

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Phylogenetic Clustering by Linear Integer Programming (PhyCLIP)

Alvin X. Han^{†,1,2,3}, Edyth Parker^{†,3,4}, Frits Scholer⁵, Sebastian Maurer-Stroh^{1,2,6}, Colin A. Russell³

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

²NUS Graduate School for Integrative Sciences and Engineering, National University of Singapore (NUS), Singapore

³Laboratory of Applied Evolutionary Biology, Department of Medical Microbiology, Academic Medical Centre, University of Amsterdam, Amsterdam, The Netherlands

⁴Department of Veterinary Medicine, University of Cambridge, Cambridge, United Kingdom

⁵Department of Medical Microbiology, Academic Medical Centre, University of Amsterdam, Amsterdam, The Netherlands

⁶Department of Biological Sciences, National University of Singapore, Singapore

[†]These authors contributed equally to this work.

Corresponding authors: Alvin X. Han (hanxc@bii.a-star.edu.sg) and Colin A. Russell (c.a.russell@amc.uva.nl)

20 **Abstract (249/250 words)**

21 Sub-species nomenclature systems of pathogens are increasingly based on sequence data. The use of
22 phylogenetics to identify and differentiate between clusters of genetically similar pathogens is
23 particularly prevalent in virology from the nomenclature of human papillomaviruses to highly pathogenic
24 avian influenza (HPAI) H5Nx viruses. These nomenclature systems rely on absolute genetic distance
25 thresholds to define the maximum genetic divergence tolerated between viruses designated as closely
26 related. However, the phylogenetic clustering methods used in these nomenclature systems are limited
27 by the arbitrariness of setting intra- and inter-cluster diversity thresholds. The lack of a consensus
28 ground truth to define well-delineated, meaningful phylogenetic subpopulations amplifies the difficulties
29 in identifying an informative distance threshold. Consequently, phylogenetic clustering often becomes
30 an exploratory, *ad-hoc* exercise.

31 Phylogenetic Clustering by Linear Integer Programming (PhyCLIP) was developed to provide a
32 statistically-principled phylogenetic clustering framework that negates the need for an arbitrarily-defined
33 distance threshold. Using the pairwise patristic distance distributions of an input phylogeny, PhyCLIP
34 parameterises the intra- and inter-cluster divergence limits as statistical bounds in an integer linear
35 programming model which is subsequently optimised to cluster as many sequences as possible. When
36 applied to the haemagglutinin phylogeny of HPAI H5Nx viruses, PhyCLIP was not only able to
37 recapitulate the current WHO/OIE/FAO H5 nomenclature system but also further delineated informative
38 higher resolution clusters that capture geographically-distinct subpopulations of viruses. PhyCLIP is
39 pathogen-agnostic and can be generalised to a wide variety of research questions concerning the
40 identification of biologically informative clusters in pathogen phylogenies. PhyCLIP is freely available at
41 <http://github.com/alvinxhan/PhyCLIP>.

42

43 **Introduction**

44 Advances in high-throughput sequencing technology and computational approaches in molecular
45 epidemiology have seen sequence data increasingly integrated into clinical care, surveillance systems
46 and epidemiological studies (Gardy and Loman 2017). Based on the growing number of available
47 pathogen sequences genomic epidemiology has yielded a wealth of information on epidemiological and
48 evolutionary questions ranging from transmission dynamics to genotype-phenotype correlations.
49 Central to all of these questions is the need for robust and consistent nomenclature systems to describe
50 and partition the genetic diversity of pathogens to meaningfully relate to epidemiological, evolutionary
51 or ecological processes. Increasingly, nomenclature systems for pathogens below the species level are

52 based on sequence information, supplementing or even displacing conventional biological properties
53 such as serology or host range (Simmonds et al. 2010; McIntyre et al. 2013). However, existing
54 sequence-based nomenclature frameworks for defining lineages, clades or clusters in pathogen
55 phylogenies are mostly based on arbitrary and inconsistent criteria.

56 Standardising the definition of a phylogenetic cluster or lineage across pathogens is complicated by
57 differences in characteristics such as genome organization and maintenance ecology. Cluster
58 definitions vary widely even between studies of the same pathogen, limiting generalisation and
59 interpretation between studies as designated clusters, clades and/or lineages carry inconsistent
60 information in the larger evolutionary context (Grabowski et al. 1904; Dennis et al. 2014; Hassan et al.
61 2017).

62 In virology, nomenclature systems are largely reliant on absolute distance thresholds that define the
63 maximum genetic divergence tolerated between viruses designated as closely related (Burk et al. 2011;
64 Van Doorslaer et al. 2011; Lauber and Gorbalenya 2012; Donald et al. 2013; Kroneman et al. 2013;
65 Poon et al. 2015; Smith et al. 2015; Poon et al. 2016; Valastro et al. 2016). Groups of closely related
66 viruses are inferred to be phylogenetic clusters when the genetic distance between them is lower than
67 the limit set on within-cluster divergence. Non-parametric distance-based clustering approaches have
68 defined the distance between sequences using pairwise genetic distances calculated directly from
69 sequence data (WHO/OIE/FAO H5N1 Evolution Working Group 2008; Aldous et al. 2012; Ragonnet-
70 Cronin et al. 2013) or pairwise patristic distances calculated from inferred phylogenetic trees (Hué et al.
71 2004; Prosperi et al. 2011; Poon et al. 2015; Pu et al. 2015; Ortiz and Neuzil 2017). Within-cluster limits
72 on tolerated divergence have been set using mean (WHO/OIE/FAO H5N1 Evolution Working Group
73 2008), median (Prosperi et al. 2011) or maximum within-cluster pairwise genetic or patristic distance
74 (Ragonnet-Cronin et al. 2013). Some methods incorporate additional criteria, such as the statistical
75 support for subtrees under consideration or minimum/maximum cluster size (Hué et al. 2004; Prosperi
76 et al. 2010; Prosperi et al. 2011; Ragonnet-Cronin et al. 2013). These genetic distance-based clustering
77 approaches are convenient, as they are rule-based and scalable, allowing for relatively easy
78 nomenclature updates. Arguably, flexibility in the distance thresholds allows researchers to curate
79 clusters based on consistency of the geographic or temporal metadata.

80 The central limitation of approaches based on pairwise genetic or patristic distance is that thresholds to
81 define meaningful within- and between-cluster diversity are arbitrary. For most pathogens, there is no
82 clear definition of a well-delineated phylogenetic unit to underlie nomenclature designation or suggest
83 what additional information would be informative to delineate subpopulations e.g. information on
84 antigenicity or geography or host range. Resultantly, there is no ground truth to optimise distance
85 thresholds when developing a nomenclature system for most pathogens. Partitioning phylogenetic trees

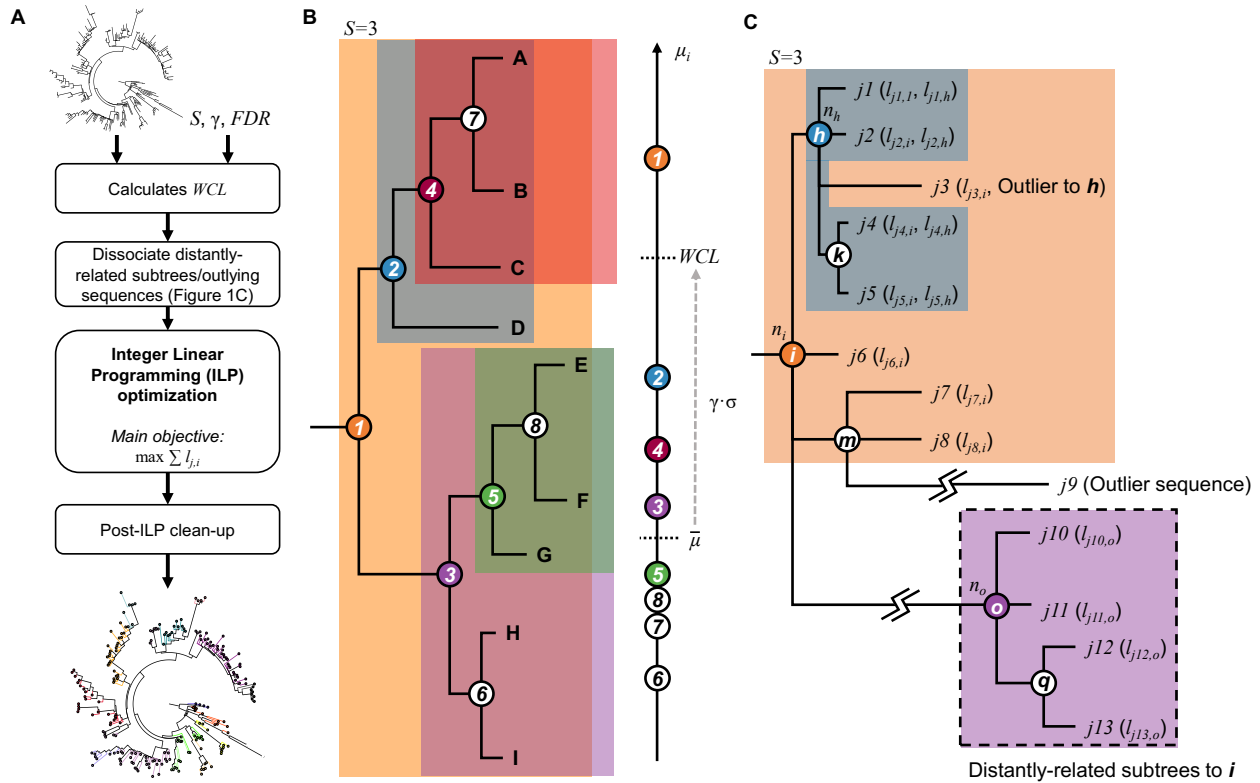
86 into meaningful subsets is therefore complex and is mostly performed *ad hoc* through exploratory
87 analyses with uninformative sensitivity analyses across thresholds. As expected, cluster membership is
88 highly sensitive to the threshold applied and therefore results can be unstable across different cluster
89 definitions (Rose et al. 2017).

90 There is a need for a consistent, automated and robust statistical framework for determining cluster-
91 defining criteria in nomenclature frameworks. Here, we describe a statistically-principled phylogenetic
92 clustering approach called PhyCLIP. PhyCLIP is based on integer linear programming (ILP)
93 optimisation, with the objective to assign statistically-principled cluster membership to as many
94 sequences as possible. We apply PhyCLIP to the haemagglutinin (HA) phylogeny of the highly
95 pathogenic avian influenza (HPAI) A/goose/Guangdong/1/1996 (Gs/GD)-like lineage of the H5Nx
96 subtype viruses, which underlies the most prominent nomenclature system for avian influenza viruses
97 and which itself is based on a genetic distance approach (WHO/OIE/FAO H5N1 Evolution Working
98 Group 2008).

99 PhyCLIP is freely available on github (<http://github.com/alvinxhan/PhyCLIP>) and documentation can be
100 found on the associated wiki page (<http://github.com/alvinxhan/PhyCLIP/wiki>).

