

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 Zaretzki⁴, Cedric Landerer^{2,3}, Juanjuan Chai^{3,5}, and Michael A. Gilchrist^{2,3,*}

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

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

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

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

⁵Current 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 indicate there is great potential for more accurate inference of phylogenetic trees and branch lengths from already existing data through the use of nested, mechanistic models. Additional parameters estimated by SelAC indicate that a large amount of non-phylogenetic, but biologically meaningful, information can be inferred from existing 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

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

34 have increased substantially over the past two decades, the realism of the mathematical models on which
35 these analyses are based has changed relatively little by comparison. The most popular models of DNA
36 substitution used in molecular phylogenetics are simple nucleotide models that date back the early 1980's
37 and 90's, e.g. F81, F84, HYK85, TN93, and GTR (see Yang (2014) for an overview), and are indifferent
38 to the type of sequences they are fitted to. For example, when evaluating protein-coding sequences these
39 models are inherently agnostic with regards to the different amino acid substitutions and their impact
40 on gene function and, as a result, cannot describe the behavior of natural selection at the amino acid or
41 protein level.

42 Two important and independent attempts to address this critical shortcoming were introduced by
43 Goldman and Yang (1994, commonly abbreviated as GY94) and Muse and Gaut (1994). These models
44 were explicitly built for protein coding data, assuming that differences in the physicochemical properties
45 between amino acids, or physicochemical distances for short, could affect substitution rates. These
46 physicochemical based codon models as originally introduced have rarely been used for empirical data.
47 Instead, these often cited models have served as the basis for an array of simpler and, in turn, more popular
48 ω models that, starting with Nielsen and Yang (1998); Yang and Nielsen (1998), typically assume an
49 equal fixation probability for *all* non-synonymous mutations. Although often attributed to GY94, these
50 later and simpler models were the first to employ the single term ω to model the differences in fixation
51 probability between nonsynonymous and synonymous changes at all sites. Since their introduction,
52 more complex models have been developed that allow ω to vary between sites or branches (as cited in
53 Anisimova, 2012) and include selection on different synonyms for the same amino acid (e.g. Yang and
54 Nielsen, 2008)

55 In Goldman and Yang (1994); Nielsen and Yang (1998); Yang and Nielsen (1998) and later studies
56 based on their work, ω is suggested to indicate whether a given site within a protein sequence is under
57 consistent 'stabilizing' ($\omega < 1$) or 'diversifying' ($\omega > 1$) selection. Contrary to popular belief, ω does not
58 describe whether a site is evolving under a constant regime of stabilizing or diversifying selection, but
59 instead how a very particular *selective environment* changes over time. Below we explain how the actual
60 behavior of these models is inconsistent with how 'stabilizing' and 'diversifying' selection are otherwise
61 defined and understood (e.g. see Pellmyr, 2002).

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

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

72 Similarly, when $\omega > 1$, synonymous substitutions have a lower substitution rate than any possible non-
73 synonymous substitutions from the resident amino acid. Again due to the model's symmetrical nature,
74 the selection *against* the resident amino acid i is equally strong relative to alternative amino acids. The
75 selection against the resident amino acid i persists until a substitution occurs at which point selection now
76 *favors* amino acid i , as well as the 19 other amino acids, to the same degree i was previously disfavored.
77 Given this behavior, ω based models are likely to only reasonably approximate a subset of scenarios
78 such as perfectly symmetrical over-/under-dominance or positive/negative frequency dependent selection
79 (Hughes and Nei, 1988; Nowak, 2006). Further, ω based models implicitly assumes the substitution is on
80 the same timescale as the shifts in the optimal (or pessimal) amino acid.

81 **New Approaches**

82 To address these fundamental shortcomings in ω based phylogenetic approaches, we present an
83 approach where selection explicitly favors minimizing the cost-benefit function η of a protein whose
84 relative performance is determined by the order and physicochemical properties of its amino acids. Our
85 approach, which we call Selection on Amino acids and Codons or SelAC, is developed in the same vein
86 as previous phylogenetic applications of the Wright-Fisher process (e.g. Dimmic et al., 2000; Halpern
87 and Bruno, 1998; Koshi and Goldstein, 1997; Koshi et al., 1999; Lartillot and Philippe, 2004; Muse and
88 Gaut, 1994; Rodrigue and Lartillot, 2014; Rodrigue et al., 2005; Thorne et al., 2012; Yang and Nielsen,
89 2008). Similar to Lartillot and Philippe (2004) and Rodrigue and Lartillot (2014), we assume there is a
90 finite set of rate matrices describing the substitution process and that each position within a protein is
91 assigned to a particular rate matrix category. Unlike that work, we assume *a priori* there are 20 different
92 families of rate matrices, one family for when a given amino acid is favored at a site. The key parameters
93 underlying these matrices are shared across genes except for gene expression. As a result, SelAC identifies
94 the amino acid at a particular position within a protein that is favored by natural selection using a simple
95 cost-benefit approach.

96 While natural selection on protein coding regions can take many forms, one general approach to
97 describing its effects is by relating a codon sequence to the cost of producing the encoded protein and
98 the functional benefit (or potential harm) from translating its sequence. The gene specific cost of protein
99 synthesis can be affected by the amino acids used, the direct and indirect costs of peptide assembly by
100 the ribosome, and the use of chaperones to aid in folding. Importantly, these costs can be computed to
101 varying degrees of realism (e.g. Lynch and Marinov, 2015; Wagner, 2005). We have previously presented
102 models of protein synthesis costs that, alternatively, take into account the cost of ribosome pausing (Shah
103 and Gilchrist, 2011) or premature termination errors (Gilchrist et al., 2009; Gilchrist, 2007; Gilchrist and
104 Wagner, 2006).

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

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

121 **Results**

122 By linking transition rates $q_{i,j}$ to gene expression in the form of protein synthesis rate ϕ , our approach
123 allows use of the same model for genes under varying degrees of stabilizing selection. Specifically, we
124 assume the strength of stabilizing selection for the optimal sequence, \bar{a}^* , is proportional to the average
125 protein synthesis rate ϕ , which we can estimate for each gene. In regards to model fit, our results clearly
126 indicated that linking the strength of stabilizing selection for the optimal sequence to gene expression

127 substantially improves our model fit. Further, including the shape parameter α_G for the random effects
128 term $G \sim \text{Gamma}(\text{shape}=\alpha_G, \text{rate}=\alpha_G)$ to allow for heterogeneity in this selection between sites within
129 a gene improves the ΔAICc of SelAC+ Γ over the simpler SelAC models by over 22,000 AIC units. Using
130 either ΔAICc or AIC_w as our measure of model support, the SelAC models fit extraordinarily better than
131 GTR + Γ , GY94, or FMutSel (Table 1). This is in spite of the need for estimating the optimal amino
132 acid at each position in each protein, which accounts for 49,881 additional model parameters. Even when
133 compared to the next most parameter rich codon model in our model set, FMutSel, SelAC+ Γ model
134 shows over 160,000 AIC unit improvement over FMutSel.

135 The analysis building upon Jhwueng et al. (2014) suggests that using the number of taxa times the
136 number of sites as the sample size performs best as a small sample size correction for estimating Kullback-
137 Liebler distance in phylogenetic models (Appendix 1). This also has an intuitive appeal. In models that
138 have at least some parameters shared across sites and some parameters shared across taxa, increasing
139 the number of sites and/or taxa should be adding more samples for the parameters to estimate. This
140 is consistent considering how likelihood is calculated for phylogenetic models: the likelihood for a given
141 site is the sum of the probabilities of each observed state at each tip, which is then multiplied across
142 sites. It is arguable that the conventional approach in comparative methods is calculating AICc in the
143 same way. That is, if only one column of data (or “site”) is examined, as remains remarkably common
144 in comparative methods, when we refer to sample size, it is technically the number of taxa multiplied by
145 number of sites, even though it is referred to simply as the number of taxa.

146 With respect to estimates of ϕ within SelAC, they were strongly correlated with both empirical
147 measurements (Pearson $r=0.33-0.48$) and theoretical predictions (Pearson $r=0.45-0.64$) of gene
148 expression (Figure 1 and Figures S1-S2, respectively). In other words, using only codon sequences, our
149 model can predict which genes have high or low expression levels. The estimate of the α_G parameter,
150 which describes the site-specific variation in sensitivity of the protein’s functionality, indicated a moderate
151 level of variation in gene expression among sites. Our estimate of $\alpha_G = 1.36$, produced a distribution
152 of sensitivity terms G ranged from 0.342-7.32, but with more than 90% of the weight for a given site-
153 likelihood being contributed by the 0.342 and 1.50 rate categories. In simulation, however, of all the
154 parameters in the model, only α_G showed a consistent bias, in that the MLE were generally lower than
155 their actual values (see Supporting Materials). Other parameters in the model, such as the Grantham
156 weights, provide an indication as to the physicochemical distance between amino acids. Our estimates of
157 these weights only strongly deviate from Grantham’s 1974 original estimates in regards to composition
158 weight, α_c , which is the ratio of non-carbon atoms in the end groups or rings to the number of

159 carbon atoms in side chains. Our estimate of the composition weighting factor of $\alpha_c=0.459$ is 1/4th
160 the value estimate by Grantham which suggests that the substitution process is less sensitive to this
161 physicochemical property when shared ancestry and variation in stabilizing selection are taken into
162 account.

