

## Phylogenetic Clustering by Linear Integer Programming (PhyCLIP)

Alvin X. Han<sup>†,1,2,3</sup>, Edyth Parker<sup>†,3,4</sup>, Frits Scholer<sup>5</sup>, Sebastian Maurer-Stroh<sup>1,2</sup>, Colin A. Russell<sup>3</sup>

<sup>1</sup>Bioinformatics Institute, Agency for Science, Technology and Research (A\*STAR), Singapore

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

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

<sup>4</sup>Department of Veterinary Medicine, University of Cambridge, Cambridge, United Kingdom

<sup>5</sup>Department of Medical Microbiology, Academic Medical Centre, University of Amsterdam, Amsterdam, The Netherlands

<sup>†</sup>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

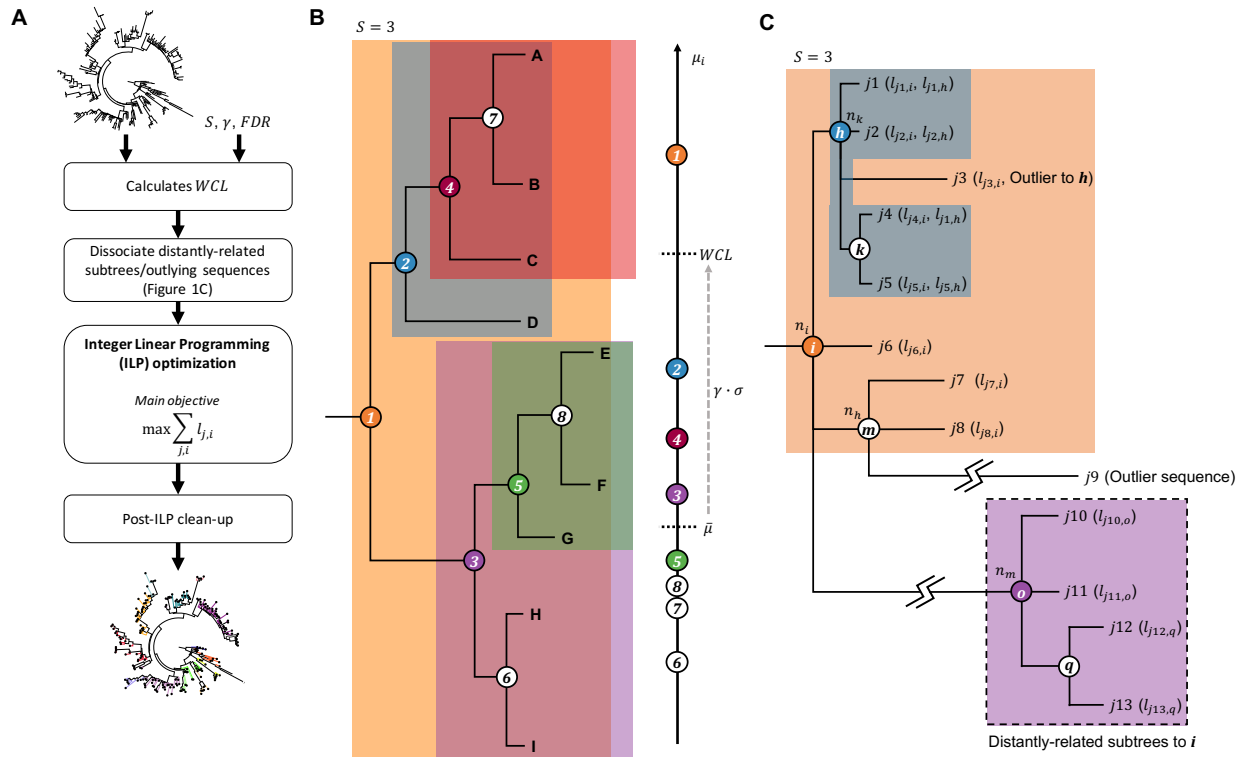
85 into meaningful subsets is therefore complex and is mostly performed *ad hoc* through exploratory  
86 analyses with uninformative sensitivity analyses across thresholds. As expected, cluster membership is  
87 highly sensitive to the threshold applied and therefore results can be unstable across 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  $S$ ,  $\gamma$  and  $FDR$  as the inputs for PhyCLIP. After determining  $WCL$ , PhyCLIP  
 105 dissociates distant subtrees and outlying sequences that skew  $\mu_i$  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  $\geq 3$  sequences are considered for clustering by the ILP model but  $WCL$   
 112 is determined from  $\mu_i$  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  $j_9$  are dissociated from subtree **i** as they are exceedingly distant from  
 114 **i**. If sequences  $j_1$ ,  $j_2$ ,  $j_4$  and  $j_5$  are clustered under subtree **h** while  $j_3$  is clustered under subtree **i** by ILP optimization, a  
 115 post-ILP clean up step will remove  $j_3$  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 ( $\gamma$ ) 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 ( $\mu_i$ ). 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  $\sigma$  is any robust estimator of scale (e.g. median absolute deviation (*MAD*) or  $Qn$ , 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  $\mu_i$  of the subtree they subtend, leading to inaccurate over-estimation of the  
145 internal divergence of the putative subtree. Similarly, distantly-related descendant subtrees can  
146 artificially inflate  $\mu_i$  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  $\mu_i$  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  $\mu_i$  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.  $\mu_i$  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 sub-clusters into their parent  
195 clusters if sequences in the descendant sub-cluster 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 sub-clusters 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 sub-clusters 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 ( $\gamma$ ) considered should be a factor of the intra- and inter-cluster distance related  
230    to the research question.

