

1           **Modeling site-specific amino-acid preferences deepens**  
2           **phylogenetic divergence estimates for viral proteins.**

3           Sarah K. Hilton<sup>1,2</sup> and Jesse D. Bloom<sup>1,2\*</sup>

<sup>1</sup>Basic Sciences and Computational Biology Program, Fred Hutchinson Cancer Research Center

<sup>2</sup>Department of Genome Sciences, University of Washington

Seattle, WA, USA

\*E-mail: [jbloom@fredhutch.org](mailto:jbloom@fredhutch.org).

4

### **Abstract**

5           Molecular phylogenetics is often used to estimate the time since the divergence of modern  
6           gene sequences. For highly diverged sequences, such phylogenetic techniques sometimes  
7           estimate surprisingly recent divergence times. In the case of viruses, independent evidence  
8           indicates that the estimates of deep divergence times from molecular phylogenetics are some-  
9           times too recent. This discrepancy is caused in part by inadequate models of purifying se-  
10          lection leading to branch-length underestimation. Here we examine the effect on branch-  
11          length estimation of using models that incorporate experimental measurements of purifying  
12          selection. We find that models informed by experimentally measured site-specific amino-acid  
13          preferences estimate longer deep branches on phylogenies of influenza virus hemagglutinin.  
14          This lengthening of branches is due to more realistic stationary states of the models, and  
15          is mostly independent of the branch-length-extension from modeling site-to-site variation in  
16          amino-acid substitution rate. The branch-length extension from experimentally informed site-  
17          specific models is similar to that achieved by other approaches that allow the stationary state  
18          to vary across sites. However, the improvements from all of these site-specific but time-  
19          homogeneous and site-independent models are limited by the fact that a protein’s amino-acid  
20          preferences gradually shift as it evolves. Overall, our work underscores the importance of  
21          modeling site-specific amino-acid preferences when estimating deep divergence times—but  
22          also shows the inherent limitations of approaches that fail to account for how these preferences  
23          shift over time.

24 **Introduction**

25 Molecular phylogenetics is commonly used to estimate the historical timing of evolutionary events  
26 ([Yang and Rannala, 2012](#)). This is done by estimating branch lengths based on the inferred number  
27 of substitutions, and then converting these branch lengths into units of time under the assumption  
28 of a molecular clock ([Zuckerkandl and Pauling, 1965](#); [Drummond et al., 2006](#)). However, phylo-  
29 genetic estimates of the divergence times of many viral lineages are clearly too recent ([Duchêne](#)  
30 [et al., 2014](#); [Ho et al., 2015](#); [Aiewsakun and Katzourakis, 2016](#)). For example, the integration of  
31 filoviruses into their host genomes indicate that Ebola and Marburg virus diverged from their com-  
32 mon ancestor 7 to 12 million years ago—but the estimate of this divergence time based on phylo-  
33 genetic analyses of the viral sequences is only ~10,000 years ago ([Carroll et al., 2013](#); [Taylor et al.,](#)  
34 [2014](#)). Similarly, the phylogenetic estimate of when major simian immunodeficiency virus groups  
35 diverged is almost 100 times more recent than the estimate based on the geographic isolation of  
36 their host species ([Wertheim and Worobey, 2009](#); [Worobey et al., 2010](#)). These examples, along  
37 with other similar discrepancies with measles virus ([Furuse et al., 2010](#)), coronavirus ([Wertheim](#)  
38 [et al., 2013](#)), and hepatitis B virus ([Fares and Holmes, 2002](#); [Holmes, 2003](#)), indicate that phylo-  
39 genetic methods have a systematic bias toward underestimation of deep branches.

40 This underestimation occurs in part because phylogenetic models do a poor job of describing  
41 the real natural selection on protein-coding genes. These genes evolve under purifying selection  
42 to maintain the structure and function of the proteins they encode. In general, these constraints  
43 are highly idiosyncratic among sites ([Echave et al., 2016](#)). However, most phylogenetic models  
44 try to account for these constraints using relatively simple approaches such as allowing the rate of  
45 substitution to vary across sites according to some statistical distribution ([Yang, 1994](#); [Yang et al.,](#)  
46 [2000](#)). These models of purifying selection are usually inadequate ([Duchêne et al., 2015b,a](#)), po-  
47 tentially causing branch lengths to be severely underestimated ([Wertheim and Kosakovsky Pond,](#)  
48 [2011](#); [Halpern and Bruno, 1998](#)).

49 More recent work has used mutation-selection models to better account for purifying selec-  
50 tion ([Halpern and Bruno, 1998](#); [Yang and Nielsen, 2008](#); [Rodrigue et al., 2010](#); [Tamuri et al.,](#)  
51 [2012](#); [McCandlish and Stoltzfus, 2014](#)). These models explicitly incorporate the fact that different  
52 protein sites prefer different amino acids, and so can improve phylogenetic estimates when there  
53 are deep branches ([Philippe and Laurent, 1998](#); [Lartillot et al., 2007](#); [Le et al., 2008](#); [Quang et al.,](#)  
54 [2008](#); [Wang et al., 2008](#)). However, these approaches require inferring the site-specific purifying  
55 selection from natural sequence data.

56 Even more recently, it has become possible to directly measure purifying selection on proteins  
57 using deep mutational scanning. This high-throughput approach involves experimentally measur-  
58 ing how each amino-acid mutation affects protein function in the lab ([Fowler and Fields, 2014](#)).

59 The resulting experimental measurements of which amino acids are preferred at each protein site  
60 can be used to inform phylogenetic substitution models (Bloom, 2014a). These experimentally  
61 informed codon models (ExpCMs) generally exhibit much better phylogenetic fit than standard  
62 substitution models (Doud et al., 2015; Hilton et al., 2017; Haddox et al., 2018; Lee et al., 2018).

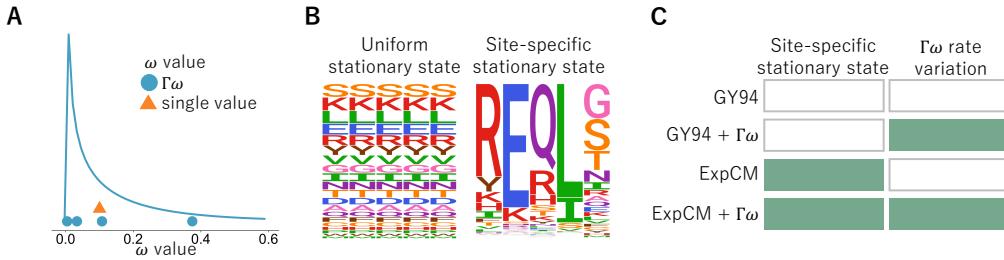
63 Here we examine how ExpCMs and other models of purifying selection estimate branch  
64 lengths on a phylogenetic tree of influenza virus hemagglutinin (HA). We find that ExpCMs esti-  
65 mate longer deep branches, and show that this extension of branch length is mostly independent  
66 and additive with that achieved by the more conventional approach of modeling rate variation. We  
67 also show that ExpCMs estimate similar branch lengths to a mutation-selection model that infers  
68 the amino-acid preferences from the natural sequence data rather than using values obtained in ex-  
69 periments. However, all of these mutation-selection models are limited by their failure to account  
70 for another feature of purifying selection: the fact that a site's amino-acid preferences shift over  
71 time due to epistasis. Therefore, truly accurate analyses of deep phylogenies need to account for  
72 the fact that amino-acid preferences vary across time as well as across sites.

## 73 Results

### 74 Different ways substitution models account for purifying selection

75 Here we consider how purifying selection is handled by codon models, which are the most so-  
76 phisticated of the three classes (nucleotide, codon, and amino acid) of phylogenetic substitution  
77 models in widespread use for protein-coding genes (Arenas, 2015). Standard codon models dis-  
78 tinguish between two types of substitutions: synonymous and nonsynonymous. The relative rate  
79 of these substitutions is referred to as dN/dS or  $\omega$ . In their simplest form, codon substitution mod-  
80 els fit a single  $\omega$  that represents the gene-wide average fixation rate of nonsynonymous mutations  
81 relative to synonymous ones. Here we will use such substitution models in the form proposed by  
82 Goldman and Yang (1994). When these models have a single gene-wide  $\omega$  they are classified as  
83 M0 by Yang et al. (2000). We will refer to M0 Goldman-Yang models simply as GY94 models  
84 (Equation 1). The gene-wide  $\omega$  is usually  $< 1$  (Murrell et al., 2015), and crudely represents the  
85 fact that many amino-acid substitutions are under purifying selection.

86 A single gene-wide  $\omega$  ignores the fact that purifying selection is heterogeneous across sites.  
87 The most common strategy to ameliorate this defect is to allow  $\omega$  to vary among sites according to  
88 some statistical distribution (Yang, 1994; Yang et al., 2000). For instance, in the M5 variant of the  
89 GY94 model (Yang et al., 2000),  $\omega$  follows a gamma distribution as shown in Figure 1A. We will  
90 denote this model as GY94+ $\Gamma\omega$ . A GY94+ $\Gamma\omega$  captures the fact that the rate of nonsynonymous  
91 substitution can vary across sites. However, these models do not capture the fact that the same  
92 amino-acid mutation can have very different effects at different sites.



**Figure 1: Different ways codon models account for purifying selection.** (A) The dN/dS parameter,  $\omega$ , can be defined as one gene-wide average (orange triangle) or allowed to vary according to some statistical distribution (blue line). For computational tractability, the distribution is discretized into  $K$  bins and  $\omega$  takes on the mean of each bin (blue circles) (Yang, 1994; Yang et al., 2000). A gamma distribution (denoted by  $\Gamma$ ) with  $K = 4$  bins is shown here. (B) A substitution model’s stationary state defines the expected sequence composition after a very long evolutionary time. Most substitution models have stationary states that are uniform across sites. However, substitution models can have site-specific stationary states. In the logo plots, each column is a site in the protein and the height of each letter is the frequency of that amino acid at stationary state. (C) Substitution models can incorporate neither, one, or both of these features. Here we will use substitution models from the Goldman-Yang (GY94; Goldman and Yang, 1994; Yang et al., 2000) and experimentally informed codon model (ExpCM; Hilton et al., 2017) families with and without gamma-distributed  $\omega$  to represent all possible combinations.

93 Mutation-selection models account for the fact that purifying selection depends idiosyncratically  
 94 on the specific amino-acid mutation at each site (Halpern and Bruno, 1998; Yang and Nielsen,  
 95 2008; Rodrigue et al., 2010; Tamuri et al., 2012; McCandlish and Stoltzfus, 2014). Here we will  
 96 consider mutation-selection models where the site-specific selection is assumed to act solely at  
 97 the protein level (different codons for the same protein are treated as selectively equivalent). Such  
 98 models explicitly define a different set of amino-acid preferences at each site in the protein. This  
 99 more mechanistic formulation results in a site-specific stationary state (Figure 1B). These models  
 100 capture the site-to-site variation in amino-acid composition that is an obvious features of real pro-  
 101 teins, and usually better describe actual evolution than models with only rate variation as assessed  
 102 by Bayesian or maximum-likelihood criteria (Lartillot and Philippe, 2004; Le et al., 2008; Quang  
 103 et al., 2008; Wang et al., 2008; Rodrigue et al., 2010; Bloom, 2014a,b; Hilton et al., 2017).