101

102 New approach



103

104 **Figure 1: Schematics of PhyCLIP workflow and inference.** **A.** Workflow of PhyCLIP. Apart from an appropriately rooted
 105 phylogenetic tree, users only need to provide S , γ and FDR as the inputs for PhyCLIP. After determining the within-cluster
 106 divergence limit (WCL), PhyCLIP dissociates distantly related subtrees and outlying sequences that inflate the mean patristic
 107 distance (μ_i) of ancestral subtrees. The integer linear programming (ILP) model is then implemented and optimised to assign
 108 cluster membership to as many sequences as possible. If a prior of cluster membership is given, this is followed by a
 109 secondary optimisation to retain as much of the prior membership as is statistically supportable within the limits of PhyCLIP.
 110 Post-ILP optimisation clean-up steps are taken before yielding finalised clustering output. **B.** PhyCLIP considers the
 111 phylogeny as an ensemble of monophyletic subtrees, each defined by an internal node (circled numbers) subtended by a set
 112 of sequences (letters encapsulated within shaded region of the same colour as the circled number). In this example, only
 113 subtrees with ≥ 3 sequences ($S = 3$) are considered for clustering by the ILP model but WCL is determined from μ_i of all
 114 subtrees, including the unshaded subtrees 6-8. Only subtrees where $\mu_i \leq WCL$ are eligible for clustering. **C.** Subtrees **o** and
 115 **q**, as well as sequence $j9$ are dissociated from subtree i as they are exceedingly distant from i . If sequences $j1$, $j2$, $j4$ and $j5$
 116 are clustered under subtree h while $j3$ is clustered under subtree i by ILP optimisation, a post-ILP clean up step will remove
 117 $j3$ from cluster i .

118

119 PhyCLIP requires an input phylogeny and three user-provided parameters:

120 (i) Minimum number of sequences (S) that should be considered a cluster.

- 121 (ii) Multiple of deviations (γ) from the grand median of the mean pairwise sequence patristic distance
122 that defines the within-cluster divergence limit (WCL).
- 123 (iii) False discovery rate (FDR) to infer that the diversity observed for every combinatorial pair of
124 output clusters is significantly distinct from one another.

125 Figure 1A shows the workflow of PhyCLIP which is further elaborated here. First, PhyCLIP considers
126 the input phylogenetic tree as an ensemble of N monophyletic subtrees (including the root) that could
127 potentially be clustered as a single phylogenetic cluster, each defined by an internal node i subtending
128 a set of sequences L_i (Figure 1B, see Methods). Consequently, as the topological structure of the
129 phylogenetic tree is incorporated in the cluster structure, it is possible to infer the evolutionary trajectory
130 of the output clusters of PhyCLIP if the tree is appropriately rooted. For clarity, we use the term *subtree*
131 to refer to the set of sequences subtended under the same node that could potentially be clustered and
132 the term *cluster* to refer to sequences that are clustered by PhyCLIP within the same subtree.

133 The within-cluster internal diversity of subtree i is measured by its mean pairwise sequence patristic
134 distance (μ_i). PhyCLIP calculates the within-cluster divergence limit (WCL), an upper bound to the
135 internal diversity of a cluster, as:

$$WCL = \bar{\mu} + (\gamma\sigma) \quad (1)$$

136 where $\bar{\mu}$ is the grand median of the mean pairwise patristic distance distribution $\{\mu_1, \mu_2, \dots, \mu_i, \dots, \mu_N\}$ and
137 σ is any robust estimator of scale (e.g. median absolute deviation (MAD) or Q_n , see Methods) that
138 quantifies the statistical dispersion of the mean pairwise patristic distance distribution for the ensemble
139 of N subtrees. In other words, only subtrees with $\mu_i \leq WCL$ will be considered for clustering by PhyCLIP
140 (Figure 1B).

141

142 **Distal dissociation**

143 The assumption that a cluster must be monophyletic can lead to incorrect assignment of cluster
144 membership to undersampled, distantly related outlying sequences that have diverged considerably
145 from the rest of the cluster (e.g. sequence $j9$ in Figure 1C). These exceedingly distant outlying
146 sequences can also inflate μ_i of the subtree they subtend, leading to inaccurate overestimation of the
147 internal divergence of the putative subtree. Similarly, distantly related descendant subtrees can
148 artificially inflate μ_i of their ancestral trunk nodes (e.g. nodes o and q in Figure 1C). In turn, historical
149 sequences immediately descending from a trunk node i will not be clustered if its μ_i exceeds WCL
150 (Figure 1C).

151 PhyCLIP dissociates any distal subtrees and/or outlying sequences from their ancestral lineage prior to
152 implementing the integer linear programming (ILP) model. For any subtree i with $\mu_i > WCL$, starting
153 from the most distant sequence to i , PhyCLIP applies a leave-one-out strategy dissociating sequences,
154 and the whole descendant subtree if every sequence subtended by it was dissociated, until the
155 recalculated μ_i without the distantly related sequences falls below WCL . For each subtree, PhyCLIP
156 also tests and dissociates any outlying sequences present. An outlying sequence is defined as any
157 sequence whose patristic distance to the node in question is $> 3 \times$ the estimator of scale away from the
158 median sequence patristic distance to the node. μ_i is recalculated for any node with changes to its
159 sequence membership L_i after dissociating these distantly related sequences. These distal dissociation
160 steps effectively offer PhyCLIP greater flexibility in its clustering construct allowing the identification of
161 paraphyletic clusters on top of monophyletic ones that may better reflect the phylogenetic relationships
162 of these sequences.

163

164 **Integer linear programming optimisation**

165 The full formulation of the ILP model is detailed in Methods. Here, we broadly describe how the
166 optimisation algorithm proceeds to delineate the input phylogeny. The primary objective of PhyCLIP is
167 to cluster as many sequences in the phylogeny as possible subject to the following constraints:

- 168 (i) All output clusters must contain $\geq S$ number of sequences.
169 (ii) All output clusters must satisfy $\mu_i \leq WCL$.
170 (iii) The pairwise sequence patristic distance distribution of every combinatorial pair of output clusters
171 must be significantly distinct from the resultant cluster if sequences from the pair of clusters were to
172 combine. This is the inter-cluster divergence constraint and herein, statistical significance is inferred
173 if the multiple-testing corrected p -value for the cluster pair is $< FDR$ (see Methods).
174 (iv) If a descendant subtree satisfies (i)-(iii) for clustering (e.g. subtree 5 in Figure 1B) and so does its
175 ancestor, which also subtends the sequences descending from the descendant, (e.g. subtree 3 in
176 Figure 1B), the leaves subtended by the descendant will be clustered under the descendant node
177 (e.g. sequences E, F and G will be clustered under cluster 5 in Figure 1B) while the direct progeny
178 of the ancestor subtree will cluster amongst themselves (e.g. sequences H and I will be clustered
179 under cluster 3 in Figure 1B).

180 The ILP model is implemented in a third-party linear programming solver fully integrated within PhyCLIP,
181 to obtain the global optimal solution. At the time of this publication, PhyCLIP supports two third-party
182 solvers:

- 183 1) Gurobi (<http://www.gurobi.com/>) is one of the fastest available commercial mathematical
184 programming solvers. Full-featured academic licenses of Gurobi are available for free to users based
185 at any academic institution.
- 186 2) GNU Linear Programming Kit (GLPK, <http://www.gnu.org/software/glpk>) is a popular, free and open-
187 source linear programming solver.

188 Based on a recent independent benchmark (<http://plato.asu.edu/talks/informs2018.pdf>), Gurobi
189 outperformed GLPK in both performance and speed (Gurobi solved all 40 Simplex LP test problems
190 while GLPK could only solve 31 of them with a geometric mean runtime that was 52 times longer than
191 Gurobi). As such, it is highly recommended that any users with internet access via an academic domain
192 run PhyCLIP with the Gurobi solver. All clustering results presented in this manuscript were obtained
193 using Gurobi.

194

195 **Post-ILP clean-up**

196 While distal dissociation prior to ILP optimisation works well for dissociating distantly related subtrees
197 and sequences, it is ineffective in identifying spurious singletons such as sequence j_3 in Figure 1C.
198 Here, in terms of sequence patristic distance, sequence j_3 is an outlying sequence to the descendant
199 node h but not so to the ancestral node i . If taxa subtended by subtree h (i.e. j_1, j_2, j_4 and j_5) were to
200 be clustered without j_3 which itself is clustered under cluster i , PhyCLIP performs a post-ILP
201 optimisation clean-up step that removes j_3 from output cluster i . This is because j_3 is clearly a
202 topologically outlying taxon to i and if unremoved, would imply that sequences clustered under cluster
203 h (i.e. j_1, j_2, j_4 and j_5) can belong to cluster i as well.

204 PhyCLIP also offers the user an optional clean-up step that subsumes subclusters into their parent
205 clusters if sequences in the descendant subcluster are still associated with the parent cluster (i.e. not
206 removed by distal dissociation) and that coalescing with the parent clusters does not lead to violation of
207 the statistical bounds that define the clustering result. This may be useful if the user prefers a relatively
208 more coarse-grained clustering (e.g. nomenclature building). As mentioned earlier, so long as a
209 statistically significant distinction could be made between a descendant subtree and its ancestral
210 lineage, the ILP model enforces the progeny sequences of the descendant subtree to cluster in the
211 descendant cluster. In turn, PhyCLIP is sensitive to the detection of clusters of highly related or identical
212 sequences that minimally satisfies the minimum cluster size (S), as their distributions are statistically
213 distinct from the rest of the population. This sensitivity may lead to over-delineation of the tree and/or
214 multiple nested clusters. Notably, these sensitivity-induced subclusters are not false-positive clusters
215 and meet the same statistical criteria as all other clusters. However, some users may want to subsume
216 these subclusters into parent clusters to facilitate higher level interpretation.

217 **Optimisation criteria**

218 PhyCLIP's user-defined parameters can be calibrated across a range of input values, optimising the
219 global statistical properties of the clustering results to select an optimal parameter set. The optimisation
220 criteria are prioritised by the research question, as the clustering resolution and cluster definition are
221 dependent on the question, and therefore the degree of information required to capture ecological,
222 epidemiological and/or evolutionary processes of interest. Users may want a high-resolution clustering
223 result, with the phylogenetic tree delineated into a large number of small, high confidence clusters with
224 very low internal divergence, tolerating a higher number of unclustered sequences. Other users may
225 want a more intermediate resolution, with more broadly defined clusters that are still well-separated but
226 encompass the majority of data in the tree (Figure S1A).

227 PhyCLIP's optimisation criteria are agnostic to the metadata of the dataset and include: 1) The grand
228 mean of the pairwise patristic distance distribution and its standard deviation. The grand mean is a
229 measure of the within-cluster divergence and can be optimised to select a clustering configuration with
230 the lowest global internal divergence. 2) The mean of the inter-cluster distance to all other clusters and
231 its standard deviation. This can be optimised to select a clustering configuration with well-separated
232 clusters. 3) The percentage of sequences clustered, which can be optimised to minimise the number of
233 unclustered sequences. 4) The total number of clusters and 5) mean or median cluster size, which can
234 be optimised to select a tolerable level of stratification of the tree.

235 The ranges of input parameters considered are also dependent on the characteristics of the dataset.
236 The minimum cluster size range considered should be a factor of the size of the phylogenetic tree,
237 whereas the multiple of deviation (γ) considered should be a factor of the intra- and inter-cluster distance
238 related to the research question.

239 Metadata can be incorporated to validate PhyCLIP's optimisation. The spatiotemporal structure of
240 phylogenies can inform clustering results if within-cluster variation in metadata such as collection times
241 or geographic origin is used as a *post-hoc* optimisation criterion. Within-cluster pairwise geographic
242 distance between the origins of sequences can act as an incomplete ground truth to determine whether
243 a clustering result delineates meaningful clusters if there is a reasonable expectation that clusters are
244 defined by spatial factors. The within-cluster deviation in collection dates can also be included as an
245 optimisation criterion if clusters are expected to be temporally structured.

