

Persistent Host Markers in Pandemic and H5N1 Influenza Viruses

David B. Finkelstein¹, Suraj Mukatira¹, Perdeep K. Mehta¹, John C. Obenauer¹,
Xiaoping Su¹, Robert G. Webster² and Clayton W. Naeve^{1,3*}

¹Hartwell Center for Bioinformatics and Biotechnology, St. Jude Children's Research Hospital, 332 North Lauderdale Street Memphis, TN 38105-2794, USA ²Department of Infectious Diseases, St. Jude Children's Research Hospital, 332 North Lauderdale Street Memphis, TN 38105-2794, USA ³Department of Pathology, University of Tennessee Health Science Center, Memphis, TN 38105, USA

Running Title: Persistent Host Markers in Influenza

Word Count: Abstract 181; Text 5,025

* **Reprint requests to:** C.W. Naeve, Ph.D.

Hartwell Center for Bioinformatics and Biotechnology

St. Jude Children's Research Hospital

332 North Lauderdale Street

Memphis, TN 38105-2794

Phone: 901-495-3689; Fax: 901-495-2945; E-mail: clayton.naeve@stjude.org

Key Words: Influenza, host, markers, H5N1, pandemic

1 **ABSTRACT**

2 Avian influenza viruses have adapted to human hosts causing pandemics in humans.
3 The key host-specific amino acid mutations required for an avian influenza virus to
4 function in humans are unknown. Through multiple sequence alignment and statistical
5 testing of each aligned amino acid we identified markers that discriminate human
6 influenza viruses from avian influenza viruses. We applied strict thresholds to select
7 only markers which are highly preserved in human influenza isolates over time. We
8 found that a subset of these persistent host markers exist in all human pandemic
9 influenza sequences from 1918, 1957 and 1968, while others are acquired as the virus
10 becomes a seasonal influenza. We also show that human H5N1 influenza viruses are
11 significantly more likely to contain the amino acid predominant in human strains for a
12 few persistent host markers when compared to avian H5N1 influenza viruses. This
13 sporadic enrichment of amino acids present in human-hosted viruses may indicate that
14 some H5N1 viruses have made modest adaptations to their new hosts in the recent past.
15 The markers reported here should be useful in monitoring potential pandemic influenza
16 viruses.

17

18 **INTRODUCTION**

19 The three best defined human influenza pandemics of the twentieth century may be
20 derived in whole or part from avian influenza viruses (2, 47, 48), although the avian
21 origin is disputed for the most deadly human pandemic known, the 1918 H1N1
22 “Spanish flu”(4, 15). This virus resulted in the deaths of millions of people worldwide
23 (47). By comparison, avian H5N1 influenza viruses have killed 172 people since 1997
24 (<http://www.who.int/>). Despite containment efforts, H5N1 influenza infections of birds
25 have spread across Asia to Europe, so that the potential for an H5N1 influenza

26 pandemic in humans still exists (21). Therefore, insight into the origin and adaptation
27 of the 1918 H1N1 virus to humans may inform our understanding of the risks posed by
28 H5N1 influenzas circulating in birds today.

29 Recently, large scale influenza genome projects have produced sufficient avian (34) and
30 human (14) sequences to address fundamental questions. Given these resources, our
31 aim was to determine precisely which amino acid changes best distinguish an avian
32 influenza virus from a human influenza virus. After successfully identifying these
33 amino acids, we used them to assess the significance of mutations in H5N1 influenza
34 viruses isolated from humans. Furthermore, we defined a subset of these key amino
35 acids which allowed us to track mutations in the H1N1 influenza lineage over time.

36
37 Although these human influenza viruses are independent isolates, they are not
38 independent of lineage. The exact number of introductions is unknown, but these three
39 influenza pandemics account for the overwhelming number of human influenzas and
40 nearly all of the readily transmissible influenzas. As a result, a distinct amino acid from
41 these three founding strains is more likely to have arisen by coincidence than the large
42 sample size would suggest. Furthermore, the host and lineage parameters are so highly
43 correlated that de-stratification methods based on PCA (39) or other methods (35) will
44 erase the host effect. These methods have proven effective in other studies (39).
45 However, our application of PCA based methods on the influenza sequence data failed
46 to resolve host from lineage. The interpretation of host and lineage are therefore
47 confounded. Specifically, we are precluded from determining whether host-
48 differentiating amino acids are new adaptations or are due to the original lineages based
49 on sequence data alone.

50

51 However, host markers that arose due to lineage may reasonably be of biological
52 importance. As the successful colonization of human beings by influenza required the
53 viruses to overcome selective pressure, even the original founding viruses of each
54 lineage may reasonably be expected to contain important adaptations. Crucially, we can
55 discern likely biologically significant host markers from those that are trivial by
56 examining conservation. Since replication in influenza relies on low fidelity RNA
57 polymerases (41), a high rate of random mutations is observed. Thus, given a large
58 number of strains, we can estimate the expected frequency of amino acid substitution at
59 a given position and compare that estimate to the observed frequency. These estimates
60 presume that the frequency of amino acid substitutions of the viruses do not vary
61 substantially within a host. Variability in the amino acid substitution frequency
62 between hosts is accounted for by our method (see Materials and Methods). Further,
63 the influence of small violations of this assumption is moderated by the averaging effect
64 of calculating the frequency of amino acid substitution across all influenzas within host.
65
66 Under this assumption, positions that are significantly more conserved than expected are
67 likely to be important. In this context, we define conserved as having a low rate of
68 amino acid substitution. Of course, we are not interested in those residues that are
69 conserved in all influenzas generally. Rather, we are interested in residues that are
70 conserved in a host-dependent manner. Thus, positions that are conserved only within
71 humans are deemed biologically significant. Thus, even if these markers arose by
72 chance or selection through the founding of the human influenza lineages, their high
73 degree of persistence despite frequent mutation is evidence of biological relevance.

74

75 MATERIALS AND METHODS

76 Publicly available DNA and protein sequences were downloaded from the Influenza
77 Virus Resource at NCBI (<http://www.ncbi.nlm.nih.gov/genomes/FLU/FLU.html>) as of
78 April 1, 2006. In addition, 847 newly sequenced avian influenza genes (Genbank
79 accessions CY014548 - CY015177) were included. Sequences were retained if they
80 began with methionine and were full length. Virus names were curated to conform to a
81 fixed vocabulary. Ambiguities were manually verified or removed. Non-structural
82 protein 2 (NS2), matrix protein 2 (M2) and polymerase basic protein 1 frame 2 product
83 (PB1-F2) sequences were derived from NS, M and PB1 nucleotide sequences and
84 translated for all downloaded sequences. NS1 sequences of lengths 217, 225, 230 and
85 237 amino acids were also included in the analysis. Qualified sequences were aligned
86 using MUSCLE (12) and classified by serotype, country of origin, host and year.
87 Bayesian analysis trees were generated to guide the manual editing of the HA and NA
88 alignments (23). A total of 9,824 avian and 13,757 human influenza sequences were
89 retained. Indonesian H5N1 human isolates (9) were also downloaded
90 (<http://flu.lanl.gov>) (29) and included in the H5N1 population tests.

