# Population Genetics Based Phylogenetics Under Stabilizing Selection for an Optimal Amino Acid Sequence: A Nested Modeling Approach

JEREMY M. BEAULIEU[1,2,3], BRIAN C. O'MEARA[2,3], RUSSELL ZARETZKI[4],

CEDRIC LANDER[2,3], JUANJUAN CHAI[2,5], AND MICHAEL A. GILCHRIST[2,3,*]

[1]Department of Biological Sciences, University of Arkansas, Fayetteville, AR 72701

[2]Department of Ecology & Evolutionary Biology, University of Tennessee, Knoxville, TN 37996-1610

[3]National Institute for Mathematical and Biological Synthesis, Knoxville, TN 37996-3410

[4]Department of Business Analytics & Statistics, Knoxville, TN   37996-0532

[5]Current address: 50 Main St, Suite 1039, White Plains, NY 10606

[*]Corresponding author. E-mail: mikeg@utk.edu

Version dated: Tuesday 3$^{\text{rd}}$ October, 2017

# Abstract

We present a phylogenetic approach rooted in the field of population genetics that more realistically models the evolution of protein-coding DNA under the assumption of stabilizing selection for a gene specific, optimal amino acid sequence. In addition to being consistent with the fundamental principles of population genetics, our new set of models, which we collectively call SelAC (Selection on Amino acids and Codons), fit phylogenetic data much better than popular models, suggesting strong potential for more accurate inference of phylogenetic trees and branch lengths. SelAC also demonstrates that a large amount of biologically meaningful information is accessible when using a nested set of mechanistic models. For example, for each position SelAC provides a probabilistic estimate of any given amino acid being optimal. SelAC also assumes the strength of selection is proportional to the expression level of a gene and, therefore, provides gene specific estimates of protein synthesis rates. Finally, because SelAC's is a nested approach based on clearly stated biological assumptions, it can be expanded or simplified as needed.

28   Phylogenetic analysis now plays a critical role in most aspects of biology,

29   particularly in the fields of ecology, evolution, paleontology, medicine, and conservation.

30   While the scale and impact of phylogenetic studies has increased substantially over the

31   past two decades, by comparison the realism of the mathematical models on which these

32   analyses are based has changed relatively little. For example, the simplest but most

33   popular models are nucleotide-based, which are naturally agnostic with regards to the

34   different amino acid substitutions and their impact on gene function (e.g. F81, F84,

35   HYK85, TN93, and GTR, see Yang (2014) for an overview).

36   Another set of models attempt to include a 'selection' term $\omega$, but the link between

37   $\omega$ and the key parameters found in standard population genetics models such as $N_e$, the

38   distribution of fitness across genotype space, and mutation bias is far from clear. For

39   instance, $\omega$ is generally interpreted as indicating whether a sequence is under 'purifying'

40   ($\omega < 1$) or 'diversifying' ($\omega > 1$) selection. However, the actual behavior of the model is

41   quite different. When $\omega < 1$ the model behaves as if the resident amino acid $i$ at a given

42   site is favored by selection since synonymous substitutions have a higher substitution rate

43   than any possible non-synonymous substitutions. Paradoxically, this selection regime for

44   the resident amino acid $i$ persists *until* a substitution for another amino acid, $j$, occurs. As

45   soon as amino acid $j$ fixes, but not before, selection now favors amino acid $j$ over all other

46   amino acids, including $i$. This is now the opposite scenario to when $i$ was the resident.

47   Similarly, when $\omega > 1$, synonymous substitutions have a lower substitution rate than any

48   possible non-synonymous substitutions the resident amino acid. In a parallel manner, this

49   selection *against* the resident amino acid $i$ persists until a substitution occurs at which

50   point selection now *favors* the former resident amino acid $i$ as well as the 18 others. Thus,

51   the simplest and most consistent interpretation of $\omega$ is that it describes the rate at which

52   the selection regime itself changes, and this change in selection perfectly coincides with the

53   fixation of a new amino acid. As a result, $\omega$ based approaches only reasonably describe a

subset of scenarios such as over/underdominance or frequency dependent selection (Hughes and Nei 1988; Nowak 2006). Because, as we show here, $\omega$ is well correlated with gene expression, its value is really an indicator of the strength of stabilizing selection on a coding sequence, rather than the 'nature' of that selection.

Given the continual growth in computational power available to researchers, it is now possible to utilize a more general set of population genetics based models for the purpose of phylogenetic analysis (e.g. Halpern and Bruno 1998; Robinson et al. 2003; Lartillot and Philippe 2004; Rodrigue and Lartillot 2014). One lesson from the field of population genetics is even when there are only a few fundamental evolutionary forces at play (mutation, drift, selection, and linkage effects), describing the evolutionary behavior of a system in which there are non-linear interactions between sites, such as epistasis, quickly becomes extremely challenging. The model formulation we evaluate here is a basic version of a more general cost-benefit model we've developed elsewhere (Gilchrist 2007; Gilchrist et al. 2009; Shah and Gilchrist 2011; Gilchrist et al. 2015). This basic version carefully avoids any non-linear interactions between evolutionary forces, resulting in simple additive effects between amino acid sites. This additivity between sites is critical to ensuring that calculation of our amino acid substitution matrix is relatively straightforward, though still computationally intensive. This additivity between sites also means our model could be generalized further and simply posed as a more generic, non-mechanistic, additive model. While often useful in the early stages of a field's development, given the maturity of the field of phylogenetics, we believe such model generalization is now counterproductive. The misinterpretation of GY94's $\omega$ we discuss above is a case in point. Another example, which we touch upon in the Discussion, is the fact that models, which include codon specific effects that presumably affect protein assembly costs and amino acid specific effects, presumably also affect protein functionality, inherently leading to epistatic interactions between sites. While this epistasis may be negligible under certain conditions, identifying

such conditions is only possible with more mechanistic models, such as the cost-benefit one we present here.

# Materials & Method

We model the substitution process as a classic Wright-Fisher process which includes the forces of mutation, selection, and drift (Fisher 1930; Kimura 1962; Wright 1969; Iwasa 1988; Berg and Lässig 2003; Sella and Hirsh 2005; McCandlish and Stoltzfus 2014). For simplicity, we ignore linkage effects and, as a result of this and other assumptions, our method behaves in a site independent manner. Our approach, which we call SelAC (Selection on Amino acids and Codons), is developed in the same vein as previous phylogenetic applications of the Wright-Fisher process (e.g. Muse and Gaut 1994; Halpern and Bruno 1998; Yang and Nielsen 2008; Rodrigue et al. 2005; Koshi and Goldstein 1997; Koshi et al. 1999; Dimmic et al. 2000; Thorne et al. 2012; Lartillot and Philippe 2004; Rodrigue and Lartillot 2014). Similar to Lartillot's work (Lartillot and Philippe 2004; Rodrigue and Lartillot 2014), we assume there is a finite set of rate matrices describing the substitution process and that each position within a protein is assigned to a particular rate matrix category. Unlike this previous work, we assume *a priori* there are 20 different families of rate matrices, one family for when a given amino acid is favored at a site. As a result, SelAC allows us to quantitatively evaluate the support for a particular amino acid being favored at a particular position within the protein encoded by a particular gene.

Because SelAC requires twenty families of $61 \times 61$ matrices, the number of parameters needed to implement SelAC would, without further assumptions, be extremely large. To reduce the number of parameters needed, while still maintaining a high degree of biological realism, we construct our gene and amino acid specific substitution matrices using a submodel nested within our substitution model, similar to approaches in Gilchrist (2007); Shah and Gilchrist (2011); Gilchrist et al. (2015).

One advantage of a nested modeling framework is that it requires only a handful of genome-wide parameters such as nucleotide specific mutation rates (scaled by effective population size $N_e$), side chain physicochemical weighting parameters, and a shape parameter describing the distribution of site sensitivities. In addition to these genome-wide parameters, SelAC requires a gene $g$ specific expression parameter $\psi_g$ which describes the average rate at which the protein's functionality is produced by the organism. Currently, $\psi$ is fixed across the phylogeny, though relaxing this assumption is a goal of future work. The gene specific parameter $\psi_g$ is multiplied by additional model terms to make a composite term $\psi'_g$ which scales the strength and efficacy of selection for the optimal amino acid sequence relative to drift. In terms of the functionality of the protein encoded, we assume that for any given gene there exists an optimal amino acid sequence $\vec{a}_*$ and that, by definition, is a complete, error free peptide consisting of $\vec{a}_*$, which provides one unit of the gene's functionality. We also assume that natural selection favors genotypes that are able to synthesize their proteome efficiently than their competitors and that each savings of an high energy phosphate bond per unit time leads to a constant proportional gain in fitness $q$. SelAC also requires the specification (as part of parameter optimization) of an optimal amino acid at each position or site within a coding sequence which, in turn, makes it the largest category of parameters we estimate. Because we use a submodel to derive our substitution matrices, SelAC requires the estimation of a fraction of the parameters required when compared to approaches where the substitution rates are allowed to vary independently (Halpern and Bruno 1998; Lartillot and Philippe 2004; Rodrigue and Lartillot 2014).

As with other phylogenetic methods, we generate estimates of branch lengths and nucleotide specific mutation rates. In addition, because the math behind our model is mechanistically derived, our method can also be used to make quantitative inferences on the optimal amino acid sequence of a given protein as well as the average synthesis rate of

131 each protein used in the analysis. The mechanistic basis of SelAC also means it can be
132 easily extended to include more biological realism and test more explicit hypotheses about
133 sequence evolution.

