# 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 Landerer[2,3], Juanjuan Chai[3,5], 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.

**Associate Editor:** TBD

## Abstract

We present a new phylogenetic approach SelAC (Selection on Amino acids and Codons), whose substitution rates are based on a nested model linking protein expression to population genetics. Unlike simpler codon models which assume a single substitution matrix for all sites, our model more realistically represents the evolution of protein coding DNA under the assumption of consistent, stabilizing selection using cost-benefit approach. This cost-benefit approach allows us generate a set of 20 optimal amino acid specific matrix families using just a handful of parameters and naturally links the strength of stabilizing selection to protein synthesis levels, which we can estimate. Using a yeast dataset of 100 orthologs for 6 taxa, we find SelAC fits the data much better than popular models by $10^4 - 10^5$ AICc units. Our results also indicated that nested, mechanistic models better predict observed data patterns highlighting the improvement in biological realism in amino acid sequence evolution that our model provides. Additional parameters estimated by SelAC indicate that a large amount of non-phylogenetic, but biologically meaningful, information can be inferred from exisiting data. For example, SelAC prediction of gene specific protein synthesis rates correlates well with both empirical ($r = 0.33 - 0.48$) and other theoretical predictions ($r = 0.45 - 0.64$) for multiple yeast species. SelAC also provides estimates of the optimal amino acid at each site. Finally, because SelAC is a nested approach based on clearly stated biological assumptions, future modifications, such as including shifts in the optimal amino acid sequence within or across lineages, are possible.

Key words: Wright-Fisher, stabilizing selection, allele substitution, protein function, gene expression

Article

# Introduction

Phylogenetic analyses plays a critical role in most aspects of biology, particularly in the fields of ecology, evolution, paleontology, medicine, and conservation. While the scale and impact of phylogenetic studies have increased substantially over the past two decades, the realism of the mathematical models on which these analyses are based has changed relatively little by comparison. The most popular models of DNA substitution used in molecular phylogenetics are simple nucleotide models that date back to the early 1980's and 90's, e.g. F81, F84, HYK85, TN93, and GTR (see **?** for an overview), and are indifferent to the type of sequences they are fitted to. For example, when evaluating protein-coding sequences these models are inherently agnostic with regards to the different amino acid substitutions and their impact on gene function and, as a result, cannot describe the behavior of natural selection at the amino acid or protein level.

Two important and independent attempts to address this critical shortcoming were introduced by **?**, commonly abbreviated as GY94 and **?**. These models were explicitly built for protein coding data, assuming that differences in the physicochemical properties between amino acids, or physicochemical distances for short, could affect substitution rates. These physicochemical based codon models as originally introduced have rarely been used for empirical data. Instead, these often cited models have served as the basis for an array of simpler and, in turn, more popular $\omega$ models that, starting with **??**, typically assume an equal fixation probability for *all* non-synonymous mutations. Although often attributed to GY94, these later and simpler models were the first to employ the single term $\omega$ to model the differences in fixation probability between nonsynonomous and synonomyous changes at all sites. Since their introduction, more complex models have been developed that allow $\omega$ to vary between sites or branches (as cited in **?**) and include selection on different synonyms for the same amino acid (e.g. **?**).

In **?**, **?** ,**?** and later studies based on their work, $\omega$ is suggested to indicate whether a given site within a protein sequence is under consistent 'stabilizing ($\omega < 1$) or 'diversifying' ($\omega > 1$) selection. Contrary to popular belief, $\omega$ does not describe whether a site is evolving under a constant regime of stabilizing or diversifying selection, but instead how a very particular *selective environment* changes over time. Below we explain how the actual behavior of these models is inconsistent with how 'stabilizing' and 'diversifying' selection are otherwise defined and understood (e.g. see **?**).

For example, when $\omega < 1$, synonymous substitutions have a higher substitution rate than any possible non-synonymous substitutions. As a result, the model behaves as if the resident amino acid $i$ at a given site is favored by natural selection. Even when $\omega$ is allowed to vary between sites, symmetrical aspects

of the model means that for any given site the strength of selection for the resident amino acid $i$ over its 19 alternatives is equally strong regardless of their physicochemical properties. Paradoxically, natural selection for amino acid $i$ persists *until* a substitution for another amino acid, $j$, occurs. As soon as amino acid $j$ fixes, but not before, selection now favors amino acid $j$ equally over all other amino acids, including amino acid $i$. This is now the opposite scenario from when $i$ was the resident. Thus, the simplest and most consistent interpretation of $\omega$ is that it represents the rate at which the selective environment itself changes, and this change in selection perfectly coincides with the fixation of a new amino acid.

Similarly, when $\omega > 1$, synonymous substitutions have a lower substitution rate than any possible non-synonymous substitutions from the resident amino acid. Again due to the model's symmetrical nature, the selection *against* the resident amino acid $i$ is equally strong relative to alternative amino acids. The selection against the resident amino acid $i$ persists until a substitution occurs at which point selection now *favors* amino acid $i$, as well as the 19 other amino acids, to the same degree $i$ was previously disfavored. Of course, in practice it is unlikely for the non-synonymous rate to be greater than the synonymous substitution rate in the absence of a particular process (e.g., antagonistic coevolution). Given these behaviors, $\omega$ based models are likely to only reasonably approximate a subset of scenarios such as perfectly symmetrical over-/under-dominance or positive/negative frequency dependent selection (**??**). Further, $\omega$ based models implicitly assumes the substitution is on the same timescale as the shifts in the optimal (or pessimal) amino acid. It is for these reasons that $\omega$ is best interpreted in a gene-wide context, as opposed to site-specific, because it averages over substitutions across many different sites.

## New Approaches

To address these fundamental shortcomings in $\omega$ based phylogenetic approaches, we present an approach where selection explicitly favors minimizing the cost-benefit function $\eta$ of a protein whose relative performance is determined by the order and physicochemical properties of its amino acids. Our approach, which we call Selection on Amino acids and Codons or SelAC, is developed in the same vein as previous phylogenetic applications of the Wright-Fisher process (e.g. **??????????**). Similar to **?** and **?**, 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 that 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. The key parameters underlying these matrices are shared across genes except for gene expression. As a result, SelAC identifies the amino acid at a particular position within a protein that is favored by natural selection using a simple cost-benefit approach.

While natural selection on protein coding regions can take many forms, one general approach to describing its effects is by relating a codon sequence to the cost of producing the encoded protein and the functional benefit (or potential harm) from translating its sequence. The gene specific cost of protein synthesis can be affected by the amino acids used, the direct and indirect costs of peptide assembly by the ribosome, and the use of chaperones to aid in folding. Importantly, these costs can be computed to varying degrees of realism (e.g. **??**). We have previously presented models of protein synthesis costs that, alternatively, take into account the cost of ribosome pausing (**?**) or premature termination errors (**???**).

Protein function or 'benefit' can be affected by the amino acids at each site and their interactions. Linking amino acid sequence to protein function is a daunting task; thus for simplicity, we assume that for any given desired biological function to be carried out by a protein, that (a) the biological importance of this protein function is invariant across the tree, (b) single optimal amino acid sequence that carries out this function best, and (c) the functionality of alternative amino acid sequences declines with their physicochemical distance from the optimum on a site by site basis.

Beyond fitting the phylogenetic data better than currently available nucleotide and codon models according to model adequacy and AICc, SelAC also makes inferences about other important biological processes. By comparing these inferences to other empirical data, such as we do with protein synthesis data, we can evaluate SelAC's performance independent of the data it is fitted to. Indeed, SelAC's assumptions lead to mechanistic and, thus, testable hypothesis about the nature of and relationships between mutation, protein function, gene expression, and rates of evolution. More importantly, alternative hypotheses could be used in place of ours and, in turn, phylogenetic and other types of data could be used to evaluate the support of these alternative models. Our hope is that by moving away from the more phenomenological models we can better connect population genetics, molecular biology, and phylogenetics allowing each area inform the others more effectively.