163 It is important to note that the nonsynonymous/synonymous mutation ratio, or ω , which we estimated
164 for each gene under the FMutSel model strongly correlated with our estimates of $\phi' = \psi' / \mathbf{B}$ where \mathbf{B}
165 depends on the sequence of each taxa. In fact, ω showed similar, though slightly reduced correlations,
166 with the same empirical estimates of gene expression described above (Figure 2) This would give the
167 impression that the same conclusions could have been gleaned using a much simpler model, both in terms
168 of the number of parameters and the assumptions made. However, as we discussed earlier, not only is
169 this model greatly restricted in terms of its biological feasibility, SelAC clearly performs better in terms
170 of its fit to the data and biological realism.

171 For example, when we simulated the sequence for *S. cerevisiae*, starting from the ancestral sequence
172 under both GTR + Γ and FMutSel, the functionality of the simulated sequence moves away from the
173 observed sequence, whereas SelAC remains near the functionality of the observed sequence (Figure 3b).
174 This is somewhat unsurprising, given that both GTR + Γ and FMutSel are agnostic to the functionality
175 of the gene, but it does highlight the improvement in biological realism in amino acid sequence evolution
176 that SelAC provides. We do note that the adequacy of the SelAC model does vary among individual
177 taxa, and does not always match the observed functionality. For instance, our simulations of *S. castellii*
178 gene function is consistently higher than estimated from the data (Figure 3c). We suspect this is an
179 indication that assuming a single set of optimal amino acid across all taxa is too simplistic. However, we
180 cannot rule out violations of SelAC's other model assumptions such as: a single set of Grantham weights,
181 a single α_G , or reductions in protein functionality \mathbf{B} being solely a function of physicochemical distances
182 d between sites.

183 Discussion

184 A central goal in evolutionary biology is to quantify the nature, strength, and, ultimately, shifts in the
185 forces of natural selection relative to genetic drift and mutation. As data set size and complexity increase,
186 so does the amount of potential information on these forces and their dynamics. As a result, there is a
187 need for more complex and realistic models to accomplish this goal (Goldman et al., 1996, 1998; Halpern
188 and Bruno, 1998; Lartillot and Philippe, 2004; Thorne et al., 1996). Although extremely popular due to
189 their elegance and computational efficiency, the utility of ω based models in helping us reach this goal is

190 substantially more limited than commonly recognized. Because these ω models use a single substitution
191 matrix, they are only applicable for situations in which the substitution process and shifts in the selective
192 environment are intrinsic to the sequence, such as with positive or negative frequency dependent selection;
193 these models do not describe stabilizing or diversifying selection as commonly envisioned (Endler, 1986;
194 Pelmyr, 2002).

195 Starting with Halpern and Bruno (1998), a number of researchers have developed methods for linking
196 site-specific selection on protein sequence and phylogenetics (e.g. Dimmic et al., 2000; Koshi and
197 Goldstein, 2000; Koshi et al., 1999; Lartillot and Philippe, 2004; Robinson et al., 2003; Rodrigue and
198 Lartillot, 2014; Thorne et al., 2012). Halpern and Bruno (1998) calculated a vector of 20 expected amino
199 acid frequencies for each amino acid site, making it the most general and most parameter rich of these
200 methods. This generality, however, comes at the cost of being purely descriptive; there is no explicit
201 biological mechanism proposed to explain the site specific amino acid frequencies estimated. By grouping
202 together amino sites with similar evolutionary behaviors, Lartillot and Philippe (2004) and Rodrigue and
203 Lartillot (2014) retained the descriptive nature of Halpern and Bruno (1998) work while greatly reduced
204 the number of model parameters needed.

205 SelAC follows in this tradition of using multiple substitution matrices, but includes some key advances.
206 First, by nesting a model of a sequence's cost-benefit function \mathbf{C}/\mathbf{B} within a broader model, SelAC allows
207 us to formulate and test a hierarchical, mechanistic models of stabilizing selection. More precisely, our
208 nested approach allows us to relax the assumption that physicochemical deviations from the optimal
209 sequence \vec{a}^* are equally disruptive at all sites within a protein. Indeed, SelAC strongly supports the
210 hypothesis that the strength of stabilizing selection against physicochemical deviations from \vec{a}^* varies
211 between sites ($\Delta\text{AICc} = 20,983$; Table1). Second, because our substitution matrices are built on a
212 formal description of a sequence's cost-benefit function \mathbf{C}/\mathbf{B} , we are able to efficiently parameterize 20
213 different matrices using a relatively small number of genome-wide parameters – e.g. our physicochemical
214 weightings, α_c , α_p , and α_v , and the shape parameter α_G for the distribution of selective strength G and
215 one gene specific expression parameter ψ . While the \mathbf{C}/\mathbf{B} function on which SelAC currently rests is
216 very simple, nevertheless, it leads to a dramatic increase in our ability to explain the sequence data
217 we analyzed. Importantly, because SelAC uses a formal description of a sequence's \mathbf{C}/\mathbf{B} , replacing our
218 assumptions with more sophisticated ones in the future is relatively straightforward. Third, our use of
219 nested models also allows us to make biologically meaningful and testable predictions. By linking a
220 gene's expression level to the strength of purifying selection it experiences, we are able to provide coarse

221 estimates of gene expression. This also suggests that ω is best explained as a proxy for gene expression,
222 rather than the nature of selection on a sequence.

223 Thus, we believe our cost-benefit approach to be a substantial advance of the more simplistic ω models,
224 is complementary to the work of others in the field (e.g. Rodrigue and Lartillot, 2014; Thorne et al., 2012),
225 and, in turn, lays the foundation for more realistic work in the future. For instance, by assuming there
226 is an optimal amino acid for each site, SelAC naturally leads to a non-symmetrical and, thus, more
227 cogent model of protein sequence evolution. Because the strength of selection depends on an additive
228 function of amino acid physicochemical properties, an amino acid more similar to the optimum has a
229 higher probability of replacing a more dissimilar amino acid than the converse situation. Further, SelAC
230 does not assume the system is always at the optimum or pessimum point of the fitness landscape, as
231 occurs when $\omega < 1$ or > 1 , respectively.

232 Importantly, the cost-benefit approach underlying SelAC allows us to link the strength of selection on a
233 protein sequence to its gene's expression level. Despite its well recognized importance in determining the
234 rate of protein evolution (e.g. Drummond et al., 2005, 2006), phylogenetic models have ignored the fact
235 that expression levels vary between genes. In order to link gene expression and the strength of stabilizing
236 selection on protein sequences, we simply assume that the strength of selection on a gene is proportional
237 to the average protein synthesis rate of the gene.

238 One possible mechanism with some theoretical and empirical support which generates a linear
239 relationship between the strength of selection and gene expression is the assumption of compensatory
240 gene expression (Allison, 2012; Allison and Goulden, 2017; Brown and Elliot, 1997; King et al., 2015;
241 Lerman et al., 2012; Thiele et al., 2012; Zanger and Schwab, 2013). That is, the assumption that any
242 reduction in protein function is compensated for by an increase in the protein's production rate and, in
243 turn, abundance. For example, a mutation which reduces the functionality of the protein to 90% of the
244 optimal protein, would require $1/0.9=1.11$ of these suboptimal proteins to be produced relative to the
245 optimal protein in order to maintain the same amount of that protein's functionality in the cell. Because
246 the energetic cost of an 11% increase in a protein's synthesis rate is proportional to its target synthesis
247 rate, our assumptions naturally link changes in protein functionality and changes in gene expression
248 and its associated costs. Under what circumstances cells actually respond in this manner, remains to be
249 determined. The fact that our method allows us to explain 13-23% of the variation in gene expression
250 measured using RNA-Seq, suggests that this assumption is a reasonable starting point.

251 Furthermore, by linking expression and selection, SelAC provides a natural framework for combining
252 information from protein coding genes with very different rates of evolution; from low expression genes

253 providing information on shallow branches to high expression genes providing information on deep
254 branches. This is in contrast to a more traditional approach of concatenating gene sequences together,
255 which is equivalent to assuming the same average functionality production rate ψ for all of the genes,
256 or more recent approaches where different models are fitted to different genes. Our results indicate that
257 including a gene specific ψ value vastly improves SelAC fits (Table 1). Perhaps more convincingly, we find
258 that the target functionality production rate ψ and the realized average protein synthesis rate $\phi = \psi/\mathbf{B}$ are
259 reasonably well correlated with laboratory measurements and theoretical predictions of gene expression
260 (Pearson $r = 0.34 - 0.64$; Figures 1, 1, and 2). The idea that quantitative information on gene expression
261 is embedded within intra-genomic patterns of synonymous codon usage is well accepted; our work shows
262 that this information can also be extracted from comparative data at the amino acid level.

263 Of course, given the general nature of SelAC and the complexity of biological systems, other biological
264 forces besides selection for reducing energy flux likely contribute to intergenic variation in the magnitude
265 of stabilizing selection. Similarly, other physicochemical properties besides composition, volume, and
266 charge likely contribute to site specific patterns of amino acid substitution. Thus, a larger and more
267 informative set of physicochemical weights might improve our model fit and reduce the noise in our
268 estimates of realized protein synthesis rates ϕ . Even if other physicochemical properties are considered,
269 the idea of a consistent, genome wide physicochemical weighting of these terms seems highly unlikely.
270 Since the importance of an amino acid's physicochemical properties likely changes with its position in a
271 folded protein, one way to incorporate such effects is to test whether the data supports multiple sets of
272 physicochemical weights for either subsets of genes or regions within genes, rather than a single set.