## *Mutation Rate Matrix μ*

135 We begin with a 4x4 nucleotide mutation matrix that defines a model for mutation rates
136 between individual bases. For our purposes, we rely on the general unrestricted
137 model(Yang 1994, UNREST) because it makes no constraint on the instantaneous rate of
138 change between any pair of nucleotides. In our view, the flexibility and potential for strong
139 asymmetries in the transition among the different nucleotide states, and ultimately among
140 the different codon states, is more consistent with our model. We note, however, that more
141 constrained models, such as the Jukes-Cantor (JC), Hasegawa-Kishino-Yano (HKY), or the
142 general time-reversible model (GTR), can also be used. The 12 parameter UNREST model
143 defines the relative rates of change between a pair of nucleotides. Thus, we arbitrarily set
144 the G→T mutation rate to 1, resulting in 11 free mutation rate parameters in the 4x4
145 mutation nucleotide mutation matrix. The nucleotide mutation matrix is also scaled by a
146 diagonal matrix $\pi$ whose entries correspond to the equilibrium frequencies of each base.
147 These equilibrium nucleotide frequencies are determined by analytically solving $\pi \times \mathbf{Q} = 0$.
148 We use this $\mathbf{Q}$ to populate a $61 \times 61$ codon mutation matrix $\mu$, whose entries $\mu_{i,j}$ describe
149 the mutation rate from codon $i$ to $j$ under a "weak mutation" assumption. That is, the rate
150 of allele fixation is much greater than $N_e\mu$ and $N_e\mu \ll 1$, such that evolution is mutation
151 limited, codon substitutions only occur one nucleotide at a time and, as a result, the rate
152 of change between any pair of codons that differ by more than one nucleotide is zero.
153 While the overall model does not assume equilibrium, we still need to scale our
154 mutation matrices $\mu$. Traditionally, it is rescaled such that at equilibrium, one unit of
155 branch length represents one expected substitution per site. Here, a scaling factor is

156  calculated as the average rate $-\sum_i \mu_i \pi_i = 1$, where $i$ indexes a particular codon in a given

157  gene. The final mutation rate matrix is the original mutation rate matrix multiplied by

158  1/scaling factor, and represents one expected *proposed mutation* per site.

## *Protein Synthesis Cost-Benefit Function $\eta$*

160  SelAC links fitness to the product of the cost-benefit function of a gene $g$, $\eta_g$, and the

161  organism's average target synthesis rate of the functionality provided by gene $g$, $\psi_g$. This is

162  because the average flux energy an organism spends to meet its target functionality

163  provided by gene $g$ is $\eta_g \times \psi_g$. In order to link genotype to our cost-benefit function

164  $\eta = \mathbf{C}/\mathbf{B}$, we begin by defining our benefit function $\mathbf{B}$.

165  *Benefit:.—* Our benefit function $\mathbf{B}$ measures the functionality of the amino acid sequence

166  $\vec{a}_i$ encoded by a set of codons $\vec{c}_i$, i.e. $a(\vec{c}_i) = \vec{a}_i$ relative to that of an optimal sequence $\vec{a}_*$.

167  By definition, $\mathbf{B}(\vec{a}_*) = 1$ and $\mathbf{B}(\vec{a}_i|\vec{a}_*) < 1$ for all other sequences. We assume all amino

168  acids within the sequence contribute to protein function and that this contribution declines

169  as an inverse function of physicochemical distance between each amino acid and the

170  optimal. Formally, we assume that

$$\mathbf{B}(\vec{a}_i|\vec{a}_*) = \left( \frac{1}{n_g} \sum_{p=1}^{n_g} \left(1 + G_p d(a_{i,p}, a_{*,p})\right) \right)^{-1} \tag{1}$$

171  where $n_g$ is the length of the protein, $d(a_{i,p}, a_{*,p})$ is a weighted physicochemical distance

172  between the amino acid encoded in gene $i$ for position $p$ and $a_{*,p}$ is the optimal amino acid

173  for that position of the protein. For simplicity, we define the distance between a stop codon

174  and a sense codon as effectively infinite and, as a result, nonsense mutations are effectively

175  lethal. The term $G_p$ describes the sensitivity of the protein's function to deviation in

176  physicochemical space. There are many possible measures for physiochemical distance; we

use (Grantham 1974) distances by default, though others may be chosen. We assume that $G_p \sim \text{Gamma}\,(\alpha = \alpha_G, \beta = \alpha_G)$ in order to ensure $\mathbb{E}(G_p) = 1$.

At the limit of $\alpha_G \to \infty$, the model collapses to a model with uniform sensitivity of $G_p = 1$ for all positions $p$. $\mathbf{B}(\vec{a}_i|\vec{a}_*)$ is inversely proportional to the average physicochemical deviation of an amino acid sequence $\vec{a}_i$ from the optimal sequence $\vec{a}_*$ weighted by each site's sensitivity to this deviation. $\mathbf{B}(\vec{a}_i|\vec{a}_*)$ can be generalized to include second and higher order terms of the distance measure $d$.

*Cost:*— Protein synthesis involves both direct and indirect assembly costs. Direct costs consist of the high energy phosphate bonds $\sim P$ of ATP or GTP's used to assemble the ribosome on the mRNA, charge tRNA's for elongation, move the ribosome forward along the transcript, and terminate protein synthesis. As a result, direct protein assembly costs are the same for all proteins of the same length. Indirect costs of protein assembly are potentially numerous and could include the cost of amino acid synthesis as well the cost and efficiency with which the protein assembly infrastructure such as ribosomes, aminoacyl-tRNA synthetases, tRNAs, and mRNAs are used. When these indirect costs are combined with sequence specific benefits, the probability of a mutant allele fixing is no longer independent of the rest of the sequence (Gilchrist et al. 2015) and, as a result, model fitting becomes substantially more complex. Thus for simplicity, in this study we ignore indirect costs of protein assembly that vary between genotypes and define,

$$\mathbf{C}(\vec{c}_i) = \text{Energetic cost of protein synthesis.} \tag{2}$$

$$= A_1 + A_2 n \tag{3}$$

where, $A_1$ and $A_2$ represent the direct cost, in high energy phosphate bonds, of ribosome initiation and peptide elongation, respectively, where $A_1 = A_2 = 4 \sim P$.

*Defining Physicochemical Distances*

Assuming that functionality declines with an amino acid $a_i$'s physicochemical distance from the optimum amino acid $a_*$ at each site provides a biologically defensible way of mapping genotype to protein function that requires relatively few free parameters. In addition, SelAC naturally lends itself to model selection since we can compare the quality of SelAC fits using different mixtures of physicochemical properties. Following Grantham (1974), we focus on using composition $c$, polarity $p$, and molecular volume $v$ of each amino acid's side chain residue to define our distance function, but the model and its implementation can flexibly handle a variety of properties. We use the Euclidian distance between residue properties where each property $c$, $p$, and $v$ has its own weighting term, $\alpha_c$, $\alpha_p$, $\alpha_v$, respectively, which we refer to as 'Grantham weights'. Because physicochemical distance is ultimately weighted by a gene's specific average protein synthesis rate $\psi$, another parameter we estimate, there is a problem with parameter identifiablity. Ultimately, the scale of gene expression is affected by how we measure physicochemical distances which, in turn, is determined by our choice of Grantham weights. As a result, by default we set $\alpha_v = 3.990 \times 10^{-4}$, the value originally estimated by Grantham, and recognize that our our estimates of $\alpha_c$ and $\alpha_p$ and $\psi$ are scaled relative to this choice for $\alpha_v$. More specifically,

$$d(a_i, a_*) = \left( \alpha_c \left[ c\left(a_i\right) - c\left(a_*\right) \right]^2 + \alpha_p \left[ p\left(a_i\right) - p\left(a_*\right) \right]^2 + \right.$$
$$\left. \alpha_v \left[ v\left(a_i\right) - v\left(a_*\right) \right]^2 \right)^{1/2}.$$

## *Linking Protein Synthesis to Allele Substitution*

Next we link the protein synthesis cost-benefit function $\eta$ of an allele with its fixation

probability. First, we assume that each protein encoded within a genome provides some

beneficial function and that the organism needs that functionality to be produced at a target average rate $\psi$. By definition, the optimal amino acid sequence for a given gene, $\vec{a}_*$, produces one unit of functionality. Second, we assume that protein expression is regulated by the organism to ensure that functionality is produced at rate $\psi$. As a result, the realized average protein synthesis rate of a gene, $\phi$, by definition, satisfies the equality $\phi = \psi/\mathbf{B}(\vec{a})$. In other words, the average production rate of a protein $\vec{a}$ with relative functionality $\mathbf{B}(\vec{a}) < 1$ must be $1/\mathbf{B}(\vec{a})$ times higher than the production rate needed if the optimal amino acid sequence $\vec{a}_*$ was encoded since, by definition, $\mathbf{B}(\vec{a}_*) = 1$. required to meet meeting the target functionality of a particular gene is the target functionality production rate $\psi$ multiplied $\eta(\vec{c})\psi$. For example, a cell with an allele $\vec{a}$ where $\mathbf{B}(\vec{a}) = 0.9$ will have to produce $1/9 \times 100\% = 11.11\%$ more proteins than a competitor cell with the optimal allele $\vec{a}_*$ at that locus Similarly, a cell with an allele $\vec{a}$ where $\mathbf{B}(\vec{a}) = 0.8$ will have to produce $1/8 \times 100\% = 12.5\%$ more proteins a cell with $\vec{a}_*$. Simply put, the fitness cost for a genotype encoding a suboptimal protein sequence stems from the need to produce suboptimal proteins at a higher rate in order to compensate for their lower functionality.

Third, we assume that every additional high energy phosphate bond, $\sim P$, spent per unit time to meet the organism's target function synthesis rate $\psi$ leads to a slight and proportional decrease in fitness $W$. This assumption, in turn, implies

$$W_i\left(\vec{c}\right) \propto \exp\left[-A_0\,\eta(\vec{c}_i)\psi\right]. \tag{4}$$

where $A_0$ describes the decline in fitness with every $\sim P$ wasted per unit time. Because $A_0$ shares the same time units as $\psi$ and $\phi$ and only occurs in SelAC in conjunction with $\psi$, we do not need to explicitly identify our time units.

