

Site-Specific Amino-Acid Preferences are Mostly Conserved in Two Closely Related Protein Homologs

Michael B. Doud^{1,2,3,✉}, Orr Ashenberg^{1,✉}, Jesse D. Bloom^{1,2,*}

1 Division of Basic Sciences and Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA

2 Department of Genome Sciences, University of Washington, Seattle, WA

3 Medical Scientist Training Program, University of Washington School of Medicine, Seattle, WA

✉These authors contributed equally to this work.

*jbloom@fredhutch.org

Abstract

Evolution drives changes in a protein's sequence over time. The extent to which these changes in sequence affect the underlying preferences for each amino acid at each site is an important question with implications for comparative sequence-analysis methods such as molecular phylogenetics. To quantify the extent that site-specific amino-acid preferences change during evolution, we performed deep mutational scanning on two homologs of human influenza nucleoprotein with 94% amino-acid identity. We found that only a small fraction of sites (14 out of 497) exhibited changes in their amino-acid preferences that exceeded the noise in our experiments. Given the limited change in amino-acid preferences between these close homologs, we tested whether our measurements could be used to build site-specific substitution models that describe the evolution of nucleoproteins from more diverse influenza viruses. We found that site-specific evolutionary models informed by our experiments greatly outperformed non-site-specific alternatives in fitting the phylogenies of nucleoproteins from human, swine, equine, and avian influenza. Combining the experimental data from both nucleoprotein homologs improved phylogenetic fit, in part because measurements in multiple genetic contexts better captured the evolutionary average of the amino-acid preferences for sites with changing preferences. Overall, our results show that site-specific amino-acid preferences are sufficiently conserved during evolution that measuring mutational effects in one protein provides information that can improve quantitative evolutionary modeling of nearby homologs.

Author Summary

Every site in a protein has inherent preferences for particular amino-acids, meaning that some mutations are better tolerated than others. These preferences result from constraints on protein structure and function, and it is hotly debated how much these preferences change as the protein's sequence evolves. Only if these preferences are substantially conserved during evolution can we measure them in one protein and extrapolate the results to related proteins. To investigate this issue, we measured the site-specific amino-acid preferences of a pair of homologous proteins from two human influenza strains separated by 30 years of evolution. We found that the amino-acid preferences are conserved for most sites in the homologs. We then used these amino-acid preferences to build site-specific models describing the evolution of homologs from more diverse influenza viruses. These models outperformed traditional non-site-specific models, suggesting that preferences are likely conserved across even longer periods of evolution than those examined in our experiments. Our results show that amino-acid preferences are mostly conserved in closely related proteins, and therefore suggest that site-independent but site-specific evolutionary models are likely to represent substantial improvements over non-site-specific ones.

Introduction

Since the first comparative analyses of homologous proteins by Zuckerkandl and Pauling fifty years ago [1], it has been obvious that different sites in proteins evolve under different constraints, with some sites substituting to a wide range of amino acids, while others are constrained to one or a few identities. Zuckerkandl and Pauling proposed [1], and decades of subsequent work have confirmed [2,3], that these constraints arise from the highly cooperative interactions among sites that shape important protein properties such as stability, folding kinetics, and biochemical function.

The complexity and among-sites cooperativity of these evolutionary constraints mean that a mutation at a single site can in principle affect the amino-acid preferences of any other site – and numerous experiments have demonstrated examples of such epistasis among sites [4–10]. However, experiments have also shown that despite such epistasis, the amino-acid preferences of many sites are similar across homologs [11–13]. For instance, protein structures themselves are highly conserved during evolution [14,15], and sites in specific structural contexts often have strong propensities for certain amino acids [16–18]. Furthermore, many of the most successful methods for identifying distant homologs (e.g. PSI-BLAST) utilize site-specific scoring models [19,20], implying that amino-acid preferences are often conserved even among homologs without strong sequence identity.

A half-century of work has therefore made it abundantly clear that site-specific amino-acid preferences can in principle shift arbitrarily during evolution, but nonetheless in practice remain somewhat conserved among homologs. The important remaining question is the *extent* to which site-specific amino-acid preferences are conserved versus shifted. This question is especially important for the development of quantitative evolutionary models for tasks such as phylogenetic inference. Initially, phylogenetic models unrealistically assumed that sites within proteins evolved both independently and under identical constraints. But more recent models have relaxed the second assumption that sites evolve identically. At first, this relaxation only allowed sites to evolve at different rates [21]. But newer models also accommodate variation in the amino-acid preferences among sites, either by inferring these preferences from naturally occurring sequences [22–25] or leveraging their direct measurement by high-throughput experiments [26,27]. Because these models retain the assumption of independence among sites, they will outperform traditional non-site-specific models only if site-specific amino-acid preferences are substantially conserved among homologs.

Here we perform the first experimental quantification of the conservation of the amino-acid preferences at all sites in two homologous proteins. We do this by using deep mutational scanning [28,29] to comprehensively measure the effects of all mutations to two homologs of influenza nucleoprotein (NP) with 94% sequence identity. We find that the amino-acid preferences are substantially conserved at most sites in the homologs. We then test whether the experimentally measured site-specific amino-acid preferences can inform site-specific phylogenetic substitution models that describe the evolution of more diverged NP homologs. We find that the experimentally informed site-specific substitution models exhibit improved fit to NP phylogenies containing diverged sequences from human, swine, equine, and avian influenza lineages. Overall, our work shows that site-specific amino-acid preferences are sufficiently conserved that measurements on one homolog can be used to improve the quantitative evolutionary modeling of moderately closely related homologs.

Results

Comparison of amino-acid preferences between two homologs

Deep mutational scanning of two influenza NP homologs Our studies focused on NP from influenza A virus. NP performs several conserved functions that are essential for the viral life cycle, including encapsidation of viral RNA into ribonucleoprotein complexes for transcription, viral-genome replication, and viral-genome trafficking [30]. NP’s structure is highly conserved in all characterized influenza strains [31,32]. Our studies compared the site-specific amino-acid preferences of NP homologs from two human influenza strains, PR/1934 (H1N1) and Aichi/1968 (H3N2) (Fig. 1). These NPs have diverged by over 30 years of evolution, and differ at 30 of their 498 residues (94% protein sequence identity).

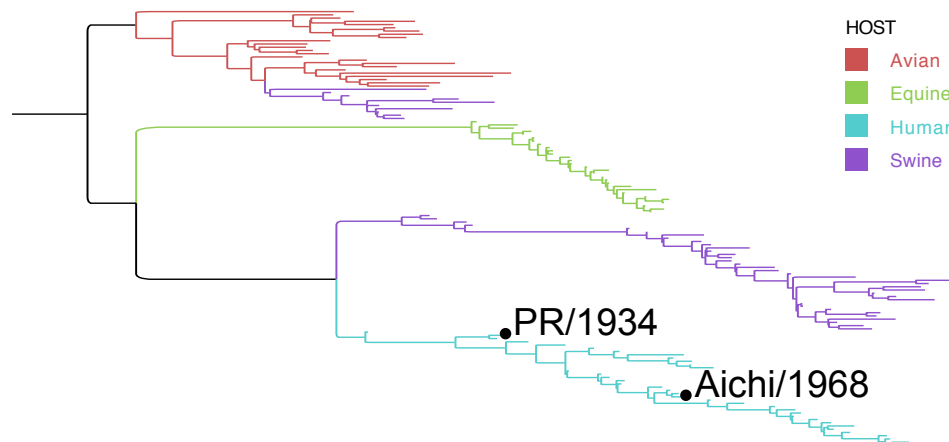


Figure 1. Phylogenetic tree of influenza NPs. The two homologs used in this work are labeled on the human influenza lineage. A diverse set of sequences was collected by sampling across years and hosts, and a maximum-likelihood tree was inferred using CodonPhyML [33] with the Goldman-Yang codon substitution model [34]. The tree was rooted using the avian clade as an outgroup.

We used our previously described approach for deep mutational scanning of influenza genes [26, 35] to measure the site-specific amino acid preferences of the PR/1934 and Aichi/1968 NPs. Briefly, this approach involved using a PCR-based technique to create mutant libraries of plasmids encoding NP genes with random codon mutations, using reverse genetics to incorporate these mutant genes into influenza viruses, and then passaging these viruses at low multiplicity of infection to select for viruses carrying functional NP variants. Deep sequencing was used to count the occurrences of each mutation before and after selection, and the amino-acid preferences for each site were inferred from these counts using `dms.tools` [36] (S1 Text, S2 Text, S2 Fig, S3 Fig). Our mutagenesis randomized 497 of the 498 codons in NP (the N-terminal methionine was not mutagenized), and so our libraries sampled all $497 \times 19 = 9,443$ amino-acid mutations at these sites. Our mutagenesis introduced an average of about two codon mutations per gene, with the number of mutations per gene following a roughly Poisson distribution (S1 Fig), and so the effect of each mutation was assayed both alone and in the background of variants that contained one or more additional mutations.

Because deep mutational scanning is subject to substantial experimental noise, we performed several full biological replicates for each NP homolog, beginning with independent creation of the plasmid mutant library. In the current work, we performed three replicates of deep mutational scanning on the PR/1934 NP and two replicates on the Aichi/1968 NP. In a previous study [26] we performed eight replicates of deep mutational scanning on Aichi/1968 NP. We will refer to these previous replicates of the Aichi/1968 NP deep mutational scanning as the *previous study*, and the two new replicates as the *current study*.

Amino-acid preferences are well correlated between homologs For each homolog we averaged the site-specific amino-acid preferences across all replicates. The mean preferences for the two NP homologs have a Pearson’s correlation coefficient of 0.78 (Fig. 2A). In comparison, the correlation between the preferences measured in the *previous study* and *current study* on the Aichi/1968 homolog is 0.83 (Fig. 2B). Therefore, the amino-acid preferences correlate nearly as well between the two homologs as they do between different experiments on the same homolog. As expected, there is no correlation between the preferences of the PR/1934 NP and a non-homologous protein (hemagglutinin, HA) for which we have previously measured the site-specific amino-acid preferences using the same approach as in this work [35] (Fig. 2C).

