# Identification of hidden population structure in time-scaled phylogenies

Erik M. Volz[1,*], Carsten Wiuf[2], Yonatan H. Grad[3], Simon D.W. Frost[4,5], Ann M. Dennis[6], and Xavier Didelot[7]

[1] *Department of Infectious Disease Epidemiology and MRC Centre for Global Infectious Disease Analysis, Imperial College London*

[2] *Department of Mathematical Sciences, University of Copenhagen*

[3] *Department of Immunology and Infectious Diseases, TH Chan School of Public Health, Harvard University*

[4] *Department of Veterinary Medicine, University of Cambridge*

[5] *The Alan Turing Institute*

[6] *School of Medicine, University of North Carolina Chapel Hill*

[7] *School of Life Sciences and Department of Statistics, University of Warwick*

[*] *Corresponding author: Norfolk Place, W2 1PG, United Kingdom; E-mail: e.volz@imperial.ac.uk*

16

## Abstract

Population structure influences genealogical patterns, however data pertaining to how populations are structured are often unavailable or not directly observable. Inference of population structure is highly important in molecular epidemiology where pathogen phylogenetics is increasingly used to infer transmission patterns and detect outbreaks. Discrepancies between observed and idealised genealogies, such as those generated by the coalescent process, can be quantified, and where significant differences occur, may reveal the action of natural selection, host population structure, or other demographic and epidemiological heterogeneities. We have developed a fast non-parametric statistical test for detection of cryptic population structure in time-scaled phylogenetic trees. The test is based on contrasting estimated phylogenies with the theoretically expected phylodynamic ordering of common ancestors in two clades within a coalescent framework. These statistical tests have also motivated the development of algorithms which can be used to quickly screen a phylogenetic tree for clades which are likely to share a distinct demographic or epidemiological history. Epidemiological applications include identification of outbreaks in vulnerable host populations or rapid expansion of genotypes with a fitness advantage. To demonstrate the utility of these methods for outbreak detection, we applied the new methods to large phylogenies reconstructed from thousands of HIV-1 partial *pol* sequences. This revealed the presence of clades which had grown rapidly in the recent past, and was significantly concentrated in young men, suggesting recent and rapid transmission in that group. Furthermore, to demonstrate the utility of these methods for the study of antimicrobial resistance, we applied the new methods to a large phylogeny reconstructed from whole genome *Neisseria gonorrhoeae* sequences. We find that population structure detected using these methods closely overlaps with the appearance and expansion of mutations conferring antimicrobial resistance.

48    Quantifying the role of population structure in shaping genetic

49  diversity is a longstanding problem in population genetics. When information

50  about how lineages are sampled is available, primarily geographic location, a

51  variety of statistics are available for describing the magnitude and role of

52  population structure (Hartl et al. 1997). In pathogen phylogenetics, such

53  geographic 'meta-data' has been instrumental in enabling the inference of

54  transmission rates over space (Dudas et al. 2017), host species (Lam et al.

55  2015), and even individual hosts (De Maio et al. 2018). Population structure

56  shapes genetic diversity, but can the existence of structure be inferred directly

57  from genetic data in the absence of structural covariates associated with each

58  lineage, such as if the geographic location or host species of a lineage is

59  unknown?

60    The problem of detecting and quantifying such 'cryptic' population

61  structure has become a pressing issue in several areas of microbial

62  phylogenetics. For example, in bacterial population genomics studies, a wide

63  diversity of methods have been recently developed to classify taxonomic units

64  based on distributions of genetic relatedness (Mostowy et al. 2017; Tonkin-Hill

65  et al. 2019, 2018; Beugin et al. 2018). In a different domain, pathogen

66  sequence data have been used for epidemiological surveillance, and 'clustering'

67  patterns of closely related sequences have been used to aid outbreak

68  investigations and prioritise public health interventions (Eyre et al. 2012;

69  Dennis et al. 2014; Miller et al. 2014; Ledda et al. 2017). In both population

70  genomics studies and outbreak investigations, a common thread is the absence

71  of variables about sampled lineages that can be correlated with phylogenetic

72  patterns. For example, in outbreak investigations, host risk behaviour and

73  transmission patterns are not usually observed and must be inferred. It is not

74  known a priori which clades are more or less likely to expand in the future,

although there is active research addressing this problem, such as to predict the emergence of strains of influenza A virus (Klingen et al. 2018) or the forecast the effect of antibiotic usage policies on the prevalence of resistant variants (Whittles et al. 2017).

In time-scaled phylogenies, the effects of population structure often appear as a difference in the distribution of branch lengths in clades circulating in different populations (Dearlove and Frost 2015). Figure 1 shows a simulated genealogy from a structured coalescent process (Notohara 1990). In two clades, the effective population size grows exponentially, and in the remaining clade, the effective size remains constant. Consequently, the lineages through time show noticeably different patterns of relatedness. For the clades with growing size, most coalescent events occur in the distant past when the size was small.
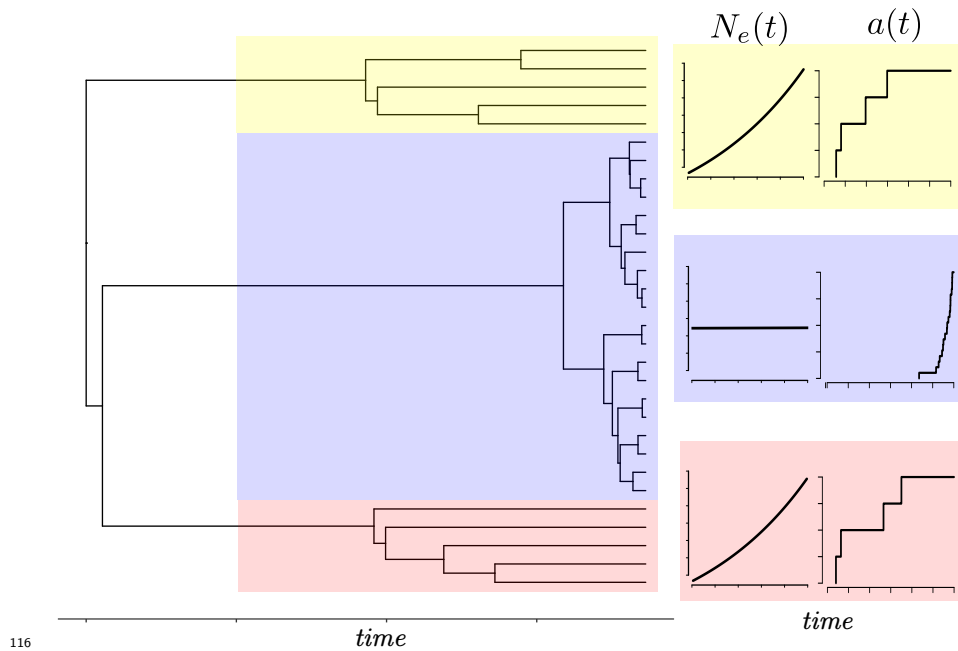
Supposing that the deme from which lineages were sampled was not observed, it is clear from visual inspection of Figure 1 which lineages were sampled from a growing population. Nevertheless, there is a paucity of objective methods readily available to automate the process of identifying temporally distinct clades. This process cannot be done manually when the differences in distributions are less obvious, and needs to be based on a theoretically grounded statistical test. Furthermore, in Figure 1, the red and yellow clades are distantly related. Their most recent common ancestor (MRCA) is at the root of the tree, but they have a very similar distribution of coalescent times suggesting that they were generated by similar demographic or epidemiological processes. For example, this can happen in infectious disease epidemics, when lineages independently colonise the same host population with greater susceptibility or higher risk behaviour (Dearlove et al. 2017). It is therefore also desirable to have an automated method for identifying polyphyletic taxonomic groups defined by shared inferred

102  population histories as opposed to genetic or phenotypic traits.

103  Here we develop a statistical test for detecting if clades within a

104  time-scaled genealogy have evidence for unobserved population structure. Our

105  approach is to develop a statistic based on an unstructured coalescent process.

106  This allows us to test a null hypothesis that two clades are both generated by

107  the same coalescent process. In this case, the coalescent model provides a

108  theoretical prediction of the order of the coalescent times between the two

109  clades in the absence of population structure. On the basis of this statistical

110  test, we also develop algorithms for systematically exploring possible partitions

111  of a genealogy into distinct sets representing evolution within latent

112  populations with different demographic or epidemic histories. Notably, these

113  algorithms not only allow us to detect outlying clades with very different

114  genealogical patterns, but also to find and classify distantly related clades

115  which likely have similar demographic or epidemic histories.

## Materials and Methods

125  As a starting point for our methodology, we assume a time-scaled phylogeny

126  has been estimated from genetic data, for example using one of the recently

127  developed fast methods (To et al. 2016; Volz and Frost 2017; Didelot et al.

128  2018; Sagulenko et al. 2018; Tamura et al. 2018; Miura et al. 2019).

129  Alternatively, summary trees obtained from full Bayesian approaches as

130  implemented in BEAST (Suchard et al. 2018; Bouckaert et al. 2014) or

131  RevBayes (Höhna et al. 2016) can be used, although these typically

132  incorporate population genetic models which presume a particular form of

133  population structure or a lack of population structure. Some precise

134  terminology and notation is required related to the structure of these

135  time-scaled trees since the basis of our approach concerns comparisons

Figure 1: A genealogy simulated from a structured coalescent process with two demes, one of which has constant effective population size (clade highlighted in blue), and the other having effective population size growing exponentially (clades highlighted in red and yellow). Migration of lineages occurs at a small constant rate in one direction from the constant size deme to the growing deme. The corresponding plots at the right show a caricature of the effective population size and lineages through time in each clade.

136   between different subsets of the tree.

## Notation

138   The tree has $n$ terminal nodes (nodes with no descendants), is rooted, and is

139   bifurcating (there are $n - 1$ internal nodes each with exactly two descendants).

