

# Deep mutational scanning of hemagglutinin helps distinguish the evolutionary fates of human H3N2 influenza virus lineages

Juhye M. Lee<sup>a,d,e,1</sup>, John Huddleston<sup>b,f,1</sup>, Michael B. Doud<sup>a,d,e</sup>, Kathryn A. Hooper<sup>a,f</sup>, Trevor Bedford<sup>b,c,2</sup>, and Jesse D. Bloom<sup>a,c,d,2</sup>

<sup>a</sup>Basic Sciences Division; <sup>b</sup>Vaccine and Infectious Diseases Division; <sup>c</sup>and Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA; <sup>d</sup>Department of Genome Sciences; <sup>e</sup>Medical Scientist Training Program; <sup>f</sup>and Molecular and Cellular Biology Program, University of Washington, Seattle, WA 98195, USA

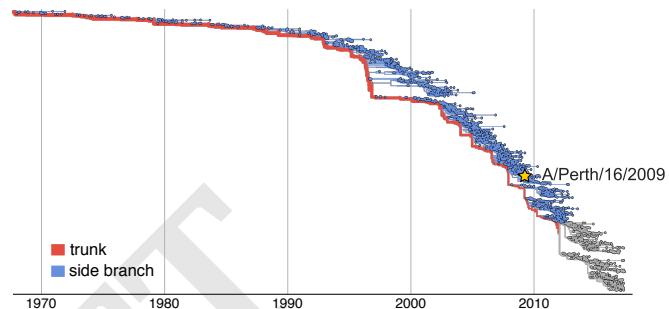
This manuscript was compiled on January 18, 2018

Human influenza virus rapidly accumulates mutations in its major surface protein hemagglutinin (HA). The evolutionary success of influenza virus lineages depends on how these mutations affect HA's functionality and antigenicity. Here we experimentally measure the effects on viral growth in cell culture of all single amino-acid mutations to the HA from a recent human H3N2 influenza virus strain. We then show that mutations in evolutionarily successful H3N2 viral lineages are measured to generally be more favorable than mutations in lineages that die out. Therefore, despite the well-known caveats about cell-culture measurements of viral fitness, such measurements can still be informative for understanding virus evolution in nature. We also compare our measurements for an H3 HA to similar data previously generated for a distantly related H1 HA, and find substantial differences in which amino acids are preferred at many sites. For instance, the H3 HA has less disparity in mutational tolerance between the head and stalk domains than the H1 HA. Overall, our work suggests that experimental measurements of mutational effects can be leveraged to help forecast the evolutionary fates of viral lineages in nature — but only when the measurements are made on a viral strain similar to the ones being studied in nature.

influenza virus | hemagglutinin | deep mutational scanning | antigenic drift | epistasis

Seasonal H3N2 influenza virus evolves rapidly, fixing 3 to 4 amino-acid mutations per year in its hemagglutinin (HA) surface protein. Many of these mutations contribute to the rapid antigenic drift that necessitates frequent updates to the annual influenza vaccine (1, 2). This evolution is further characterized by strain competition and frequent population turnover (3–8), producing a spindly phylogenetic tree with a persistent trunk lineage and short-lived side branches (9) (Figure 1). Several lines of evidence indicate that the trunk lineage outcompetes other strains at least in part because it has higher fitness (4–7). A key challenge in the study of H3N2 evolution is to forecast which strain will dominate the upcoming influenza season, as this information can guide selection of a vaccine strain (10). In order to make these forecasts, it is important to identify the features that enable the trunk strain to persist as other strains die out.

Two main characteristics distinguish the evolutionarily successful trunk from its competitors: greater antigenic change, and efficient viral growth and transmission. In principle, experiments could be informative for identifying how mutations affect these features. Most work on influenza evolution to date has utilized experimental data primarily to assess the antigenicity of circulating strains (11–13). While the anti-



**Fig. 1.** Human H3N2 HA phylogeny from 1968–2017. The trunk is shown in red, and side branches are shown in blue. The gray branches represent the part of the tree for which we cannot yet distinguish the trunk from side branches. The Perth/2009 strain is indicated with a star.

genic features of a virus certainly contribute to its success, the non-antigenic effects of mutations also play an important role (4, 7, 14–16). Specifically, due to influenza virus's high mutation rate (17–19) and lack of intra-segment recombination (20), deleterious mutations become linked to beneficial ones. The resulting accumulation of deleterious mutations can affect non-antigenic properties that are central to viral fitness (16). However, while extensive work has gone into characterizing how mutations affect HA antigenicity (21–26), we lack any large-scale quantitative characterization of how mutations to H3N2 HA affect non-antigenic properties crucial to viral fitness.

It is now possible to use deep mutational scanning (27)

## Significance Statement

A key goal in the study of influenza virus evolution is to forecast which viral strains will persist and which ones will die out. Here we experimentally measure the effects of all amino-acid mutations to the hemagglutinin protein from a human H3N2 influenza strain on viral growth in cell culture. We show that these measurements have utility for distinguishing among viral strains that do and do not succeed in nature. Overall, our work suggest that new high-throughput experimental approaches may be useful for evolutionary forecasting.

Please provide details of author contributions here.

Please declare any conflict of interest here.

<sup>1</sup>J.M.L. and J.H. contributed equally to this work.

<sup>2</sup>To whom correspondence should be addressed: trevor@bedford.io, jbloom@fredhutch.org

to measure the functional effects of large numbers of mutations to viral proteins (28–32). However, the only HA for which such measurements have previously been made is from the highly lab-adapted A/WSN/1933 (H1N1) strain (28–30). Here we measure the effects on viral growth in cell culture of all mutations to the HA from a recent human H3N2 strain, A/Perth/16/2009. We then examine whether these measurements are useful for distinguishing between H3N2 trunk and side branch lineages, and find that mutations on the trunk are measured to be more favorable than those on the side branches, both at epitope and non-epitope sites. However, the experiments are useful for assessing the effects of mutations only in closely related viral strains, since we show that mutations often have disparate effects on H3 and H1 HAs. Our work highlights the potential for using high-throughput experimental measurements of mutational effects to inform evolutionary forecasting of human seasonal influenza virus.

## Results

**Deep mutational scanning of HA from a recent strain of human H3N2 influenza virus.** We performed a deep mutational scan to measure the effects of all amino-acid mutations to HA from the A/Perth/16/2009 (H3N2) strain on viral growth in cell culture. This strain was the H3N2 component of the influenza vaccine from 2010–2012 (33, 34). Relative to the consensus sequence for this HA in Genbank, we used a variant with two mutations that enhanced viral growth in cell culture, G78D and T212I (Figure S1 and Dataset S1). The G78D mutation occurs at low frequency in natural H3N2 sequences, and T212 is a site where a mutation to Ala rose to fixation in human influenza in ~2011.

We mutagenized the entire HA coding sequence at the codon level to create mutant plasmid libraries harboring an average of ~1.4 codon mutations per clone (Figure S2). We then generated mutant virus libraries from the mutant plasmids using a helper-virus system that enables efficient generation of complex influenza virus libraries (30) (Figure 2A). These mutant viruses derived all their non-HA genes from the lab-adapted A/WSN/1933 strain. Using WSN/1933 for the non-HA genes reduces biosafety concerns, and also helped increase viral titers. To further increase viral titers, we used MDCK-SIAT1 cells that we had engineered to constitutively express the TMPRSS2 protease, which cleaves the HA precursor to activate it for membrane fusion (35, 36).

After generating the mutant virus libraries, we passaged them at low MOI in cell culture to create a genotype-phenotype link and select for functional HA variants (Figure 2A). All experiments were completed in full biological triplicate (Figure 2B). We also passaged and deep sequenced library 3 in duplicate (denoted as library 3-1 and 3-2) to gauge the experimental noise occurring *within* a single biological replicate. As a control to measure sequencing and mutational errors, we used the unmutated HA gene to generate and passage viruses carrying wildtype HA.

Deep sequencing of the initial plasmid mutant libraries and the passaged mutant viruses revealed selection for functional HA mutants. Specifically, stop codons were purged to 20–45% of their initial frequencies after correcting for error rates estimated by sequencing the wildtype controls (Figure 2C). The incomplete purging of stop codons is likely because genetic complementation due to co-infection (37) enabled the

persistence of some virions with nonfunctional HAs. We also observed selection against many nonsynonymous mutations (Figure 2C), with their frequencies falling to 30–40% of their initial values after error correction.

We next quantified the reproducibility of our deep mutational scanning measurements across biological and technical replicates. We first used the deep sequencing data for each replicate estimate the preference of each site in HA for all 20 amino acids using the method described in (38). Because there are 567 residues in HA, there are  $567 \times 20 = 11,340$  estimated amino-acid preferences, corresponding to  $567 \times 19 = 10,773$  distinct measurements (the 20 preferences at each site are constrained to sum to one (38)). The correlations of the amino-acid preferences between pairs of replicates are shown in Figure 2D. The biological replicates are fairly well-correlated, with Pearson's *R* ranging from 0.69 to 0.78. Replicate 1 exhibited the weakest correlation with the other replicates; this replicate also showed the weakest selection against stop and nonsynonymous mutations (Figure 2C), perhaps indicating more experimental noise. The two technical replicates 3-1 and 3-2 were only slightly more correlated than pairs of biological replicates, suggesting that bottlenecking of library diversity after the reverse-genetics step contributes most of the experimental noise.

**Our measurements are consistent with existing knowledge about HA's evolution and function.** How do the HA amino-acid preferences measured in our experiments relate to the evolution of H3N2 influenza virus in nature? This question can be addressed by evaluating how well an experimentally informed codon substitution model (ExpCM) using our measurements describe H3N2 evolution compared to standard phylogenetic substitution models (39, 40). Table 1 shows that the ExpCM using the across-replicate average of our measurements greatly outperforms conventional substitution models. This result indicates that our experiments authentically capture some of the constraints on HA evolution. The relative rate of nonsynonymous to synonymous substitutions ( $dN/dS$  or  $\omega$ ) is  $\ll 1$  for conventional substitution models (Table 1). However, the relative rate of nonsynonymous to synonymous substitutions after accounting for the amino-acid constraints measured in our experiments (the  $\omega$  for the ExpCM) is close to one (Table 1), indicating that the deep mutational scanning captures much of the purifying selection on HA. ExpCM fit a stringency parameter that relates the selection in the experiments to that in nature (39, 40). The stringency parameter for our HA ExpCM is 2.44 (Table 1), indicating that natural selection favors the same amino acids as the experiments, but with greater stringency. Throughout the rest of this paper, we use experimental measurements re-scaled (39, 40) by this stringency parameter. These re-scaled preferences are shown in Figure 3.