We also asked if the site-specific amino-acid preferences from each replicate showed the same pattern of correlation between homologs that we observed when comparing mean preferences. We again found that

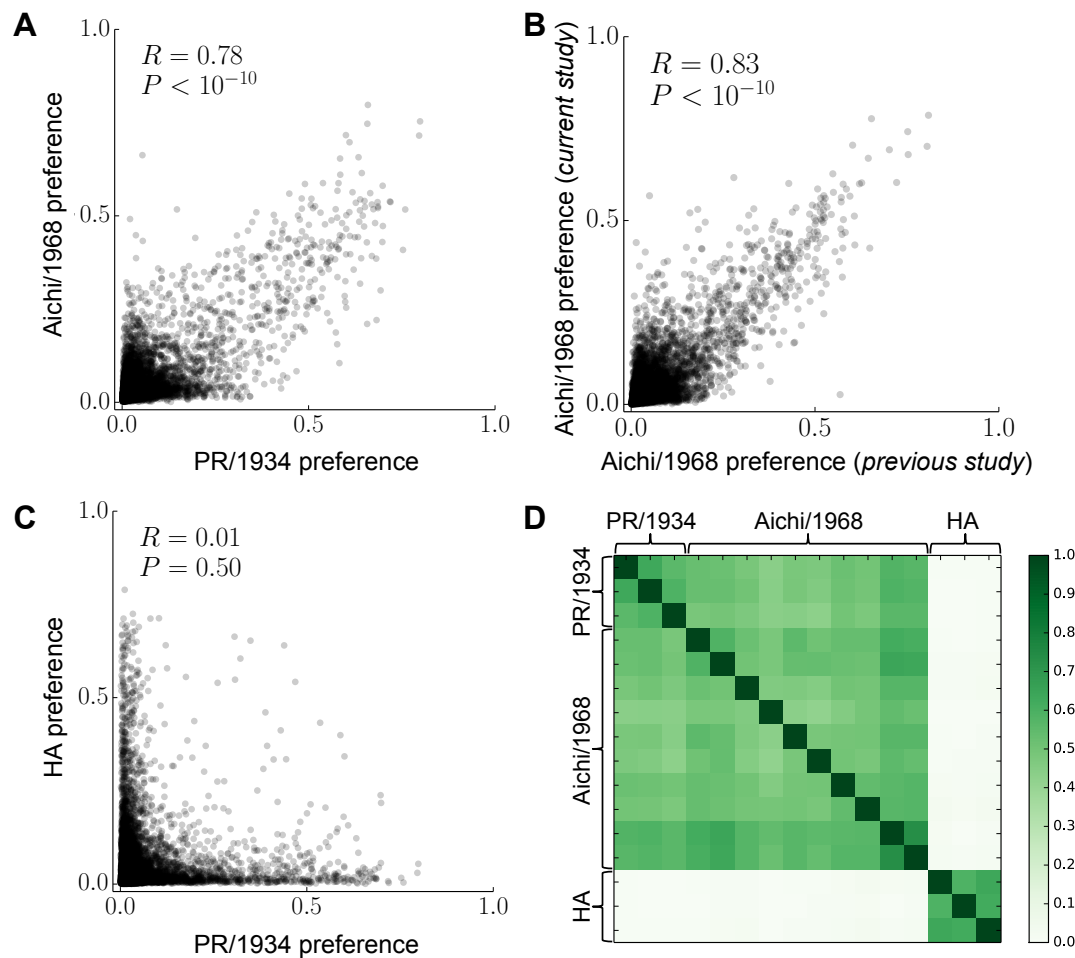


Figure 2. Site-specific amino-acid preferences correlate nearly as well between NP homologs as between replicate measurements on the same homolog. (A), (B) The correlation between the mean of the preferences taken over all replicates on each NP homolog is nearly as large as that between the preferences measured in the *current study* and *previous study* on the Aichi/1968 NP. (C) However, there is no correlation between the preferences measured for NP and the non-homologous protein HA. Each data point in (A)-(C) is the preference for one of the twenty amino acids at one of the 497 sites in NP. R is the Pearson correlation coefficient. (D) The Pearson correlations between the preferences measured in all pairs of individual replicates.

correlation coefficients are just as high between NP homologs as they are between replicate measurements on the same homolog, and that there is no correlation between the preferences for NP and the non-homologous protein HA (Fig. 2D). Overall, these results indicate that at the vast majority of sites, any differences in the amino-acid preferences between NP homologs are smaller than the noise in our experimental measurements.

Shifts in amino-acid preferences are small for most sites The previous section shows that any widespread shifts in site-specific amino-acid preferences are smaller than the noise in our experiments. However, it remains possible that a small subset of sites show substantial shifts in their amino-acid preferences that are masked by examining all sites together. We therefore performed an analysis to identify specific sites that show shifted amino-acid preferences between homologs.

This analysis needed to account for the fact that experimental noise induced variation in the preferences measured in each replicate. Fig. 3 shows replicate measurements for both homologs at several sites in NP. At

many sites, such as 298, all replicate measurements yielded highly reproducible amino-acid preferences both between and within homologs. At many other sites, such as 309, replicate measurements were quite variable both between and within homologs, probably due to fairly weak selection at that site. Finally, at a few sites, such as 470, replicate measurements were highly reproducible within each homolog but clearly differed between homologs. We therefore developed a quantitative measure of the shift in preferences between homologs that accounts for this site-specific experimental noise.





Site	Biological Replicates	$RMSD_{\text{between}}$	$RMSD_{\text{within}}$	$RMSD_{\text{corrected}}$
309		0.58	0.58	0.00
298		0.32	0.30	0.02
408		0.69	0.44	0.25
470		0.71	0.26	0.45
	PR/1934 Aichi/1968			

Figure 3. Replicate measurements quantify the shift in amino-acid preferences between homologs after correcting for experimental noise. The amino-acid preferences measured in multiple replicates of deep mutational scanning of both homologs are shown for selected sites. $RMSD_{\text{between}}$ (the average difference between the two homologs) and $RMSD_{\text{within}}$ (the average variation within replicates of each homolog) are shown to the right. $RMSD_{\text{corrected}}$ is calculated by subtracting $RMSD_{\text{within}}$ from $RMSD_{\text{between}}$.

We used the Jensen-Shannon distance metric (the square root of the Jensen-Shannon divergence) to quantify the distance between pairs of replicate measurements at a site. This distance ranges from 0 (identical amino-acid preferences) to 1 (completely different amino-acid preferences). To quantify experimental noise at a site, we calculated the root-mean-square of the Jensen-Shannon distance for all pairwise comparisons among replicate measurements on the same homolog, and termed this quantity $RMSD_{\text{within}}$. Sites with large $RMSD_{\text{within}}$ have high experimental noise (Fig. 3). We defined an analogous statistic, $RMSD_{\text{between}}$, to quantify the distance in preferences between homologs by calculating the root-mean-square of the Jensen-Shannon distance for all pairwise comparisons between replicates of PR/1934 and replicates of Aichi/1968.

The fact that we had data from two independent sets of experiments on the Aichi/1968 NP (the *current study* and *previous study*) enabled us to perform a control analysis by calculating $RMSD_{\text{between}}$ and $RMSD_{\text{within}}$ for the replicates from these two experiments. As an additional control to gauge the extent of amino-acid preference differences between non-homologous proteins, we also calculated $RMSD_{\text{between}}$ and $RMSD_{\text{within}}$ for our experiments on PR/1934 NP and our previous study on HA.

The relationship between $RMSD_{\text{between}}$ and $RMSD_{\text{within}}$ for all sites is shown in Fig. 4A-C. When comparing two independent experiments on the same NP (Fig. 4A) or comparing experiments on two homologs of NP (Fig. 4B), the relationship between $RMSD_{\text{between}}$ and $RMSD_{\text{within}}$ is approximately linear, indicating that the difference in amino-acid preferences between homologs at a given site is usually comparable to the experimental noise. Deviations from this linear relationship are more frequent in the comparison between PR/1934 and Aichi/1968 (Fig. 4B) than in the comparison between the two studies of Aichi/1968 (Fig. 4A). These deviations mostly arise from sites that have larger $RMSD_{\text{between}}$ than $RMSD_{\text{within}}$, indicating that these sites have differences in their amino-acid preferences between homologs that exceed the experimental noise. These results comparing NP homologs are in stark contrast with the

$RMSD_{between}$ and $RMSD_{within}$ calculated when comparing NP to HA (Fig. 4C), where the difference in preferences between proteins is almost always substantially greater than the experimental noise.

To quantify the extent of amino-acid preference change between the two homologs in a way that corrects for the experimental noise, we defined another statistic, $RMSD_{corrected}$, by subtracting $RMSD_{within}$ from $RMSD_{between}$ (Fig. 3, S1 Table). Sites with changes in amino-acid preferences greater than the experimental noise will have $RMSD_{corrected} > 0$. However, we also expect many sites to have positive $RMSD_{corrected}$ values due to statistical noise. In order to determine the distribution of $RMSD_{corrected}$ values expected due to such statistical noise alone, we generated null distributions of $RMSD_{corrected}$ by performing simulations under a model where there are no differences in the underlying preferences between the two homologs, but varying levels of noise for each replicate.

In order to generate a null distribution of the $RMSD_{corrected}$ values expected under the null hypothesis that there is no difference in the preferences between homologs, we simulated amino-acid preferences at each site by drawing from a Dirichlet distribution, which is well-suited for this purpose because its support is a normalized vector of values, in this case corresponding to the vector of amino-acid preferences at a site. Our null hypothesis is that the amino-acid preferences are the same for both homologs, so we performed simulations assuming that the *true* vector of amino-acid preferences at a site is equal to the average of our experimental measurements for both homologs. We simulated the amino-acid preferences for each replicate by drawing from a Dirichlet distribution centered on this vector of assumed true preferences. The extent to which any given sample drawn from this Dirichlet distribution differs from the true vector can be tuned with a single scaling parameter (the concentration parameter). We identified a value for the concentration parameter for each experiment (Aichi/1968 *current study*, Aichi/1968 *previous study*, and PR/1934) that resulted in correlation coefficients between replicates that matched those in the actual experiment. We performed 1000 replicate simulations and combined the calculated $RMSD_{corrected}$ values from all simulations to build the null distribution.

The distribution of $RMSD_{corrected}$ obtained experimentally closely overlaps the simulated distribution of $RMSD_{corrected}$ when comparing the two independent Aichi/1968 experiments (Fig. 4D). The close overlap between the simulated and experimental $RMSD_{corrected}$ is consistent with the hypothesis that the true amino-acid preferences are the same in both experiments on the Aichi/1968 NP, and that the observed differences are attributable to noise. In contrast, when comparing PR/1934 to Aichi/1968, some $RMSD_{corrected}$ values are shifted in the positive direction (Fig. 4E), indicating larger differences in preferences than can be explained by experimental noise alone. However, this shift in preferences is subtle compared to the shift observed when comparing PR/1934 to the non-homologous HA (Fig. 4F), where nearly all sites have differences in preferences that are inconsistent with the null hypothesis that the only differences are due to experimental noise.

As a second null model, we also used exact randomization testing to shuffle the replicate datasets between the two homologs. For each of the three PR/1934 and ten Aichi/1968 replicate datasets, we randomized which dataset was placed with which homolog. For every possible shuffling of datasets, we computed $RMSD_{corrected}$ at every site and combined the results across all shufflings. The distributions of $RMSD_{corrected}$ obtained by randomizing the datasets are in S5 Fig and are similar to those obtained from simulated data.

Features of the few sites with changed amino-acid preferences Using either of the two null distributions of $RMSD_{corrected}$, we were able to identify specific sites with $RMSD_{corrected}$ values significantly larger than expected due to experimental noise alone (S1 Table). These are sites for which we can reject the null hypothesis that there is no difference in amino-acid preferences. The set of significant sites identified using the simulated data contained all three of the sites identified using randomization, plus an additional eleven. Among these fourteen sites with significant changes in amino-acid preferences, five are sites with different wild-type amino-acid identity between PR/1934 and Aichi/1968 (Fig. 5A). Therefore, of sites with variable amino-acid identity between the two homologs, 17% have detectable shifts in preference. In contrast, 2% of conserved sites have detectable shifts in preferences.

