

1 RH: BEAULIEU ET AL.— Pop. Gen. Based Phylo.

2 **Population Genetics Based Phylogenetics Under**
3 **Stabilizing Selection for an Optimal Amino Acid**
4 **Sequence: A Nested Modeling Approach**

5 JEREMY M. BEAULIEU^{1,2,3}, BRIAN C. O'MEARA^{2,3}, RUSSELL ZARETZKI⁴,
6 CEDRIC LANDER^{2,3}, JUAN JUAN CHAI^{2,5}, AND MICHAEL A. GILCHRIST^{2,3,*}

7 ¹Department of Biological Sciences, University of Arkansas, Fayetteville, AR 72701

8 ²Department of Ecology & Evolutionary Biology, University of Tennessee, Knoxville, TN
9 37996-1610

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

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

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

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

Version dated: Friday 24th March, 2017

14

ABSTRACT

15 We present a phylogenetic approach rooted in the field of population genetics that more
16 realistically models the evolution of protein-coding DNA under the assumption of
17 stabilizing selection for a gene specific, optimal amino acid sequence. In addition to being
18 consistent with the fundamental principles of population genetics, our new set of models,
19 which we collectively call SelAC, fit phylogenetic data astronomically better than popular
20 models, suggesting strong potential for more accurate inference of phylogenetic trees and
21 branch lengths. SelAC also demonstrates that a large amount of biologically meaningful
22 information is accessible when using a nested set of mechanistic models. For example, for
23 each position SelAC provides a probabilistic estimate of any given amino acid being
24 optimal. Because SelAC assumes the strength of selection is proportional to the expression
25 level of a gene, SelAC provides gene specific estimates of protein synthesis rates. Finally,
26 because SelAC's is a nested approach based on clearly stated biological assumptions, it can
27 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 agnostic with regards to the different amino acid substitutions and
34 their impact on gene function (e.g. F81, F84, HYK85, TN93, and GTR, see Yang (2014)
35 for an overview).

36 Another set of models attempt to include a 'selection' term ω , but the link between
37 ω 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 are far from clear. For
39 instance, ω 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 as 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 ω 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, ω based approaches only reasonably describe a

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

58 Fortunately, given the continual growth in computational power available to
59 researchers, it is now possible to utilize a more general set of population genetics based
60 models for the purpose of phylogenetic analysis (e.g. Halpern and Bruno 1998; Robinson
61 et al. 2003; Lartillot and Philippe 2004; Rodrigue and Lartillot 2014). One lesson from the
62 field of population genetics is even when there are only a few fundamental evolutionary
63 forces at play (mutation, drift, selection, and linkage effects), describing the evolutionary
64 behavior of a system in which there are non-linear interactions between sites, such as
65 epistasis, quickly becomes extremely challenging. Fortunately, under the simplifying
66 assumptions of additivity between sites and alleles, calculating stationary and substitution
67 probabilities are relatively straightforward, making fitting additive models of the
68 evolutionary process to sequence data computationally feasible.

69 MATERIALS & METHODS

70 *Overview*

71 We model the substitution process as a classic Wright-Fisher process which includes
72 the forces of mutation, selection, and drift (Fisher 1930; Kimura 1962; Wright 1969; Iwasa
73 1988; Berg and Lässig 2003; Sella and Hirsh 2005; McCandlish and Stoltzfus 2014). For
74 simplicity, we ignore linkage effects and, as a result of this and other assumptions, our
75 method behaves in a site independent manner. Our approach, which we call SelAC
76 (Selection on Amino acids and Codons), is developed in the same vein as previous

77 phylogenetic applications of the Wright-Fisher process (e.g. Muse and Gaut 1994; Halpern
78 and Bruno 1998; Yang and Nielsen 2008; Rodrigue et al. 2005; Koshi and Goldstein 1997;
79 Koshi et al. 1999; Dimmic et al. 2000; Thorne et al. 2012; Lartillot and Philippe 2004;
80 Rodrigue and Lartillot 2014). Similar to Lartillot’s work (Lartillot and Philippe 2004;
81 Rodrigue and Lartillot 2014), we assume there is a finite set of rate matrices describing the
82 substitution process and that each position within a protein is assigned to a particular rate
83 matrix category. Unlike this previous work, we assume *a priori* there are 20 different
84 families of rate matrices, one family for when a given amino acid is favored at a site. As a
85 result, SelAC allows us to quantitatively evaluate the support for a particular amino acid
86 being favored at a particular position within the protein encoded by a particular gene.

87 Because SelAC requires twenty families of 61×61 matrices, the number of
88 parameters needed to implement SelAC would, without further assumptions, be extremely
89 large. To reduce the number of parameters needed while still maintaining a high degree of
90 biological realism, we construct our gene and amino acid specific substitution matrices
91 using a submodel nested within our substitution model, similar to approaches in Gilchrist
92 (2007); Shah and Gilchrist (2011); Gilchrist et al. (2015).

93 One advantage of a nested modeling framework is that it requires only a handful of
94 genome wide parameters such as nucleotide specific mutation rates (scaled by effective
95 population size N_e), side chain physicochemical weighting parameters, and a shape
96 parameter describing the distribution of site sensitivities. In addition to these genome wide
97 parameters, SelAC requires a gene g specific expression parameter ψ_g which describes the
98 average rate at which the protein’s functionality is produced by the organism. Currently, ψ
99 is fixed across the phylogeny, though relaxing this assumption is a goal of future work. The
100 gene specific parameter ψ_g is multiplied by additional model terms to make a composite
101 term ψ'_g which scales the strength and efficacy of selection for the optimal amino acid
102 sequence relative to drift. In terms of the functionality of the protein encoded, we assume

103 that for any given gene there exists an optimal amino acid sequence \vec{a}_* and that, by
104 definition, a complete, error free peptide consisting of a_* provides one unit of the gene's
105 functionality. We also assume that natural selection favors genotypes that are able to
106 synthesize their proteome efficiently than their competitors and that each savings of an
107 high energy phosphate bond per unit time leads to a constant proportional gain in fitness
108 q . SelAC also requires the specification (as part of parameter optimization) of an optimal
109 amino acid at each position or site within a coding sequence which, in turn, makes it the
110 largest category of parameters we estimate. Because we use a submodel to derive our
111 substitution matrices, SelAC requires the estimation of a fraction of the parameters
112 required when compared to approaches where the substitution rates are allowed to vary
113 independently (Halpern and Bruno 1998; Lartillot and Philippe 2004; Rodrigue and
114 Lartillot 2014).

115 As with other phylogenetic methods, we generate estimates of branch lengths and
116 nucleotide specific mutation rates. In addition, because the math behind our model is
117 mechanistically derived, our method can also be used to make quantitative inferences on
118 the optimal amino acid sequence of a given protein as well as the average synthesis rate of
119 each protein used in the analysis. The mechanistic basis of SelAC also means it can be
120 easily extended to include more biological realism and test more explicit hypotheses about
121 sequence evolution.

122 *Mutation Rate Matrix μ*

123 We begin with a 4x4 nucleotide mutation matrix that defines a model for mutation rates
124 between individual bases. For our purposes, we rely on the general unrestricted
125 model (Yang 1994, UNREST) because it makes no constraint on the instantaneous rate of
126 change between any pair of nucleotides. We note, however, that more constrained models,
127 such as the Jukes-Cantor (JC), Hasegawa-Kishino-Yano (HKY), or the general time

128 reversible model (GTR), can also be used. The 12 parameter UNREST model defines the
129 relative rates of change between a pair of nucleotides. Thus, we arbitrarily set the G→T
130 mutation rate to 1, resulting in 11 free mutation rate parameters in the 4x4 mutation
131 nucleotide mutation matrix. The nucleotide mutation matrix is also scaled by a diagonal
132 matrix π whose entries correspond to the equilibrium frequencies of each base. These
133 equilibrium nucleotide frequencies are determined by analytically solving $\pi \times \mathbf{Q} = 0$. We
134 use this \mathbf{Q} to populate a 61×61 codon mutation matrix μ , whose entries $\mu_{i,j}$ describe the
135 mutation rate from codon i to j under a "weak mutation" assumption. That is, the rate of
136 allele fixation is much greater than $N_e\mu$ and $N_e\mu \ll 1$, such that evolution is mutation
137 limited, codon substitutions only occur one nucleotide at a time and, as a result, the rate
138 of change between any pair of codons that differ by more than one nucleotide is zero.

139 While the overall model does not assume equilibrium, we still need to scale our
140 mutation matrices μ . Traditionally, it is rescaled such that at equilibrium, one unit of
141 branch length represents one expected substitution per site. Here, a scaling factor is
142 calculated as the average rate $-\sum_i \mu_i \pi_i = 1$, where i indexes a particular codon in a given
143 gene. The final mutation rate matrix is the original mutation rate matrix multiplied by
144 $1/\text{scaling factor}$.

145 *Protein Synthesis Cost-Benefit Function η*

146 SelAC links fitness to the product of the cost-benefit function of a gene g , η_g , and the
147 organism's average target synthesis rate of the functionality provided by gene g , ψ_g . This is
148 because the average flux energy an organism spends to meet its target functionality
149 provided by gene g is $\eta_g \times \psi_g$. In order to link genotype to our cost-benefit function
150 $\eta = \mathbf{B}/\mathbf{C}$, we begin by defining our benefit function \mathbf{B} .

151 *Benefit*.— Our benefit function \mathbf{B} measures the functionality of the amino acid sequence
152 \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}_* .

