

1

## 2 **Supplementary Information for**

### 3 **Deep mutational scanning of hemagglutinin helps predict evolutionary fates of human H3N2** 4 **influenza variants**

5 **JM Lee, J Huddleston, MB Doud, KA Hooper, NC Wu, T Bedford, JD Bloom**

6 **Trevor Bedford**  
7 **tbedford@fredhutch.org**  
8 **Jesse D. Bloom**  
9 **jbloom@fredhutch.org**

#### 10 **This PDF file includes:**

- 11     Supplementary text
- 12     Figs. S1 to S11
- 13     Captions for Databases S1 to S4
- 14     References for SI reference citations

#### 15 **Other supplementary materials for this manuscript include the following:**

- 16     Databases S1 to S4

## 17 Supporting Information Text

18 **Codon mutagenesis.** The codon-mutant libraries were generated in the Perth/2009 HA-G78D-T212I background using the  
19 PCR-based approach described in (1) with the primer melting-temperature modifications described in (2), using two rounds of  
20 mutagenesis. The script to design the mutagenesis primers is at <https://github.com/jbloombloom/CodonTilingPrimers>. We created  
21 three independent libraries, one for each biological replicate. The mutant variants were then cloned at high efficiency into the  
22 pICR2 (3) vector using digestion with BsmBI, ligation with T4 DNA ligase, and electroporation into ElectroMAX DH10B  
23 competent cells (Invitrogen 18290015). We obtained >6 million transformants for each replicate. We scraped the plates,  
24 expanded the cultures in liquid LB (Luria-Bertani Broth) + ampicillin at 37°C for 3 h with shaking, and then maxiprepmed.  
25 We randomly chose 31 clones to Sanger sequence to evaluate the mutation rate (see SI Appendix, Figure S2).

26 **Barcoded-subamplicon sequencing.** We generated the HA PCR amplicons for the three plasmid libraries, the three virus  
27 libraries, the wildtype plasmid control, and the wildtype virus control using KOD Hot Start Master Mix (EMD Millipore  
28 71842) using the PCR reaction mixture and cycling conditions described in (1) and the P09-HA-For and P09-HA-Rev primers.  
29 We prepared the sequencing libraries using a barcoded-subamplicon strategy (4) to increase the accuracy from deep sequencing.  
30 The exact details of this approach are described in (5) (also see [https://jbloomlab.github.io/dms\\_tools2/bcsubamp.html](https://jbloomlab.github.io/dms_tools2/bcsubamp.html)). The  
31 primers used to generate the subamplicons are in the SI Appendix (Dataset S2). We performed deep sequencing on a lane of  
32 an Illumina HiSeq 2500 using  $2 \times 250$  bp paired-end reads in rapid-run mode.

33 **Phylogenetic model comparison and fitting of a stringency parameter.** For the analysis in Table 1, we downloaded all full-length  
34 H3 HA sequences from the Influenza Virus Resource (6), and randomly subsampled two sequences per year. These sequences  
35 were aligned using MAFFT (7) and used to infer a phylogenetic tree using RAxML (8) with a GTRCAT model of nucleotide  
36 substitution. We then used `phydms` (9) (<https://github.com/jbloombloom/phydms>, version 2.2.2) to fit the substitution models listed  
37 in Table 1.

38 The amino-acid preferences were re-scaled by the stringency parameter using the approach described in (9). Note that the  
39 re-scaling simply puts the amino-acid preferences on a useful scale for visualization in logo plots, and has no effect on any  
40 of the quantitative conclusions relating the deep mutational scanning to natural evolution. The reason is that all of these  
41 conclusions use the effects of mutations as calculated using Equation 1, and the re-scaling simply acts as a constant multiplier  
42 on all mutational effects (e.g., a scale factor) when re-scaled preferences are converted to mutational effects.

43 The phylogenetic tree of HA subtypes in Figure 7A was generated as described in (10).

44 **Inference of human H3N2 phylogenetic tree and calculation of maximum mutation frequencies.** To generate the tree shown in  
45 the SI Appendix (Figure S8), we applied Nextstrain’s augur pipeline (11) (<https://github.com/nextstrain/augur>; commit 006896d) to  
46 publicly available H3N2 HA sequences from GISAID (12) (see the SI Appendix, Dataset S4), sampling six viruses per month over  
47 the time interval of January 1, 1968 to February 1, 2018. We aligned the resulting 2,189 HA sequences with MAFFT v7.310 (7)  
48 and constructed a maximum likelihood phylogeny from this alignment with RAxML 8.2.10 (8). Ancestral state reconstruction  
49 and branch length timing were performed with TreeTime (13). The phylogenetic tree is available as a JSON file on GitHub  
50 at [https://github.com/jbloombloom/Perth2009-DMS-Manuscript/blob/master/analysis\\_code/data/flu\\_h3n2\\_ha\\_1968\\_2018\\_6v\\_tree.json.gz](https://github.com/jbloombloom/Perth2009-DMS-Manuscript/blob/master/analysis_code/data/flu_h3n2_ha_1968_2018_6v_tree.json.gz).  
51 The tree was visualized using BALTIC (<https://github.com/blab/baltic>).

