

Inferring putative transmission clusters with Phydely

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

¹Protein Sequence Analysis Group, Bioinformatics Institute, Agency for Science, Technology and Research (A*STAR), 30 Biopolis Street, 138671 Singapore, ²NUS Graduate School for Integrative Sciences and Engineering, National University of Singapore (NUS), 21 Lower Kent Ridge, 119077 Singapore, ³Laboratory of Applied Evolutionary Biology, Department of Medical Microbiology, Academic Medical Centre, Meibergdreef 9, 1105 AZ Amsterdam-Zuidoost, The Netherlands, ⁴Department of Veterinary Medicine, University of Cambridge, Madingley Rd, Cambridge CB3 0ES, UK and ⁵Department of Biological Sciences, National University of Singapore, 16 Science Drive 4, 117558 Singapore

[†]<http://orcid.org/0000-0001-6281-8498>

*Corresponding authors: E-mail: hanxc@bii.a-star.edu.sg; c.a.russell@amsterdamumc.nl

Abstract

Current phylogenetic clustering approaches for identifying pathogen transmission clusters are limited by their dependency on arbitrarily defined genetic distance thresholds for within-cluster divergence. Incomplete knowledge of a pathogen's underlying dynamics often reduces the choice of distance threshold to an exploratory, ad hoc exercise that is difficult to standardise across studies. Phydely is a new tool for the identification of transmission clusters in pathogen phylogenies. It identifies groups of sequences that are more closely related than the ensemble distribution of the phylogeny under a statistically principled and phylogeny-informed framework, without the introduction of arbitrary distance thresholds. Relative to other distance threshold- and model-based methods, Phydely outputs clusters with higher purity and lower probability of misclassification in simulated phylogenies. Applying Phydely to empirical datasets of hepatitis B and C virus infections showed that Phydely identified clusters with better correspondence to individuals that are more likely to be linked by transmission events relative to other widely used non-parametric phylogenetic clustering methods without the need for parameter calibration. Phydely is generalisable to any pathogen and can be used to identify putative direct transmission events. Phydely is freely available at <https://github.com/alvinxhan/Phydely>.

Key words: phylogenetic clustering; molecular epidemiology; transmission

1. Introduction

Recent advances in high-throughput sequencing technologies have led to the widespread use of sequence data in infectious disease epidemiology (Gardy and Loman 2018). In particular, epidemiologically relevant information such as the structure of transmission networks and infection source identification are increasingly inferred from virus phylogenies, especially for measurably evolving viral pathogens like HIV-1 and hepatitis C viruses

(Ambrosioni et al. 2012; Bezemer et al. 2015; de Oliveira et al. 2017; Matsuo et al. 2017; Charre et al. 2018). Non-parametric phylogenetic-based clustering tools operate on the assumption that pathogens in a transmission cluster are linked by transmission events rapid enough that molecular evolution between the transmitted pathogens is minimal, and thus genetically more similar amongst themselves than to the ensemble of input isolates (Prosperi et al. 2011; Ragonnet-Cronin et al. 2013). This assumption is generally valid for rapidly evolving pathogens such as RNA

viruses as genetic changes between sequences sampled from transmission pairs are generally low (Campbell et al. 2018).

Non-parametric phylogenetic clustering methods typically measure the genetic divergence of sequence pairs either by their genetic distances that are computed from the sequence data directly (Aldous et al. 2012; Ragonnet-Cronin et al. 2013) or by their patristic distances from the inferred phylogenetic tree (i.e. the sum of the inferred phylogenetic branch lengths linking the two sequences; Brenner et al. 2007; Prospero et al. 2011). The divergence of a cluster can be defined as the median (Prospero et al. 2011) or largest (Ragonnet-Cronin et al. 2013) pairwise distance between member sequences of the cluster. To define transmission clusters, an upper divergence threshold is implemented either as an absolute distance limit (Ragonnet-Cronin et al. 2013) or as a percentile of the distribution of pairwise sequence distances (Prospero et al. 2011). A fundamental limitation of these non-parametric phylogenetic clustering tools is the need to define this arbitrary absolute transmission cluster divergence thresholds (termed as ‘cutpoints’ by Villandre et al. 2016). The lack of a consensus definition of a phylogenetic transmission cluster (Grabowski and Redd 2014) coupled with incomplete knowledge of a pathogen’s underlying epidemiological dynamics often reduces the choice of cutpoints to an ad hoc exploratory exercise resulting in subjective cluster definitions.

Phydelity is a novel phylogenetic clustering tool designed to negate the need for arbitrarily defined cluster divergence thresholds. Requiring only the phylogenetic tree as input, Phydelity infers putative transmission clusters through the identification of groups of sequences that are more closely related to one another than the ensemble distribution under a statistically principled framework. Phydelity, like another phylogenetic clustering tool that we recently developed, PhyCLIP, is based on integer linear programming (ILP) optimisation (Han et al. 2019). However, the two clustering tools are substantially different in their approaches and ILP models such that their clustering results have entirely distinct interpretations. PhyCLIP uses the divergence information of the entire phylogenetic tree to inclusively assign statistically supported cluster membership to as many sequences in the tree as possible that putatively capture variant ecological, evolutionary or epidemiological processes. To this end, PhyCLIP is useful for sub-species nomenclature development. Phydelity, on the other hand, exclusively distinguishes closely related pathogens with pairwise sequence divergence that are significantly more likely to be drawn from the same low divergence distribution than that of the ensemble. As such, while PhyCLIP’s designated clusters are underpowered to be interpreted as sequences linked by transmission events, clusters inferred by Phydelity can be interpreted as putative transmission clusters (see Supplementary Data).

To demonstrate the utility of Phydelity in identifying putative transmission clusters, the algorithm underlying Phydelity is first presented in detail. The clustering tool is then applied to both simulated and empirical datasets, including outbreaks of hepatitis B and C viruses as well as seasonal A/H3N2 influenza virus infections, and compared against results generated by existing phylogenetic clustering methods. Phydelity is freely available at <http://github.com/alvinxhan/Phydelity>.

2. Method

2.1 Clustering algorithm

Figure 1a shows the overall workflow of Phydelity. First, Phydelity considers the input phylogeny as an ensemble of

putative clusters, each consisting of an internal node i and the leaves it subtends. The within-cluster diversity of node i is measured by its mean pairwise patristic distance (μ_i). The patristic distance between two nodes, which can be any sequence tips or internal nodes in the phylogeny, refers to the sum of branch lengths linking those two nodes. Sequences subtended by i (i.e. all descendant tree tips of node i) are considered for clustering if μ_i is less than the maximal patristic distance limit (MPL), under which sequences are considered more closely related to one another than the ensemble distribution (Fig. 1b).

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

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

The user can opt to either input the desired k parameter or allow Phydelity to automatically scale k to the value that yields the supremum k th core distance distribution with the lowest overall divergence (i.e. the largest possible k that still yields the lowest overall divergence between k -neighbouring tips). This is done by testing if \mathcal{D}_{k+1} and \mathcal{D}_k are statistically distinct ($P < 0.01$) using Kuiper’s test (see Supplementary Data). All clustering results of Phydelity presented in this work were generated using the autoscaled value of k .