273 Both of these points highlight the advantage of the detailed, mechanistic modeling approach underlying
274 SelAC. Because there is a clear link between protein expression, synthesis cost, and functionality, SelAC
275 can be extended by increasing the realism of the mapping between these terms and the coding sequences
276 being analyzed. For example, SelAC currently assumes the optimal amino acid for any site is fixed along
277 all branches. This assumption can be relaxed by allowing the optimal amino acid to change during the
278 course of evolution along a branch. From a computational standpoint, the additive nature of selection
279 between sites is desirable because it allows us to analyze sites within a gene largely independently of
280 each other. From a biological standpoint, this additivity between sites ignores any non-linear interactions
281 between sites, such as epistasis, or between alleles, such as dominance. Thus, our work can be considered
282 a first step to modeling these more complex scenarios.

283 For example, our current implementation ignores any selection on synonymous codon usage bias (CUB)
284 (c.f. Pouyet et al., 2016; Yang and Nielsen, 2008). Including such selection is tricky because introducing

285 the site-specific cost effects of CUB, which is consistent with the hypothesis that codon usage affects
286 the efficiency of protein assembly or **C**, into a model where amino acids affect protein function or **B**,
287 results in a cost-benefit ratio **C/B** with epistatic interactions between all sites. These epistatic effects
288 can likely be ignored under certain conditions or reasonably approximated based on an expectation of
289 codon specific costs (e.g. Kubatko et al., 2016). Nevertheless, it is difficult to see how one could identify
290 such conditions without modeling the way in which codon and amino acid usage affects **C/B**.

291 This work also points out the potential importance of further investigation into model choice in
292 phylogenetics. For likelihood models, use of AICc has become standard. However, how one determines the
293 appropriate number of data points in a model is more complicated than generally recognized. Common
294 sense suggests that dataset size is increased by adding taxa and/or sites. In other words, a dataset of 1000
295 taxa and 100 sites must have more information on substitution models than a dataset of 4 taxa and 100
296 sites. Our simple analyses agree that the number of observations in a dataset (number of sites \times number
297 of taxa) should be taken as the sample size for AICc, but this conclusion likely only applies when there
298 is sufficient independence between taxa. For instance, one could imagine a phylogeny where one taxon is
299 sister to a polytomy of 99 taxa that have zero length terminal branches. Absent measurement error or
300 other intraspecific variation, one would have 100 species but only two unique trait values, and the only
301 information about the process of evolution comes from what happens on the path connecting the lone
302 taxon to the polytomy. Although this is a rather extreme example, it seems prudent for researchers to
303 use a simulation based approach similar to the one we take here to determine the appropriate means for
304 calculating the effective number of data points in their data.

305 There are still significant shortcomings in the approach outlined here. Most worrisome are biological
306 oversimplifications in SelAC. For example, at its heart, SelAC assumes that suboptimal proteins can
307 be compensated for, at a cost, simply by producing more of them. However, this is likely only true
308 for proteins reasonably close to the optimal sequence. Different enough proteins will fail to function
309 entirely: the active site will not sufficiently match its substrates, a protein will not properly pass through
310 a membrane, and so forth. Yet, in our model, even random sequences still permit survival, just requiring
311 more protein production. Like the other oversimplifications previously discussed, these assumptions can be
312 relaxed through further extension of our model.

313 There are also deficiencies in our implementation. Though reasonable to use for a given topology with
314 a modest number of species, it is currently too slow for practical use for tree search. Our work serves
315 as a proof of concept, or of utility for targeted questions where a more realistic model may be of use
316 (placement of particular taxa, for example). Future work will encode SelAC models into a variety of

317 mature, popular tree-search programs. SelAC also represents a challenging optimization problem: the
318 nested models reduce parameter complexity vastly, but there are still numerous parameters to optimize,
319 including the discrete parameter of the optimal amino acid at each site. One way to avoid the use of
320 discrete parameters at the expense of more of them would be to have SelAC estimate the optimum
321 physicochemical values on a per site basis rather than a specific amino acid. While this would increase
322 the number of parameters estimated, it would have the practical advantage of continuous parameter
323 optimization rather than discrete, and biologically would be more realistic (as it is the properties that
324 selection “sees”, not the identity of the amino acid itself).

325 In spite of these difficulties, SelAC represents an important step in uniting phylogenetic and population
326 genetic models. For example, while Dimmic et al. (2000); Koshi and Goldstein (2000); Koshi et al. (1999);
327 Lartillot and Philippe (2004); Robinson et al. (2003); Rodrigue and Lartillot (2014); Thorne et al. (2012)
328 are all models of constant, stabilizing selection, SelAC can be generalized further to include diversifying
329 selection. Specifically, by letting SelAC’s sensitivity term G , which we now assume is ≥ 0 , to take on
330 negative values, SelAC will behave as if there is a pessimal, rather than optimal, amino acid for the given
331 site. In this diversifying selection scenario, amino acids with physicochemical qualities more dissimilar to
332 the pessimal amino acid are increasingly favored, potentially resulting in multiple fitness peaks.

333 Because SelAC infers the optimal amino acid for each site, it is substantially more parameter rich than
334 more commonly used models such as GTR+ Γ , GY94, and FMutSel. Despite this increase in number of
335 model parameters, SelAC drastically outperforms these models with AICc values on the order of 10,000s
336 to 100,000s. We predict that SelAC’s performance could be improved even further if we use a hierarchical
337 approach where the optimal amino acid is not estimated on a per site basis, but rather as a vector of
338 probability an amino acid is optimal at the gene level.

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

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

358 **Materials & Methods**

359 **Overview**

360 We model the substitution process as a classic Wright-Fisher process which includes the forces
361 of mutation, selection, and drift (Berg and Lässig, 2003; Fisher, 1930; Iwasa, 1988; Kimura, 1962;
362 McCandlish and Stoltzfus, 2014; Sella and Hirsh, 2005; Wright, 1969). For simplicity, we ignore linkage
363 effects and, as a result of this and other assumptions, sequences evolve in a site independent manner.

364 Because SelAC requires twenty families of 61×61 matrices, the number of parameters needed to
365 implement SelAC would, without further assumptions, be extremely large (i.e. on the order of 74,420
366 parameters). To reduce the number of parameters needed, while still maintaining a high degree of
367 biological realism, we construct our gene and amino acid specific substitution matrices using a submodel
368 nested within our substitution model, similar to approaches in Gilchrist (2007); Gilchrist et al. (2015);
369 Shah and Gilchrist (2011).

370 One advantage of a nested modeling framework is that it requires only a handful of genome-
371 wide parameters such as nucleotide specific mutation rates (scaled by effective population size N_e),
372 amino acid side chain physicochemical weighting parameters, and a shape parameter describing the
373 distribution of site sensitivities. In addition to these genome-wide parameters, SelAC requires a gene g
374 specific functionality expression parameter ψ_g which describes the average rate at which the protein's
375 functionality is produced by the organism or a gene's 'average functionality production rate' for short (for
376 notational simplicity, we will ignore the gene specific indicator $_g$, unless explicitly needed). Currently, ψ
377 is fixed across the phylogeny, though relaxing this assumption is a goal of future work. The gene specific
378 parameter ψ is multiplied by additional model terms to make a composite term ψ' which scales the
379 strength and efficacy of selection for the optimal amino acid sequence relative to drift (see Implementation
380 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 (Halpern and Bruno, 1998; Lartillot and Philippe, 2004; Rodrigue and Lartillot, 2014).

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 μ

We begin with a 4x4 nucleotide mutation matrix μ that describes mutation rates between different bases and, in turn, different codons. For our purposes, we rely on the general unrestricted model (UNREST from Yang, 1994) 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→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 π whose entries, $\pi_{i,i}$, correspond to the equilibrium frequencies of each base. These equilibrium nucleotide frequencies are determined by analytically solving $\pi \times \mathbf{Q} = 0$. We use this \mathbf{Q} to populate a 61 × 61 codon mutation matrix μ , 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.

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

417 **Protein Synthesis Cost-Benefit Function η**

418 SelAC links fitness to the product of the cost-benefit function of a gene η and the organism's average
 419 target synthesis rate of the functionality provided by gene ψ . As a result, the average flux energy an
 420 organism spends to meet its target functionality provided by the gene is $\eta \times \psi$. Compensatory changes
 421 that allow an organism to maintain functionality even with loss of one or both copies of a gene are
 422 widespread. There is evidence of compensation for protein function. Metabolism with gene expression
 423 models (ME-models) link those factors to successfully make predictions about response to perturbations
 424 in a cell (King et al., 2015; Lerman et al., 2012). For example, an ME-model for *E. coli* successfully
 425 predicted gene expression levels in vivo (Thiele et al., 2012). Here we assume that for finer scale problems
 426 than entire loss (for example, a 10% loss of functionality) the compensation is more production of the
 427 protein. The particular type of dosage compensation assumed by SelAC in response to stress (e.g.
 428 reduced functionality) is commonly assumed in microbial ecology (Allison, 2012; Allison and Goulden,
 429 2017). Our assumption is also consistent with the Michaelis-Menten enzyme kinetics. Moreover, there is
 430 evidence that mutations can influence expression level, though this does not always match our expression
 431 compensation assumption (Brown and Elliot, 1997; Zanger and Schwab, 2013). In order to link genotype
 432 to our cost-benefit function $\eta = \mathbf{C}/\mathbf{B}$, we begin by defining our benefit function \mathbf{B} .