153 By definition, $\mathbf{B}(\vec{a}_*) = 1$ and $\mathbf{B}(\vec{a}_i|\vec{a}_*) < 1$ for all other sequences. We assume all amino
154 acids within the sequence contribute to protein function and that this contribution declines
155 as an inverse function of physicochemical distance between each amino acid and the
156 optimal. Formally, we assume that

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

157 where n_g is the length of the protein, $d(a_{i,p}, a_{*,p})$ is a weighted physicochemical distance
158 between the amino acid encoded in gene i for position p and $a_{*,p}$ is the optimal amino acid
159 for that position of the protein. For simplicity, we define the distance between a stop codon
160 and a sense codon as effectively infinite and, as a result, nonsense mutations are effectively
161 lethal. The term G_p describes the sensitivity of the protein's function to deviation in
162 physicochemical space. There are many possible measures for physicochemical distance; we
163 use (Grantham 1974) distances by default, though others may be chosen. We assume that
164 $G_p \sim \text{Gamma}(\alpha = \alpha_G, \beta = \alpha_G)$ in order to ensure $\mathbb{E}(G_p) = 1$.

165 At the limit of $\alpha_G \rightarrow \infty$, the model collapses to a model with uniform sensitivity of
166 $G_p = 1$ for all positions p . $\mathbf{B}(\vec{a}_i|\vec{a}_*)$ is inversely proportional to the average physicochemical
167 deviation of an amino acid sequence \vec{a}_i from the optimal sequence \vec{a}_* weighted by each sites
168 sensitivity to this deviation. $\mathbf{B}(\vec{a}_i|\vec{a}_*)$ can be generalized to include second and higher order
169 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 any indirect costs of protein assembly that vary between genotypes and define,

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

$$= A_1 + A_2 n \quad (3)$$

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

172

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, α_c , α_p , α_v , respectively, which we refer to as 'Grantham weights'. Because physicochemical distance is ultimately weighted by a gene's specific average protein synthesis rate ψ , another

parameter we estimate, there is a problem with parameter identifiability. 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 estimates of α_c and α_p and ψ are scaled relative to this choice for α_v . More specifically,

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

173

Linking Protein Synthesis to Allele Substitution

174

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

175

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

176

beneficial function and that the organism needs that functionality to be produced at a

177

target average rate ψ . By definition, the optimal amino acid sequence for a given gene, \vec{a}_* ,

178

produces one unit of functionality. Second, we assume that protein expression is regulated

179

by the organism to ensure that functionality is produced at rate ψ . As a result, the realized

180

average protein synthesis rate of a gene, ϕ , is equal to $\psi/\mathbf{B}(\vec{a})$ and the total energy flux

181

allocated towards meeting the target functionality of a particular gene is $\eta(\vec{c})\psi$. The fitness

182

cost for a genotype encoding a suboptimal protein sequence stems from the need to

183

produce $1/\mathbf{B}(\vec{a})$ proteins in order to produce the equivalent functionality of one protein

184

consisting of the optimal amino acid sequence a_* . For example, a protein encoding allele

185

which has a 10% reduction in functionality relative to the optimal sequence,

186

i.e. $\mathbf{B}(\vec{a}) = 0.9$, will have the same energetic burden and selective cost relative to its

187

optimal sequence as a protein encoding allele of similar length which has a 20% reduction

188

in functionality but whose target synthesis rate is 1/2 of the first protein.

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 ψ leads to a slight and proportional decrease in fitness W . This assumption, in turn, implies

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

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

Correspondingly, the ratio of fitness between two genotypes is,

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

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

$$\frac{W_i}{W_j} = \exp \left[-\frac{A_0 (A_1 + A_2 n)}{n} \right] \quad (5)$$

$$\times \sum_{p \in \mathbb{P}} [d(a_{i,p}, a_{*,p}) - d(a_{j,p}, a_{*,p})] \psi \quad (6)$$

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,

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

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,

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

192 where, given our weak mutation assumption, $\mu_{i,j} = 0$ when two codons differ by more than
193 one nucleotide. In the end, each optimal amino acid has a separate 64 x 64 substitution
194 rate matrix \mathbf{Q}_a , which incorporates selection for the amino acid (and the fixation rate
195 matrix this creates) as well as the common mutation parameters across optimal amino
196 acids. This results in the creation of 20 \mathbf{Q}_a matrices, one for each amino acid, with up to
197 26,880 unique rates, based on few parameters (one to 11 mutation rates, two free
198 Grantham weights, the cost of protein assembly, A_1 and A_2 , the gene specific target
199 functionality synthesis rate ψ , and optimal amino acid at each position p , $a_{*,p}$), which can
200 either be specified *a priori* or estimated from the data. SelAC can be generalized to allow
201 transitions between optimal amino acids as well as between codons, which would result in a
202 $(20 \times 64) \times (20 \times 64) = 1344 \times 1344$ matrix.

203 Finally, given our assumption of independent evolution among sites, the probability
204 of the whole data set is the product of the probabilities of observing the data at each

205 individual site. Thus, the log likelihood \mathcal{L} of amino acid a being optimal at a given site
206 position p is calculated as

$$\mathcal{L}(\mathbf{Q}_a|\mathbf{D}_p, \mathbf{T}) \propto \mathbf{P}(\mathbf{D}_p|\mathbf{Q}_a, \mathbf{T}) \quad (7)$$

207 In this case, the data, \mathbf{D}_p , are the observed codon states at position p for the tips of the
208 phylogenetic tree with topology \mathbf{T} . For our purposes we take \mathbf{T} as given but it could be
209 estimated as well. The pruning algorithm of Felsenstein (1981) is used to calculate $\mathcal{L}(\mathbf{Q}_a)$.
210 The log likelihood is maximized by estimating the genome scale parameters which consist
211 of 11 mutation parameters which are implicitly scaled by $2N_e/b$, and two Grantham
212 distance parameters, α_c and α_p , and the sensitivity distribution parameter α_G . Because A_0
213 and ψ_g always co-occur and are scaled by N_e , for each gene g we estimate a composite term
214 $\psi'_g = \psi_g A_0 b N_e$ and the optimal amino acid for each position $a_{*,p}$ of protein. When
215 estimating α_G , the likelihood then becomes the average likelihood which we calculate using
216 the generalized Laguerre quadrature with $k = 4$ points (Felsenstein 2001).

217 *Implementation*

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

227 For the SelAC models, the starting guess for the optimal amino acid at a site comes
228 from ‘majority’ rule, where the initial optimum is the most frequently observed amino acid

229 at a given site (ties resolved randomly). Our optimization routine then proceeds by cycling
230 through multiple phases. The first phase optimizes the branch lengths while holding the
231 model parameters constant. The second phase optimizes the gene specific composite
232 parameter $\psi' = A_0\psi N_e$ across genes, while holding constant both the branch lengths and
233 the model parameters shared across the genome (i.e., α_c and α_p , and the sensitivity
234 distribution parameter α_G). This is followed by a third phase that optimizes the
235 parameters across the genome, while keeping the branch lengths and the composite
236 parameters constant. Finally, the fourth phase estimates the optimal amino acid at each
237 site while keeping the branch lengths and all model parameters at their current values.
238 This entire procedure is repeated six times. For optimization of a given set of parameters,
239 we rely on a bounded subplex routine (Rowan 1990) in the package `NLopt` (Johnson 2012)
240 to maximize the log-likelihood function. To help the optimization navigate through local
241 peaks, we perform a set of independent analyses with different sets of naive starting points
242 with respect to the gene specific composite ψ' parameters, α_c , and α_p . Confidence in the
243 parameter estimates can be generated by an 'adaptive search' procedure that we
244 implemented to provide an estimate of the parameter space that is some pre-defined
245 likelihood distance (e.g., 2 lnL units) from the maximum likelihood estimate (MLE), which
246 follows Beaulieu and OMeara (2016); Edwards (1984).

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

255

Simulations

256 We evaluated the performance of our codon model by simulating datasets and estimating
257 the bias of the inferred model parameters from these data. Our 'known' parameters under
258 a given generating model were based on fitting SelAC to the 106 gene data set and
259 phylogeny of Rokas et al. (2003). The tree used in these analyses is outdated with respect
260 to the current hypothesis of relationships within Saccharomyces, but we rely on it simply as
261 a training set that is separate from our empirical analyses (see section on Analyzing Yeast
262 Genome). Bias in the model parameters were assessed under two generating models: one
263 where we assumed a model of SelAC assuming $\alpha_G = \infty$, and one where we estimated α_G
264 from the data. Under each of these two scenarios, we used parameter estimates from the
265 corresponding empirical analysis and simulated 50 five-gene data sets. For the gene specific
266 composite parameter ψ'_g the 'known' values used for the simulation were five evenly spaced
267 points along the rank order of the estimates across the 106 genes. The MLE estimate for a
268 given replicate were taken as the fit with the highest log-likelihood after running five
269 independent analyses with different sets of naive starting points with respect to the
270 composite ψ'_g parameter, α_c , and α_p . All analyses were carried out in our `selac` R package.

271

Analysis of yeast genome and tests of model adequacy

272 We focus our empirical analyses on the large yeast data set and phylogeny of Salichos and
273 Rokas (2013). The yeast genome is an ideal system to examine our phylogenetic estimates
274 of gene expression and its connection to real world measurements of these data within
275 individual taxa. The complete data set of Salichos and Rokas (2013) contain 1070
276 orthologs, where we selected 100 at random for our analyses. We also focus our analyses
277 only on Saccharomyces *sensu stricto*, including their sister taxon *Candida glabrata*, and we
278 rely on the phylogeny depicted in Fig. 1 of Salichos and Rokas (2013) for our fixed tree.