A closer examination of Figure 3 reveals that the experimentally measured amino-acid preferences generally agree with existing knowledge about HA's structure and function. For instance, sites that form structurally important disulfide bridges (sites 52 & 277, 64 & 76, 97 & 139, 281 & 305, 14 & 137-HA2, 144-HA2 & 148-HA2) (43) possess high preference for cysteine. At residues involved in receptor binding, there are strong preferences for the amino acids that are known to be involved in binding sialic acid, such Y98, D190, W153, and S228 (44–47). A positively charged amino acid at site 329 is



**Fig. 2. Deep mutational scanning of the Perth/2009 H3 HA.** (A) We generated mutant virus libraries using a helper-virus approach (30), and passed the libraries at low MOI to establish a genotype-phenotype linkage and to select for functional HA variants. Deep sequencing of the variants before and after selection allowed us to estimate each site's amino-acid preferences. (B) The experiments were performed in full biological triplicate. We also passedaged and deep sequenced library 3 in duplicate. (C) Frequencies of nonsynonymous, stop, and synonymous mutations in the mutant plasmid DNA, the passedaged mutant viruses, and wildtype DNA and virus controls. (D) The Pearson correlations among the amino-acid preferences estimated in each replicate.

**Table 1. Substitution models informed by the experiments describe H3N2's natural evolution better than traditional substitution models.**

| Model            | $\Delta\text{AIC}$ | $\text{LnL}$ | Stringency | $\omega$          |
|------------------|--------------------|--------------|------------|-------------------|
| ExpCM            | 0.0                | -8439        | 2.44       | 0.91              |
| GY94 M5          | 2166               | -9516        | –          | 0.36 (0.30, 0.84) |
| ExpCM, site avg. | 2504               | -9691        | 0.68       | 0.32              |
| GY94 M0          | 2608               | -9738        | –          | 0.31              |

Maximum likelihood phylogenetic fit to an alignment of human H3N2 influenza HAs using ExpCM (40), ExpCM in which the experimental measurements are averaged across sites (site avg.), and the M0 and M5 versions of the Goldman-Yang (GY94) model (41). Models are compared by AIC (42) computed from the log likelihood ( $\text{LnL}$ ) and number of model parameters. The  $\omega$  parameter is dN/dS for the Goldman-Yang models, and the relative dN/dS after accounting for the measurements for the ExpCM. For the M5 model, we give the mean followed by the shape and rate parameters of the gamma distribution over  $\omega$ .

important for cleavage of the HA0 precursor into the mature form (48, 49), and this site strongly prefers arginine. However, a notable exception occurs at the start codon at position -16, which does not show a strong preference for methionine. This codon is part of the signal peptide and is cleaved from the mature HA protein. One possible reason that our experiments do not show a strong preference for methionine at this site could be alternative translation-initiation at a downstream or upstream start site, as has been described for other HAs (50).

**There is less difference in mutational tolerance between the HA head and stalk domains for H3 than for H1.** Our experiments measure which amino acids are tolerated at each HA site under selection for protein function. We can therefore use our experimentally measured amino-acid preferences to calculate the inherent mutational tolerance of each site, which we quantify as the Shannon entropy of the re-scaled preferences. In prior mutational studies of the WSN/1933 H1 HA, the stalk domain was found to be substantially less mutationally tolerant than the globular head (28–30).

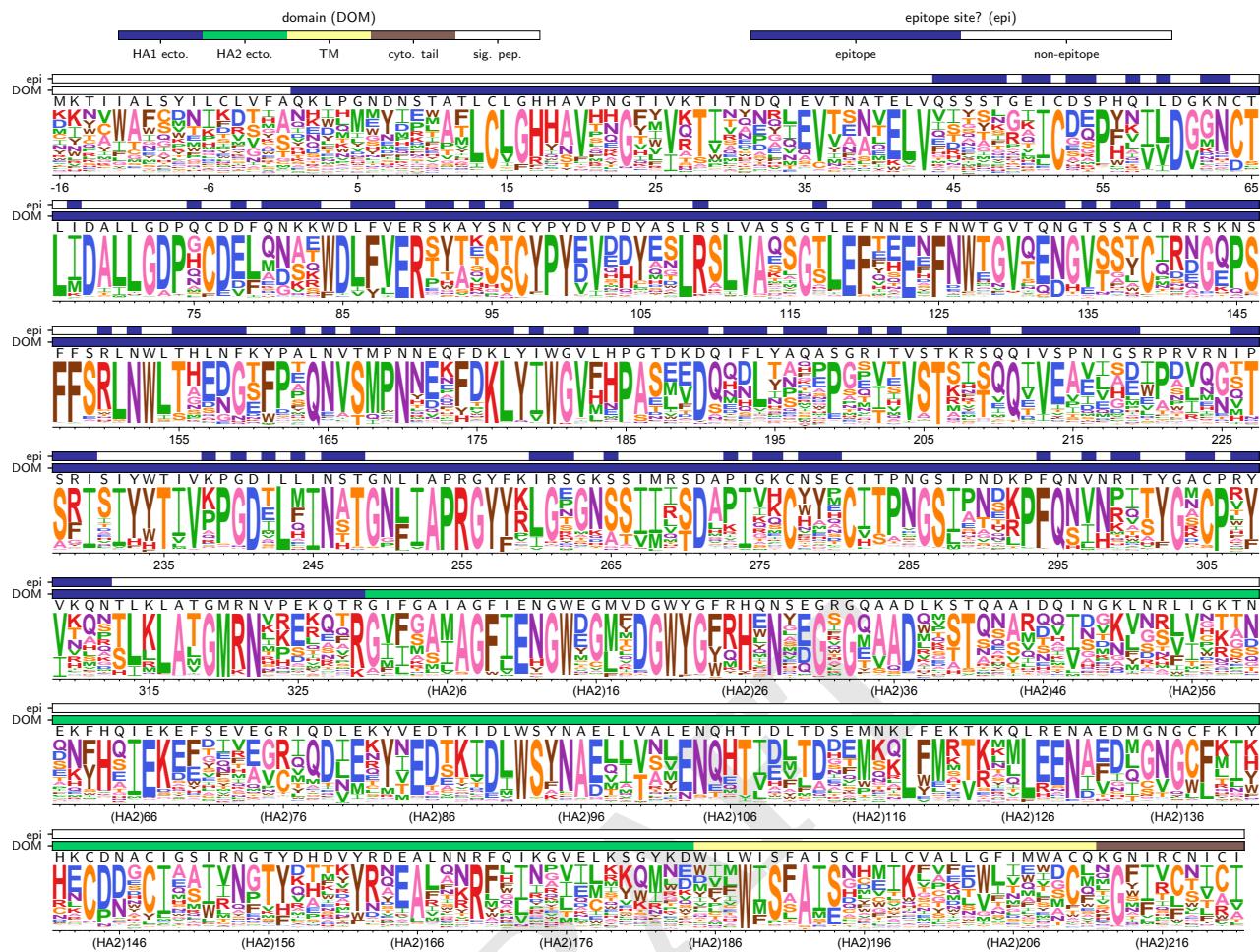
We performed a similar analysis using the current data for the Perth/2009 H3 HA. Surprisingly, there was much less contrast in mutational tolerance between the stalk and head domains for the H3 HA than for the H1 (Figure 4). For instance, in the H3 HA, the short helix A in the stalk is very mutationally tolerant. Interestingly, there are more reports of readily selecting escape mutants from broadly neutralizing anti-stalk antibodies in H3 (52–55) than in H1 (56–58) HAs.

We also see high mutational tolerance in many of the known antigenic regions of H3 HA (23). For instance, in recent H3N2 strains, antigenic region B is immunodominant, and most recent major antigenic drift mutations have occurred in this region (24, 25, 59). We find that the most distal portion of the globular head near the 190-helix, which is part of antigenic region B, is highly tolerant of mutations (Figure 4). Antigenic region C is also notably mutationally tolerant.

Many residues inside HA's receptor binding pocket are known to be highly functionally constrained (45, 60), and our data indicates that these sites are relatively mutationally intolerant in both H3 and H1 HAs. In contrast, the residues surrounding the receptor binding pocket are fairly mutationally tolerant, which may contribute to the rapidity of influenza's antigenic evolution, since mutations at these sites can have large effects on antigenicity (23, 24).

**The experimental measurements can help discriminate between successful and unsuccessful influenza virus lineages.** A major goal in the study of rapidly evolving viruses such as influenza is to forecast evolutionary trajectories (10, 63). For instance, forecasting which viral lineages will dominate the upcoming influenza season is an integral part of vaccine-strain selection (10, 63, 64). Evolutionary forecasts must ultimately distinguish between successful and unsuccessful viral lineages, which in the case of human influenza virus means distinguishing between the trunk and side branches of the phylogenetic tree (Figure 1).

To investigate whether our experiments can aid in distinguishing trunk and side branch lineages, we calculated the experimentally measured effects of all mutations in a maximum-likelihood reconstruction of the phylogeny in Figure 1. The



**Fig. 3.** The site-specific amino-acid preferences of the Perth/2009 HA measured in our experiments. The height of each letter is the preference for that amino acid, after taking the average over experimental replicates and re-scaling (40) by the stringency parameter in Table 1. The sites are in H3 numbering. The top overlay bar indicates whether or not a site is in the set of epitope residues delineated in (51). The bottom overlay bar indicates the HA domain (sig. pep. = signal peptide, HA1 ecto. = HA1 ectodomain, HA2 ecto. = HA2 ectodomain, TM = transmembrane domain, cyto. tail. = cytoplasmic tail). The letters directly above each logo stack indicate the wildtype amino acid at that site.



**Fig. 4.** Mutational tolerance of each site in H3 and H1 HAs. Mutational tolerance as measured in the current study is mapped onto the structure of the H3 trimer (PDB 4O5N; (61)). Mutational tolerance of the WSN/1933 H1 HA as measured in (30) is mapped onto the structure of the H1 trimer (PDB 1RVX; (62)). Different color scales are used because measurements are comparable among sites within the same HA, but not necessarily across HAs. Both trimers are shown in approximately the same orientation. For each HA, the structure at left shows a surface representation of the full trimer, while the structure at right side shows a ribbon representation of just one monomer. The sialic acid receptor is shown in black sticks.

mutations that occurred on the trunk of the tree were consistently more beneficial for viral growth according to our experimental measurements (Figure 5A). This difference in the effects of mutations between the trunk and side branch lineages was statistically significant when taken across the entire phylogeny (Figure 5B). Some influenza sequences are derived from egg- or cell-passaged isolates that contain lab-adaptation mutations (65–67). These mutations mostly occur on the terminal side branches that lead to tip nodes on the phylogenetic tree. We therefore re-calculated the statistics separating side branches to tip and internal nodes, and again found that the difference between the trunk and both sets of side branches was statistically significant (Figure 5B). Therefore, strains with mutations that we experimentally measure to be favorable for viral growth tend to do better in nature than strains with mutations that we measure to be less favorable.

