

1 **Development of avian influenza A(H5) virus datasets for Nextclade enables rapid**
2 **and accurate clade assignment**

3
4 Jordan T. Ort¹, Samuel S. Shepard², Sonja A. Zolnoski¹, Tommy T.-Y. Lam^{3,4}, Todd
5 Davis², Richard Neher⁵, Louise H. Moncla^{6,*}

6
7 ¹Department of Microbiology, Perelman School of Medicine, University of Pennsylvania,
8 Philadelphia, PA, USA

9 ²Influenza Division, National Center for Immunizations and Respiratory Disease,
10 Centers for Disease Control and Prevention, Atlanta, Georgia, USA

11 ³State Key Laboratory of Emerging Infectious Diseases, HKU-Pasteur Research Pole,
12 School of Public Health, The University of Hong Kong, Hong Kong SAR, China

13 ⁴Centre for Immunology and Infection, Hong Kong Science and Technology Park, Hong
14 Kong SAR, China.

15 ⁵Biozentrum, University of Basel, Basel, Switzerland

16 ⁶Department of Pathobiology, School of Veterinary Medicine, University of
17 Pennsylvania, Philadelphia, PA, USA

18
19 *Correspondence should be addressed to Louise Moncla, lhmoncla@upenn.edu

20 **Abstract**

21

22 The ongoing panzootic of highly pathogenic avian influenza (HPAI) A(H5) viruses is the
23 largest in history, with unprecedented transmission to multiple mammalian species.
24 Avian influenza A viruses of the H5 subtype circulate globally among birds and are
25 classified into distinct clades based on their hemagglutinin (HA) genetic sequences.
26 Thus, the ability to accurately and rapidly assign clades to newly sequenced isolates is
27 key to surveillance and outbreak response. Co-circulation of endemic, low pathogenic
28 avian influenza (LPAI) A(H5) lineages in North American and European wild birds
29 necessitates the ability to rapidly and accurately distinguish between infections arising
30 from these lineages and epizootic HPAI A(H5) viruses. However, currently available
31 clade assignment tools are limited and often require command line expertise, hindering
32 their utility for public health surveillance labs. To address this gap, we have developed
33 datasets to enable A(H5) clade assignments with Nextclade, a drag-and-drop tool
34 originally developed for SARS-CoV-2 genetic clade classification. Using annotated
35 reference datasets for all historical A(H5) clades, clade 2.3.2.1 descendants, and clade
36 2.3.4.4 descendants provided by the Food and Agriculture Organization/World Health
37 Organization/World Organisation for Animal Health (FAO/WHO/WOAH) H5 Working
38 Group, we identified clade-defining mutations for every established clade to enable tree-
39 based clade assignment. We then created three Nextclade datasets which can be used
40 to assign clades to A(H5) HA sequences and call mutations relative to reference strains
41 through a drag-and-drop interface. Nextclade assignments were benchmarked with
42 19,834 unique sequences not in the reference set using a pre-released version of
43 LABEL, a well-validated and widely used command line software. Prospective
44 assignment of new sequences with Nextclade and LABEL produced very well-matched
45 assignments (match rates of 97.8% and 99.1% for the 2.3.2.1 and 2.3.4.4 datasets,
46 respectively). The all-clades dataset also performed well (94.8% match rate) and
47 correctly distinguished between all HPAI and LPAI strains. This tool additionally allows
48 for the identification of polybasic cleavage site sequences and potential N-linked
49 glycosylation sites. These datasets therefore provide an alternative, rapid method to
50 accurately assign clades to new A(H5) HA sequences, with the benefit of an easy-to-use
51 browser interface.

52

53 Introduction

54

55 Highly pathogenic avian influenza (HPAI) A viruses of the H5 subtype continue to pose
56 risks to human and animal health, with transmission in wild and domestic birds seeding
57 periodic cross-species transmission and outbreaks [1, 2]. First detected in China in
58 1996 with the A/Goose/Guangdong/1/1996 (GsGd) strain, these viruses have now
59 spread worldwide and are currently causing the largest panzootic of HPAI in history [3,
60 4]. Highly pathogenic A(H5Nx) viruses—those carrying an H5 hemagglutinin (HA) and
61 various subtypes of NA (e.g. N1, N2, N3, N5, N6, and N8)—have greatly diversified
62 since the emergence of the GsGd lineage, a result of both substitutions in HA and
63 reassortment with other gene segments [5]. Through this accumulation of HA mutations,
64 A(H5Nx) viruses have evolved into distinct phylogenetic groups, known as clades,
65 based on their HA sequence [6].