Having identified evolutionarily variable sites as enriched for shifts in amino-acid preferences, we next looked at sites with other special structural or functional properties. One group of functionally important

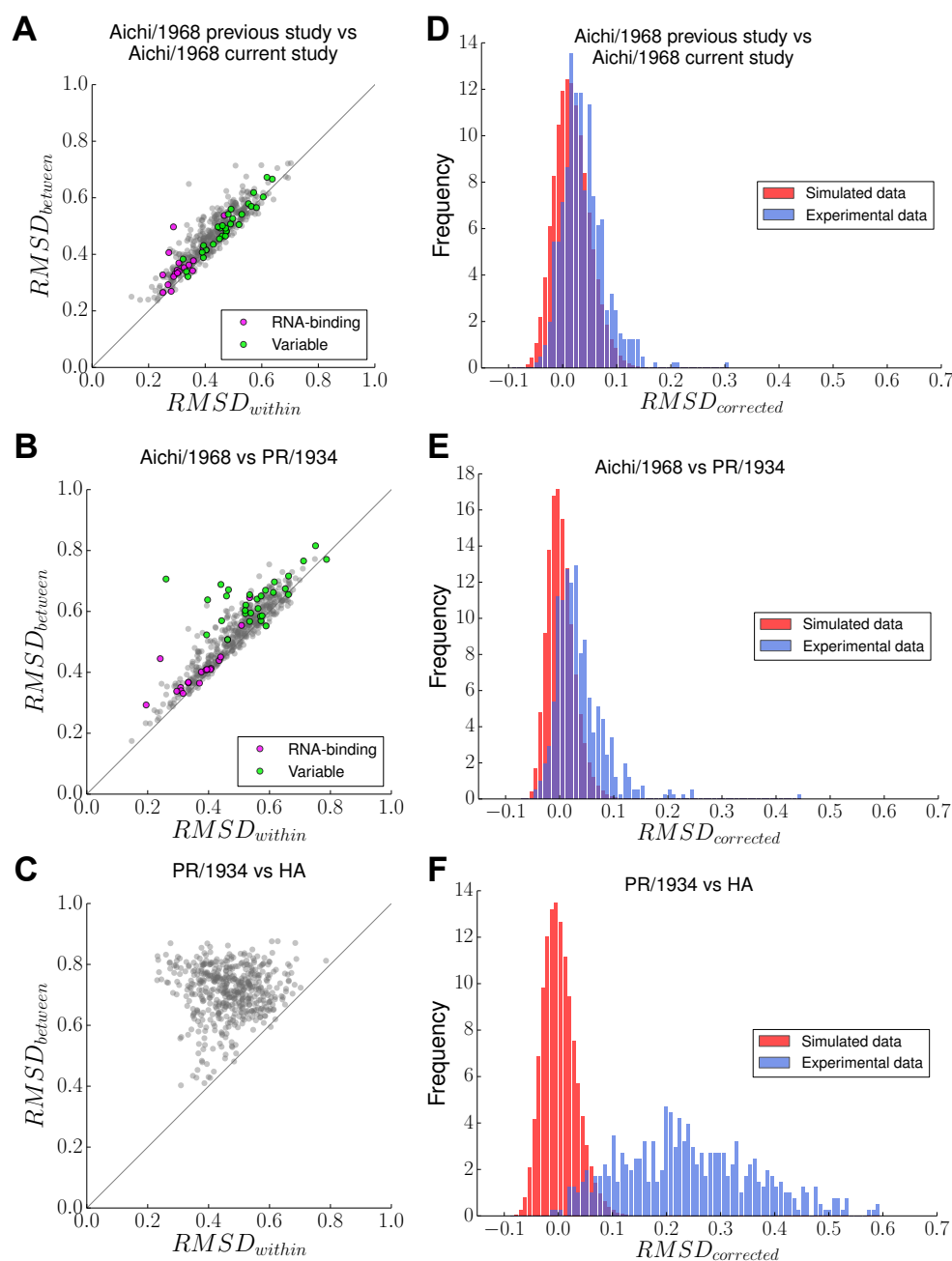


Figure 4. Identification of sites with shifts in amino-acid preferences. (A)-(C) Each plot shows statistics calculated for a comparison between two groups of replicate experiments. Each point represents a site in NP. $RMSD_{within}$ quantifies the average difference in amino-acid preferences within each of the two groups (experimental noise), and $RMSD_{between}$ quantifies the average difference in preferences between the two groups. Sites in the RNA-binding groove are in purple; sites that have different wild-type identities in PR/1934 and Aichi/1968 are in green. (D)-(F) The actual distribution of $RMSD_{corrected}$ values is shown in blue, and the distribution of $RMSD_{corrected}$ from data simulated under a null hypothesis of no difference in preferences is shown in red. Comparisons are made between the two studies on Aichi/1968 (A, D), between Aichi/1968 and PR/1934 (B, E), and between PR/1934 and the non-homologous HA (C, F).

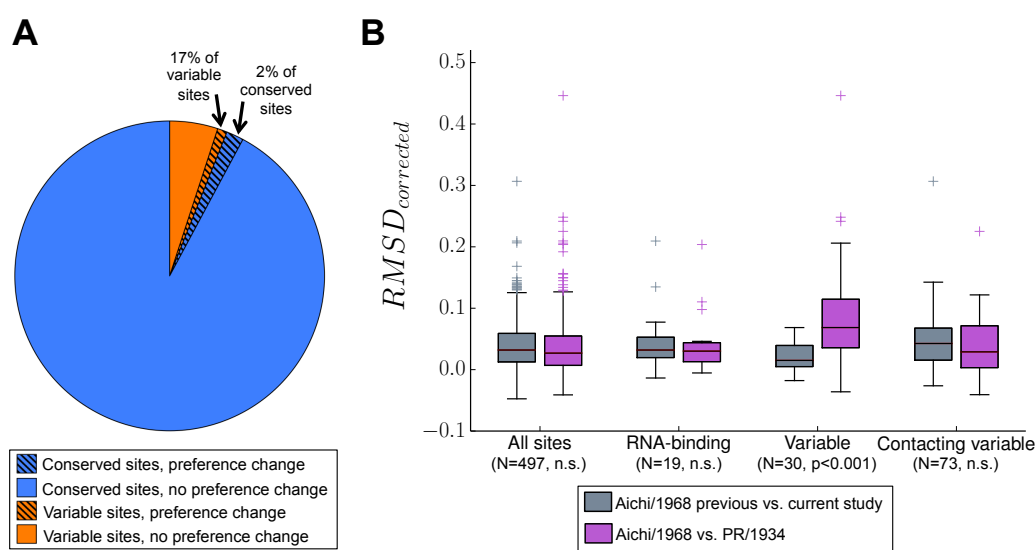


Figure 5. Evolutionarily variable sites are enriched for changes in amino-acid preference. (A) Sites with shifts in amino-acid preferences were identified by $RMSD_{corrected}$ values greater than expected under a null model assuming no difference between homologs ($p \leq 0.05$ after Bonferroni correction for multiple hypothesis testing). *Variable* sites have different wild-type residues in the two NP homologs. (B) The distributions of $RMSD_{corrected}$ for various groups of sites. The median is marked by a horizontal line, boxes extend from 25th to 75th percentile, and whiskers extend to data points within 1.5 times the interquartile range. Outliers are marked with crosses. *Contacting variable* sites are conserved sites with side-chain atoms within 4.5 Å of a variable side-chain atom. $RMSD_{corrected}$ distributions for each group of sites are shown for two comparisons: one comparing two independent experiments on Aichi/1968, and one comparing Aichi/1968 to PR/1934. P-values were determined using the Mann-Whitney U test and adjusted using the Bonferroni correction.

sites are those that comprise the RNA-binding groove of NP. These RNA-binding sites have low $RMSD_{within}$ (Fig. 4A and B), indicating below-average noise among replicates. RNA-binding sites also have low $RMSD_{corrected}$ (Fig. 5B, Fig. 6A). These results are consistent with the expectation that RNA-binding sites in NP are under strong and conserved functional constraint, since RNA binding is essential for viral genome packing, transcription, and replication.

We also hypothesized that sites in structural proximity to evolutionarily variable sites may experience shifts in amino-acid preferences due to changes in the surrounding biochemical environment. We identified sites contacting the evolutionarily variable sites and found that they do not exhibit any shift in $RMSD_{corrected}$ (Fig. 5B). Therefore, the underlying amino-acid preferences at these sites are not detectably affected by changes in nearby amino acids.

The fourteen sites with detectable shifts in amino-acid preferences are distributed throughout the surface of NP in the body, head, and tail loop domains (Fig. 6). Six of these fourteen sites are located in the flexible tail loop which inserts into a neighboring NP monomer during NP oligomerization. Seven of the fourteen sites are in contact with each other. The thirty evolutionarily variable sites are also mostly found on the surface of NP (Fig. 6). Although the variable sites are enriched for shifts in preferences relative to other sites, most variable sites still have low $RMSD_{corrected}$ values.

Overall, these results indicate that sites with evolutionarily variable amino-acid identity are more likely than conserved sites to exhibit detectable shifts in amino-acid preferences. However, the majority of sites

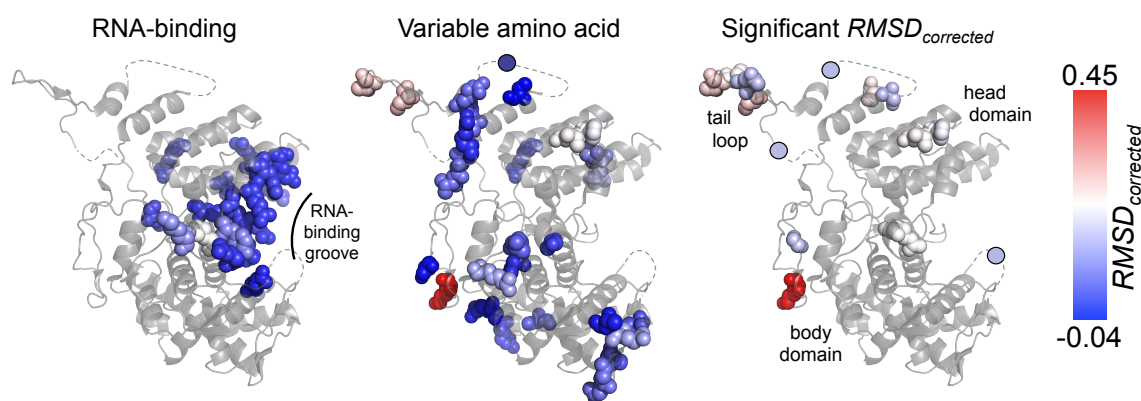


Figure 6. Magnitude of the shift in amino-acid preferences mapped on the NP structure. $RMSD_{corrected}$ values for each site are used to color space-filling models for the indicated sites in the NP crystal structure (PDB ID: 2IQH, chain C [31]). Sites are shown as circles when in regions that are not present in the crystal structure (dashed lines). Blue represents small shifts in amino-acid preferences between PR/1934 and Aichi/1968; red represents large shifts. *Variable amino acid* refers to sites where the wild-type residue differs between PR/1934 and Aichi/1968 NP.

with variable identity still exhibit no detectable shift in amino-acid preference, and overall, only 3% of sites in NP have shifts in amino-acid preferences that are sufficiently large to justify rejecting the null hypothesis that the preferences are identical between homologs.

