

# Application of mechanistic models to separate the effects of mutation, selection, and drift on protein sequence evolution

A Dissertation Presented for the  
Doctor of Philosophy  
Degree

The University of Tennessee, Knoxville

Cedric Lars Florian Landerer

December 2018

© by Cedric Lars Florian Landerer, 2018  
All Rights Reserved.

*To my mother*

## Acknowledgments

I am grateful for the many people at the University of Tennessee and in Knoxville that that made my time here such a pleasure. First and foremost I want to thank my Adviser, Dr. Michael Gilchrist for his long lasting patience, his availability and his teachings; always sharpening my focus and providing a new angle to a problem. Great thanks also goes to my committee Dr. Benjamin Fitzpatrick, Dr. Brian O'Meara, and Dr. Russel Zaretzki as they were always available for questions and discussions and for their great guidance. In particular Brian O'Meara who always had an open door and tolerated my frequent visits. None of the work presented in this dissertation would have been possible without their great guidance. For many great discussions and never a dull moment in the office I also have to thank my labmate Alex Cope. I also have to thank the faculty and students in Ecology and Evolutionary Biology, allowing me to broaden my knowledge and insights with always stimulating discussions and for moral support. Specially John Reese, Cassie Dresser, Liam Muller, Athmanathan Senthilnathan, Harmony Yomai, and Jim Fordyce. Thanks also goes to my former roommate Cassie Watters without whom my stay in Knoxville would have been a lot less exciting.

*Nothing in Biology Makes Sense Except in the Light of Evolution.*

-Theodosius Dobzhansky

*Nothing in evolutionary biology makes sense except in the light of population genetics*

-Michael Lynch

## Abstract

Mathematical and statistical models are useful for describing and understanding observations in genetics and genomics. These models have to constantly be updated to reflect current biological understanding. As opposed to descriptive and phenomenological models, mechanistic models allow for the extraction of more biologically relevant information based on underlying principles. Mutation, selection, and genetic drift are the three forces guiding evolution. Mechanistic models rooted in population genetics principles allow us to determine how these forces shape observed data. I demonstrate the usage of mechanistic models to relate protein coding sequences to their fitness landscapes and the evolutionary forces shaping them. Using the yeast *L. kluyveri*, I show the increased cost of protein synthesis due to a large scale introgression with mismatched codon usage. Furthermore, I analyze site-specific selection on amino acids in the beta-lactamase protein TEM, which confers antibiotic resistance in *E. coli* and related species.

## Table of Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Cost: Decomposing Codon Usage . . . . .	3
1.2	Benefit: Selection on Amino acids . . . . .	5
<b>2</b>	<b>AnaCoDa: Analyzing Codon Data with Bayesian mixture models</b>	<b>7</b>
2.1	Abstract . . . . .	8
2.2	Introduction . . . . .	10
2.3	Features . . . . .	11
2.4	Appendix: Supplementary Material . . . . .	14
2.4.1	The AnaCoDa framework . . . . .	14
2.4.2	AnaCoDa setup . . . . .	15
2.4.3	File formats . . . . .	22
2.4.4	Analyzing and Visualizing results . . . . .	24
<b>3</b>	<b>Decomposing mutation and selection to identify mismatched codon usage</b>	<b>34</b>
3.1	Abstract . . . . .	35
3.2	Introduction . . . . .	37
3.3	Results . . . . .	39
3.3.1	The Signatures of two Cellular Environments within <i>L. kluveri</i> 's Genome . . . . .	39
3.3.2	Comparing Differences in the Endogenous and Exogenous Codon Usage	40
3.3.3	Determining Source of Exogenous Genes . . . . .	42

3.3.4	Estimating Introgression Age . . . . .	43
3.3.5	Genetic Load due to Mismatching Codon Usage of the Exogenous Genes	44
3.4	Discussion . . . . .	45
3.5	Materials and Methods . . . . .	49
3.5.1	Separating Endogenous and Exogenous Genes . . . . .	49
3.5.2	Model Fitting with ROC SEMPPR . . . . .	49
3.5.3	Comparing Codon Specific Parameter Estimates . . . . .	50
3.5.4	Syntenly Comparison . . . . .	50
3.5.5	Estimating Age of Introgression . . . . .	50
3.6	Acknowledgments . . . . .	53
3.7	Appendix: Supplementary Material . . . . .	54
<b>4</b>	<b>Phylogenetic model of stabilizing selection is more informative about site specific selection than extrapolation from laboratory estimates</b>	<b>64</b>
4.1	Abstract . . . . .	65
4.2	Introduction . . . . .	66
4.3	Results . . . . .	68
4.3.1	Site Specific Stabilizing Selection on Amino Acids Improves Model Fit	68
4.3.2	Laboratory Inferences Inconsistent with Observed Sequences. . . . .	72
4.3.3	Stabilizing Selection for Optimal Physicochemical Properties Improves Model Adequacy . . . . .	73
4.3.4	Estimating Site Specific Selection on Amino Acids . . . . .	76
4.4	Discussion . . . . .	79
4.5	Materials and Methods . . . . .	83
4.5.1	Phylogenetic Inference and Model selection . . . . .	83
4.5.2	Sequence Simulation . . . . .	84
4.5.3	Estimating site specific efficacy of selection $G$ . . . . .	84
4.5.4	Estimating site specific fitness values $w_i$ . . . . .	84



4.5.5	Model Adequacy . . . . .	85
4.6	Acknowledgments . . . . .	86
4.7	Appendix: Supplementary Material . . . . .	87
<b>5</b>	<b>Conclusion</b>	<b>98</b>
5.1	Summary . . . . .	98
5.1.1	The Value of Mechanistic Models . . . . .	99
5.1.2	Mechanistic Models Supplement Experiments . . . . .	100
5.2	Estimating Protein Functional and Fitness Landscape . . . . .	100
5.2.1	The Importance of Translation Errors . . . . .	100
5.2.2	Homogeneous Selection . . . . .	102
	<b>Bibliography</b>	<b>103</b>
	<b>Vita</b>	<b>118</b>

## List of Tables

3.1	Model selection of the two competing hypothesis. Reported are the log-likelihood, $\log(\mathcal{L})$ , the number of parameters estimated $n$ , AIC, and $\Delta\text{AIC}$ values. . . . .	39
3.2	Synonymous codon preference in the various data sets based on our estimates of $\Delta M$ . . . . .	54
3.3	Synonymous codon preference in the various data sets based on our estimates of $\Delta\eta$ . . . . .	55
3.4	. . . . .	56
4.1	Model selection, shown are the three models of stabilizing site specific amino acid selection ( <i>SelAC</i> , <i>SelAC</i> +DMS, <i>phydms</i> ) and the best performing codon and nucleotide model (GOLDMAN and YANG, 1994; ZHARKIKH, 1994). Reported are the log-likelihood $\log(\mathcal{L})$ , the number of parameters estimated $n$ , AIC, $\Delta\text{AIC}$ , $\text{AICc}$ , and $\Delta\text{AICc}$ values. See Table 4.3 for results from all models we tested. . . . .	69
4.2	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. . . . .	77
4.3	Model selection of 230 models of nucleotide and codon evolution. . . . .	87