140   Being rooted implies there is one node with no ancestor. Mathematically we

141   describe this tree as a node-labelled directed acyclic graph:

$$\mathcal{G} = (\mathcal{N}, \mathcal{E}, \tau)$$

142   where $\mathcal{N}$ is a set of $2n - 1$ nodes, $\mathcal{E} \subseteq \{(u,v)|u,v \in \mathcal{N}^2\}$ is the set of $2n - 2$

143   edges or 'lineages', and $\tau \colon \mathcal{N} \to \mathbb{R}_{\geq 0}$ defines the time of each node. With

144   reference to an edge $(u,v) \in \mathcal{E}$ we say that $u$ is the 'direct ancestor' and $v$ is

145   the 'direct descendant' and we require $\tau(u) < \tau(v)$. Nodes are further

146   classified into two sets: 'tips' (terminal nodes) denoted $\mathcal{T}$ with no descendants

147   and internal nodes denoted $\mathcal{I}$ with exactly two direct descendants. The trees

148   may be heterochronous, meaning that tips of the tree can represent samples

149   taken at different time points.

150        For a node $u \in \mathcal{N}$ we define the clade $C_u$ to be the set of nodes

151   descending from $u$, that is, the node $u$ and all $v \in \mathcal{N}$ such that there is a

152   directed path of edges from $u$ to $v$. We say that nodes $v$ in $C_u$ are 'descended

153   from' $u$. We will also have occasion to define clades 'top down' in terms of a

154   subset of tips in the tree. For this, we define the most recent common ancestor

155   MRCA($X$) of a set $X \subseteq \mathcal{T}$ to be the most recent node $u$ such that $X \subseteq C_u$,

156   that is, all other nodes $v$ with $X \subseteq C_v$ have $\tau(v) < \tau(u)$. Then we let the

157   top-down clade $B_X$ be defined as

$$B_X = \{u \in \mathcal{N} | C_u \cap X \neq \emptyset\}.$$

Note that $B_X$ includes the tips $X$ as well as some nodes ancestral to MRCA($X$).

In general $B_X \neq C_{\mathrm{MRCA}(X)}$ since $X$ does not necessarily include all tips descending from MRCA($X$). We will also need to refer to the nodes corresponding to coalescent events among lineages of the set $X$ only, excluding those between lineages of $X$ and lineages of the complement of $X$,

$$D_X = X \cup \{u \in B_X | \exists (u,v), (u,w) \in \mathcal{E}, v \neq w, C_v \cap X \neq \emptyset, C_w \cap X \neq \emptyset\},$$

Figure 2A illustrates a tree and the sets $B_X, D_X,$ and $C_{\mathrm{MRCA}(X)}$.

Since each node has a time, we can define the set of 'extant' lineages $\mathcal{A}(t)$ at a particular time $t$ to be the set of nodes occurring after time $t$ with a direct ancestor before time $t$,

$$\mathcal{A}(t) = \{v \in \mathcal{N} | \exists (u,v) \in \mathcal{E}, \tau(u) < t \leq \tau(v)\}.$$

We might also refer to the number of extant lineages at time $t$, $a(t) = |\mathcal{A}(t)|$, and if considering the number of extant lineages within a particular clade ancestral to (and including) $X$ we write

$$a_X(t) = |\mathcal{A}(t) \cap B_X|.$$

## Non-parametric test for a given pair of clades

With the above notation, the rank-sum statistic can now be defined which will form the basis for subsequent statistical tests and can be used to compare any pair of clades in the tree.

Let $X$ and $Y$ represent disjoint sets of tips as represented in Figure 2B-D. Having sorted the nodes according to time and assigned a corresponding rank to each internal node, this statistic computes the sum of ranks in a given clade in comparison to a different clade:

$$\rho(X|Y) = \sum_{i=1}^{K} i \, \mathbf{1}_{D_X \setminus D_Y}(w_i), \tag{1}$$

where $S_{X,Y} = (w_1, w_2, \ldots, w_K)$ is the sequence of internal nodes in $D_X \cup D_Y$ sorted by time (present to past), and $\mathbf{1}_A(u)$ is an indicator that takes the value 1 if $u \in A$ and is zero otherwise. Note that $\rho(X|Y)$ is asymmetric in $X$ and $Y$. Also note that $\rho(X|Y)$ makes use of $D_X, D_Y$ and not $B_X, B_Y$ because we are interested in the relative ordering of coalescent events among lineages of $X$ and $Y$. Only the ordering of the events matter, the absolute times are immaterial to the test.

Under a neutral coalescent process, the distribution of coalescent times in two clades ancestral to $X$ and $Y$ will depend on the number of extant lineages through time in both clades and on the effective population size $N_e(t)$ (Wakeley 2009). However, the distribution of the relative ordering of coalescent times only depends on the sizes of the clades. This distribution can be computed rapidly by Monte-Carlo simulation as shown below, provided that we know the probability that the next coalescent will be in $X$ or $Y$ as a function of the number of lineages ancestral to $X$ and $Y$, given by $a_X(t)$ and
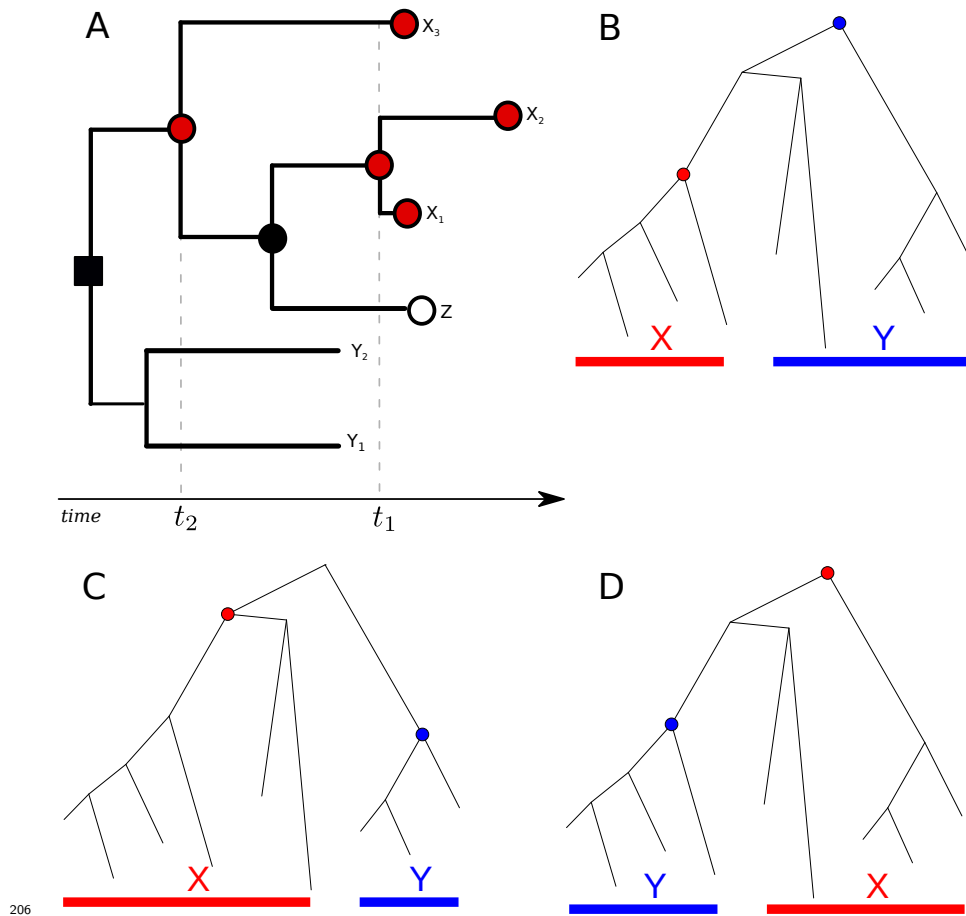
194   $a_Y(t)$. We here provide new theoretical results on the distribution of the

195   relative ordering of coalescence times under the null hypothesis that both $B_X$

196   and $B_Y$ are clades within a single tree generated by a neutral unstructured

197   coalescent process. In the following we consider three different scenarios.

198   **Event $E_1$.** Suppose that a clade $B_X$ has a MRCA before any tip of $X$ shares

199   a common ancestor with the clade of another set of tips $Y$, disjoint to $X$.

200   After lineages in $X$ have found a common ancestor, the MRCA of $X$ may or

201   may not coalesce with lineages in $B_Y$ before $Y$ has found a common ancestor.

202   Figures 2B-C illustrates trees that satisfies this condition. Note that in Figure

203   2B, a lineage in $Y$ coalesces with the MRCA of $X$ before lineages in $Y$ find a

204   MRCA and in Figure B, both $X$ and $Y$ have a common ancestor before they

205   find a common ancestor with one another.

217         Observing a taxonomic pattern such as shown in Figures 2B-C is a

218   random event in a stochastic unstructured coalescent process, and we denote

219   this event by $E_1$ (suppressing $X$ and $Y$ for convenience). Wiuf and Donnelly

220   (Wiuf and Donnelly 1999) showed that the probability of observing $E_1$, given

221   the state of the tree at a particular time $t$, only depends on the number of

222   lineages $z = a_X(t)$ and $w = a_Y(t)$,

$$Q_1(z, w) = \frac{2(z-1)!w!}{(z+w-1)!(z+1)}, \quad z, w \geq 1. \tag{2}$$

223         The numbers of extant lineages in $B_X$ (or its complement) following

224   each coalescent event conditional on $E_1$ is a Markov chain. The transition

225   probabilities of this chain are exactly those needed to simulate the null

226   distribution of the test statistic $\rho(X|Y)$. The probability that the next

227   coalescent event is among lineages in the clade $B_X$ given $E_1$ (starting at a

228   particular time $t$) and the current ancestral number of lineages of $X$, say $z$,

Figure 2: Coalescent trees for illustrating taxonomic relationships and notation used throughout the text. In panel A, the shape and colour of nodes correspond to to variables $B_X, D_X$, and $C_{\mathrm{MRCA}(X)}$ in relation to the set of tips $X = \{x_1, x_2, x_3\}$. All circles regardless of colour correspond to $C_{\mathrm{MRCA}(X)}$. All filled shapes (red or black, square or circle) correspond to $B_X$. Note that this includes nodes ancestral to the MRCA of $X$. All red filled circles correspond to $D_X$. Two coalescent events occur among nodes in $D_X$ at times $t_1$ and $t_2$. Panels B-D show a coalescent tree and examples of potential taxonomic relationships between two clades. Prior knowledge of taxonomic relationships between $X$ and $Y$ influences the probability that the next coalescent event will be observed in clade $X$.