66

67 A(H5) clades are defined by the Food and Agriculture Organization/World Health
68 Organization/World Organisation for Animal Health (FAO/WHO/WOAH) H5 Working
69 Group, who periodically reviews global A(H5) virus genetic sequence data to determine
70 whether there is sufficient intra-clade diversity to support designating sublineages as
71 new clades [7–11]. Under this nomenclature system, clades are defined as
72 monophyletic groups having <1.5% within-clade average pairwise nucleotide distance,
73 >1.5% distance between other clades, and at least 60% bootstrap support [7]. Currently,
74 high levels of circulating diversity within the 2.3.2.1c and 2.3.4.4 clades necessitate the
75 division of each, yielding subclades 2.3.2.1d–g and 2.3.4.4a–h. In addition to these
76 HPAI clades, there are two lineages of primarily low pathogenic avian influenza (LPAI)
77 viruses that remain endemic in wild birds and are denoted American and Eurasian non-
78 GsGd (Am-nonGsGd and EA-nonGsGd, respectively). This nomenclature system allows
79 for consistent classification and interpretation of surveillance data across laboratories
80 [7]. The ability to classify A(H5) HA sequences provides a framework for studying
81 A(H5Nx) virus evolution, including changes in antigenicity or pathogenicity, host
82 adaptation, and reassortment [7, 12–14]. By enabling researchers to rapidly distinguish
83 the frequency of lineages that are currently circulating among birds and other animals in
84 different regions, this system is critical to viral surveillance and public health monitoring.
85 Clade classifications are therefore crucial for accurately tracking the spread of A(H5Nx)
86 viruses, identifying outbreaks, and improving public health outcomes through more
87 informed decision-making and risk assessment.

88

89 Despite the importance of this nomenclature system, A(H5) clade annotation tools are
90 limited. One popular tool for A(H5) clade assignment is LABEL [15]
91 (<https://wonder.cdc.gov/amd/flu/label/>), which is a command-line software that utilizes a
92 hidden Markov model to classify HA sequences without an alignment step. While using
93 an alignment-free approach allows for rapid clade assignments, this methodology
94 obscures the ability to identify amino acid substitutions and key features of HA—such as
95 the cleavage site sequence—and to assess sequence quality based on the presence of
96 indicators like frameshifts, indels, and premature stop codons. Additionally, LABEL is a
97 command-line software without a graphical user interface, posing a barrier to
98 researchers in public health and surveillance laboratories without bioinformatics

99 expertise or command line access. Alternative web-based tools are available, such as
100 those from the Bacterial and Viral Bioinformatics Resource Center (BV-BRC;
101 <https://www.bv-brc.org>) and The University of Hong Kong
102 (https://tipars.hku.hk/reference/Influenza-A-H5_HA_Tree), which can classify A(H5) HA
103 sequences into clades and perform phylogenetic placement [16]. However, these
104 cannot be used to assess sequence quality, identify amino acid motifs, or call mutations
105 relative to reference strains.
106

107 To address these gaps, we have developed datasets to allow for A(H5) clade
108 assignment with Nextclade, a drag-and-drop browser tool originally developed for
109 SARS-CoV-2 classification that can be run via the command line or via a web interface
110 [17]. Each Nextclade dataset is composed of a reference phylogeny in which each tip
111 and internal node has been annotated with a clade assignment. To perform clade
112 assignments, Nextclade first aligns and trims each query sequence relative to the
113 dataset-specific reference sequence. Next, mutations in each query sequence are
114 compared to the sets of mutations possessed by each node and tip on the tree. After
115 identifying the node or tip with the most similar set of mutations (see
116 <https://docs.nextstrain.org/projects/nextclade/en/stable/user/algorithm/index.html> for
117 algorithm details), its reference clade is assigned to the query sequence and the query
118 sequence is phylogenetically placed relative to this node or tip. This process allows for
119 the underlying reference phylogeny to be utilized in clade assignments, better matching
120 the methodology used in A(H5) clade nomenclature. Nextclade also reports the query
121 sequences' sets of mutations, which can be useful in studying viral evolution—for
122 instance, to identify arising variants or to understand if mammalian-adaptive
123 substitutions are present. Further, sequence quality control scores are generated based
124 on the aforementioned indicators—a feature useful to researchers generating and
125 analyzing their own A(H5) HA sequence data.

