

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

Sequence amplification via cell passaging creates spurious signals of positive adaptation in influenza H3N2 hemagglutinin

Claire D. McWhite^{1,2}, Austin G. Meyer^{1,3,4}, Claus O. Wilke^{1,3,4*}

¹Center for Systems and Synthetic Biology and Institute for Cellular and Molecular Biology, The University of Texas at Austin, Austin, TX 78712

²Department of Molecular Biosciences, The University of Texas at Austin, Austin, TX 78712

³Center for Computational Biology and Bioinformatics, The University of Texas at Austin, Austin, TX 78712

⁴Department of Integrative Biology, The University of Texas at Austin, Austin, TX 78712

Address correspondence to wilke@austin.utexas.edu

23 Clinical influenza A isolates are frequently not sequenced directly. Instead, a majority of
24 these isolates (~70% in 2015) are first subjected to passaging for amplification, most
25 commonly in non-human cell culture. Here, we find that this passaging leaves distinct
26 signals of adaptation in the viral sequences, which can confound evolutionary analyses
27 of the viral sequences. We find distinct patterns of adaptation to Madin-Darby (MDCK)
28 and monkey cell culture absent from unpassaged hemagglutinin sequences. These
29 patterns also dominate pooled datasets not separated by passaging type, and they
30 increase in proportion to the number of passages performed. By contrast, MDCK-SIAT1
31 passaged sequences seem mostly (but not entirely) free of passaging adaptations.
32 Contrary to previous studies, we find that using only internal branches of the influenza
33 phylogenetic trees is insufficient to correct for passaging artifacts. These artifacts can
34 only be safely avoided by excluding passaged sequences entirely from subsequent
35 analysis. We conclude that future influenza evolutionary analyses should appropriately
36 control for potentially confounding effects of passaging adaptations.

37

38

39

40

41

42

43

44

45 1. INTRODUCTION

46 The routine sequencing of clinical isolates has become a critical component of global
47 seasonal influenza surveillance (World Health Organization Global influenza
48 surveillance network, 2011). Analysis of these viral sequences informs the selection of
49 future vaccine strains (Stöhr et al., 2012; WHO Writing Group et al., 2012), and a wide
50 variety of computational methods have been developed to identify sites under selection
51 or immune-escape mutations (Blackburne et al., 2008; Koelle et al., 2006; Nelson et al.,
52 2006; Suzuki, 2008; Wolf et al., 2006), or to predict the short-term evolutionary future of
53 influenza virus (Łuksza and Lässig, 2014; Neher et al., 2014). However, sites that
54 appear positively selected in sequence analysis frequently do not agree with sites
55 identified experimentally in hemagglutination inhibition assays (Kratsch et al., 2016;
56 Meyer and Wilke, 2015; Tusche et al., 2012), and the origin of this discrepancy is
57 unclear. Here, we argue that a major cause of this discrepancy is widespread passaging
58 of influenza virus before sequencing.

59

60 Clinical isolates are often passaged in culture one or more times to amplify viral copy
61 number, as well as to introduce virus into a living system for testing strain features such
62 as vaccine response, antiviral response, and replication efficiency (Kumar and
63 Henrickson, 2012; World Health Organization Global influenza surveillance network,
64 2011). A variety of culture systems are used for virus amplification. Cell cultures derived
65 from Madin-Darby canine kidney (MDCK) cells are by far the most widely used system,
66 with the majority of sequences in influenza repositories deriving from virus that has
67 been passaged through an MDCK or modified MDCK cell culture (Balish et al., 2005;

68 Bogner et al., 2006). Influenza virus may also be passaged through monkey kidney
69 (RhMK or TMK) cell culture or injected directly into egg amniotes. Alternatively,
70 complete influenza genomes can be obtained from PCR-amplified influenza samples
71 without intermediate passaging (Katz et al., 1990; Lee et al., 2013a).

72

73 Several experiments have demonstrated that influenza virus hemagglutinin (HA)
74 accumulates mutations following rounds of passaging in both cell (Ilyushina et al., 2012;
75 Lee et al., 2013b; Wyde et al., 1977) and egg culture (Robertson et al., 1993). The
76 decreased number of mutations in MDCK-based cell culture is the main argument for
77 use of this system over egg amniotes in vaccine production (Katz and Webster, 1989),
78 with MDCK cells expressing human SIAT1 having the highest fidelity to the original
79 sequence and reduced host adaptation (Hamamoto et al., 2013). Viral adaptations to
80 eggs were recently linked to reduced vaccine efficacy (Skowronski et al., 2014; Xie et
81 al., 2015) and were implicated as potentially contributing to reduced efficacy of 2014-
82 2015 seasonal H3N2 influenza vaccination in the World Health Organization's
83 recommendations for 2015-2016 vaccine strains (The World Health Organization,
84 2015). As the majority of influenza vaccines worldwide are produced in eggs, vaccine
85 strain selection is limited to virus with the ability to replicate rapidly in this system (World
86 Health Organization Global influenza surveillance network, 2011).

87

88 Although egg-passaged sequences are increasingly excluded from influenza virus
89 phylogenetic analysis (see e.g. the NextFlu tracker (Neher and Bedford, 2015)), due to
90 the known high host-specific substitution rates, cell culture is generally not thought to be

91 sufficiently selective to produce a discernable evolutionary signal. One of few existing
92 evolutionary analyses of passaging effects on influenza virus (Bush et al., 2000) found
93 that passaging caused no major changes in clade structure between egg and cell
94 passaged viruses. However, several studies have recommended the use of internal
95 branches in the phylogenetic tree to reduce passaging effects in evolutionary analysis of
96 Influenza A (Bush et al., 2001; Suzuki, 2006). Another study discovered egg culture to
97 be the cause of misidentification of several sites under positive selection in Influenza B
98 (Gatherer, 2010), but this study was limited to comparing egg-cultured to cell-cultured
99 virus. As the availability of unpassaged influenza sequences has dramatically increased
100 over the past ten years, we can now perform a direct comparison of passaged to
101 circulating virus.

102

103 Here, we compare patterns of adaptation in North American seasonal H3N2 influenza
104 HA sequences derived from passaged and unpassaged virus. We divide viral
105 sequences by their passaging history, distinguishing between unpassaged clinical
106 samples, egg amniotes, RhMK (monkey) cell culture, and generic/MDCK-based cell
107 culture. For the latter, we also distinguish between virus passaged in MDCK-SIAT1 cell
108 culture (SIAT1) and in unmodified MDCK or unspecified cell culture (non-SIAT1). We
109 find clear signals of adaptation to the various passaging conditions, and demonstrate
110 that passaging artifacts become more severe with additional rounds of passaging.
111 These signals are strongly present in the tip branches of the phylogenetic trees but can
112 also be detected in internal and trunk branches. We demonstrate the accumulation of
113 these passaging artifacts with additional rounds of serial passaging in non-SIAT1 cells.

114 Finally, we demonstrate that the identification of antigenic escape sites from sequence
115 data has been confounded by passaging adaptations, and that the exclusion of
116 passaged sequences allows us to use sequence and structural data to highlight regions
117 involved in antigenic escape.

118

119 **2. METHODS**

120 **2.1 Influenza sequence data**

121 Non-laboratory strain H3N2 hemagglutinin (HA) sequences collected in North America
122 were downloaded from The Global Initiative for Sharing Avian Influenza Data (GISAID)
123 (Bogner et al., 2006) for the 1968–2015 influenza seasons. We used exclusively North
124 American sequences to reduce regional variation between influenza virus strains. Non-
125 complete HA sequences were excluded. Sequences were trimmed to open reading
126 frames, filtered to remove redundancies, and aligned by translation–alignment–back-
127 translation using MAFFT (Kato and Standley, 2013) for the alignment step. Sequence
128 headers of FASTA files were standardized into an uppercase text format with non-
129 alphanumeric characters replaced by underscores. As H3N2 strains have experienced
130 no persistent insertion or deletion events, we deleted sequences that introduced gaps to
131 the alignment. To ascertain overall data quality, we built a phylogenetic tree of the entire
132 sequence set (using FastTree 2.0 compiled for short branch lengths (Price et al., 2010))
133 and checked for any abnormal clades or other unexpected tree features. We found one
134 abnormal clade of approximately 20 sequences with an exceptionally long branch length
135 (> 0.01) and removed the sequences in that clade from further analysis. Our final
136 dataset consisted of 6873 sequences from 2005–2015 as well as one outgroup of 45

137 sequences from 1968–1977 (not considered for further analysis). We did not consider
138 sequences collected from 1978-2004.

139

140 **2.2 Identification of passage history and evolutionary-rate calculations**

141 We divided sequences into groups by their passage-history annotation and collection
142 year, determining passage history by parsing with regular expressions for keywords in
143 FASTA headers (Table 1). We classified 1133 sequences with indeterminate or missing
144 passage histories, or passage through multiple categories of hosts (i.e. both egg and
145 cell), as “other”. The final datasets for individual passage groups contained between 79
146 and 3041 sequences (Table 1).

147

148 We additionally divided passage groups into singly and serially passaged subgroups.
149 Sequences matching the regular expression
150 "2|3|4|5|6|1_C|X_|X_C|AND_MV1|X_S|1_S|CX|MX|C1S1|EX|X_E|MIX_RHMK|RII" were
151 classified as having been passaged two or more times. All remaining sequences were
152 passaged only once.

153

154 To determine the number of times a sequence was passaged, we used a different set of
155 regular expressions, which we applied only to non-SIAT1 cell-passaged sequences. We
156 first excluded sequences with an indeterminate number of passages by excluding
157 sequences whose record IDs matched the regular expression
158 "^X_\$|DETAILS__MDCK|MX_C|^X_C1_|^MX_\$|CX_C1|X_C1|DETAILS__ND". We
159 then collected multiply passaged sequences using the regular expression

160 "3|1_C2|2_C1|M1M1_C1|1_MDCK2|4|2_C2|3_C1|1_C3|5|2_C3|3_C2" and doubly
161 passaged sequences using the regular expression "2|1_C1". The remaining sequences
162 were only passaged a single time.

163

164 We next constructed phylogenetic trees for each passage group as well as one tree for
165 a pooled dataset combining all individual passage groups and other sequences. All
166 phylogenetic trees were constructed using FastTree 2.0 (Price et al., 2010). We
167 calculated site-specific dN/dS values using a one-rate SLAC (Single-Likelihood
168 Ancestor Counting) model implemented in HyPhy (Pond et al., 2005). One rate models
169 fit a site-specific dN and a global dS , where the global dS is the mean site-wise dS for a
170 given condition (Spielman et al., 2015). Among different one-rate, site-specific models,
171 SLAC performs nearly identically to other approaches (Spielman et al., 2015), and it
172 was chosen here due to its speed and ease of extracting dN/dS estimates along internal
173 and tip branches. To obtain internal and tip branch-specific estimates, we extracted the
174 dN/dS values calculated by the SLAC algorithm. We manually calculated dN/dS along
175 trunk branches by counting the number of synonymous and non-synonymous
176 substitutions at each trunk site. We defined the trunk as the sequence of branches from
177 the root to the penultimate node before a randomly chosen terminal sequence from the
178 most recent year represented in the tree.

179

180 We chose sequences from 2005-2015 as our sample set due the low number of
181 available sequences prior to this period. As dN/dS estimates can be confounded by
182 sample size (Spielman et al., 2015), we sought to limit this effect by down-sampling

183 each experimental set to match the number of sequences in the smallest group being
184 considered in a particular analysis (Table 1). To reduce season-to-season variation in
185 the comparison of unpassaged, SIAT1, and non-SIAT1 cell culture, we performed one
186 analysis with sequences from only 2014, which is the year that maximizes sequences
187 available from all three conditions ($n = 249$ each).

188

189 **2.3 Geometric analysis of dN/dS distributions**

190 For each amino acid site i in HA, we computed the inverse Euclidean distance to each
191 amino acid site j ($j \neq i$) in the 3D crystal structure. For each site i , we then correlated the
192 list of inverse distances to sites j with site-wise dN/dS values at sites j . This procedure
193 resulted in a correlation coefficient for each site i , and we mapped these correlation
194 coefficients onto the corresponding sites i in the 3D structure model of HA. In this
195 analysis, sites spatially closest to positively selected regions in the protein yielded the
196 highest correlation coefficients. Thus, this approach allowed us to visualize regions of
197 increased positive selection. As the correlation coefficient for site $i = 224$ is consistently
198 highest for sets of sequences undivided by passage history, we chose this site as a
199 reference to compare passage conditions. See Meyer and Wilke (2015) for additional
200 discussion of this approach.

201

202 We processed the HA PDB structure to allow for easy alignment with site-wise
203 measures as discussed previously (Meyer and Wilke, 2015). We provided a
204 renumbered and formatted H3N2 HA structure derived from PDBID:2YP7
205 (2YP7clean.pdb) (Lin et al., 2012) with our data analysis code (see below). Noting that

206 the hemagglutinin protein and gene numbering is offset by 16, all site numbering in this
207 manuscript refers to the protein site position. The alignment of gene and protein
208 numbering schemes to amino acid sequence is available in each supplementary data
209 file.

210

211 **2.4 Local Branching Index analysis**

212 We used the framework and code (<https://github.com/rneher/FitnessInference>) from
213 Neher et al. (2014) to rank sequences according to their Local Branching Index (LBI), a
214 metric that uses branching density to predict progenitor lineages. To build our sample
215 set we divided sequences by year and passage history. We then down-sampled
216 alignments to 70% of the available sequences, up to a maximum of 100 sequences. We
217 repeated each down-sampling fifty times for each condition. We then ranked
218 sequences in each sample according to the LBI algorithm (script rank_sequences.py
219 available from <https://github.com/rneher/FitnessInference>) and calculated the hamming
220 distance of the top ranked sequence from each condition to the ancestrally
221 reconstructed root sequence of the following year's unpassaged and pooled trees. The
222 hamming distance derived from the top ranked sequence was divided by the hamming
223 distance of a randomly chosen sequence from the same condition, to assess if a
224 predicted progenitor was better than a randomly chosen sequence. These ratios were
225 averaged over all possible choices of the randomly chosen sequence and over the fifty
226 trials, to yield the mean ratio score for a particular year and passage condition. A mean
227 ratio score < 1 indicates that the LBI algorithm performs better than random chance.