246

247 Results

248 To evaluate the utility of PhyCLIP we compared its clustering of the global HPAI H5Nx virus data against
249 the WHO/OIE/FAO nomenclature (WHO/OIE/FAO HN Evolution Working Gr 2009; Smith et al. 2015).
250 The WHO/OIE/FAO H5 nomenclature has been updated progressively since its development in 2007
251 as new sequences are added to the global phylogeny including updates in 2009 and 2015. The primary
252 analysis of PhyCLIP's performance was assessed with the full dataset of H5N1 haemagglutinin (HA)
253 sequences included in the WHO/OIE/FAO H5 nomenclature update of 2015 (n=4357), with comparison
254 to the WHO/OIE/FAO clade designation. PhyCLIP was run with different combinations of the parameters
255 varied over the following ranges: a minimum cluster size of 2-10, a multiple of deviation (γ) of 1-3, and
256 an FDR of 0.05, 0.1, 0.15 or 0.2. The optimisation criteria were prioritised as follows: 1) percentage of
257 sequences clustered, 2) grand mean of within-cluster patristic distance distribution, 3) mean within-
258 cluster geographic distance and 4) mean of the inter-cluster distances.

259 The percentage of sequences clustered was prioritised as the primary optimisation criterion to ensure
260 that the maximum number of sequences were assigned a nomenclature identifier. Mean within-cluster
261 geographic distance was included as a *post-hoc* optimisation criterion as many avian influenza viruses
262 cluster with high spatial consistency owing to their transmission dynamics in localised avian populations.
263 For influenza viruses endemic to poultry such as H5Nx, this is likely owing to increased local
264 transmission during outbreaks in large poultry populations, as well as the associated sampling biases
265 (Smith et al. 2015). Within-cluster genetic divergence was optimised with higher priority than within-
266 cluster mean geographic distance, as the use of phylogenetic geographic structure as a ground truth for
267 avian influenza viruses is restricted by the long-distance dissemination of related viruses through
268 mechanisms such as the poultry trade or migration of wild birds (WHO/OIE/FAO H5N1 Evolution
269 Working Group 2014a; Smith et al. 2015a). The within-cluster geographic distance was calculated for
270 each cluster in each clustering result as the mean within-cluster pairwise Vicenty distance in miles.

271 The temporal consistency of clusters can also be used as optimisation criteria for measurably evolving
272 viruses such as Influenza A virus (Drummond et al. 2003). Results ranking the grand mean within-cluster
273 standard deviation in sampling dates of each clustering result as the fourth optimisation criterium, with
274 mean of the inter-cluster distance in fifth, were identical to those only including the above mentioned
275 four optimisation criteria.

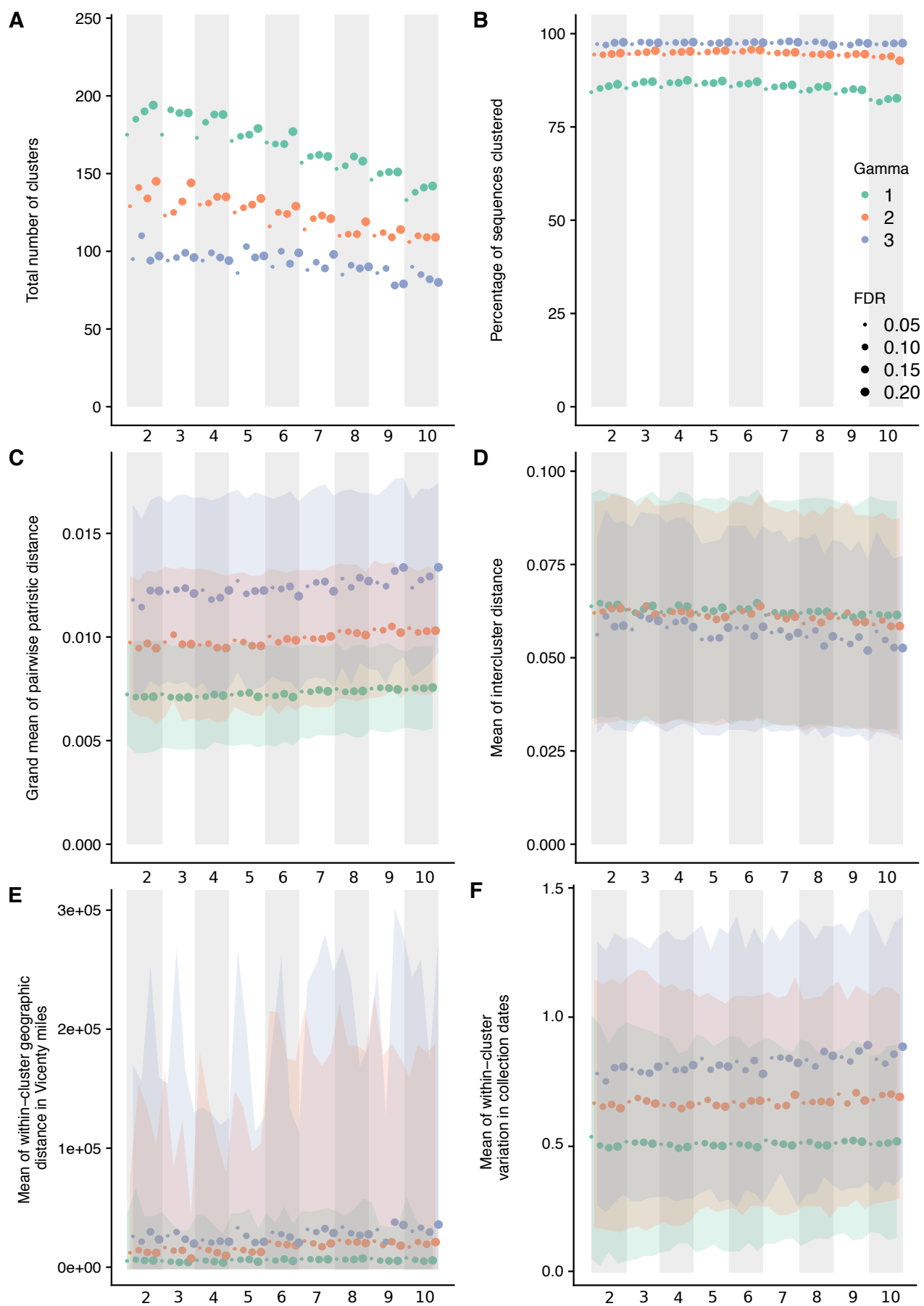
276 As PhyCLIP incorporates topological information of the phylogeny into the clustering construct, non-
277 terminal internal nodes with zero branch lengths can lead to erroneous clustering and over-delineation
278 (Figure S1B). Such internal nodes are usually found in bifurcating trees as representations of
279 polytomies, arising from a lack of phylogenetic signal among the sequences subtended by the node to
280 resolve them into dichotomies. As such, prior to implementing PhyCLIP, all non-terminal, zero branch

281 length nodes in the input phylogenetic trees were collapsed into polytomies, which more accurately
282 depicts the relationship between identical/indiscernible sequences and/or ancestral states. In the H5Nx
283 analysis, all subclusters were subsumed if the statistical requisites of the parent clade were maintained,
284 to aid in easing the interpretation of the nomenclature designation (as discussed in the New Approach
285 section).

286

287 **Influence of the parameters**

288 The influence of the parameters on PhyCLIP's clustering properties was assessed with the 2015-update
289 H5 phylogeny. Lower multiples of deviation (γ) define a more conservative expected range for tolerated
290 within-cluster divergence, informed by the global pairwise patristic distance distribution (Figure S2). As
291 a result, clusters designated at a γ of 1 have the lowest internal divergence, measured by the grand
292 mean of the pairwise patristic distance distribution (Figure 2C). These clusters are expected to be highly
293 related, with low variation in clustered sequence spatiotemporal metadata (Figure 2E-F). More
294 conservative ranges of tolerated within-cluster divergence result in a higher clustering resolution with a
295 greater number of clusters, lower mean cluster sizes and a higher percentage of sequences unclustered
296 (Figure 2A-B). A higher γ increases the limit of tolerated within-cluster divergence, resulting in a lower
297 clustering resolution that coalesces smaller clusters into larger, more internally-divergent clusters. The
298 collapsing of the smaller clusters decreases the total number of clusters while concurrently increasing
299 the percentage of sequences clustered and mean cluster size. The influence of γ is less pronounced for
300 the mean inter-cluster distance, with no apparent distinction between $\gamma = 1$ and 2. The total number of
301 clusters decreases approximately linearly as the minimum cluster size (S) increases from two towards
302 ten (Figure 2A). Lower FDRs are more conservative in designating the pairwise patristic distance
303 distributions of two clusters as statistically distinct. A higher or less conservative FDR therefore
304 designates more similar distributions as distinct from one another, increasing the number of clusters
305 (Figure 2A). The effect of FDR is muted at a higher minimum cluster size or higher γ , as these
306 parameters designate larger clusters, which limits the number of clustering configurations available.



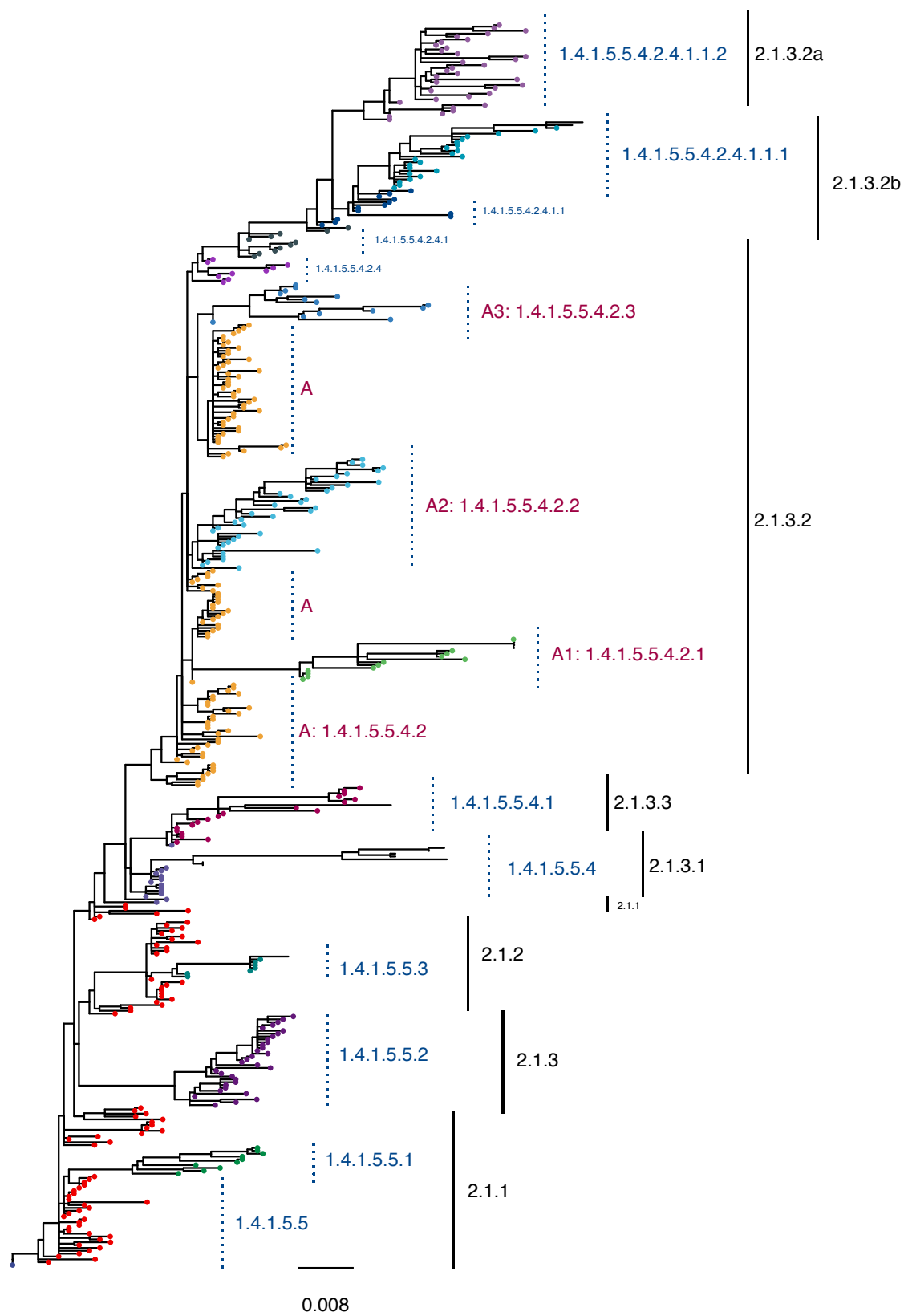
308 **Figure 2: Influence of parameters on the clustering properties of PhyCLIP in the WHO/OIE/FAO 2015-update**
309 **phylogeny.** Figures A-F have the parameter set combinations ordered according to minimum cluster size, FDR and γ on the
310 x axis. The banded background and x-axis subscript numbering indicate the minimum cluster size of the parameter set. Marker
311 colour and size is indicative of the γ and the FDR respectively of the parameter set as indicated by the legend in Figure B. **A.**
312 Total number of clusters. **B.** Percentage of sequences clustered. **C.** Grand mean of the pairwise patristic distance distribution.
313 **D.** Mean of the inter-cluster distance to all other clusters. **E.** Mean within-cluster geographic distance calculated in Vicenty
314 miles. **F.** Mean within-cluster standard deviation in collection dates.

