

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

Phylogenetic Clustering by Linear Integer Programming (PhyCLIP)

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

[†]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)

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

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

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

41

42 **Introduction**

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

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

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

61 In virology, nomenclature systems are largely reliant on absolute distance thresholds that define the
62 maximum genetic divergence tolerated between viruses designated as closely related (Smith et al.; Burk
63 et al. 2011; Van Doorslaer et al. 2011; Lauber and Gorbalenya 2012; Kroneman et al. 2013; Poon et al.
64 2015; Smith, Donis, and WHO/OIE/FAO H5 Evolution Working Group 2015; Poon et al. 2016; Valastro
65 et al. 2016). Groups of closely related viruses are inferred to be phylogenetic clusters when the genetic
66 distance between them is lower than the limit set on within-cluster divergence. Non-parametric distance-
67 based clustering approaches have defined the distance between sequences using pairwise genetic
68 distances calculated directly from sequence data (WHO/OIE/FAO H5N1 Evolution Working Group 2008;
69 Aldous et al. 2012; Ragonnet-Cronin et al. 2013) or pairwise patristic distances calculated from inferred
70 phylogenetic trees (Hu   et al. 2004; Prosperi et al. 2011; Poon et al. 2015; Pu et al. 2015; Ortiz and
71 Neuzil 2017). Within-cluster limits on tolerated divergence have been set using mean (WHO/OIE/FAO
72 H5N1 Evolution Working Group 2008), median (Prosperi et al. 2011) or maximum within-cluster pairwise
73 genetic or patristic distance (Ragonnet-Cronin et al. 2013). Some methods incorporate additional
74 criteria, such as the statistical support for subtrees under consideration or minimum/maximum cluster
75 size (Hu   et al. 2004; Prosperi et al. 2010; Prosperi et al. 2011; Ragonnet-Cronin et al. 2013). These
76 genetic distance-based clustering approaches are convenient, as they are rule-based and scalable,
77 allowing for relatively easy nomenclature updates. Arguably, flexibility in the distance thresholds allows
78 researchers to curate clusters based on consistency of the geographic or temporal metadata.

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

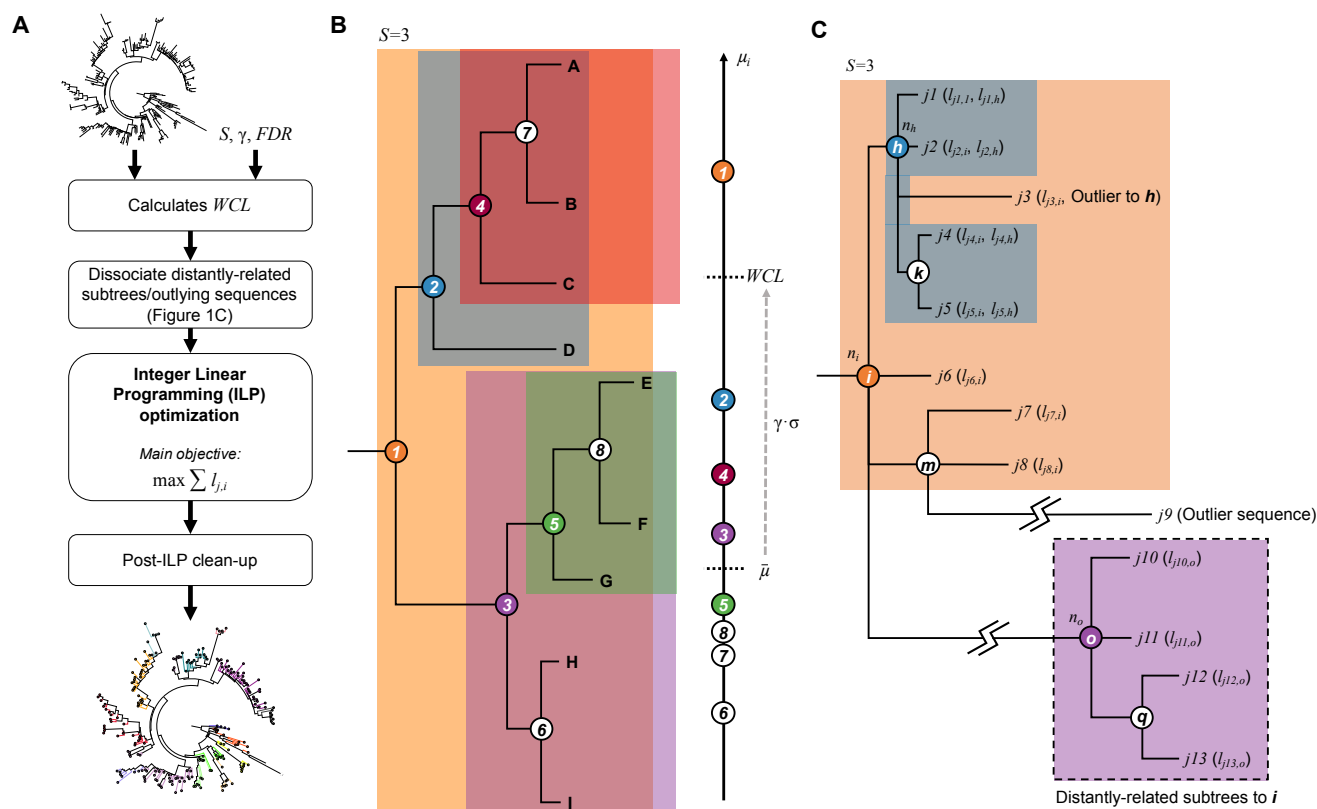
85 into meaningful subsets is therefore complex and is mostly performed *ad hoc* through exploratory
86 analyses with uninformative sensitivity analyses across thresholds. As expected, cluster membership is
87 highly sensitive to the threshold applied and therefore results can be unstable across different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

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

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

100

101 New approach



102

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

116

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

- 118 (i) Minimum number of sequences (S) that should be considered a cluster.
 119 (ii) Multiple of deviations (γ) from the grand median of the mean pairwise sequence patristic distance
 120 that defines the within-cluster divergence limit (WCL).

121 (iii) False discovery rate (*FDR*) to infer that the diversity observed for every combinatorial pair of
122 output clusters is significantly distinct from one another.

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

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

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

134 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
135 σ is any robust estimator of scale (e.g. median absolute deviation (*MAD*) or Q_n , see Methods) that
136 quantifies the statistical dispersion of the mean pairwise patristic distance distribution for the ensemble
137 of N subtrees. In other words, only subtrees with $\mu_i \leq WCL$ will be considered for clustering by PhyCLIP
138 (Figure 1B).

139

140 **Distal dissociation**

141 The assumption that a cluster must be monophyletic can lead to incorrect assignment of cluster
142 membership to undersampled, distantly related outlying sequences that have diverged considerably
143 from the rest of the cluster (e.g. sequence j_9 in Figure 1C). These exceedingly distant outlying
144 sequences can also skew μ_i of the subtree they subtend, leading to inaccurate overestimation of the
145 internal divergence of the putative subtree. Similarly, distantly related descendant subtrees can
146 artificially inflate μ_i of their ancestral trunk nodes (e.g. nodes o and q in Figure 1C). In turn, historical
147 sequences immediately descending from a trunk node i will never be clustered if its μ_i exceeds *WCL*
148 (Figure 1C).