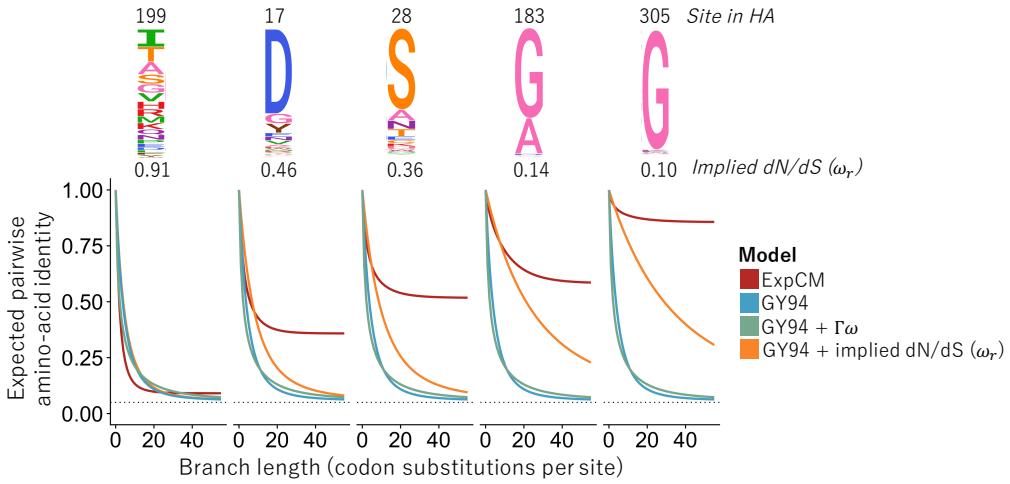
104 However, the increased realism of mutation-selection models comes at the cost of an increased  
 105 number of parameters. Codon substitution models with uniform stationary states have only a mod-  
 106 est number of parameters that must be fit from the phylogenetic data. For instance, a GY94+ $\Gamma\omega$   
 107 model with the commonly used F3X4 stationary state has 12 parameters: two describing the shape  
 108 of the gamma distribution over  $\omega$ , a transition-transversion rate, and nine parameters describing  
 109 the nucleotide composition of the stationary state. However, mutation-selection models must ad-

ditionally specify 19 parameters defining the amino-acid preferences for *each* site (there are 20 amino acids whose preferences are constrained to sum to one). This corresponds to  $19 \times L$  parameters for a protein of length  $L$ , or 9,500 parameters for a 500-residue protein. It is challenging to obtain values for these amino-acid preference parameters in a maximum-likelihood framework without overfitting the data (Rodrigue, 2013). Here we will primarily use experimentally informed codon models (ExpCMs), which define the site-specific amino-acid preference parameters *a priori* from deep mutational scanning experiments so that they do not need to be fit from phylogenetic data (see Methods and Bloom, 2014a; Hilton et al., 2017; Bloom, 2017). Because the amino-acid preference parameters in an ExpCM are obtained from experiments, the number of ExpCM free parameters is similar to a non-site-specific substitution model. An alternative strategy of obtaining the amino-acid preference parameters via Bayesian inference (Lartillot and Philippe, 2004; Rodrigue and Lartillot, 2014) is discussed in the last section of the Results.

Importantly, these two strategies for modeling purifying selection are not mutually exclusive. Mutation-selection models such as an ExpCM can still incorporate an  $\omega$  parameter, which now represents the relative rate of nonsynonymous to synonymous substitution *after* accounting for the constraints due to the site-specific amino-acid preferences (Bloom, 2017; Rodrigue and Lartillot, 2017). This  $\omega$  parameter for an ExpCM can be drawn from a statistical distribution (e.g., a gamma distribution) just like for GY94-style models (Rodrigue and Lartillot, 2014; Haddox et al., 2018). We will denote such models as ExpCM+ $\Gamma\omega$ . Figure 1C shows the full spectrum of models that incorporate all combinations of gamma-distributed  $\omega$  and site-specific stationary states.

### Effect of stationary state and rate variation on branch-length estimation

Given a single branch, a substitution model transforms sequence identity into branch length. Under a molecular-clock assumption, this branch length is proportional to time. The transformation from sequence identity to branch length is trivial when the sequence identity is high. For instance, when there has only been one substitution, then the sequence identity will simply be  $\frac{L-1}{L}$  for a gene of  $L$  sites, and even a simple exponential model (Zuckerkandl and Pauling, 1965) will correctly infer the short branch length of  $1/L$  substitutions per site. However, as substitutions accumulate it becomes progressively more likely for multiple changes to occur at the same site. In this regime, the accuracy of the substitution model becomes critical for transforming sequence identity into branch length. Any time-homogenous substitution model predicts that after a very large number of substitutions, two related sequences will approach some asymptotic amino-acid sequence identity. For instance, if all 20 amino acids are equally likely in the stationary state, then this asymptotic sequence identity will be  $\frac{1}{20} = 0.05$ . If the substitution model underestimates the asymptotic sequence identity then it will also underestimate long branch lengths, since it will predict that sequences that have evolved for a very long time should be more diverged than is



**Figure 2: Effect of stationary state and  $\Gamma\omega$  rate variation on predicted asymptotic sequence divergence.** The logo plots at top show the amino-acid preferences for some sites in an H1 influenza hemagglutinin protein as experimentally measured by Doud and Bloom (2016). The graphs show the expected amino-acid identity at that site for two sequences separated by a branch of the indicated length (Equation 9). For the GY94 model, the graphs are identical for all sites since this model does not have site-specific parameters; the same is true for GY94+ $\Gamma\omega$ . The graphs do differ among sites if we calculate a different  $\omega_r$  for each site  $r$  in the GY94 model using the amino-acid preferences (Equation 7; Spielman and Wilke, 2015b). However, all GY94 models, including the one with site-specific  $\omega_r$  values, approach the same asymptote since they all have the same stationary state. The ExpCM has different asymptotes for different sites since it accounts for how amino-acid preferences lead to site-specific stationary states.

145 actually the case.

146 Figure 2 shows how different substitution models predict amino-acid sequence identity to  
 147 decrease as a function of branch length using model parameters fit to a phylogeny of H1 influenza  
 148 hemagglutinin (HA) genes. The GY94 model predicts the same behavior for all sites, since it does  
 149 not have any site-specific parameters, with an asymptotic sequence identity of 0.062. While this  
 150 predicted sequence identity is higher than  $\frac{1}{20} = 0.05$  due to redundant codon and nucleotide biases  
 151 favoring certain amino acids, it is much lower than the pairwise identity of even the most diverged  
 152 HAs in nature. While it is of course possible that the identity of HAs in nature would become  
 153 even lower given more time, it seems biochemically improbable that it would ever become as low  
 154 as 0.062. The reason is that like many proteins HA has a highly conserved structure and function  
 155 that imposes constraints that cause many sites to sample only a small subset of the 20 amino acids  
 156 among all known HA homologs (Nobusawa et al., 1991).

157 Accounting for site-to-site dN/dS rate variation in GY94 models affects the rate at which the  
 158 asymptotic sequence identity is approached, but not the actual value of this asymptote. For in-

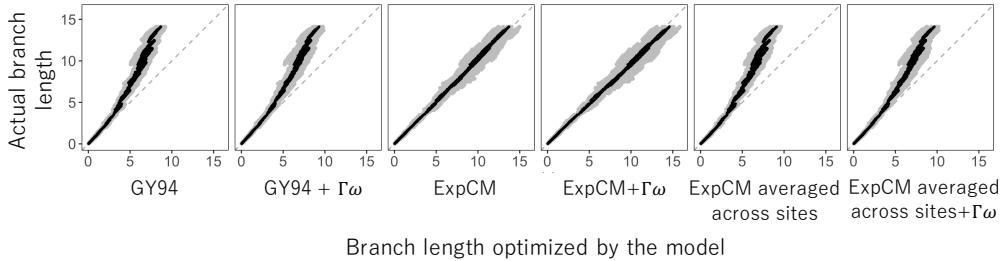
stance, [Figure 2](#) shows that the GY94+ $\Gamma\omega$  model takes longer to reach the asymptote than GY94, but that the asymptote is identical for both models. This fact holds true even if we use experimental measurements of HA’s site-specific amino-acid preferences ([Doud and Bloom, 2016](#)) to calculate a different  $\omega_r$  value for each site using the method of [Spielman and Wilke \(2015b\)](#) (see [Equation 7](#)). Specifically, this GY94+ $\omega_r$  model predicts that different sites will approach the asymptote at different rates, but the asymptote is always the same ([Figure 2](#)). The invariance of the asymptotic sequence identity under different schemes for modeling  $\omega$  is a fundamental feature of the mathematics of reversible substitution models. These models are reversible stochastic matrices, which can be decomposed into stationary states and symmetric exchangeability matrices ([Nielsen, 2006](#)). The stationary state is invariant with respect to multiplication of the symmetric exchangeability matrix by any non-zero number. Different schemes for modeling  $\omega$  only multiply elements of the symmetric exchangeability matrix. Therefore, no matter how “well” a model accounts for site-to-site variation in  $\omega$ , it will always have the same stationary state as a simple GY94 model.

However, mutation-selection models such as ExpCMs have site-specific stationary states. They predict that different sites will have different asymptotic sequence identities ([Figure 2](#))—a prediction that accords with the empirical observation that some sites are much more variable than others in alignments of highly diverged sequences. For instance, [Figure 2](#) shows that at sites such as 183 and 305 in the H1 HA, an ExpCM but not a GY94-style model predicts that the identity will always be relatively high. When sites with highly constrained amino-acid preferences such as these are common, an ExpCM can estimate a long branch length at modest sequence identities that a GY94 model might attribute to a shorter branch.

## Simulations demonstrate how failure to model site-specific amino-acid preferences leads to branch-length underestimation.

To directly demonstrate the effect of stationary state and  $\Gamma\omega$  rate variation on branch-length estimation, we tested the ability of a variety of models to accurately infer branch lengths on simulated data ([Figure 3](#)). Specifically, we simulated alignments of sequences along the HA phylogenetic tree in [Figure 4](#) using an ExpCM parameterized by the amino-acid preferences of H1 HA as experimentally measured by deep mutational scanning ([Doud and Bloom, 2016](#)). We then estimated the branch lengths from the simulated sequences using all the substitution models in [Figure 1C](#), and compared these estimates to the actual branch lengths used in the simulations. Note that these simulations closely parallel those performed by [Halpern and Bruno \(1998\)](#) and [Wertheim and Kosakovsky Pond \(2011\)](#).

The models with a uniform stationary state underestimated the lengths of long branches on the phylogenetic tree of the simulated sequences ([Figure 3](#)). The GY94 model estimated branch lengths that are ~60% of the true values for the longest branches. Accounting for site-to-site



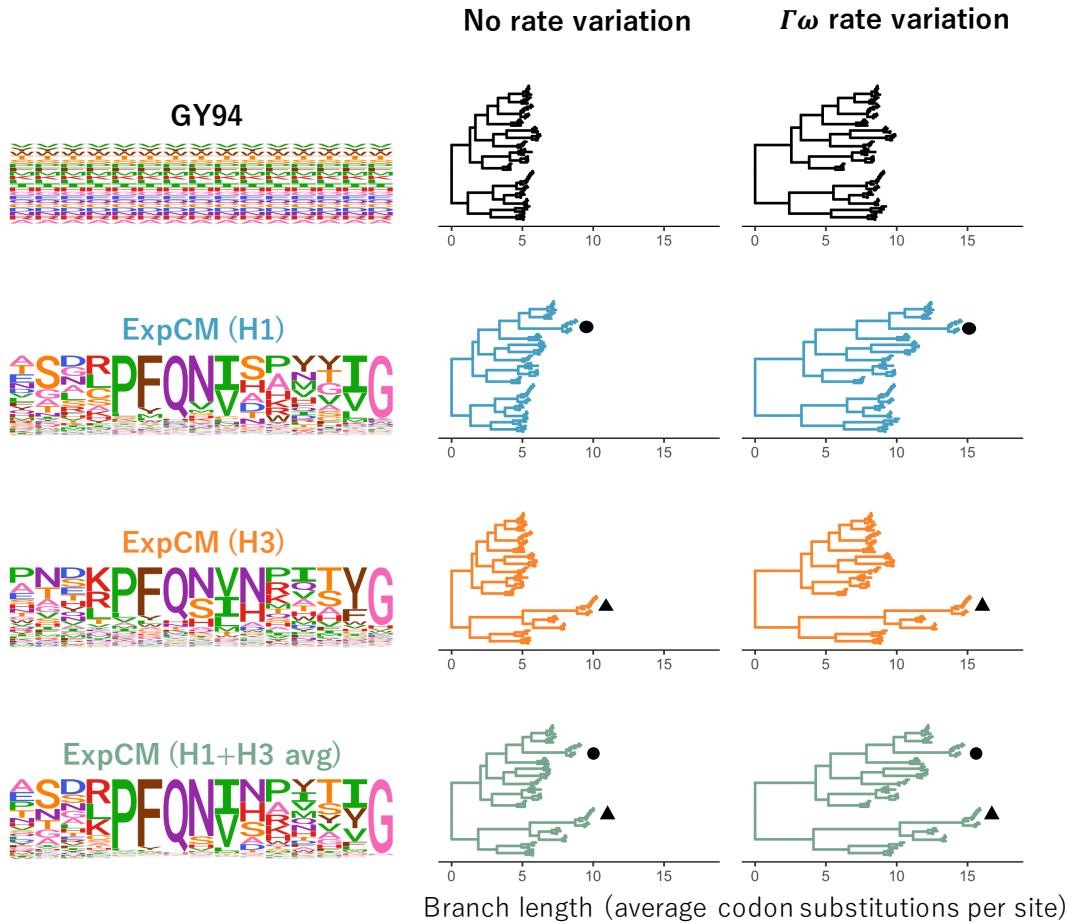
**Figure 3: Branch lengths inferred on data simulated under a model with site-specific amino-acid preferences.** We simulated alignments along a HA phylogenetic tree (see Figure 4) using an ExpCM parameterized by the actual site-specific amino-acid preferences for an H1 HA (Doud and Bloom, 2016). We then inferred the branch lengths of this tree from the simulated alignments. The inferred branch lengths for various models are plotted on the x-axis, and the actual branch lengths used in the simulations are on the y-axis. We performed 10 simulations and inferences, and gray points show each inferred branch length from each simulation, and black points show the average of each branch length across simulations. The grey dashed line at  $y = x$  represents what would be seen if the inferred branch lengths exactly matched those used in the simulations.