315
316 **Optimal PhyCLIP clustering result for HPAI avian H5 viruses**

317 For the full phylogeny of Gs/GD-like H5 viruses from the 2015 nomenclature update, the optimal
318 parameter set combined a minimum cluster size of 7, an FDR of 0.15 and a γ of 3. The optimal clustering
319 configuration clustered 98% of the sequences into a total of 89 clusters with a median cluster size of 21
320 sequences. The topology of the optimal clustering result yields informative source-sink trajectories that
321 are supported by previously reported phylogenetic and phylogeographic evidence of the global
322 panzootic of the Gs/GD-like H5N1 lineage (Duan et al. 2008; Wang et al. 2008; Smith et al. 2015; The
323 Global Consortium for H5N8 and Related Influenza Viruses 2016).

324 Principally, pathogen nomenclature systems should delineate population structure, highlighting the
325 underlying population dynamics that may be informative about the evolutionary trajectory of pathogen
326 variants. The distal dissociation approach of PhyCLIP produces a clustering topology where divergent
327 subclusters nest within a larger cluster structure termed a supercluster, as exemplified with
328 WHO/OIE/FAO clade 2.1x viruses in Figure 3. Sufficiently diverse subclusters are dissociated from the
329 ancestral trunk node of a putative cluster. This enables the remaining sequences that meet the statistical
330 criteria to cluster with the ancestral node based on their pairwise patristic distance, as the divergent
331 subcluster is no longer inflating the ancestral node's mean pairwise patristic distance above the within-
332 cluster limit. Cluster A in Figure 3 depicts the supercluster topology: the source population viruses (tips
333 in yellow) are annotated as A, and the divergent descendant subclusters are annotated as A.1, A.2 and
334 A.3 respectively. This approach captures source-sink ecological dynamics: the supercluster acts as a
335 putative source population to its subclusters, reflecting the clear evolutionary divergence and trajectory
336 of descendants of the source population (sub-lineages). The nomenclature system algorithmically
337 imposed on PhyCLIP's clustering for avian influenza is designed to enhance the evolutionary information
338 in the clustering (see Methods).

339



341 **Figure 3: Phylogeny of the Clade 2.1x viruses circulating in Indonesia.** The WHO/OIE/FAO H5 nomenclature is annotated
342 in black. PhyCLIP's cluster designation is indicated in blue, corresponding to tip colour. PhyCLIP's supercluster topology is
343 exemplified by Cluster A. The source population of the supercluster is annotated as A in pink, with tips coloured yellow. The
344 divergent descendant clusters are annotated as A.1, A.2 and A.3 respectively here. The letter A here is shorthand for its
345 nomenclature address, 1.4.1.5.5.4.2. This nomenclature address indicates that supercluster A is the second descendant of
346 cluster 1.4.1.5.5.4 (indicated in light purple), which in turn is the fourth descendant of the source supercluster 1.4.1.5.5, indicated
347 in red. See Methods sections for full explanation of nomenclature addresses.

348

349 PhyCLIP's optimal cluster designation delineated the spatiotemporal structure of the phylogeny at high
350 resolution (Figure S3). Viruses circulating in south, central and northeast China and Hong Kong in 1996-
351 2003 acted as the source population for the emergence of the classical viruses, seeding four lineages
352 (cluster 1, seeding cluster 1.1-1.4, Table S1). The second supercluster captures the first major wave of
353 expansion into neighbouring countries in east and southeast Asia in the early 2000s, with a source
354 population of viruses circulating in south central, east and north China, Viet Nam and Hong Kong in
355 2000-2003 (1.4 and 1.4.1 and their descendant lineages). The third supercluster captures the second
356 major wave of expansion of the Gs/GD-like H5 viruses, characterised by global spread (cluster 1.4.1.5
357 and its descendants). The source population of viruses from east, south central and southwest China,
358 Hong Kong and Viet Nam circulated from 2002-2005, giving rise to diverse and distinct viral lineages in
359 different regions globally (1.4.1.5.1-6). The supercluster topology highlights single lineage introductions
360 for countries with endemic circulation such as Indonesia and Egypt, but delineates multiple co-circulating
361 lineages structured over time. The clustering topology also highlights multiple incursions of diverse
362 viruses into countries such as South Korea and Japan (Table S3).

363 In addition to source-sink dynamics, distal dissociation also identifies probable outlying sequences,
364 defined as sequences more than 3 times the estimator of scale away from the median patristic distance
365 to the node. For example, PhyCLIP identifies seven outliers in its delineation of WHO/OIE/FAO clade
366 2.3.2.1c in the 2015 phylogeny (indicated by red tip-points in Figure 4). These sequences may represent
367 under-sampled populations with unobserved diversity, introductions from otherwise unsampled
368 populations or lower quality sequences introducing error into phylogenetic reconstruction.

369

370 **Comparison to the WHO/OIE/FAO H5 nomenclature**

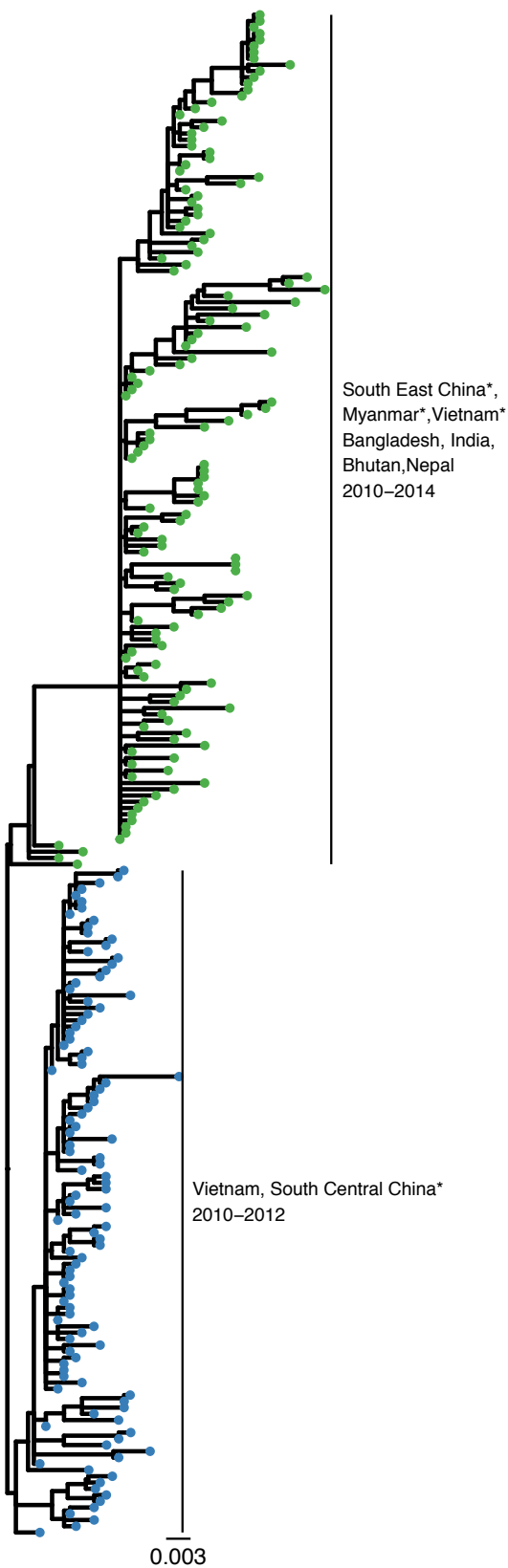
371 The current WHO/OIE/FAO nomenclature system designates 43 different clades and 7 clade-like
372 groupings for the full H5 phylogeny as of the 2015 update (Smith, Donis, and WHO/OIE/FAO H5
373 Evolution Working Group 2015) (Table S2). PhyCLIP recovers the current WHO/OIE/FAO H5
374 nomenclature with varying degrees of agreement across parameter sets, as measured by the variation
375 of information (VI) between the clustering partitions (Figure S4). VI is an information theoretic criterion

376 for comparing partitions of the same data set, based on the information lost and gained when moving
377 between partitions (Meilă 2007). A lower VI indicates more similar partitions. Parameter sets with a γ of
378 3 consistently had the lowest VI compared to the WHO/OIE/FAO system, indicating that the
379 WHO/OIE/FAO nomenclature system has the highest agreement with PhyCLIP clustering results that
380 tolerate higher within-cluster divergence.