The MPL is then calculated by

$$\text{MPL} = \bar{\mu} + \sigma$$

where $\bar{\mu}$ is the median pairwise distance of \mathcal{D}_k and σ is the corresponding robust estimator of scale without assuming symmetry about $\bar{\mu}$ using the Qn method (see Supplementary Data; Rousseeuw and Croux 1993; Fig. 1b).

This is then followed by dissociation of distantly related descendant subtrees/sequences to all putative nodes for clustering, thereby facilitating identification of both monophyletic as well as nested paraphyletic clusters (Fig. 1c; see Supplementary Data). Phydelity filters outlying tips from putative clusters under the assumption that viruses infecting individuals in a transmission chain coalesce to the same most recent common ancestor (MRCA). Additionally, Phydelity requires any clonal ancestors in between the MRCA and tips of a putative cluster to be as genetically similar to each other as they are to the tips of the cluster. As such, for a putative transmission cluster, the mean pairwise nodal distance between all internal and tip nodes of a cluster must also be $\leq \text{MPL}$ (Fig. 1c).

An ILP model is implemented and optimised under the objective to assign cluster membership to sequences satisfying the aforementioned relatedness criteria within the least number of clusters. In other words, Phydelity uses ILP optimisation to search for the clustering configuration that favours the designation of larger clusters of closely related sequences which are likely linked by transmission events. Any topologically outlying singletons that were spuriously clustered are removed. Finally, it is important to note that a transmission cluster identified by Phydelity should only be interpreted as a fully connected network of likely transmission pairs without implying any underlying transmission directionality. The full algorithm description

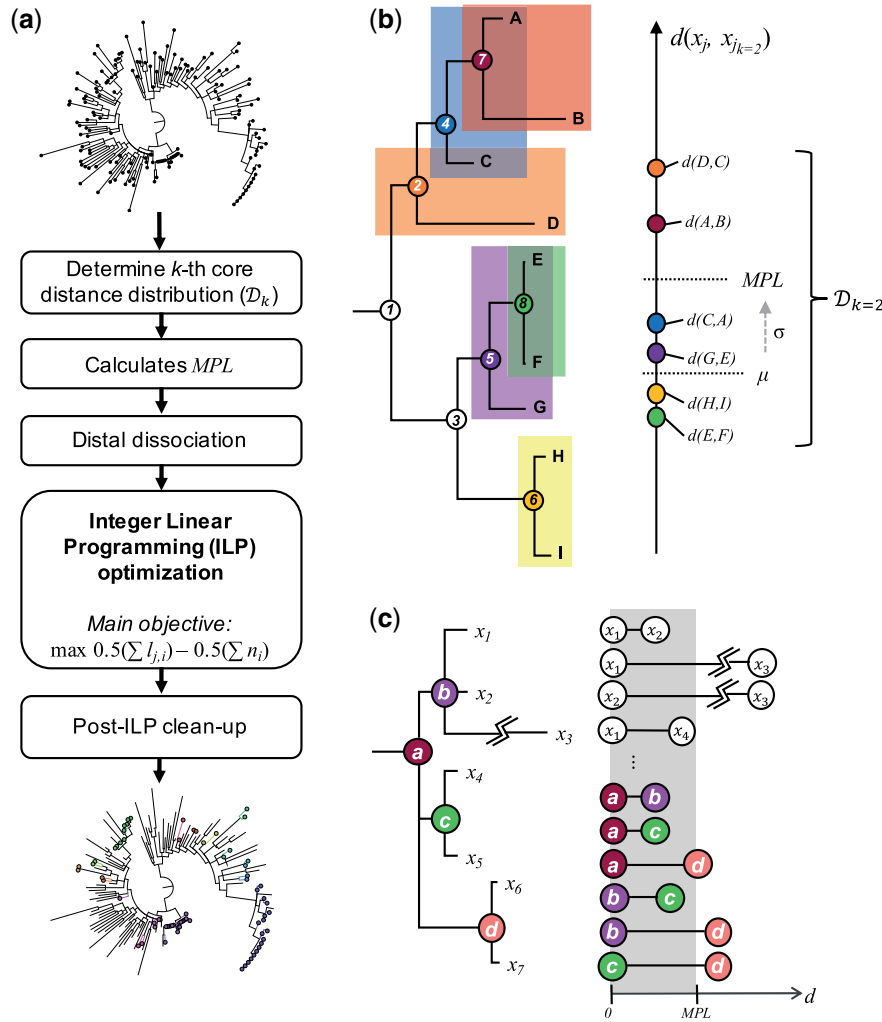


Figure 1. (a) Phydely algorithm pipeline. Phydely considers the input phylogenetic tree as a collection of putative clusters each defined by an internal node i and tips j that it subtends. The algorithm first infers the k th core distance distribution (\mathcal{D}_k) from the pairwise patristic distances of the closest k -neighbouring tips. k can be defined by the user or scaled by Phydely to obtain the supremum \mathcal{D}_k with the lowest divergence. \mathcal{D}_k is then used to compute the maximal patristic distance limit (MPL) under which tips are considered to be more closely related than to the ensemble. Dissociation of distally related subtrees/sequences (c) ensues such that both monophyletic and paraphyletic clustering structures can be identified. Phydely then incorporates the distance and topological information of the remaining nodes and tips into an integer linear programming (ILP) model to be optimised by clustering all tips that satisfy the relatedness constraints within the least number of clusters. Finally, post-ILP steps are implemented to remove any tips that may have been spuriously clustered. (b) Determination of the maximal patristic distance limit (MPL) using the median (μ) and robust estimator of scale (σ) based on the k th core distance distribution (\mathcal{D}_k) of every sequence x_j and its k -closest neighbours ($d(x_j, x_{j+k})$; $k=2$ in this case as shown by the pairs of sequences highlighted with distinct colours). (c) Distal dissociation of a putative transmission cluster subtended by internal node a . If a sequence tip has a pairwise sequence distance that is greater than MPL , it will be dissociated and not be clustered under the internal node of interest (i.e. internal node a). In this case, sequence x_3 is dissociated from the putative cluster a due to its exceedingly long branch length violating the MPL threshold (i.e. $d(x_3, x_{3_2}) > MPL$). Additionally, whole subtrees subtended by the internal node of interest will be dissociated if any of its inter-nodal patristic distance exceeds MPL . Here, subtree d and its descending sequences (i.e. x_6 and x_7) will be dissociated from a as its inter-nodal distances with internal nodes b and c are both larger than MPL .

and mathematical formulation of Phydely are detailed in [Supplementary Data](#).

2.2 Assessing clustering results of simulated epidemics

Phydely was evaluated on phylogenetic trees derived from simulated HIV epidemics of a hypothetical men who have sex with men (MSM) sexual contact network (C-type networks in [Villandre et al. 2016](#)). The simulated sexual contact network comprised of 100 subnetworks (communities) sampled from an empirical distribution obtained from the Swiss HIV Cohort Study. All communities were linked in a chain initially and additional connections between any two communities were generated at a probability of 0.00075. Subjects in the network could either be in the ‘susceptible’,

‘infected’ or ‘removed’ (i.e. individual was diagnosed and sampled) state. Transmission clusters were attributed to sexual contact among individuals belonging to the same community.

