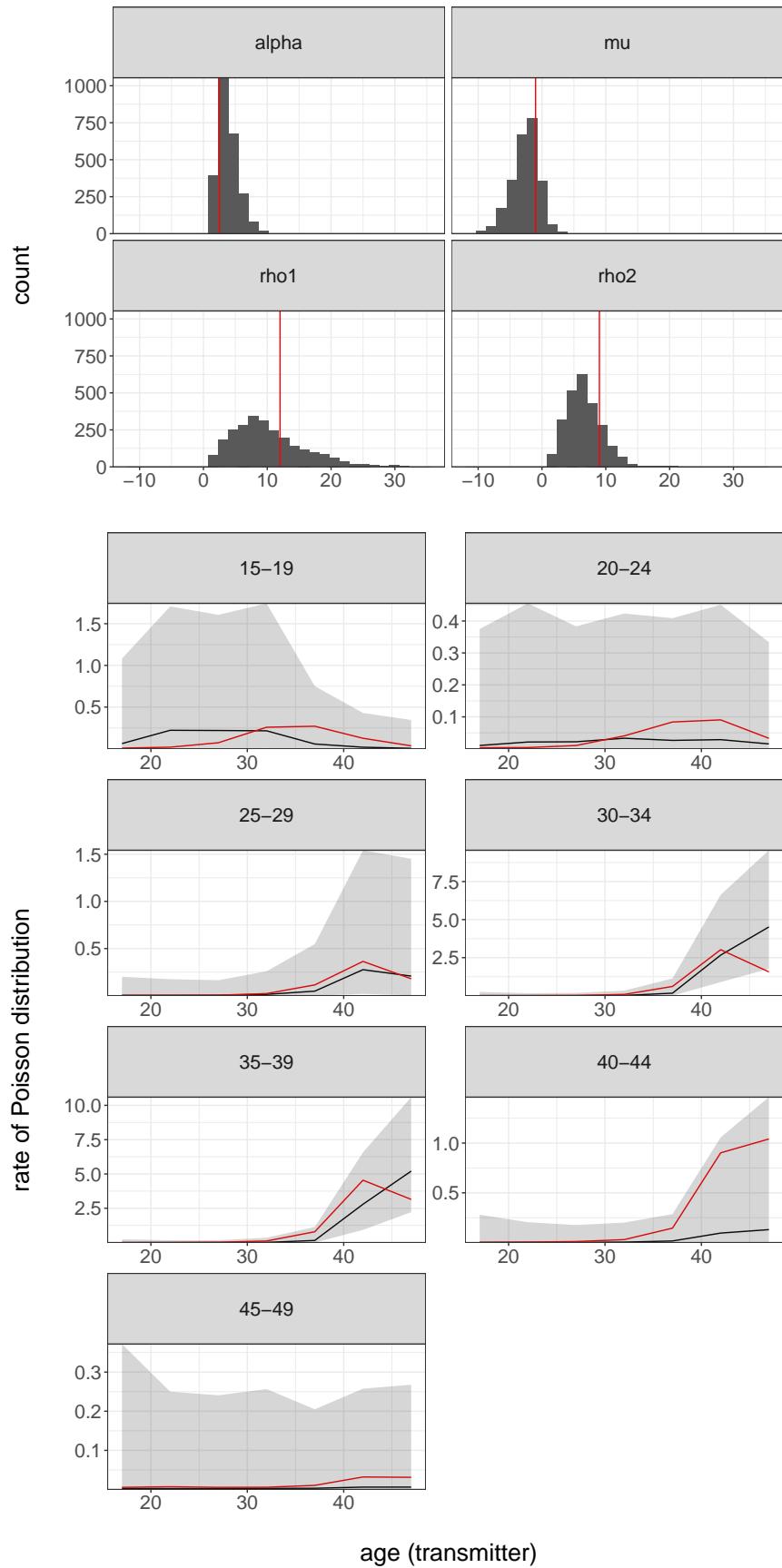


Fig. S2. Demonstration of Gaussian process regression. Panel A compares posterior distributions of model parameters and true values. We simulated transmission count data on the assumption of Gaussian process prior (16) with fixed hyperparameters (vertical red lines). Hyperparameters were re-estimated and posterior distributions are plotted in dark grey area. Panel B compares estimated (black lines and grey shaded areas) and simulated poisson rates (red lines).

```

866     }
867     fi1 = fi[,1];
868     for (i in 2:c){
869         fi1 = fi1 .* fi[,i];
870     }
871     return fi1;
872 }
873 }
874 data {
875     int<lower=1> N_xi; \\ number of sampling categories
876     int<lower=1> N; \\ number of inputs
877     int<lower=1> D; \\ dimension of inputs
878     real L[D]; \\ domain
879     int<lower=1> M; \\ number of basis functions in 1D
880     int<lower=1> M_nD; \\ number of basis functions
881     int indices[M_nD,D]; \\ index matrix S
882     int xi_id[N,2]; \\ sampling category indices
883     matrix[N,D] x; \\ inputs
884     int<lower=0> y[N]; \\ outputs
885     matrix[N_xi,2] shape; \\ parameters of sampling rates
886 }
887 transformed data {
888     matrix[N,M_nD] PHI; \\ eigenfunctions
889     for (m in 1:M_nD){
890         PHI[,m] = phi_nD(L, indices[m,], x);
891     }
892 }
893 parameters {
894     row_vector<lower=0>[D] rho; \\ lengthscale
895     real<lower=0> alpha; \\ marginal standard deviation
896     real mu; \\ baseline
897     vector[M_nD] beta;
898     vector<lower=0,upper=1>[N_xi] xi; \\ sampling proportions
899 }
900 transformed parameters {
901     vector[N] f; \\ gaussian process
902     vector[N] xi1; \\ sampling proportions for transmitters