Experimentally informed site-specific substitution models describe vast swaths of nucleoprotein evolution

We next quantitatively assessed how well our experimentally measured amino-acid preferences reflected the actual constraints on NP evolution. To do so, we used the amino-acid preferences to inform site-specific phylogenetic substitution models. We have previously shown that substitution models informed by experimentally measured site-specific amino-acid preferences greatly outperform common non-site-specific codon-substitution models [26,27,35].

In the prior work, site-specific amino-acid preferences were experimentally measured in a single sequence context. Here, we asked whether combining the preferences measured in the two different sequence contexts of Aichi/1968 and PR/1934 would more accurately describe NP sequence evolution. Any improvement could be due to two effects: First, a combined substitution model might better reflect the evolutionary average of the amino-acid preferences at sites with significant changes in preferences over time. Second, combining data from multiple experiments should reduce noise and yield more accurate site-specific amino-acid preferences.

Combining deep mutational scanning datasets from nucleoprotein homologs improves phylogenetic fit To compare the performance of different substitution models, we used a likelihood-based framework. We first built a maximum-likelihood tree for NP sequences from human influenza using CodonPhyML [33] with the Goldman-Yang (GY94) codon-substitution model [34] (Fig. 1). We fixed this tree topology and used HyPhy to optimize branch lengths and model parameters for each substitution model by maximum likelihood. The relative fits of the substitution models were evaluated using the Akaike information criterion (AIC) [37].

We tested experimentally informed substitution models derived from the Aichi/1968 and PR/1934 mutational scans either alone or in combination. The Aichi/1968 model used amino-acid preferences averaged across the *current study* and the *previous study*. To build a combined substitution model based on both NP homologs, we averaged the amino-acid preferences for the Aichi/1968 and PR/1934 homologs (Aichi/1968 + PR/1934). Each substitution model had five free parameters: four nucleotide mutation rates

and a stringency parameter β that accounts for the possibility of a different strength of selection in natural sequence evolution compared to the mutational-scanning experiments [27]. These parameters were fit by maximum-likelihood optimization, and their values are listed in Table 1.

Table 1. Combining experimental data improves phylogenetic fit to NPs from human influenza

model	ΔAIC	log likelihood	parameters (optimized + empirical)	optimized parameters
Aichi/1968 + PR/1934	0.0	-4395.8	5 (5 + 0)	$R_{A \rightarrow G} = 4.6$, $R_{A \rightarrow T} = 0.8$, $R_{C \rightarrow A} = 1.4$, $R_{C \rightarrow G} = 0.1$, $\beta = 3.0$
PR/1934	322.3	-4556.9	5 (5 + 0)	$R_{A \rightarrow G} = 4.9$, $R_{A \rightarrow T} = 0.8$, $R_{C \rightarrow A} = 1.4$, $R_{C \rightarrow G} = 0.1$, $\beta = 2.1$
Aichi/1968	485.7	-4638.6	5 (5 + 0)	$R_{A \rightarrow G} = 4.8$, $R_{A \rightarrow T} = 0.7$, $R_{C \rightarrow A} = 1.4$, $R_{C \rightarrow G} = 0.1$, $\beta = 2.4$
GY94, gamma ω , gamma rates	2582.3	-5678.9	13 (4 + 9)	$\kappa = 6.2$, ω shape = 0.1, mean $\omega = 0.1$, rate shape = 2.4

Substitution models are sorted by ΔAIC , and the corresponding log likelihoods, number of free parameters, and values of optimized parameters are shown. Log likelihoods for each model were calculated through maximum-likelihood optimization of branch lengths and model parameters given the fixed tree topology of human NPs shown with blue lines in Fig. 1. The only parameters in the experimentally informed models are the four nucleotide mutation rates and the stringency parameter β . The non-site-specific GY94 model [34] has nine empirical nucleotide equilibrium frequencies [38], and optimized parameters describing the transition-transversion ratio (κ), the gamma distribution of the nonsynonymous-synonymous ratio (ω) [39], and the gamma distribution of substitution rates [21]. In the Aichi/1968 model, the preferences from *current study* and *previous study* have been averaged.

As a comparison to the experimentally informed substitution models, we also tested the non-site-specific GY94 model. Relative to the experimentally informed substitution models, the GY94 model includes many more free parameters including nucleotide equilibrium frequencies, a transition-transversion ratio, and parameters describing gamma distributions of the nonsynonymous-synonymous ratio and substitution rate across sites [21, 39].

The Aichi/1968 and PR/1934 experimentally informed models described the human NP phylogeny far better than the non-site-specific GY94 model (Table 1). Strikingly, combining amino-acid preferences from both NP homologs (Aichi/1968 + PR/1934) resulted in a greatly improved substitution model (Table 1). For each experimentally informed model, the stringency parameter β fit with value greater than 1, consistent with the idea that selection during natural evolution is more stringent than our laboratory selection.

Experimentally informed models also describe the evolution of more diverged non-human influenza strains Given the success of the experimentally informed substitution models in describing the human NP phylogeny, we asked whether these models could be extended to more diverged NPs from non-human influenza strains. We expect these models to exhibit good fit if the NP site-specific amino-acid preferences are mostly conserved across these viral strains. We examined NPs from influenza strains from three hosts: swine, equine, and avian. The average protein-sequence identity between human NPs and swine, equine, and avian NPs was 91%, 91%, and 93% respectively.

We built a phylogenetic tree of NPs of influenza viruses from human, swine, equine, and avian hosts (Fig. 1, S4 Text, S5 Text). As previously reported, the avian sequences could be divided into western and eastern hemispheric clades, and the swine sequences consisted of the North American Classical H1N1 clade and the more recent Eurasian H1N1 clade [40]. Using this tree, we performed a phylogenetic analysis similar to that described above for human influenza NPs.

Again, the experimentally informed models greatly outperformed the non-site-specific GY94 model, and combining the Aichi/1968 and PR/1934 models resulted in a far superior model (Table 2). Since the amino-acid preferences were experimentally measured for human NP, we wanted to ensure that this superior performance was not driven solely by the human clade of the tree. We separately fit subtrees consisting only of swine, equine, or avian NP sequences (S2 Table, S3 Table, S4 Table). Each subtree showed the same trend as the full tree: the experimentally informed models were superior to the GY94 model, and combining data from the two NP homologs resulted in large improvements in likelihood. Therefore, substitution models informed by deep mutational scanning of human influenza NP homologs can be extended to the NPs of

influenza from other hosts, suggesting that the site-specific amino-acid preferences are substantially conserved among all influenza A viruses.

Table 2. Combining experimental data improves phylogenetic fit to NPs from human, swine, equine, and avian influenza

model	Δ AIC	log likelihood	parameters (optimized + empirical)	optimized parameters
Aichi/1968 + PR/1934	0.0	-17507.9	5 (5 + 0)	$R_{A \rightarrow G} = 6.0, R_{A \rightarrow T} = 1.0, R_{C \rightarrow A} = 1.4, R_{C \rightarrow G} = 0.1, \beta = 2.7$
PR/1934	700.2	-17858.0	5 (5 + 0)	$R_{A \rightarrow G} = 6.3, R_{A \rightarrow T} = 1.0, R_{C \rightarrow A} = 1.4, R_{C \rightarrow G} = 0.1, \beta = 2.1$
Aichi/1968	1030.2	-18023.0	5 (5 + 0)	$R_{A \rightarrow G} = 6.2, R_{A \rightarrow T} = 0.9, R_{C \rightarrow A} = 1.4, R_{C \rightarrow G} = 0.1, \beta = 2.3$
GY94, gamma ω , gamma rates	1784.7	-18392.2	13 (4 + 9)	$\kappa = 6.9, \omega \text{ shape} = 0.3, \text{mean } \omega = 0.1, \text{rate shape} = 3.1$

This table differs from Table 1 in that it fits the combined tree of human, swine, equine, and avian NPs in Fig. 1.

Combining data from NP homologs improves phylogenetic fit to sites with shifted preferences

The results above show that the experimentally informed substitution models improved phylogenetic fit relative to the non-site-specific model, and that combining data from two NP homologs resulted in the best model. This increased performance when combining data may come from more accurate measurement of amino-acid preferences due to more replicates, or from averaging amino-acid preferences over multiple sequence contexts. To examine these possible explanations, we analyzed which sites in NP were more accurately modeled when the Aichi/1968 and PR/1934 experimental models were combined. This analysis was performed using the full phylogenetic tree of NP sequences (Fig. 1).

While fixing the branch lengths and model parameters to their maximum-likelihood values for each model, we calculated for each site the difference in likelihoods (Δ log-likelihood) when the site was modeled using the combined Aichi/1968 + PR/1934 model compared to using the Aichi/1968 model. We binned sites into quintiles of Δ log-likelihood. Sites in the top quintile had the greatest increases in likelihood when the Aichi/1968 and PR/1934 models were combined. Overall 67% of sites in NP had increased likelihoods under the Aichi/1968 + PR/1934 model.

To determine whether these improved likelihoods came from lower noise in the combined experimental model, we used the $RMSD_{within}$ statistic. Sites with greater variance in amino-acid preferences across experimental replicates have higher $RMSD_{within}$ scores. We analyzed the distribution of the $RMSD_{within}$ scores for sites within each quintile (Fig. 7). The top and bottom quintiles did not have significantly different $RMSD_{within}$ distributions, indicating that sites prone to experimental noise contributed both positively and negatively to the tree likelihood when experimental datasets were combined. Thus, the improved modeling with the combined dataset was not chiefly due to reduced experimental noise.

Next, to determine whether the improved likelihoods were driven by sites with different preferences between the two NP homologs, we used the $RMSD_{corrected}$ statistic (Fig. 7). If the improvements under the combined model came from sites with different amino-acid preferences between Aichi/1968 and PR/1934, then we would expect that the sites with the greatest increases in likelihood would also have the greatest $RMSD_{corrected}$ values. This was indeed the case, as sites in the top quintile of log-likelihoods had the highest median $RMSD_{corrected}$. The $RMSD_{corrected}$ scores in the top quintile were significantly different from those in the lower quintiles (Mann-Whitney U with Bonferroni correction $p < 0.002$), whereas there were no significant differences in the $RMSD_{corrected}$ scores when comparing the lower quintiles. Therefore, improvements in the combined model were partly due to better describing those sites that had the largest shifts in amino-acid preferences over evolutionary time.

