

1 **Inferring putative transmission clusters with Phydelity**

2

3 Alvin X. Han^{1,2,3*}, Edyth Parker^{3,4}, Sebastian Maurer-Stroh^{1,5} and Colin A. Russell^{3*}

4

5 ¹Bioinformatics Institute, Agency for Science, Technology and Research (A*STAR), 30

6 Biopolis Street, Singapore 138671

7 ²NUS Graduate School for Integrative Sciences and Engineering, National University of

8 Singapore (NUS), 21 Lower Kent Ridge, Singapore 119077

9 ³Laboratory of Applied Evolutionary Biology, Department of Medical Microbiology,

10 Academic Medical Centre, Meibergdreef 9, 1105 AZ Amsterdam-Zuidoost, The Netherlands,

11 ⁴Department of Veterinary Medicine, University of Cambridge, Madingley Rd, Cambridge

12 CB3 0ES, United Kingdom

13 ⁵Department of Biological Sciences, National University of Singapore, 16 Science Drive 4,

14 Singapore 117558

15

16 *Corresponding authors: hanxc@bii.a-star.edu.sg or c.a.russell@amc.uva.nl

17

18 **Abstract:** Current phylogenetic clustering approaches for identifying pathogen transmission

19 clusters are centrally limited by their dependency on arbitrarily-defined genetic distance

20 thresholds for within-cluster divergence. Incomplete knowledge of a pathogen's underlying

21 dynamics often reduces the choice of distance threshold to an exploratory, ad-hoc exercise

22 that is difficult to standardise across studies. Phydelity is a new tool for the identification of

23 transmission clusters in pathogen phylogenies. It identifies groups of sequences that are more

24 closely-related than the ensemble distribution of the phylogeny under a statistically-

25 principled and phylogeny-informed framework, without the introduction of arbitrary distance

26 thresholds. Relative to other distance threshold-based and model-based methods, Phydelity

27 outputs clusters with higher purity and lower probability of misclassification in simulated

28 phylogenies. Applying Phydelity to empirical datasets of hepatitis B and C virus infections

29 showed that Phydelity identified clusters with better correspondence to individuals that are

30 more likely to be linked by rapid transmission events relative to other widely-used

31 phylogenetic clustering methods without the need for parameter calibration. Phydelity is

32 generalisable to any pathogen and can be used to reliably identify putative direct transmission

33 events. Phydelity is freely available at <https://github.com/alvinxhan/Phydelity>.

34

35 **Introduction**

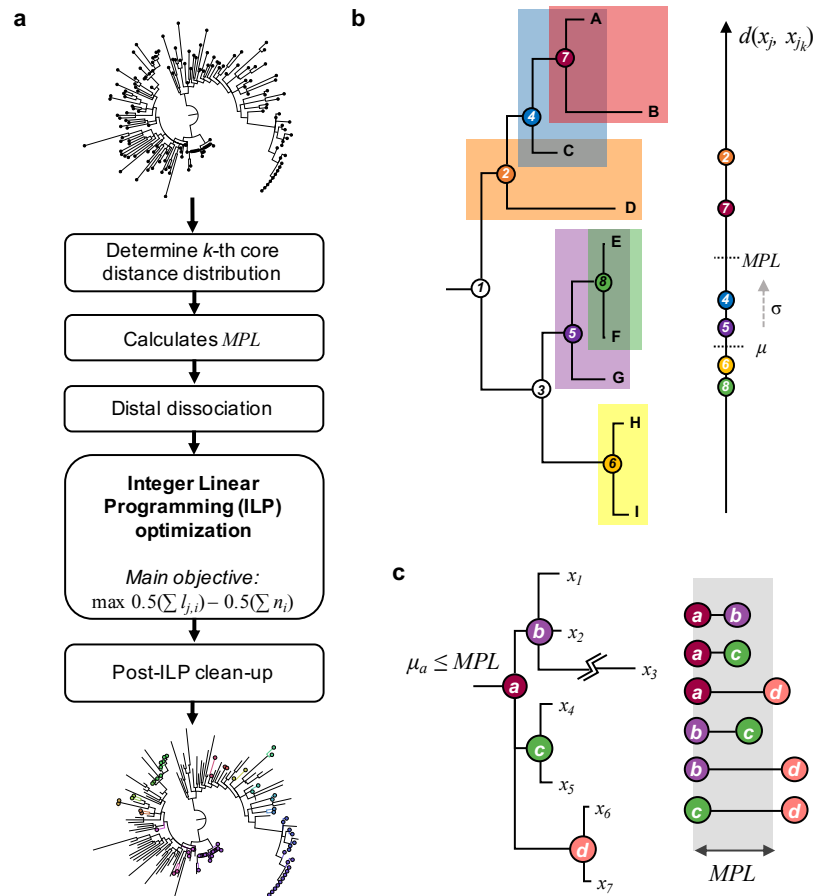
36

37 Recent advancements in high-throughput sequencing technologies have led to the widespread
38 use of sequence data in infectious disease epidemiology (Gardy and Loman 2017). In
39 particular, phylogenetics is frequently used to infer genetic clusters underlying the structure
40 of transmission networks (Ambrosioni et al. 2012; Bezemer et al. 2015; Matsuo et al. 2017;
41 de Oliveira et al. 2017; Charre et al. 2018). Current phylogenetic approaches for inferring
42 transmission clusters (primarily ‘cutpoint-based’ methods) are centrally limited by the need
43 to define arbitrary, absolute cluster divergence thresholds (Prosperi et al. 2011; Ragonnet-
44 Cronin et al. 2013). The lack of a consensus definition of a phylogenetic transmission cluster
45 (Grabowski and Redd 2014) coupled with incomplete knowledge of a pathogen’s underlying
46 epidemiological dynamics often reduces the choice of cutpoints to an *ad hoc* exploratory
47 exercise resulting in subjective cluster definitions.

48

49 Phydelity is a new tool for inferring putative transmission clusters through the identification
50 of groups of sequences that are more closely-related than the ensemble distribution under a
51 statistically-principled framework. Notably, Phydelity only requires a phylogeny as input,
52 negating the need to define arbitrary cluster divergence thresholds, and also only has a single
53 parameter that can either be user defined or determined directly by Phydelity. Phydelity is
54 freely available at <http://github.com/alvinxhan/Phydelity>.

55



56

57 **Figure 1. (a)** Phydality algorithm pipeline. Phydality considers the input phylogenetic tree as a collection of
 58 putative clusters each defined by an internal node i and tips j that it subtends. The algorithm first infers the k -th
 59 core distance distribution (\mathcal{D}_k) from the pairwise patristic distances of the closest k -neighbouring tips. k can be
 60 defined by the user or scaled by Phydality to obtain the supremum \mathcal{D}_k with the lowest divergence. \mathcal{D}_k is then
 61 used to compute the maximal patristic distance limit (MPL) under which tips are considered to be more closely-
 62 related than to the ensemble. Dissociation of distally related subtrees/sequences (Figure 1c) ensues such that
 63 both monophyletic and paraphyletic clustering structures can be identified. Phydality then incorporates the
 64 distance and topological information of the remaining nodes and tips into an integer linear programming (ILP)
 65 model to be optimised by clustering all tips that satisfy the relatedness constraints within the least number of
 66 clusters. Finally, post-ILP steps are implemented to remove any tips that may have been spuriously clustered.
 67 **(b)** Determination of the maximal patristic distance limit (MPL) using the median (μ) and robust estimator of
 68 scale (σ) based on the k -th core distance distribution (\mathcal{D}_k) of every sequence x_j and its k -closest neighbours
 69 ($d(x_j, x_{j_k}); k=2$ in this case). **(c)** Distal dissociation of a putative cluster subtended by internal node a where
 70 $\mu_a \leq MPL$. Sequence x_3 is dissociated from the putative cluster a due to its exceedingly long branch length
 71 violating the MPL threshold (i.e. $d(x_3, x_{3_k}) > MPL$). Additionally, subtree d is also dissociated from a as its
 72 inter-nodal distance with internal nodes b and c exceeds MPL .