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

These simulated datasets were tested by Villandre et al. (2016) to compare the outputs of four ‘cutpoint-based’ phylogenetic clustering methods where the arbitrary distance threshold defining a transmission cluster (i.e. cutpoint) was computed as the: (1) absolute patristic distance threshold between any two tips (Brenner et al. 2007); (2) standardised number of nucleotide changes (i.e. ClusterPicker; Ragonnet-Cronin et al. 2013); (3) percentile of the phylogeny’s pairwise sequence patristic distance distribution (i.e. PhyloPart; Prosperi et al. 2011) and (4) height of an ultrametric tree obtained using the weighted pair-group method of analysis (WPGMA). For each method, Villandre et al. varied the corresponding cutpoint parameter over an equivalent range of thresholds. Comparing the output clusters generated by the four methods at their respective optimal cutpoint by adjusted rand index (ARI) (see below), it was found that the WPGMA method tended to produce clusters with better correspondence to the underlying sexual contact structure. As such, clustering results from Phydely were compared with those obtained by Villandre et al. using the WPGMA method. Additionally, Phydely was also compared with the multi-state birth–death (MSBD) method which inferred transmission clusters on the same simulated datasets by detecting significant changes in transmission rates (Barido-Sottani, Vaughan, and Stadler 2018).

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

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

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

where C^* is the set of clusters defined to have correctly classified sequences attributed to community t (i.e. any cluster that constitutes the largest proportion of sequences from community t at both the cluster and the community label levels), p_t is the probability of sequence from community t and $p(c|t)$ refers to the probability that a sequence is clustered under cluster c conditional of it being from community t . If output clusters perfectly align with the underlying sexual contact network (i.e. one cluster only constitute one class of community), $I_G = 0$. Conversely, if clustering results are completely random, $I_G = 1$.

- iii. Purity measures the average extent that the output clusters contain only a single class (i.e. a particular sexual contact community; Manning, Raghavan, and Schütze 2008):

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

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

- iv. Normalised mutual information (NMI) trades off the output clustering quality against the number of clusters (Manning, Raghavan, and Schütze 2008):

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

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

2.3 Empirical datasets

Phydely was also tested on three empirical datasets—acute hepatitis C virus infections among MSM (Charre et al. 2018), hepatitis B viruses (HBVs) collected from members of the same families (Matsuo et al. 2017) as well as A/H3N2 influenza viruses collected from a community-based cohort of households during the 2014/2015 season (McCrone et al. 2018). All phylogenetic trees were reconstructed using RAxML (v8.2.12) under the GTRGAMMA model (Stamatakis 2014).

2.4 Comparisons to ClusterPicker and PhyloPart

ClusterPicker (Ragonnet-Cronin et al. 2013) and PhyloPart (Prosperi et al. 2011), two non-parametric phylogenetic clustering tools that are methodologically comparable to Phydely, were also applied to the hepatitis C and B virus datasets for comparisons. Either clustering tool has been previously applied to multiple studies involving different pathogens (Prosperi et al. 2011; Jacka et al. 2014; Bezemer et al. 2015; Bartlett et al. 2016; Coll et al. 2017; de Oliveira et al. 2017; Charre et al. 2018). Other than the phylogenetic tree, both ClusterPicker and PhyloPart also require users to input an arbitrarily defined genetic distance threshold (as an absolute distance limit for ClusterPicker and percentile of the global pairwise patristic distance for PhyloPart). As such, a range of distance limits (PhyloPart: 0.5–10th percentile; ClusterPicker: 0.005–0.1 nucleotide/site) were applied to both tools. No bootstrap support threshold was implemented for comparability to Phydely.

The lowest optimal threshold for the distance range tested was found by maximisation of the mean Silhouette index (SI) for both ClusterPicker and PhyloPart. The SI measures how similar an item is to members of its own cluster as opposed to the nearest neighbouring clusters—i.e. a larger mean SI indicates

that items of the same cluster are more closely related amongst themselves than to its neighbours (Rousseeuw 1987). No parameter optimisation was required for Phydelyty.

3. Results

3.1 Simulated HIV epidemics

Phydelyty was applied to simulated HIV epidemics among MSM belonging to a hypothetical sexual contact network structures where transmission clusters were attributed to transmission by sexual contact among individuals belonging to the same sub-network (see Section 2; Villandre et al. 2016). These simulations were originally used to assess the performance of ‘cutpoint-based’ clustering tools, including ClusterPicker, PhyloPart as well as the WPGMA that generally attained the highest ARI score across all simulations when calibrating their respective cutpoint thresholds against the ground-truth. Phylogenetic trees generated from these simulations were also tested by the MSBD method (Barido-Sottani, Vaughan, and Stadler 2018).

Clustering results from Phydelyty were compared with outputs from the MSBD method and those from the WPGMA method achieving the best ARI scores. The purity, modified Gini index (I_C) and NMI measures were also used to provide a more comprehensive assessment of the clustering results (Fig. 2; Supplementary Fig. S3 and Table S1; see Section 2).

The phylogenetic trees generated from the simulations had a large number of clusters that were relatively small in size (i.e. percentage of sequences that were part of ground truth clusters with sizes < 8 tips = 33.9% (weight of inter-community transmission rates, $w = 25\%$); 55.5% ($w = 100\%$); see Barido-Sottani, Vaughan, and Stadler (2018) for more details). Furthermore, these ground truth clusters were not all monophyletic (Fig. 2c). As a result, while Phydelyty and WPGMA yielded comparable ARI scores (Phydelyty: 0.44–0.45 (SD = 0.05); WPGMA: 0.44–0.56 (SD = 0.05–0.05); Supplementary Table S1), Phydelyty’s output clusters, which allows paraphyletic clusters (Fig. 2c), are substantially purer (mean purity; Phydelyty: 0.81–0.88 (SD = 0.03); WPGMA: 0.67–0.74 (SD = 0.06–0.06)) and have a lower probability of misclassification when compared with WPGMA which assumes clusters are strictly monophyletic (mean I_C ; Phydelyty: 0.27–0.28 (SD = 0.04–0.05); WPGMA: 0.33–0.40 (SD = 0.04–0.05)). Coverage of sequences clustered by Phydelyty lies between 58.2 per cent and 61.6 per cent.

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

3.2 HBV transmission between family members

Phydelyty was tested on empirical datasets to demonstrate its applicability on real-world data, including HBVs collected from residents in the Binh Thuan Province of Vietnam (Matsuo et al. 2017). In such highly endemic regions, HBV is commonly transmitted either vertically from mothers to children during the perinatal period or horizontally between cohabitants of the same household (Matsuo et al. 2017). As complete genome nucleotide sequences were not available for all individuals, a phylogenetic tree was reconstructed using the viral polymerase sequences collected from forty-one patients, of which twelve of

them were confirmed to be members of three families (i.e. denoted as F2, F3 and F4) by a family survey as well as mitochondrial analyses. Besides Phydelyty, the resulting phylogeny was also implemented in ClusterPicker and PhyloPart.

Phydelyty identified three likely transmission clusters that distinguish between the separate family households (Fig. 3). At their respective optimal distance thresholds by mean SI (see Section 2), ClusterPicker and PhyloPart achieved similar clustering results. Importantly, Phydelyty was able to obtain the same optimal clustering results without optimisation and implementation of a hard-to-interpret distance parameter.

