

# The relationship between scaled selection coefficients and $dN/dS$

Stephanie J. Spielman\* and Claus O. Wilke\*

\*Department of Integrative Biology, Center for Computational Biology and Bioinformatics, and Institute of Cellular and Molecular Biology. The University of Texas at Austin, Austin, TX 78712, USA.

Submitted to Proceedings of the National Academy of Sciences of the United States of America

Two broad classes of Markov-process models have been developed to describe the strength of natural selection in protein-coding sequences in a phylogenetic context. The first and most widely-used modeling framework estimates the evolutionary rate ratio  $dN/dS$ , which represents the ratio of nonsynonymous to synonymous substitution rates, in a maximum likelihood (ML) framework. These models have been developed to a high level of sophistication and are a staple of modern-day comparative sequence analysis. The second class of models, known as mutation-selection-balance (MutSel) models, explicitly consider the dynamic interplay between mutation and selection. Instead of  $dN/dS$  estimates, these models calculate site-specific amino-acid and/or codon scaled selection coefficients, which indicate the extent to which natural selection acts on all possible mutations. However, the extent to which these two modeling frameworks relate to each other is unknown. Do  $dN/dS$  estimates reveal similar or distinct information from scaled selection coefficients? To answer this question, we derive a formal mathematical relationship between these two quantities. We find that scaled selection coefficients can precisely predict  $dN/dS$ , indicating that  $dN/dS$  and MutSel models are in complete agreement. We additionally show that  $dN/dS$  calculated from selection coefficients is strictly less than 1, and thus MutSel models inherently cannot describe positive selection and/or adaptive evolution. However, if synonymous mutations are not neutral, it is actually possible for  $dN/dS$  to be greater than 1, even though positive selection is not occurring. Finally, we find that ML  $dN/dS$  models produce systematically biased  $dN/dS$  estimates when nucleotide mutation rates are asymmetric, demonstrating that these models do not properly account for nucleotide compositional bias.

$dN/dS$  | mutation-selection-balance models | scaled selection coefficients | Markov models of sequence evolution

## Introduction

The oldest and most-widely used method to infer selection pressure in protein-coding genes calculates the evolutionary rate ratio  $dN/dS$ , which represents the ratio of nonsynonymous to synonymous substitution rates. This metric indicates how quickly a protein's constituent amino acids change, and it is commonly used to identify proteins or protein sites that experience negative selection ( $dN/dS < 1$ ), evolve neutrally ( $dN/dS \approx 1$ ), or experience positive, diversifying selection ( $dN/dS > 1$ ) [1, 2, 3, 4]. Frameworks for calculating  $dN/dS$  have broadly fallen into two camps: heuristic counting methods [5, 6, 7, 8, 9] and maximum likelihood (ML) methods [10, 11, 1, 12]. The latter variety assume an explicit Markov-process model of sequence evolution to yield maximum likelihood estimates (MLEs) of the parameter  $\omega$ , which represents the quantity  $dN/dS$  (although we note that other styles of these models use separate parameters for nonsynonymous and synonymous substitution rates [11, 13]). These  $dN/dS$  models have become a staple of comparative sequence analysis since their introduction in the 1990s (see ref [14] for a comprehensive review), and we will refer to them as  $\omega$  models throughout this paper.

A second class of models, known as mutation-selection-balance (MutSel) models, are increasingly being viewed as

a viable alternative to  $\omega$  models. While  $\omega$  models describe the how quickly a protein's constituent amino acids change, MutSel models assess the strength of natural selection operating on specific amino-acid or codon changes. In particular, the MutSel framework, couched firmly in population genetics theory, considers the specific selective responses to of all site-wise mutations in a protein-coding sequence [15, 16]. MutSel models yield estimates of site-wise scaled selection coefficients  $S = 2N_e s$ , which indicate the extent to which natural selection favors, or disfavors, particular codon or amino acid changes [15, 17, 18, 19]. Although first introduced over 15 years ago [15], MutSel models have seen little use due to their high computational expense. Recently, however, several computationally tractable model implementations have emerged [20, 21], allowing for the first time the potential for widespread adoption.

$\omega$  models have undergone rigorous development in their 20 years of existence and have advanced to high levels of sophistication. These models can accommodate a variety of evolutionary scenarios, including synonymous rate variation [11, 13], episodic [22, 23] and/or lineage-specific selection [24, 25, 26], and they can also incorporate information regarding protein structure and/or epistatic interactions [27, 28, 29, 30, 31]. This flexibility, along with accessible software implementations [32, 33, 34], make  $\omega$  models a very attractive modeling choice. On the other hand, some have argued that MutSel models, given their explicit consideration of population genetics theory and attention to site-specific amino acid fitness differences, offer a more fine-grained approach to studying protein evolution than do  $dN/dS$  models [15, 18, 19, 16]. Recent phylogenetic studies have also demonstrated that evolutionary models which account for amino acid fitness values dramatically outperform other  $\omega$  models, suggesting that MutSel models may more aptly represent the evolutionary process [35, 36].

Although both  $\omega$  and MutSel models describe the same fundamental process of protein-coding sequence evolution along a phylogeny, it is unknown how these two modeling classes relate to one another. In particular, as these inference methods have been developed independently, it remains an open question whether or not parameter estimates from one model are comparable to those of the other model. As a consequence, although certain rhetorical arguments may be made

## Reserved for Publication Footnotes

in favor of using one method over another, there is currently no formalized, concrete rationale to guide researchers in their methodological choices. Elucidating the relationship between these competing modeling frameworks will more precisely reveal under which circumstances the use of these models is justified.

Here, we formalize the relationship between  $\omega$  and MutSel models by examining the extent to which their respective focal parameters,  $dN/dS$  and scaled selection coefficients, yield overlapping information about the evolutionary process. To this end, we derive a mathematical framework to calculate  $dN/dS$  values from scaled selection coefficients. We find that  $dN/dS$  values can be precisely calculated from scaled selection coefficients, and that  $dN/dS$  and the distribution of scaled selection coefficients are strongly related. Furthermore, we prove that, when synonymous mutations are neutral,  $dN/dS$  calculated from selection coefficients is necessarily less than 1. This proof demonstrates that MutSel models are inherently only able to model purifying selection, and therefore would be an inappropriate model choice if positive selection is expected. However, we also find that, when synonymous codons have different fitnesses, it is possible to recover  $dN/dS$  values above 1, even though no positive selection is occurring.

Finally, we are able to use this robust relationship to assess the performance of  $\omega$  models. If  $\omega$  models are behaving as expected, we expect that they will yield the same  $dN/dS$  estimates as our calculations. We find that, in the absence of nucleotide compositional bias,  $dN/dS$  values inferred in an ML framework agree precisely with those calculated from scaled selection coefficients, meaning that MutSel and  $dN/dS$  models are in complete agreement. However, we found that, in the presence of mutational or nucleotide biases, ML inference frameworks produce systematically biased  $dN/dS$  estimates.

## Results and Discussion

**Theoretical model.** We model sequence evolution as a continuous-time Markov process [12] under the assumptions of a fixed effective population size  $N_e$  and constant selection pressure over time. This process is governed by the  $61 \times 61$  transition matrix  $P(t) = e^{Qt}$ , where the corresponding instantaneous rate matrix  $Q$  gives the instantaneous substitution probabilities between all 61 sense codons. We further assume that only single nucleotide changes occur instantaneously. We adopt the Halpern-Bruno [15, 17, 19, 16] MutSel modeling framework, which models the evolutionary process with explicit population genetics theory.