126 **Methods**

127

128 **Generation of reference phylogenies**

129 Three reference datasets containing HPAI A(H5) HA sequences with clade annotations
130 were obtained from the WHO/FAO/WOAH H5 Working Group: one that included
131 sequences from each of the 55 historic HPAI A(H5) clades (n=407), one that included
132 sequences from clade 2.3.2.1 and its descendants (2.3.2.1-like and 2.3.2.1a–g;
133 n=2,422), and one that included sequences from clade 2.3.4.4 and its descendants
134 (2.3.4.4-like and 2.3.4.4a–h; n=10,079). To enable assignment of LPAI lineages in
135 addition to established HPAI A(H5) clades, we supplemented this dataset with an
136 additional 10 Am-nonGsGd and 8 EA-nonGsGd sequences from GenBank (accessions
137 available in **Supplementary Table S1**). We next built 3 guide trees from these
138 reference sequences using the Nextstrain pipeline [18]. For the all-clades dataset, we
139 used all sequences in the reference set. For the 2.3.2.1 and 2.3.4.4 datasets, we
140 randomly sampled up to 50 sequences per clade per year and supplemented this
141 dataset with one randomly chosen sequence from each unrepresented clade to serve
142 as outgroup sequences. These outgrouped sequences are included to prevent
143 erroneous Nextclade assignments to the most basal clade (i.e., 2.3.2.1 or 2.3.4.4) when
144 query sequences from other clades are analyzed (see detailed description below).
145 Sequences for each dataset were aligned with MAFFT [19] and divergence phylogenies
146 were inferred using IQ-TREE 2 [20]. TreeTime [21] was used to re-root the phylogeny
147 and to perform ancestral state reconstruction of nucleotide and amino acid sequences
148 across every branch of the tree. The final trees contain 432 (all-clades), 1,654 (2.3.2.1),
149 and 1,900 (2.3.4.4) sequences, including those from outgroups.

150

151 **Determination of clade-defining mutations**

152 To assign query sequences to clades, Nextclade compares each query sequence to
153 every tip and internal node on reference phylogeny to find the closest sequence match.
154 To enable this for our A(H5) datasets, we next inferred clade-defining mutations (i.e.,
155 mutations that are unique to each clade) and used these sets of mutations to assign
156 clades to both internal nodes and tips using Augur clades [22]. For each non-outgroup
157 clade on the reference phylogeny, the last common ancestor (LCA) of all tips belonging
158 to that clade was determined and this branch was used to define the start of the clade.
159 Clade inheritance was established by assessing whether an LCA of one clade
160 descended from the LCA of another clade; this allows for the clade-defining mutations of
161 a parental clade to be inferred for its descendants. Then, sets of clade-defining
162 mutations were determined starting at the root of the tree and working towards the tips.
163 As an LCA was encountered, the set of mutations on that branch was assigned to its
164 respective clade. If a mutation arose on the path to a descendant LCA and was located
165 at the same site as a clade-defining mutation for a parental clade, it was also assigned
166 to the descendent clade to overwrite the inferred parental mutation. For the all-clades
167 tree, which does not include any basal, unassigned outgroup sequences, an amino acid
168 present in the root sequence (HA 17D) was manually assigned to the basal EA-
169 nonGsGd lineage to allow for clade assignment beginning at the root of the tree.

170

171 **Generation and benchmarking of Nextclade datasets**

172 To build the Nextclade datasets, reference phylogenies were paired with appropriate
173 reference HA sequences, including gene maps denoting their coding regions, to be
174 used for alignment of query sequences and subsequent mutation calling. For the all-
175 clades dataset, the ancestral HPAI A/Goose/Guangdong/1/1996 strain was chosen,
176 while the candidate vaccine viruses A/duck/Vietnam/NCVD-1584/2012 and
177 A/Astrakhan/3212/2020 were chosen for the 2.3.2.1 and 2.3.4.4 datasets, respectively
178 [23]. In a config file for each dataset, amino acid motifs of interest were specified to
179 allow for determination of potential N-linked glycosylation sites (PNGS) and cleavage
180 site sequences. Cutoff values for quality control metrics were also specified here, which
181 allow for an estimate of sequence quality based on factors such as the number of
182 frameshifts or private mutations (mutations that separate a query sequence from its
183 attachment site on the tree).

184
185 To benchmark our Nextclade datasets, we acquired a pre-released version (now v0.6.5)
186 of LABEL—a widely used and well-validated command line software for A(H5) clade
187 assignment—that includes the updated, preliminary clade splits. We then generated
188 testing datasets using all unique A(H5) HA sequences available in GISAID and FluDB
189 (accessed April 3, 2024; accessions available in **Supplementary Table S2**) by
190 assigning a clade to each test sequence with both Nextclade and LABEL. For the all-
191 clades dataset, each unique sequence that was not included in the reference dataset
192 was included, while for the 2.3.2.1 and 2.3.4.4 datasets, each unique sequence from the
193 appropriate clades that does not appear on the reference tree was included. These
194 sequences were analyzed with the appropriate Nextclade dataset and assignments
195 were compared to those from LABEL. To determine private mutation quality control
196 values for each dataset, medians and 97.5th percentiles of private mutation counts from
197 these analyses were calculated and used as the ‘typical’ and ‘cutoff’ values,
198 respectively. Additionally, to compare runtimes of the two tools, three randomly selected
199 sets of 100, 1,000, and 10,000 sequences from the all-clades testing dataset were
200 analyzed using the command-line interface for the all-clades Nextclade dataset and
201 LABEL. For each of the three sets of sequences, runtimes were determined in duplicate
202 and the fold-change in mean runtime (LABEL / Nextclade) was calculated. This runtime
203 benchmarking was performed on a 2021 MacBook Pro with a 10-core M1 Pro CPU.