3.3 Hepatitis C virus transmission among MSM

Incidence of HCV infections among HIV-negative MSM has been relatively limited as compared with their HIV-positive counterparts. However, the recent uptake of pre-exposure prophylaxis (PrEP) among HIV-negative individuals to prevent HIV infection could pose higher risk of sexually transmitted HCV infections (Volk et al. 2015; Charre et al. 2018). In a study on HIV-positive and HIV-negative MSM patients in Lyon, 108 cases of acute HCV infections (80 primary infections; 28 reinfections) were reported between 2014 and 2017 among 96 MSM (72 HIV-positive; 24 HIV-negative, of which 16 (67%) of them were on PrEP; Charre et al. 2018). Separate phylogenetic analyses were performed on a subset of 89 (68 HIV-positive; 21 HIV-negative) HCV isolates belonging to Genotypes 1a and 4d based on their NS5B sequences. Additionally, 25 HCV sequences from HIV-infected MSM collected before 2014 were included along with 60 control HCV sequences derived from HIV-negative, non-MSM patients residing in the same geographical area as controls. All sequences collected from MSM patients were given strain names in the format of ‘MAH(ID)_accession’ while control sequences from non-HIV, non-MSM patients were denoted as ‘NCH(ID)_accession’ (Fig. 4). Phydelyty as well as ClusterPicker and PhyloPart were applied to the reconstructed phylogenies, with the latter calibrated over a range of distance thresholds. Again, only clustering results based on the lowest distance threshold maximising the mean SI for ClusterPicker and PhyloPart were compared with Phydelyty’s output clusters (see Section 2).

Generally, membership of the MSM transmission clusters and pairs identified by Phydelyty across both genotypes were strictly limited to sequences derived from MSM patients. Relaxing the monophyletic assumption by dissociating distantly related tips from putative monophyletic clusters (see Section 2) enables Phydelyty to identify likely outlying sequences as evidenced by their relatively longer branch lengths from the cluster ensemble (Table 1; Fig. 4; Genotype 1a: Cluster C1—MAH66 and Cluster C3—MAH31, MAH62 and MAH72; Genotype 4d: Cluster C3—MAH24 and MAH08). In particular, for Genotype 1a, even though the mean pairwise distance of MAH72 to members of Cluster C3 is within a standard deviation of the latter’s within-cluster diversity, its distance to the more distant members (e.g. MAH15 and MAH40; Fig. 4) violated the inferred MPL (Table 1). Additionally, as a result of distal dissociation, Phydelyty distinguishes clusters that are genetically more alike amongst themselves than to those phylogenetically ancestral to it (e.g. Cluster C1.1 that is ‘nested’ within Cluster C1 for Genotype 1a; Fig. 4a).

For both genotypes, Phydelyty found multiple clusters that included both HIV-positive and HIV-negative MSM patients (i.e. Genotype 1a: Clusters C2 and C3, Fig. 4a; Genotype 4d: Clusters C2 and C2.2, as well as pair P2, Fig. 4b). While it is not clear which of the HIV-negative patients were on PrEP (information

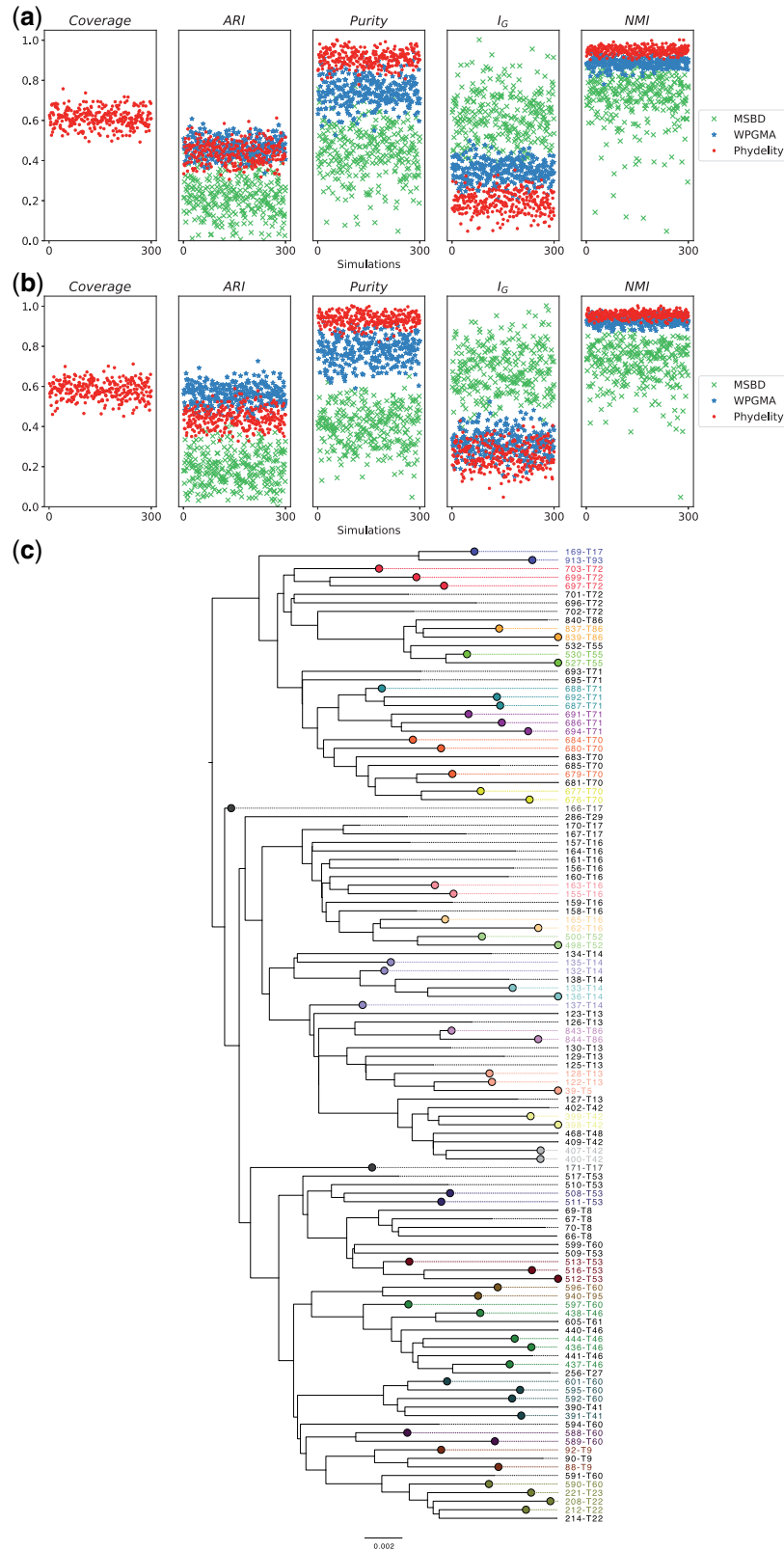


Figure 2. Clustering results of simulated HIV epidemics in a hypothetical MSM sexual contact network. (a) Clustering metrics for clustering algorithms (Phylometry, weighted pair-group method of analysis (WPGMA) and multi-state birth–death (MSBD) methods) applied simulated phylogenies with inter-communities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