73

74 Method

75 Clustering Algorithm

76 Phydality considers the input phylogeny as an ensemble of putative clusters, each consisting
 77 of an internal node i and the leaves it subtends. The within-cluster diversity of node i is
 78 measured by its mean pairwise patristic distance (μ_i) Sequences subtended by i are

79 considered for clustering if μ_i is less than the maximal patristic distance limit (*MPL*), under
80 which sequences are considered more closely-related to one another than the ensemble
81 distribution (Figure 1).

82

83 Phydelyity computes the *MPL* by first calculating the pairwise patristic distance distribution
84 (i.e. k -th core distance distribution, \mathcal{D}_k) of closely-related tips comprising the pairwise
85 patristic distances of sequence x_j to the closest k -neighbouring tips (i.e. $d(x_j, x_{j_k}) = d_l$)
86 wherein their closest k -neighbours include sequence x_j as well. Additionally, \mathcal{D}_k is
87 incrementally sorted ($d_l \leq d_{l+1}$) and truncated up to d_L if the log difference between d_L and
88 d_{L+1} is more than zero:

$$89 \quad \mathcal{D}_k = \left\{ d_1, \dots, d_l, d_{l+1}, \dots, d_L \mid d_l \leq d_{l+1}, \lg\left(\frac{d_{l+1} - d_l}{d_l}\right) \leq 0 \right\}$$

90 The user can opt to either input the desired k parameter or allow Phydelyity to automatically
91 scale k to the value that yields the supremum k -th core distance distribution with the lowest
92 overall divergence. This is done by testing if \mathcal{D}_{k+1} and \mathcal{D}_k are statistically distinct ($p < 0.01$)
93 using the Kuiper's test (see Supplementary Materials). All clustering results of Phydelyity
94 presented in this work were performed using the autoscaled value of k .

95

96 The *MPL* is then calculated by:

$$97 \quad MPL = \bar{\mu} + \sigma$$

98 where $\bar{\mu}$ is the median pairwise distance of \mathcal{D}_k and σ is the corresponding robust estimator of
99 scale without assuming symmetry about $\bar{\mu}$ (Figure 1b, see Supplementary Materials).

100

101 This is then followed by dissociation of distantly-related descendant subtrees/sequences to all
102 putative nodes where $\mu_i > MPL$, thereby facilitating identification of both monophyletic as
103 well as nested, paraphyletic clusters (Figure 1c; see Supplementary Materials). Phydelyity
104 filters outlying tips from putative clusters under the assumptions that viruses infecting
105 individuals in a rapid transmission chain likely coalesce to the same most recent common
106 ancestor (MRCA). Additionally, Phydelyity requires any clonal ancestors in between the
107 MRCA and tips of a putative cluster to be as genetically similar to each other as they are to
108 the tips of the cluster. As such, for a putative transmission cluster, the mean pairwise nodal
109 distance between all internal and tip nodes of a cluster must also be $\leq MPL$.

110

111 An integer linear programming (ILP) model is implemented and optimised under the
112 objective to assign cluster membership to sequences satisfying the aforementioned
113 relatedness criteria within the least number of clusters. In other words, Phydelity uses ILP
114 optimisation to search for the clustering configuration that favours the designation of larger
115 clusters of closely-related sequences which are likely linked by rapid transmission events.
116 Lastly, topologically outlying singletons that were spuriously clustered are removed. The full
117 algorithm description and mathematical formulation of Phydelity is detailed in
118 Supplementary Materials.

119

120 *Assessing clustering results of simulated epidemics*

121 Phydelity was evaluated on phylogenetic trees derived from simulated HIV epidemics of a
122 hypothetical men who have sex with men (MSM) sexual contact network (C-type networks in
123 Villandre *et al.*, 2016). The simulated sexual contact network comprised 100 subnetworks
124 (communities) sampled from an empirical distribution obtained from the Swiss HIV Cohort
125 Study. All communities were linked in a chain initially and additional connections were
126 generated at a probability of 0.00075. Subjects in the network can either be in the
127 “susceptible”, “infected” or “removed” (i.e. individual is diagnosed and sampled) state.
128 Quick transmission chains (i.e. transmission clusters) were attributed to sexual contact among
129 individuals belonging to the same community.

130

131 300 epidemics were simulated for four different weights of inter-community transmission
132 rates (i.e. $w = 25\%$, 50% , 75% or 100% of the within-community rate). Two infected
133 individuals were randomly introduced in any of the 100 communities. Transmission time
134 along an edge followed an exponential distribution with rates directly proportional to the
135 associated weights. Time until removal was based on a shifted exponential distribution with
136 the shift representing the minimum amount of time required for a virus to be transmitted to
137 susceptible neighbours. The simulation ended once 200 individuals were in the “removed”
138 state.

139

140 These simulated datasets were tested by Villandre *et al.* (2016) to compare the outputs of four
141 “cutpoint-based” phylogenetic clustering methods where the arbitrary distance threshold
142 defining a transmission cluster (i.e. cutpoint) was computed as the: (i) absolute patristic
143 distance threshold between any two tips (Brenner *et al.* 2007); (ii) standardised number of
144 nucleotide changes (i.e. ClusterPicker, Ragonnet-Cronin *et al.*, 2013); (iii) percentile of the

145 phylogeny's pairwise sequence patristic distance distribution (i.e. PhyloPart, Prosperi *et al.*,
146 2011) and (iv) height of an ultrametric tree obtained using the weighted pair-group method of
147 analysis (WPGMA). For each method, Villandre *et al.* varied the corresponding cutpoint
148 parameter over an equivalent range of thresholds. Comparing the output clusters generated by
149 the four methods at their respective optimal cutpoint by adjusted rand index (see below), it
150 was found that the WPGMA method tended to produce clusters with better correspondence to
151 the underlying sexual contact structure. As such, clustering results from Phydelyty were
152 compared to those obtained by Villandre *et al.* using the WPGMA method. Additionally,
153 Phydelyty was also compared to the multi-state birth-death (MSBD) method which inferred
154 transmission clusters on the same simulated datasets by detecting significant changes in
155 transmission rates (Barido-Sottani et al. 2018).

156

157 To assess and compare the output clusters from Phydelyty and the aforementioned clustering
158 methods that had been tested on these networks previously, several metrics were used to
159 measure how well the clustering results corresponded with the known sexual contact
160 network:

- 161 i. Adjusted rand index (*ARI*) measures the accuracy of the clustering results by computing
162 the frequencies whereby a pair of sequences of the identical (or distinct) subnetwork(s)
163 was assigned to the same (or different) cluster(s) (Hubert and Arabie 1985). *ARI* ranges
164 between -1 (matching between output clusters and community labels is worse than
165 random clustering) and 1 (perfect match between output clusters and ground truth).
- 166 ii. Modified Gini index (I_G). Gini impurity, commonly used in decision tree learning,
167 refers to the probability of a randomly selected item from a set of classes would be
168 incorrectly labelled if it was randomly labelled by the distribution of occurrences in the
169 class set (Breiman et al. 1984). Here, I_G measures how often a randomly selected
170 sequence from the given network would be incorrectly clustered by the inferred
171 clusters. For a sexual contact network with T communities (i.e. $t \in \{1, 2, \dots, T\}$), I_G is
172 computed as:

$$173 \quad I_G = \sum_{t=1}^T \left[p_t \left(1 - \sum_{c=1}^{c^*} p(c|t) \right) \right]$$

174 where C^* is the set of clusters defined to have correctly classified sequences attributed
175 to community t (i.e. any cluster that constitutes the largest proportion of sequences
176 from community t at both the cluster and the community label levels), p_t is the

177 probability of sequence from community t and $p(c|t)$ refers to the probability that a
178 sequence is clustered under cluster c conditional of it being from community t . If output
179 clusters perfectly align with the underlying sexual contact network (i.e. one cluster only
180 constitute one class of community), $I_G = 0$. Conversely, if clustering results are
181 completely random, $I_G = 1$.