## List of Figures

1.1	ROC SEMPFR model behavior for Isoleucine. The proportion of each codon observed changes with protein synthesis rate. Mutation is dominant when protein synthesis rate is low, mutationally favored codons are observed with the highest frequency. With the increase of protein synthesis rate, the influence of selection increases until the system is dominated by selection. The selectively favored codon is observed with the highest frequency. . . . .	3
1.2	Decline in fitness with distance in physicochemical space from the optimal amino acid. Fitness decline of amino acids (black dots) relative to optimal amino acid (Alanine). Weighting of physicochemical properties according to GRANTHAM (1974). The full fitness surface can be described but only 20 discrete amino acid states are available for selection to act on. . . . .	5
2.1	Distribution of $s$ for codon GCA for amino acid alanine. Dashed line indicates the CAI weight for GCA. The comparison provides a more nuanced picture as we can see that the selection on GCA varies across the genome. . . . .	27
2.2	Trace plot showing the traces of all 40 codon specific selection parameters $\Delta\eta$ organized by amino acid. . . . .	29
2.3	Trace plot showing the protein synthesis trace $\phi$ for gene 669. . . . .	30
2.4	Trace plot showing the $\log(\text{posterior})$ trace for the current model fit. Window inset shows the last 1.000 samples . . . . .	31

2.5	Fit of the ROC model for a random yeast. The solid line represent the model fit from the data, showing how synonymous codon frequencies change with gene expression. The points are the observed mean frequencies of a codon in that synthesis rate bin and the whisks indicate the standard deviation within the bin. The codon favored by selection is indicated by a ”*”. The bottom right panel shows how many genes are contained in each bin . . . . .	32
2.6	Comparison of the selection parameter of seven yeast species estimated with ROC-SEMPPR. . . . .	33
3.1	Comparison of predicted protein synthesis rate $\phi$ to microarray data from <i>TSANKOV et al. (2010)</i> for (a) the combined genome and (b) the separated endogenous and exogenous genes. Endogenous genes are displayed in black and exogenous genes in red. Black line indicates type II regression line ( <i>SOKAL and ROHLF, 1981</i> ). . . . .	40
3.2	Comparison of (a) mutation bias $\Delta M$ and (b) selection bias $\Delta\eta$ parameters for endogenous and exogenous genes. Estimates are relative to the mean for each codon family. Black dots indicate $\Delta M$ or $\Delta\eta$ parameters with the same sign for the endogenous and exogenous genes, red dots indicate parameters with different signs. Black line shows the type II regression line ( <i>SOKAL and ROHLF, 1981</i> ). Dashed lines mark quadrants. . . . .	41
3.3	Correlation coefficients of $\Delta M$ and $\Delta\eta$ of the exogenous genes with 38 examined yeast lineages. Dots indicate the correlation of $\Delta M$ and $\Delta\eta$ of the lineages with the endogenous and exogenous parameter estimates. All regressions were performed using a type II regression ( <i>SOKAL and ROHLF, 1981</i> ). . . . .	43
3.4	Genetic load $s = \Delta\eta\phi$ (a) at the time of introgression ( $\kappa = 5$ ), and (b) currently ( $\kappa = 1$ ). . . . .	45

3.5	Correlation coefficient of $\Delta M$ and $\Delta\eta$ of the endogenous genes with 38 examined yeast lineages. Dots indicate the correlation of $\Delta M$ and $\Delta\eta$ of the lineages with the endogenous and exogenous parameter estimates. All regressions were performed using a type II regression line (SOKAL and ROHLF, 1981).	57
3.6	Comparison of (a) mutation bias $\Delta M$ and (b) selection bias $\Delta\eta$ parameters for endogenous genes and combined gene sets. Estimates are relative to the mean for each codon family. Black dots indicate $\Delta M$ or $\Delta\eta$ parameters with the same sign for the endogenous and exogenous genes, red dots indicate parameters with different signs. Black line shows the type II regression line (SOKAL and ROHLF, 1981). Dashed lines mark quadrants.	58
3.7	Synteny relationship of <i>E. gossypii</i> and the exogenous genes. Indicated is the GC content along the introgression.	59
3.8	Amount of synteny for each species in units of standard deviations for selected species.	60
3.9	Genetic load (left) without scaling of $\phi$ per gene, and change of total genetic load with scaling $\kappa$ between <i>E. gossypii</i> and <i>L. kluyveri</i> (right)	61
3.10	Total amount of adaptation estimated to have occurred between time of introgression and currently observed per gene.	62
3.11	Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dashed line indicates the exogenous codon usage, dotted line indicates the combined codon usage.	63

4.1	Phylogenies resulting from <i>SelAC</i> , <i>SelAC</i> +DMS, <i>phydms</i> , and <i>GY94</i> . As <i>SelAC</i> is currently too slow for the inference of topologies, the topology for the <i>SelAC</i> phylogenies was inferred using the codon model of KOSIOL <i>et al.</i> (2007).	71
4.2	Alignment of TEM optimal and simulated sequences. Indicated is the percentage identity at each site. . . . .	72
4.3	Sequences simulated from the ancestral state under the site specific selection on amino acids estimated using deep mutation scanning. (left) Sequence similarity to the observed consensus sequence at various times for a range of values of $N_e$ . (right) Genetic load of the simulated sequences at various times for a range of values of $N_e$ . Time is given in number of expected mutations per site, which equals the substitution rate of a neutral mutation. Points indicate sample means and vertical bars indicate standard deviations. Initial sequence is the inferred ancestral state of the TEM variants and indicated by a black triangle. . . . .	74
4.4	Sequences simulated from the ancestral state under the site specific selection on amino acids estimated using <i>SelAC</i> . (left) Sequence similarity to the observed consensus sequence at various times for a range of values of $N_e$ . (right) Genetic load of the simulated sequences at various times for a range of values of $N_e$ . Time is given in number of expected mutations per site, which equals the substitution rate of a neutral mutation. Points indicate sample means and vertical bars indicate standard deviations. Initial sequence is the inferred ancestral state of the TEM variants and indicated by a black triangle.	75
4.5	Distribution of average site specific genetic load in TEM over all observed TEM variants. Average site specific genetic load is indicated by the black line. Light gray bars indicate where helices are found, and dark gray bars indicate $\beta$ -sheets. The three residues forming the binding site and the two residues forming the active are indicated by triangles at the top of the plot. .	78