229 and $Y$, say $w$, was found by Wiuf and Donnelly (Wiuf and Donnelly 1999):

$$(z, w) \mapsto (z - 1, w) \quad \text{with probability} \quad \frac{z + 1}{z + w}. \tag{3}$$

230 **Event $E_2$.** We further derive analogous probabilities under slightly different

231 conditions. Suppose we have disjoint sets of tips, $X$ and $Y$. Let all lineages in

232 $X$ share a common ancestor before any share a common ancestor with $Y$ *and*

233 vice versa, all lineages in $Y$ share a common ancestor before any share a

234 common ancestor with tips in $X$. Figure 2C illustrates a tree and two clades

235 that satisfy this condition, which we denote by $E_2$. As before, the number of

236 ancestors in $B_X$ and $B_Y$ will form a Markov chain, conditional on $E_2$.

237 The probability that the next coalescent event is among lineages in

238 the clade $B_X$ given $E_2$ at a particular time $t$ and the current ancestral number

239 of lineages of $X$, $z = a_X(t)$, and $Y$, $w = a_Y(t)$, can be given as:

$$(z, w) \mapsto (z - 1, w) \quad \text{with probability} \quad \frac{z - 1}{z + w - 2}, \quad z, w \geq 1. \tag{4}$$

240 To see this, note that without conditioning on $E_2$, the probability that

241 the next coalescent is among ancestral nodes in $B_X$ is

$$\frac{z(z - 1)}{(z + w)(z + w - 1)}.$$

242 This is simply the ratio of the coalescent rates in $B_X$, which is $\binom{z}{2}/N_e(t)$, to

243 the rate in $B_X \cup B_Y$, which is $\binom{z+w}{2}/N_e(t)$. The effective population size is

244 homogenous through the tree by hypothesis of the statistical test, and it

245 cancels out in this ratio. The probability that the coalescent event would be

246 between the clades ancestral to $X$ and $Y$ would be

$$\frac{2zw}{(z+w)(z+w-1)}.$$

247 The probability $Q_2(z,w)$ of the event $E_2$ must fulfil the recursion,

$$(z+w)(z+w-1)Q_2(z,w)$$
$$= \quad z(z-1)Q_2(z-1,w) + w(w-1)Q_2(z,w-1), \qquad (5)$$

248 where $z, w \geq 1$. If there is exactly one lineage in both $B_X$ and $B_Y$, then

249 $Q_2(1,1) = 1$. If there is one lineage remaining in $B_X$ and $w > 1$ in $V_Y$, then

250 $Q_2(1,w)$ is the probability that the next $w-1$ coalescent events only occur

251 between lineages in $B_Y$ and do not include the single lineage ancestral to $X$.

252 The probability of the next coalescent event being in $B_Y$ is the probability of

253 not selecting the $B_X$ lineage when sampling two extant lineages without

254 replacement:

$$Q_2(1,w) = \prod_{j=2}^{w} \left(\frac{j}{j+1}\right)\left(\frac{j-1}{j}\right)$$
$$= \frac{2}{w(w+1)}, \quad w \geq 1. \qquad (6)$$

255 Similarly, $Q_2(z,1) = \frac{2}{z(z+1)}, z \geq 1$. This recursion can be solved explicitly to

256 give

$$Q_2(z,w) = \frac{2z!w!}{(z+w)!(z+w-1)}, \quad z,w \geq 1. \qquad (7)$$

257 Now the transition probability (Equation 4) can be defined in terms of the

rate of coalescence in $B_X$ and $B_Y$ and the probability of $E_2$ being satisfied following the coalescent event:

$$(z, w) \mapsto (z - 1, w) \quad \text{with probability}$$
$$\frac{z(z-1)Q_2(z-1, w)}{z(z-1)Q_2(z-1, w) + w(w-1)Q_2(z, w-1)} = \frac{z-1}{z+w-2}. \quad (8)$$

**Event $E_3$.** Finally, we consider an event that is a combination of the scenarios described by events $E_1$ and $E_2$. We denote $E_3$ to be the event that all $X$ have a MRCA before sharing a common ancestor with lineages of $Y$ and/or *vice versa*, all lineages in $Y$ have a MRCA before sharing an ancestor with lineages of $X$. The trees in Figures 2B-D satisfy this condition.

The probability of the event $E_3$ might be defined in terms of $Q_1$ and $Q_2$ given previously:

$$Q_3(z, w) = Q_1(z, w) + Q_1(w, z) - Q_2(z, w)$$
$$= \frac{2z!w!}{(z+w-1)!} \left( \frac{1}{z(z+1)} + \frac{1}{w(w+1)} - \frac{1}{(z+w)(z+w-1)} \right), \quad (9)$$

with $z = a_X(t)$ and $w = a_Y(t)$ being sample sizes at a particular time $t$, as before. Note that $Q_2$ is subtracted once in this equation because the taxonomic relationship described by $E_2$ is already included in $E_1$. The function $Q_3$ satisfies the same recursion as above (Equation 5) with slightly different boundary conditions:

$$Q_3(1, w) = Q_3(z, 1) = 1, \quad z, w \geq 1.$$

Transition probabilities can be derived as above by substituting $Q_3$ for $Q_2$ in Equation 8. The probability that the next coalescent event is among lineages

in $D_X$ conditional on $E_3$ is

$$(z, w) \mapsto (z - 1, w) \quad \text{with probability} \quad \frac{(z - 1)R_{z-1,w}}{(z - 1)R_{z-1,w} + (w - 1)R_{w,z-1}}, \quad (10)$$

where

$$R_{z,w} = \frac{1}{z(z + 1)} + \frac{1}{w(w + 1)} - \frac{1}{(z + w)(z + w - 1)}, \quad z, w \geq 1. \quad (11)$$

## Algorithms for detecting population structure

The null distribution of the test statistic $\rho(X, Y)$ can be computed by Monte-Carlo simulation using Equations 3, 4 or 10 depending on the taxonomic constraints to be conditioned on. This can be computed given any pair of disjoint clades $X$ and $Y$. Algorithm 1 in the Supporting Information provides the simulation procedure for computing the two-sided p-values of an empirical measurement $\hat{R} = \rho(X, Y)$, and we denote these p-values $\xi(X, Y, R)$. The algorithm works by simulating many replicates of the rank-sum statistic conditional on the sets $X$, $Y$, and the taxonomic relationship between these clades. Furthermore, the order of sampling events and coalescent events is part of the data within a time-scaled phylogeny. Thus the simulation procedure does not simulate coalescent trees per se, but rather the number of lineages through time $a_X(t)$ and $a_Y(t)$ by proceeding from the most recent sample back to the MRCA of clades $X$ and $Y$. Upon visiting a node in the ordered sequence of coalescent events, the algorithm selects at random a clade $D_X$ or $D_Y$ for this event using the transition probabilities from Equations 3, 4 or 10. Upon visiting a coalescent event, $a_X(t)$ or $a_Y(t)$ is incremented using the observed clade membership of the sample at that time. The end result of of this simulation procedure is a large set of replicate rank-sum statistics which serves as a null distribution for comparison with the value computed from the

295   time-scaled phylogeny.

296        While in principle this test allows comparison of any pair of disjoint

297   clades, the number of possible comparisons is vast, and deriving a useful

298   summary of taxonomic structure requires additional heuristic algorithms.

299   These algorithms are designed to stratify clades into self-similar sets and to do

300   so in a computationally efficient manner. Algorithm 2 in the Supporting

301   Information identifies 'cladistic outliers', which are clades that have a

302   coalescent pattern that is different from the remainder of the tree. It performs

303   a single pre-order traversal of the tree and greedily adds clades to the partition

304   with the most outlying values of the test statistic. At each node $u$ visited in

305   pre-order traversal, Algorithm 2 examines all descendants $v$ in $C_u$ and

306   compares $C_v$ with to $C_u \setminus C_v$. If no outliers are found, the algorithm will desist

307   from searching $C_u$ and the set of tips $C_u \cap \mathcal{T}$ will be added to the partition. If

308   at least one outlier is found in $C_u$, a search will begin on the biggest outlier

309   (smallest p-value computed using Algorithm 1). The final result of this

310   algorithm is a partition of $m$ non-overlapping clades $M = \{X_1, \cdots, X_m\}$.

311        In practice, it is often desirable to not compare very small clades

312   against one another or much larger clades, so additional parameters are

313   available to desist the pre-order traversal upon reaching a clade with few

314   descendants. It is also often of practical interest to only compare clades that

315   overlap in time to a significant extent, so yet another parameter is available to

316   desist from comparing a pair of clades if few lineages in the pair ever coexist at

317   any time.

318        Additional algorithms are required to detect polyphyletic relationships

319   as depicted in Figure 1 which arise if, for example, distantly related lineages

320   colonise the same area and have similar population dynamics or if

321   near-identical fitness-enhancing mutations occur independently on different

lineages. Figure 1 depicts two distantly related clades (yellow and red) with similar population dynamics, and it is desirable to classify these as a single deme based on shared population dynamic history. Algorithm 2 will partition tips of the tree into distinct clades with monophyletic or paraphyletic relationships, however an approach based on pre-order traversal of the tree can not on its own arrive at a polyphyletic partition of the tree. Therefore we can implement a final hierarchical clustering step in order to group similar clades as follows:

1. For each distinct pair of clades $X$ and $Y$ in partition $M$, compute
   $q_{XY} = \xi(X, Y, \hat{R}_{XY})$.

2. Convert the p-value into a measure of distance between all clades:
   $d_{XY} = |F^{-1}(q_{XY})|$, where $F^{-1}$ is the inverse Gaussian cumulative distribution function (quantile function). Set $d_{XX} = 0$ for all $X$.