Correspondingly, the ratio of fitness between two genotypes is,

$$W_i/W_j = \exp\left[-A_0\,\eta(\vec{c}_i)\psi\right] / \exp\left[-A_0\,\eta(\vec{c}_j)\psi\right] \tag{5}$$

$$= \exp\left[-A_0\left(\eta(\vec{c}_i) - \eta(\vec{c}_j)\right)\psi\right] \tag{6}$$

$$\tag{7}$$

Given our formulations of $\mathbf{C}$ and $\mathbf{B}$, the fitness effects between sites are multiplicative and, therefore, the substitution of an amino acid at one site can be modeled independently of the amino acids at the other sites within the coding sequence. As a result, the fitness ratio for two genotypes differing at a single site $p$ simplifies to

$$\frac{W_i}{W_j} = \exp\left[-\frac{A_0\left(A_1 + A_2 n_g\right)}{n_g}\right.$$
$$\left. \times \sum_{p \in \mathbb{P}} \left[d\left(a_{i,p}, a_{*,p}\right) - d\left(a_{j,p}, a_{*,p}\right)\right] G_p \psi \right]$$

where $\mathbb{P}$ represents the codon positions in which $\vec{c}_i$ and $\vec{c}_j$ differ. Fourth, we make a weak mutation assumption, such that alleles can differ at only one position at any given time, i.e. $|\mathbb{P}| = 1$, and that the population is evolving according to a Fisher-Wright process. As a result, the probability a new mutant, $j$, introduced via mutation into a resident population $i$ with effective size $N_e$ will go to fixation is,

$$u_{i,j} = \frac{1 - (W_i/W_j)^b}{1 - (W_i/W_j)^{2N_e}}$$
$$= \frac{1 - \exp\left\{-\frac{A_0}{n_g}\left(A_1 + A_2 n_g\right)\left[d\left(a_i, a_*\right) - d\left(a_j, a_*\right)\right] G_p \psi\, b\right\}}{1 - \exp\left\{-\frac{A_0}{n_g}\left(A_1 + A_2 n_g\right)\left[d\left(a_i, a_*\right) - d\left(a_j, a_*\right)\right] G_p \psi\, 2N_e\right\}}$$

where $b = 1$ for a diploid population and 2 for a haploid population (Kimura 1962; Wright 1969; Iwasa 1988; Berg and Lässig 2003; Sella and Hirsh 2005). Finally, assuming a

constant mutation rate between alleles $i$ and $j$, $\mu_{i,j}$, the substitution rate from allele $i$ to $j$ can be modeled as,

$$q_{i,j} = \frac{2}{b}\mu_{i,j}N_e u_{i,j}.$$

208  where, given our weak mutation assumption, $\mu_{i,j} = 0$ when two codons differ by more than

209  one nucleotide. In the end, each optimal amino acid has a separate 64 x 64 substitution

210  rate matrix $\mathbf{Q}_a$, which incorporates selection for the amino acid (and the fixation rate

211  matrix this creates) as well as the common mutation parameters across optimal amino

212  acids. This results in the creation of 20 $\mathbf{Q}_a$ matrices, one for each amino acid, with up to

213  26,880 unique rates, based on few parameters (one to 11 mutation rates, two free

214  Grantham weights, the cost of protein assembly, $A_1$ and $A_2$, the gene specific target

215  functionality synthesis rate $\psi$, and optimal amino acid at each position $p$, $a_{*,p}$), which can

216  either be specified *a priori* or estimated from the data. SelAC can be generalized to allow

217  transitions between optimal amino acids as well as between codons, which would result in a

218  $(20 \times 64) \times (20 \times 64) = 1344 \times 1344$ matrix.

219      Given our assumption of independent evolution among sites, it follows that the

220  probability of the whole data set is the product of the probabilities of observing the data at

221  each individual site. Thus, the likelihood $\mathcal{L}$ of amino acid $a$ being optimal at a given site

222  position $p$ is calculated as

$$\mathcal{L}\left(\mathbf{Q}_a|\mathbf{D}_p, \mathbf{T}\right) \propto \mathbf{P}\left(\mathbf{D}_p|\mathbf{Q}_a, \mathbf{T}\right) \tag{8}$$

223  In this case, the data, $\mathbf{D}_p$, are the observed codon states at position $p$ for the tips of the

224  phylogenetic tree with topology $\mathbf{T}$. For our purposes we take $\mathbf{T}$ as given but it could be

225  estimated as well. The pruning algorithm of Felsenstein (1981) is used to calculate $\mathcal{L}(\mathbf{Q}_a)$.

226  The log of the likelihood is maximized by estimating the genome scale parameters which

227  consist of 11 mutation parameters which are implicitly scaled by $2N_e/b$, and two Grantham

distance parameters, $\alpha_c$ and $\alpha_p$, and the sensitivity distribution parameter $\alpha_G$. Because $A_0$ and $\psi_g$ always co-occur and are scaled by $N_e$, for each gene $g$ we estimate a composite term $\psi_g' = \psi_g A_0 b N_e$ and the optimal amino acid for each position $a_{*,p}$ of protein. When estimating $\alpha_G$, the likelihood then becomes the average likelihood which we calculate using the generalized Laguerre quadrature with $k = 4$ points (Felsenstein 2001).

Finally, we note that because we infer the ancestral state of the system, our approach does not rely on any assumptions of model stationary. Nevertheless, as our branch lengths grow the probability of observing a particular amino acid $a$ at a given site approaches a stationary value proportional to $W(a)^{2N_e - b}$ (Sella and Hirsh 2005).

# *Implementation*

All methods described above are implemented in the new R package, `selac` available through GitHub (`https://github.com/bomeara/selac`) [it will be uploaded to CRAN once peer review has completed]. Our package requires as input a set of fasta files that contain each coding sequence for a set of taxa, and the phylogeny depicting the hypothesized relationships among them. In addition to the SelAC models, we implemented the GY94 codon model of Goldman and Yang (1994), the FMutSel0 mutation-selection model of Yang and Nielsen (2008), and the standard general time-reversible nucleotide model that allows for $\Gamma$ distributed rates across sites. These likelihood-based models represent a sample of the types of popular models often fit to codon data.

For the SelAC models, the starting guess for the optimal amino acid at a site comes from 'majority' rule, where the initial optimum is the most frequently observed amino acid at a given site (ties resolved randomly). Our optimization routine then proceeds by cycling through multiple phases. The first phase optimizes the branch lengths while holding the model parameters constant. The second phase optimizes the gene specific composite parameter $\psi' = A_0 \psi N_e$ across genes, while holding constant both the branch lengths and

the model parameters shared across the genome (i.e., $\alpha_c$ and $\alpha_p$, and the sensitivity distribution parameter $\alpha_G$). This is followed by a third phase that optimizes the parameters across the genome, while keeping the branch lengths and the composite parameters constant. Finally, the fourth phase estimates the optimal amino acid at each site while keeping the branch lengths and all model parameters at their current values. This entire procedure is repeated six times. For optimization of a given set of parameters, we rely on a bounded subplex routine (Rowan 1990) in the package `NLopt` (Johnson 2012) to maximize the log-likelihood function. To help the optimization navigate through local peaks, we perform a set of independent analyses with different sets of naive starting points with respect to the gene specific composite $\psi'$ parameters, $\alpha_c$, and $\alpha_p$. Confidence in the parameter estimates can be generated by an 'adaptive search' procedure that we implemented to provide an estimate of the parameter space that is some pre-defined likelihood distance (e.g., 2 lnL units) from the maximum likelihood estimate (MLE), which follows Beaulieu and OMeara (2016); Edwards (1984).

We note that our current implementation is painfully slow, and is particularly suited for smaller data sets in terms of numbers of taxa. This is largely due to the size and quantity of matrices we create and manipulate just to calculate the log-likelihood of an individual site. We have parallelized operations wherever possible, but the fact remains that, long term, this model may not be well-suited for R. Ongoing work will address the need for speed, with the eventual goal of implementing the model in popular phylogenetic inference toolkits, such as RevBayes (Hhna et al. 2016), PAML (Yang 2007) and RAxML (Stamatakis 2006).

## *Simulations*

We evaluated the performance of our codon model by simulating datasets and estimating the bias of the inferred model parameters from these data. Our 'known' parameters under

278 a given generating model were based on fitting SelAC to the 106 gene data set and

279 phylogeny of Rokas et al. (2003). The tree used in these analyses is outdated with respect

280 to the current hypothesis of relationships within Saccharomyces, but we rely on it simply as

281 a training set that is separate from our empirical analyses (see section on Analyzing Yeast

282 Genome). Bias in the model parameters were assessed under two generating models: one

283 where we assumed a model of SelAC assuming $\alpha_G = \infty$, and one where we estimated $\alpha_G$

284 from the data. Under each of these two scenarios, we used parameter estimates from the

285 corresponding empirical analysis and simulated 50 five-gene data sets. For the gene specific

286 composite parameter $\psi'_g$ the 'known' values used for the simulation were five evenly spaced

287 points along the rank order of the estimates across the 106 genes. The MLE estimate for a

288 given replicate were taken as the fit with the highest log-likelihood after running five

289 independent analyses with different sets of naive starting points with respect to the

290 composite $\psi'_g$ parameter, $\alpha_c$, and $\alpha_p$. All analyses were carried out in our `selac` R package.