52 The frequency trajectory of each individual mutation on the phylogeny is estimated following Nextstrain’s augur pipeline  
53 and as first implemented in Nextflu (14). Herein, mutation frequency dynamics are modeled according to a Brownian motion  
54 diffusion process discretized to one-month intervals. The number of viruses sampled in each interval determines the denominator  
55 of the mutation frequency calculations. Relative to a simple Brownian motion, the expectation includes an “inertia” term  $\epsilon$   
56 that adds velocity to the diffusion and the variance includes a term  $x(1-x)$  to scale variance according to frequency following  
57 a Wright-Fisher population genetic process. This results in the following diffusion process

$$58 \quad x(t + dt) = \mathcal{N}(x(t) + \epsilon dx, dt \sigma^2 x(t) (1 - x(t))), \quad [1]$$

59 with ‘volatility’ parameter  $\sigma^2$ . The term  $dx$  is the increment in the previous timestep, so that  $dx = x(t) - x(t - dt)$ . We used  
60  $\epsilon = 0.7$  and  $\sigma^2 = 0.05$  to maximize fit to empirical trajectory behavior.

61 We also include an Bernoulli observation model for mutation presence / absence among sampled viruses at timestep  $t$ . This  
62 observation model follows

$$63 \quad f(x, t) = \prod_{v \in V} x(t) \prod_{v \notin V} (1 - x(t)), \quad [2]$$

64 where  $v \in V$  represents the set of viruses that have the mutation and  $v \notin V$  represents the set of viruses that do not have the  
65 mutation. Each frequency trajectory is estimated by simultaneously maximizing the likelihood of the process model and the  
66 likelihood of the observation model via adjusting frequency trajectory  $\mathbf{x} = (x_1, \dots, x_n)$ .

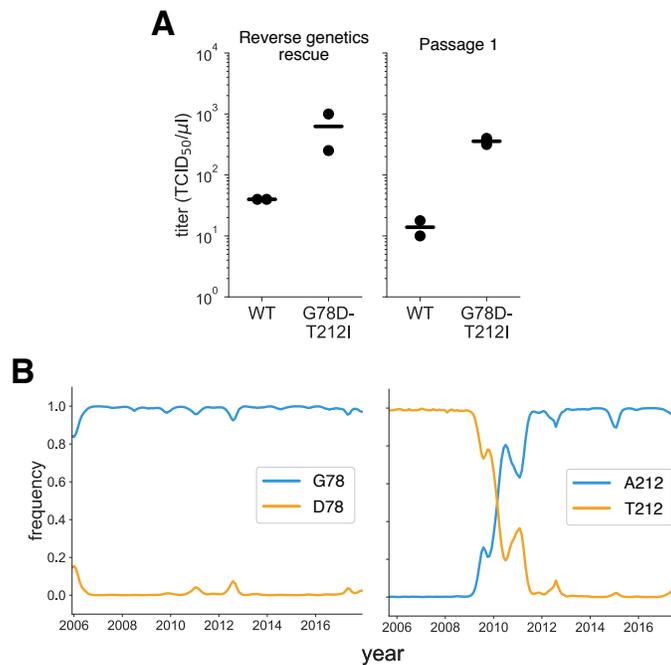
67 We also repeated the above analyses using only viruses that were sequenced directly without passaging. Routine direct  
68 sequencing did not begin until the early 2000s (15). To construct a tree with a similar number of viruses as the original  
69 analysis, we sampled 30 viruses per month between January 1, 2000 and April 1, 2018, producing a tree with 2,374 unpassaged  
70 viruses with augur (commit: 6d9f708). We included the passaged DMS strain, A/Perth/16/2009, in the resulting tree to  
71 enable comparison between pre-Perth and post-Perth clades.

72 **Analysis of mutational shifts.** To compare the Perth/2009 H3 and WSN/1933 H1 HA preferences, we first aligned the wildtype  
73 HA sequences using MAFFT (7). To quantify the shifts in preference for every alignable site while accounting for experimental  
74 noise, we used the approach described in (16) and used the  $\text{RMSD}_{\text{corrected}}$  values as our quantification of the extent of each  
75 shift.

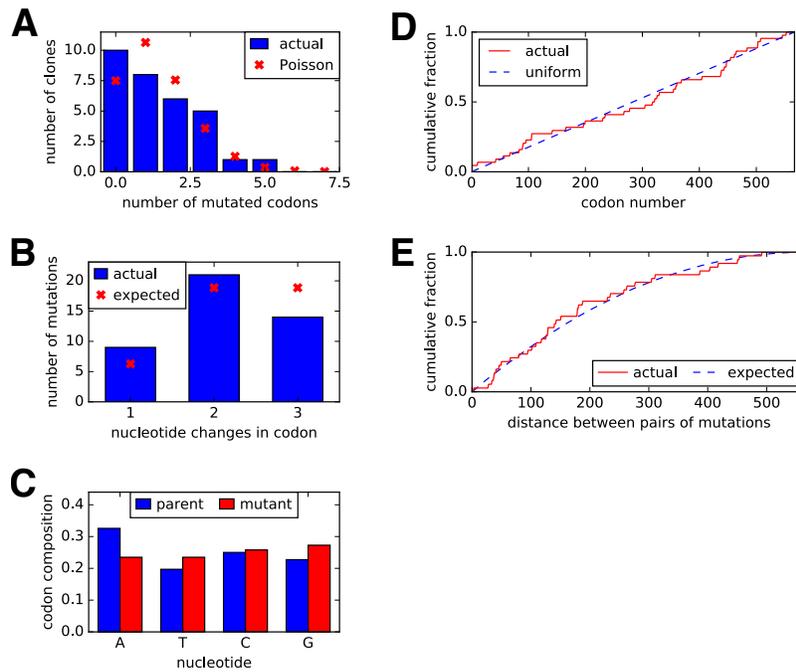
76 For the plots shown in Figure 8B, any residues falling between Cys-52 and Cys-277 were defined as the head domain, and all  
77 other residues were defined as the stalk domain. We used the multiple sequence alignment of the HA subtype sequences from  
78 (10) to identify sites that are absolutely conserved across all subtypes, or in the different clades described in Figure 8.