149 PhyCLIP dissociates any distal subtrees and/or outlying sequences from their ancestral lineage prior to
150 implementing the integer linear programming (ILP) model. For any subtree i with $\mu_i > WCL$, starting
151 from the most distant sequence to i , PhyCLIP applies a leave-one-out strategy dissociating sequences,

152 and the whole descendant subtree if every sequence subtended by it was dissociated, until the
153 recalculated μ_i without the distantly related sequences falls below WCL . For each subtree, PhyCLIP
154 also tests and dissociates any outlying sequences present. An outlying sequence is defined as any
155 sequence whose patristic distance to the node in question is $> 3 \times$ the estimator of scale away from the
156 median sequence patristic distance to node. μ_i is recalculated for any node with changes to its sequence
157 membership L_i after dissociating these distantly related sequences. These distal dissociation steps
158 effectively offer PhyCLIP greater flexibility in its clustering construct allowing the identification of
159 paraphyletic clusters on top of monophyletic ones that may better reflect the phylogenetic relationships
160 of these sequences.

161

162 **Integer linear programming optimisation**

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

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

178 The ILP model is implemented in Gurobi (<http://www.gurobi.com/>), a third-party commercial linear
179 programming solver fully integrated within PhyCLIP, to obtain the global optimal solution. At the time of
180 this publication, Gurobi is one of the fastest available mathematical programming solvers (2018
181 benchmark tests of popular linear programming solvers by Hans Mittelmann,
182 <http://plato.asu.edu/ftp/lpsimp.html>). Full-featured academic licenses of Gurobi are available for free to
183 users based at any academic institution.

184

185 **Post-ILP clean-up**

186 While distal dissociation prior to ILP optimisation works well for dissociating distantly related subtrees
187 and sequences, it is ineffective in identifying spurious singletons such as sequence $j3$ in Figure 1C.
188 Here, in terms of sequence patristic distance, sequence $j3$ is an outlying sequence to the descendant
189 node h but not so to the ancestral node i . If taxa subtended by subtree h (i.e. $j1$, $j2$, $j4$ and $j5$) were to
190 be clustered without $j3$ which itself is clustered under cluster i , PhyCLIP performs a post-ILP
191 optimisation clean-up step that removes $j3$ from output cluster i . This is because $j3$ is clearly a
192 topologically outlying taxon to i and if unremoved, would also suggest fuzzy clustering for the sequences
193 clustered under cluster h .

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

207

208

209 **Optimisation criteria**

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

219 PhyCLIP's generated optimisation criteria is agnostic to the metadata of the dataset and includes: 1)
220 The grand mean of the pairwise patristic distance distribution and its standard deviation. The grand
221 mean is a measure of the within-cluster divergence and can be optimised to select a clustering
222 configuration with the lowest global internal divergence. 2) The mean of the inter-cluster distance to all
223 other clusters and its standard deviation. This can be optimised to select a clustering configuration with
224 well-separated clusters. 3) The percentage of sequences clustered, which can be optimised to minimise
225 the number of unclustered sequences. 4) The total number of clusters and 5) mean or median cluster
226 size, which can be optimised to select a tolerable level of stratification of the tree.

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

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

238

239 Results

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

252 The percentage of sequences clustered was prioritised as the primary optimisation criterion to ensure
253 that the maximum number of sequences were assigned a nomenclature identifier. Mean within-cluster
254 geographic distance was included as a *post-hoc* optimisation criterion as many avian influenza viruses
255 cluster with high spatiotemporal consistency owing to their transmission dynamics in localised avian
256 populations. For influenza viruses endemic to poultry such as H5Nx, this is likely owing to increased
257 local transmission during outbreaks in large poultry populations, as well as the associated sampling bias
258 (Smith, Donis, and WHO/OIE/FAO H5 Evolution Working Group 2015). Within-cluster genetic
259 divergence was optimised with higher priority than within-cluster mean geographic distance, as the use
260 of phylogenetic geographic structure as a ground truth for avian influenza viruses is restricted by the
261 long-distance dissemination of related viruses through mechanisms such as the poultry trade or
262 migration of wild birds (WHO/OIE/FAO H5N1 Evolution Working Group 2014; Smith, Donis, and
263 WHO/OIE/FAO H5 Evolution Working Group 2015). The within-cluster geographic distance was
264 calculated for each cluster in each clustering result as the mean within-cluster pairwise Vicenty distance
265 in miles.

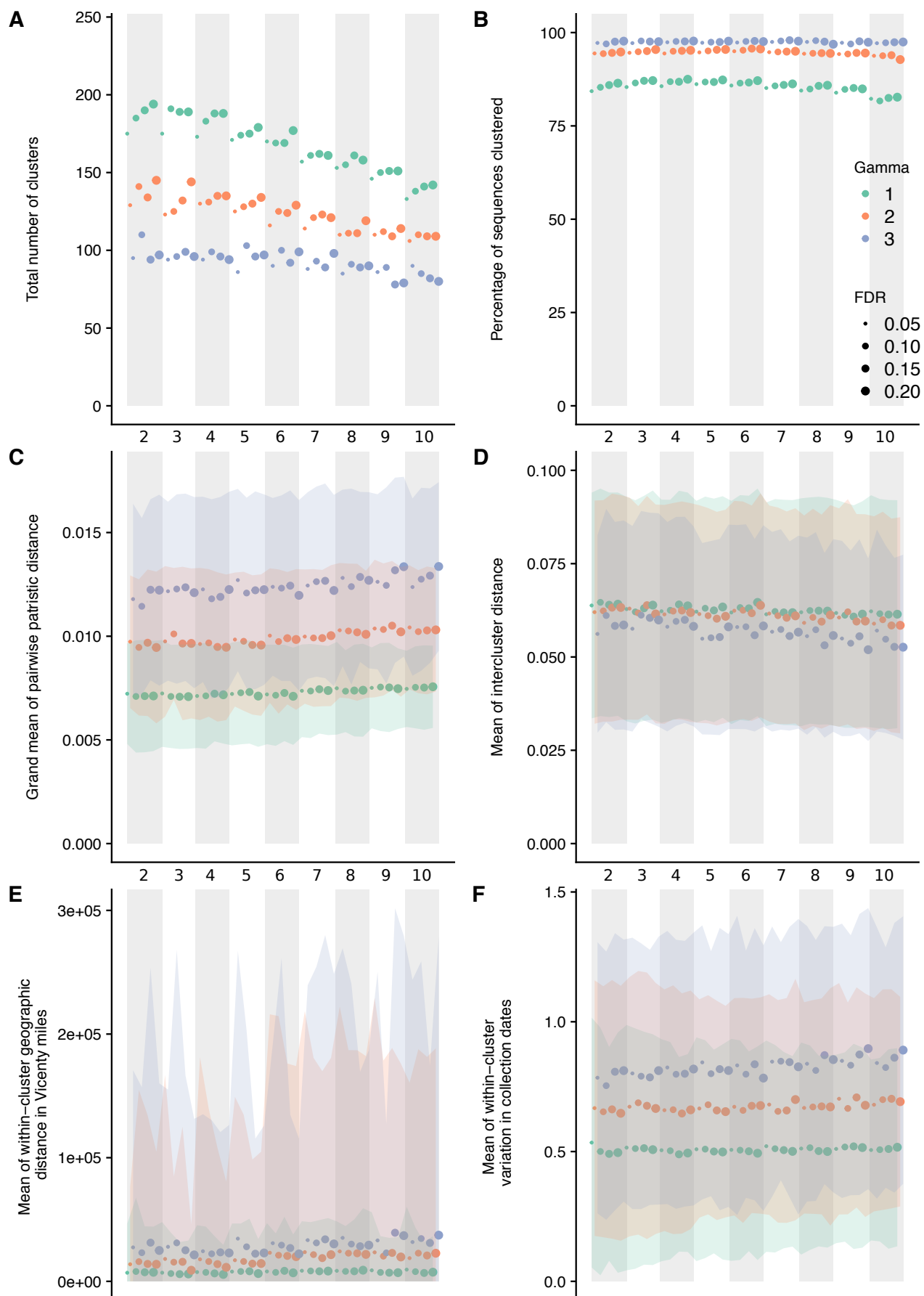
266 As PhyCLIP incorporates topological information of the phylogeny into the clustering construct, non-
267 terminal internal nodes with zero branch lengths can lead to erroneous clustering and over-delineation
268 (Figure S1). Such internal nodes are usually found in bifurcating trees as representations of polytomies,
269 arising from a lack of phylogenetic signal among the sequences subtended by the node to resolve them
270 into dichotomies. As such, prior to implementing PhyCLIP, all non-terminal, zero branch length nodes
271 in the input phylogenetic trees were collapsed into polytomies, which more accurately depicts the
272 relationship between identical/indiscernible sequences and/or ancestral states. In the H5Nx analysis, all

273 subclusters were subsumed if the statistical requisites of the parent clade were maintained, to aid in
274 easing the interpretation of the nomenclature designation (as discussed in the New Approach section).