279 We fit both the new models described in this paper, as well as two codon models, GY94
280 and FMutSel0, and a standard GTR + Γ nucleotide model. The FMutSel0 model, which
281 assumes that the amino acid frequencies are determined by functional requirements of the
282 protein, is the most similar to our model. In all cases, we assumed that the model was
283 partitioned by gene, but with branch lengths linked across genes.

284 For SelAC, we compared our estimates of $\phi' = \psi'/\mathbf{B}$, which represents the average
285 protein synthesis rate of a gene, to estimates of gene expression from empirical data.
286 Specifically, we obtained expression data for five of the six species used - four species were
287 measured during log-growth phase, whereas the other was measured at the beginning of the
288 stationary phase (*S. kudriavzevii*) from the Gene Expression Omnibus (GEO). Gene
289 expression was measured using either Microarray chips (*C. glabrata*, *S. castellii*, and *S.*
290 *kudriavzevii*) or RNA-Seq (*S. paradoxus*, *S. mikatae*, and *S. cerevisiae*). For further
291 comparison, we also predicted protein synthesis rate (ϕ) by analyzing gene and
292 genome-wide patterns of synonymous codon usage using ROC-SEMPPR (Gilchrist et al.
293 2015) for each individual genome. While, like SelAC, ROC-SEMPPR uses codon level
294 information, it does not rely on any inter-specific comparisons and, unlike selac, assumes
295 selection on synonymous codon usage is contributing to these patterns. Nevertheless,
296 ROC-SEMPPR predictions of gene expression ϕ correlates strongly ($r = 0.53 - 0.74$) with
297 a wide range of laboratory measurements of gene expression (Gilchrist et al. 2015).

298 While one of our main objectives was to determine the improvement of fit that
299 SelAC has with respect to other standard phylogenetic models, we also evaluated the
300 adequacy of SelAC. Model fit, measured with assessments such as the Akaike Information
301 Criterion (AIC), can tell which model is least bad as an approximation for the data, but it
302 does not reveal whether a model is actually doing a good job of representing the biological
303 processes. An adequate model does the latter, one measure of which is that data generated
304 under the model resemble real data (Goldman 1993). For example, Beaulieu et al. (2013)

305 assessed whether parsimony scores and the size of monomorphic clades of empirical data
306 were within the distributions of simulated under a new model and the best standard
307 model; if the empirical summaries were outside the range for each, it would have suggested
308 that neither model was adequately modeling this part of the biology. For a given gene we
309 first remove a particular taxon from the data set and the phylogeny. A marginal
310 reconstruction of the likeliest sequence across all remaining nodes is conducted under the
311 model, including where the attachment point of pruned taxon to the tree. The marginal
312 probabilities of each site are used to sample and assemble the starting coding sequence.
313 This sequence is then evolved along the branch, periodically being sampled and its current
314 functionality assessed. We repeat this process 100 times and compare the distribution of
315 trajectories against the observed functionality calculated for the gene. For comparison, we
316 also conducted the same test, by simulating the sequence under the standard GTR + Γ
317 nucleotide model, which is often used on these data but does not account for the fact that
318 the sequence codes for a specific protein, and under FMutSel0, which includes selection on
319 codons but in a fundamentally different way as our model.

320 RESULTS

321 By linking transition rates $q_{i,j}$ to gene expression ψ , our approach allows use of the same
322 model for genes under varying degrees of stabilizing selection. Specifically, we assume the
323 strength of stabilizing selection for the optimal sequence, \vec{a}_* , is proportional to the average
324 protein synthesis rate ϕ , which we can estimate for each gene. In regards to model fit, our
325 results clearly indicated that linking the strength of stabilizing selection for the optimal
326 sequence to gene expression substantially improves our model fit. Further, including the
327 single random effects term $G \sim \text{Gamma}(\alpha_G, \beta_g)$ to allow for heterogeneity in this selection
328 between sites within a gene, improves the ΔAIC of SelAC+ Γ score over the simpler SelAC

329 model by over 23,000 AIC units. Using either ΔAIC or AIC_w as our measure of model
330 support, the SelAC models fit extraordinarily better than GTR + Γ , GY94, or FMutSel0
331 (Table 1). This is in spite of the need for estimating the optimal amino acid at each
332 position in each protein, which accounts for more than 47,000 additional model parameters.
333 Even when compared to the next most parameter rich codon model in our model set,
334 FMutSel0, SelAC+ Γ model shows nearly 400,000 AIC unit improvement over FMutSel0.

335 With respect to estimates of ϕ within SelAC, they were strongly correlated with two
336 separate measures of gene expression, one empirical (See Figure S1), and one model-based
337 prediction that does not account for shared ancestry (Figure S1-S2). In other words, using
338 only codon sequences our model can predict which genes have high or low expression levels.
339 The estimate of the α_G parameter, which describes the site-specific variation in sensitivity
340 of the protein's functionality, indicated a moderate level of variation in gene expression
341 among sites. Our estimate of $\alpha_G = 1.40$, produced a distribution of sensitivity terms G
342 ranged from 0.344-7.16, but with nearly 90% of the weight for a given site-likelihood being
343 contributed by the 0.344 and 1.48 rate categories. In simulation, however, of all the
344 parameters in the model, only α_G showed a consistent bias, in that the estimates were
345 generally underestimated (see Supporting Materials). Other parameters in the model, such
346 as the Grantham weights, provide an indication as to the physicochemical distance between
347 amino acids. Our estimates of these weights only strongly deviate from Grantham's 1974
348 original estimates in regards to composition weight, α_c , which is the ratio of noncarbon
349 elements in the end groups to the number of side chains. Our estimate of the composition
350 weighting factor of $\alpha_c=0.484$ is 1/4th the value estimate by Grantham which suggests that
351 the substitution process is less sensitive to this physicochemical property when shared
352 ancestry and variation in stabilizing selection are taken into account.

353 It is important to note that the nonsynonymous/synonymous mutation ratio, or ω ,
354 which we estimated for each gene under the FMutSel0 model strongly correlated with our

355 estimates of $\phi' = \psi'/\mathbf{B}$ where \mathbf{B} depends on the sequence of each taxa. In fact, ω showed
356 similar, though slightly reduced correlations, with the same empirical estimates of gene
357 expression described above (See Figure 2). This would give the impression that the same
358 conclusions could have been gleaned using a much simpler model, both in terms of the
359 number of parameters and the assumptions made. However, as we discussed earlier, not
360 only is this model greatly restricted in terms of its biological feasibility, SelAC clearly
361 performs better in terms of its fit to the data and biological realism. For example, when we
362 simulated the sequence for *S. cerevisiae*, starting from the ancestral sequence under both
363 GTR + Γ and FMutSel0, the functionality of the simulated sequence moves away from the
364 observed sequence, whereas SelAC remains near the functionality of the observed sequence
365 (Figure 3b). In a way, this is somewhat unsurprising, given that both GTR + Γ and
366 FMutSel0 are agnostic to the functionality of the gene, but it does highlight the
367 improvement in biological realism in amino acid sequence evolution that SelAC provides.
368 We do note that the adequacy of the SelAC model does vary among individual taxa, and
369 does not always perfectly match the observed functionality. For instance, *S. castellii* is
370 simulated with consistently higher functionality than observed (Figure 3c). We suspect this
371 is an indication that assuming a single set of optimal amino acid across all taxa may be too
372 simplistic, but we cannot also rule out other potential simplifying assumptions in our
373 model, such as a single set of Grantham weights and α_G values or the simple, inverse
374 relationship between physicochemical distance d and benefit \mathbf{B} .

375 DISCUSSION

376 The work presented here contributes to the field of phylogenetics and molecular
377 evolution in a number of ways. First, SelAC provides an complementary example to
378 Thorne et al. (2012) studies of how models of molecular and evolutionary scales can be

379 combined together in a nested manner. While the mapping between genotype and
380 phenotype is more abstract than Thorne et al. (2012), SelAC has the advantage of not
381 requiring knowledge of a protein's native folding. Second, our use of model nesting also
382 allows us to formulate and test specific biological hypotheses. For example, we are able to
383 compare a model formulation which assumes that physiochemical deviations from the
384 optimal sequence are equally disruptive at all sites within a protein to one which assumes
385 the effect of deviation from the optimal amino acid's physicochemical properties on protein
386 function varies between sites. By linking the strength of stabilizing selection for an optimal
387 amino acid sequence to gene expression, we can weight the historical information encoded
388 in genes evolving at vastly different rates in a biologically plausible manner while
389 simultaneously estimating their expression levels. Finally, because our fitness functions are
390 well defined, we can provide estimates of key evolutionary statistics such as the distribution
391 of effects on fitness and genetic load.

392 As phylogenetic methods become ever more ubiquitous in biology, and data set size
393 and complexity increase, there is a need and an opportunity for more complex and realistic
394 models (Goldman et al. 1996; Thorne et al. 1996; Goldman et al. 1998; Halpern and Bruno
395 1998; Lartillot and Philippe 2004). Despite their widespread use, phylogenetic models
396 based on purifying and diversifying selection, i.e. Goldman and Yang (1994) and
397 extensions, are very narrow categories of selection that mostly apply to cases of positive
398 and negative frequency dependent selection at the level of a particular amino acid, not for
399 tree inference itself.