## Results

The SelAC model requires the construction of gene and amino acid specific substitution matrices that uses a submodel nested within our substitution model. This requires only a handful of genome-wide parameters such as nucleotide specific mutation rates $\mu_{i,j}$ , which are scaled by effective population size Ne, amino acid side chain physicochemical weighting parameters $\alpha_c$, $\alpha_p$, and $\alpha_v$, and a gamma distribution shape parameter $\alpha_G$ describing the distribution of site sensitivities $G$. In addition to these genome-wide parameters, the model requires a gene-specific functionality expression parameter that describes the average rate at which the protein's functionality is produced by the organism or a gene's

4

<sup>125</sup> 'average functionality production rate' $\psi$. By linking transition rates $q_{i,j}$ to gene expression in the form

<sup>126</sup> of protein synthesis rate $\phi$, our approach allows use of the same model for genes under varying degrees of

<sup>127</sup> stabilizing selection. Specifically, we assume the strength of stabilizing selection for an optimal sequence,

<sup>128</sup> $\vec{a}^*$, is proportional to $\psi$, which we can estimate for each gene.

<sup>129</sup> We first evaluated the performance of our codon model by simulating datasets and estimating the

<sup>130</sup> bias of the inferred model parameters from these data. Overall, the simulation results indicated that

<sup>131</sup> our SelAC model can reasonably recover the known values of the generating model (Figure 1S3). This

<sup>132</sup> includes not only the parameters in SelAC, but also the optimal amino acids for a given sequence as well

<sup>133</sup> as the estimates of the branch lengths. There are, however, a few observations to note. First, the ability to

<sup>134</sup> accurately recover the true optimal amino acid sequence $\vec{a}^*$ will largely depend on the magnitude of the

<sup>135</sup> realized average protein synthesis rate of the gene $\phi$, which is the target functionality rate $\psi$ divided by

<sup>136</sup> the functionality of the observed amino acid sequence $\mathbf{B}(\vec{a})$. This is, of course, intuitive, given that $\psi$ sets

<sup>137</sup> the strength of stabilizing selection towards an optimal amino acid at a site. However, the inclusion of

<sup>138</sup> between site variation in selection via the shape parameter $\alpha_G$ into SelAC generally increase our estimates

<sup>139</sup> of $\psi$ as well as improve our ability to recover the optimal amino acids $\vec{a}^*$. This is true even for the gene

<sup>140</sup> with the lowest baseline $\psi$. Second, we found a strong downward bias in estimates of $\alpha_G$, which actually

<sup>141</sup> translates to greater variation among the rate categories. The choice of a gamma distribution to represent

<sup>142</sup> site-specific variation in sensitivity was based on mathematical convenience and convention, rather than

<sup>143</sup> on biological reality. Further, given the fact that the density of the gamma distribution is infinite at

<sup>144</sup> $G\!=\!0$ when $\alpha_G\!<\!1$, imputing site specific $G$ values will be an issue in these scenarios. Nevertheless, we

<sup>145</sup> suspect that this downward estimation bias of $\alpha_G$ is in large part due to the difficulty in determining

<sup>146</sup> the baseline $\psi$ for a given gene and the value of $\alpha_G$ that globally satisfies the site-specific variation in

<sup>147</sup> sensitivity across all genes, as indicated by the slight upward bias in estimates of $\psi$ (see Figure S5).

In regards to model fit in an empirical setting, 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 shape parameter $\alpha_G$ for the random effects term $G \sim$

Gamma(shape$=\alpha_G$,rate$=\alpha_G$) to allow for heterogeneity in this selection between sites within a gene

improves the $\Delta$AICc of SelAC+$\Gamma$ over the simpler SelAC models by over 22,000 AIC units. Using either

$\Delta$AICc or AIC$_{\mathrm{w}}$ as our measure of model support, the SelAC models fit extraordinarily better than

GTR + $\Gamma$, GY94, or FMutSel (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, FMutSel (with 178 parameters),

5

SelAC+$\Gamma$ model shows over 160,000 AIC unit improvement over FMutSel. SelAC models also appeared to outperform, based on likelihood, and reported AIC and AICc from each program, the codon models in IQtree. Supplemental table $_{comparesthese.}$

We note our use of AICc, as opposed to the standard AIC, in the above model comparisons. At the outset of our study it was unclear what the appropriate sample size, $n$, when comparing models of sequence evolution. Building upon the work of **?**, our simulations suggest that using the number of taxa times the number of sites as the sample size correction performed best as a small sample size correction for estimating Kullback-Liebler distance in phylogenetic models (Supporting Materials). This also has an intuitive appeal. In models that have at least some parameters shared across sites and some parameters shared across taxa, increasing the number of sites and/or taxa should be adding more samples for the parameters to estimate. This is consistent considering how likelihood is calculated for phylogenetic models: the likelihood for a given site is the sum of the probabilities of each observed state at each tip, which is then multiplied across sites. It is arguable that the conventional approach in comparative methods is calculating AICc in the same way. That is, if only one column of data (or "site") is examined, as remains remarkably common in comparative methods, when we refer to sample size, it is technically the number of taxa multiplied by number of sites, even though it is referred to simply as the number of taxa.

With respect to estimates of $\phi$ within SelAC, they were strongly correlated with both empirical measurements (Pearson $r=0.33-0.48$) and theoretical predictions (Pearson $r=0.45-0.64$) of gene expression (Figure 2 and Figures S1-S2, respectively). These correlations are remarkable given that they were uncovered using only codon sequences. The estimate of the $\alpha_G$ parameter, which controls the shape of the site-specific, gamma-distributed, variation in sensitivity of the protein's functionality, indicated a moderate level of variation in gene expression among sites. Our estimate of $\alpha_G = 1.36$, produced a distribution of sensitivity terms $G$ ranged from 0.342-7.32, but with more than 90% of the weight for a given site-likelihood being contributed by the 0.342 and 1.50 rate categories. In simulation, however, of all the parameters in the model, only $\alpha_G$ showed a consistent bias, in that the MLE were generally lower than their actual values (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 **?** original estimates in regards to composition weight, $\alpha_c$, which is the ratio of non-carbon atoms in the end groups or rings to the number of carbon atoms in side chains. Our estimate of the composition weighting factor of $\alpha_c$=0.459 is 1/4th the value estimate by Grantham which suggests that the substitution process is less sensitive to

6

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 FMutSel 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 (Figure 3) 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 FMutSel, the functionality of the simulated sequence, defined as protein function in relation to the physicochemical distance from each amino acid to the optimal, moves away from the observed sequence. By contrast, SelAC remains near the functionality of the observed sequence (Figure 4b). This is somewhat unsurprising, given that both GTR $+$ $\Gamma$ and FMutSel are agnostic to the functionality of the gene, but it does highlight the improvement in biological realism in amino acid sequence evolution that SelAC provides. We do note that the adequacy of the SelAC model does vary among individual taxa, and does not always match the observed functionality. For instance, our simulations of *S. castellii* gene function is consistently higher than estimated from the data (Figure 4c). We suspect this is an indication that assuming a single set of optimal amino acid across all taxa is too simplistic. However, we cannot rule out violations of SelAC's other model assumptions such as, a single set of Grantham weights, a single $\alpha_G$, or reductions in protein functionality $\mathbf{B}$ being solely a function of physicochemical distances $d$ between sites.