4.6	Distribution of genetic load in TEM mapped on its structure (1xpb). Average genetic load over all observed TEM variants is indicated by the color, blue low, red medium, yellow high genetic load. Active site is indicated in black. .	93
4.7	Distribution of genetic load in SHV. Average genetic load over all observed SHV variants is indicated by the black line. Light gray bars indicate where helices are found, and dark gray bars indicate $\beta$ -sheets. The three residues forming the binding site and the two residues forming the active are indicated by triangles at the top of the plot. . . . .	94
4.8	Sequences simulated from a random codon sequence under the site specific selection on amino acids estimated using <i>SelAC</i> . (left) Sequence similarity to the observed consensus sequence at various times for a range of values of $N_e$ . (right) Genetic load of the simulated sequences at various times for a range of values of $N_e$ . Time is given in number of expected mutations per site, which equals the substitution rate of a neutral mutation. Points indicate sample means and vertical bars indicate standard deviations. Initial sequence is the inferred ancestral state of the TEM variants and indicated by a black triangle.	95
4.9	Sequences simulated from a random codon sequence under the site specific selection on amino acids estimated using deep mutation scanning. (left) Sequence similarity to the observed consensus sequence at various times for a range of values of $N_e$ . (right) Genetic load of the simulated sequences at various times for a range of values of $N_e$ . Time is given in number of expected mutations per site, which equals the substitution rate of a neutral mutation. Points indicate sample means and vertical bars indicate standard deviations. Initial sequence is the inferred ancestral state of the TEM variants and indicated by a black triangle. . . . .	96

4.10 Comparison of selection related parameters between TEM and SHV. (left)	
Estimated site specific efficacy of selection $G$ . (right) Selection related pa-	
rameter estimates. Protein functionality production rate $\psi$ , physicochemical	
weight for amino acid composition $\alpha_c$ , physicochemical weight for amino acid	
polarity $\alpha_p$ , and the parameter describing the distribution of $G$ , $\alpha_G$ estimated	
by <i>SelAC</i> . . . . .	97



# Chapter 1

## Introduction

Protein synthesis is the most costly metabolic process a cell performs (REEDS *et al.*, 1985; WATERLOW and MILLWARD, 1989; BUTTGEREIT and BRAND, 1995; WARNER, 1999; AKASHI and GOJOBORI, 2002; LINDQVIST *et al.*, 2018) causing selection to maximize the benefit of protein synthesis and performing it as efficiently as possible. Studying the ratio of cost to benefit of protein synthesis is, therefore, important to understand the evolution of protein coding sequences (GILCHRIST *et al.*, 2009; SHAH and GILCHRIST, 2011a; GILCHRIST *et al.*, 2015; BEAULIEU *et al.*, in review). However, the strength of selection varies greatly between genes, from low expression genes with codon usage dominated by mutation bias between nucleotides over highly expressed genes reflecting the dominance of selection for efficient translation of the mRNA, to selection on the amino acid composition required for the function of the protein.

We can formalize the cost and benefit of a protein coding sequence and formulate mathematical models. Mathematical and statistical models have long been used to describe or summarize observations in genetics and genomics. Often without addressing the underlying biological mechanisms - mutation, selection, and genetic drift - shaping DNA sequences, but as phenomenological descriptions. As researchers learn more about the underlying processes and more genetic and genomic data is available, the mathematical models that allow for the extraction of information from this data have to keep up. For example, after the unraveling of the degenerate genetic code by MATTHAEI and NIERENBERG

(1961); NIERENBERG and MATTHAEI (1961); MAXWELL (1962); LEDER and NIERENBERG (1964), and many others, researchers noticed that synonymous codons are not found in uniform proportions (FITCH, 1976; GRANTHAM *et al.*, 1980; IKEMURA, 1981; GRANTHAM *et al.*, 1981; SHARP *et al.*, 1988). Models of codon usage, however, were long purely descriptive and heuristic (IKEMURA, 1981; BENNETZEN and HALL, 1982; SHARP, 1987; WRIGHT, 1990). Similarly, phylogenetic models have long been phenomological (JUKES and CANTOR, 1969; DAYHOFF *et al.*, 1978; KIMURA, 1980; FELSENSTEIN, 1981; ALTSCHUL, 1991), describing the rate of change between states without regards for the forces guiding evolution, mutation, selection, and genetic drift. ZUCKERKANDL and PAULING (1962) proposed that the evolution of proteins is constant over time and between lineages before the genetic code was fully deciphered and at a time where protein synthesis was barely understood based on their observation that similarity on hemoglobin is correlated with divergence time. This work is therefore focused on the application of mechanistic models rooted in first principles and their application to protein coding sequences

Mechanistic models are used throughout biology (GOLDMAN and YANG, 1994; LAUREAU, 1998; DAVIS and PELSOR, 2001; DORON-FAIGENBOIM and PUPKO, 2007; MCGILL *et al.*, 2007). By modeling the process underlying the observed data mechanistic models provide insights into the processes and estimates of parameters shaping the data (LIBERLES *et al.*, 2013). A wide variety of information is stored in protein and protein coding sequences, e.g. structure (ANFINSEN, 1973), mutation bias (SHAH and GILCHRIST, 2011a; GILCHRIST *et al.*, 2015), protein synthesis rate (GILCHRIST, 2007; GILCHRIST *et al.*, 2015). Mechanistic models can be used to extract these informations and to study the relative strength of mutation, selection, and genetic drift leading to the observed sequences.

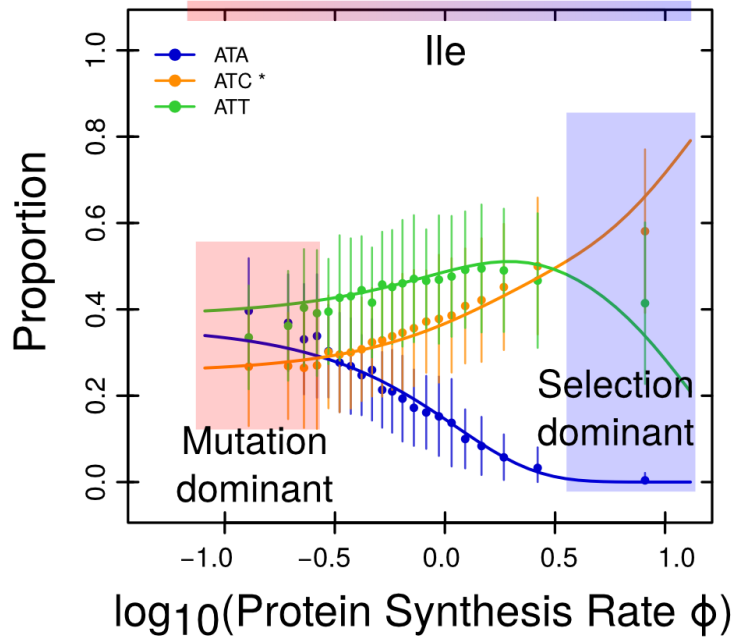


Figure 1.1: ROC SEMPPR model behavior for Isoleucine. The proportion of each codon observed changes with protein synthesis rate. Mutation is dominant when protein synthesis rate is low, mutationally favored codons are observed with the highest frequency. With the increase of protein synthesis rate, the influence of selection increases until the system is dominated by selection. The selectively favored codon is observed with the highest frequency.