400 Instead of heuristically extending population genetic models of neutral evolution for
401 use in phylogenetics, it makes sense to derive these extensions from population genetic
402 models that *explicitly* include the fundamental forces of mutation, drift, and natural
403 selection. Starting with Halpern and Bruno (1998), a number of researchers have developed
404 methods for linking site-specific selection on protein sequence and phylogenetics(e.g. Koshi

405 et al. 1999; Dimmic et al. 2000; Koshi and Goldstein 2000; Robinson et al. 2003; Lartillot
406 and Philippe 2004; Thorne et al. 2012; Rodrigue and Lartillot 2014). Our work follows this
407 tradition, but includes some key advances. For instance, even though SelAC requires a
408 large number of matrices, because of our assumption about protein functionality and
409 physicochemical distance from the optimum, we are able to parameterize our substitution
410 matrices using a relatively small number of genome-wide parameters and one gene specific
411 parameter. We show that all of these parameters can be estimated simultaneously with
412 branch lengths from the data at the tips of the tree.

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

429 Of course, given the general nature of SelAC and the complexity of biological
430 systems, other biological forces besides selection for reducing energy flux likely contribute

431 intergenic variation in the magnitude of stabilizing selection. Similarly, other
432 physicochemical properties besides composition, volume, and charge likely contribute to
433 site specific patterns of amino acid substitution. Thus, a larger and more informative set of
434 Grantham weights might improve our model fit and reduce the noise in our estimates of ϕ .
435 Even if other physicochemical properties are considered, the idea of a consistent, genome
436 wide Grantham weighting of these terms seems highly unlikely. Since the importance of an
437 amino acid's physicochemical properties likely changes with where it lies in a folded
438 protein, one way to incorporate such effects is to test whether the data supports multiple
439 sets of Grantham weights for either subsets of genes or regions within genes, rather than a
440 single set.

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

448 From a computational standpoint, the additive nature of selection between sites is
449 desirable because it allows us to analyze sites within a gene largely independently of each
450 other. From a biological standpoint, this additivity between site ignores any non-linear
451 interactions between sites, such as epistasis, or between alleles, such as dominance. Thus,
452 our work can be considered a first step to modeling to these more complex scenarios. For
453 example, our current implementation ignores any selection on synonymous codon usage bias
454 (CUB). Including such selection is tricky because introducing the site specific cost effects of
455 CUB leads to non-additive (i.e. epistatic) interactions between sites. Relative to stabilizing
456 selection on amino acid sequence, selection on CUB is thought to be substantially weaker.

457 As a result, epistatic effects due to synonymous codon specific differences in assembly costs
458 can likely ignored and selection on CUB incorporated into our current framework.

459 There are still significant deficiencies in the approach outlined here. Most worrisome
460 are biological flaws in the model. For example, at its heart, the model assumes that
461 suboptimal proteins can be compensated for, at a cost, simply by producing more of them.
462 However, this is likely only true for proteins reasonably close to the optimal sequence.
463 Different enough proteins will fail to function entirely: the active site will not sufficiently
464 match its substrates, a protein will not properly pass through a membrane, and so forth.
465 Yet, in our model, even random sequences still permit survival, just requiring more protein
466 production.

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

479 Overall, SelAC represents an important step in uniting phylogenetic and population
480 genetic models. It allows biologically relevant population genetic parameters to be
481 estimated from phylogenetic information, while also dramatically improving fit and
482 accuracy of phylogenetic models. Moreover, it demonstrates that there remains

483 substantially more information in the coding sequences used for phylogenetic analysis than
484 other methods acknowledge.

485 ACKNOWLEDGEMENTS

486 This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ) and
487 DEB-1355033 (BO, MAG, and RZ) with additional support from The University of
488 Tennessee Knoxville and University of Arkansas (JB). JB was supported, in part, as a
489 Postdoctoral Fellow at the National Institute for Mathematical and Biological Synthesis,
490 an Institute sponsored by the National Science Foundation through NSF Award
491 DBI-1300426, with additional support from UTK.

*

492

493

REFERENCES

- 494 Beaulieu, J. M., B. C. O'Meara, and M. J. Donoghue. 2013. Identifying Hidden Rate
495 Changes in the Evolution of a Binary Morphological Character: The Evolution of Plant
496 Habit in Campanulid Angiosperms. *Systematic Biology* 62:725–737.
- 497 Beaulieu, J. M. and B. C. O'Meara. 2016. Detecting Hidden Diversification Shifts in Models
498 of Trait-Dependent Speciation and Extinction. *Systematic Biology* 65:583–601.
- 499 Berg, J. and M. Lässig. 2003. Stochastic Evolution and Transcription Factor Binding Sites.
500 *Biophysics* 48:S36–S44.
- 501 Dimmic, M. W., D. P. Mindell, and R. A. Goldstein. 2000. Modeling evolution at the
502 protein level using an adjustable amino acid fitness model. *Pacific Symposium on*
503 *Biocomputing* 5:18–29.
- 504 Edwards, A. 1984. *Likelihood*. Cambridge science classics Cambridge University Press.
- 505 Felsenstein, J. 1981. Evolutionary trees from DNA-sequences - a maximum-likelihood
506 approach. *Journal of Molecular Evolution* 17:368–376.
- 507 Felsenstein, J. 2001. Taking Variation of Evolutionary Rates Between Sites into Account in
508 Inferring Phylogenies. *Journal of Molecular Evolution* 53:447–455.
- 509 Fisher, S., Ronald A. 1930. *The Genetical Theory of Natural Selection*. Oxford University
510 Press, Oxford.
- 511 Gilchrist, M. A. 2007. Combining Models of Protein Translation and Population Genetics
512 to Predict Protein Production Rates from Codon Usage Patterns. *Molecular Biology and*
513 *Evolution* 24:2362–2373.

- 514 Gilchrist, M. A., W.-C. Chen, P. Shah, C. L. Landerer, and R. Zaretzki. 2015. Estimating
515 Gene Expression and Codon-Specific Translational Efficiencies, Mutation Biases, and
516 Selection Coefficients from Genomic Data Alone. *Genome Biology and Evolution*
517 7:1559–1579.
- 518 Goldman, N. 1993. Statistical tests of models of DNA substitution. *Journal of molecular*
519 *evolution* 36:182–198.
- 520 Goldman, N., J. L. Thorne, and D. T. Jones. 1996. Using Evolutionary Trees in Protein
521 Secondary Structure Prediction and Other Comparative Sequence Analyses. *Journal of*
522 *Molecular Biology* 263:196 – 208.
- 523 Goldman, N., J. L. Thorne, and D. T. Jones. 1998. Assessing the Impact of Secondary
524 Structure and Solvent Accessibility on Protein Evolution. *Genetics* 149:445–458.
- 525 Goldman, N. and Z. H. Yang. 1994. Codon-based model of nucleotide substitution for
526 protein-coding DNA-sequences. *Molecular Biology and Evolution* 11:725–736.
- 527 Grantham, R. 1974. Amino acid difference formula to help explain protein evolution.
528 *Science* 185:862–864.
- 529 Halpern, A. L. and W. J. Bruno. 1998. Evolutionary distances for protein-coding sequences:
530 Modeling site-specific residue frequencies. *Molecular Biology And Evolution* 15:910–917.
- 531 Hughes, A. L. and M. Nei. 1988. Pattern of nucleotide substitution at major
532 histocompatibility complex class-i loci reveals overdominant selection. *Nature*
533 335:167–170.
- 534 Hhna, S., M. J. Landis, T. A. Heath, B. Boussau, N. Lartillot, B. R. Moore, J. P.
535 Huelsenbeck, and F. Ronquist. 2016. RevBayes: Bayesian Phylogenetic Inference Using

536 Graphical Models and an Interactive Model-Specification Language. *Systematic Biology*
537 65:726.

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

540 Johnson, S. G. 2012. The NLOpt nonlinear-optimization package. Version 2.4.2 – Released
541 20 May 2014.

542 Kimura, M. 1962. on the probability of fixation of mutant genes in a population. *Genetics*
543 47:713–719.

544 Koshi, J. M. and R. A. Goldstein. 1997. Mutation matrices and physical-chemical
545 properties: Correlations and implications. *Proteins-Structure Function And Genetics*
546 27:336–344.

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

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

552 Lartillot, N. and H. Philippe. 2004. A Bayesian mixture model for across-site
553 heterogeneities in the amino-acid replacement process. *Molecular Biology And Evolution*
554 21:1095–1109.

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

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

- 559 Muse, S. V. and B. S. Gaut. 1994. A likelihood approach for comparing synonymous and
560 nonsynonymous nucleotide substitution rates, with application to the chloroplast
561 genome. *Molecular Biology and Evolution* 11:715–724.
- 562 Nowak, M. A. 2006. *Evolutionary Dynamics: Exploring the Equations of Life*. Belknap of
563 Harvard University Press, Cambridge, MA.
- 564 Robinson, D. M., D. T. Jones, H. Kishino, N. Goldman, and J. L. Thorne. 2003. Protein
565 evolution with dependence among codons due to tertiary structure. *Molecular Biology
566 And Evolution* 20:1692–1704.
- 567 Rodrigue, N. and N. Lartillot. 2014. Site-heterogeneous mutation-selection models within
568 the PhyloBayes-MPI package. *Bioinformatics* 30:1020–1021.
- 569 Rodrigue, N., N. Lartillot, D. Bryant, and H. Philippe. 2005. Site interdependence
570 attributed to tertiary structure in amino acid sequence evolution. *Gene* 347:207–217.
- 571 Rokas, A., B. L. Williams, N. King, and S. B. Carroll. 2003. Genome-scale approaches to
572 resolving incongruence in molecular phylogenies. *Nature* 425:798–804.
- 573 Rowan, T. 1990. *Functional Stability Analysis of Numerical Algorithms*. Ph.D. thesis
574 University of Texas, Austin.
- 575 Salichos, L. and A. Rokas. 2013. Inferring ancient divergences requires genes with strong
576 phylogenetic signals. *Nature* 497:327–331.
- 577 Sella, G. and A. E. Hirsh. 2005. The application of statistical physics to evolutionary
578 biology. *Proceedings of the National Academy of Sciences of the United States of
579 America* 102:9541–9546.