## Discussion

A central goal in evolutionary biology is to quantify the nature, strength, and, ultimately, shifts in the forces of natural selection relative to genetic drift and mutation. As data set size and complexity increase, so does the amount of potential information on these forces and their dynamics. As a result, there is a need for more complex and realistic models to accomplish this goal (**?????**). Although extremely popular due to their elegance and computational efficiency, the utility of $\omega$ based models in helping us reach this goal is substantially more limited than commonly recognized. Because these $\omega$ models use a single substitution matrix, they are only applicable for situations in which the substitution process and shifts in the selective

7

environment are intrinsic to the sequence, such as with positive or negative frequency dependent selection; these models do not describe stabilizing or diversifying selection as commonly envisioned (**??**).

Starting with **?**, a number of researchers have developed methods for linking site-specific selection on protein sequence and phylogenetics (e.g. **???????**). **?** calculated a vector of 20 expected amino acid frequencies for each amino acid site, making it the most general and most parameter rich of these methods. This generality, however, comes at the cost of being purely descriptive; there is no explicit biological mechanism proposed to explain the site specific amino acid frequencies estimated. By grouping together amino sites with similar evolutionary behaviors, **?** and **?** retained the descriptive nature of **?** work while greatly reduced the number of model parameters needed.

SelAC follows in this tradition of using multiple substitution matrices, but includes some key advances. First, by nesting a model of a sequence's cost-benefit function $\mathbf{C}/\mathbf{B}$ within a broader model, SelAC allows us to formulate and test a hierarchical, mechanistic models of stabilizing selection. More precisely, our nested approach allows us to relax the assumption that physicochemical deviations from the optimal sequence $\vec{a}^*$ are equally disruptive at all sites within a protein. Indeed, SelAC results are consistent with the idea that the strength of stabilizing selection against physicochemical deviations from $\vec{a}^*$ varies between sites ($\Delta$AICc = 20,983; Table1). Second, because our substitution matrices are built on a formal description of a sequence's cost-benefit function $\mathbf{C}/\mathbf{B}$, we are able to efficiently parameterize 20 different matrices using a relatively small number of genome-wide parameters – e.g. our physicochemical weightings, $\alpha_c$, $\alpha_p$, and $\alpha_v$, and the shape parameter $\alpha_G$ for the distribution of selective strength $G$ and one gene specific expression parameter $\psi$. While the $\mathbf{C}/\mathbf{B}$ function on which SelAC currently rests is very simple, nevertheless, it leads to a dramatic increase in our ability to explain the sequence data we analyzed. Importantly, because SelAC uses a formal description of a sequence's $\mathbf{C}/\mathbf{B}$, replacing our assumptions with more sophisticated ones in the future is relatively straightforward. Third, our use of nested models also allows us to make biologically meaningful and testable predictions. By linking a gene's expression level to the strength of purifying selection it experiences, we are able to provide coarse estimates of gene expression. This also suggests that the anticorrelation between $\omega$ and gene expression is a proxy for the nature of selection on a sequence.

Thus, we believe our cost-benefit approach to be a substantial advance of the more simplistic $\omega$ models, is complementary to the work of others in the field (e.g. **??**), and, in turn, lays the foundation for more realistic work in the future. For instance, by assuming there is an optimal amino acid for each site, SelAC naturally leads to a non-symmetrical and, thus, more cogent model of protein sequence evolution. Because the strength of selection depends on an additive function of amino acid physicochemical properties, an

amino acid more similar to the optimum has a higher probability of replacing a more dissimilar amino acid than the converse situation. Further, SelAC does not assume the system is always at the optimum or pessimum point of the fitness landscape, as occurs when $\omega < 1$ or $> 1$, respectively.

Importantly, the cost-benefit approach underlying SelAC allows us to link the strength of selection on a protein sequence to its gene's expression level. Despite its well recognized importance in determining the rate of protein evolution (e.g. **??**), phylogenetic models have ignored the fact that expression levels vary between genes. In order to link gene expression and the strength of stabilizing selection on protein sequences, we simply assume that the strength of selection on a gene is proportional to the average protein synthesis rate of the gene.

One possible mechanism with some theoretical and empirical support which generates a linear relationship between the strength of selection and gene expression is the assumption of compensatory gene expression (**???????**). That is, the assumption that any reduction in protein function is compensated for by an increase in the protein's production rate and, in turn, abundance. For example, a mutation which reduces the functionality of the protein to 90% of the optimal protein, would require $1/0.9 = 1.11$ of these suboptimal proteins to be produced relative to the optimal protein in order to maintain the same amount of that protein's functionality in the cell. Because the energetic cost of an 11% increase in a protein's synthesis rate is proportional to its target synthesis rate, our assumptions naturally link changes in protein functionality and changes in gene expression and its associated costs. Under what circumstances cells actually respond in this manner, remains to be determined. The fact that our method allows us to explain 13-23% of the variation in gene expression measured using RNA-Seq, suggests that this assumption is a reasonable starting point.

Furthermore, by linking expression and selection, SelAC provides a natural framework for combining information from protein coding genes with very different rates of evolution; from low expression genes providing information on shallow branches to high expression genes providing information on deep branches. This is in contrast to a more traditional approach of concatenating gene sequences together, which is equivalent to assuming the same average functionality production 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 functionaly production rate $\psi$ and the realized average protein synthesis rate $\phi = \psi/\mathbf{B}$ are reasonably well correlated with laboratory measurements and theoretical predictions of gene expression (Pearson $r = 0.34 - 0.64$; Figures 2, 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 to 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 physicochemical weights might improve our model fit and reduce the noise in our estimates of realized protein synthesis rates $\phi$. Even if other physicochemical properties are considered, the idea of a consistent, genome wide physicochemical weighting of these terms seems highly unlikely. Since the importance of an amino acid's physicochemical properties likely changes with its position in a folded protein, one way to incorporate such effects is to test whether the data supports multiple sets of physicochemical 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 sites 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 these more complex scenarios.

For example, our current implementation ignores any selection on synonymous codon usage bias (CUB) (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 function 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. **?**). Nevertheless, it is difficult to see how one could identify such conditions without modeling the way in which codon and amino acid usage affects $\mathbf{C}/\mathbf{B}$.

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. However, how one determines the

10

appropriate number of data points in a model is more complicated than generally recognized. 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 × number of taxa) should be taken as the sample size for AICc, but this conclusion likely only applies when there is sufficient independence between taxa. 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 what happens on the path connecting the lone taxon to the polytomy. Although this is a rather extreme example, it seems prudent for researchers to use a simulation based approach similar to the one we take here to determine the appropriate means for calculating the effective number of data points in their data.

There are still significant shortcomings in the approach outlined here. Most worrisome are biological oversimplifications in SelAC. For example, at its heart, SelAC 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. Like the other oversimplifications previously discussed, these assumptions can be relaxed through further extension of our model.

A deeper potential issue comes from the nature of model fitting itself. Our model assumes the world works a particular way (selection for an optimal amino acid rather than on properties themselves, independence between sites, the cell's response to inefficient proteins, and more). Some of these assumptions, and the parameters based on them, are undoubtedly wrong. However, as we see in other models (Goldberg and Rabosky, Hisse),

There are also deficiencies in our implementation. Though reasonable to use for a given topology with a modest number of species, it is currently too slow for practical use for tree search. Our work 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 challenging optimization problem: the nested models reduce parameter complexity vastly, but there are still numerous parameters to optimize, including the discrete parameter of the optimal amino acid at each site. One way to avoid the use of

discrete parameters at the expense of more of them would be to have SelAC estimate the optimum physicochemical values on a per site basis rather than a specific amino acid. While this would increase the number of parameters estimated, it 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).

