

**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}

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

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

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

Version dated: October 5, 2018

Introduction

- Incorporating selection into phylogenetic frameworks is already a long lasting endeavor.
 - Phylogenetic inference of sequence relationship was long focused on rates of substitutions.
 - Focus has shifted towards site specific equilibrium frequencies (HB98, Bloom2014, ...) in the last 20 years.
 - Such models however, tend to be unfeasible as they are very parameter rich.
 - The type of selection on a protein is not always clear, or differs between proteins phylogenetic models also have to make generalizing assumptions.
 - Incorporating selection from experimental sources therefore seems like an attractive option.
 - Incorporating empirical fitness has some important features.
 - * It allows for site specific amino acid preferences, acknowledging the heterogeneity of selection along the protein sequence.
 - * It greatly reduces the number of parameters that have to be estimated from the data.
 - * It allows for the fitting more complex models
 - However, the incorporation of empirical fitness also has some important shortcomings.
 - * Loss of generality.
 - * DMS experiments are limited to proteins and organisms that can be manipulated under laboratory conditions.
 - * But even in the case of TEM, the applied selection pressure is limited to the defense against a specific antibiotic.

- * TEM, however, has evolved to compete against conspecifics and other microbes using secreted metabolites to gain an advantage.
 - * Furthermore, DMS relies on a library of mutants and therefore on a heterogeneous population with competing genotypes.
 - * Therefore, it is important to ask how adequate such experiments reflect natural evolution.
- In this study we will assess how adequate DMS inference of site specific selection on amino acids, using TEM and provide an alternative, more generally applicable solution.
 - Simulations using DMS inferred site specific selection on amino acids show that observed TEM variants are unexpected; revealing the inadequacy of DMS.
 - Models fits achieved by the incorporation of DMS experiments can be improved upon using a hierarchical phylogenetic framework of stabilizing selection: SelAC.
 - Extrapolating site specific selection on amino acids between sequences (TEM and SHV) with related function can be inadequate.

Results

Site Specific Selection on Amino Acids Improves Model Fit

We compared the models *phydms* and *SelAC*, models of stabilizing site specific amino acid selection, to 281 other codon and nucleotide models by fitting them to the β -lactamase TEM. Models with site specific selection on amino acids improved model fits by at least 917 AICc units over codon or nucleotide models without site specific selection (Table ??tab:AICc)). *SelAC* estimates 263 site specific parameters, $\sim 5\%$ of the 4997 parameters necessary to estimate the site specific selection on amino acids. In contrast, *phydms* utilizes site specific selection on amino acids estimated from deep mutation scanning experiments.

Model	L	n	AIC	Δ AIC
SelAC	-1498	374	3744	0
SelAC+DMS	-1768	111	3758	14
phyDMS	-2060	105	4331	586

Table 1: L , number of model parameters n , AIC, and Δ AIC., Full table has > 200 models

Laboratory inferences of selection are inconsistent with observed sequences.

Improved model fits with phydms are deceiving. The site specific selection inferred by the deep mutation scanning experiment is inconsistent with the observed TEM sequences. We find that the sequence of selectively favored amino acids has only 49 % sequence similarity with the observed consensus sequence (Figure 3).

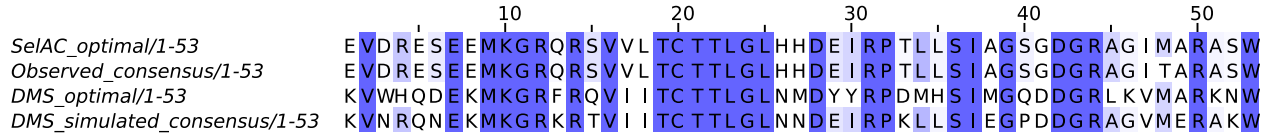


Figure 1: Every 5th residue. DMS and simulation based on DMS do not reflect natural sequences

Simulating under the experimentally inferred fitness landscape reveals that the observed sequences are unlikely to occur. Sequences simulated using a wide range of effective population sizes N_e show that we expect a sequence similarity of $\sim 70\%$ (Figure 4a). Similarly, the observed substitutional load is a twice the substitutional load of the sequences simulated under the experimentally inferred fitness landscape (Figure 4b).