903     vector[N] xi2; \\ sampling proportions for recipients
904     vector[M_nD] diagSPD; \\ square root of eigenvalues
905     vector[M_nD] SPD_beta; \\ spectral density evaluations
906     for(m in 1:M_nD){
907         diagSPD[m] = sqrt(spd_nD(alpha, rho, sqrt(lambda_nD(L, indices[m,], D)), D));
908     }
909     SPD_beta = diagSPD .* beta;
910     f= PHI * SPD_beta; \\ compute (S16)
911     for (n in 1:N){
912         xi1[n] = xi[xi_id[n,1]]; \\ take sampling proportions for transmitters
913     }
914     for (n in 1:N){
915         xi2[n] = xi[xi_id[n,2]]; \\ take sampling proportions for recipients
916     }
917 }
918 model {
919     for (i in 1:N_xi){
920         xi[i] ~ beta(shape[i,1],shape[i,2]); \\ sampling proportions priors
921     }
922     beta ~ normal(0,1); \\ normal random variable

```

```

930 rho ~ inv_gamma(2.22121,7.04478); \\ lengthscale
931 alpha ~ normal(0, 10); \\ marginal standard deviation
932 mu ~ normal(0, 10); \\ baseline
933 y ~ poisson_log(f + mu + log(xi1) + log(xi2)); \\ poisson model
934 }
```

937 **S2.3. Maximum likelihood estimates of the model**

938 Recall the likelihood of Bayesian flow model is

$$\begin{aligned} l &= \Pi_{ab} \text{Poisson}(n_{ab}; \lambda_{ab} \xi_a^T \xi_b^R) \Pi_{ap} (N_a^{Te}, N_a^{Ts} | \xi_a^T) \Pi_{ap} (N_a^{Re}, N_a^{Rs} | \xi_a^R) \\ &= \Pi_{ab} \text{Poisson}(n_{ab}; \pi_{ab} \eta \xi_a^T \xi_b^R) \Pi_{ap} (N_a^{Te}, N_a^{Ts} | \xi_a^T) \Pi_{ap} (N_a^{Re}, N_a^{Rs} | \xi_a^R) \end{aligned} \quad (\text{S18})$$

939 Taking the derivative of log-likelihood to zero gives maximum likelihood estimates.

$$\begin{aligned} \pi_{ab} \eta &= \frac{n_{ab}}{\xi_a^T \xi_b^R} \\ \xi_a^T &= \frac{N_a^{Ts}}{N_a^{Te}} \\ \xi_a^R &= \frac{N_a^{Rs}}{N_a^{Re}} \end{aligned} \quad (\text{S19})$$