433 *Benefit:* Our benefit function \mathbf{B} measures the functionality of the amino acid sequence \vec{a}_i encoded by a
 434 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
 435 $\mathbf{B}(\vec{a}_i|\vec{a}^*) < 1$ for all other sequences. We assume all amino acids within the sequence contribute to protein
 436 function and that this contribution declines as an inverse function of physicochemical distance from each
 437 amino acid to the optimal one. Formally, we assume that

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

438 where n is the length of the protein, $d(a_p, a_p^*)$ is a weighted physicochemical distance between the amino
 439 acid encoded at a given position p and a_p^* is the optimal amino acid for that position. There are many
 440 possible measures for physicochemical distance; we use Grantham (1974) distances by default, though

441 others may be chosen. For simplicity, we assume all nonsense mutations are lethal by defining the the
442 physicochemical distance between a stop codon and a sense codon as ∞ . The term G_p describes the
443 sensitivity of the protein's function to physicochemical deviation from the optimum at site position p .
444 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
445 the Gamma distribution, the variance in G_p is equal to $\text{shape}/\text{rate}^2=1/\alpha_G$. We note that at the limit of
446 $\alpha_G \rightarrow \infty$, the model becomes equivalent to assuming uniform site sensitivity where $G_p=1$ for all positions
447 p . Further, $\mathbf{B}(\vec{a}_i|\vec{a}^*)$ is inversely proportional to the average physicochemical deviation of an amino acid
448 sequence \vec{a}_i from the optimal sequence \vec{a}^* weighted by each site's sensitivity to this deviation. $\mathbf{B}(\vec{a}_i|\vec{a}^*)$
449 can be generalized to include second and higher order terms of the distance measure d .

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

$$\begin{aligned} \mathbf{C}(\vec{c}_i) &= \text{Direct energetic cost of protein synthesis.} \\ &= A_1 + A_2 n \end{aligned}$$

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

463 **Defining Physicochemical Distances**

464 Assuming that functionality declines with an amino acid a_i 's physicochemical distance from the
465 optimum amino acid a^* at each site provides a biologically defensible way of mapping genotype to protein
466 function that requires relatively few free parameters. In addition, SelAC naturally lends itself to model
467 selection since one could compare the quality of SelAC fits using different mixtures of physicochemical
468 properties. Following (Grantham, 1974), we focus on using composition c , polarity p , and molecular
469 volume v of each amino acid's side chain residue to define our distance function, but the model and

470 its implementation can flexibly handle a variety of properties. We use the Euclidian distance between
 471 residue properties where each property c , p , and v has its own weighting term, α_c , α_p , α_v , respectively,
 472 which we refer to as ‘Grantham weights’. Because physicochemical distance is ultimately weighted by a
 473 gene’s specific average protein synthesis rate ψ , another parameter we estimate, there is a problem with
 474 parameter identifiability. The scale of gene expression is affected by how we measure physicochemical
 475 distances which, in turn, is determined by our choice of Grantham weights. As a result, by default we
 476 set $\alpha_v = 3.990 \times 10^{-4}$, the value originally estimated by Grantham, and recognize that our estimates of
 477 α_c and α_p and ψ are scaled relative to this choice for α_v . More specifically,

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

478 **Linking Protein Synthesis to Allele Substitution**

479 Next we link the protein synthesis cost-benefit function η of an allele with its fixation probability.
 480 First, we assume that each protein encoded within a genome provides some beneficial function and that
 481 the organism needs that functionality to be produced at a target average rate ψ . Again, by definition,
 482 the optimal amino acid sequence for a given gene, \vec{a}^* , produces one unit of functionality, i.e. $\mathbf{B}(\vec{a}^*) = 1$.
 483 Second, we assume that the actual average rate a protein is synthesized ϕ is regulated by the organism
 484 to ensure that functionality is produced at rate ψ . As a result, it follows that $\phi = \psi / \mathbf{B}(\vec{a} | \vec{a}^*)$ and the
 485 energetic burden of a suboptimal amino acid increases the more it decreases the protein’s functionality,
 486 \mathbf{B} . In other words, the average production rate of a protein \vec{a} with relative functionality $\mathbf{B}(\vec{a}) < 1$ must
 487 be $1/\mathbf{B}(\vec{a} | \vec{a}^*)$ times higher than the production rate needed if the optimal amino acid sequence \vec{a}^* was
 488 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
 489 produce the protein at rate $\phi = 10/9 \times \psi = 1.11\psi$. Similarly, a cell with an allele \vec{a} where $\mathbf{B}(\vec{a} | \vec{a}^*) = 1/2$
 490 will have to produce the protein at $\phi = 2\psi$. In contrast, a cell with the optimal allele \vec{a}^* would have to
 491 produce the protein at rate $\phi = \psi$.

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

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

495 where A_0 , again, describes the proportional decline in fitness with every $\sim P$ wasted per unit time.
 496 Because A_0 shares the same time units as ψ and ϕ and only occurs in SelAC in conjunction with ψ , we

497 do not need to explicitly identify our time units. Instead, we recognize that our estimates of ψ share an
 498 unknown scaling term.

499 Correspondingly, the ratio of fitness between two genotypes is,

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

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

$$W_i/W_j = \exp \left[- \left(\frac{A_0(A_1 + A_2n_g)}{n_g} \right) \sum_{p \in \mathbb{P}} [d(a_{i,p}, a_p^*) - d(a_{j,p}, a_p^*)] G_p \psi \right]$$

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

$$\begin{aligned} u_{i,j} &= \frac{1 - (W_i/W_j)^b}{1 - (W_i/W_j)^{2N_e}} \\ &= \frac{1 - \exp \left\{ - \frac{A_0}{n_g} (A_1 + A_2n_g) [d(a_i, a^*) - d(a_j, a^*)] G_p \psi b \right\}}{1 - \exp \left\{ - \frac{A_0}{n_g} (A_1 + A_2n_g) [d(a_i, a^*) - d(a_j, a^*)] G_p \psi 2N_e \right\}} \end{aligned}$$

508 where $b=1$ for a diploid population and 2 for a haploid population (Berg and Lässig, 2003; Iwasa, 1988;
 509 Kimura, 1962; Sella and Hirsh, 2005; Wright, 1969). Finally, assuming a constant mutation rate between
 510 alleles i and j , $\mu_{i,j}$, the substitution rate from allele i to j can be modeled as,

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

511 where, given the substitution model's weak mutation assumption, $N_e\mu \ll 1$. In the end, each optimal
 512 amino acid has a separate 61×61 substitution rate matrix \mathbf{Q}_a , which incorporates selection for the
 513 amino acid (and the fixation rate matrix this creates) as well as the common mutation parameters across
 514 optimal amino acids. This results in the creation of 20 \mathbf{Q} matrices, one for each amino acid and each with
 515 3,721 entries which are based on a relatively small number of model parameters (one to 11 mutation rates,
 516 two free Grantham weights, the cost of protein assembly, A_1 and A_2 , the gene specific target functionality
 517 synthesis rate ψ , and optimal amino acid at each position p , a_p^*). These model parameters can either be
 518 specified *a priori* and/or estimated from the data.

519 Given our assumption of independent evolution among sites, it follows that the probability of the
520 whole data set is the product of the probabilities of observing the data at each individual site. Thus, the
521 likelihood \mathcal{L} of amino acid a being optimal at a given site position p is calculated as

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

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

531 Finally, we note that because we infer the ancestral state of the system, our approach does not rely
532 on any assumptions of model stationarity. Nevertheless, as our branch lengths grow the probability
533 of observing a particular amino acid a at a given site approaches a stationary value proportional to
534 $W(a)^{2N_e-b}$ and any effects of mutation bias (Sella and Hirsh, 2005).

535 **Implementation**

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

544 For the SelAC models, the starting guess for the optimal amino acid at a site comes from ‘majority’
545 rule, where the initial optimum is the most frequently observed amino acid at a given site (ties resolved
546 randomly). Our optimization routine utilizes a four stage hill climbing approach. More specifically, within
547 each stage a block of parameters are optimized while the remaining parameters are held constant. The
548 first stage optimizes the block of branch length parameters. The second stage optimizes the block of

549 gene specific composite parameters $\psi'_g = A_0\psi_g N_e b$. The third stage optimizes SelAC's parameters shared
550 across the genome α_c and α_p , and the sensitivity distribution parameter α_G . The fourth stage estimates
551 the optimal amino acid at each site a^* . This entire four stage cycle is repeated six more times, using the
552 estimates from the previous cycle as the initial conditions for the new one. The search is terminated when
553 the improvement in the log-likelihood between cycles is less than 10^{-8} at which point we consider the
554 ML solution found and the search is terminated. For optimization of a given set of parameters, we rely
555 on a bounded subplex routine (Rowan, 1990) in the package `NLoptR` (Johnson, 2012) to maximize the
556 log-likelihood function. To ensure the robustness of our results, we perform a set of independent analyses
557 with different sets of naive starting points with respect to the gene specific composite ψ' parameters, α_c ,
558 and α_p and were able to repeatedly reach the same log-likelihood (lnL) peak in our parameter space.
559 Confidence in the parameter estimates can be generated by an 'adaptive search' procedure that we
560 implemented to provide an estimate of the parameter space that is some pre-defined likelihood distance
561 (e.g., 2 lnL units) from the maximum likelihood estimate (MLE), which follows Beaulieu and O'Meara
562 (2016) and Edwards (1984).