79 **Validation of individual point mutants.** To validate the viral growth of Perth/2009 HA point mutants M(-16)K, C52A, C52C,  
80 T24F, T40V, S287A, and C199(HA2)K, we used site-directed mutagenesis to introduce the amino-acid mutation into the  
81 Perth/2009 HA bidirectional reverse genetics plasmid, and verified the sequence of two clones for each mutant by Sanger  
82 sequencing. We generated these individual mutant viruses carrying GFP in the PB1 segment using a protocol described  
83 in (17, 18), and the PB2, PA, NP, NA, M, and NS segments of the A/WSN/1933 (H1N1) strain. To rescue each mutant  
84 GFP-carrying virus in duplicate, we transfected a co-culture of  $4 \times 10^5$  293T-CMV-PB1 and  $0.5 \times 10^5$  MDCK-SIAT1-CMV-  
85 PB1-TMPRSS2 cells with the eight reverse genetics plasmids and the pHAGE2-EF1aInt-TMPRSS2-IRES-mCherry-W plasmid.  
86 Each well received a transfection mixture of 100  $\mu\text{L}$  DMEM, 3  $\mu\text{L}$  BioT transfection reagent, and 250 ng of each plasmid. We  
87 changed the media in each well with 2 mL IGM eight hours post-transfection. At 53 hours post-transfection, transfection  
88 supernatants were harvested, clarified by centrifugation at  $2000 \times g$  for five minutes, aliquoted, and frozen at  $-80^\circ\text{C}$ .

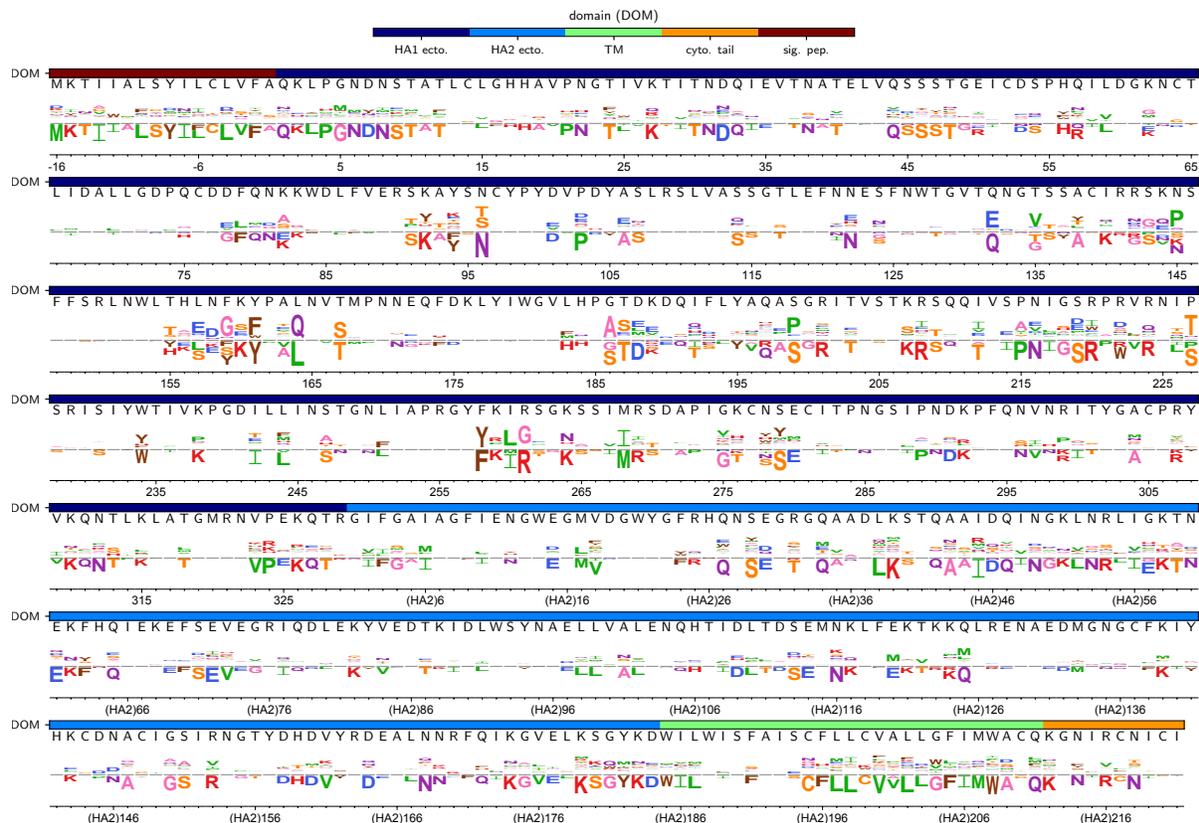
89 To titer the GFP-carrying viruses, we plated  $1 \times 10^5$  MDCK-SIAT1-CMV-PB1-TMPRSS2 cells per well in 12-well plates in  
90 IGM. Four hours after plating, we infected cells with dilutions of viral supernatant. At 16 hours post-infection, we selected  
91 wells with 1 to 10% of cells that were GFP-positive and analyzed the fraction of GFP-positive cells by flow cytometry to  
92 calculate the titer of infectious particles per  $\mu\text{L}$ .