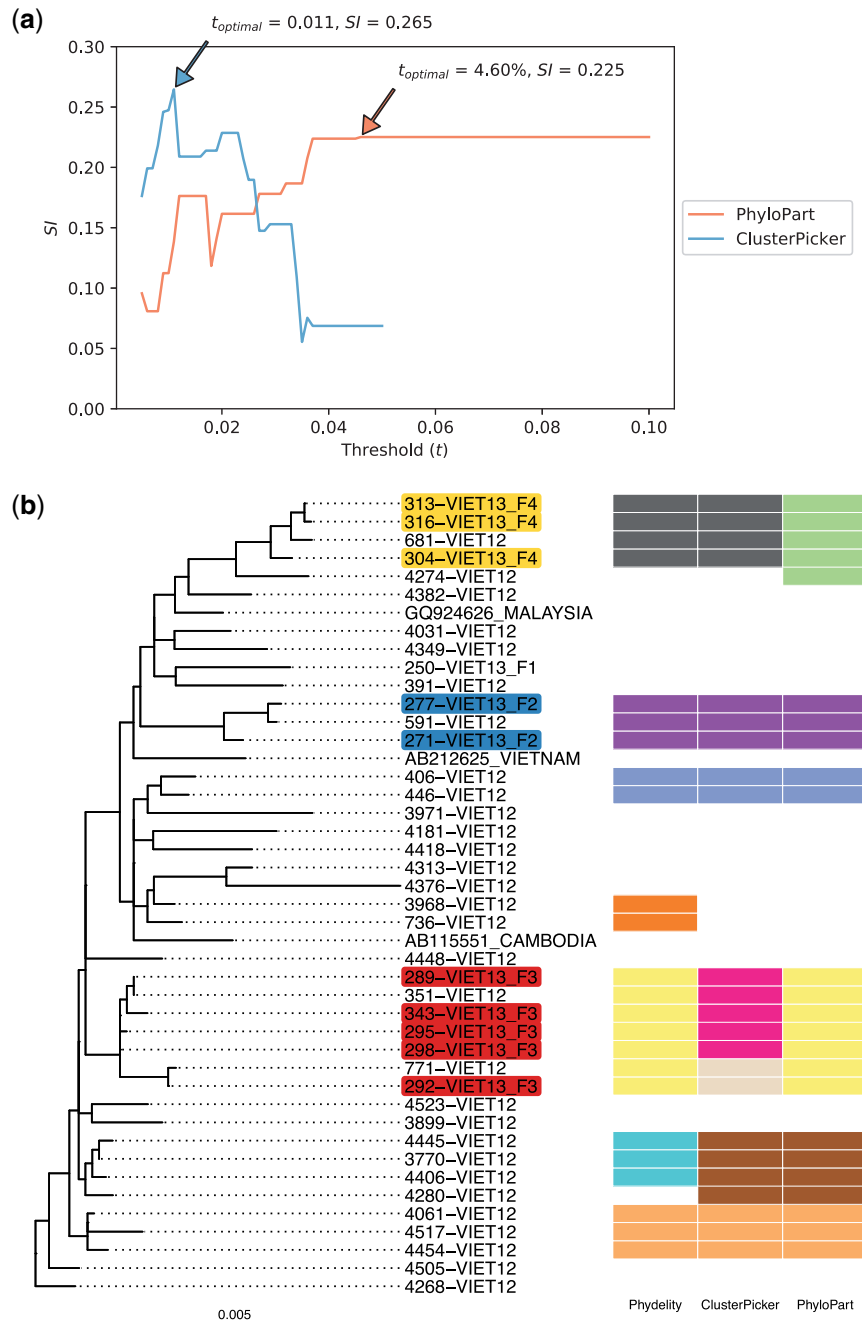


Figure 3. Clustering results of hepatitis B viruses (HBV) collected from residents in the Binh Thuan Province of Vietnam. (a) Plots of mean Silhouette index (SI) computed for the range of genetic distance thresholds implemented in ClusterPicker and PhyloPart. Clustering results from the lowest optimal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

not supplied in the original paper), the clustering results from Phydelyity were in line with the findings by Charre *et al.* that acute HCV infections among HIV-negative MSM were likely sourced from their HIV-positive counterparts.

While ClusterPicker managed to consolidate all of the MSM Genotype 4d sequences into a single monophyletic cluster (Fig. 4b), its clustering of Genotype 1a was problematic as a large

number of non-MSM control sequences were clustered together with those from MSM patients (Fig. 4a). PhyloPart's optimal clustering output was consistent with Phydelyity's for Genotype 1a. However, the larger number of identical sequences in the Genotype 4d tree skewed the optimal distance parameter (expressed as xth percentile of the pairwise patristic distribution of the entire phylogeny) to only cluster these identical sequences.

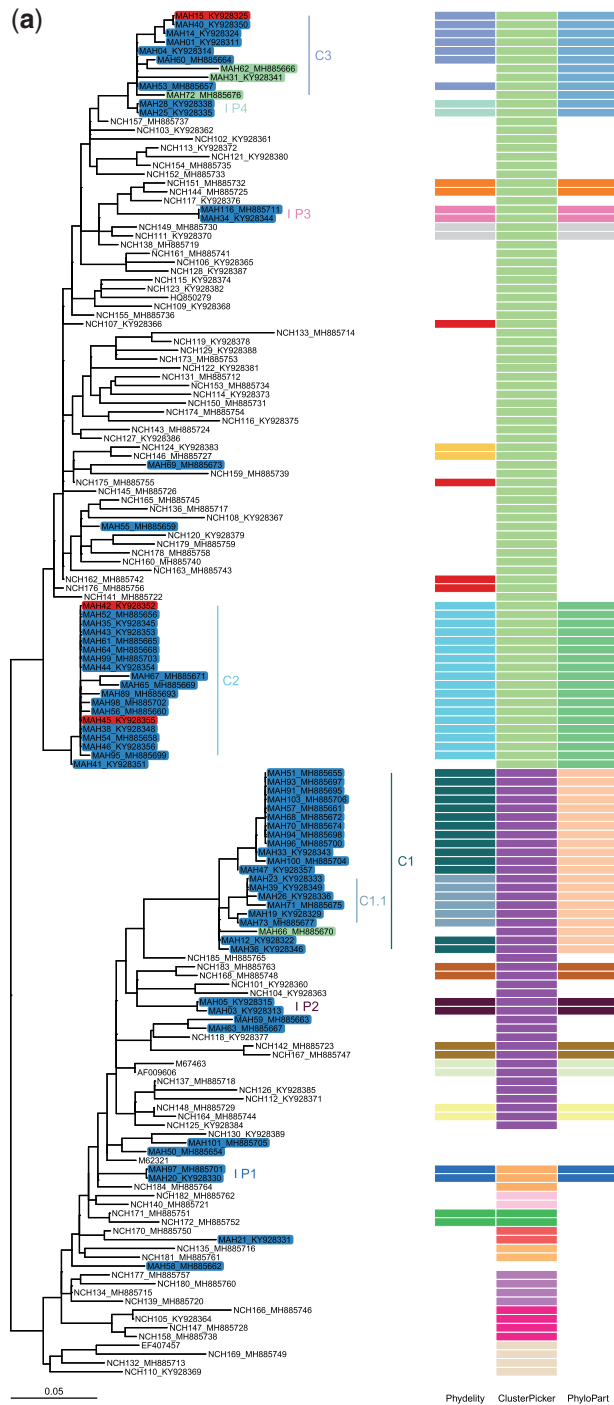


Figure 4. Maximum likelihood phylogeny and clustering results of hepatitis C viruses (HCV) obtained from men who have sex with men (MSM) in Lyon, France. All highlighted tip names denoted in the format ‘MAH(ID)_accession’ were samples from MSM patients (blue: HIV-positive, red: HIV-negative, green: HIV-positive and considered as outlying sequences by Phydelity). Non-highlighted tips were collected from non-HIV, non-MSM patients residing in the same geographic region and time period. Clustering results from Phydelity, ClusterPicker and PhyloPart are depicted as a heatmap. Each distinct colour refers to a different cluster. Similar to the Vietnamese hepatitis B empirical viral datasets (Fig. 3a; Supplementary Fig. 5a), mean Silhouette index was used as the optimality criterion to determine the optimal absolute distance threshold for ClusterPicker and PhyloPart. Only results based on the optimised thresholds are shown here for ClusterPicker and PhyloPart. No parameter optimisation is required for Phydelity. (a) Genotype 1a and (b) Genotype 4d.