## 1.1 Cost: Decomposing Codon Usage

Mutation bias on codon usage is a reflection of the cellular environment while selection on codon usage allows us to make inferences about the cellular and external environment a genome has evolved in. The relative strength of mutation and selection on individual genes varies, allowing us to separate mutation bias and selection, specifically selection against translation overhead cost (GILCHRIST, 2007; SHAH and GILCHRIST, 2011a; GILCHRIST *et al.*, 2015). Genes with low protein synthesis rates are thought to be under weak selection for codon usage and their codon usage is therefore dominated by mutation bias. In contrast, genes with high protein synthesis rates are thought to be under strong selection and their codon usage is therefore dominated by selection. However, mutation bias and selection can differ within the genome.

For example, strand specific mutation bias (LAFAY *et al.*, 1999; ROMERO *et al.*, 2000), differences in the tRNA pool throughout life stages (SAGI *et al.*, 2016), or introgressions and horizontal gene transfer (MDIGUE *et al.*, 1991; LAWRENCE and OCHMAN, 1997) can produce multiple genomic environments. Chapter 2 extends the mechanistic model ROC SEMPPR GILCHRIST *et al.* (2015) to allow for a mixture distribution of mutation and selection parameters LANDERER *et al.* (2018) and provides researchers with a software tool to address intra genomic variation in codon usage. However, there is a significant difference to classical mixture approaches. In addition to gene population specific parameters, ROC SEMPPR also estimates a gene specific parameter (protein synthesis rate). Therefore, the protein synthesis rate for each gene has to be estimated assuming that the a gene is in each gene population. This can provide additional insight into the adaptiveness of a gene to alternative genomic environments. Figure 1.1 illustrates how the proportions of synonymous codons change with increasing protein synthesis rate. When the protein synthesis rate is low, mutation bias between codons dominates the proportions of synonymous codons while increasing protein synthesis increases the strength of selection (see GILCHRIST *et al.* (2015) for details).

In chapter 3, I apply AnaCoDa to analyze the synonymous codon usage of the yeast *L. kluyveri* which experienced a large scale introgression replacing the whole left arm of chromosome C (FRIEDRICH *et al.*, 2015). I studied the differences in the parameters describing codon usage between the endogenous *L. kluyveri* genes and the introgressed exogenous genes. Recognizing the differences in codon usage between the endogenous and exogenous genes allowed me to improve prediction of protein synthesis rate, and separate the effects of mutation bias and selection in the endogenous *L. kluyveri* genes and the introgressed exogenous genes. This information was used to determine a potential donor lineage in *E. gossypii*, estimate the time since introgression, and estimate the genetic load of the introgression.

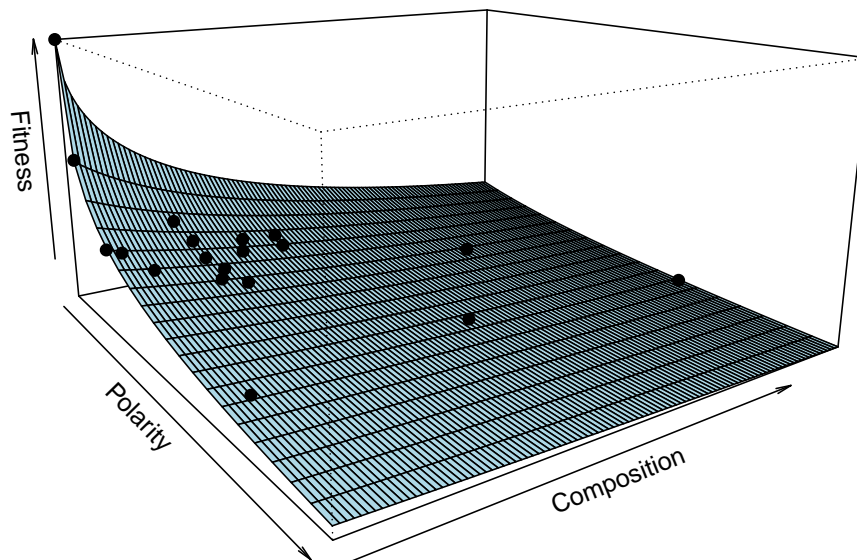


Figure 1.2: Decline in fitness with distance in physicochemical space from the optimal amino acid. Fitness decline of amino acids (black dots) relative to optimal amino acid (Alanine). Weighting of physicochemical properties according to [GRANTHAM \(1974\)](#). The full fitness surface can be described but only 20 discrete amino acid states are available for selection to act on.

## 1.2 Benefit: Selection on Amino acids

Genes are evolving with natural selection favoring proteins that encode their function optimally, with mutations and genetic drift reducing functionality. Amino acid preference and the relative strength of mutation, selection, and genetic drift usually varies between sites along the protein sequence. The number of parameters required to describe protein fitness increases exponentially with the length of the protein if interactions between sites are accounted for. Attempts to incorporate selection into phylogenetic approaches are, therefore, limited to site specific selection. The goal of chapter 4 is to estimate the strength of site specific selection on amino acids from protein coding sequences in a phylogenetic framework.

Ignoring interactions between sites allows to describe the site specific fitness landscape of a protein. Some approaches rely on the description of the full fitness landscape and therefore require  $19 \times L$ , where  $L$  is the length of the peptide in amino acids, parameters ([LARTILLOT and PHILIPPE, 2004](#); [LE \*et al.\*, 2008](#); [WANG \*et al.\*, 2008](#); [HOLDER \*et al.\*,](#)

2008; WU *et al.*, 2013; TAMURI *et al.*, 2014). As this is still a large number of parameters the incorporation of experimentally determined site specific selection on amino acids is an attractive alternative (BLOOM, 2014; THYAGARAJAN and BLOOM, 2014; BLOOM, 2017). Alternatively, assumptions about the nature of selection can reduce the number of parameters required. For example, negative frequency dependent selection (GOLDMAN and YANG, 1994; MUSE and GAUT, 1994; THORNE *et al.*, 1996) or stabilizing selection (BEAULIEU *et al.*, in review) allow for a reduction in fitness of amino acids with distance in physicochemical space.

*SelAC* (BEAULIEU *et al.*, in review), a model of stabilizing selection, assesses the fitness of each amino acid relative to the fitness peak (Figure 1.2). Fitness is assumed to decline exponentially with distance in physicochemical space to the optimal amino acid. In chapter 4 I apply *SelAC* to the  $\beta$ -lactamase TEM and estimate site specific selection on amino acids and compare the inferred fitness landscape to empirical estimates from deep mutation scanning experiments (STIFFLER *et al.*, 2016). I find that experimentally informed amino acid preferences improve model fit but do not accurately reflect the evolution of TEM. Furthermore, I show that the information on site specific selection on amino acids can be extracted from protein coding sequences by models rooted in first principles like *SelAC*.