194 variation in  $\omega$  did not fix the fundamental problem: the GY94+ $\Gamma\omega$  did slightly better, but still sub-  
 195 substantially underestimated the longest branches. However, there was no systematic underestimation  
 196 of long branches by the ExpCM and ExpCM+ $\Gamma\omega$  models. The improved performance of the Ex-  
 197 pCMs is due to their modeling of the site-specific amino-acid preferences: if we parameterize  
 198 ExpCMs by amino-acid preferences that have been averaged across HA sites (and so are no longer  
 199 site-specific), then they perform no better than GY94 models (Figure 3). Therefore, models with  
 200 uniform stationary states underestimate the length of long branches in phylogenies of sequences  
 201 that have evolved under strong site-specific amino-acid preferences.

## 202 Experimentally informed site-specific models estimate longer branches on real data.

203 The foregoing section shows the superiority of ExpCMs to GY94 models for estimating long  
 204 branches on phylogenies simulated with ExpCMs. But how do these models perform on real data?  
 205 Real genes do evolve under functional constraint, but these constraints are almost certainly more  
 206 complex than what is modeled by an ExpCM. However, if ExpCMs do a substantially better job  
 207 than GY94 models of capturing the true constraints, then we might still expect them to estimate  
 208 more accurate branch lengths.

209 To test the models on real data, we used actual sequences of influenza HA. The topology of  
 210 HA phylogenetic trees makes these sequences an interesting test case for branch-length estima-  
 211 tion. HA consists of a number of different subtypes. Sequences within a subtype have >68%



**Figure 4: Effect of site-specific amino-acid preferences and  $\Gamma\omega$  rate variation on HA branch length estimation.** The branch lengths of the HA tree were optimized using the indicated ExpCM or GY94 model. The amino-acid preferences defining the model (ExpCM) or implied by the model (GY94) are shown as logo plots for 15 sites in HA; the full set of experimentally measured amino-acid preferences are in [Supplementary figure 1](#), [Supplementary figure 2](#), and [Supplementary figure 3](#). The ExpCMs use amino-acid preferences measured in deep mutational scanning of an H1 HA ([Doud and Bloom, 2016](#)), an H3 HA ([Lee et al., 2018](#)), or the average of the measurements for these two HAs. Circle denotes the H1 clade and triangle denotes the H3 clade. The root of each tree is placed where it would fall if the tree was midpoint rooted using the branch lengths inferred by RAxML using the GTRCAT model.

212 amino-acid identity, but sequences in different subtypes have as little as 38% identity. However,  
 213 HA proteins from all subtypes have a highly conserved structure that performs a highly conserved  
 214 function ([Ha et al., 2002](#); [Russell et al., 2004](#)). We used RAxML ([Stamatakis, 2006](#)) with a nu-  
 215 cleotide substitution model (GTRCAT) to infer a phylogenetic tree for 92 HA sequences drawn

**Table 1: Fitting of substitution models to the HA phylogenetic tree.** All ExpCMs describe the evolution of HA better than the GY94 models, as evaluated by the Akaike information criteria ( $\Delta\text{AIC}$ , [Posada and Buckley, 2004](#)). The models fit here are the same ones in [Figure 4](#). The  $\omega$  value for each of the  $K = 4$  bins is shown for the models with  $\Gamma\omega$  rate variation. All ExpCMs fit a stringency parameter  $> 1$ .

Model	$\Delta\text{AIC}$	Log Likelihood	$\omega$	Stringency parameter ( $\beta$ )
ExpCM (H1+H3 avg) + $\Gamma\omega$	0	-51083	0.19, 0.50, 0.91, 1.86	1.69
ExpCM (H1+H3 avg)	1063	-51616	0.14	1.77
ExpCM (H1) + $\Gamma\omega$	1321	-51744	0.12, 0.42, 0.89, 2.13	1.11
ExpCM (H3) + $\Gamma\omega$	1777	-51972	0.10, 0.36, 0.76, 1.84	1.28
ExpCM (H1)	2670	-52419	0.12	1.21
ExpCM (H3)	3377	-52773	0.12	1.43
GY94 + $\Gamma\omega$	4817	-53487	0.00, 0.03, 0.08, 0.24	-
GY94	7892	-55025	0.07	-

from 15 of the 18 subtypes (we excluded bat influenza and one other rare subtype). For the rest of this paper, we fix the tree topology to this RAxML-inferred tree. Although the nucleotide model used with RAxML to infer this tree topology is probably less accurate than codon models, the modular subtype structure of the HA phylogeny means that most of the phylogenetic uncertainty lies in the length of the long branches separating the subtypes rather than in the tree topology itself.

Deep mutational scanning has been used to measure the amino-acid preferences of all sites in two different HAs. One scan measured the preferences of an H1 HA ([Doud and Bloom, 2016](#)) and the other measured the preferences of an H3 HA ([Lee et al., 2018](#)). The amino-acid preferences measured for these two HAs are shown in [Supplementary figure 1](#) and [Supplementary figure 2](#). The H1 and H3 HAs have only  $\sim 42\%$  amino-acid identity, and so are separated by a large distance on the phylogenetic tree (see triangle and circle on [Figure 4](#)). As described in [Lee et al. \(2018\)](#), the amino-acid preferences clearly differ between the H1 and H3 HA at a substantial number of sites (these differences are apparent in a simple visual comparison of [Supplementary figure 1](#) and [Supplementary figure 2](#); see site 33 as an example). Therefore, we also created a third set of amino-acid preferences by averaging the measurements for the H1 and H3 HAs, under the conjecture that these averaged preferences might better describe the “average” constraint on sites across the full HA tree ([Supplementary figure 3](#)). These three sets of HA amino-acid preferences define three different ExpCMs.

We fit the GY94 model and each of the three ExpCMs to the fixed HA tree topology estimated using RAxML, and also tested a version of each model with  $\Gamma\omega$  rate variation. [Table 1](#) shows that all ExpCMs fit the actual data much better than the GY94 models. The best fit was for the ExpCM informed by the average of the H1 and H3 deep mutational scans. For all models, incorporating

238  $\Gamma\omega$  rate variation improved the fit, although even ExpCMs without  $\Gamma\omega$  greatly outperformed the  
239 GY94+ $\Gamma\omega$  model (Table 1). As mentioned in the previous section,  $\omega$  is generally  $< 1$  when a  
240 single value is fit to all sites in a gene (Murrell et al., 2015), and this is the case for all the models we  
241 tested (Table 1). However, the ExpCMs always fit an  $\omega$  greater than the GY94 model, suggesting  
242 that the site-specific amino-acid preferences capture some of the purifying selection that the GY94  
243 models can represent only via a small  $\omega$ . Among the models with  $\Gamma\omega$ , the GY94+ $\Gamma\omega$  model fits  
244 all four  $\omega$  categories to values  $\ll 1$ , but the ExpCM+ $\Gamma\omega$  models fit one of the  $\omega$  categories to a  
245 value  $> 1$ . This increase in  $\omega$  values makes sense given the different interpretation of  $\omega$  for each  
246 family of models. The ExpCM  $\omega$  is the relative rate of fixation of nonsynonymous to synonymous  
247 mutations *after* accounting for the functional constraints described by the amino-acid preferences.  
248 This more realistic null model gives ExpCMs enhanced power to detect diversifying selection for  
249 amino-acid change (Bloom, 2017; Rodrigue and Lartillot, 2017), which is known to occur at some  
250 sites in HA due to immune selection (Bedford et al., 2014).

251 Importantly, models that account for purifying selection via either  $\Gamma\omega$  rate variation or site-  
252 specific amino-acid preferences do not just exhibit better fit—they also estimate longer branches  
253 on the HA tree. Figure 4 shows the branch lengths optimized by each model on a common scale.  
254 The tree’s deepest branches are shortest when they are optimized by the GY94 model, which lacks  
255 both  $\Gamma\omega$  and site-specific amino-acid preferences. Adding either  $\Gamma\omega$  rate variation or site-specific  
256 amino-acid preferences increases the length of the deep branches. Specifically, the tree’s diameter  
257 (the distance between the two most divergent tips) for the GY94+ $\Gamma\omega$  model is 159% of the GY94  
258 model tree diameter (Supplementary table 1). The tree diameter is 122% and 135% of the GY94  
259 model tree diameter for ExpCMs informed by H1 or H3 amino-acid preferences, respectively, and  
260 160% of the GY94 model for the ExpCM informed by the average of the H1 and H3 preferences  
(Supplementary table 1).

262 The deepening of branch lengths that results from the  $\Gamma\omega$  and site-specific amino-acid pref-  
263 erence approaches to modeling purifying selection are largely independent. This can be seen by  
264 examining the ExpCM+ $\Gamma\omega$  models, which combine  $\Gamma\omega$  rate variation with site-specific amino-acid  
265 preferences. As shown in Figure 4, these ExpCM+ $\Gamma\omega$  models estimate longer branches than mod-  
266 els with just  $\Gamma\omega$  rate variation (GY94+ $\Gamma\omega$ ) or just site-specific amino-acid preferences (ExpCMs).  
267 The near independence of these effects is quantified in Supplementary table 1, which shows that  
268 76% of the tree diameter extension of ExpCM(H1+H3 avg)+ $\Gamma\omega$  versus can be explained by sim-  
269 plly adding the extension from incorporating  $\Gamma\omega$  (GY94+ $\Gamma\omega$  versus GY94) to the extension from  
270 incorporating site-specific amino-acid preferences (compare ExpCM(H1+H3 avg) to GY94).

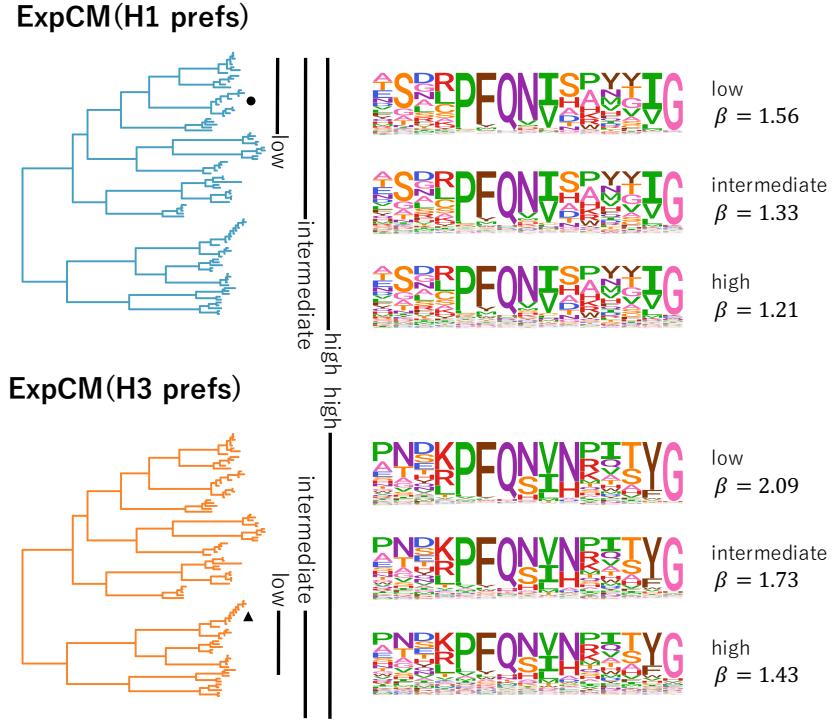
271 However, while adding  $\Gamma\omega$  rate variation increases the length of deep branches in a roughly  
272 uniform fashion across the tree, the branch lengthening from adding site-specific amino-acid pref-  
273 erences is not uniform across the tree (Figure 4) [I think we should make a Figure 5 that looks

274 similar to the current Figure 6B, that would show GY94 versus each ExpCM, GY94+gamma ver-  
275 sus each ExpCM+gamma, and each model versus it's gamma. These could be panels (one or two  
276 rows) and panel B.] Instead, the increase in branch length is most pronounced on branches leading  
277 to the HA sequence that was used in the deep mutational scanning experiment that informed the  
278 ExpCM. For instance, the ExpCM informed by the H1 data most dramatically lengthens branches  
279 near the H1 clade of the tree, while the ExpCM informed by the H3 data has the largest effect on  
280 branches near the H3 clade. The ExpCM informed by the average of the H1 and H3 data has a  
281 more uniform effect across the tree, but still most strongly extends branches leading to either the  
282 H1 or H3 clade. Therefore, Figure 4 shows that ExpCMs estimate longer branches, but that the  
283 effect is shaped by the set of amino-acid preferences used to inform the model.