In spite of these difficulties, SelAC represents an important step in uniting phylogenetic and population genetic models. For example, while **???????** are all models of constant, stabilizing selection, SelAC can be generalized further to include diversifying selection. Specifically, by letting SelAC's sensitivity term $G$, which we now assume is $\geq 0$, to take on negative values, SelAC will behave as if there is a pessimal, rather than optimal, amino acid for the given site. In this diversifying selection scenario, amino acids with physicochemical qualities more dissimilar to the pessimal amino acid are increasingly favored, potentially resulting in multiple fitness peaks.

The ability to extend our model and, in turn, sharpen our thinking about the nature of natural selection on amino acid sequences illustrates the value of moving from descriptive to more mechanistic models in general and phylogenetics in particular. How frequently diversifying selection of this nature occurs is an open, but addressable, question. Regardless of the frequency at which diversifying selection occurs, another question of interest to evolutionary biologists is, "How often does the optimal/pessimal amino sequence change along any given branch?" Due to its mechanistic nature, SelAC can also be extended to include changes in the optimal/pessimal sequence over a phylogeny using a hidden Markov modelling approach **???**. Extending SelAC in these ways, will allow researchers to explicitly model shifts in selection on protein sequences and, in turn, quantify their frequency and magnitude thus deepening our understanding of biological evolution.

In summary, SelAC allows biologically relevant population genetic parameters to be estimated from phylogenetic information, while also dramatically improving fit and accuracy of phylogenetic models. By explicitly modeling the optimal/pessimal sequence of a gene, SelAC can be extended to include shifts in the optimal/pessimal sequence over evolutionary time. Moreover, it demonstrates that there remains substantially more information in the coding sequences used for phylogenetic analysis than other methods can access. Given the enormous amount of efforts expended to generate sequence datasets, it makes sense for researchers to continue developing more realistic models of sequence evolution in order to extract the biological information embedded in these datasets. The cost-benefit model we develop here is just one of many possible paths of mechanistic model development.

12

# Materials & Methods

## Overview

We model the substitution process as a classic Wright-Fisher process which includes the forces of mutation, selection, and drift (???????). For simplicity, we ignore linkage effects and, as a result of this and other assumptions, sequences evolve in a site independent manner.

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 (i.e. on the order of 74,420 parameters). 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 ???.

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 $N_e$), amino acid 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 functionality expression parameter $\psi_g$ which describes the average rate at which the protein's functionality is produced by the organism or a gene's 'average functionality production rate' for short (for notational simplicity, we will ignore the gene specific indicator $_g$, unless explicitly needed). Currently, $\psi$ is fixed across the phylogeny, though relaxing this assumption is a goal of future work. The gene specific parameter $\psi$ is multiplied by additional model terms to make a composite term $\psi'$ which scales the strength and efficacy of selection for the optimal amino acid sequence relative to drift (see Implementation below). 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, a complete, error free peptide consisting of $\vec{a}^*$ provides one unit of the gene's functionality. We also assume that natural selection favors genotypes that are able to synthesize their proteome more 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 $A_0$. SelAC also requires the specification (as part of parameter optimization) of an optimal amino acid $a^*$ at each position within a coding sequence. This requirement of one $a^*$ per site makes our $\vec{a}^*$ the largest category of parameters SelAC estimates. Despite the need to specify $a^*$ for each site, because we use a submodel to derive our substitution matrices, SelAC estimates a relatively small number of the parameters when compared to more general approaches where the fitness of each amino acid is allowed to vary freely of any physicochemical properties (???).

13

As with other phylogenetic methods, SelAC generates estimates of branch lengths and nucleotide specific mutation rates. In addition, the method can also be used to make quantitative inferences on the optimal amino acid sequence of a given protein as well as the realized average synthesis rate of each protein used in the analysis. The mechanistic basis of SelAC also means it can be easily extended to include more biological realism and test more explicit hypotheses about sequence evolution.

**Mutation Rate Matrix $\boldsymbol{\mu}$**

We begin with a 4x4 nucleotide mutation matrix $\boldsymbol{\mu}$ that describes mutation rates between different bases and, in turn, different codons. For our purposes, we rely on the general unrestricted model (UNREST from **?**) because it imposes no constraints on the instantaneous rate of change between any pair of nucleotides. More constrained models, such as the Jukes-Cantor (JC), Hasegawa-Kishino-Yano (HKY), or the general time-reversible model (GTR), could also be used.

The 12 parameter UNREST model defines the relative rates of change between a pair of nucleotides. Thus, we arbitrarily set the G$\rightarrow$T mutation rate to 1, resulting in 11 free mutation rate parameters in the 4x4 mutation nucleotide mutation matrix. The nucleotide mutation matrix is also scaled by a diagonal matrix $\boldsymbol{\pi}$ whose entries, $\pi_{i,i}$, correspond to the equilibrium frequencies of each base. These equilibrium nucleotide frequencies are determined by analytically solving $\boldsymbol{\pi} \times \mathbf{Q} = 0$. We use this $\mathbf{Q}$ to populate a $61 \times 61$ codon mutation matrix $\boldsymbol{\mu}$, whose entries $\mu_{i,j}$ $i \neq j$ describes the mutation rate from codon $i$ to $j$ and $\mu_{i,i} = -\sum_j \mu_{i,j}$. We generate this matrix using a "weak mutation" assumption, such that evolution is mutation limited, codon substitutions only occur one nucleotide at a time. As a result, the rate of change between any pair of codons that differ by more than one nucleotide is zero.

While the overall model does not assume equilibrium, we still need to scale our mutation matrices $\mu$ by a scaling factor $S$. As traditionally done, we rescale our time units such that at equilibrium, one unit of branch length represents one expected mutation per site (which equals the substitution rate under neutrality). More explicitly, $S = -\left(\sum_{i \in \text{codons}} \mu_{i,i} \pi_{i,i}\right)$ where the final mutation rate matrix is the original mutation rate matrix multiplied by $1/S$.

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

SelAC links fitness to the product of the cost-benefit function of a gene $\eta$ and the organism's average target synthesis rate of the functionality provided by gene $\psi$. As a result, the average flux energy an organism spends to meet its target functionality provided by the gene is $\eta \times \psi$. Compensatory changes that allow an organism to maintain functionality even with loss of one or both copies of a gene are widespread. There is evidence of compensation for protein function. Metabolism with gene expression models (ME-models) link those factors to successfully make predictions about response to perturbations

14

in a cell (**??**). For example, an ME-model for *E. coli* successfully predicted gene expression levels in vivo (**?**). Here we assume that for finer scale problems than entire loss (for example, a 10% loss of functionality) the compensation is more production of the protein. The particular type of dosage compansation assumed by SelAC in respondse to stress (e.g. reduced functionality) is commonly assumed in microbial ecology (**??**). Our assumption is also consistent with the Michaelis-Menten enzyme kinetics. Moreover, there is evidence that mutations can influence expression level, though this does not always match our expression compensation assumption (**??**). In order to link genotype to our cost-benefit function $\eta = \mathbf{C}/\mathbf{B}$, we begin by defining our benefit function $\mathbf{B}$.