Discussion

Determining the extent to which site-specific amino-acid preferences shift during evolution is important for evaluating how well experimental measurements can be extrapolated across homologs, and for guiding the development of site-specific phylogenetic substitution models. We have performed the first comprehensive

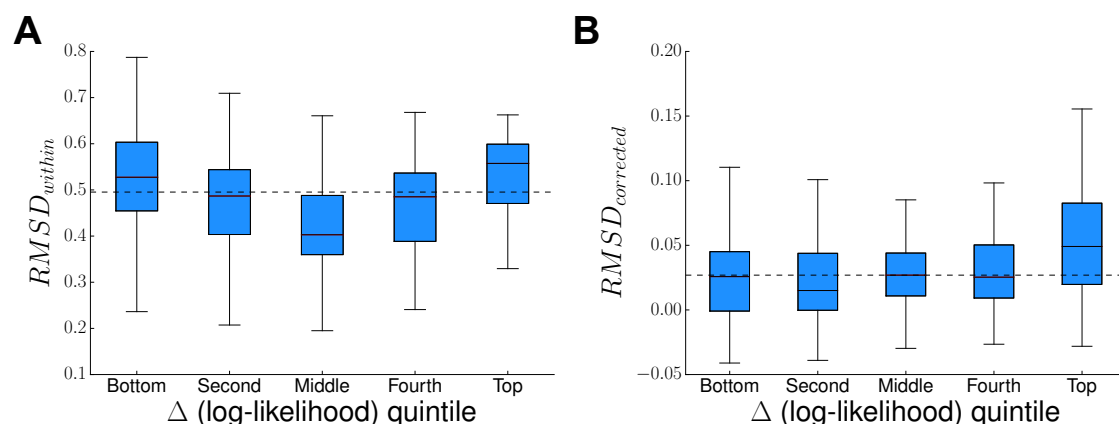


Figure 7. NP sites that are better described by combining data from both homologs have shifted amino-acid preferences. The change in per-site likelihood in going from the Aichi/1968 model to the Aichi/1968 + PR/1934 model was plotted against the per-site $RMSD_{within}$ (A) or per-site $RMSD_{corrected}$ (B). Sites were ranked by $\Delta(\log\text{-likelihood})$, divided into quintiles, and the per-site $RMSD_{within}$ or per-site $RMSD_{corrected}$ for sites in each quintile was displayed as a box and whisker plot. Outlier sites beyond the interquartile range are omitted. Quintiles are ordered left to right from least improved likelihoods to most improved likelihoods under the combined model. The median $RMSD_{within}$ or $RMSD_{corrected}$ is shown as a horizontal, dashed line. Sites with the most improved likelihoods did not have significantly higher variation in amino-acid preferences (high $RMSD_{within}$) across replicate measurements on the same homolog. However, these sites did have significantly higher differences in amino-acid preferences between Aichi/1968 and PR/1934 (high $RMSD_{corrected}$).

assessment of the conservation of site-specific amino-acid preferences by using deep mutational scanning to measure the effects of all mutations on two closely related homologs of influenza NP.

We found that for the vast majority of sites, any shift in amino-acid preferences between homologs was smaller than the noise in our experiments. Only for 14 of the 497 sites could we reject the null hypothesis that the amino-acid preferences were identical between homologs. Sites that varied in amino-acid identity between the two homologs were more likely to have a detectable shift in amino-acid preferences (5 of 30 variable sites had shifted preferences, versus only 9 of 467 non-variable ones) – but even among variable sites, there was usually no detectable shift. Admittedly, our experiments had substantial noise, so it is likely that other sites have more subtle shifts below our limit of detection. However, the fact that the preferences we measured for the two NP homologs were substantially correlated with each other but completely uncorrelated with those for the non-homologous hemagglutinin shows that the site-specific amino-acid preferences of homologs are tremendously more similar than those of unrelated proteins.

This broad conservation of site-specific amino-acid preferences does not imply an absence of epistasis during NP's evolution. Even small shifts in mutational effects can induce evolutionarily important epistasis, since natural selection is highly discerning. Indeed, we have previously demonstrated epistasis among mutations to NP [8], indicating that NP is no different than the many other proteins for which evolutionarily relevant epistasis has been identified [4–7, 9, 10]. The key result is not that epistasis is absent, but rather that its frequency and magnitude are sufficiently low that the amino-acid preferences for most sites are largely conserved between homologs.

The implications of this finding are illustrated by the second part of our study, which shows that the experimentally measured site-specific amino-acid preferences can inform phylogenetic substitution models that greatly outperform traditional non-site-specific models even for more diverged NP homologs. It is well known that the actual constraints on protein evolution involve cooperative interactions among sites [1–3], and so substitution models that treat sites either independently or identically are obviously imperfect. But computational biology must balance realism with tractability. Site-independent but site-specific substitution models are becoming feasible for real-world datasets [22–27], but approaches that relax the assumption of

independence among sites remain in their infancy [41,42]. Are amino-acid preferences sufficiently conserved for site-independent but site-specific models to represent substantial improvements over existing alternatives? Both our experimental and computational results answer this question with a resounding yes.

Why are the site-specific amino-acid preferences mostly conserved? As is the case for virtually all proteins [14,15], the structure of NP is highly conserved among homologs [31,32], and sites in specific structural contexts often have propensities for certain amino acids [16–18]. In addition, selection for protein stability is a major constraint on evolution [2,43], and experiments on NP [12] and other proteins [11,13] have shown that the effects of mutations on stability are similar among homologs. Therefore, conserved structural and stability constraints probably naturally lead to substantial conservation of site-specific amino-acid preferences. We refer the reader to an excellent recent study by Sanchez-Ruiz and coworkers for a more biophysically nuanced discussion of these issues [11].

The extent to which site-specific amino-acid preferences will be conserved among more distant homologs remains an open question. Computational simulations of the divergence of distant homologs have been used to argue that preferences shift substantially [44], but the reliability of such simulations is unclear since computational predictions of the effects of even single amino-acid mutations are only modestly accurate [45,46]. The only direct experimental data come from a study showing that the effects of a handful of mutations on stability are mostly conserved among homologs with about 50% protein-sequence identity [11]. More comprehensive determination of the relationship between sequence divergence and shifts in site-specific amino-acid preferences therefore remains an important topic for future work.

Materials and Methods

Availability of data and computer code FASTQ files can be accessed at the Sequence Read Archive (SRA Accession SRP056028). The computer code and an accompanying description sufficient to reproduce all the analysis in this work is available at <https://github.com/mbdoud/Compare-NP-Preferences>.

Deep mutational scanning of two influenza nucleoprotein homologs

We performed deep mutational scanning of influenza nucleoprotein (NP) in three biological replicates for A/PR/1934 (H1N1) and two biological replicates for A/Aichi/1968 (H3N2) (termed here as Aichi/1968 *current study*). We broadly followed the methods used for mutagenesis, viral rescue, deep sequencing, and inference of amino-acid preferences from sequence data described in [26], with the following notable changes to the protocol.

Codon mutagenesis For each replicate mutant library, we followed the mutagenesis protocol as previously described [26], but performed two rounds of mutagenesis instead of three to decrease the average number of mutations per clone. After ligation of mutagenized PCR products to the pHW2000 [47] plasmid backbone, multiple parallel transformations and platings were combined to ensure that each replicate library contained more than 10^6 unique transformants. Sanger sequencing of 30 clones from each homolog revealed that the number of mutations per clone was approximately Poisson distributed with an average of 1.7 mutations per clone for the PR/1934 libraries and 2.1 mutations per clone for the Aichi/1968 libraries, with mutations distributed uniformly across the length of the gene.

Growth of mutant virus libraries We used reverse genetics [47] to rescue viruses carrying mutant NP genes. Co-cultures of 293T and MDCK-SIAT1 cells were plated 16 hours prior to transfection in D10 media (DMEM supplemented with 10% FBS, 100 U/mL of penicillin, 100 μ g/mL of streptomycin, and 2 mM L-glutamine) at cell densities of 3×10^5 293T/mL and 2.5×10^4 MDCK-SIAT1/mL. Co-cultures were transfected using BioT transfection reagent (Bioland Scientific) with a mixture of 250 ng of each of the eight reverse genetics plasmids per well in 6-well plates. In order to circumvent the possibility of rare mutants with exceptional replication fitness growing to high frequencies and limiting the growth of other mutants, we divided each transfection into multiple tissue-culture wells.

For the PR/1934 libraries, we rescued viruses containing the mutagenized PR/1934 NP with the seven remaining PR/1934 viral gene segments, and each replicate mutant library was transfected into the twelve wells of two 6-well plates. For the Aichi/1968 libraries, we used a viral rescue protocol that increases the number of parallel transfections and uses 293T cells that constitutively express protein V from hPIV2. This protein targets STAT1 for degradation, thereby inhibiting type I interferon signaling [48]. We rescued these Aichi/1968 virus libraries by transfecting the Aichi/1968 NP mutant library along with PB1/PB2/PA from Nanchang/933/1995 (using the plasmids in [8] and HA/NA/M/NS from WSN/1933 into 48 wells of eight 6-well plates. For both homologs, in parallel, we performed similar transfections using the corresponding unmutated NP genes to grow unmutated virus.

At 24 hours after transfection, co-culture media was aspirated, cells were rinsed with PBS, and the media was changed to influenza growth media (OptiMEM I media (Gibco) supplemented with 0.01% FBS, 0.3% BSA, 100 U/mL of penicillin, 100 μ g/mL of streptomycin, 100 μ g/mL calcium chloride, and 3 μ g/mL TPCK-trypsin). Co-culture supernatant was collected 72 hours after transfection, clarified by centrifugation at 2,000 \times g for 5 min, aliquoted and stored at -80° C.

Since many of the virions obtained from transfection with mutant NP library plasmids are likely to have originated in cells that contained more than one mutant NP gene and therefore might carry NP genes and NP proteins with different mutations, we passaged viruses in MDCK-SIAT1 cells at a low multiplicity of infection (MOI) to enforce genotype-phenotype linkage. We titrated viruses from thawed transfection supernatant aliquots for each replicate virus library using the TCID₅₀ protocol described in [35]. We then passaged viral libraries in MDCK-SIAT1 cells. Cells were plated in D10 media at 2 \times 10⁵ cells/mL. After 16 hours, the media was changed to influenza growth media containing diluted transfection supernatant virus. PR/1934 libraries were each passaged in 20 wells of 6-well dishes at an MOI of 0.05 TCID₅₀/cell, and Aichi/1968 libraries were each passaged in eight 10-cm dishes at an MOI of 0.1 TCID₅₀/cell. After 48 hours, supernatant was clarified by centrifugation at 2,000 \times g for 5 min, aliquoted and stored at -80° C.