91 The data were reformatted for statistical testing such that each aligned position was in
92 its own separate column and each row was from a single strain. Initial surveys revealed
93 repeatedly sequenced examples of the same strain or of viruses from the same outbreak.
94 To reduce this sampling bias, a representative set of genomes was selected. First, all
95 sequences were classed by "outbreak", defined here as the set of all viruses with the
96 same year, host, country and serotype. Next, the sequence which most closely matched
97 the consensus of each outbreak was selected as the representative strain. This was
98 repeated for each outbreak and for each gene. The resulting representative data set
99 contained 6,561 protein sequences (Supplemental Information, data file).

100 We reasoned that H5N1, H7N7, and H9N2 isolates, recently introduced into humans,
101 are not fully adapted, and thus may lack the persistent host markers we were seeking.
102 Therefore these strains were excluded from defining host-specific residues, but were
103 used to validate the results. Next, we identified the most frequent amino acid at each
104 position in each gene of the avian and human isolates among the representative genes.
105 Host specificity was tested at residues where the most frequent in avian influenza
106 isolates differed from that in human influenza isolates. These positions were
107 statistically tested (chi squared test using STATA 9.2/SE) for host specificity and the p-
108 values were adjusted for multiple comparisons using the Bonferroni method. Next, the
109 Euclidean distance between hosts at each position was calculated based on the
110 frequency of each observed amino acid. A vector of amino acid frequencies from avian
111 hosts was compared to a vector of amino acid frequencies from human hosts at each
112 site. The Euclidean distance between these two vectors was then calculated for each
113 site. This number was then divided by the square root of two, the maximum Euclidean
114 distance possible, to create a proportion from zero to one that reflects the percentage of
115 representative isolates that differ between host classes. A minimum proportional
116 Euclidean distance of 0.95 was used as a threshold so that each host marker
117 differentiated at least 95% of the representative isolates.

118

119 For each aligned position, the conservation of the most frequent amino acid found in
120 human influenza isolates was calculated for both avian and human influenzas.
121 Conservation was measured for all human H1N1, H2N2 and H3N2 influenza sequences,
122 and for all avian sequences, not just the representative strains. Next, the conservation
123 frequencies in human influenzas (x) were regressed against the conservation frequencies
124 in avian influenzas (y) and standardized residuals were calculated. The regression line
125 finds the overall difference in variability of influenzas between hosts. Extreme cases of

126 host dependent conservation have extreme standardized residuals of large absolute
127 value. In this regression, large negative standardized residuals are those positions where
128 conservation within human isolates is much larger than observed in birds. Based on z
129 tables, the probability of a single marker having a standardized residual of less than -4
130 is .000032. Given 4,728 positions and 61 discoveries, a false discovery rate (7) of 1%
131 was then calculated for positions which are below -4 standardized residuals. A position
132 was deemed of interest if it passed the standardized residual threshold, had a 99% level
133 of conservation in human influenza viruses and had a proportional Euclidean of 95% or
134 greater. These positions were then designated as persistent host markers.

135 All full-length human influenza virus sequences of the correct serotype in the pandemic
136 year were classed as pandemic viral sequences. For comparison purposes, an H5N1
137 consensus sequence was constructed by surveying all H5N1 viruses. The most frequent
138 amino acid at each position in the pandemic strains and the H5N1 strains were then
139 found. Pandemic markers are those host-differentiating sites where the most frequent
140 amino acid is the same in all three pandemics.

141

142 **RESULTS**

143 **Persistent host markers.** We surveyed 9,824 avian influenza sequences including 847
144 novel avian genes sequenced by this laboratory and 13,757 human influenza sequences.
145 We reduced sample bias by producing a representative data set of 6,561 sequences,
146 minimized false positives by selecting markers with significant host-dependent
147 conservation, and measured the persistence of changes over time (details in Methods
148 and Supplemental Data). Using this approach we identified 32 amino acids, from 4,728
149 aligned positions, that distinguish avian and human influenza virus populations and that
150 met all standards of host differentiation and host dependent conservation. Given a false

151 discovery rate of 1%, we expect that no more than 1 marker persists in human influenza
152 viruses by chance. Given 611 sites that discriminate host to any degree, the median
153 discriminator of host varies 205 times in 1000 human isolates. By comparison, each site
154 we selected varies 10 times or less in 1000 human strains and simultaneously
155 differentiates avian influenza viruses from human influenza viruses with 95% success.

156 These host markers are in five of the 11 proteins tested: RNA polymerase basic protein
157 2 (PB2), RNA polymerase acidic protein (PA), nucleoprotein (NP), matrix protein
158 (M1), and the non-structural protein (NS1). The distribution of these residues among
159 avian and human virus populations and the H1, H2 and H3 pandemic strains isolated
160 during the first year of their respective pandemic is shown in Figure 1. The early 1918
161 (H1) pandemic isolates contain only 13 markers while the subsequent 1957 (H2) and
162 1968 (H3) pandemic strains contain all 32. This is likely due to the fact that these
163 genes/proteins were derived from pre-existing human strains in the H2 and H3
164 pandemic viruses. A fourteenth marker, V100A in the PA protein, is shared by all
165 pandemic strains but one and is thus not 100% conserved. H5N1 isolates, as a
166 population, do not contain any of these markers although isolated cases do exist. By our
167 stringent criteria, there are no persistent host markers in the surface glycoproteins,
168 hemagglutinin (HA) and neuraminidase (NA), or in the RNA polymerase basic protein
169 1 (PB1). This may seem surprising since avian derived HA, NA and PB1 genes have
170 been identified in H2 and H3 pandemic isolates and one might expect to find host
171 markers associated with these proteins, and particularly those residues in HA and NA
172 selected to escape immune surveillance mechanisms (2). However, our methods are
173 designed to identify strictly conserved residues that persist over time and will not
174 capture seasonal changes or even changes between pandemic isolates (see Supplemental
175 text).

176 Remarkably, 26 of the 32 markers (81%) are found in three of the four proteins that
177 form the viral RNA replication complex (NP, PB2, and PA) (24, 44). Fourteen of these
178 markers may be directly associated with the formation of the RNA replication complex
179 as they fall in regions where NP, PB1 and PB2 are known to interact (37, 38) (Figure 2).
180 Six markers in NP fall within known PB2 binding regions (38) and eight markers in
181 PB2 are in regions of the molecule known to bind to either NP or PB1 (37). Two
182 additional PB2 markers may influence RNA replication indirectly. The residue at 475 in
183 the polymerase gene PB2 is predominately leucine (L) in avian isolates and methionine
184 (M) human strains. This marker, L475M, is in a domain necessary for nuclear
185 importation (30) and residue D567N is in the RNA cap binding region (20). The
186 implication is that nuclear importation and formation of the RNA polymerase complex
187 is influenced by the host environment. Less clear is the role PA plays in RNA
188 replication and, consequently, the functional significance of the 10 markers in PA. One
189 host marker, S225C, is located in a region involved in nuclear localization (32). The
190 remaining 9 markers in PA are in regions of unknown or ambiguous functional
191 importance.