231    Meta-data 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 ( $\gamma$ ) 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.

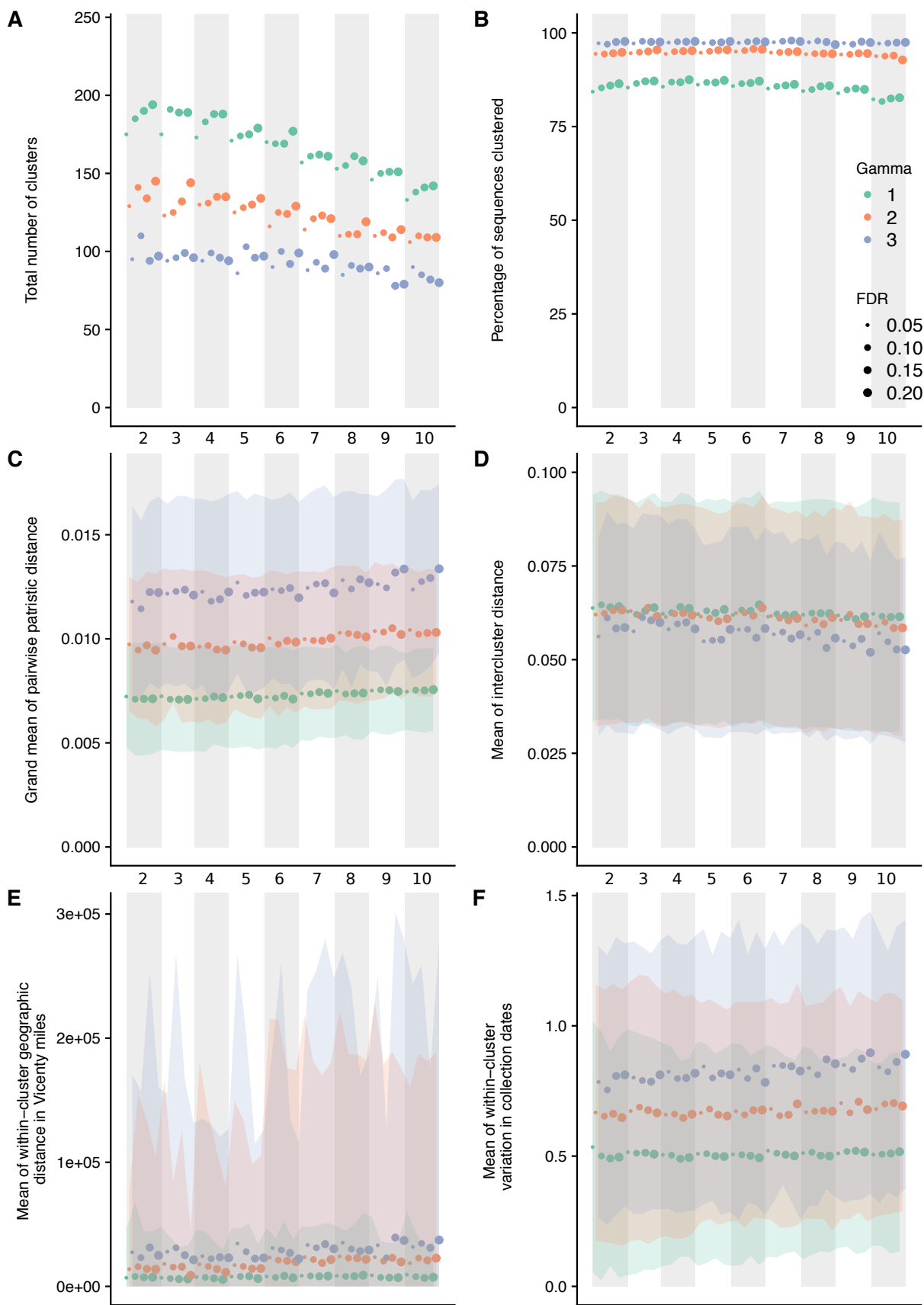
266 As PhyCLIP incorporates topological information of the phylogeny into the clustering construct, non-  
267 terminal internal nodes with zero branch lengths can lead to erroneous clustering and over-delineation  
268 (Figure S1). Such internal nodes are usually found in bifurcating trees as representations 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,