Our experiments were performed on the Perth/2009 HA, but we are scoring mutations on a phylogenetic tree that extends from 1968 to 2017. Because the effects of mutations can vary with genetic background (68–72), it is possible that the effects of mutations that we measured in the Perth/2009 HA might have shifted somewhat in other related HAs. To

explore this question, we scored the complete HA sequence  $\mathbf{s}$  of every node in the phylogenetic tree by quantifying its average per-site sequence preference as

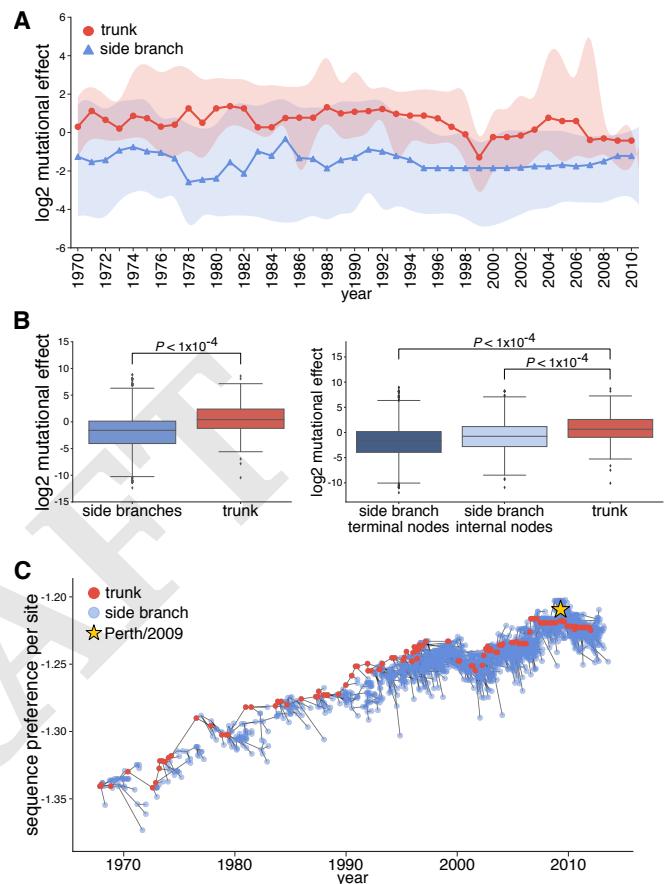
$$F(\mathbf{s}) = \frac{1}{L} \sum_{r=1}^L \ln \pi_{r,s_r}, \quad [1]$$

where  $\pi_{r,s_r}$  is the preference measured in our experiments (e.g., Figure 3) for the amino acid  $s_r$  at site  $r$  in HA sequence  $\mathbf{s}$ , and  $L$  is the length of the sequence. Figure 5C shows that the sequences of trunk nodes tend to have higher preference than those of side branch nodes across the entire timespan, consistent with the finding that trunk mutations are generally more favorable than side branch mutations. However, the average per-site sequence preference increases as the nodes approach the Perth/2009 strain. The Perth/2009 itself has among the highest per-site sequence preference of the entire tree, despite falling on a side branch. We do not think that this result means that Perth/2009 is actually higher fitness than the other strains. Rather, if there are modest shifts in the effects of mutations as HA evolves, then sequences more distant from Perth/2009 could be scored as having lower preference simply because these shifts start to degrade the accuracy of our measurements. But importantly, our experiments consistently show the trunk to be more favorable than the side branches regardless of whether we look at sequences that precede Perth/2009 (most of Figure 5C) or sequences that occur after Perth/2009 has diverged from the trunk (Figure S4). Therefore, our experimental measurements consistently reveal the selective advantage of the trunk over more than a half-century of viral evolution.

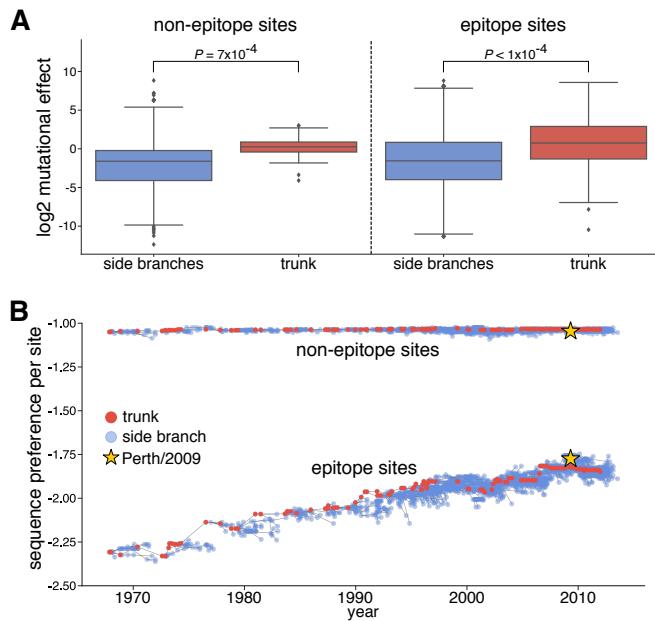
**Our experiments suggest different patterns of evolution at epitope and non-epitope sites.** HA is under selection both to maintain its essential function in viral growth and to escape pre-existing immunity. Most immune selection is focused on a subset of so-called “epitope sites” that are targets of the immunodominant antibody response. Although both epitope and non-epitope sites evolve rapidly, immune selection drives a higher rate of evolution at epitope sites (2, 9, 51). There are various classifications of HA epitope sites; here we will use the classification of Wolf *et al* (51), which defines 129 epitope sites among the 567 residues in HA. In the timeframe from 1968 to 2012, the trunk of the tree in Figure 1 fixed 0.84 epitope mutations per site, but only 0.07 non-epitope mutations per site. Since our experiments only measure how mutations affect viral growth, they might be expected to describe the evolution of epitope and non-epitope sites differently.

Mutations that occur on the trunk are scored as significantly more favorable in our experiments than side-branch mutations at both epitope and non-epitope sites (Figure 6A). Therefore, our experiments can distinguish evolutionarily favorable mutations both at sites that are predominantly under functional constraint and at sites that are also under immune pressure. The fact that our experiments can discriminate between the trunk and side branches at epitope sites despite only assaying for viral growth underscores the fact that most epitope sites are still important for HA function (24, 73, 74).

However, there are differences in how the average per-site sequence preference changes over time at the epitope and non-epitope sites (Figure 6B). The per-site sequence preference at epitope sites increases over time, whereas at non-epitope sites



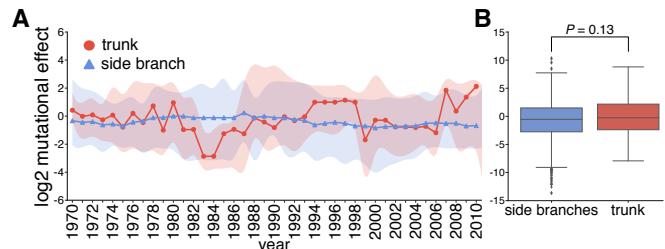
**Fig. 5. Mutations on the trunk are more favorable than those on the side branches according to our experimental measurements.** (A) We used our experiments to calculate the  $\log_2$  mutational effect for all trunk and side branch mutations in five-year windows from 1968 to 2013. The central year in each window is denoted on the x-axis. Median mutational effects in each window are shown as circles for the trunk and triangles for side branches. Shaded regions demarcate the interquartile range. Negative numbers signify mutations to less preferred amino acids, while positive numbers signify mutations to more preferred mutations. (B) The  $\log_2$  mutational effect for all side branch and trunk mutations (left panel), and the same data but separating side branches to terminal and internal nodes.  $P$ -values were computed by randomizing the experimental measurements among sites 10,000 times, and determining how often the difference in median mutational effect between the trunk and side branches exceeded the actual value. (C) The average per-site sequence preference of every node in the phylogenetic tree in Fig. 1. Trunk nodes are in red, side branch nodes in blue, and Perth/2009 is marked with a yellow star. More preferred sequences have larger values.



**Fig. 6.** Effects of mutations at epitope and non-epitope sites during HA evolution. We partitioned HA into epitope and non-epitope sites (51). (A) The  $\log_2$  mutational effects for side branch and trunk mutations at non-epitope (left) and epitope (right) sites.  $P$ -values were computed as in Figure 5 but only performing the randomizations over the appropriate set of sites (non-epitope or epitope). (B) The average per-site sequence preference for all nodes in the phylogenetic tree, calculated separately for each set of sites. The sequence preferences for epitope and non-epitope sites on separate y-axes is shown in Figure S5.

it remains relatively constant (Figure 6B and Figure S5). In addition, the per-site sequence preference is consistently higher at non-epitope than epitope sites, perhaps because immune pressure sometimes drives the fixation of functionally less favorable mutations at epitope sites. The fact that only epitope sites exhibit an increase in sequence preference over time can also be rationalized in terms of prior knowledge about protein evolution. As proteins evolve, there can be epistatic shifts in the effects of mutations (68–72, 75). Such epistasis has been experimentally demonstrated for HA (74, 76), including at some epitope sites in H3 HA (77). It is also known that immune pressure can increase the amount of epistasis among substitutions by selecting for combinations of functionally deleterious immune-escape mutations and counterbalancing secondary mutations (78). We hypothesize that immune pressure increases the accumulation of epistatic interactions involving substitutions at epitope sites. Over time, these epistatic interactions could manifest as shifts in the effects of mutations, leading to an apparent increase in per-site sequence preference over time when mutational effects are measured in the Perth/2009 HA (Figure 6B and Figure S5). Such a trend is not apparent for non-epitope sites, probably because they are under less selection to fix adaptive mutations.