192 The remaining six persistent host markers are found in M1 and NS1 proteins and are
193 located in regions generally associated with binding to host cell proteins. All three host
194 markers in M1 are in the C-terminal half of the molecule known to bind heat shock
195 protein Hsc70 in host cells (49). Hsc70 has been shown to enhance viral replication
196 through interaction with M1 (49). All three markers in the non-structural protein NS1
197 are located in regions of the molecule with known host cell binding functions. The N-
198 terminal domain of NS1 binds to the 30 kDa subunit of the cleavage and
199 polyadenylation specificity factor (CPSF) and to eukaryotic translation initiation factor
200 4 gamma 1 (eIF4G1) (5, 8) and contains the host marker I81M. In the course of
201 identifying these markers in NS1 we identified a previously unreported SRC-homology
202 3 (SH3) motif. One of our host markers, P215T, is in this SH3 recognition motif where

203 the PPLPP motif is preserved in avian influenzas and is altered to PPLTP in human
204 influenzas. The third persistent host marker in NS1 is at residue R227E in the PDZ
205 (postsynaptic density, PSD-95; discs large, Dlg; zonula occludens-1, ZO-1) binding
206 domain we previously identified at the C-terminus of the molecule and demonstrated to
207 interact with numerous human PDZ domains *in vitro* (34). This region is also known to
208 bind poly-A binding protein II, PABPII (11, 26). Given the NS1 protein's role in
209 suppressing the host cell immune response through binding host proteins and host RNA,
210 it may be that markers in this molecule point to key mutations needed to improve the
211 immune suppression function of NS1 and enhance viral replication (25, 27).

212 Overall, persistent host markers are typically found in RNA replication complex
213 proteins and are often located in known protein binding domains. These regions may
214 directly influence the RNP replication complex or they may enhance replication through
215 the interaction with host factors.

216 **Persistent host markers in pandemic isolates.** We next focused on the early isolates
217 of pandemic influenza viruses to determine which markers they might have acquired.
218 We found that 13 of our 32 host markers (Figure 1, arrows) are absolutely conserved
219 (100%) in the influenza viruses that caused the 1918, 1957 and 1968 pandemics and are
220 distributed among four viral genes; PB2, PA, NP and M1 (Table S1). Again the
221 majority of these markers reside in RNA replication complex proteins. We should
222 emphasize that it is unlikely that all 13 pandemic markers must be acquired to gain any
223 single phenotypic trait of pandemic influenza viruses such as efficient replication, tissue
224 tropism or transmissibility. Further, we cannot estimate how long it would take an
225 avian virus such as H5N1 to acquire these traits.

226 **DISCUSSION**

227 **Pandemic versus seasonal influenza.** Although we cannot determine the rate at which
228 avian isolates would acquire the 13 “pandemic” host markers, we can look at historical
229 data to determine whether early pandemic isolates acquired additional markers over
230 time. We can do this only for H1N1 isolates as they represent the introduction of all 8
231 influenza virus genes from an avian precursor, they have circulated in humans for 88
232 years, and many isolates have been sequenced. In contrast, subsequent pandemics
233 involved the introduction of only HA, NA, and PB1 genes from an avian isolate into a
234 pre-existing human strain, none of which carry host markers as defined by our criteria.
235 If we plot the proportion of host markers in M1, NP, NS1, PA and PB2 proteins over
236 time (Figure 3) we see that early H1N1 isolates, already containing 13 of the amino
237 acids prevalent in human-hosted viruses, acquire the remaining 19 within 10-20 years
238 depending on the protein. The stepped progression of PA and PB2 markers suggests
239 that the H1N1 pandemic influenza adapted to human hosts in stages. In contrast, NS1
240 marker acquisition appears to be more abrupt. However, this is likely a sampling
241 artefact as there are no H1N1 human influenza sequences available during the years
242 1919 to 1932. Further this abruptness may be due to the relatively few number of
243 markers in NS1 and M1. Unlike these other genes, the host markers in proteins NP and
244 M1 do not appear to be stably preserved, despite passing our 99% persistence criteria
245 (see methods). The instability of these markers in these genes may be due to the
246 reintroduction of H1N1 viruses from swine or birds or to seasonal variation in the
247 human host.

248 The progressive changes seen in H1N1 human influenza isolates implies that these
249 viruses gradually acquired mutations that confer the phenotypic traits of seasonal
250 viruses and that these additional sites are not required for an influenza virus to cause a
251 pandemic. Rather it is likely that these additional mutations are associated with the
252 traits of seasonal influenza viruses such as low mortality. Over time, through
253 successive rounds of transmission and selection, we would expect avian influenzas, like

254 H5N1, introduced into humans to acquire all 32 persistent host markers seen in seasonal
255 influenza viruses.

256 **Persistent host markers in H5N1 viruses.** We examined H5N1 influenza sequences
257 from avian hosts and compared them to H5N1 influenzas isolated from humans,
258 focusing on the 32 persistent host markers. We included 7 H5N1 strains recently
259 reported to transmit within a family in Indonesia (9). Although the predominant amino
260 acid found in H5N1 isolates is consistent with avian influenzas at most marker
261 locations, in a fraction of H5N1 isolates, the amino acid prevalent in human-hosted
262 viruses has been acquired. We found four sites that are significantly enriched ($p <$
263 0.0001) in human H5N1 isolates (Table 1). Three of the four host markers that are
264 enriched in H5N1 are also 100% conserved in human pandemic isolates. These three
265 host markers are in PB2, one of which is the well known marker E627K. This mutation
266 was seen in all seven of the putatively human-transmissible Indonesian H5N1 viruses
267 (9) and in the 2003 H7N7 outbreak in the Netherlands (13). Four of the seven
268 Indonesian strains also have the PB2 host marker K702R. This novel marker site is
269 adjacent to a known high-pathogenicity site in PB2 at residue D701N (27). The
270 enrichment of four host markers in H5N1 isolates suggests H5N1 influenza can adapt to
271 human hosts. However, no single H5N1 virus sequenced contains more than two of
272 these four sites.

273 The polymerase protein PB2 appears critical to adaptation of avian viruses to humans,
274 based on this and other studies (9, 13). Significantly, we identify 10 PB2 host markers
275 here (Table A1). These are all high quality discriminators of host (95% or greater) and
276 all of these sites are preserved in 99% of human H1N1, H2N2 and H3N2 sequences
277 over time. Not only does PB2, along with PA, have the most persistent host markers, it
278 also has A199S, E627K, and K702R. These residues are the only host markers that are
279 absolutely (100%) conserved in all pandemic influenza isolates we surveyed (Table A1)

280 and are also enriched in the population of human H5N1 isolates (Table 1). We suspect
281 that acquisition of the amino acids that are prevalent in humans are required for the
282 evolution of an avian influenza virus like H5N1 into a virus that is capable of causing a
283 human pandemic. Here we must note that the sporadic and modest acquisition of
284 markers in H5N1 human isolates and the stability of the H5N1 avian isolates indicate
285 that currently circulating H5N1 viruses are no more adapted to human hosts today than
286 they were in the past. What has changed is the geographic dispersion of the H5N1 virus
287 and thus the size of the population at risk. Therefore the current risk of an H5N1
288 influenza pandemic in humans is due to an increased frequency of human exposure to
289 the H5N1 virus from birds, rather than to a human adapted H5N1 virus.