275

276 **Influence of the parameters**

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



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

304

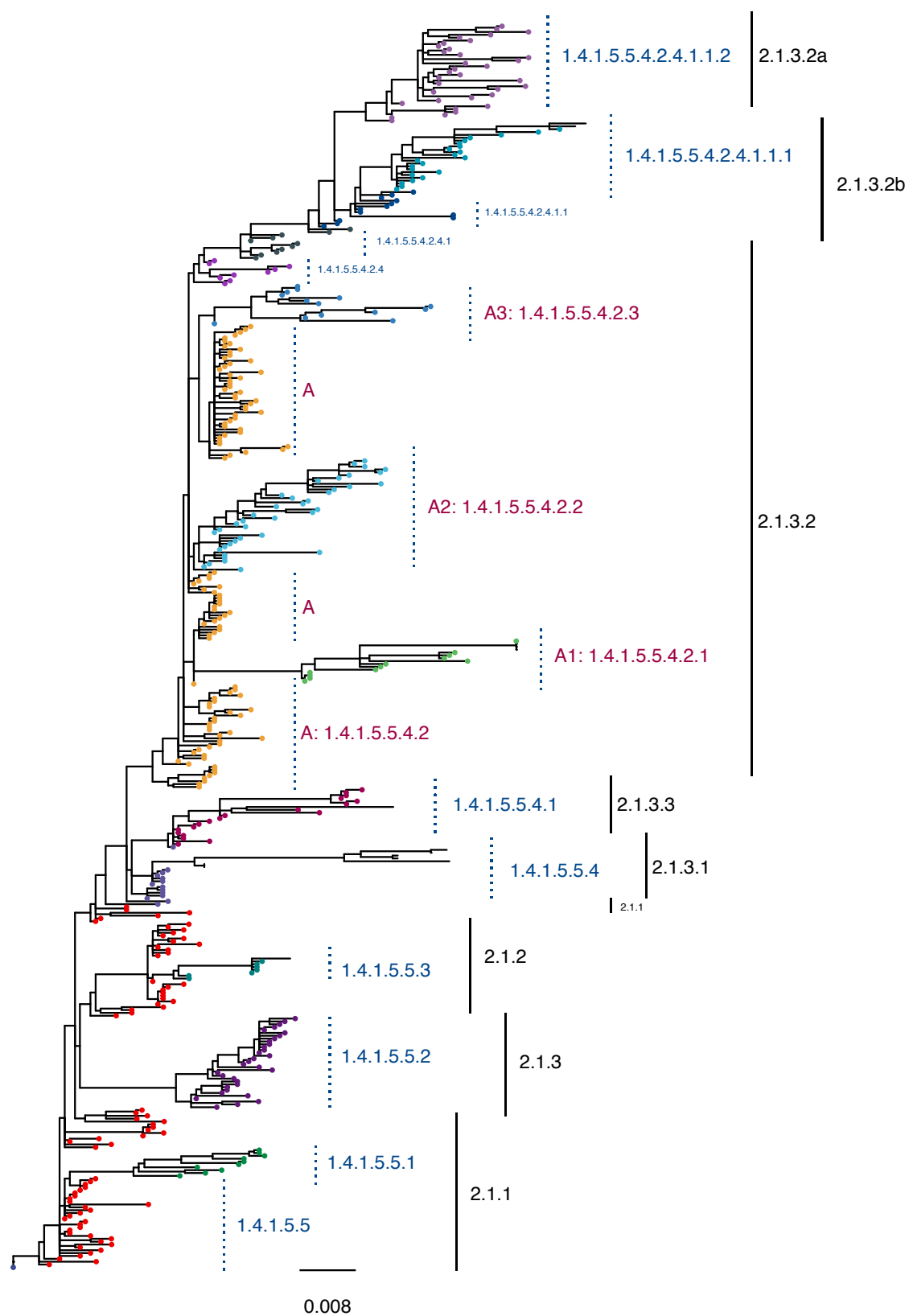
305 **Optimal PhyCLIP clustering result for HPAI avian H5 viruses**

306 For the full phylogeny of Gs/GD-like H5 viruses from the 2015 nomenclature update, the optimal
307 parameter set combined a minimum cluster size of 7, an FDR of 0.15 and a γ of 3. The optimal clustering
308 configuration clustered 98% of the sequences into a total of 89 clusters with a median cluster size of 21
309 sequences.

310 The topology of the optimal clustering result yields informative source-sink trajectories that are
311 supported by previously reported phylogenetic and phylogeographic evidence of the global panzootic
312 of the Gs/GD-like H5N1 lineage (Duan et al. 2008; Wang et al. 2008; Smith, Donis, for Animal
313 Health/Food, et al. 2015; The Global Consortium for H5N8 and Related Influenza Viruses 2016).

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

329



331 **Figure 3: Phylogeny of the Clade 2.1x viruses circulating in Indonesia.** The WHO/OIE/FAO H5 nomenclature is annotated
332 in black. PhyCLIP's cluster designation is indicated in blue, corresponding to tip colour. PhyCLIP's supercluster topology is
333 exemplified by Cluster A. The source population of the supercluster is annotated as A in pink, with tips coloured yellow. The
334 divergent descendant clusters are annotated as A.1, A.2 and A.3 respectively here. The letter A here is shorthand for its
335 nomenclature address, 1.4.1.5.5.4.2. This nomenclature address indicates that supercluster A is the second descendant of
336 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
337 in red. See Methods sections for full explanation of nomenclature addresses.

338

339 PhyCLIP's optimal cluster designation delineated the spatiotemporal structure of the phylogeny at high
340 resolution (Figure S3). Viruses circulating in south, central and northeast China and Hong Kong in 1996-
341 2003 acted as the source population for emergence of the classical viruses, seeding four lineages
342 (cluster 1, seeding cluster 1.1-1.4, Table S1). The second supercluster captures the first major wave of
343 expansion into neighbouring countries in east and southeast Asia in the early 2000's, with a source
344 population of viruses circulating in south central, east and north China, Viet Nam and Hong Kong in
345 2000-2003 (1.4 and 1.4.1 and their descendant lineages). The third supercluster captures the second
346 major wave of expansion of the Gs/GD-like H5 viruses, characterised by global spread (cluster 1.4.1.5
347 and its descendants). The source population of viruses from east, south central and southwest China,
348 Hong Kong and Viet Nam circulated from 2002-2005, giving rise to diverse and distinct viral lineages in
349 different regions globally (1.4.1.5.1-6). The supercluster topology highlights single lineage introductions
350 for countries with endemic circulation such as Indonesia and Egypt, but delineates multiple co-circulating
351 lineages structured over time. The clustering topology also highlights multiple incursions of diverse
352 viruses into countries such as South Korea and Japan (Table S3).