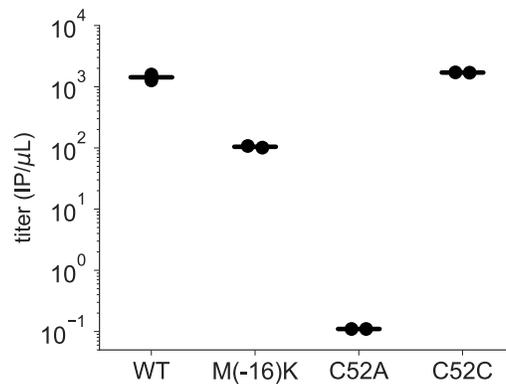
**Fig. S1. Characterization of the G78D-T212I Perth/2009 HA variant.** (A) The G78D-T212I Perth/2009 HA variant supports better viral growth than the wildtype Perth/2009 HA. Viruses were generated in duplicate by reverse genetics with the Perth/2009 NA and WSN internal genes, and passaged once at MOI = 0.01 in MDCK-SIAT1-TMPRSS2 cells. The rescue and passage viral supernatants were collected at 72 hours post-transfection and 44 hours post-infection, respectively, and titered in MDCK-SIAT1-TMPRSS2 cells. The points mark each duplicate and the bar marks the mean. (B) The D78 variant remained at a low frequency in natural human H3N2 sequences over the past ~10 years. The A212 variant rose to fixation in ~2011, replacing the T212 variant.



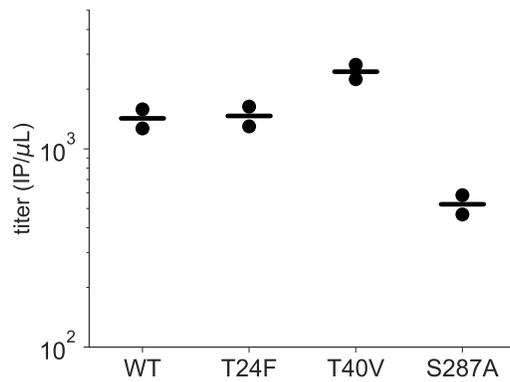
**Fig. S2. Sanger sequencing of 31 randomly chosen clones from the mutant plasmid libraries.** (A) There were an average of  $\sim 1.4$  codon mutations per clone across the three plasmid mutant libraries. (B) A mixture of one-, two-, and three-nucleotide mutations were present, with slightly fewer triple-nucleotide changes than expected. (C) Nucleotide frequencies were uniform in the codon mutations. (D) The mutations were distributed relatively evenly across the length of the HA coding sequence. (E) We calculated the pairwise distances between mutations for clones carrying more than one mutation. The cumulative distribution of these distances is shown in the red line. The blue line indicates the expected distribution if mutations in multiply mutated genes are randomly dispersed along the sequence.



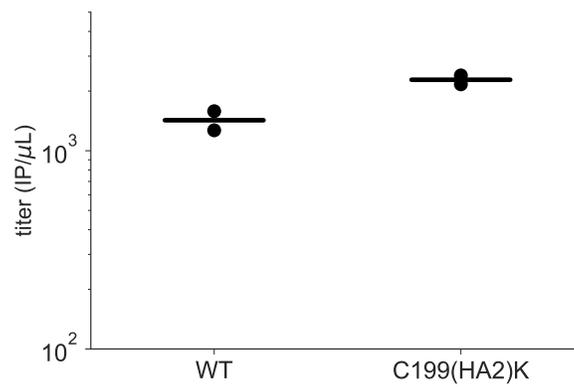
**Fig. S3.** Sites where there are strong differences between our experimental measurements and the amino-acid frequencies among natural HA sequences. We calculated the distance between our H3 measurements and the alignment frequencies using the same approach as in Figure 7, but using the alignment frequencies in place of the H1 preferences. For each site, the height of each letter above or below the line indicates how much more or less preferred that amino acid is in our experiments compared to its frequency nature. The overlays show the same information as in Figure 2 (domain and wildtype amino acid). The sites are in H3 numbering. The HA alignment used to calculate the natural frequencies is the same one used in Table 1.