### *Analysis of yeast genome and tests of model adequacy*

292 We focus our empirical analyses on the large yeast data set and phylogeny of Salichos and

293 Rokas (2013). The yeast genome is an ideal system to examine our phylogenetic estimates

294 of gene expression and its connection to real world measurements of these data within

295 individual taxa. The complete data set of Salichos and Rokas (2013) contain 1070

296 orthologs, where we selected 100 at random for our analyses. We also focus our analyses

297 only on Saccharomyces *sensu stricto*, including their sister taxon *Candida glabrata*, and we

298 rely on the phylogeny depicted in Fig. 1 of Salichos and Rokas (2013) for our fixed tree. We

299 fit the two SelAC models described above (i.e., SelAC and SelAC+Γ), as well as two codon

300 models, GY94 and FMutSel0, and a standard GTR + Γ nucleotide model. The FMutSel0

301 model, which assumes that the amino acid frequencies are determined by functional

302 requirements of the protein, is the most similar to our model. In all cases, we assumed that

303 the model was partitioned by gene, but with branch lengths linked across genes.

304 For SelAC, we compared our estimates of $\phi' = \psi'/\mathbf{B}$, which represents the average

305 protein synthesis rate of a gene, to estimates of gene expression from empirical data.

306 Specifically, we obtained expression data for five of the six species used - four species were

307 measured during log-growth phase, whereas the other was measured at the beginning of the

308 stationary phase (*S. kudriavzevii*) from the Gene Expression Omnibus (GEO). Gene

309 expression was measured using either Microarray chips (*C. glabrata*, *S. castellii*, and *S.*

310 *kudriavzevii*) or RNA-Seq (*S. paradoxus*, *S. mikatae*, and *S. cerevisiae*). For further

311 comparison, we also predicted protein synthesis rate ($\phi$) by analyzing gene and

312 genome-wide patterns of synonymous codon usage using ROC-SEMPPR (Gilchrist et al.

313 2015) for each individual genome. While, like SelAC, ROC-SEMPPR uses codon level

314 information, it does not rely on any inter-specific comparisons and, unlike SelAC, assumes

315 selection on synonymous codon usage is contributing to these patterns. Nevertheless,

316 ROC-SEMPPR predictions of gene expression $\phi$ correlates strongly ($r = 0.53 - 0.74$) with

317 a wide range of laboratory measurements of gene expression (Gilchrist et al. 2015).

318 While one of our main objectives was to determine the improvement of fit that

319 SelAC has with respect to other standard phylogenetic models, we also evaluated the

320 adequacy of SelAC. Model fit, measured with assessments such as the Akaike Information

321 Criterion (AIC), can tell which model is least bad as an approximation for the data, but it

322 does not reveal whether a model is actually doing a good job of representing the biological

323 processes. An adequate model does the latter, one measure of which is that data generated

324 under the model resemble real data (Goldman 1993). For example, Beaulieu et al. (2013)

325 assessed whether parsimony scores and the size of monomorphic clades of empirical data

326 were within the distributions of simulated under a new model and the best standard

327 model; if the empirical summaries were outside the range for each, it would have suggested

328  that neither model was adequately modeling this part of the biology.

329       For a given gene we first remove a particular taxon from the data set and the

330  phylogeny. A marginal reconstruction of the likeliest sequence across all remaining nodes is

331  conducted under the model, including where the attachment point of pruned taxon to the

332  tree. The marginal probabilities of each site are used to sample and assemble the starting

333  coding sequence. This sequence is then evolved along the branch, periodically being

334  sampled and its current functionality assessed. We repeat this process 100 times and

335  compare the distribution of trajectories against the observed functionality calculated for

336  the gene. For comparison, we also conducted the same test, by simulating the sequence

337  under the standard GTR $+ \Gamma$ nucleotide model, which is often used on these data but does

338  not account for the fact that the sequence codes for a specific protein, and under FMutSel0,

339  which includes selection on codons but in a fundamentally different way as our model.

## The appropriate estimator of bias for AIC

341  As part of the model set described above, we also included a reduced form of each of the

342  two SelAC models. Specifically, rather than optimizing the amino acid at any given site, we

343  assume the optimal is the most abundantly observed amino acid at each site, which greatly

344  decreased the computation time. We refer to these 'majority rule' models as $SelAC_M$ and

345  $SelAC_M + \Gamma$. Since these models assume that the optimal amino acids are known prior to

346  fitting of our model, it is tempting to reduce the number of parameters in the model by the

347  number of total sites being analyzed. Despite having become standard behavior in the field

348  of phylogenetics, this reduction is statistically inappropriate due to the fact that

349  determination of the majority rule amino acid is made by examining the data prior to the

350  fitting of the model. Because the difference in $K$ when counting or not counting number of

351  nucleotide sites drops out when comparing nucleotide models with AIC, this statistical

352  issue does not apply to nucleotide models. It does, however, matter for AICc, where the

number of parameters, $K$, and the sample size, $n$, combine in the penalty term. This also matters in our case, where the number of estimated parameters for the majority rule estimation differs based on whether one is looking at codons or single nucleotides.

In phylogenetics two variants of AICc are used. In comparative methods (e.g. Butler and King 2004; O'Meara et al. 2006; Beaulieu et al. 2013) the number of data points, $n$, is taken as the number of taxa. More taxa allow the fitting of more complex models, given more data. However, in DNA evolution, which is effectively the same as a discrete character model used in comparative methods, the $n$ is taken as the number of sites. Obviously, both cannot be correct. The original derivation of AICc by Hurvich and Tsai (1989) assumed a regression model, where the true model was in the set of examined models, as well as approximations in the derivation itself. It might not be an appropriate approximation for phylogenetic data, where data points are not independent of each other. In any case, we argue that for phylogenetic data, a good estimate of data set size is number of taxa multiplied by number of sites. First of all, this is what is conventionally seen as the size of the dataset in the field. Second, when considering how likelihood is calculated, the likelihood for a given site is the sum of the probabilities of each observed state at each tip, and this is then multiplied across sites. It is arguable that the conventional approach in comparative methods is calculating AICc in this way: number of taxa multiplied by number of sites equals the number of taxa, if only one site is examined, as remains remarkably common in comparative methods. An notable exception is the program SURFACE implemented by Ingram and Mahler (2013), which uses multiple characters and taxa. While its default is to use AIC to compare models, if one chooses to use AICc, the number of samples is taken as the product of number of sites and number of taxa. Recently, Jhwueng et al. (2014) performed an analysis that investigated what variant of AIC and AICc worked best as an estimator, but the results were inconclusive. Here, we have adopted and extended the simulation approach of Jhwueng et al. (2014) in order to

examine a large set of different penalty functions and how well they approximate the

remaining portion of Kullback-Liebler (KL) divergence between two models after

accounting for the deviance (i.e., $-2\mathcal{L}$). Please see Appendix 1 for details.

# RESULTS

By linking transition rates $q_{i,j}$ to gene expression $\psi$, our approach allows use of the same

model for genes under varying degrees of stabilizing selection. Specifically, we assume the

strength of stabilizing selection for the optimal sequence, $\vec{a}_*$, is proportional to the average

protein synthesis rate $\phi$, which we can estimate for each gene. In regards to model fit, our

results clearly indicated that linking the strength of stabilizing selection for the optimal

sequence to gene expression substantially improves our model fit. Further, including the

single random effects term $G \sim \text{Gamma}(\alpha_G, \beta_g)$ to allow for heterogeneity in this selection

between sites within a gene, improves the $\Delta$AICc of SelAC+$\Gamma$ (i.e., includes a shape

parameter on the gamma distributed site-specific sensitivity to protein function) score over

the simpler SelAC models by over 22,000 AIC units. Using either $\Delta$AICc or AIC$_\text{w}$ as our

measure of model support, the SelAC models fit extraordinarily better than GTR + $\Gamma$,

GY94, or FMutSel0 (Table 1). This is in spite of the need for estimating the optimal amino

acid at each position in each protein, which accounts for 49,881 additional model

parameters. Even when compared to the next most parameter rich codon model in our

model set, FMutSel0, SelAC+$\Gamma$ model shows nearly 180,000 AIC unit improvement over

FMutSel0.

With respect to estimates of $\phi$ within SelAC, they were strongly correlated with two

separate measures of gene expression, one empirical (See Figure 1), and one model-based

prediction that does not account for shared ancestry (Figure S1-S2). In other words, using

only codon sequences our model can predict which genes have high or low expression levels.

The estimate of the $\alpha_G$ parameter, which describes the site-specific variation in sensitivity of the protein's functionality, indicated a moderate level of variation in gene expression among sites. Our estimate of $\alpha_G = 1.40$, produced a distribution of sensitivity terms $G$ ranged from 0.344-7.16, but with nearly 90% of the weight for a given site-likelihood being contributed by the 0.344 and 1.48 rate categories. In simulation, however, of all the parameters in the model, only $\alpha_G$ showed a consistent bias, in that the estimates were generally underestimated (see Supporting Materials). Other parameters in the model, such as the Grantham weights, provide an indication as to the physicochemical distance between amino acids. Our estimates of these weights only strongly deviate from Grantham's 1974 original estimates in regards to composition weight, $\alpha_c$, which is the ratio of noncarbon elements in the end groups to the number of side chains. Our estimate of the composition weighting factor of $\alpha_c$=0.484 is 1/4th the value estimate by Grantham which suggests that the substitution process is less sensitive to this physicochemical property when shared ancestry and variation in stabilizing selection are taken into account.