940 As $\sum_{ab} \pi_{ab} = 1$,

$$\begin{aligned} \eta &= \sum_{ab} \frac{n_{ab}}{\xi_a^T \xi_b^R} \\ \pi_{ab} &= \frac{n_{ab}}{\xi_a^T \xi_b^R} / \sum_{ab} \frac{n_{ab}}{\xi_a^T \xi_b^R} \end{aligned} \quad (\text{S20})$$

941 **S3. Simulation experiments**

942 *S3.1. Overall simulation strategy*

943 Simulations were implemented to validate our bias-adjusted method and investigate
 944 errors brought by sampling bias in a classical SIR-type epidemics and Rakai Community
 945 Cohort Study. Our strategy is to simulate transmission counts, re-estimate transmission
 946 flows and compare the estimated flows with the true flows. This Supplementary Text
 947 provides full details on the simulation experiments.

948 *S3.2. 2×2 simulations*

949 The first experiment was a minimal example to assess sampling differences. We con-
 950 sidered transmission flows between two population groups $A = (a, b)$, which for ease of
 951 illustration we refer to as individuals living in rural areas or small communities (group
 952 a), and individuals living in large communities (group b). Assume 60% individuals are
 953 in group b and both groups have the same male and female populations. We simulated
 954 stochastic epidemics between four compartments specified by genders and community
 955 types through structured SIR model (4). ODE parameters were set to make prevalence
 956 approximately 60% and run till equilibrium. Assume people in this epidemics were sam-
 957 pled at the end of the simulation period. Assume the sampling intensity in a was 60%,
 958 and the sampling fraction in b ranged from 60% to 35%, in order to assess the impact of
 959 sampling differences of 0%, 5%, 10% (baseline), 15%, 20%, 25%, from which the number
 960 of sampled individuals could be calculated. We considered transmission events in the
 961 past ten years as they could be reconstructed from sampled individuals. Simulated epi-
 962 demics gave the true values for the 4-dimensional flow vector, $\boldsymbol{\pi}^0 = (\pi_{aa}, \pi_{ab}, \pi_{ba}, \pi_{bb})$,
 963 and observed transmission counts between and within a and b , $\mathbf{n} = (n_{aa}, n_{ab}, n_{ba}, n_{bb})$.
 964 Population sizes were calibrated to ensure roughly 300 transmission events observed.
 965 The vector $\boldsymbol{\pi}$ was inferred when the prior distribution is constructed by Equation (??)
 966 (adjustment for sampling differences) or Beta(25 + 0.5, 50 - 25 + 0.5) (no adjustment
 967 for sampling differences). These prior distributions, together with the assumed true
 968 sampling fractions, were displayed in Figure S3.

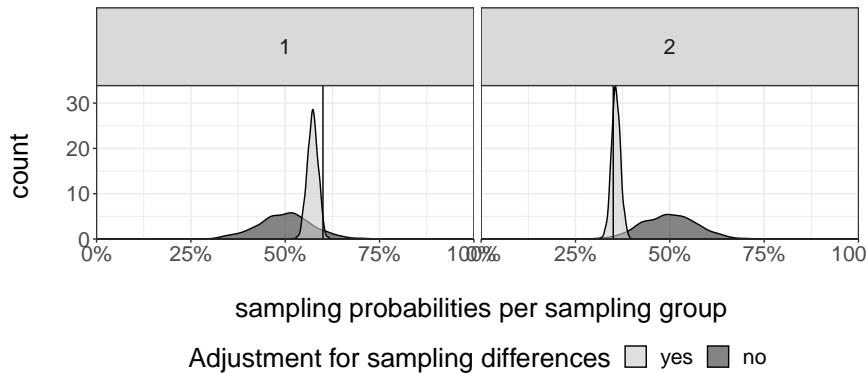


Fig. S3. Prior distributions for the 2×2 simulation. We re-estimated transmission flows with and without adjustment for sampling differences in the 2×2 simulation, where the adjustment is through the specification of prior distributions. Light grey areas represent prior distributions of ξ_a (left panel) and ξ_b (right panel) without adjustment for sampling differences, i.e. Beta (25.5, 25.5). Dark grey areas refer to prior distributions of ξ_a (left) and ξ_b (right) with adjustment for sampling differences, i.e. Equation (??) where X_a^p and X_a^s are informed by sampling data. The vertical line shows the true sampling fractions, $\xi_a = 0.6$ and $\xi_b = 0.35$.