To being, let  $f_i$  be the fitness of codon  $i$ , and thus the selection coefficient acting on a mutation from codon  $i$  to codon  $j$  is  $s_{ij} = f_j - f_i$  [37, 17]. The fixation probability for a mutation from codon  $i$  to codon  $j$  is given by

$$P_{\text{fix}}(i \rightarrow j) = \frac{2s_{ij}}{1 - e^{-2N_e s_{ij}}} \approx \frac{1}{N_e} \frac{2N_e s_{ij}}{1 - e^{-2N_e s_{ij}}} \quad [1]$$

[38, 15, 17]. We further define  $S_{ij} = 2N_e s_{ij}$  as the scaled selection coefficient for this mutation. We model the substitution as the product of fixation and mutation rates,  $\mu$ . Therefore, the substitution probability for codon  $i$  to codon  $j$  is

$$q_{ij} = N_e \mu_{ij} P_{\text{fix}}(i \rightarrow j) = \mu_{ij} \frac{S_{ij}}{1 - e^{-S_{ij}}}, \quad [2]$$

and this expression corresponds to the instantaneous matrix element  $Q_{ji}$ . [15, 37]. Given detailed balance (reversibility), we have

$$q_{ij}p_i = q_{ji}p_j, \quad [3]$$

where  $p_i$  is the stationary frequency of codon  $i$ . From equations [2] and [3], we can write the ratio of substitution probabilities as

$$\frac{q_{ij}}{q_{ji}} = \frac{p_i \mu_{ij} S_{ij} (1 - e^{-S_{ji}})}{p_j \mu_{ji} S_{ji} (1 - e^{-S_{ij}})} \quad [4]$$

Given that  $S_{ij} = -S_{ji}$ , we can simplify equation [4] to show that  $q_{ij}/q_{ji} = e^{S_{ij}}$ , and we therefore find that

$$S_{ij} = \ln \left( \frac{p_j \mu_{ji}}{p_i \mu_{ij}} \right). \quad [5]$$

These equations establish a relationship between scaled selection coefficients and the stationary codon frequencies of the Markov model. Moreover, we note that in the specific case of symmetric mutation rates  $\mu_{ij} = \mu_{ji}$ , we have  $S_{ij} = \ln \left( \frac{p_j}{p_i} \right)$  [37].

**Mathematical relationship between scaled selection coefficients and  $dN/dS$ .** Using the theory laid out in the previous subsection, we can calculate an evolutionary rate by summing over all substitution probabilities weighted by the frequency of the originating codon. Further, we can establish specific expressions for nonsynonymous and synonymous evolutionary rates, and then divide them in order to obtain a value for the evolutionary rate ratio  $dN/dS$ .

To begin, we can write the nonsynonymous rate  $K_N$  as

$$K_N = N_e \sum_i \sum_{j \in \mathcal{N}_i} p_i P_{\text{fix}}(i \rightarrow j) \mu_{ij}, \quad [6]$$

where  $\mathcal{N}_i$  is the set of codons that are nonsynonymous to codon  $i$  and differ from it by one nucleotide. To normalize  $K_N$ , we divide it by the number of nonsynonymous sites, which we calculate according to the mutational opportunity definition of a site [10, 12] as

$$L_N = \sum_i \sum_{j \in \mathcal{N}_i} p_i \mu_{ij}, \quad [7]$$

and thus we find that

$$dN = \frac{K_N}{L_N} = \frac{N_e \sum_i \sum_{j \in \mathcal{N}_i} p_i P_{\text{fix}}(i \rightarrow j) \mu_{ij}}{\sum_i \sum_{j \in \mathcal{N}_i} p_i \mu_{ij}}. \quad [8]$$

Similarly, for  $dS$ , the synonymous evolutionary rate  $K_S$  per synonymous site  $L_S$ , we find

$$dS = \frac{K_S}{L_S} = \frac{N_e \sum_i \sum_{j \in \mathcal{S}_i} p_i P_{\text{fix}}(i \rightarrow j) \mu_{ij}}{\sum_i \sum_{j \in \mathcal{S}_i} p_i \mu_{ij}}, \quad [9]$$

where  $\mathcal{S}_i$  is the set of codons that are synonymous to codon  $i$  and differ from it by one nucleotide substitution. The quantities  $K_S$  and  $L_S$  are defined as in Eqs. [??] and [7] but sum over  $j \in \mathcal{S}_i$  instead of  $j \in \mathcal{N}_i$ . Moreover, we note that, if we make the dual assumptions that nucleotide mutation rates are symmetric and that all synonymous codons have equal fitness (i.e. synonymous mutations are neutral), the synonymous fixation rate  $P_{\text{fix}}(i \rightarrow j) = 1/N_e$  [39]. Under this circumstance, the value for  $dS$  reduces to 1.

**$dN/dS$  accurately reflects selection strength.** Using the theoretical framework established in equations [1] - [9], we can examine the relationship between the  $dN/dS$  values corresponding to different distributions of scaled selection coefficients. To this end, we generated 200 distinct distributions of amino acid fitness values  $f_a$ . We drew these 20 amino acid fitness values from a normal distribution  $\mathcal{N}(0, \sigma^2)$ , where  $\sigma^2 \sim \mathcal{U}(0, 4)$ . Here,  $\sigma^2$  effectively represents the strength of

natural selection; higher  $\sigma^2$  correspond to larger fitness difference among amino acids, prompting selection to act more strongly against nonsynonymous changes. In other words, high  $\sigma^2$  values indicate strong purifying selection, while lower values indicate weaker purifying selection. We additionally note that these amino acid fitness quantities correspond to the amino acid propensity parameters estimated by currently available MutSel inference methods [20, 21].

We then converted each distribution of amino acid fitnesses to a corresponding set of codon fitnesses. For 100 of the distributions, we assumed that synonymous changes were neutral, and thus we directly assigned each codon the same fitness  $f_i = f_a$ . For the other 100 sets of fitnesses, we allowed synonymous codons to have different fitness values. In these circumstances, we randomly selected a preferred codon for each amino acid, and we assigned the preferred codon the fitness of  $f_i = f_a + \lambda$  and all non-preferred codons the fitness  $f_i = f_a - \lambda$ . We drew a unique  $\lambda$  for each fitness distribution from  $\mathcal{U}[0, 2]$ . We refer to first set of codon selection coefficients as “no codon bias,” and the second set as “codon bias.”

Finally, using equations [1] - [9], we computed  $dN/dS$  for each distribution of codon fitnesses. For these calculations, we also need to select mutation rates. We set the mutation rate for transitions as  $\mu\kappa$ , and the rate for all transversions as  $\mu$ . We use the value  $\mu = 10^{-6}$  for all  $dN/dS$  calculations, and we draw a unique value for  $\kappa$  from  $\mathcal{U}[1, 6]$  for each set of codon fitnesses.

We found that  $dN/dS$  values scale excellently with the variance ( $\sigma^2$ ) of the distribution of amino-acid scaled selection coefficients (Figure 1). As Figure 1 shows,  $dN/dS$  and  $\sigma^2$  are strongly negatively correlated; when fitness differences among amino acids are very high,  $dN/dS$  takes on lower values, properly reflecting stronger purifying selection. Furthermore, as expected, this trend is much stronger for alignments without codon bias (Figure 1A,  $r^2 = 0.83$ ) than for alignments with codon bias (Figure 1B,  $r^2 = 0.45$ ). The weakened relationship for alignments with codon bias emerges from the fact that fitness differences among synonymous codons will obscure underlying amino acid fitness differences. Even so, the presence of codon bias does not remove the significant negative correlation between  $dN/dS$  and selection strength.

Importantly, Figure 1A demonstrates that, in the limiting case when  $\sigma^2$  approaches 0, and thus all codons have virtually the same fitness,  $dN/dS$  converges to 1. More precisely, the largest  $dN/dS$  value recovered for alignments without codon bias was 0.997, and this alignment featured a  $\sigma^2 = 0.08$ . This result properly reflects the case of neutral evolution. In fact, in **PROOF**, we prove that, when synonymous changes are neutral and mutation is symmetric (i.e.  $\mu_{xy} = \mu_{yx}$ ),  $dN/dS$  is necessarily always less than or equal to 1. We have also proven that, when synonymous changes are neutral and mutation rates are symmetric,  $dN/dS$  as calculated from scaled selection coefficients will always be less than 1. This proof formalizes the MutSel model underlying assumption that selection pressure is constant over the phylogeny and confirms that MutSel models are inherently unable to describe positive, diversifying selection. Although this proof assumes symmetric nucleotide mutation rates, we do not expect that deviations from this assumption will have dramatic effects on  $dN/dS$  estimates. Therefore, we conclude that the MutSel framework is an inappropriate model when positive selection is expected, as the model may yield spurious and misleading results.