3. Perform a conventional hierarchical clustering using a threshold distance $F^{-1}(1 - \alpha/2)$ for confidence level $\alpha$. Various clustering algorithms can be used at this point, and our software has implemented the 'complete linkage' algorithm (Everitt et al. 2001).

Algorithms 1 and 2 as well as the final hierarchical clustering step are implemented as an open source R package called *treestructure* available at https://github.com/emvolz-phylodynamics/treestructure. The R package supports parallelisation and includes facilities for tree visualisation using the *ggtree* package (Yu et al. 2017). The package provides convenience functions to output cluster and partition assignment for downstream statistical analysis in R.
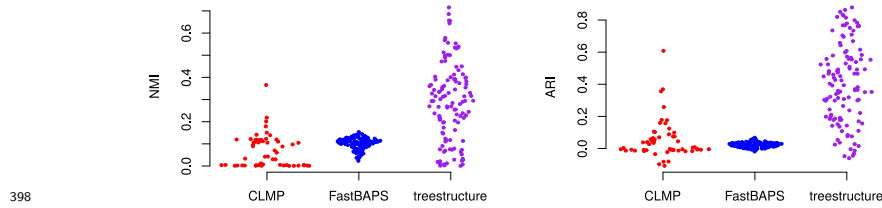
## Simulation studies

To evaluate the potential for *treestructure* to detect outbreaks we applied the new method to phylogenies estimated from newly simulated data using a structured coalescent model, as well as previously published simulation data based on a discrete-event branching process (McCloskey and Poon 2017).

The structured coalescent simulation was based on a model with two demes: a large deme with constant effective population size and a smaller deme which grows exponentially up to the time of sampling. Migration occurs at a constant rate in both directions between the growing and constant-size demes, and equal proportions of these two demes are sampled. Coalescent simulations were implemented using the *phydynR* package http://github.com/emvolz-phylodynamics/phydynR. All genealogies simulated from this model were comprised of 1000 tips with 200 of these sampled from the growing deme. Each of 100 simulations were based on different parameters such that there was a spectrum of difficulty identifying population structure from the trees. The sample proportion was chosen uniformly between 5% and 75% and, the growth rate in the growing deme was chosen uniformly between 5% and 100% per year. Bidirectional migration between demes was fixed at 5% per year. While most tips were sampled at a single time point, 50 tips from the constant-size deme were distributed uniformly through time in order to facilitate molecular clock dating. Multiple sequence alignments were simulated based on trees using seq-gen (Rambaut and Grass 1997). Each sequence comprised 1000 nucleotides from a HKY model with a substitution rate of $10^{-3}$ per site per year. A neighbour joining tree was estimated from each alignment and dated phylogenies estimated using the *treedater* R package (Volz and Frost 2017) with a strict molecular clock. The *treestructure* algorithm was applied to each phylogeny using the default
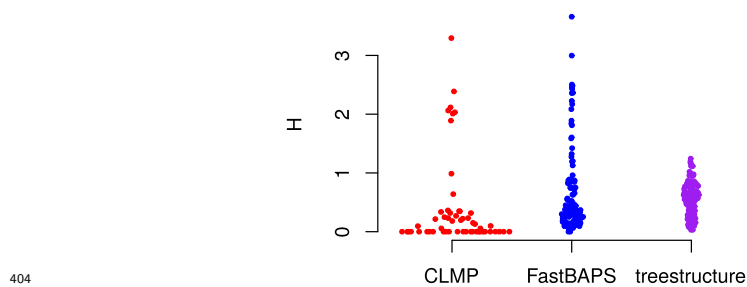
$\alpha = 1\%$ threshold.

Previously, McCloskey et al. simulated 100 genealogies from a discrete-event birth-death process (McCloskey and Poon 2017; Vaughan and Drummond 2013). These simulations were based on a process with heterogeneous classes of individuals with different birth rates. With some probability, lineages migrate to a class with higher birth rates. This could represent a generic outbreak scenario such as a set of individuals with higher risk behaviour or other exposures. In a separate set of simulations, the outbreak population differs from the main population along multiple dimensions: the birth rate and the sampling rate are both increased by a common factor ($5\times$). 100 genealogies were simulated under both scenarios and the *treestructure* algorithm was applied to each. To create more challenging conditions for the method and to evaluate the sensitivity of the method to sample coverage, we also applied the method to genealogies based on subsampled lineages with a frequency of 25%. Complete descriptions of parameters and simulation methods can be found in (McCloskey and Poon 2017).

The performance of *treestructure* was evaluated using the normalised mutual information (NMI) statistic and adjusted Rand index (ARI) computed using the *aricode* R package `https://github.com/jchiquet/aricode` (Vinh et al. 2010). Both statistics quantify the strength of association between the estimated and actual structure of the tree, with larger values corresponding to higher quality reconstructions.

Figure 3: The normalised mutual information (NMI) and adjusted Rand index (ARI) as a function of classifications from several tree-partitioning algorithms and membership of lineages in outbreaks or a constant-size reservoir. Each point corresponds to a structured coalescent simulation where 20% of tips are sampled from an exponentially growing outbreak.



Figure 4: Entropy ($H$) of classification from several tree partitioning algorithms applied to the structured coalescent simulations but only counting lineages sampled from the exponentially growing outbreak.
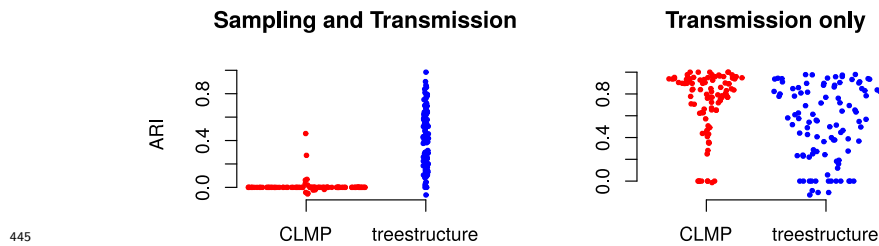
# Results

## Simulation studies

The *treestructure* algorithm achieves relatively high fidelity of classifications in comparison to other methods in the structured coalescent simulations which included 20% of samples from a rapidly growing outbreak. Figure 3 compares three methods according to the NMI and ARI statistics. In these figures, the partition of the tree computed by each method is compared to the true membership of each sampled lineage in outbreak or in the constant-size reservoir population. Across 100 simulations, *treestructure* has mean ARI of 41% (IQR: 20-57%). The FastBAPS method (Tonkin-Hill et al. 2019) has mean ARI of 2.3% (IQR:1.2-3.3%) and the CLMP method (McCloskey and Poon 2017) has mean ARI 5.2% (IQR:-1-7.5%).

The relatively lower performance of CLMP and FastBAPS in these comparisons is largely a consequence of false-positive partitioning of samples from the reservoir population, but CLMP and FastBAPS usually correctly identify a clade that closely corresponds to the outbreak. In contrast, the *treestructure* method seldom sub-divides clades from the reservoir. Figure 4 compares the entropy of partition assignments only within lineages sampled from the outbreak. This shows that all methods are assigning outbreak lineages to a small number of partitions and no method is clearly superior by this metric. The CLMP method has the lowest entropy (mean 40%) but also several large outliers. *treestructure* has higher entropy (mean 57%) but few outliers, and FastBAPS is intermediate.

The performance of all methods depended on the sample density and growth rate of the outbreak. Fast growing outbreaks are easier to detect by all methods but the role of sample density is more ambiguous. The Pearson correlation of ARI with growth rate is $\rho =53\%$, 27%, 71% for *treestructure*,

Figure 5: The adjusted Rand index for 100 previously published simulations (McCloskey and Poon 2017). This describes accuracy of classification of tips into outbreaks using the *treestructure* method and CLMP. Results on left were based on simulations where both transmission and sampling rates varied in the outbreak cluster, whereas simulations on the right only allowed transmission rates to vary.

CLMP and FastBAPS respectively. Not all methods are equally sensitive to these parameters however and FastBAPS is especially sensitive to growth and sample density. The growth rate and sample density collectively explain 41%, 28% and 60% of variance of ARI in *treestructure*, CLMP and FastBAPS respectively.

We also performed analyses with Phydelity (Han et al. 2018) (results not shown), a recently proposed method for transmission cluster identification. This tended to generate a very large number of clusters, both within and outside of the outbreak demes, reflecting a different emphasis of this method on finding closely related clusters rather than addressing differences in macro-level population structure. Thus, results with Phydelity and other clustering methods were not easily comparable to *treestructure*.

Figure 5 shows performance of *treestructure* on previously published tree simulations (McCloskey and Poon 2017). These simulations differ from the structured coalescent simulations because both the reservoir and outbreak demes are growing exponentially at different rates. The birth rate in the

456 outbreak deme is 5-fold the birth rate in the reservoir, but in one set of

457 simulations, both the birth rate and sampling rate in the outbreak was also

458 increased 5-fold. In these simulations, the performance of *treestructure*

459 (median ARI 56%) is slightly lower than the CLMP method (McCloskey and

460 Poon 2017) (median ARI 83%) when only the birth-rate differs in the

461 outbreak deme. However *treestructure* maintains good performance when

462 death and sampling rates also differ. In that case, *treestructure* has median

463 ARI 42% and CLMP has median ARI 0%. The difficulty of detecting

464 outbreaks with different sampling patterns was previously highlighted as a

465 challenge for CLMP (McCloskey and Poon 2017).

