# Phylogenetic model of stabilizing selection is more informative about site specific selection than extrapolation from laboratory estimates

CEDRIC LANDERER[1,2,*], BRIAN C. OMEARA[1,2], AND MICHAEL A. GILCHRIST[1,2]

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

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

[*]Corresponding author. E-mail: cedric.landerer@gmail.com

Version dated: November 28, 2018

# Introduction

- Numerous attempts to incorporate selection into phylogenetic models have been made.

  - Phylogenetic inference of sequence relationship was long only focused on substitution rates and fixation probabilities.

  - However, the importance of site specific equilibrium frequencies has long been noted.

  - Models of site specific equilibrium frequencies tend to be unfeasible as they are very parameter rich.

  - Independent fitness estimates have potential to greatly reduce number of parameters estimated from phylogenetic data.

  - Incorporating selection from experimental sources therefore seems like an attractive option.

    * site specific amino acid preferences acknowledge the heterogeneity of selection along the protein sequence.

    * It allows for the fitting complex site specific models to smaller data sets.

    * DMS allows to estimate empirical selection on amino acids for large amount of mutations on a single experiment.

  - Value of DMS depends on many factors like initial library of mutants and applied selection.

    * Extensive mutation libraries lead to heterogeneous competing population.

    * DMS experiments are limited to proteins and organisms that can be manipulated under laboratory conditions.

    * Greatly limiting application of experimentally informed phylogenetic models.

- Even when empirical selection estimates are available, their application for phylogenetic inference is questionable.

- In this study, we compared experimentally inferred site specific selection to inform phylogenetic models to a site specific model of stabilizing selection.

- We assessed model fit of two codon models of site specific stabilizing selection with ($SelAC$+DMS, $phydms$) and without ($SelAC$) experimentally inferred selection and compared the models fits to 227 other codon and nucleotide models.

  - We used the class A $\beta$-lactamase TEM found in gram-negative bacteria like *E. coli* for which empirical selection estimates are available.

  - The applied selection pressure was limited to ampicillin and focused on the sequence variant TEM-1.

  - TEM can confer resistance to a wide range of antibiotics.

  - Models fits informed by experimentally inferred selection improve model fit over conventional codon and nucleotide models but can be improved upon using a hierarchical phylogenetic framework of stabilizing selection: $SelAC$.

  - Simulations highlight the inadequacy of experimentally inferred selection.

  - Comparison between $SelAC$ and empirical estimates of selection show that they are comparable when site specific selection is captured adequately by the experiment.

  - Furthermore, we show that extrapolating experimentally inferred selection between homologous proteins (TEM and SHV) can be inadequate.

# Results

## Site Specific Stabilizing Selection on Amino Acids Improves Model Fit

- We evaluated here model fits of site specific selection and 227 other codon and nucleotide models to 49 observed TEM sequences.

    - All three models of site specific selection improved model fit.

    - Number of parameters estimated from phylogenetic data differs between *SelAC*, and *SelAC*+DMS and *phydms*, resulting in slightly worse AICc for *SelAC*.

    - However, *SelAC* outperforms *phydms* (Table 1).

Table 1: Model selection, shown are the three models of stabilizing site specific amino acid selection (*SelAC*, *SelAC*+DMS, *phydms*) and the best performing codon and nucleotide model (**??**). Reported are the log-likelihood $\log(L)$, the number of parameters estimated $n$, AIC, $\Delta$AIC, AICc, and $\Delta$AICc values. See Table X for results from all models we tested.

| Model | $\log(L)$ | n | AIC | $\Delta$AIC | AICc | $\Delta$AICc |
|---|---|---|---|---|---|---|
| *SelAC*+DMS | -1768 | 111 | 3758 | 14 | 3760 | 0 |
| *SelAC* | -1498 | 374 | 3744 | 0 | 3766 | 6 |
| *phydms* | -2061 | 102 | 4326 | 582 | 4328 | 568 |
| *SYM*+R2 | -2230 | 102 | 4663 | 919 | 4694 | 934 |
| *GY94* +F1X4+R2 | -2243 | 102 | 4690 | 946 | 4821 | 1061 |

- We observe differences in topology.

    - *SelAC* is to slow for a topology search, therefore unclear if the difference in topology can be attributes to the experimentally inferred selection.

    - *GY94* is outperformed by several nucleotide model e.g. *SYM*+R2, potentially indicating that negative frequency dependent selection is inappropriate for TEM.

    - Results indicate shift in the evolution from the tips (*SelAC*) to internal branches (*SelAC*+DMS, *phydms*, *GY94*).