182 iii. Purity measures the average extent that the output clusters contain only a single class
183 (i.e. a particular sexual contact community; Manning et al. 2008):

$$184 \quad Purity = \sum_{c=1}^C \frac{1}{N_c} \left(\frac{\max_t \{N_{c,t}\}}{N_c} \right)$$

185 where N_c is the size of cluster c , $N_{c,t}$ is the number of tips from community t clustered
186 under cluster c and C is the set of all output clusters. Note that purity (as well as I_G) can
187 be inflated if the total number of clusters is large (i.e. if each tip is assigned to a unique
188 cluster, purity = 1 and $I_G = 0$).

189 iv. Normalised mutual information (*NMI*) trades off the output clustering quality against
190 the number of clusters (Manning et al. 2008):

$$191 \quad NMI = \frac{I(T, C)}{[H(T) + H(C)]/2}$$

192 where $H(T)$ and $H(C)$ are the respective entropies of the network communities and
193 output clusters, and $I(T, C)$ is the mutual information between them. If clustering is
194 random with respect to the network community labels, $I(T, C) = 0$ (i.e. $NMI = 0$). On
195 the other hand, maximum mutual information is achieved (i.e. $I(T, C) = I(T, C)_{max}$)
196 either when the output clusters map the sexual contact network perfectly or all clusters
197 have one member only. Hence, to penalise large cardinalities while normalising $I(T, C)$
198 between 0 and 1, *NMI* is calculated since (a) entropy increases with increasing number
199 of clusters and (b) $[H(T) + H(C)]/2$ is a tight upper bound to $I(T, C)$.

200

201 *Comparison to ClusterPicker and PhyloPart*

202 Phydelity was also tested on two empirical datasets – acute hepatitis C infections among men
203 who have sex with men and hepatitis B viruses collected from members of the same families.
204 Phylogenetic trees for both datasets were reconstructed using RAxML under the
205 GTRGAMMA model (Stamatakis 2014).

206

207 ClusterPicker (Ragonnet-Cronin et al. 2013) and PhyloPart (Prosperi et al. 2011), two
208 popular phylogenetic clustering tools that are methodologically comparable to Phydelyity,
209 were also applied to the same datasets for comparisons. However, other than the phylogenetic
210 tree, both ClusterPicker and PhyloPart also require users to input an arbitrarily-defined
211 genetic distance threshold (as an absolute distance limit for ClusterPicker and percentile of
212 the global pairwise patristic distance for PhyloPart). As such, a range of distance limits
213 (PhyloPart: 0.5-10th percentile; ClusterPicker: 0.005-0.1 nucleotide/site) were applied to both
214 tools. No bootstrap support threshold were implemented for comparability to Phydelyity.

215
216 The lowest optimal threshold for the distance range tested was found by maximisation of the
217 mean silhouette index (*SI*) for both ClusterPicker and PhyloPart. The Silhouette index
218 measures how similar an item is to members of its own cluster as opposed to the nearest
219 neighbouring clusters - i.e. a larger mean silhouette index indicates that items of the same
220 cluster are more closely related amongst themselves than to its neighbours (Rousseeuw,
221 1987). No parameter optimisation was required for Phydelyity.

222

223 **Results**

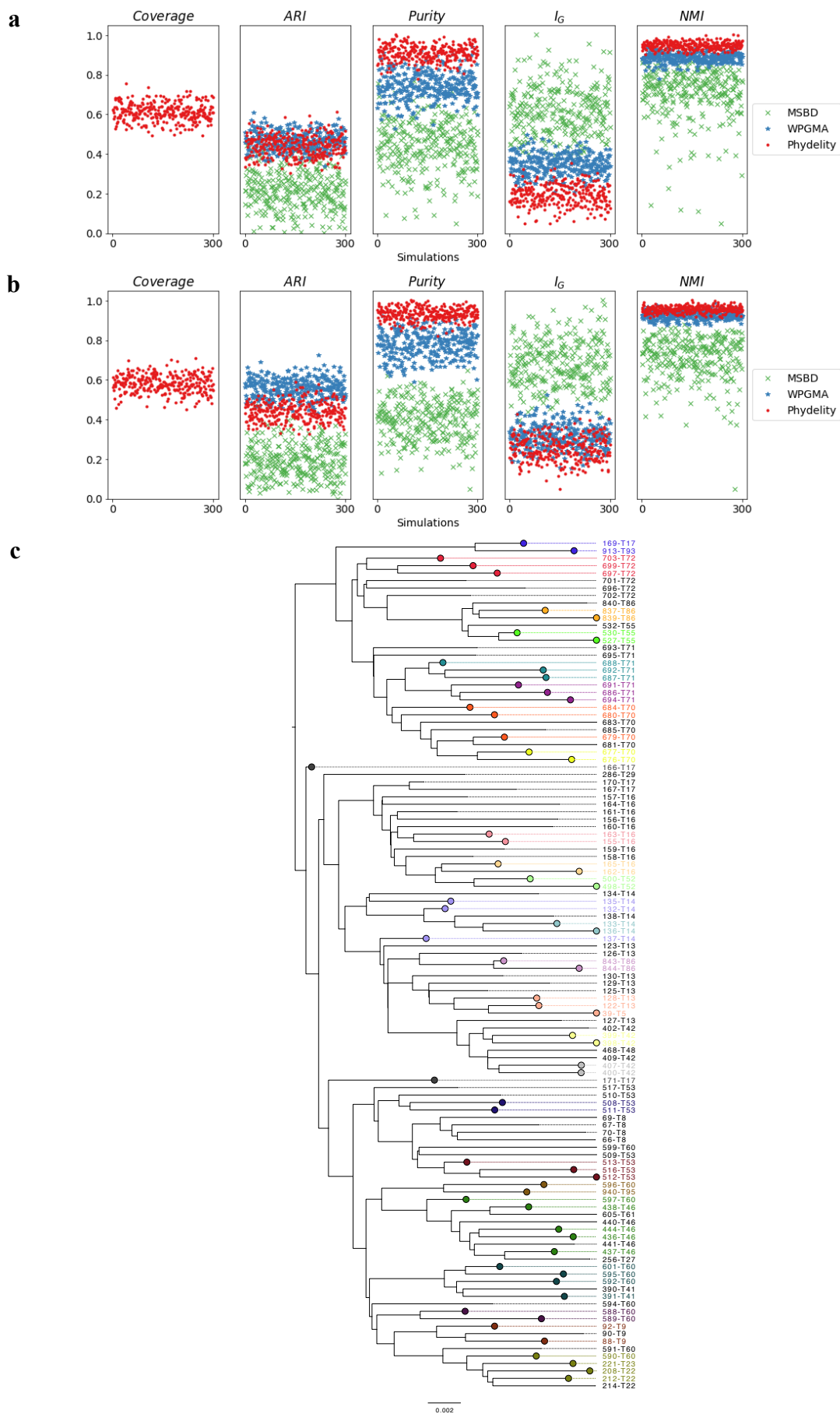
224 *Simulated HIV epidemics*

225 Phydelyity was applied to simulated HIV epidemics among men who have sex with men
226 (MSM) belonging to a hypothetical sexual contact network structures where transmission
227 clusters were attributed to quick transmission chains due to sexual contact among individuals
228 belonging to the same subnetwork (see Methods; Villandre *et al.*, 2016). These simulations
229 were originally used to assess the performance of “cutpoint-based” clustering tools, including
230 ClusterPicker, PhyloPart as well as the weighted pair-group method of analysis (WPGMA)
231 which generally attained the highest adjusted rand-index (ARI) score across all simulations
232 when calibrating their respective cutpoint thresholds against the ground-truth. Phylogenetic
233 trees generated from these simulations were also tested by the multi-state birth death (MSBD)

234 method (Barido-Sottani et al. 2018).

235

236 Clustering results from Phydelyity were compared to outputs from the MSBD method and
237 those achieving the best ARI scores based on the WPGMA method. The purity, modified
238 Gini index (I_G) and normalised mutual information (*NMI*) measures were also used to
239 provide a more comprehensive assessment of the clustering results (Figure 2, Supplementary
240 Figure 3 and Supplementary Table 1; see Methods).