## 466 Clonal expansion of drug-resistant *N. gonorrhoeae*

467 We examined the role of evolution of antimicrobial resistance in shaping the

468 phylogenetic structure of *N. gonorrhoeae* using 1102 previously described

469 whole genome sequences (Grad et al. 2016). These isolates were collected from

470 multiple sites in the United States between 2000 and 2013 and featured clonal

471 expansion of lineages with antibiotic resistance to different classes of

472 antibiotics. We estimated a maximum likelihood tree using *PhyML*(Guindon

473 et al. 2010) and corrected for the distorting effect of recombination using

474 *ClonalFrameML* (Didelot and Wilson 2015). We estimated a rooted

475 time-scaled phylogeny using *treedater* (Volz and Frost 2017). A relaxed clock

476 model was inferred, with a mean rate of $4.6 \times 10^{-6}$ substitutions per site per

477 year. *BactDating* (Didelot et al. 2018) was also applied for the same purpose

478 and found to give very similar estimates for the clock rate and dating of clades.

479 We focus on the origin and expansion of two clades which

480 independently developed resistance to cefixime (CFX) by acquiring the mosaic

481 *penA* XXXIV allele (Grad et al. 2016). These clades are indicated in Figure 6.

482 Note, however, that the level of susceptibility to CFX varies, particularly in

483 the larger of these two clades. In one lineage within this clade, the mosaic

484 *penA* XXXIV allele was replaced by recombination with an allele associated

485 with susceptibility. Other isolates within this clade gained mutations that

486 further modified the extent of resistance. The larger of these two clades

487 emerged on a genomic background that is resistant to ciprofloxacin (CIP), so

488 that it has reduced susceptibility to both CIP and CFX. The smaller of the

489 two clades is resistant to CFX but not CIP. We therefore extracted a tree with

490 just 576 tips, representing the genomes from these two CFX-resistant clades as

491 well as genomes from the two clades that are most closely related to the two

492 CFX-resistant clades. The output of *treestructure* is shown in Figure 6, using

493 unique colours to highlight each of the 11 clusters that were identified with

494 $\alpha = 1\%$. The clusters reported by *treestructure* are highly correlated with

495 CFX resistance. Among all distinct pairs of sampled isolates, 84% share the

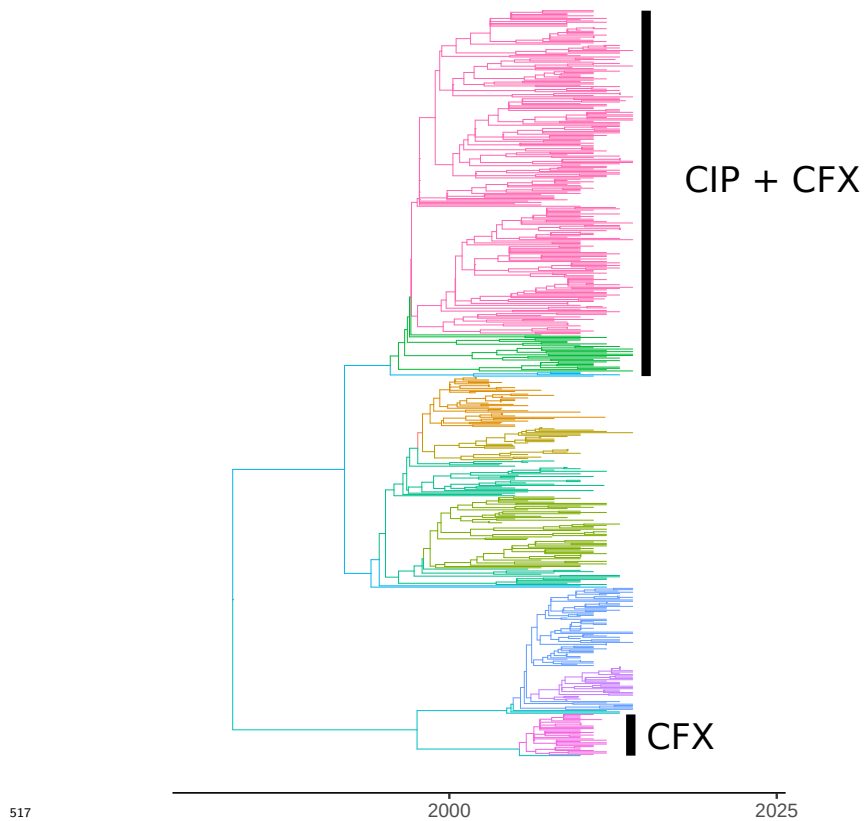496 same resistance profile and cluster membership.

497 We compared *treestructure* with a different method for detecting

498 community structure, FastBAPS (Tonkin-Hill et al. 2019), since BAPS models

499 are often applied to bacterial pathogens. We applied FastBAPS using the

500 same time-scaled phylogeny described previously and using a trimmed

501 sequence alignment consisting of 38830 polymorphic sites and removing sites

502 with many gaps. This produced a similar partition of the tree (Figure S2)

503 with a few differences: FastBAPS clusters overlap exactly with the clade

504 featuring dual resistance (CIP and CFX), whereas *treestructure* classified a

505 small number of deep-splitting lineages into a different cluster. Note however

506 that this behaviour is not necessarily problematic, and may represent a

507 progressive increase in fitness following the acquisition of resistance through

508 the evolution of compensatory mutations (Didelot et al. 2016). On the other

509 hand, FastBAPS failed to identify the smaller clade with resistance to CFX

510 and not CIP and instead grouped that clade with its drug-sensitive sister

511 clade. In general, *treestructure* found many more clusters within the sister

512 clades and FastBAPS tended to group these together. We also applied the

513 much more computationally intensive RhierBAPS method (Tonkin-Hill et al.

514 2018), and obtained almost identical results to FastBAPS. Overall, BAPS

515 methods appear to give greater weight to long internal branches when

516 identifying clusters than *treestructure*.

## Epidemiological transmission patterns of HIV-1

524 We reanalysed a time-scaled phylogeny reconstructed from 2068 partial *pol*

525 HIV-1 subtype B sequences collected from Tennessee between 2001 and 2015

526 (Dennis et al. 2018). Each lineage within this phylogeny corresponds to a

527 single HIV patient sampled at a single time point, and various clinical and

528 demographic covariate data concerning these patients can be associated with

529 each lineage. In the original study, these sequence data were used to show high

530 rates of transmission among young (age $< 26.4$) men who have sex with men

531 (MSM) (Dennis et al. 2018). Clustering by threshold genetic distance is often

532 used in HIV epidemiology (Dennis et al. 2014) and indicated that young white

533 MSM had the highest odds of clustering.

534 We applied the *treestructure* algorithm with default settings to the

535 time-scaled tree which yielded ten partitions with sizes ranging from 58 to 398.

536 The tree and partitions are shown in Figure 7 where partitions are labeled

537 according to the median year of birth among patients in each partition. Many

538 of these partitions were polyphyletic, suggesting possible multiple importations

539 of lineages to specific risk groups. We then compared the estimated partition

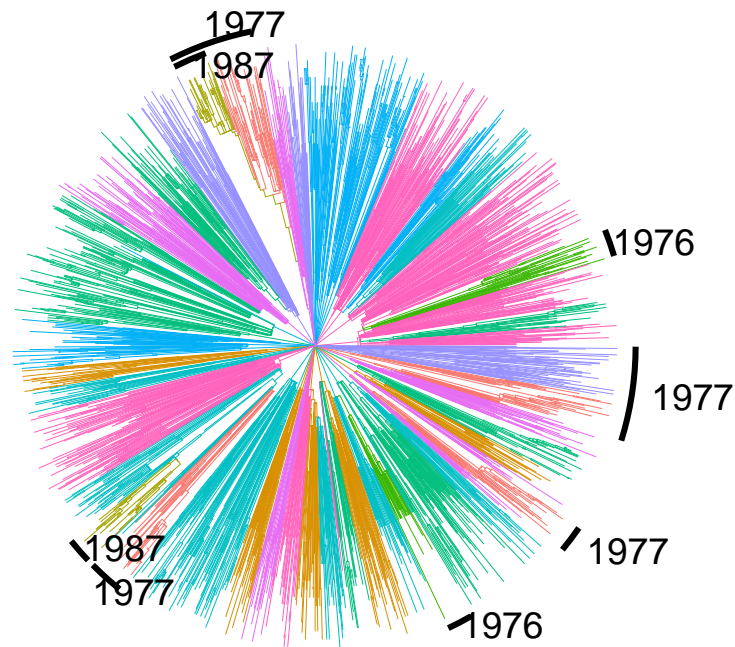540 of the tree with patient covariates. A particular partition stands out along

Figure 6: A time-scaled phylogeny based on 576 whole genomes of *N. gonorrhoeae*, comprising two clades with reduced susceptibility to cefixime (CFX) and their two sister clades. The top clade also has resistance to ciprofloxacin (CIP). Different colours on the tree represent the partition detected using the *treestructure* algorithm.

541 multiple dimensions: it is the smallest (size 58), polyphyletic, arose in the

542 recent past, and is characterised by very young MSM. The median year of

543 birth in this partition is 1987, in stark contrast to the rest of the sample with

544 year of birth in the 1970s. Clades within this young partition are also nested

545 paraphyletically under other relatively young partitions (cf. Figure 7).

546 We did not find a significant association between the tree partition

547 and residential postal codes (Tukey analysis of variance, $p = 0.097$). This is in

548 agreement with the original study which found minimal impact of geography

549 on genetic clustering in this sample, however this is largely a consequence of

550 the highly concentrated nature of the sample around Nashville. The ethnicity

551 of patients (black, white, and other) was strongly associated with the

552 estimated partition. Black MSM were strongly concentrated in the 1987

553 partition in particular (83% in contrast to 26-38% in all other partitions). The

554 odds ratio of black ethnicity given membership in the 1987 partition was 9.7

555 (95% CI:5.2-19.8).