284 **Shifting amino-acid preferences limit the benefits of models with site-specific station-  
285 ary states for estimating long branch lengths.**

286 The fact that an ExpCM leads to the most profound increase in branch length leading to the se-  
287 quence used in the experiment can be rationalized in terms of existing knowledge about epistasis  
288 during protein evolution. Each ExpCM is informed by a single set of experimentally measured  
289 amino-acid preferences. But in reality, the effect of a mutation at one site in a protein can depend  
290 on the amino-acid identities of other sites in the protein (Ortlund et al., 2007; Gong et al., 2013;  
291 Harms and Thornton, 2014; Tufts et al., 2014; Starr et al., 2018). This epistasis can lead to shifts in  
292 a protein's amino-acid preferences over evolutionary time (Pollock et al., 2012; Doud et al., 2015;  
293 Shah et al., 2015; Bazykin, 2015; Haddox et al., 2018). Because the deep mutational scanning  
294 experiments that inform our ExpCMs were each performed in the context of a single HA genetic  
295 background, their measurements do not account for the accumulation of epistatic shifts in amino-  
296 acid preferences as HA evolves. Therefore, an ExpCM is expected to most accurately describe the  
297 evolution of sequences closely related to the one used in the experiment.

298 We can observe how shifting amino-acid preferences degrade the accuracy of an ExpCM by  
299 fitting the model to trees containing increasingly diverged sequences. For both H1 and H3 HAs, we  
300 created three phylogenetic trees (Supplementary figure 5): a “low” divergence tree that contains  
301 sequences with  $\geq 59\%$  amino-acid identity to the HA used in the experiment, an “intermediate”  
302 divergence tree that contains sequences with  $\geq 46\%$  amino-acid identity to the HA in the exper-  
303 iment, and a “high” divergence tree that contains all HAs (which have as little as 38% identity  
304 to the HA in the experiment). Figure 5 shows the subtrees containing each of these sets of HA  
305 sequences. For each subtree, we examined the congruence between site-specific natural selection  
306 and the amino-acid preferences measured in the deep mutational scanning experiment using the  
307 ExpCM stringency parameter  $\beta$  (Bloom, 2014b; Hilton et al., 2017). Values of  $\beta$  that are  $> 1$   
308 indicate that natural selection prefers the same amino acids as the experiments but with a greater



**Figure 5: The congruence between natural selection and the deep mutational scanning measurements decreases with sequence divergence.** We fit an ExpCM informed by the H1 or H3 deep mutational scanning experiments to trees spanning sequences with low, intermediate, and high divergence from the sequence used in the experiment. The ExpCM stringency parameter ( $\beta$ ) is a measure of the congruence between natural selection and the experimental measurements (Bloom, 2014b; Hilton et al., 2017). Larger values of  $\beta$  indicate that natural selection prefers the same amino acids as the experiments but with greater stringency. As divergence increases between the HA used in the experiment and the other sequences in the tree, the  $\beta$  value decreases and the amino-acid preference “flatten.” Therefore, the preferences measured in each experiment are progressively less congruent with natural selection as we include increasingly diverged sequences.

stringency, suggesting strong congruence between natural selection and the experimental preferences. In contrast, values of  $\beta$  that are  $<1$  flatten the preferences, suggesting that they provide a relatively poor description of natural selection on the protein.

The value of  $\beta$  decreases as the divergence from the sequence used in the deep mutational scan increases Figure 5. This inverse relationship between  $\beta$  and overall divergence is seen for the ExpCMs informed by both the H1 and H3 experiments. As  $\beta$  decreases, the preferences “flatten” and so the ExpCM draws less information from the experiment. At the most extreme value of  $\beta = 0$ , the preferences would be perfectly uniform and look similar to the GY94 preferences in Figure 4. In reality,  $\beta$  never reaches a value this low, indicating the deep mutational scanning

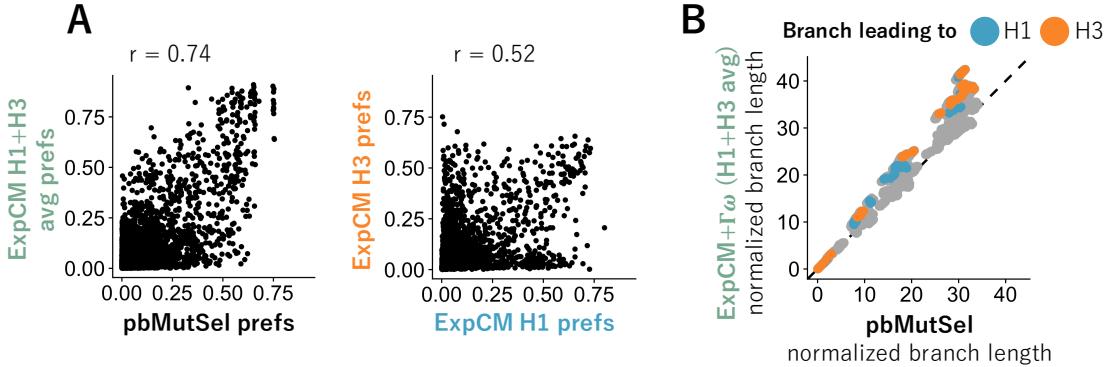
318 experiments remain somewhat informative about real natural selection across the entire swath of  
319 HAs. However, [Figure 5](#) shows that the amino-acid preferences clearly become less informative  
320 about natural selection as we move away from the experimental sequence on the tree. This shifting  
321 of amino-acid preferences helps explain why the ExpCM informed by the average of the H1 and  
322 H3 experiments performs best ([Table 1](#) and [Figure 4](#)): averaging the measurements across these  
323 two HAs is a heuristic method of accounting for shifts in preferences during HA evolution.

324 The fact that amino-acid preferences shift as a protein evolves leaves us with an inherent  
325 tension: models with site-specific amino-acid preferences only become important for accurate  
326 branch-length estimation as sequences become increasingly diverged, but this same divergence  
327 degrades the accuracy of extrapolating the amino-acid preferences from any given experiment  
328 across the phylogenetic tree. Crucially, this problem is more fundamental than the inability of a  
329 single deep mutational scanning experiment to measure amino-acid preferences in more than one  
330 genetic background. If amino-acid preferences shift during evolution, there simply will not be any  
331 single set of time-homogeneous site-specific preferences that accurately describes evolution along  
332 the entirety of a phylogenetic tree that covers a wide span of sequences.

333 **A model with amino-acid preferences estimated from natural sequences gives similar  
334 results to an ExpCM**

335 The previous sections used ExpCMs, which are mutation-selection models that use site-specific  
336 amino-acid preferences that have been measured by experiments. However, there are other math-  
337 ematically similar implementations of mutation-selection models that infer the amino-acid pref-  
338 erences directly from the natural sequence data. When these models are designed for use in phy-  
339 logenetic inference, they are generally implemented in a Bayesian framework, which avoids the  
340 overfitting problems associated with trying to make maximum-likelihood estimates of the thou-  
341 sands of amino-acid preference parameters ([Lartillot, 2014](#)). (Note that the maximum-likelihood  
342 implementations of [Tamuri et al. \(2012, 2014\)](#) are designed for estimating the amino-acid pref-  
343 erences, *not* for phylogenetic inference.) The model most comparable to our ExpCMs is the  
344 codon mutation-selection model implemented in [PhyloBayes-MPI](#), which we will refer to as  
345 pbMutSel ([Rodrigue and Lartillot, 2014](#)). In the pbMutSel model, the amino-acid preferences are  
346 modeled using Dirichlet processes rather than derived from experiments. However, like an Ex-  
347 pCM, a pbMutSel model still assumes a single set of time-homogeneous site-specific amino-acid  
348 preferences for the entire tree.

349 Comparing ExpCM and pbMutSel models can help determine the ultimate limits of mutation-  
350 selection models that assign each site a single set of amino-acid preferences. If the limitations  
351 of ExpCMs described above arise simply because the deep mutational scanning experiments do  
352 not correctly measure the “true” amino-acid preferences across the entirety of a highly diverged



**Figure 6: Models inferred from natural sequences have similar stationary states to models defined by experimental preferences and estimate similar branch lengths.** We fit an  $\text{ExpCM}(\text{H1+H3 avg})+\Gamma\omega$  and a pbMutSel to the full HA tree in Figure 4. The pbMutSel amino-acid preferences are inferred from the natural HA sequences, while the ExpCM amino-acid preferences are experimentally measured and then rescaled by the stringency parameter in Table 1. (A) The pbMutSel preferences are more correlated with the re-scaled average of the H1 and H3 deep mutational scanning preferences than the individual re-scaled H1 and H3 deep mutational scanning preferences are to each other (Pearson’s  $r$ : 0.74 versus 0.52). (B) The  $\text{ExpCM}(\text{H1+H3 avg})+\Gamma\omega$  and pbMutSel models estimated similar branch lengths when fit to the entire HA tree. Points denote branch lengths between all pairs of tips on the tree. Blue and orange denote branches that lead to the H1 and H3 deep mutational scanning reference sequences respectively. The phydms program implementing ExpCMs and the PhyloBayes-MPI program implementing pbMutSel models give branch lengths in different units, so to facilitate direct comparison between the models, we have normalized all branch lengths returned by each program by the length of the branches separating the earliest (A/South Carolina/1918) and latest (A/Solomon Islands/2006) seasonal human H1 sequences on the tree.

353 phylogenetic tree, then we would expect the pbMutSel models (which infer these preferences from  
 354 the entire tree) to perform better. On the other hand, if the major limitation is that no single set  
 355 of time-homogenous amino-acid preferences can fully describe evolution over the entire tree, then  
 356 we would expect ExpCM and pbMutSel models to perform similarly.

357 We fit a pbMutSel model to the entire HA phylogenetic tree, and compared the results to those  
 358 from analyzing the same tree with the best ExpCM, which is the  $\text{ExpCM}(\text{H1+H3 avg})+\Gamma\omega$  vari-  
 359 ant. This is a direct apple-to-apples comparison, since the pbMutSel model also draws  $\omega$  from  
 360 a gamma-distribution (Rodrigue and Lartillot, 2014). First, we compared the amino-acid prefer-  
 361 ences inferred by the pbMutSel model to the preferences measured in the experiments. Figure 6A  
 362 shows that the preferences inferred by pbMutSel are quite similar to the (H1+ H3 avg) obtained  
 363 by averaging the deep mutational scanning measurements for the H1 and H3 HAs. Notably, the  
 364 amino-acid preferences from the pbMutSel model are more correlated with the (H1+ H3 avg) than

365 the H1 and H3 measurements are with each other (Figure 6A). This strong correlation indicates  
366 that the ExpCM(H1+H3 avg)+ $\Gamma\omega$  is unlikely to be much different than a pbMutSel model that is  
367 parameterized only using the natural sequence data.

368 We next compared the branch lengths estimated by using the ExpCM(H1+H3 avg)+ $\Gamma\omega$  and  
369 pbMutSel models. As shown in Figure 6B, these two models estimated similar branch lengths  
370 across the entire HA phylogenetic tree. However, the estimates are not identical, and the tension  
371 between local and global accuracy of the amino-acid preferences is still apparent. Specifically, the  
372 branches leading to the H1 or H3 sequences used in the experiments were estimated to be slightly  
373 longer by the ExpCM, while some other branches were estimated to be slightly longer by the pb-  
374 MutSel model. The relatively longer branches leading to the experimental sequences when using  
375 the ExpCM(H1+H3 avg)+ $\Gamma\omega$  suggests that the “tree average” amino-acid preferences inferred by  
376 the pbMutSel model are not as accurate as the preferences from the deep mutational scanning for  
377 sequences close to those used in the experiments. However, for sequences distant from those used  
378 in the experiments, the “tree average” preferences inferred by the pbMutSel model appear to be  
379 slightly better than the experimental values. Therefore, while the ExpCM and pbMutSel models  
380 differ slightly in the extent to which they lengthen different branches, neither model can avoid the  
381 tension between the local and global accuracy of amino-acid preferences.

## 382 Discussion

383 We examined how estimates of deep branch lengths on phylogenetic trees are affected by account-  
384 ing for the fact that proteins prefer specific amino acids at specific sites. We did this by compar-  
385 ing inferences from models informed by experimental measurements of site-specific amino-acid  
386 preferences with more conventional codon substitution models, as well as with models that infer  
387 the amino-acid preferences from the natural sequences. We found that models that account for  
388 site-specific amino-acid preferences estimated deeper long branches, regardless of whether these  
389 preferences are measured experimentally or inferred from the sequence alignment. Additionally,  
390 we showed that the extension in branch length from site-specific amino-acid preferences is mostly  
391 independent of the extension that results from simply modeling rate variation.