290 Interestingly, two of these persistent host markers in PB2 occur in a unique set of four
291 H5N1 human isolates from Indonesia. These Indonesian influenza isolates are
292 distinguished from nearly all other human H5N1 isolates in that they may be acquired
293 by human to human transmission rather than by avian to human transmission (9).
294 Although the numbers are too small to allow a valid statistical test, these H5N1 isolates
295 from a single Indonesian family appear more adapted to humans than the other H5N1
296 human isolates presumably acquired directly from birds. The residue A199S in PB2 is
297 the only marker that is absolutely conserved in the seasonal human influenza isolates we
298 surveyed (Table A1).

299 In summary, we have examined large collections of both avian and human influenza
300 protein sequences and identified persistent host markers across the influenza proteome.
301 By minimizing false positives and by focusing on those sites preserved over time in a
302 host dependent manner, we have identified a set of 32 amino acids which are persistent
303 host markers. These include both well known and novel sites, including a potential SH3
304 binding motif in NS1. By tracking the acquisition of these sites over time, we observed
305 evidence of progressive adaptation of the avian H1N1 virus to human hosts. We show

306 that 13 of the 32 persistent host markers are 100% conserved amino acid changes in
307 pandemic viruses and suggest these are likely important in the evolution of pandemic
308 influenza. Further, we show that a small fraction of the population of H5N1 isolates
309 from humans have acquired four of these 32 markers, although no single H5N1 isolate
310 surveyed contains more than two markers and current H5N1 viruses are no more
311 adapted to humans today than they were in the past (Table 1).

312 APPENDIX

313 Details concerning bioinformatic and statistical methods are provided in this
314 supplement. The persistence percentages, and proportional Euclidean distances are
315 given in Table A1. A summary graphic of the proportional Euclidean distances by
316 genes is provided in Figure A1 and a frequency table of a key HA amino acid is shown
317 in Table A2. Further discussion of the statistical methods and results is provided
318 including a summary table of sample sizes in the representative set by gene. A brief
319 discussion of the protein interaction regions and the appropriate references are also
320 included here. The supplemental information includes an Excel file listing the
321 accession numbers, year, serotype and country of origin for all 6,561 representative
322 proteins.

323

324 **Multiple Sequence Alignments.** The protein sequences were aligned using the
325 MUSCLE (12) program. The MUSCLE program was chosen due to its performance,
326 flexibility, and the speed with which it aligns a large set of sequences. The protein
327 alignments were manually inspected and edited using BioEdit. Nucleotide sequences
328 were then aligned based on the protein alignments using the tranalign program in the
329 EMBOSS package. We found that protein-guided alignments of nucleotide sequences
330 produced better alignments than aligning the nucleotides directly. After generating

331 maximum likelihood trees, the sequences in each multiple alignment were re-ordered to
332 match the ordering in the trees for easy visual comparison. The clade-guided re-
333 alignment of nucleotide and protein sequences helped further improve the quality of
334 alignments by manual checking and editing especially for HA and NA genes in the
335 highly variable regions. Custom Perl scripts and additional EMBOSS tools were used
336 to facilitate this process.

337 **Potential for false negatives.** One might expect, *a priori*, to find host markers in the
338 surface glycoproteins HA and NA because of immune pressure and because of the
339 receptor specificity of the HA receptor binding site (16, 40, 45, 50) or in the polymerase
340 protein PB1 because of its association with HA/NA in the H2 and H3 pandemic strains.
341 However, in this study, as a result of stringent criteria designed to eliminate false
342 positives, authentic host adaptations may have been lost. As noted in the text, there are
343 no host markers in the surface glycoproteins HA and NA or in the polymerase protein
344 PB1. All amino acid markers from the genes HA, NA, and PB1 as well as the alternate
345 transcripts NS2, M2 and PB1-F2 were either poor quality host discriminators
346 (proportional Euclidean distance < 0.95) or were not preserved in human strains over
347 time (persistence < 99%). Host specific residues in HA have been reported elsewhere
348 (40), but HA residues do not differentiate more than 72% of viruses by host in this
349 broad study (proportional Euclidean distance of 0.72). Two studies have also reported
350 host-specific M2 sites (10, 28), however, these sites failed to pass thresholds used in this
351 study. The best M2 site V86A (V is the predominant avian residue and A is the
352 predominant human residue), did pass the Euclidean distance test, but failed the 99%
353 persistence test. Thus, by our stringent criteria, this M2 site is a valid host-specific
354 marker, but was excluded as a persistent host marker because it was not sufficiently
355 preserved in human influenzas over time.

356 In addition, our methods test each residue separately, so that if host specific pressure
357 can be relieved at any number of sites, then the pressure to conserve a given site is
358 reduced as is the high degree of differentiation at that site. Direct evidence indicates
359 that HA receptor specificity can be altered by mutations at any one of several sites(16,
360 40, 50). Furthermore our survey of all HA sequences indicates that the amino acids at
361 key sites such as 226 in the HA receptor binding site are well preserved among avian
362 influenza isolates, but are not well preserved among human influenza isolates (Table
363 A2). We recognize that accurate alignments of HA and NA are hampered by high
364 variability and despite the care taken in manual editing, false negative errors may occur
365 due to alignment errors. While it might be possible to improve these alignments by
366 adding structural data, this data only exists for portions of each protein and for only a
367 few serotypes. Finally, while a lack of markers in HA and NA proteins are a concern,
368 we note that there is also a lack of markers in PB1, which was trivial to align due to
369 high conservation. Thus, it may be that residues in HA, NA, PB1 and PB1-F2 are
370 simply less host-differentiating than are other genes, as we have observed.

371 **Statistical Tests.** All statistical tests in this paper are performed on categorical data.
372 For each position we compared the frequency of amino acid categories across host using
373 a two sided chi-squared test. This test assumed independence of the categories and is in
374 common usage. For each position, the table size varied in accordance with the number
375 of amino acid types. For the host test, we decided not to fix the table size at 2 by 20 to
376 minimize table sparsity and to avoid false discoveries due to excessive degrees of
377 freedom. We relied on the strictness of the Bonferroni correction and the application of
378 absolute quality metrics to minimize false discoveries.

379 In total there were 4,728 aligned positions. Of these 611 positions passed the initial
380 screening. Including the initial screens as informal hypothesis tests there were 4,728
381 tests. The Bonferroni threshold at the 0.05 alpha level was 1.06 e-05. There were 599

382 positions that passed Bonferroni criteria. As described in the methods section,
383 applying the regression standardized residual threshold at -4 reduced this set to 61. This
384 list was further reduced by the use of absolute metrics, proportional Euclidean distance
385 and percent persistence. These metrics were chosen arbitrarily to guarantee minimum
386 quality standards. Sample sizes were dependent on the gene of interest (Table A3).

