

1 **Inferring putative transmission clusters with Phydely**

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. Phydely 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. In simulated phylogenies, Phydely achieves higher rates of correspondence to
27 ground-truth clusters than current model-based methods, and comparable results to
28 parametric methods without the need for parameter calibration.

29

30 **Availability and implementation:** Phydely is available at

31 <http://github.com/alvinxhan/Phydely>.

32

33

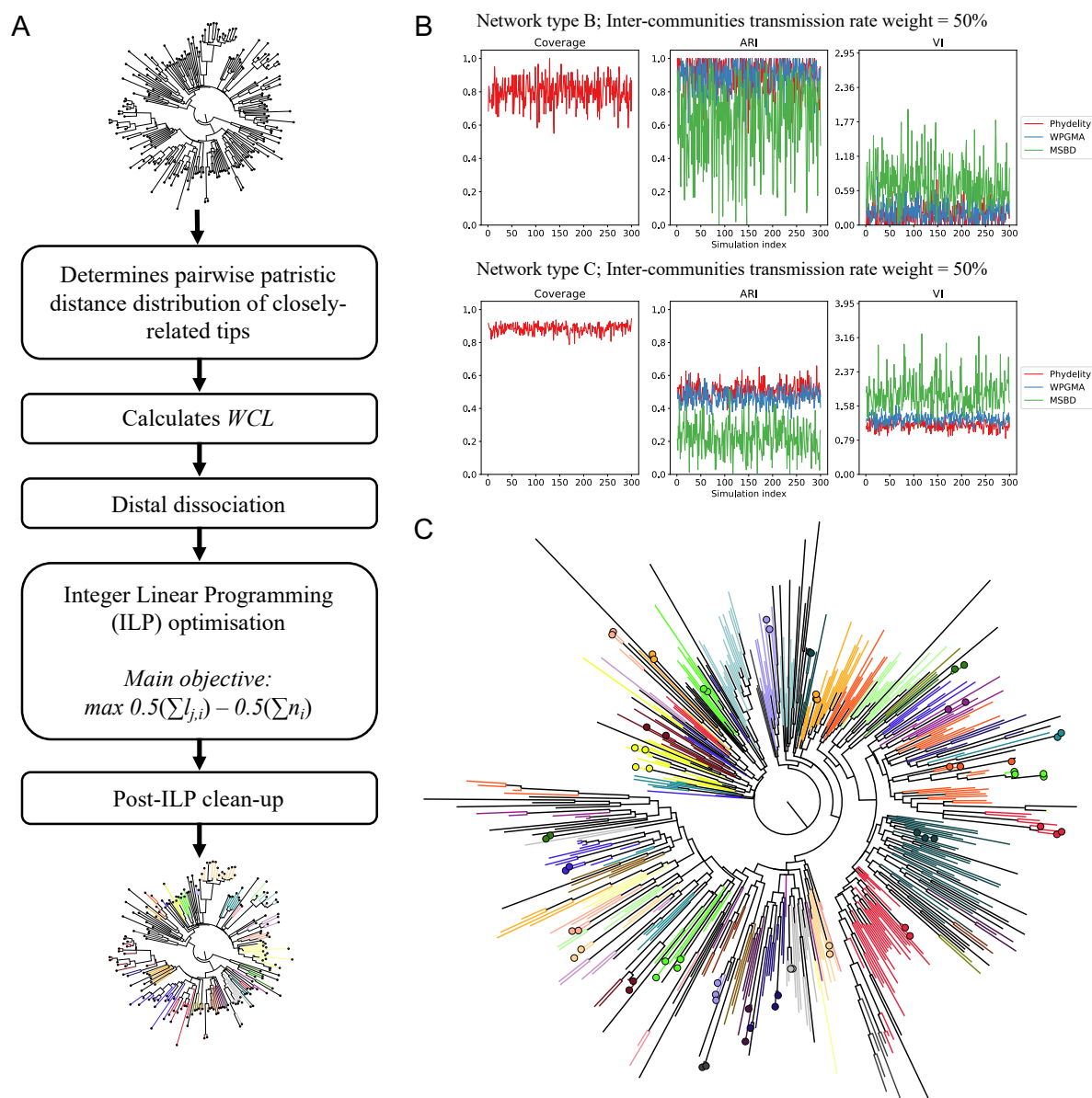
34 **Introduction**

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

46

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

53



54

55 **Figure 1.** (A) Phydelyty algorithm pipeline. (B). Clustering correspondence metrics for clustering algorithms
 56 (Phydelyty, WPGMA and MSBD) applied to phylogenies generated from simulations of HIV epidemics of
 57 hypothetical MSM sexual contact network types B and C, with inter-communities transmission rates weighted at
 58 50% of within-community rate (Villandre *et al.*, 2016). (C) Clustering results of Phydelyty on HIV-1 subtype A
 59 *env* sequences collected from the Rakai Community Cohort Study (Grabowski *et al.*, 2014). Tips are coloured
 60 by Phydelyty clusters and those marked with ● were phylogenetic clusters identified by Grabowski *et al.*