392 Overall, our results underscore the importance of modeling purifying selection in a way that  
393 is more nuanced than simply allowing the substitution rate to vary across sites. Protein sites do  
394 not simply differ in their rates of substitution—different sites also prefer different amino acids.  
395 There are now two ways to account for this fact: using models informed by deep mutational  
396 scanning experiments, or using models that infer site-specific amino-acid preferences from the  
397 natural sequence alignment. Combining either type of model with rate variation increases the  
398 inferred length of deep branches relative to models that only incorporate rate variation.

399     However, assuming a single set of site-specific amino-acid preferences is still an imperfect  
400 way to model evolution over a highly diverged phylogenetic tree. In the case of the experimentally  
401 informed models, it is fairly obvious why this is true: the amino-acid preferences are measured in  
402 just one genetic background, and therefore provide only a single snapshot of preferences that shift  
403 over evolutionary time due to epistasis (Pollock et al., 2012; Shah et al., 2015; Bazykin, 2015;  
404 Haddox et al., 2018; Starr et al., 2018; Doud et al., 2015). As a result, experimentally measured  
405 amino-acid preferences are most accurate for sequences similar to the one used in the experiment,  
406 and so cause the largest increases in branch length in that region of the phylogenetic tree. However,  
407 this limitation is not unique to experimentally informed models, but is a general limitation of  
408 describing purifying selection using a single set of site-independent and time-homogenous amino-  
409 acid preferences. For instance, we showed that averaging experimental measurements on two  
410 protein homologs does a somewhat better job of capturing the “average” constraint across the tree,  
411 and performs similarly to approaches that infer the “average” preferences from natural sequence  
412 data (Rodrigue et al., 2010; Rodrigue and Lartillot, 2014). But even these “average” preferences  
413 exhibit a tradeoff between local and global accuracy for the inference of deep branch lengths.

414     So while modeling site-specific amino-acid preferences is a clear improvement over most con-  
415 ventional models, the next step towards greater accuracy will require relaxing the assumption that  
416 these preferences are time homogeneous and site independent. Of course, many authors have  
417 pointed out the shortcomings of models that fail to account for the full site-interdependent com-  
418 plexity of purifying selection (Rodrigue et al., 2005; Choi et al., 2007; Pollock et al., 2012; Gold-  
419 stein and Pollock, 2017). However, the challenge is to overcome these shortcomings with models  
420 that are tractable for real phylogenetic questions. There are two main issues: first, the Felsenstein  
421 pruning algorithm (Felsenstein, 1981) that is typically used to evaluate phylogenetic likelihoods  
422 breaks down when sites are no longer treated independently. Some alternative algorithms have  
423 been proposed (Bordner and Mittelmann, 2013; Rodrigue et al., 2009, 2005; Choi et al., 2007),  
424 but they are still in their infancy. Second, site-interdependent models require a realistic “fitness  
425 function” that describes the interactions among sites. It appears that typical structural modeling  
426 programs are insufficient for this purpose (Rodrigue et al., 2009). But hope comes from experi-  
427 mental progress in measuring actual site-interdependent constraints on proteins (Olson et al., 2014;  
428 Wu et al., 2016; Steinberg and Ostermeier, 2016; Li et al., 2016), combined with new methods for  
429 using these measurements to parameterize fitness functions (Sailer and Harms, 2017; Otwinowski  
430 et al., 2018; Otwinowski, 2018). Perhaps some day such truly realistic models might be useful  
431 for phylogenetic inference. Until that day, our work shows that modeling a single set of time  
432 homogenous amino-acid preferences provides at least some improvement.

433 **Methods**

434 **Substitution models**

435 All of the substitution models used in this paper have been described previously. However, here  
 436 we briefly recap their exact mathematical implementations.

437 **GY94 model**

The GY94 model is M0 variant of the Goldman-Yang model described by [Yang et al. \(2000\)](#). Specifically, the substitution rate  $P_{xy}$  from codon  $x$  to codon  $y$  is

$$P_{xy} = \begin{cases} 0 & \text{if } x \text{ and } y \text{ differ by more than one nucleotide,} \\ \Phi_y & \text{if } \mathcal{A}(x) = \mathcal{A}(y) \text{ and } x \text{ is converted to } y \text{ by a single-nucleotide transversion,} \\ \omega\Phi_y & \text{if } \mathcal{A}(x) \neq \mathcal{A}(y) \text{ and } x \text{ is converted to } y \text{ by a single-nucleotide transversion,} \\ \kappa\Phi_y & \text{if } \mathcal{A}(x) = \mathcal{A}(y) \text{ and } x \text{ is converted to } y \text{ by a single-nucleotide transition,} \\ \omega\kappa\Phi_y & \text{if } \mathcal{A}(x) \neq \mathcal{A}(y) \text{ and } x \text{ is converted to } y \text{ by a single-nucleotide transition,} \\ -\sum_{z \neq x} P_{xz} & \text{if } x = y, \end{cases} \quad (\text{Equation 1})$$

438 where  $\mathcal{A}(x)$  is the amino-acid encoded by codon  $x$ ,  $\kappa$  is the transition-transversion rate,  $\Phi_y$  is the  
 439 equilibrium frequency of codon  $y$ , and  $\omega$  is the relative rate of nonsynonymous and synonymous  
 440 substitutions. We define the codon frequency parameters,  $\Phi_y$ , using the “corrected F3X4” method  
 441 from [Pond et al. \(2010\)](#). There are nine parameters describing the nucleotide frequencies at each  
 442 codon site (the nucleotides are constrained to sum to one at each codon position), and these pa-  
 443 rameter values are calculated from the empirical alignment frequencies. The “corrected F3X4”  
 444 method calculates the  $\Phi_y$  values from these nucleotide frequencies but corrects for the exclusion  
 445 of sequences with premature stop codons from the analysis.

The frequency  $p_x$  of codon  $x$  in the stationary state of a GY94 model is simply

$$p_x = \Phi_x. \quad (\text{Equation 2})$$

446 Overall, a GY94 model has 11 free parameters:  $\kappa$ ,  $\omega$ , and the 9 nucleotide frequency parameters  
 447 used to define  $\Phi_y$ .

<sup>448</sup> **Experimentally Informed Codon Model (ExpCM)**

The ExpCM models used in this paper are the ones described in Bloom (2017). Briefly, the rate of substitution  $P_{r,xy}$  of site  $r$  from codon  $x$  to  $y$  is

$$P_{r,xy} = Q_{xy} \times F_{r,xy} \quad (\text{Equation 3})$$

<sup>449</sup> where  $Q_{xy}$  is proportional to the rate of mutation from  $x$  to  $y$ ,  $F_{r,xy}$  is proportional to the probability that this mutation fixes, and the diagonal elements  $P_{xx}$  are set by  $P_{xx} = -\sum_{z \neq x} P_{xz}$ .

The rate of mutation  $Q_{xy}$  is assumed to be uniform across sites, and takes an HKY85-like (Hasegawa et al., 1985) form as

$$Q_{xy} = \begin{cases} 0 & \text{if } x \text{ and } y \text{ differ by more than one nucleotide,} \\ \phi_w & \text{if } x \text{ can be converted to } y \text{ by a transversion of a nucleotide to } w, \\ \kappa \times \phi_w & \text{if } x \text{ can be converted to } y \text{ by a transition of a nucleotide to } w \end{cases} \quad (\text{Equation 4})$$

<sup>451</sup> where  $\phi_w$  is the nucleotide frequency of nucleotide  $w$  and  $\kappa$  is the transition-transversion rate.

The deep mutational scanning amino-acid preferences are incorporated into the ExpCM via the  $F_{r,xy}$  terms. The experiments measure the preference  $\pi_{r,a}$  of every site  $r$  for every amino-acid  $a$ .  $F_{r,xy}$  is defined in terms of these experimentally measured amino-acid preferences as

$$F_{r,xy} = \begin{cases} 1 & \text{if } \mathcal{A}(x) = \mathcal{A}(y), \\ \omega \times \frac{\ln[(\pi_{r,\mathcal{A}(y)} / \pi_{r,\mathcal{A}(x)})^\beta]}{1 - (\pi_{r,\mathcal{A}(x)} / \pi_{r,\mathcal{A}(y)})^\beta} & \text{if } \mathcal{A}(x) \neq \mathcal{A}(y), \end{cases} \quad (\text{Equation 5})$$

<sup>452</sup> where  $\beta$  is the stringency parameter (Bloom, 2014b; Hilton et al., 2017) and  $\omega$  is the relative rate of nonsynonymous to synonymous substitutions after accounting for the amino-acid preferences.

The stationary state of an ExpCM is

$$p_{r,x} = \frac{\phi_{x_1} \phi_{x_2} \phi_{x_3} (\pi_{r,A(x)})^\beta}{\sum_z \phi_{z_1} \phi_{z_2} \phi_{z_3} (\pi_{r,A(z)})^\beta} \quad (\text{Equation 6})$$

<sup>454</sup> where  $\phi_{x_1}$ ,  $\phi_{x_2}$ , and  $\phi_{x_3}$  are the nucleotides at position 1, 2, and 3 of codon  $x$ .

<sup>455</sup> An ExpCM has five free parameters:  $\kappa$ ,  $\omega$ , and the three independent  $\phi_x$  values. The amino-acid preferences  $\pi_{r,a}$  are *not* free parameters since they are determined *a priori* by an experiment independent of the sequence alignment being analyzed.

458  **$\Gamma\omega$  rate variation**

459 The GY94+ $\Gamma\omega$  is equivalent to the M5 model in [Yang et al. \(2000\)](#) with  $\omega$  drawn from  $K = 4$   
 460 categories. The ExpCM+ $\Gamma\omega$  similarly draws  $\omega$  from a  $\Gamma$  distribution discretized into  $K = 4$  bins.  
 461 Each bin is equally weighted and  $\omega$  takes on the mean value of the bin. Because the  $\Gamma$  distribution  
 462 is defined by two parameters, adding  $\Gamma\omega$  to a model with a single  $\omega$  adds one free parameter.  
 463 Therefore, the GY94+ $\Gamma\omega$  model has 12 free parameters, and the ExpCM+ $\Gamma\omega$  model has 6 free  
 464 parameters.

465 **GY94 with  $\omega_r$**

In [Figure 2](#), we describe GY94 models where each site  $r$  has its own  $\omega_r$  value that is calculated  
 from the amino-acid preferences using the relationship described by [Spielman and Wilke \(2015b\)](#).  
 This relationship defines the expected rate of nonsynonymous to synonymous substitutions given  
 the amino-acid preferences. We first fit an ExpCM to the “low divergence” H1 subtree (parameter  
 values in [Supplementary table 2](#)), which allows us to calculate  $P_{r,xy}$  ([Equation 3](#)),  $Q_{xy}$  ([Equation 4](#)), and  $p_{r,x}$  ([Equation 6](#)). We then calculated  $\omega_r$  using the equation of [Spielman and Wilke \(2015b\)](#), normalizing by the gene-wide  $\omega$  fit by the ExpCM:

$$\omega_r = \frac{\sum_x \sum_{y \in N_x} p_{r,x} \times \frac{P_{r,xy}}{\omega}}{\sum_x \sum_{y \in N_x} p_{r,x} \times Q_{xy}}, \quad (\text{Equation 7})$$

466 where  $N_x$  is the set of codons that are nonsynonymous to codon  $x$  and differ from codon  $x$  by only  
 467 one nucleotide.

468 **HA amino-acid preferences from deep mutational scanning experiments**

We used amino-acid preferences measured in deep mutational scans of the A/WSN/1933 H1  
 HA ([Doud and Bloom, 2016](#)) and the A/Perth/2009 H3 HA ([Lee et al., 2018](#)) to define the amino-  
 acid preferences that inform the ExpCMs. We only used sites that can be unambiguously aligned  
 in these H1 and H3 HAs. These alignable sites and their mapping to sequential numbering of the  
 HA sequences used in the deep mutational scanning experiments are in [Supplementary file 1](#). The  
 experimentally measured amino-acid preferences masked to just include these alignable sites are  
 in [Supplementary file 2](#) and [Supplementary file 3](#). For the average preference set, we took the  
 pairwise average of the H1 and H3 preferences. The preference for every amino acid  $a$  at every  
 site  $r$  in the average preference set is

$$\pi_{r,a,(\text{H1+H3 avg})} = \frac{\pi_{r,a,\text{H1}} + \pi_{r,a,\text{H3}}}{2} \quad (\text{Equation 8})$$

469 **HA sequences and tree topology**

470 We downloaded all full-length, coding sequences for 15 of the 18 influenza A virus HA subtypes  
471 from the Influenza Virus Resource Database (Bao et al., 2008) in June of 2017. We excluded rare  
472 subtypes 15, 17, and 18, which have few sequences in the database. We filtered and aligned  
473 the sequences using phydms\_prepalignment (Hilton et al., 2017). Specifically, we used  
474 phydms\_prepalignment with the flag --minidentity 0.3 to remove sequences with  
475 ambiguous nucleotides, premature stops, or frameshift mutations as well as redundant sequences.  
476 We also removed all codon sites which are not alignable between the H1 HA and H3 HA  
477 used in the deep mutational scanning experiments (these alignable sites are listed in [Supplementary file 1](#)). We subsampled the remaining sequences to five per subtype with  $\leq 1$  sequence per  
478 year per subtype. We also included a small number of sequences from the major human and equine  
479 influenza lineages to ensure representation of these well-studied lineages. The resulting alignment  
480 contains 92 sequences, and is provided in [Supplementary file 4](#).