**There are large differences between H3 and H1 HAs in the amino-acid preferences of many sites.** The results above show that measurements of mutational effects in the Perth/2009 HA are informative about the evolutionary fate of other human H3N2 HAs. How broadly can such measurements be generalized across HAs? To address this question, we repeated the trunk versus side-branch analyses using amino-acid preferences



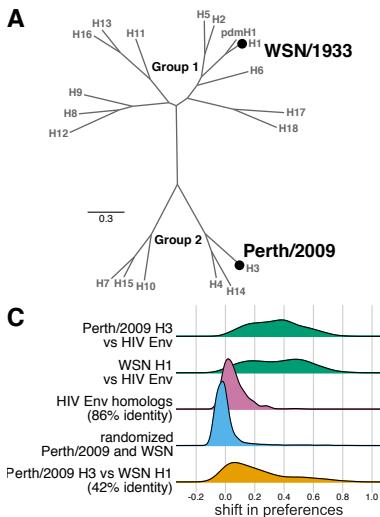
**Fig. 7.** Experimental measurements on an H1 HA do not provide evolutionarily relevant information for H3 HAs. (A), (B) This figure is analogous to that in Figure 5 except that it scores the H3 sequences using experimental measurements made on the lab-adapted WSN/1933 HA (30).

from our prior deep mutational scanning of the WSN/1933 H1 HA, which is highly diverged from the Perth/2009 H3 HA. Figure 7 shows that the H1 measurements are not informative for distinguishing the trunk and side branches of an H3 phylogeny. This fact indicates that the utility of an experiment for discriminating successful and unsuccessful strains degrades over sufficiently long evolutionary distances.

An obvious hypothesis for why the H1 deep mutational scanning is not informative about the evolutionary fate of H3 strains is that the effect of the same mutation will often be very different between these two HA subtypes. To determine if this is the case, we can examine how much the amino-acid preferences of homologous sites have shifted between H3 and H1 HAs. Prior experiments have found only modest shifts in amino-acid preferences between two variants of influenza nucleoprotein with 94% amino-acid identity (79) and variants of HIV envelope (Env) with 86% amino-acid identity (75). However, the H1 and H3 HAs are far more diverged, with only 42% amino-acid identity (Figure 8A). One simple way to investigate the extent of shifts in amino-acid preferences is to correlate measurements from independent deep mutational scanning replicates on H1 and H3 HAs. Figure 8B shows that replicate measurements on the same HA variant are more correlated than those on different HA variants.

To more rigorously quantify shifts in amino-acid preferences after correcting for experimental noise, we used the statistical approach in (75, 79). Figure 8C shows the distribution of shifts in amino-acid preferences between H3 and H1 HAs after correcting for experimental noise. Although some sites have small shifts near zero, many sites have large shifts. These shifts between H3 and H1 are much larger than expected from the null distribution that would be observed purely from experimental noise. They are also much larger than the shifts previously observed between two HIV Envs with 86% amino-acid identity (75). However, the shifts are still smaller than those observed when comparing HA to the non-homologous HIV Env protein. Therefore, there are very substantial shifts in mutational effects between highly diverged HA homologs, although the effects of mutations remain more similar than for non-homologous proteins.

**Properties associated with the shifts in amino-acid preferences between H3 and H1 HAs.** What features distinguish the sites with shifted amino-acid preferences between H3 and H1 HAs? The sites of large shifts do not obviously map to one specific region of HA's structure (Figure 9A). However, at the domain level, sites in HA's stalk tend to have smaller shifts than sites



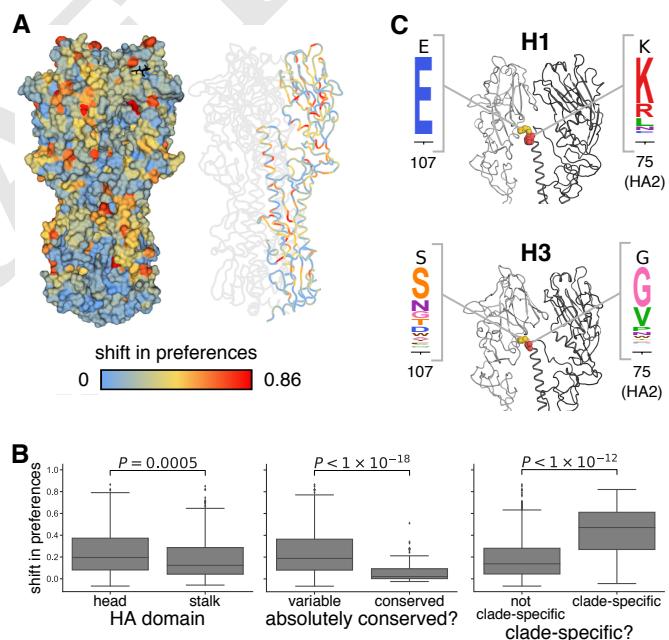
**Fig. 8.** There are substantial differences in the effects of mutations between H1 and H3 HAs. (A) Phylogenetic tree of HA subtypes, with the WSN/1933 H1 and Perth/2009 H3 HAs labeled. These HAs have 42% amino-acid identity. (B) All pairwise correlations of the amino-acid preferences measured in the three individual deep mutational scanning replicates in the current study and the three replicates in prior deep mutational scanning of an H1 HA (30). Comparisons between H3 replicates are in purple, those between H1 replicates are in brown, and those across H1 and H3 replicates are in gray.  $R$  indicates the Pearson correlation coefficient. (C) We calculated the shift in amino-acid preferences at each site between H3 and H1 HAs using the method in (75), and plotted the distribution of shifts for all sites. The shifts between H3 and H1 (yellow) are much larger than the null distribution (blue) expected if all differences are due to experimental noise. The shifts are also much larger than those previously observed between two variants of HIV envelope protein (Env) that share 86% amino-acid identity (pink). However, the shifts between H3 and H1 are still less than the differences between either HA and HIV Env (green).

in HA's globular head (Figure 9B). The HA stalk domain is also more conserved in sequence (80–82), suggesting that conservation of amino-acid sequence tends to be correlated with conservation of amino-acid preferences. Consistent with this idea, sites that are absolutely conserved across all 18 HA subtypes are significantly less shifted than sites that are variable across HA subtypes (Figure 9B). Presumably these sites are under consistent functional constraint across all HAs.

Despite their high sequence divergence, H1 and H3 adopt very similar protein folds (83, 84). However, there are differences in the rotation and upward translation of the globular head subdomains relative to the central stalk domain among different HA subtypes (83, 84). Previous work has defined clades of structurally related HA subtypes (83, 84). One such clade includes H1, H2, H5, and H6, whereas another clade includes H3, H4, and H14 HAs (Figure 8). Sites that are conserved at different amino-acid identities in these two clades tend to have exceptionally large shifts in amino acid preferences (Figure 9B). The clade containing H1 has an upward shift of the globular head relative to the clade containing H3. This structural shift has been attributed largely to the interaction between sites 107 and 75(HA2) (83, 84). Specifically, the clade containing H1 has a taller turn in the interhelical loop connecting helix A and helix B in the stalk domain, and this tall turn is stabilized by a hydrogen bond between Glu-107 and Lys-75(HA2) (Figure 9C). In deep mutational scanning of the H1 HA, site 107 has a high preference for Glu and 75(HA2) strongly prefers positively charged Lys and Arg. In contrast, the interhelical loop in H3 HA makes a sharper and shorter turn which is facilitated by a Gly at 75(HA2). In the deep mutational scanning of the Perth/2009 H3 HA, site 75(HA2) prefers Gly and to a lesser extent Val, while site 107 is fairly tolerant of mutations. Therefore, some of the shifts in HA amino-acid preferences can be directly rationalized in terms of changes in HA structure.

## Discussion

We have measured the effect of all possible single amino-acid mutations to the Perth/2009 H3 HA on viral growth in cell



**Fig. 9.** Sites with strongly shifted amino-acid preferences between H3 and H1 HAs. (A) The shift in amino-acid preferences between the H3 and H1 HA at each site as calculated in Figure 8C is mapped onto the structure of the H3 HA. (B) Amino-acid preferences of sites in the stalk domain are less shifted than those in the head domain. Sites absolutely conserved in all 18 HA subtypes are less shifted than other sites. Sites with one amino-acid identity in the clade containing H1, H2, H5, and H6 and another identity in the clade containing H3, H4, and H14 are more shifted than other sites. (C) Sites 107 and 75(HA2) help determine the different orientation of the globular head domain in H1 versus H3 HAs. These sites are shown in spheres on the structure of H1 and H3 and colored as in panel (A), and the experimentally measured amino-acid preferences in the H1 and H3 HAs are shown. One monomer is in dark gray, while the HA1 domain of the neighboring monomer is in lighter gray.

culture, and demonstrated that these measurements are informative for distinguishing the evolutionary fate of human H3N2 viral strains in nature. Specifically, HA sequences on the trunk of the phylogenetic tree tend to have mutations that our experiments measure to be more beneficial for viral growth than mutations that occur in side-branch HA sequences. The fact that our experiments can help distinguish between trunk and side-branch lineages suggests that they might inform evolutionary forecasting. In their landmark paper introducing predictive viral fitness models that accounted for both antigenic and non-antigenic mutations, Luksza and Lassig (16) noted that the models could in principle be improved by integrating “diverse genotypic and phenotypic data” that more realistically represented the effects of specific mutations. Our work suggests that deep mutational scanning can provide such data.

It is important to emphasize that measurements of viral growth in cell culture do *not* represent true fitness in nature. Indeed, a vast amount of work in virology has chronicled the many ways in which experiments can select for lab artifacts or fail to capture important pressures that are relevant in nature (65, 85–87). Mutations in viral genes other than HA are also certainly important in determining strain success (88, 89). Given these caveats, it might seem surprising that simply measuring viral growth in cell culture can be informative about the success of viral strains in nature. Yet, prior to our work, there were no comprehensive studies of the functional effects of mutations to H3 HA on any property that even resembled viral fitness in nature, and modeling work has either omitted the non-antigenic effects of mutations (11–13) or simply assumed that all non-epitope mutations had equivalent deleterious effects (16). The strength of our measurements are not that they perfectly capture fitness in nature, but that they are systematic and quantitative — and represent a vast improvement over no information at all. That this is true is demonstrated by the fact that these measurements have utility for distinguishing successful and unsuccessful viral lineages over ~50 years of H3N2 evolution. We suspect that performing the experiments using more realistic and complex selections (e.g., ferrets or primary human airway cultures (?)) might further improve their utility.

We measured the effects of all single amino-acid mutations to a specific HA, and then generalized these measurements to other H3N2 HAs from a 50-year timespan. These generalizations will only be valid to the extent that the effects of mutations are conserved during HA’s evolution. Extensive work on protein evolution has shown that epistasis can shift the effects of mutations (68–72, 75), and the extent to which this occurs is an important question in its own right. The fact that our measurements score trunk mutations more favorably than side branch ones across a 50-year timespan suggests that the mutational effects have not shifted to dramatically among these relatively closely related sequences. On the other hand, when we compared our measurements for an H3 HA to prior measurements on H1 HA, we found substantial shifts at many sites — much greater than those observed in prior protein-wide comparisons of more closely related homologs (75, 79). Further investigation of how mutational effects shift as proteins diverge will be important for determining how broadly any given experiment can be generalized when attempting to make evolutionary forecasts.