- 580 Shah, P. and M. A. Gilchrist. 2011. Explaining complex codon usage patterns with
581 selection for translational efficiency, mutation bias, and genetic drift. *Proceedings of the*
582 *National Academy of Sciences of the United States of America* 108:10231–10236.
- 583 Stamatakis, A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses
584 with thousands of taxa and mixed models. *Bioinformatics* 22:2688–2690.
- 585 Thorne, J. L., N. Goldman, and D. T. Jones. 1996. Combining protein evolution and
586 secondary structure. *Molecular Biology and Evolution* 13:666–673.
- 587 Thorne, J. L., N. Lartillot, N. Rodrigue, and S. C. Choi. 2012. Codon models as a vehicle
588 for reconciling population genetics with inter-specific sequence data. *Codon Evolution:*
589 *Mechanisms And Models Pages 97–110 D2* 10.1093/acprof:osobl/9780199601165.001.0001
590 ER.
- 591 Wright, S. 1969. *Evolution and the genetics of populations. Vol. 2. The theory of gene*
592 *frequencies. vol. 2. University of Chicago Press.*
- 593 Yang, Z. 2014. *Molecular Evolution: A Statistical Approach. Oxford University Press, New*
594 *York.*
- 595 Yang, Z. H. 1994. Maximum-likelihood phylogenetic estimation from DNA-sequences with
596 variable rates over sites - approximate methods. *Journal Of Molecular Evolution*
597 39:306–314.
- 598 Yang, Z. H. 2007. PAML 4: Phylogenetic analysis by maximum likelihood. *Molecular*
599 *Biology And Evolution* 24:1586–1591.
- 600 Yang, Z. H. and R. Nielsen. 2008. Mutation-selection models of codon substitution and
601 their use to estimate selective strengths on codon usage. *Molecular Biology and*
602 *Evolution* 25:568–579.

TABLE

Model	logLik	Parameters Estimated	AIC	Δ AIC	Model Weight
GTR+ Γ	-655166.4	610	1,311,553	504,151	<0.001
GY94	-612121.5	210	1,224,663	417,261	<0.001
FMutSel0	-598848.2	2810	1,203,316	395,914	<0.001
SelAC	-465616.7	50,004	831,226	23,824	<0.001
SelAC+ Γ	-453706.0	50,005	807,402	0	0.999

Table 1: Comparison of model fits using Δ AIC.

FIGURES

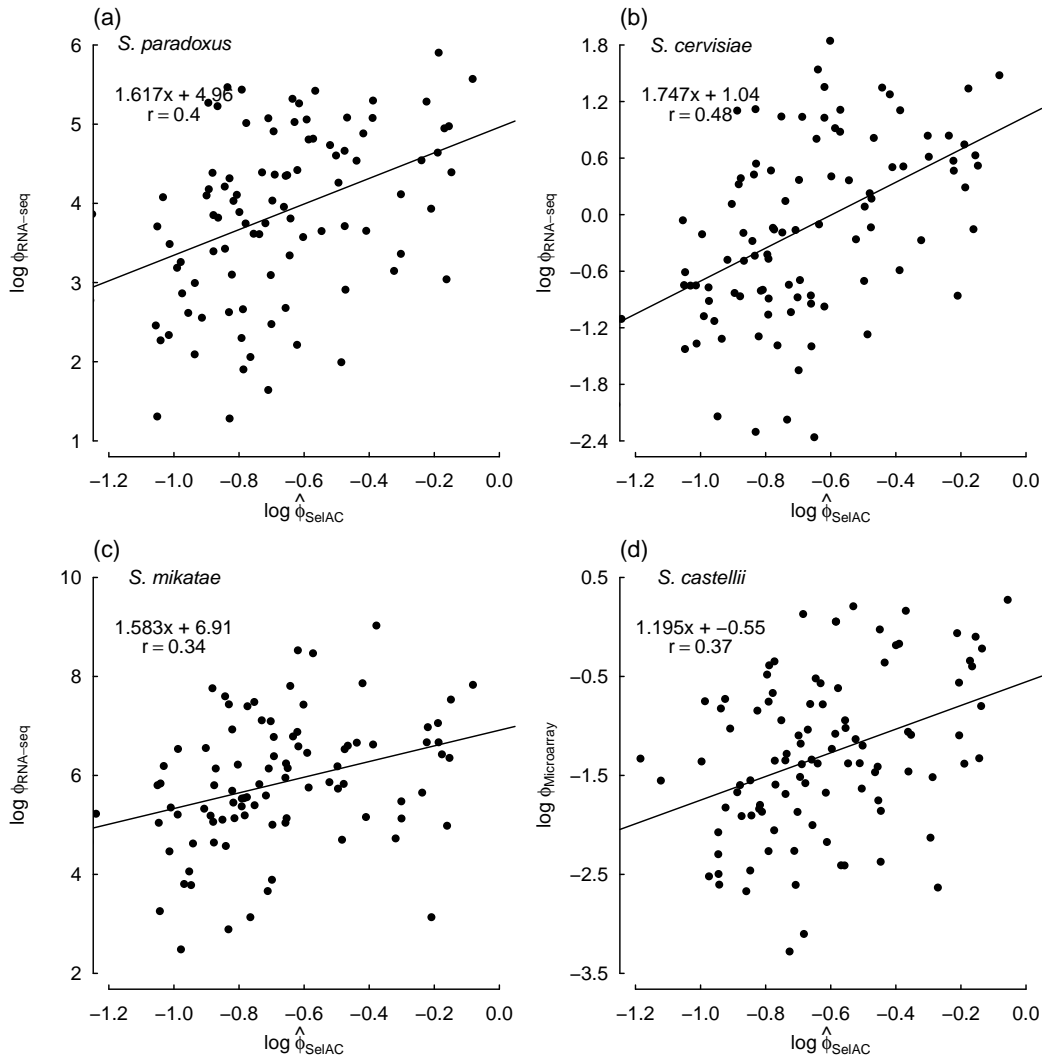


Figure 1: Comparisons between estimates of ϕ obtained from SelAC+ Γ and direct measurements of expression for individual yeast taxa across the 100 selected genes from Salichos and Rokas (2013). Estimates of ϕ were obtained by solving for ψ based on estimates of ψ' , and then dividing by $\mathbf{B}(\vec{a}_i | \vec{a}_*)$. 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 .