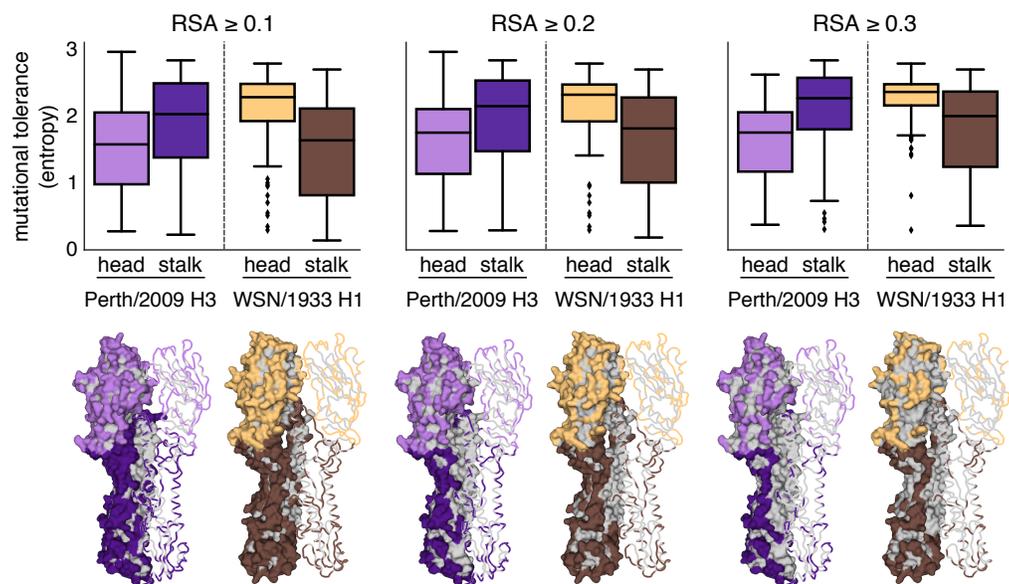
**Fig. S4. Validation that the Perth/2009 HA is somewhat tolerant to mutation of the canonical start codon.** We created variants of the Perth/2009 HA in which the canonical start codon was mutated (amino-acid mutation M(-16)K, codon mutation ATG→AAA), a conserved cysteine was mutated to alanine (C52A, codon mutation TGC→GCA), and the same cysteine was synonymously mutated (C52C, codon mutation TGC→TGT). We selected these mutations because our deep mutational scanning results in Figure 2 surprisingly suggest that the start codon (position -16) is fairly tolerant of mutations, but that site 52 is highly intolerant of mutation to anything other than cysteine. The synonymous C52C mutation is a negative control mutation that is not expected to have any effect. We then used reverse genetics to generate viruses carrying the wildtype HA or each point mutant, with GFP packaged in the PB1 segment to enable easy titring by flow cytometry (17, 18). Each variant was generated in duplicate, and the plots show the viral titer in the supernatant at 53 hours post-transfection. As expected, the C52A virus is essentially non-viable, while the C52C mutation is roughly equivalent to wildtype. The M(-16)K mutation is only moderately attenuated, validating that the canonical start codon is not completely essential for viral growth.



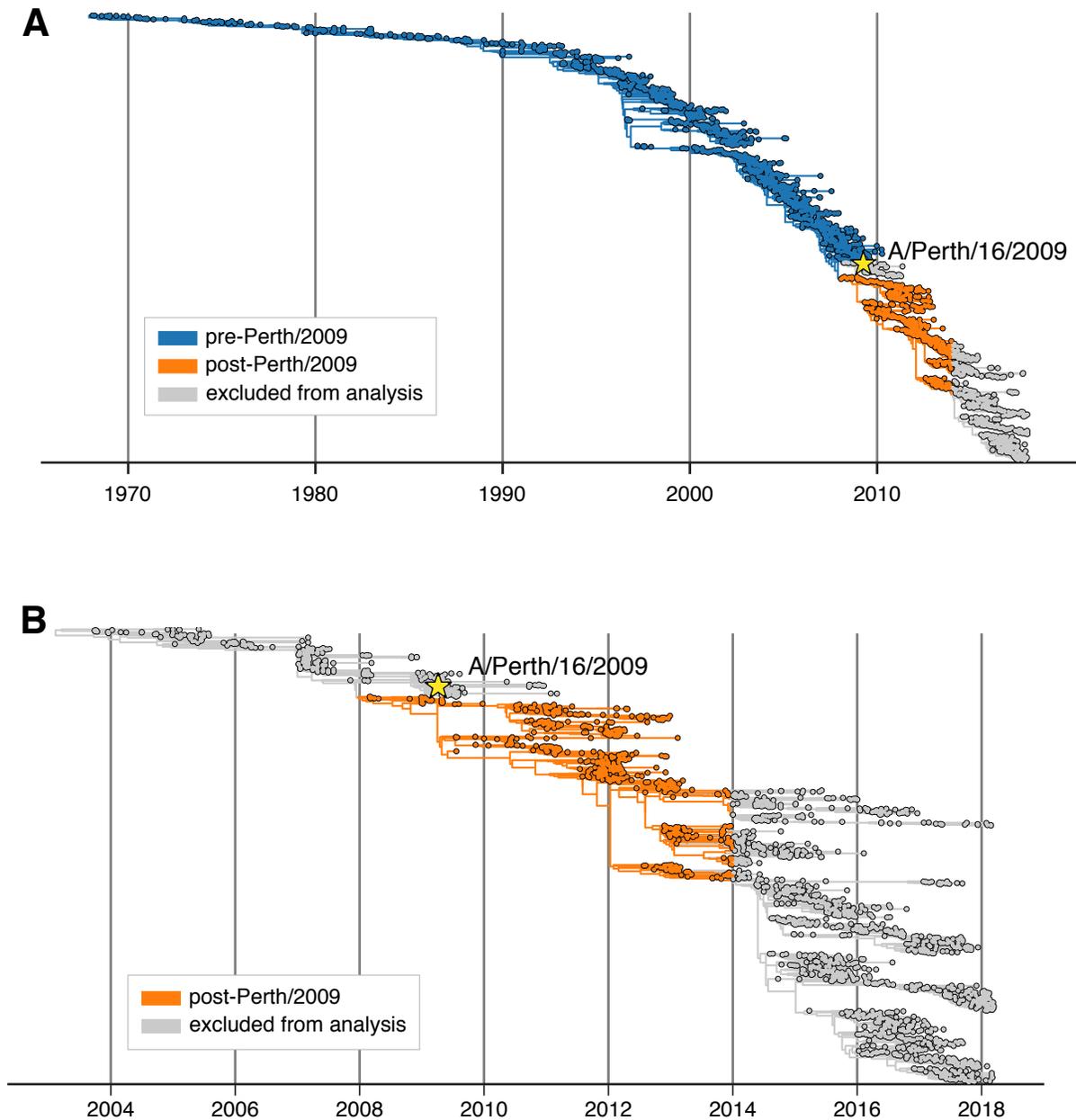
**Fig. S5. Validation of viral variants with mutations at N-linked glycosylation motifs.** We created variants of the Perth/2009 HA with mutations that eliminated N-linked glycosylation motifs (Asn-Xaa-Ser/Thr) at asparagine residues 22, 38, and 285 (these are mutations T24F, T40V, and S287A, respectively). The codon mutations were ACG→TTC, ACT→GTA, and AGC→GCG, respectively. We then used reverse genetics to generate viruses carrying the wildtype HA or each of these mutants. Each variant was generated in duplicate, and the plots show the viral titer in the supernatant at 53 hours post-transfection. The viruses with mutations at the motifs at residues 22 and 38 reached titers at least as high as wildtype, whereas the virus with a mutation to the motif at residue 285 was modestly attenuated.



