

Accurate detection of HIV transmission clusters from phylogenetic trees using a multi-state birth-death model

Joëlle Barido-Sottani^{1,*} and Tanja Stadler¹

¹*Department of Biosystems Science and Engineering, ETH Zürich, Basel, Switzerland*

**Correspondence to be addressed to: joelle.barido-sottani@m4x.org*

Abstract

HIV transmission networks are highly clustered, and accurate identification of these clusters is essential for effective targeting of public health interventions. This clustering affects the transmission dynamics of the HIV epidemic, which affects the pathogen phylogenies reconstructed from patient samples. We present a new method for identifying transmission clusters by detecting the changes in transmission rate provoked by the introduction of the epidemic into a new cluster. The method employs a multi-state birth-death (MSBD) model where each state represents a cluster. Transmission rates in each cluster decrease exponentially over time, simulating susceptible depletion in the cluster. This model is fitted to the pathogen phylogeny using a Maximum Likelihood approach. Using simulated datasets we show that the MSBD method is able to reliably infer both the cluster repartition and the transmission parameters from a pathogen phylogeny. In contrast to existing cutpoint-based methods for cluster identification, which are dependent on a parameter set by the user, the MSBD method is consistently reliable. It also performs better on phylogenies containing nested clusters. We present an application of our method to the inference of transmission clusters using sequences obtained from the Swiss HIV Cohort Study. The MSBD method is available as an R package.

23 1 Background

24 Basic epidemiologic models rest on the random mixing assumption (1; 2). In the presence
25 of random mixing, each individual in a population has an equal probability of coming into
26 contact with any other individual, which can lead to very quick epidemic spread. The random
27 mixing assumption may be appropriate for airborne diseases in small communities. For sexually-
28 transmitted infections (STIs) such as HIV-1 however, the random mixing hypothesis does not
29 hold: STIs spread within sexual contact networks that limit their propagation to a specific subset
30 of possible transmission events.

31 Identifying the structure in the sexual contact network has multiple applications, for instance
32 allowing public health officials to target the populations most vulnerable to infection. One
33 particular aim is to identify communities in the sexual contact network. These communities,
34 or clusters, are defined as sets of nodes in the sexual contact network such that most or all
35 nodes are connected within a cluster but few links exist between clusters (3). These clusters
36 will affect the dynamics of an epidemic: at first the infection will spread quickly in the cluster
37 where it has been introduced. The rate of transmission will then go down as the population of
38 susceptibles in the cluster is progressively exhausted (2). Eventually a new introduction event
39 may occur, where an individual from a previously uninfected contact cluster will be infected
40 through one of the inter-cluster connections. Since the newly infected cluster is completely
41 susceptible, the rate of transmission will then go up suddenly as new transmission routes open.
42 Thus the cluster structure of the sexual contact network shapes transmission dynamics and thus
43 may leave a detectable footprint in the phylogeny reconstructed from an epidemic. In what
44 follows, we always consider phylogenies on the epidemic level, i.e. phylogenies obtained from
45 pathogen genetic sequences of different infected individuals within an epidemic; thus each tip in
46 the phylogeny represents a unique infected host.

47 Previous studies have found varying degrees of influence of the contact network on the phy-
48 logeny. (4) found almost no influence of the clustering coefficient of a network and the shape of
49 transmission trees when the degree distribution of the network was kept constant, (5) found a
50 modest effect of the degree distribution in the network on the shape of phylogenies reconstructed
51 from simulated genetic data, whereas (6) found that the variance in degree distribution and the

52 mean path length of the contact network could significantly affect the shapes of phylogenies.
53 The link between network structures and phylogenies is also affected by viral characteristics such
54 as within-host viral diversity (7). Several methods have been proposed to identify structural
55 characteristics, such as connectivity and clustering coefficient, of the population network from a
56 viral phylogeny (8; 9).

57 A number of methods have been proposed which exploit the effects that contact networks have
58 on phylogenies to identify HIV transmission clusters from those phylogenies. In this paper we
59 will focus on the methods evaluated in (10), which we will refer to as “cutpoint-based” methods.
60 These methods differ in how they define the distance between two tips of the tree, but they have
61 two major features in common: first, they require an ad hoc cutpoint to be specified by the user
62 ; second, they assume that the clusters are monophyletic in the phylogeny or monophyletic in
63 a tree obtained from hierarchical clustering (Def. 4 in (10)), i.e that the most recent common
64 ancestor of all tips belonging to a given cluster has no other descending tips. As (10) found,
65 both features have a strong impact on the quality of the recovered clusters. Thus there is a need
66 for a method which does not have these limitations.

67 Multi-state birth-death models have been widely used to model population structure and
68 analyze phylogenies built from individuals in a structured population, both in epidemiological
69 and macroevolutionary applications. Thus in principle such a model may be used to study the
70 sexual contact network. In this context the aim is to infer which tips in a phylogeny belong to
71 which cluster of an unknown contact network. Clusters differ by having different transmission
72 dynamics through time, meaning different birth rates, so each cluster corresponds to a state in
73 a multi-state birth-death model.

74 The Binary State Speciation and Extinction (BiSSE, (11)) and its extension to multiple states
75 MuSSE, included in the package Diversitree (12), were the first efforts to infer state-specific birth
76 and death rates from ultrametric phylogenies where each tip is assigned to a state. In (13), these
77 approaches were extended to non-ultrametric trees. More recently the Beast2 package BDMM
78 (14) allowed the joint reconstruction of a phylogeny and quantification of the parameters of an
79 underlying multi-state birth-death model. These approaches require the user to specify how many
80 states the model contains and to which state each tip of the phylogeny belongs. An exception to
81 the latter is (13), which can integrate over tip states, but does not assign states to tips.

82 However we cannot readily use any of the above approaches to infer transmission clusters, for
83 two reasons. First, the state of tips, i.e which cluster they belong to, is not known prior to the
84 analysis. Second, integrating over the tip states instead explicitly assigning states to tips means
85 that the repartition of tips into clusters cannot be inferred.

86 The method Bayesian Analysis of Macroevolutionary Mixtures (BAMM, (15)) addresses these
87 issues and is able to infer the number of clusters and assign each tip to a cluster. Further, the
88 birth- and death rate parameters associated with each cluster are quantified. However, it was
89 designed to be used with macroevolutionary datasets, meaning at the time of this writing it
90 can only analyze ultrametric trees, i.e with all tips sampled at the same point in time. For
91 epidemiological datasets, we have non-ultrametric trees as samples are collected through time.
92 Furthermore, its results have been called into question (16).

93 In this paper, we present a new method to identify clusters of transmission in a phylogeny
94 built from viral sequences, by detecting 'jumps' in transmission rate. We associate these jumps
95 with introduction events into previously untouched clusters. From the detected jumps, we can
96 readily read off the partition of the tips of our phylogeny into distinct clusters. Our method
97 uses the multi-state birth-death (MSBD) model with allowing decreasing transmission rates
98 within clusters to account for the depletion of susceptibles. In particular, it does not require
99 prior knowledge on the number of clusters or the tip repartition in clusters. We evaluate the
100 performance of this new method on the simulated dataset of (10) and compare it to cutpoint-
101 based methods. We then apply it to a published HIV phylogeny (9) which was obtained based
102 on 192 sequences from the Swiss HIV Cohort Study. Finally we discuss the limitations of the
103 method and planned future work.