562 Finally, we applied phylodynamic analysis methods to see if the

563 partition structure supported the previously published findings that young

564 MSM were transmitting at a higher rate (Dennis et al. 2018). To estimate $N_e$

565 through time, we used the nonparametric *skygrowth* R package (Volz and

566 Didelot 2018). We estimated $N_e(t)$ for each partition individually using a

567 range of precision parameters which control the smoothness ($\tau$) of the

568 estimated trajectories since we lack a priori information about volatility of

569 these trajectories. Figure 8 shows $N_e(t)$ for each partition with $\tau = 10$ and

570 supporting Figures S3 and S4 show results using different values of $\tau$. The

571 1987 partition again stands out as the only group which shows evidence of

572 recent and rapid population growth. Less dramatic recent periods of growth

573 are also noticeable for other partitions with young patients. The current

556

Figure 7: A time-scaled phylogeny estimated from HIV-1 *pol* sequences in Tennessee (Dennis et al. 2018). The colours correspond to the ten partitions identified using the *treestructure* algorithm. Several partitions are annotated with the median year of birth of HIV patients from whom sequences were sampled. Partitions lacking annotation had years of birth 1969-1972.

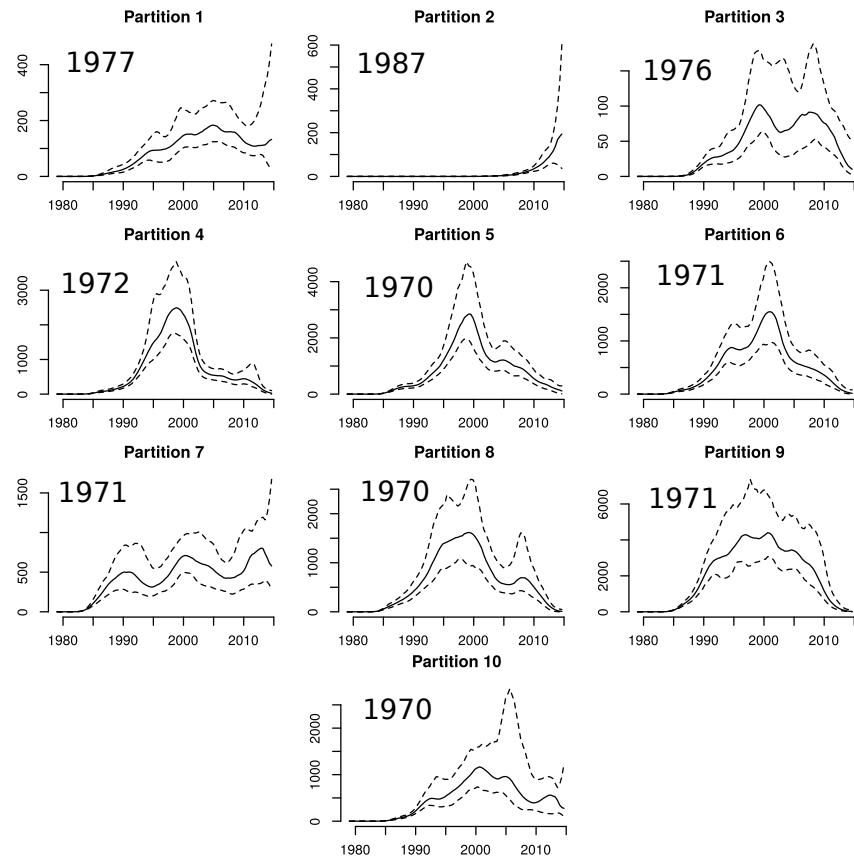574  exponential growth in the 1987 partition is not consistent across all analyses,

575  but when $\tau < 10$ we find $N_e(t)$ drops precipitously in 2014-2015 (Figure S3).

576  However, this could also be an artefact of non-random sampling and inclusion

577  of transmission pairs within the sample.

583  This analysis supports the hypothesis that there has been a recent and

584  rapid rise of HIV transmissions among young MSM in Tennessee and in

585  particular among young black MSM. This interpretation is mostly in

586  agreement with the original study (Dennis et al. 2018), but we find that black

587  MSM are a group at greater risk than young white MSM.

# Discussion

589  Contrasting the distribution of ordering of nodes provides a natural criterion

590  for distinguishing clades within a time-scaled phylogeny which are shaped by

591  different evolutionary or demographic processes. The non-parametric nature of

592  this classification method imposes minimal assumptions on the mechanisms

593  that generate phylogenetic patterns. Thus, we have found this method

594  maintains good performance over a diverse range of situations where

595  phylogenetic structure is produced, including differential transmission rates,

596  epidemiological outbreaks, evolution of beneficial mutations, and differential

597  sampling patterns. Our work is related to the research on species delimitation

598  methods (see for example (Zhang et al. 2013)) although targeted at

599  within-species variation, and is also related to recent work on methods for

600  detecting co-diversification of species(Oaks et al. 2019). This method appears

601  relatively robust compared to other methods against false-positive

602  identification of phylogenetic structure, but nevertheless has good sensitivity

603  for detecting structure in most situations.

604  There are many immediate applications of this method in the area of

578

Figure 8: Estimated effective population size through time for each partition in the Tennessee HIV-1 phylogeny. Each panel is annotated with the median year of birth among HIV patients in each partition. $N_e(t)$ was estimated using the *skygrowth* method (Volz and Didelot 2018) with precision parameter $\tau = 10$.

pathogen evolution where time-scaled phylogenetics is increasingly used in epidemiological investigations (Biek et al. 2015). We have demonstrated the role of natural selection in shaping phylogenetic structure of *N. gonorrhoeae*, and our method clearly identifies clades which expanded in the recent past due to acquisition of antimicrobial resistance. We have demonstrated the role of human demography and transmission patterns in shaping the evolution of HIV-1, and our method has shown distinct outbreaks of HIV-1 in specific groups defined by age, race, and behaviour. Furthermore, we have shown how clades detected by this method can be analysed using phylodynamic methods that can yield additional insights into recent outbreaks or the mechanisms which generated phylogenetic structure. For example, we have applied non-parametric methods to estimate the effective population size through time in HIV outbreaks detected using *treestructure* which highlighted particular groups that appear to be at higher risk of transmission. Such analyses would be more problematic using other partitioning or clustering algorithms because phylogenetic clusters can appear by chance in homogeneous populations of neutrally evolving pathogens, and this can give the false appearance of recent growth (Dearlove et al. 2017). This application of phylodynamics analysis methods is possible because the statistical test used in *treestructure* provides theoretical justification for treating each partition as a separate unstructured population.

Applications of the *treestructure* algorithms are scalable to relatively large phylogenies. The main algorithms require only a single pre-order traversal of the tree and all of the computations presented here required less than one minute to run. The method is based on a time-scaled phylogeny, and the computational burden of this preliminary step is typically higher than that of running *treestructure*, even though significant progress has been made

632  recently is this area (Volz and Frost 2017; Didelot et al. 2018; Sagulenko et al.

633  2018; Tamura et al. 2018; Miura et al. 2019). Future developments of

634  *treestructure* and other methods post-processing time-scaled phylogenies (Volz

635  and Didelot 2018; Didelot et al. 2017) should address the uncertainty in the

636  input phylogeny, for example by accounting for bootstrap or Bayesian support

637  values for phylogenetic splits, or by summarising results from multiple trees.

# 646  References

647  Beugin, M. P., T. Gayet, D. Pontier, S. Devillard, and T. Jombart. 2018. A

648  fast likelihood solution to the genetic clustering problem. Methods Ecol.

649  Evol. 9:1006–1016.

650  Biek, R., O. G. Pybus, J. O. Lloyd-Smith, and X. Didelot. 2015. Measurably

651  evolving pathogens in the genomic era. Trends Ecol. Evol. 30:306–313.

652  Bouckaert, R., J. Heled, D. Kühnert, T. Vaughan, C.-H. Wu, D. Xie, M. A.

653  Suchard, A. Rambaut, and A. J. Drummond. 2014. Beast 2: a software

654  platform for bayesian evolutionary analysis. PLoS Comput. Biol.

655  10:e1003537.

656  De Maio, N., C. J. Worby, D. J. Wilson, and N. Stoesser. 2018. Bayesian

reconstruction of transmission within outbreaks using genomic variants. PLoS Comput. Biol. 14:e1006117.

Dearlove, B. L. and S. D. W. Frost. 2015. Measuring Asymmetry in Time-Stamped Phylogenies. PLoS Comput. Biol. 11:e1004312.

Dearlove, B. L., F. Xiang, and S. D. Frost. 2017. Biased phylodynamic inferences from analysing clusters of viral sequences. Virus Evolution 3.

Dennis, A. M., J. T. Herbeck, A. L. Brown, P. Kellam, T. de Oliveira, D. Pillay, C. Fraser, and M. S. Cohen. 2014. Phylogenetic studies of transmission dynamics in generalized HIV epidemics: an essential tool where the burden is greatest? J. Acquir. Immune Defic. Syndr. 67:181–195.

Dennis, A. M., E. Volz, S. D. Frost, M. Hossain, A. F. Poon, P. F. Rebeiro, S. H. Vermund, T. R. Sterling, and M. L. Kalish. 2018. Hiv-1 transmission clustering and phylodynamics highlight the important role of young men who have sex with men. AIDS Research and Human Retroviruses 34:879–888.

Didelot, X., N. J. Croucher, S. D. Bentley, S. R. Harris, and D. J. Wilson. 2018. Bayesian inference of ancestral dates on bacterial phylogenetic trees. Nucleic Acids Res. 46:e134.

Didelot, X., C. Fraser, J. Gardy, and C. Colijn. 2017. Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks. Mol. Biol. Evol. 34:997–1007.

Didelot, X., A. S. Walker, T. E. Peto, D. W. Crook, and D. J. Wilson. 2016. Within-host evolution of bacterial pathogens. Nat. Rev. Microbiol. 14:150–162.

681    Didelot, X. and D. J. Wilson. 2015. ClonalFrameML: Efficient Inference of

682    Recombination in Whole Bacterial Genomes. PLoS Comput. Biol.

683    11:e1004041.