However, this restriction of  $dN/dS < 1$  does not hold when synonymous changes are not neutral, as seen in Figure 1B. In other words, even though the underlying evolutionary model assumes evolutionary equilibrium, i.e. selection pressures remain constant over time and hence there is no positive, diver-

sifying selection,  $dN/dS$  can readily be greater than 1. Indeed, it is theoretically possible to achieve arbitrarily high  $dN/dS$  values when there are fitness differences among synonymous codons; in the most extreme case of codon bias, where only a single codon per amino acid is selectively tolerated, the number of synonymous sites  $L_S = 0$ , and thus the value for  $dN/dS$  approaches infinity. Given that the MutSel model framework assumes an overarching regime of purifying selection, this finding might seem paradoxical. However, the logical argument that  $dN/dS > 1$  represents positive, diversifying selection assumes that the rate of synonymous change may be used as a neutral benchmark, an assumption which codon bias clearly violates. Thus, in theory, what is classically termed positive selection can result simply from strong synonymous fitness differences.

We acknowledge that it is unlikely that this result will have a strong influence in real analyses, as selection on synonymous codons is likely relatively weak in most taxa [40]. For instance, experimental evidence from the yeast Hsp90 protein suggests that, while there are some fitness differences among synonymous codons, these differences are exceedingly minimal compared to fitness differences among amino acids [41, 42]. Moreover, our implementation of codon bias explicitly assumed that selection alone, and not mutation, was the sole source of codon bias. This implementation might not be entirely biologically realistic, as both mutational and selective forces likely contribute to codon bias in real genomes [43, 44, 40, 45, 46]. Even so, the fact that  $dN/dS$  can theoretically bear the hallmark of positive selection, even when purifying selection alone is occurring, is an important insight. It is, therefore, possible that estimates of positive selection in species with high levels of codon bias driven in part by selection, such as bacterial, *Drosophila*, or certain mammalian species [44, 47, 46], may not be true cases of positive selection, but rather signals of strong codon bias.

**$dN/dS$  calculated from ssc’s as a model benchmark.** As a consequence of the relationship between  $dN/dS$  and scaled selection coefficients, we have a unique opportunity to assess the robustness of  $dN/dS$  inference methods. It is conventional practice in model development to benchmark the model against data simulated according to the model itself. While this strategy is crucial for testing whether a model was implemented correctly, it is inherently incapable of discerning limitations and properties of the inference framework under more general conditions. Moreover, it cannot confirm that the underlying model accurately represents the evolutionary process. Therefore, we apply a novel benchmarking approach which uses the theoretical relationship among modeling frameworks to assess the accuracy and specific utility of those models. This approach, outlined in Figure 2A, entails comparing  $dN/dS$  values calculated from selection coefficients to those inferred by an  $\omega$ -based model, in order to benchmark the model’s accuracy.

Using the selection coefficients and mutation rates derived in the previous subsection, we simulated alignments using standard methods [12] according the Halpern-Bruno MutSel model [15]. We then inferred  $dN/dS$  for each alignment using the M0 model [10, 2], as implemented in the HyPhy batch language [32]. This model uses the GY94 instantaneous rate matrix, which includes parameters for transition/transversion bias ( $\kappa$ ), equilibrium codon frequencies ( $\pi$ ), and finally the  $dN/dS$  rate ratio ( $\omega$ ) (see equation [13]). Throughout the remaining text, we refer to  $dN/dS$  as inferred by M0 as  $\omega$ , and to  $dN/dS$  computed using equations [1] - [9] simply as  $dN/dS$ .

We found that  $dN/dS$  values agree nearly perfectly with  $\omega$  values (Figure 2B). This agreement was neither influenced by the presence of codon bias, nor by nucleotide compositional bias; indeed, simulated alignments featured a wide range (0.21–0.89) of GC contents. Additionally, in Figure 2C, we demonstrate that  $\omega$  converges to the true  $dN/dS$  value as the size of the data set, represented by simulated alignment length, increases. These results unequivocally show that the  $dN/dS$  quantity is fully contained within MutSel model parameters, and importantly that ML  $dN/dS$  inference methods are robust indeed reveal that ML  $dN/dS$  inference frameworks behave exactly as expected, yielding precise  $dN/dS$  estimates. This finding has important implications for modeling choices; although the MutSel framework might model the sequence evolution in a way that more mechanistically matches the evolutionary process,  $\omega$ -based models do not dramatically suffer from any modeling limitations.

**$dN/dS$  inference with realistic data.** Having confirmed that  $\omega$ -based and MutSel models agree under general conditions, we sought to test the accuracy of  $\omega$ -based models using more realistic data. To this end, we made use of realistic amino acid fitness and nucleotide mutation rate parameters to construct a series of experimental evolutionary models. In particular, we used influenza nucleoprotein (NP) site-specific amino acid preference values, given by ref. [35], which consisted of experimentally-determined fitness values for each individual amino acid across all sites in NP, yielding 498 distinct amino acid propensity distributions. We combined these experimental fitness parameters with three sets of experimentally-determined mutation rates, either for NP [35], yeast [48], or polio virus [49]. While each of these mutation matrices is asymmetric, they feature differing degrees of asymmetry, with NP mutation rates being the most symmetric and polio mutation rates the most asymmetric. More precisely, in the absence of amino-acid level selection, the GC contents that the NP, yeast, and polio mutation rates would generate are 0.518, 0.336, and 0.192, respectively. Finally, we built a unique experimentally-informed evolutionary model for all combinations of amino acid fitness distributions and mutation rates using the approach outlined in refs. [35, 36]. We calculated stationary codon frequencies for each experimental model, and used these values in combination with their corresponding mutation rates to simulate sequences and infer  $dN/dS$  quantities.

As different M0 model parameterizations are known to yield distinct  $\omega$  estimates [50, 12], we inferred  $\omega$  using the commonly-used frequency estimators F61 [10], F3x4 [11], and CF3x4 [51]. The F61 estimator approximates these model parameters using an alignment’s empirical codon frequencies, while the F3x4 and CF3x4 estimators approximate codon frequency parameters using positional nucleotide frequencies. We calculated values for codon frequency parameters using two different ways. First, we used the empirical codon frequencies found in the simulated alignments; as we simulated unique alignments for all site-specific amino acid preferences, we pooled all alignments for a given mutation rate specification to determine the global empirical frequencies, and subsequently calculate values for F61, F3x4 and CF3x4 frequency parameterizations. This strategy is analogous to how these frequency parameters are estimated in common practice. Second, we computed frequency parameter values using the codon frequencies which would exist in the absence of natural selection, but would arise strictly from mutational processes. As these values correspond precisely to those intended for these

model parameters [10, 11, 9, 12], we term these the “True Frequencies.”

Overall, we found that  $\omega$ -based models clearly before best when mutation rates symmetric.

We found that these results, while still pretty good, were not as good as under symmetric mutation rates. NP, which is the closest to symmetric, yielded very good results, with high correlations and minimal bias. On the other hand, yeast and polio go nuts. Presence of asymmetric mutation rates led all M0 model parameterizations to systematically infer underestimated values. We found that the three frequency estimators performed comparably,  $\omega$  estimates are systematically biased downwards as mutation rates become increasingly asymmetric. This trend generally holds regardless of the frequency estimator used (with a single exception, which we discuss below).  $\omega$  inferences on alignments with NP mutation rates have the least amount of bias, followed by yeast and finally polio mutation rates, and a parallel trend exists for the correlation strengths. However, using true frequencies (Figure 3C–D), as opposed to empirical (Figure 3A–B) frequencies, did reduce the bias (mean absolute decrease of 4.47%) for simulations which used polio mutation rates, but bias was largely unchanged for NP and yeast mutation rates. Additionally, correlations between  $\omega$  and  $dN/dS$  were generally higher when empirical frequencies were used instead of true frequencies. Taken together, these results suggest that current codon frequency estimators cannot fully account for strongly asymmetric mutation rates and/or even nucleotide composition introduced by mutation.