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

276

277 **Influence of the parameters**

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



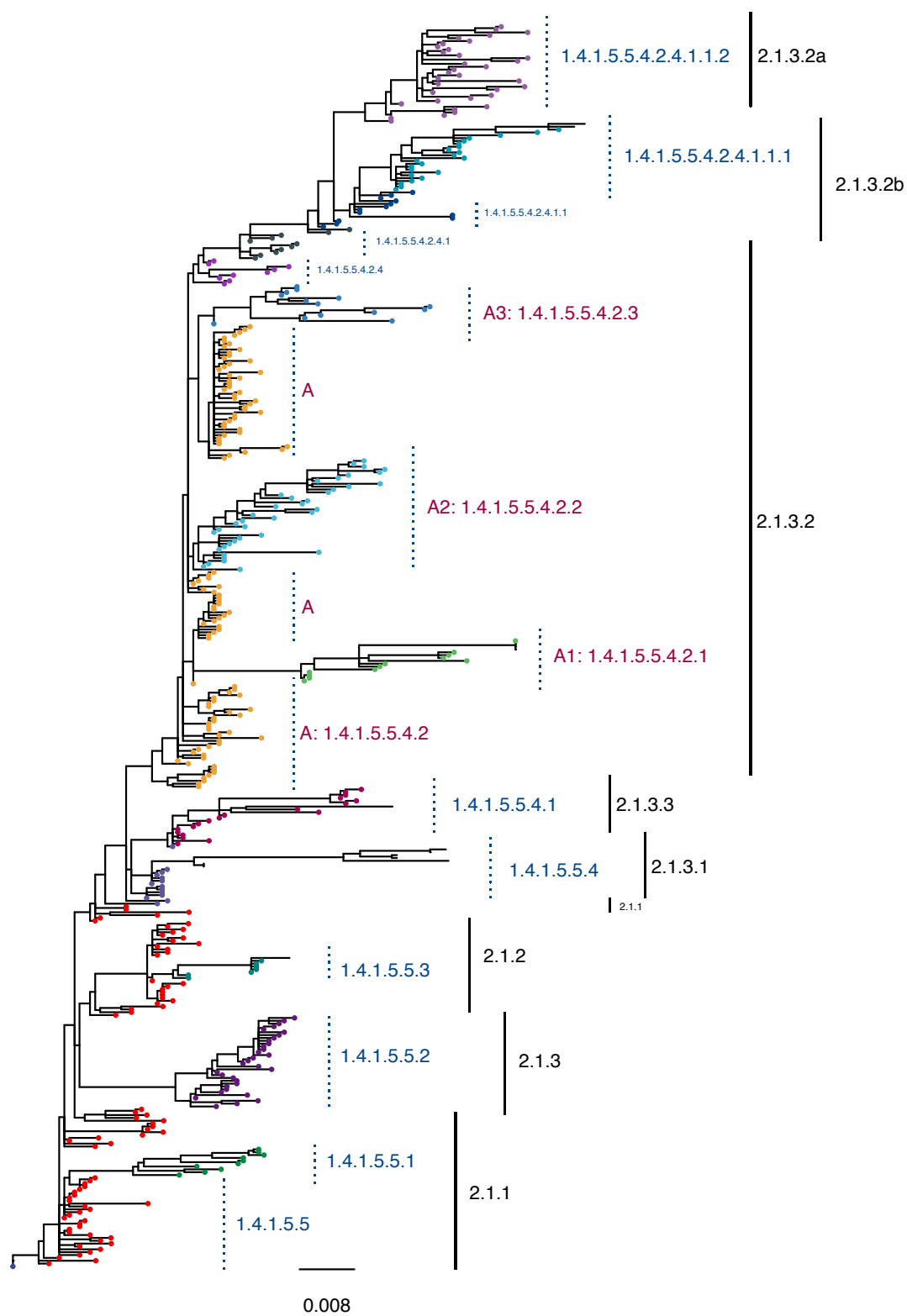
**Figure 2: Influence of parameters on the clustering properties of PhyCLIP in the WHO/OIE/FAO 2015-update phylogeny.** Figure A-F have the parameter set combinations ordered according to minimum cluster size, FDR and gamma on the x axis. The banded background and x-axis subscript numbering indicate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