3.4 Seasonal A/H3N2 influenza virus infections within a community and the effects of sampling

Phydelity was also applied to A/H3N2 influenza viruses collected from a community-based cohort of 340 households (1431 participants) in Southeastern Michigan, USA during the 2014/

2015 season (McCrone et al. 2018). Of the influenza-positive cases, 206 virus samples were collected from 166 individuals that belonged to 81 households and sequenced. As concurrent infections among individuals within the same household do not necessarily imply transmission, McCrone et al. implemented

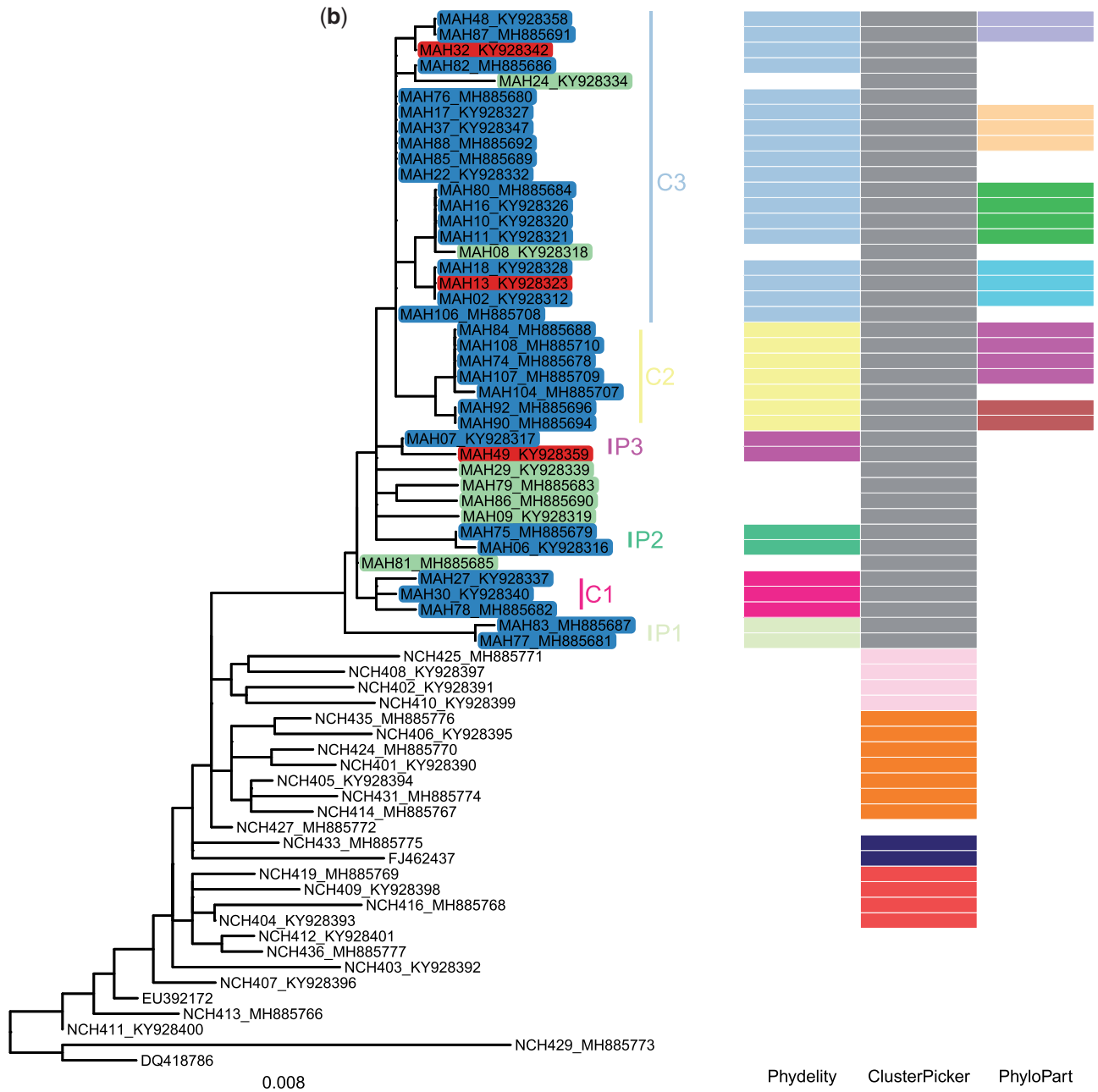


Figure 4. continued

stringent epidemiological as well as genetic distance constraints to identify transmission pairs: (1) the donor and recipient of a transmission pair were of the same household with onset of illness symptoms occurring within 7 days of each other, with the donor having the earlier symptom onset date; (2) there must be no other potential donors with the same symptom onset date; (3) symptom onset dates of donor and recipient should not be on the same day unless they were index cases; and (4) genetic distance between the within-host viral populations of donor and recipient must be below the 5th percentile of the distance distribution of random pairs of infected individuals from the community (McCrone et al. 2018). In total, fifty virus isolates constituting thirty-two high-quality transmission pairs

were identified. Consolidating transmission pairs with overlapping donors and recipients into clusters, there were twenty-two genetically validated transmission clusters, comprising of sixteen pairs and six trios in total.

Using the phylogeny constructed from the consensus whole genome sequences of all 206 viruses, Phydelity was able to identify 20 of the 22 high-quality transmission clusters as distinct clusters (Supplementary Fig. S5). Applying the same metrics used to assess clustering performance of the simulated dataset earlier and using the high-quality transmission cluster labels as ground truth, Phydelity was able to produce highly pure clusters (97.8%), with a low probability of misclassification ($I_C = 0.022$) and good accuracy (ARI = 0.962), even after accounting for the

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

| Genotype | MPL | 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) |

Table 2. Clustering performance of Phydely on seasonal A/H3N2 influenza viruses collected by McCrone et al. (2018).

| Basis | n_{sample} | % _{trans} | Purity | I_G | ARI | NMI | |
|------------------------------------|---------------------|--------------------|-------------|-------------|-------------|------|------|
| High-quality transmission clusters | All | 0.98 | 0.02 | 0.96 | 0.99 | | |
| | | 52 | 25% | 0.87 | 0.06 | 0.72 | 0.93 |
| | | | 45% | 0.87 | 0.04 | 0.74 | 0.95 |
| | 93 | 25% | 0.85 | 0.07 | 0.76 | 0.94 | |
| | | | 45% | 0.94 | 0.03 | 0.88 | 0.98 |
| | | | 45% | 0.94 | 0.03 | 0.90 | 0.98 |
| Household | All | 0.89 | 0.08 | 0.79 | 0.96 | | |
| | | 52 | 25% | 0.56 | 0.29 | 0.35 | 0.82 |
| | | | 45% | 0.73 | 0.16 | 0.56 | 0.90 |
| | 93 | 25% | 0.82 | 0.11 | 0.74 | 0.93 | |
| | | | 45% | 0.75 | 0.16 | 0.64 | 0.92 |
| | | | 45% | 0.87 | 0.11 | 0.80 | 0.95 |