Our work did not characterize the antigenic effects of mutations, which also play an important role in determining strain success in nature (?). Indeed, some of our results (e.g., Figure 6) suggest as much by showing that the experiments better capture selection at non-epitope than epitope sites. However, our basic selection and deep-sequencing approach can be harnessed to completely map how mutations affect antibody recognition (57, 90). But so far, experiments using this approach have either failed to examine antibodies or sera that are relevant to driving the evolution of H3N2 influenza (57, 90), or have used relevant sera but examined a non-comprehensive set of mutations (26). Future experiments that completely map how HA mutations affect recognition by human sera seem likely to be especially fruitful for informing viral forecasting.

## Materials and Methods

Deep sequencing data are available from the Sequence Read Archive under BioSample accessions SAMN08102609 and SAMN08102610. Computer code used to analyze the data and produce the results in the paper are in [link to repository].

**HA numbering.** Unless otherwise indicated, all sites are in H3 numbering, with the signal peptide in negative numbers, the HA1 subunit in plain numbers, and the HA2 subunit denoted with "(HA2)". The conversion between sequential numbering of the A/Perth/16/2009 HA and H3 numbering was performed using an HA numbering Python script (available at [https://github.com/jbloomlab/HA\\_numbering](https://github.com/jbloomlab/HA_numbering)).

**Creation of MDCK-SIAT1-TMPRSS2 cell line.** When growing influenza virus, TPCK-trypsin must normally be added in order to cleave HA into its mature form. In order to obviate the need for trypsin, we engineered an MDCK-SIAT1 cell line to constitutively express the TMPRSS2 protease. This protease cleaves and activates HA and is found endogenously in the human airway (35, 36). The human TMPRSS2 cDNA ORF was ordered from OriGene (NM\_005656), PCR amplified, and cloned into a pHAGE2 lentiviral vector under an EF1α-Int promoter followed by an IRES driving expression of mCherry. We used the lentiviral vector to transduce MDCK-SIAT1 cells, and sorted an intermediate mCherry population by flow cytometry. We refer to the resulting bulk population of sorted cells as MDCK-SIAT1-TMPRSS2 cells. There is no selectable marker for the TMPRSS2; however, we generally maintain the cells at low passage number, and have seen no indication that lose their ability to support the growth of viruses with H3 HAs in the absence of exogenous trypsin.

**Generation of HA codon mutant plasmid libraries.** Recombinant A/Perth/16/2009 (HA, NA) × A/Puerto Rico/8/1934 influenza virus, NIB-64, NR-41803 was ordered from BEI Resources, NIAID, NIH. Bulk RNA from the viral sample was extracted using the QIAamp Viral RNA Mini Kit (QIAGEN) according to manufacturer’s instructions. The Perth/2009 HA and NA genes were then reverse transcribed using the SuperScript III First-Strand Synthesis SuperMix Kit (Invitrogen), PCR amplified, and cloned into the pHW2000 (92) and pICR2 (93) plasmid backbones by In-Fusion.

To generate wildtype Perth/2009 H3N2 virus, we transfected in triplicate co-cultures of 293T and MDCK-SIAT1-TMPRSS2 cells with reverse genetics plasmids encoding the Perth/2009 HA and NA, WSN/1933 internal genes, and the TMPRSS2 lentiviral vector. [add more details about reverse genetics?] The viral supernatant was harvested at 72 hours post-transfection, clarified at 1200 rpm for 5 min, and 0.01 to 1.0 ul of the supernatant was passaged onto  $4 \times 10^5$  MDCK-SIAT1-TMPRSS2 cells in six-well plates for a total of six serial passages. After the final passage, RNA was extracted from 140 ul of the viral supernatant for each replicate using the QIAamp Viral RNA Mini Kit. The HA and NA genes were reverse transcribed, PCR amplified, and sequenced by Sanger sequencing. Two mutations, G78D and T212I, arose to fixation in one of the replicates after serial passaging. The Perth/2009 HA carrying these

two mutations were cloned into the pHW2000 (92) and pICR2 (93) plasmid backbones.

To validate that the variant Perth/2009 HA could support viral growth, we generated in duplicate by reverse genetics viruses with the original Perth/2009 HA or the Perth/2009 HA-G78D-T212I variant on a background of Perth/2009 NA + WSN internal genes. The transfection supernatant was collected 72 hours post-transfection and passaged onto a monolayer of MDCK-SIAT1-TMPRSS2 cells at an MOI = 0.01. The viral supernatant was collected at 44 hours post-infection. Both the transfection and passage supernatants were titered by TCID<sub>50</sub> in MDCK-SIAT1-TMPRSS2 cells.

The codon-mutant libraries were generated in the Perth/2009 HA-G78D-T212I background using a PCR-based approach described in (91, 94) using two rounds of mutagenesis. The script used to design the mutagenesis primers is on <https://github.com/jbloomlab/CodonTilingPrimers>. The codon mutations were introduced into three independent replicates of the Perth/2009 HA variant. The mutant variants were then cloned into the pICR2 vector by digestion with BsmBI restriction sites, ligation by T4 DNA ligase, and electroporation and transformation into ElectroMAX DH10B competent cells (Invitrogen; 18290015). We pooled over 6 million transformants for each replicate, cultured in liquid LB + ampicillin at 37°C for 3 h with shaking, and maxiprepped. We randomly chose 31 clones to Sanger sequence, and summary statistics from sequencing are shown in Figure S2.

**Generation and passaging of mutant viruses.** The mutant virus libraries were generated and passaged using the approach described in (30) with several modifications. We used the HA-deficient WSN/1933 helper virus generated in (30). Briefly, these helper viruses were generated by transfecting seven WSN/1933 non-HA reverse-genetics plasmids plus a WSN/1933 HA protein expression plasmid into a coculture of 293T and MDCK-SIAT1-EF1α-WSN-HA cells engineered to constitutively express the WSN/1933 HA. The transfection supernatant was then expanded by passaging onto a monolayer of MDCK-SIAT1-EF1α-WSN-HA cells, and the viral supernatant was then harvested and aliquots were frozen at -80°C.

To generate the mutant virus libraries, we transfected  $5 \times 10^5$  MDCK-SIAT1-TMPRSS2 cells in suspension with 937.5 ng each of four protein expression plasmids encoding the ribonucleoprotein complex (HDM-Nan95-PA, HDM-Nan95-PB1, HDM-Nan95-PB2, and HDM-Aichi68-NP) (68), and 1250 ng of one of the three pICR2-mutant-HA libraries (or the wildtype pICR2-Perth2009-HA-G78D-T212I control) using Lipofectamine3000 (ThermoFisher; L3000008). We allowed the transfected cells to adhere in 6-well plates, then four hours later changed the media to D10 media (DMEM supplemented with 10% heat-inactivated FBS, 2 mM L-glutamine, 100 U of penicillin/mL, and 100 μg of streptomycin/mL). Eighteen hours after transfection, we infected the cells with WSN/1933 HA-deficient helper virus by preparing an inoculum of 500 TCID<sub>50</sub> per μL in influenza growth media (Opti-MEM supplemented with 0.01% heat-inactivated FBS, 0.3% BSA, 100 U of penicillin/mL, 100 μg of streptomycin/mL, and 100 μg of calcium chloride/mL), aspirating the D10 media from the cells, and adding 2 mL of the helper-virus inoculum to each well. After three hours, we changed the media to fresh influenza growth media. Twenty-four hours after helper-virus infection, we harvested the viral supernatants for each replicate, froze aliquots at -80°C, and titered in MDCK-SIAT1-TMPRSS2 cells.

We passaged over  $9 \times 10^5$  TCID<sub>50</sub> of the transfection supernatants at an MOI of 0.0035 TCID<sub>50</sub> per cell. To passage, we plated  $4.6 \times 10^6$  MDCK-SIAT1-TMPRSS2 cells per dish in 15-cm dishes in D10 media, and allowed the cells to grow for 24 hours, at which the cells reached a density of  $\sim 1.7 \times 10^7$  cells per dish. We then replaced the media of each dish with 25 mL of an inoculum of 2.5 TCID<sub>50</sub> of virus per μL. We collected viral supernatant for sequencing 48 hours post-infection.

**Barcoded subamplicon sequencing.** To extract viral RNA from the three replicate HA virus libraries and the wildtype HA virus, we first ultracentrifuged 24 mL of the supernatant at 22,000 rpm for 1.5 h at 4°C in a Beckman Coulter SW28 rotor. We then extracted the RNA using the Qiagen RNeasy Mini Kit by resuspending the viral pellet in 400 μL of buffer RLT supplemented with β-mercaptoethanol, pipetting 30 times, transferring the liquid to a microcentrifuge

tube, adding 600 μL 70% ethanol, and proceeding with the rest of the extraction according to manufacturer's instructions. The HA gene was then reverse-transcribed with AccuScript Reverse Transcriptase (Agilent 200820) using the primers P09-HA-For (5'-AGCAAAAGCAGGGATAATTCTATTAAATC-3') and P09-HA-Rev (5'-AGTAGAACAAAGGGTGTTTAATTACTAATACAC-3').

We generated the HA PCR amplicons for each of the three plasmid libraries, the three virus libraries, one wildtype plasmid, and one wildtype virus samples using KOD Hot Start Master Mix (EMD Millipore; 71842) according to the PCR reaction mixture and cycling conditions described in (94) and the P09-HA-For and P09-HA-Rev primers. We prepared the sequencing libraries using a barcoded-subamplicon strategy (29) to increase the accuracy from deep sequencing. The approach we used is described in (30) (also see [https://jbloomlab.github.io/dms\\_tools2/bcsubamp.html](https://jbloomlab.github.io/dms_tools2/bcsubamp.html)). The primers we used to generate the subamplicons are in Dataset S2. We performed sequencing on a lane of a flow cell of an Illumina HiSeq 2500 using 2 × 250 bp paired-end reads in rapid-run mode.

**Analysis of deep sequencing data.** We used the dms\_tools2 software package (38) ([https://github.com/jbloomlab/dms\\_tools2](https://github.com/jbloomlab/dms_tools2)) to analyze the deep sequencing data. The algorithm used to estimate the site-specific amino-acid preferences from the deep sequencing counts is described in (38). The amino-acid preferences for each replicate and for the re-scaled, cross-replicate average are provided in Datasets [suppfile].