353 In addition to source-sink dynamics, distal dissociation also identifies probable outlying sequences,
354 defined as sequences more than 3 times the estimator of scale away from the median patristic distance
355 to the node. For example, PhyCLIP identifies seven outliers in its delineation of WHO/OIE/FAO clade
356 2.3.2.1c in the 2015 phylogeny (indicated by red tip-points in Figure 4). These sequences may represent
357 under-sampled populations with unobserved diversity, introductions from otherwise unsampled
358 populations or lower quality sequences introducing error into phylogenetic reconstruction.

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

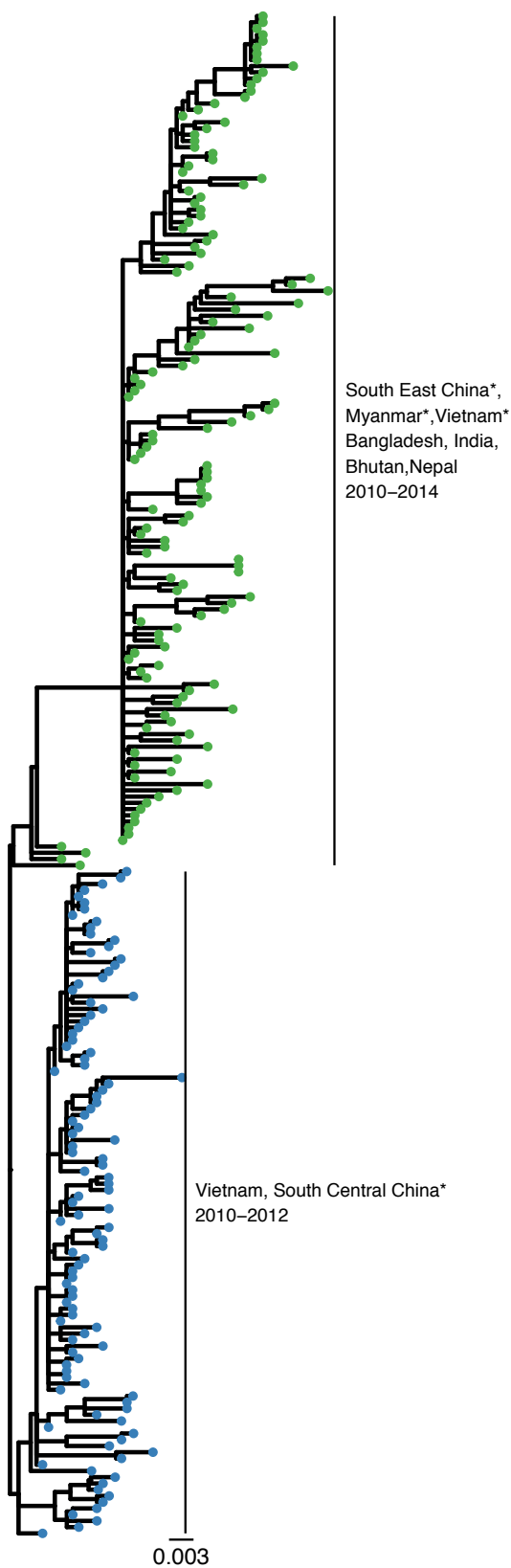
360 The current WHO/OIE/FAO nomenclature system designates 43 different clades and 7 clade-like
361 groupings for the full H5 phylogeny as of the 2015 update (Smith, Donis, and WHO/OIE/FAO H5
362 Evolution Working Group 2015) (Table S2). PhyCLIP recovers the current WHO/OIE/FAO H5
363 nomenclature with varying degrees of agreement across parameter sets, as measured by the variation
364 of information (VI) between the clustering partitions (Figure S4). VI is an information theoretic criterion
365 for comparing partitions of the same data set, based on the information lost and gained when moving

366 between partitions (Meilă 2007). A lower VI indicates more similar partitions. Parameter sets with a γ of
367 3 consistently had the lowest VI compared to the WHO/OIE/FAO system, indicating that the
368 WHO/OIE/FAO nomenclature system has the highest agreement with PhyCLIP clustering results that
369 tolerate higher within-cluster divergence.

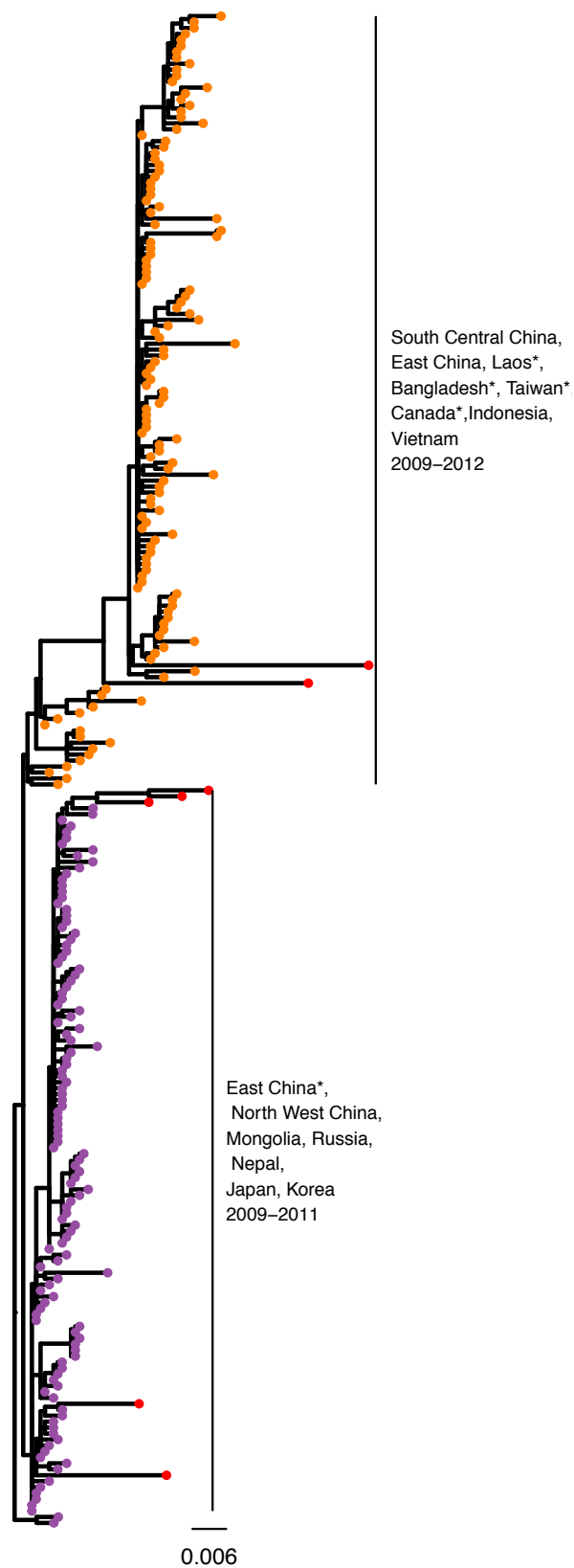
370 In the optimal clustering result, PhyCLIP delineates the spatiotemporal structure of the phylogeny with
371 a higher resolution than the WHO/OIE/FAO nomenclature system (89 vs 50 phylogenetic units, Figure
372 S3). The supercluster structure of the PhyCLIP clustering topology recapitulates the hierarchical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

383

Clade 2.3.2.1a



Clade 2.3.2.1c



385 **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
386 according to PhyCLIP's cluster designation. The tips coloured in red in B are viruses that were designated as outliers by
387 PhyCLIP's outlier detection. Countries represented by single viruses in the cluster are indicated with an asterisk.