104 **2 Methods**

105 **2.1 Model**

106 We use a multi-state birth-death model similar to the model used in the BDMM package (14).
107 The birth-death process starts with one infected individual at time τ in the past in an ancestral
108 state and is stopped at present time 0. This means that we measure time in the backwards

109 direction, increasing from the present to the root. State changes happen in each individual
110 through time with a rate γ . Our MSBD model contains an unknown number of states n^* ,
111 corresponding to n^* clusters in the underlying population network. We assume that all states
112 are equally likely to transition to, so that the state change rate between any state i and j is as
113 follows:

$$m_{i,j} = \frac{\gamma}{n^* - 1} \quad \forall i, j \neq i$$

114 Each state i is characterized by a specific initial transmission rate $\lambda_{0,i}$, a transmission decay
115 rate z_i , and a removal rate μ_i . Each individual produces an additional individual with a state-
116 and time-dependent transmission rate $\lambda_i(t)$ (function of $\lambda_{0,i}, z_i$ as defined below), and is removed
117 with a state-dependent removal rate μ_i corresponding to the rate of “becoming non-infectious”.

118 The depletion of the susceptible population is modeled by the exponential decay of the trans-
119 mission rates in the process. Each state is associated with a specific initial transmission rate $\lambda_{0,i}$
120 and a transmission decay rate z_i . Equation 1 shows the transmission rate for a lineage in state
121 i at time t before the present, where $t_{0,i}$ is the time of introduction into state i . Since time is
122 backwards, we impose $z_i \geq 0$, so that the transmission rate decreases as the process progresses
123 towards the present.

$$\lambda_i(t) = \lambda_{0,i} \times e^{z_i(t-t_{0,i})} \quad (1)$$

124 The infected individuals are sampled upon removal with a probability σ . This birth-death
125 model produces a tree on all infected individuals together with position and times of rate changes
126 on the tree, and we obtain the phylogeny by considering the subtree spanned by the sampled
127 infected individuals. The phylogeny contains information about the transmission and removal
128 times of the sampled individuals, as well as the positions and times of the rate changes, as
129 shown in Figure 1. We assume that the state changes correspond to introduction events in newly
130 infected clusters, so that all tips inferred to be in the same state belong to the same cluster in
131 the original transmission network.

132 We refer to a node in the phylogeny being either a branching event, a tip, or a state change
133 event. Edges in the phylogeny connect any two nodes, and so any edge belongs to only one state.

134 2.2 Likelihood function

135 We now derive the probability density of a phylogeny (including the state change times) given
 136 the MSBD parameters, i.e. we derive the likelihood of the parameters given a phylogeny.

137 2.2.1 Differential equations

138 Following (13; 14), the likelihood function of the model parameters given the phylogeny can be
 139 calculated from the differential equations below. Eqn. (2) describes the probability $p_i(t)$ of a
 140 lineage in state i at time t not producing any sampled offspring until the present (referred to
 141 extinction probability below). Eqn. (3) describes the probability density $q_{i,N}(t)$ of an edge N in
 142 state i at time t evolving according to the phylogeny in time interval $[t, 0]$.

$$\begin{aligned} \frac{dp_i}{dt}(t) &= -(\gamma + \lambda_i(t) + \mu_i)p_i(t) + \mu_i(1 - \sigma) + \lambda_i(t)p_i(t)^2 + \gamma \sum_{j \neq i} p_j(t), \\ p_i(0) &= 1, \end{aligned} \quad (2)$$

143

$$\begin{aligned} \frac{dq_{i,N}}{dt}(t) &= -(\gamma + \lambda_i(t) + \mu_i)q_{i,N}(t) + 2\lambda_i(t)q_{i,N}(t)p_i(t), \\ q_{i,N}(t_s) &= \mu_i\sigma && \text{if } N \text{ leads to a tip at time } t_s, \\ q_{i,N}(t_t) &= \lambda_i(t_t)q_{i,N'}(t_t)q_{i,N''}(t_t) && \text{if } N \text{ undergoes transmission at } t_t, \text{ leading to } N' \text{ and } N'', \\ q_{i,N}(t_c) &= \frac{\gamma}{n^* - 1}q_{j,N}(t_c) && \text{if } N \text{ changes to state } j \text{ at } t_c. \end{aligned} \quad (3)$$

144 The probability of a phylogeny starting at root time τ with initial state I is $q_{I,N}(\tau)$ so the full
 145 likelihood can be calculated from Eq 3. Rather than writing it recursively as in Eq 3, it can be
 146 written as a closed form equation by defining the edge likelihood function $f_N = \frac{q_{i,N}(t_b)}{q_{i,N}(t_e)}$ for an
 147 edge N in state i with start time t_b and end time t_e . f_N follows the differential equation in Eq
 148 3 with initial condition $f_N(t_e) = 1$. The full likelihood of the model M given the phylogeny T
 149 is then obtained by multiplying the likelihoods of all edges as shown in Equation (4), where n is
 150 the number of states (including the root state) in the tree, N_i is the set of edges in state i , T_i

151 the set of transmission events in state i and S_i the set of tips in state i .

$$L(M|T) = q_{I,N}(\tau) = \prod_i \left[\prod_{N \in N_i} f_N \times \prod_{t \in T_i(T)} \lambda_i(t_t) \times \prod_{s \in S_i(T)} \sigma \mu_i \right] \times \left(\frac{\gamma}{n^* - 1} \right)^{n-1} \quad (4)$$

152 This likelihood function can be applied to trees with or without a root edge, i.e trees starting
153 with one lineage or two at time τ .

154 2.3 Approximations to the likelihood function

155 2.3.1 Simplifying the number of states

156 Since the real number of clusters in the underlying network n^* is unknown, we need to estimate
157 it. However this parameter only appears in the likelihood in the factor $\left(\frac{\gamma}{n^*-1}\right)^{n-1}$ so maximizing
158 the likelihood is equivalent to minimizing n^* . We further assume that each migration enters a
159 previously not visited state, i.e. $n^* \geq n$. Together, the maximum likelihood estimate will always
160 be $n^* = n$. Thus we fix $n^* = n$ in the inference.

161 2.3.2 Ignoring state changes in unsampled subtrees

162 The equations for p and f_N do not have an analytical solution. Numerical integration is compu-
163 tationally expensive and can be unstable for certain parameters, so we make the assumption that
164 no state changes happen in the unsampled parts of the tree, meaning we observe all state changes
165 in the final tree. With this assumption, the master equation for $p_i(t)$ changes to Equation (5),

$$\begin{aligned} \frac{dp_i}{dt}(t) &= -(\gamma + \lambda_i(t) + \mu_i)p_i(t) + \mu_i(1 - \sigma) + \lambda_i(t)p_i(t)^2, \\ p_i(0) &= 1. \end{aligned} \quad (5)$$

166 These equations have an analytical solution for constant transmission and removal rates, but not
167 necessarily for time-dependent rates. To obtain a closed form solution, we use time discretization
168 and assume that in each time step the transmission rate can be considered constant, as described
169 in the next section.

170 2.3.3 Time discretization

171 We discretize the time-dependent transmission rates by assuming that they can be considered
 172 locally constant on small enough intervals. The grid size used for the discretization is fixed across
 173 the tree and needs to be specified by the user. A smaller size will improve the accuracy of the
 174 likelihood calculation but also increase the computational cost.

175 Time discretization for p

176 A closed form of the extinction probability and the likelihood function can be obtained for
 177 piecewise constant transmission and removal rates. Assuming constant rates in Eqn 5, and a
 178 generic initial condition $p_i(t_{IC}) = V_{IC}$ (rather than the initial condition $p_i(0) = 1$), we obtain
 179 an analytic solution of Eqn 5,

$$\begin{aligned}
 p_i(t) &= -\frac{1}{\lambda_i} \frac{(y_i + \lambda_i V_{IC})x_i e^{-ct} - y_i(x_i + \lambda_i V_{IC})e^{-ct_{IC}}}{(y_i + \lambda_i V_{IC})e^{-ct} - (x_i + \lambda_i V_{IC})e^{-ct_{IC}}} \\
 c &= \sqrt{(\gamma + \lambda_i + \mu_i)^2 - 4\mu_i(1 - \sigma)\lambda_i} \\
 x_i &= \frac{-(\gamma + \lambda_i + \mu_i) - c}{2} \quad \text{and} \quad y_i = \frac{-(\gamma + \lambda_i + \mu_i) + c}{2}
 \end{aligned} \tag{6}$$

180 This solution can be verified by differentiating the solution and substituting the result into Eqn
 181 5.

182 To obtain $p_i(t)$ using this time discretization, we divide the time interval $[\tau; 0]$ into a grid.
 183 Starting with $p_i(0) = 1$, we can then evaluate p_i using Eq 6 in each grid interval going backwards
 184 in time, using as initial value the solution of the previous grid interval.