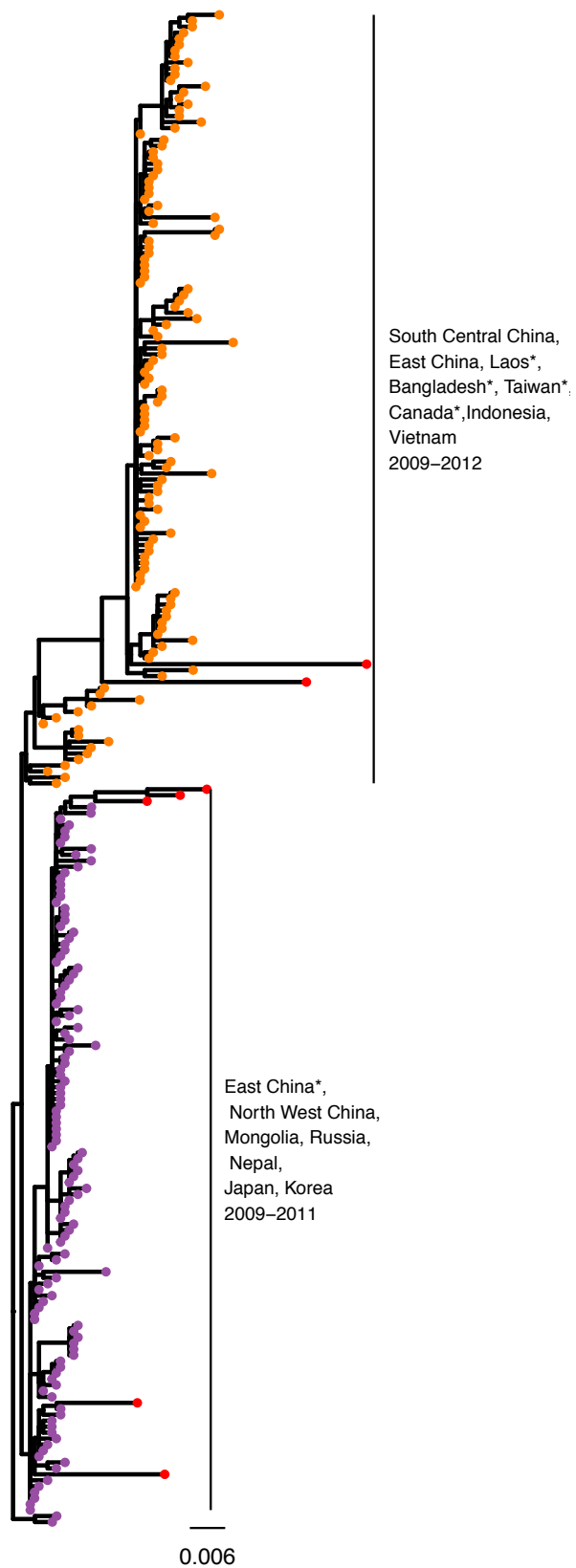
381 In the optimal clustering result, PhyCLIP delineates the spatiotemporal structure of the phylogeny with
382 a higher resolution than the WHO/OIE/FAO nomenclature system (89 vs 50 phylogenetic units, Figure
383 S3). The supercluster structure of the PhyCLIP clustering topology recapitulates the hierarchical
384 structure of the WHO/OIE/FAO nomenclature (Figure 3). Simultaneously, PhyCLIP's clustering captures
385 clear lineage distinctions for viruses from different geographic regions and years in several
386 WHO/OIE/FAO demarcated clades. For example, PhyCLIP delineates clade 2.3.2.1a into two separate
387 clusters: 1) a cluster that circulated in Viet Nam in 2011-2012, with sporadic detection in south central
388 China and 2) a cluster that circulated largely in Bangladesh, India, Bhutan and Nepal from 2010 to 2014,
389 with single viruses detected in south east China, Viet Nam and Myanmar (Figure 4A). PhyCLIP also
390 delineates clade 2.3.2.1c into two clusters: 1) a cluster that captures the expansion of viruses from north
391 west and east China into Mongolia, Russia, Nepal, Japan and Korea for the period 2009-2011, and 2)
392 a cluster that predominantly circulates in China, Viet Nam and Indonesia for 2009-2012, with single
393 viruses from Lao PDR, Bangladesh and Taiwan (Figure 4B).

394

Clade 2.3.2.1a



Clade 2.3.2.1c



396 **Figure 4: PhyCLIP's delineation of WHO/OIE/FAO demarcated clades 2.3.2.1a (A) and 2.3.2.1c (B).** Tips are coloured
397 according to PhyCLIP's cluster designation. The tips coloured in red in B are viruses that were designated as outliers by
398 PhyCLIP's outlier detection. Countries represented by single viruses in the cluster are indicated with an asterisk.

399 **Impact of sampling**

400 PhyCLIP's clustering results are sensitive to the diversity in the input population that informs the global
401 distribution and resultant sampling. The influence of sampling was assessed by comparing the optimal
402 clustering result of the phylogeny underlying the WHO/OIE/FAO H5 2015 nomenclature (n=4357) to the
403 phylogeny underlying the 2009 nomenclature update (n=1224), a subset nested in the 2015-update
404 phylogeny. The WHO/OIE/FAO 2009 nomenclature update was performed after the geographic
405 expansion and divergence of clade 2.2, which necessitated further delineation into clade 2.2.1. It
406 designated 20 clades, including 8 third order clades (WHO/OIE/FAO HN Evolution Working Gr 2009).
407 The WHO/OIE/FAO 2015 nomenclature update includes approximate 3.5-times the number of
408 sequences as the 2009 nomenclature update, and includes novel clade designation to the fourth and
409 fifth order WHO/OIE/FAO H5 Evolution Working Group 2015). The optimal PhyCLIP parameter set for
410 the 2009 WHO/OIE/FAO nomenclature system combines a minimum cluster size of 3, a FDR of 0.2 and
411 a γ of 3. In the 2009 tree, this clustered 98% of the n=1224 viruses into 39 clusters, with a median
412 cluster size of 12 (Figure S5).

413 Overall, the source-sink inference of PhyCLIP's clustering topology is largely consistent between the
414 WHO/OIE/FAO 2009 and 2015 update phylogeny optimal clustering results (Table S1). The optimal
415 result for the 2009 update phylogeny captures a similar topology and source population for the South
416 East Asian (clusters 1.3.1 and 1.3.1.1) and the post-2005 global wave of expansion (cluster
417 1.3.1.1.2.2.2) compared to the optimal 2015 clustering, with substantial overlap between the source
418 populations identified (100% and 83% for source populations for southeast Asia wave and global wave
419 respectively).

420 Changes in the clustering topology between the 2009 and 2015 update phylogenies are expected as
421 the underlying datasets are substantially different. More than 3000 viruses were added to the tree in the
422 six years between nomenclature updates. The Gs/GD-like H5 viruses evolved significantly in the
423 intervening period owing to genetic drift and reassortment. The addition of a large number of divergent
424 viruses to the underlying dataset fundamentally alters the ensemble statistical properties of the tree,
425 driving changes in the clustering configuration by changes in the global patristic distance distribution,
426 topology and statistical power between datasets. As a result, the ecological inferences drawn from the
427 2015 clustering topology are different from that of the 2009 phylogeny (Table S1).

428 Primarily, the addition of a set of highly divergent sequences increases the spread of the global pairwise
429 patristic distance distribution (Figure S2). The within-cluster limit it informs increases concurrently,

430 increasing the tolerance of allowable within-cluster divergence. In the distal dissociation approach,
431 increased tolerance of divergence would allow for the incorporation of more distant trunk viruses into
432 supercluster source populations if the enclosed viruses are sufficiently distinct to be dissociated as
433 independent clusters (Figure S6). If the within-cluster limit is lowered, inclusion of the considered trunk
434 viruses will violate the within-cluster limit. Resultantly, these trunk viruses and their descendants will be
435 assessed for clustering as independent subtrees.

436 Clustering changes between 2009 and 2015 update phylogenies are also induced by the local effects
437 of the addition of multiple lineages to the 2015 phylogeny within clusters defined in 2009 owing to their
438 continued circulation and diversification post-2009. Notably, many distinct clusters in the 2009
439 phylogeny are structured as source populations in superclusters in the 2015 phylogeny (Figure S7).
440 Here, PhyCLIP identifies that the statistical properties of these divergent post-2009 lineages are distinct
441 enough to reliably dissociate them from the ancestral node and delineate them as separate clusters.
442 The viruses present in the 2009 phylogeny that these divergent lineages descend from meet the within-
443 cluster limit after the dissociation and are structured as the source population to the post-2009 nested
444 diversity.

445 Topological differences between phylogenetic trees built from different underlying datasets can also
446 drive changes in PhyCLIP's clustering, as observed for the classical clade 0 viruses (Figure S6). The
447 source population of the classical clade viruses for both the 2009 and 2015 updates optimal clustering
448 result is estimated to have originated from south central and east China and Hong Kong in 1997-2003.
449 However, the 2015 cluster designation resolves an additional seed lineage within the 2009-source
450 population (Figure S6). In the 2009 phylogeny, this additional cluster forms part of the source population
451 as it is part of the trunk of the tree. The equivalent cluster does not form part of the trunk of the tree in
452 the 2015 phylogeny and is dissociated as a statistically distinct cluster. Moreover, the substantial
453 increase in the number of viruses between the 2009 and 2015 datasets along with the increase in
454 diversity results in more statistical power to delineate among groups of viruses resulting in a higher
455 clustering resolution for the 2015 phylogeny.

456

457 **Comparison of optimal to suboptimal clustering results**

458 So far, we have focused our interpretation on the optimal PhyCLIP clustering. To ensure that our results
459 were robust across similarly optimal PhyCLIP parameter sets we compared the optimal set against the
460 next four similarly optimal sets. Comparing the top 5 clustering results ranked by the optimality criterion
461 (in order of greatest number of sequences clustered, lowest internal genetic and geographic divergence,
462 and greatest average between-cluster distance), the clustering result from the optimal parameters set
463 of the 2015 phylogeny was generally consistent with those generated from the four highest-ranked

464 suboptimal parameter sets (see Figure S8). Each of the top four suboptimal clustering was found to
465 have low VI (0.817-0.984) relative to the optimal clustering, with a large proportion (74.4%-82.7%) of
466 viruses clustered in the same corresponding clusters. The supercluster source populations leading to
467 the early 2000 expansion into east and southeast Asia as well as the global expansion in 2005 were
468 similarly found in all suboptimal results.

469 However, changes to parameter sets fundamentally changed the statistical constraints defining the
470 clustering solution space and in turn, altered the partitions between resultant clusters. Specifically, in
471 this case where $\gamma = 3$ in all five optimal/suboptimal parameter sets, varying minimum cluster size not
472 only changed the distribution of putative subtrees for clustering but the distribution of inter-cluster
473 divergence p -values for multiple-testing correction as well. As such, while the global superclusters were
474 largely recapitulated in the suboptimal results, local partitions of co-circulating viruses descending from
475 these supercluster sources, and consequently the inferences of source-sink dynamics, varied amongst
476 the different parameter sets.

477

478 **PhyCLIP clustering of the 1996-2018 H5Nx phylogeny**

479 In recent years the Gs/GD-lineage of H5 viruses has undergone substantial evolution, with viruses from
480 WHO/OIE/FAO clade 2.3.4.4 reassorting with co-circulating viruses to give rise to multiple H5Nx
481 subtypes including H5N2, H5N5, H5N6 and H5N8. We applied PhyCLIP to a phylogeny representing
482 the Gs/GD-lineage up to and including early 2018 to investigate how the global expansion of the H5Nx
483 viruses changes clustering inference ($n=7898$) (Figure S9, S10). Applying the same optimisation
484 approach described above, the optimal parameter set for the 2018 phylogeny combines a minimum
485 cluster size of 4, a FDR of 0.2 and a γ of 3. This parameter set clustered 97% of the viruses into 135
486 clusters, with a median cluster size of 23 (Figure S11).

487 The addition of the H5Nx viruses collected from 2014-2018 to the 2015 phylogeny changed the
488 distribution in two ways: 1. it added diversity to the right tail of the distribution, owing to the increased
489 divergence of the H5Nx viruses compared to the H5N1 viruses; 2. it increased the number of putative
490 clusters with low internal divergence, as a large amount of the H5Nx viruses possess highly similar HA
491 genes owing to both sampling biases during outbreaks and the relative short circulation time following
492 their emergence. This shift in the distribution reduced the within-cluster limit compared to that of the
493 2015 dataset (Figure S2).

494 Filtering the 2015-update and 2018 datasets (see Methods) resulted in changes in tree topology and
495 overall sequence diversity, and consequently altered the ecological inference of source-sink clusters
496 circulating from 1997-2005 (Table S1). However, the ecological inferences of the second major wave of
497 expansion, the post-2005 global expansion characterised by cluster 1.2.1.1.1.3.2 and its descendants

498 1.2.1.1.1.3.2.1-8, were largely consistent across the 2009 (cluster 1.3.1.1.2.2.2), 2015 (cluster 1.4.1.5)
499 and 2018 (cluster 1.2.1.1.1.3.2) trees, including a shared core source population (Table S1).

