

2 **Predicting amino acid functionality from sequence**
3 **data in a phylogenetic framework.**

4 **Abstract**

5
6 CEDRIC LANDERER^{1,2,*}, BRIAN C. OMEARA^{1,2}, AND MICHAEL
7 A. GILCHRIST^{1,2}

8 ¹Department of Ecology & Evolutionary Biology, University of Tennessee, Knoxville, TN 37996-
9 1610

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

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

Version dated: June 4, 2018

Outline

Introduction

- The introduction of selection into phylogenetic frameworks has been a long effort, with limited success as these frameworks are very parameter rich.
- Another shortcoming of these models is the uniform stationary distribution of amino acids, making each one equally likely at a position.
- Insert brief review of methods, what distinguishes them and what do they share, that is relevant to DMS/SelAC.
- The most popular tools, however, are still based purely on the mutation process (RaxML, RevBayes).
- A novel take on the incorporation of selection on a protein is the independent estimation of fitness effects.
- DMS experiments provide site specific fitness values on synonymous and non-synonymous mutations (focus on non-synonymous).
- This limits the number of estimated parameters greatly and allows for computationally feasible models.
- However, the information on selection gained by DMS experiments is limited to single proteins of organisms that can be manipulated in the laboratory with short generation times.
- SelAC on the other hand, is a mechanistic model with relatively few parameters, explicitly modeling the functionality of a gene by estimating site specific amino acid preferences from the sequence data.

- SelAC has multiple advantages over other models incorporating selection into a phylogenetic framework, as it does not assume a uniform stationary amino acid distribution which allows for the estimation of the preferred amino acid at a site, estimates relatively few parameters and does not depend on experimental data.
- In this study, we compare the quality of phylogenetic estimates obtained utilizing DMS experiments to estimates from SelAC.
- We utilize DMS experiments from Firnberg (2014) and Stiffler (2016) for the TEM β -lactamase of *E. coli*.
- We compare model fit and adequacy of the DMS and SelAC amino acid preferences using SelAC and phydms.
- phydms is a tool explicitly designed to utilize selection information from DMS experiments.
- We show that DMS experiments do not accurately reflect natural evolution of protein sequences.
- We find that amino acid preferences estimated with SelAC provides better model fit and higher model adequacy than DMS experiments.
- We show that phylogenetic models can extract information on amino acid preference and do not require it as input.

Results

- Compare DMS from Firnberg and Stiffler to SelAC and majority under SelAC and phydms
 - Comparison of Firnberg under SelAC for TEM and SHV (three sequences: DMS, Majority, SelAC)

- Comparison of Frinberg under phydms for TEM and SHV (three sequences: DMS, Majority, SelAC)
- Comparison of Stiffler under SelAC for TEM and SHV (three sequences: DMS, Majority, SelAC)
- Comparison of Stiffler under phydms for TEM and SHV (three sequences: DMS, Majority, SelAC)
- Comparison of preferred sequence
 - Simulations of sequences under each preferred sequence.
 - Only majority rule (duh) and SelAC agree with observed sequences.
- SelAC is dependent on choice of PC properties to produce amino acid rankorder and assumes stabilizing selection.
 - Rankorder of certain sites can not be produced by any of the PC checked (no combination checked)

Discussion

- SelAC sequence outperforms DMS experiments, reflecting evolution better than DMS sequences under artificial selection pressure.
- SelAC only uses preferred state as input, no information about 2nd or third preferred amino acid.
- The reduction of a DMS experiment to this state might be considered an unfair comparison, however, we tested the sequences under phydms (no reduction of information), with the same result.
- This also means that SelAC produces the same information a DMS experiment would, but for naturally evolving sequences and can be applied to any sequence.

- TEM/SHV have not evolved to combat specific human developed antibiotics, but as means of "warfare" between bacteria (need more reading here).
- This could be the cause for the great difference between DMS and observed sequences.
- SelAC, however can not provide any information about antibiotic resistency, making DMS very valuable, but not for phylogenetics.
- but additional tip information could be combined with SelAC to get at this information (out of scope? future directions?).

Introduction

Materials & Methods

Results

Discussion

Supplemental Material