228

229 **2.5 Statistical analysis and data availability**

230 Raw influenza sequences used in this analysis are available for download from GISAID
231 (<http://gisaid.org>) using the parameters “North America”, “H3N2”, “1976 – 2015”.

232 Acknowledgements for sequences used in this study are available in Supplementary
233 File 11. The complete, processed dataset used in our statistical analysis is available in
234 Supplementary Data 10, including protein and gene numbering, computed evolutionary
235 rates, relative solvent accessibility for the hemagglutinin trimer, and site-wise distance
236 to protein site 224. Relative solvent accessibility of the hemagglutinin trimer was taken
237 from Meyer and Wilke (2015). Site-wise Euclidean distances between all amino acids in
238 the HA structure PDBID:2YP7 were recalculated from structural coordinates using the
239 script `distances.py` from:

240 https://github.com/wilkelab/influenza_H3N2_passaging/tree/master/scripts. Statistical
241 analysis was performed using R (Ihaka and Gentleman, 1996), and all graph figures
242 were drawn with the R package `ggplot2` (Wickham, 2009). Throughout this work, *
243 denotes a significance of $0.01 \leq P < 0.05$, ** denotes a significance of $0.001 \leq P < 0.01$,
244 and *** denotes a significance of $P < 0.001$.

245
246 Linear models between site-wise dN/dS and RSA or inverse distance were fit using the
247 `lm()` function in R. Correlations were calculated using the R function `cor()` and
248 significance determined using `cor.test()`.

249

250 Our entire analysis pipeline, instructions for running analyses, and raw data (except
251 initial sequence data per the GISAID user agreement) are available at the following
252 Github project repository:

253 https://github.com/wilkelab/influenza_H3N2_passaging.

254

255 **3. RESULTS**

256 Many influenza-virus samples collected from patients are first passaged through one or
257 more culturing systems (Table 1) prior to PCR amplification and sequencing (Figure
258 1A). Samples may be passaged either once or serially (Table 1), even though a single
259 passage is generally sufficient to obtain adequate amounts of viral DNA for sequencing.
260 Reconstructed trees of influenza evolution contain a mixture of passage histories at
261 their tips (Figure 1B). During passaging, influenza genomes accumulate adaptive
262 mutations, and the effect of these mutations on evolutionary analyses of influenza
263 sequences is not well understood.

264

265 **3.1 Site-wise evolutionary rate patterns differ between passage groups**

266 To quantify any evolutionary signal that may be introduced by passaging, we
267 assembled, from the GISAID database (Bogner et al., 2006), a set of North American
268 human influenza H3N2 hemagglutinin sequences collected between 2005 and 2015.
269 We initially sorted these sequences into groups by their passage history: (1)
270 unpassaged, (2) egg-passaged, (3) generic cell-passaged, and (4) monkey cell-
271 passaged (Table 1). To assess evolutionary variation at individual sites, we calculated
272 site-specific dN/dS (Echave et al., 2016), using Single Likelihood Ancestor Counting

273 (SLAC). Specifically, we calculated one-rate dN/dS estimates, i.e., site-specific dN
274 values normalized by a global dS value (see Methods for details). In addition to
275 considering groups of sequences with specific passage histories, we also calculated
276 dN/dS values by pooling all sequences into one combined analysis. This pooled group
277 corresponds to a typical influenza evolutionary analysis for which passage history has
278 not been accounted.

279
280 We first correlated the site-wise dN/dS values we obtained for virus sequences derived
281 from different passage histories. If passage history did not matter, then the dN/dS
282 values obtained from different sources should have correlated strongly with each other,
283 with r approaching 1. Instead, we found that correlation coefficients ranged from 0.68 to
284 0.87, depending on which specific comparison we made (Figure 2A). In this analysis,
285 and throughout this work, we down-sampled alignments to the smallest number of
286 sequences available for any of the conditions compared, to keep the samples as
287 comparable as possible overall. The analysis of Figure 2 used $n = 917$ randomly drawn
288 sequences for each condition. Unpassed dN/dS correlated more strongly with cell
289 and pooled dN/dS (correlations of 0.77 and 0.79, respectively) than with monkey-cell
290 dN/dS (0.68). Note that the dN/dS values from the pooled group, which corresponds to
291 a typical dataset used in a phylogenetic analysis of influenza, more closely correlated
292 with the dN/dS values from the generic cell group ($r = 0.87$) than from the unpassed
293 group ($r = 0.79$). Egg-derived sequences were excluded from this analysis due to low
294 sequence numbers ($n = 79$), however evolutionary rates from this condition correlated
295 particularly poorly with those of random draws of 79 unpassed sequences

296 (Supplementary Figure 1). This result is consistent with previous conclusions (Bush et
297 al., 2000; Gatherer, 2010; Suzuki, 2006) that egg-derived sequences show specific
298 adaptations not found otherwise in influenza sequences.

299
300 Because the common ancestor of any two passaged influenza viruses is a virus that
301 replicated in humans, we expected that any adaptations introduced during passaging
302 would not extend into the internal branches of a reconstructed tree. Therefore, we
303 additionally subdivided phylogenetic trees into internal branches and tip branches, and
304 calculated site-specific dN/dS values separately for these two sets of branches. In fact,
305 Bush et al. (2000) recommended the use of internal branches to reduce variation seen
306 between egg and cell culture-passaged virus. As expected, we found that when dN/dS
307 calculations were restricted to the internal branches, the correlations between the
308 passage groups increased overall (Figure 2B), even though distinct differences between
309 the passage groups remained. Conversely, when we only considered tip branches,
310 correlations among most groups were relatively low (Figure 2C), with the exception of
311 cell-passaged sequences compared to the pooled sequences. This finding emphasizes
312 once again that the pooled sample is most similar to the cell-passaged sample. We
313 conclude that different passaging histories leave distinct, evolutionary signatures of
314 adaptation to the passaging environment.

315
316 In aggregate, these results show that both generic-cell-passaged sequences and
317 monkey-cell-passaged sequences yield different site-wise dN/dS patterns relative to
318 unpassaged sequences (Fig. 2A-C), with dN/dS values derived from monkey-cell-

319 passaged sequences being the least similar to dN/dS from unpassaged sequences (Fig.
320 2A–C). The pooled group of sequences, which corresponds to a typical dataset used in
321 evolutionary analyses of influenza virus, describes evolutionary rates of specifically cell-
322 passaged virus and poorly matches evolutionary rates of unpassaged virus.

323

324 **3.2 Adaptations to cell and monkey-cell passage display characteristic patterns of** 325 **site variation**

326 We next asked whether adaptations to passage history were located in specific regions
327 of the hemagglutinin (HA) protein. To address this question, we employed the geometric
328 model of HA evolution we recently introduced (Meyer and Wilke, 2015), where structural
329 measurements explain variation in dN/dS . For H3N2 HA, this model explains over 30%
330 of the variation in dN/dS using two simple physical measures, the relative solvent
331 accessibility (RSA) of individual residues in the structure (Tien et al., 2013) and the
332 inverse linear distance in 3D space from each residue to protein site 224 in the
333 hemagglutinin monomer. Notably, the geometric model was previously applied to a
334 pooled sequence set including sequences of various passaging histories. To what
335 extent it carries over to sequences with specific passaging histories is not known.

336

337 We first considered the correlation between dN/dS and RSA (Figure 3A). We found that
338 for all passage groups, R^2 values ranged from 0.10 to 0.16 in the full tree, consistent
339 with our earlier work (Meyer and Wilke, 2015). The high congruence among R^2 values
340 for internal branches and all branches suggests that RSA imposes a pervasive selection
341 pressure on HA, independent of passaging adaptations. Thus, RSA represents a useful

342 structural measure of a persistent effect of dN/dS with stronger correlations in the full
343 tree and internal branches than in tip branches.

344

345 Next we considered the correlation between dN/dS and the inverse distance to site 224
346 to each site in the HA structure (Figure 3B). In contrast to RSA, correlations here were
347 systematically higher in tip branches, suggesting a recent adaptive signal. We found
348 virtually no correlation for unpassaged sequences, while a low correlation existed for
349 monkey-cell cultured sequences and a higher correlation for cell-passaged and pooled
350 sequences. To confirm that differences between pooled and unpassaged correlations
351 were not simply due to variation from random sampling, we created a null distribution of
352 correlations from 200 random draws of pooled sequences. The correlation for
353 unpassaged sequences was significantly lower than it was for pooled sequences ($z = -$
354 4.22 , $p = 1.2 \times 10^{-5}$, Supplementary Figure 2). Correlations from pooled sequences
355 mirrored cell-culture correlations and persisted through internal branches. Thus, the
356 correlation of dN/dS with the inverse distance to site 224 seems to be primarily an
357 artifact of cell passage, even though its effect can be seen along internal branches as
358 well. This cell-specific signal dominates the pooled dataset. Further, this cell-specific
359 signal is partially attenuated along internal branches and amplified along tip branches,
360 as we would expect from a signal caused by recent host-specific adaptation. Even
361 though this signal is a true predictor of influenza evolutionary rates for virus grown in
362 cell culture, it does not transfer to unpassaged sequences and therefore has no
363 relevance for the circulating virus. This finding serves as a strong demonstration of

364 passage history as a confounder in analysis of hemagglutinin evolution, not just for egg
365 passage as previously demonstrated, but also for cell and monkey-cell passage.

366

367 Surprisingly, the correlation we found here between dN/dS and inverse distance to site
368 224 for pooled sequences ($R^2 = 0.067$) was less than half of the value previously
369 reported (Meyer and Wilke, 2015) (Fig. 3B). However, using a dataset of sequences
370 more temporally matched to the previously published analysis (2005–2014 instead of
371 2005–2015), we recovered the earlier higher correlation. This finding suggests that
372 there is some feature in the additional 2015 sequences that changes the pooled data's
373 relationship with inverse distance to site 224. In 2015, unpassed and SIAT1
374 sequences each doubled in number compared to 2014, while the number of non-SIAT1
375 cell cultured sequences dropped dramatically (Table 1). SIAT1, an MDCK cell line which
376 overexpresses human-like 6-linked sialic acids over native 3-linked sialic acids
377 (Matrosovich et al., 2003) has higher sequence fidelity than unmodified MDCK
378 (Hamamoto et al., 2013). Therefore, we next investigated whether the drop in
379 correlation from 2014 to 2015 could be attributed to the recent reduction in cell culture
380 using non-SIAT1 cells.

381

382 **3.3 Adaptation to passage in SIAT1 cells is weak or absent**

383 In the preceding analyses, we lumped all cell cultures except monkey cells into the
384 same category. However, there are more subtle distinctions in cell passaging systems,
385 and they can exert differential selective pressures on human adapted virus (Hamamoto
386 et al., 2013; Oh et al., 2008). As our generic cell culture group was composed of a

387 mixture of wild type MDCK, SIAT1, and unspecified cell cultures, we next investigated
388 whether any one culture type was the source of the high cell-culture signal seen in
389 Figure 3B.

390
391 SIAT1 is currently the dominant system for passaging of influenza virus in North
392 America, with approximately half of the 2015 influenza sequences currently available
393 from GISAID deriving from serial passaging through SIAT1 cells. Experimental analysis
394 of SIAT1 demonstrates improved sequence fidelity and reduced positive selection over
395 unmodified MDCK cell culture (Hamamoto et al., 2013; Oh et al., 2008). We sought to
396 determine if the apparently cell-culture-specific correlation of site-wise evolutionary
397 rates and inverse distance to site 224 extended to SIAT1 cell culture. To compare cell-
398 culture varieties, we created sample-size matched groups of non-SIAT1 cell culture,
399 SIAT1 cell culture, and unpassaged sequences collected between 2005 and 2015 ($n =$
400 1046), excluding sequences that had been passaged through both a non-SIAT1 and a
401 SIAT1 cell culture.

402
403 All groups showed similar correlations between dN/dS and RSA, regardless of whether
404 dN/dS was calculated for the entire tree, for internal branches only, or for tip branches
405 only (Figure 4A). By contrast, inverse distance to site 224 uniquely correlated with
406 dN/dS from non-SIAT1-cultured virus (Figure 4B). This effect was strongest along tip
407 branches ($R^2 = 0.139$), but it was almost as strong along the entire tree ($R^2 = 0.129$).
408 The correlation was reduced, though still significant, among internal branches ($R^2 =$
409 0.075). Thus, we conclude that the correlation between dN/dS and the inverse distance

410 to site 224 represents a unique signal of adaptation to passaging in non-SIAT1 cells. In
411 other words, a non-SIAT1-specific signal can completely dominate all signals of positive
412 adaptation when a dataset contains a sufficiently high number of sequences passaged
413 in non-SIAT1 cells. In our analysis (Figure 3B), the high correlation of non-SIAT1 cell
414 dN/dS with inverse distance to site 224 is suppressed in the pooled condition because
415 the number of unpassaged and SIAT1-passaged sequences grew substantially in 2015.
416 This difference in sample composition explains the lower than expected correlations in
417 Figure 3B for pooled dN/dS .