242 **Figure 2.** Clustering results of simulated HIV epidemics in a hypothetical MSM sexual contact network. **(a)**
243 Clustering metrics for clustering algorithms (Phydelity, weighted pair-group method of analysis (WPGMA) and
244 multi-state birth death (MSBD) methods) applied simulated phylogenies with inter-communities transmission
245 rates weighted at half of within-community rates (i.e. $w = 0.5$). Coverage refers to the proportion of tips
246 clustered by Phydelity. Adjusted rand index (*ARI*) measures how accurate the output clusters corresponded with
247 the community labels. *Purity* gives the average extent clusters contain only a single class of community.
248 Modified Gini index (I_G) is the probability that a randomly selected sequence would be incorrectly clustered.
249 Normalised mutual information (*NMI*) accounts for the tradeoff between clustering quality and number of
250 clusters. **(b)** Results for simulations where inter-communities transmission rates were identical to within-
251 community rates (i.e. $w = 1.0$). **(c)** Sample output clusters of Phydelity for a subtree of an example simulation (w
252 = 0.5). Tips that were clustered by Phydelity are distinctly coloured according to their cluster membership. By
253 relaxing the monophyletic assumption, Phydelity is capable of detecting paraphyletic clusters (e.g. transmission
254 pair 166-T17 and 171-T17 and cluster subtending 132-T14, 135-T14 and 137-14).

255

256 The phylogenetic trees generated from the simulations had a large number of clusters that
257 were relatively small in size (i.e. percentage of sequences that were part of ground truth
258 clusters with sizes < 8 tips = 33.9% ($w = 25\%$); 55.5% ($w = 100\%$); see Barido-Sottani *et al.*
259 (2018) for more details). Furthermore, these ground truth clusters were not all monophyletic
260 (Figure 2c). As a result, while Phydelity and WPGMA yielded comparable ARI scores
261 (Phydelity: 0.44-0.45 (s.d. = 0.05); WPGMA: 0.44-0.56 (s.d = 0.05-0.05); Supplementary
262 Table 1), the output clusters of Phydelity, which can be paraphyletic (Figure 2c), are purer
263 (mean purity; Phydelity: 0.81-0.88 (s.d. = 0.03); WPGMA: 0.67-0.74 (s.d. = 0.06-0.06)) and
264 have a lower probability of misclassification when compared to WPGMA which assumes
265 clusters are strictly monophyletic (mean I_G ; Phydelity: 0.27-0.28 (s.d. = 0.04-0.05);
266 WPGMA: 0.33-0.40 (s.d. = 0.04-0.05)). Coverage of sequences clustered by Phydelity lies
267 between 58.2% and 61.6%.

268

269 The clustering results from WPGMA presented in this work were based on the optimal
270 distance threshold derived by calibration against the simulated ground-truth. Notably,
271 Phydelity's auto-scaling mitigates the need for threshold calibration and enables application
272 to empirical datasets where ground truth clustering is unavailable, as is largely the case for
273 epidemiological studies.

274

275 *Hepatitis B virus transmission between family members*

276 Phydelity was tested on empirical datasets to demonstrate its applicability on real-world data,
277 including hepatitis B viruses (HBV) collected from residents in the Binh Thuan Province of
278 Vietnam. In such highly endemic regions, HBV is commonly transmitted either vertically

279 from mothers to children during the perinatal period or horizontally between cohabitants of
280 the same household (Matsuo et al. 2017).

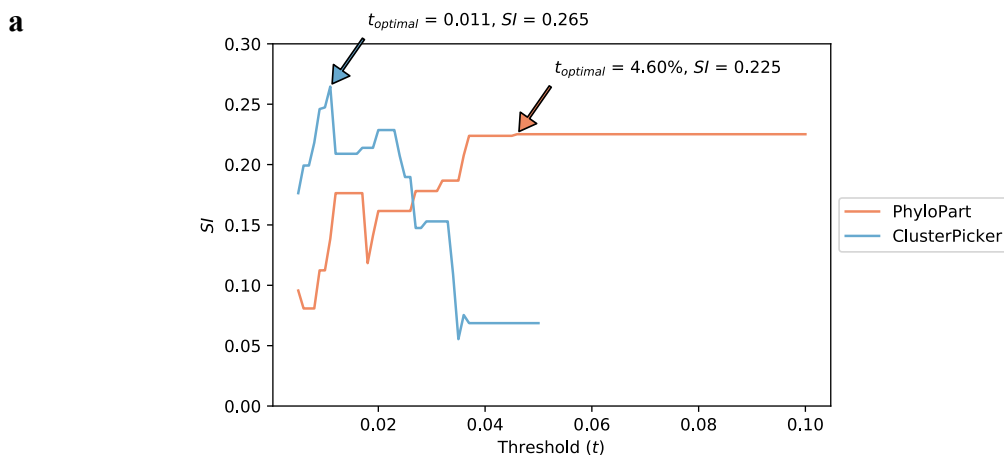
281

282 As complete genome nucleotide sequences were not available for all individuals, a
283 phylogenetic tree was reconstructed using the viral polymerase sequences collected from 41
284 patients, of which 12 of them were confirmed to be members of three families (i.e. denoted as
285 F2, F3 and F4) by a family survey as well as mitochondrial analyses. Besides Phydelyty, the
286 resulting phylogeny was also implemented in ClusterPicker and PhyloPart for comparison.
287 While WPGMA performed better in the simulations by Villandre *et al.*, ClusterPicker and
288 PhyloPart are arguably the more widely-used phylogenetic clustering tools to date.

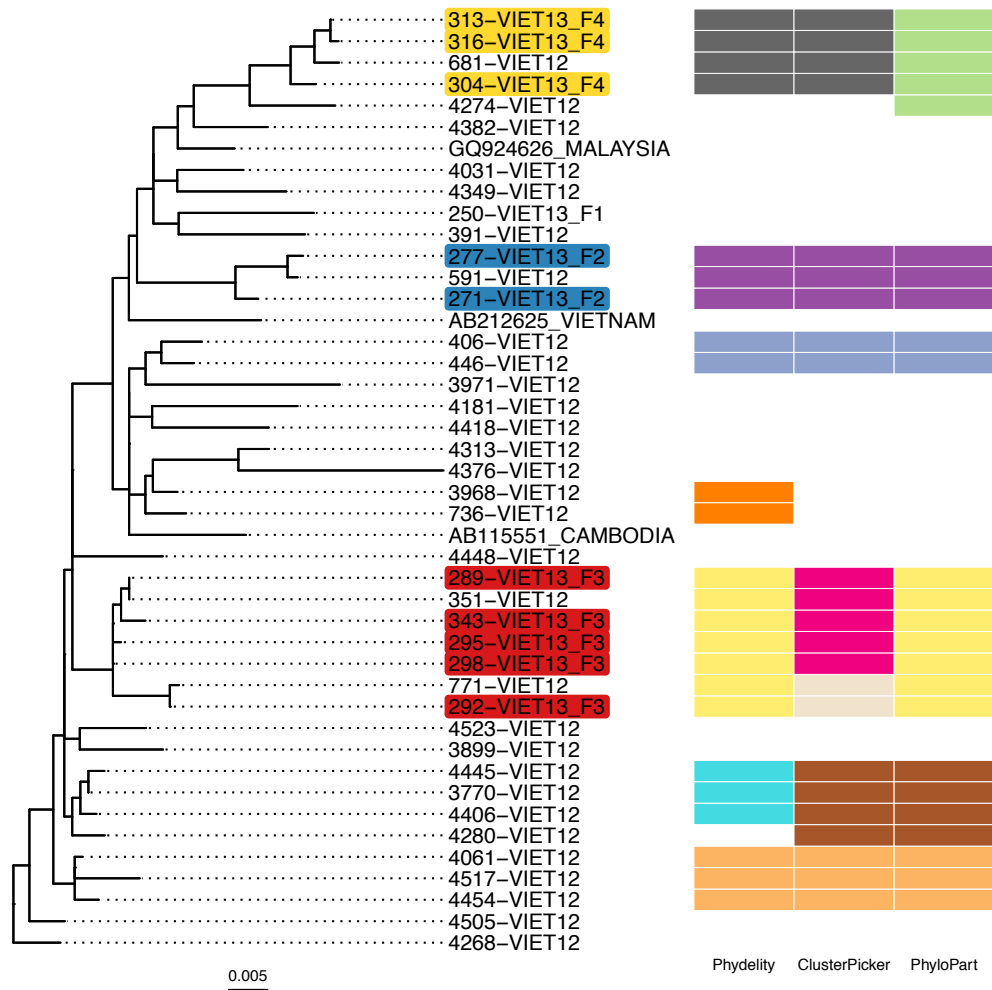
289

290 Phydelyty identified three likely transmission clusters that distinguish between the separate
291 family households (Figure 3). At their respective optimal distance thresholds by mean
292 Silhouette index (see Methods), ClusterPicker and PhyloPart achieved similar clustering
293 results as well. Importantly, however, Phydelyty was able to obtain the same optimal
294 clustering results without having to optimise and implement a hard-to-interpret distance
295 parameter.