61

62 **Method**

63 Phydelyty considers the input phylogeny as an ensemble of putative clusters, each consisting
 64 of an internal node i and the leaves it subtends. The within-cluster diversity of node i is
 65 measured by its mean pairwise patristic distance (μ_i). Phydelyty then determines the pairwise
 66 patristic distance distribution of closely-related tips, which comprises the pairwise distances
 67 of sequence j and its closest k -neighbouring tips, where the closest k -neighbours includes
 68 sequence j . The user can input the desired k parameter or Phydelyty can automatically scale k

69 to the value that yields the supremum distribution with the lowest overall divergence. All
70 tests of Phydelity presented in this work were performed using the autoscaled value of k .

71

72 Regardless of how the distance distribution of closely-related tips is determined, Phydelity
73 uses this distribution to calculate the within-cluster divergence limit (WCL), an upper bound
74 to μ_i of putative clusters:

$$75 \quad WCL = \bar{\mu} + \sigma$$

76 where $\bar{\mu}$ is the median pairwise distance of the closely-related tips distance distribution and σ
77 is the corresponding robust estimator of scale without assuming symmetry about $\bar{\mu}$.

78

79 This is followed by distal dissociation of distantly-related descendant subtrees/sequences to
80 any ancestral node with $\mu_i > WCL$, thereby facilitating identification of both monophyletic
81 as well as nested, paraphyletic clusters (Han *et al.*, 2018). Phydelity filters outlying tips from
82 putative clusters under the assumption that viruses infecting individuals in a quick
83 transmission chain are ultimately descended from the same source and are highly similar
84 genetically. An outliers is defined by a node-to-tip distance more than three deviations from
85 the median distance. An integer linear programming model is implemented and optimised
86 under a blended objective of equal weights to maximise the number of sequences clustered
87 within the lowest number of clusters. Lastly, clean-up steps are taken to remove any
88 topologically outlying singletons that were spuriously clustered as described in Han *et al.*
89 (2018). The full algorithm description and mathematical formulation of Phydelity is detailed
90 in Supplementary Materials.

91

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

95

96 **Results**

97 Phydelity was evaluated on phylogenetic trees derived from simulated HIV epidemics of two
98 hypothetical MSM sexual contact network types (network types B & C) produced by
99 Villandre *et al.* (2016), wherein quick transmission chains (i.e. transmission clusters) could
100 be attributed to sexual contact among individuals belonging to the same community. Network
101 type B, which corresponded best with the assumption of monophyletic clusters, consisted of a

102 main contact network of 60 individuals with single linkages to 25 disjoint subnetworks of 20
103 subjects. Conversely, the more realistic network type C included 100 communities of sizes
104 sampled from an empirical distribution obtained from the Swiss HIV Cohort Study. For both
105 network types, inter-community transmission rates were weighted at 25%, 50%, 75% or
106 100% of the within-community rate. These simulated datasets were also tested by Barido-
107 Sottani *et al* using their multi-state birth-death (MSBD) method which infers transmission
108 clusters by detecting significant changes in transmission rates (Barido-Sottani *et al.*, 2018).
109 More information on the simulated epidemics are included in Supplementary Materials.

110

111 Clustering results from Phydelyty were compared to those generated by the MSBD method
112 and a cutpoint method based on the weighted pair-group method of analysis (WPGMA)
113 (Villandre *et al.*, 2016) (Figure 1B, Supplementary Fig. 1 and Supplementary Table 1). Both
114 the adjusted rand index (ARI) and variation of information (VI) were calculated to quantify
115 the correspondence between the actual network communities and clustering results. Villandre
116 *et al.* (2016) assessed four different commonly used cutpoint-based methods, including
117 arbitrarily varying patristic distance thresholds between any two tips (Brenner *et al.*, 2007),
118 ClusterPicker (varying standardised number of nucleotide changes; Ragonnet-Cronin *et al.*,
119 2013), PhyloPart (changing arbitrary percentile of pairwise patristic distance distribution;
120 Proserpi *et al.*, 2011) and agglomerative hierarchical clustering methods such as the
121 WPGMA method. WPGMA methods achieved the best overall ARI. As such, only WPGMA
122 clustering results derived from the optimal cutpoint parameter (i.e. maximum ARI) were
123 compared.