204

205 **Data and code availability**

206 Accessions for the LPAI GenBank sequences used in the reference phylogenies are
207 provided in **Supplementary Table S1** and those for the GISAID sequences used in the
208 testing datasets are provided in **Supplementary Table S2**. Scripts and other files used
209 in the generation of the reference phylogenies and plots are available at
210 <https://github.com/moncla-lab/h5-nextclade>. Additionally, the three Nextclade datasets
211 can be accessed at
212 [https://github.com/nextstrain/nextclade_data/tree/master/data/community/moncla-](https://github.com/nextstrain/nextclade_data/tree/master/data/community/moncla-lab/iav-h5/ha)
213 [lab/iav-h5/ha](https://github.com/nextstrain/nextclade_data/tree/master/data/community/moncla-lab/iav-h5/ha).

214 Results

215

216 Generation of Nextclade datasets with A(H5)-specific features

217 Using WHO/FAO/WOAH-annotated reference sequence datasets—for all A(H5) clades,
218 for clade 2.3.2.1 and its descendants, and for clade 2.3.4.4 and its descendants—we
219 constructed three reference phylogenies. Additionally, for the 2.3.2.1 and 2.3.4.4 trees,
220 one randomly chosen sequence from each unrepresented clade was added to serve as
221 an outgroup sequence (to prevent erroneous downstream assignments to the most
222 basal clade). While the all-clades dataset includes annotation of 2.3.4.4 and 2.3.2.1
223 subclades, using additional datasets specific to these allows for better resolution of the
224 phylogeny for these currently circulating clades. As Nextclade compares query
225 sequences to both tips and internal nodes of the reference phylogeny, we next
226 determined clade-defining mutations based on the common ancestor of all tips of a
227 given clade. These mutations were used to assign clades to internal nodes, yielding
228 reference phylogenies with clade annotations on all tips and nodes, except for the basal
229 outgroup sequences where applicable. Given the high circulating diversity of A(H5)
230 viruses and the desire for these Nextclade datasets to be used on sequences from any
231 time, quality control cutoffs for private mutations (mutations that separate a query
232 sequence from its placement on the reference tree) had to be adjusted for these
233 datasets. Based on analysis of unique sequences from GISAID, the median number of
234 private mutations was used as the ‘typical’ value and the 97.5th percentile as the ‘cutoff’
235 value (24 and 97 for all-clades; 5 and 25 for 2.3.2.1; 6 and 18 for 2.3.4.4)

236 (**Supplementary Figure S1**). If private mutation counts exceed the ‘typical’ value,
237 Nextclade will flag this QC parameter with an increasing score, with a ceiling of 100
238 when the value exceeds the ‘cutoff’ value.

239

240 After alignment, Nextclade will report all nucleotide mutations and amino acid
241 substitutions relative to the dataset’s reference sequence. Using this set of mutations,
242 the most closely related tip or node on the reference tree is determined, and the clade
243 annotation of the tip/node is assigned to the query sequence. Further, the tool will
244 search for N–X–S/T (where X is any amino acid except proline) amino acid motifs in the
245 HA ectodomain of query sequences and annotate each as a PNGS. Additionally, we
246 added in a feature to annotate amino acid sequences at sites 341–346 (all-clades) or
247 341–345 (2.3.2.1 and 2.3.4.4) of HA—corresponding to the polybasic cleavage site
248 motif in each reference sequence—to produce an annotation for the cleavage site, and
249 important determinant of pathogenicity for A(H5Nx) viruses. If four or more arginines
250 and/or lysines are present in the motif, it is also marked as being a polybasic cleavage
251 site.