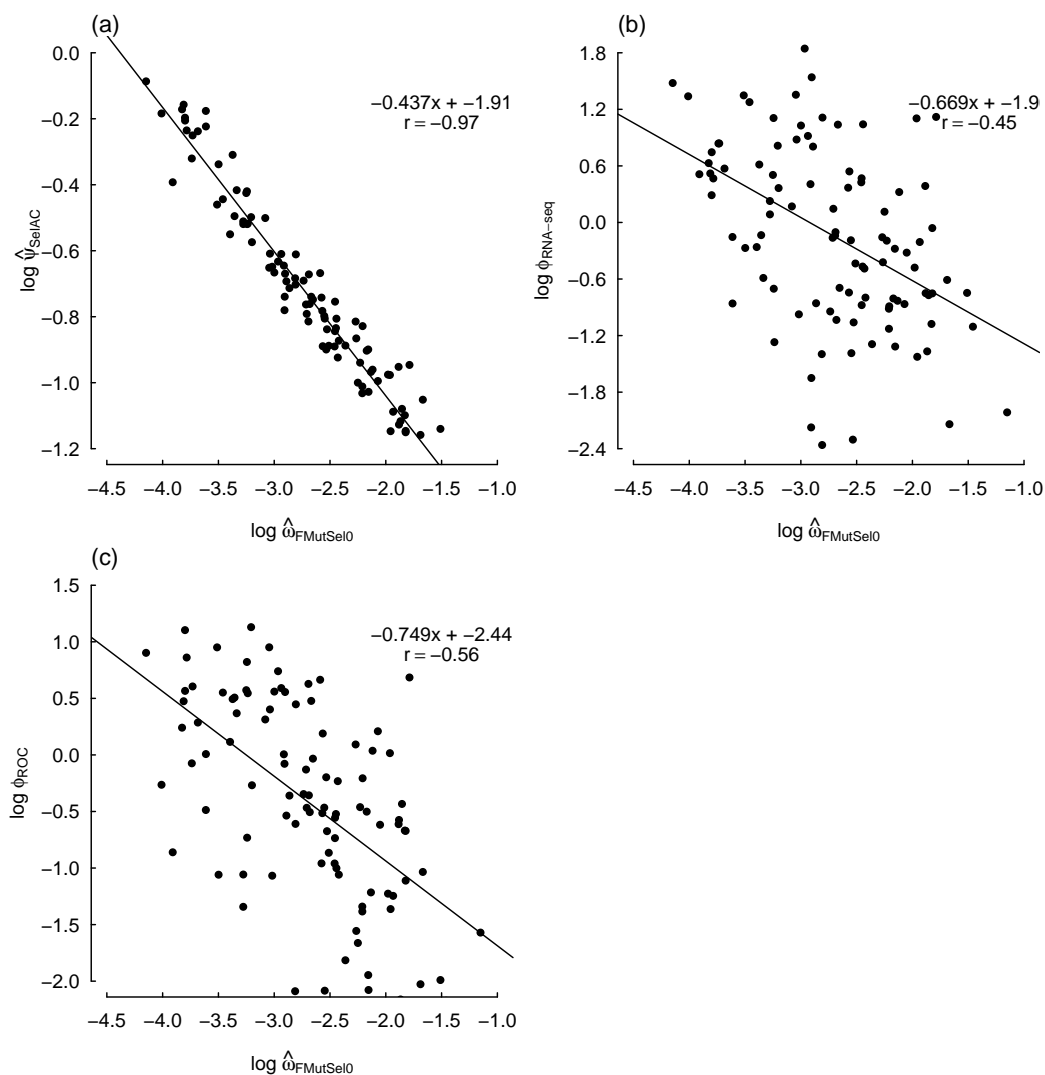


Figure 2: Comparisons between ω , which is the nonsynonymous/synonymous mutation ratio in FMutSel0, ψ obtained from SelAC+ Γ (a), a direct measurement of expression (b), and a model-based prediction of gene expression that does not account for ancestry (c), for *S. cerevisiae* across the 100 selected 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 ψ were solved from estimates of ψ' .

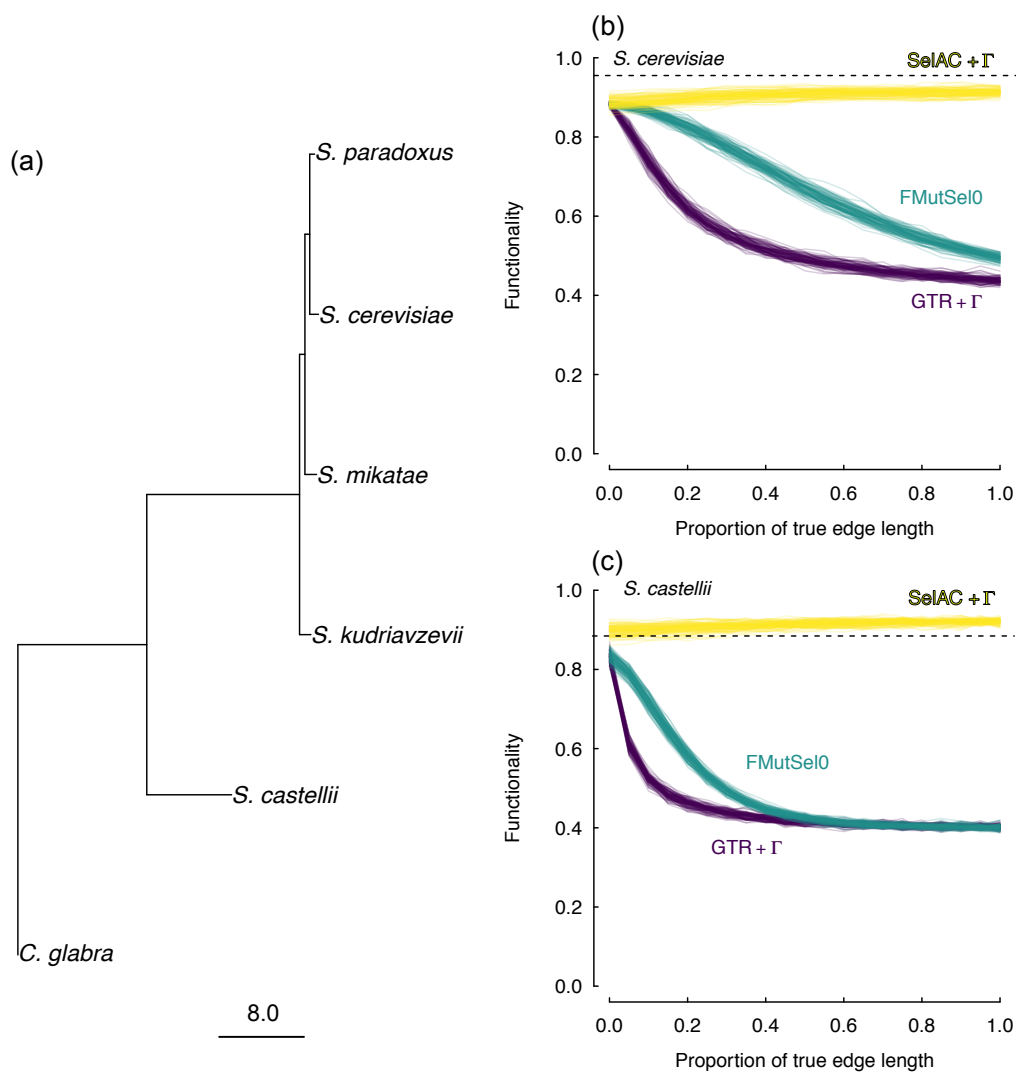


Figure 3: (a) Maximum likelihood estimates of branch lengths under SelAC+ Γ for 100 selected genes from Salichos and Rokas (2013). Tests of model adequacy for *S. cerevisiae* (b) and *S. castellii* (c) indicated that, when these taxa are removed from the tree, and their sequences are simulated, the parameters 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 .

605 **Part I**

606 **Supporting Materials**

607 **SUPPORTING MATERIALS**

608 *Comparisons of SelAC gene expression estimates with empirical*
609 *measurements*

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

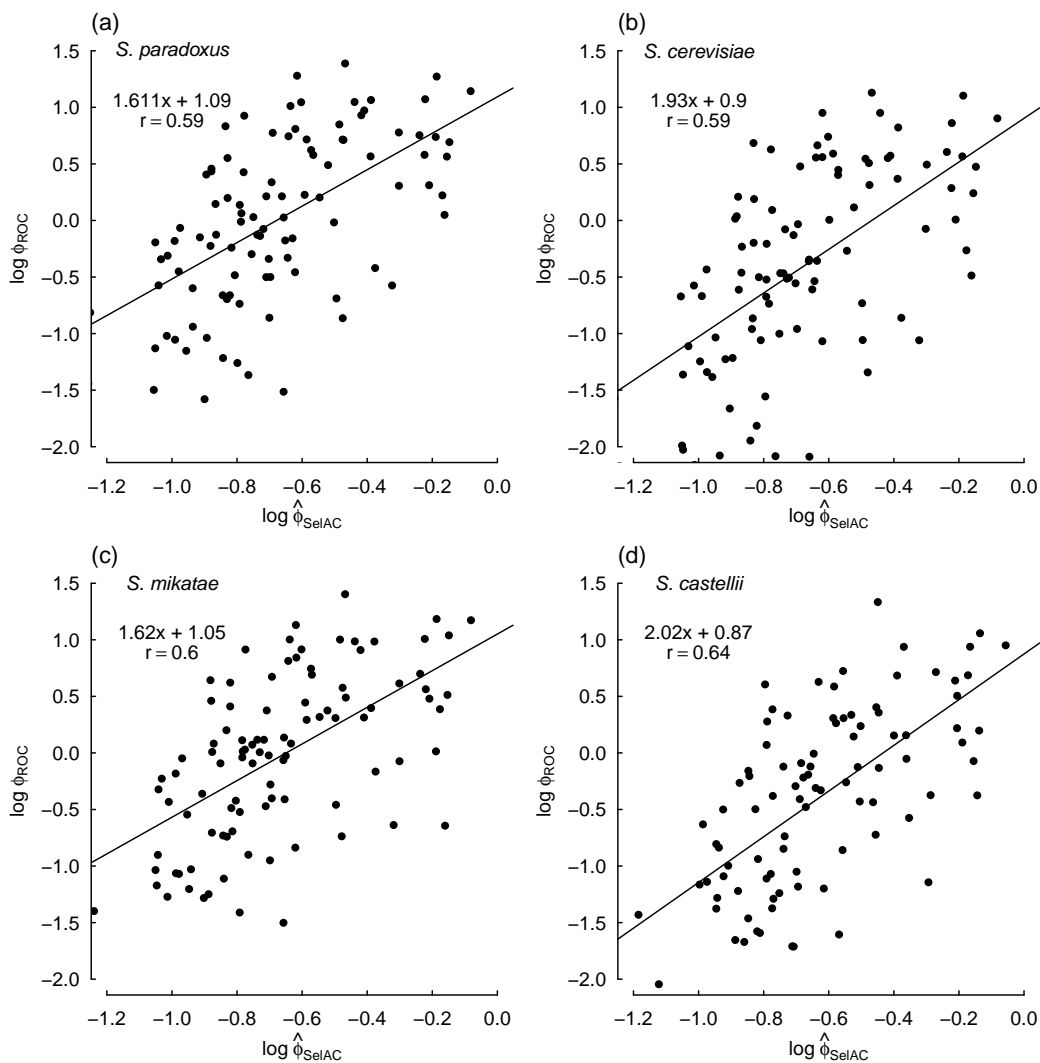


Figure S1: Comparisons between estimates of ϕ obtained from SelAC+ Γ and the predicted gene expression from the ROC SEMPER model (Gilchrist et al. (2015)) for individual yeast taxa across the 100 selected genes from Salichos and Rokas (2013). As with figures in the main text, estimates of ϕ were obtained by solving for ψ based on estimates of ψ' , 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.