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

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

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

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



332 **Figure 3: Phylogeny of the Clade 2.1x viruses circulating in Indonesia.** The WHO/OIE/FAO H5 nomenclature is annotated  
333 in black. PhyCLIP's cluster designation is indicated in blue, corresponding to tip colour. PhyCLIP's supercluster topology is  
334 exemplified by Cluster A. The source population of the supercluster is annotated as A in pink, with tips coloured yellow. The  
335 divergent descendant clusters are annotated as A.1, A.2 and A.3 respectively here. The letter A here is shorthand for its  
336 nomenclature address, 1.4.1.5.4.2. This nomenclature address indicates that supercluster A is the second descendant of  
337 cluster 1.4.1.5.4 (indicated in light purple), which in turn is the forth descendant of the source supercluster 1.4.1.5.5, indicated  
338 in red. See Methods sections for full explanation of nomenclature addresses.

339

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

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

360

## 361 **Comparison to the WHO/OIE/FAO H5 nomenclature**

362 The current WHO/OIE/FAO nomenclature system designates 43 different clades and 7 clade-like  
363 groupings for the full H5 phylogeny as of the 2015 update (Smith, Donis, and WHO/OIE/FAO H5  
364 Evolution Working Group 2015) (Table S2). PhyCLIP recovers the current WHO/OIE/FAO H5  
365 nomenclature with varying degrees of agreement across parameter sets, as measured by the variation

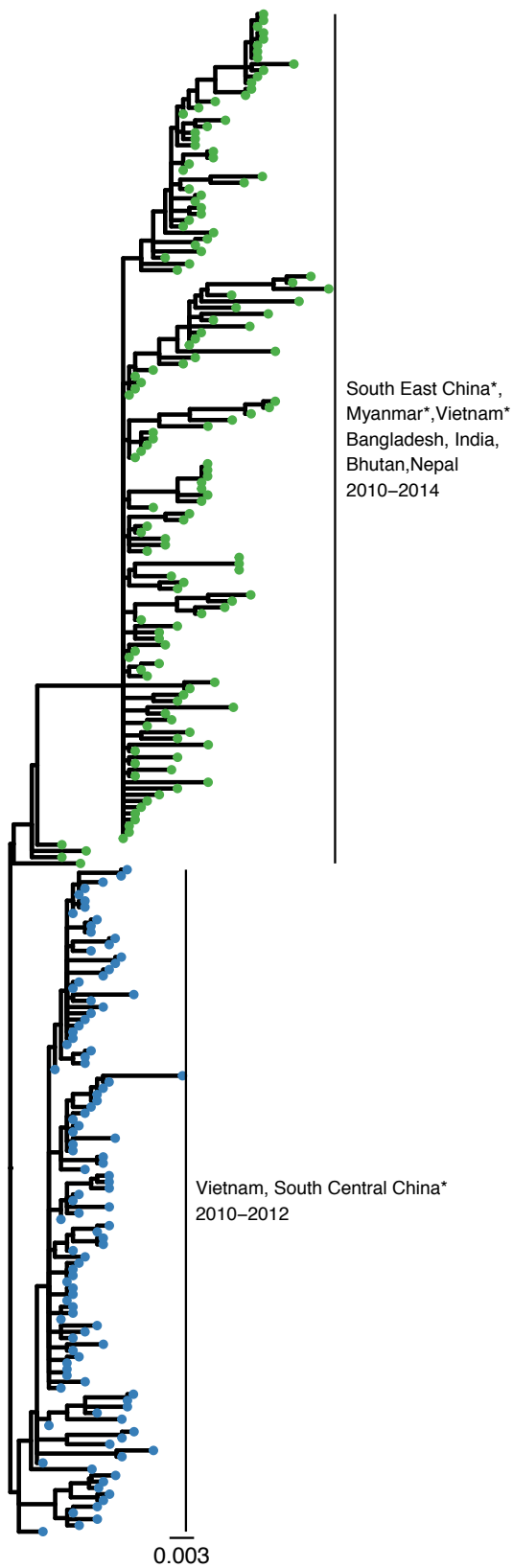
366 of information (VI) between the clustering partitions (Figure S4). VI is an information theoretic criterion  
367 for comparing partitions of the same data set, based on the information lost and gained when moving  
368 between partitions (Meilă 2007). A lower VI indicates more similar partitions. Parameter sets with a  $\gamma$  of  
369 3 consistently had the lowest VI compared to the WHO/OIE/FAO system, indicating that the  
370 WHO/OIE/FAO nomenclature system has the highest agreement with PhyCLIP clustering results that  
371 tolerate higher within-cluster divergence.