4

# Assessing Adequacy of Laboratory and *SelAC* Inferences of Site Specific Selection

- Assessing model adequacy as sequence similarity sequence of selectively favored amino acids and observed consensus sequence.

    - Experimentally inferred selection is inconsistent with observed sequences.

    - Experimentally inferred sequence of selectively favored amino acids has only 52% sequence similarity with the observed consensus sequence.

    - *SelAC* inferred sequence of selectively favored amino acids has 99% sequence similarity with the observed consensus sequence.

    - The average sequence similarity between the 49 observed sequences is 98%.

- Assessing model adequacy as genetic load.

    - Simulations under experimentally and *SelAC* inferred selection were used to establish a baseline expectation.

    - Assuming the site specific selection estimated by DMS, the observed TEM sequences represent an average sequence specific genetic load of 17.88 and an average site specific load of 0.065.

    - Simulated sequences showed an average sequence specific load of 6.68 and an average site specific genetic load of 0.025

    - Assuming the site specific selection estimated by *SelAC*, the observed TEM sequences represent an average sequence specific genetic load of $6.4 \times 10^{-5}$ and an average site specific load of $2.4 \times 10^{-7}$.

    - Simulated sequences showed an average sequence specific load of $1.3 \times 10^{-5}$ and an average site specific genetic load of $4.8 \times 10^{-8}$.

## Comparing Laboratory and *SelAC* Inferences of Site Specific Selection

- Distribution of genetic load differs between experimentally inferred site specific selection and *SelAC* inferred site specific selection.

    - Assuming the site specific selection estimated by DMS, 111 sites do not carry a genetic load.

    - Assuming the site specific selection estimated by *SelAC*, 207 sites do not carry a genetic load.

    - The selection estimates from DMS and *SelAC* agree for 107 sites that no genetic load is carried.

    - Thus, for 100 sites *SelAC* does not estimate a genetic load but DMS does, while the inverse is true for four sites.

    - For the 52 sites where both, DMS and *SelAC*, estimate a non-zero genetic load we find a correlation of $\rho = 0.247$, explaining 6% of the variation in the empirical selection estimates.

## Comparing *SelAC* Inferences of Site Specific Selection for Homologous Sequences TEM and SHV

- Site specific genetic load for TEM and SHV is not correlated ($\rho = 0.006$) and , despite similar $\alpha_G$

    - Excluding site with no genetic load and calculating the correlation on the log scale lead to a correlation coefficient of $\rho = 0.22$.

- Greatest difference is observed in the physicochemical properties, specifically $\alpha$.

6

- No significant differences are observed in average genetic load between secondary structure elements (Table 2).

Table 2: Efficacy of selection ($G$) and genetic load for TEM and SHV, and separated by secondary structure. $G$ was estimated as a truncated variable with an upper bound of 300.

| Protein | Secondary Structure | # Residues | $G$ Mean | SE | Genetic Load $L_i$ Mean | SE |
|---|---|---|---|---|---|---|
| TEM | | 263 | 219.3 | 7.5 | $15.9 \times 10^{-8}$ | $6.5 \times 10^{-8}$ |
| | Helix | 113 | 206.1 | 12.4 | $17.5 \times 10^{-8}$ | $13.1 \times 10^{-8}$ |
| | $\beta$-Sheet | 48 | 238.6 | 15.8 | $6.8 \times 10^{-8}$ | $2.9 \times 10^{-8}$ |
| | Unstructured | 102 | 224.8 | 11.4 | $18.6 \times 10^{-8}$ | $8.1 \times 10^{-8}$ |
| | Active/Binding Sites | 5 | 202.6 | 62.2 | $0.01 \times 10^{-8}$ | $0.01 \times 10^{-8}$ |
| SHV | | 263 | 244.9 | 6.8 | $4.0 \times 10^{-8}$ | $1.9 \times 10^{-8}$ |
| | Helix | 102 | 234.6 | 11.5 | $7.3 \times 10^{-8}$ | $4.8 \times 10^{-8}$ |
| | $\beta$-Sheet | 66 | 253.1 | 12.8 | $2.1 \times 10^{-8}$ | $1.1 \times 10^{-8}$ |
| | Unstructured | 95 | 224.7 | 11.4 | $1.5 \times 10^{-8}$ | $0.6 \times 10^{-8}$ |
| | Active/Binding Sites | 5 | 239.9 | 60.0 | $1.5 \times 10^{-8}$ | $1.5 \times 10^{-8}$ |

# Discussion

- We evaluated how well experimental selection estimates from DMS experiments explain natural sequence evolution and compared it to a novel phylogenetic framework, SelAC.