252

253 Dataset benchmarking and performance

254 To assess the accuracy of clade assignments generated by these Nextclade datasets,
255 we utilized LABEL v0.6.5 for benchmarking. Given that this is a pre-released version and
256 the new subclade designations are not yet finalized by the WHO/FAO/WOAH, it is
257 possible that there are inaccuracies in its clade assignments; however, it allows us to
258 compare assignments for any A(H5) HA sequences, including those not officially
259 annotated, rather than a small subset of officially annotated sequences. Using all

260 available unique HA sequences that were not found in the reference dataset,
261 comparisons were made between Nextclade and LABEL clade assignments. The all-
262 clades dataset performed well, with a 94.8% match rate (n=19,833) (**Figure 1A**); of note,
263 no mismatches between HPAI A(H5) clades and LPAI A(H5) lineages were observed.
264 Performance was best for clades with many available sequences, presumably due to
265 better representation of within-clade diversity on the reference tree. The 2.3.2.1 and
266 2.3.4.4 datasets produced very well-matched assignments, with match rates of 97.8%
267 and 99.1%, respectively (**Figure 1B, 1C**). Similarly to the all-clades dataset, these
268 generally performed best for clades that were well represented on the reference tree.
269 Additionally, poor resolution of the 2.3.2.1c-like and 2.3.4.4-like clades, which act as
270 outlier groups and represent sequences that do not fall within the newly designated
271 subclades, appears to play a large role in the mismatches seen in these datasets. Given
272 the role of the “like” clades (i.e., to capture sequences that do not neatly cluster within a
273 defined subclade), different approaches vary in the stringency with which sequences may
274 be allocated into these clades, which we expect may account for many mismatches
275 observed between LABEL and Nextclade. To compare runtimes, randomly selected sets
276 of 100, 1,000, and 10,000 sequences were analyzed with the all-clades Nextclade dataset
277 (via the command-line interface) and LABEL on a laptop computer, and the total
278 execution times were determined. Nextclade very quickly analyzed large queries,
279 performing 10,000 assignments in under 20 seconds (**Figure 2A**). With these rapid
280 execution times, Nextclade greatly outperformed LABEL for each of the testing sets, with
281 the runtime of LABEL being approximately 100× longer than Nextclade when 1,000 or
282 more sequences were analyzed (**Figure 2B**).

283 Discussion

284

285 We have generated three Nextclade datasets that can be used to assign clades to both
286 HPAI and LPAI A(H5) sequences. Using WHO/FAO/WOAH-annotated reference
287 sequences, we generated reference phylogenies and determined clade-defining
288 mutations for each represented clade. The all-clades dataset can be used for all
289 historical clades and will call mutations relative to the ancestral HPAI GsGd strain, while
290 the 2.3.2.1 and 2.3.4.4 datasets are specific to currently circulating clade groups and
291 will call mutations relative to candidate vaccine viruses. Nextclade will also
292 automatically call the cleavage site sequence—and determine if a polybasic cleavage
293 site is present—as well as locate potential N-linked glycosylation sites in the HA
294 ectodomain. These features provide users with further details about their sequences,
295 including those useful in identifying potential alterations to pathogenicity or antigenicity
296 [12–14, 24, 25]. While other available clade assignment tools can provide clade
297 assignments, we believe that Nextclade's user-friendly interface, rapid execution, and
298 additional features will enable a broad audience to fully leverage their A(H5) sequence
299 data. We demonstrate that tree-based classification of avian influenza viruses is
300 possible with Nextclade, thus future work should include the expansion of this tool to
301 accommodate other avian influenza clade systems, such as A(H7) and A(H9).
302 Additionally, updating these A(H5) datasets upon the designation of new clades will be
303 critical to their continued utility. As the reference datasets can be easily amended to
304 reflect clade updates—through the reassignment of clades to reference sequences or
305 the addition of newly annotated sequences—Nextclade is poised to rapidly respond to
306 nomenclature changes as they arise.

307

308 Acknowledgements

309

310 We gratefully acknowledge labs that have generated sequence data deposited into
311 public databases, and the WHO/FAO/WOAH H5 Working Group for continuous updates
312 of clade designations. The findings and conclusions in this report are those of the
313 authors and do not necessarily represent the views of the Centers for Disease Control
314 and Prevention or the Agency for Toxic Substances and Disease Registry.