969 **S3.3. Sensitivity to the sample size**

970 The population size was changed to make the number of transmission events to be around
 971 100 and 600 for the baseline scenario in order to assess the robustness of the method to
 972 studies with varying transmission events observed. Figure ?? shows the median WCE
 973 when sampling differences are ignored (light grey bars) compared to when they are not
 974 ignored (dark grey bars) in three studies where 100, 300, 600 transmission events were
 975 observed.

976 **S3.4. Rakai-type simulations**

The second experiment mimicked the substantial sampling heterogeneity observed in the Rakai case study in simulations. We considered the flow vector of length 576 between 24 subpopulations defined by three age brackets, genders, immigration status and community types. We used data on sampling procedure, $N_a^e, N_a^p, N_a^n, N_a^s$, in the Rakai case study and took the estimated flow vector in Section 3.3 as the true flow π_0 . We simulated the observed transmission counts between population groups following the inference procedure. For a fixed sample size $N = 300$, the total number of transmissions Z were simulated from

$$Z \sim \text{Poisson}(N/\bar{\xi})$$

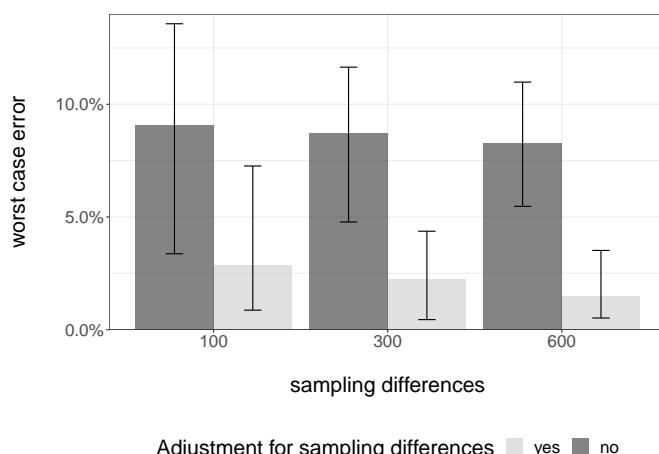


Fig. S4. Sensitivity of bias-adjusted method to varying sample size. The Bayesian flow model (13) was used to re-estimate transmission flows within and between two population groups in simulations, in which the observed flow counts were subject to varying (100, 300, 600) number of observed transmission events in x-axis. In the first set of inference runs (light grey), group-differences in sampling were not adjusted for, by using the same sampling prior density for both groups, $p(\xi_a) = \text{Beta}(25.5, 25.5)$ and $p(\xi_b) = \text{Beta}(25.5, 25.5)$. In the second set of inferences (dark grey), group-differences in sampling were adjusted based on counts of sampled and total infected individuals in both populations, see Equation (??). The worst case error between the true transmission flows and median posterior estimates was calculated on 100 simulations in each scenario, and is shown on the y-axis. With increasing differences in sample sizes, we found increasing bias in flow estimates when sampling differences were not adjusted for, and decreasing bias in flow estimates when sampling differences were adjusted for. Despite capability to reduce bias through bias-adjusted method in all scenarios, the algorithm performs better for larger sample sizes.

where $\bar{\xi}$ was the average sampling probability in the population. The actual transmission flows between population groups followed from

$$\mathbf{z} \sim \text{Multinomial}(Z, \boldsymbol{\pi}_0).$$

The observed transmission counts \mathbf{n} were then simulated under

$$n_{ab} \sim \text{Binomial}(z_{ab}, \xi_a \xi_b)$$

for any population group a and b . The quantity we are interested in is the flow vector between and within inland and fishing communities. The vector $\boldsymbol{\pi}$ was inferred when the prior distribution is constructed by Equation (??) (adjustment for sampling differences) or Beta(25+0.5, 50–25+0.5) (no adjustment for sampling differences), and aggregated to the target flow vector. The difference between sampling intensities in inland and fishing communities is about 10%. We adjusted sampling proportions to make the sampling difference range from 5% to 25% while keeping sampling structure unchanged by using sine-type function. The observed transmission counts \mathbf{n} were then simulated for each case, and used to make inference on the flow vector.