387 The 32 remaining positions were then tested in H5N1 isolates for enrichment in human
388 versus avian hosts. As this test generated two-by-two tables, Fisher's exact test was
389 used. This two tailed test gives the most accurate p-value available under the
390 independence assumptions and is computationally feasible for small tables. Again a
391 Bonferroni threshold at the 0.05 alpha level was applied for these 32 tests. All results in
392 Table 1 pass the Bonferroni threshold of 0.0015625. Sample size for the PB2 tests in
393 Table 1 was 214 and the sample size for the PA test was 196.

394 **Regions of Protein Interaction.** Known regions of protein interaction in Figure 2 in
395 the main text were derived from the literature. The M1 interaction regions are based on
396 crystal structure (1) and viral assembly studies (3, 6, 18, 19, 31, 46). M1 interacts with
397 other influenza viral proteins PA, NP, HA, PB1, and NS2 (18). M1 Also interacts with
398 Hsc70 (36). NS1 interacts with binds eukaryotic initiation factor 4GI at the N terminus
399 (26), and has several protein binding domains (11, 25, 26), including a SH3 domain
400 reported here for the first time (33). NP, PB1 (not pictured in Figure 2) and PB2 bind
401 each other in RNA polymerase complex formation (37, 38). PA has casein kinase II
402 sites (42), and binds to the RNA cap (17) and host proteins (22). PA also has a large
403 proteolytic region (43) (from 1-247) that may interact with M1 (18); however, specific
404 sites in PA are not known so this region was not included in the figure.

405 **Accession Numbers.** For Genbank or Genpept accession numbers of the genes
406 included in the representative set are available on as a supplement.

407

408 **ACKNOWLEDGEMENTS**

409 This work was supported by the American Lebanese Syrian Associated Charities
410 (ALSAC).

411 The authors gratefully acknowledge the technical comments and text editing of Caroline
412 Obert, the editorial assistance of Geoff Neale, Yiping Fan, and Jinhua Wang and the
413 statistical advice of Stan Pounds.

414

415 **REFERENCES**

- 416 1. **Akarsu, H., W. P. Burmeister, C. Petosa, I. Petit, C. W.**
417 **Muller, R. W. Ruigrok, and F. Baudin.** 2003. Crystal
418 structure of the M1 protein-binding domain of the influenza
419 A virus nuclear export protein (NEP/NS2). *Embo J* **22**:4646-
420 55.
- 421 2. **Alexander, D. J.** 2006. Avian influenza viruses and human
422 health. *Dev Biol (Basel)* **124**:77-84.
- 423 3. **Ali, A., R. T. Avalos, E. Ponimaskin, and D. P. Nayak.**
424 2000. Influenza virus assembly: effect of influenza virus
425 glycoproteins on the membrane association of M1 protein. *J*
426 *Virology* **74**:8709-19.
- 427 4. **Antonovics, J., M. E. Hood, and C. H. Baker.** 2006.
428 Molecular virology: was the 1918 flu avian in origin? *Nature*
429 **440**:E9; discussion E9-10.

- 430 5. **Aragon, T., S. de la Luna, I. Novoa, L. Carrasco, J.**
431 **Ortin, and A. Nieto.** 2000. Eukaryotic translation initiation
432 factor 4GI is a cellular target for NS1 protein, a
433 translational activator of influenza virus. *Mol Cell Biol*
434 **20:6259-68.**
- 435 6. **Barman, S., A. Ali, E. K. Hui, L. Adhikary, and D. P.**
436 **Nayak.** 2001. Transport of viral proteins to the apical
437 membranes and interaction of matrix protein with
438 glycoproteins in the assembly of influenza viruses. *Virus*
439 *Res* **77:61-9.**
- 440 7. **Benjamini, Y., and Y. Hochberg.** 1995. Controlling the
441 False Discovery Rate - a Practical and Powerful Approach to
442 Multiple Testing. *Journal of the Royal Statistical Society*
443 *Series B-Methodological* **57:289-300.**
- 444 8. **Burgui, I., T. Aragon, J. Ortin, and A. Nieto.** 2003.
445 PABP1 and eIF4GI associate with influenza virus NS1
446 protein in viral mRNA translation initiation complexes. *J*
447 *Gen Virol* **84:3263-74.**
- 448 9. **Butler, D.** 2006. Family tragedy spotlights flu mutations.
449 *Nature* **442:114-5.**
- 450 10. **Chen, G. W., S. C. Chang, C. K. Mok, Y. L. Lo, Y. N.**
451 **Kung, J. H. Huang, Y. H. Shih, J. Y. Wang, C. Chiang, C.**
452 **J. Chen, and S. R. Shih.** 2006. Genomic signatures of
453 human versus avian influenza A viruses. *Emerg Infect Dis*
454 **12:1353-60.**

- 455 11. **Chen, Z., Y. Li, and R. M. Krug.** 1999. Influenza A virus
456 NS1 protein targets poly(A)-binding protein II of the
457 cellular 3'-end processing machinery. *Embo J* **18**:2273-83.
- 458 12. **Edgar, R. C.** 2004. MUSCLE: multiple sequence alignment
459 with high accuracy and high throughput. *Nucleic Acids Res*
460 **32**:1792-7.
- 461 13. **Fouchier, R. A., P. M. Schneeberger, F. W. Rozendaal, J.**
462 **M. Broekman, S. A. Kemink, V. Munster, T. Kuiken, G. F.**
463 **Rimmelzwaan, M. Schutten, G. J. Van Doornum, G. Koch,**
464 **A. Bosman, M. Koopmans, and A. D. Osterhaus.** 2004.
465 Avian influenza A virus (H7N7) associated with human
466 conjunctivitis and a fatal case of acute respiratory distress
467 syndrome. *Proc Natl Acad Sci U S A* **101**:1356-61.
- 468 14. **Ghedini, E., N. A. Sengamalay, M. Shumway, J. Zaborsky,**
469 **T. Feldblyum, V. Subbu, D. J. Spiro, J. Sitz, H. Koo, P.**
470 **Bolotov, D. Dernovoy, T. Tatusova, Y. Bao, K. St George,**
471 **J. Taylor, D. J. Lipman, C. M. Fraser, J. K.**
472 **Taubenberger, and S. L. Salzberg.** 2005. Large-scale
473 sequencing of human influenza reveals the dynamic nature
474 of viral genome evolution. *Nature* **437**:1162-6.
- 475 15. **Gibbs, M. J., and A. J. Gibbs.** 2006. Molecular virology:
476 was the 1918 pandemic caused by a bird flu? *Nature* **440**:E8;
477 discussion E9-10.
- 478 16. **Glaser, L., J. Stevens, D. Zamarin, I. A. Wilson, A.**
479 **Garcia-Sastre, T. M. Tumpey, C. F. Basler, J. K.**