500 The WHO/OIE/FAO clade 2.3.4.4 viruses are of interest owing to their reassortment promiscuity and
501 rapid global expansion. PhyCLIP delineates the clade 2.3.4.4 viruses into two distinct lineages, seeded
502 from a source population of viruses circulating in east and south-central China and Malaysia in 2005-
503 2010 (cluster 7.8, Table S1). The first lineage circulated in east, south central and northeast China from
504 2008 to 2011 (7.8.2, Figure S11, Table S1). The second lineage (7.8.3) circulated in south central and
505 east China in 2008-2012 and seeded six distinct sub-lineages: Lineage 7.8.3.1 circulated in China from
506 2010 to 2014 before expanding to Viet Nam and circulating there for 2014-2015. Lineage 7.8.3.2
507 captures the global expansion of viruses from 2009 onwards. This includes the early subclade of H5N8
508 viruses described in Lycett et al (The Global Consortium for H5N8 and Related Influenza Viruses 2016).
509 Lineage 7.8.3.3 was restricted to China and was detected in 2013-2016. Lineage 7.8.3.4 also captures
510 a pan-national lineage that was detected from 2014 to 2016, and captures the more recent H5N8
511 subclade described in Lycett et al (The Global Consortium for H5N8 and Related Influenza Viruses
512 2016). Lineage 7.8.3.5 circulated in east and southeast Asia from 2013 to 2017. Lineage 7.8.3.6 is
513 seeded from a source population of viruses circulating in east and southeast Asia, expanding into
514 multiple co-circulating H5N6 southeast Asian lineages from 2013 onwards (Table S1).

515

516 **Benchmarking against other phylogenetic clustering tools**

517 PhyCLIP was benchmarked for performance against two open-source non-parametric clustering tools,
518 PhyloPart (Prosperi et al. 2011) and ClusterPicker (Ragonnet-Cronin et al. 2013). Both tools require a
519 phylogenetic tree as input, as well as a user-specified distance threshold and minimum statistical node-
520 support level. Additionally, both algorithms carry out a depth-first traversal of the tree, considering
521 subtrees as putative clusters if the node support is above the user-defined level. In PhyloPart, the user
522 specifies a percentile of the global pairwise patristic distance distribution as a threshold. If the median
523 of the pairwise patristic distances of the putative cluster is below the percentile threshold, a cluster is
524 designated. ClusterPicker requires a user-defined maximum pairwise genetic distance (calculated as p-
525 distance directly from the sequences) threshold for cluster designation. In both tools, a subtree is
526 designated as a cluster if it meets the respective clustering criteria. If the subtree violates the clustering
527 criteria, the algorithm tests the children of the subtree as potential clusters until a leaf is reached, when
528 no cluster is designated in the path.

529 In contrast, traversal order has no bearing on the clustering outcomes of PhyCLIP. Although PhyCLIP
530 parses the input phylogeny by level-order, prior to ILP optimisation, PhyCLIP dissociates outlying taxa
531 if $\mu_i < WCL$ and proceeds with full distal dissociation heuristics described in the New Approach section

532 if otherwise for every internal node i in the input tree. In both cases, tip dissociation is performed by
533 ranking taxa based on their patristic distance to node i (i.e. the common ancestor) without consideration
534 of their topological placement. Finally, all putative subtrees (i.e. tree nodes) after distal dissociation are
535 given equal consideration by ILP optimisation to maximally assign cluster membership to all tips (see
536 New Approach). In doing so, not only does PhyCLIP allow for paraphyletic clustering, tree traversal
537 order does not affect clustering results.

538 Accepted practice for these tools is to incorporate previous knowledge of sequence divergence into a
539 distance threshold or to calibrate the threshold over a tolerable range with metadata or expert
540 consensus. The two methods were applied to the 2009-update phylogeny (n=1224 sequences) with
541 thresholds ranging from 0.005 to 0.05 substitutions/site. For PhyloPart, the respective percentile of the
542 global pairwise patristic distance distribution was chosen to match the distance threshold. Required
543 bootstrap support level was set to 0 in both methods to make it comparable to PhyCLIP, which lacks
544 node-support criteria. The optimal threshold was selected by maximisation of the mean silhouette index
545 across the clustering partitions (see Methods). All programs were run on the Ubuntu 16.04 LTS
546 operating system with an Intel Core i7-4790 3.60 GHz CPU.

547 The optimal thresholds and clustering statistics for each of the methods are reported in Table S4. A
548 direct comparison of cluster inference between PhyCLIP and the other methods is difficult owing to
549 notable differences in cluster definitions, as these methods were largely designed to detect highly
550 related clusters of sequences linked by direct transmission events. The optimal clustering result for
551 ClusterPicker by silhouette maximisation had a very low maximum genetic distance threshold at 0.5%.
552 (Figure S12). This resulted in a highly stratified tree with 246 small, highly-related clusters and 33.8%
553 outliers, compared to PhyCLIP's 39 clusters and 2% outliers (VI to PhyCLIP of 2.7) (Figure S13, Table
554 S4).

555 Clustering results between PhyCLIP and PhyloPart's optimal results showed better correspondence,
556 with PhyloPart designating 37 clusters to PhyCLIP's 39 (VI to PhyCLIP of 0.64, Figure S13, Table S4).
557 However, the cluster delineations and inferences drawn are substantially different between the two
558 methods (Table S5). The nomenclature scheme developed for PhyCLIP was applied to PhyloPart's
559 optimal clustering result for a more meaningful comparison. PhyCLIP's distal dissociation approach
560 allows for the identification of paraphyletic clusters, forming supercluster topologies throughout the tree
561 (as discussed above). Notably, PhyloPart's depth-first approach and monophyletic cluster criteria
562 prevent it from designating paraphyletic clusters, obscuring the suggestive source-sink dynamics of
563 H5N1's expansion wave identified by PhyCLIP's distal dissociation approach (Table S5). Concurrently,
564 PhyloPart is unable to identify hierarchical clusters, which PhyCLIP identifies as divergent trajectories
565 nested in larger clusters (Figure S13).

566 PhyCLIP is appreciably more computationally intensive than PhyloPart and ClusterPicker as it not only
567 parses the global pairwise patristic distance distribution of the phylogeny but recursively recalculates
568 the distribution for subtrees in the distal dissociation approach, performs hypothesis testing across every
569 combinatorial pair of subtrees to test their inter-cluster divergence, as well as optimise the ILP model.
570 To relieve some of the computational cost, PhyCLIP is written in Python 2.7 employing multiprocessing
571 modules to parallelise the computational tasks involved resulting in ~3.2x times speedup with 8 CPU
572 cores relative to a single core run (Table 1).

573 Despite the differences in computation time, the principal advantage of PhyCLIP is its use of the
574 background genetic diversity to inform its within-cluster limit without the need to arbitrarily define it or
575 calibrate it over a range of thresholds. This is especially helpful as there is typically a lack of prior
576 knowledge on meaningful delineation of phylogenetic units for most pathogens to recommend a range
577 of distance thresholds. Additionally, PhyCLIP's distal dissociation and outlier detection approaches are
578 capable of identifying informative paraphyletic and hierarchical clusters, unlike the other tools.

579

Approach	Time to completion	Peak memory usage	Number of CPUs
PhyCLIP	1 hour 4 minutes	2.0 GB	8
	3 hours 25 minutes	1.7 GB	1
ClusterPicker	2.8 minutes	0.3 GB	1
PhyloPart	10.6 minutes	4.1 GB	8

580 Table 1: Benchmarking the performance of PhyCLIP against widely-used phylogenetic clustering tools

581

582 Discussion

583 PhyCLIP provides a statistically-principled, phylogeny-informed framework to assign cluster
584 membership to taxa in phylogenetic trees without the introduction of arbitrary distance thresholds for
585 cluster designation. PhyCLIP uses the pairwise patristic distance distribution of the entire tree to inform
586 its limit on within-cluster internal divergence against the background genetic diversity of the population
587 included in the phylogeny. Testing against the global background genetic diversity indicates whether
588 the putative clustered sequences are sufficiently more related to one another than to the rest of the
589 dataset to be designated a distinct cluster.

590 PhyCLIP's cluster assignment is agnostic to metadata but is capable of capturing the geographic and
591 temporal structure of the H5 phylogeny informatively. PhyCLIP recovers the overall structure of the
592 current WHO/OIE/FAOH5 nomenclature developed on a sequence divergence threshold but delineates

593 more informative, higher resolution clusters that capture geographically-distinct subpopulations.
594 PhyCLIP therefore plausibly provides the foundation for an alternative nomenclature that minimises the
595 limitations of currently employed approaches.

596 PhyCLIP's clustering is expected to improve with the addition of new sequences to the tree as new
597 information about the genetic diversity and evolutionary trajectory of the pathogen becomes known and
598 can be incorporated into the background diversity of the tree that informs the algorithm. Additionally,
599 topological information that captures how sequences are related by common ancestors is inherently
600 incorporated in PhyCLIP owing to its distal dissociation approach. The distal dissociation approach also
601 does not assume all clusters are monophyletic as the most recent common ancestor of all tips in a
602 cluster is not assumed to have any descendants. As such, PhyCLIP can identify nested clusters both
603 as clusters with sufficiently high information content to meet the statistical requirements of cluster
604 designation or sufficiently diverse clusters that are dissociated from their ancestral nodes. The
605 designation of divergent descendant clusters nested within a supercluster suggestively captures source-
606 sink population dynamics that may be informative about the evolutionary trajectory of the clustered
607 sequences. At the same time, users could also opt for PhyCLIP to subsume subclusters that do not
608 violate the statistical criteria of the parent clusters into the latter, aiding higher level interpretation.
609 Importantly, the distal dissociation approach also identifies highly divergent outlying sequences that may
610 be indicative of under-sampled diversity.

611 For pathogens that evolve more rapidly than they spread geographically, it is expected that clusters of
612 related sequences would be temporally structured. However, it is important to consider the distribution
613 of sampling times, which can drive clustering artificially. This is especially pertinent for transmission
614 dynamic studies, where clustering is often driven by heterogeneity in sampling rates across
615 subpopulations rather than heterogeneity in transmission rates (Poon 2016; McCloskey and Poon
616 2017). PhyCLIP can be applied to time-resolved phylogenies in heterochronous datasets. However,
617 molecular clock analyses make strong biological assumptions and require sufficient temporal signal to
618 inform the model reconstructing the statistical relationship between genetic divergence and time
619 (Rambaut et al. 2016). These models rely on high-quality sampling dates and alignments free of
620 sequence error and laboratory-altered strains or recombinant viruses to reconstruct valid and unbiased
621 time-scaled phylogenies (Rambaut et al. 2016). As PhyCLIP centrally operates on the branch lengths
622 of the phylogeny, we recommend it is only applied to robust time-resolved phylogenies after a thorough
623 investigation of the temporal signal as well as a rigorous assessment of model and prior assumptions
624 (Boskova et al. 2018).

625 PhyCLIP's methodology has limitations. Notably, PhyCLIP is tree-based and is therefore subject to error
626 in phylogenetic reconstruction. PhyCLIP does not include criteria for the statistical support of nodes
627 under consideration, which omits uncertainty in phylogenetic reconstruction. However, high statistical

628 support for a node does not necessarily indicate that all sequences subtended by it are highly related
629 but merely reflects the statistical support of the bipartition to the exclusion of other sequences.
630 Additionally, the relationship between the statistical significance of internal nodes and population
631 dynamics is unresolved as is an appropriate definition of a robustly supported node (Zharkikh and Li
632 1992; Susko 2009; Anisimova et al. 2011; Kumar et al. 2012; Volz et al. 2012). There is often less
633 phylogenetic signal to resolve internal nodes subtending small subtrees in measurably evolving
634 populations, increasing uncertainty in the arrangement of the internal structure of smaller subtrees. If a
635 statistical support threshold is set for nodes, these viruses will consistently be left unclustered or will be
636 forced to coalesce with more ancestral nodes subtending larger clusters, which would violate PhyCLIP's
637 statistical framework.

638 As with any phylogenetic clustering methods, PhyCLIP is also sensitive to variation in sampling rates
639 (Volz et al. 2012). There is a significant surveillance bias towards certain pathogens (e.g. HPAI H5)
640 owing to their consequences for animal and human health. The evolution and divergence of these
641 pathogens are currently captured in surveillance data as a more accurate approximation to a continuum
642 of evolution. PhyCLIP's clustering is strongly influenced by the diversity in the input population it tests
643 against and will perform best when the background diversity of the phylogeny is complete or
644 representative.