It is important to note that the nonsynonymous/synonymous mutation ratio, or $\omega$, which we estimated for each gene under the FMutSel0 model strongly correlated with our estimates of $\phi' = \psi'/\mathbf{B}$ where $\mathbf{B}$ depends on the sequence of each taxa. In fact, $\omega$ showed similar, though slightly reduced correlations, with the same empirical estimates of gene expression described above (See Figure 2). This would give the impression that the same conclusions could have been gleaned using a much simpler model, both in terms of the number of parameters and the assumptions made. However, as we discussed earlier, not only is this model greatly restricted in terms of its biological feasibility, SelAC clearly performs better in terms of its fit to the data and biological realism. For example, when we simulated the sequence for *S. cervisieae*, starting from the ancestral sequence under both GTR + $\Gamma$ and FMutSel0, the functionality of the simulated sequence moves away from the observed sequence, whereas SelAC remains near the functionality of the observed sequence

429 (Figure 3b). In a way, this is somewhat unsurprising, given that both GTR + $\Gamma$ and

430 FMutSel0 are agnostic to the functionality of the gene, but it does highlight the

431 improvement in biological realism in amino acid sequence evolution that SelAC provides.

432 We do note that the adequacy of the SelAC model does vary among individual taxa, and

433 does not always perfectly match the observed functionality. For instance, *S. castellii* is

434 simulated with consistently higher functionality than observed (Figure 3c). We suspect this

435 is an indication that assuming a single set of optimal amino acid across all taxa may be too

436 simplistic, but we cannot also rule out other potential simplifying assumptions in our

437 model, such as a single set of Grantham weights and $\alpha_G$ values or the simple, inverse

438 relationship between physicochemical distance $d$ and benefit $\mathbf{B}$.

439       Finally, we note that our simulation analysis suggested that the best measure of

440 dataset size for AICc uses a scaled value of the product of number of sites and number of

441 characters was the best at estimating KL distance. The model comparison approach

442 described above included this assumption. For more details on the simulation approach, see

443 Appendix 1.


# DISCUSSION

445 The work presented here contributes to the field of phylogenetics and molecular evolution

446 in a number of ways. First, SelAC provides an complementary example to Thorne et al.

447 (2012) studies of how models of molecular and evolutionary scales can be combined

448 together in a nested manner. While the mapping between genotype and phenotype is more

449 abstract than Thorne et al. (2012), SelAC has the advantage of not requiring knowledge of

450 a protein's native folding. Second, our use of model nesting also allows us to formulate and

451 test specific biological hypotheses. For example, we are able to compare a model

452 formulation which assumes that physiochemical deviations from the optimal sequence are

equally disruptive at all sites within a protein to one which assumes the effect of deviation from the optimal amino acid's physicochemical properties on protein function varies between sites. By linking the strength of stabilizing selection for an optimal amino acid sequence to gene expression, we can weight the historical information encoded in genes evolving at vastly different rates in a biologically plausible manner while simultaneously estimating their expression levels. Finally, because our fitness functions are well defined, we can provide estimates of key evolutionary statistics such as the distribution of the effects of amino acid substitutions on fitness and genetic load.

As phylogenetic methods become ever more ubiquitous in biology, and data set size and complexity increase, there is a need and an opportunity for more complex and realistic models (Goldman et al. 1996; Thorne et al. 1996; Goldman et al. 1998; Halpern and Bruno 1998; Lartillot and Philippe 2004). Despite their widespread use, phylogenetic models based on purifying and diversifying selection, i.e. Goldman and Yang (1994) and extensions, are very narrow categories of selection that mostly apply to cases of positive and negative frequency dependent selection at the level of a particular amino acid, not for tree inference itself.

Instead of heuristically extending population genetic models of neutral evolution for use in phylogenetics, it makes sense to derive these extensions from population genetic models that *explicitly* include the fundamental forces of mutation, drift, and natural selection. Starting with Halpern and Bruno (1998), a number of researchers have developed methods for linking site-specific selection on protein sequence and phylogenetics(e.g. Koshi et al. 1999; Dimmic et al. 2000; Koshi and Goldstein 2000; Robinson et al. 2003; Lartillot and Philippe 2004; Thorne et al. 2012; Rodrigue and Lartillot 2014). Our work follows this tradition, but includes some key advances. For instance, even though SelAC requires a large number of matrices, because of our assumption about protein functionality and physicochemical distance from the optimum, we are able to parameterize our substitution

matrices using a relatively small number of genome-wide parameters and one gene specific parameter. We show that all of these parameters can be estimated simultaneously with branch lengths from the data at the tips of the tree.

By assuming fitness declines with extraneous energy flux, SelAC explicitly links the variation in the strength of stabilizing selection for the optimal protein sequence among genes, to the variation among genes in their target expression levels $\psi$. Furthermore, by linking expression and selection, SelAC provides a natural framework for combining information from protein coding genes with very different rates of evolution with the low expression genes providing information on shallow branches and the high expression genes providing information on deep branches. This is in contrast to more traditional approach of concatenating gene sequences together, which is equivalent to assuming the same average protein synthesis rate $\psi$ for all of the genes, or more recent approaches where different models are fitted to different genes. Our results indicate that including a gene specific $\psi$ value vastly improves SelAC fits (Table 1). Perhaps more convincingly, we find that the target expression level $\psi$ and realized protein synthesis rate $\phi$ are reasonably well correlated with laboratory measurements of gene expression ($r = 0.34 - 0.65$; Figures 1, S1, and S2). The idea that quantitative information on gene expression is embedded within intra-genomic patterns of synonymous codon usage is well accepted; our work shows that this information can also be extracted from comparative data at the amino acid level.

Of course, given the general nature of SelAC and the complexity of biological systems, other biological forces besides selection for reducing energy flux likely contribute intergenic variation in the magnitude of stabilizing selection. Similarly, other physicochemical properties besides composition, volume, and charge likely contribute to site specific patterns of amino acid substitution. Thus, a larger and more informative set of Grantham weights might improve our model fit and reduce the noise in our estimates of $\phi$. Even if other physicochemical properties are considered, the idea of a consistent, genome

wide Grantham weighting of these terms seems highly unlikely. Since the importance of an amino acid's physicochemical properties likely changes with where it lies in a folded protein, one way to incorporate such effects is to test whether the data supports multiple sets of Grantham weights for either subsets of genes or regions within genes, rather than a single set.

Both of these points highlight the advantage of the detailed, mechanistic modeling approach underlying SelAC. Because there is a clear link between protein expression, synthesis cost, and functionality, SelAC can be extended by increasing the realism of the mapping between these terms and the coding sequences being analyzed. For example, SelAC currently assumes the optimal amino acid for any site is fixed along all branches. This assumption can be relaxed by allowing the optimal amino acid to change during the course of evolution along a branch.

From a computational standpoint, the additive nature of selection between sites is desirable because it allows us to analyze sites within a gene largely independently of each other. From a biological standpoint, this additivity between site ignores any non-linear interactions between sites, such as epistasis, or between alleles, such as dominance. Thus, our work can be considered a first step to modeling to these more complex scenarios. For example, our current implementation ignores any selection on synonymous codon usage bias (CUB) (Yang and Nielsen 2008; Pouyet et al. 2016, c.f. ). Including such selection is tricky because introducing the site specific cost effects of CUB, which is consistent with the hypothesis that codon usage affects the efficiency of protein assembly or $\mathbf{C}$, into a model where amino acids affect protein funcdtion or $\mathbf{B}$, results in a cost-benefit ratio $\mathbf{C}/\mathbf{B}$ with epistatic interactions between all sites. These epistatic effects can likely be ignored under certain conditions or reasonably approximated based on an expectation of codon specific costs (e.g. Kubatko et al. 2016). Nevertheless, it is difficult to see how one could identify such conditions without modeling the way in which codon and amino acid usage affects

**C/B**.

This work also points out the potential importance of further investigation into model choice in phylogenetics. For likelihood models, use of AICc has become standard, but, in our view, the conflict in how number of data points is calculated remains. Common sense suggests that dataset size is increased by adding taxa and/or sites. In other words, a dataset of 1000 taxa and 100 sites must have more information on substitution models than a dataset of 4 taxa and 100 sites. Our simple analyses agree that the number of observations in a dataset (number of sites times number of taxa) should be taken as the sample size for AICc, but we suspect this is still too simple. For instance, one could imagine a phylogeny where one taxon is sister to a polytomy of 99 taxa that have zero length terminal branches. Absent measurement error or other intraspecific variation, one would have 100 species but only two unique trait values, and the only information about the process of evolution comes from one happens on the path connecting the lone taxon to the polytomy. A fully resolved 100-taxon tree with nonzero branch lengths would have much more information than the above tree, despite the same number of taxa. This is just one scenario, we can envision many more, and so a deeper investigation of this issue is clearly warranted.

There are still significant deficiencies in the approach outlined here. Most worrisome are biological flaws in the model. For example, at its heart, the model assumes that suboptimal proteins can be compensated for, at a cost, simply by producing more of them. However, this is likely only true for proteins reasonably close to the optimal sequence. Different enough proteins will fail to function entirely: the active site will not sufficiently match its substrates, a protein will not properly pass through a membrane, and so forth. Yet, in our model, even random sequences still permit survival, just requiring more protein production. Other oversimplifications include the assumption of no selection on codon usage, no change of optimal amino acids through time, and no change of the effect of

physiochemical properties on fitness through time. However, many of these can be relaxed through further elaborations of the model.

There are also deficiencies in our implementation. Though reasonable to use for a given topology with a modest number of species, it is too slow for practical use for tree search. It thus serves as a proof of concept, or of utility for targeted questions where a more realistic model may be of use (placement of particular taxa, for example). Future work will encode SelAC models into a variety of mature, popular tree-search programs. SelAC also represents a hard optimization problem: the nested models reduce parameter complexity vastly, but there are still numerous parameters to optimize, including the discrete parameter of optimal amino acid at each site. A different implementation, more parameter-rich, would optimize values of three (or more) physiochemical properties per site. This would have the practical advantage of continuous parameter optimization rather than discrete, and biologically would be more realistic (as it is the properties that selection "sees", not the identity of the amino acid itself).

Overall, SelAC represents an important step in uniting phylogenetic and population genetic models. It allows biologically relevant population genetic parameters to be estimated from phylogenetic information, while also dramatically improving fit and accuracy of phylogenetic models. Moreover, it demonstrates that there remains substantially more information in the coding sequences used for phylogenetic analysis than other methods acknowledge.