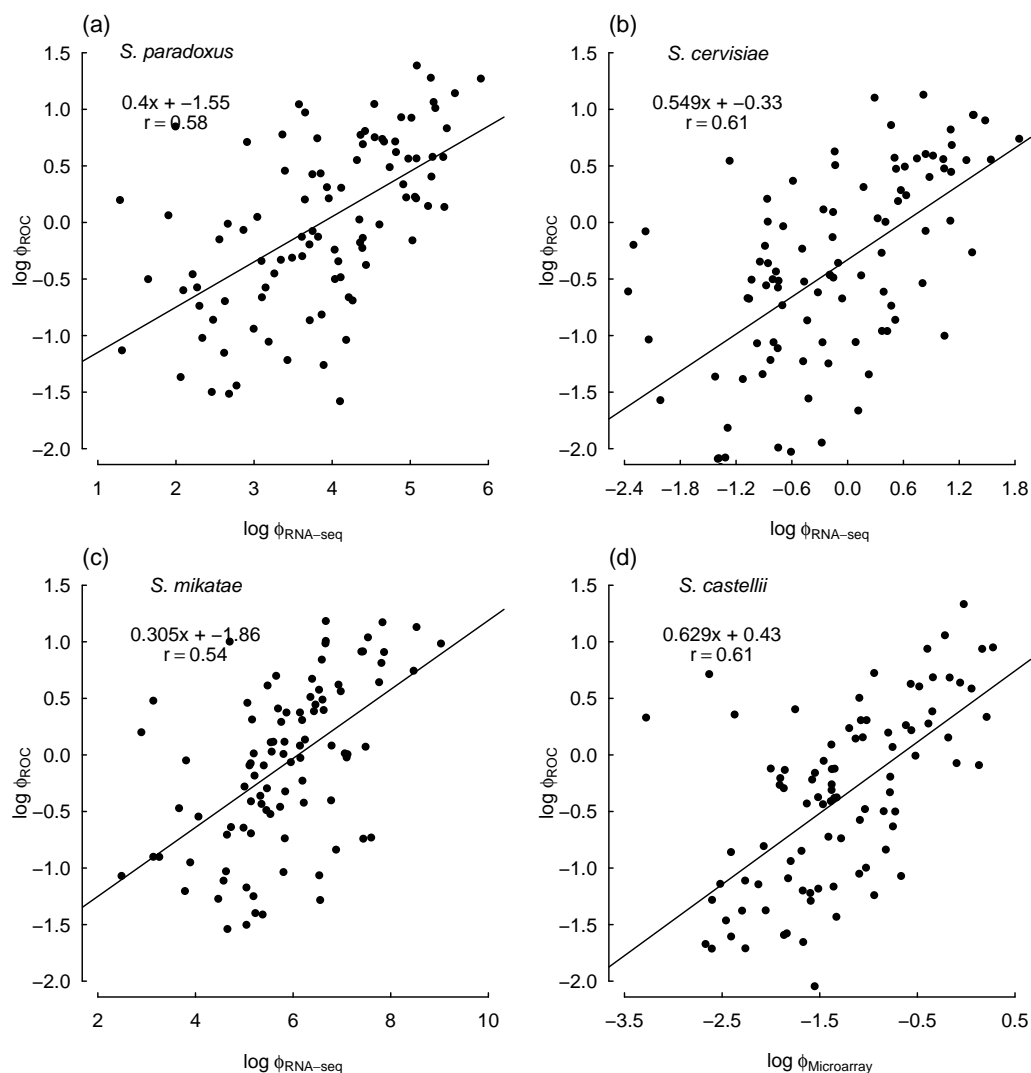


Figure S2: Comparisons of predicted gene expression from the ROC SEMPER model (Gilchrist et al. (2015)) and direct measurements of expression from RNA-Seq or Microarray data for individual yeast taxa across the 100 selected genes from Salichos and Rokas (2013). The equations in the upper left hand corner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

618

Simulations

619

Overall, the simulation results indicate that SelAC model can reasonably recover

620

the known values of the generating model (Figure S3 - S6). This includes not only the

621 parameters in the model, but also the optimal amino acids for a given sequence as well as
622 the estimates of the branch lengths. There are a few observations to note. First, the ability
623 to accurately recover the true optimal amino acid sequence will largely depend on the
624 magnitude of ϕ . This is, of course, intuitive, given that ϕ sets the strength of stabilizing
625 selection towards an optimal amino acid at a site. However, the inclusion of α_G into the
626 model, appears to generally increase values of ϕ and generally improves the ability to
627 recover the optimal amino acids even for the gene with the lowest baseline ϕ . Second, we
628 found a strong downward bias in estimates of α_G , which actually translates to greater
629 variation among the rate categories. The choice of a gamma distribution to represent
630 site-specific variation in sensitivity was based on mathematical convenience and
631 convention, rather than on biological reality. Nevertheless, we suspect that this bias is in
632 large part due to the difficulty in determining the baseline ψ for a given gene and the value
633 of α_G that globally satisfies the site-specific variation in sensitivity across all genes, as
634 indicated by the slight upward bias in estimates of ψ . It has been suggested, in studies of
635 the behavior of the the gamma distribution in applications of nucleotide substitution
636 model, that increasing the number of rate categories can often improve accuracy of the
637 shape parameter (Mayrose et al. (2005)). Future work will address this issue.

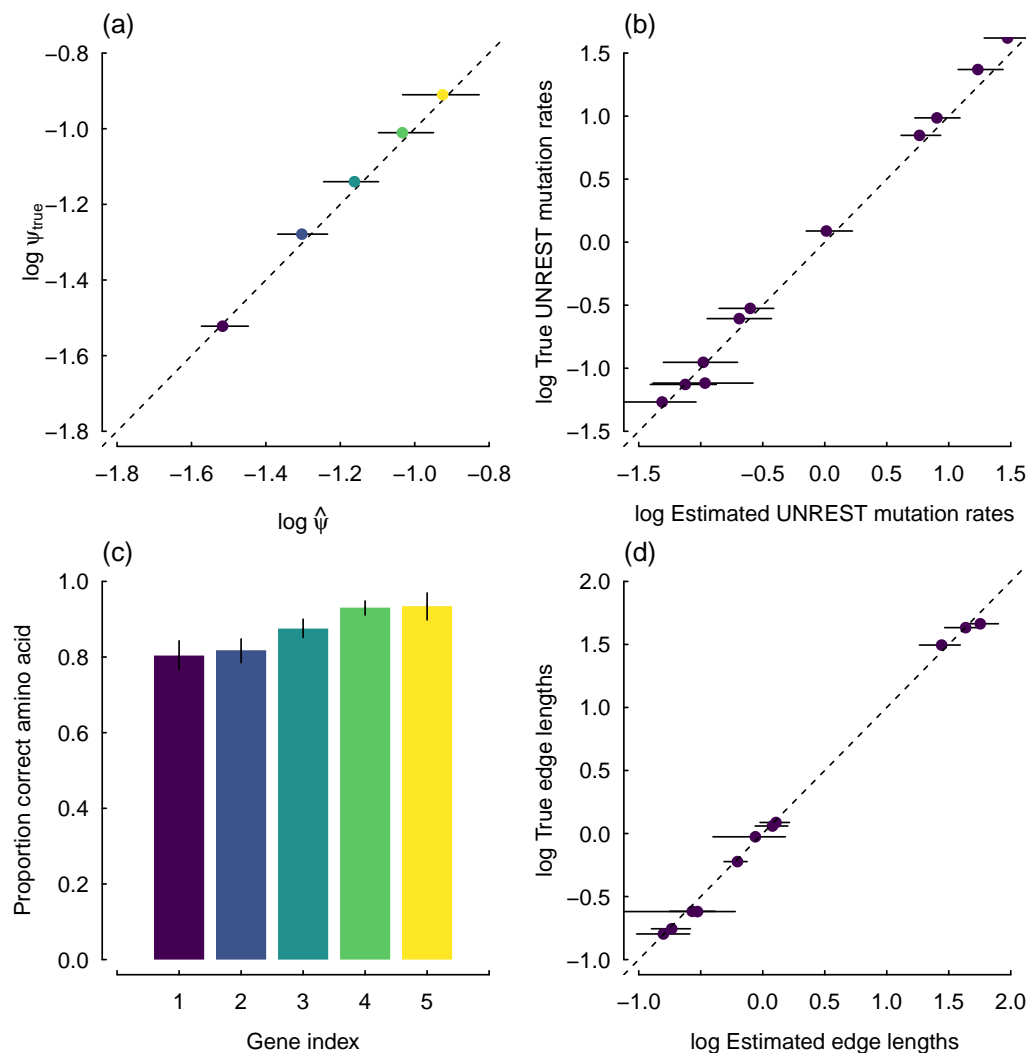


Figure S3: Summary a 5-gene simulation for a SelAC model where we assume $\alpha_G = \infty$, and thus, no site-specific sensitivity in the generating model. The 'known' parameters were based on fitting the same SelAC to the 106 gene data set and phylogeny of Rokas et al. (2003), with gene choice being based on five evenly spaced points along the rank order of the gene specific composite parameter ψ'_g . The points and associated uncertainty in the estimates of the gene-specific average protein synthesis rate, or ψ (calculated from ψ')(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.