418
419 As these three conditions were somewhat temporally separated (most non-SIAT1 cell
420 culture sequences were pre-2015, and most unpassaged and SIAT1 culture sequences
421 were post-2014), we controlled for season-to-season variation by drawing 249
422 sequences from each group from 2014. We again considered site-wise dN/dS
423 correlations among passaging groups, and we found that overall, unpassaged and
424 SIAT1-passaged sequences appeared the most similar (Supplementary Figure 3A–C).

425

426 **3.4 Signals of passaging adaptation accumulate with additional rounds of** 427 **passaging in non-SIAT1 cells**

428 Having identified non-SIAT1 cell culture as the source of the contaminating signal in
429 analyses of inverse distance to site 224, we next investigated the source of this signal at
430 the single amino acid level. We expected that a signal of adaptation to a passaging
431 system would strengthen with additional exposure to that system. Thus, we compared
432 the magnitude of the site-wise dN/dS values in sequences that had never been

433 passaged, had been passaged once, passaged twice, or passaged three to five times
434 (Figure 5A). For this analysis, we only considered passage in non-SIAT1 cells. This
435 analysis revealed distinct regions of increasing positive selection along the
436 hemagglutinin molecule (arrows in Figure 5A) and a strong relationship between the
437 magnitude of these signals and the number of times influenza viruses were passaged.
438 Further, we found an overall increase in dN/dS with increased numbers of passages in
439 non-SIAT1 cells (Figure 5B). The strongly selected sites 221 and 225 are adjacent to
440 site 224, explaining the specific relationship between dN/dS calculated from non-SIAT1
441 sequences and the inverse distance in 3D space to this site. The correlation between
442 dN/dS and inverse distance increased in strength with increasing numbers of passages
443 (Figure 5C), even though it was observable after a single passage in non-SIAT1 cells.
444 Mapping the raw dN/dS values onto the hemagglutinin structure showed how specific
445 sites light up as passage numbers increase (Figure 5D).

446

447 **3.5 Evolutionary variation in sequences from unpassaged virus predicts regions** 448 **involved in antigenic escape**

449 Our preceding analyses might suggest that the inverse distance metric for describing
450 regions of selection only captures effects of adaptation to non-SIAT1 cell culture.
451 However, this is not necessarily the case. Importantly, inverse distance needs to be
452 calculated relative to a specific reference point. Site 224 was previously used as the
453 reference point because it yielded the highest correlation for the dataset analyzed
454 (Meyer and Wilke, 2015). For a different dataset, one that doesn't carry the signal of

455 adaptation to non-SIAT1 cell culture, a different reference point may be more
456 appropriate.

457

458 We thus repeated the inverse distance analysis of Meyer and Wilke (2015) for a size-
459 matched sample of 1046 sequences from non-SIAT1, pooled, SIAT1, and unpassaged
460 virus collected between 2005 and 2015 (Figure 6). In brief, for each possible reference
461 site in the hemagglutinin structure, we measured the inverse distance in 3D space from
462 that site to every other site in the structure (see Methods for details). We then correlated
463 the inverse distances with the dN/dS values at each site, resulting in one correlation
464 coefficient per reference site. Finally, we mapped these correlation coefficients onto the
465 HA structure, coloring each reference site by its associated correlation coefficient. If
466 inverse distances measured from a particular reference amino acid have higher
467 correlation with the site-wise dN/dS values, then this reference site will appear
468 highlighted on the structure.

469

470 For non-SIAT1-passaged and pooled virus, this analysis recovered the finding of Meyer
471 and Wilke (2015) that the loop containing site 224 appeared strongly highlighted (Figure
472 7A). However, this signal was entirely absent in unpassaged and SIAT1 passaged virus
473 (Figure 7B), with no sites in that loop working well as a reference point. These results
474 suggested that this loop was specifically involved in adaptation of hemagglutinin to non-
475 SIAT1 cell culture, explaining the non-SIAT1-specific signal shown in Figure 4A.
476 Globally, the pattern of correlations from pooled sequences strongly resembled the non-
477 SIAT1 pattern, in contrast to the resemblance of SIAT1 to unpassaged. Thus, the

478 inverse distance metric is useful for differentiating regions of selection particular to
479 different experimental groups.

480

481 Therefore, we next asked what residual patterns of positive selection remained once the
482 adaptation to non-SIAT1 cells was removed. Even though site-wise correlations are
483 relatively low for unpassaged virus compared to the ones observed for non-SIAT1-
484 passaged virus, we could still recover relevant patterns of HA adaptation after rescaling
485 our coloring. In particular, we found that sites opposite to the loop-containing site 224 lit
486 up in our analysis of unpassaged sequences (Figure 7A). Sites in this region are known
487 to be involved in antigenic escape. In fact, many of the highlighted regions contain
488 amino acid positions where substitution led to antigenic change (Table 2). We found a
489 similar pattern of concordance with antigenic sites when mapping dN/dS values directly
490 onto the structure (Figure 7B). The inverse-distance correlations, however, performed
491 better at identifying antigenic residues than did raw dN/dS values. When considering the
492 90th percentile (top 10% highest scored sites) by either metric, the inverse-distance
493 correlations recovered 5 of 7 sites while dN/dS alone recovered only 1 of 7 sites (Table
494 2). Additionally, while several sites involved in antigenic change had very low dN/dS , all
495 had inverse-distance correlations above the 86th percentile.

496

497 **3.6 Passaging artifacts extend deep into reconstructed trees**

498 Passaging adaptations could reasonably be expected to only affect peripheral clades of
499 influenza virus evolutionary trees, as they represent recent signals of adaptation that
500 should not penetrate far into the tree or significantly affect tree structure (Kryazhimskiy

501 and Plotkin, 2008; Strelkova and Lässig, 2012). Surprisingly, however, we found
502 signals of passage adaptations in season-to-season fixed mutations and branch
503 density.

504

505 To capture mutations that became fixed across seasons, we calculated dN/dS along the
506 trunks of trees constructed from sequences of difference passage histories (Figure 8A).
507 To time-calibrate our trunk, we limited this analysis to passage types that had
508 sequences at least in 2005 and in 2015, which excludes SIAT1 sequences. Trunk
509 dN/dS measures season-to-season adaptation and might be expected to be robust to
510 the effects of passaging. However, recurring adaptations to passaging conditions as
511 samples were processed across seasons could falsely appear as trunk mutations.

512

513 For pooled sequences, trunk dN/dS appeared to be generally free of artifacts,
514 resembling the trunk dN/dS of unpassaged sequences (Figure 8A). In contrast, for
515 datasets composed entirely of passaged sequences we found artifacts extending into
516 the trunk. When trees were constructed from only cell-passaged sequences or only non-
517 SIAT1 sequences, we observed a general inflation in dN/dS as well as several spurious
518 sites of high dN/dS that do not occur in the unpassaged condition. Together, this result
519 shows that the relative proportion of passaged to unpassaged sequences in a sample
520 matters; when sequences with passaging artifacts are overrepresented compared to
521 unpassaged sequences, there is a risk that a spurious signal will be found in the trunk.
522 For example, a trunk dN/dS analysis of mainly non-SIAT1 sequences would direct
523 attention to site 261, even though this site does not appear to be positively selected on

524 the unpassaged tree trunk. This analysis demonstrates the ability of sequences
525 containing major passaging artifacts to confound both deep and peripheral analyses of
526 influenza virus evolution.

527

528 We next investigated the effect of passaging on a tree-topology based metric, Local
529 Branching Index (LBI) (Neher et al., 2014). Notably, this metric is entirely independent
530 of dN/dS values. The LBI algorithm uses the degree of local branching around a
531 terminal node to predict sequences similar to progenitors of the following season's
532 strain. As a read-out of the algorithm's performance, we calculated the mean ratio
533 score, as described (Neher et al., 2014) (see also Methods). Mean ratio scores below 1
534 indicate that the algorithm performs better than random chance. The lowest possible
535 mean ratio score possible corresponds to the mean ratio score of the sequence with the
536 lowest hamming distance to the following year's progenitor (i.e., the theoretical best
537 possible prediction from sequences in a condition).

538

539 We saw clear differences in accuracy of predictions made using trees composed of non-
540 SIAT1, pooled, or unpassaged sequences (Figure 8B, C). We could not examine SIAT1
541 patterns, as these sequences were not consistently available across seasons until
542 2013. As a general trend, predictions from unpassaged sequences seemed to be more
543 accurate (both less likely to exceed 1 and more likely to be closer to the best possible
544 prediction) than predictions from either passaged or pooled sequences.

545

546 **4. DISCUSSION**

547 We find that serial passaging of influenza virus introduces a measurable signal of
548 adaptation into the evolutionary analysis of natural influenza sequences. There are
549 unique, characteristic patterns of adaptation to egg passage, monkey cell passage, and
550 non-SIAT1 cell passage. Monkey-cell-derived sequences show different molecule-wide
551 evolutionary rate patterns. Non-SIAT1 cell-derived sequences instead display a hotspot
552 of positive selection in a loop underneath the sialic-acid binding region. This hotspot has
553 been previously noted (Meyer and Wilke, 2015) but no explanation for its origin was
554 available. Additional passages in non-SIAT1 cell strengthen this artifact. Further, we find
555 that virus passaged in SIAT1 cells seems to accumulate only minor passaging artifacts.
556 Throughout our analyses, we find limited utility in subdividing phylogenetic trees to
557 internal and terminal branches. While signals of passage adaptation are consistently
558 elevated along terminal branches and attenuated along internal branches, evolutionary
559 rates along internal branches remain confounded by passaging artifacts. Additionally,
560 passage adaptation can resemble fixed season-to-season mutation along trunk
561 branches and alter topology-based predictions of sequence fitness. Finally, we can
562 accurately recover the experimentally determined antigenic regions of hemagglutinin
563 from evolutionary-rate analysis by using a dataset consisting of only unpassaged viral
564 sequences.

565

566 Previous studies (Bush et al., 2001; Suzuki, 2006) suggest the use of internal branches
567 to alleviate passage adaptations. However, we find here that this strategy is insufficient,
568 because the evolutionary signal of passage adaptations can often be detected along
569 internal branches. This finding may seem counterintuitive, as internal nodes should

570 exclusively represent human-adapted virus. We suggest that passaging adaptations in
571 internal branches may be homoplasies caused by convergent evolution; if different
572 clinical isolates converge onto the same adaptive mutations under passaging, then
573 these mutations may incorrectly be placed along internal branches under phylogenetic
574 tree reconstruction. Additionally, although the use of only internal branches removes
575 some differences between the passage groups, the exclusion of terminal sequences
576 can obscure recent natural adaptations and thus obscure actual sites under positive
577 selection. Therefore, analysis of internal branches is not only insufficient for eliminating
578 artifacts from passaging adaptations but also suboptimal for detecting positive selection
579 in seasonal H3N2 influenza.

580

581 The safest route to avoid passaging artifacts is to limit sequence datasets to only
582 unpassaged virus, although this approach limits sequence numbers. The human-like 6-
583 linked sialic acids in SIAT1 (Matrosovich et al., 2003) greatly reduce observed cell
584 culture-specific adaptations, particularly in the loop of hemagglutinin which contains site
585 224. This lack of selection concords with multiple experiments finding low levels of
586 adaptation in this cell line (Hamamoto et al., 2013; Oh et al., 2008). As our analysis only
587 detects minor differences between unpassaged and SIAT1 passaged virus, we posit
588 that this passage condition is an acceptable substitute for unpassaged clinical samples.
589 Even so, our findings do not preclude the existence of SIAT1-specific adaptations that
590 may confound specific analyses.

591

592 Over half of the passaged hemagglutinin sequences in the GISAID database from 2005-
593 2015 were passaged more than once. Multiple passages cause increasing
594 accumulation of passaging artifacts in non-SIAT1 cells, and we predict that any yet
595 unknown passaging effects would accumulate similarly. However, even a single
596 passage in non-SIAT1 cells introduces noticeable artifacts of adaptation. Influenza virus
597 is often passaged multiple times to improve viral titers for hemagglutination inhibition
598 assays, and thus, we expect that multiply passaged viruses will continue to be
599 deposited for the foreseeable future. We recommend that such viruses be used with
600 care when studying the evolutionary dynamics of influenza strains circulating in the
601 human population.

602

603 Although the majority of the sequences from the year 2015 are SIAT1-passaged or
604 unpassaged, several hundred sequences from that year derive from monkey cell
605 culture. The use of monkey cell culture surged in 2014 and 2015 compared to previous
606 years. We recommend that these recently collected sequences be excluded from
607 influenza rate analysis, in favor of the majority of unpassaged and SIAT1-passaged
608 sequences. As passaging is a useful and cost effective method for amplification of
609 clinically collected virus, unpassaged viral sequences are unlikely to completely
610 dominate influenza sequence databases in the near future. However, new human
611 epithelial cell culture systems for influenza passaging (Ilyushina et al., 2012) could soon
612 provide an ideal system that both amplifies virus and protects it from non-human
613 selective pressures.