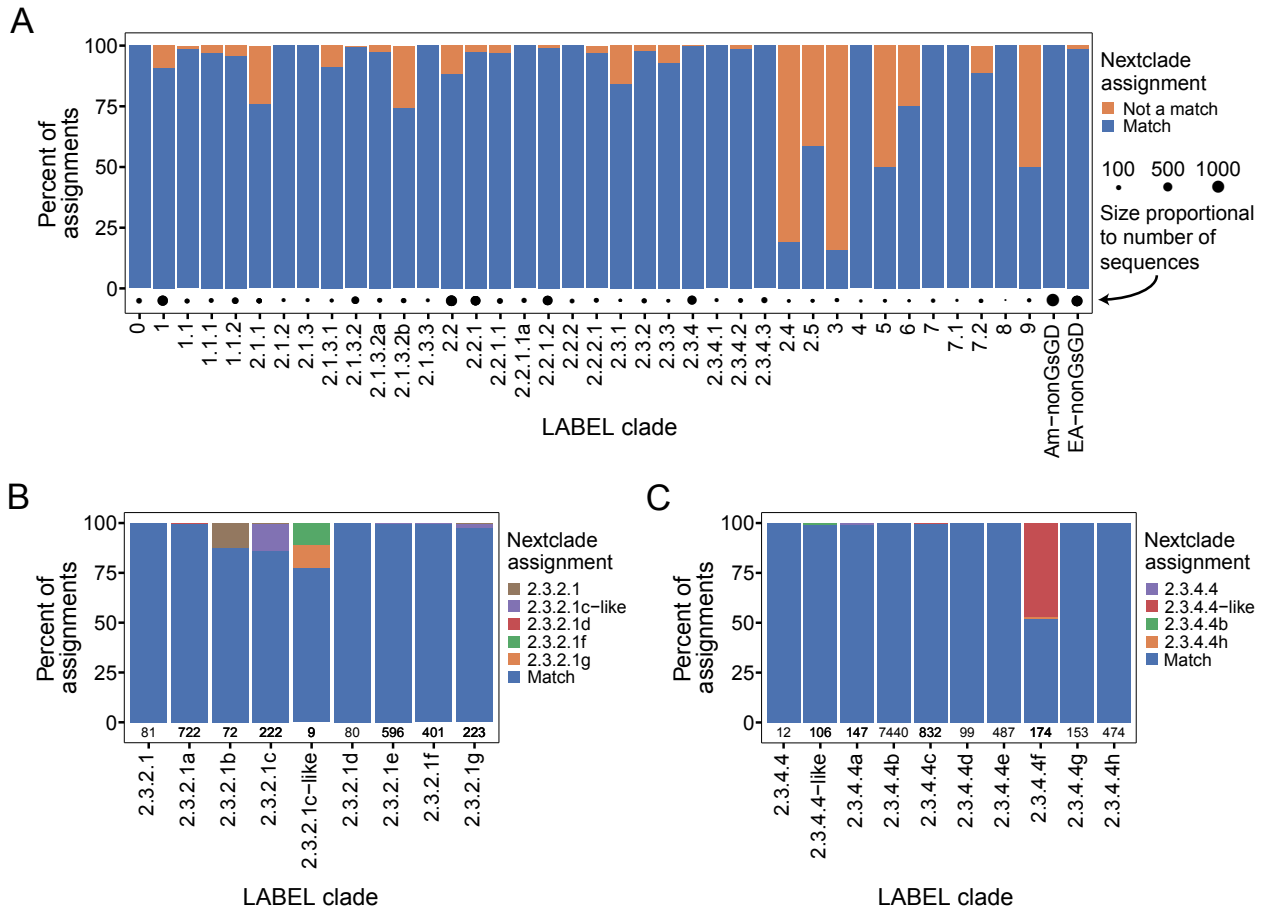
315

316 Funding

317

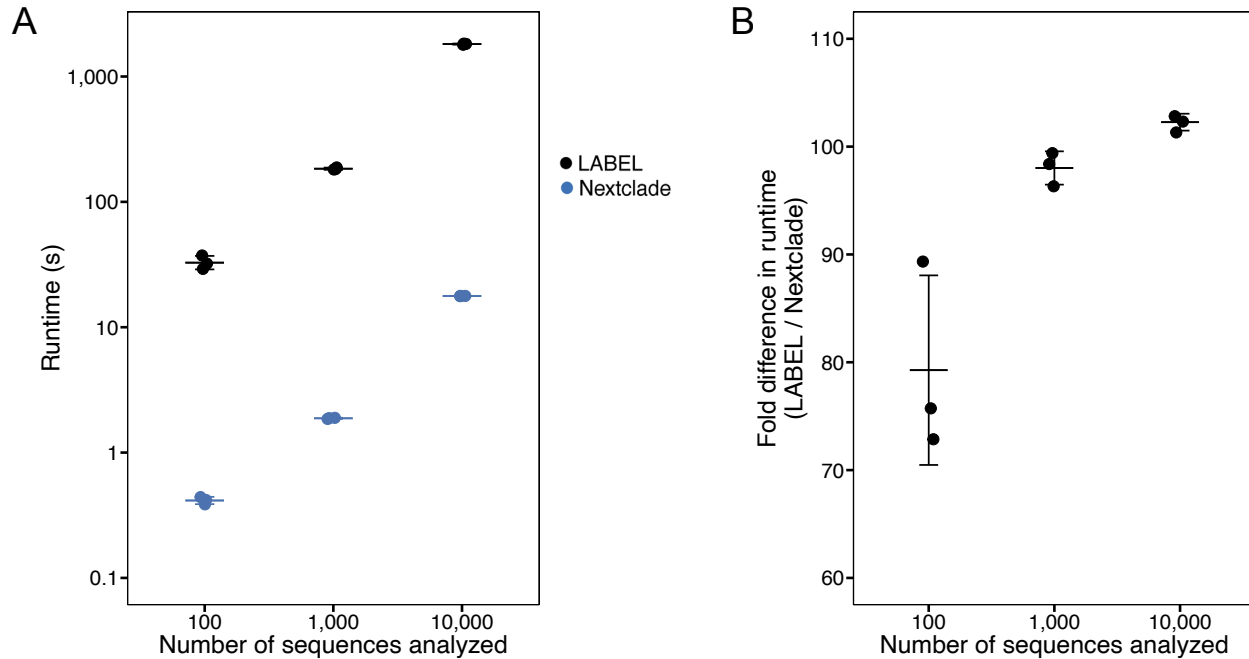
318 This work was supported by the National Institute for Allergy and Infectious Diseases at
319 the National Institutes of Health [grant number R00-AI147029-05], and by the Centers of
320 Excellence for Influenza Research and Response Computational Modeling Core,
321 funded by the National Institute for Allergy and Infectious Diseases at the National
322 Institutes of Health [grant number 75N93021C00015]. LHM is a Pew Biomedical
323 Scholar and is supported by the National Institutes of Health [grant number R00-
324 AI147029-05]. JTO is supported by the National Institutes of Health [grant number
325 75N93021C00015]. TTYL is supported by InnoHK funding from Innovation and
326 Technology Commission of Hong Kong SAR Government and Seed Funding for
327 Strategic Interdisciplinary Research Scheme from University Research Committee
328 (URC) [project number 102010190].