185 Time discretization for f_N

186 A closed form solution of the edge likelihood function f_N can now be calculated, for a small
 187 time interval $[t_l; t_{l-1}]$ on an edge N in state i . This expression uses the value of $p_i(t_{l-1})$, which
 188 can be calculated as explained in “Time discretization for p ”. We define $f_N(t_l, t_{l-1}) = \frac{q_{i,N}(t_l)}{q_{i,N}(t_{l-1})}$,
 189 and obtain

$$f_N(t_l, t_{l-1}) = e^{c(t_{l-1}-t_l)} \left(\frac{y_i - x_i}{(y_i + \lambda_i p_i(t_{l-1}))e^{-c(t_l-t_{l-1})} - (x_i + \lambda_i p_i(t_{l-1}))} \right)^2 \tag{7}$$

190 This expression for $f_N(t_l, t_{l-1})$ is a solution of the differential equation 3 with $f_N(t_{l-1}) = 1$,

191 assuming the rates are constant in interval $[t_l, t_{l-1}]$ and using the approximate function $p_i(t)$
192 from Eq. 6. This can be easily verified by differentiating Eq. 6 and Eq. 7 and substituting the
193 resulting expressions $\frac{d}{dt}p_i(t)$ and $\frac{d}{dt}f_N(t_l, t_{l-1})$ into the differential equations 5 and 3. Equations
194 6 and 7 are identical to the expressions used in the birth-death skyline model and a full derivation
195 of them can be found in (17).

196 We now describe how to calculate f_N and obtain an evaluation of the likelihood provided in
197 Eq. 4, using Eqn 7. Values of p_i for all branching times and state change times are precomputed
198 to avoid the repetition of those calculations for multiple edges. For edge N in state i starting at
199 time t_b and ending at time t_e (i.e. $t_b < t_e$), we aim to calculate $f_N(t_b, t_e)$. Thus we aim to solve,
200 using the time discretization, the differential equation in Eqn 3 with initial value $f(t_e, t_e) = 1$:

- 201 1. Fetch the precomputed value of $p_i(t_e)$.
- 202 2. Divide the interval $[t_b, t_e]$ in k equidistant intervals $[t_k, t_{k-1}], [t_{k-1}, t_{k-2}], \dots, [t_1, t_0]$ with
203 $t_0 = t_e$ and $t_k = t_b$.
- 204 3. For each step $l \in [1..k]$ do the following:
 - 205 (a) calculate $\lambda_{i,l}$ the mean of $\lambda_i(t)$ on the interval $[t_l, t_{l-1}]$, then
 - 206 (b) calculate $p_i(t_l)$ and $f_N(t_l, t_{l-1})$ by using the constant rates solutions provided in Eqn
207 6 for p and in Eqn 7 for f with $\lambda_i = \lambda_{i,l}$, based on the value $p_i(t_{l-1})$ given by the
208 precomputed value if $l = 1$ and by the previous step $l - 1$ otherwise.
- 209 4. Finally, compute $f_N(t_b, t_e) = \prod_{l=1}^k f_N(t_l, t_{l-1})$.

210 2.4 Algorithm

211 We now present an algorithm which identifies the state change configuration and associated
212 parameters that maximize the likelihood in Eq. 4 for a particular phylogeny T .

213 2.4.1 Initial condition

214 The first step of the algorithm is to infer the most likely parameters for a constant rate birth-
215 death model given the tree. These parameters will be used as starting values for the optimization

216 in further steps. The initial values used in the optimization can have a great impact on the entire
217 inference: if they are too distant from the optimal values, it can happen that the constant rates
218 optimization finds only a local optima, and this will in turn affect all subsequent steps of the
219 inference. Our method avoids this issue by applying an initial coarse-grained optimization step
220 prior to the main optimization algorithm. Initial values are tested until no further improvement
221 of the optima found by the optimization can be obtained. This optima will then be accepted
222 as the global optima for the constant rates model. The user-provided starting values define the
223 order of magnitude of the values tested in this phase.

224 **2.4.2 Maximum likelihood search**

225 We then use a greedy approach to add state changes until no further improvement of the likelihood
226 can be obtained. New maximum likelihood estimates are obtained for all transmission, decay,
227 removal and state change rates each time a new state change is tested, but the positions and
228 times of previous state changes are fixed.

229 Once a configuration has been found in which no more state changes can be added to improve
230 the likelihood, we will attempt to recursively remove all the states from this configuration. This
231 step is designed to compensate partly for the fact that the greedy approach never goes back on
232 previous state change assignments, and so can end up in sub-optimal configurations.

233 Once no further improvements of the likelihood can be obtained by either adding or removing
234 a state, the method will return the best fitting model found, including the state configuration
235 and the maximum likelihood estimates for all parameters.

236 The full algorithm, including the initial coarse-grained search phase, is as follows:

- 237 1. Find the most likely parameters for a one-state birth-death model (i.e with identical birth
238 and death rates across the tree).
- 239 2. For all edges in the tree:
 - 240 (a) add a state change on this edge, then
 - 241 (b) find the most likely parameters for this state configuration, then
 - 242 (c) keep the edge as candidate if it is the most likely found so far.

- 243 3. If a configuration with $n+1$ states was found that is more likely than the configuration
244 with n states, keep it and go back to step 2.
- 245 4. For each state change in the configuration:
 - 246 (a) remove this state change.
 - 247 (b) find the most likely parameters for this state configuration, and
 - 248 (c) if the configuration without this state was more likely than the previous configuration,
249 keep it.
- 250 5. If at least one state was removed, go back to step 4.
- 251 6. Otherwise, end and record the most likely model.

252 **2.5 Implementation**

253 The likelihood calculation and Maximum Likelihood inference are implemented as a publicly
254 available R package. Partial results of the inference are automatically saved after each opti-
255 mization step, so that an interrupted run can be resumed at any point. The full results returned
256 include the best estimates for the number and positions of states, as well as all initial transmission
257 rates, transmission decay rates and removal rates of each state. An estimation of the uncertainty
258 around the result is provided by the maximum likelihood values found for each number of states
259 n up to $\tilde{n} + 1$ where \tilde{n} is the maximum likelihood inferred number of states.

260 All analysis, pre- and post-processing of the datasets were done using custom R scripts,
261 included in the Supplementary Materials.

262 **2.5.1 Time positions of state changes**

263 The model and the likelihood function allow for state changes to be placed anywhere on an edge.
264 The implementation of the algorithm allows for the time positions of changes to be estimated as
265 additional parameters, but this is computationally expensive especially when the number of state
266 changes grow. As a consequence we also provide the option to limit the positioning of changes
267 to predetermined positions on edges: they can be positioned at either 10%, 50% or 90% of the

268 length of the edge they are on. An intermediate option is also available, which will test all three
269 predetermined options and keep the most likely.

270 **2.5.2 Speed improvement option**

271 The algorithm as presented in the previous sections is fast at the beginning of the inference but
272 will progressively slow down as more states are added, due to the increase in the number of
273 parameters that need to be optimized.