Therefore, we sought to ameliorate this systematic bias by introducing a novel model parameterization, which we term “Fnuc.” Fnuc merges the equilibrium nucleotide and codon frequencies to reveal a distinct target nucleotide frequency for each specific instantaneous change. Therefore, this parameterization takes into account both mutation at the nucleotide level as well as overarching codon frequencies. A derivation of this parameterization is given in **Appendix 2**. We used the Fnuc parameterization to again infer  $\omega$ , using both empirical and true frequencies. We found that Fnuc substantially reduces bias while simultaneously increasing precision in  $\omega$  estimates, in particular for the yeast mutation rates. While empirical Fnuc did not dramatically improve bias for the polio mutation rate data set, it did offer improvements over all other frequency estimators. We additionally note that the polio mutation rates were strongly asymmetric, and in the absence of selection would lead to the strongly biased A/T content of roughly 80%. Thus, we recommend that future implementations of  $\omega$ -based models consider the Fnuc parameterization, as it has the ability to reduce misleading signal induced by asymmetric mutation rates.

## Future Directions

We will conclude with insights gained from our study and recommendations for using  $\omega$ -based and MutSel modeling frameworks going forward. We have shown that  $dN/dS$  be accurately calculated from selection coefficients, revealing that  $\omega$ -based and MutSel models yield consistent and overlapping information about the strength natural selection. Importantly, our proof that  $dN/dS \leq 1$ , when calculated from selection coefficients and when synonymous mutations are neutral, indicates that the use of MutSel models is only justified under conditions of strictly purifying selection, or neutral evolution. This restriction is in part indicated by the basic MutSel model assumption of constant selection pressures over time, or in other words a static fitness landscape [15, 28, ?, 16].

Thus, if the aim is to identify positive selection, only  $\omega$ -based models, of the two frameworks examined here, are justified. However, there are still some theoretical issues with the approach that these models take, namely in their use of equilibrium codon frequency parameters. We expect that, if positive selection is occurring, there is necessarily a shift in the underlying fitness landscape causing different amino acids to become preferred.

$\omega$ -based models are an apt model choice for examining positive, diversifying selection. However, if one desires site-specific point estimates of  $dN/dS$ ,  $\omega$ -based models

2. if and when you use omega models, you absolutely must parameterize them properly, otherwise dnds is a meaningless quantity. this seems difficult to do. there is an internal tension in omega models, wrt to equilibrium codon frequencies. If codon frequencies are assumed to be at equilibrium, how can we ever properly account for adaptive processes, such as positive selection? moreover, if we screw up any model parameter, dnds may be wrong. suggests that point estimates for dnds from these models are not ideal, see that one plos response. precise point estimates can, however, be calculated from a mutsel model, provided purifying selection. but if detecting positive selection is your goal, seriously do NOT use mutsel models. future work should investigate modeling frameworks which fully account for nonequilibrium evolutionary processes.

We additionally emphasize that improper model parameterizations lead to spurious  $\omega$  MLEs which do not accurately represent  $dN/dS$ . If other model parameters ( $\kappa$  and equilibrium codon frequencies) are specified incorrectly, or inadvertently contain information about amino-acid level natural selection, the resulting  $\omega$  MLE will not represent the true  $dN/dS$  evolutionary rate ratio. Only by ensuring that  $\omega$  is the only model parameter which contains information about natural selection will it assuredly represent  $dN/dS$ .

Taken together, these results strongly suggest that the MC model's codon frequency parameters are ill-suited to accommodate compositional biases which result from forces other than amino-acid level selection. We therefore suggest that future work investigate the utility of novel parameters for MC models which better account for asymmetry in the mutational process.

In sum, we have garnered several important insights into the behavior of MC and MutSel models, as well as the  $dN/dS$  metric. These results were only made possible through establishing a formal mathematical relationship between distinct modeling frameworks. We believe that the approach presented in this paper represents a promising future avenue for methodological benchmarking. Typically, researchers assess the performance of a given inference framework through simulations which adhere to the underlying model's assumptions (with a notable exception of ref. [52]). While this strategy is critical for testing whether a model implementation behaves as expected, it is innately incapable of assessing the limitations and properties of the inference framework under more general conditions, and it cannot confirm that the underlying model accurately represents the evolutionary process. Therefore, we suggest an alternate approach to benchmark inference methods: assessing the extent to which distinct models agree may serve as a novel, robust strategy to determine the accuracy and specific utility of different modeling frameworks. As we have shown here, this approach has great potential to reveal previously unrecognized model properties or biases and will help ensure robust model development going forward.

We additionally emphasize that improper model parameterizations lead to spurious  $\omega$  MLEs which do not accurately represent  $dN/dS$ . If other model parameters ( $\kappa$  and equilibrium codon frequencies) are specified incorrectly, or inad-

vertently contain information about amino-acid level natural selection, the resulting  $\omega$  MLE will not represent the true  $dN/dS$  evolutionary rate ratio. Only by ensuring that  $\omega$  is the only model parameter which contains information about natural selection will it assuredly represent  $dN/dS$ .

In sum, we have garnered several important insights into the behavior of MC and MutSel models, as well as the  $dN/dS$  metric. These results were only made possible through establishing a formal mathematical relationship between distinct modeling frameworks. We believe that the approach presented in this paper represents a promising future avenue for methodological benchmarking. Typically, researchers assess the performance of a given inference framework through simulations which adhere to the underlying model's assumptions (with a notable exception of ref. [52]). While this strategy is critical for testing whether a model implementation behaves as expected, it is innately incapable of assessing the limitations and properties of the inference framework under more general conditions, and it cannot confirm that the underlying model accurately represents the evolutionary process. Therefore, we suggest an alternate approach to benchmark inference methods: assessing the extent to which distinct models agree may serve as a novel, robust strategy to determine the accuracy and specific utility of different modeling frameworks. As we have shown here, this approach has great potential to reveal previously unrecognized model properties or biases and will help ensure robust model development going forward.

## Methods

**Simulation of scaled selection coefficients** We first examined the relationship between  $dN/dS$  and scaled selection coefficients by simulating 200 distributions of amino acid scaled fitness values,  $F_a = 2Nf_a$ , from a normal distribution  $\mathcal{N}(0, \sigma^2)$ , where a unique  $\sigma^2$  was drawn from  $\mathcal{U}(0, 4)$  for each fitness distribution. We converted these amino acid fitnesses to codon fitnesses,  $F_i$ . For 100 of the fitness distributions, we directly assigned all codons within a given amino acid family the fitness  $f_a$ , giving all synonymous codons the same fitness. For the other 100 fitness distributions, we assigned synonymous codons different fitnesses by randomly selected a preferred codon for each amino acid. This preferred codon was assigned the fitness of  $F_i = F_a + \lambda$ , and all non-preferred codons were given the fitness  $F_i = F_a - \lambda$ . We drew a unique  $\lambda$  for each fitness distribution from  $\mathcal{U}[0, 2]$ . We then computed stationary codon frequencies as

$$p_i = \frac{e^{F_i}}{\sum e^{F_k}}, \quad [10]$$

where the sum in the denominator runs over all 61 sense codons [37]. Equation [10] gives the analytically precise stationary frequencies for a MutSel model, under the assumption of symmetric nucleotide mutation rates, i.e. where  $\mu_{xy} = \mu_{yx}$  [37]. For each resulting set of stationary codon frequencies, we used equations [6] - [9] to compute a  $dN/dS$  value. For these calculations, we set the mutation rate for transitions as  $\mu\kappa$ , and the rate for all transversions as  $\mu$ . We used the value  $\mu = 10^{-6}$  for all  $dN/dS$  calculations, and we drew a unique value for  $\kappa$  from  $\mathcal{U}[1, 6]$  for each set of codon frequencies.