## Chapter 2

**AnaCoDa: Analyzing Codon Data with Bayesian mixture models**

This chapter is a lightly revised version of a paper by the same name published in Bioinformatics and co-authored with Alexander Cope, Russell Zaretzki, and Michael A. Gilchrist.

C. Landerer, A. Cope, R. Zaretzki, M.A. Gilchrist, AnaCoDa: analyzing codon data with Bayesian mixture models, Bioinformatics, 34, 2018, 2496-2498

## 2.1 Abstract

**AnaCoDa** is an R package for estimating biologically relevant parameters of mixture models, such as selection against translation inefficiency, nonsense error rate, and ribosome pausing time, from genomic and high throughput datasets. **AnaCoDa** provides an adaptive Bayesian MCMC algorithm, fully implemented in C++ for high performance with an ergonomic R interface to improve usability. **AnaCoDa** employs a generic object-oriented design to allow users to extend the framework and implement their own models. Current models implemented in **AnaCoDa** can accurately estimate biologically relevant parameters given either protein coding sequences or ribosome foot-printing data. Optionally, **AnaCoDa** can utilize additional data sources, such as gene expression measurements, to aid model fitting and parameter estimation. By utilizing a hierarchical object structure, some parameters can vary between sets of genes while others can be shared. Genes may be assigned to clusters or membership may be estimated by **AnaCoDa**. This flexibility allows users to estimate the same model parameter under different biological conditions and categorize genes into different sets based on shared model properties embedded within the data. **AnaCoDa** also allows users to generate simulated data which can be used to aid model development and model analysis as well as evaluate model adequacy. Finally, **AnaCoDa** contains a set of



visualization routines and the ability to revisit or re-initiate previous model fitting, providing researchers with a well rounded easy to use framework to analyze genome scale data.

## **Availability:**

**AnaCoDa** is freely available under the Mozilla Public License 2.0 on CRAN (<https://cran.r-project.org/web/packages/AnaCoDa/>).

## 2.2 Introduction

**AnaCoDa** is an open-source software implemented in R ([R CORE TEAM, 2015](#)) that allows researchers to analyze genome-scale data like coding sequences and ribosome footprinting data using evolutionary or analytical models in a Bayesian framework. **AnaCoDa** was developed to analyze selection on synonymous codon usage in the form of ribosome overhead cost ([GILCHRIST \*et al.\*, 2015](#); [WALLACE \*et al.\*, 2013](#); [SHAH and GILCHRIST, 2011b](#)). However, other codon metrics like the codon adaptation index ([SHARP, 1987](#)) or the effective number of codons ([WRIGHT, 1990](#)) are also provided as reference. In addition, three currently unpublished models to analyze coding sequences for evidence of selection against nonsense errors and estimate ribosome pausing times from ribosome footprinting data are included. **AnaCoDa** implements an adaptive Gibbs sampler within a Metropolis-Hastings Monte Carlo Markov Chain (MCMC). This allows for the incorporation of prior knowledge such as observed gene expression levels and easy sampling from the posterior distribution to estimate parameter values and quantify degree of uncertainty. **AnaCoDa** provides a mixture distribution option to all implemented models, combining genes into sets by estimating the posterior probabilities of set membership based on gene-set specific parameters shared by all genes assigned to a given set. **AnaCoDa** provides a generic, mixture distribution option to all implemented models, allowing for the estimation of condition specific parameters or the automatic categorization of data into different sets based on differences in their posterior probabilities of set membership. In addition to the four models provided, **AnaCoDa** provides a modular infrastructure such that additional genome scale or even phylogenetic models can be integrated.

The **AnaCoDa** framework works with **AnaCoDa** requires gene specific data such as codon frequencies obtained from coding sequences or position specific footprint counts. Conceptually, **AnaCoDa** allows for three different types of parameters. The first type are gene specific parameters such as protein synthesis rate or relative functionality. The second type are gene-set specific parameters, such as mutation bias terms or translation error rates.

These parameters are shared across genes within a set and can be exclusive to a single set or shared with other sets. While the number of gene sets must be pre-defined by the user, set assignment of genes can be pre-defined or estimated as part of the model fitting. Estimation of the set assignment provides the probability of a gene being assigned to a set allowing the user to assess the uncertainty in each assignment. The third type are hyperparameters allowing for the construction and analysis of hierarchical model. Hyperparameters control the prior distribution for gene and gene-set specific parameters such as mutation bias or protein synthesis rate.

## 2.3 Features

**AnaCoDa** provides an interface written in R, a freely available programming language noted for its ease of use for even inexperienced programmers. As a result, **AnaCoDa** is accessible to researchers with minimal computational experience.

The interface of **AnaCoDa** is designed for quick and efficient data analysis. Generally, the only input needed for fitting a model to the data are protein-coding codon sequences in the form of a FASTA file or a flat-file containing codon counts obtained from ribosome foot-printing experiments. **AnaCoDa** also provides visualization functionality, including plotting functions to compare parameter estimates for different mixture distributions and display codon usage patterns. In addition, diagnostic functions such as those for calculating and visualizing the degree of autocorrelation in the parameter traces are provided.

### Robust and efficient model fitting

**AnaCoDa** has built-in features designed to improve the robustness and performance of the implemented MCMC approach. For example, the implemented MCMC automatically adapts the proposal width for sampled parameters such that a user defined acceptance range is met, improving sampling efficiency of the MCMC and computational performance. Even though

**AnaCoDa** is written in C++, analysis of large datasets and/or complex models can be very computationally intensive. To protect users from computer failures or aid in the collection of additional MCMC samples, **AnaCoDa** can periodically produce output checkpoint files, which can be used to restart an MCMC chain from a previous time point. In addition, **AnaCoDa** automatically thins all parameter traces - meaning only every  $k^{th}$  sample is kept - increasing the effective number of samples and reducing its memory footprint.

Although **AnaCoDa** is provided as an R package, the main computational work is implemented in C++. Because R does not provide native C++ support, Rcpp was employed to expose whole C++ classes as modules to R (EDELBUETTEL and FRANCOIS, 2011). Using Rcpp eliminates time consuming data transfers between the R environment and the C++ core during model fitting, resulting in improved computational performance and allows for a fully object-oriented code design (BOOCH, 1993). As expected, the runtime of **AnaCoDa** scales linearly with genome size and number of iterations, and scales polynomially with the number of mixture distributions in the data set. The polynomial increase in runtime with the number of mixture distributions is due to the necessity to condition the gene assignment on the estimation of gene specific parameters, such as, protein synthesis rate.