614

615 Passage history should routinely be considered as a potential confounding variable in
616 future analyses of influenza evolutionary rates. Future studies should be checked
617 against unpassaged samples to ensure that conclusions are not based on adaptation to
618 non-human hosts. We recommend the exclusion of viral sequences that derive from
619 serial passage in egg amniotes, monkey kidney cell culture, and any unspecified cell
620 culture. Prior work that did not consider passaging history may be confounded by
621 passaging adaptations, as occurred in our previous publication (Meyer and Wilke,
622 2015). In that manuscript, we concluded that sites under positive selection differ from
623 sites involved in immune escape. Here, we find that the origin of this positive selection
624 is adaptation to the non-human passaging host, not immune escape in or adaptation to
625 humans. In particular, we suggest that the evolutionary markers of influenza virus
626 determined in (Belanov et al., 2015) be reevaluated to ensure these sites are not
627 artifacts of viral passaging. Similarly, many of the earlier studies (Bush et al., 1999;
628 Meyer and Wilke, 2013, 2015; Pan and Deem, 2011; Shih et al., 2007; Suzuki, 2006,
629 2008; Tusche et al., 2012) performing site-specific evolutionary analysis of
630 hemagglutinin likely contain some conclusions that can be traced back to passaging
631 artifacts. Additionally, even though passage artifacts do not appear to be sufficiently
632 strong to affect clade-structure reconstruction (Bush et al., 2000), they do have the
633 potential to cause artificially long branch lengths, due to dN/dS inflation, or misplaced
634 branches, due to convergent evolution under passaging. We find samples composed of
635 non-SIAT1 appear to behave differently than unpassaged samples under the Local
636 Branching Index metric. Thus, future phylogenetic predictive models of influenza fitness
637 and antigenicity, as in (Bedford et al., 2014; Łuksza and Lässig, 2014; Neher et al.,

638 2014), should also be checked for robustness to passage-related signals. Finally, while
639 it is beyond the scope of this work to investigate passage history effects in other
640 viruses, we suspect that passage-derived artifacts could be a factor in their phylogenetic
641 analyses as well. The use of datasets free of passage adaptations will likely bring
642 computational predictions of influenza positive selection more in line with corresponding
643 experimental results.

644

645 Sequences without passage annotations are inadequate for reliable evolutionary
646 analysis of influenza virus. Yet, passage annotations are often completely missing from
647 strain information, and, when present, are often inconsistent; there is currently no
648 standardized language to represent number and type of serial passage. We note,
649 however, that passage annotations from the 2015 season are greatly improved when
650 compared to previous seasons. Several major influenza repositories, including the
651 Influenza Research Database (Squires et al., 2012) and the NCBI Influenza Virus
652 Resource (Bao et al., 2008), do not provide any passaging annotations at all.
653 Additionally, passage history is not required for new sequence submissions to the NCBI
654 Genbank (Benson et al., 2012). The EpiFlu database maintained by the Global Initiative
655 for Sharing Avian Influenza Data (GISAID) (Bogner et al., 2006) and OpenFluDB
656 (Liechti et al., 2010), however, stand apart by providing passage history annotations for
657 the majority of their sequences. Of these, only the OpenFluDB repository allows filtering
658 of sequences by passage history during data download. Our results demonstrate the
659 strength of passaging artifacts in evolutionary analysis of influenza. The lack of a
660 universal standard for annotation of viral passage histories and a universal standard for

661 serial passage experimental conditions complicate the analysis and mitigation of
662 passaging effects.

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684 **REFERENCES**

- 685 Balish, A.L., Katz, J.M., and Klimov, A.I. (2005). Influenza: Propagation, Quantification, and Storage. In
686 *Current Protocols in Microbiology*, (John Wiley & Sons, Inc.), p.
- 687 Bao, Y., Bolotov, P., Dernovoy, D., Kiryutin, B., Zaslavsky, L., Tatusova, T., Ostell, J., and Lipman, D.
688 (2008). The Influenza Virus Resource at the National Center for Biotechnology Information. *J. Virol.* **82**,
689 596–601.
- 690 Bedford, T., Suchard, M.A., Lemey, P., Dudas, G., Gregory, V., Hay, A.J., McCauley, J.W., Russell, C.A.,
691 Smith, D.J., and Rambaut, A. (2014). Integrating influenza antigenic dynamics with molecular evolution.
692 *eLife* **3**, e01914.
- 693 Belanov, S.S., Bychkov, D., Benner, C., Ripatti, S., Ojala, T., Kankainen, M., Lee, H.K., Tang, J.W.-T., and
694 Kainov, D.E. (2015). Genome-wide analysis of evolutionary markers of human influenza A(H1N1)pdm09
695 and A(H3N2) viruses may guide selection of vaccine strain candidates. *Genome Biol. Evol.* **evv240**.
- 696 Benson, D.A., Karsch-Mizrachi, I., Clark, K., Lipman, D.J., Ostell, J., and Sayers, E.W. (2012). GenBank.
697 *Nucleic Acids Res.* **40**, D48–D53.
- 698 Blackburne, B.P., Hay, A.J., and Goldstein, R.A. (2008). Changing Selective Pressure during Antigenic
699 Changes in Human Influenza H3. *PLoS Pathog* **4**, e1000058.
- 700 Bogner, P., Capua, I., Lipman, D.J., Cox, N.J., and others (2006). A global initiative on sharing avian flu
701 data. *Nature* **442**, 981–981.
- 702 Bush, R.M., Fitch, W.M., Bender, C.A., and Cox, N.J. (1999). Positive selection on the H3 hemagglutinin
703 gene of human influenza virus A. *Mol. Biol. Evol.* **16**, 1457–1465.
- 704 Bush, R.M., Smith, C.B., Cox, N.J., and Fitch, W.M. (2000). Effects of passage history and sampling bias on
705 phylogenetic reconstruction of human influenza A evolution. *Proc. Natl. Acad. Sci.* **97**, 6974–6980.
- 706 Bush, R.M., Fitch, W.M., Smith, C.B., and Cox, N.J. (2001). Predicting influenza evolution: the impact of
707 terminal and egg-adapted mutations. *Int. Congr. Ser.* **1219**, 147–153.
- 708 Echave, J., Spielman, S.J., and Wilke, C.O. (2016). Causes of evolutionary rate variation among protein
709 sites. *Nat. Rev. Genet.*
- 710 Gatherer, D. (2010). Passage in egg culture is a major cause of apparent positive selection in influenza B
711 hemagglutinin. *J. Med. Virol.* **82**, 123–127.
- 712 Hamamoto, I., Takaku, H., Tashiro, M., and Yamamoto, N. (2013). High Yield Production of Influenza
713 Virus in Madin Darby Canine Kidney (MDCK) Cells with Stable Knockdown of IRF7. *PLoS ONE* **8**, e59892.
- 714 Ihaka, R., and Gentleman, R. (1996). R: A Language for Data Analysis and Graphics. *J. Comput. Graph.*
715 *Stat.* **5**, 299–314.

- 716 Ilyushina, N.A., Ikizler, M.R., Kawaoka, Y., Rudenko, L.G., Treanor, J.J., Subbarao, K., and Wright, P.F.
717 (2012). Comparative study of influenza virus replication in MDCK cells and in primary cells derived from
718 adenoids and airway epithelium. *J. Virol.* *86*, 11725–11734.
- 719 Katoh, K., and Standley, D.M. (2013). MAFFT Multiple Sequence Alignment Software Version 7:
720 Improvements in Performance and Usability. *Mol. Biol. Evol.* *30*, 772–780.
- 721 Katz, J.M., and Webster, R.G. (1989). Efficacy of Inactivated Influenza A Virus (H3N2) Vaccines Grown in
722 Mammalian Cells or Embryonated Eggs. *J. Infect. Dis.* *160*, 191–198.
- 723 Katz, J.M., Wang, M., and Webster, R.G. (1990). Direct sequencing of the HA gene of influenza (H3N2)
724 virus in original clinical samples reveals sequence identity with mammalian cell-grown virus. *J. Virol.* *64*,
725 1808–1811.
- 726 Koelle, K., Cobey, S., Grenfell, B., and Pascual, M. (2006). Epochal evolution shapes the phylodynamics of
727 interpandemic influenza A (H3N2) in humans. *Science* *314*, 1898–1903.
- 728 Kratsch, C., Klingen, T.R., Mümken, L., Steinbrück, L., and McHardy, A.C. (2016). Determination of
729 antigenicity-altering patches on the major surface protein of human influenza A/H3N2 viruses. *Virus*
730 *Evol.* *2*, vev025.
- 731 Kryazhimskiy, S., and Plotkin, J.B. (2008). The Population Genetics of dN/dS. *PLOS Genet* *4*, e1000304.
- 732 Kumar, S., and Henrickson, K.J. (2012). Update on Influenza Diagnostics: Lessons from the Novel H1N1
733 Influenza A Pandemic. *Clin. Microbiol. Rev.* *25*, 344–361.
- 734 Lee, H.K., Tang, J.W.-T., Kong, D.H.-L., and Koay, E.S.-C. (2013a). Simplified Large-Scale Sanger Genome
735 Sequencing for Influenza A/H3N2 Virus. *PLoS ONE* *8*.
- 736 Lee, H.K., Tang, J.W.-T., Kong, D.H.-L., Loh, T.P., Chiang, D.K.-L., Lam, T.T.-Y., and Koay, E.S.-C. (2013b).
737 Comparison of Mutation Patterns in Full-Genome A/H3N2 Influenza Sequences Obtained Directly from
738 Clinical Samples and the Same Samples after a Single MDCK Passage. *PLoS ONE* *8*, e79252.
- 739 Liechti, R., Gleizes, A., Kuznetsov, D., Bougueleret, L., Mercier, P.L., Bairoch, A., and Xenarios, I. (2010).
740 OpenFluDB, a database for human and animal influenza virus. *Database* *2010*, baq004.
- 741 Lin, Y.P., Xiong, X., Wharton, S.A., Martin, S.R., Coombs, P.J., Vachieri, S.G., Christodoulou, E., Walker,
742 P.A., Liu, J., Skehel, J.J., et al. (2012). Evolution of the receptor binding properties of the influenza
743 A(H3N2) hemagglutinin. *Proc. Natl. Acad. Sci.* *109*, 21474–21479.
- 744 Łuksza, M., and Lässig, M. (2014). A predictive fitness model for influenza. *Nature* *507*, 57–61.
- 745 Matrosovich, M., Matrosovich, T., Carr, J., Roberts, N.A., and Klenk, H.-D. (2003). Overexpression of the
746 alpha-2,6-sialyltransferase in MDCK cells increases influenza virus sensitivity to neuraminidase
747 inhibitors. *J. Virol.* *77*, 8418–8425.
- 748 Meyer, A.G., and Wilke, C.O. (2013). Integrating Sequence Variation and Protein Structure to Identify
749 Sites under Selection. *Mol. Biol. Evol.* *30*, 36–44.

- 750 Meyer, A.G., and Wilke, C.O. (2015). Geometric Constraints Dominate the Antigenic Evolution of
751 Influenza H3N2 Hemagglutinin. *PLoS Pathog.* *11*.
- 752 Neher, R.A., and Bedford, T. (2015). nextflu: real-time tracking of seasonal influenza virus evolution in
753 humans. *Bioinformatics* *btv381*.
- 754 Neher, R.A., Russell, C.A., and Shraiman, B.I. (2014). Predicting evolution from the shape of genealogical
755 trees. *eLife* *3*, e03568.
- 756 Nelson, M.I., Simonsen, L., Viboud, C., Miller, M.A., Taylor, J., George, K.S., Griesemer, S.B., Ghedin, E.,
757 Sengamalay, N.A., Spiro, D.J., et al. (2006). Stochastic Processes Are Key Determinants of Short-Term
758 Evolution in Influenza A Virus. *PLoS Pathog* *2*, e125.
- 759 Oh, D.Y., Barr, I.G., Mosse, J.A., and Laurie, K.L. (2008). MDCK-SIAT1 cells show improved isolation rates
760 for recent human influenza viruses compared to conventional MDCK cells. *J. Clin. Microbiol.* *46*, 2189–
761 2194.
- 762 Pan, K., and Deem, M.W. (2011). Quantifying selection and diversity in viruses by entropy methods, with
763 application to the haemagglutinin of H3N2 influenza. *J. R. Soc. Interface* *8*, 1644–1653.
- 764 Pond, S.L.K., Frost, S.D.W., and Muse, S.V. (2005). HyPhy: hypothesis testing using phylogenies.
765 *Bioinformatics* *21*, 676–679.
- 766 Price, M.N., Dehal, P.S., and Arkin, A.P. (2010). FastTree 2 – Approximately Maximum-Likelihood Trees
767 for Large Alignments. *PLoS ONE* *5*, e9490.
- 768 Robertson, J.S., Nicolson, C., Major, D., Robertson, E.W., and Wood, J.M. (1993). The role of amniotic
769 passage in the egg-adaptation of human influenza virus is revealed by haemagglutinin sequence
770 analyses. *J. Gen. Virol.* *74* (Pt 10), 2047–2051.
- 771 Shih, A.C.-C., Hsiao, T.-C., Ho, M.-S., and Li, W.-H. (2007). Simultaneous amino acid substitutions at
772 antigenic sites drive influenza A hemagglutinin evolution. *Proc. Natl. Acad. Sci.* *104*, 6283–6288.
- 773 Skowronski, D.M., Janjua, N.Z., De Serres, G., Sabaiduc, S., Eshaghi, A., Dickinson, J.A., Fonseca, K.,
774 Winter, A.-L., Gubbay, J.B., Kraiden, M., et al. (2014). Low 2012–13 Influenza Vaccine Effectiveness
775 Associated with Mutation in the Egg-Adapted H3N2 Vaccine Strain Not Antigenic Drift in Circulating
776 Viruses. *PLoS ONE* *9*, e92153.
- 777 Spielman, S., Wan, S., and Wilke, C.O. (2015). One-rate models outperform two-rate models in site-
778 specific dN/dS estimation. *bioRxiv* 32805.
- 779 Squires, R.B., Noronha, J., Hunt, V., García-Sastre, A., Macken, C., Baumgarth, N., Suarez, D., Pickett, B.E.,
780 Zhang, Y., Larsen, C.N., et al. (2012). Influenza Research Database: an integrated bioinformatics resource
781 for influenza research and surveillance. *Influenza Other Respir. Viruses* *6*, 404–416.
- 782 Stöhr, K., Bucher, D., Colgate, T., and Wood, J. (2012). Influenza Virus Surveillance, Vaccine Strain
783 Selection, and Manufacture. In *Influenza Virus*, Y. Kawaoka, and G. Neumann, eds. (Humana Press), pp.
784 147–162.

