

# Introduction, Rough

Over the years, a variety of models have been proposed to describe the effects of natural selection on protein-coding sequences, in a phylogenetic context. Traditionally, the focus has been on mechanistic codon-substitution models (see ref. [1] for a comprehensive review). Since their introduction in the 1990s, these models have seen great success in inferring protein evolutionary rates, or the nonsynonymous/synonymous rate ratio ( $dN/dS$ ). This metric indicates how quickly a protein’s constituent amino acids change [2–4], allowing for the identification of positively-selected regions in protein sequences [4, 5].

More recently, a second class of models, known as mutation-selection-balance (MutSel) models, has emerged as a popular alternative to  $dN/dS$  models. Unlike mechanistic codon models, MutSel models explicitly model the dynamic balance between mutation and selection, rather than merely the final outcome (e.g. substitution) of this process [6–9]. Moreover, these models yield estimates of amino acid selection coefficients, which indicate the extent to which natural selection favors, or disfavors, particular amino acids at protein positions. These selection coefficients, which can in turn be scaled relative to a focal amino acid, the primary parameters of interest that MutSel models produce. Although MutSel models were first introduced over 15 years ago [6], they have seen virtually no use due to their high computational expense. However, recently, several computationally tractable model implementations have emerged [10, 11], allowing for the first time the potential for widespread use.

Some have argued that MutSel models are more robust than  $dN/dS$  and can better describe the evo process owing to their treatment of amino acid identities. More fine-grained modeling results than a  $dN/dS$  analysis would yield. However, it is virtually unknown how these models really relate, so it remains unclear whether one model should be preferred over another. We don’t know how parameter estimates even relate.

Although both  $dN/dS$  and MutSel models describe the same fundamental process of protein evolution along a phylogeny, the relationship between these models is largely unknown. These two classes of models have largely been developed independently, and as a consequence we do not know whether parameters estimated from a  $dN/dS$  model are similar, distinct, or even contradictory to those estimated from from a MutSel model of coding sequence evolution. Here, we aim to formalize the relationship between  $dN/dS$  and MutSel models by examining the extent to which their focal parameters,  $dN/dS$  and scaled amino acid selection coefficients, yield overlapping information

about the evolutionary process. To this end, we derive a mathematical relationship between these two parameter classes, and we demonstrate that MutSel models fully embody the  $dN/dS$  values. Using a simulation approach, we show that we can accurately estimate  $dN/dS$  values using MutSel model parameter estimates, and these estimates correspond precisely to those inferred using a traditional  $dN/dS$  maximum likelihood inference approach. Importantly, we additionally show that this relationship holds only under regimes of purifying selection or neutral evolution ( $dN/dS \leq 1$ ). Therefore, MutSel models are inherently unable to describe protein evolution under a regime of positive selection, in which  $dN/dS > 1$ .

## Mathematical relationship between selection coefficients and omega

We describe here how to calculate  $dN/dS$  from the parameters of a MutSel model. We assume the following: (i) the mutational process is symmetric, such that  $\mu_{xy} = \mu_{yx}$  for all nucleotide pairs  $xy$ ; (ii) all synonymous codons for a given amino acid have the same fitness; there is no synonymous rate variation or codon bias.

In the framework of a MutSel model, we can write the steady-state frequency of codon  $i$  as

$$f_i = e^{s_i} / \sum_k e^{s_k}, \quad (1)$$

where the sum in the denominator runs over all 61 sense codons [12]. Here,  $s_i$  is the *scaled selection coefficient* for codon  $i$ ; larger  $s_i$  correspond to higher frequencies of codon  $i$ . The fixation probability for a mutation from  $i$  to  $j$  is [6, 12]

$$\pi_{i \rightarrow j} = \frac{1 - (f_i/f_j)^{1/N_e}}{1 - f_i/f_j} \approx \frac{1}{N_e} \frac{\ln f_j - \ln f_i}{1 - f_i/f_j}, \quad (2)$$

where  $N_e$  is the effective population size. We can calculate an evolutionary rate by summing over all fixation probabilities weighted by the frequency of the originating codon. For example, we can write the synonymous rate  $K_S$  as