329 **Figures**
330



331
332

333 **Figure 1. Clade assignments are well-matched between our Nextclade datasets**
334 **and LABEL.** Results for benchmarking against LABEL are shown for each dataset. (A)
335 Results for the all-clades dataset (94.8% match rate, n=19,833; results for 2.3.2.1 and
336 2.3.4.4 descendants not shown), with blue indicating matched and orange indicating
337 mismatched assignments. (B) Results for the clade 2.3.2.1 dataset (97.8% match rate,
338 n=2,406), with blue indicating matched assignments and mismatched assignments
339 shown colored by the Nextclade assignment (2.3.2.1g in orange, 2.3.2.1f in green,
340 2.3.2.1d in red, 2.3.2.1c-like in purple, and 2.3.2.1 in brown). (C) Results for the clade
341 2.3.4.4 dataset (99.1% match rate, n=9,924), with blue indicating matched assignments
342 and mismatched assignments shown colored by the Nextclade assignment (2.3.4.4h in
343 orange, 2.3.4.4b in green, 2.3.4.4-like in red, and 2.3.4.4 in purple).



344
345

346 **Figure 2. Nextclade analyzes sequences more rapidly than LABEL.** Three randomly
347 selected sets of 100, 1,000, and 10,000 sequences were annotated via the command-
348 line interface of the all-clades Nextclade dataset and LABEL in duplicate. (A) Mean
349 absolute runtimes were determined in seconds (s) and (B) fold differences in means
350 were calculated, with each dot representing the value for a unique set of sequences and
351 the mean and standard deviation shown by error bars.