Ground truth used for clustering assessment was either based on the identities of genetically validated, high-quality transmission clusters as defined by McCrone et al. or by the patients' households. Besides analysing all of the viruses collected (bolded results), Phydely was also applied to downsampled datasets consisting of different sample size (n_{sample}) and proportion of sequences derived from the aforementioned high-quality transmission pairs (%_{trans}). Adjusted rand index (ARI) measures how accurate the output clusters corresponded with the ground truth labels. Purity gives the average extent clusters contain only a single class. Modified Gini index (I_G) is the probability that a randomly selected sequence would be incorrectly clustered. Normalised mutual information (NMI) accounts for the trade-off between clustering quality and number of clusters (see Section 2).

number of predicted clusters (NMI = 0.993; Table 2). As transmission events defined by McCrone et al. were based on highly conservative criteria imposed on deep sequencing datasets, Phydely, which operates at the consensus sequence level, could also cluster viruses that did not satisfy these constraints but were still linked epidemiologically by their household identities. As such, Phydely's clustering results were assessed based on the household association of the clustered individuals as well, yielding slightly diminished but nonetheless high quality performance (Purity = 0.894, I_G = 0.081, ARI = 0.791, NMI = 0.964; Supplementary Fig. S5; Table 2).

The full A/H3N2 sequence dataset was then randomly sampled to smaller pools of fifty-two (25%) as well as ninety-three (45%) isolates to assess how low sampling rates might affect Phydely's performance. To ensure that sequences involved in high-quality transmission pairs were also sampled, such isolates would constitute different proportions (either 25% or 45%; as well as 70% for pools of fifty-two sequences only) of the downsampled datasets. Ten distinct downsamples were generated for each sample size/high-quality transmission sequence combination and the average results were tabulated (Table 2).

As the MPL is informed by the phylogenetic tree, clustering results will consequently be sensitive to the diversity of closely

related tips within the input phylogeny. Specifically, the closely related sequences that constitute the k th core patristic distance distribution (\mathcal{D}_k) must be homogenous (i.e. similar difference between consecutive distances when \mathcal{D}_k is sorted; see Section 2) but sufficiently distinct from the background diversity of the phylogeny. This was demonstrated by the improved clustering results with respect to household identities with greater proportional inclusion of genetically similar, high-quality transmission pairs in the downsampled dataset (Table 2). Furthermore, erroneous clustering of distantly related tips can be obtained if \mathcal{D}_k has a similar distance distribution relative to the entire tree due to insufficient divergence information from reduced sampling. This is evident from the general decrease in the clustering performance of all downsampled data. In particular, clustering closely related, high-quality transmission clusters was worse off with a lower sample size.

3.5 Computational performance

For computational performance, Phydely can process a phylogeny of 1,000 tips, on an Ubuntu 16.04 LTS operating system with an Intel Core i7-4790 3.60 GHz CPU, in ~3 min using a single CPU core and 253 MB of peak memory usage.

4. Discussion

Phydely is a statistically principled tool capable of identifying putative transmission clusters from pathogen phylogenies without the need to introduce arbitrary distance thresholds. Instead, Phydely infers the maximal patristic distance limit (MPL) for cluster designation using the pairwise patristic distance distribution of closely related tips in the input phylogenetic tree. Unlike other cutpoint-based methods, Phydely does not assume clusters are strictly monophyletic and can identify paraphyletic clustering owing to its distal dissociation approach. For datasets that span extended periods of time, multiple introductions within the same contact network and concurrent onward transmissions to other communities can result in 'nested' introduction events that would go undetected by monophyletic clustering (Barido-Sottani, Vaughan, and Stadler 2018). By relaxing this assumption, not only can Phydely pick up these 'nested' events, it tends to produce clusters that are purer with a lower chance of misclassification while excluding putative outlying tips that are exceedingly distant from the inferred cluster.

Even though there are algorithmic similarities between PhyCLIP (Han et al. 2019) and Phydely, clustering results generated by PhyCLIP should not be interpreted as sequences linked by transmission events. For instance, when applied to the HCV Genotype 1a NS5B dataset, PhyCLIP clustered 131 of the 155 input sequences into seven clades, all of which encompasses

genetically similar viruses of both MSM and non-MSM origins that were endemic in Lyon during a specific period in time. In contrast, Phydelyty assigned seventy-three sequences into twelve transmission pairs and five transmission clusters that distinguished the underlying MSM transmission events from non-MSM ones (Supplementary Fig. S1). A detailed comparison between Phydelyty and PhyCLIP can be found in Supplementary Data.

One of the key assumptions made by Phydelyty is that the transmitted pathogens coalesce to the same MRCA and that the pairwise genetic distance of internal nodes found between the MRCA and the tips of the cluster to be bounded below MPL. While Phydelyty does not explicitly equate the inferred phylogeny to a transmission tree, imposing a distance threshold between the internal nodes within a phylogenetic cluster may be construed as an implicit assumption that the internal nodes are representative of transmission events. There are important differences in the interpretation of phylogenetic and transmission trees. The former depicts the shared ancestry between the sampled tips while the latter represents the true transmission history between the transmitted pathogens (Pybus and Rambaut 2009; Ypma, van Ballegooijen, and Wallinga 2013). It should be noted that Phydelyty neither attributes any interpretation of transmission events to the internal nodes nor does it relate branch lengths of the phylogenetic tree, which correlate with the timing of coalescence, to transmission times. Restricting the distances between internal nodes below the MPL is strictly meant to increase conservatism in identifying clusters that are as closely related as possible.

There have also been criticisms that non-parametric cluster identification by genetic similarity is biased towards the detection of recent infections as opposed to discerning variations in transmission rates between different subpopulations, which can be further exacerbated by oversampling (Poon 2016; Dearlove, Xiang, and Frost 2017; Le Vu et al. 2018). While this caveat limits the interpretation of the inferred transmission clusters, it does not render all phylogenetic clustering tools obsolete. Phylogenetic clustering tools supplemented by epidemiological meta-data can still be used to systematically identify infection trends, potential risk factors and target subpopulations, as demonstrated by multiple epidemiological studies of different measurably evolving pathogens (de Oliveira et al. 2017; Matsuo et al. 2017; Charre et al. 2018).

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

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

Acknowledgements

We would like to thank Frits Scholer for his assistance with optimising Phydelyty, as well as Jelle Koopsen and Velislava Petrova for their intellectual contributions.

Funding

A.X.H. was supported by the A*STAR Graduate Scholarship programme from A*STAR to carry out his PhD work via collaboration between Bioinformatics Institute (A*STAR) and NUS Graduate School for Integrative Sciences and Engineering from the National University of Singapore. E.P. was funded by the Gates Cambridge Trust (Grant number: OPP1144). S.M.S. was supported by the A*STAR HEIDI programme (Grant number: H1699f0013) and Bioinformatics Institute (A*STAR).