274 We have thus added a so-called ‘fast optimization’ option, which limits the number of pa-
275 rameters which are allowed to change during one step of the maximum likelihood optimization.
276 In practice, when adding the n -th state change, only the parameters $\lambda_{0,n+1}$, $\lambda_{0,a}$, z_{n+1} , z_a , μ_{n+1}
277 and μ_a are optimized, where a is the state ancestral to the new state change. All other parame-
278 ters are fixed to the values inferred when adding the n -th state. Thus this option results in each
279 step of the algorithm having a constant cost instead of a cost dependent on n , however it will
280 lose some precision by fixing parameters.

281 It is to be noted that it is possible to run the normal analysis for the early steps of the
282 algorithm and turn on the fast optimization afterwards.

283 **3 Results**

284 **3.1 Cluster inference on simulated data**

285 **3.1.1 Dataset**

286 We use a simulated dataset produced by (10). This dataset contains simulated epidemics on three
287 different types of networks, A, B and C. The network structure A is composed of 13 communities
288 of 20 subjects each, with each community being a fully-connected graph and one bridge linking
289 any two communities.

290 The network structure B consists of one central community of size 60, representing a main
291 sexual contact network, connected by single bridges to 25 communities of size 20. Each small
292 community is a fully-connected graph. The set of small communities represents disjoint sexual
293 contact subnetworks in a population of interest.

294 The C networks are made of 100 communities each. The size of those clusters was sampled
295 from a distribution obtained from a phylogeny of the Swiss HIV Cohort Study (SHCS) dataset
296 (see (10) for details). To ensure that all communities are accessible, they are first linked in a
297 chain. Additional bridges are then created by connecting any two vertices belonging to different
298 communities with probability 0.00075.

299 In all types of networks, edges between communities are weighted, with the weight value 0.25,
300 0.5, 0.75, or 1. This means that the rate of transmission on these edges is respectively 25%, 50%,
301 75 % and 100% of the transmission rate on within-community edges.

302 Epidemics were simulated on these networks starting from one random introduction in A
303 networks, one random introduction in the main community in B networks, and two random
304 introductions in C networks. All infected individuals were sampled upon removal and a trans-
305 mission tree was built from the sampled tips. Thus there is no phylogenetic uncertainty in this
306 dataset: the tree represents exactly the progress of the simulated epidemic. For each type of net-
307 work (A,B,C) and each weighting scheme ($w=0.25,0.5,0.75$ or 1), 300 epidemics were simulated,
308 for a total dataset of 3600 trees.

309 Network structure B was designed to correspond best to the monophyletic assumption of
310 the cutpoint-based clustering methods: the epidemic starts in the main cluster and the smaller
311 islands are not connected with each other so all infections originating from the same population
312 cluster will be grouped in a single clade. Network structure A, on the other hand, allows for
313 the possibility of multiple introductions in the same population cluster and nested clusters, thus
314 breaking some of the assumptions of the cutpoint-based methods.

315 Various features of the A,B,C networks and the resulting simulated trees are shown in table
316 1. Networks A and B are very similar both in the size of their trees and in the cluster partition
317 inside trees. Network C, on the other hand, contains a large number of fairly small clusters.
318 Even though C trees are much larger on average, the clusters they contain are very small on
319 average and 34% of them include only 1 or 2 tips of the tree. These very small clusters contain
320 very little signal from the underlying contact network, and thus are not expected to be detected
321 by the method.

322 3.1.2 Comparison with cutpoint-based methods

323 We ran our maximum likelihood inference on the trees and inferred clusters by considering all
324 tips in the same state to be coming from the same community. In accordance with the simulation
325 conditions we set $\sigma = 1$ in the inference. The removal rates μ_i are assumed independent of the
326 population cluster, and so they are set to the same value μ for all states. The time positions of
327 the state changes were fixed using the intermediate option of testing positions at 10%, 50% and
328 90% of the length of the edges the state changes were on.

329 The correspondance between the real network communities and the clusters inferred from the
330 tree was assessed using the Adjusted Rand Index (ARI) (18; 19). We compare the results from
331 our method to the results obtained by (10) using cutpoint-based clustering methods.

332 Figure 2 shows the scores obtained by our MSBD method on the simulated A,B,C networks
333 compared to the scores of the cutpoint-based clustering methods. All methods used the same
334 cutpoints values, except for the method based on Definition 3 (Def3). Data corresponding to
335 this method was rescaled to fit in the same figure. As shown in (10), the results of the cutpoint-
336 based methods are highly variable and good scores can only be obtained from a narrow range of
337 cutpoints. In addition, the best cutpoint value is highly dependent on the underlying network
338 structure: in methods other than Def3, the best scores are obtained for a cutpoint of $c \approx 0.15$
339 for networks A, $c \approx 0.03$ for networks B and $c \approx 0.02$ for networks C. For Def3, the best score is
340 obtained for $c \approx 0.05$ for networks A, $c \approx 0.16$ for networks B and $c \approx 0.04$ for networks C. We
341 define the “peak range” of cutpoints for each method, network structure and weighting scheme
342 as the range of cutpoints which give a score which is at least 75% of the best score obtained for
343 any cutpoint. With this definition the peak ranges are very narrow, with an average length of
344 respectively 0.008, 0.015 and 0.016 for networks A, B and C in methods other than Def3. The
345 peak ranges obtained with Def3 are much wider, but a direct comparison is difficult due to the
346 different definition used for the cutpoint. In all methods the peak ranges for networks A and C
347 on one hand, and B on the other hand have very little overlap and the best cutpoint for C is
348 never found in the peak range of either A or B, and vice-versa. In conclusion it is impossible to
349 get good results from all network types with any single cutpoint value.

350 In addition the cutpoint-based methods are sensitive to network features and in particular to

351 the non-respect of the monophyletic assumption. In both the A and C networks, the best score
352 obtained by any cutpoint-based method is ≈ 0.45 for the weighting scheme $w = 0.25$ and ≈ 0.55
353 for $w = 1$, whereas it goes up to ≈ 0.85 and ≈ 0.9 , respectively, in networks B.

354 In comparison, the MSBD method performs less well on B networks, with an average score of
355 0.73 for $w = 0.25$ and 0.49 for $w = 1$. However, it performs much better on A networks, with an
356 average score of 0.64 for $w = 0.25$ and 0.53 for $w = 1$. The worst results are obtained on the C
357 networks, where the average score is ≈ 0.2 for all weights, less than half the best scores obtained
358 by cutpoint-based methods.

359 The low scores obtained on the C networks point to a potential limitation of our method on
360 the number of clusters that can be inferred from a tree. As seen in the network features, the
361 trees simulated on the C networks contain clusters which have less elements on average, and a
362 higher proportion of very small clusters. These clusters may be harder to detect due to their low
363 signal. To confirm this hypothesis we calculated the scores obtained by the MSBD method when
364 excluding all tips that belonged to a cluster with strictly less than 8 tips. The results are shown
365 in figure 2 (dotted line). The proportion of tips excluded by applying this criteria is shown in
366 table 2. The scores of all network structures and all weighting schemes improved when applying
367 this criteria. The improvement increased with the proportion of tips belonging to the excluded
368 clusters, supporting our hypothesis that the MSBD method has difficulty identifying them. In
369 particular, the MSBD scores on the C network structure for weight ≥ 0.5 increase to a level on
370 par with the scores obtained by cutpoint-based methods.

371 **3.2 Quality of the parameter inference**

372 To evaluate the performance of our MSBD method beyond cluster identification, we simulated
373 several datasets of 200 trees each under the multi-state birth-death process, with various param-
374 eter combinations. Simulations were done using Gillespie's algorithm. Tips were sampled upon
375 removal and the process was ran until the tree reached 50 sampled tips. Since these trees were not
376 built from network simulations we did not try to assess the quality of the cluster inference, but
377 we focused on the quality of the parameter inference and on whether our method can adequately
378 distinguish between trees that contain several clusters and trees that do not.