563 We note that our current implementation of SelAC is painfully slow, and is best suited for data sets
564 with relatively few number of taxa (i.e. < 10). This limitation is largely due to the size and quantity of
565 matrices we create and manipulate to calculate the log-likelihood of an individual site. Ongoing work
566 will address the need for speed, with the eventual goal of implementing SelAC in popular phylogenetic
567 inference toolkits, such as RevBayes (Hhna et al., 2016), PAML (Yang, 2007) and RAxML (Stamatakis,
568 2006).

569 **Simulations**

570 We evaluated the performance of our codon model by simulating datasets and estimating the bias of the
571 inferred model parameters from these data. Our 'known' parameters under a given generating model were
572 based on fitting SelAC to the 106 gene data set and phylogeny of Rokas et al. (2003). The tree used in
573 these analyses is outdated with respect to the current hypothesis of relationships within *Saccharomyces*,
574 but we rely on it simply as a training set that is separate from our empirical analyses (see section below).
575 Bias in the model parameters were assessed under two generating models: one where we assumed a model
576 of SelAC assuming uniform sensitivity across sites (i.e. $G_p = 1$ for all sites, i.e. $\alpha_G = \infty$), and one where
577 we used the Gamma distribution joint shape and rate parameter α_G estimated from the empirical data.
578 Under each of these two scenarios, we used parameter estimates from the corresponding empirical analysis
579 and simulated 50 five-gene data sets. For the gene specific composite parameter ψ'_g the 'known' values
580 used for the simulation were five evenly spaced points along the rank order of the estimates across the

581 106 genes. The MLE estimate for a given replicate were taken as the fit with the highest log-likelihood
582 after running five independent analyses with different sets of naive starting points with respect to the
583 composite ψ'_g parameter, α_c , and α_p . All analyses were carried out in our `selac` R package.

584 **Analysis of yeast genomes & tests of model adequacy**

585 We focus our empirical analyses on the large yeast data set and phylogeny of Salichos and Rokas
586 (2013). As a model system, the yeast genome is an ideal system to examine our phylogenetic estimates
587 of gene expression and its connection to real world measurements of these data within individual taxa.
588 The complete data set of Salichos and Rokas (2013) contain 1070 orthologs, where we selected 100 at
589 random for our analyses. We also focus our analyses on *Saccharomyces sensu stricto* and their sister
590 taxon *Candida glabrata*, and we used the phylogeny depicted in Fig. 1 of Salichos and Rokas (2013) for
591 our fixed tree. We fit the two SelAC models described above (i.e., SelAC and SelAC+ Γ), as well as two
592 codon models, GY94 and FMutSel, and a standard GTR + Γ nucleotide model. The FMutSel model
593 assumes that the amino acid frequencies are determined by functional requirements of the protein while
594 the other models make no assumptions about amino acid frequencies. In all cases, we assumed that the
595 model was partitioned by gene, but with branch lengths linked across genes.

596 For SelAC, we compared our estimates of $\phi' = \psi' / \mathbf{B}$, which represents the average protein synthesis
597 rate of a gene, to estimates of gene expression from empirical data. Specifically, we examined gene
598 expression data for five of the six species measured during log-growth phase. Gene expression in this
599 context corresponds to mRNA abundances, which were measured using either microarrays (*C. glabrata*
600 and *S. castellii*, or RNA-Seq (*S. paradoxus*, *S. mikatae*, and *S. cerevisiae*). We obtained expression data
601 for the remaining species, *S. kudriavzevii*, which was measured at the beginning of the stationary phase
602 from the Gene Expression Omnibus (GEO). *Saccharomyces*, however, only enter the stationary growth
603 phase in response to severe stress, such as starvation. In addition, only 56 % of the genes examined with
604 SelAC had expression measurements available. For these reasons, we excluded *S. kudriavzevii* from our
605 comparisons of empirical gene expression.

606 For further comparison, we also predicted the average protein synthesis rate for each gene ϕ by analyzing
607 gene and genome-wide patterns of synonymous codon usage using ROC-SEMPPR (Gilchrist et al., 2015)
608 for each individual genome. While, like SelAC, ROC-SEMPPR uses codon level information, it does not
609 rely on any interspecific comparisons and, unlike SelAC, uses only the intra- and inter-genic frequencies
610 of synonymous codon usage as its data. Nevertheless, ROC-SEMPPR predictions of gene expression
611 ϕ correlates strongly (Pearson $r = 0.53 - 0.74$) with a wide range of laboratory measurements of gene
612 expression (Gilchrist et al., 2015).

613 While one of our main objectives was to determine the improvement of fit that SelAC has with respect
614 to other standard phylogenetic models, we also evaluated the adequacy of SelAC. Model fit, measured
615 with assessments such as the Akaike Information Criterion (AIC), can tell which model is least bad
616 as an approximation for the data, but it does not reveal whether a model is actually doing a good
617 job of representing the data. An adequate model does the latter, one measure of which is that data
618 generated under the model resemble real data (Goldman, 1993). For example, Beaulieu et al. (2013)
619 assessed whether parsimony scores and the size of monomorphic clades of empirical data were within
620 the distributions of simulated data under a new model and the best standard model; if the empirical
621 summaries were outside the range for each, it would have suggested that neither model was adequately
622 modeling this part of the biology.

623 In order to test adequacy for a given gene we first remove a particular taxon from the data set
624 and the phylogeny. A marginal reconstruction of the likeliest sequence across all remaining nodes is
625 conducted under the model, including the node where the pruned taxon attached to the tree. The
626 marginal probabilities of each site are used to sample and assemble the starting coding sequence. This
627 sequence is then evolved along the branch, periodically being sampled and its current functionality
628 assessed. We repeat this process 100 times and compare the distribution of trajectories against the
629 observed functionality calculated for the gene. For comparison, we also conducted the same test, by
630 simulating the sequence under the standard GTR + Γ nucleotide model, which is often used on these
631 data but does not account for the fact that the sequences are protein coding, and under FMutSel, which
632 includes selection on codons but in a fundamentally different way as our model.

633 **The appropriate estimator of bias for AIC**

634 As part of the model set described above, we also included a reduced form of each of the two SelAC
635 models, SelAC and SelAC+ Γ . Specifically, rather than optimizing the amino acid at any given site, we
636 assume the the most frequently observed amino acid at each site is the optimal amino acid a^* . We refer to
637 these ‘majority rule’ models as SelAC_M and SelAC_M+ Γ and note that these majority rule formulations
638 greatly accelerate model fitting.

639 Since these majority rule models assume that the optimal amino acids are known prior to fitting of
640 our model, it is tempting to reduce the count of estimated parameters in the model by the number of
641 parameters estimated using majority rule. While using majority rule does not necessarily provide the
642 most likely parameter estimate, it nevertheless uses the data to generate the estimate and, , represents
643 a parameter estimated from the data. Thus, despite having become standard behavior in the field
644 of phylogenetics, this reduction is statistically inappropriate. Because the difference in the number of

645 parameters K when counting or not counting the number of nucleotide sites drops out when comparing
646 nucleotide models with AIC, this statistical issue does not apply to nucleotide models. It does, however,
647 matter for AICc, where K and the sample size n combine in the penalty term. This also matters in our
648 case, where the number of estimated parameters for the majority rule estimation differs based on whether
649 one is looking at codons or single nucleotides.

650 In phylogenetics two variants of AICc are used. In comparative methods (e.g. Beaulieu et al., 2013;
651 Butler and King, 2004; O'Meara et al., 2006) the number of data points, n , is taken as the number of
652 taxa. More taxa allow the fitting of more complex models, given more data. However, in DNA evolution,
653 which is effectively the same as a discrete character model used in comparative methods, the n is taken
654 as the number of sites. Obviously, both cannot be correct. This uncertainty was highlighted by Posada
655 and Buckley (2004): they chose to use number of sites, but mentioned in their discussion that sample size
656 also depends on the number of taxa. Sullivan and Joyce (2005) also mention that while the number of
657 sites is often taken as sample size, whether that is appropriate in phylogenetics is not entirely clear. One
658 approach incorporating both number of taxa and sites in calculating AICc is the program SURFACE
659 implemented by Ingram and Mahler (2013), which uses multiple characters and taxa. While its default
660 is to use AIC to compare models, if one chooses to use AICc, the number of samples is taken as the
661 product of number of sites and number of taxa.

662 Recently, Jhwueng et al. (2014) performed an analysis that investigated what variant of AIC and AICc
663 worked best as an estimator, but the results were inconclusive. Here, we have adopted and extended the
664 simulation approach of Jhwueng et al. (2014) in order to examine a large set of different penalty functions
665 and how well they approximate the remaining portion of the Kullback-Liebler (KL) divergence between
666 two models after accounting for the deviance (i.e., $-2\mathcal{L}$) (see Appendix 1 for more details).

667 **Acknowledgements**

668 This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ) and DEB-1355033
669 (BCO, MAG, and RZ) with additional support from The University of Tennessee Knoxville and University
670 of Arkansas (JMB). JJC and JMB received support as Postdoctoral Fellows and CL received support as a
671 Graduate Student Fellow at the National Institute for Mathematical and Biological Synthesis, an Institute
672 sponsored by the National Science Foundation through NSF Award DBI-1300426, with additional support
673 from UTK. The authors would like to thank Premal Shah, Todd Oakley, and two anonymous reviewers
674 for their helpful criticisms and suggestions for this work.