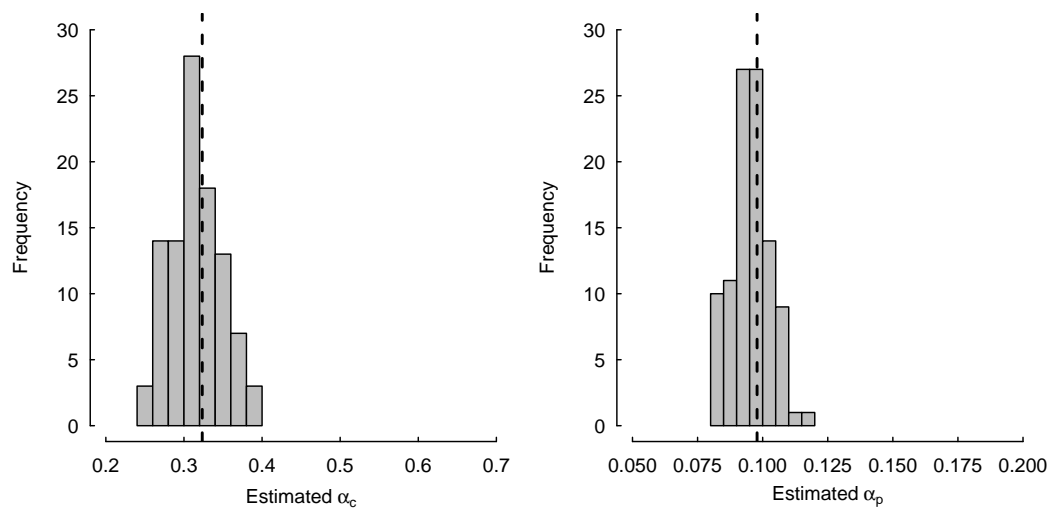


Figure S4: The distribution of estimates of the Grantham weights, α_c and α_p , in a SelAC model, where we assume $\alpha_G = \infty$, and thus no site-specific sensitivity in the generating model. The dashed line represents the value used in the generating model.

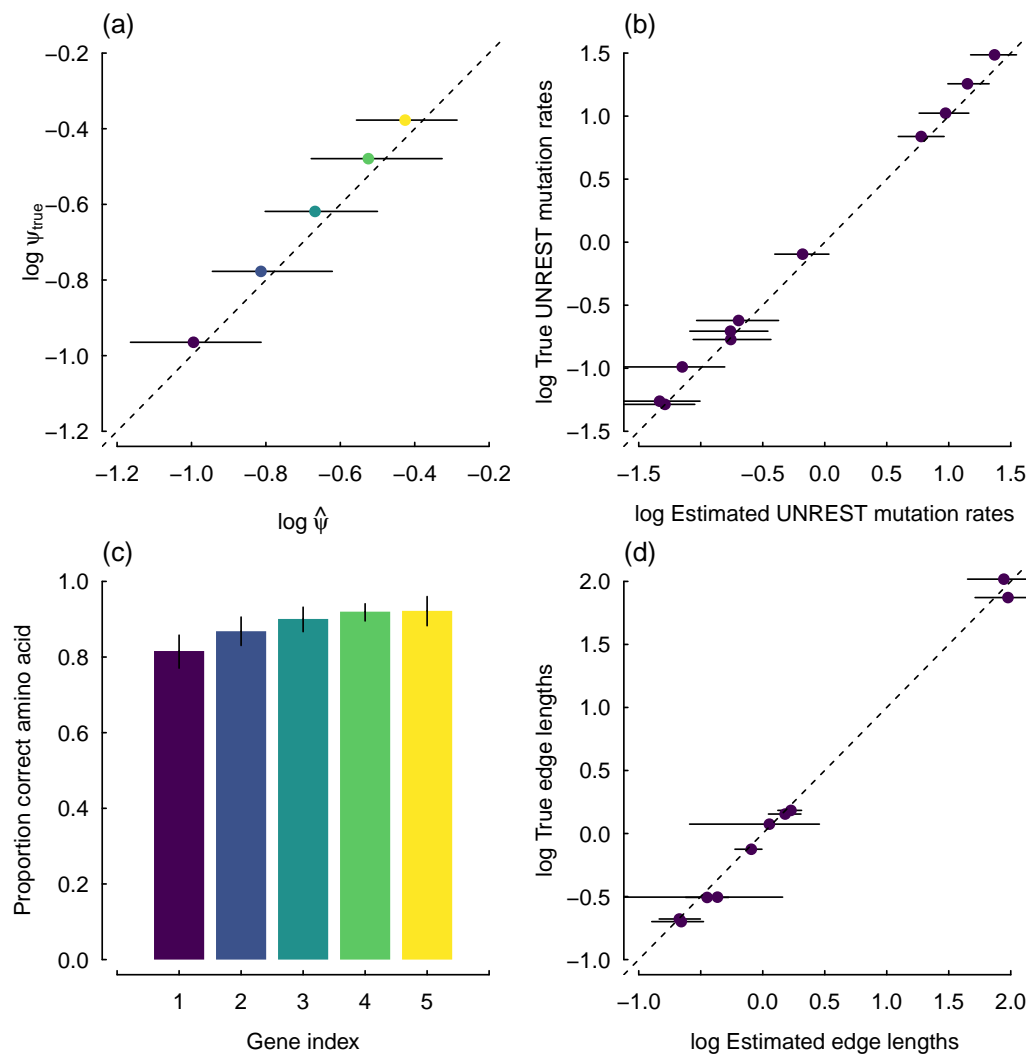


Figure S5: Same figure as in Figure S3, except the generating model includes site-specific sensitivity in the generating model (i.e., α_G).

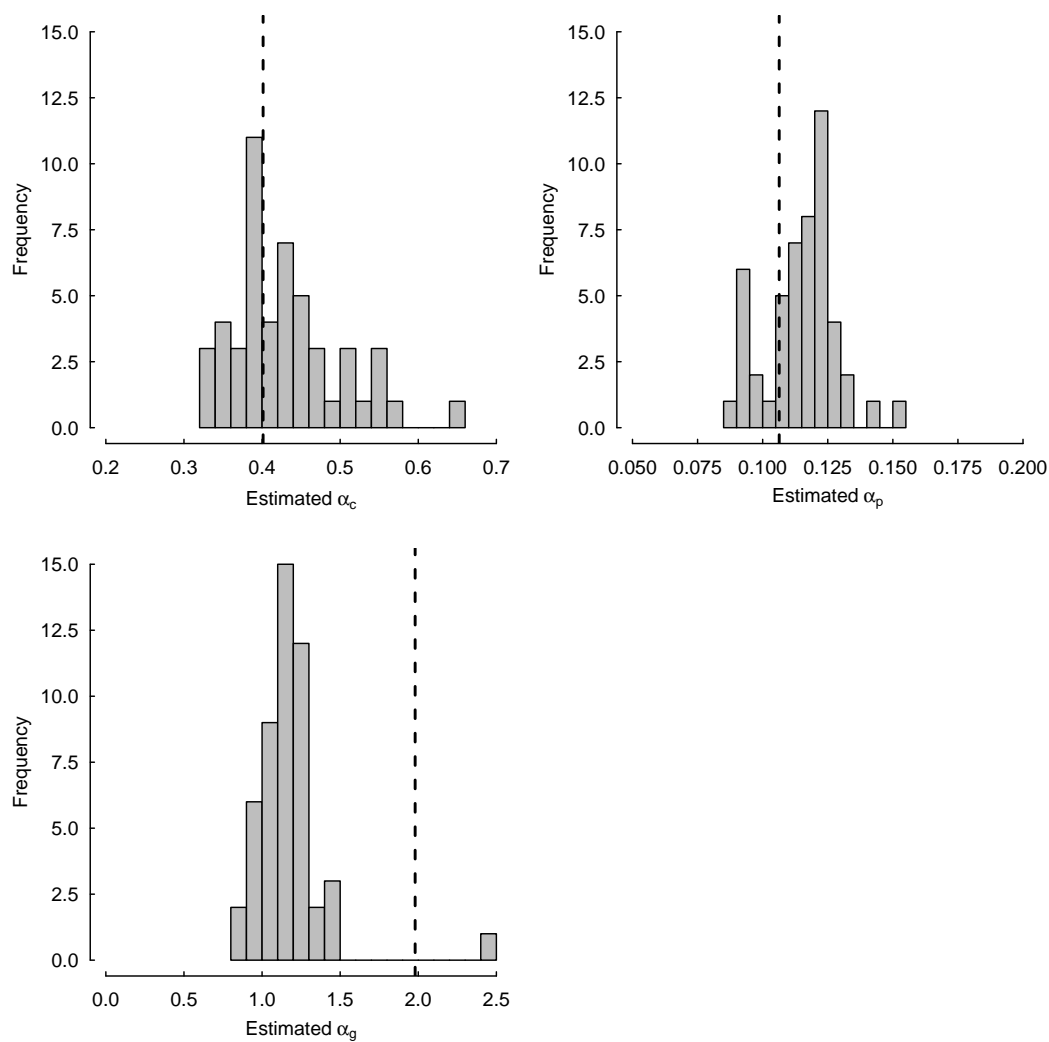


Figure S6: Same figure as in Figure S4, except the generating model includes site-specific sensitivity in the generating model (i.e., α_G). Unlike, Grantham weights, which showed no systematic bias, there is a downward bias in estimates of α_G .