379 The results are summarized in Table 3. We can see that although the MSBD method is
380 able to consistently infer multiple clusters when they are present, it will also wrongly detect one
381 additional cluster in around 25% of the trees that only contain one cluster. This may be a problem
382 of noise, where due to the stochasticity of the simulation one subtree is slightly more likely to be
383 attributed different rates than the rest of the tree. This problem can be alleviated by looking at
384 the difference in the inferred transmission rates of each cluster, which are also outputted by our
385 method: a smaller difference is more likely to be indicative of noise. As previously noted, the
386 method also tends to underestimate the number of clusters in multi-cluster trees, mostly because
387 it cannot detect clusters below a certain size.

388 Regarding the parameter inference, the method has a slight bias towards overestimating
389 the transmission rate and underestimating the removal rate. This is potentially due to our
390 simulation process being conditioned on reaching 50 tips, which could bias datasets in favour of
391 trees showing apparent higher diversification rates. Overall, the absolute error on the inferred
392 parameters remain low compared to the true values, both in datasets with one cluster and in
393 datasets with multiple clusters.

394 In conclusion, the parameter inference from the MSBD method is reliable, although it suffers
395 from noise when applied to trees which contain only one cluster.

396 **3.3 Speed improvement option**

397 In this section we compared the performance of the “fast optimization” option and the regular
398 algorithm. We used a dataset of 300 trees of average size 200 tips, on which a partial inference
399 had already been performed, so the algorithm started from a saved state in which multiple state
400 changes had already been found. One optimization step of the algorithm was then performed, i.e
401 the inference added a state change on one edge of the tree. As shown in figure 3, we measured
402 both the speed-up resulting from using the faster option and the difference in the maximum log
403 likelihood found.

404 As expected the speed-up achieved increases with the number of states already present in the
405 tested configuration. At 5 state changes, the fast optimization is on average 10 times faster than
406 the regular one, with a number of outliers with speed-ups of up to 50 times. At 15 state changes

407 the speed-up is of 70 on average, a considerable improvement. The difference in the maximum
408 log-likelihood obtained using the less-precise fast option also increases with the number of state
409 changes, although the difference remains small compared to the log-likelihood value, which is on
410 average -1690 for the regular optimization across all categories. The runtimes for one edge are
411 on average 170s at 5 state changes and 1250s at 15 state changes for the regular optimization.
412 Since every step of the algorithm involves testing all edges of the tree, the “fast” option is thus
413 necessary to ensure completion of the inference on trees with more than 10 clusters.

414 **3.4 Cluster inference on HIV dataset**

415 In this section we analyze a tree used in another study of the correlation between sexual net-
416 works and tree features, (9). HIV-1 subtype B pol sequences were obtained from the Swiss
417 HIV Cohort Study 192 (SHCS). While the Swiss epidemic includes a mixture of population risk
418 groups including heterosexuals, injection drug users and MSM, only viral samples from MSM
419 were analyzed. A large cluster including 200 sampled individuals who predominantly lived or
420 sought treatment in the Zürich area was identified from a maximum likelihood (ML) phylogeny
421 of the complete dataset. The phylogeny of this cluster was then obtained by fitting a SIR-type
422 pairwise epidemic model to this sub-epidemic while simultaneously inferring the tree from the
423 sequence data in BEAST2. We re-analyze the tree provided for that cluster in the Supplement
424 of (9), this is a random tree from the posterior sample.

425 The results of the MSBD analysis on cluster 581 are shown in figure 4, part A. Three sub-
426 clusters are identified in the tree, one with a higher base transmission rate than in the backbone
427 of the tree, and two with similar base transmission rates which are lower than in the backbone
428 of the tree.

429 We compare our results to results obtained using the software Cluster Picker (20), which
430 detects clusters based on a combination of genetic distance between tip sequences and bootstrap
431 support at the nodes. It relies on two user-defined thresholds for both these measures, and so it
432 is a cutpoint-based method. Genetic sequences were generated for that tree using the software
433 SeqGen (21), using a GTR model with a gamma distribution with 4 rate categories and invariant
434 sites. The parameters of the molecular evolution model were set to the estimates obtained by

435 (9) when inferring the tree, which are shown in table 4.

436 As with other cutpoint-based methods, the results depend strongly on the user-defined values.
437 We used three different cutpoint values for the genetic distance: 1.5%, 4.5% and 8%. 4.5% is the
438 default value proposed by Cluster Picker and is the higher bound of the range recommended by
439 Cluster Picker for HIV data, whereas 1.5% is the lower bound of the recommended range. For the
440 bootstrap support threshold we used the value 0.0. With this cutpoint the bootstrap support is
441 disregarded entirely, which mimicks the behaviour of the methods studied by (10). The results
442 are shown in figure 4. We can see that the number of identified clusters is strongly dependent
443 on the cutpoint values, in keeping with the results obtained by (10). The size of the identified
444 clusters varies also widely, even within the bounds of the recommended range of cutpoints. With
445 the default setting of 4.5%, the clusters identified by MSBD are also recovered with Cluster
446 Picker, although one of the clusters is split in two in the clustering by Cluster Picker.

447 4 Discussion

448 We have introduced a novel method of identifying transmission clusters from a phylogeny, based
449 on a multi-state birth-death (MSBD) model. Our likelihood function makes two important
450 assumptions: the first one is that all the clusters in the transmission network appear in the
451 tree, and the second one is that unsampled subtrees, i.e subtrees that do not appear in the
452 reconstructed phylogeny, do not contain new introduction events. The implementation also relies
453 on a time discretization which approximates all transmission rates as locally constant on small
454 time intervals. A similar discretization can be applied to extend our method to time-dependent
455 removal rates, although the current implementation only allows removal rates to vary with state,
456 not through time.

457 This new method has a few key differences compared to the cutpoint-based clustering meth-
458 ods. Firstly, it is not restricted to monophyletic clades and can find clusters that are nested
459 within one another in the phylogeny. As a result our method clearly outperformed the others
460 on trial networks which were designed specifically to violate the monophyletic assumption. Sec-
461 ondly, as the MSBD method is model-based it does not rely on an arbitrary cutpoint ; this is
462 particularly important as (10) found that the quality of the clusters obtained by the cutpoint-

463 based clustering methods was extremely dependent on the value of this cutpoint, on all network
464 types. As seen in the results it is not possible to define a single cutpoint value as adequate
465 for all network types, which limits the usefulness of cutpoint-based methods in the absence of
466 prior information on the transmission network. The chosen cutpoint value is strongly linked with
467 the number of clusters inferred by cutpoint-based methods, thus obtaining the correct clusters
468 requires prior knowledge on the true number of clusters. Overall, while our method may not
469 perform as well on certain types of network as cutpoint-based methods, it is more reliable and
470 consistent in its results and does not require additional information from the user to get optimal
471 results.

472 As seen from the low scores obtained on the more fragmented trial networks and the improve-
473 ments obtained by limiting the size of the clusters to be detected, the MSBD method has a strong
474 limitation on the size of clusters that can be inferred from a tree. Contrary to the cutpoint-based
475 methods, which can handle arbitrary numbers and sizes of clusters, our method can only add
476 clusters when there is a strong signal for them and thus performs poorly in datasets with many
477 small clusters. As before though it should be noted that this low performance is compared to
478 the optimal results obtained by cutpoint-based methods, which require reliable information on
479 the expected number of clusters.