SelAC better explains observed Sequences

Site specific estimates of Selection on Amino Acids

- Model selection shows that DMS can improve phylogenetic inference.
 - phyDMS improved model fit to 49 TEM sequences by 917 AICc units

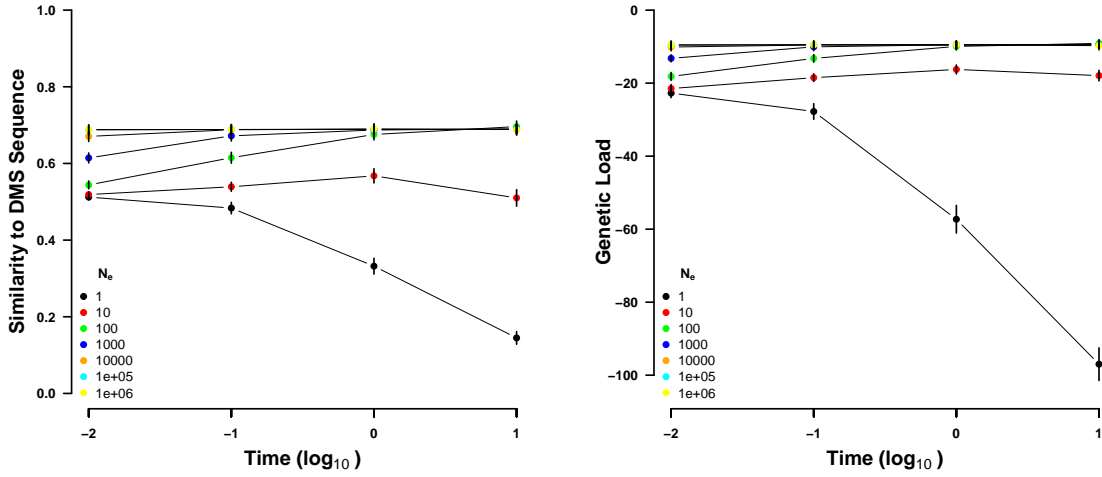


Figure 2: Sequences simulated under various values of N_e and for various times.

- Number of parameters estimated from data comparable to GY94 and others despite complex description of fitness landscape because of experimental estimates.
- Lab inferences of selection (DMS) are inconsistent with observed sequences.
 - The inferred fitness landscape is inconsistent with observed sequences.
 - * The optimal amino acid sequence inferred by DMS only shows 49% sequence similarity with the observed sequences (Figure 3).
 - Observed sequences unlikely under the lab inferred fitness landscape (Figure 4a,b).
 - * We would expect about half of the observed fitness burden.
 - * Sequence similarity is expected to be about $\sim 70\%$.
- SelAC better explains observed sequences than DMS and other models.
 - Model selection shows that SelAC outperforms phydms (Table 2).
 - Model adequacy assessed by sequence similarity shows that SelAC better represents the observed sequences.
- Application of SelAC to TEM.

- Site specific estimates of aa fitness (Figure 5).
 - * Fitness burden based on SSE.
 - * Fitness burden at binding sites
 - * Most sites show the estimated optimal amino acid.
 - * We find that selection against used amino acids is clustered and locally confined.
 - * Role of secondary structures?
- Application of SelAC to TEM and comparison to TEM
 - * Site specific G terms for TEM and SHV are only weakly correlated ($\rho = 0.17$), despite similar α_G (Figure 6a).
 - * Greatest difference is observed in the physicochemical properties, specifically α (which PC is that?) (Figure 6b).

Discussion

- Given
- 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 2).
- 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 3).
- 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 4).
- 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).

- 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 be expected to be near the optimum.
 - We find the majority of sequences near the optimum, therefore the SelAC estimates 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 alignment.
 - This removes the need for extrapolation e.g. from TEM to SHV.