388 **Impact of sampling**

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

402 Overall, the source-sink inference of PhyCLIP's clustering topology is largely consistent between the
403 WHO/OIE/FAO 2009 and 2015 update phylogeny optimal clustering results (Table S1). The optimal
404 result for the 2009 update phylogeny captures a similar topology and source population for the South
405 East Asian (clusters 1.3.1 and 1.3.1.1) and the post-2005 global wave of expansion (cluster
406 1.3.1.1.2.2.2) compared to the optimal 2015 clustering, with substantial overlap between the source
407 populations identified (100% and 83% for source populations for southeast Asia wave and global wave
408 respectively).

409 Changes in the clustering topology between the 2009 and 2015 update phylogenies are expected as
410 the underlying datasets are substantially different. More than 3000 viruses were added to the tree in the
411 six years between nomenclature updates. The Gs/GD-like H5 viruses evolved significantly in the
412 intervening period owing to genetic drift and reassortment. The addition of a large number of divergent
413 viruses to the underlying dataset fundamentally alters the ensemble statistical properties of the tree,
414 driving changes in the clustering configuration by changes in the global patristic distance distribution,
415 topology and statistical power between datasets. As a result, the ecological inferences drawn from the
416 2015 clustering topology are different from that of the 2009 phylogeny (Table S1).

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

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

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

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

445

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

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

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

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

466

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

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

476 The addition of the H5Nx viruses collected from 2014-2018 to the 2015 phylogeny changed the
477 distribution in two ways: 1. it added diversity to the right tail of the distribution, owing to the increased
478 divergence of the H5Nx viruses compared to the H5N1 viruses; 2. it increased the number of putative
479 clusters with low internal divergence, as a large amount of the H5Nx viruses possess highly similar HA
480 genes owing to both sampling biases during outbreaks and the relative short circulation time following
481 their emergence. This shift in the distribution reduced the within-cluster limit compared to that of the
482 2015 dataset (Figure S2).

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

487 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)
488 and 2018 (cluster 1.2.1.1.1.3.2) trees, including a shared core source population (Table S1).

489 The WHO/OIE/FAO clade 2.3.4.4 viruses are of interest owing to their reassortment promiscuity and
490 rapid global expansion. PhyCLIP delineates the clade 2.3.4.4 viruses into two distinct lineages, seeded
491 from a source population of viruses circulating in east and south-central China and Malaysia in 2005-
492 2010 (cluster 7.8, Table S1). The first lineage circulated in east, south central and northeast China in
493 2008 to 2011 (7.8.2, SFigure 11, Table S1). The second lineage (7.8.3) circulated in south central and
494 east China in 2008-2012 and seeded six distinct sub-lineages: Lineage 7.8.3.1 circulated in China from
495 2010 to 2014 before expanding to Viet Nam and circulating there for 2014-2015. Lineage 7.8.3.2
496 captures the global expansion of viruses from 2009 onwards. This includes the early subclade of H5N8
497 viruses described in Lycett et al (The Global Consortium for H5N8 and Related Influenza Viruses 2016).
498 Lineage 7.8.3.3 was restricted to China and was detected in 2013-2016. Lineage 7.8.3.4 also captures
499 a pan-national lineage that was detected from 2014 to 2016, and captures the more recent H5N8
500 subclade described in Lycett et al (The Global Consortium for H5N8 and Related Influenza Viruses
501 2016). Lineage 7.8.3.5 circulated in east and southeast Asia from 2013 to 2017. Lineage 7.8.3.6 is
502 seeded from a source population of viruses circulating in east and southeast Asia, expanding into
503 multiple co-circulating H5N6 southeast Asian lineages from 2013 onwards (Table S1).

504

505 **Benchmarking against other phylogenetic clustering tools**

506 PhyCLIP was benchmarked for performance against PhyloPart and ClusterPicker, two popular open-
507 source non-parametric phylogenetic clustering tools based on distance thresholds (Prosperi et al. 2011;
508 Ragonnet-Cronin et al. 2013). Both tools require a phylogenetic tree as input, as well as a user-specified
509 distance threshold and minimum statistical node-support level. Both tools carry out a depth-first traversal
510 of the tree, considering subtrees as putative clusters if the node support is above the user-defined level.
511 In PhyloPart, the user specifies a percentile of the global pairwise patristic distance distribution as a
512 threshold. If the median of the pairwise patristic distances of the putative cluster is below the percentile
513 threshold, a cluster is designated. ClusterPicker requires a user-defined maximum pairwise genetic
514 distance (calculated as p-distance directly from the sequences) threshold for cluster designation.

515 Accepted practice for these tools is to incorporate previous knowledge of sequence divergence into a
516 distance threshold or to calibrate the threshold over a tolerable range with metadata or expert
517 consensus. A direct comparison of PhyCLIP's clustering to PhyloPart or ClusterPicker is difficult owing
518 differences in generating within-cluster limits and a lack of prior knowledge of a meaningful delineation
519 of phylogenetic units for avian influenza to recommend a range of distance thresholds. PhyloPart and
520 ClusterPicker were applied to the 2009-update phylogeny (n=1224 sequences), with their input distance

521 thresholds, the within-cluster median pairwise patristic distance and within-cluster maximum genetic
522 distance respectively, set to match PhyCLIP's within-cluster limit for the optimal clustering result of the
523 2009-update phylogeny. Clustering results between PhyCLIP and PhyloPart showed high
524 correspondence (VI to PhyCLIP of 0.76, Figure S12), whereas the absolute maximum genetic distance
525 threshold of ClusterPicker lead to a highly stratified tree (VI to PhyCLIP of 1.87).

526 PhyCLIP is appreciably more computationally intensive than PhyloPart and ClusterPicker as it not only
527 has to parse the global pairwise patristic distance distribution of the phylogeny, but recursively
528 recalculate the distribution for subtrees in the distal dissociation approach, perform hypothesis testing
529 across every combinatorial pair of subtrees to test their inter-cluster divergence, as well as optimise the
530 ILP model. To relieve some of the computational cost, PhyCLIP is written in Python 2.7 employing
531 multiprocessing modules to parallelise the computational tasks involved resulting in ~3.2x times
532 speedup with 8 CPU cores relative to a single core run (Table 1).

533

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	14 seconds	0.6 GB	1
Phylopart	54 seconds	2.6 GB	1

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

535

536 Discussion

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

544 PhyCLIP's cluster assignment is agnostic to metadata but is capable of capturing the geographic and
545 temporal structure of the H5 phylogeny informatively. PhyCLIP recovers the overall structure of the
546 current WHO/OIE/FAOH5 nomenclature developed on a sequence divergence threshold, but delineates
547 more informative, higher resolution clusters that capture geographically-distinct subpopulations.

548 PhyCLIP therefore plausibly provides the foundation for an alternative nomenclature that minimizes the
549 limitations of currently employed approaches.