The second simulation experiment closely mimicked sampling heterogeneity as seen in the data from the RCCS observational study (Figure 3B). For a 10% sampling difference between high- and low-prevalence communities, the median WCE was reduced from 6.7%(2.8%–9.8%) to 2.4%(1.1%–4.5%) by considering sampling bias when the sampling difference between inland and fishing communities was 10%.

991 **S4. Numerical performance**

992 This supplementary text presents the numerical performance of the bias-adjusted method.

993 *S4.1. Transmission flows between areas with high and low disease prevalence*

994 The algorithm takes 2.8 hours. The effective sample sizes range from 49910.32 to
 995 105115.13. The acceptance rates from 0.84 to 0.97 for Metropolis-within-Gibbs steps.
 996 Trace plots of three parameters with the the worst performance in terms of effective
 997 sample size are in Figure S6A.

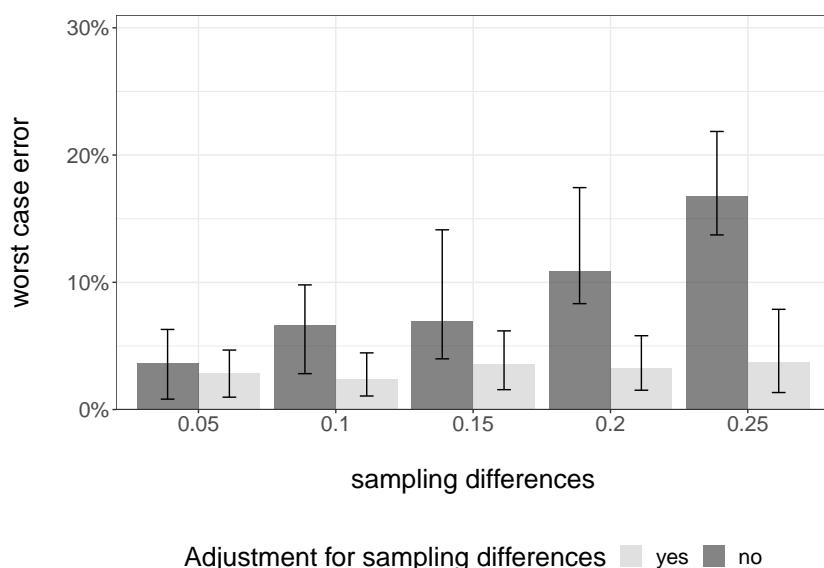


Fig. S5. Comparing errors in estimates of transmission flow obtained with and without adjustment for sampling differences. We assessed bias in estimating transmission flows with model (13) before and after accounting for sampling differences. Transmission flows were simulated between population groups in Rakai analysis. The sampling difference between inland and fishing sites is 10% in Rakai analysis, and was adjusted to range from 5% to 25% in x-axis. The worst case error between the transmission flows in RCCS and median posterior estimates was calculated on 20 simulations in each sampling scenario, and is shown on the y-axis. With increasing differences in sampling differences, we found increasing bias in flow estimates when sampling differences were not adjusted for. When sampling differences were adjusted for, the median worst case error in flow estimates remained on average below 5%.