Sample preparation and deep sequencing For each virus sample to be sequenced, 10 mL of clarified viral passage supernatant was centrifuged at 64,000 \times g for 1.5 hours to pellet viruses. RNA was extracted using the Qiagen RNeasy kit by lysing viral pellets in buffer RLT and following the manufacturer's recommended protocol. The NP gene was reverse transcribed using AccuScript High-Fidelity Reverse Transcriptase (Agilent Technologies) from both positive-sense and negative-sense viral RNA templates using the primers PR8-NP-RT-F (5'-agcaaaagcagggtagataatcactcactgagtgac-3') and PR8-NP-RT-R (5'-agtagaaacaagggtatttttcttta-3') for PR/1934 viruses or the primers 5'-BsmBI-Aichi68-NP (5'-catgatcgtctcaggagcaaaagcagggtagataatcactcacag-3') and 3'-BsmBI-Aichi68-NP (5'-catgatcgtctcgtattagtagaacaagggtatttttcttta-3') for Aichi/1968 viruses.

To ensure a sufficiently large number of unique RNA molecules were reverse transcribed, we used qPCR (SYBR Green Real-Time PCR Master Mix, Life Technologies) using primers qWSN-NP-for (5'-ACGGCTGGTCTGACTCACAT-3') and qPR8-NP-rev (5'-TCCATTCCGGTGCGAACAAG-3') to quantify the concentration of first-strand cDNA molecules against a standard curve of linear NP amplicons quantified by Quant-iT PicoGreen dsDNA Assay Kit (Life Technologies). We then made PCR amplicons with KOD DNA Polymerase (Merck Millipore) using at least 1 \times 10⁹ first-strand cDNA molecules as template in each reaction for viral gene sequencing. We also made PCR amplicons using 10 ng of the indicated plasmids for plasmid sequencing. For each biological replicate, we generated these PCR amplicons with 25 cycles of amplification using unmutated NP plasmid, mutated NP plasmid, NP cDNA from unmutated virus, and NP cDNA from mutated virus as template for the **DNA**, **mutDNA**, **virus**, and **mutvirus** samples, respectively.

To reduce the sequencing error rate, we developed a sequencing sample preparation protocol that results in sequencing libraries with inserts approximately 150 bp long. This allowed us to use paired-end 150 bp sequencing to achieve mostly overlapping reads so that sequencing errors resulting in mismatches between the two reads could be identified and ignored during data analysis. To make these sequencing libraries, we gel-purified the **DNA**, **mutDNA**, **virus**, and **mutvirus** PCR amplicons and sheared 1 μ g of each amplicon using Covaris to a median size of approximately 150 bp. We followed the modified Illumina paired-end library preparation protocol provided in [49] for end repair, 3' A overhang, and adapter ligation steps, using

Zymo DNA Clean & Concentrator columns (Zymo Research) or Ampure XP (Beckman Coulter) magnetic beads for DNA clean-up after shearing, end repair, and 3' A overhang steps. Barcoded Y-adapters were made by annealing 10 μ L of 100 μ M PAGE purified universal adapter (5'-AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTC TTCCGATC*T-3', where * indicates phosphorothioate bond) to 10 μ L of 100 μ M PAGE purified barcoded adapter (5'-PGATCGGAAGAGCACACGTCTGAACTC CAGTCACNNNNNATCTCGTATGCCGTCTTCTGCTT*G-3', where P indicates 5' phosphorylation, * indicates phosphorothioate bond, and NNNNNN indicates sample-specific barcode). Each 20 μ L mixture (one mixture for each barcode sequence) was annealed by heating to 95° C for 5 minutes and cooling at 0.3° C/second to 4° C. The resulting Y-adapters were diluted to 25 μ M by adding 20 μ L 10 mM Tris pH 7.5 and stored in 4 μ L aliquots at -20° C. Y-adapters with unique barcodes (ATCACG, ACTTGA, TAGCTT, GGCTAC, TTAGGC, GATCAG, ACTGAT, CGTACG, CGATGT, TGACCA, CAGATC, and CCGTCC) were ligated to samples derived from each biological replicate of each amplicon and ligation products were purified using 0.8X bead-to-sample ratio Ampure XP.

Purified adapter-ligated products for each sample were quantified by Quant-iT PicoGreen dsDNA Assay Kit (Life Technologies) and 25 ng was used as template for a 4-cycle PCR using Phusion High-Fidelity Polymerase (Thermo Scientific) to amplify inserts with adapters properly ligated on both sides. This amplification step was performed with the following components: 25 ng template DNA, 5 μ L 5X Phusion buffer, 2.5 μ L mixture of each dNTP at 2.5 mM, 2 μ L forward primer at 10 μ M (5'-AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGA-3'), 2 μ L reverse primer at 10 μ M (5'-CAAGCAGAAGACGGCATACGAGAT-3'), and 0.25 μ L Phusion polymerase in a final reaction volume of 25 μ L. PCR products were purified using 1.0X bead-to-sample ratio Ampure XP and quantified using PicoGreen. Samples were pooled in equal amounts and size-selected on a 2.0% agarose gel for fragments between 240 bp and 300 bp, which contain sequencing inserts in the size range of 120-180 bp. The size-selected sample was then sequenced at the Fred Hutchinson Genomics Core on an Illumina HiSeq 2500 using a paired-end 150 bp sequencing strategy in rapid run mode.

Analysis of deep sequencing data Sequencing data processing was performed using the software package `mapmut` [26]. Briefly, for each replicate sample of **DNA**, **mutDNA**, **virus**, and **mutvirus**, paired reads were stripped of any adapter sequence and aligned to each other. Read pairs were discarded if any of the following criteria were met: less than 100 bp of overlap between reads, average Q-score less than 25 across either read, more than 5 ambiguous nucleotides (N nucleotides) in either read, or more than 1 mismatch in the overlap between reads. Retained read pairs were then aligned to the appropriate reference NP gene sequence for PR/1934 or Aichi/1968 NP, and read pairs with more than 10 mismatches to the reference sequence or with any gaps or insertions were discarded. Once aligned to the reference sequence, codon identities at every position were called only if all three nucleotides in the codon matched unambiguously in both reads. The total number of codon identities at every codon position in the coding region were totaled for each sample (**DNA**, **mutDNA**, **virus**, and **mutvirus**), separately for each biological replicate.

Inference of amino-acid preferences We specify that at every site r in the protein, there is an inherent preference $\pi_{r,a}$ for every amino acid a , and we specify that $\sum_a \pi_{r,a} = 1$. The preference $\pi_{r,a}$ can be considered to be the expected frequency of amino acid a at site r in a mutant virus library after viral growth from a starting plasmid mutant library that contains equal numbers of every amino acid encoded at site r . Thus, mutations to amino acids with high preferences are beneficial and will be selected for during viral growth, and mutations to amino acids with low preferences will inhibit viral growth and will be selected against. Since the plasmid mutant libraries we generated contain on average more than one mutation per clone, the amino-acid preferences we measure represent an average preference in a variety of genetic backgrounds very similar to the starting sequence.

Let $\mathcal{A}(x)$ represent the amino acid encoded by codon x and let \mathcal{C} represent the set of all codons. The effect of the preference $\pi_{r,\mathcal{A}(x)}$ on the frequency f of observing codon x at site r in the mutant virus library

sample **mutvirus** is given by:

$$f_{r,x}^{\text{mutvirus}} = \epsilon_{r,x} + \rho_{r,x} + \frac{\mu_{r,x} \times \pi_{r,\mathcal{A}(x)}}{\sum_{y \in \mathcal{C}} \mu_{r,y} \times \pi_{r,\mathcal{A}(y)}} \quad (1)$$

where $\epsilon_{r,x}$ is the rate of PCR and sequencing errors at site r resulting in codon x , $\rho_{r,x}$ is the rate of reverse transcription errors at site r resulting in codon x , and $\mu_{r,x}$ is the frequency of codon x at site r in the plasmid mutant library **mutDNA**.

We inferred the amino-acid preferences independently for each biological replicate using the Bayesian algorithm described in [36] as implemented in **dms_tools** where codon counts in the **DNA**, **virus**, and **mutDNA** samples are used to infer the unknown parameters ϵ , ρ , and μ at each site.

Amino-acid preferences for Aichi/1968 NP were previously published in [26], where 8 biological replicates of the entire experiment were performed. In this work we report two additional biological replicates of the deep mutational scanning experiment for Aichi/1968. We will distinguish the two data sets when they are used separately for comparison as Aichi/1968 *previous study* and Aichi/1968 *current study*, and we will call the combined dataset of all 10 biological replicates for this homolog *Aichi/1968*.

Comparison of site-specific amino-acid preferences between homologs

Quantifying the magnitude of amino-acid preference difference between homologs At every site in the protein, each replicate deep mutational scanning experiment allows for the inference of an amino-acid preference distribution $\vec{\pi}$ that provides the preference at that site for all 20 amino acids. We used the Jensen-Shannon distance metric (the square root of the Jensen-Shannon divergence) to quantify the distance d between two amino-acid preference distributions:

$$d(\vec{\pi}_1, \vec{\pi}_2) = \sqrt{H\left(\frac{\vec{\pi}_1 + \vec{\pi}_2}{2}\right) - \frac{H(\vec{\pi}_1) + H(\vec{\pi}_2)}{2}} \quad (2)$$

where $H(\vec{\pi})$ is the Shannon entropy of the amino-acid preference distribution $\vec{\pi}$. The Jensen-Shannon distance metric quantifies the similarity between two amino-acid preference distributions, ranging from 0 (identical distributions) to 1 (completely dissimilar distributions). The average distance d between amino-acid preferences inferred from replicate experiments in the same homolog varies across sites. In other words, at some sites in the protein $\vec{\pi}$ is measured with greater precision than others. We therefore sought to develop, for every site r , a quantitative measure of the magnitude of change in $\vec{\pi}$ between homologs that corrects for the variation in $\vec{\pi}$ within replicate experiments of the same homolog.

For two groups of replicate mutational-scanning experiments A and B done in different homologs, each containing several replicate inferences of $\vec{\pi}$ for every site, we calculate the root-mean-square distance at site r over all pairwise comparisons of $\vec{\pi}$ measured in replicate experiments i (from group A) and j (from group B):

$$RMSD_{r,between} = \sqrt{\frac{1}{N_{A,B}} \sum_{i \in A} \sum_{j \in B} d(\vec{\pi}_{r,i}, \vec{\pi}_{r,j})^2} \quad (3)$$

where $N_{A,B}$ is the total number of non-redundant pairwise comparisons between replicate preferences measured from groups A and B . At the same site, to estimate the amount of experimental noise within replicates of the same homolog, we calculate the root-mean-square distance over all pairwise comparisons of $\vec{\pi}$ *within* the same group of replicate experiments, and average this site-specific noise estimate across the two groups:

$$RMSD_{r,within} = \frac{1}{2} \sqrt{\frac{1}{N_{A,A}} \sum_{i,j \in A, i < j} d(\vec{\pi}_{r,i}, \vec{\pi}_{r,j})^2} + \frac{1}{2} \sqrt{\frac{1}{N_{B,B}} \sum_{i,j \in B, i < j} d(\vec{\pi}_{r,i}, \vec{\pi}_{r,j})^2} \quad (4)$$

where $N_{A,A}$ and $N_{B,B}$ are the number of non-redundant pairwise comparisons between replicates within groups A and B , respectively. We then subtract the magnitude of the noise at this site observed *within*

groups from our measurement of the difference in amino-acid preferences seen *between* groups to obtain a corrected value for the change in $\vec{\pi}$ at site r between homologs:

$$RMSD_{r,corrected} = RMSD_{r,between} - RMSD_{r,within} \quad (5)$$

Identifying sites with statistically significant changes in amino-acid preference To determine whether site-specific $RMSD_{corrected}$ values are significantly larger than expected if amino-acid preferences are unchanged between homologs, we applied two methods to generate null distributions of $RMSD_{corrected}$ values. First, we used exact randomization testing to make all possible shuffles of the replicate homolog datasets into the two groups A and B . For each permutation, we calculated the $RMSD_{corrected}$ at every site, and the results are combined for all permutations. If there are no differences in preferences between homologs, the distribution of scores generated through randomization should be similar to the distribution of scores from the actual experiment.