480 Another limitation of the current implementation is its computational cost, which limits the
481 size of the trees that can be analyzed in a reasonable time to a few hundred tips. Improving the
482 speed was the reason for several approximations such as the limitations on the positions of state
483 changes and the ‘fast optimization’ option, however these options necessarily limit the precision
484 of the results. Future work will focus on implementing the algorithm in parallel in order to
485 address this limitation.

486 Finally as a result of using a Maximum Likelihood framework, estimating the uncertainty
487 around the various estimated parameters is problematic, in particular for the positions and
488 number of state changes. The current implementation will estimate the uncertainty around the
489 number of states n by returning the best likelihood values for each n tested, although we expect
490 that better estimates could be achieved using a Bayesian framework.

491 **5 Acknowledgements**

492 We thank Dr. Timothy Vaughan for his comments and suggestions on the manuscript and Dr.
493 Luc Villandr e for his help with using the network simulation dataset.

494 **6 Funding statement**

495 JBS and TS are supported in part by the European Research Council under the Seventh Frame-
496 work Programme of the European Commission (PhyPD: grant agreement number 335529).

497 **7 Author contributions**

498 JBS implemented the model, performed the simulations, analysed the data and drafted the
499 manuscript. TS conceived of the study, designed the study, coordinated the study and helped
500 draft the manuscript. All authors gave final approval for publication.

501 References

- 502 [1] Allen LJS. In: Brauer F, van den Driessche P, Wu J, editors. An Introduction to Stochastic
503 Epidemic Models. Berlin, Heidelberg: Springer Berlin Heidelberg; 2008. p. 81–130. Available
504 from: http://dx.doi.org/10.1007/978-3-540-78911-6_3.
- 505 [2] Leventhal GE, Günthard HF, Bonhoeffer S, Stadler T. Using an Epidemiological Model for
506 Phylogenetic Inference Reveals Density Dependence in HIV Transmission. *Molecular Biol-*
507 *ogy and Evolution*. 2014;31(1):6. Available from: [+http://dx.doi.org/10.1093/molbev/](http://dx.doi.org/10.1093/molbev/mst172)
508 [mst172](http://dx.doi.org/10.1093/molbev/mst172).
- 509 [3] Girvan M, Newman MEJ. Community structure in social and biological networks. *Pro-*
510 *ceedings of the National Academy of Sciences*. 2002;99(12):7821–7826. Available from:
511 <http://www.pnas.org/content/99/12/7821.abstract>.
- 512 [4] Welch D. Is Network Clustering Detectable in Transmission Trees? *Viruses*. 2011;3(6):659–
513 676. Available from: <http://www.mdpi.com/1999-4915/3/6/659>.
- 514 [5] Robinson K, Fyson N, Cohen T, Fraser C, Colijn C. How the Dynamics and Structure of
515 Sexual Contact Networks Shape Pathogen Phylogenies. *PLOS Computational Biology*. 2013
516 06;9(6):1–15. Available from: <http://dx.doi.org/10.1371/journal.pcbi.1003105>.
- 517 [6] Leventhal GE, Kouyos R, Stadler T, von Wyl V, Yerly S, Böni J, et al. Inferring Epidemic
518 Contact Structure from Phylogenetic Trees. *PLOS Computational Biology*. 2012 03;8(3):1–
519 10. Available from: <http://dx.doi.org/10.1371/journal.pcbi.1002413>.
- 520 [7] Giardina F, Romero-Severson EO, Albert J, Britton T, Leitner T. Inference of Transmis-
521 sion Network Structure from HIV Phylogenetic Trees. *PLOS Computational Biology*. 2017
522 01;13(1):1–22. Available from: <https://doi.org/10.1371/journal.pcbi.1005316>.
- 523 [8] McCloskey RM, Liang RH, Poon AFY. Reconstructing contact network parameters from
524 viral phylogenies. *Virus Evolution*. 2016;2(2):vew029. Available from: [http://dx.doi.org/](http://dx.doi.org/10.1093/ve/vew029)
525 [10.1093/ve/vew029](http://dx.doi.org/10.1093/ve/vew029).

- 526 [9] Rasmussen DA, Kouyos R, Günthard HF, Stadler T. Phylodynamics on local sexual contact
527 networks. *PLOS Computational Biology*. 2017 03;13(3):1–23. Available from: [http://dx.
528 doi.org/10.1371%2Fjournal.pcbi.1005448](http://dx.doi.org/10.1371/journal.pcbi.1005448).
- 529 [10] Villandre L, Stephens DA, Labbe A, Günthard HF, Kouyos R, Stadler T, et al. Assessment
530 of overlap of phylogenetic transmission clusters and communities in simple sexual contact
531 networks: Applications to HIV-1. *PLoS ONE*. 2016;11(2):1–18.
- 532 [11] Maddison WP, Midford PE, Otto SP. Estimating a binary character’s effect on speciation
533 and extinction. *Systematic biology*. 2007 oct;56(5):701–10. Available from: [http://sysbio.
534 oxfordjournals.org/content/56/5/701.full](http://sysbio.oxfordjournals.org/content/56/5/701.full).
- 535 [12] FitzJohn RG. Diversitree : comparative phylogenetic analyses of diversification in R.
536 *Methods in Ecology and Evolution*. 2012 dec;3(6):1084–1092. Available from: [http:
537 //doi.wiley.com/10.1111/j.2041-210X.2012.00234.x](http://doi.wiley.com/10.1111/j.2041-210X.2012.00234.x).
- 538 [13] Stadler T, Bonhoeffer S. Uncovering epidemiological dynamics in heterogeneous host popu-
539 lations using phylogenetic methods. *Philosophical Transactions of the Royal Society B: Bio-
540 logical Sciences*. 2013;368(1614). Available from: [http://rstb.royalsocietypublishing.
541 org/content/368/1614/20120198](http://rstb.royalsocietypublishing.org/content/368/1614/20120198).
- 542 [14] Kühnert D, Stadler T, Vaughan TG, Drummond AJ. Phylodynamics with Migration: A
543 Computational Framework to Quantify Population Structure from Genomic Data. *Molecular
544 biology and evolution*. 2016 aug;33(8):2102–16.
- 545 [15] Rabosky DL, Santini F, Eastman J, Smith SA, Sidlauskas B, Chang J, et al. Rates of
546 speciation and morphological evolution are correlated across the largest vertebrate radiation.
547 *Nature communications*. 2013 jan;4:1958.
- 548 [16] Moore BR, Höhna S, May MR, Rannala B, Huelsenbeck JP. Critically evaluating the theory
549 and performance of Bayesian analysis of macroevolutionary mixtures. *Proceedings of the
550 National Academy of Sciences*. 2016 aug;p. 201518659. Available from: [http://www.pnas.
551 org/lookup/doi/10.1073/pnas.1518659113](http://www.pnas.org/lookup/doi/10.1073/pnas.1518659113).

- 552 [17] Stadler T, Kühnert D, Bonhoeffer S, Drummond AJ. Birth-death skyline plot reveals tem-
553 poral changes of epidemic spread in HIV and hepatitis C virus (HCV). Proceedings of the
554 National Academy of Sciences of the United States of America. 2013 jan;110(1):228–33.
- 555 [18] Rand WM. Objective Criteria for the Evaluation of Clustering Methods. Journal of the
556 American Statistical Association. 1971;66(336):846–850.
- 557 [19] Hubert L, Arabie P. Comparing partitions. Journal of Classification. 1985 Dec;2(1):193–218.
558 Available from: <https://doi.org/10.1007/BF01908075>.
- 559 [20] Ragonnet-Cronin M, Hodcroft E, Hué S, Fearnhill E, Delpech V, Brown AJL, et al. Au-
560 tomated analysis of phylogenetic clusters. BMC Bioinformatics. 2013;14(1):317. Available
561 from: <http://dx.doi.org/10.1186/1471-2105-14-317>.
- 562 [21] Rambaut A, Grass NC. Seq-Gen: an application for the Monte Carlo simulation of DNA
563 sequence evolution along phylogenetic trees. Bioinformatics. 1997;13(3):235. Available from:
564 [+http://dx.doi.org/10.1093/bioinformatics/13.3.235](http://dx.doi.org/10.1093/bioinformatics/13.3.235).