- 480 **Taubenberger, and P. Palese.** 2005. A single amino acid
481 substitution in 1918 influenza virus hemagglutinin changes
482 receptor binding specificity. *J Virol* **79**:11533-6.
- 483 17. **Hara, K., F. I. Schmidt, M. Crow, and G. G. Brownlee.**
484 2006. Amino acid residues in the N-terminal region of the
485 PA subunit of influenza A virus RNA polymerase play a
486 critical role in protein stability, endonuclease activity, cap
487 binding, and virion RNA promoter binding. *J Virol* **80**:7789-
488 98.
- 489 18. **Hara, K., M. Shiota, H. Kido, K. Watanabe, K. Nagata,**
490 **and T. Toyoda.** 2003. Inhibition of the protease activity of
491 influenza virus RNA polymerase PA subunit by viral matrix
492 protein. *Microbiol Immunol* **47**:521-6.
- 493 19. **Harris, A., F. Forouhar, S. Qiu, B. Sha, and M. Luo.**
494 2001. The crystal structure of the influenza matrix protein
495 M1 at neutral pH: M1-M1 protein interfaces can rotate in the
496 oligomeric structures of M1. *Virology* **289**:34-44.
- 497 20. **Honda, A., K. Mizumoto, and A. Ishihama.** 1999. Two
498 separate sequences of PB2 subunit constitute the RNA cap-
499 binding site of influenza virus RNA polymerase. *Genes Cells*
500 **4**:475-85.
- 501 21. **Horimoto, T., and Y. Kawaoka.** 2005. Influenza: lessons
502 from past pandemics, warnings from current incidents. *Nat*
503 *Rev Microbiol* **3**:591-600.

- 504 22. **Huarte, M., J. J. Sanz-Ezquerro, F. Roncal, J. Ortin, and**
505 **A. Nieto.** 2001. PA subunit from influenza virus polymerase
506 complex interacts with a cellular protein with homology to a
507 family of transcriptional activators. *J Virol* **75**:8597-604.
- 508 23. **Huelsenbeck, J. P., and F. Ronquist.** 2001. MRBAYES:
509 Bayesian inference of phylogenetic trees. *Bioinformatics*
510 **17**:754-5.
- 511 24. **Kawaguchi, A., T. Naito, and K. Nagata.** 2005.
512 Involvement of influenza virus PA subunit in assembly of
513 functional RNA polymerase complexes. *J Virol* **79**:732-344.
- 514 25. **Krug, R. M., W. Yuan, D. L. Noah, and A. G. Latham.**
515 2003. Intracellular warfare between human influenza viruses
516 and human cells: the roles of the viral NS1 protein. *Virology*
517 **309**:181-9.
- 518 26. **Li, Y., Y. Yamakita, and R. M. Krug.** 1998. Regulation of
519 a nuclear export signal by an adjacent inhibitory sequence:
520 the effector domain of the influenza virus NS1 protein. *Proc*
521 *Natl Acad Sci U S A* **95**:4864-9.
- 522 27. **Li, Z., H. Chen, P. Jiao, G. Deng, G. Tian, Y. Li, E.**
523 **Hoffmann, R. G. Webster, Y. Matsuoka, and K. Yu.** 2005.
524 Molecular basis of replication of duck H5N1 influenza
525 viruses in a mammalian mouse model. *J Virol* **79**:12058-64.
- 526 28. **Liu, W., P. Zou, J. Ding, Y. Lu, and Y. H. Chen.** 2005.
527 Sequence comparison between the extracellular domain of
528 M2 protein human and avian influenza A virus provides new

- 529 information for bivalent influenza vaccine design. *Microbes*
530 *Infect* **7**:171-7.
- 531 29. **Macken, C., H. Lu, J. Goodman, and L. Boykin.** 2001. The
532 value of a database in surveillance and vaccine selection, p.
533 103-106. *In* A. D. Osterhaus, N. J. Cox, and A. W. Hampson
534 (ed.), *Options for the Control of Influenza IV*. Elsevier
535 Science, Amsterdam.
- 536 30. **Mukaigawa, J., and D. P. Nayak.** 1991. Two signals
537 mediate nuclear localization of influenza virus (A/WSN/33)
538 polymerase basic protein 2. *J Virol* **65**:245-53.
- 539 31. **Nayak, D. P., E. K. Hui, and S. Barman.** 2004. Assembly
540 and budding of influenza virus. *Virus Res* **106**:147-65.
- 541 32. **Nieto, A., S. de la Luna, J. Barcena, A. Portela, and J.**
542 **Ortin.** 1994. Complex structure of the nuclear translocation
543 signal of influenza virus polymerase PA subunit. *J Gen Virol*
544 **75 (Pt 1)**:29-36.
- 545 33. **Obenauer, J. C., L. C. Cantley, and M. B. Yaffe.** 2003.
546 Scansite 2.0: Proteome-wide prediction of cell signaling
547 interactions using short sequence motifs. *Nucleic Acids Res*
548 **31**:3635-41.
- 549 34. **Obenauer, J. C., J. Denson, P. K. Mehta, X. Su, S.**
550 **Mukatira, D. B. Finkelstein, X. Xu, J. Wang, J. Ma, Y.**
551 **Fan, K. M. Rakestraw, R. G. Webster, E. Hoffmann, S.**
552 **Krauss, J. Zheng, Z. Zhang, and C. W. Naeve.** 2006.

- 553 Large-scale sequence analysis of avian influenza isolates.
554 Science **311**:1576-80.
- 555 35. **Patterson, N., A. L. Price, and D. Reich.** 2006. Population
556 Structure and Eigenanalysis. PLoS Genet **2**:e190.
- 557 36. **Perez-Gonzalez, A., A. Rodriguez, M. Huarte, I. J.**
558 **Salanueva, and A. Nieto.** 2006. hCLE/CGI-99, a human
559 protein that interacts with the influenza virus polymerase, is
560 a mRNA transcription modulator. J Mol Biol **362**:887-900.
- 561 37. **Poole, E., D. Elton, L. Medcalf, and P. Digard.** 2004.
562 Functional domains of the influenza A virus PB2 protein:
563 identification of NP- and PB1-binding sites. Virology
564 **321**:120-33.
- 565 38. **Portela, A., and P. Digard.** 2002. The influenza virus
566 nucleoprotein: a multifunctional RNA-binding protein
567 pivotal to virus replication. J Gen Virol **83**:723-34.
- 568 39. **Price, A. L., N. J. Patterson, R. M. Plenge, M. E.**
569 **Weinblatt, N. A. Shadick, and D. Reich.** 2006. Principal
570 components analysis corrects for stratification in genome-
571 wide association studies. Nat Genet **38**:904-9.
- 572 40. **Rogers, G. N., J. C. Paulson, R. S. Daniels, J. J. Skehel, I.**
573 **A. Wilson, and D. C. Wiley.** 1983. Single amino acid
574 substitutions in influenza haemagglutinin change receptor
575 binding specificity. Nature **304**:76-8.
- 576 41. **Sallie, R.** 2005. Replicative homeostasis II: influence of
577 polymerase fidelity on RNA virus quasispecies biology:

578 implications for immune recognition, viral autoimmunity and
579 other "virus receptor" diseases. *Virology* **2**:70.