$$K_S = N_e \sum_i \sum_{j \in \mathcal{S}_i} f_i \pi_{i \rightarrow j} \mu_{ij}, \quad (3)$$

where  $\mathcal{S}_i$  is the set of codons that are synonymous to codon  $i$  and differ from it by one nucleotide substitution. To normalize  $K_S$ , we divide it by the number of synonymous sites  $L_S$ , which we can calculate as

$$L_S = \sum_i \sum_{j \in \mathcal{S}_i} f_i. \quad (4)$$

Under the assumption that all synonymous codons have equal fitness (all synonymous mutations are neutral), we have  $\pi_{i \rightarrow j} = 1/N_e$  [13], and thus we find for  $dS$ , the synonymous rate per synonymous site,

$$dS = \frac{K_S}{L_S} = \frac{\sum_i \sum_{j \in \mathcal{S}_i} f_i \mu_{ij}}{\sum_i \sum_{j \in \mathcal{S}_i} f_i}. \quad (5)$$

Similarly, for  $dN$ , the non-synonymous rate per non-synonymous site, we find

$$dN = \frac{K_N}{L_N} = \frac{N_e \sum_i \sum_{j \in \mathcal{N}_i} f_i \pi_{i \rightarrow j} \mu_{ij}}{\sum_i \sum_{j \in \mathcal{N}_i} f_i}, \quad (6)$$

where  $\mathcal{N}_i$  is the set of codons that are not synonymous to codon  $i$  and differ from it by one nucleotide substitution. The quantities  $K_N$  and  $L_N$  are defined as in Eqs. (3) and (4) but summing over  $j \in \mathcal{N}_i$  instead of  $j \in \mathcal{S}_i$ .

Equations (1)–(6) establish a connection between the scaled selection coefficients  $s_i$  (i.e., the primary parameters of a MutSel model) and the evolutionary rate ratio  $dN/dS$ .

## Methods, INCREDIBLY ROUGH

### Sequence simulation and omega inference

We simulated protein-coding sequences according as a continuous-time Markov process according to the MutSel model proposed by [6]. This model's instantaneous rate matrix  $Q = q_{ij}$ , which describes the probability of substitution of from codon  $i$  to codon  $j$ , is given by

$$Q_{ij} = \begin{cases} 0 & \text{multiple nucleotide changes} \\ \mu_{ij} f_{ij} & \text{single nucleotide transversion} \\ \kappa \mu_{ij} f_{ij} & \text{single nucleotide transition} \end{cases}, \quad (7)$$

where  $\mu_{ij}$  is the symmetric nucleotide mutation rate and  $f_{ij}$  is the fixation probability from codon  $i$  to  $j$ . The fixation probability is defined as

$$f_{ij} = \ln \left( \frac{\pi_j \mu_{ij}}{\pi_i \mu_{ji}} \right) / \left( 1 - \frac{\pi_i \mu_{ji}}{\pi_j \mu_{ij}} \right), \quad (8)$$

where  $\pi_i$  is the equilibrium frequency of codon  $i$ .

For each simulation, we derived Boltzmann distributed values for the steady-state amino acid frequencies, where

$$F(a) = \frac{e^{s_a \beta}}{\sum_b e^{s_b \beta}} \quad (9)$$

, and the denominator sums over all 20 amino acids.  $s_a$  represents the scaled selection coefficient for amino acid  $a$ , analogous to the primary parameter given by a MutSel model. We sampled values for selection coefficients from a normal distribution. We fixed  $\beta$  to 3(??).

Once amino acid frequencies were determined, we assigned them to individual amino acids as follows. There was a certain number of values, out of 20, greater than 0.05 (the freq expected by random chance). For this number, we got a set of reasonably co-occurring amino acids based on mean pair-wise grantham scores. We selected a random group of aa's such that the mean distance was less than or equal to 100. The remaining frequencies were assigned randomly to the remaining amino acids.

Simulation parameters - 2 taxon tree with branch lengths fixed at 0.005.  $\mu$ , the symmetric base mutation rate was fixed at  $1e-6$ , and we varied  $\kappa$  from time to time. We evolved each position according to the same model; no site-wise variation.