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

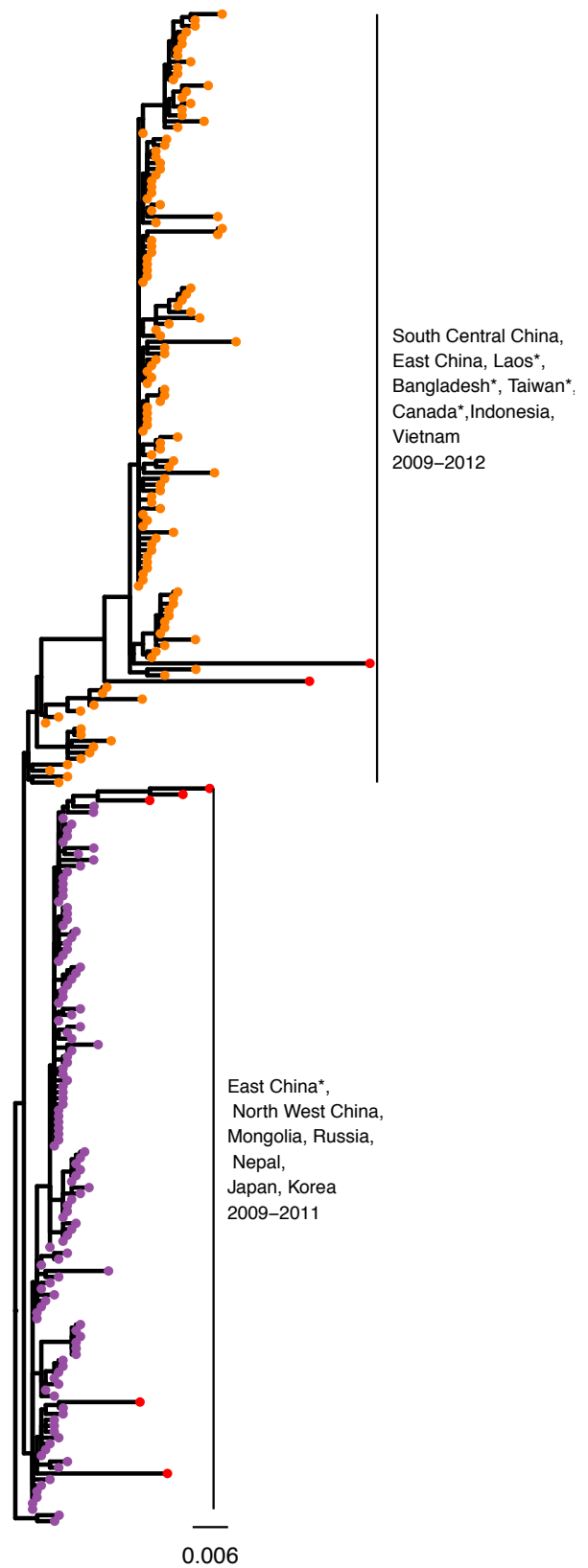
385



Clade 2.3.2.1a



Clade 2.3.2.1c



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

390 **Impact of sampling**

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

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

411 Changes in the clustering topology between the 2009 and 2015 update phylogenies are expected as  
412 the underlying datasets are substantially different. More than 3000 viruses were added to the tree in the  
413 six years between nomenclature updates. The Gs/GD-like H5 viruses evolved significantly in the  
414 intervening period owing to genetic drift and reassortment. The addition of a large number of divergent  
415 viruses to the underlying dataset fundamentally alters the ensemble statistical properties of the tree,  
416 driving changes in the clustering configuration by changes in the global patristic distance distribution,  
417 topology and statistical power between datasets. As a result, the ecological inferences drawn from the  
418 2015 clustering topology are different from that of the 2009 phylogeny (Table S1).

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

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

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

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

447

## 448 **Comparison of optimal to suboptimal clustering results**

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

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

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

468

#### PhyCLIP clustering of the 1996-2018 H5Nx phylogeny

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

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

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

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

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

## **Benchmarking against other phylogenetic clustering tools**

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

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

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

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

535

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

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

537

538 **Discussion**

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

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

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

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

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

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

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

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

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

598

## 599 **Materials and methods**

### 600 **Robust estimator of scale (deviation)**

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

608

### 609 **Integer linear programming model**

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



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

617

618 subject to the following constraints:

619

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

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