550 PhyCLIP's clustering is expected to improve with the addition of new sequences to the tree as new
551 information about the genetic diversity and evolutionary trajectory of the pathogen becomes known and
552 can be incorporated into the background diversity of the tree that informs the algorithm. Additionally,
553 topological information that capture how sequences are related by common ancestors is inherently
554 incorporated in PhyCLIP owing to its distal dissociation approach. The distal dissociation approach also
555 does not assume all clusters are monophyletic as the most recent common ancestor of all tips in a
556 cluster is not assumed to have no other descendants. As such, PhyCLIP can identify nested clusters
557 both as clusters with sufficiently high information content to meet the statistical requirements of cluster
558 designation or sufficiently diverse clusters that are dissociated from their ancestral nodes. The
559 designation of divergent descendant clusters nested within a supercluster suggestively captures source-
560 sink population dynamics that may be informative about the evolutionary trajectory of the clustered
561 sequences. At the same time, users could also opt for PhyCLIP to subsume subclusters that do not
562 violate the statistical criteria of the parent clusters into the latter, aiding higher level interpretation.
563 Importantly, the distal dissociation approach also identifies highly divergent outlying sequences that may
564 be indicative of under-sampled diversity.

565 PhyCLIP's methodology has limitations. Notably, PhyCLIP is tree-based and is therefore subject to error
566 in phylogenetic reconstruction. PhyCLIP does not include criteria for the statistical support of nodes
567 under consideration, which omits uncertainty in phylogenetic reconstruction. However, high statistical
568 support for a node does not necessarily indicate that all sequences subtended by it are highly related
569 but merely reflects the statistical support of the bipartition to the exclusion of other sequences.
570 Additionally, the relationship between the statistical significance of internal nodes and population
571 dynamics is unresolved as is an appropriate definition of a robustly supported node (Zharkikh and Li
572 1992; Susko 2009; Anisimova et al. 2011; Kumar et al. 2012; Volz et al. 2012). There is often less
573 phylogenetic signal to resolve internal nodes subtending small subtrees in measurably evolving
574 populations, increasing uncertainty in the arrangement of the internal structure of smaller subtrees. If a
575 statistical support threshold is set for nodes, these viruses will consistently be left unclustered or will be
576 forced to coalesce with more ancestral nodes subtending larger clusters, which would violate PhyCLIP's
577 statistical framework.

578 As with any phylogenetic clustering methods, PhyCLIP is also sensitive to variation in sampling rates
579 (Volz et al. 2012). There is a significant surveillance bias towards certain pathogens (e.g. HPAI H5)
580 owing to their consequences for animal and human health. The evolution and divergence of these
581 pathogens is currently captured in surveillance data as a more accurate approximation to a continuum
582 of evolution. PhyCLIP's clustering is strongly influenced by the diversity in the input population it tests

583 against, and will perform best when the background diversity of the phylogeny is complete or
584 representative.

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

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

596

597 **Materials and methods**

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

599 PhyCLIP computes the robust estimator of scale (σ) either as the median absolute deviation (MAD) or
600 Qn . Note that MAD may not suitably account for any potential skewness of the pairwise sequence
601 patristic distance distribution as it inherently assumes symmetry about the median ($\bar{\mu}$). On the contrary,
602 Qn , an alternative estimator of scale proposed by Rousseeuw & Croux (1993), is as robust as MAD (i.e.
603 50% breakdown point), calculated solely using the differences between the values in the distribution
604 without needing a location estimate, and has been proven to be statistically more efficient in both
605 Gaussian and non-Gaussian distributions relative to MAD .

606

607 **Integer linear programming model**

608 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
609 variables indicating if subtree i satisfies the conditions for clustering as a clade ($n_i = 1$ if it does and
610 $n_i = 0$ vice versa, Figure 2C). Each sequence j subtended by subtree i is also assigned a binary variable
611 $l_{j,i}$ indicating if the sequence is clustered under subtree i ($l_{j,i} = 1$ if j is clustered under node i and $l_{j,i} =$
612 0 vice versa, Figure 2C). PhyCLIP then formulates the phylogenetic clustering problem as an integer
613 linear programming (ILP) model with the objective to maximize the number of sequences assigned with
614 cluster membership:

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

615

616 subject to the following constraints:

617

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

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

620

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

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

625

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

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

628

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

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

631

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

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

633

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

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

640 Although both tests are nonparameteric, the Kuiper's test statistic incorporates both the greatest positive
641 and negative deviations between the two distributions whereas the KS test statistic is defined only by
642 their maximum difference. As a result, the Kuiper's test becomes equally sensitive to differences to the
643 tails as well as the median of the distributions but the KS test works best when the distributions differ
644 mostly at the median. In other words, the KS test is good at detecting *shifts* between the distributions
645 but lacks the sensitivity to uncover *spreads* between the distributions characterized by changes in their
646 tails. Kuiper's test is, however, sensitive to detect both types of changes in distributions.

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

- 648 (i) Subtree i is ancestral to k . The hypothesis test assumes the null hypothesis that the pairwise
649 sequence patristic distance distribution of subtree k is statistically identical to the pairwise
650 sequence patristic distance distribution of its ancestor i .
- 651 (ii) Neither subtree i nor k is an ancestor of the other. In this case, two hypothesis tests are carried
652 out comparing the distribution of each subtree to the distribution of pairwise sequence patristic
653 distance should both subtrees be combined as a single cluster and we take the more
654 conservative $q_{i,k} = \max\{q_{i,combined}, q_{k,combined}\}$.

655

656 Nomenclature

657 Traversing the output clusters of PhyCLIP by pre-order of the input phylogeny, a unique number is
658 assigned to any cluster with no immediate ancestral supercluster precursor to it (i.e. parent node of the
659 cluster node is not part of any PhyCLIP clusters). Otherwise, the descendant cluster in question is
660 designated as a *child cluster* should its membership size be $>25^{\text{th}}$ percentile of PhyCLIP's output cluster
661 size distribution (i.e. for having proliferated in numbers substantial enough to be deemed a progeny
662 cluster). Every child cluster of a supercluster is assigned a progeny number separated by a decimal
663 point (e.g. 1.2 refers to the second child cluster of supercluster 1). On other hand, descendant clusters
664 that fall below the cluster size cut-off are distinguished from child clusters as *nested clusters*, each
665 assigned an address in the form of a parenthesized letter, alphabetised by tree traversal order, prefixed

666 by its parent supercluster nomenclature (e.g. 1.1(c) refers to the third nested cluster of supercluster 1.1).
667 Nested clusters in superclusters fundamentally have different properties from the sensitivity-induced
668 nested clusters discussed in New Approach section and cannot be subsumed as it will violate the within-
669 cluster limit of the parent supercluster. The structure of the resultant clustering topology is highlighted
670 in Figure 3.

671

672 **Phylogenetic analyses**

673 PhyCLIP's performance was evaluated on an empirical dataset. The sequence datasets used to
674 construct the haemagglutinin (HA) gene phylogenetic trees underlying the WHO/OIE/FAO nomenclature
675 for the A/goose/Guangdong /1/1996 (Gs/GD/96)-like H5 avian influenza viruses were downloaded from
676 GISAID (Anon 2008; WHO/OIE/FAO H5N1 Evolution Working Group 2012; WHO/OIE/FAO H5N1
677 Evolution Working Group 2014; Smith, Donis, and WHO/OIE/FAO H5 Evolution Working Group 2015).
678 The primary analysis is based on the full dataset included in the 2009(n=1224) and 2015(n=4357)
679 nomenclature updates. Viruses that were inconsistently included across WHO/OIE/FAO updates were
680 followed up and included (WHO/OIE/FAO HN Evolution Working Gr 2009; Smith, Donis,
681 andWHO/OIE/FAO H5 Evolution Working Group 2015). Sequences were curated based on criteria
682 defined by the H5 nomenclature: sequences with more than 5 ambiguous nucleotides, with a sequence
683 length shorter than 60% of the alignment, or with frameshifts or duplicated by name were removed. For
684 the 2018 phylogeny, all avian and human viruses from the Gs/GD-like H5 lineage were downloaded
685 from GISAID up to April 2018, including H5Nx subtypes H5N2, H5N3, H5N5, H5N6 and H5N8. An
686 alternative filtering approach compared to the published WHO nomenclature approach was applied to
687 ensure a dataset of high-quality sequences that would be robust to error in phylogenetic reconstruction
688 as PhyCLIP is inherently sensitive to topological information. In this approach, duplicate sequences and
689 sequences with a length below 95% of the full HA sequence or more than 1% ambiguous nucleotides
690 were discarded. Sequences were aligned with MAFFT v7.397 and trimmed to the start of the mature
691 protein (Kato et al. 2002). Each sequence set was annotated with the WHO/OIE/FAOH5 nomenclature
692 using LABEL(v0.5.2), and the version of the module corresponding to the nomenclature update of the
693 dataset (e.g. H5v2015 module for the full tree from the nomenclature update in 2015) (Shepard et al.
694 2014). Maximum likelihood phylogenetic trees were constructed for each dataset with RAxML 8.2.12
695 under the GTR+GAMMA substitution model, and rooted to Gs/GD/96 (Stamatakis 2014). Phylogenetic
696 trees were visualized using Figtree (<http://tree.bio.ed.ac.uk/software/figtree/>) and ggtree (Yu et al. 2017).