For each simulated sequence set, we inferred  $dN/dS$  values using two approaches: the derivation presented in this paper (see Results) and using the standard maximum likelihood GY94 model [2] within the HyPhy package [14]. This model includes two primary parameters,  $\omega$  and  $\kappa$ . For each inference, we fixed  $\kappa$  to the known simulated value, and we provided HyPhy with equal equilibrium codon frequencies, such that each codon had a frequency of  $1/61$ . This was necessary to achieve accurate  $dN/dS$  estimates, and is discussed more in depth in Results. Using ML, we estimated a single average  $\omega$  value for each alignment.

Several things emerge from this analysis - 1. If you fit a mutsel model and calculate dnds from it, it's just as good as if you had used a mech codon model. 2. we prove that dnds, when calc'd from mutsel, must be less than 1. although generally acknowledged that purifying selection is a feature, we demonstrate it precisely and reveal that, in cases of positive selection, mutsel models are likely not appropriate. 3. we also have a more realistic way to assess performance of dnds models. bias is introduced by codon frequencies used by the model. otherwise, freqs will likely capture selective pressure, rendering omega estimates bonkers. typically highly elevated. More work should be done in this field, for instance, in dnds analyses in which codon frequencies are somehow constrained (eg membrane protein in which hydrophobic enriched) may result in artifactually high dnds. Discussion points: - Important insight is that  $dN/dS$  inherently cannot be described by mutation-selection models, and thus at those sites its results may be misleading. In particular, possibly a confounding factor in the Rodrigue implementation, as positively selected sites in the alignment could introduce bias. - Importance of examining intersections between models. Must understand how estimates from one relate to another. Helps to ensure robust results; model agreement is key, so we must formulate explicit relationships among them to systematically assess agreement. -

Future directions: 1. Consider which info the models share and which they don't. For purifying selection, params a dNdS model would yield are fully contained within a mutsel model. On the other hand, the model overlap disappears under positive selection, when steady-state equilibrium is violated. Thus, mutsel models are sufficient for purifying, but not at all useful for positive selection.

dNdS limitations: One such limitation is that the  $dN/dS$  parameter ignores the influence of site-specific amino acid propensities. It is universally recognized that a particular position in a protein will only tolerate certain amino acids, due either to functional or structural constraints. However,  $dN/dS$  ignores this key aspect of protein evolution and considers all nonsynonymous changes, regardless of which amino acid was substituted, as having equal weight on protein fitness, a biologically implausible assumption. An additional limitation is that codon-substitution models merely describe the result of the evolutionary process, rather than the explicit underlying mechanism producing those results. More precisely, substitutions in protein-coding sequences result from an ongoing mutation-selection balance. When a mutation occurs in a DNA sequence, natural selection must act on this mutation, either by disfavoring it, resulting in the mutation's removal from the population, or favoring it, ultimately yielding an amino acid replacement or substitution. By overlooking this underlying mechanism and focusing only on whether substitutions have occurred, codon substitution models are unable to capture the full extent of the evolutionary process.

## References

- [1] Anisimova M, Kosiol C (2009) Investigating protein-coding sequence evolution with probabilistic codon substitution models. *Mol Biol Evol* 26: 255–271.
- [2] Goldman N, Yang Z (1994) A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Mol Biol Evol* 11: 725–736.
- [3] 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.
- [4] 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.
- [5] Yang ZH, Nielsen R, Goldman N, Pedersen AMK (2000) Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics* 155: 431–449.
- [6] Halpern AL, Bruno WJ (1998) Evolutionary distances for protein-coding sequences: modeling site-specific residue frequencies. *Mol Biol Evol* 15: 910–917.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] Rodrigue N, Lartillot N (2014) Site-heterogeneous mutation-selection models within the PhyloBayes-MPI package. *Bioinformatics* : 1020–1021.
- [11] 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.

- [12] Sella G, Hirsh AE (2005) The application of statistical physics to evolutionary biology. *Proc Natl Acad Sci USA* 102: 9541–9546.
- [13] Crow JF, Kimura M (1970) *An Introduction to Population Genetics Theory*. California: Burgess Pub. Co.
- [14] Kosakovsky Pond SL, Frost SDW, Muse SV (2005) HyPhy: hypothesis testing using phylogenetics. *Bioinformatics* 21: 676–679.