645 Clusters identified by PhyCLIP should not be interpreted as sequences linked by rapid direct
646 transmission events. Transmission dynamic studies aim to integrate epidemiological clustering with
647 phylogenetic clusters to study transmission chains or local outbreak networks by assuming putative
648 transmission links between highly related sequences (Hassan et al. 2017). Datasets from transmission
649 dynamic studies are likely to be sampled from localised outbreaks over a very specific period of time.
650 The global distribution generated from the resulting phylogenetic trees will not contain sufficient
651 information or power to meaningfully compare subpopulations to identify high confidence transmission
652 clusters.

653 In conclusion, PhyCLIP provides an automated, statistically-principled framework for phylogenetic
654 clustering that can be generalised to research questions concerning the identification of biologically
655 informative clusters in pathogen phylogenies.

656

657 **Materials and methods**

658 **Robust estimator of scale (deviation)**

659 PhyCLIP computes the robust estimator of scale (σ) either as the median absolute deviation (*MAD*) or
660 Q_n . Note that *MAD* may not suitably account for any potential skewness of the pairwise sequence

661 patristic distance distribution as it inherently assumes symmetry about the median ($\bar{\mu}$). On the contrary,
662 Q_n , an alternative estimator of scale proposed by Rousseeuw & Croux (1993), is as robust as MAD (i.e.
663 50% breakdown point), calculated solely using the differences between the values in the distribution
664 without needing a location estimate, and has been proven to be statistically more efficient in both
665 Gaussian and non-Gaussian distributions relative to MAD .

666

667 Integer linear programming model

668 Here, we fully elaborate the ILP model underlying PhyCLIP. Let $n_1, n_2, \dots, n_i, \dots, n_N$ be the set of binary
669 variables indicating if subtree i satisfies the conditions for clustering as a clade ($n_i = 1$ if it does and
670 $n_i = 0$ vice versa, Figure 2C). Each sequence j subtended by subtree i is also assigned a binary variable
671 $l_{j,i}$ indicating if the sequence is clustered under subtree i ($l_{j,i} = 1$ if j is clustered under node i and $l_{j,i} =$
672 0 vice versa, Figure 2C). PhyCLIP then formulates the phylogenetic clustering problem as an integer
673 linear programming (ILP) model with the objective to maximise the number of sequences assigned with
674 cluster membership:

$$\max \sum_{j,i} l_{j,i} \quad (2)$$

675

676 subject to the following constraints:

677

$$l_{j,i} \leq n_i \quad \forall j \in L_i, i \quad (3)$$

678 Constraint (3) stipulates that sequence j can be clustered under subtree i if and only if subtree i is a
679 potential clade ($n_i = 1$).

680

$$l_{j,i} \leq 2 - n_i - n_k \quad \forall j \in \{L_i, L_k\}, k; i < k \quad (4)$$

681 If sequence j is subtended by subtrees i and k , wherein i is ancestral to k and both nodes are potential
682 clusters ($n_i = n_k = 1$), constraints (3) and (4) stipulate sequence j will not be clustered under the
683 ancestor node i . Implementing these constraints across all pairwise combinations of subtrees
684 subtending sequence j in turn constrains j to be clustered under the most descendant node k possible.

685

$$\sum_i l_{j,i} \leq 1 \quad \forall j \quad (5)$$

686 Constraint (5) stipulates that each sequence can only be clustered under a single subtree, hence
687 abrogating any fuzzy clustering.

688

$$C(n_i - 1) \leq \sum_j l_{j,i} - S \quad \forall i \quad (6)$$

689 where C is any arbitrarily large positive constant. Constraint (6) requires all clusters to contain at least S
690 number of taxa as defined by the user (Figures 1B and C).

691

$$C(n_i - 1) \leq WCL - \mu_i \quad \forall i \quad (7)$$

692 Constraint (7) ensures that μ_i of all clades fall below the stipulated WCL limit.

693

$$C(2 - n_i - n_k) \geq q_{i,k} - FDR \quad \forall i, k \neq i \quad (8)$$

694 where $q_{i,k}$ is the Benjamini-Hochberg corrected p -value testing if subtrees i and k are significantly
695 divergent from one and another under the user-defined significance level, FDR . Constraint (8) is the
696 inter-cluster divergence constraint. Inter-cluster divergence between subtrees i and k is tested under
697 the null hypothesis that the pairwise sequence distance distributions of i and k are empirically equivalent
698 to that if the two subtrees were clustered together. This can be done either by the putative Kolmogorov-
699 Smirnov (KS) test or Kuiper's test.

700 Although both tests are nonparametric, the Kuiper's test statistic incorporates both the greatest positive
701 and negative deviations between the two distributions whereas the KS test statistic is defined only by
702 their maximum difference. As a result, the Kuiper's test becomes equally sensitive to differences to the
703 tails as well as the median of the distributions but the KS test works best when the distributions differ
704 mostly at the median. In other words, the KS test is good at detecting *shifts* between the distributions
705 but lacks the sensitivity to uncover *spreads* between the distributions characterised by changes in their
706 tails. Kuiper's test is, however, sensitive to detect both types of changes in distributions.

707 There are two scenarios under which $q_{i,k}$ may be calculated:

- 708 (i) Subtree i is ancestral to k . The hypothesis test assumes the null hypothesis that the pairwise
709 sequence patristic distance distribution of subtree k is statistically identical to the pairwise
710 sequence patristic distance distribution of its ancestor i .

711 (ii) Neither subtree i nor k is an ancestor of the other. In this case, two hypothesis tests are carried
712 out comparing the distribution of each subtree to the distribution of pairwise sequence patristic
713 distance should both subtrees be combined as a single cluster and we take the more
714 conservative $q_{i,k} = \max\{q_{i,combined}, q_{k,combined}\}$.

715

716 **Nomenclature**

717 Traversing the output clusters of PhyCLIP by pre-order of the input phylogeny, a unique number is
718 assigned to any cluster with no immediate ancestral supercluster precursor to it (i.e. parent node of the
719 cluster node is not part of any PhyCLIP clusters). Otherwise, the descendant cluster in question is
720 designated as a *child cluster* should its membership size be $>25^{\text{th}}$ percentile of PhyCLIP's output cluster
721 size distribution (i.e. for having proliferated in numbers substantial enough to be deemed a progeny
722 cluster). Every child cluster of a supercluster is assigned a progeny number separated by a decimal
723 point (e.g. 1.2 refers to the second child cluster of supercluster 1). On other hand, descendant clusters
724 that fall below the cluster size cut-off are distinguished from child clusters as *nested clusters*, each
725 assigned an address in the form of a parenthesized letter, alphabetised by tree traversal order, prefixed
726 by its parent supercluster nomenclature (e.g. 1.1(c) refers to the third nested cluster of supercluster 1.1).
727 Nested clusters in superclusters fundamentally have different properties from the sensitivity-induced
728 nested clusters discussed in New Approach section and cannot be subsumed as it will violate the within-
729 cluster limit of the parent supercluster. The structure of the resultant clustering topology is highlighted
730 in Figure 3.

731

732 **Phylogenetic analyses**

733 PhyCLIP's performance was evaluated on an empirical dataset. The sequence datasets used to
734 construct the haemagglutinin (HA) gene phylogenetic trees underlying the WHO/OIE/FAO nomenclature
735 for the A/goose/Guangdong /1/1996 (Gs/GD/96)-like H5 avian influenza viruses were downloaded from
736 GISAID (Anon 2008; WHO/OIE/FAO H5N1 Evolution Working Group 2012; WHO/OIE/FAO H5N1
737 Evolution Working Group 2014; Smith, Donis, and WHO/OIE/FAO H5 Evolution Working Group 2015).
738 The primary analysis is based on the full dataset included in the 2009 (n=1224) and 2015 (n=4357)
739 nomenclature updates. Viruses that were inconsistently included across WHO/OIE/FAO updates were
740 followed up and included (WHO/OIE/FAO HN Evolution Working Gr 2009; Smith, Donis,
741 and WHO/OIE/FAO H5 Evolution Working Group 2015). Sequences were curated based on criteria
742 defined by the H5 nomenclature: sequences with more than 5 ambiguous nucleotides, with a sequence
743 length shorter than 60% of the alignment, or with frameshifts or duplicated by name were removed. For

744 the 2018 phylogeny, all avian and human viruses from the Gs/GD-like H5 lineage were downloaded
745 from GISAID up to April 2018, including H5Nx subtypes H5N2, H5N3, H5N5, H5N6 and H5N8. An
746 alternative filtering approach compared to the published WHO nomenclature approach was applied to
747 ensure a dataset of high-quality sequences that would be robust to error in phylogenetic reconstruction
748 as PhyCLIP is inherently sensitive to topological information. In this approach, duplicate sequences and
749 sequences with a length below 95% of the full HA sequence or more than 1% ambiguous nucleotides
750 were discarded. Sequences were aligned with MAFFT v7.397 and trimmed to the start of the mature
751 protein (Kato et al. 2002). Each sequence set was annotated with the WHO/OIE/FAOH5 nomenclature
752 using LABEL(v0.5.2), and the version of the module corresponding to the nomenclature update of the
753 dataset (e.g. H5v2015 module for the full tree from the nomenclature update in 2015) (Shepard et al.
754 2014). Maximum likelihood phylogenetic trees were constructed for each dataset with RAxML 8.2.12
755 under the GTR+GAMMA substitution model, and rooted to Gs/GD/96 (Stamatakis 2014). Phylogenetic
756 trees were visualised using Figtree (<http://tree.bio.ed.ac.uk/software/figtree/>) and ggtree (Yu et al. 2017).

757

758 **Silhouette index**

759 The silhouette index is based on the distance, here patristic distance, of each cluster member to other
760 cluster members compared to the distance to its nearest neighbours (Rousseeuw 1987). Silhouette
761 values approaching one indicate that the cluster member is correctly assigned, whereas values close to
762 zero indicate that the sequence is equally matched to its neighbouring cluster. A negative Silhouette
763 index indicates that the sequence is more closely related to the neighbouring cluster than to its fellow
764 cluster members. Calculation of the silhouette index was performed in R (R Core Team 2016).

765

766 **Code availability**

767 PhyCLIP is freely available on github (<http://github.com/alvinxhan/PhyCLIP>) and documentation can be
768 found on the associated wiki page (<http://github.com/alvinxhan/PhyCLIP/wiki>).

769

770 **Acknowledgments**

771 We thank the GISAID Initiative and the influenza surveillance and research groups that openly shared
772 the genetic sequence data that made this work possible (full acknowledgement table is available as
773 supplementary). A.X.H. was supported by the A*STAR Graduate Scholarship programme from A*STAR
774 to carry out his PhD work via collaboration between Bioinformatics Institute (A*STAR) and NUS
775 Graduate School for Integrative Sciences and Engineering from the National University of Singapore.

776 E.P. was funded by the Gates Cambridge Trust (Grant number OPP1144). S.M.S. was supported by
777 the A*STAR HEIDI programme (Grant number: H1699f0013) and Bioinformatics Institute (A*STAR).
778 C.A.R. was supported by University Research Fellowship from the Royal Society.

779

780 **References**

781 Aldous JL, Pond SK, Poon A, Jain S, Qin H, Kahn JS, Kitahata M, Rodriguez B, Dennis AM, Boswell
782 SL, et al. 2012. Characterizing HIV Transmission Networks Across the United States. *Clin. Infect. Dis.*
783 55:1135–1143.