- 785 Strelkova, N., and Lässig, M. (2012). Clonal Interference in the Evolution of Influenza. *Genetics* *192*,
786 671–682.
- 787 Suzuki, Y. (2006). Natural selection on the influenza virus genome. *Mol. Biol. Evol.* *23*, 1902–1911.
- 788 Suzuki, Y. (2008). Positive selection operates continuously on hemagglutinin during evolution of H3N2
789 human influenza A virus. *Gene* *427*, 111–116.
- 790 The World Health Organization (2015). Recommended composition of influenza virus vaccines for use in
791 the 2015- 2016 northern hemisphere influenza season.
- 792 Tien, M.Z., Meyer, A.G., Sydykova, D.K., Spielman, S.J., and Wilke, C.O. (2013). Maximum allowed solvent
793 accessibilities of residues in proteins. *PLoS One* *8*, e80635.
- 794 Tusche, C., Steinbrück, L., and McHardy, A.C. (2012). Detecting patches of protein sites of influenza A
795 viruses under positive selection. *Mol. Biol. Evol.* *29*, 2063–2071.
- 796 WHO Writing Group, Ampofo, W.K., Baylor, N., Cobey, S., Cox, N.J., Daves, S., Edwards, S., Ferguson, N.,
797 Grohmann, G., Hay, A., et al. (2012). Improving influenza vaccine virus selection Report of a WHO
798 informal consultation held at WHO headquarters, Geneva, Switzerland, 14–16 June 2010. *Influenza*
799 *Other Respir. Viruses* *6*, 142–152.
- 800 Wickham, H. (2009). *ggplot2: elegant graphics for data analysis* (Springer New York).
- 801 Wolf, Y.I., Viboud, C., Holmes, E.C., Koonin, E.V., and Lipman, D.J. (2006). Long intervals of stasis
802 punctuated by bursts of positive selection in the seasonal evolution of influenza A virus. *Biol. Direct* *1*,
803 34.
- 804 World Health Organization Global influenza surveillance network (2011). *Manual for the laboratory*
805 *diagnosis and virological surveillance of influenza* (Geneva, Switzerland).
- 806 Wyde, P.R., Couch, R.B., Mackler, B.F., Cate, T.R., and Levy, B.M. (1977). Effects of Low- and High-
807 Passage Influenza Virus Infection in Normal and Nude Mice. *Infect. Immun.* *15*, 221–229.
- 808 Xie, H., Wan, X.-F., Ye, Z., Plant, E.P., Zhao, Y., Xu, Y., Li, X., Finch, C., Zhao, N., Kawano, T., et al. (2015).
809 H3N2 Mismatch of 2014–15 Northern Hemisphere Influenza Vaccines and Head-to-head Comparison
810 between Human and Ferret Antisera derived Antigenic Maps. *Sci. Rep.* *5*.

811

812

813

814

815

816

817

818 **ACKNOWLEDGEMENTS**

819 We would like to thank Sebastian Maurer-Stroh for help with interpreting passaging
820 annotations in GISAID. This work was supported in part by NIH grant no. R01
821 GM088344, DTRA grant no. HDTRA1-12-C-0007, and NSF Cooperative agreement no.
822 DBI-0939454 (BEACON Center). The funders had no role in study design, data
823 collection and analysis, decision to publish, or preparation of the manuscript.

824

825 **DATA AVAILABILITY**

826 Processed data are available as supplementary material. Sequence data are available
827 from GISAID as detailed in Methods. All analysis code used to generate the processed
828 data is available at: https://github.com/wikelab/influenza_H3N2_passaging

829

830 **AUTHOR CONTRIBUTIONS**

831 Conceived and designed the experiments: CDM COW. Wrote scripts and analytic tools:
832 CDM AGM. Performed the experiments: CDM. Analyzed the data: CDM COW. Wrote
833 the paper: CDM AGM COW.

834

835 **COMPETING FINANCIAL INTERESTS STATEMENT**

836 The authors declare no competing financial interests.

837

838 **TABLES**

839 **Table 1. Parsing of passage-annotated FASTA sequences into passage history**

840 **groups.** For each passage group, we defined a regular expression that could reliably
 841 identify sequences with that passage history. Regular expressions were applied through
 842 the built-in python library “re”. SIAT1 and non-SIAT1 cell culture regular expressions
 843 were applied to the subset of sequences identified as generic cell culture sequences.
 844 The three middle columns list the number of sequences we identified for each passage
 845 group, for years 2014 only, 2015 only, and 2005–2015. The two furthest right columns
 846 list the number of single and multiply passaged sequences from 2005-2015 for each
 847 condition.

848

Passage group	Regular expression	Number of sequences			Rounds of passage	
		2014	2015	2005–2015	Single	Multiple
Chicken egg amniotes	AM[1-9] E[1-7] AMNIOTIC EGG EX AM_[1-9]	6	0	79	1	79
Monkey cell culture	TMK RMK RHMK RII PMK R[1-9] RX	366	290	917	904	13
Generic cell culture	S[1-9] SX SIAT MDCK C[1-9] CX C_[1-9] M[1-9] MX X[1-9] ^X_\$	794	787	3041	867	2158
SIAT1	^S[1-9]_\$ ^SX_\$ SIAT2_SIAT1 SIAT3_SIAT1	389	626	1046	459	587
Non-SIAT1 cell culture	not SIAT SX S[1-9]	297	56	1755	408	1331
Unpassaged	LUNG P0 OR_ ORIGINAL CLINICAL DIRECT	249	506	1703	N/A	N/A
Pooled		1508	1601	6873	1772	2317

849

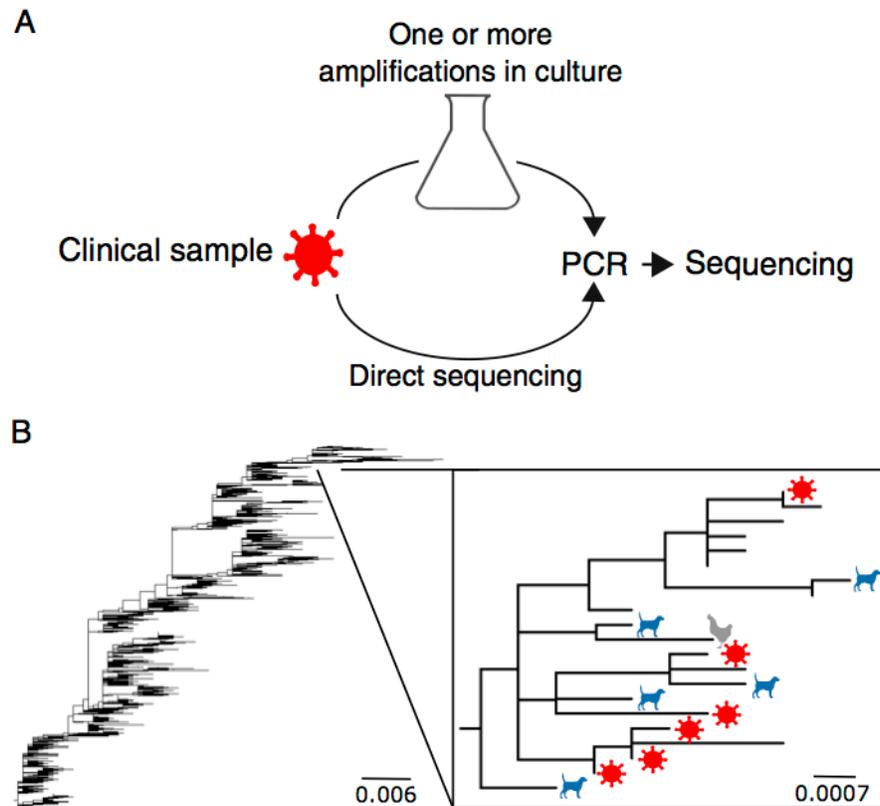
850 **Table 2. Evolutionary rates and inverse distance correlations of residues**
851 **responsible for antigenic change.** For each site, we determined dN/dS and the
852 correlation between dN/dS and inverse distance for unpassaged sequences collected
853 between 2005 and 2015 ($n = 1703$). 5/7 residues linked to antigenic changes have
854 inverse-distance correlations above the 90th percentile, while only 1/7 have dN/dS
855 values above the 90th percentile. Sites were experimentally determined by Koel et al.,
856 (2013).

857

Site		Raw dN/dS		Inv.-dist. correlation	
Gene	Protein	dN/dS	percentile	r	percentile
161	145	0.672	0.823	0.082	0.883
171	155	0	0	0.077	0.867
172	156	0.672	0.832	0.1317	0.971
174	158	1.36	0.958	0.1797	0.996
175	159	0.49	0.75	0.1837	1
205	189	0.474	0.763	0.0887	0.905
209	193	0.672	0.845	0.098	0.936

858

859 **FIGURES**



860

861 **Figure 1. Schematic of influenza A virus sequence collection and analysis. (A)**

862 Typical processing steps of influenza A virus clinical isolates. Virus collected from

863 patients may be passaged a single time or multiple times prior to PCR amplification and

864 sequencing in a variety of different environments (Ex. canine cell culture, monkey cell

865 culture, egg amniotes). However, some clinical virus is not passaged and is sequenced

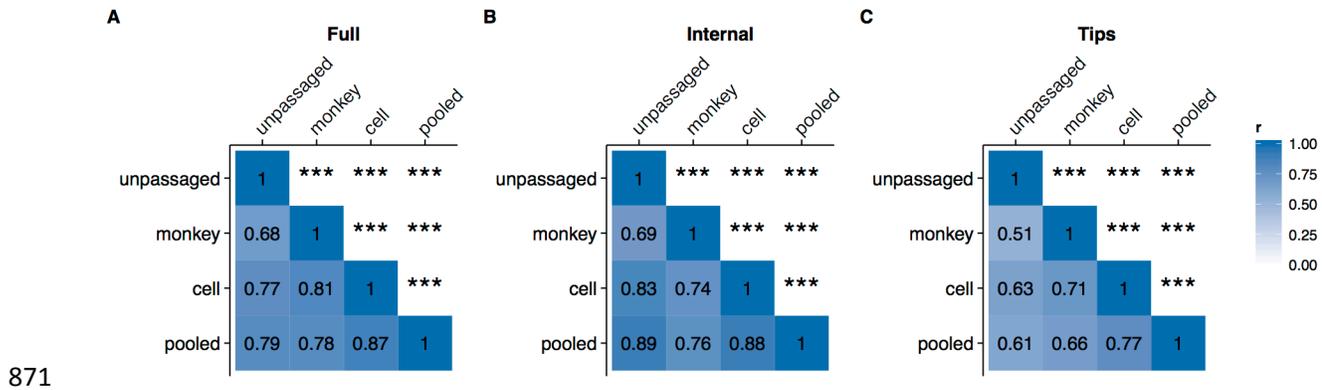
866 directly. (B) Phylogenetic tree of H3N2 HA sequences from the 2005-2015 seasons.

867 The inset shows a small clade of sequences from the 2006/2007 season, with colored

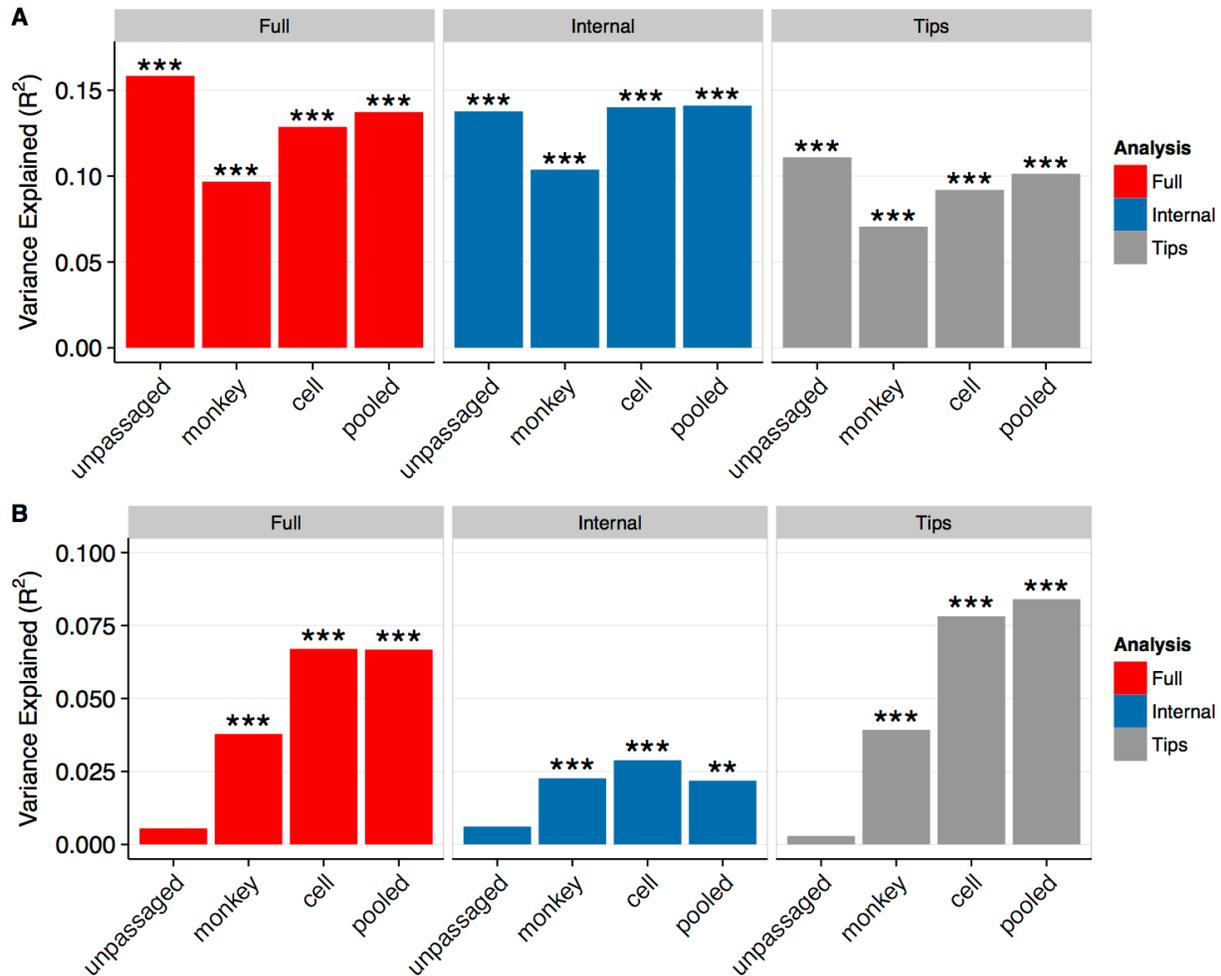
868 dots representing sequences with passage annotations (red virion: unpassed, blue

869 dog: canine cell culture, gray hen: egg amniote, unlabeled: missing or unclear passage-

870 history annotation).