  - Previous work has shown that DMS selection estimates can improve model fit over classical approaches like GY94 and our work confirms this.

  - Model selection favored the SelAC model fit and the corresponding fitness estimates over the DMS estimates using both, SelAC and phyDMS (Table 1).

- Adequacy of the DMS selection has previously not been assessed.

  - The amino acid with the cumulative highest fitness experimentally estimated with DMS only has 49% concordance with the observed alignment.

  - In contrast, the SelAC estimate has 99% concordance (Figure ??).

- Estimates of selection coefficients do not represent evolution.

    * Due to artificial selection environment; Heterogeneous population, very large $s$.

    * Only one antibiotic used, maybe a mixture of antibiotics would better reflect natural evolution.

    * Lack of repeatability between labs introduces further problems (Firnberg et al 2014 vs. Stifler et al. 2016).

- Assuming that the DMS selection inference adequately reflects natural evolution, the observed TEM sequences are either mal-adapted or where unable to reach a fitness peak.

    - *E. coli* has a large effective population size, estimates are on the order of $10^8$ to $10^9$ (Ochman and Wilson 1987, Hartl et al 1994).

    - The large $N_e$ would allow *E. coli* to effectively "explore" the sequence space, thus suggesting that the TEM sequences are mal-adapted according to the DMS estimates.

    - Our simulations of sequence evolution with various $N_e$ values and the DMS fitness values in contrast show that we would expect higher adaptation even with much smaller $N_e$ (Figure **??**).

- Estimates of selection coefficients do not represent evolution.

    - Due to artificial selection environment; Heterogeneous population, very large $s$.

    - Only one antibiotic used, maybe a mixture of antibiotics would better reflect natural evolution.

    - Lack of repeatability between labs introduces further problems (Firnberg et al 2014 vs. Stifler et al. 2016).

– Still better than models without site specific equilibrium frequencies.

• DMS estimates of the observed TEM variants predict them to be mal-adapted while
SelAC predicts most TEM variants to be well adapted.

– Given *E. coli*'s large effective population size, the efficacy of selection should be
very large.

– We therefore expect the observed sequence variants to be at the selection-mutation-
drift barrier, which in turn can expected to be near the optimum.

– We find the majority of sequences near the optimum, therefore the SelAC esti-
mates are consistent with theoretical population genetics results.

– In contrast, finding strong selection against the observed TEM variants indicates
that DMS is not consistent with theoretical population genetics expectations.

– This is consistent when thinking about that DMS only reflects the selection on
the TEM sequence with regards to one antibiotic, which seems appropriate to
model selection in modern hospital environments but not when the interest lies
in the natural evolution of TEM.

• We find that SelAC produces similar selection against the observed TEM variants if
we assume the fitness peaks (optimal AA) that are estimated by DMS.

– This shows that DMS and SelAC can provide consistent estimates of selection
against amino acids.

– SelAC has the advantage that it can be applied to any protein coding sequence
alignment.

– This removes the need for extrapolation e.g. from TEM to SHV.

• SelAC has the advantage that it can be applied to any protein coding sequence align-
ment.

177     − This removes the need for extrapolation e.g. from TEM to SHV.

178 • Difference in selection parameters between TEM and SHV indicate that extrapolation
179     is not a good idea.

180     − The difference in the site specific strength of selection shows that TEM and SHV
181       are facing different selection pressures.

182     − this is also highlighted by the differences in physicochemical weightings between
183       the two proteins.

184 • SelAC outperforms DMS, but is not without flaws itself

185     − Like DMS and most phylogenetic models, SelAC assumes site independence.

186     − SelAC is a model of stabilizing selection, in contrast to e.g. GY94 which is a
187       model of frequency dependent selection.

188       * Since TEM plays a role in the chemical warfare with conspecifics and other
189         microbes, some sites may be under negative frequency dependent selection.

190     − SelAC assumes the same G distribution across all sites.

191       * Different G distribution for each type of secondary structure

192       * active sites may not follow distribution.

193     − SelAC assumes that selection is proportional to distance in physicochemical space.

194       * We used Grantham (1974) properties, however many other distances are avail-
195         able which may an even better model fit.

196 • Low sequence variation in the TEM may be cause for concern as it could be misinter-
197     preted by the model as stabilizing selection because of the short branches.

198     − However, provided our simulations support that TEM is actually under stabilizing
199       selection

- In conclusion, DMS experiments have been proposed to supplement information on selection on amino acids in phylogenetic studies.

  - This study shows that information on selection can be extracted from alignments of protein coding sequences.

  - This highlights the limitations of DMS to explain natural evolution.