**Alignment simulations.** We simulated protein-coding sequences as a continuous-time Markov process using standard methods [12] according to the Halpern-Bruno MutSel model [15]. In simplified form, this model's instantaneous rate ma-

trix  $Q$  is given by

$$Q_{ji} = \begin{cases} \mu_{ij} \frac{S_{ij}}{1-1/S_{ij}} & \text{single nucleotide change} \\ 0 & \text{multiple nucleotide changes} \end{cases}, \quad [11]$$

for a mutation from codon  $i$  to  $j$ , where  $\mu_{ij}$  is the mutation rate,  $p_i$  is the stationary frequency for codon  $i$ , and the scaled selection coefficient  $S_{ij}$  is defined in equation [5]. All alignments presented here were simulated along a 4-taxon phylogeny, beginning with a root sequence selected from stationary codon frequencies. Unless otherwise stated, all simulated alignments contained 500,000 codon positions. A single evolutionary model was applied to all positions in the simulated sequences. While this lack of site-wise heterogeneity is unrealistic for real sequence evolution, it allowed us to verify our derived relationship between scaled selection coefficients and  $dN/dS$  with a sufficiently sized data set.

**Computation of stationary frequencies for experimental data sets.** We used experimentally-determined site-specific amino acid fitness parameters ( $F_a$ ) for influenza nucleoprotein (NP), from Bloom 2014 [35], in combination with experimental nucleotide mutation rates for either NP [35], yeast [48], or polio virus [49] to derive realistic distributions of stationary codon frequencies. Bloom 2014 reported 498 distinct site-wise amino acid preference distributions for NP [35]. We combined these 498 amino acid preference sets with each of the three mutation rate matrices sets to construct a total of  $498 \times 3 = 1494$  unique experimental evolutionary Markov models, using the approach in refs. [35, 36]. The instantaneous matrix for these experimental models is given by

$$Q_{ji} = \begin{cases} \frac{F_j}{F_i} \mu_{ij} & \text{single nucleotide change, where } F_j \geq F_i \\ \mu_{ij} & \text{single nucleotide change, where } F_j < F_i \\ 0 & \text{multiple nucleotide changes} \end{cases}, \quad [12]$$

where  $F_i$  is the fitness of codon  $i$  [35, 36]. We calculated  $F_i$  values by simply assigning a given amino acid's experimental fitness  $F_a$  to each of its constituent codons; thus, all synonymous changes are neutral. We determined the stationary codon frequencies for each resulting experimental model from the matrix's eigenvector corresponding to the eigenvalue 0. Finally, we simulated alignments for each set of stationary frequencies and corresponding mutation rates according to equation [11].

**Maximum likelihood inference of  $dN/dS$ .** We inferred  $\omega$  for all simulated alignments using the M0 model. For the 200 alignments simulated with symmetric mutation rates, we inferred  $dN/dS$  using the M0 model [2], as implemented in the HyPhy batch language [32]. The M0 model uses the GY94 instantaneous rate matrix,

$$Q_{ji} = \begin{cases} \pi_j & \text{synonymous transversion} \\ \kappa \pi_j & \text{synonymous transition} \\ \omega \pi_j & \text{nonsynonymous transversion} \\ \omega \kappa \pi_j & \text{nonsynonymous transition} \\ 0 & \text{multiple nucleotide changes} \end{cases}, \quad [13]$$

where  $\kappa$  is the transition-transversion bias,  $\pi_j$  is the equilibrium frequency of the target codon  $j$ , and  $\omega$  represents  $dN/dS$  [10, 1]. Importantly, this model's  $\pi$  parameters are intended to represent those codon frequencies which would exist in absence of selection pressure, but those which would result from mutation alone [10, 11, 9, 12]. Therefore, when inferring  $\omega$  for alignments simulated with symmetric mutation rates, used the F\_equal frequency parameterization, which assigns equal

values of  $1/61$  for all  $\pi_i$  [12]. F\_equal gives the exact codon frequencies expected under symmetric mutation, in the absence of selection. Alternatively, when inferring  $\omega$  for alignments simulated with experimental mutation rates, we used 5 different sets of equilibrium frequency parameterizations. First, we inferred  $\omega$  by specifying codon frequencies which would arise strictly from mutational processes in the absence of natural selection. We computed these codon frequency values using the same approach as we did in calculating the true steady-state codon frequencies, except instead of using the experimental amino acid preference data, we assigned all amino acids the same preference value of 0.05. This strategy eliminated amino-acid level selection and allowed mutation rates alone to determine equilibrium codon frequencies. We term this frequency parameterization "Ftrue." Finally, we used the common frequency estimators F61 [10], F3x4 [11], and CF3x4 [51]. As typical analyses consider model frequency parameters as protein-wide (not site-specific) parameters, we computed these parameter values by pooling, for each set of mutation rates, all 498 steady-state codon frequencies to derive average codon frequencies. This approach yielded a set of global equilibrium frequencies for each set of mutation rates, and we calculated the F61, F3x4, and CF3x4 frequencies from these distributions. Finally, we used  $\omega$  with the F\_equal parameterization.

## Appendix 1

To prove that  $dN/dS \leq 1$ , when computed from scaled selection coefficients, we note that

We assume that  $P_i \leq P_j$  and that  $P_j > 0$ .

To this end, we introduce the function

$$F(x, y) = x + y - \frac{2xy[\log(x) - \log(y)]}{x - y} \quad [14]$$

, and we will now show that  $F(x, y) \geq 0$  for  $x \leq y$  and  $y \geq 0$ . It is straightforward to show this for  $x = y$ . For  $x < y$ , we show that the first derivative of equation [14] is negative throughout  $x \in (0, y)$ , which proves that the function monotonically decreases, and thus  $F(x, y) > 0$ , in this interval. We calculate the first derivative as

$$\frac{\partial F(x, y)}{\partial x} = \frac{[(x - 3y)(x - y) + 2y^2(\log x - \log y)]}{(x - y)^2}. \quad [15]$$

We can further rewrite (expand?) the expression  $\log x - \log y$  as a series, thus yielding

$$\frac{\partial F(x, y)}{\partial x} = \frac{(x - 3y)(x - y) + 2y^2 \log(\frac{x}{y})}{(x - y)^2}. \quad [16]$$

If we then take only the first two terms of the series, we find that the expression simplifies to 0. Therefore, we can write

$$\frac{\partial F(x, y)}{\partial x} = \frac{-2y^2 \sum_{n=3}^{\infty} \left(1 - \frac{x}{y}\right)^n}{(x - y)^2}, \quad [17]$$

which is clearly negative.

## Appendix 2

To derive the Fnuc model parameterization, we begin with a modified version of the MG94 [11] instantaneous rate matrix,

$$Q_{ji} = \begin{cases} \pi_n^j & \text{synonymous transversion} \\ \kappa \pi_n^j & \text{synonymous transition} \\ \omega \pi_n^j & \text{nonsynonymous transversion} \\ \omega \kappa \pi_n^j & \text{nonsynonymous transition} \\ 0 & \text{multiple nucleotide changes} \end{cases}, \quad [18]$$

where  $\pi_n^j$  is the equilibrium frequency of the target nucleotide which experienced a mutation from codon  $i$  to codon  $j$ . The primary difference between equation [18] and the original MG94 rate matrix is that we consider a single parameter,  $\omega$

for the nonsynonymous/synonymous rate ratio, whereas the original formulation used parameters  $\beta$  and  $\alpha$  to describe separately nonsynonymous and synonymous rates, respectively. As Muse and Gaut [11] noted, the stationary frequency for a codon comprised of nucleotides  $x$ ,  $y$ , and  $z$  (in that order) is

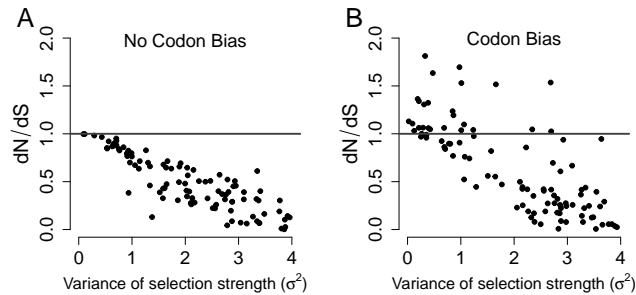
$$\pi_x \pi_y \pi_z \frac{1}{C}, \quad [19]$$

where  $C = \Pi_{\text{stop}} = \pi_T \pi_A \pi_G + \pi_T \pi_G \pi_A + \pi_T \pi_A \pi_A$ . This expression is better known as either the frequency estimator F1x4, if four parameters for nucleotide frequencies are used, or F3x4, if twelve parameters representing positional nucleotide frequencies are used. Using these equilibrium codon frequencies, we can derive precise values for target nucleotide frequencies  $\pi_n$  used in [18]. We note that  $p_{in} =$