# ACKNOWLEDGEMENTS

*

## References

Beaulieu, J. M., B. C. O'Meara, and M. J. Donoghue. 2013. Identifying Hidden Rate Changes in the Evolution of a Binary Morphological Character: The Evolution of Plant Habit in Campanulid Angiosperms. Systematic Biology 62:725–737.

Beaulieu, J. M. and B. C. OMeara. 2016. Detecting Hidden Diversification Shifts in Models of Trait-Dependent Speciation and Extinction. Systematic Biology 65:583–601.

Berg, J. and M. Lässig. 2003. Stochastic Evolution and Transcription Factor Binding Sites. Biophysics 48:S36–S44.

Butler, M. A. and A. A. King. 2004. Phylogenetic comparative analysis: a modeling approach for adaptive evolution. American Naturalist 164:683–695.

Dimmic, M. W., D. P. Mindell, and R. A. Goldstein. 2000. Modeling evolution at the protein level using an adjustable amino acid fitness model. Pacific Symposium on Biocomputing 5:18–29.

Edwards, A. 1984. Likelihood. Cambridge science classics Cambridge University Press.

Felsenstein, J. 1981. Evolutionary trees from DNA-sequences - a maximum-likelihood approach. Journal of Molecular Evolution 17:368–376.

Felsenstein, J. 2001. Taking Variation of Evolutionary Rates Between Sites into Account in Inferring Phylogenies. Journal of Molecular Evolution 53:447–455.

Fisher, S., Ronald A. 1930. The Genetical Theory of Natural Selection. Oxford University Press, Oxford.

Gilchrist, M., P. Shah, and R. Zaretzki. 2009. Measuring and detecting molecular adaptation in codon usage against nonsense errors during protein translation. Genetics 183:1493–1505.

Gilchrist, M. A. 2007. Combining Models of Protein Translation and Population Genetics to Predict Protein Production Rates from Codon Usage Patterns. Molecular Biology and Evolution 24:2362–2373.

Gilchrist, M. A., W.-C. Chen, P. Shah, C. L. Landerer, and R. Zaretzki. 2015. Estimating Gene Expression and Codon-Specific Translational Efficiencies, Mutation Biases, and Selection Coefficients from Genomic Data Alone. Genome Biology and Evolution 7:1559–1579.

Goldman, N. 1993. Statistical tests of models of DNA substitution. Journal of molecular evolution 36:182–198.

Goldman, N., J. L. Thorne, and D. T. Jones. 1996. Using Evolutionary Trees in Protein Secondary Structure Prediction and Other Comparative Sequence Analyses. Journal of Molecular Biology 263:196 – 208.

Goldman, N., J. L. Thorne, and D. T. Jones. 1998. Assessing the Impact of Secondary Structure and Solvent Accessibility on Protein Evolution. Genetics 149:445–458.

Goldman, N. and Z. H. Yang. 1994. Codon-based model of nucleotide substitution for protein-coding DNA-sequences. Molecular Biology and Evolution 11:725–736.

Grantham, R. 1974. Amino acid difference formula to help explain protein evolution. Science 185:862–864.

Halpern, A. L. and W. J. Bruno. 1998. Evolutionary distances for protein-coding sequences: Modeling site-specific residue frequencies. Molecular Biology And Evolution 15:910–917.

629 Hughes, A. L. and M. Nei. 1988. Pattern of nucleotide substitution at major
630     histocompatibility complex class-i loci reveals overdominant selection. Nature
631     335:167–170.

632 Hurvich, C. M. and C.-L. Tsai. 1989. Regression and time series model selection in small
633     samples. Biometrika 76:297–307.

634 Hhna, S., M. J. Landis, T. A. Heath, B. Boussau, N. Lartillot, B. R. Moore, J. P.
635     Huelsenbeck, and F. Ronquist. 2016. RevBayes: Bayesian Phylogenetic Inference Using
636     Graphical Models and an Interactive Model-Specification Language. Systematic Biology
637     65:726.

638 Ingram, T. and D. L. Mahler. 2013. SURFACE: detecting convergent evolution from data
639     by fitting Ornstein-Uhlenbeck models with stepwise Akaike Information Criterion.
640     Methods in ecology and evolution 4:416–425.

641 Iwasa, Y. 1988. Free fitness that always increases in evolution. Journal of Theoretical
642     Biology 135:265–281.

643 Jhwueng, D.-C., H. Snehalata, B. C. O'Meara, and L. Liu. 2014. Investigating the
644     performance of AIC in selecting phylogenetic models. Statistical applications in genetics
645     and moleculr biology 13:459–475.

646 Johnson, S. G. 2012. The NLopt nonlinear-optimization package. Version 2.4.2 – Released
647     20 May 2014.

648 Kimura, M. 1962. on the probability of fixation of mutant genes in a population. Genetics
649     47:713–719.

650 Koshi, J. M. and R. A. Goldstein. 1997. Mutation matrices and physical-chemical

properties: Correlations and implications. Proteins-Structure Function And Genetics 27:336–344.

Koshi, J. M. and R. A. Goldstein. 2000. Analyzing site heterogeneity during protein evolution. Pages 191–202 *in* Biocomputing 2001. World Scientific.

Koshi, J. M., D. P. Mindell, and R. A. Goldstein. 1999. Using physical-chemistry-based substitution models in phylogenetic analyses of HIV-1 subtypes. Molecular biology and evolution 16:173–179.

Kubatko, L., P. Shah, R. Herbei, and M. A. Gilchrist. 2016. A codon model of nucleotide substitution with selection on synonymous codon usage. Molecular Phylogenetics and Evolution 94:290 – 297.

Lartillot, N. and H. Philippe. 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. Molecular Biology And Evolution 21:1095–1109.

Mayrose, I., N. Friedman, and T. Pupko. 2005. A Gamma mixture model better accounts for among site rate heterogeneity. Bioinformatics 21:ii151–ii158.

McCandlish, D. M. and A. Stoltzfus. 2014. Modeling evolution using the probability of fixation: History and implications. The Quarterly Review of Biology 89:225–252.

Muse, S. V. and B. S. Gaut. 1994. A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. Molecular Biology and Evolution 11:715–724.

Nowak, M. A. 2006. Evolutionary Dynamics: Exploring the Equations of Life. Belknap of Harvard University Press, Cambridge, MA.

673  O'Meara, B. C., C. Ane, M. J. Sanderson, and W. P.C. 2006. Testing for different rates of

674    continuous trait evolution using likelihood. Evolution 60:922–933.

675  Pouyet, F., M. Bailly-Bechet, D. Mouchiroud, and L. Guguen. 2016. SENCA: A

676    Multilayered Codon Model to Study the Origins and Dynamics of Codon Usage. Genome

677    Biology and Evolution 8:2427–2441.

678  Robinson, D. M., D. T. Jones, H. Kishino, N. Goldman, and J. L. Thorne. 2003. Protein

679    evolution with dependence among codons due to tertiary structure. Molecular Biology

680    And Evolution 20:1692–1704.

681  Rodrigue, N. and N. Lartillot. 2014. Site-heterogeneous mutation-selection models within

682    the PhyloBayes-MPI package. Bioinformatics 30:1020–1021.

683  Rodrigue, N., N. Lartillot, D. Bryant, and H. Philippe. 2005. Site interdependence

684    attributed to tertiary structure in amino acid sequence evolution. Gene 347:207–217.

685  Rokas, A., B. L. Williams, N. King, and S. B. Carroll. 2003. Genome-scale approaches to

686    resolving incongruence in molecular phylogenies. Nature 425:798–804.

687  Rowan, T. 1990. Functional Stability Analysis of Numerical Algorithms. Ph.D. thesis

688    University of Texas, Austin.

689  Salichos, L. and A. Rokas. 2013. Inferring ancient divergences requires genes with strong

690    phylogenetic signals. Nature 497:327–331.

691  Sella, G. and A. E. Hirsh. 2005. The application of statistical physics to evolutionary

692    biology. Proceedings of the National Academy of Sciences of the United States of

693    America 102:9541–9546.

Shah, P. and M. A. Gilchrist. 2011. Explaining complex codon usage patterns with selection for translational efficiency, mutation bias, and genetic drift. Proceedings of the National Academy of Sciences of the United States of America 108:10231–10236.

Stamatakis, A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. Bioinformatics 22:2688–2690.

Thorne, J. L., N. Goldman, and D. T. Jones. 1996. Combining protein evolution and secondary structure. Molecular Biology and Evolution 13:666–673.

Thorne, J. L., N. Lartillot, N. Rodrigue, and S. C. Choi. 2012. Codon models as a vehicle for reconciling population genetics with inter-specific sequence data. Codon Evolution: Mechanisms And Models Pages 97–110 D2 10.1093/acprof:osobl/9780199601165.001.0001 ER.

Wright, S. 1969. Evolution and the genetics of populations. Vol. 2. The theory of gene frequencies. vol. 2. University of Chicago Press.

Yang, Z. 2014. Molecular Evolution: A Statistical Approach. Oxford University Press, New York.

Yang, Z. H. 1994. Maximum-likelihood phylogenetic estimation from DNA-sequences with variable rates over sites - approximate methods. Journal Of Molecular Evolution 39:306–314.

Yang, Z. H. 2007. PAML 4: Phylogenetic analysis by maximum likelihood. Molecular Biology And Evolution 24:1586–1591.

Yang, Z. H. and R. Nielsen. 2008. Mutation-selection models of codon substitution and their use to estimate selective strengths on codon usage. Molecular Biology and Evolution 25:568–579.