*Benefit:* Our benefit function $\mathbf{B}$ measures the functionality of the amino acid sequence $\vec{a}_i$ 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 $\vec{a}^*$. By definition, $\mathbf{B}(\vec{a}^*|\vec{a}^*) = 1$ and $\mathbf{B}(\vec{a}_i|\vec{a}^*) < 1$ for all other sequences. We assume all amino acids within the sequence contribute to protein function and that this contribution declines as an inverse function of physicochemical distance from each amino acid to the optimal one. Formally, we assume that

$$\mathbf{B}(\vec{a}|\vec{a}^*) = \left(\frac{1}{n}\sum_{p=1}^{n}\left(1 + G_p d(a_p, a_p^*)\right)\right)^{-1} \tag{1}$$

where $n$ is the length of the protein, $d(a_p, a_p^*)$ is a weighted physicochemical distance between the amino acid encoded at a given position $p$ and $a_p^*$ is the optimal amino acid for that position. There are many possible measures for physiochemical distance; we use **?** distances by default, though others may be chosen. For simplicity, we assume all nonsense mutations are lethal by defining the the physicochemical distance between a stop codon and a sense codon as $\infty$. The term $G_p$ describes the sensitivity of the protein's function to physicochemical deviation from the optimimum at site position $p$. We assume that $G_p \sim \text{Gamma}(\text{shape} = \alpha_G, \text{rate} = \alpha_G)$ in order to ensure $\mathbb{E}(G_p) = 1$. Given the definition of the Gamma distribution, the variance in $G_p$ is equal to shape/rate$^2 = 1/\alpha_G$. We note that at the limit of $\alpha_G \to \infty$, the model becomes equivalent to assuming uniform site sensitivity where $G_p = 1$ for all positions $p$. Further, $\mathbf{B}(\vec{a}_i|\vec{a}^*)$ 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 ATPs or GTPs 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

15

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 (**?**) 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,

$$\mathbf{C}(\vec{c}_i) = \text{Direct energetic cost of protein synthesis.}$$

$$= A_1 + A_2 n$$

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 one could compare the quality of SelAC fits using different mixtures of physicochemical properties. Following (**?**), 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. 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 estimates of $\alpha_c$ and $\alpha_p$ and $\psi$ are scaled relative to this choice for $\alpha_v$. More specifically,

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

**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

16

the organism needs that functionality to be produced at a target average rate $\psi$. Again, by definition, the optimal amino acid sequence for a given gene, $\vec{a}^*$, produces one unit of functionality, i.e. $\mathbf{B}(\vec{a}^*) = 1$. Second, we assume that the actual average rate a protein is synthesized $\phi$ is regulated by the organism to ensure that functionality is produced at rate $\psi$. As a result, it follows that $\phi = \psi/\mathbf{B}(\vec{a}|\vec{a}^*)$ and the energetic burden of a suboptimal amino acid increases the more it decreases the protein's functionality, $\mathbf{B}$. In other words, the average production rate of a protein $\vec{a}$ with relative functionality $\mathbf{B}(\vec{a}) < 1$ must be $1/\mathbf{B}(\vec{a}|\vec{a}^*)$ times higher than the production rate needed if the optimal amino acid sequence $\vec{a}^*$ was encoded since $\mathbf{B}(\vec{a}^*|\vec{a}^*) = 1$. For example, a cell with an allele $\vec{a}$ where $\mathbf{B}(\vec{a}|\vec{a}^*) = 9/10$ would have to produce the protein at rate $\phi = 10/9 \times \psi = 1.11\psi$. Similarly, a cell with an allele $\vec{a}$ where $\mathbf{B}(\vec{a}|\vec{a}^*) = 1/2$ will have to produce the protein at $\phi = 2\psi$. In contrast, a cell with the optimal allele $\vec{a}^*$ would have to produce the protein at rate $\phi = \psi$.

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

$$W_i(\vec{c}) \propto \exp[-A_0 \eta(\vec{c}_i)\psi].$$

where $A_0$, again, describes the proportional 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. Instead, we recognize that our estimates of $\psi$ share an unknown scaling term.

Correspondingly, the ratio of fitness between two genotypes is,

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

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

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 multiple sites simplifies to

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

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 Wright-Fisher process. As a result, the probability a new mutant,

510    $j$, introduced via mutation into a resident population $i$ with effective size $N_e$ will go to fixation is,

$$u_{i,j} = \frac{1 - (W_i/W_j)^b}{1 - (W_i/W_j)^{2N_e}}$$

$$= \frac{1 - \exp\left\{-\frac{A_0}{n_g}(A_1 + A_2 n_g)[d(a_i, a^*) - d(a_j, a^*)]G_p \psi b\right\}}{1 - \exp\left\{-\frac{A_0}{n_g}(A_1 + A_2 n_g)[d(a_i, a^*) - d(a_j, a^*)]G_p \psi 2N_e\right\}}$$

511    where $b = 1$ for a diploid population and 2 for a haploid population (**?????**). Finally, assuming a constant

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

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

513    where, given the substitution model's weak mutation assumption, $N_e \mu \ll 1$. In the end, each optimal

514    amino acid has a separate $61 \times 61$ substitution rate matrix $\mathbf{Q}_a$, which incorporates selection for the

515    amino acid (and the fixation rate matrix this creates) as well as the common mutation parameters across

516    optimal amino acids. This results in the creation of 20 $\mathbf{Q}$ matrices, one for each amino acid and each with

517    3,721 entries which are based on a relatively small number of model parameters (one to 11 mutation rates,

518    two free Grantham weights, the cost of protein assembly, $A_1$ and $A_2$, the gene specific target functionality

519    synthesis rate $\psi$, and optimal amino acid at each position $p$, $a_p^*$). These model parameters can either be

520    specified *a priori* and/or estimated from the data.

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

522    whole data set is the product of the probabilities of observing the data at each individual site. Thus, the

523    likelihood $\mathcal{L}$ of amino acid $a$ being optimal at a given site position $p$ is calculated as

$$\mathcal{L}(\mathbf{Q}_a | \mathbf{D}_p, \mathbf{T}) \propto \mathbf{P}(\mathbf{D}_p | \mathbf{Q}_a, \mathbf{T}) \tag{2}$$

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

525    with topology $\mathbf{T}$. For our purposes we take $\mathbf{T}$ as given, but it could be estimated as well. The pruning

526    algorithm of **?** is used to calculate $\mathcal{L}(\mathbf{Q}_a | \mathbf{D}_p, \mathbf{T})$. The log of the likelihood is maximized by estimating

527    the genome scale parameters which consist of 11 mutation parameters, which are implicitly scaled by

528    $2N_e/b$, and two Grantham distance parameters, $\alpha_c$ and $\alpha_p$, and the sensitivity distribution parameter

529    $\alpha_G$. Because $A_0$ and $\psi_g$ always co-occur and are scaled by $N_e$, for each gene $g$ we estimate a composite

530    term $\psi_g' = \psi_g A_0 b N_e$ and the optimal amino acid for each position $a_p^*$ of the protein. When estimating $\alpha_G$,

531    the likelihood then becomes the average likelihood which we calculate using the generalized Laguerre

532    quadrature with $k = 4$ points (**?**).

533    Finally, we note that because we infer the ancestral state of the system, our approach does not rely

534    on any assumptions of model stationarity. Nevertheless, as our branch lengths grow the probability

18

of observing a particular amino acid $a$ at a given site approaches a stationary value proportional to $W(a)^{2N_e - b}$ and any effects of mutation bias (**?**).

## Implementation

All methods described above are implemented in the new R package, `selac` available through GitHub (`https://github.com/bomeara/selac`) which will be uploaded to CRAN once peer review has completed. Our package requires as input a set of fasta files that each contain an alignment of 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 **?**, the FMutSel mutation-selection model of **?**, 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 utilizes a four stage hill climbing approach. More specifically, within each stage a block of parameters are optimized while the remaining parameters are held constant. The first stage optimizes the block of branch length parameters. The second stage optimizes the block of gene specific composite parameters $\psi'_g = A_0 \psi_g N_e b$. The third stage optimizes SelAC's parameters shared across the genome $\alpha_c$ and $\alpha_p$, and the sensitivity distribution parameter $\alpha_G$. The fourth stage estimates the optimal amino acid at each site $a^*$. This entire four stage cycle is repeated six more times, using the estimates from the previous cycle as the initial conditions for the new one. The search is terminated when the improvement in the log-likelihood between cycles is less than $10^{-8}$ at which point we consider the ML solution found and the search is terminated. For optimization of a given set of parameters, we rely on a bounded subplex routine (**?**) in the package `NLoptR` (**?**) to maximize the log-likelihood function. To ensure the robustness of our results, 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 $\alpha_p$ and were able to repeatedly reach the same log-likelihood (lnL) peak in our parameter space. 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 **?** and **?**.