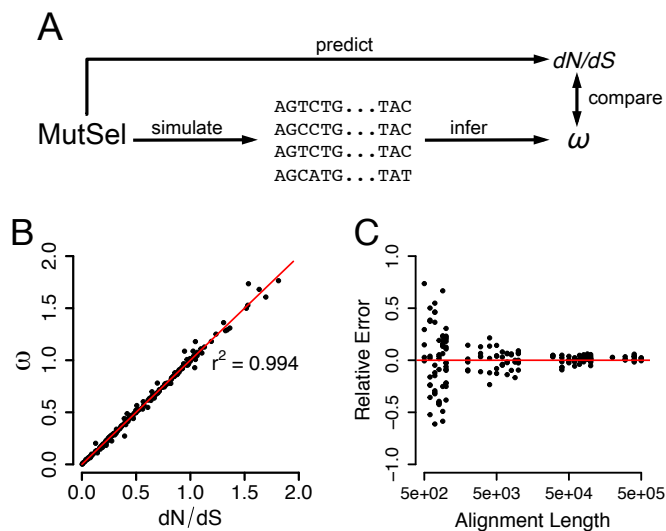
**ACKNOWLEDGMENTS.** This work was supported by the army and by NIH.

- Nielsen R, Yang Z (1998) Likelihood models for detecting positive selected amino acid sites and applications to the HIV-1 envelope gene. *Genetics* 148:929–936.
- Yang ZH, Nielsen R, Goldman N, Pedersen AMK (2000) Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics* 155:431–449.
- Kosakovsky Pond S, Frost SD (2005) Not so different after all: A comparison of methods for detecting amino acid sites under selection. *Mol Biol Evol* 22:1208–1222.
- Huelsenbeck JP, Jain S, Frost SWD, Kosakovsky Pond SL (2006) A Dirichlet process model for detecting positive selection in protein-coding DNA sequences. *Proc Natl Acad Sci USA* 103:6263–6268.
- Li WH, Wu CI, Luo CC (1985) A new method for estimating synonymous and nonsynonymous rates of nucleotide substitution consider the relative likelihood of nucleotide and codon changes. *Mol Biol Evol* 2:150–174.
- Nei M, Gojobori T (1986) Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Mol Biol Evol* 3:418–426.
- Pamilo P, Bianchi NO (1993) Evolution of the Zfx and Zfy genes: rates and interdependence between the genes. *Mol Biol Evol* 10:271–281.
- Ina Y (1995) New methods for estimating the numbers of synonymous and nonsynonymous substitutions. *Journal of Molecular Evolution* 40:190–226.
- Yang Z, Nielsen R (2000) Estimating synonymous and nonsynonymous substitution rates under realistic evolutionary models. *Mol Biol Evol* 17:32–42.
- Goldman N, Yang Z (1994) A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Mol Biol Evol* 10:725–736.
- Muse SV, Gaut BS (1994) A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. *Mol Biol Evol* 11:715–724.
- Yang Z (2006) *Computational Molecular Evolution* (Oxford University Press).
- Kosakovsky Pond S, Muse SV (2005) Site-to-site variation of synonymous substitution rates. *Mol Biol Evol* 22:2375–2385.
- Anisimova M, Kosiol C (2009) Investigating protein-coding sequence evolution with probabilistic codon substitution models. *Mol Biol Evol* 26:255–271.
- Halpern AL, Bruno WJ (1998) Evolutionary distances for protein-coding sequences: modeling site-specific residue frequencies. *Mol Biol Evol* 15:910–917.
- Thorne JL, Lartillot N, Rodrigue N, Choi SC (2012) Codon models as vehicles for reconciling population genetics with inter-specific data. In Cannarozzi G, Schneider A, eds., *Codon evolution: mechanisms and models* (Oxford University Press, New York).
- Yang Z, Nielsen R (2008) Mutation-selection models of codon substitution and their use to estimate selective strengths on codon usage. *Mol Biol Evol* 25:568–579.
- Rodrigue N, Philippe H, Lartillot N (2010) Mutation-selection models of coding sequence evolution with site-heterogeneous amino acid fitness profiles. *Proc Natl Acad Sci USA* 107:4629–4634.
- Tamuri AU, dos Reis M, Goldstein RA (2012) Estimating the distribution of selection coefficients from phylogenetic data using sitewise mutation-selection models. *Genetics* 190:1101–1115.
- Rodrigue N, Lartillot N (2014) Site-heterogeneous mutation-selection models within the PhyloBayes-MPI package. *Bioinformatics* 30:1020–1021.
- Tamuri AU, Goldman N, dos Reis M (2014) A penalized-likelihood method to estimate the distribution of selection coefficients from phylogenetic data. *Genetics* 197:257–271.
- Kosakovsky Pond S, et al. (2011) A random effects branch-site model for detecting episodic diversifying selection. *Mol Biol Evol* 28:3033–3043.
- Murrell B, et al. (2012) Detecting individual sites subject to episodic diversifying selection. *PLoS Genet* 8:e1002764.
- Yang Z, Nielsen R (2002) Codon-substitution models for detecting molecular adaptation at individual sites along specific lineages. *Mol Biol Evol* 19:908–917.
- Zhang J, Nielsen R, Yang Z (2005) Evaluation of an improved branch-site likelihood method for detecting positive selection at the molecular level. *Mol Biol Evol* 22:2472–2479.
- Kosakovsky Pond S, Frost S (2005) A genetic algorithm approach to detecting lineage-specific variation in selection pressure. *Mol Biol Evol* 22:478–485.
- Robinson DM, Jones DT, Kishino H, Goldman N, Thorne JL (2003) Protein evolution with dependence among codons due to tertiary structure. *Mol Biol Evol* 20:1692–1704.
- Thorne J, Choi S, Yu J, Higgs P, Kishino H (2007) Population genetics without intraspecific data. *Mol Biol Evol* 24:1667–1677.
- Rodrigue N, Kleinman C, Philippe H, Lartillot N (2000) Computational methods for evaluating phylogenetic models of codon sequence evolution with dependence between codons. *Mol Biol Evol* 26:1663–1676.
- Scherrer MP, Meyer AG, Wilke CO (2012) Modeling coding-sequence evolution within the context of residue solvent accessibility. *BMC Evol Biol* 12:179.
- Meyer AG, Wilke CO (2012) Integrating sequence variation and protein structure to identify sites under selection. *Mol Biol Evol* 30:36–44.
- Kosakovsky Pond SL, Frost SDW, Muse SV (2005) HyPhy: hypothesis testing using phylogenetics. *Bioinformatics* 21:676–679.
- Yang Z (2007) PAML 4: Phylogenetic analysis by maximum likelihood. *Molecular Biology and Evolution* 24:1586–1591.
- Delport W, Poon A, Frost S, Pond S (2010) Datamonkey 2010: a suite of phylogenetic analysis tools for evolutionary biology. *Bioinformatics* 26:2455–2457.
- Bloom JD (2014) An experimentally determined evolutionary model dramatically improves phylogenetic fit. *Mol Biol Evol* :To appear.
- Bloom JD (2014) An experimentally informed evolutionary model improves phylogenetic fit to divergent lactamase homologs. *Mol Biol Evol* 31:1956–1978.
- Sella G, Hirsh AE (2005) The application of statistical physics to evolutionary biology. *Proc Natl Acad Sci USA* 102:9541–9546.
- Kimura M (1962) On the probability of fixation of mutant genes in a population. *Genetics* 4:713–719.
- Crow JF, Kimura M (1970) *An Introduction to Population Genetics Theory* (Burgess Pub. Co., California).
- Hershberg R, Petrov D (2008) Selection on codon bias. *Annu Rev Genet* 42.
- Hietpas RT, Jensen JD, Bolon DNA (2011) Experimental illumination of a fitness landscape. *Proc Natl Acad Sci USA* 108:7896–7901.
- Hietpas RT, Bank C, Jensen JD, Bolon DNA (2013) Shifting fitness landscapes in response to altered environments. *Evolution* 67:3512–3522.
- Blumer M (1991) The selection-mutation-drift theory of synonymous codon usage. *Genetics* 129:897–907.
- Duret L (2002) Evolution of synonymous codon usage in metazoans. *Curr Opin Genet Dev* 12:640–649.
- Chen SL, Lee W, Hottes AK, Shapiro L, McAdams HH (2009) Codon usage between genomes is constrained by genome-wide mutational processes. *Proc Natl Acad Sci USA* 106:3480–3485.
- Plotkin JB, Kudla G (2011) Synonymous but not the same: the causes and consequences of codon bias. *Nature Rev Genet* 12:32–42.
- Chamary JV, Parmley JL, Hurst LD (2006) Hearing silence: non-neutral evolution at synonymous sites in mammals. *Nature Rev Genet* 7:98–108.
- Zhu YO, Siegal ML, Hall DW, Petrov DA (2014) Precise estimates of mutation rate and spectrum in yeast. *Proc Natl Acad Sci USA* :FORTHCOMING.
- Acevedo A, Brodsky L, Andino R (2014) Mutational and fitness landscapes of an RNA virus revealed through population sequencing. *Nature* 505:686 – 690.
- Zhang Z, Yu J (2006) Evaluation of six methods for estimating synonymous and nonsynonymous substitution rates. *Geno Prot Bioinfo* 4:173–181.
- Kosakovsky Pond SL, Delport W, Muse SV, Scheffler K (2010) Correcting the bias of empirical frequency parameter estimators in codon models. *PLoS One* 5:e11230.
- Holder M, Zwickl D, Dessimoz C (2008) Evaluating the robustness of phylogenetic methods to among-site variability in substitution processes. *Phil Trans R Soc B* 363:4013–4021.