- 352 1. Krammer F, Schultz-Cherry S. We need to keep an eye on avian influenza. *Nat Rev*
353 *Immunol.* 2023;23:267–8.
- 354 2. The Global Consortium for H5N8 and Related Influenza Viruses. Role for migratory
355 wild birds in the global spread of avian influenza H5N8. *Science.* 2016;354:213–7.
- 356 3. Xu X, Subbarao K, Cox NJ, Guo Y. Genetic Characterization of the Pathogenic
357 Influenza A/Goose/Guangdong/1/96 (H5N1) Virus: Similarity of Its Hemagglutinin Gene
358 to Those of H5N1 Viruses from the 1997 Outbreaks in Hong Kong. *Virology.*
359 1999;261:15–9.
- 360 4. Ramey AM, Hill NJ, DeLiberto TJ, Gibbs SEJ, Camille Hopkins M, Lang AS, et al.
361 Highly pathogenic avian influenza is an emerging disease threat to wild birds in North
362 America. *J Wildl Manag.* 2022;86:e22171.
- 363 5. El-Shesheny R, Franks J, Turner J, Seiler P, Walker D, Friedman K, et al. Continued
364 Evolution of H5Nx Avian Influenza Viruses in Bangladeshi Live Poultry Markets:
365 Pathogenic Potential in Poultry and Mammalian Models. *J Virol.* 2020;94:e01141-20.
- 366 6. Cui P, Shi J, Wang C, Zhang Y, Xing X, Kong H, et al. Global dissemination of H5N1
367 influenza viruses bearing the clade 2.3.4.4b HA gene and biologic analysis of the ones
368 detected in China. *Emerging Microbes & Infections.* 2022;11:1693–704.
- 369 7. WHO/OIE/FAO H5N1 Evolution Working Group. Toward a Unified Nomenclature
370 System for Highly Pathogenic Avian Influenza Virus (H5N1). *Emerg Infect Dis.*
371 2008;14:e1–e1.
- 372 8. Who/Oie/Fao H5N1 Evolution Working Group. Continuing progress towards a unified
373 nomenclature for the highly pathogenic H5N1 avian influenza viruses: divergence of
374 clade 2·2 viruses. *Influenza Resp Viruses.* 2009;3:59–62.
- 375 9. WHO/OIE/FAO H5N1 Evolution Working Group. Continued evolution of highly
376 pathogenic avian influenza A (H5N1): updated nomenclature. *Influenza Resp Viruses.*
377 2012;6:1–5.
- 378 10. World Health Organization/World Organisation for Animal Health/Food and
379 Agriculture Organization (WHO/OIE/FAO) H5N1 Evolution Working Group. Revised and
380 updated nomenclature for highly pathogenic avian influenza A (H5N1) viruses. *Influenza*
381 *Resp Viruses.* 2014;8:384–8.
- 382 11. WHO/OEI/FAO H5 Evolution Working Group. Evolution of the influenza A(H5)
383 haemagglutinin: WHO/OIE/FAO H5 Working Group reports a new clade designated
384 2.3.4.4. 2015.
- 385 12. Herfst S, Schrauwen EJA, Linster M, Chutinimitkul S, De Wit E, Munster VJ, et al.
386 Airborne Transmission of Influenza A/H5N1 Virus Between Ferrets. *Science.*
387 2012;336:1534–41.
- 388 13. Imai M, Watanabe T, Hatta M, Das SC, Ozawa M, Shinya K, et al. Experimental
389 adaptation of an influenza H5 HA confers respiratory droplet transmission to a
390 reassortant H5 HA/H1N1 virus in ferrets. *Nature.* 2012;486:420–8.
- 391 14. Kong H, Burke DF, Da Silva Lopes TJ, Takada K, Imai M, Zhong G, et al. Plasticity
392 of the Influenza Virus H5 HA Protein. *mBio.* 2021;12:e03324-20.
- 393 15. Shepard SS, Davis CT, Bahl J, Rivaille P, York IA, Donis RO. LABEL: Fast and
394 Accurate Lineage Assignment with Assessment of H5N1 and H9N2 Influenza A
395 Hemagglutinins. *PLoS ONE.* 2014;9:e86921.

- 396 16. Ye Y, Shum MH, Tsui JL, Yu G, Smith DK, Zhu H, et al. Robust expansion of
397 phylogeny for fast-growing genome sequence data. *PLoS Comput Biol*.
398 2024;20:e1011871.
- 399 17. Aksamentov I, Roemer C, Hodcroft E, Neher R. Nextclade: clade assignment,
400 mutation calling and quality control for viral genomes. *JOSS*. 2021;6:3773.
- 401 18. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al. Nextstrain:
402 real-time tracking of pathogen evolution. *Bioinformatics*. 2018;34:4121–3.
- 403 19. Katoh K. MAFFT: a novel method for rapid multiple sequence alignment based on
404 fast Fourier transform. *Nucleic Acids Research*. 2002;30:3059–66.
- 405 20. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Von Haeseler A,
406 et al. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the
407 Genomic Era. *Molecular Biology and Evolution*. 2020;37:1530–4.
- 408 21. Sagulenko P, Puller V, Neher RA. TreeTime: Maximum-likelihood phylodynamic
409 analysis. *Virus Evolution*. 2018;4.
- 410 22. Huddleston J, Hadfield J, Sibley T, Lee J, Fay K, Ilcisin M, et al. Augur: a
411 bioinformatics toolkit for phylogenetic analyses of human pathogens. *JOSS*.
412 2021;6:2906.
- 413 23. World Health Organization. Genetic and antigenic characteristics of zoonotic
414 influenza A viruses and development of candidate vaccine viruses for pandemic
415 preparedness. 2023.
- 416 24. Luczo JM, Stambas J, Durr PA, Michalski WP, Bingham J. Molecular pathogenesis
417 of H5 highly pathogenic avian influenza: the role of the haemagglutinin cleavage site
418 motif. *Reviews in Medical Virology*. 2015;25:406–30.
- 419 25. Dadonaite B, Ahn JJ, Ort JT, Yu J, Furey C, Dosey A, et al. Deep mutational
420 scanning of H5 hemagglutinin to inform influenza virus surveillance. 2024.
421