**Alignments and phylogenetic analysis of HA sequences.** To build the ExpCM using our measurements, we first downloaded all full-length H3 HA sequences from the Influenza Virus Resource (95), and randomly subsampled two sequences per year. The subsampled sequences were then aligned using MAFFT (96) and used to build a phylogenetic tree using RAxML (97). We then used the phydms program (40) (<https://github.com/jbloomlab/phydms>) to fit phylogenetic substitution models using the aligned subsampled sequences.

A description of the sequences and methods used to build the phylogenetic tree of HA subtypes shown in Figure 8A is provided in (57).

**Inference of human H3N2 phylogenetic tree.** [John to add methods here. We downloaded X sequences from the Influenza Virus Resource ?.... etc. subsampled how? inferred the tree, ancestral state reconstruction, visualized the tree...]. To parse out trunk mutations from side branch mutations, we first defined a set of recent nodes sampled on or after Jan. 1, 2017, and traced these nodes back to their most recent common ancestor. We defined all branches ancestral to this most recent common ancestor of recently sampled nodes as the trunk. All nodes descended from the trunk were defined as side branch nodes. Nodes descended from the most recent common ancestor were defined as unresolved.

**Quantification of mutational effects and sequence preferences from an H3N2 phylogeny.** For a given mutation a1-r-a2, we calculated the mutational effect as:

$$\log_2 \frac{\pi_{r,a2}}{\pi_{r,a1}} \quad [2]$$

where  $\pi_{r,a1}$  and  $\pi_{r,a2}$  are the preferences for amino-acid a1 or a2 at site r, and the preferences are re-scaled by the stringency parameter provided in Table 1 and averaged across replicates. The WSN/1933 H1 HA amino-acid preferences were also re-scaled 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)) and averaged across replicates. The mutational effects were also calculated for five-year windows from 1968 to 2012, and windows were incremented by one year.

To calculate significance, we randomized the amino-acid preferences across all sites, epitope sites, or non-epitope sites, and quantified the median difference in trunk and side branch mutational effects. We repeated this randomization for a total of 10,000 times and calculated how frequently the difference in median mutational effect between trunk and side branches exceeded the true difference to obtain a P-value.

Sequence preference per site for every node on the phylogenetic tree shown in Figure 1 was calculated using Equation 1 for all sites, or the set of epitope and non-epitope sites defined in (51).

**Analysis of mutational shifts.** To compare the Perth/2009 H3 and WSN/1933 H1 HA preferences, we first aligned the wildtype HA sequences using MAFFT (96). To quantify the shifts in preference for every alignable site while accounting for experimental noise, we used the approach described in (75) and based off the method presented in (79). Briefly, for each site we calculated half of the sum of the absolute value of the difference between preferences for each amino acid. The RMSD<sub>corrected</sub> value signifies the shift in preference.

For the plots shown in Figure 9B, any residues falling between Cys-52 and Cys-277 were defined as the head domain, and all other residues were defined as the stalk domain. We used a multiple sequence alignment of the HA subtype sequences provided in (57) to list sites that are absolutely conserved across all subtypes. To find clade-specific sites, we listed sites with a single amino-acid identity in all of H1, H2, H5, and H6, and a different identity in H3, H4, and H14.

**ACKNOWLEDGMENTS.** We thank Sarah Hilton, Hugh Haddox, and Sidney Bell for helpful discussions about data analysis. We thank the Fred Hutch Genomics Core for performing the Illumina deep sequencing. This work was supported by...

## References.

- Smith DJ, et al. (2004) Mapping the antigenic and genetic evolution of influenza virus. *Science* 305(5682):371–376.
- Bhatt S, Holmes EC, Pybus OG (2011) The genomic rate of molecular adaptation of the human influenza A virus. *Molecular Biology and Evolution* 28(9):2443.
- Fitch WM, Bush RM, Bender CA, Cox NJ (1997) Long term trends in the evolution of H(3) HA1 human influenza type A. *Proc. Natl. Acad. Sci. USA* 94(15):7712–7718.
- Strelkowa N, Lässig M (2012) Clonal interference in the evolution of influenza. *Genetics* 192(2):671–682.
- Bedford T, Cobey S, Pascual M (2011) Strength and tempo of selection revealed in viral gene genealogies. *BMC evolutionary biology* 11(1):220.
- Neher RA, Russell CA, Shraiman BI (2014) Predicting evolution from the shape of genealogical trees. *Elife* 3:e03568.
- Koelle K, Rasmussen DA (2015) The effects of a deleterious mutation load on patterns of influenza A/H3N2's antigenic evolution in humans. *Elife* 4:e07361.
- Bedford T, et al. (2015) Global circulation patterns of seasonal influenza viruses vary with antigenic drift. *Nature* 523(7559):217–220.
- Fitch WM, Leiter J, Li Y, Palese P (1991) Positive Darwinian evolution in human influenza A viruses. *Proc. Natl. Acad. Sci. USA* 88(10):4270–4274.
- Morris DH, et al. (2017) Predictive modeling of influenza shows the promise of applied evolutionary biology. *Trends in Microbiology*.
- Sun H, et al. (2013) Using sequence data to infer the antigenicity of influenza virus. *MBio* 4(4):e00230–13.
- Harvey WT, et al. (2016) Identification of low-and high-impact hemagglutinin amino acid substitutions that drive antigenic drift of influenza A (H1N1) viruses. *PLoS Pathogens* 12(4):e100526.
- Neher RA, Bedford T, Daniels RS, Russell CA, Shraiman BI (2016) Prediction, dynamics, and visualization of antigenic phenotypes of seasonal influenza viruses. *Proceedings of the National Academy of Sciences* 113(12):E1701–E1709.
- Pybus OG, et al. (2007) Phylogenetic evidence for deleterious mutation load in RNA viruses and its contribution to viral evolution. *Molecular Biology and Evolution* 24(3):845–852.
- Kucharski A, Gog JR (2011) Influenza emergence in the face of evolutionary constraints. *Proceedings of the Royal Society of London B: Biological Sciences* p. rspb20111168.
- Łukasz M, Lässig M (2014) A predictive fitness model for influenza. *Nature* 507(7490):57–61.
- Holland J, et al. (1982) Rapid evolution of RNA genomes. *Science* 215(4540):1577–1585.
- Steinhardt D, Holland J (1987) Rapid evolution of RNA viruses. *Annual Reviews in Microbiology* 41(1):409–431.
- Lauring AS, Andino R (2010) Quasispecies theory and the behavior of RNA viruses. *PLoS Pathogens* 6(7):e1001005.
- Boni MF, Zhou Y, Taubenberger JK, Holmes EC (2008) Homologous recombination is very rare or absent in human influenza A virus. *Journal of Virology* 82(10):4807–4811.
- Laver W, et al. (1979) Antigenic drift in type A influenza virus: sequence differences in the hemagglutinin of hong kong (H3N2) variants selected with monoclonal hybridoma antibodies. *Virology* 98(1):226–237.
- Webster R, Laver W (1980) Determination of the number of nonoverlapping antigenic areas on Hong Kong (H3N2) influenza virus hemagglutinin with monoclonal antibodies and the selection of variants with potential epidemiological significance. *Virology* 104(1):139–148.
- Wiley D, Wilson I, Skehel J, , et al. (1981) Structural identification of the antibody-binding sites of hong kong influenza haemagglutinin and their involvement in antigenic variation. *Nature* 289(5796):373–378.
- Koel BF, et al. (2013) Substitutions near the receptor binding site determine major antigenic change during influenza virus evolution. *Science* 342(6161):976–979.
- Chambers BS, Parkhouse K, Ross TM, Alby K, Hensley SE (2015) Identification of hemagglutinin residues responsible for H3N2 antigenic drift during the 2014–2015 influenza season. *Cell Reports* 12(1):1–6.
- Li C, et al. (2016) Selection of antigenically advanced variants of seasonal influenza viruses. *Nature Microbiology* 1:16058.
- Fowler DM, Fields S (2014) Deep mutational scanning: a new style of protein science. *Nature methods* 11(8):801–807.
- Thyagarajan B, Bloom JD (2014) The inherent mutational tolerance and antigenic evolvability of influenza hemagglutinin. *eLife* 3:e03300.
- Wu NC, et al. (2014) High-throughput profiling of influenza A virus hemagglutinin gene at single-nucleotide resolution. *Scientific Reports* 4:4942.
- Doud MB, Bloom JD (2016) Accurate measurement of the effects of all amino-acid mutations to influenza hemagglutinin. *Viruses* 8:155.
- Haddox HK, Dingens AS, Bloom JD (2016) Experimental estimation of the effects of all amino-acid mutations to HIV's envelope protein on viral replication in cell culture. *PLoS Pathogens* 12(12):e1006114.
- Qi H, Wu NC, Du Y, Wu TT, Sun R (2015) High-resolution genetic profile of viral genomes: why it matters. *Current Opinion in Virology* 14:62–70.
- WHO (2010) Recommended viruses for influenza vaccines for use in the 2010–2011 northern hemisphere influenza season. [http://www.who.int/influenza/vaccines/virus/recommendations/201002\\_Recommendation.pdf?ua=1](http://www.who.int/influenza/vaccines/virus/recommendations/201002_Recommendation.pdf?ua=1).
- WHO (2011) Recommended composition of influenza virus vaccines for use in the 2011–2012 northern hemisphere influenza season. [http://www.who.int/influenza/vaccines/2011\\_02\\_recommendation.pdf?ua=1](http://www.who.int/influenza/vaccines/2011_02_recommendation.pdf?ua=1).
- Böttcher E, et al. (2006) Proteolytic activation of influenza viruses by serine proteases TP-PRSS2 and HAT from human airway epithelium. *Journal of Virology* 80:9896–9898.
- Böttcher-Friebertshäuser, E, et al. (2010) Cleavage of influenza virus hemagglutinin by airway proteases TMPRSS2 and HAT differs in subcellular localization and susceptibility to protease inhibitors. *Journal of Virology* 81:5605–5614.
- Marshall N, Priyamvada L, Ende Z, Steel J, Lowen AC (2013) Influenza virus reassortment occurs with high frequency in the absence of segment mismatch. *PLoS Pathogens* 9(6):e1003421.
- Bloom JD (2015) Software for the analysis and visualization of deep mutational scanning data. *BMC Bioinformatics* 16(1):1.
- Bloom JD (2017) Identification of positive selection in genes is greatly improved by using experimentally informed site-specific models. *Biology Direct* 12(1):1.
- Hilton SK, Doud MB, Bloom JD (2017) phdms: Software for phylogenetic analyses informed by deep mutational scanning. *PeerJ* 5:e3657.
- Yang Z, Nielsen R, Goldman N, Pedersen AMK (2000) Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics* 155(1):431–449.
- Posada D, Buckley TR (2004) Model selection and model averaging in phylogenetics: advantages of Akaike information criterion and Bayesian approaches over likelihood ratio tests. *Systematic Biology* 53(5):793–808.
- Waterfield M, Scrafe G, Skehel J (1981) Disulphide bonds of haemagglutinin of Asian influenza virus. *Nature* 289:422–424.
- Weis W, et al. (1988) Structure of the influenza virus haemagglutinin complexed with its receptor, sialic acid. *Nature* 333:426–431.
- Martin J, et al. (1998) Studies of the binding properties of influenza hemagglutinin receptor-site mutants. *Virology* 241(1):101–111.
- Nobusawa E, Ishihara H, Morishita T, Sato K, Nakajima K (2000) Change in receptor-binding specificity of recent human influenza A viruses (H3N2): a single amino acid change in hemagglutinin altered its recognition of sialyloligosaccharides. *Virology* 278(2):587–596.
- Yang H, et al. (2015) Structure and receptor binding preferences of recombinant human A (H3N2) virus hemagglutinins. *Virology* 477:18–31.
- Kido H, et al. (1992) Isolation and characterization of a novel trypsin-like protease found in rat bronchiolar epithelial Clara cells. A possible activator of the viral fusion glycoprotein. *J Biol Chem* 267:13573–13579.
- Stech J, Garn H, Wegmann M, Wagner R, Klenk H (2005) A new approach to an influenza live vaccine: modification of the cleavage site of hemagglutinin. *Nature Medicine* 11(6):683–689.
- Girard G, Gulyaev A, Olsthoorn R (2011) Upstream start codon in segment 4 of North American H2 avian influenza A viruses. *Infect. Genet. Evol.* 11:489–495.
- Wolf Y, Viboud C, Holmes E, Koonin E, Lipman D (2006) Long intervals of stasis punctuated by bursts of positive selection in the seasonal evolution of influenza A virus. *Biology Direct* 1:34.
- Ekiert DC, et al. (2011) A highly conserved neutralizing epitope on group 2 influenza A viruses. *Science* 333(6044):843–850.
- Friesen R, et al. (2014) A common solution to group 2 influenza virus neutralization. *Proc. Natl. Acad. Sci. USA* 111:445–450.
- Chai N, et al. (2016) Two escape mechanisms of influenza A virus to a broadly neutralizing stalk-binding antibody. *PLoS Pathogens* 12(6):e1005702.
- Yamayoshi S, et al. (2017) Human protective monoclonal antibodies against the HA stem of group 2 HAs derived from an H3N2 virus-infected human. *Journal of Infection*.
- Okuno Y, Isegawa Y, Sasao F, Ueda S (1993) A common neutralizing epitope conserved between the hemagglutinins of influenza A virus H1 and H2 strains. *Journal of Virology* 67(5):2552–2558.
- Doud MB, Lee JM, Bloom JD (2017) Quantifying the ease of viral escape from broad and narrow antibodies to influenza hemagglutinin. *bioRxiv* p. 210468.
- Anderson CS, et al. (2017) Natural and directed antigenic drift of the H1 influenza virus hemagglutinin stalk domain. *Scientific Reports* 7(1):14614.
- Popova L, et al. (2012) Immunodominance of antigenic site B over site A of hemagglutinin of recent H3N2 influenza viruses. *PLoS One* 7(7):e41895.
- Wilson I, Skehel J, Wiley D (1981) Structure of the haemagglutinin membrane glycoprotein of influenza virus at 3 Å resolution. *Nature* 289:366–373.
- Lee PS, et al. (2014) Receptor mimicry by antibody F045-092 facilitates universal binding to the H3 subtype of influenza virus. *Nat Commun* 5:3614.
- Gamblin S, et al. (2004) The structure and receptor binding properties of the 1918 influenza hemagglutinin. *Science* 303(5665):1838–1842.