872 **Figure 2. Comparison of sitewise dN/dS values among sequences with differing**
873 **passage histories.** Pearson correlations between sitewise dN/dS values for HA
874 sequences derived from passaged and unpassaged influenza virus collected between
875 2005 and 2015 (downsampled to $n = 917$ in all groups). Correlations were calculated
876 separately for dN/dS estimated from complete trees (A), internal branches only (B), and
877 tip branches only (C). Asterisks denote significance of correlations (* $0.01 \leq P < 0.05$,
878 ** $0.001 \leq P < 0.01$, *** $P < 0.001$). Data used to generate this figure are available in
879 Supplementary Data 1.



880

881 **Figure 3. Percent variance in dN/dS explained by relative solvent accessibility (A)**

882 **and by inverse distance to protein site 224 (B).** (A) Relative solvent accessibility

883 (RSA) explains ~10%–16% of the variation in dN/dS for all sequences. (B) Inverse

884 distance to site 224 explains ~7% of the variation in dN/dS for cell-passaged sequences

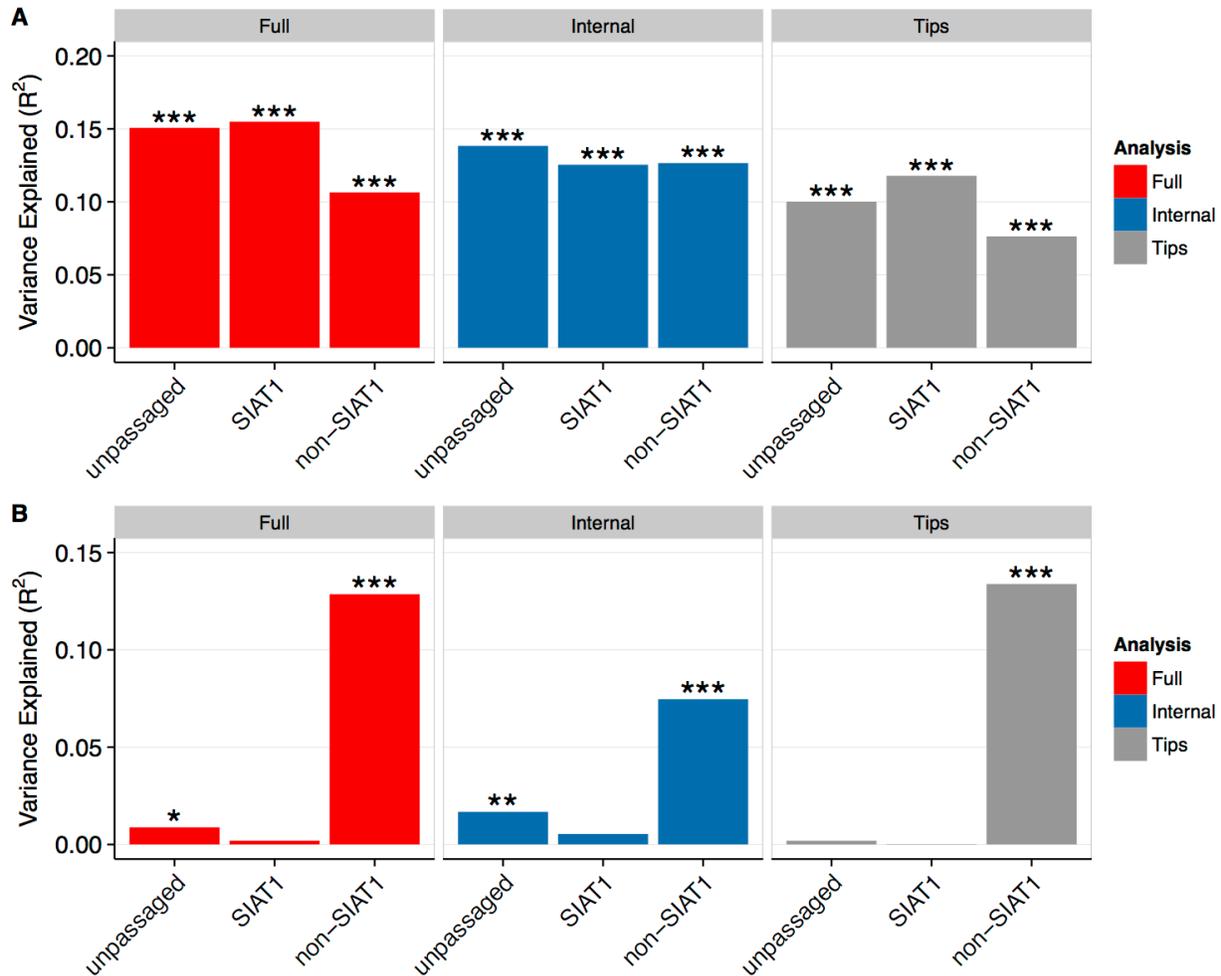
885 and for all sequences (pooled), however it explains virtually no variation for unpassaged

886 sequences. Asterisks denote significance of correlations (* $0.01 \leq P < 0.05$, ** $0.001 \leq P$

887 < 0.01 , *** $P < 0.001$). Data used to generate this figure are available in Supplementary

888 Data 1.

889



890

891 **Figure 4. Virus passaged in non-SIAT1 cells carries unique adaptations not**

892 **present in unpassaged or SIAT1-passaged virus. (A) The correlation between dN/dS**

893 **and RSA is weakened for virus passaged in non-SIAT1 cells. (B) The correlation**

894 **between dN/dS and inverse distance to site 224, representing a positive-selection**

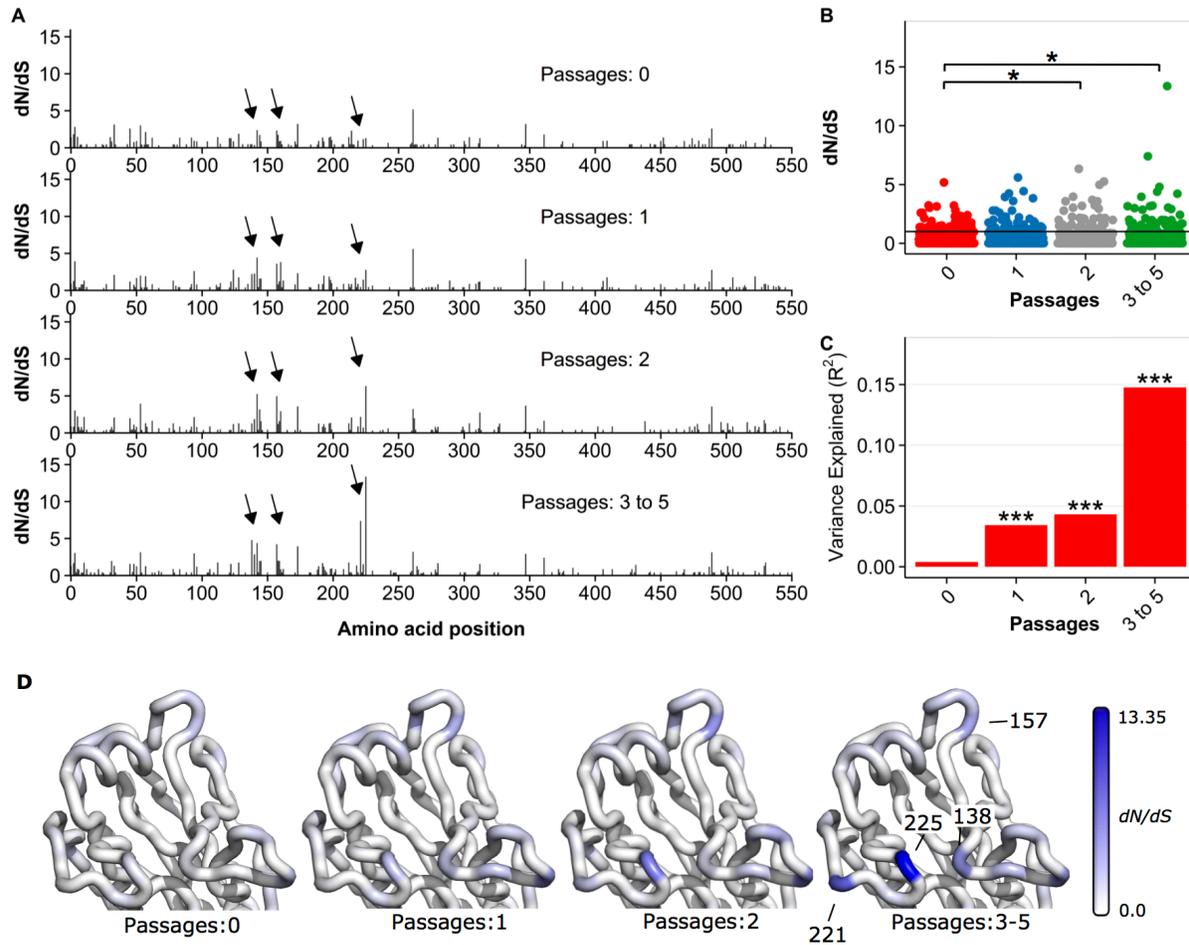
895 **hotspot in the vicinity of that site, is only present in virus passaged in non-SIAT1 cells.**

896 **Asterisks denote significance levels (* $0.01 \leq P < 0.05$, ** $0.001 \leq P < 0.01$, *** $P < 0.001$).**

897 **Sequences analyzed were collected between 2005 and 2015. Alignments were**

898 **randomly down-sampled to yield identical numbers of sequences in each alignment ($n =$**

899 **1046). Data used to generate this figure are available in Supplementary Data 4.**



900

901 **Figure 5. Accumulation of passaging artifacts with increasing numbers of serial**

902 **passages in non-SIAT1 cell culture. (A) Sitewise dN/dS values for virus which was**

903 **not passaged, passaged once, passaged twice, and passaged three to five times in**

904 **non-SIAT1 cell culture ($n = 304$ for each group). Arrows highlight regions of increased**

905 **dN/dS in passaged virus. Notably, dN/dS inflation is increased with increasing rounds of**

906 **passaging. (B) dN/dS values vs. number of passages. dN/dS values are significantly**

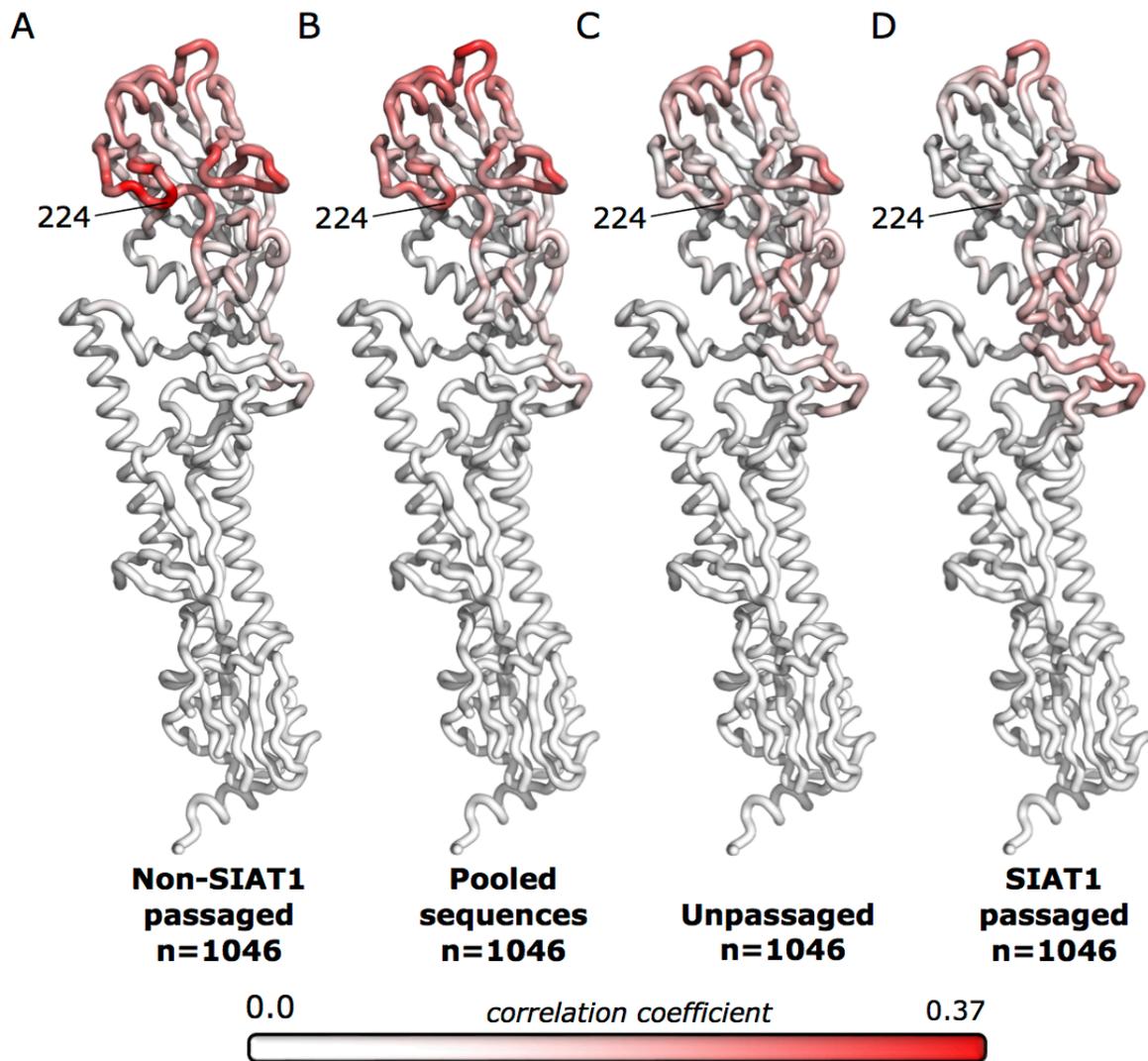
907 **elevated after two or more passages in non-SIAT1 cell culture, relative to unpassaged**