684    Dudas, G., L. M. Carvalho, T. Bedford, A. J. Tatem, G. Baele, N. R. Faria,

685    D. J. Park, J. T. Ladner, A. Arias, D. Asogun, F. Bielejec, S. L. Caddy,

686    M. Cotten, J. D'Ambrozio, S. Dellicour, A. Di Caro, J. W. Diclaro,

687    S. Duraffour, M. J. Elmore, L. S. Fakoli, O. Faye, M. L. Gilbert, S. M.

688    Gevao, S. Gire, A. Gladden-Young, A. Gnirke, A. Goba, D. S. Grant, B. L.

689    Haagmans, J. A. Hiscox, U. Jah, J. R. Kugelman, D. Liu, J. Lu, C. M.

690    Malboeuf, S. Mate, D. A. Matthews, C. B. Matranga, L. W. Meredith,

691    J. Qu, J. Quick, S. D. Pas, M. V. T. Phan, G. Pollakis, C. B. Reusken,

692    M. Sanchez-Lockhart, S. F. Schaffner, J. S. Schieffelin, R. S. Sealfon,

693    E. Simon-Loriere, S. L. Smits, K. Stoecker, L. Thorne, E. A. Tobin, M. A.

694    Vandi, S. J. Watson, K. West, S. Whitmer, M. R. Wiley, S. M. Winnicki,

695    S. Wohl, R. Wölfel, N. L. Yozwiak, K. G. Andersen, S. O. Blyden, F. Bolay,

696    M. W. Carroll, B. Dahn, B. Diallo, P. Formenty, C. Fraser, G. F. Gao, R. F.

697    Garry, I. Goodfellow, S. Günther, C. T. Happi, E. C. Holmes, B. Kargbo,

698    S. Keïta, P. Kellam, M. P. G. Koopmans, J. H. Kuhn, N. J. Loman,

699    N. Magassouba, D. Naidoo, S. T. Nichol, T. Nyenswah, G. Palacios, O. G.

700    Pybus, P. C. Sabeti, A. Sall, U. Ströher, I. Wurie, M. A. Suchard, P. Lemey,

701    and A. Rambaut. 2017. Virus genomes reveal factors that spread and

702    sustained the ebola epidemic. Nature 544:309–315.

703    Everitt, B., S. Landau, and M. Leese. 2001. Cluster Analysis. Wiley New York.

704

705    Eyre, D. W., T. Golubchik, N. C. Gordon, R. Bowden, P. Piazza, E. M. Batty,

706    C. L. C. Ip, D. J. Wilson, X. Didelot, L. O'Connor, R. Lay, D. Buck, A. M.

707    Kearns, A. Shaw, J. Paul, M. H. Wilcox, P. J. Donnelly, T. E. A. Peto, A. S.

708    Walker, and D. W. Crook. 2012. A pilot study of rapid benchtop sequencing

709    of Staphylococcus aureus and Clostridium difficile for outbreak detection

710    and surveillance. BMJ Open 2:e001124.

711  Grad, Y. H., S. R. Harris, R. D. Kirkcaldy, A. G. Green, D. S. Marks, S. D.

712    Bentley, D. Trees, and M. Lipsitch. 2016. Genomic epidemiology of

713    gonococcal resistance to extended-spectrum cephalosporins, macrolides, and

714    fluoroquinolones in the united states, 2000–2013. The Journal of Infectious

715    Diseases 214:1579–1587.

716  Guindon, S., J.-F. Dufayard, V. Lefort, M. Anisimova, W. Hordijk, and

717    O. Gascuel. 2010. New algorithms and methods to estimate

718    maximum-likelihood phylogenies: assessing the performance of PhyML 3.0.

719    Systematic Biology 59:307–21.

720  Han, A., E. Parker, S. Maurer-Stroh, and C. Russell. 2018. Inferring putative

721    transmission clusters with phydelity. bioRxiv Page 477653.

722  Hartl, D. L., A. G. Clark, and A. G. Clark. 1997. Principles of population

723    genetics vol. 116. Sinauer associates Sunderland, MA.

724  Höhna, S., M. J. Landis, T. A. Heath, B. Boussau, N. Lartillot, B. R. Moore,

725    J. P. Huelsenbeck, and F. Ronquist. 2016. Revbayes: Bayesian phylogenetic

726    inference using graphical models and an interactive model-specification

727    language. Systematic Biology 65:726–736.

728  Klingen, T. R., S. Reimering, C. A. Guzmán, and A. C. McHardy. 2018. In

729    silico vaccine strain prediction for human influenza viruses. Trends in

730    microbiology 26:119–131.

731  Lam, T. T.-Y., B. Zhou, J. Wang, Y. Chai, Y. Shen, X. Chen, C. Ma,

732    W. Hong, Y. Chen, Y. Zhang, L. Duan, P. Chen, J. Jiang, Y. Zhang, L. Li,

L. L. M. Poon, R. J. Webby, D. K. Smith, G. M. Leung, J. S. M. Peiris, E. C. Holmes, Y. Guan, and H. Zhu. 2015. Dissemination, divergence and establishment of H7N9 influenza viruses in china. Nature 522:102–105.

Ledda, A., J. R. Price, K. Cole, M. J. Llewelyn, A. M. Kearns, D. W. Crook, J. Paul, and X. Didelot. 2017. Re-emergence of methicillin susceptibility in a resistant lineage of Staphylococcus aureus. J. Antimicrob. Chemother. 72:1285–1288.

McCloskey, R. M. and A. F. Poon. 2017. A model-based clustering method to detect infectious disease transmission outbreaks from sequence variation. PLoS Comput. Biol. 13:e1005868.

Miller, R., J. Price, E. Batty, X. Didelot, D. Wyllie, T. Golubchik, D. W. Crook, J. Paul, T. E. A. Peto, D. J. Wilson, M. Cule, C. Ip, N. Day, C. Moore, R. Bowden, and M. Llewelyn. 2014. Healthcare-associated outbreak of meticillin-resistant Staphylococcus aureus bacteraemia: role of a cryptic variant of an epidemic clone. J. Hosp. Infect. 86:83–89.

Miura, S., K. Tamura, S. L. K. Pond, L. A. Huuki, J. Priest, J. Deng, and S. Kumar. 2019. A new method for inferring timetrees from temporally sampled molecular sequences. BioRxiv Page 620187.

Mostowy, R., N. J. Croucher, C. P. Andam, J. Corander, W. P. Hanage, and P. Marttinen. 2017. Efficient inference of recent and ancestral recombination within bacterial populations. Mol. Biol. Evol. 34:1167–1182.

Notohara, M. 1990. The coalescent and the genealogical process in geographically structured population. J. Math. Biol. 29:59–75.

Oaks, J. R., N. LBahy, and K. A. Cobb. 2019. Insights from a general, full-likelihood bayesian approach to inferring shared evolutionary events

758    from genomic data: Inferring shared demographic events is challenging.

759    bioRxiv Page 679878.

760    Rambaut, A. and N. C. Grass. 1997. Seq-Gen: an application for the Monte

761    Carlo simulation of DNA sequence evolution along phylogenetic trees.

762    Bioinformatics 13:235–238.

763    Sagulenko, P., V. Puller, and R. A. Neher. 2018. Treetime:

764    Maximum-likelihood phylodynamic analysis. Virus Evolution 4:vex042.

765    Suchard, M. A., P. Lemey, G. Baele, D. L. Ayres, A. J. Drummond, and

766    A. Rambaut. 2018. Bayesian phylogenetic and phylodynamic data

767    integration using beast 1.10. Virus Evolution 4:vey016.

768    Tamura, K., Q. Tao, and S. Kumar. 2018. Theoretical foundation of the

769    RelTime method for estimating divergence times from variable evolutionary

770    rates. Mol. Biol. Evol. 35:1770–1782.

771    To, T.-H., M. Jung, S. Lycett, and O. Gascuel. 2016. Fast dating using

772    Least-Squares criteria and algorithms. Systematic Biology 65:82–97.

773    Tonkin-Hill, G., J. A. Lees, S. D. Bentley, S. D. W. Frost, and J. Corander.

774    2018. RhierBAPS: An R implementation of the population clustering

775    algorithm hierBAPS. Wellcome Open Res 3:93.

776    Tonkin-Hill, G., J. A. Lees, S. D. Bentley, S. D. W. Frost, and J. Corander.

777    2019. Fast hierarchical Bayesian analysis of population structure. Nucleic

778    Acids Res. Pages 1–11.

779    Vaughan, T. G. and A. J. Drummond. 2013. A stochastic simulator of

780    birth–death master equations with application to phylodynamics. Molecular

781    biology and evolution 30:1480–1493.

Vinh, N. X., J. Epps, and J. Bailey. 2010. Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance. Journal of Machine Learning Research 11:2837–2854.

Volz, E. M. and X. Didelot. 2018. Modeling the growth and decline of pathogen effective population size provides insight into epidemic dynamics and drivers of antimicrobial resistance. Systematic Biology 67:719–728.

Volz, E. M. and S. D. W. Frost. 2017. Scalable relaxed clock phylogenetic dating. Virus Evolution 3.

Wakeley, J. 2009. Coalescent theory: an introduction. Greenwood Village: Roberts & Company Publishers.

Whittles, L. K., P. J. White, and X. Didelot. 2017. Estimating the fitness benefit and cost of cefixime resistance in Neisseria gonorrhoeae to inform prescription policy: A modelling study. PLoS Med. 14:e1002416.

Wiuf, C. and P. Donnelly. 1999. Conditional genealogies and the age of a neutral mutant. Theor. Popul. Biol. 56:183–201.

Yu, G., D. K. Smith, H. Zhu, Y. Guan, and T. T.-Y. Lam. 2017. ggtree: an r package for visualization and annotation of phylogenetic trees with their covariates and other associated data. Methods in Ecology and Evolution 8:28–36.

Zhang, J., P. Kapli, P. Pavlidis, and A. Stamatakis. 2013. A general species delimitation method with applications to phylogenetic placements. Bioinformatics 29:2869–2876.