622

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

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

627

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

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

630

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

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

633

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

634 Constraint (7) ensures that  $\mu_i$  of all clades fall below the stipulated  $WCL$  limit.

635

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

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

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

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

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

657

## 658 Nomenclature

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

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

673

## 674 **Phylogenetic analyses**

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

699

## 700 **Benchmarking**

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

**Code availability**

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

**Acknowledgments**

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

**References**

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

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

Anon. 2008. Toward a Unified Nomenclature System for Highly Pathogenic Avian Influenza Virus

732 (H5N1). *Emerg. Infect. Dis.* 14:e1–e1.

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

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

740 Van Doorslaer K, Bernard H-U, Chen Z, de Villiers E-M, zur Hausen H, Burk RD. 2011.  
 741 Papillomaviruses: evolution, Linnaean taxonomy and current nomenclature. *Trends Microbiol.* 19:49-  
 742 50; author reply 50-1.

743 Duan L, Bahl J, Smith GJD, Wang J, Vijaykrishna D, Zhang LJ, Zhang JX, Li KS, Fan XH, Cheung CL,  
 744 et al. 2008. The development and genetic diversity of H5N1 influenza virus in China, 1996–2006.  
 745 *Virology* [Internet] 380:243–254. Available from:  
 746 <https://www.sciencedirect.com/science/article/pii/S0042682208004856?via%3Dihub>

747 Gardy JL, Loman NJ. 2017. Towards a genomics-informed, real-time, global pathogen surveillance  
 748 system. *Nat. Rev. Genet.* 19:9–20.

749 Grabowski MK, Herbeck JT, Poon AFY. 1904. Genetic Cluster Analysis for HIV Prevention.

750 Hassan AS, Pybus OG, Sanders EJ, Albert J, Esbjörnsson J. 2017. Defining HIV-1 transmission  
 751 clusters based on sequence data. *AIDS* 31:1211–1222.

752 Hué S, Clewley JP, Cane PA, Pillay D. 2004. HIV-1 pol gene variation is sufficient for reconstruction of  
 753 transmissions in the era of antiretroviral therapy. *AIDS* 18:719–728.

754 Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple sequence  
 755 alignment based on fast Fourier transform. *Nucleic Acids Res.* 30:3059–3066.

756 Kroneman A, Vega E, Vennema H, Vinjé J, White PA, Hansman G, Green K, Martella V, Katayama K,  
 757 Koopmans M. 2013. Proposal for a unified norovirus nomenclature and genotyping. *Arch. Virol.*  
 758 158:2059–2068.

759 Kumar S, Filipski AJ, Battistuzzi FU, Kosakovsky Pond SL, Tamura K. 2012. Statistics and Truth in  
 760 Phylogenomics. *Mol. Biol. Evol.* 29:457–472.