## Data Simulation

In addition to fitting the models to datasets, **AnaCoDa** can be used to generate simulated data sets as well. On their own, simulated datasets are useful for model development and analysis. Simulating data under different conditions allows the user to explore model behavior and explore theoretical scenarios. Different conditions can include the addition or elimination of parameters, or simply allowing a set of parameter values to vary. Fitting models to simulated data can provide insight into potential pitfalls or shortcomings when fitting observational data and can serve as the basis for evaluating model adequacy of a model fit to observational data (MI *et al.*, 2015). Significant differences between simulated

and observational data suggests the current set of parameters or the model as a whole fail to include or adequately represent biological mechanisms underlying the observed data.

## Available models

**AnaCoDa** currently provides codon models for analyzing genome scale data. The ROC model implements and extends the codon usage bias (CUB) models developed by [GILCHRIST \*et al.\* \(2015\)](#); [WALLACE \*et al.\* \(2013\)](#); [SHAH and GILCHRIST \(2011b\)](#), which can reliably estimate the strength of selection on ribosome overhead cost, mutation bias and allows for the inference of protein synthesis rates. This model allows for the separation of effects of mutation and selection based on gene ordering by protein synthesis rate, and the addition of a mixture distribution allows for gene clustering based on mutation bias and selection for translation efficiency. In addition to identifying the most efficient codons, ROC provides estimates of mutation bias allowing the approximation of mutation ratios between codons ([GILCHRIST \*et al.\*, 2015](#); [WALLACE \*et al.\*, 2013](#)).

The ability to estimate protein synthesis rates in the absence of empirical data is useful for investigating CUB of non-model organisms for which such data is lacking and enables the usage of protein synthesis rate in comparative frameworks or other analyses requiring protein synthesis rate information ([DUNN \*et al.\*, 2018](#)). Use of the mixture model allows for the investigation of CUB heterogeneity at the genome or gene level. Following the same framework, additional models included in **AnaCoDa** provide estimates of codon-specific nonsense errors rates (FONSE) and ribosome pausing times (PA and PANSE).

Parameters estimated with the evolutionary models ROC and FONSE represent evolutionary averages and do not depend on experimental conditions. In contrast, PA and PANSE estimate the distribution of biologically relevant parameters like ribosome pausing times along a gene from experimental data such as ribosome footprinting data. The distribution can be dependent (PANSE) or independent (PA) of evidence for nonsense errors in the data.

## 2.4 Appendix: Supplementary Material

AnaCoDa allows for the estimation of biologically relevant parameters like mutation bias or ribosome pausing time, depending on the model employed. Bayesian estimation of parameters is performed using an adaptive Metropolis-Hasting within Gibbs sampling approach. Models implemented in AnaCoDa are currently able to handle gene coding sequences and ribosome footprinting data.

### 2.4.1 The AnaCoDa framework

The AnaCoDa framework works with gene specific data such as codon frequencies or position specific footprint counts. Conceptually, AnaCoDa uses three different types of parameters.

- The first type of parameters are **gene specific parameters** such as gene expression level or functionality. Gene-specific parameters are estimated separately for each gene and can vary between potential gene categories or sets.
- The second type of parameters are **gene-set specific parameters**, such as mutation bias terms or translation error rates. These parameters are shared across genes within a set and can be exclusive to a single set or shared with other sets. While the number of gene sets must be pre-defined by the user, set assignment of genes can be pre-defined or estimated as part of the model fitting. Estimation of the set assignment provides the probability of a gene being assigned to a set allowing the user to assess the uncertainty in each assignment.
- The third type of parameters are **hyperparameters**, such as parameters controlling the prior distribution for mutation bias or error rate. Hyperparameters can be set specific or shared across multiple sets and allow for the construction and analysis of hierarchical models, by controlling prior distributions for gene or gene-set specific parameters.

## Analyzing protein coding gene sequences

AnaCoDa always requires the following four objects:

- **Genome** contains the codon data read from a fasta file as well as empirical protein synthesis rate in the form of a comma separated (.csv) ID/Value pairs.
- **Parameter** represents the parameter set (including parameter traces) for a given genome. The parameter object also hold the mapping of parameters to specified sets.
- **Model** allows you to specify which model should be applied to the genome and the parameter object.
- **MCMC** specifies how many samples from the posterior distribution of the specified model should be stored to obtain parameter estimates.

### 2.4.2 AnaCoDa setup

#### Application of codon model to single genome

In this example we are assuming a genome with only one set of gene-set specific parameters, hence `num.mixtures` = 1. We assign all genes the same gene-set, and provide an initial value for the hyperparameter  $s_\phi$ .  $s_\phi$  controls the lognormal prior distribution on the gene specific parameters like the protein synthesis rate  $\phi$ . To ensure identifiability the expected value of the prior distribution is assumed to be 1.

$$E[\phi] = \exp\left(m_\phi + \frac{s_\phi^2}{2}\right) = 1 \quad (2.1)$$

Therefore the mean  $m_\phi$  is set to be  $-\frac{s_\phi^2}{2}$ . For more details see [GILCHRIST \*et al.\* \(2015\)](#).

After choosing the model and specifying the necessary arguments for the MCMC routine, the MCMC is run

```
genome <- initializeGenomeObject(file = "genome.fasta")
```

```

parameter <- initializeParameterObject(genome = genome, sphl = 1,
                                     num.mixtures = 1,
                                     gene.assignment = rep(1, length(genome)))
model <- initializeModelObject(parameter = parameter, model = "R0C")
mcmc <- initializeMCMCObject(samples = 5000, thinning = 10,
                             adaptive.width=50)
runMCMC(mcmc = mcmc, genome = genome, model = model)

```

`runMCMC` does not return a value, the results of the MCMC are stored automatically in the `mcmc` and `parameter` objects created earlier.

**Please note that AnaCoDa utilizes C++ object orientation and therefore employs pointer structures. This means that no return value is necessary for such objects as they are modified within the the `runMCMC` routine. You will find that after a completed run, the `parameter` object will contain all necessary information without being directly passed into the MCMC routine. This might be confusing at first as it is not default R behavior.**

## Application of codon model to a mixture of genomes

This case applies if we assume that parts of the genome differ in their gene-set specific parameters. This could be due to introgression events or strand specific mutation difference, horizontal gene transfers or other reasons. We make the assumption that all sets of genes are independent of one another. For two sets of gene-set specific parameter with a random gene assignment we can use:

```

parameter <- initializeParameterObject(genome = genome,
                                     sphl = c(0.5, 2), num.mixtures = 2,
                                     gene.assignment = sample.int(2,
                                                                    length(genome), replace = T))
gene.assignment = sample.int(2, length(genome), replace = T)

```

To accommodate for this mixing we only have to adjust `sphl`, which is now a vector of length 2, `num.mixtures`, and `gene.assignment`, which is chosen at random here.