63. Lässig M, Mustonen V, Walczak AM (2017) Predicting evolution. *Nature Ecology & Evolution* 1:0077.
64. Neher RA, Bedford T (2015) nextflu: real-time tracking of seasonal influenza virus evolution in humans. *Bioinformatics* 31(21):3546–3548.
65. Wu N, et al. (2017) A structural explanation for the low effectiveness of the seasonal influenza H3N2 vaccine. *PLoS Pathogens* 13(10):e1006682.
66. McWhite C, Meyer A, Wilke C (2016) Sequence amplification via cell passaging creates spurious signals of positive adaptation in influenza virus H3N2 hemagglutinin. *Virus Evolution* 2:eve026.
67. Skowronski D, et al. (2016) Mutations acquired during cell culture isolation may affect antigenic characterisation of influenza A(H3N2) clade 3C.2a viruses. *Euro Surveill.* 21:30112.
68. Gong Li, Suchard MA, Bloom JD (2013) Stability-mediated epistasis constrains the evolution of an influenza protein. *eLife* 2:e00631.
69. Natarajan C, et al. (2013) Epistasis among adaptive mutations in deer mouse hemoglobin. *Science* 340(6138):1324–1327.
70. Harms MJ, Thornton JW (2014) Historical contingency and its biophysical basis in glucocorticoid receptor evolution. *Nature* 512(7513):203–207.
71. Starr TN, Thornton JW (2016) Epistasis in protein evolution. *Protein Science* 25(7):1204–1218.
72. Starr TN, Picton LK, Thornton JW (2017) Alternative evolutionary histories in the sequence space of an ancient protein. *Nature* 549(7672):409–413.
73. Nakajima K, Nobusawa E, Tonegawa K, Nakajima S (2003) Restriction of amino acid change in influenza A virus H3HA: comparison of amino acid changes observed in nature and in vitro. *Journal of Virology* 77(18):10088–10098.
74. Das SR, et al. (2013) Defining influenza A virus hemagglutinin antigenic drift by sequential monoclonal antibody selection. *Cell Host & Microbe* 13(3):314–323.
75. Haddox HK, Dingens AS, Hilton SK, Overbaugh J, Bloom JD (2017) Mapping mutational effects along the evolutionary landscape of HIV envelope. *bioRxiv* p. 235630.
76. Myers JL, et al. (2013) Compensatory hemagglutinin mutations alter antigenic properties of influenza viruses. *Journal of virology* 87(20):11168–11172.
77. Wu NC, et al. (2017) Diversity of functionally permissive sequences in the receptor-binding site of influenza hemagglutinin. *Cell Host & Microbe* 21(6):742–753.
78. Gong Li, Bloom JD (2014) Epistemically interacting substitutions are enriched during adaptive protein evolution. *PLoS genetics* 10(5):e1004328.
79. Doud MB, Ashenberg O, Bloom JD (2015) Site-specific amino acid preferences are mostly conserved in two closely related protein homologs. *Mol. Biol. Evol.* 32:2944–2960.
80. Nobusawa E, et al. (1991) Comparison of complete amino acid sequences and receptor-binding properties among 13 serotypes of hemagglutinins of influenza A viruses. *Virology* 182(2):475–485.
81. Hai R, et al. (2012) Influenza viruses expressing chimeric hemagglutinins: globular head and stalk domains derived from different subtypes. *Journal of Virology* 86(10):5774–5781.
82. Mallajosyula VV, et al. (2014) Influenza hemagglutinin stem-fragment immunogen elicits broadly neutralizing antibodies and confers heterologous protection. *Proc. Natl. Acad. Sci. USA* 111(25):E2514–E2523.
83. Ha Y, Stevens DJ, Skehel JJ, Wiley DC (2002) H5 avian and H9 swine influenza virus haemagglutinin structures: possible origin of influenza subtypes. *The EMBO Journal* 21(5):865–875.
84. Russell R, et al. (2004) H1 and H7 influenza haemagglutinin structures extend a structural classification of haemagglutinin subtypes. *Virology* 325(2):287–296.
85. Daniels R, et al. (1985) Fusion mutants of the influenza virus hemagglutinin glycoprotein. *Cell* 40(2):431–439.
86. Sun X, Longping VT, Ferguson AD, Whittaker GR (2010) Modifications to the hemagglutinin cleavage site control the virulence of a neurotropic H1N1 influenza virus. *Journal of virology* 84(17):8683–8690.
87. Lee HK, et al. (2013) Comparison of mutation patterns in full-genome A/H3N2 influenza sequences obtained directly from clinical samples and the same samples after a single MDCK passage. *PLoS One* 8(11):e79252.
88. Memoli MJ, et al. (2009) Recent human influenza A/H3N2 virus evolution driven by novel selection factors in addition to antigenic drift. *The Journal of Infectious Diseases* 200(8):1232–1241.
89. Raghwan J, Thompson RN, Koelle K (2017) Selection on non-antigenic gene segments of seasonal influenza A virus and its impact on adaptive evolution. *Virus Evolution* 3(2).
90. Doud MB, Hensley SE, Bloom JD (2017) Complete mapping of viral escape from neutralizing antibodies. *PLoS Pathogens* 13(3):e1006271.
91. Dingens AS, Haddox HK, Overbaugh J, Bloom JD (2017) Comprehensive mapping of HIV-1 escape from a broadly neutralizing antibody. *Cell Host & Microbe* 21:777–787.
92. Hoffmann E, Neumann G, Kawaoka Y, Hobom G, Webster RG (2000) A DNA transfection system for generation of influenza A virus from eight plasmids. *Proc. Natl. Acad. Sci. USA* 97:6108–6113.
93. Ashenberg O, Padmakumar J, Doud MB, Bloom JD (2017) Deep mutational scanning identifies sites in influenza nucleoprotein that affect viral inhibition by MxA. *PLoS pathogens* 13(3):e1006288.
94. Bloom JD (2014) An experimentally determined evolutionary model dramatically improves phylogenetic fit. *Molecular Biology and Evolution* 31:1956–1978.
95. Bao Y, et al. (2008) The Influenza Virus Resource at the National Center for Biotechnology Information. *J. Virol.* 82:596–601.
96. Katoh K, Standley DM (2013) MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Molecular Biology and Evolution* 30(4):772–780.
97. Stamatakis A (2006) RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. *Bioinformatics* 22(21):2688–2690.

**Supporting Information (SI).** The main text of the paper must stand on its own without the SI. Refer to SI in the manuscript at an appropriate point in the text. Number supporting figures and tables starting with S1, S2, etc. Authors are limited to no more than 10 SI files, not including movie files. Authors who place detailed materials and methods in SI must provide sufficient detail in the main text methods to enable a reader to follow the logic of the procedures and results and also must reference the online methods. If a paper is fundamentally a study of a new method or technique, then the methods must be described completely in the main text. Because PNAS edits SI and composes it into a single PDF, authors must provide the following file formats only.