Network type	A	B	C
Number of clusters in the network	13	26	100
Number of elements per cluster	20	21.54	9.84
Number of tips per tree	52	60	196
Number of clusters per tree	5.95	6.63	39.10
Number of elements per cluster in the tree	9.45	9.57	5.17
Proportion of small clusters (<3 elements) in trees	21%	14%	34%

Table 1: General features of the A, B, C networks. All numbers are averages over the 4 weighting schemes, i.e averages over all 1200 trees in each network.

Network type	A	B	C
$w = 0.25$	9.02	9.87	33.9
$w = 0.5$	15.3	16	44.6
$w = 0.75$	21.7	20.8	51.4
$w = 1$	25.1	22.2	55.5

Table 2: Proportion (%) of tips belonging to clusters with strictly less than 8 tips, per network structure and weighting scheme.

Dataset parameters		$\lambda_0 = 25, z = 12,$ $\mu = 1, \gamma = 0$	$\lambda_0 = 25, z = 15,$ $\mu = 1, \gamma = 0$	$\lambda_0 = 10, z = 1,$ $\mu = 5, \gamma = 0.5$	$\lambda_0 = 10, z = 2,$ $\mu = 5, \gamma = 0.5$
Average number of clusters	simulated	1	1	4.95	6.38
	> 5 individuals, simulated	1	1	1.92	2.49
	inferred	1.22	1.25	2.43	2.65
Average transmission rate	simulated	1.09	0.86	6.95	5.40
	inferred	1.54	1.38	7.52	6.20
	median absolute error	0.37	0.49	0.75	0.78
Average removal rate	simulated	1.0	1.0	5.0	5.0
	inferred	0.88	0.91	4.64	4.50
	median absolute error	0.21	0.20	0.73	0.71

Table 3: Parameter inference on several datasets. Each dataset contains 200 trees of 50 tips each, simulated under a multi-state birth-death process using Gillespie’s algorithm. Transmission rates are averaged over the entire tree.

Parameter	Value used
Proportion of invariant sites	48%
Frequency of A	0.38
Frequency of C	0.16
Frequency of G	0.22
Frequency of T	0.24
Shape of gamma heterogeneity	0.57
Substitution rate	0.0015
Transition rate $A \rightarrow C$	0.23
Transition rate $A \rightarrow G$	1.12
Transition rate $A \rightarrow T$	0.09
Transition rate $C \rightarrow G$	0.14
Transition rate $C \rightarrow T$	1.0
Transition rate $G \rightarrow T$	0.11

Table 4: Parameter values used to simulate sequences with SeqGen.

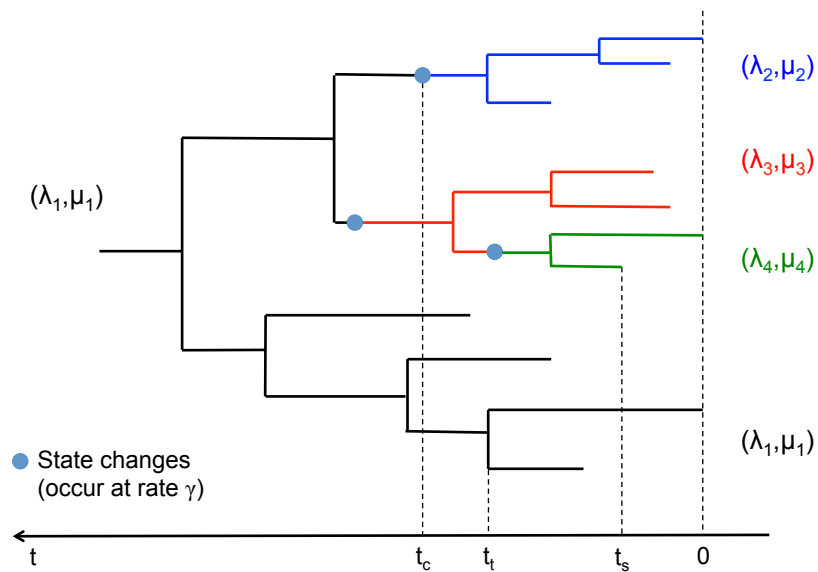


Figure 1: Visual representation of the phylogeny under a MSBD model. Each state is represented by a colour: the ancestral state, in black, starts at the root and represent the first cluster infected. The other states, in blue, red and green, start at change points along the tree. These states represent the clusters infected later in the course of the epidemic and the state change points represent the introduction event for each associated cluster.