908 **virus (paired t test). Asterisks denote significance levels ($*0.01 \leq P < 0.05$, $**0.001 \leq P <$**

909 **0.01 , $***P < 0.001$). (C) The correlation between dN/dS and inverse distance to site 224**

910 **increases with the number of passages. (D) Mapping dN/dS values onto the**

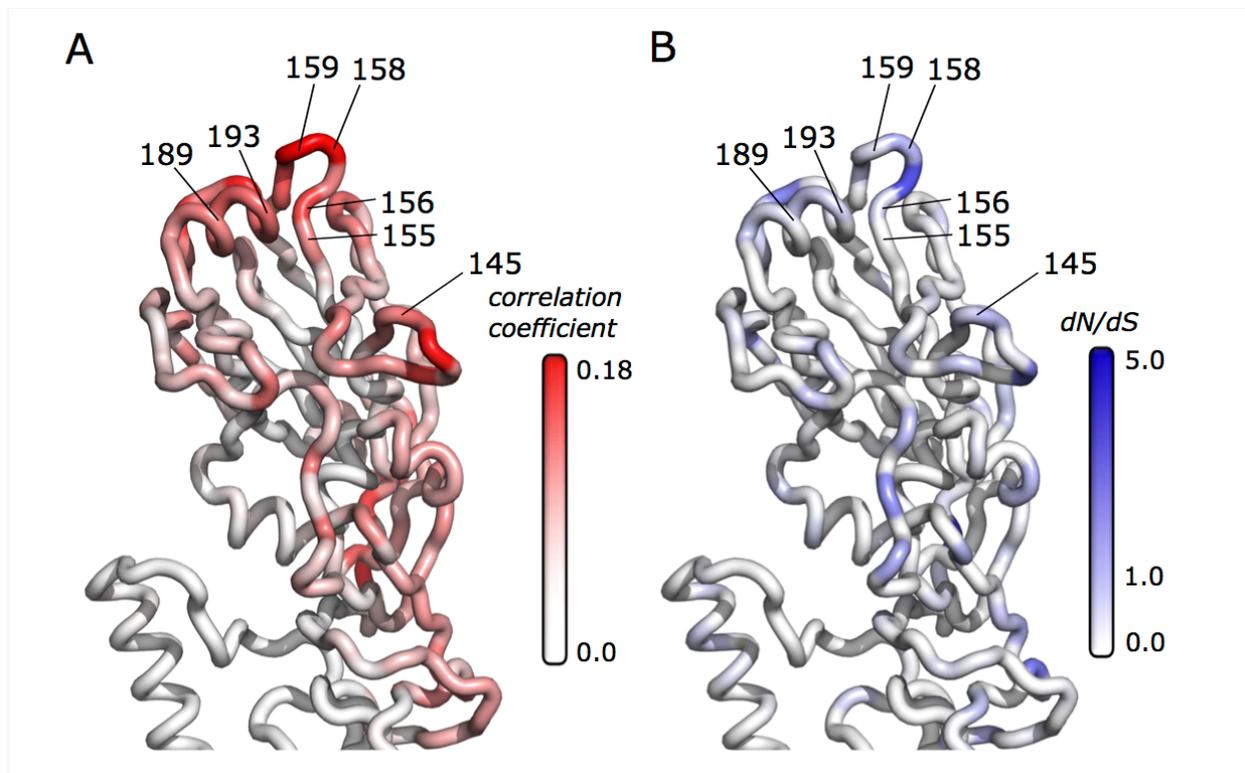
911 hemagglutinin head structure demonstrates the accumulation of passage adaptations
912 with increasing rounds of passages. Labeled sites correspond to regions denoted with
913 arrows in (A). Data used to generate this figure are available in Supplementary Data 6.



914

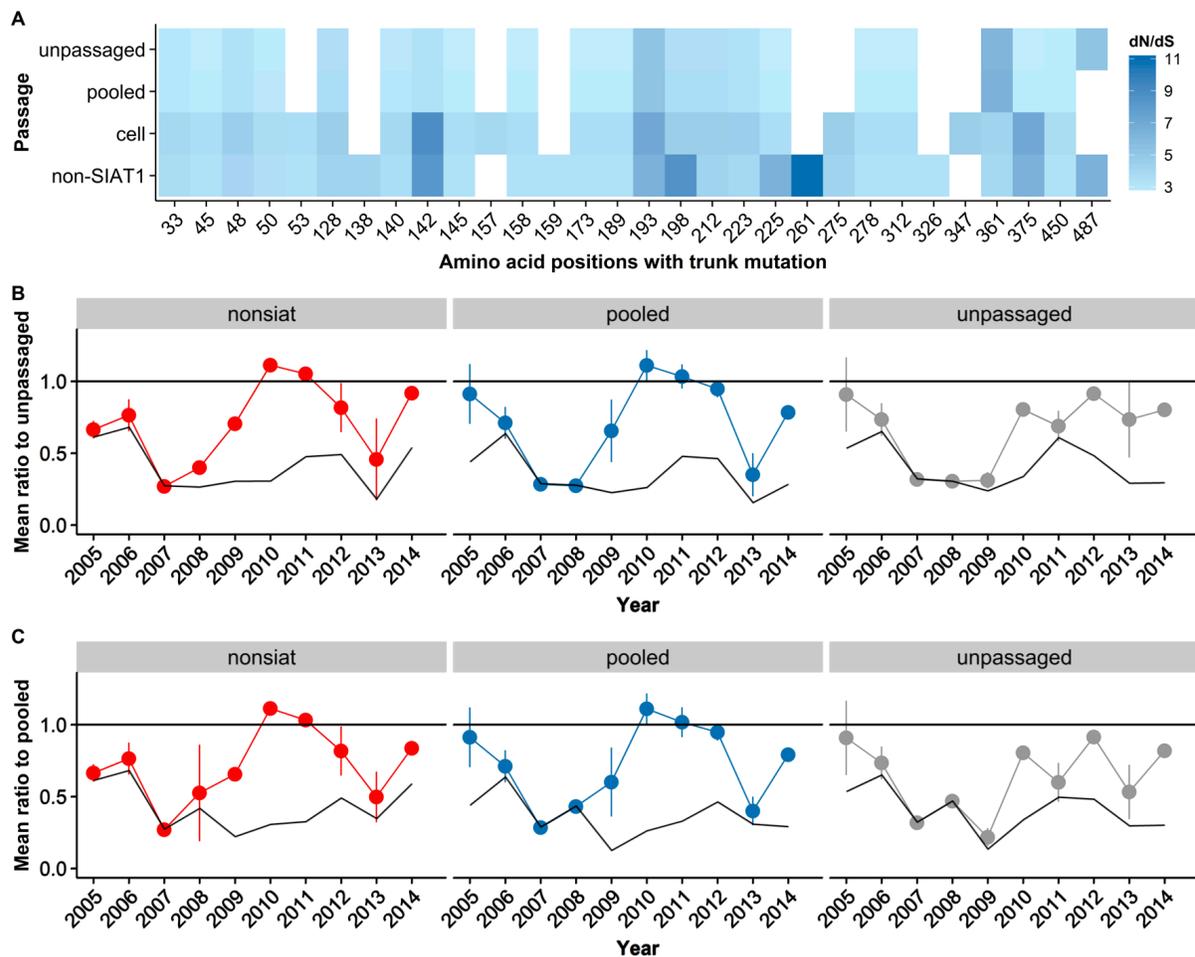
915 **Figure 6. Correlations of dN/dS with inverse distances, mapped onto the**
916 **hemagglutinin structure for non-SIAT1-passaged, pooled, unpassed, and SIAT1**
917 **passed sequences.** The correlation between dN/dS and inverse distance for each
918 reference site was mapped onto the hemagglutinin structure for (A) non-SIAT1
919 sequences, (B) pooled sequences, (C) unpassed sequences, and (D) SIAT1
920 passed sequences. Sequences analyzed were collected between 2005 and 2015.
921 Alignments were randomly down-sampled to yield identical numbers of sequences in

922 each alignment ($n = 1046$). Red coloring represents positive correlations, while white
923 represents zero or negative correlations. The four conditions group into two distinct
924 correlation patterns, non-SIAT1/pooled and unpassaged/SIAT1. In particular, the loop
925 containing site 224 lights up strongly for non-SIAT1 and pooled sequences but not for
926 unpassaged and SIAT1 sequences. Data used to generate this figure are available in
927 Supplementary Data 7.



928

929 **Figure 7. Unpassed sequences allow recovery of antigenic regions from**
930 **positive-selection analysis.** For each site, the correlation between dN/dS and inverse
931 distance (A) or dN/dS directly (B) were mapped onto the hemagglutinin structure, for
932 dN/dS derived from unpassed sequences collected between 2005 and 2015 ($n =$
933 1703). Red coloring represents higher correlation; blue coloring represents higher
934 dN/dS . Highlighted regions contain residues (labeled with protein site number) which
935 experimentally determined to cause antigenic change by Koel et al., 2013. Correlations
936 and dN/dS for antigenic residues are given in Table 2. Data used to generate this figure
937 are available in Supplementary Data 7.



938

939 **Figure 8. Passage artifacts affect trunk dN/dS and topology-based predictions. (A)**

940 Sitewise trunk dN/dS values for passage groups ($n = 1703$). Only sites with at least one

941 non-synonymous mutation along the trunk were included. Many sites appear under

942 positive selection only in the trunk reconstructed from passaged sequences, e.g. 159,

943 261, 275. (B, C) Prediction of future dominant clade by Local Branching Index (LBI)

944 depends on passaging history. LBI was calculated for trees derived from non-SIAT1,

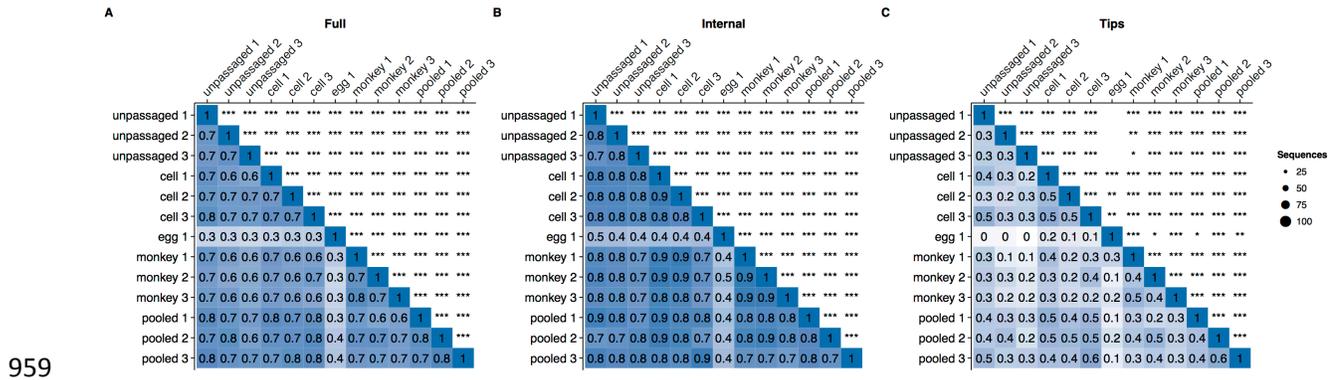
945 pooled, and unpassed sequences using a maximum of 100 sequences per condition.