580 42. **Sanz-Ezquerro, J. J., J. Fernandez Santaren, T. Sierra, T.**
581 **Aragon, J. Ortega, J. Ortin, G. L. Smith, and A. Nieto.**
582 1998. The PA influenza virus polymerase subunit is a
583 phosphorylated protein. *J Gen Virol* **79 (Pt 3)**:471-8.

584 43. **Sanz-Ezquerro, J. J., T. Zurcher, S. de la Luna, J. Ortin,**
585 **and A. Nieto.** 1996. The amino-terminal one-third of the
586 influenza virus PA protein is responsible for the induction
587 of proteolysis. *J Virol* **70**:1905-11.

588 44. **Sidorenko, Y., and U. Reichl.** 2004. Structured model of
589 influenza virus replication in MDCK cells. *Biotechnol*
590 *Bioeng* **88**:1-14.

591 45. **Stevens, J., O. Blixt, L. Glaser, J. K. Taubenberger, P.**
592 **Palese, J. C. Paulson, and I. A. Wilson.** 2006. Glycan
593 microarray analysis of the hemagglutinins from modern and
594 pandemic influenza viruses reveals different receptor
595 specificities. *J Mol Biol* **355**:1143-55.

596 46. **Takeuchi, H., A. Okada, and T. Miura.** 2003. Roles of the
597 histidine and tryptophan side chains in the M2 proton
598 channel from influenza A virus. *FEBS Lett* **552**:35-8.

599 47. **Taubenberger, J. K., A. H. Reid, R. M. Lourens, R. Wang,**
600 **G. Jin, and T. G. Fanning.** 2005. Characterization of the
601 1918 influenza virus polymerase genes. *Nature* **437**:889-93.

- 602 48. **Tumpey, T. M., A. Garcia-Sastre, J. K. Taubenberger, P.**
603 **Palese, D. E. Swayne, M. J. Pantin-Jackwood, S. Schultz-**
604 **Cherry, A. Solorzano, N. Van Rooijen, J. M. Katz, and C.**
605 **F. Basler.** 2005. Pathogenicity of influenza viruses with
606 genes from the 1918 pandemic virus: functional roles of
607 alveolar macrophages and neutrophils in limiting virus
608 replication and mortality in mice. *J Virol* **79**:14933-44.
- 609 49. **Watanabe, K., T. Fuse, I. Asano, F. Tsukahara, Y. Maru,**
610 **K. Nagata, K. Kitazato, and N. Kobayashi.** 2006.
611 Identification of Hsc70 as an influenza virus matrix protein
612 (M1) binding factor involved in the virus life cycle. *FEBS*
613 *Lett* **580**:5785-90.
- 614 50. **Yamada, S., Y. Suzuki, T. Suzuki, M. Q. Le, C. A. Nidom,**
615 **Y. Sakai-Tagawa, Y. Muramoto, M. Ito, M. Kiso, T.**
616 **Horimoto, K. Shinya, T. Sawada, M. Kiso, T. Usui, T.**
617 **Murata, Y. Lin, A. Hay, L. F. Haire, D. J. Stevens, R. J.**
618 **Russell, S. J. Gamblin, J. J. Skehel, and Y. Kawaoka.**
619 2006. Haemagglutinin mutations responsible for the binding
620 of H5N1 influenza A viruses to human-type receptors.
621 *Nature* **444**:378-82.

622
623
624

625 **TABLE 1.** Host markers are enriched in human H5N1 influenza viruses.

Gene	Mutation	Strains with a Human Adaptation Marker				p-value
		Avian H5N1 frequency %		Human H5N1 frequency %		
<i>PB2</i>	A199S	0/177	0%	7/37	19%	2.79E-06
	E627K	22/177	12%	20/37	54%	1.62E-07
	K702R	0/177	0%	6/37	16%	1.87E-05
<i>PA</i>	S409N	5/162	3%	11/34	32%	2.20E-06

626 P values are from Fisher's exact test.

627

ACCEPTED

628 **TABLE A1.** The 32 persistent host markers.

	Gene	Position	Distance	Human	H1N1 1918	H2N2 1957	H3N2 1968	H5N1	Avian	Persistence
<i>MI</i>		115	0.967	I	V	I	I	V	V	99.27%
		121	0.962	A	A	A	A	T	T	99.92%
		137	0.958	A	T	A	A	T	T	99.10%
<i>NP</i>		16	0.953	D	D	D	D	G	G	99.41%
		61	0.973	L	I	L	L	I	I	99.32%
		283	0.981	P	P	P	P	L	L	99.24%
		305	0.96	K	R	K	K	R	R	99.07%
		313	0.973	Y	Y	Y	Y	F	F	99.32%
<i>NS</i>		357	0.964	K	K	K	K	Q	Q	99.32%
		81	0.958	M	I	M	M	deleted	I	99.32%
		215	0.955	T	P	T	T	P	P	99.74%
<i>PA</i>		227	0.966	R	K	R	R	E	E	99.40%
		28	0.988	L	L	L	L	P	P	99.45%
		55	0.984	N	N	N	N	D	D	99.73%
		57	0.958	Q	R	Q	Q	R	R	99.27%
		100	0.955	A	A	A	A	V	V	99.54%
		225	0.969	C	S	C	C	S	S	99.36%
		268	0.951	I	L	I	I	L	L	99.09%
		337	0.978	S	A	S	S	A	A	99.91%
		404	0.967	S	A	S	S	A	A	99.27%
		409	0.959	N	S	N	N	S	S	99.54%
		552	0.999	S	S	S	S	T	T	99.91%
	<i>PB2</i>		44	0.966	S	A	S	S	A	A
		64	0.954	T	M	T	T	I	M	99.82%
		199	0.997	S	S	S	S	A	A	100.00%
		271	0.958	A	T	A	A	T	T	99.38%
		475	0.994	M	M	M	M	L	L	99.91%
		567	0.977	N	N	N	N	D	D	99.29%
		588	0.971	I	A	I	I	A	A	99.38%
		627	0.977	K	K	K	K	E	E	99.82%
	674	0.969	T	A	T	T	A	A	99.47%	
	702	0.955	R	R	R	R	K	K	99.38%	

629 Note each column has the most frequently occurring amino acid by class. Position is the
 630 location in the protein sequence. Distance refers to the proportional Euclidean distance
 631 of amino acid frequency between human and avian hosted viruses.

632 **TABLE A2.** The frequency of residues at position 226 varies by host in influenza A
 633 haemagglutinin.

Amino Acid	Avian virus	Human virus	Row total
I	2	220	222
L	49	125	174
M	1	0	1
P	1	0	1
Q	1007	298	1305
R	0	1	1
V	2	626	628
ambiguity	0	2	2
deletion	1	0	1
Column total	1063	1272	2335

634

635 **TABLE A3.** Sample sizes by gene for statistical testing.