**Fig. S6. Validation of the mutational tolerance of a site in the transmembrane domain.** We created a variant of the Perth/2009 HA with a transmembrane domain mutation, C199(HA2)K (codon mutation TGT→AAG), at a site that our deep mutational scanning suggests should be highly mutationally tolerant (Figure 2). We then used reverse genetics to generate viruses carrying the wildtype HA or this mutant. Each variant was generated in duplicate, and the plots show the viral titer in the supernatant at 53 hours post-transfection. The virus with the mutation reached titers comparable to wildtype.



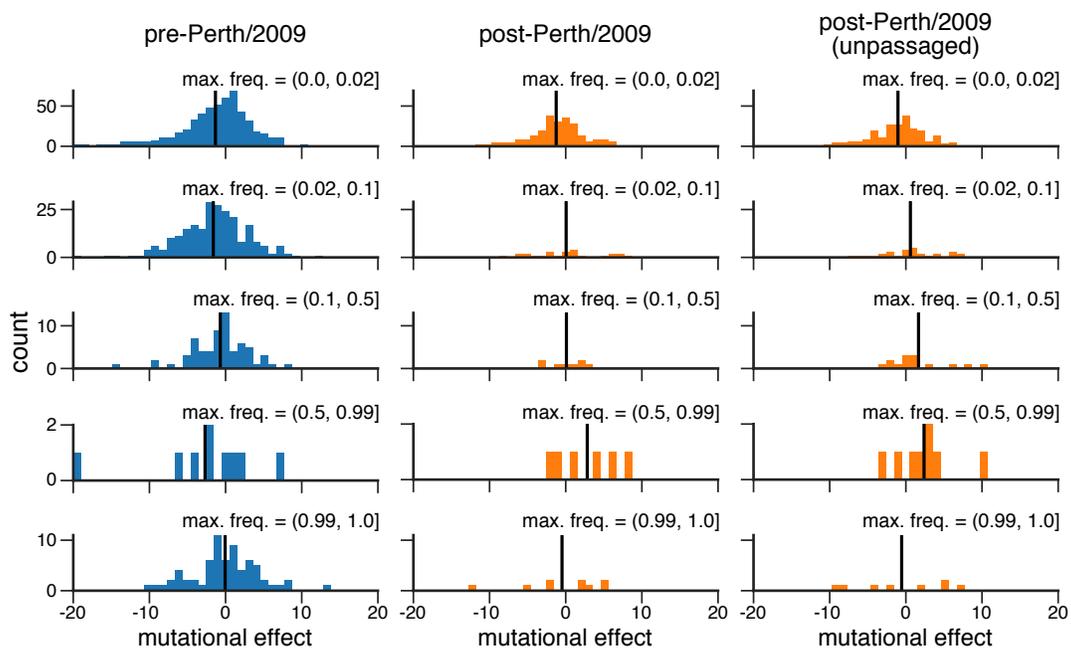
**Fig. S7. Mutational tolerances of the head and stalk domains at various relative solvent accessibility cutoffs.** The mutational tolerances of the head and stalk domains show less disparity for the Perth/2009 H3 HA compared to those for the WSN/1933 H1 HA. We used relative solvent accessibility (RSA) cutoffs of 0.1, 0.2, and 0.3 to define solvent-exposed residues and plotted the mutational tolerances (Shannon entropy of re-scaled preferences) of these residues in the head and stalk domains for the Perth/2009 H3 HA (purple) and the WSN/1933 H1 HA (brown). Residues falling in between the two cysteines at sites 52 and 277 were defined as belonging to the head domain, while all other residues were defined as the stalk domain. The HA structures color the residues that are defined as solvent exposed at a given RSA cutoff. One monomer is shown in surface representation and another monomer shown in ribbon representation. Residues in lighter shades of purple or brown are in the head domain, while residues in darker shades are in the stalk domain. Note that the mutational tolerance values are not comparable between the two HAs.