References

- 675
676 Allison, S. 2012. A trait-based approach for modelling microbial litter decomposition. *Ecology Letters* 15:1058–1070.
- 677 Allison, S. and M. Goulden. 2017. Consequences of drought tolerance traits for microbial decomposition in the DEMENT
678 model. *Soil Biology & Biochemistry* 107:104–113.
- 679 Anisimova, M. 2012. Parametric models of codon evolution. Pages 12–33 in *Codon Evolution: Mechanisms and Models*
680 (G. M. Cannarozzi and A. Schneider, eds.). Oxford University Press, Oxford, UK.
- 681 Beaulieu, J. M. and B. C. O’Meara. 2016. Detecting Hidden Diversification Shifts in Models of Trait-Dependent Speciation
682 and Extinction. *Systematic Biology* 65:583–601.
- 683 Beaulieu, J. M., B. C. O’Meara, and M. J. Donoghue. 2013. Identifying Hidden Rate Changes in the Evolution of a Binary
684 Morphological Character: The Evolution of Plant Habit in Campanulid Angiosperms. *Systematic Biology* 62:725–737.
- 685 Berg, J. and M. Lässig. 2003. Stochastic Evolution and Transcription Factor Binding Sites. *Biophysics* 48:S36–S44.
- 686 Brown, L. and T. Elliot. 1997. Mutations That Increase Expression of the rpoS Gene and Decrease Its Dependence on hfq
687 Function in *Salmonella typhimurium*. *J. Bacteriol.* 179:656–662.
- 688 Butler, M. A. and A. A. King. 2004. Phylogenetic comparative analysis: a modeling approach for adaptive evolution.
689 *American Naturalist* 164:683–695.
- 690 Dimmic, M. W., D. P. Mindell, and R. A. Goldstein. 2000. Modeling evolution at the protein level using an adjustable amino
691 acid fitness model. *Pacific Symposium on Biocomputing* 5:18–29.
- 692 Drummond, D. A., J. D. Bloom, C. Adami, C. O. Wilke, and F. H. Arnold. 2005. Why highly expressed proteins evolve
693 slowly. *Proceedings of the National Academy of Sciences of the United States of America* 102:14338–14343.
- 694 Drummond, D. A., A. Raval, and C. O. Wilke. 2006. A single determinant dominates the rate of yeast protein evolution.
695 *Molecular Biology and Evolution* 23:327–337.
- 696 Edwards, A. 1984. *Likelihood*. Cambridge science classics Cambridge University Press.
- 697 Endler, J. A. 1986. Natural Selection in the Wild Pages 16–17. No. 21 in *Monographs in Population Biology* Princeton
698 University Press, Princeton, NJ reference for definition of diversifying selection.
- 699 Felsenstein, J. 1981. Evolutionary trees from DNA-sequences - a maximum-likelihood approach. *Journal of Molecular*
700 *Evolution* 17:368–376.
- 701 Felsenstein, J. 2001. Taking Variation of Evolutionary Rates Between Sites into Account in Inferring Phylogenies. *Journal*
702 *of Molecular Evolution* 53:447–455.
- 703 Fisher, S., Ronald A. 1930. *The Genetical Theory of Natural Selection*. Oxford University Press, Oxford.
- 704 Gilchrist, M., P. Shah, and R. Zaretzki. 2009. Measuring and detecting molecular adaptation in codon usage against nonsense
705 errors during protein translation. *Genetics* 183:1493–1505.
- 706 Gilchrist, M. A. 2007. Combining Models of Protein Translation and Population Genetics to Predict Protein Production
707 Rates from Codon Usage Patterns. *Molecular Biology and Evolution* 24:2362–2373.
- 708 Gilchrist, M. A., W.-C. Chen, P. Shah, C. L. Landerer, and R. Zaretzki. 2015. Estimating Gene Expression and Codon-
709 Specific Translational Efficiencies, Mutation Biases, and Selection Coefficients from Genomic Data Alone. *Genome Biology*
710 *and Evolution* 7:1559–1579.
- 711 Gilchrist, M. A. and A. Wagner. 2006. A model of protein translation including codon bias, nonsense errors, and ribosome
712 recycling. *Journal of Theoretical Biology* 239:417–434.

- 713 Goldman, N. 1993. Statistical tests of models of DNA substitution. *Journal of molecular evolution* 36:182–198.
- 714 Goldman, N., J. L. Thorne, and D. T. Jones. 1996. Using Evolutionary Trees in Protein Secondary Structure Prediction
715 and Other Comparative Sequence Analyses. *Journal of Molecular Biology* 263:196 – 208.
- 716 Goldman, N., J. L. Thorne, and D. T. Jones. 1998. Assessing the Impact of Secondary Structure and Solvent Accessibility
717 on Protein Evolution. *Genetics* 149:445–458.
- 718 Goldman, N. and Z. H. Yang. 1994. Codon-based model of nucleotide substitution for protein-coding DNA-sequences.
719 *Molecular Biology and Evolution* 11:725–736.
- 720 Grantham, R. 1974. Amino acid difference formula to help explain protein evolution. *Science* 185:862–864.
- 721 Halpern, A. L. and W. J. Bruno. 1998. Evolutionary distances for protein-coding sequences: Modeling site-specific residue
722 frequencies. *Molecular Biology And Evolution* 15:910–917.
- 723 Hughes, A. L. and M. Nei. 1988. Pattern of nucleotide substitution at major histocompatibility complex class-i loci reveals
724 overdominant selection. *Nature* 335:167–170.
- 725 Hhna, S., M. J. Landis, T. A. Heath, B. Boussau, N. Lartillot, B. R. Moore, J. P. Huelsenbeck, and F. Ronquist. 2016.
726 RevBayes: Bayesian Phylogenetic Inference Using Graphical Models and an Interactive Model-Specification Language.
727 *Systematic Biology* 65:726.
- 728 Ingram, T. and D. L. Mahler. 2013. SURFACE: detecting convergent evolution from data by fitting Ornstein-Uhlenbeck
729 models with stepwise Akaike Information Criterion. *Methods in ecology and evolution* 4:416–425.
- 730 Iwasa, Y. 1988. Free fitness that always increases in evolution. *Journal of Theoretical Biology* 135:265–281.
- 731 Jhwueng, D.-C., H. Snehalata, B. C. O’Meara, and L. Liu. 2014. Investigating the performance of AIC in selecting
732 phylogenetic models. *Statistical applications in genetics and molecular biology* 13:459–475.
- 733 Johnson, S. G. 2012. The NLOpt nonlinear-optimization package. Version 2.4.2 – Released 20 May 2014.
- 734 Kimura, M. 1962. on the probability of fixation of mutant genes in a population. *Genetics* 47:713–719.
- 735 King, Z. A., C. J. Lloyd, A. M. Feist, and B. O. Palsson. 2015. Next-generation genome-scale models for metabolic
736 engineering. *Current Opinion in Biotechnology* 35:23 – 29 chemical biotechnology Pharmaceutical biotechnology.
- 737 Koshi, J. M. and R. A. Goldstein. 1997. Mutation matrices and physical-chemical properties: Correlations and implications.
738 *Proteins-Structure Function And Genetics* 27:336–344.
- 739 Koshi, J. M. and R. A. Goldstein. 2000. Analyzing site heterogeneity during protein evolution. Pages 191–202 *in*
740 *Biocomputing 2001*. World Scientific.
- 741 Koshi, J. M., D. P. Mindell, and R. A. Goldstein. 1999. Using physical-chemistry-based substitution models in phylogenetic
742 analyses of HIV-1 subtypes. *Molecular biology and evolution* 16:173–179.
- 743 Kubatko, L., P. Shah, R. Herbei, and M. A. Gilchrist. 2016. A codon model of nucleotide substitution with selection on
744 synonymous codon usage. *Molecular Phylogenetics and Evolution* 94:290 – 297.
- 745 Lartillot, N. and H. Philippe. 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement
746 process. *Molecular Biology And Evolution* 21:1095–1109.
- 747 Lerman, J. A., D. R. Hyduke, H. Latif, V. A. Portnoy, N. E. Lewis, J. D. Orth, A. C. Schrimpe-Rutledge, R. D. Smith, J. N.
748 Adkins, K. Zengler, and B. O. Palsson. 2012. In silico method for modelling metabolism and gene product expression at
749 genome scale. *Nature Communications* 3:929 EP – article.
- 750 Lynch, M. and G. K. Marinov. 2015. The bioenergetic costs of a gene. *Proceedings Of The National Academy Of Sciences*
751 *Of The United States Of America* 112:15690–15695.