784 Anisimova M, Gil M, Dufayard J-F, Dessimoz C, Gascuel O. 2011. Survey of branch support methods
785 demonstrates accuracy, power, and robustness of fast likelihood-based approximation schemes. *Syst.*
786 *Biol.* 60:685–699.

787 Boskova V, Stadler T, Magnus C. 2018. The influence of phylodynamic model specifications on
788 parameter estimates of the Zika virus epidemic. *Virus Evol.* 4:1–14.

789 Burk RD, Chen Z, Harari A, Smith BC, Kocjan BJ, Maver PJ, Poljak M. 2011. Classification and
790 nomenclature system for human Alphapapillomavirus variants: general features, nucleotide landmarks
791 and assignment of HPV6 and HPV11 isolates to variant lineages. *Acta dermatovenerologica Alpina,*
792 *Pannonica, Adriat.* 20:113–123.

793 Dennis AM, Herbeck JT, Brown AL, Kellam P, de Oliveira T, Pillay D, Fraser C, Cohen MS. 2014.
794 Phylogenetic studies of transmission dynamics in generalized HIV epidemics: an essential tool where
795 the burden is greatest? *J. Acquir. Immune Defic. Syndr.* 67:181–195.

796 Donald BS, Jens B, Carla K, Scott MA, M. RC, Simmonds P, T. SJ. 2013. Expanded classification of
797 hepatitis C virus into 7 genotypes and 67 subtypes: updated criteria and assignment web resource.
798 *Hepatology* 59:318–327.

799 Van Doorslaer K, Bernard H-U, Chen Z, de Villiers E-M, zur Hausen H, Burk RD. 2011.
800 Papillomaviruses: evolution, Linnaean taxonomy and current nomenclature. *Trends Microbiol.* 19:49-
801 50; author reply 50-1.

802 Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG. 2003. Measurably evolving
803 populations. *Trends Ecol. Evol.* 18:481–488.

804 Duan L, Bahl J, Smith GJD, Wang J, Vijaykrishna D, Zhang LJ, Zhang JX, Li KS, Fan XH, Cheung CL,
805 et al. 2008. The development and genetic diversity of H5N1 influenza virus in China, 1996-2006.
806 *Virology* 380:243–254.

807 Gardy JL, Loman NJ. 2017. Towards a genomics-informed, real-time, global pathogen surveillance

- 808 system. *Nat. Rev. Genet.* 19:9–20.
- 809 Grabowski MK, Herbeck JT, Poon AFY. 1904. Genetic Cluster Analysis for HIV Prevention.
- 810 Hassan AS, Pybus OG, Sanders EJ, Albert J, Esbjörnsson J. 2017. Defining HIV-1 transmission
811 clusters based on sequence data. *AIDS* 31:1211–1222.
- 812 Hué S, Clewley JP, Cane PA, Pillay D. 2004. HIV-1 pol gene variation is sufficient for reconstruction of
813 transmissions in the era of antiretroviral therapy. *AIDS* 18:719–728.
- 814 Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple sequence
815 alignment based on fast Fourier transform. *Nucleic Acids Res.* 30:3059–3066.
- 816 Kroneman A, Vega E, Vennema H, Vinjé J, White PA, Hansman G, Green K, Martella V, Katayama K,
817 Koopmans M. 2013. Proposal for a unified norovirus nomenclature and genotyping. *Arch. Virol.*
818 158:2059–2068.
- 819 Kumar S, Filipski AJ, Battistuzzi FU, Kosakovsky Pond SL, Tamura K. 2012. Statistics and Truth in
820 Phylogenomics. *Mol. Biol. Evol.* 29:457–472.
- 821 Lauber C, Gorbalenya AE. 2012. Toward Genetics-Based Virus Taxonomy: Comparative Analysis of a
822 Genetics-Based Classification and the Taxonomy of Picornaviruses. *J. Virol.* 86:3905–3915.
- 823 McCloskey RM, Poon AFY. 2017. A model-based clustering method to detect infectious disease
824 transmission outbreaks from sequence variation. Kosakovsky Pond SL, editor. *PLOS Comput. Biol.*
825 13:e1005868.
- 826 McIntyre CL, Knowles NJ, Simmonds P. 2013. Proposals for the classification of human rhinovirus
827 species A, B and C into genotypically assigned types. *J. Gen. Virol.* 94:1791–1806.
- 828 Meilă M. 2007. Comparing clusterings—an information based distance. *J. Multivar. Anal.* 98:873–895.
- 829 Ortiz JR, Neuzil KM. 2017. Influenza immunization of pregnant women in resource-constrained
830 countries: an update for funding and implementation decisions. *Curr. Opin. Infect. Dis.* 30:455–462.
- 831 Poon AFY. 2016. Impacts and shortcomings of genetic clustering methods for infectious disease
832 outbreaks. *Virus Evol.* 2:vew031.
- 833 Poon AFY, Gustafson R, Daly P, Zerr L, Demlow SE, Wong J, Woods CK, Hogg RS, Kraiden M,
834 Moore D, et al. 2016. Near real-time monitoring of HIV transmission hotspots from routine HIV
835 genotyping: an implementation case study. *Lancet HIV* 3:e231–e238.
- 836 Poon AFY, Joy JB, Woods CK, Shurgold S, Colley G, Brumme CJ, Hogg RS, Montaner JSG, Harrigan
837 PR. 2015. The Impact of Clinical, Demographic and Risk Factors on Rates of HIV Transmission: A
838 Population-based Phylogenetic Analysis in British Columbia, Canada. *J. Infect. Dis.* 211:926–935.

- 839 Prosperi MCF, Ciccozzi M, Fanti I, Saladini F, Pecorari M, Borghi V, Di Giambenedetto S, Bruzzone B,
840 Capetti A, Vivarelli A, et al. 2011. A novel methodology for large-scale phylogeny partition. *Nat.*
841 *Commun.* 2:321.
- 842 Prosperi MCF, De Luca A, Di Giambenedetto S, Bracciale L, Fabbiani M, Cauda R, Salemi M. 2010.
843 The Threshold Bootstrap Clustering: A New Approach to Find Families or Transmission Clusters
844 within Molecular Quasispecies. Poon AFY, editor. *PLoS One* 5:e13619.
- 845 Pu J, Wang S, Yin Y, Zhang G, Carter RA, Wang J, Xu G, Sun H, Wang M, Wen C, et al. 2015.
846 Evolution of the H9N2 influenza genotype that facilitated the genesis of the novel H7N9 virus. *Proc.*
847 *Natl. Acad. Sci. U. S. A.* 112:548–553.
- 848 R Core Team. 2016. R: A Language and Environment for Statistical Computing. Available from:
849 <https://www.r-project.org/>
- 850 Ragonnet-Cronin M, Hodcroft E, Hué S, Fearnhill E, Delpech V, Brown AJ, Lycett S, Holmes E, Nee
851 S, Rambaut A, et al. 2013. Automated analysis of phylogenetic clusters. *BMC Bioinformatics* 14:317.
- 852 Rambaut A, Lam TT, Max Carvalho L, Pybus OG. 2016. Exploring the temporal structure of
853 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol.* 2:vew007.
- 854 Rose R, Lamers SL, Dollar JJ, Grabowski MK, Hodcroft EB, Ragonnet-Cronin M, Wertheim JO, Redd
855 AD, German D, Laeyendecker O. 2017. Identifying Transmission Clusters with Cluster Picker and HIV-
856 TRACE. *AIDS Res. Hum. Retroviruses* 33:211–218.
- 857 Rousseeuw PJ. 1987. Silhouettes: A graphical aid to the interpretation and validation of cluster
858 analysis. *J. Comput. Appl. Math.* 20:53–65.
- 859 Rousseeuw PJ, Croux C. 1993. Alternatives to the Median Absolute Deviation. *J. Am. Stat. Assoc.*
860 88:1273–1283.
- 861 Shepard SS, Davis CT, Bahl J, Rivaille P, York IA, Donis RO. 2014. LABEL: Fast and Accurate
862 Lineage Assignment with Assessment of H5N1 and H9N2 Influenza A Hemagglutinins. Woo PCY,
863 editor. *PLoS One* 9:e86921.
- 864 Simmonds P, McIntyre C, Savolainen-Kopra C, Tapparel C, Mackay IM, Hovi T. 2010. Proposals for
865 the classification of human rhinovirus species C into genotypically assigned types. *J. Gen. Virol.*
866 91:2409–2419.
- 867 Smith GJD, Donis RO, World Health Organization/World Organisation for Animal Health/Food and
868 Agriculture Organization (WHO/OIE/FAO) H5 Evolution Working Group WHOO for AH and AO
869 (WHO/OIE/FAO) HEW. 2015. Nomenclature updates resulting from the evolution of avian influenza
870 A(H5) virus clades 2.1.3.2a, 2.2.1, and 2.3.4 during 2013-2014. *Influenza Other Respi. Viruses* 9:271–

- 871 276.
- 872 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
873 phylogenies. *Bioinformatics* 30:1312–1313.
- 874 Susko E. 2009. Bootstrap Support Is Not First-Order Correct. *Syst. Biol.* 58:211–223.
- 875 The Global Consortium for H5N8 and Related Influenza Viruses. 2016. Role for migratory wild birds in
876 the global spread of avian influenza H5N8. *Science* (80-.). 354:213–217.
- 877 Valastro V, Holmes EC, Britton P, Fusaro A, Jackwood MW, Cattoli G, Monne I. 2016. S1 gene-based
878 phylogeny of infectious bronchitis virus: An attempt to harmonize virus classification. *Infect. Genet.*
879 *Evol.* 39:349–364.
- 880 Volz EM, Koopman JS, Ward MJ, Brown AL, Frost SDW. 2012. Simple Epidemiological Dynamics
881 Explain Phylogenetic Clustering of HIV from Patients with Recent Infection. Fraser C, editor. *PLoS*
882 *Comput. Biol.* 8:e1002552.
- 883 Wang J, Vijaykrishna D, Duan L, Bahl J, Zhang JX, Webster RG, Peiris JSM, Chen H, Smith GJD,
884 Guan Y. 2008. Identification of the Progenitors of Indonesian and Vietnamese Avian Influenza A
885 (H5N1) Viruses from Southern China. *J. Virol.* 82:3405–3414.
- 886 WHO/OIE/FAO H5N1 Evolution Working Group. 2008. Toward a Unified Nomenclature System for
887 Highly Pathogenic Avian Influenza Virus (H5N1). *Emerg. Infect. Dis.* 14:e1–e1.
- 888 WHO/OIE/FAO H5N1 Evolution Working Group WHEW. 2012. Continued evolution of highly
889 pathogenic avian influenza A (H5N1): updated nomenclature. *Influenza Other Respi. Viruses* 6:1–5.
- 890 WHO/OIE/FAO HN Evolution Working Gr. 2009. Continuing progress towards a unified nomenclature
891 for the highly pathogenic H5N1 avian influenza viruses: divergence of clade 2?2 viruses. *Influenza*
892 *Other Respi. Viruses* 3:59–62.
- 893 World Health Organization/World Organisation for Animal Health/Food and Agriculture Organization
894 (WHO/OIE/FAO) H5N1 Evolution Working Group. 2014. Revised and updated nomenclature for highly
895 pathogenic avian influenza A (H5N1) viruses. *Influenza Other Respi. Viruses* 8:384–388.
- 896 Yu G, Smith DK, Zhu H, Guan Y, Lam TTY. 2017. Ggtree: an R Package for Visualization and
897 Annotation of Phylogenetic Trees With Their Covariates and Other Associated Data. *Methods Ecol.*
898 *Evol.* 8:28–36.
- 899 Zharkikh A, Li WH. 1992. Statistical properties of bootstrap estimation of phylogenetic variability from
900 nucleotide sequences. I. Four taxa with a molecular clock. *Mol. Biol. Evol.* 9:1119–1147.
- 901