## Figures

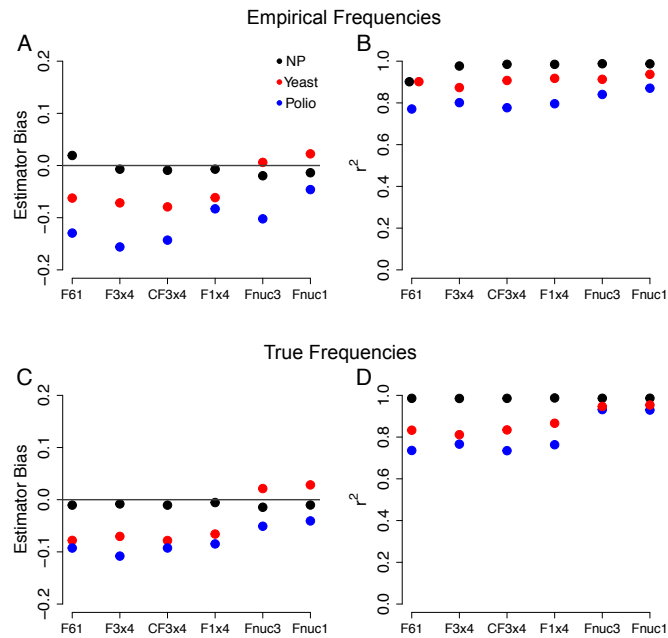


**Fig. 1.**  $dN/dS$  decreases in proportion to amino-acid level selection strength.  $dN/dS$  is plotted against the  $\sigma^2$  of the simulated distribution of amino-acid scaled selection coefficients. Higher values of  $\sigma^2$  indicate larger fitness differences among amino acids, whereas the limiting value of  $\sigma^2 = 0$  means that all amino acids have the same fitness. (A) Synonymous codons have equal fitness values ( $r^2 = 0.83$ ). (B) Synonymous codons have different fitness values ( $r^2 = 0.45$ ). Importantly, (B), but not (A) shows  $dN/dS$  values greater than 1, in spite of the steady-state evolutionary process.



**Fig. 2.** Regressions between  $dN/dS$  values as calculated from scaled selection coefficients and as inferred using the M0 mechanistic codon model. Each point corresponds to a single simulated alignment. All  $\omega$  values shown here were inferred by parameterizing the M0 model with  $\kappa$  fixed to its true, simulated value as well as the Fequal codon frequency specification [12]. The red line in panels (A-B) is the  $x = y$  line. (A) Synonymous codons have equal fitness ( $r^2 = 0.997$ ). (B) Synonymous codons have different fitness values ( $r^2 = 0.992$ ). (C) Convergence of  $\omega$  MLEs to the true  $dN/dS$  value. The y-axis indicates the relative error of the maximum likelihood  $dN/dS$  estimate, and the x-axis indicates the number of positions in the simulated alignment. As the number of positions, and hence the size of the data set, increases, the maximum likelihood estimates converge to the  $dN/dS$  values calculated using equations [??]-[9]. The red line in panel (C) is the line  $y = 0$ , indicating no error.





**Fig. 3.** (A) Estimator bias and (B)  $r^2$  values between  $dN/dS$  and  $\omega$  MLEs across M0 codon frequency parameterizations, for each set of nucleotide mutation rates. Note that negative biases indicate  $\omega$  values that are, on average, lower than  $dN/dS$ . All bias and  $r^2$  values are highly statistically significant, with all  $p < 10^{-12}$ . In this figure, we see that  $\omega$ -based models tend to systematically underestimate  $dN/dS$ , across all codon frequency parameterizations. Generally, F\_equal features the least amount of bias, and has very high  $r^2$  values for both NP and yeast mutation rates. Although Fequal yields lower  $r^2$  values for polio mutation rates than do F61, F3x4, and CF3x4, the latter three estimators also have relatively high biases, demonstrating that they systematically underestimate  $dN/dS$ . That Ftrue, which assigns codon frequencies to those which would exist in the absence of amino-acid level selection, also underestimates  $dN/dS$  implies that codon frequency parameters are ill-suited to accommodate mutation-induced nucleotide compositional bias.

Supplementary Information

Table S1. Estimator bias between  $\omega$  MLEs and the expected, true  $dN/dS$  values, for all mutation rates and M0 codon frequency parameterizations examined. Negative bias values indicate that  $\omega$  MLEs are, on average, lower than  $dN/dS$ . All biases are statistically significant, with all  $p < 2 \times 10^{-16}$ , except yeast fnuc data (9.23e-5) .

Mutation rate	F61	F3x4	CF3x4	Fnuc
Empirical Frequencies				
NP	0.019	-0.007	-0.009	-0.020
Yeast	-0.062	-0.072	-0.079	0.006
Polio	-0.129	-0.156	-0.143	-0.102
True Frequencies				
NP	-0.010	-0.008	-0.010	-0.015
Yeast	-0.078	-0.070	-0.078	0.021
Polio	-0.093	-0.108	-0.093	-0.051

Table S2. NYP  $r^2$  values between  $\omega$  MLEs and  $dN/dS$ , for all mutation rates and M0 codon frequency parameterizations examined. All  $r^2$  values are statistically significant, with all  $p < 10^{-15}$ .