**Data:**  1) Disjoint sets of tips $X$ and $Y$
2) Empirical value of test statistic $\hat{R}$
3) Number of simulations $n_{\text{sim}}$
4) Taxonomic condition $E$ (see Equations 3, 4 or 10)
**Result:** Two-sided p-value denoted $q = \xi(X, Y, \hat{R})$.
Initialisation;
Form a time-ordered sequence of nodes

$$U = (u_1, \cdots, u_{|D_X|+|D_Y|}) | u_i \in (D_X \cup D_Y), \tau(u_i) \geq \tau(u_{i+1})$$

Form a corresponding numeric sequence:
$\Upsilon = (v_1, \cdots, v_{|D_X|+|D_Y|})$ where

$$v_i = \begin{cases} 1 & \text{if } u_i \in X \\ -1 & \text{if } u_i \in Y \\ 0 & \text{if } u_i \in (D_X \cup D_Y) \cap \mathcal{I} \end{cases}$$

**for** $k = 1$ to $n_{sim}$ **do**
    $z \leftarrow 0$ (simulated lineages through time in clade $X$)
    $w \leftarrow 0$ (simulated lineages through time in clade $Y$)
    $r_{\text{sim}} \leftarrow 0$ (simulated rank-sum statistic)
    $c \leftarrow 0$ (number of coalescent events simulated)
    **for** $i = 1$ to $|D_X| + |D_Y|$ **do**
        **if** $v_i = 1$ **then**
            Account for sample in $X$: $z \leftarrow z + 1$ ;
        **if** $v_i = -1$ **then**
            Account for sample in $Y$: $w \leftarrow w + 1$ ;
        **if** $W_i = 0$ **then**
            Increment coalescent counter: $c \leftarrow c + 1$ ;
            Compute probability $\tilde{p} = \tilde{Q}_E(z, w)$ that next coalescent is in
              $D_X$ or $D_Y$ using Equation 3, 4 or 10;
            Draw a random uniform variable $\omega \leftarrow \text{Unif}(0, 1)$ ;
            **if** $\omega < \tilde{p}$ **then**
                $z \leftarrow z - 1$
                $r_{\text{sim}} \leftarrow r_{\text{sim}} + c$
            **else**
                $w \leftarrow w - 1$
    **end**
    Record simulated statistic:
    $R_k \leftarrow r_{\text{sim}}$
**end**
Compute number of simulations more and less than empirical value:

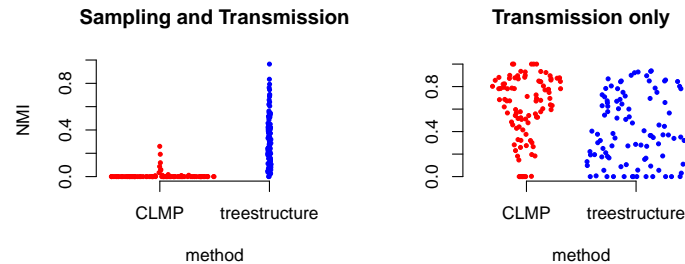$$m_+ \leftarrow |\{r' \in R_k | r' > \hat{R}\}|$$

$$m_- \leftarrow |\{r' \in R_k | r' < \hat{R}\}|$$

Return $\min(\frac{m_+}{n_{\text{sim}}}, \frac{m_-}{n_{\text{sim}}})$.
**Algorithm 1:** Algorithm for computing the null distribution and associated p-value of the test-statistic for cladistic outliers.
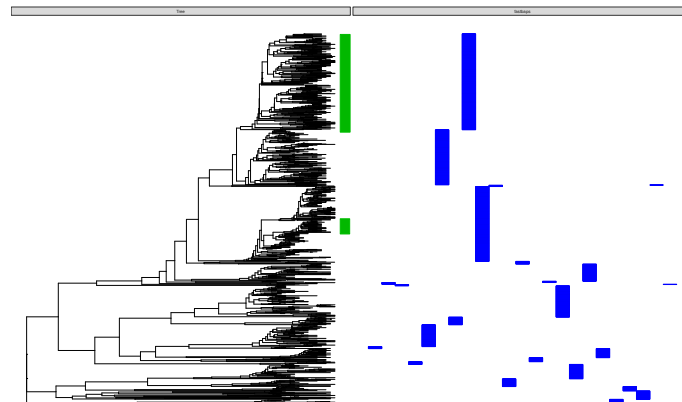
804

**Data:** Time-scale genealogy $\mathcal{G}$

**Result:** Partition of tips of tree, denoted $M$.

Initialise 'active set' to consist of root node: $\Omega \leftarrow \{\text{root}\}$ ;

Initialise partition: $M \leftarrow \emptyset$ ;

**for** $u \in \mathcal{I}$ *(internal nodes)* **do**
  | Initialise $\tilde{C}_u \leftarrow C_u$ ;
**end**

**while** $|\Omega| > 0$ **do**
  | Initialise $\Omega' \leftarrow \Omega$ ;
  | **for** $u \in \Omega$ **do**
  | | Find biggest outlier descended from $u$:
  | | $v^* \leftarrow \text{argmax}_{v \in C_u} f(v) = \xi(\tilde{C}_u \setminus \tilde{C}_v, \tilde{C}_v)$ (Algorithm 1);
  | | $q \leftarrow \xi(\tilde{C}_u, \tilde{C}_{v^*})$ ;
  | | **if** $q < \alpha$ **then**
  | | | $\Omega' \leftarrow \Omega' \cup v^*$ ;
  | | | $\tilde{C}_u \leftarrow \tilde{C}_u \setminus C_{v^*}$ ;
  | | **else**
  | | | No significant outliers, so remove $u$ from active sets:
  | | | $\Omega' \leftarrow \Omega' \setminus u$ ;
  | | | Add the clade descended from $u$ to the partition:
  | | | $M \leftarrow M \cup \{(\mathcal{T} \cap \tilde{C}_u)\}$ ;
  | **end**
  | $\Omega \leftarrow \Omega'$.
**end**

Return $M$.

**Algorithm 2:** Algorithm for detecting cladistic outliers.
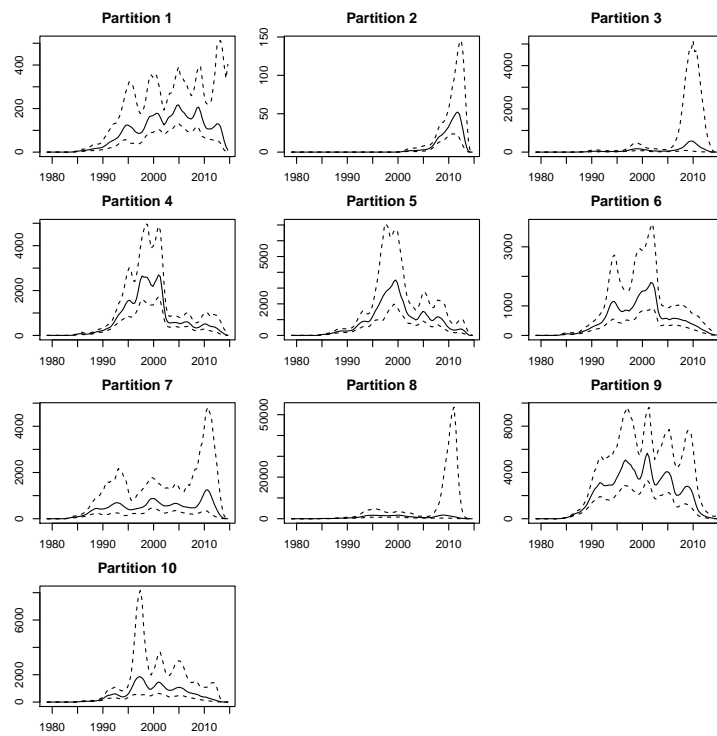
Figure S1: The normalised mutual information (NMI) for 100 previously published simulations (McCloskey and Poon 2017). This describes accuracy of classification of tips into outbreaks using the *treestructure* method and CLMP (McCloskey and Poon 2017). Results on left were based on simulations where both transmission and sampling rates varied in the outbreak cluster, whereas simulations on the right only allowed transmission rates to vary.
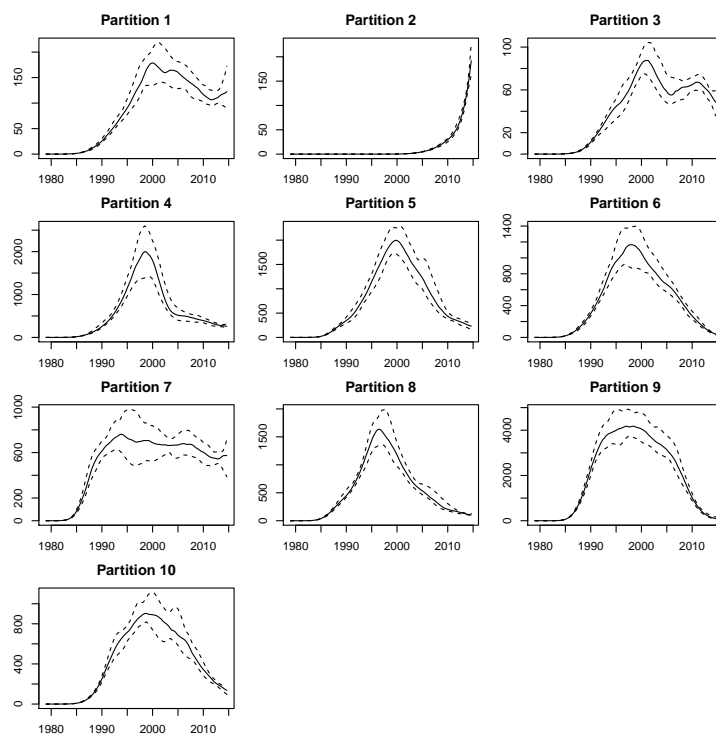


Figure S2: The output of FastBAPS classification applied to 1102 *N. gonorrhoeae* isolates described in the main text. Clades indicated in green have CFX resistance.

817

818 Figure S3: Estimated effective population size through time for each partition

819 in the Tennessee HIV-1 phylogeny. $N_e(t)$ was estimated using the *skygrowth*

820 method (Volz and Didelot 2018) with precision parameter $\tau = 1$.

Figure S4: Estimated effective population size through time for each partition in the Tennessee HIV-1 phylogeny. $N_e(t)$ was estimated using the *skygrowth* method (Volz and Didelot 2018) with precision parameter $\tau = 100$.