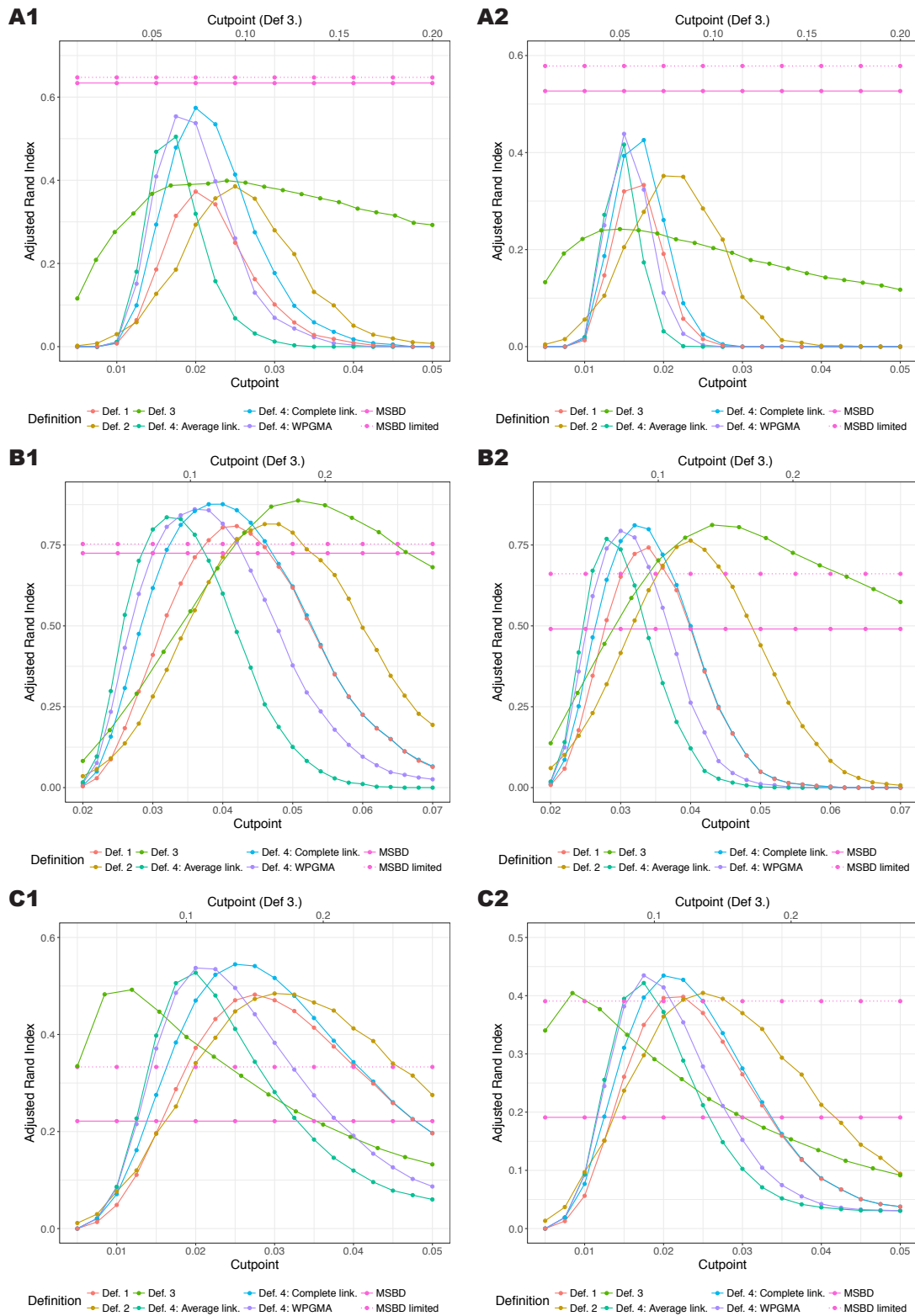


Figure 2: Comparison of the average ARI obtained by the different clustering methods in function of the set cutpoint on networks A (parts A1,A2), B (parts B1,B2) and C (parts C1,C2). For each network the first column (part 1) shows the results for weight $w = 0.25$ and the second column (part 2) for $w = 1$. Our proposed MSBD method is not dependent on a cutpoint. The cutpoint range for Definition 3. is shown on the x-axis on the top, the cutpoint range for all other definitions are shown on the x-axis at the bottom.

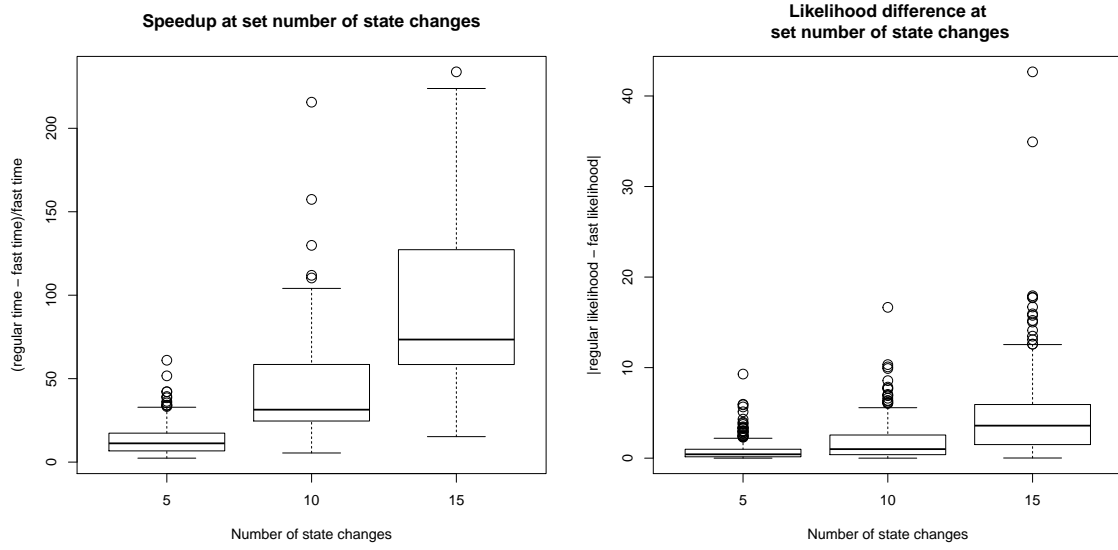


Figure 3: Box plots representing the speed-up (A) and likelihood difference (B) on one step of the algorithm when using the ‘fast optimization’ option compared to the default settings. The dataset used was divided in three categories based on the number of state changes already found by the inference before the test was run.

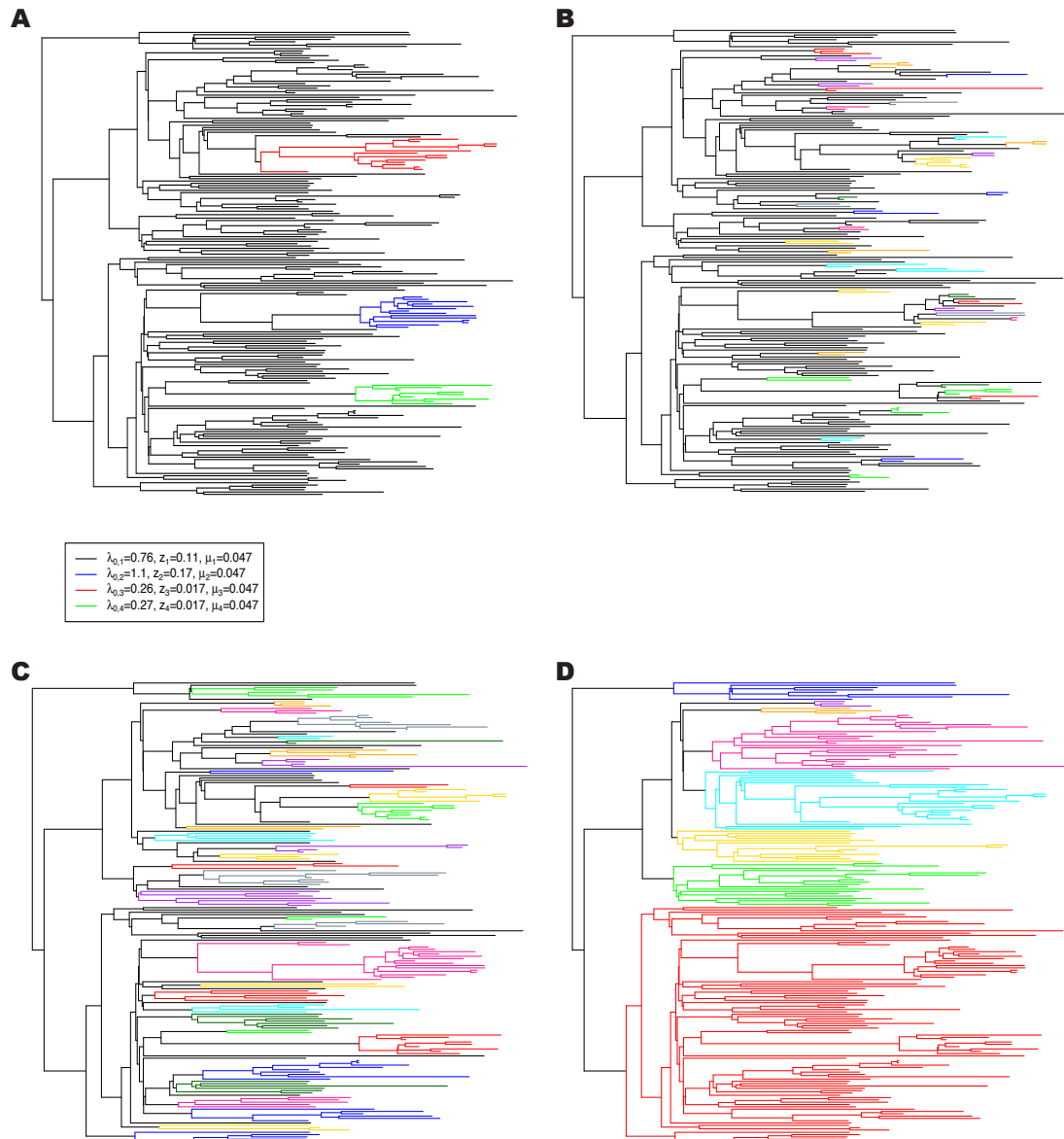


Figure 4: Comparison of the clusters obtained with MSBD (part A) or with Cluster Picker with a bootstrap threshold of 0.0 and a genetic distance threshold of 1.5% (part B), 4.5% (part C) and 8% (part D). Cluster Picker only identifies monophyletic clades as clusters, so each coloured clade is a separate cluster.