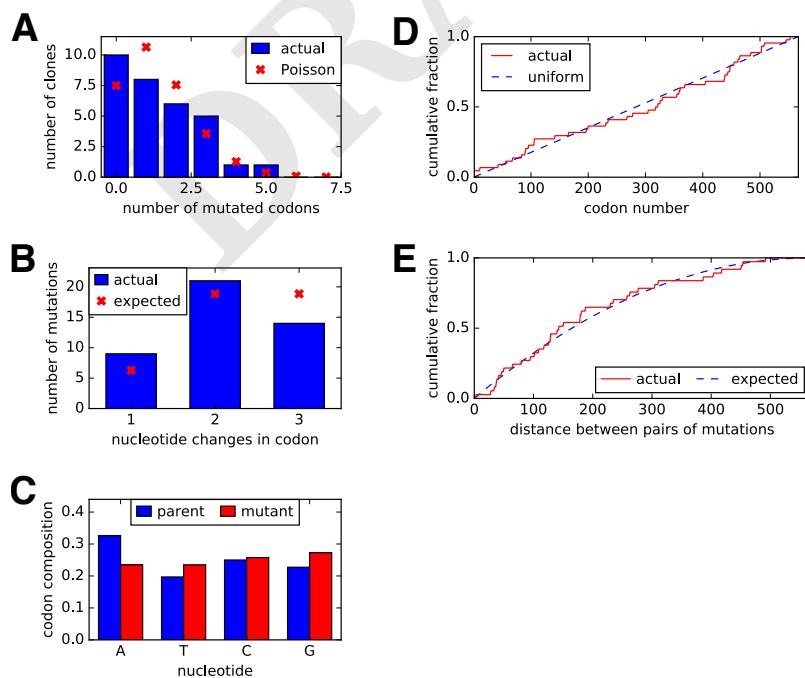
**SI Text.** Supply Word, RTF, or LaTeX files (LaTeX files must be accompanied by a PDF with the same file name for visual reference).

**SI Figures.**

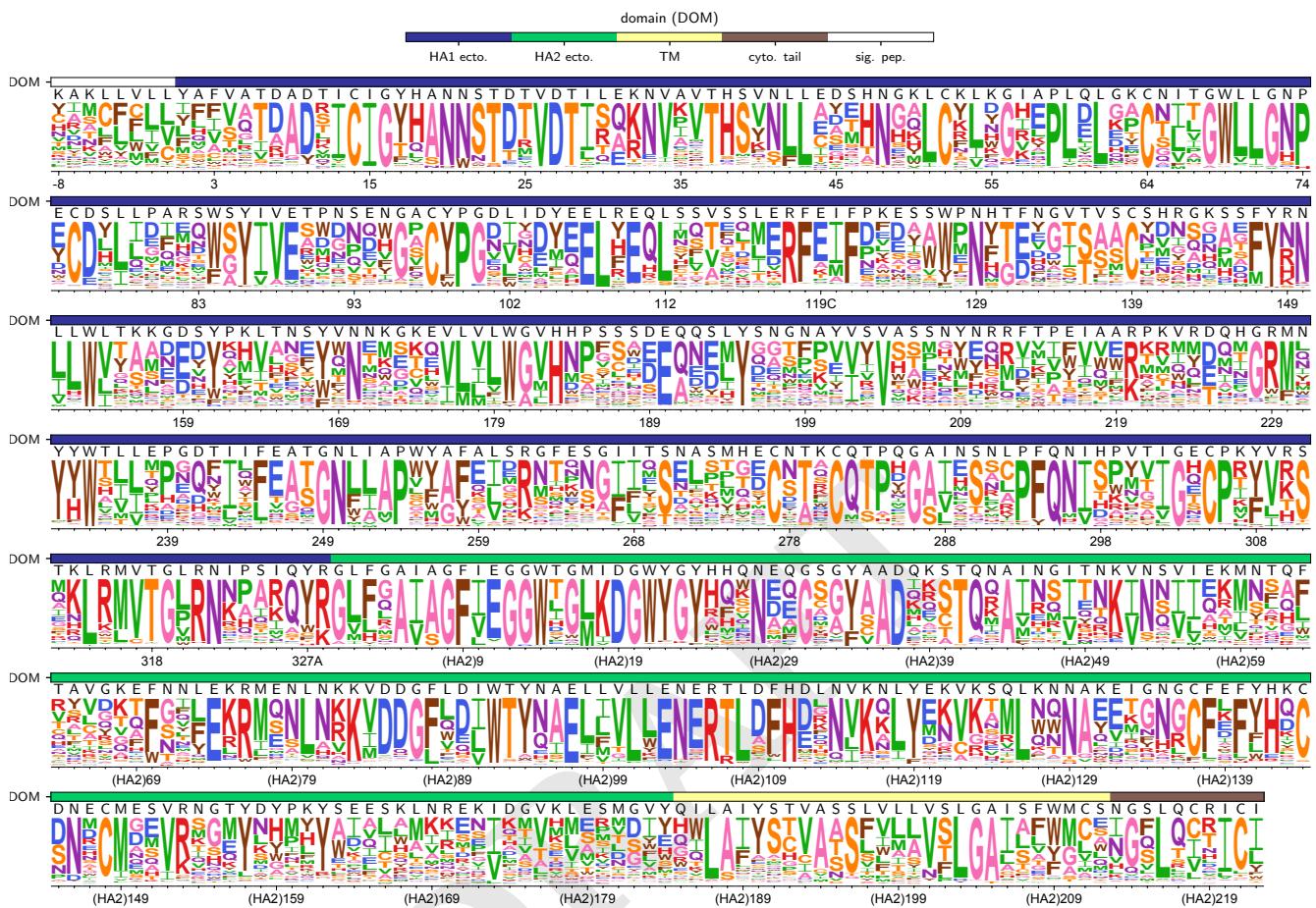
DRAFT



**Figure 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 titrated 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.



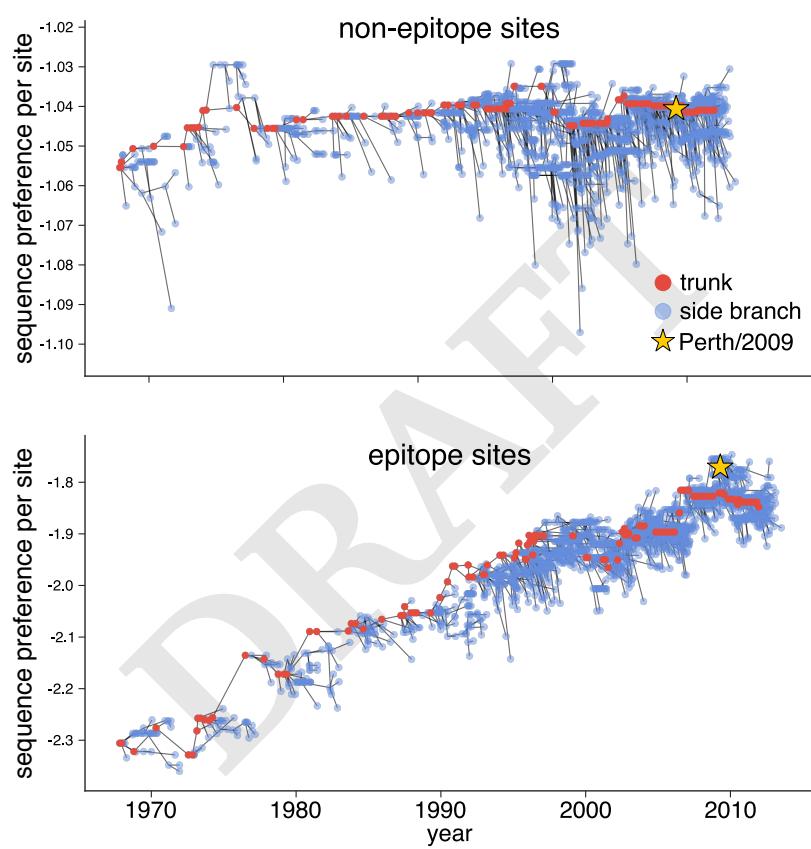
**Figure S2. Sanger sequencing of 31 randomly chosen clones from the mutant plasmid libraries.** (A) There were an average of ~1.4 codon mutations per clone across the three plasmid mutant libraries. (B) A mixture of one-, two-, and three-nucleotide mutations were present, with slightly fewer triple-nucleotide 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.



**Figure S3.** The site-specific amino-acid preferences of the WSN/1933 H1 HA. The amino-acid preferences from (30) after taking the average of the experimental replicates and re-scaling (40) 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 sites are in H3 numbering. The overlays show the same information as in Figure 3 (domain and wildtype amino acid).

[add figure]

**Figure S4.** Our experiments show a selective advantage for the trunk even if we limit the analysis to sequences that occur after the split of Perth/2009 off the trunk. (A) Phylogenetic tree showing all nodes that branch from the trunk before or contemporaneously with Perth/2009 (orange), all nodes that occur after the branch to Perth/2009 for which we can resolve the trunk (red) or side branches (blue), and all nodes for which it is not yet clear which sequences will be on the trunk or side branches. In this figure, we restrict the analysis to the nodes in red or blue. (B) Average per-site sequence preference (similar to Figure 5C) for all trunk and side-branch nodes that occur after the branch to Perth/2009. (C) The sequence preferences of nodes on the trunk generally exceed those of nodes on the side branches along the post-Perth/2009 portion of the tree. Note that because there are relatively few post-Perth/2009 sequences for which we can currently distinguish trunk from side branch, the number of points in this plot is small, precluding meaningful statistical analyses of the type performed in Figure 5B.



**Figure S5.** The per-site sequence preference at epitope and non-epitope sites. The per-site sequence preferences shown in Figure 6B, but on separate y-axes.

DRAFT

**Dataset S1.** Genbank file giving the full sequence of the bidirectional reverse-genetics plasmid pHW-Perth2009-HA-G78D-T212I, which encodes the wildtype HA sequence used in this study.

**Dataset S2.** Excel file providing the primers used to generate the barcoded subamplicons for Perth/2009 HA deep sequencing.

**Dataset S3.** Excel file giving the amino-acid preferences in sequential 1, 2, ... numbering of the Perth/2009 HA. The unscaled preferences 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 re-scaled amino-acid preferences in sequential numbering and in H3 numbering as shown in Figure 3. There are also tabs that give the conversion from sequential to H3 numbering, and the epitope sites in H3 numbering as defined by Wolf et al and used in this study. Each tab can simply be exported to CSV for computational analyses.