# TABLE

| Model | logLik | Parameters Estimated | AIC | AICc | ΔAICc | Model Weight |
|---|---|---|---|---|---|---|
| GTR+Γ | -655166.4 | 610 | 1,311,553 | 1,311,554 | 287,415 | <0.001 |
| GY94 | -612121.5 | 210 | 1,224,663 | 1,224,663 | 200,524 | <0.001 |
| FMutSel0 | -598848.2 | 2810 | 1,203,316 | 1,203,362 | 179,223 | <0.001 |
| SelAC$_M$ | -478282.7 | 50,004 | 1,056,573 | 1,073,290 | 49,151 | <0.001 |
| SelAC | -465616.7 | 50,004 | 1,031,241 | 1,047,958 | 23,819 | <0.001 |
| SelAC$_M$ + Γ | -465089.7 | 50,005 | 1,030,189 | 1,046,906 | 22,767 | <0.001 |
| SelAC+Γ | -453706.0 | 50,005 | 1,007,422 | 1,024,139 | 0 | >0.999 |

Table 1: Comparison of model fits using AIC, AICc, and AIC$_w$. Note the subscripts $M$ indicate model fits where the most common or 'majority rule' amino acid was fixed as the optimal amino acid $a_*$ for each site. As discussed in text, despite the fact that $a_*$ for each site was not fitted by our algorithm, its value was determined by examining the data and, as a result, represent an additional parameter estimated from the data and are accounted for in our table.

# FIGURES



Figure 1: Comparisons between estimates of $\phi$ obtained from SelAC+$\Gamma$ and direct measurements of expression for individual yeast taxa across the 100 selected genes from Salichos and Rokas (2013). Estimates of $\phi$ were obtained by solving for $\psi$ based on estimates of $\psi'$, and then dividing by $\mathbf{B}(\vec{a}_i|\vec{a}_*)$. Gene expression was measured using either RNA-Seq (a-c) or Microarray chips (d), and the equations in the upper left hand corner of each panel represent the regression fit and correlation coefficient $r$.

Figure 2: Comparisons between $\omega$, which is the nonsynonymous/synonymous mutation ratio in FMutSel0, $\psi$ obtained from SelAC+$\Gamma$ (a), a direct measurement of expression (b), and a model-based prediction of gene expression that does not account for ancestry (c), for *S. cerevisiae* across the 100 selected genes from Salichos and Rokas (2013). As in Figure 1, the equations in the upper left hand corner of each panel provide the regression fit and correlation coefficient. Estimates of $\psi$ were solved from estimates of $\psi'$.

Figure 3: (a) Maximum likelihood estimates of branch lengths under SelAC+Γ for 100 selected genes from Salichos and Rokas (2013). Tests of model adequacy for *S. cerevisiae* (b) and *S. castellii* (c) indicated that, when these taxa are removed from the tree, and their sequences are simulated, the parameters of SelAC+Γ exhibit functionality that is far closer to the observed (dashed black line) than data sets produced from parameters of either FMutSel0 or GTR + Γ.

# Part I

# Supporting Materials

### *Comparisons of SelAC gene expression estimates with empirical measurements*

In our model, the parameter $\phi$ measures the realized average protein synthesis rate of a gene. We compared our estimates of $\phi$ to two separate measures of gene expression, one empirical (See Figure S1), and one model-based prediction that does not account for shared ancestry, for individual yeast taxa across the same set of genes. Our estimates of $\phi$ are positively correlated both measures, which are also strongly correlated with each other (Figure 1 - S2) On the whole, these comparisons indicate not only a high degree of consistency among all three measures, but also, importantly, that estimates of $\phi$ obtained from SelAC provide real biological insight into the expression level of a gene.

Figure S1: Comparisons between estimates of $\phi$ obtained from SelAC+$\Gamma$ and the predicted gene expression from the ROC SEMPER model (Gilchrist et al. (2015)) for individual yeast taxa across the 100 selected genes from Salichos and Rokas (2013). As with figures in the main text, estimates of $\phi$ were obtained by solving for $\psi$ based on estimates of $\psi'$, and then dividing by $\mathbf{B}(\vec{a}_i|\vec{a}_*)$. The equations in the upper left hand corner of each panel represent the regression fit and correlation coefficient.

Figure S2: Comparisons of predicted gene expression from the ROC SEMPER model (Gilchrist et al. (2015)) and direct measurements of expression from RNA-Seq or Microarray data for individual yeast taxa across the 100 selected genes from Salichos and Rokas (2013). The equations in the upper left hand corner of each panel represent the regression fit and correlation coefficient.

## Simulations

Overall, the simulation results indicate that SelAC model can reasonably recover the known values of the generating model (Figure S3 - S6). This includes not only the