We note that our current implementation of SelAC is painfully slow, and is best suited for data sets with relatively few number of taxa (i.e. $< 10$). This limitation is largely due to the size and quantity of matrices we create and manipulate to calculate the log-likelihood of an individual site. Ongoing work

19

<sup>567</sup> will address the need for speed, with the eventual goal of implementing SelAC in popular phylogenetic

<sup>568</sup> inference toolkits, such as RevBayes (**?**), PAML (**?**) and RAxML (**?**).

<sup>569</sup> **Simulations**

<sup>570</sup> We evaluated the performance of our codon model by simulating datasets and estimating the bias of

<sup>571</sup> the inferred model parameters from these data. Our 'known' parameters under a given generating model

<sup>572</sup> were based on fitting SelAC to the 106 gene data set and phylogeny of **?**. The tree used in these analyses

<sup>573</sup> is outdated with respect to the current hypothesis of relationships within *Saccharomyces*, but we rely on

<sup>574</sup> it simply as a training set that is separate from our empirical analyses (see section below). Bias in the

<sup>575</sup> model parameters were assessed under two generating models: one where we assumed a model of SelAC

<sup>576</sup> assuming uniform sensitivity across sites (i.e. $G_p = 1$ for all sites, i.e. $\alpha_G = \infty$), and one where we used

<sup>577</sup> the Gamma distribution joint shape and rate parameter $\alpha_G$ estimated from the empirical data. Under

<sup>578</sup> each of these two scenarios, we used parameter estimates from the corresponding empirical analysis and

<sup>579</sup> simulated 50 five-gene data sets. For the gene specific composite parameter $\psi'_g$ the 'known' values used for

<sup>580</sup> the simulation were five evenly spaced points along the rank order of the estimates across the 106 genes.

<sup>581</sup> The MLE estimate for a given replicate were taken as the fit with the highest log-likelihood after running

<sup>582</sup> five independent analyses with different sets of naive starting points with respect to the composite $\psi'_g$

<sup>583</sup> parameter, $\alpha_c$, and $\alpha_p$. All analyses were carried out in our `selac` R package.

<sup>584</sup> **Analysis of yeast genomes & tests of model adequacy**

<sup>585</sup> We focus our empirical analyses on the large yeast data set and phylogeny of **?**. As a model system,

<sup>586</sup> the yeast genome is an ideal system to examine our phylogenetic estimates of gene expression and its

<sup>587</sup> connection to real world measurements of these data within individual taxa. The complete data set of **?**

<sup>588</sup> contain 1070 orthologs, where we selected 100 at random for our analyses. We also focus our analyses on

<sup>589</sup> *Saccharomyces sensu stricto* and their sister taxon *Candida glabrata*, and we used the phylogeny depicted

<sup>590</sup> in Fig. 1 of **?** for our fixed tree. We fit the two SelAC models described above (i.e., SelAC and SelAC+Γ),

<sup>591</sup> as well as two codon models, GY94 and FMutSel, and a standard GTR + Γ nucleotide model. The

<sup>592</sup> FMutSel model assumes that the amino acid frequencies are determined by functional requirements of

<sup>593</sup> the protein while the other models make no assumptions about amino acid frequencies. In all cases, we

<sup>594</sup> assumed that the model was partitioned by gene, but with branch lengths linked across genes.

We also compared SelAC models with 195 codon models in IQtree. This is popular software

implementing recent (i.e.,

<sub>) models of evolution, including codon models. Likelihoods appear generally comparable between models in selac and models in IQtree (though details of their</sub>

<sup>595</sup> For SelAC, we compared our estimates of $\phi' = \psi'/\mathbf{B}$, which represents the average protein synthesis

<sup>596</sup> rate of a gene, to estimates of gene expression from empirical data. Specifically, we examined gene

20

expression data for five of the six species measured during log-growth phase. Gene expression in this context corresponds to mRNA abundances, which were measured using either microarrays (*C. glabrata* and *S. castellii*, or RNA-Seq (*S. paradoxus*, *S. mikatae*, and *S. cerevisiae*). We obtained expression data for the remaining species, *S. kudriavzevii*, which was measured at the beginning of the stationary phase from the Gene Expression Omnibus (GEO). Saccharomyces, however, only enter the stationary growth phase in response to severe stress, such as starvation. In addition, only 56 % of the genes examined with SelAC had expression measurements available. For these reasons, we excluded *S. kudriavzevii* from our comparisons of empirical gene expression.

For further comparison, we also predicted the average protein synthesis rate for each gene $\phi$ by analyzing gene and genome-wide patterns of synonymous codon usage using ROC-SEMPPR (**?**) for each individual genome. While, like SelAC, ROC-SEMPPR uses codon level information, it does not rely on any interspecific comparisons and, unlike SelAC, uses only the intra- and inter-genic frequencies of synonymous codon usage as its data. Nevertheless, ROC-SEMPPR predictions of gene expression $\phi$ correlates strongly (Pearson $r = 0.53 - 0.74$) with a wide range of laboratory measurements of gene expression (**?**).

While one of our main objectives was to determine the improvement of fit that SelAC has with respect to other standard phylogenetic models, we also evaluated the adequacy of SelAC. Model fit, measured with assessments such as the Akaike Information Criterion (AIC), can tell which model is least bad as an approximation for the data, but it does not reveal whether a model is actually doing a good job of representing the data. An adequate model does the latter, one measure of which is that data generated under the model resemble real data (**?**). For example, **?** assessed whether parsimony scores and the size of monomorphic clades of empirical data were within the distributions of simulated data under a new model and the best standard model; if the empirical summaries were outside the range for each, it would have suggested that neither model was adequately modeling this part of the biology.

In order to test adequacy for a given gene we first remove a particular taxon from the data set and the phylogeny. A marginal reconstruction of the likeliest sequence across all remaining nodes is conducted under the model, including the node where the pruned taxon attached to the tree. The marginal probabilities of each site are used to sample and assemble the starting coding sequence. This sequence is then evolved along the branch, periodically being sampled and its current functionality assessed. We repeat this process 100 times and compare the distribution of trajectories against the observed functionality calculated for the gene. For comparison, we also conducted the same test, by simulating the sequence under the standard GTR + $\Gamma$ nucleotide model, which is often used on these

21

data but does not account for the fact that the sequences are protein coding, and under FMutSel, which includes selection on codons but in a fundamentally different way as our model.

**The appropriate estimator of bias for AIC**

As part of the model set described above, we also included a reduced form of each of the two SelAC models, SelAC and SelAC+$\Gamma$. Specifically, rather than optimizing the amino acid at any given site, we assume the the most frequently observed amino acid at each site is the optimal amino acid $a^*$. We refer to these 'majority rule' models as SelAC$_M$ and SelAC$_M$+$\Gamma$ and note that these majority rule formulations greatly accelerate model fitting.