- Difference in selection parameters between TEM and SHV indicate that extrapolation is not a good idea.
 - The difference in the site specific strength of selection shows that TEM and SHV are facing different selection pressures.
 - this is also highlighted by the differences in physicochemical weightings between the two proteins.
- SelAC outperforms DMS, but is not without flaws itself
 - Like DMS and most phylogenetic models, SelAC assumes site independence.
 - SelAC is a model of stabilizing selection, in contrast to e.g. GY94 which is a model of frequency dependent selection.
 - * Since TEM plays a role in the chemical warfare with conspecifics and other microbes, some sites may be under negative frequency dependent selection.
 - SelAC assumes the same G distribution across all sites.
 - * Different G distribution for each type of secondary structure
 - * active sites may not follow distribution.
 - SelAC assumes that selection is proportional to distance in physicochemical space.
 - * We used Grantham (1974) properties, however many other distances are available which may an even better model fit.
- Low sequence variation in the TEM may be cause for concern as it could be misinterpreted by the model as stabilizing selection because of the short branches.
 - However, provided our simulations support that TEM is actually under stabilizing selection
- In conclusion, DMS experiments have been proposed to supplement information on selection on amino acids in phylogenetic studies.

Model	L	n	AIC	Δ AIC
SelAC	-1498	374	3744	0
SelAC+DMS	-1768	111	3758	14
phyDMS	-2060	105	4331	586

Table 2: L , number of model parameters n , AIC, and Δ AIC., Full table has > 200 models

- 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.

Tables

Figures

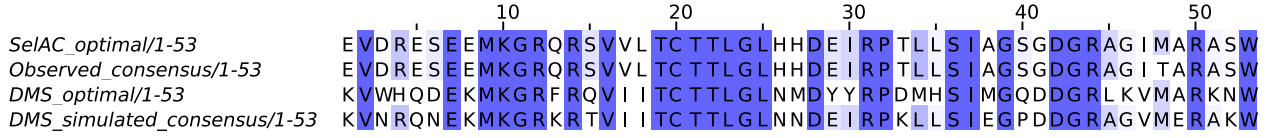


Figure 3: Every 5th residue. DMS and simulation based on DMS do not reflect natural sequences

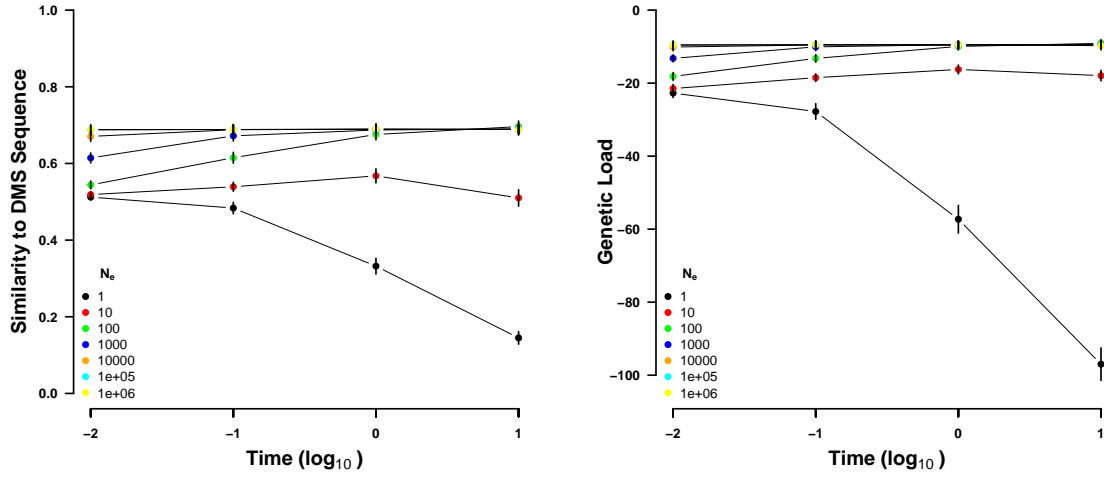


Figure 4: Sequences simulated under various values of N_e and for various times.