697

698 **Benchmarking**

699 PhyCLIP was benchmarked for performance against two other non-parametric clustering methods,
700 ClusterPicker (Ragonnet-Cronin et al. 2013) and PhyloPart (Prosperi et al. 2011). PhyloPart and
701 ClusterPicker were applied to the WHO/OIE/FAO 2009-update phylogeny, with their input distance
702 thresholds, the within-cluster median pairwise patristic distance and within-cluster maximum genetic
703 distance respectively, set to match PhyCLIP's within-cluster limit for the optimal clustering result of the
704 2009-update phylogeny. Required bootstrap support level was set to 0 in both PhyloPart and
705 ClusterPicker to make it comparable to PhyCLIP, which lacks node-support criteria. All programs were
706 run on the Ubuntu 16.04 LTS operating system with an Intel Core i7-4790 3.60 GHz CPU.

707

708 **Code availability**

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

711

712 **Acknowledgments**

713 We thank the GISAID Initiative and the influenza surveillance and research groups that openly shared
714 the genetic sequence data that made this work possible. A.X.H. was supported by the A*STAR Graduate
715 Scholarship programme from A*STAR to carry out his PhD work via collaboration between
716 Bioinformatics Institute (A*STAR) and NUS Graduate School for Integrative Sciences and Engineering
717 from the National University of Singapore. E.P. was funded by the Gates Cambridge Trust (Grant number
718 OPP1144). S.M.S. was supported by the A*STAR HEIDI programme (Grant number: H1699f0013) and
719 Bioinformatics Institute (A*STAR). C.A.R. was supported by University Research Fellowship from the
720 Royal Society.

721

722 **References**

723

724 Aldous JL, Pond SK, Poon A, Jain S, Qin H, Kahn JS, Kitahata M, Rodriguez B, Dennis AM, Boswell
725 SL, et al. 2012. Characterizing HIV Transmission Networks Across the United States. *Clin. Infect. Dis.*
726 55:1135–1143.

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

- 730 Anon. 2008. Toward a Unified Nomenclature System for Highly Pathogenic Avian Influenza Virus
731 (H5N1). *Emerg. Infect. Dis.* 14:e1–e1.
- 732 Burk RD, Chen Z, Harari A, Smith BC, Kocjan BJ, Maver PJ, Poljak M. 2011. Classification and
733 nomenclature system for human Alphapapillomavirus variants: general features, nucleotide landmarks
734 and assignment of HPV6 and HPV11 isolates to variant lineages. *Acta dermatovenerologica Alpina,*
735 *Pannonica, Adriat.* 20:113–123.
- 736 Dennis AM, Herbeck JT, Brown AL, Kellam P, de Oliveira T, Pillay D, Fraser C, Cohen MS. 2014.
737 Phylogenetic studies of transmission dynamics in generalized HIV epidemics: an essential tool where
738 the burden is greatest? *J. Acquir. Immune Defic. Syndr.* 67:181–195.
- 739 Van Doorslaer K, Bernard H-U, Chen Z, de Villiers E-M, zur Hausen H, Burk RD. 2011.
740 Papillomaviruses: evolution, Linnaean taxonomy and current nomenclature. *Trends Microbiol.* 19:49-
741 50; author reply 50-1.
- 742 Duan L, Bahl J, Smith GJD, Wang J, Vijaykrishna D, Zhang LJ, Zhang JX, Li KS, Fan XH, Cheung CL,
743 et al. 2008. The development and genetic diversity of H5N1 influenza virus in China, 1996–2006.
744 *Virology [Internet]* 380:243–254. Available from:
745 <https://www.sciencedirect.com/science/article/pii/S0042682208004856?via%3Dihub>
- 746 Gardy JL, Loman NJ. 2017. Towards a genomics-informed, real-time, global pathogen surveillance
747 system. *Nat. Rev. Genet.* 19:9–20.
- 748 Grabowski MK, Herbeck JT, Poon AFY. 1904. Genetic Cluster Analysis for HIV Prevention.
- 749 Hassan AS, Pybus OG, Sanders EJ, Albert J, Esbjörnsson J. 2017. Defining HIV-1 transmission
750 clusters based on sequence data. *AIDS* 31:1211–1222.
- 751 Hué S, Clewley JP, Cane PA, Pillay D. 2004. HIV-1 pol gene variation is sufficient for reconstruction of
752 transmissions in the era of antiretroviral therapy. *AIDS* 18:719–728.
- 753 Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple sequence
754 alignment based on fast Fourier transform. *Nucleic Acids Res.* 30:3059–3066.
- 755 Kroneman A, Vega E, Vennema H, Vinjé J, White PA, Hansman G, Green K, Martella V, Katayama K,
756 Koopmans M. 2013. Proposal for a unified norovirus nomenclature and genotyping. *Arch. Virol.*
757 158:2059–2068.
- 758 Kumar S, Filipowski AJ, Battistuzzi FU, Kosakovsky Pond SL, Tamura K. 2012. Statistics and Truth in
759 Phylogenomics. *Mol. Biol. Evol.* 29:457–472.
- 760 Lauber C, Gorbalenya AE. 2012. Toward Genetics-Based Virus Taxonomy: Comparative Analysis of a
761 Genetics-Based Classification and the Taxonomy of Picornaviruses. *J. Virol.* 86:3905–3915.