Since these majority rule models assume that the optimal amino acids are known prior to fitting of our model, it is tempting to reduce the count of estimated parameters in the model by the number of parameters estimated using majority rule. While using majority rule does not necessarily provide the most likely parameter estimate, it nevertheless uses the data to generate the estimate and represents a parameter estimated from the data. Thus, despite having become standard behavior in the field of phylogenetics, this reduction is statistically inappropriate. Because the difference in the number of parameters $K$ when counting or not counting the number of nucleotide sites drops out when comparing nucleotide models with AIC, this statistical issue does not apply to nucleotide models. It does, however, matter for AICc, where $K$ 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. **???**) 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. This uncertainty was highlighted by **?**: they chose to use number of sites, but mentioned in their discussion that sample size also depends on the number of taxa. **?** also mention that while the number of sites is often taken as sample size, whether that is appropriate in phylogenetics is not entirely clear. One approach incorporating both number of taxa and sites in calculating AICc is the program SURFACE implemented by **?**, 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, **?** 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

22

of ? in order to examine a large set of different penalty functions and how well they approximate the remaining portion of the Kullback-Liebler (KL) divergence between two models after accounting for the deviance (i.e., $-2\mathcal{L}$) (see Appendix 1 for more details).

## Acknowledgements

23

**Table**

|  |  | Parameters |  |  |  | Model |
| Model | logLik | Estimated | AIC | AICc | $\Delta$AICc | Weight |
| --- | --- | --- | --- | --- | --- | --- |
| SelAC+$\Gamma$ | -453,620.8 | 50,005 | 1,007,252 | 1,027,314 | 0 | >0.999 |
| SelAC | -464,114.8 | 50,004 | 1,028,238 | 1,048,299 | 20,985 | <0.001 |
| SelAC$_M$+$\Gamma$ | -465,106.9 | 50,005 | 1,030,224 | 1,050,286 | 22,972 | <0.001 |
| SelAC$_M$ | -478,302.4 | 50,004 | 1,056,613 | 1,076,674 | 49,360 | <0.001 |
| FMutSel | -597,140.7 | 178 | 1,194,637 | 1,194,638 | 167,324 | <0.001 |
| GY94 | -612,670.4 | 111 | 1,225,563 | 1,225,563 | 198,249 | <0.001 |
| GTR+$\Gamma$ | -655,166.4 | 610 | 1,311,553 | 1,311,554 | 284,240 | <0.001 |

**Table 1.** Comparison of model fits using AIC, AICc, and AIC$_w$ from analyses of 100 selected genes from 6 yeast taxa **?**. 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. Also, the sample size used in the calculation of AICc is assumed to be equal to the size of the matrix (number of taxa x number of sites, which is 6 x 49,881).

**Figures**



**FIG. 1.** Summary of 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 SelAC model to the 106 gene data set and phylogeny of **?**, 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 all 50 simulated datasets (d). Gene index on the x-axis refers to the arbitrary number assigned to the simulated gene.

**FIG. 2.** Comparisons between estimates of average protein translation rate $\hat{\phi}_{\mathrm{SelAC}}$ obtained from SelAC+Γ and direct measurements of expression for individual yeast taxa across the 100 selected genes from **?** measured during log-growth phase. Estimates of $\hat{\phi}_{\mathrm{SelAC}}$ were generated by dividing the composite term $\psi'$ by $\mathbf{B}(\vec{a}_i|\vec{a}^*)$. Gene expression was measured using either RNA-Seq (a)-(c) or microarray (d)-(e). The equations in the upper left hand corner of each panel represent the regression fit and the Pearson correlation coefficient $r$.

**FIG. 3.** Comparisons between $\omega_{\text{FMutSel}}$, which is the nonsynonymous/synonymous mutation ratio in FMutSel, SelAC+$\Gamma$ estimates of protein functionality production rates $\hat{\psi}_{\text{SelAC}}$ (a), RNA-Seq based measurements of mRNA abundance $\phi_{\text{RNA-seq}}$ (b), and ROC-SEMPPER's estimates of protein translation rates $\phi_{\text{ROC}}$, which are based solely on *S. cerevisiae*'s patterns of codon usage bias (c), for *S. cerevisiae* across the 100 selected genes from **?**. As in Figure 2, the equations in the upper right hand corner of each panel provide the regression fit and correlation coefficient.

**FIG. 4.** (a) Maximum likelihood estimates of branch lengths under SelAC+Γ for 100 selected genes from **?**. 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 $\mathbf{B}(\vec{a}_{\text{obs}}|\vec{a}^*)$ that is far closer to the observed (dashed black line) than data sets produced from parameters of either FMutSel or GTR + Γ.

## Supporting Materials

Supporting Materials for *Population Genetics Based Phylogenetics Under Stabilizing Selection for an Optimal Amino Acid Sequence: A Nested Modeling Approach* by Beaulieu *et al.* (In Review).

**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 (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 with both measures, which are also reasonably well correlated with each other (Figure 2 - 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.

**FIG. S1.** Comparisons between estimates of $\phi$ obtained from SelAC+$\Gamma$ and the predicted gene expression from the ROC SEMPER model (**?**) for individual yeast taxa across the 100 selected genes from **?**. 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.

**FIG. S2.** Comparisons of predicted gene expression from the ROC SEMPER model (**?**) and direct measurements of expression from RNA-Seq or microarray data for individual yeast taxa across the 100 selected genes from **?**. The equations in the upper left hand corner of each panel represent the regression fit and correlation coefficient.