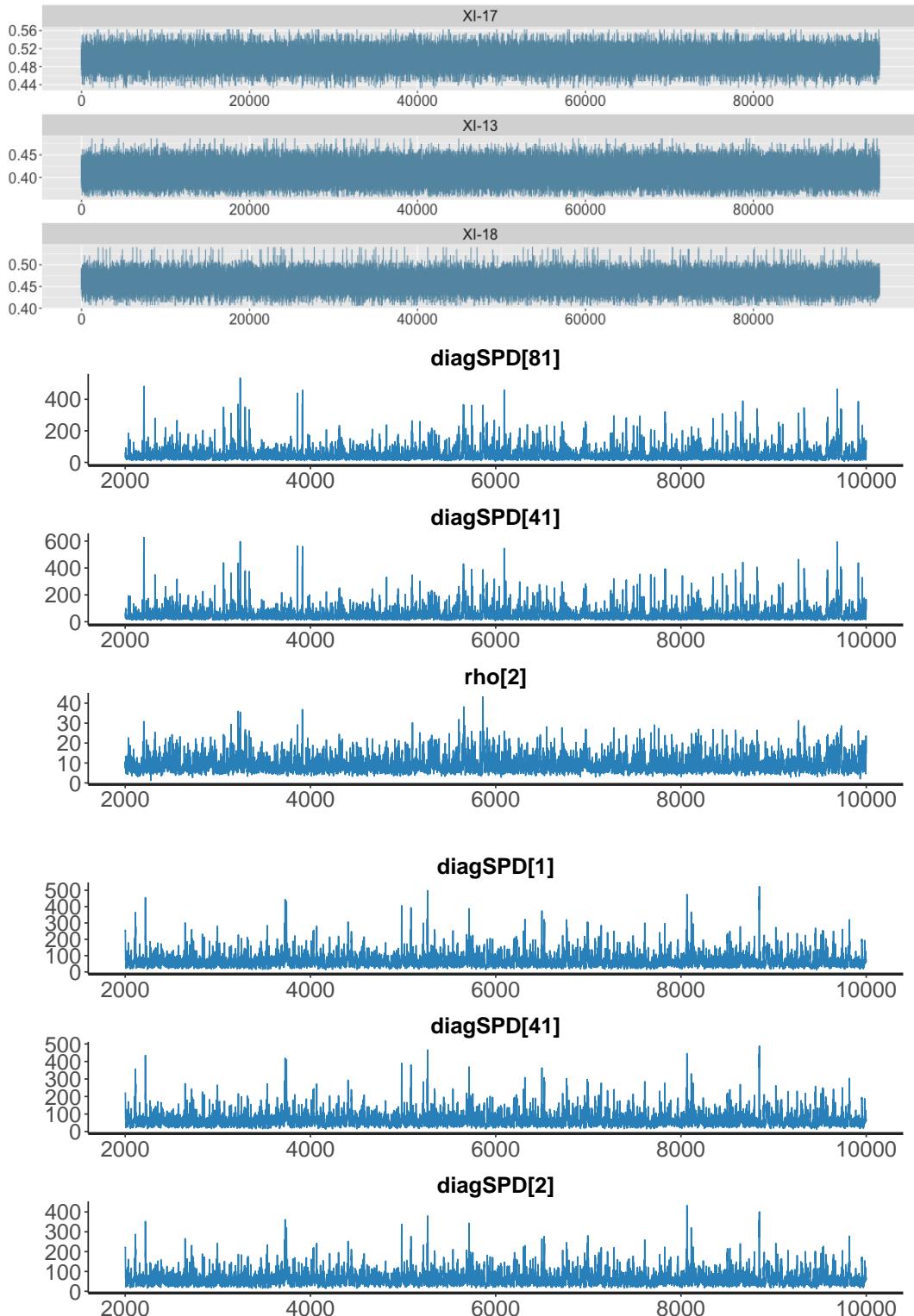


Fig. S6. numerical performance in RCCS. Panel A/B/C provides trace plots of three parameters with lowest effective sample sizes when analysing transmission flows between areas with high and low disease prevalence, flows between age groups from females to males and flows between age groups from males to females respectively.

998 *S4.2. Transmission flows between age groups*

999 The algorithm takes 25.3 hours for transmission pairs from women to men. The effective
1000 sample sizes range from 1620.61 to 28485.75. Trace plots of three parameters with the
1001 the worst performance in terms of effective sample size are in Figure S6B.

1002 The algorithm takes 41.4 hours for transmission pairs from men to women. The effective
1003 sample sizes range from 1486.12 to 20623.32. Trace plots of three parameters with the
1004 the worst performance in terms of effective sample size are in Figure S6C.