296



b



297 **Figure 3.** Clustering results of hepatitis B viruses (HBV) collected from residents in the Binh Thuan Province
 298 of Vietnam. **(a)** Plots of mean Silhouette index (SI) computed for the range of genetic distance thresholds
 299 implemented in ClusterPicker and PhyloPart. Clustering results from the lowest optimal distance threshold
 300 ($t_{optimal}$) with the highest SI value for each method were compared to Phydelity as depicted in **b** (ClusterPicker:
 301 $t_{optimal} = 0.011$ nucleotide/site, $SI = 0.265$; PhyloPart: $t_{optimal} = 4.60\%$, $SI = 0.225$). Plot for ClusterPicker is
 302 truncated at ~ 0.05 nucleotide/site as the entire tree collapsed to single cluster after this threshold. **(b)** Maximum
 303 likelihood phylogeny of HBV polymerase sequences derived from viruses collected from 41 patients. 12
 304 patients were confirmed to be members of three separate family households (F2, F3 and F4; tip names shaded
 305 with a distinct colour for each family). Clustering results from Phydelity are depicted as a heatmap alongside
 306 outputs from ClusterPicker and PhyloPart based on their respective $t_{optimal}$. Each distinct colour of the heatmap
 307 cells denotes a different cluster.

308

309 *Hepatitis C virus transmission among MSM*

310 Incidence of HCV infections among HIV-negative MSM has been relatively limited as
 311 compared to their HIV-positive counterparts. However, the recent uptake of pre-exposure
 312 prophylaxis (PrEP) among HIV-negative individuals to prevent HIV infection could pose
 313 higher risk of sexually transmitted HCV infections (Volk et al. 2015; Charre et al. 2018). In a

314 study on HIV-positive and HIV-negative MSM patients in Lyon, 108 cases of acute HCV
315 infections (80 primary infections; 28 reinfections) were reported between 2014 and 2017
316 among 96 MSM (72 HIV-positive; 24 HIV-negative, of which 16 (67%) of them were on
317 PrEP). Separate phylogenetic analyses were performed on a subset of 89 (68 HIV-positive;
318 21 HIV-negative) HCV isolates belonging to genotypes 1a and 4d based on their NS5B
319 sequences. Additionally, 25 HCV sequences from HIV-infected MSM collected before 2014
320 were included along with 60 control HCV sequences derived from HIV-negative, non-MSM
321 patients residing in the same geographical area as controls. All sequences collected from
322 MSM patients were given strain names in the format of “MAH(ID)_accession” while control
323 sequences from non-HIV, non-MSM patients were denoted as “NCH(ID)_accession” (Figure
324 4). Phydelity as well as ClusterPicker and PhyloPart were applied to the reconstructed
325 phylogenies, with the latter calibrated over a range of distance thresholds. Again, only
326 clustering results based on the lowest distance threshold maximising the mean Silhouette
327 index for ClusterPicker and PhyloPart were compared to Phydelity’s output clusters (see
328 Methods).

329

330 Generally, membership of the MSM transmission clusters and pairs identified by Phydelity
331 across both genotypes were strictly limited to sequences derived from MSM patients.
332 Relaxing the monophyletic assumption by dissociating distantly-related tips from putative
333 monophyletic clusters (see Methods) enables Phydelity to identify likely outlying sequences
334 as evidenced by their relatively longer branch lengths from the cluster ensemble (Table 1 and
335 Figure 4; Genotype 1a: cluster C1 – MAH66 and cluster C3 – MAH31, MAH62 and
336 MAH72; Genotype 4d: cluster C3 – MAH24 and MAH08). In particular, for genotype 1a,
337 even though the mean pairwise distance of MAH72 to members of cluster C3 is within a
338 standard deviation of the latter’s within-cluster diversity, its distance to the more distant
339 members (e.g. MAH15 and MAH40, Figure 4) violated the inferred *MPL* (Table 1).
340 Additionally, as a result of distal dissociation, Phydelity distinguishes clusters that are
341 genetically more alike amongst themselves than to those phylogenetically ancestral to it (e.g.
342 cluster C1.1 from C1 for genotype 1a; Figure 4a).

343

344 For both genotypes, Phydelity found multiple clusters that included both HIV-positive and
345 HIV-negative MSM patients (i.e. Genotype 1a: clusters C2 and C3; Genotype 4d: clusters C2
346 and C2.2, as well as pair P2). While it is not clear which of the HIV-negative patients were
347 on PrEP (information is not given in the original paper), the clustering results from Phydelity

348 were in line with the findings by Charre *et al.* that acute HCV infections among HIV-
 349 negative MSM were likely sourced from their HIV-positive counterparts.

350

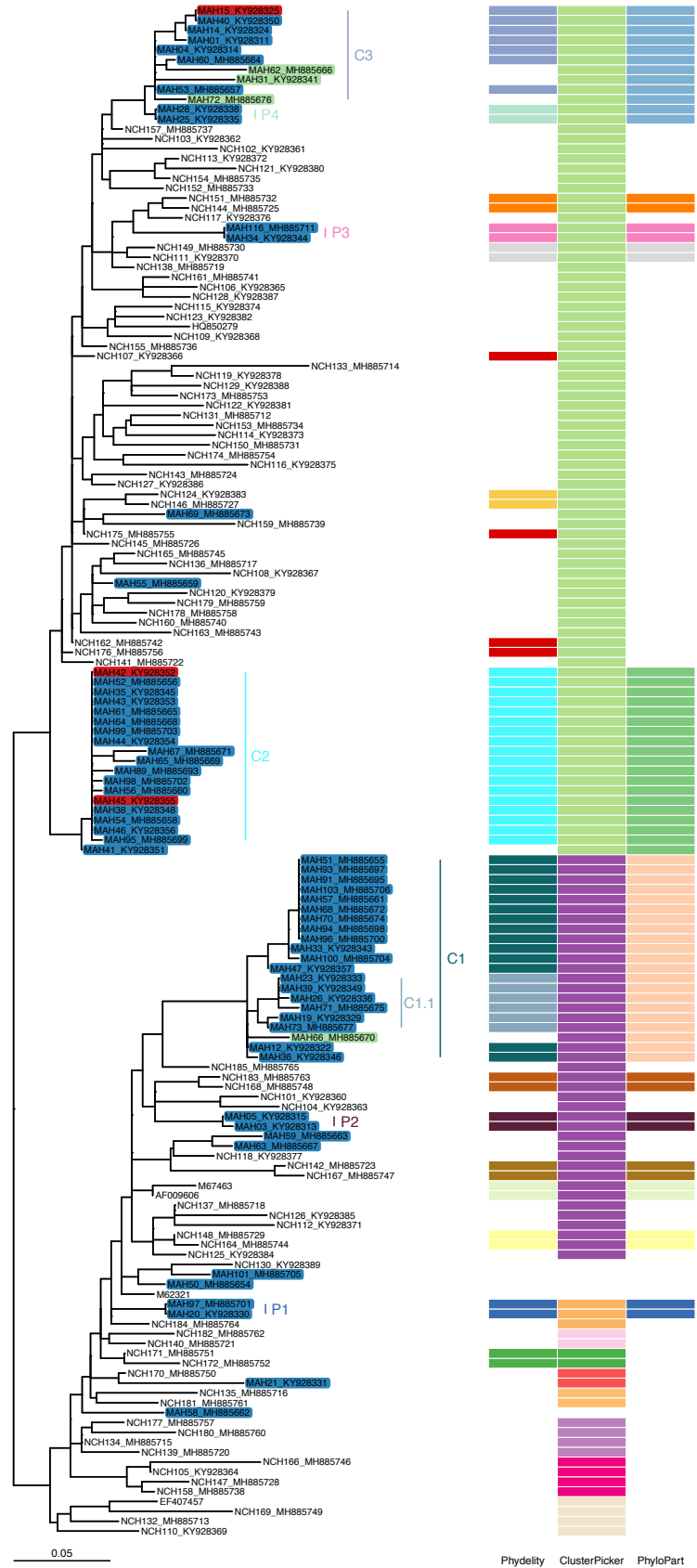
351 While ClusterPicker did manage to consolidate all of the MSM genotype 4d sequences into a
 352 single monophyletic cluster, its clustering of genotype 1a was clearly problematic as a large
 353 number of non-MSM control sequences were clustered together with those from MSM
 354 patients. PhyloPart's optimal clustering output was consistent Phydely's for genotype 1a.
 355 However, the larger number of identical sequences in the genotype 4d tree skewed the
 356 optimal distance parameter (expressed as x -th percentile of the pairwise patristic distribution
 357 of the entire phylogeny) to only cluster these identical sequences.