761 Lauber C, Gorbalenya AE. 2012. Toward Genetics-Based Virus Taxonomy: Comparative Analysis of a  
 762 Genetics-Based Classification and the Taxonomy of Picornaviruses. *J. Virol.* 86:3905–3915.

763 McIntyre CL, Knowles NJ, Simmonds P. 2013. Proposals for the classification of human rhinovirus  
 764 species A, B and C into genotypically assigned types. *J. Gen. Virol.* 94:1791–1806.

765 Meilä M. 2007. Comparing clusterings—an information based distance. *J. Multivar. Anal.* 98:873–895.

766 Ortiz JR, Neuzil KM. 2017. Influenza immunization of pregnant women in resource-constrained  
 767 countries: an update for funding and implementation decisions. *Curr. Opin. Infect. Dis.* 30:455–462.

768 Poon AFY, Gustafson R, Daly P, Zerr L, Demlow SE, Wong J, Woods CK, Hogg RS, Krajden M,  
 769 Moore D, et al. 2016. Near real-time monitoring of HIV transmission hotspots from routine HIV  
 770 genotyping: an implementation case study. *Lancet HIV* 3:e231–e238.

771 Poon AFY, Joy JB, Woods CK, Shurgold S, Colley G, Brumme CJ, Hogg RS, Montaner JSG, Harrigan  
 772 PR. 2015. The Impact of Clinical, Demographic and Risk Factors on Rates of HIV Transmission: A  
 773 Population-based Phylogenetic Analysis in British Columbia, Canada. *J. Infect. Dis.* 211:926–935.

774 Prosperi MCF, Ciccozzi M, Fanti I, Saladini F, Pecorari M, Borghi V, Di Giambenedetto S, Bruzzone B,  
 775 Capetti A, Vivarelli A, et al. 2011. A novel methodology for large-scale phylogeny partition. *Nat.*  
 776 *Commun.* 2:321.

777 Prosperi MCF, De Luca A, Di Giambenedetto S, Bracciale L, Fabbiani M, Cauda R, Salemi M. 2010.  
 778 The Threshold Bootstrap Clustering: A New Approach to Find Families or Transmission Clusters  
 779 within Molecular Quasispecies. Poon AFY, editor. *PLoS One* 5:e13619.

780 Pu J, Wang S, Yin Y, Zhang G, Carter RA, Wang J, Xu G, Sun H, Wang M, Wen C, et al. 2015.  
 781 Evolution of the H9N2 influenza genotype that facilitated the genesis of the novel H7N9 virus. *Proc.*  
 782 *Natl. Acad. Sci. U. S. A.* 112:548–553.

783 Ragonnet-Cronin M, Hodcroft E, Hué S, Fearnhill E, Delpech V, Brown AJ, Lycett S, Holmes E, Nee  
 784 S, Rambaut A, et al. 2013. Automated analysis of phylogenetic clusters. *BMC Bioinformatics* 14:317.

785 Rose R, Lamers SL, Dollar JJ, Grabowski MK, Hodcroft EB, Ragonnet-Cronin M, Wertheim JO, Redd  
 786 AD, German D, Laeyendecker O. 2017. Identifying Transmission Clusters with Cluster Picker and HIV-  
 787 TRACE. *AIDS Res. Hum. Retroviruses* 33:211–218.

788 Rousseeuw PJ, Croux C. 1993. Alternatives to the Median Absolute Deviation. *J. Am. Stat. Assoc.*  
 789 [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)  
 790 [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)

791 Shepard SS, Davis CT, Bahl J, Rivaller P, York IA, Donis RO. 2014. LABEL: Fast and Accurate  
 792 Lineage Assignment with Assessment of H5N1 and H9N2 Influenza A Hemagglutinins. Woo PCY,  
 793 editor. *PLoS One* 9:e86921.

794 Simmonds P, McIntyre C, Savolainen-Kopra C, Tapparel C, Mackay IM, Hovi T. 2010. Proposals for

795 the classification of human rhinovirus species C into genotypically assigned types. *J. Gen. Virol.*  
796 91:2409–2419.

797 Smith DB, Bukh J, Kuiken C, Muerhoff AS, Rice CM, Stapleton JT, Simmonds P. Expanded  
798 Classification of Hepatitis C Virus Into 7 Genotypes and 67 Subtypes: Updated Criteria and Genotype  
799 Assignment Web Resource.