482 We created four subalignments with “low” and “intermediate” divergence from either the H1  
483 or the H3 deep mutational scanning reference sequence for the analysis in [Figure 5](#). The “low di-  
484 vergence” alignments had  $\geq 59\%$  amino-acid identity to the sequence used in the deep mutational  
485 scanning, and the “intermediate divergence” alignments had  $\geq 46\%$  identity from the reference  
486 sequence ([Supplementary figure 5](#)).

487 We inferred the tree topology of each alignment using RAxML (Stamatakis, 2006) and the  
488 GTRCAT model. We estimated the branch lengths of this fixed topology using each ExpCM and  
489 GY94 models with phydms\_comprehensive (Hilton et al., 2017).

490 **Asymptotic amino-acid sequence identity**

491 For the analysis in [Figure 2](#), we fit models to the “low divergence” H1 subtree. This gave the  
492 parameter values in [Supplementary table 2](#).

For each model, we calculated the expected amino-acid sequence identity for two sequences  
separated by a branch length of  $t$  as

$$\sum_a \sum_{x \in a} p_{r,x} \sum_{y \in a} [e^{tP_r}]_{xy} \quad (\text{Equation 9})$$

493 where  $a$  ranges over all 20 amino acids,  $x \in a$  indicates that  $x$  ranges over all codons that encode  
494 amino-acid  $a$ ,  $p_{r,x}$  is the stationary state of the model at site  $r$  and codon  $x$  (given by [Equation 2](#)  
495 for GY94-family models, and [Equation 6](#) for ExpCM-family models), and  $[e^{tP_r}]_{xy}$  is the value in  
496 row  $x$  and column  $y$  of the matrix obtained by exponentiating the product of  $t$  and the substitution  
497 matrix  $P_r$  for site  $r$  (defined by [Equation 1](#) for GY94-family models and [Equation 3](#) for ExpCM-

498 family models).

499 **Simulations**

500 For Figure 3, we simulated sequences using `pyvolve` (Spielman and Wilke, 2015a) along the  
501 full HA tree using an ExpCM defined by parameters fit to the “low divergence” H1 subtree (Sup-  
502 [plementary table 2](#)). We performed 10 replicate simulations and estimated the branch lengths for  
503 each replicate using `phydms_comprehensive` (Hilton et al., 2017).

504 **pbMutSel inference with PhyloBayes-MPI.**

505 For Figure 6, we fit a pbMutSel model to the full HA tree. We ran one chain for 5500 steps, saved  
506 every sample, and discarded the first 550 samples as a burnin. We used PhyloBayes-MPI  
507 program `readpb_mpi` to compute the majority-rule consensus tree and the posterior average  
508 site-specific amino-acid preferences. Convergence was assessed visually using Tracer (<http://tree.bio.ed.ac.uk/software/tracer/>) and by the correlation of amino-acid pref-  
509 erences inferred by two independent chains ( $r=0.996$ ).

510 In order to make the branch lengths in Figure 6 comparable between the pbMutSel tree re-  
511 turned by PhyloBayes-MPI and the other trees returned by `phydms`, we normalized the  
512 branch lengths on the pbMutSel consensus tree and the ExpCM(H1+H3 avg)+ $\Gamma\omega$  by dividing each  
513 branch by the length from A/South Carolina/1/1918 and A/Solomon Islands/3/2006. These two  
514 H1 sequences are early and late representatives of the longest known human influenza lineage, and  
515 are of sufficiently high identity that different ExpCM and GY94 substitution models all estimate  
516 nearly identical branch lengths separating them.

518 **Software versions and computer code**

519 All code used for the analyses in this paper is available at [https://github.com/jbloomlab/divergence\\_timing\\_manuscript](https://github.com/jbloomlab/divergence_timing_manuscript). The external computer programs that we used were

- 521 • `phydms` (Hilton et al., 2017) version 2.2.2 (available at [github.com/jbloomlab/phydms](https://github.com/jbloomlab/phydms)) to fit the ExpCM and GY94 models.
- 523 • `pyvolve` (Spielman and Wilke, 2015a) version 0.8.7 (available at <https://github.com/sjspielman/pyvolve>) to simulate the sequences.
- 525 • PhyloBayes-MPI (Rodrigue and Lartillot, 2014) version 1.8 (available at <https://github.com/bayesiancook/pbmpi>) to fit the pbMutSel model.

- 527     ● RAxML ([Stamatakis, 2006](#)) version 8.2.11 (available at <https://github.com/stamatak/standard-RAxML>) to infer tree topology.
- 528
- 529     ● We used ggplot2 ([Wickham, 2016](#)), ggtree ([Yu et al., 2017](#)), and ggseqlogo ([Wagih, 2017](#)) for visualization of the results.
- 530
- 531     ● snakemake ([Köster and Rahmann, 2012](#)) version 3.11.2 (available at <https://snakemake.readthedocs.io/en/stable/>) to run the pipelines.
- 532

533 **Acknowledgments**

534 We thank Erick Matsen and Trevor Bedford for helpful comments about the project and manuscript.

535 SKH is supported in part by training grant T32AI083203 from the NIAID of the National Institutes

536 of Health. This work was supported by the NIAID and NIGMS of the NIH under grant numbers

537 R01AI127893 and R01GM102198. JDB is supported in part by a Faculty Scholars grant from the

538 Howard Hughes Medical Institute and the Simons Foundation. The funders had no role in study

539 design, data collection and analysis, decision to publish, or preparation of the manuscript.

540 **References**

- 541 Aiewsakun P, Katzourakis A. 2016. Time-dependent rate phenomenon in viruses. *Journal of*  
542 *Virology*. 90:7184–7195.
- 543 Arenas M. 2015. Trends in substitution models of molecular evolution. *Frontiers in Genetics*.  
544 6:319.
- 545 Bao Y, Bolotov P, Dernovoy D, Kiryutin B, Zaslavsky L, Tatusova T, Ostell J, Lipman D. 2008.  
546 The influenza virus resource at the National Center for Biotechnology Information. *Journal of*  
547 *Virology*. 82:596–601.
- 548 Bazykin GA. 2015. Changing preferences: deformation of single position amino acid fitness  
549 landscapes and evolution of proteins. *Biology Letters*. 11:20150315.
- 550 Bedford T, Suchard MA, Lemey P, Dudas G, Gregory V, Hay AJ, McCauley JW, Russell CA,  
551 Smith DJ, Rambaut A. 2014. Integrating influenza antigenic dynamics with molecular evolution.  
552 *eLife*. 3:e01914.
- 553 Bloom JD. 2014a. An experimentally determined evolutionary model dramatically improves phy-  
554 logenetic fit. *Molecular Biology and Evolution*. 31:1956–1978.
- 555 Bloom JD. 2014b. An experimentally informed evolutionary model improves phylogenetic fit to  
556 divergent lactamase homologs. *Mol. Biol. Evol.* 31:2753–2769.
- 557 Bloom JD. 2017. Identification of positive selection in genes is greatly improved by using experi-  
558 mentally informed site-specific models. *Biology Direct*. 12:1.
- 559 Bordner AJ, Mittelmann HD. 2013. A new formulation of protein evolutionary models that ac-  
560 count for structural constraints. *Molecular Biology and Evolution*. 31:736–749.
- 561 Carroll SA, Towner JS, Sealy TK, McMullan LK, Khristova ML, Burt FJ, Swanepoel R, Rollin  
562 PE, Nichol ST. 2013. Molecular evolution of viruses of the family filoviridae based on 97  
563 whole-genome sequences. *Journal of Virology*. 87:2608–2616.
- 564 Choi SC, Hobolth A, Robinson DM, Kishino H, Thorne JL. 2007. Quantifying the impact of  
565 protein tertiary structure on molecular evolution. *Molecular Biology and Evolution*. 24:1769–  
566 1782.
- 567 Doud MB, Ashenberg O, Bloom JD. 2015. Site-specific amino acid preferences are mostly con-  
568 served in two closely related protein homologs. *Mol. Biol. Evol.* 32:2944–2960.

- 569 Doud MB, Bloom JD. 2016. Accurate measurement of the effects of all amino-acid mutations to  
570 influenza hemagglutinin. *Viruses*. 8:155.
- 571 Drummond AJ, Ho SY, Phillips MJ, Rambaut A. 2006. Relaxed phylogenetics and dating with  
572 confidence. *PLoS Biology*. 4:e88.
- 573 Duchêne DA, Duchêne S, Holmes EC, Ho SY. 2015a. Evaluating the adequacy of molecular clock  
574 models using posterior predictive simulations. *Molecular Biology and Evolution*. 32:2986–  
575 2995.
- 576 Duchêne S, Di Giallondo F, Holmes EC. 2015b. Substitution model adequacy and assessing  
577 the reliability of estimates of virus evolutionary rates and time scales. *Molecular Biology and*  
578 *Evolution*. 33:255–267.
- 579 Duchêne S, Holmes EC, Ho SY. 2014. Analyses of evolutionary dynamics in viruses are hindered  
580 by a time-dependent bias in rate estimates. *Proc. R. Soc. B*. 281:20140732.
- 581 Echave J, Spielman SJ, Wilke CO. 2016. Causes of evolutionary rate variation among protein  
582 sites. *Nature Reviews Genetics*. .
- 583 Fares MA, Holmes EC. 2002. A revised evolutionary history of hepatitis b virus (HBV). *Journal*  
584 *of Molecular Evolution*. 54:807–814.
- 585 Felsenstein J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach.  
586 *J. Mol. Evol.* 17:368–376.
- 587 Fowler DM, Fields S. 2014. Deep mutational scanning: a new style of protein science. *Nature*  
588 *Methods*. 11:801–807.
- 589 Furuse Y, Suzuki A, Oshitani H. 2010. Origin of measles virus: divergence from rinderpest virus  
590 between the 11th and 12th centuries. *Virology journal*. 7:52.
- 591 Goldman N, Yang Z. 1994. A codon-based model of nucleotide substitution for protein-coding  
592 dna sequences. *Molecular Biology and Evolution*. 11:725–736.
- 593 Goldstein RA, Pollock DD. 2017. Sequence entropy of folding and the absolute rate of amino acid  
594 substitutions. *Nature Ecology & Evolution*. 1:1923.
- 595 Gong LI, Suchard MA, Bloom JD. 2013. Stability-mediated epistasis constrains the evolution of  
596 an influenza protein. *eLife*. 2:e00631.
- 597 Ha Y, Stevens DJ, Skehel JJ, Wiley DC. 2002. H5 avian and H9 swine influenza virus haemagglu-  
598 tinin structures: possible origin of influenza subtypes. *The EMBO journal*. 21:865–875.