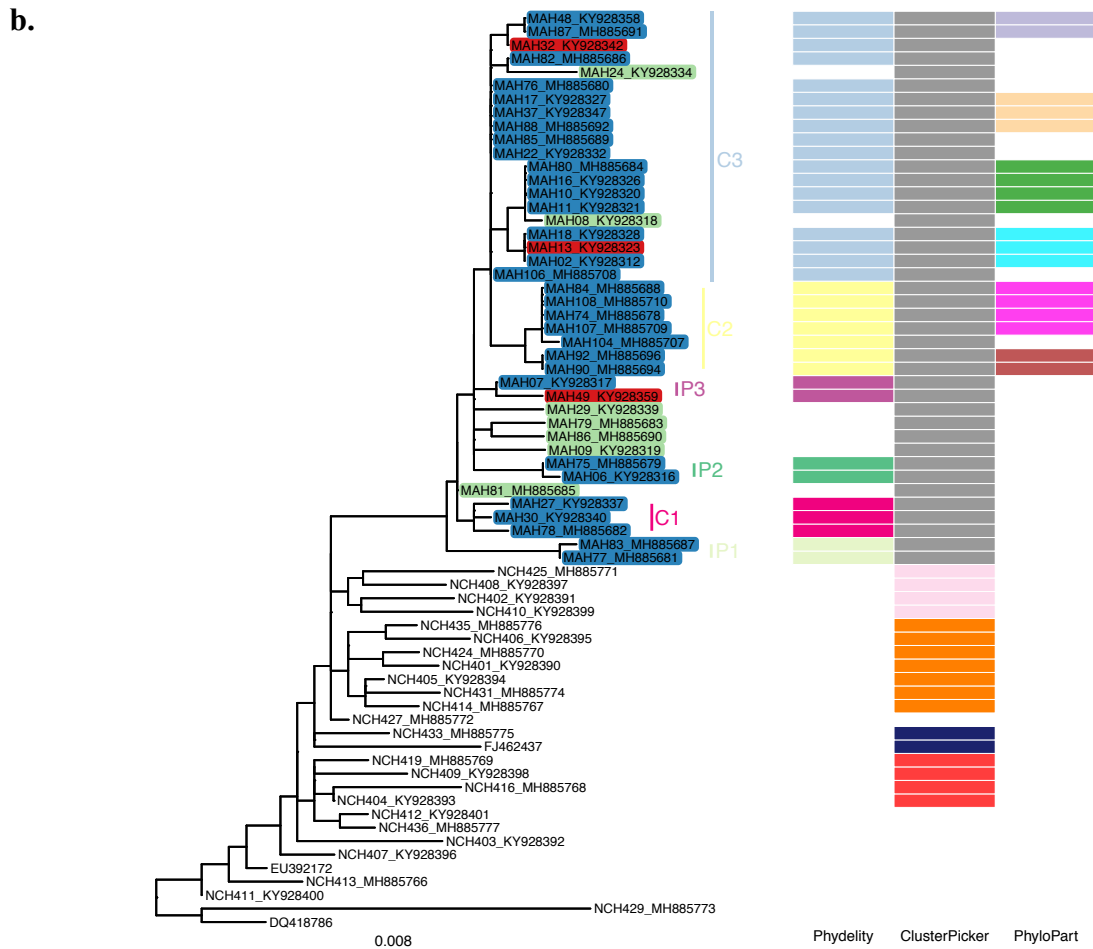
358

Genotype	<i>MPL</i>	Cluster	Mean pairwise patristic distance of cluster (σ)	Outlier	Mean pairwise patristic distance of outliers to cluster members (σ)
1a	0.029	C1	0.011 (0.012)	MAH66	0.043 (0.009)
		C3	0.016 (0.009)	MAH62	0.045 (0.027)
				MAH31	0.041 (0.025)
				MAH72	0.022 (0.015)
4d	0.010	C1	0.006 (0.004)	MAH24	0.019 (0.006)
				MAH08	0.009 (0.005)

359 **Table 1:** Comparing the genetic distance between outlying tips and the clusters they coalesce to with the
 360 genetic diversity of those clusters.

a.





361 **Figure 4.** Maximum likelihood phylogeny and clustering results of hepatitis C viruses (HCV) obtained from
 362 men who have sex with men (MSM) in Lyon, France. All highlighted tip names denoted in the format
 363 “MAH(ID)_accession” were samples from MSM patients (blue: HIV-positive, red: HIV-negative, green: HIV-
 364 positive and considered as outlying sequences by Phydelyity). Non-highlighted tips were collected from non-
 365 HIV, non-MSM patients residing in the same geographic region and time period. Clustering results from
 366 Phydelyity, ClusterPicker and PhyloPart are depicted as a heatmap. Each distinct colour refers to a different
 367 cluster. Similar to the Vietnamese hepatitis B empirical viral datasets (Figure 3a and Supplementary Figure 4),
 368 mean Silhouette index was used as the optimality criterion to determine the optimal absolute distance threshold
 369 for ClusterPicker and PhyloPart. Only results based on the optimised thresholds are shown here for
 370 ClusterPicker and PhyloPart. No parameter optimisation is required for Phydelyity. **(a)** Genotype 1a. **(b)**
 371 Genotype 4d.

373 *Computational performance*

374 For computational performance, Phydelyity can process a phylogeny of 1000 tips, on an
 375 Ubuntu 16.04 LTS operating system with an Intel Core i7-4790 3.60 GHz CPU, in ~3
 376 minutes using a single CPU core and 253 MB of peak memory usage.

378 **Discussion**

379 Phydelyity is a statistically-principled tool capable of identifying putative transmission clusters
 380 from pathogen phylogenies without the need to introduce arbitrary distance thresholds.

381 Instead, Phydelyity infers the maximal patristic distance limit (*MPL*) for cluster designation
382 using the pairwise patristic distance distribution of closely-related tips in the input
383 phylogenetic tree. Additionally, unlike other cutpoint-based methods, Phydelyity does not
384 assume clusters are strictly monophyletic and can identify paraphyletic clustering owing to its
385 distal dissociation approach. For datasets that span extended periods of time, multiple
386 introductions within the same contact network and concurrent onward transmissions to other
387 communities can result in “nested” introduction events that would go undetected by
388 monophyletic clustering (Barido-Sottani et al. 2018). By relaxing this assumption, not only
389 can Phydelyity pick up these “nested” events, it tends to produce clusters that are purer with a
390 lower chance of misclassification while excluding putative outlying tips that are exceedingly
391 distant from the inferred cluster.