636

Gene	Sample size
HA	606
M1	697
M2	690
NA	683
NP	573
NS1	681
NS2	699
PA	481
PB1	490
PB1F2	481
PB2	480

637 **FIGURE LEGENDS**

638 **Fig 1. Host differentiating sites are compared to pandemic strains.** Each of the 32
639 host-differentiating sites are displayed and color-coded by host. Avian is in blue, human
640 in yellow. The intensity of each position is determined by the proportional Euclidean
641 distance between hosts. Positions where the consensus residue of each pandemic strain
642 agrees with the most frequent human amino acid are boxed. The 13 positions where all
643 pandemic isolates surveyed absolutely agree with the most frequent human amino acid
644 are denoted by a black arrow. Wherever the most frequently observed amino acid is
645 neither the avian nor human consensus residue, it appears in gray. Position numbers for
646 the markers in each protein are given in the final row.

647 **Fig. 2. Persistent host markers occur in known protein binding domains.** Blue
648 squares denote regions where the named protein is known to bind to a specific protein
649 (Supplemental Information, text) or the novel SH3 domain. Red lines denote host
650 markers found in this study.

651 **Fig. 3. The preservation of host markers increases over time in human H1N1**
652 **viruses.** The panels above plot the proportion of host markers acquired over time by
653 H1N1 influenza isolated from human hosts. All 32 markers are 99% persistent.
654 Position numbered in red are referred to in the text.

Fig. 1

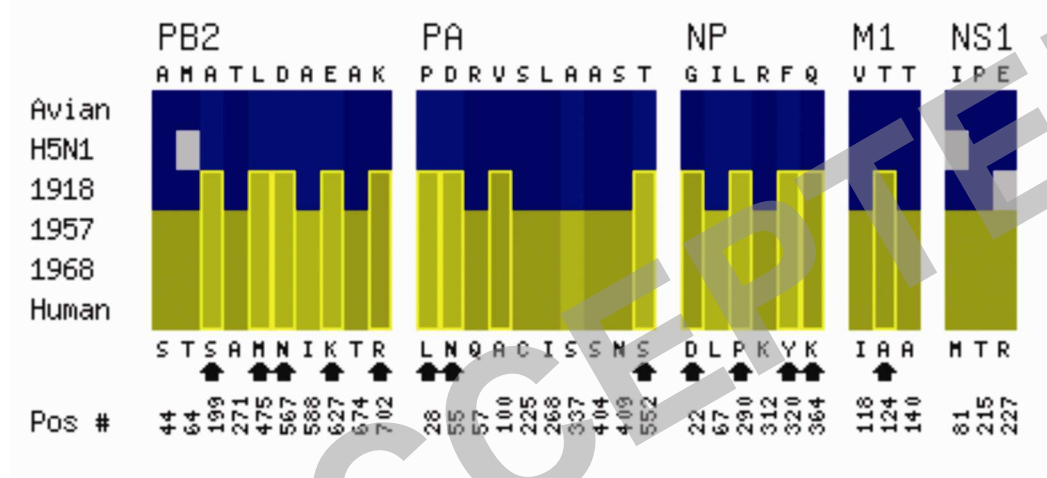


Fig. 2

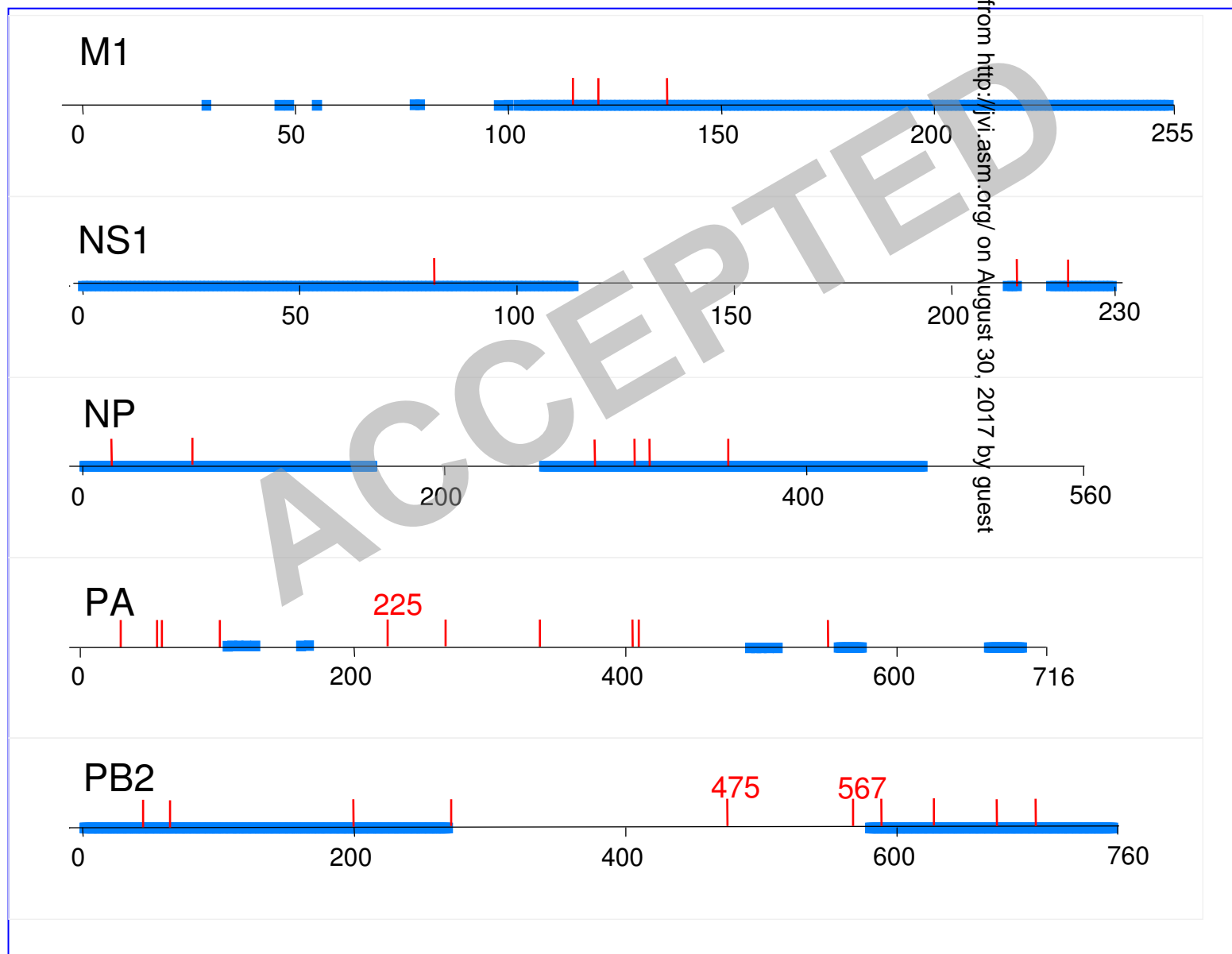


Fig. 3