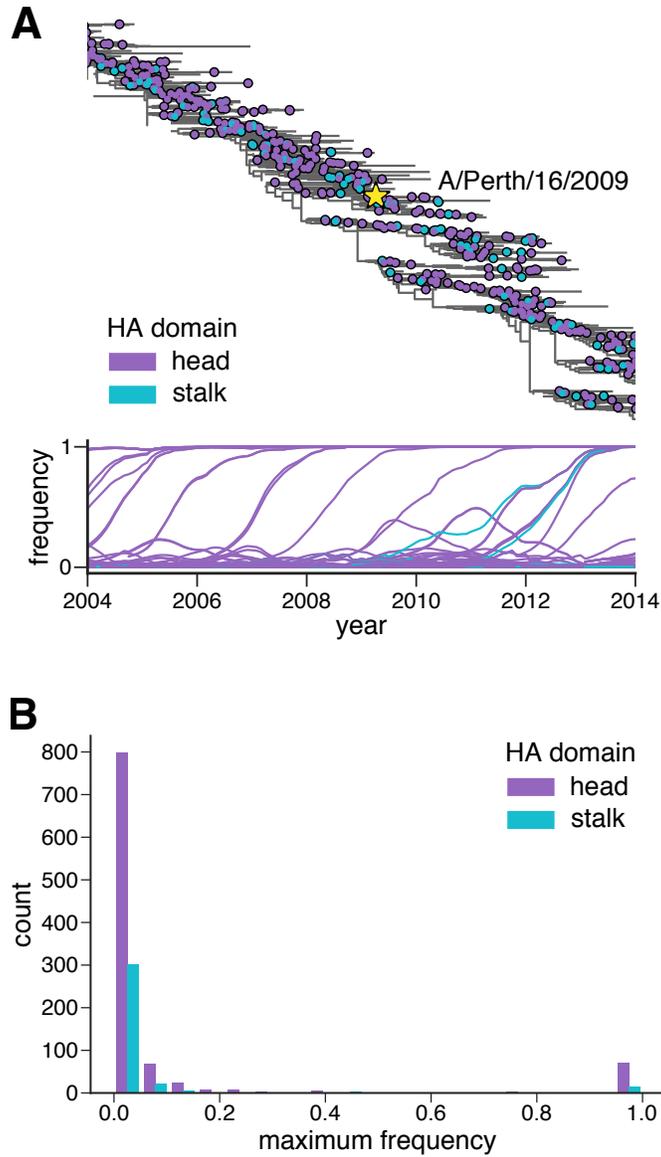
**Fig. S8. A phylogenetic tree of all HA sequences used in our analysis of mutation frequencies.** (A) HA sequences were sampled at a rate of six viruses per month from January 1, 1968 through February 1, 2018. The Perth/2009 strain used in our experiments is indicated. The rest of the tree is partitioned into nodes that preceded the split of the Perth/2009 strain from the trunk of the tree (blue) and nodes that branched off the trunk after the clade containing Perth/2009 (orange). In Figure 5, these two partitions of the tree are analyzed separately. Nodes in the clade containing the Perth/2009 strain and nodes sampled in 2014 or after were excluded from our analyses. The Perth/2009 strain was excluded to avoid artifacts related to mutations that occurred on the branches leading to the HA sequence used in the experiment. The post-2014 nodes were excluded because the evolutionary fates of many sequences after this date are not yet fully resolved. (B) The post-Perth/2009 partition of the tree containing only sequences from unpassed isolates.



**Fig. S9. The site-specific amino-acid preferences of the WSN/1933 H1 HA.** The amino-acid preferences of the WSN/1933 H1 HA from (5) after taking the average of the experimental replicates and re-scaling (9) by a stringency parameter of 2.05 (see [https://github.com/jbloomlab/dms\\_tools2/blob/master/examples/Doud2016/analysis\\_notebook.ipynb](https://github.com/jbloomlab/dms_tools2/blob/master/examples/Doud2016/analysis_notebook.ipynb)). The overlays show the same information as in Figure 2 (domain and wildtype amino acid). Note that (5) used libraries in which all codons were mutagenized *except* for the one encoding N-terminal methionine. Therefore, no data is shown for the first codon in the gene. The sites are in H3 numbering.



**Fig. S10.** The distribution of mutational effects measured in H1 HA among H3N2 mutations binned by the maximum frequency that they reach. This figure repeats the analysis of the H3N2 mutation frequencies in Figure 5B, but uses the deep mutational scanning data for an H1 HA as measured in (5).