392

393 There are algorithmic overlaps between Phydelyity and PhyCLIP, which is also a statistically-
394 principled phylogenetic clustering algorithm based on integer linear programming
395 optimisation (Han et al. 2019). However, the two clustering tools have substantially different
396 approaches that recover clusters with distinctly different interpretations. PhyCLIP was
397 developed to identify statistically-supported subpopulations in pathogen phylogenies that
398 putatively capture variant ecological, evolutionary or epidemiological processes that could
399 underlie sub-species nomenclature development. As such, PhyCLIP’s designated clusters
400 should not be interpreted as sequences linked by rapid transmission events. For instance,
401 when applied to the HCV genotype 1a NS5B dataset, PhyCLIP clustered 131 of the 155 input
402 sequences into seven clades, all of which encompasses genetically similar viruses of both
403 MSM and non-MSM origins that were endemic in Lyon during a specific period in time. In
404 contrast, Phydelyity assigned 73 sequences into 12 transmission pairs and 5 transmission
405 clusters that distinguished the underlying MSM transmission events from non-MSM ones
406 (Supplementary Figure 1). A detailed comparison between Phydelyity and PhyCLIP can be
407 found in Supplementary Materials.

408

409 There are two key assumptions underlying Phydelyity’s clustering algorithm. First, Phydelyity
410 expects transmissions to be linked by events rapid enough such that molecular evolution
411 between the transmitted pathogens is minimal, and thus genetically more similar amongst
412 themselves than to the given ensemble. At least for rapidly evolving pathogens such as RNA
413 viruses, genetic changes between sequences sampled from transmission pairs were found to
414 be generally low (Campbell et al. 2018).

415

416 Phydelity also assumes that the transmitted pathogens coalesce to the same most recent
417 common ancestor (MRCA) and that the pairwise genetic distance of internal nodes found
418 between the MRCA and the tips of the cluster to be bounded below *MPL*. Even though
419 Phydelity does not explicitly equate the inferred phylogeny to a transmission tree, imposing a
420 distance threshold between the internal nodes within a phylogenetic cluster may be construed
421 as an implicit assumption that the internal nodes are representative of transmission events.
422 There are central differences in the interpretation of phylogenetic and transmission trees
423 respectively. The former depicts the shared ancestry between the sampled tips while the latter
424 represents the true transmission history between the transmitted pathogens (Pybus and
425 Rambaut 2009; Ypma et al. 2013). It should be noted that Phydelity does not attribute any
426 interpretation of transmission events to the internal nodes and does not relate branch lengths
427 of the phylogenetic tree, which correlates with the timing of coalescence, to transmission
428 times. Restricting the distances between internal nodes below the *MPL* is strictly meant to
429 increase conservatism in identifying clusters that are as closely-related as possible. Any tips
430 clustered within the same cluster should be interpreted as a network of undirected
431 transmission pairs.

432

433 Furthermore, incorporating the aforementioned assumptions also means that Phydelity is not
434 exempt from well-documented pitfalls associated with other non-parametric, phylogenetic
435 clustering methods. Firstly, these clustering algorithms largely operate with conservatively
436 low thresholds. Resultantly, cluster identification is biased towards recent infections as
437 opposed to detecting differences in transmission rates between subpopulations. This bias
438 could be further worsened if oversampling occurs (Poon 2016; Dearlove et al. 2017; Le Vu et
439 al. 2018). While this caveat limits the interpretation of phylogenetic clusters, it does not
440 render phylogenetic clustering tools obsolete. As demonstrated by the HBV and HCV
441 empirical studies above, with meta-data associated with the individuals clustered,
442 phylogenetic clustering can still be used to identify infection trends as well as potential risk
443 factors and/or target subpopulations in retrospective studies.

444

445 There are a few more limitations that should be noted when using Phydelity. As the *MPL* is
446 wholly informed by the phylogenetic tree, clustering results will consequently be sensitive to
447 the diversity of closely-related tips within the input phylogeny. Specifically, the closely-
448 related sequences that constitute the k -th core patristic distance distribution (\mathcal{D}_k) must be

449 homogenous (i.e. similar difference between consecutive distances when \mathcal{D}_k is sorted; see
450 Methods) but sufficiently distinct from the background diversity of the phylogeny. Two
451 scenarios can arise if this is not the case: 1) few to no tips will be clustered if \mathcal{D}_k is not
452 homogenous. This may arise if sampling is so scattered such that few to no transmission pairs
453 are sampled; 2) potentially erroneous clustering of distantly-related tips may be obtained if
454 \mathcal{D}_k has a similar distance distribution relative to the entire tree. This could be possible if
455 sampling rate is too low relative to the mutation rate of the pathogen.

456

457 Additionally, constructing a phylogenetic tree can be a computational bottleneck for large
458 sequence datasets. As an alternative, genetic distance-based clustering algorithms such as
459 HIV-TRACE (Kosakovsky Pond et al. 2018) which negate the need to build a phylogenetic
460 tree have becoming increasingly popular. However, HIV-TRACE still requires users to
461 specify an arbitrary absolute distance threshold. Additionally, while it performed better than
462 other phylogenetic clustering method, HIV-TRACE did not preclude problems with bias
463 towards higher sampling rates (Poon 2016).

464

465 Despite the limitations discussed above, clustering results generated by Phydelyity for the
466 simulation and empirical datasets in this study demonstrate its superior performance over
467 current widely used phylogenetic clustering methods. Importantly, Phydelyity obviates the
468 need for users to define or optimise non-biologically-informed distance thresholds. Phydelyity
469 is fast, generalisable, and freely available at <https://github.com/alvinxhan/Phydelyity>.

470

471 **Acknowledgements**

472 We would like to thank Frits Scholer for his help in writing the program, as well as Jelle
473 Koopsen and Velislava Petrova for their key intellectual contributions.

474

475 **Data availability**

476 Phydelyity is freely available on <https://github.com/alvinxhan/Phydelyity>. All simulated
477 datasets were downloaded from Villandre et al. (2016). Genbank accession numbers of HBV
478 polymerase sequences: AB212625, GQ924626, AB115551, LC57377-LC57378, LC60789-
479 LC60790, LC63767, LC64366-LC64378, LC64380-LC64381, LC80779-LC80783,
480 LC80785, LC80787-LC80800, and LC80802-LC80804. Genbank accession numbers of
481 HCV NS5B sequences: AF9606, EF407457, HQ850279, EU392172, FJ462437, DQ418786,

482 M62321, MH885654-MH885777, and KY928311-KY928401. Jupyter notebooks used to
483 analyse both simulated and empirical datasets can be found in the same aforementioned
484 github repository.

485

486 **Funding**

487 A.X.H. was supported by the A*STAR Graduate Scholarship programme from A*STAR to
488 carry out his PhD work via collaboration between Bioinformatics Institute (A*STAR) and
489 NUS Graduate School for Integrative Sciences and Engineering from the National University
490 of Singapore. E.P. was funded by the Gates Cambridge Trust (Grant numer OPP1144).

491 S.M.S. was supported by the A*STAR HEIDI programme (Grant number: H1699f0013) and
492 Bioinformatics Institute (A*STAR).