We next observed that the overall correlation of amino-acid preferences across all sites between replicates can vary between experiments. For instance, the average Pearson’s correlation between PR/1934 replicates is 0.59, the correlation between Aichi/1968 replicates in the *previous study* is 0.50, and the correlation between Aichi/1968 replicates in the *current study* is 0.74. We considered whether the varying precision between homologs might lead to biases in the calculated $RMSD_{corrected}$.

To test this, we generated a second null distribution of $RMSD_{corrected}$ under the hypothesis that the “true” amino-acid preferences are the same for both homologs and can be approximated by averaging the mean observed preferences for each homolog:

$$\langle\langle\vec{\pi}_r\rangle\rangle = \frac{\langle\vec{\pi}_{r,homolog\ A}\rangle + \langle\vec{\pi}_{r,homolog\ B}\rangle}{2} \quad (6)$$

Under this hypothesis, the observed differences in amino-acid preferences between homologs is solely due to the different amounts of experimental noise between replicates of each homolog. To model the effects of this noise on our analysis, we drew replicate simulated amino-acid preferences at each site r from a Dirichlet distribution with mean centered on the “true” amino-acid preferences:

$$\vec{\pi}_{r,simulated\ A} = Dir(\langle\langle\vec{\pi}_r\rangle\rangle \times \sigma_A) \quad (7)$$

where σ_A is a scaling factor that is chosen to yield simulated replicate preferences across the entire protein that have an average Pearson’s correlation between replicates equal to the correlation between experimental replicates. In other words, we simulate replicate amino-acid preference measurements with noise tuned to match the actual noise in each experiment. For each simulated experiment, we simulated the same number of replicates that were performed experimentally, and calculated $RMSD_{corrected}$ for all sites. We ran the entire simulation 1000 times, combining all $RMSD_{corrected}$ values to obtain a null distribution.

We then separately used the two null distributions (generated through randomization or simulation) to assign p-values to site-specific $RMSD_{corrected}$ at each site r :

$$p_r = \frac{\text{number of scores in null distribution} \geq RMSD_{r,corrected}}{\text{number of scores in null distribution}} \quad (8)$$

We used the Bonferroni correction to set the threshold for significant p-values to account for the 497 sites we tested for significance.

Phylogenetic analysis

Experimental substitution model overview We used a previously described approach to build site-specific substitution models for influenza nucleoprotein [26,27]. Briefly, this approach calculates the codon-substitution rate at each site in nucleoprotein based on the rate at which nucleotide mutations arise and the level of selection acting on these new mutations. The rate of codon substitution, $P_{r,xy}$, at site r of codon x to a different codon y is described as,

$$P_{r,xy} = Q_{xy} \times F_{r,xy} \quad (9)$$

where Q_{xy} is the rate of mutation from x to y, and $F_{r,xy}$ is the probability that a mutation from x to y at site r is selected and reaches fixation. In this equation, the mutation rates Q_{xy} are assumed to be identical across sites whereas the selection is modeled as site-specific and site-independent. The site-specific fixation probabilities $F_{r,xy}$ were calculated from the experimentally measured amino-acid preferences using the relationship proposed by Halpern and Bruno [27, 50]. The four mutation rate free parameters and the stringency parameter were defined as in [27].

We then calculated the phylogenetic likelihood of the observed nucleoprotein sequences given the resulting experimental substitution model $P_{r,xy}$, the nucleoprotein phylogenetic tree, and the model parameters. The tree consisted of influenza nucleoproteins from either human, swine, equine, or avian hosts. While holding the tree topology fixed, tree branch lengths, and any other model parameters (discussed below), were optimized by maximum likelihood.

To compare overall phylogenetic likelihoods calculated under various substitution models, we calculated the difference in the Akaike Information Criteria (ΔAIC) between models. We compared site-specific models derived from experimentally determined amino-acid preferences to a non-site-specific model. We tested separate site-specific models using the amino-acid preferences from PR/1934 and Aichi/1968. The Aichi/1968 preferences were an average of the amino-acid preferences from the *current study* and *previous study*. In addition, we tested a site-specific model where we combined data from the separate Aichi/1968 and PR/1934 mutational-scanning experiments, by averaging amino-acid preferences for each amino acid at each site across the two homologs, weighting each homolog equally.

The non-site-specific model used the Goldman-Yang (GY94) codon substitution model [34], with nucleotide equilibrium frequencies calculated by the CF3x4 method [38]. In this model, the transition-transversion ratio was optimized by maximum likelihood, along with the mean and shape parameters describing gamma distributions of the nonsynonymous-synonymous ratios [39] and the substitution rates [21] across sites. Each gamma distribution was discretized with four categories. In previous comparisons of non-site-specific models, this non-site-specific model performed better than other variants of the GY94 model [26, 27]. All analyses were performed using the software packages *phyloExpCM* [26] and *HyPhy* [51], and the data, scripts, and descriptions to replicate the results in this article are available at <https://github.com/mbdoud/Compare-NP-Preferences>.

Phylogenetic trees for different influenza hosts We built phylogenetic trees for nucleoprotein coding sequences from strains of human influenza, swine influenza, equine influenza, and avian influenza. Full-length nucleoprotein sequences were downloaded from the Influenza Virus Resource [52], and for each host, a small number of unique sequences per year per influenza subtype were retained. For human influenza, we retained one sequence every other year from each of the H1N1, H2N2, and H3N2 lineages. For swine influenza, we retained one sequence per year from either the North American Classical H1N1 lineage or the Eurasian H1N1 lineage. For equine influenza, we retained one sequence per year from the H3N8 lineage. For avian influenza, one sequence every other year per subtype was retained, and the examined hosts were further restricted to only duck species, to make a sequence set with a size manageable for phylogenetic modeling.

Sequences from each host were aligned by EMBOSS *needle* [53], and maximum-likelihood trees were built by RAxML [54]. Using these trees and the program *Path-O-Gen* (<http://tree.bio.ed.ac.uk/software/pathogen/>), we identified and removed any sequences that were noticeable outliers from the molecular clock. The final tree contained 37, 46, 29, and 24 sequences from human, swine, equine, and avian hosts respectively.

Maximum-likelihood phylogenetic trees were then built from the nucleoprotein sequence alignment using *codonPhyML* [33]. The GY94 model [34] was run using the CF3x4 nucleotide equilibrium frequencies [38] along with maximum-likelihood optimization of a transition-transversion ratio and of a mean and shape parameter describing a gamma distribution of nonsynonymous-synonymous ratios [39]. This gamma distribution was discretized with four categories. The final, unrooted tree was visualized with *FigTree* (<http://tree.bio.ed.ac.uk/software/figtree/>) and rooted using the avian clade [40].

Supporting Information

S1 Text

Mean amino-acid preferences for PR/1934 NP. This text file lists the mean amino-acid preferences for sites 2 through 498 in PR/1934 NP. The amino-acid preferences inferred from three biological replicates for PR/1934 NP were averaged at each site.

S2 Text

Mean amino-acid preferences for Aichi/1968 NP. This text file lists the mean amino-acid preferences for sites 2 through 498 in Aichi/1968 NP. The average is taken from the average across the previous study replicates and the average from the current study replicates.

S3 Text

Mean amino-acid preferences for combined PR/1934+Aichi/1968 NP. This text file lists the amino-acid preferences averaged evenly across the two homologs.

S4 Text

NP sequence alignment used to build phylogenetic tree. The alignment consists of human, swine, equine, and avian NP-coding sequences.

S5 Text

Phylogenetic tree of NP. The tree contains human, swine, equine, and avian NPs.

S1 Fig

Characterization of plasmid mutant libraries generated by codon mutagenesis. The distributions of number of mutated codons per clone (A, C) and number of nucleotide changes per codon mutation (B, D) were determined by full-length Sanger sequencing of individual clones. A-B: PR/1934 libraries, C-D: Aichi/1968 libraries.

S2 Fig

Logoplot of amino-acid preferences for PR/1934 NP. The mean preferences for sites 2 through 498 of PR/1934 are represented in a sequence logo-like visualization created with the program `dms_logoplot`. The height of each letter is proportional to the preference for that amino-acid at that site.

S3 Fig

Logoplot of amino-acid preferences for Aichi/1968 NP. The mean preferences for sites 2 through 498 of Aichi/1968 are represented in a sequence logo-like visualization created with the program `dms_logoplot`. The height of each letter is proportional to the preference for that amino-acid at that site.

S4 Fig

Logoplot of amino-acid preferences for combined PR/1934+Aichi/1968 NP. The mean preferences for sites 2 through 498 of the combined Aichi/1968 + PR/1934 model are represented in a sequence logo-like visualization created with the program `dms_logoplot`. The height of each letter is proportional to the preference for that amino-acid at that site.

S5 Fig

Null distributions of RMSDcorrected generated by randomization. The null distributions generated by randomization are shown in red; experimental distributions are shown in blue.

S1 Table

Amino-acid preference RMSD calculations for PR/1934 vs. Aichi/1968 NP. This table lists summary findings for the comparison of amino-acid preferences between PR/1934 and Aichi/1968 NP. For sites 2 through 498 (the initiating methionine was not mutagenized in our experiments) the amino-acid identity is noted to either be conserved or variable between PR/1934 and Aichi/1968 NP homologs. For variable sites, the PR/1934 amino-acid identity is listed first. The $RMSD_{between}$, $RMSD_{within}$, and $RMSD_{corrected}$ statistics calculated for the comparison between 3 replicates of PR/1934 and 10 replicates of Aichi/1968 are listed in the following columns. The next column marks sites with a corrected p-value of 0.05 or less as significant with an asterisk. P-values were calculated as the fraction of all simulated $RMSD_{corrected}$ values that are greater than the experimental value at the site, with the Bonferroni correction applied to control for 497 hypothesis tests. The final two columns list the mean amino-acid preferences for the most preferred amino acids as measured in each homolog. Amino acids are listed in their order of preference to account for the top 65% of preferences (preferences sum to 1 at each site), with up to three listed for each homolog; full amino-acid preference data is available in S1 Text, S2 Text, and S3 Text.