**Fig. S11. Frequency trajectories of head and stalk domain mutations.** (A) This figure repeats the analysis of the H3N2 mutation frequencies in Figure 4, but colors amino-acid mutations by whether they occur in the head (purple) or the stalk (blue) domain. (B) Histogram of mutation maximum frequencies by the number of mutations in the head and stalk domains. It is clear that mutations in the head domain are more numerous than those in the stalk, particularly among mutations that reach high frequencies.

93 **Additional data table S1 (pHW\_Perth09-HA-G78D-T212I.txt)**

94 Genbank file giving the full sequence of the bidirectional reverse-genetics plasmid pHW-Perth09-HA-G78D-T212I, which  
95 encodes the wildtype HA sequence used in this study.

96 **Additional data table S2 (Perth2009\_subamplicon\_primers.xlsx)**

97 Excel file providing the primers used to generate the barcoded subamplicons for Perth/2009 HA deep sequencing.

98 **Additional data table S3 (DatasetS3.xlsx)**

99 Excel file giving the amino-acid preferences in sequential 1, 2, ... numbering of the Perth/2009 HA. The unscaled preferences  
100 for replicates 1, 2, 3-1, and 3-2 are each in a separate tab of the file. Additional tabs give the across-replicate averaged and  
101 re-scaled amino-acid preferences in sequential numbering and in H3 numbering as shown in Figure 2. There are also tabs  
102 that give the conversion from sequential to H3 numbering, and a list of the mutations in the 31 clones we Sanger sequenced  
103 to evaluate the mutation rate of the mutant plasmid libraries. Each tab can simply be exported to CSV for computational  
104 analyses.

105 **Additional data table S4 (gisaid\_acknowledgment\_table.xls)**

106 Excel file providing acknowledgments and accessions for sequences downloaded from GISAID.

107 **References**

- 108 1. Bloom JD (2014) An experimentally determined evolutionary model dramatically improves phylogenetic fit. *Molecular*  
109 *Biology and Evolution* 31:1956–1978.
- 110 2. Dingens AS, Haddock HK, Overbaugh J, Bloom JD (2017) Comprehensive mapping of HIV-1 escape from a broadly  
111 neutralizing antibody. *Cell Host & Microbe* 21:777–787.
- 112 3. Ashenberg O, Padmakumar J, Doud MB, Bloom JD (2017) Deep mutational scanning identifies sites in influenza  
113 nucleoprotein that affect viral inhibition by MxA. *PLoS Pathogens* 13(3):e1006288.
- 114 4. Wu NC, et al. (2014) High-throughput profiling of influenza A virus hemagglutinin gene at single-nucleotide resolution.  
115 *Scientific Reports* 4:4942.
- 116 5. Doud MB, Bloom JD (2016) Accurate measurement of the effects of all amino-acid mutations to influenza hemagglutinin.  
117 *Viruses* 8:155.
- 118 6. Bao Y, et al. (2008) The Influenza Virus Resource at the National Center for Biotechnology Information. *J. Virol.*  
119 82:596–601.
- 120 7. Katoh K, Standley DM (2013) MAFFT multiple sequence alignment software version 7: improvements in performance  
121 and usability. *Molecular Biology and Evolution* 30(4):772–780.
- 122 8. Stamatakis A (2006) RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and  
123 mixed models. *Bioinformatics* 22(21):2688–2690.
- 124 9. Hilton SK, Doud MB, Bloom JD (2017) phydms: Software for phylogenetic analyses informed by deep mutational scanning.  
125 *PeerJ* 5:e3657.
- 126 10. Doud MB, Lee JM, Bloom JD (2018) How single mutations affect viral escape from broad and narrow antibodies to H1  
127 influenza hemagglutinin. *Nature Communications* DOI 10.1038/s41467-018-03665-3.
- 128 11. Hadfield J, et al. (2017) Nextstrain: real-time tracking of pathogen evolution. *bioRxiv*.
- 129 12. Shu Y, McCauley J (2017) GISAID: Global initiative on sharing all influenza data—from vision to reality. *Eurosurveillance*  
130 22(13).
- 131 13. Sagulenko P, Puller V, Neher RA (2018) TreeTime: Maximum-likelihood phylodynamic analysis. *Virus Evolution*  
132 4(1):vex042.
- 133 14. Neher RA, Bedford T (2015) nextflu: real-time tracking of seasonal influenza virus evolution in humans. *Bioinformatics*  
134 31(21):3546–3548.
- 135 15. McWhite C, Meyer A, Wilke C (2016) Sequence amplification via cell passaging creates spurious signals of positive  
136 adaptation in influenza virus H3N2 hemagglutinin. *Virus Evolution* 2:vew026.
- 137 16. Haddock HK, Dingens AS, Hilton SK, Overbaugh J, Bloom JD (2018) Mapping mutational effects along the evolutionary  
138 landscape of HIV envelope. *eLife* 7:e34420.
- 139 17. Bloom JD, Gong LI, Baltimore D (2010) Permissive secondary mutations enable the evolution of influenza oseltamivir  
140 resistance. *Science* 328:1272–1275.
- 141 18. Hooper KA, Bloom JD (2013) A mutant influenza virus that uses an N1 neuraminidase as the receptor-binding protein.  
142 *Journal of Virology* 87(23):12531–12540.