493

494 **References**

495 Ambrosioni J, Junier T, Delhumeau C, Calmy A, Hirschel B, Zdobnov E, Kaiser L, Yerly S,
496 Study the SHIVC. 2012. Impact of highly active antiretroviral therapy on the molecular
497 epidemiology of newly diagnosed HIV infections. *AIDS* [Internet] 26. Available from:
498 https://journals.lww.com/aidsonline/Fulltext/2012/10230/Impact_of_highly_active_antiretroviral_therapy_on.10.aspx

500 Barido-Sottani J, Vaughan TG, Stadler T. 2018. Detection of HIV transmission clusters from
501 phylogenetic trees using a multi-state birth–death model. *J. R. Soc. Interface* [Internet] 15.
502 Available from: <http://rsif.royalsocietypublishing.org/content/15/146/20180512.abstract>

503 Bezemer D, Cori A, Ratmann O, van Sighem A, Hermanides HS, Dutilh BE, Gras L,
504 Rodrigues Faria N, van den Hengel R, Duits AJ, et al. 2015. Dispersion of the HIV-1
505 Epidemic in Men Who Have Sex with Men in the Netherlands: A Combined Mathematical
506 Model and Phylogenetic Analysis. *PLOS Med.* [Internet] 12:e1001898. Available from:
507 <https://doi.org/10.1371/journal.pmed.1001898>

508 Breiman L, Friedman J, Stone CJ, Olshen RA. 1984. *Classification and regression trees*.
509 Florida.

510 Brenner BG, Roger M, Routy J-P, Moisi D, Ntemgwa M, Matte C, Baril J-G, Thomas R,
511 Rouleau D, Bruneau J, et al. 2007. High Rates of Forward Transmission Events after
512 Acute/Early HIV-1 Infection. *J. Infect. Dis.* [Internet] 195:951–959. Available from:
513 <http://dx.doi.org/10.1086/512088>

514 Campbell F, Strang C, Ferguson N, Cori A, Jombart T. 2018. When are pathogen genome
515 sequences informative of transmission events? Parrish C, editor. *PLOS Pathog.* [Internet]
516 14:e1006885. Available from: <https://dx.plos.org/10.1371/journal.ppat.1006885>

517 Charre C, Cotte L, Kramer R, Mialhes P, Godinot M, Koffi J, Scholtès C, Ramière C. 2018.
518 Hepatitis C virus spread from HIV-positive to HIV-negative men who have sex with
519 men. Shoukry NH, editor. *PLoS One* [Internet] 13:e0190340. Available from:

- 520 <https://dx.plos.org/10.1371/journal.pone.0190340>
- 521 Dearlove BL, Xiang F, Frost SDW. 2017. Biased phylodynamic inferences from analysing
522 clusters of viral sequences. *Virus Evol.* [Internet] 3:vex020. Available from: +
- 523 Gardy JL, Loman NJ. 2017. Towards a genomics-informed, real-time, global pathogen
524 surveillance system. *Nat. Rev. Genet.* 19:9–20.
- 525 Grabowski MK, Redd AD. 2014. Molecular tools for studying HIV transmission in sexual
526 networks. *Curr. Opin. HIV AIDS* [Internet] 9. Available from: [https://journals.lww.com/co-](https://journals.lww.com/co-hivandaids/Fulltext/2014/03000/Molecular_tools_for_studying_HIV_transmission_in.6.aspx)
527 [hivandaids/Fulltext/2014/03000/Molecular_tools_for_studying_HIV_transmission_in.6.aspx](https://journals.lww.com/co-hivandaids/Fulltext/2014/03000/Molecular_tools_for_studying_HIV_transmission_in.6.aspx)
- 528 Han AX, Parker E, Scholer F, Maurer-Stroh S, Russell CA. 2019. Phylogenetic Clustering by
529 Linear Integer Programming (PhyCLIP). Shapiro B, editor. *Mol. Biol. Evol.* [Internet].
530 Available from: [https://academic.oup.com/mbe/advance-](https://academic.oup.com/mbe/advance-article/doi/10.1093/molbev/msz053/5373046)
531 [article/doi/10.1093/molbev/msz053/5373046](https://academic.oup.com/mbe/advance-article/doi/10.1093/molbev/msz053/5373046)
- 532 Hubert L, Arabie P. 1985. Comparing partitions. *J. Classif.* [Internet] 2:193–218. Available
533 from: <http://link.springer.com/10.1007/BF01908075>
- 534 Kosakovsky Pond SL, Weaver S, Leigh Brown AJ, Wertheim JO. 2018. HIV-TRACE
535 (TRANsmiSSion Cluster Engine): A tool for large scale molecular epidemiology of HIV-1 and
536 other rapidly evolving pathogens. Shapiro B, editor. *Mol. Biol. Evol.* 35:1812–1819.
- 537 Manning CD, Raghavan P, Schütze H. 2008. *Introduction to Information Retrieval*. New
538 York, NY, USA: Cambridge University Press
- 539 Matsuo J, Do SH, Yamamoto C, Nagashima S, Chuon C, Katayama K, Takahashi K, Tanaka
540 J. 2017. Clustering infection of hepatitis B virus genotype B4 among residents in Vietnam,
541 and its genomic characters both intra- and extra-family. Chemin I, editor. *PLoS One* [Internet]
542 12:e0177248. Available from: <https://dx.plos.org/10.1371/journal.pone.0177248>
- 543 de Oliveira T, Kharsany ABM, Gräf T, Cawood C, Khanyile D, Grobler A, Puren A, Madurai
544 S, Baxter C, Karim QA, et al. 2017. Transmission networks and risk of HIV infection in
545 KwaZulu-Natal, South Africa: a community-wide phylogenetic study. *Lancet HIV* [Internet]
546 4:e41–e50. Available from: [https://doi.org/10.1016/S2352-3018\(16\)30186-2](https://doi.org/10.1016/S2352-3018(16)30186-2)
- 547 Poon AFY. 2016. Impacts and shortcomings of genetic clustering methods for infectious
548 disease outbreaks. *Virus Evol.* 2:vew031.
- 549 Prospero MCFF, Ciccozzi M, Fanti I, Saladini F, Pecorari M, Borghi V, Di Giambenedetto S,
550 Bruzzone B, Capetti A, Vivarelli A, et al. 2011. A novel methodology for large-scale
551 phylogeny partition. *Nat. Commun.* [Internet] 2:321. Available from:
552 <https://doi.org/10.1038/ncomms1325>
- 553 Pybus OG, Rambaut A. 2009. Evolutionary analysis of the dynamics of viral infectious
554 disease. *Nat. Rev. Genet.* [Internet] 10:540–550. Available from:
555 <http://www.nature.com/articles/nrg2583>
- 556 Ragonnet-Cronin M, Hodcroft E, Hué S, Fearnhill E, Delpech V, Brown AJLAAL, Lycett
557 SS, Holmes E, Nee S, Rambaut A, et al. 2013. Automated analysis of phylogenetic clusters.
558 *BMC Bioinformatics* [Internet] 14:317. Available from: <https://doi.org/10.1186/1471-2105->

559 14-317

560 Rousseeuw PJ. 1987. Silhouettes: A graphical aid to the interpretation and validation of
561 cluster analysis. *J. Comput. Appl. Math.* 20:53–65.

562 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
563 large phylogenies. *Bioinformatics* [Internet] 30:1312–1313. Available from:
564 <http://dx.doi.org/10.1093/bioinformatics/btu033>

565 Villandre L, Stephens DA, Labbe A, Günthard HF, Kouyos R, Stadler T, Study TSHIVC.
566 2016. Assessment of Overlap of Phylogenetic Transmission Clusters and Communities in
567 Simple Sexual Contact Networks: Applications to HIV-1. *PLoS One* [Internet] 11:e0148459.
568 Available from: <https://doi.org/10.1371/journal.pone.0148459>

569 Volk JE, Marcus JL, Phengrasamy T, Hare CB. 2015. Incident Hepatitis C Virus Infections
570 Among Users of HIV Preexposure Prophylaxis in a Clinical Practice Setting. *Clin. Infect.*
571 *Dis.* [Internet] 60:1728–1729. Available from: [https://academic.oup.com/cid/article-](https://academic.oup.com/cid/article-lookup/doi/10.1093/cid/civ129)
572 [lookup/doi/10.1093/cid/civ129](https://academic.oup.com/cid/article-lookup/doi/10.1093/cid/civ129)

573 Le Vu S, Ratmann O, Delpech V, Brown AE, Gill ON, Tostevin A, Fraser C, Volz EM.
574 2018. Comparison of cluster-based and source-attribution methods for estimating
575 transmission risk using large HIV sequence databases. *Epidemics* [Internet] 23:1–10.
576 Available from:
577 <https://www.sciencedirect.com/science/article/pii/S1755436517301159#bib0030>

578 Ypma RJF, van Ballegooijen WM, Wallinga J. 2013. Relating phylogenetic trees to
579 transmission trees of infectious disease outbreaks. *Genetics* [Internet] 195:1055–1062.
580 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24037268>

581