- 599 Haddox HK, Dingens AS, Hilton SK, Overbaugh J, Bloom JD. 2018. Mapping mutational effects  
600 along the evolutionary landscape of HIV envelope. *eLife*. 7:e34420.
- 601 Halpern AL, Bruno WJ. 1998. Evolutionary distances for protein-coding sequences: modeling  
602 site-specific residue frequencies. *Molecular Biology and Evolution*. 15:910–917.
- 603 Harms MJ, Thornton JW. 2014. Historical contingency and its biophysical basis in glucocorticoid  
604 receptor evolution. *Nature*. 512:203–207.
- 605 Hasegawa M, Kishino H, Yano Ta. 1985. Dating of the human-ape splitting by a molecular clock  
606 of mitochondrial DNA. *Journal of Molecular Evolution*. 22:160–174.
- 607 Hilton SK, Doud MB, Bloom JD. 2017. phydms: Software for phylogenetic analyses informed by  
608 deep mutational scanning. *PeerJ*. 5:e3657.
- 609 Ho SY, Duchêne S, Molak M, Shapiro B. 2015. Time-dependent estimates of molecular evolu-  
610 tionary rates: evidence and causes. *Molecular Ecology*. 24:6007–6012.
- 611 Holmes EC. 2003. Molecular clocks and the puzzle of RNA virus origins. *Journal of Virology*.  
612 77:3893–3897.
- 613 Köster J, Rahmann S. 2012. Snakemake—a scalable bioinformatics workflow engine. *Bioinfor-  
614 matics*. 28:2520–2522.
- 615 Lartillot N. 2014. The Bayesian Kitchen: overcoming the fear of over-  
616 parameterization. [http://bayesiancook.blogspot.com/2014/01/  
617 the-myth-of-over-parameterization.html](http://bayesiancook.blogspot.com/2014/01/the-myth-of-over-parameterization.html). Last accessed: March-12-2018.
- 618 Lartillot N, Brinkmann H, Philippe H. 2007. Suppression of long-branch attraction artefacts in the  
619 animal phylogeny using a site-heterogeneous model. *BMC Evolutionary Biology*. 7:S4.
- 620 Lartillot N, Philippe H. 2004. A Bayesian mixture model for across-site heterogeneities in the  
621 amino-acid replacement process. *Molecular Biology and Evolution*. 21:1095–1109.
- 622 Le SQ, Lartillot N, Gascuel O. 2008. Phylogenetic mixture models for proteins. *Phil. Trans. R.  
623 Soc. B*. 363:3965–3976.
- 624 Lee JM, Huddleston J, Doud MB, Hooper K, Wu NC, Bedford T, Bloom JD. 2018. Deep mu-  
625 tational scanning of hemagglutinin helps predict evolutionary fates of human H3N2 influenza  
626 variants. *bioRxiv*. DOI: 10.1101/298364.
- 627 Li C, Qian W, Maclean CJ, Zhang J. 2016. The fitness landscape of a trna gene. *Science*. 352:837–  
628 840.

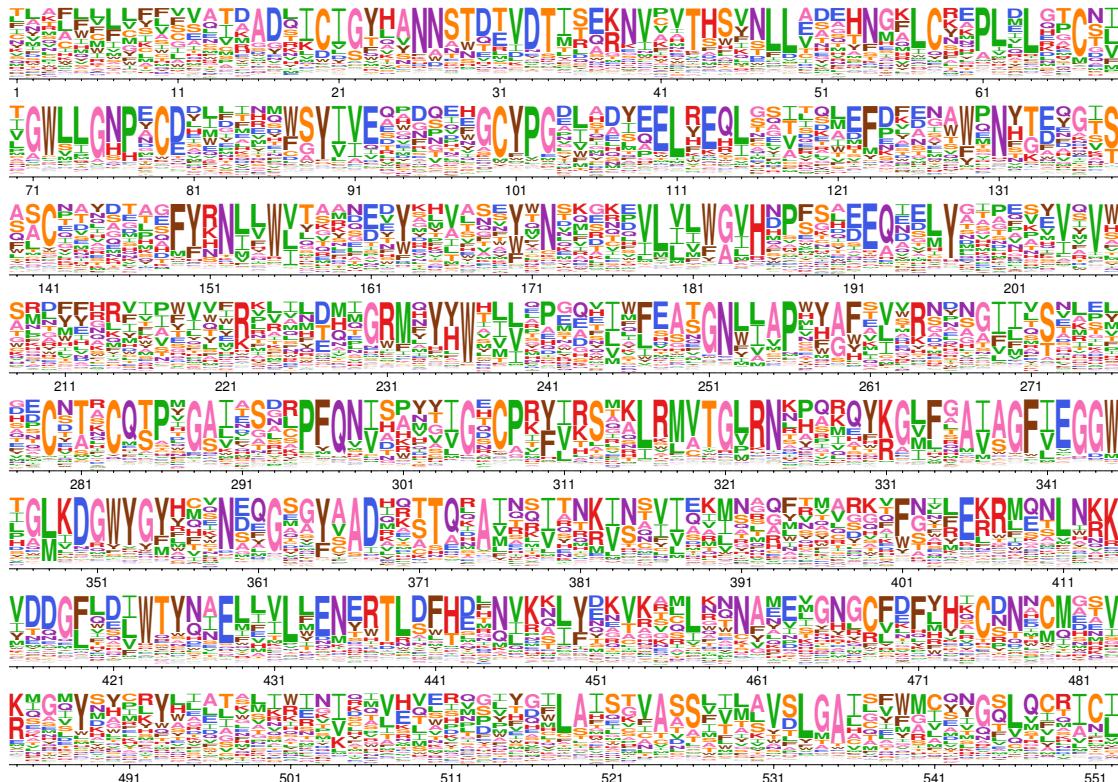
- 629 McCandlish DM, Stoltzfus A. 2014. Modeling evolution using the probability of fixation: history  
630 and implications. *The Quarterly Review of Biology*. 89:225–252.
- 631 Murrell B, Weaver S, Smith MD, et al. (11 co-authors). 2015. Gene-wide identification of episodic  
632 selection. *Molecular Biology and Evolution*. 32:1365–1371.
- 633 Nielsen R. 2006. Statistical methods in molecular evolution. Springer.
- 634 Nobusawa E, Aoyama T, Kato H, Suzuki Y, Tateno Y, Nakajima K. 1991. Comparison of complete  
635 amino acid sequences and receptor-binding properties among 13 serotypes of hemagglutinins  
636 of influenza A viruses. *Virology*. 182:475–485.
- 637 Olson CA, Wu NC, Sun R. 2014. A comprehensive biophysical description of pairwise epistasis  
638 throughout an entire protein domain. *Current Biology*. 24:2643–2651.
- 639 Ortlund EA, Bridgham JT, Redinbo MR, Thornton JW. 2007. Crystal structure of an ancient  
640 protein: evolution by conformational epistasis. *Science*. 317:1544–1548.
- 641 Otwinowski J. 2018. Inferring protein stability and function from a high-throughput assay. *arXiv*.  
642 1802.08744.
- 643 Otwinowski J, McCandlish DM, Plotkin J. 2018. Inferring the shape of global epistasis. *bioRxiv*.  
644 DOI: 10.1101/278630.
- 645 Philippe H, Laurent J. 1998. How good are deep phylogenetic trees? *Current Opinion in Genetics  
& Development*. 8:616–623.
- 647 Pollock DD, Thiltgen G, Goldstein RA. 2012. Amino acid coevolution induces an evolutionary  
648 stokes shift. *Proc. Natl. Acad. Sci. USA*. 109:E1352–E1359.
- 649 Pond SK, Delport W, Muse SV, Scheffler K. 2010. Correcting the bias of empirical frequency  
650 parameter estimators in codon models. *PLoS One*. 5:e11230.
- 651 Posada D, Buckley TR. 2004. Model selection and model averaging in phylogenetics: advantages  
652 of Akaike information criterion and Bayesian approaches over likelihood ratio tests. *Systematic  
Biology*. 53:793–808.
- 654 Quang SL, Gascuel O, Lartillot N. 2008. Empirical profile mixture models for phylogenetic re-  
655 construction. *Bioinformatics*. 24:2317–2323.
- 656 Rodrigue N. 2013. On the statistical interpretation of site-specific variables in phylogeny-based  
657 substitution models. *Genetics*. 193:557–564.

- 658 Rodrigue N, Kleinman CL, Philippe H, Lartillot N. 2009. Computational methods for evaluating  
659 phylogenetic models of coding sequence evolution with dependence between codons. *Mol. Biol.*  
660 *Evol.* 26:1663–1676.
- 661 Rodrigue N, Lartillot N. 2014. Site-heterogeneous mutation-selection models within the  
662 PhyloBayes-MPI package. *Bioinformatics*. 30:1020–1021.
- 663 Rodrigue N, Lartillot N. 2017. Detecting adaptation in protein-coding genes using a Bayesian site-  
664 heterogeneous mutation-selection codon substitution model. *Molecular Biology and Evolution*.  
665 34:204–214.
- 666 Rodrigue N, Lartillot N, Bryant D, Philippe H. 2005. Site interdependence attributed to tertiary  
667 structure in amino acid sequence evolution. *Gene*. 347:207–217.
- 668 Rodrigue N, Philippe H, Lartillot N. 2010. Mutation-selection models of coding sequence evolu-  
669 tion with site-heterogeneous amino acid fitness profiles. *Proceedings of the National Academy*  
670 *of Sciences*. 107:4629–4634.
- 671 Russell R, Gamblin S, Haire L, Stevens D, Xiao B, Ha Y, Skehel J. 2004. H1 and H7 influenza  
672 haemagglutinin structures extend a structural classification of haemagglutinin subtypes. *Virol-*  
673 *ogy*. 325:287–296.
- 674 Sailer ZR, Harms MJ. 2017. Detecting high-order epistasis in nonlinear genotype-phenotype  
675 maps. *Genetics*. 205:1079–1088.
- 676 Shah P, McCandlish DM, Plotkin JB. 2015. Contingency and entrenchment in protein evolu-  
677 tion under purifying selection. *Proceedings of the National Academy of Sciences*. 112:E3226–  
678 E3235.
- 679 Spielman SJ, Wilke CO. 2015a. Pyvolve: a flexible Python module for simulating sequences along  
680 phylogenies. *PLoS One*. 10:e0139047.
- 681 Spielman SJ, Wilke CO. 2015b. The relationship between dN/dS and scaled selection coefficients.  
682 *Molecular Biology and Evolution*. 32:1097–1108.
- 683 Stamatakis A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with  
684 thousands of taxa and mixed models. *Bioinformatics*. 22:2688–2690.
- 685 Starr TN, Flynn JM, Mishra P, Bolon DNA, Thornton JW. 2018. Pervasive contingency and  
686 entrenchment in a billion years of Hsp90 evolution. *Proceedings of the National Academy of*  
687 *Sciences*. DOI: 10.1073/pnas.1718133115.

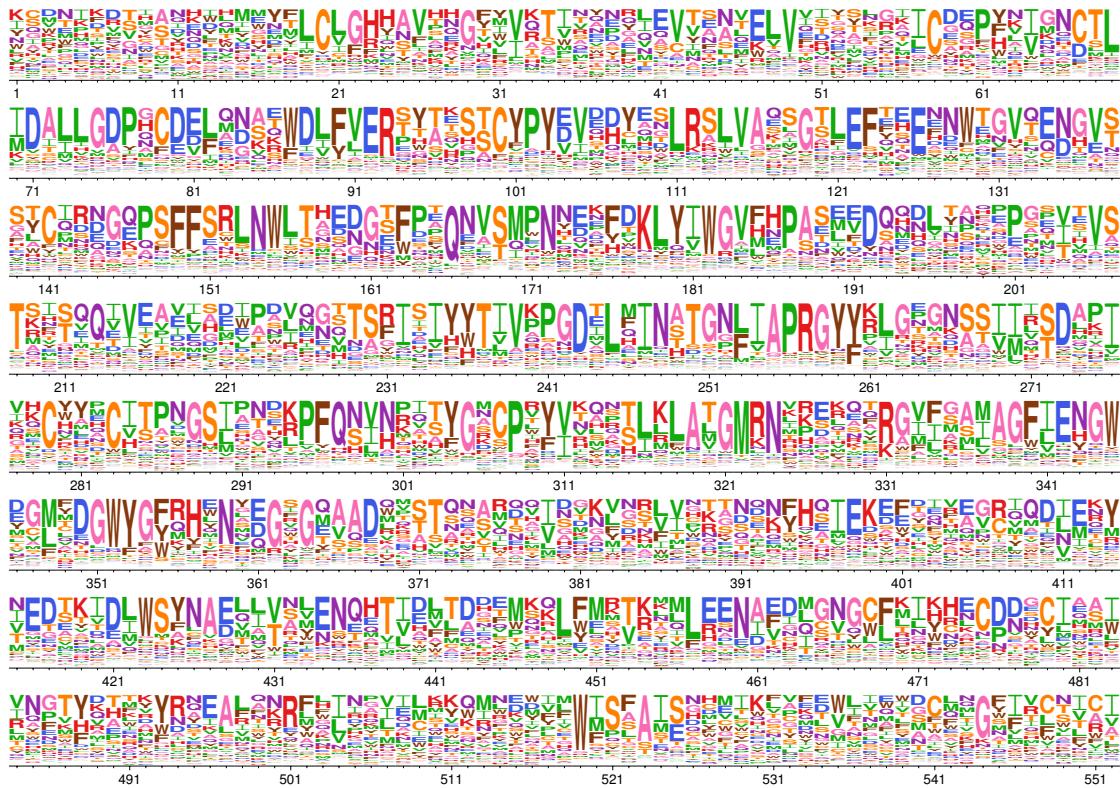
- 688 Steinberg B, Ostermeier M. 2016. Shifting fitness and epistatic landscapes reflect trade-offs along  
689 an evolutionary pathway. *Journal of Molecular Biology*. 428:2730–2743.
- 690 Tamuri AU, dos Reis M, Goldstein RA. 2012. Estimating the distribution of selection coefficients  
691 from phylogenetic data using sitewise mutation-selection models. *Genetics*. 190:1101–1115.
- 692 Tamuri AU, Goldman N, dos Reis M. 2014. A penalized likelihood method for estimating the  
693 distribution of selection coefficients from phylogenetic data. *Genetics*. pp. genetics–114.
- 694 Taylor DJ, Ballinger MJ, Zhan JJ, Hanzly LE, Bruenn JA. 2014. Evidence that Ebolaviruses and  
695 Cuevaviruses have been diverging from Marburgviruses since the Miocene. *PeerJ*. 2:e556.
- 696 Tufts DM, Natarajan C, Revsbech IG, Projecto-Garcia J, Hoffmann FG, Weber RE, Fago A,  
697 Moriyama H, Storz JF. 2014. Epistasis constrains mutational pathways of hemoglobin adap-  
698 tation in high-altitude pikas. *Molecular Biology and Evolution*. 32:287–298.
- 699 Wagih O. 2017. ggseqlogo: a versatile R package for drawing sequence logos. *Bioinformatics*.  
700 33:3645–3647.
- 701 Wang HC, Li K, Susko E, Roger AJ. 2008. A class frequency mixture model that adjusts for  
702 site-specific amino acid frequencies and improves inference of protein phylogeny. *BMC Evolu-  
703 tionary Biology*. 8:331.
- 704 Wertheim JO, Chu DK, Peiris JS, Pond SLK, Poon LL. 2013. A case for the ancient origin of  
705 coronaviruses. *Journal of Virology*. 87:7039–7045.
- 706 Wertheim JO, Kosakovsky Pond SL. 2011. Purifying selection can obscure the ancient age of viral  
707 lineages. *Molecular Biology and Evolution*. 28:3355–3365.
- 708 Wertheim JO, Worobey M. 2009. Dating the age of the SIV lineages that gave rise to HIV-1 and  
709 HIV-2. *PLoS Computational Biology*. 5:e1000377.
- 710 Wickham H. 2016. *ggplot2: elegant graphics for data analysis*. Springer.
- 711 Worobey M, Telfer P, Souquière S, et al. (11 co-authors). 2010. Island biogeography reveals the  
712 deep history of siv. *Science*. 329:1487–1487.
- 713 Wu NC, Dai L, Olson CA, Lloyd-Smith JO, Sun R. 2016. Adaptation in protein fitness landscapes  
714 is facilitated by indirect paths. *eLife*. 5:e16965.
- 715 Yang Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable  
716 rates over sites: approximate methods. *J. Mol. Evol.* 39:306–314.