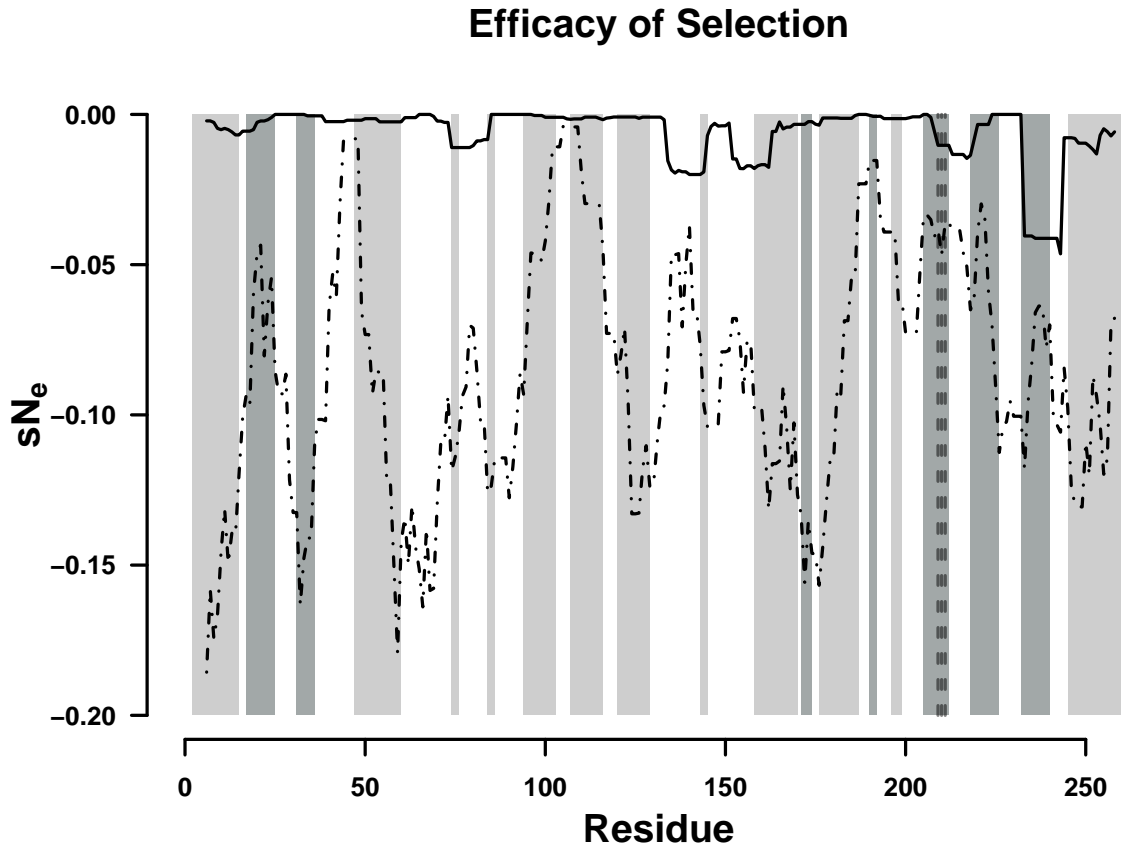


Figure 5: TEM, bars are different secondary structure elements. Dashed dotted line is DMS, solid is SelAC sN_e , all lines are means of all sequences, sliding window of 10 sites. vertical lines are active/binding sites.,