## Empirical protein synthesis rate values

To use empirical values as prior information one can simply specify an `observed.expression.file` when initializing the genome object.

```
genome <- initializeGenomeObject(file = "genome.fasta",  
                                observed.expression.file = "synthesis_values.csv")
```

These observed expression or synthesis values ( $\Phi$ ) are independent of the number of gene-sets. The error in the observed  $\Phi$  values is estimated and described by `sepsilon` ( $s_\epsilon$ ). The csv file can contain multiple observation sets separated by comma. For each set of observations an initial  $s_\epsilon$  has to be specified.

```
# One case of observed data  
sepsilon <- 0.1  
  
# Two cases of observed data  
sepsilon <- c(0.1, 0.5)  
  
# ...  
  
# Five cases of observed data  
sepsilon <- c(0.1, 0.5, 1, 0.8, 3)  
  
parameter <- initializeParameterObject(genome = genome, sphi = 1,  
                                       num.mixtures = 1,  
                                       gene.assignment = rep(1, length(genome)),  
                                       init.sepsilon = sepsilon)
```

In addition one can choose to keep the noise in the observations ( $s_\epsilon$ ) constant by using the `fix.observation.noise` flag in the model object.

```
model <- initializeModelObject(parameter = parameter, model = "ROC",  
                               fix.observation.noise = TRUE)
```

## Fixing parameter types

It can sometime be advantages to fix certain parameters, like the gene specific parameters. For example in cases where only few sequences are available but gene expression measurements are at hand we can fix the gene specific parameters to increase confidence in our estimates of gene-set specific parameters.

We again initialize the **genome**, **parameter**, and **model** objects.

```
genome <- initializeGenomeObject(file = "genome.fasta")
parameter <- initializeParameterObject(genome = genome, phi = 1,
                                       num.mixtures = 1,
                                       gene.assignment = rep(1, length(genome)))
model <- initializeModelObject(parameter = parameter, model = "ROC")
```

To fix gene specific parameters we will set the **est.expression** flag to **FALSE**. This will estimate only gene-set specific parameters, hyperparameters, and the assignments of genes to various sets.

```
mcmc <- initializeMCMCObject(samples, thinning=1,
                             adaptive.width=100, est.expression=FALSE,
                             est.csp=TRUE, est.hyper=TRUE, est.mix=TRUE)
```

If we would like to fix gene-set specific parameters we instead disable the **est.csp** flag.

```
mcmc <- initializeMCMCObject(samples, thinning=1,
                             adaptive.width=100, est.expression=TRUE,
                             est.csp=FALSE, est.hyper=TRUE, est.mix=TRUE)
```

The same applies to the hyper parameters (**est.hyper**),

```
mcmc <- initializeMCMCObject(samples, thinning=1,
                             adaptive.width=100, est.expression=TRUE,
                             est.csp=TRUE, est.hyper=FALSE, est.mix=TRUE)
```

and gene set assignment (**est.mix**).

```
mcmc <- initializeMCMCObject(samples, thinning=1,
                             adaptive.width=100, est.expression=TRUE,
                             est.csp=TRUE, est.hyper=TRUE, est.mix=FALSE)
```

We can use these flags to fix parameters in any combination.

### Combining various gene-set specific parameters to a gene-set description.

We distinguish between three simple cases of gene-set descriptions, and the ability to customize the parameter mapping. The specification is done when initializing the parameter object with the **mixture.definition** argument.

We encounter the simplest case when we assume that all gene sets are independent.

```
parameter <- initializeParameterObject(genome = genome,
                                       sphi = c(0.5, 2), num.mixtures = 2,
                                       gene.assignment = sample.int(2,
                                                                    length(genome), replace = T),
                                       mixture.definition = "allUnique")
```

The **allUnique** keyword allows each type of gene-set specific parameter to be estimated independent of parameters describing other gene sets.

In case we want to share mutation parameter between gene sets we can use the keyword **mutationShared**

```
parameter <- initializeParameterObject(genome = genome,
                                       sphi = c(0.5, 2), num.mixtures = 2,
                                       gene.assignment = sample.int(2,
                                                                    length(genome), replace = T),
                                       mixture.definition = "mutationShared")
```

This will force all gene sets to share the same mutation parameters.

The same can be done with parameters describing selection, using the keyword **selectionShared**

```
parameter <- initializeParameterObject(genome = genome,
                                       sphi = c(0.5, 2), num.mixtures = 2,
                                       gene.assignment = sample.int(2,
                                                                    length(genome), replace = T),
                                       mixture.definition = "selectionShared")
```

For more intricate compositions of gene sets, one can specify a custom  $n \times 2$  matrix, where  $n$  is the number of gene sets, to describe how gene-set specific parameters should be shared. Instead of using the **mixture.definition** argument one uses the **mixture.definition.matrix** argument.

The matrix representation of **mutationShared** can be obtained by

```
# [,1] [,2]
# [1,] 1 1
# [2,] 1 2
# [3,] 1 3
defMatrix <- matrix(c(1,1,1,1,2,3), ncol=2)
parameter <- initializeParameterObject(genome = genome,
                                       sphi = c(0.5, 2, 1), num.mixtures = 3,
                                       gene.assignment = sample.int(3,
                                                                    length(genome), replace = T),
                                       mixture.definition.matrix = defMatrix)
```

Columns represent mutation and selection, while each row represents a gene set. In this case we have three gene sets, each sharing the same mutation category and three different selection categories. In the same way one can produce the matrix for three independent gene sets equivalent to the **allUnique** keyword.

```
# [,1] [,2]
```

```
#[1,] 1 1
#[2,] 2 2
#[3,] 3 3
defMatrix <- matrix(c(1,2,3,1,2,3), ncol=2)
```

We can also use this matrix to produce more complex gene set compositions.

```
# [,1] [,2]
#[1,] 1 1
#[2,] 2 1
#[3,] 1 2
defMatrix <- matrix(c(1,2,1,1,1,2), ncol=2)
```

In this case gene set one and three share their mutation parameters, while gene set one and two share their selection parameters.

## Checkpointing

AnaCoDa does provide checkpointing functionality in case runtime has to be restricted. To enable checkpointing, one can use the function **setRestartSettings**.

```
# writing a restart file every 1000 samples
setRestartSettings(mcmc, "restart_file", 1000, write.multiple=TRUE)
# writing a restart file every 1000 samples
# but overwriting it every time
setRestartSettings(mcmc, "restart_file", 1000, write.multiple=FALSE)
```

To re-initialize a parameter object from a restart file one can simply pass the restart file to the initialization function:

```
initializeParameterObject(init.with.restart.file="restart_file.rst")
```

## Load and save parameter objects

AnaCoDa is based on C++ objects using the Rcpp ([EDELBUETTEL and FRANCOIS, 2011](#)). This comes with the problem that C++ objects are by default not serializable and can therefore not be saved/loaded with the default R save/load functions.