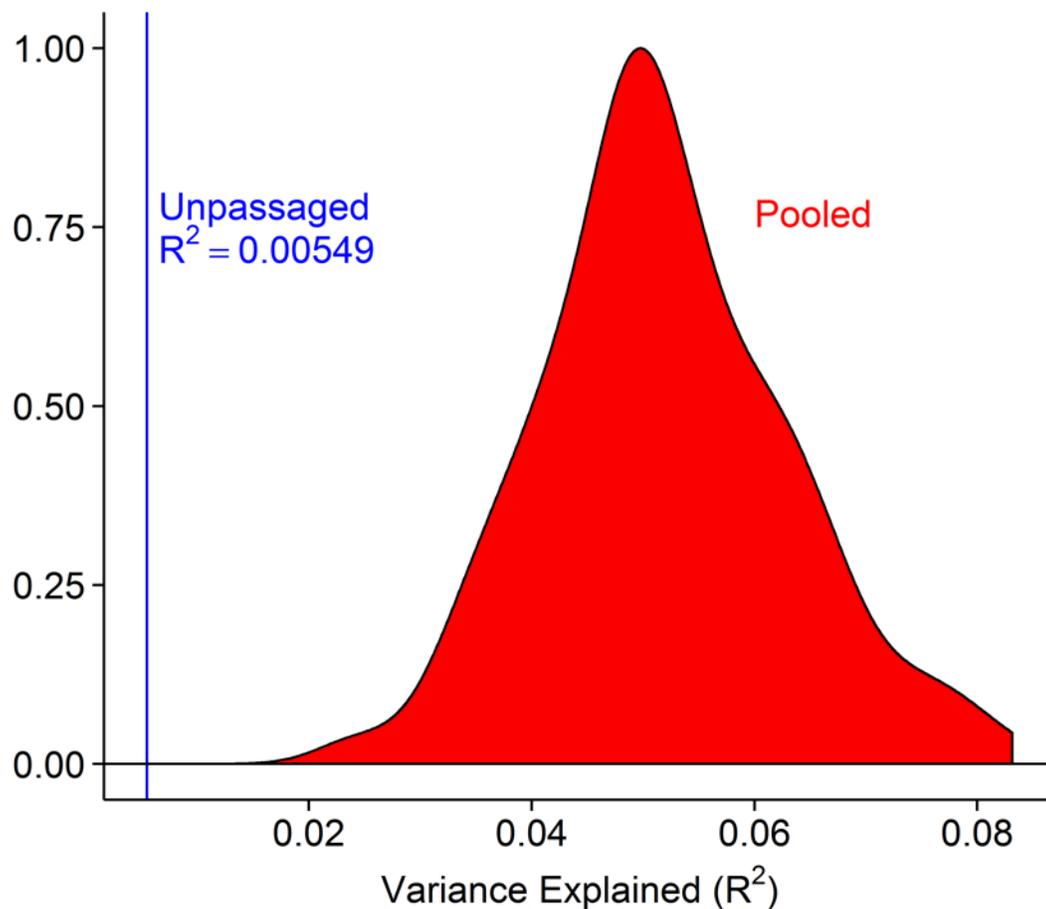
946 The mean ratio estimates the quality of the prediction (lower is better), and values < 1

947 indicate the prediction performs better than random. Error bars indicate the standard

948 deviation of the mean, estimated from resampling (Methods). The solid black lines
949 represent the best possible prediction in each year. Predictions were evaluated relative
950 to the following year's ancestrally reconstructed root sequence obtained from (A)
951 unpassaged and (B) pooled sequences. Non-SIAT1, pooled, and unpassaged
952 sequences have divergent prediction quality across years, and on average, unpassaged
953 sequences seem to perform better than passaged or pooled sequences. (The distances
954 between dots and the solid black lines are, on average, the smallest for predictions
955 derived from unpassaged sequences.) Data used to generate this figure are available in
956 Supplementary Dataset 8 (A) and Supplementary Dataset 9 (B and C).
957

958 **SUPPLEMENTARY FIGURES AND DATA**





971

972 **Supplementary Figure 2. Differences in correlation strength not due to variance**

973 **from random sampling.** Percent variance in dN/dS explained by inverse distance to

974 site 224 for unpassaged and pooled sequences ($n = 917$). For pooled sequences, sets

975 of 917 sequences were randomly drawn 200 times, with replacement. For each set, we

976 calculated dN/dS and correlated with inverse distance. The R^2 for unpassaged

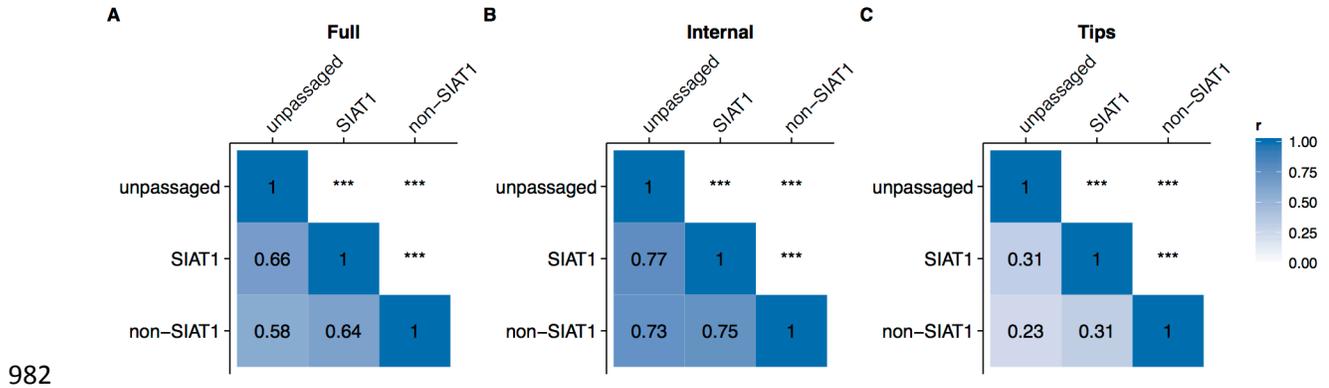
977 sequences falls well outside the distribution of R^2 values for pooled sequences ($z =$

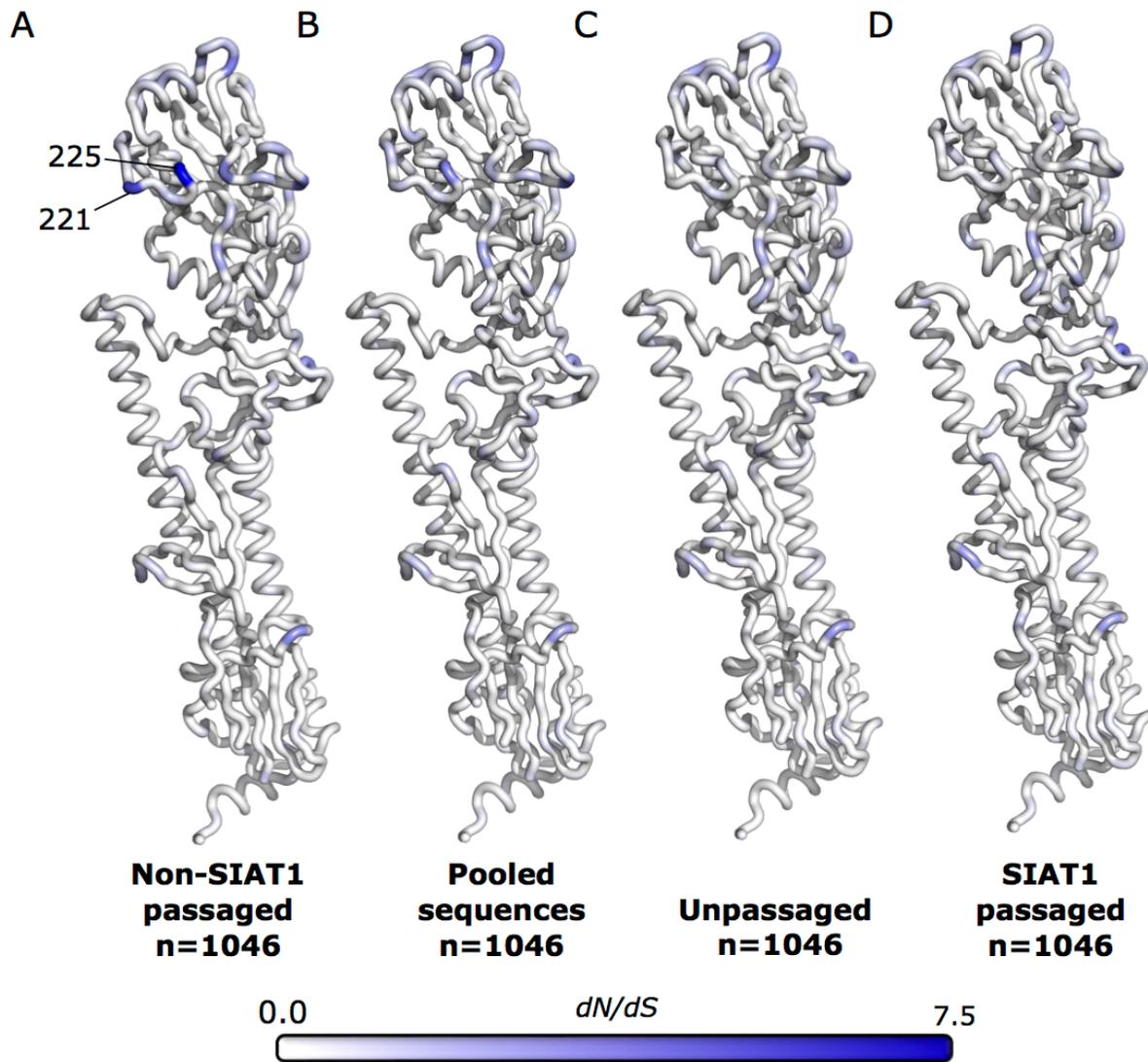
978 -4.22 , $p = 1.2 \times 10^{-5}$). Data used to generate this figure are available in Supplementary

979 Data 3.

980

981





988

989 **Supplementary Figure 4. dN/dS mapped onto the hemagglutinin structure for non-**

990 **SIAT1-passaged, pooled, unpassaged, and SIAT1 passaged sequences. dN/dS**

991 was mapped onto the hemagglutinin structure for (A) non-SIAT1 sequences, (B) pooled

992 sequences, (C) unpassaged sequences, and (D) SIAT1 passaged sequences collected

993 between 2005 and 2015 ($n = 1046$). dN/dS from pooled sequences appears to be

994 intermediate between dN/dS from non-SIAT1-passaged and unpassaged sequences,

995 particularly at sites 221 and 225. Data used to generate this figure are available in
996 Supplementary Data 7.

997 **Supplementary Data 1.** Data used to generate Figures 2 and 3. This file includes 1)
998 sitewise dN/dS values of random draws of 917 unpassaged, generic cell cultured,
999 monkey cell cultured, and the pooled group sequences collected between 2005 and
1000 2015, 2) protein and gene numbering, 3) PDB:2YP7 sequence, 4) relative solvent
1001 accessibilities of the hemagglutinin trimer, 5) linear distances to protein site 224.

1002

1003 **Supplementary Data 2.** Data used to generate Supplementary Figure 1. This file
1004 includes 1) sitewise dN/dS values of random draws of 97 unpassaged, generic cell
1005 cultured, egg cultured, monkey cell cultured, and the pooled group sequences collected
1006 between 2005 and 2015, 2) protein and gene numbering, 3) PDB:2YP7 sequence.

1007

1008 **Supplementary Data 3.** Data used to generate Figure Supplementary Figure 2. This
1009 file includes 1) sitewise dN/dS values of 500 random draws of 917 pooled sequences
1010 and one column of sitewise dN/dS from 917 unpassaged sequences collected between
1011 2005 and 2015, 2) protein and gene numbering, 3) PDB:2YP7 sequence 4) linear
1012 distances to protein site 224.

1013

1014 **Supplementary Data 4.** Data used to generate Figure 4. This file includes 1) sitewise
1015 dN/dS values of random draws of 1046 unpassaged, SIAT1, and non-SIAT1 cell culture
1016 sequences collected between 2005 and 2015, 2) protein and gene numbering, 3)
1017 PDB:2YP7 sequence, 4) relative solvent accessibilities of the hemagglutinin trimer, and
1018 5) linear distances to protein site 224.

1019

1020 **Supplementary Data 5.** Data used to generate Supplementary Figure 3. This file
1021 includes 1) sitewise dN/dS values of random draws of 249 unpassaged, SIAT1 cultured,
1022 and non-SIAT1 cell cultured sequences collected in 2014, 2) protein and gene
1023 numbering, 3) PDB:2YP7 sequence.

1024

1025 **Supplementary Data 6** Data used to generate Figure 5. This file includes 1) sitewise
1026 dN/dS values of random draws of 304 unpassaged and 304 non-SIAT1 cell cultured
1027 sequences passaged once, twice, or 3-5 times collected between 2005 and 2015, 2)
1028 protein and gene numbering, 3) PDB:2YP7 sequence.

1029

1030 **Supplementary Data 7.** Data used to generate Figures 6, 7, and Supplementary Figure
1031 4. This file includes 1) sitewise dN/dS values of random draws of 1703 unpassaged and
1032 non-SIAT1 cell cultured sequences collected between 2005 and 2015, 2) protein and
1033 gene numbering, 3) PDB:2YP7 sequence, 4) sitewise inverse distance correlations.

1034

1035 **Supplementary Data 8.** Data used to generate Figure 8A. This file includes 1) trunk
1036 sitewise dN/dS values of all available unpassaged, SIAT1 cultured, non-SIAT1 cell,
1037 generic cell culture, monkey cell cultured, and pooled sequences collected from 2005-
1038 2015, 2) protein and gene numbering, 3) PDB:2YP7 sequence.

1039

1040 **Supplementary Data 9.** Data used to generate Figure 8B and C. This file includes
1041 output of Local-Branching-Index analyses.

1042

1043 **Supplementary Data 10.** This file includes 1) all sitewise dN/dS values used to
1044 generate figures, 2) protein and gene numbering, 3) PDB:2YP7 sequence, 4) relative
1045 solvent accessibilities of the hemagglutinin trimer, and 5) linear distances to protein site
1046 224.

1047

1048 **Supplementary File 11.** GISAID acknowledgements for hemagglutinin sequences
1049 collected between 1968 and 2015.