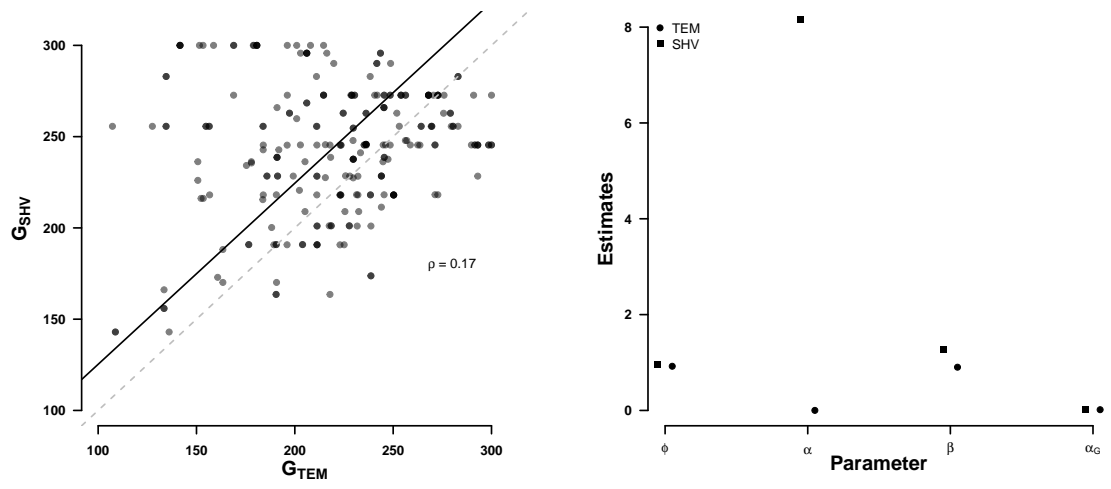


Figure 6: Comparisson of selection related parameters between TEM and SHV.

Suppl. Figures

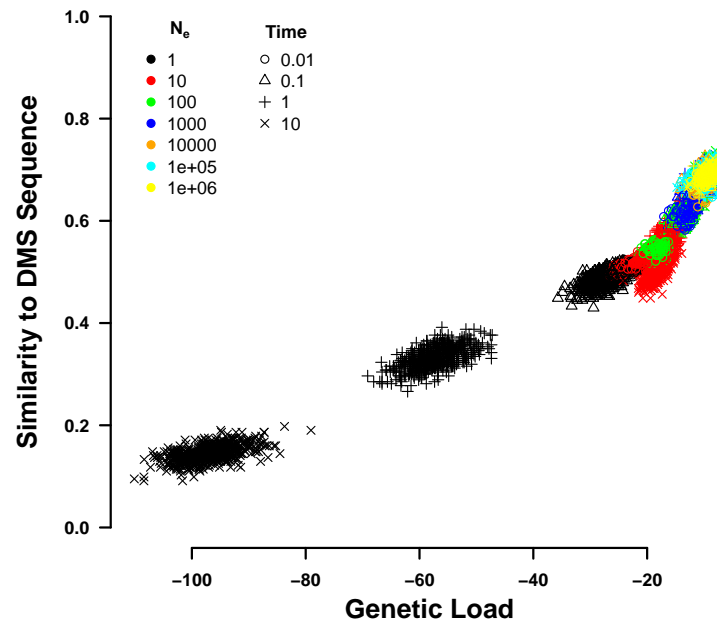


Figure 7: Suppl: Sequences simulated under various values of N_e and for various times.
TODO: replace clouds by mean+sd bars

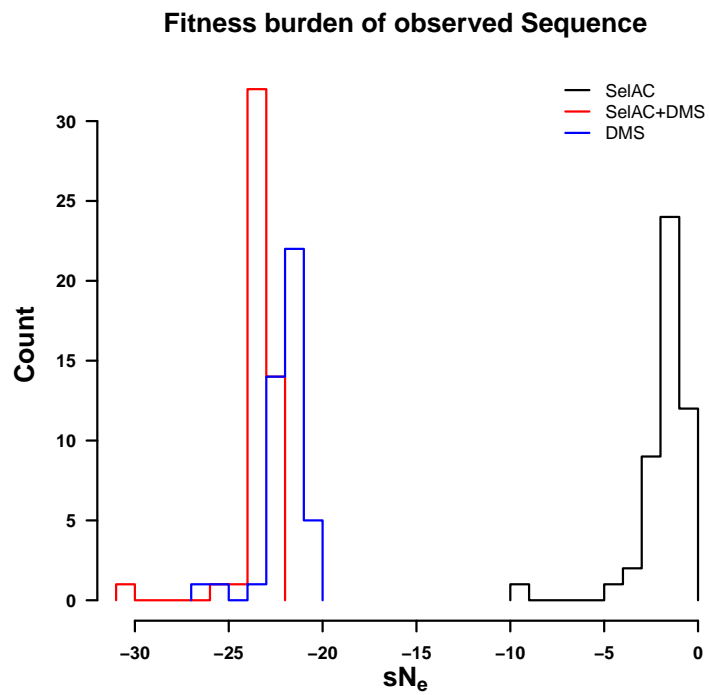


Figure 8: Suppl: sN_e of whole sequence, variation across tips. TEM