- 752 Mayrose, I., N. Friedman, and T. Pupko. 2005. A Gamma mixture model better accounts for among site rate heterogeneity.
753 *Bioinformatics* 21:ii151–ii158.
- 754 McCandlish, D. M. and A. Stoltzfus. 2014. Modeling evolution using the probability of fixation: History and implications.
755 *The Quarterly Review of Biology* 89:225–252.
- 756 Muse, S. V. and B. S. Gaut. 1994. A likelihood approach for comparing synonymous and nonsynonymous nucleotide
757 substitution rates, with application to the chloroplast genome. *Molecular Biology and Evolution* 11:715–724.
- 758 Nielsen, R. and Z. H. Yang. 1998. Likelihood models for detecting positively selected amino acid sites and applications to
759 the HIV-1 envelope gene. *Genetics* 148:929–936.
- 760 Nowak, M. A. 2006. *Evolutionary Dynamics: Exploring the Equations of Life*. Belknap of Harvard University Press,
761 Cambridge, MA.
- 762 O’Meara, B. C., C. Ane, M. J. Sanderson, and W. P.C. 2006. Testing for different rates of continuous trait evolution using
763 likelihood. *Evolution* 60:922–933.
- 764 Pellmyr, O. 2002. Microevolution. Pages 731–732 *in* *Encyclopedia of Evolution* (M. Pagel, ed.). Oxford University Press,
765 Oxford, UK.
- 766 Pellmyr, O. 2002. Microevolution. Pages 731–732 *in* *Encyclopedia of Evolution* (M. Pagel, ed.) vol. 2. Oxford University
767 Press, Oxford, UK.
- 768 Posada, D. and T. R. Buckley. 2004. Model Selection and Model Averaging in Phylogenetics: Advantages of Akaike
769 Information Criterion and Bayesian Approaches over Likelihood Ratio Tests. *Systematic Biology* 53:793–808.
- 770 Pouyet, F., M. Bailly-Bechet, D. Mouchiroud, and L. Guguen. 2016. SENCA: A Multilayered Codon Model to Study the
771 Origins and Dynamics of Codon Usage. *Genome Biology and Evolution* 8:2427–2441.
- 772 Robinson, D. M., D. T. Jones, H. Kishino, N. Goldman, and J. L. Thorne. 2003. Protein evolution with dependence among
773 codons due to tertiary structure. *Molecular Biology And Evolution* 20:1692–1704.
- 774 Rodrigue, N. and N. Lartillot. 2014. Site-heterogeneous mutation-selection models within the PhyloBayes-MPI package.
775 *Bioinformatics* 30:1020–1021.
- 776 Rodrigue, N., N. Lartillot, D. Bryant, and H. Philippe. 2005. Site interdependence attributed to tertiary structure in amino
777 acid sequence evolution. *Gene* 347:207–217.
- 778 Rokas, A., B. L. Williams, N. King, and S. B. Carroll. 2003. Genome-scale approaches to resolving incongruence in molecular
779 phylogenies. *Nature* 425:798–804.
- 780 Rowan, T. 1990. *Functional Stability Analysis of Numerical Algorithms*. Ph.D. thesis University of Texas, Austin.
- 781 Salichos, L. and A. Rokas. 2013. Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature*
782 497:327–331.
- 783 Sella, G. and A. E. Hirsh. 2005. The application of statistical physics to evolutionary biology. *Proceedings of the National*
784 *Academy of Sciences of the United States of America* 102:9541–9546.
- 785 Shah, P. and M. A. Gilchrist. 2011. Explaining complex codon usage patterns with selection for translational efficiency,
786 mutation bias, and genetic drift. *Proceedings of the National Academy of Sciences of the United States of America*
787 108:10231–10236.
- 788 Stamatakis, A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed
789 models. *Bioinformatics* 22:2688–2690.

- 790 Sullivan, J. and P. Joyce. 2005. Model Selection in Phylogenetics. *Annual Review of Ecology, Evolution, and Systematics*
791 36:445–466.
- 792 Thiele, I., R. M. T. Fleming, R. Que, A. Bordbar, D. Diep, and B. O. Palsson. 2012. Multiscale Modeling of Metabolism
793 and Macromolecular Synthesis in *E. coli* and Its Application to the Evolution of Codon Usage. *PLOS ONE* 7:1–18.
- 794 Thorne, J. L., N. Goldman, and D. T. Jones. 1996. Combining protein evolution and secondary structure. *Molecular Biology*
795 and Evolution 13:666–673.
- 796 Thorne, J. L., N. Lartillot, N. Rodrigue, and S. C. Choi. 2012. Codon models as a vehicle for reconciling
797 population genetics with inter-specific sequence data. *Codon Evolution: Mechanisms And Models* Pages 97–110 D2
798 10.1093/acprof:osobl/9780199601165.001.0001 ER.
- 799 Wagner, A. 2005. Energy constraints on the evolution of gene expression. *Molecular Biology and Evolution* 22:1365–1374.
- 800 Wright, S. 1969. *Evolution and the genetics of populations. Vol. 2. The theory of gene frequencies.* vol. 2. University of
801 Chicago Press.
- 802 Yang, Z. 2014. *Molecular Evolution: A Statistical Approach.* Oxford University Press, New York.
- 803 Yang, Z. H. 1994. Maximum-likelihood phylogenetic estimation from DNA-sequences with variable rates over sites -
804 approximate methods. *Journal Of Molecular Evolution* 39:306–314.
- 805 Yang, Z. H. 2007. PAML 4: Phylogenetic analysis by maximum likelihood. *Molecular Biology And Evolution* 24:1586–1591.
- 806 Yang, Z. H. and R. Nielsen. 1998. Synonymous and nonsynonymous rate variation in nuclear genes of mammals. *Journal Of*
807 *Molecular Evolution* 46:409–418.
- 808 Yang, Z. H. and R. Nielsen. 2008. Mutation-selection models of codon substitution and their use to estimate selective
809 strengths on codon usage. *Molecular Biology and Evolution* 25:568–579.
- 810 Zanger, U. and M. Schwab. 2013. Cytochrome P450 enzymes in drug metabolism: Regulation of gene expression, enzyme
811 activities, and impact of genetic variation. *Pharmacology & Therapeutics* 138:103–141.

812 **Table**

Model	logLik	Parameters			Model	
		Estimated	AIC	AICc	Δ AICc	Weight
SelAC+ Γ	-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 + Γ	-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+ Γ	-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. 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).

813 **Figures**

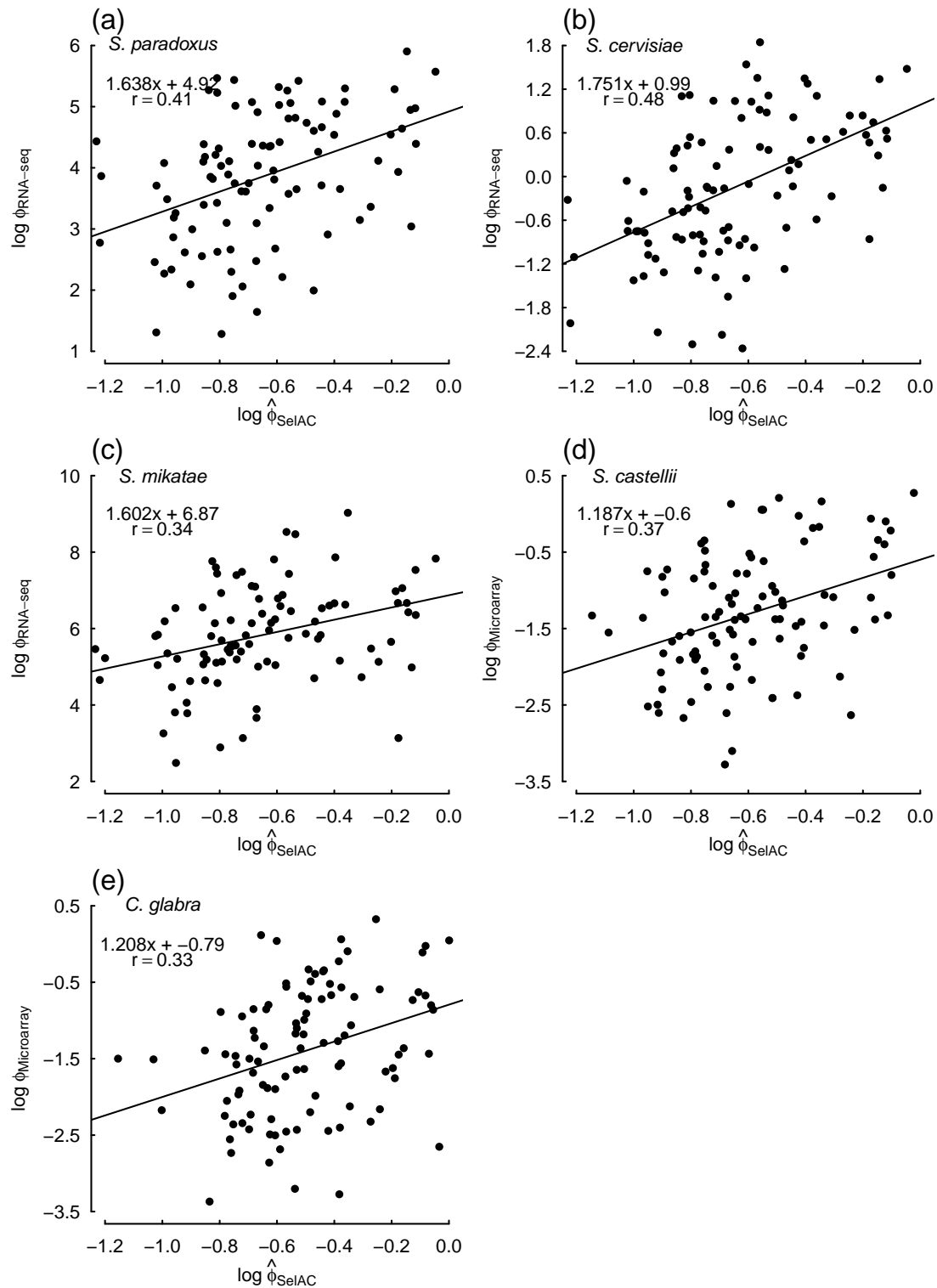


FIG. 1. Comparisons between estimates of average protein translation rate $\hat{\phi}_{\text{SelAC}}$ obtained from SelAC+ Γ and direct measurements of expression for individual yeast taxa across the 100 selected genes from Salichos and Rokas (2013) measured during log-growth phase. Estimates of $\hat{\phi}_{\text{SelAC}}$ were generated by dividing the composite term ψ^l by $\mathbf{B}(\bar{a}_i|\bar{a}^*)$. 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 .