124

125 Owing to Phydelyty's definition of a closely-related neighbourhood, its distal dissociation
126 approach and outlier detection, its mean coverage of sequences clustered ranges from 70.8–
127 80.0% for B networks and 87.5–87.9% for C networks. Phydelyty (mean ARI = 0.90-0.91,
128 mean VI = 0.17-0.18) consistently performed as well as optimised WPGMA (mean ARI =
129 0.88-0.96, mean VI = 0.09-0.25) for B networks. Phydelyty was the best performing method
130 for C networks (Phydelyty: mean ARI = 0.49-0.59, mean VI = 0.95-1.18; WPGMA: mean
131 ARI = 0.44-0.56, mean VI = 1.07-1.33; MSBD: mean ARI = 0.19-0.22, mean VI = 1.80-2.02;
132 Supplementary Table 1). Notably, even though results generated by WPGMA are comparable
133 to those from Phydelyty, this was only possible for WPGMA when the optimal cutpoint could
134 be determined by calibration with the simulated ground-truth. However, ground truth

135 clustering is largely unavailable in epidemiological studies. Phydelity, on the other hand,
136 does not require this calibration step.

137

138 Phydelity was also tested on an empirical dataset of HIV-1 subtype A *env* sequences obtained
139 from the Rakai Community Cohort Study (Grabowski *et al.*, 2014). Grabowski *et al.*
140 identified 35 clusters by phylogenetic analysis, of which 18 constituted individuals of the
141 same community and 12 clusters were confirmed to be from the same household. As
142 community and household information was blinded for privacy reasons, the available
143 sequence data could not be matched for the exact clusters identified by Grabowski *et al.*
144 However, the stringent cluster definitions used by Grabowski *et al.* ($\geq 90\%$ bootstrap support
145 (1000 iterations), median genetic pairwise distance $\leq 2.6\%$) restricts most of these clusters to
146 pairs, with non-pair clusters made up of no more than 5 individuals. As such, we found the
147 same number of clusters of similar size distribution through visual inspection when we
148 recapitulated the phylogeny using the same methods (GTR+I+G substitution model, Garli)
149 and cluster definition as Grabowski *et al.* Phydelity identified 33 out of the 35
150 epidemiologically-identified clusters as distinct transmission clusters (Figure 1C).

151

152 **Conclusion**

153 Phydelity is a statistically-principled and phylogeny-informed tool capable of identifying
154 putative transmission clusters in pathogen phylogenies without the introduction of arbitrary
155 distance thresholds. It is fast, generalizable, and freely available at
156 <https://github.com/alvinxhan/Phydelity>.

157

158 **Acknowledgements**

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

161

162 **Funding**

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

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

169

170 **References**

171 Ambrosioni,J. *et al.* (2012) Impact of highly active antiretroviral therapy on the molecular
172 epidemiology of newly diagnosed HIV infections. *AIDS*, **26**.

173 Barido-Sottani,J. *et al.* (2018) Detection of HIV transmission clusters from phylogenetic trees
174 using a multi-state birth–death model. *J. R. Soc. Interface*, **15**.

175 Bezemer,D. *et al.* (2015) Dispersion of the HIV-1 Epidemic in Men Who Have Sex with Men
176 in the Netherlands: A Combined Mathematical Model and Phylogenetic Analysis. *PLOS*
177 *Med.*, **12**, e1001898.

178 Brenner,B.G. *et al.* (2007) High Rates of Forward Transmission Events after Acute/Early
179 HIV-1 Infection. *J. Infect. Dis.*, **195**, 951–959.

180 Gardy,J.L. and Loman,N.J. (2017) Towards a genomics-informed, real-time, global pathogen
181 surveillance system. *Nat. Rev. Genet.*, **19**, 9–20.

182 Grabowski,M.K. *et al.* (2014) The Role of Viral Introductions in Sustaining Community-
183 Based HIV Epidemics in Rural Uganda: Evidence from Spatial Clustering, Phylogenetics,
184 and Egocentric Transmission Models. *PLOS Med.*, **11**, e1001610.

185 Grabowski,M.K. and Redd,A.D. (2014) Molecular tools for studying HIV transmission in
186 sexual networks. *Curr. Opin. HIV AIDS*, **9**.

187 Han,A.X. *et al.* (2018) Phylogenetic Clustering by Linear Integer Programming (PhyCLIP).
188 *bioRxiv*.

189 de Oliveira,T. *et al.* (2017) Transmission networks and risk of HIV infection in KwaZulu-
190 Natal, South Africa: a community-wide phylogenetic study. *Lancet HIV*, **4**, e41–e50.

191 Prosperi,M.C.F. *et al.* (2011) A novel methodology for large-scale phylogeny partition. *Nat.*
192 *Commun.*, **2**, 321.

193 Ragonnet-Cronin,M. *et al.* (2013) Automated analysis of phylogenetic clusters. *BMC*
194 *Bioinformatics*, **14**, 317.

195 Villandre,L. *et al.* (2016) Assessment of Overlap of Phylogenetic Transmission Clusters and
196 Communities in Simple Sexual Contact Networks: Applications to HIV-1. *PLoS One*, **11**,
197 e0148459.

198