Data availability

Phydelyty is freely available on <https://github.com/alvinxhan/Phydelyty>. All simulated datasets were downloaded from Villandre et al. (2016). Genbank accession numbers of HBV polymerase sequences: AB212625, GQ924626, AB115551, LC57377-LC57378, LC60789-LC60790, LC63767, LC64366-LC64378, LC64380-LC64381, LC80779-LC80783, LC80785, LC80787-LC80800 and LC80802-LC80804. Genbank accession numbers of HCV NS5B sequences: AF9606, EF407457, HQ850279, EU392172, FJ462437, DQ418786, M62321, MH885654-MH885777 and KY928311-KY928401. The A/H3N2 influenza virus consensus sequences were downloaded from https://github.com/lauringlab/Host_level_IAV_evolution. Jupyter notebooks used to analyse both simulated and empirical datasets can be downloaded from <https://github.com/alvinxhan/Phydelyty/tree/master/manuscript>.

Supplementary data

Supplementary data are available at *Virus Evolution* online.

Conflict of interest: None declared.

References

- Aldous, J. L. et al. (2012) 'Characterizing HIV Transmission Networks across the United States', *Clinical Infectious Diseases: An Official Publication of the Infectious Diseases Society of America*, 55: 1135–43.
- Ambrosioni, J. et al. (2012) 'Impact of Highly Active Antiretroviral Therapy on the Molecular Epidemiology of Newly Diagnosed HIV Infections', *AIDS*, 26: 2079–86.
- Barido-Sottani, J., Vaughan, T. G., and Stadler, T. (2018) 'Detection of HIV Transmission Clusters from Phylogenetic Trees Using a Multi-State Birth–Death Model', *Journal of the Royal Society Interface*, 15: 20180512.
- Bartlett, S. R. et al. (2016) 'HIV Infection and Hepatitis C Virus Genotype 1a Are Associated with Phylogenetic Clustering among People with Recently Acquired Hepatitis C Virus Infection', *Infection, Genetics and Evolution*, 37: 252–8.
- Bezemer, D. et al. (2015) 'Dispersion of the HIV-1 Epidemic in Men Who Have Sex with Men in The Netherlands: A Combined Mathematical Model and Phylogenetic Analysis', *PLoS Medicine*, 12: e1001898.
- Breiman, L. et al. (1984) *Classification and Regression Trees*. Monterey, CA: Wadsworth & Brooks.
- Brenner, B. G. et al. (2007) 'High Rates of Forward Transmission Events after Acute/Early HIV-1 Infection', *The Journal of Infectious Diseases*, 195: 951–9.

- Campbell, F. et al. (2018) 'When Are Pathogen Genome Sequences Informative of Transmission Events?', *PLoS Pathogens*, 14: e1006885.
- Charre, C. et al. (2018) 'Hepatitis C Virus Spread from HIV-Positive to HIV-Negative Men Who Have Sex with Men', *PLoS One*, 13: e0190340.
- Coll, F. et al. (2017) 'Longitudinal Genomic Surveillance of MRSA in the UK Reveals Transmission Patterns in Hospitals and the Community', *Science Translational Medicine*, 9: eaak9745.
- Dearlove, B. L., Xiang, F., and Frost, S. (2017) 'Biased Phylodynamic Inferences from Analysing Clusters of Viral Sequences', *Virus Evolution*, 3: vex020.
- Gardy, J. L., and Loman, N. J. (2018) 'Towards a Genomics-Informed, Real-Time, Global Pathogen Surveillance System', *Nature Reviews Genetics*, 19: 9–20.
- Grabowski, M. K., and Redd, A. D. (2014) 'Molecular Tools for Studying HIV Transmission in Sexual Networks', *Current Opinion in HIV and AIDS*, 9: 126–33.
- Han, A. X. et al. (2019) 'Phylogenetic Clustering by Linear Integer Programming (PhyCLIP)', *Molecular Biology and Evolution*, 36: 1580–95.
- Hubert, L., and Arabie, P. (1985) 'Comparing Partitions', *Journal of Classification*, 2: 193–218.
- Jacka, B. et al. (2014) 'Phylogenetic Clustering of Hepatitis C Virus among People Who Inject Drugs in Vancouver, Canada', *Hepatology (Baltimore, Md.)*, 60: 1571–80.
- Kosakovsky Pond, S. L. et al. (2018) 'HIV-TRACE (TRANSMISSION Cluster Engine): a Tool for Large Scale Molecular Epidemiology of HIV-1 and Other Rapidly Evolving Pathogens', *Molecular Biology and Evolution*, 35: 1812–9.
- Manning, C. D., Raghavan, P., and Schütze, H. (2008) *Introduction to Information Retrieval*. New York, NY: Cambridge University Press.
- Matsuo, J. et al. (2017) 'Clustering Infection of Hepatitis B Virus Genotype B4 among Residents in Vietnam, and Its Genomic Characters Both Intra- and Extra-Family', *PLoS One*, 12: e0177248.
- McCrone, J. T. et al. (2018) 'Stochastic Processes Constrain the within and between Host Evolution of Influenza Virus', *eLife*, 7: e35962.
- de Oliveira, T. et al. (2017) 'Transmission Networks and Risk of HIV Infection in KwaZulu-Natal, South Africa: A Community-Wide Phylogenetic Study', *The Lancet HIV*, 4: e41–e50.
- Poon, A. (2016) 'Impacts and Shortcomings of Genetic Clustering Methods for Infectious Disease Outbreaks', *Virus Evolution*, 2: vew031.
- Prosperi, M. et al. (2011) 'A Novel Methodology for Large-Scale Phylogeny Partition', *Nature Communications*, 2: 321.
- Pybus, O. G., and Rambaut, A. (2009) 'Evolutionary Analysis of the Dynamics of Viral Infectious Disease', *Nature Reviews Genetics*, 10: 540–50.
- Ragonnet-Cronin, M. et al. (2013) 'Automated Analysis of Phylogenetic Clusters', *BMC Bioinformatics*, 14: 317.
- Rousseeuw, P. J. (1987) 'Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis', *Journal of Computational and Applied Mathematics*, 20: 53–65.
- , and Croux, C. (1993) 'Alternatives to the Median Absolute Deviation', *Journal of the American Statistical Association*, 88: 1273–83.
- Stamatakis, A. (2014) 'RAxML Version 8: A Tool for Phylogenetic Analysis and Post-Analysis of Large Phylogenies', *Bioinformatics (Oxford, England)*, 30: 1312–3.
- Villandre, L. et al. (2016) 'Assessment of Overlap of Phylogenetic Transmission Clusters and Communities in Simple Sexual Contact Networks: Applications to HIV-1', *PLoS One*, 11: e0148459.
- Volk, J. E. et al. (2015) 'Incident Hepatitis C Virus Infections among Users of HIV Preexposure Prophylaxis in a Clinical Practice Setting', *Clinical Infectious Diseases*, 60: 1728–9.
- Le Vu, S. et al. (2018) 'Comparison of Cluster-Based and Source-Attribution Methods for Estimating Transmission Risk Using Large HIV Sequence Databases', *Epidemics*, 23: 1–10.
- Ypma, R. J. F., van Ballegooijen, W. M., and Wallinga, J. (2013) 'Relating Phylogenetic Trees to Transmission Trees of Infectious Disease Outbreaks', *Genetics*, 195: 1055–62.