- 717 Yang Z, Nielsen R. 2008. Mutation-selection models of codon substitution and their use to estimate  
718 selective strengths on codon usage. *Molecular Biology and Evolution*. 25:568–579.
- 719 Yang Z, Nielsen R, Goldman N, Pedersen AMK. 2000. Codon-substitution models for heteroge-  
720 neous selection pressure at amino acid sites. *Genetics*. 155:431–449.
- 721 Yang Z, Rannala B. 2012. Molecular phylogenetics: principles and practice. *Nature Reviews  
722 Genetics*. 13:303.
- 723 Yu G, Smith DK, Zhu H, Guan Y, Lam TTY. 2017. ggtree: an R package for visualization and  
724 annotation of phylogenetic trees with their covariates and other associated data. *Methods in  
725 Ecology and Evolution*. 8:28–36.
- 726 Zuckerkandl E, Pauling L. 1965. Evolutionary divergence and convergence in proteins. In: Evolv-  
727 ing genes and proteins. New York, NY: Academic Press, pp. 97–166.

728 **Supplemental Information**



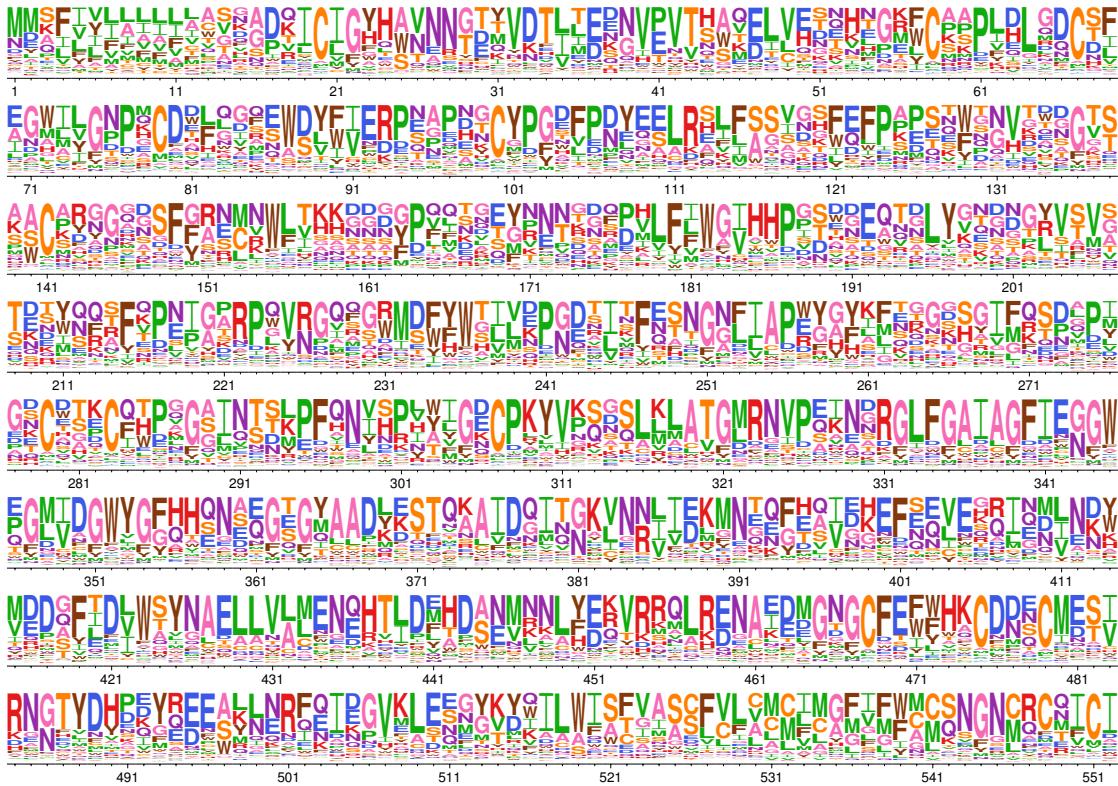
**Supplementary figure 1: H1 HA amino-acid preferences measured by deep mutational scanning.** Each column represents a site in the HA protein, and the height of each letter is proportional to the preference for the amino acid measured by Doud and Bloom (2016) and then re-scaled by the stringency parameter in Table 1. The plot only shows sites that are alignable between the H1 and H3 HAs, and these alignable sites are numbered sequentially starting from 1. The conversion between the numbering scheme in this figure and sequential numbering of the H1 HA reference sequence is in [Supplementary file 1](#).



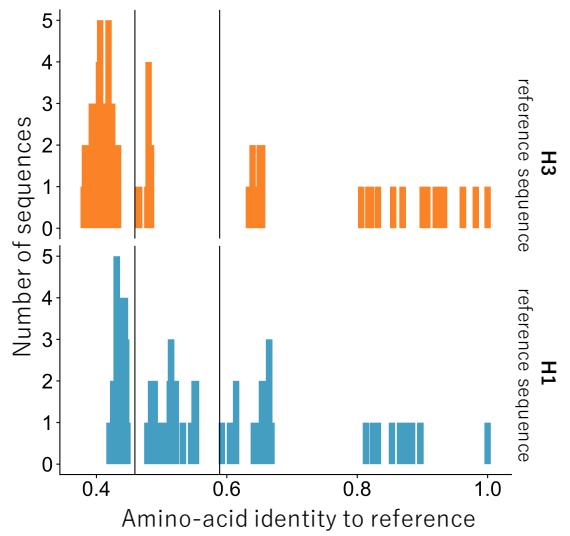
**Supplementary figure 2: H3 HA amino-acid preferences measured by deep mutational scanning.** Similar to [Supplementary figure 1](#) but shows the re-scaled preferences for the H3 HA as measured by [Lee et al. \(2018\)](#).



**Supplementary figure 3: Average of H1 HA and H3 HA amino-acid preferences measured by deep mutational scanning.** Similar to [Supplementary figure 1](#) but shows the re-scaled average of the preferences for the H1 and H3 HAs.



**Supplementary figure 4: Amino-acid preferences inferred by the pbMutSel model.** Similar to [Supplementary figure 1](#), but shows the preferences inferred by fitting the pbMutSel model to the full HA tree.



**Supplementary figure 5: Overall divergence for subtrees.** We created two subalignments for each HA used in the deep mutational scanning experiments. The “low divergence” alignments had  $\geq 59\%$  amino-acid identity to either the H1 or H3 reference sequence. The “intermediate divergence” alignments had  $\geq 46\%$  amino-acid identity to the reference sequences.

Model	Tree diameter (average codon substitutions per site)	Percentage of GY94 tree diameter
GY94	12.04	100%
ExpCM(H1)	14.70	122%
ExpCM(H3)	16.28	135%
ExpCM(H1+H3 avg)	19.21	160%
GY94 + $\Gamma\omega$	19.15	159%
ExpCM(H1) + $\Gamma\omega$	24.75	206%
ExpCM(H3) + $\Gamma\omega$	25.03	208%
ExpCM(H1+H3 avg) + $\Gamma\omega$	30.78	256%

**Supplementary table 1: Branch length extension as measured by tree diameter.** We calculated the tree diameter, the distance between the two most divergent tips, for the trees in [Figure 4](#). For each tree, the diameter is reported as a raw value and as a percentage of the GY94 model tree, the smallest of the eight trees.

Model	Parameters
GY94	$\kappa = 3.17, \omega = 0.10,$ $\phi_{1,A} = 0.32, \phi_{1,C} = 0.14, \phi_{1,G} = 0.28,$ $\phi_{2,A} = 0.38, \phi_{2,C} = 0.18, \phi_{2,G} = 0.20,$ $\phi_{3,A} = 0.36, \phi_{3,C} = 0.19, \phi_{3,G} = 0.21$
GY94 + $\Gamma\omega$	$\alpha_\omega = 0.51, \beta_\omega = 3.92, \kappa = 3.49,$ $\phi_{1,A} = 0.32, \phi_{1,C} = 0.14, \phi_{1,G} = 0.28,$ $\phi_{2,A} = 0.38, \phi_{2,C} = 0.18, \phi_{2,G} = 0.20,$ $\phi_{3,A} = 0.36, \phi_{3,C} = 0.19, \phi_{3,G} = 0.21$
ExpCM(H1)	$\beta = 1.56, \kappa = 3.64, \omega = 0.24,$ $\phi_A = 0.378, \phi_C = 0.17, \phi_G = 0.23$

**Supplementary table 2: Model parameters fit to a low divergence tree.** We fit GY94 models and an ExpCM defined by H1 deep mutational scanning preferences to the “low divergence from H1” tree in [Figure 5](#). We used these model parameters calculate the expected pairwise sequence identity in [Figure 2](#) and simulate the sequences in [Figure 3](#).

**Supplementary file 1:** List of alignable sites between H1 HA and H3 HA. This files provides a conversion between the numbering scheme we use in the paper (sequential numbering of just the alignable sites) to sequential numbering of the H1 HA reference sequence A/Wilson Smith/1933 and the H3 HA reference sequence A/Perth/2009.

**Supplementary file 2:** Amino acid preferences measured by the deep mutational scanning of the H1 HA strain A/WSN/1933 ([Doud and Bloom, 2016](#)). This file only contains measurements for the alignable sites between H1 and H3 HAs. Conversion from this numbering scheme to sequential numbering of A/WSN/1933 is in [Supplementary file 1](#).

**Supplementary file 3:** Amino acid preferences measured by the deep mutational scanning of the H3 HA strain A/Perth/2009 ([Lee et al., 2018](#)). This file only contains measurements for the alignable sites between H1 and H3 HAs. Conversion from this numbering scheme to sequential numbering of A/Perth/2009 is in [Supplementary file 1](#).

**Supplementary file 4:** The HA sequences for the full HA tree. The sequences in this alignment contain only the alignable sites between H1 and H3 HAs. Conversion from this numbering scheme to sequential numbering of A/Perth/2009 is in [Supplementary file 1](#).