S2 Table

Combining experimentally informed substitution models for swine influenza NP. This table differs from Table 1 in that the phylogenetic fit is for the tree of swine NPs shown in Figure 1.

S3 Table

Combining experimentally informed substitution models for equine influenza NP. This table differs from Table 1 in that the phylogenetic fit is for the tree of equine NPs shown in Figure 1.

S4 Table

Combining experimentally informed substitution models for avian influenza NP. This table differs from Table 1 in that the phylogenetic fit is for the tree of avian NPs shown in Figure 1.

Acknowledgments

We thank Hugh Haddock, Alistair Russell, and Heather Machkovech for critical reading of the manuscript and Trevor Bedford for helpful discussions about statistical analysis. We thank the Summer Institute in Statistics and Modeling in Infectious Diseases at the University of Washington for helpful instruction and the Genomics Shared Resource at the Fred Hutchinson Cancer Research Center for performing high-throughput sequencing.

References

1. Zuckerkandl E, Pauling L. Evolutionary divergence and convergence in proteins. In: *Evolving genes and proteins*. New York, NY: Academic Press; 1965. p. 97–166.
2. DePristo MA, Weinreich DM, Hartl DL. Missense meanderings in sequence space: a biophysical view of protein evolution. *Nature Reviews Genetics*. 2005;6(9):678–687.
3. Harms MJ, Thornton JW. Evolutionary biochemistry: revealing the historical and physical causes of protein properties. *Nature Reviews Genetics*. 2013;14(8):559–571.

4. Weinreich DM, Delaney NF, DePristo MA, Hartl DL. Darwinian evolution can follow only very few mutational paths to fitter proteins. *Science*. 2006;312(5770):111–114.
5. Ortlund EA, Bridgham JT, Redinbo MR, Thornton JW. Crystal structure of an ancient protein: evolution by conformational epistasis. *Science*. 2007;317:1544–1548.
6. da Silva J, Coetzer M, Nedellec R, Pastore C, Mosier DE. Fitness epistasis and constraints in adaptation in a human immunodeficiency virus type 1 protein region. *Genetics*. 2010;185:293–303.
7. Lunzer M, Golding GB, Dean AM. Pervasive cryptic epistasis in molecular evolution. *PLoS Genetics*. 2010;6(10):e1001162.
8. Gong LI, Suchard MA, Bloom JD. Stability-mediated epistasis constrains the evolution of an influenza protein. *eLife*. 2013;2:e00631.
9. Natarajan C, Inoguchi N, Weber RE, Fago A, Moriyama H, Storz JF. Epistasis among adaptive mutations in deer mouse hemoglobin. *Science*. 2013;340(6138):1324–1327.
10. Podgornaia A, Laub M. Protein evolution. Pervasive degeneracy and epistasis in a protein-protein interface. *Science*. 2015;347(6222):673–677.
11. Risso VA, Manssour-Triedo F, Delgado-Delgado A, Arco R, Barroso-delJesus A, Ingles-Prieto A, et al. Mutational Studies on Resurrected Ancestral Proteins Reveal Conservation of Site-Specific Amino Acid Preferences throughout Evolutionary History. *Molecular Biology and Evolution*. 2015;32(2):440–455.
12. Ashenberg O, Gong LI, Bloom JD. Mutational effects on stability are largely conserved during protein evolution. *Proceedings of the National Academy of Sciences*. 2013;110(52):21071–21076.
13. Serrano L, Day AG, Fersht AR. Step-wise mutation of barnase to binase: a procedure for engineering increased stability of proteins and an experimental analysis of the evolution of protein stability. *J Mol Biol*. 1993;233:305–312.
14. Chothia C, Lesk AM. The relation between the divergence of sequence and structure in proteins. *EMBO J*. 1986;5:823–826.
15. Sander C, Schneider R. Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins: Structure, Function, and Bioinformatics*. 1991;9(1):56–68.
16. Chou PY, Fasman GD. Conformational parameters for amino acids in helical, β -sheet, and random coil regions calculated from proteins. *Biochemistry*. 1974;13(2):211–222.
17. Richardson JS, Richardson DC. Amino acid preferences for specific locations at the ends of alpha helices. *Science*. 1988;240(4859):1648–1652.
18. Lim WA, Sauer RT. The role of internal packing interactions in determining the structure and stability of a protein. *Journal of molecular biology*. 1991;219(2):359–376.
19. Henikoff S, Henikoff JG. Embedding strategies for effective use of information from multiple sequence alignments. *Protein Science*. 1997;6(3):698–705.
20. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic acids research*. 1997;25(17):3389–3402.
21. Yang Z. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *J Mol Evol*. 1994;39(3):306–314.
22. Lartillot N, Philippe H. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Molecular biology and evolution*. 2004;21(6):1095–1109.

23. Le SQ, Lartillot N, Gascuel O. Phylogenetic mixture models for proteins. *Phil Trans R Soc B*. 2008;363(1512):3965–3976.
24. Wang HC, Li K, Susko E, Roger AJ. A class frequency mixture model that adjusts for site-specific amino acid frequencies and improves inference of protein phylogeny. *BMC evolutionary biology*. 2008;8(1):331.
25. Rodrigue N, Philippe H, Lartillot N. Mutation-selection models of coding sequence evolution with site-heterogeneous amino acid fitness profiles. *Proceedings of the National Academy of Sciences*. 2010;107(10):4629–4634.
26. Bloom JD. An experimentally determined evolutionary model dramatically improves phylogenetic fit. *Molecular Biology and Evolution*. 2014;30:1956–1978.
<http://mbe.oxfordjournals.org/content/31/8/1956>.
27. Bloom JD. An experimentally informed evolutionary model improves phylogenetic fit to divergent lactamase homologs. *Molecular Biology and Evolution*. 2014;31:2753–2769.
<http://mbe.oxfordjournals.org/content/31/10/2753>.
28. Fowler DM, Araya CL, Fleishman SJ, Kellogg EH, Stephany JJ, Baker D, et al. High-resolution mapping of protein sequence-function relationships. *Nat Methods*. 2010;7(9):741–746.
29. Fowler DM, Fields S. Deep mutational scanning: a new style of protein science. *Nature methods*. 2014;11(8):801–807.
30. Eisefeld AJ, Neumann G, Kawaoka Y. At the centre: influenza A virus ribonucleoproteins. *Nature Reviews Microbiology*. 2015;13(1):28–41.
31. Ye Q, Krug RM, Tao YJ. The mechanism by which influenza A virus nucleoprotein forms oligomers and binds RNA. *Nature*. 2006;444(7122):1078–1082.
32. Das K, Aramini JM, Ma LC, Krug RM, Arnold E. Structures of influenza A proteins and insights into antiviral drug targets. *Nature structural & molecular biology*. 2010;17(5):530–538.
33. Gil M, Zanetti MS, Zoller S, Anisimova M. CodonPhyML: Fast Maximum Likelihood Phylogeny Estimation under Codon Substitution Models. *Mol Biol Evol*. 2013;30(6):1270–1280.
34. Goldman N, Yang Z. A codon-based model of nucleotide substitution probabilities for protein-coding DNA sequences. *Mol Biol Evol*. 1994;11:725–736.
35. Thyagarajan B, Bloom JD. The inherent mutational tolerance and antigenic evolvability of influenza hemagglutinin. *eLife*. 2014;3:e03300.
36. Bloom JD. Software for the analysis and visualization of deep mutational scanning data. *bioRxiv*. 2015;p. doi:10.1101/013623. <http://dx.doi.org/10.1101/013623>.
37. Posada D, Buckley TR. Model selection and model averaging in phylogenetics: advantages of Akaike information criterion and Bayesian approaches over likelihood ratio tests. *Systematic Biology*. 2004;53(5):793–808.
38. Pond SK, Delpont W, Muse SV, Scheffler K. Correcting the bias of empirical frequency parameter estimators in codon models. *PLoS One*. 2010;5(7):e11230.
39. Yang Z, Nielsen R, Goldman N, Pedersen AMK. Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics*. 2000;155(1):431–449.
40. Worobey M, Han GZ, Rambaut A. A synchronized global sweep of the internal genes of modern avian influenza virus. *Nature*. 2014;508(7495):254.

41. Choi SC, Hobolth A, Robinson DM, Kishino H, Thorne JL. Quantifying the impact of protein tertiary structure on molecular evolution. *Molecular biology and evolution*. 2007;24(8):1769–1782.
42. Bordner AJ, Mittelman HD. A new formulation of protein evolutionary models that account for structural constraints. *Molecular biology and evolution*. 2014;31(3):736–749.
43. Bloom JD, Silberg JJ, Wilke CO, Drummond DA, Adami C, Arnold FH. Thermodynamic prediction of protein neutrality. *Proc Natl Acad Sci USA*. 2005;102:606–611.
44. Pollock DD, Thiltgen G, Goldstein RA. Amino acid coevolution induces an evolutionary Stokes shift. *Proceedings of the National Academy of Sciences*. 2012;109(21):E1352–E1359.
45. Kellogg EH, Leaver-Fay A, Baker D. Role of conformational sampling in computing mutation-induced changes in protein structure and stability. *Proteins: Structure, Function, and Bioinformatics*. 2011;79(3):830–838.
46. Potapov V, Cohen M, Schreiber G. Assessing computational methods for predicting protein stability upon mutation: good on average but not in the details. *Prot Eng Des Sel*. 2009;22(9):553–560.
47. Hoffmann E, Neumann G, Kawaoka Y, Hobom G, Webster RG. A DNA transfection system for generation of influenza A virus from eight plasmids. *Proceedings of the National Academy of Sciences*. 2000;97(11):6108–6113.
48. Andrejeva J, Young D, Goodbourn S, Randall R. Degradation of STAT1 and STAT2 by the V proteins of simian virus 5 and human parainfluenza virus type 2, respectively: consequences for virus replication in the presence of alpha/beta and gamma interferons. *Journal of virology*. 2002;76(5):2159–2167.
49. Henikoff JG, Belsky JA, Krassovsky K, MacAlpine DM, Henikoff S. Epigenome characterization at single base-pair resolution. *Proceedings of the National Academy of Sciences*. 2011;108(45):18318–18323.
50. Halpern AL, Bruno WJ. Evolutionary distances for protein-coding sequences: modeling site-specific residue frequencies. *Mol Biol Evol*. 1998;15(7):910–917.
51. Pond SL, Frost SD, Muse SV. HyPhy: hypothesis testing using phylogenies. *Bioinformatics*. 2005;21(5):676–679.
52. Bao Y, Bolotov P, Dernovoy D, Kiryutin B, Zaslavsky L, Tatusova T, et al. The Influenza Virus Resource at the National Center for Biotechnology Information. *J Virol*. 2008;82:596–601.
53. Rice P, Longden I, Bleasby A. EMBOSS: the European molecular biology open software suite. *Trends in Genetics*. 2000;16(6):276–277.
54. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. *Bioinformatics*. 2006;22(21):2688–2690.