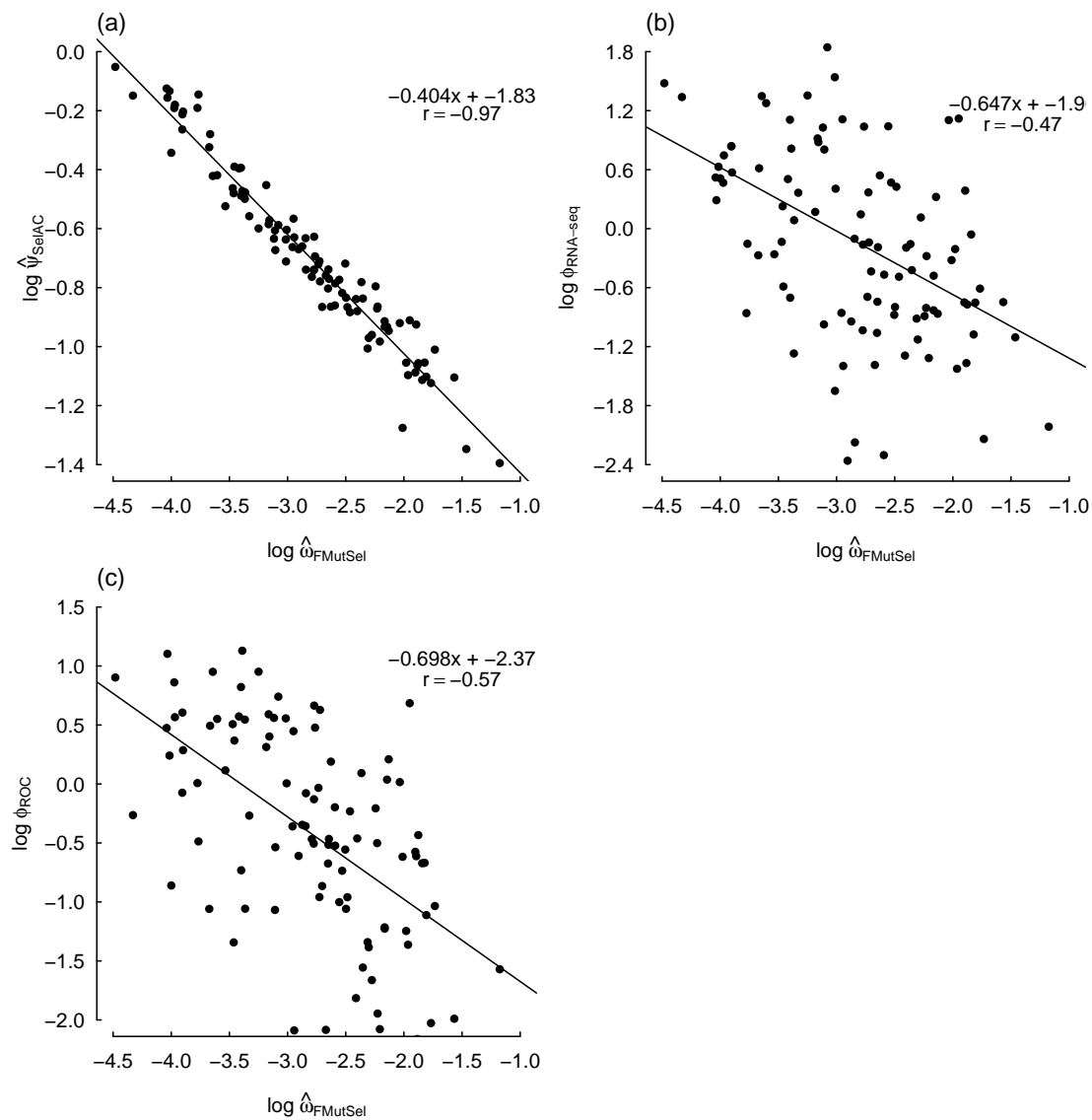


FIG. 2. Comparisons between ω_{FMutSel} , which is the nonsynonymous/synonymous mutation ratio in FMutSel, SelAC+ Γ estimates of protein functionality production rates $\hat{\psi}_{\text{SelAC}}$ (a), RNA-Seq based measurements of mRNA abundance $\phi_{\text{RNA-seq}}$ (b), and ROC-SEMPER's estimates of protein translation rates ϕ_{ROC} , which are based solely on *S. cerevisiae*'s patterns of codon usage bias (c), for *S. cerevisiae* across the 100 selected genes from Salichos and Rokas (2013). As in Figure 1, the equations in the upper right hand corner of each panel provide the regression fit and correlation coefficient.

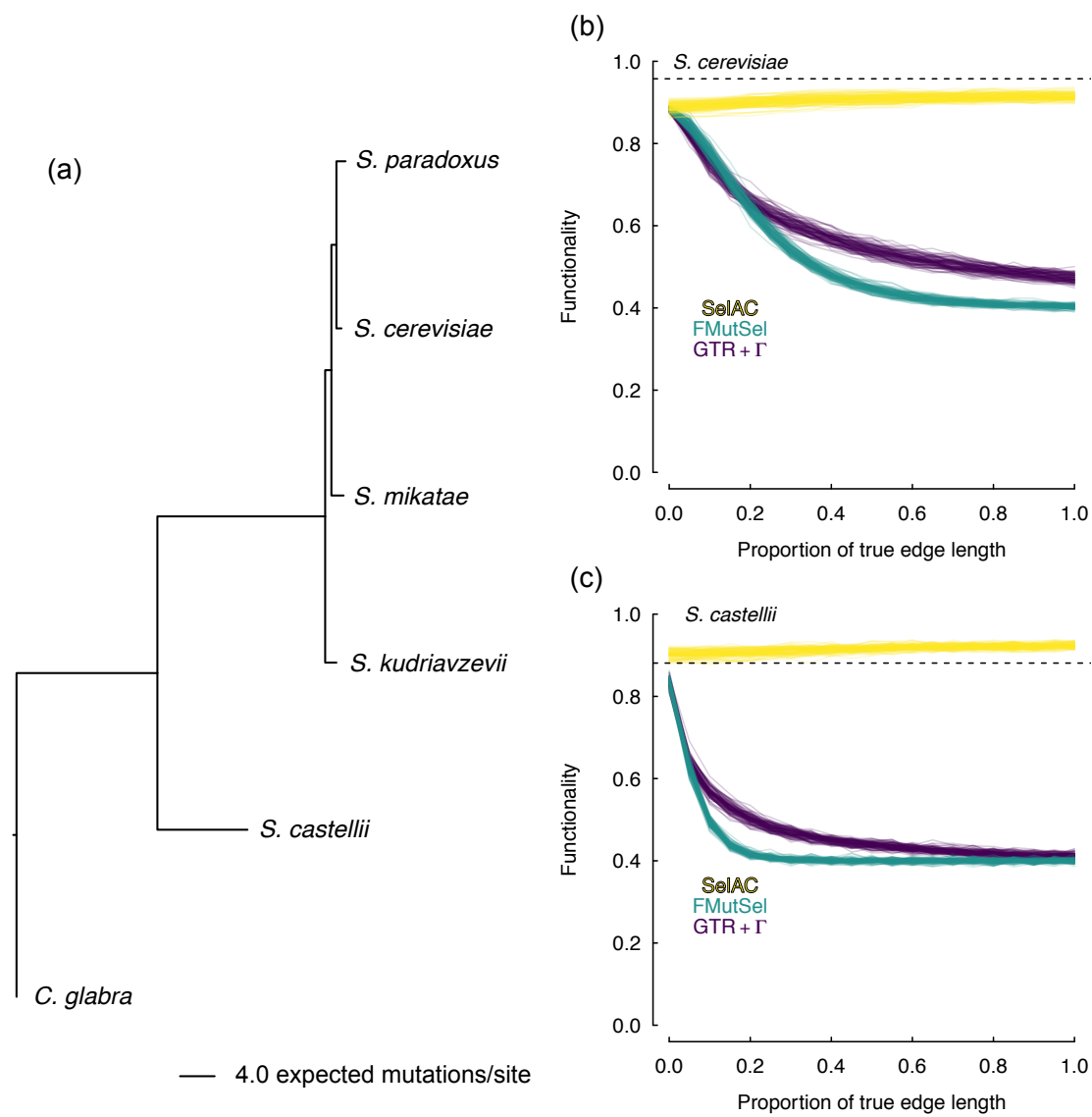


FIG. 3. (a) Maximum likelihood estimates of branch lengths under SelAC+ Γ for 100 selected genes from Salichos and Rokas (2013). Tests of model adequacy for *S. cerevisiae* (b) and *S. castellii* (c) indicated that, when these taxa are removed from the tree, and their sequences are simulated, the parameters of SelAC+ Γ exhibit functionality $\mathbf{B}(\vec{a}_{\text{obs}}|\vec{a}^*)$ that is far closer to the observed (dashed black line) than data sets produced from parameters of either FMutSel or GTR + Γ .