800 Smith GJD, Donis RO, for Animal Health/Food WHOO, Group AO (WHO/OIE/FAO) HEW. 2015.  
801 Nomenclature updates resulting from the evolution of avian influenza A(H5) virus clades 2.1.3.2a,  
802 2.2.1, and 2.3.4 during 2013-2014. *Influenza Other Respi. Viruses* [Internet] 9:271–276. Available  
803 from: <http://dx.doi.org/10.1111/irv.12324>

804 Smith GJD, Donis RO, World Health Organization/World Organisation for Animal Health/Food and  
805 Agriculture Organization (WHO/OIE/FAO) H5 Evolution Working Group WHOO for AH and AO  
806 (WHO/OIE/FAO) HEW. 2015. Nomenclature updates resulting from the evolution of avian influenza  
807 A(H5) virus clades 2.1.3.2a, 2.2.1, and 2.3.4 during 2013-2014. *Influenza Other Respi. Viruses* 9:271–  
808 276.

809 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
810 phylogenies. *Bioinformatics* 30:1312–1313.

811 Susko E. 2009. Bootstrap Support Is Not First-Order Correct. *Syst. Biol.* 58:211–223.

812 The Global Consortium for H5N8 and Related Influenza Viruses. 2016. Role for migratory wild birds in  
813 the global spread of avian influenza H5N8. *Science* (80-. ). [Internet] 354:213 LP-217. Available from:  
814 <http://science.sciencemag.org/content/354/6309/213.abstract>

815 Valastro V, Holmes EC, Britton P, Fusaro A, Jackwood MW, Cattoli G, Monne I. 2016. S1 gene-based  
816 phylogeny of infectious bronchitis virus: An attempt to harmonize virus classification. *Infect. Genet.*  
817 *Evol.* 39:349–364.

818 Volz EM, Koopman JS, Ward MJ, Brown AL, Frost SDW. 2012. Simple Epidemiological Dynamics  
819 Explain Phylogenetic Clustering of HIV from Patients with Recent Infection. Fraser C, editor. *PLoS*  
820 *Comput. Biol.* 8:e1002552.

821 Wang J, Vijaykrishna D, Duan L, Bahl J, Zhang JX, Webster RG, Peiris JSM, Chen H, Smith GJD,  
822 Guan Y. 2008. Identification of the Progenitors of Indonesian and Vietnamese Avian Influenza A  
823 (H5N1) Viruses from Southern China. *J. Virol.* [Internet] 82:3405 LP-3414. Available from:  
824 <http://jvi.asm.org/content/82/7/3405.abstract>

825 WHO/OIE/FAO H5N1 Evolution Working Group. 2008. Toward a Unified Nomenclature System for  
826 Highly Pathogenic Avian Influenza Virus (H5N1). *Emerg. Infect. Dis.* 14:e1–e1.

827 WHO/OIE/FAO H5N1 Evolution Working Group WHEW. 2012. Continued evolution of highly  
828 pathogenic avian influenza A (H5N1): updated nomenclature. *Influenza Other Respi. Viruses* 6:1–5.

829 WHO/OIE/FAO HN Evolution Working Gr. 2009. Continuing progress towards a unified nomenclature  
830 for the highly pathogenic H5N1 avian influenza viruses: divergence of clade 2?2 viruses. *Influenza*  
831 *Other Respi. Viruses* 3:59–62.

832 World Health Organization/World Organisation for Animal Health/Food and Agriculture Organization  
833 (WHO/OIE/FAO) H5N1 Evolution Working Group. 2014. Revised and updated nomenclature for highly  
834 pathogenic avian influenza A (H5N1) viruses. *Influenza Other Respi. Viruses* 8:384–388.

835 Yu G, Smith DK, Zhu H, Guan Y, Lam TTY. 2017. Ggtree: an R Package for Visualization and  
836 Annotation of Phylogenetic Trees With Their Covariates and Other Associated Data. *Methods Ecol.*  
837 *Evol.* 8:28–36.

838 Zharkikh A, Li WH. 1992. Statistical properties of bootstrap estimation of phylogenetic variability from  
839 nucleotide sequences. I. Four taxa with a molecular clock. *Mol. Biol. Evol.* 9:1119–1147.

840