Mutation rate	F61	F3x4	CF3x4	Fnuc
Empirical Frequencies				
NP	0.901	0.977	0.985	0.988
Yeast	0.902	0.874	0.908	0.913
Polio	0.771	0.801	0.777	0.841
True Frequencies				
NP	0.986	0.985	0.986	0.986
Yeast	0.833	0.812	0.834	0.947
Polio	0.736	0.766	0.735	0.932

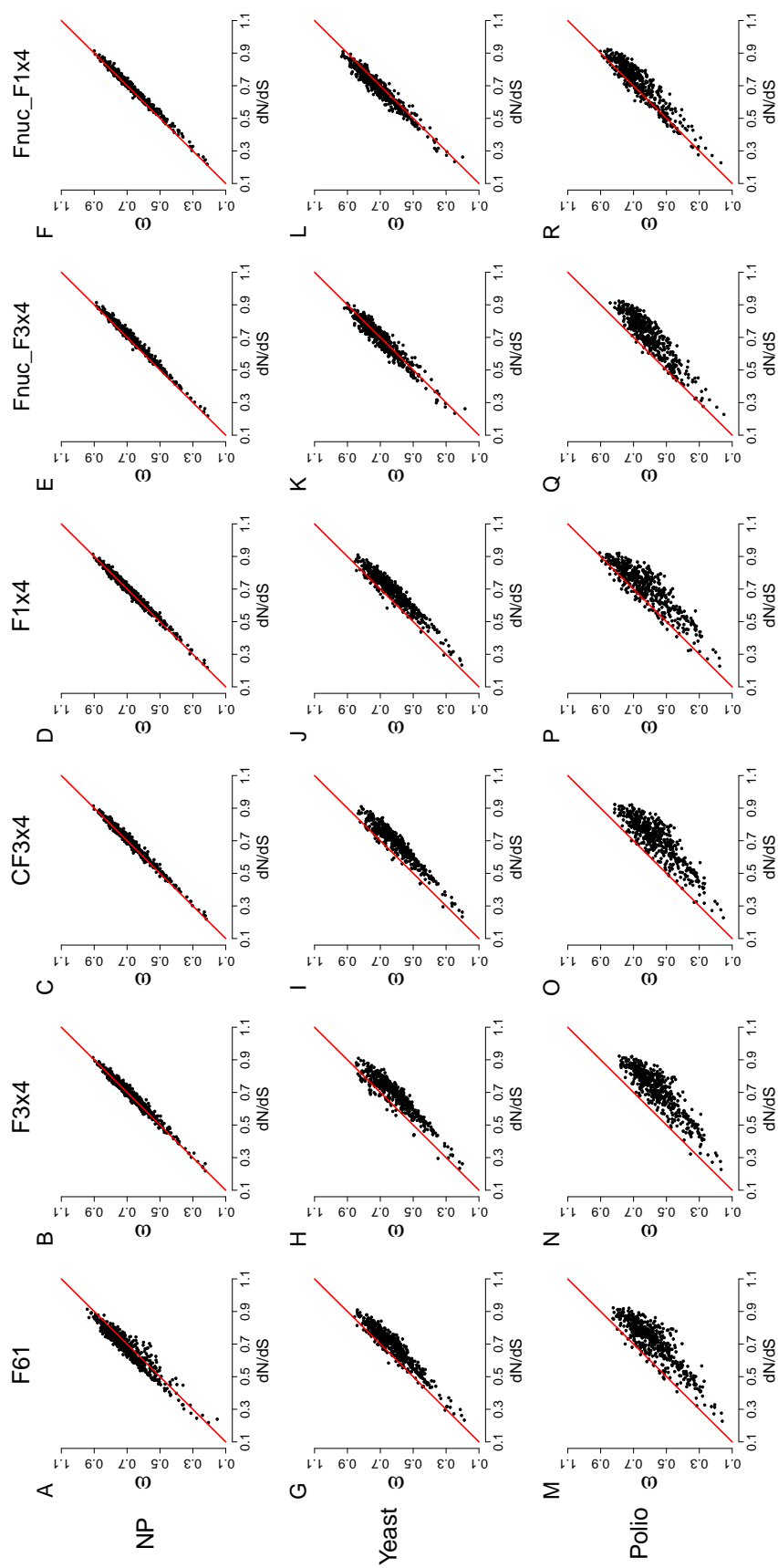


Figure S1. REWRITE: Regression plots for  $\omega$  MLEs versus  $dN/dS$  values computed from scaled selection coefficients, for each set of nucleotide mutation rates and all M0 codon frequency parameterizations. Each point represents a single simulated alignment, and the red lines correspond to  $x = y$ . (A) Simulations which assigned equal fitness values to synonymous codons. (B) Simulations which allowed different fitness values among synonymous codons.

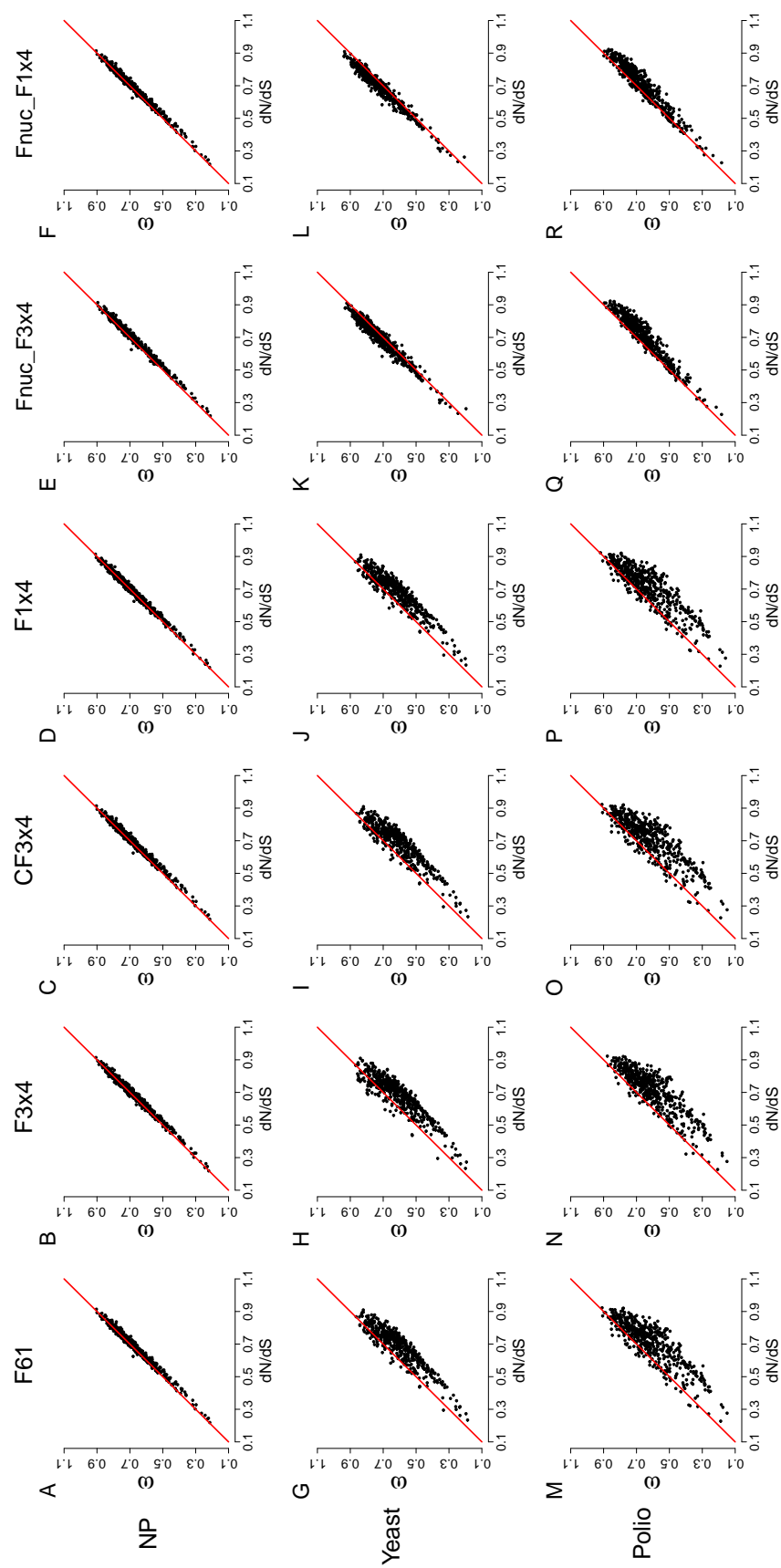


Figure S2. REWRITE: Regression plots for  $\omega$  MLEs versus  $dN/dS$  values computed from scaled selection coefficients, for each set of nucleotide mutation rates and all M0 codon frequency parameterizations. Each point represents a single simulated alignment, and the red lines correspond to  $x = y$ . (A) Simulations which assigned equal fitness values to synonymous codons. (B) Simulations which allowed different fitness values among synonymous codons.