- 762 McIntyre CL, Knowles NJ, Simmonds P. 2013. Proposals for the classification of human rhinovirus
763 species A, B and C into genotypically assigned types. *J. Gen. Virol.* 94:1791–1806.
- 764 Meilä M. 2007. Comparing clusterings—an information based distance. *J. Multivar. Anal.* 98:873–895.
- 765 Ortiz JR, Neuzil KM. 2017. Influenza immunization of pregnant women in resource-constrained
766 countries: an update for funding and implementation decisions. *Curr. Opin. Infect. Dis.* 30:455–462.
- 767 Poon AFY, Gustafson R, Daly P, Zerr L, Demlow SE, Wong J, Woods CK, Hogg RS, Krajden M,
768 Moore D, et al. 2016. Near real-time monitoring of HIV transmission hotspots from routine HIV
769 genotyping: an implementation case study. *Lancet HIV* 3:e231–e238.
- 770 Poon AFY, Joy JB, Woods CK, Shurgold S, Colley G, Brumme CJ, Hogg RS, Montaner JSG, Harrigan
771 PR. 2015. The Impact of Clinical, Demographic and Risk Factors on Rates of HIV Transmission: A
772 Population-based Phylogenetic Analysis in British Columbia, Canada. *J. Infect. Dis.* 211:926–935.
- 773 Prosperi MCF, Ciccozzi M, Fanti I, Saladini F, Pecorari M, Borghi V, Di Giambenedetto S, Bruzzone B,
774 Capetti A, Vivarelli A, et al. 2011. A novel methodology for large-scale phylogeny partition. *Nat.*
775 *Commun.* 2:321.
- 776 Prosperi MCF, De Luca A, Di Giambenedetto S, Bracciale L, Fabbiani M, Cauda R, Salemi M. 2010.
777 The Threshold Bootstrap Clustering: A New Approach to Find Families or Transmission Clusters
778 within Molecular Quasispecies. Poon AFY, editor. *PLoS One* 5:e13619.
- 779 Pu J, Wang S, Yin Y, Zhang G, Carter RA, Wang J, Xu G, Sun H, Wang M, Wen C, et al. 2015.
780 Evolution of the H9N2 influenza genotype that facilitated the genesis of the novel H7N9 virus. *Proc.*
781 *Natl. Acad. Sci. U. S. A.* 112:548–553.
- 782 Ragonnet-Cronin M, Hodcroft E, Hué S, Fearnhill E, Delpech V, Brown AJ, Lycett S, Holmes E, Nee
783 S, Rambaut A, et al. 2013. Automated analysis of phylogenetic clusters. *BMC Bioinformatics* 14:317.
- 784 Rose R, Lamers SL, Dollar JJ, Grabowski MK, Hodcroft EB, Ragonnet-Cronin M, Wertheim JO, Redd
785 AD, German D, Laeyendecker O. 2017. Identifying Transmission Clusters with Cluster Picker and HIV-
786 TRACE. *AIDS Res. Hum. Retroviruses* 33:211–218.
- 787 Rousseeuw PJ, Croux C. 1993. Alternatives to the Median Absolute Deviation. *J. Am. Stat. Assoc.*
788 [Internet] 88:1273–1283. Available from: [https://www.tandfonline-](https://www.tandfonline-com.libproxy1.nus.edu.sg/doi/pdf/10.1080/01621459.1993.10476408?needAccess=true)
789 [com.libproxy1.nus.edu.sg/doi/pdf/10.1080/01621459.1993.10476408?needAccess=true](https://www.tandfonline-com.libproxy1.nus.edu.sg/doi/pdf/10.1080/01621459.1993.10476408?needAccess=true)
- 790 Shepard SS, Davis CT, Bahl J, Rivaller P, York IA, Donis RO. 2014. LABEL: Fast and Accurate
791 Lineage Assignment with Assessment of H5N1 and H9N2 Influenza A Hemagglutinins. Woo PCY,
792 editor. *PLoS One* 9:e86921.
- 793 Simmonds P, McIntyre C, Savolainen-Kopra C, Tapparel C, Mackay IM, Hovi T. 2010. Proposals for

- 794 the classification of human rhinovirus species C into genotypically assigned types. *J. Gen. Virol.*
795 91:2409–2419.
- 796 Smith DB, Bukh J, Kuiken C, Muerhoff AS, Rice CM, Stapleton JT, Simmonds P. Expanded
797 Classification of Hepatitis C Virus Into 7 Genotypes and 67 Subtypes: Updated Criteria and Genotype
798 Assignment Web Resource.
- 799 Smith GJD, Donis RO, for Animal Health/Food WHOO, Group AO (WHO/OIE/FAO) HEW. 2015.
800 Nomenclature updates resulting from the evolution of avian influenza A(H5) virus clades 2.1.3.2a,
801 2.2.1, and 2.3.4 during 2013-2014. *Influenza Other Respi. Viruses* [Internet] 9:271–276. Available
802 from: <http://dx.doi.org/10.1111/irv.12324>
- 803 Smith GJD, Donis RO, World Health Organization/World Organisation for Animal Health/Food and
804 Agriculture Organization (WHO/OIE/FAO) H5 Evolution Working Group WHOO for AH and AO
805 (WHO/OIE/FAO) HEW. 2015. Nomenclature updates resulting from the evolution of avian influenza
806 A(H5) virus clades 2.1.3.2a, 2.2.1, and 2.3.4 during 2013-2014. *Influenza Other Respi. Viruses* 9:271–
807 276.
- 808 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
809 phylogenies. *Bioinformatics* 30:1312–1313.
- 810 Susko E. 2009. Bootstrap Support Is Not First-Order Correct. *Syst. Biol.* 58:211–223.
- 811 The Global Consortium for H5N8 and Related Influenza Viruses. 2016. Role for migratory wild birds in
812 the global spread of avian influenza H5N8. *Science* (80-.). [Internet] 354:213 LP-217. Available from:
813 <http://science.sciencemag.org/content/354/6309/213.abstract>
- 814 Valastro V, Holmes EC, Britton P, Fusaro A, Jackwood MW, Cattoli G, Monne I. 2016. S1 gene-based
815 phylogeny of infectious bronchitis virus: An attempt to harmonize virus classification. *Infect. Genet.*
816 *Evol.* 39:349–364.
- 817 Volz EM, Koopman JS, Ward MJ, Brown AL, Frost SDW. 2012. Simple Epidemiological Dynamics
818 Explain Phylogenetic Clustering of HIV from Patients with Recent Infection. Fraser C, editor. *PLoS*
819 *Comput. Biol.* 8:e1002552.
- 820 Wang J, Vijaykrishna D, Duan L, Bahl J, Zhang JX, Webster RG, Peiris JSM, Chen H, Smith GJD,
821 Guan Y. 2008. Identification of the Progenitors of Indonesian and Vietnamese Avian Influenza A
822 (H5N1) Viruses from Southern China. *J. Virol.* [Internet] 82:3405 LP-3414. Available from:
823 <http://jvi.asm.org/content/82/7/3405.abstract>
- 824 WHO/OIE/FAO H5N1 Evolution Working Group. 2008. Toward a Unified Nomenclature System for
825 Highly Pathogenic Avian Influenza Virus (H5N1). *Emerg. Infect. Dis.* 14:e1–e1.

- 826 WHO/OIE/FAO H5N1 Evolution Working Group WHEW. 2012. Continued evolution of highly
827 pathogenic avian influenza A (H5N1): updated nomenclature. *Influenza Other Respi. Viruses* 6:1–5.
- 828 WHO/OIE/FAO HN Evolution Working Gr. 2009. Continuing progress towards a unified nomenclature
829 for the highly pathogenic H5N1 avian influenza viruses: divergence of clade 2?2 viruses. *Influenza*
830 *Other Respi. Viruses* 3:59–62.
- 831 World Health Organization/World Organisation for Animal Health/Food and Agriculture Organization
832 (WHO/OIE/FAO) H5N1 Evolution Working Group. 2014. Revised and updated nomenclature for highly
833 pathogenic avian influenza A (H5N1) viruses. *Influenza Other Respi. Viruses* 8:384–388.
- 834 Yu G, Smith DK, Zhu H, Guan Y, Lam TTY. 2017. Ggtree: an R Package for Visualization and
835 Annotation of Phylogenetic Trees With Their Covariates and Other Associated Data. *Methods Ecol.*
836 *Evol.* 8:28–36.
- 837 Zharkikh A, Li WH. 1992. Statistical properties of bootstrap estimation of phylogenetic variability from
838 nucleotide sequences. I. Four taxa with a molecular clock. *Mol. Biol. Evol.* 9:1119–1147.
- 839