parameters in the model, but also the optimal amino acids for a given sequence as well as the estimates of the branch lengths. There are a few observations to note. First, the ability to accurately recover the true optimal amino acid sequence will largely depend on the magnitude of $\phi$. This is, of course, intuitive, given that $\phi$ sets the strength of stabilizing selection towards an optimal amino acid at a site. However, the inclusion of $\alpha_G$ into the model, appears to generally increase values of $\phi$ and generally improves the ability to recover the optimal amino acids even for the gene with the lowest baseline $\phi$. Second, we found a strong downward bias in estimates of $\alpha_G$, which actually translates to greater variation among the rate categories. The choice of a gamma distribution to represent site-specific variation in sensitivity was based on mathematical convenience and convention, rather than on biological reality. Nevertheless, we suspect that this bias is in large part due to the difficulty in determining the baseline $\psi$ for a given gene and the value of $\alpha_G$ that globally satisfies the site-specific variation in sensitivity across all genes, as indicated by the slight upward bias in estimates of $\psi$. A reviewer pointed out that it may also be difficulty for the model to account for changing amino-acid, which we agree may also play a role. It has been suggested, in studies of the behavior of the gamma distribution in applications of nucleotide substitution model, that increasing the number of rate categories can often improve accuracy of the shape parameter (Mayrose et al. (2005)). Future work will address this issue.
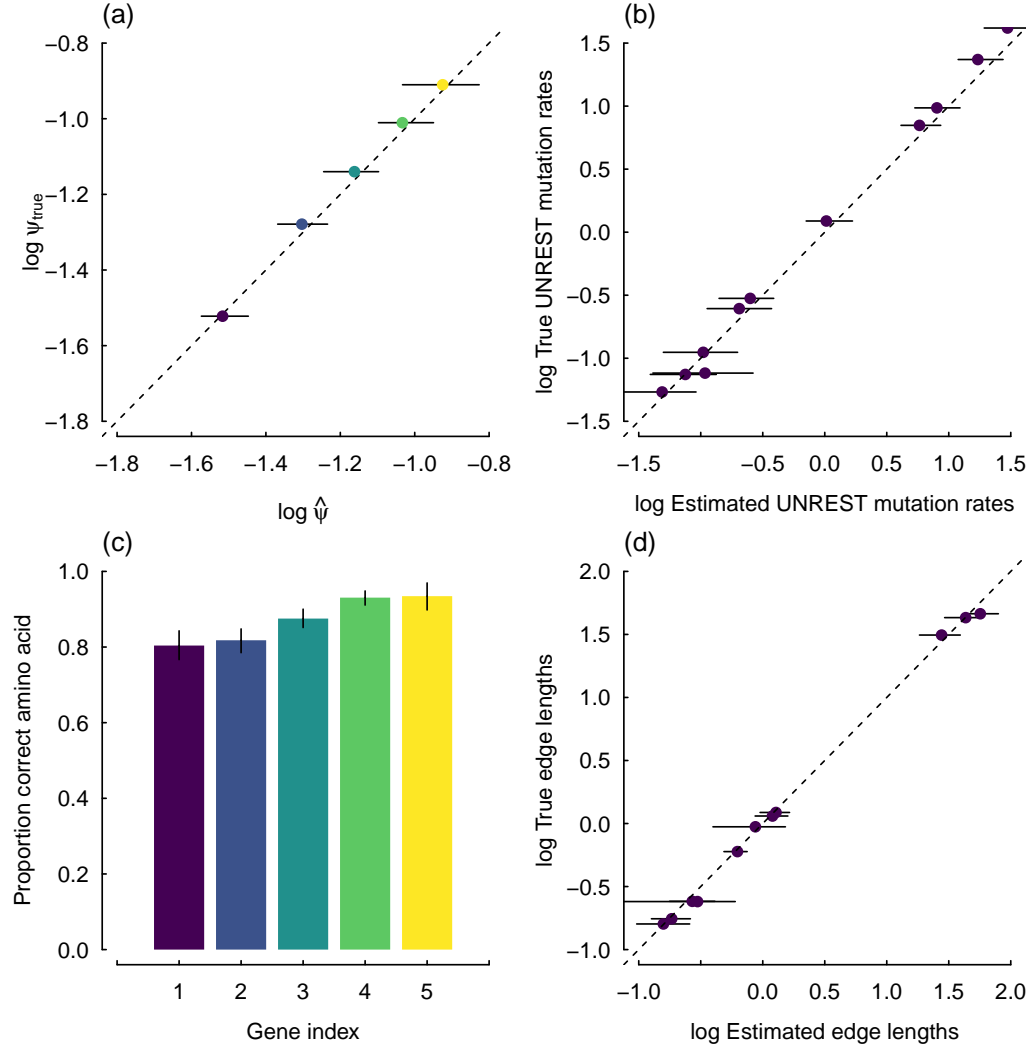
Figure S3: Summary a 5-gene simulation for a SelAC model where we assume $\alpha_G = \infty$, and thus, no site-specific sensitivity in the generating model. The 'known' parameters were based on fitting the same SelAC to the 106 gene data set and phylogeny of Rokas et al. (2003), with gene choice being based on five evenly spaced points along the rank order of the gene specific composite parameter $\psi'_g$. The points and associated uncertainty in the estimates of the gene-specific average protein synthesis rate, or $\psi$ (calculated from $\psi'$)(a), nucleotide mutation rates under the UNREST model (b), proportion of correct optimal amino acids for a given gene (c), and estimates of the individual edge lengths are based the mean and 2.5% and 97.5% quantiles across on 50 simulated datasets (d). Gene index on the x-axis refers to the arbitrary number assigned to the simulated gene.
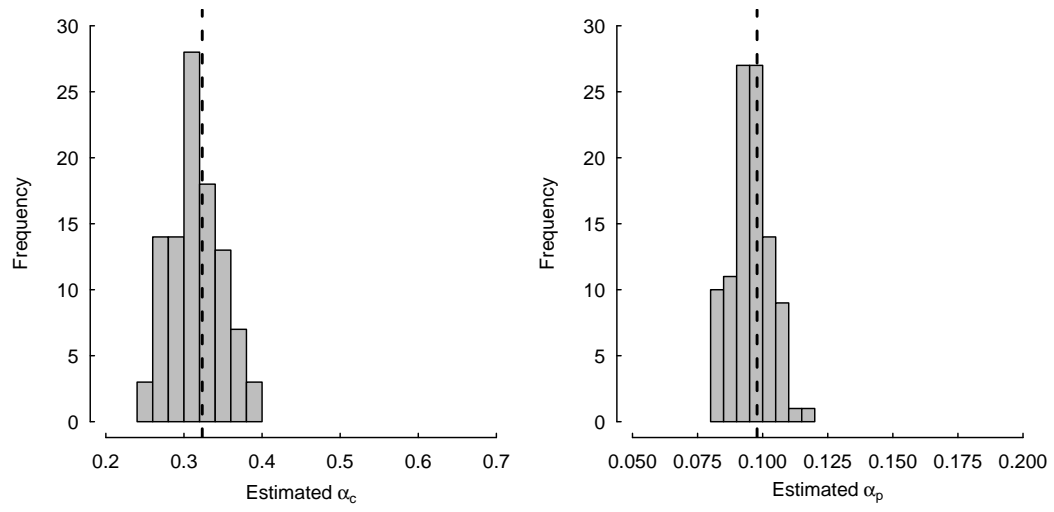
Figure S4: The distribution of estimates of the Grantham weights, $\alpha_c$ and $\alpha_p$, in a SelAC model, where we assume $\alpha_G = \infty$, and thus no site-specific sensitivity in the generating model. The dashed line represents the value used in the generating model.
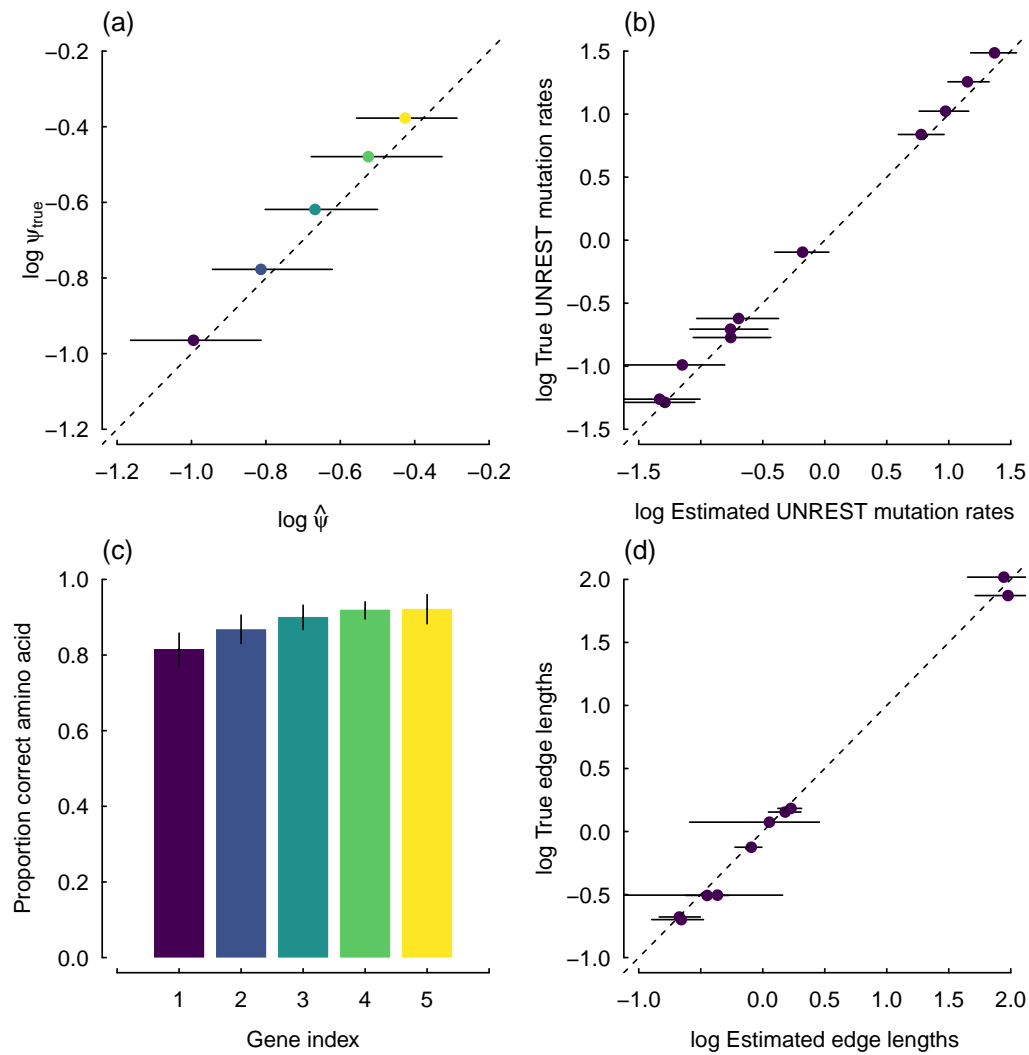
Figure S5: Same figure as in Figure S3, except the generating model includes site-specific sensitivity in the generating model (i.e., $\alpha_G$).
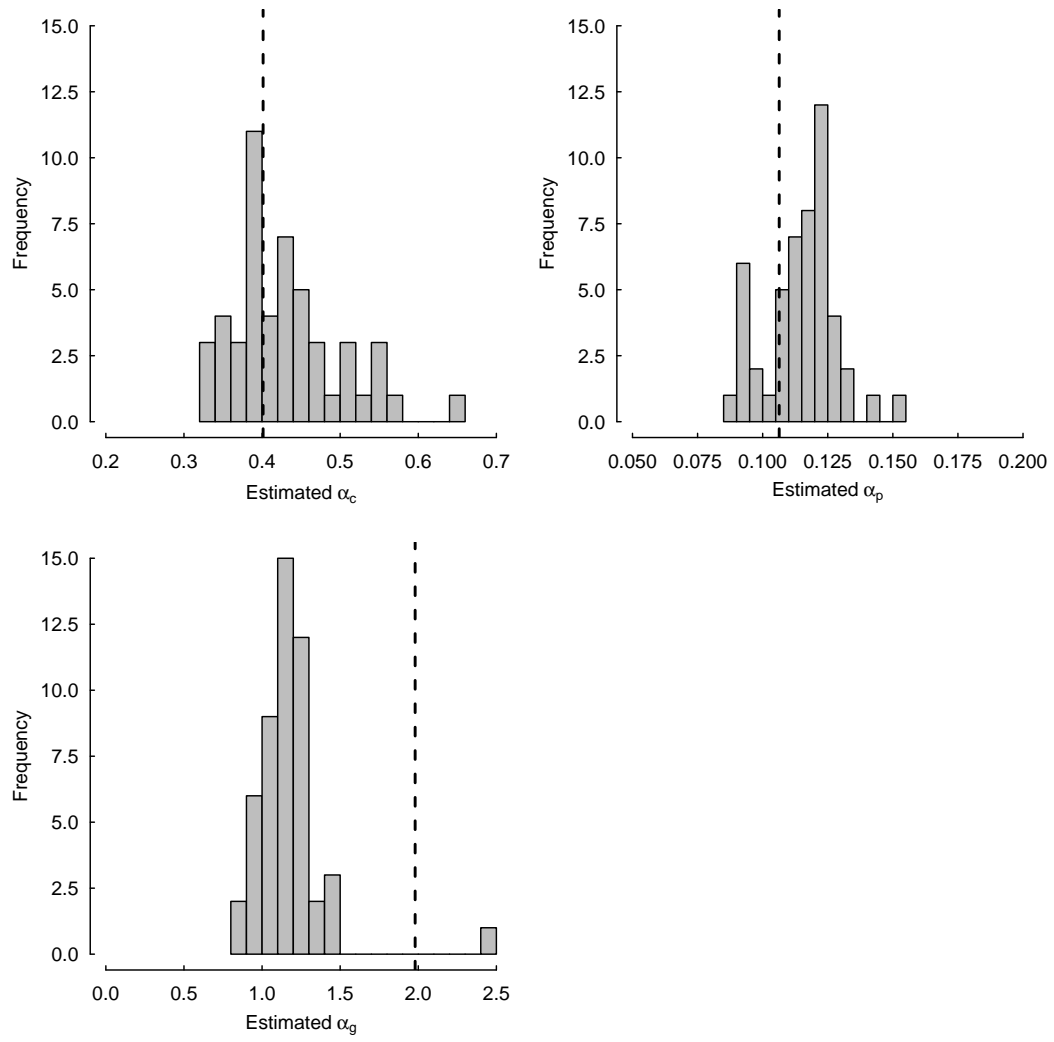
Figure S6: Same figure as in Figure S4, except the generating model includes site-specific sensitivity in the generating model (i.e., $\alpha_G$). Unlike, Grantham weights, which showed no systematic bias, there is a downward bias in estimates of $\alpha_G$.