## Simulations

As stated in the main text, overall, the simulation results indicate that the SelAC model can reasonably recover the known values of the generating model (Figure S3 - S5). This includes not only the parameters
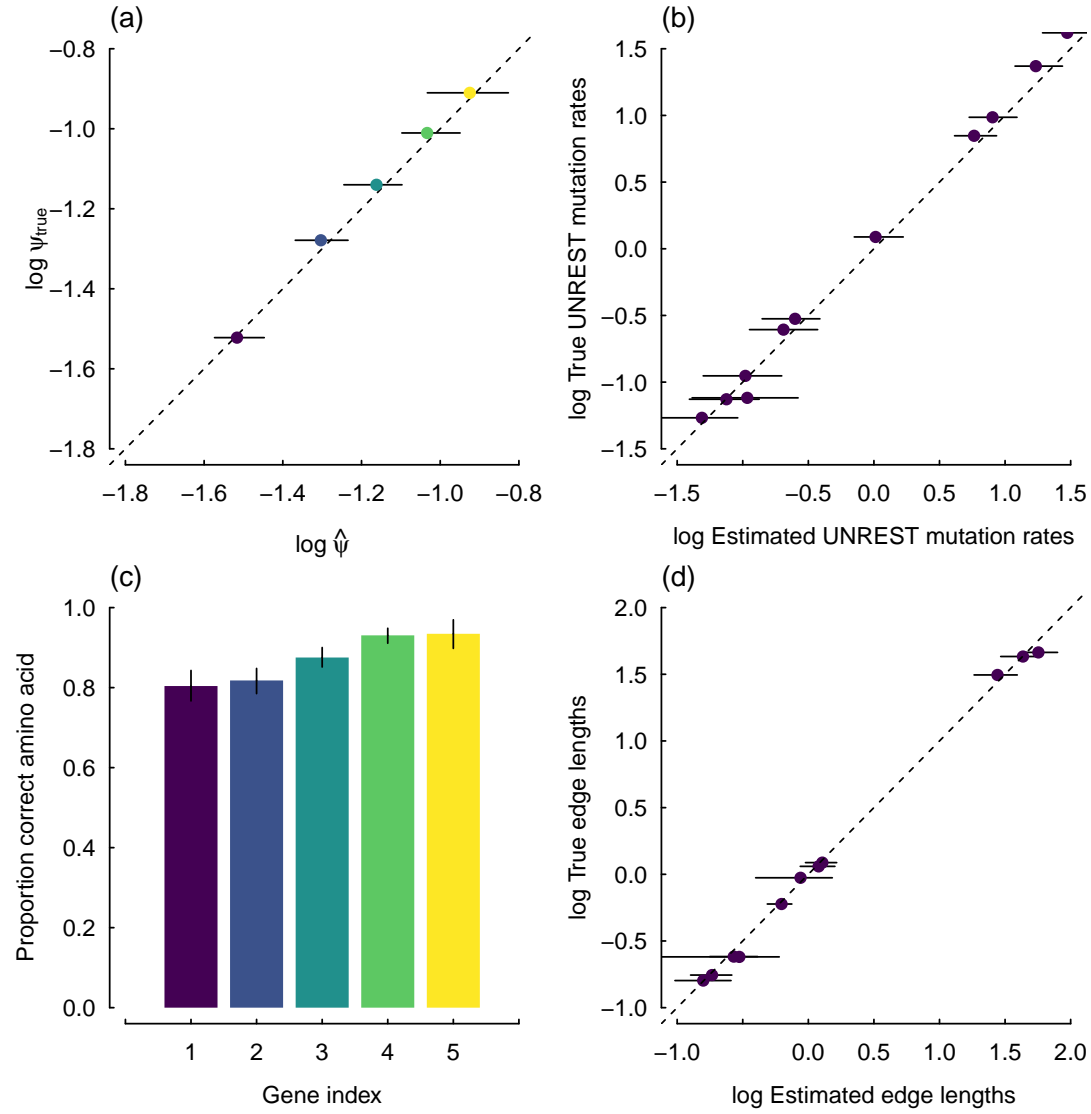
**FIG. S3.** Same figure as in Figure 1 in the main text, except the generating model does not include a site-specific sensitivity in the generating model (i.e., $\alpha_G = \infty$).
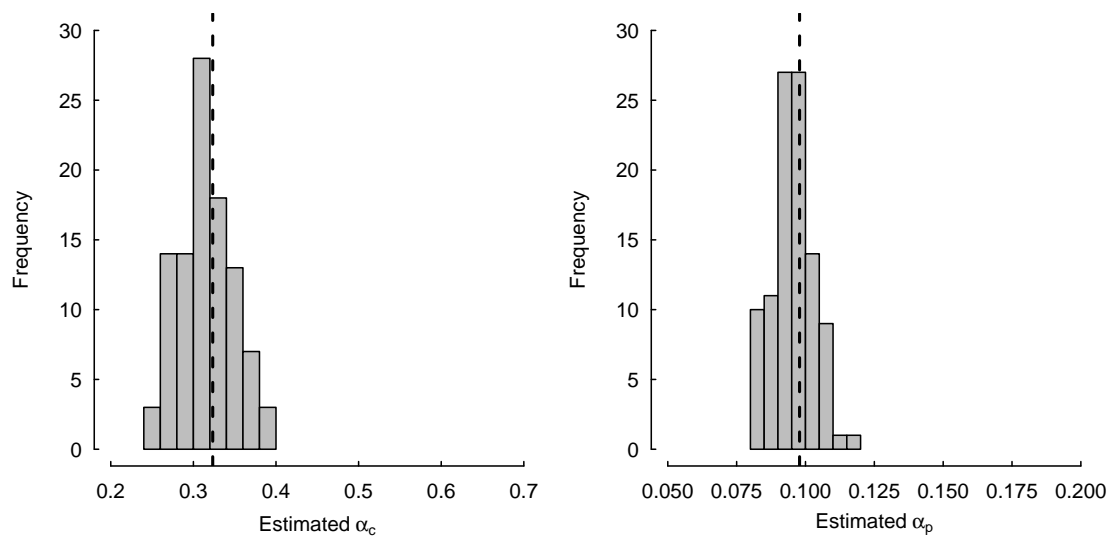
**FIG. S4.** The distribution of estimates of the Grantham weights, $\alpha_c$ and $\alpha_p$, in a SelAC model, where we assume $\alpha_G = \infty$, and thus no site-specific sensitivity in the generating model. The dashed line represents the value used in the generating model.
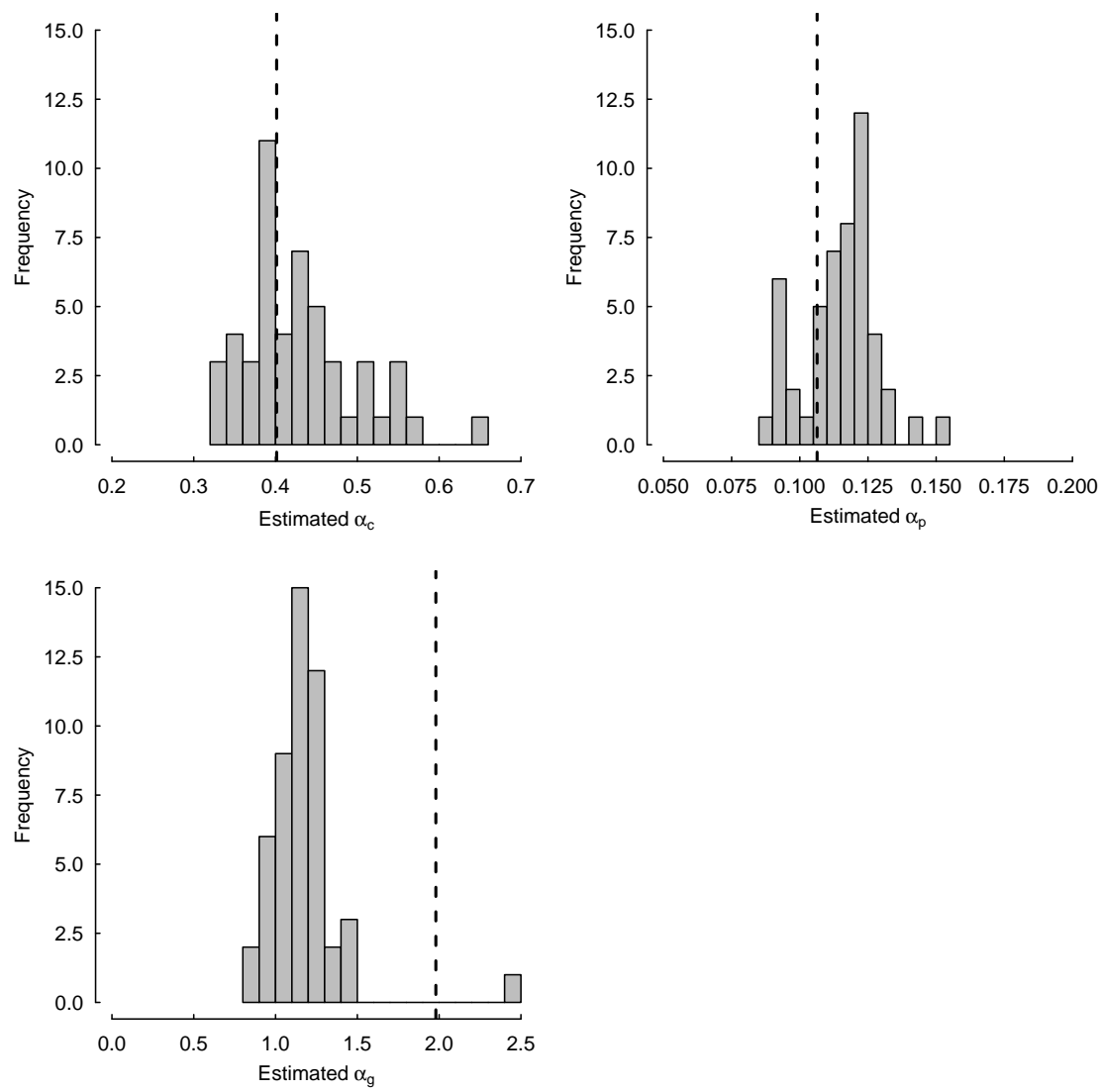
**FIG. S5.** Same figure as in Figure S4, except the generating model includes site-specific sensitivity in the generating model (i.e., $\alpha_G$). Unlike, Grantham weights, which showed no systematic bias, there is a downward bias in estimates of $\alpha_G$.