AnaCoDa however, does provide functions to load and save parameter and mcmc objects. These are the only two objects that store information during a run.

```
#save objects after a run

runMCMC(mcmc = mcmc, genome = genome, model = model)

writeParameterObject(parameter = parameter, file = "parameter.Rda")

writeMCMCObject(mcmc = mcmc, file = "mcmc_out.Rda")
```

As **genome**, and **model** objects are purely storage containers, no save/load function is provided at this point, but will be added in the future.

```
#load objects

parameter <- loadParameterObject(file = "parameter.Rda")

mcmc <- loadMCMCObject(file = "mcmc_out.Rda")
```

## 2.4.3 File formats

### Protein coding sequence

Protein coding sequences are provided by fasta file with the default format. One line containing the sequence id starting with > followed by the id and one or more lines containing the sequence. The sequences are expected to have a length that is a multiple of three. If a codon can not be recognized (e.g AGN) it is ignored.

```
>YAL001C
TTGGTTCTGACTCATTAGCCAGACGAACTGGTTCAA
CATGTTTCTGACATTCATTCTAACATTGGCATTTCAT
```

```

ACTCTGAACCAACTGTAAGACCATTCTGGCATTTAG
>YAL002W
TTGGAACAAAACGGCCTGGACCACGACTCACGCTCT
TCACATGACACTACTCATAACGACACTCAAATTACT
TTCCTGGAATTCCGCTCTTAGACTCAACTGTCAGAA

```

## Empirical gene expression

Empirical expression or gene specific parameters are provided in a csv file format. The first line is expected to be a header describing each column. The first column is expected to be the gene id, and every additional column is expected to be represent a measurement. Each row corresponds to one gene and contains all measurements for that gene, including missing values.

```

>YAL001C
ORF,DATA_1,DATA_2,...DATA_N
YAL001C,0.254,0.489,...,0.156
YAL002W,1.856,1.357,...,2.014
YAL003W,10.45,NA,...,9.564
YAL005C,0.556,0.957,...,0.758

```

## Ribosome foot-printing counts

Ribosome foot-printing (RFP) counts are provided in a csv file format. The first line is expected to be a header describing each column. The columns are expected in the following order gene id, position, codon, rfpcount. Each row corresponds to a single codon with an associated number of ribosome footprints.

```

GeneID,Position,Codon,rfpCount
YBR177C, 0, ATA, 8

```

YBR177C, 1, CGG, 1

YBR177C, 2, GTT, 8

YBR177C, 3, CGC, 1

## 2.4.4 Analyzing and Visualizing results

### Parameter estimates

After we have completed the model fitting, we are interested in the results. AnaCoDa provides functions to obtain the posterior estimate for each parameter. For gene-set specific parameters or codon specific parameters we can use the function **getCSPEstimates**. Again we can specify for which mixture we would like the posterior estimate and how many samples should be used. **getCSPEstimates** has an optional argument `filename` which will cause the routine to write the result as a csv file instead of returning a **data.frame**.

```
cspMat <- getCSPEstimates(parameter = parameter, CSP="Mutation",
                           mixture = 1, samples = 1000)

head(cspMat)
# AA Codon Posterior 0.025% 0.975%
#1 A GCA -0.2435340 -0.2720696 -0.2165220
#2 A GCC 0.4235546 0.4049132 0.4420680
#3 A GCG 0.7004484 0.6648690 0.7351707
#4 C TGC 0.2016298 0.1679025 0.2387024
#5 D GAC 0.5775052 0.5618199 0.5936979
#6 E GAA -0.4524295 -0.4688044 -0.4356677

getCSPEstimates(parameter = parameter, filename = "mutation.csv",
                 CSP="Mutation", mixture = 1, samples = 1000)
```

To obtain posterior estimates for the gene specific parameters, we can use the function **getExpressionEstimatesForMixture**. In the case below we ask to get the gene specific parameters for all genes, and under the assumption each gene is assigned to mixture 1.



```

phiMat <- getExpressionEstimates(parameter = parameter,
                                gene.index = 1:length(genome),
                                samples = 1000)

head(phiMat)
# PHI log10.PHI Std.Error log10.Std.Error 0.025 0.975 log10.025 ...
#[1,] 0.2729446 -0.6188447 0.0001261525 2.362358e-04 0.07331819 ...
#[2,] 1.4221716 0.1498953 0.0001669425 5.194123e-05 1.09593642 ...
#[3,] 0.7459888 -0.1512764 0.0002313539 1.529267e-04 0.31559618 ...
#[4,] 0.6573082 -0.2030291 0.0001935466 1.400333e-04 0.31591233 ...
#[5,] 1.6316901 0.2098120 0.0001846631 4.986347e-05 1.28410352 ...
#[6,] 0.6179711 -0.2286806 0.0001744928 1.374863e-04 0.28478950 ...

```

However we can decide to only obtain certain gene parameters. in the first case we sample 100 random genes.

```

# sampling 100 genes at random
phiMat <- getExpressionEstimates(parameter = parameter,
                                gene.index = sample(1:length(genome), 100),
                                samples = 1000)

```

Furthermore, AnaCoDa allows to calculate the selection coefficient  $s$  for each codon and each gene. We can use the function **getSelectionCoefficients** to do so. Please note, that this function returns the  $\log(sN_e)$ .

**getSelectionCoefficients** returns a matrix with  $\log(sN_e)$  relative to the most efficient synonymous codon.

```

selectionCoefficients <- getSelectionCoefficients(genome = genome,
                                                  parameter = parameter, samples = 1000)
head(selectionCoefficients)
# GCA GCC GCG GCT TGC TGT GAC GAT ...
#SAKLOA00132g -0.1630284 -0.008695144 -0.2097771 0 -0.1014373 ...

```

```
#SAKLOA00154g -0.8494558 -0.045305847 -1.0930388 0 -0.5285367 ...
#SAKLOA00176g -0.4455753 -0.023764823 -0.5733448 0 -0.2772397 ...
#SAKLOA00198g -0.3926068 -0.020939740 -0.5051875 0 -0.2442824 ...
#SAKLOA00220g -0.9746002 -0.051980440 -1.2540685 0 -0.6064022 ...
#SAKLOA00242g -0.3691110 -0.019686586 -0.4749542 0 -0.2296631 ...
```

We can compare these values to the weights from the codon adaptation index (CAI) [citepSharp1987](#) or effective number of codons ( $N_c$ ) ([WRIGHT, 1990](#)) by using the functions `getCAIweights` and `getNcAA`.

```
caiWeights <- getCAIweights(referenceGenome = genome)
head(caiWeights)
# GCA GCC GCG GCT TGC TGT
#0.7251276 0.6282192 0.2497737 1.0000000 0.6222628 1.0000000
nc.per.aa <- getNcAA(genome = genome)
head(nc.per.aa)
# A C D E F G ...
#SAKLOA00132g 3.611111 1.000000 2.200000 2.142857 1.792453 ...
#SAKLOA00154g 1.843866 2.500000 2.035782 1.942505 1.986595 ...
#SAKLOA00176g 5.142857 NA 1.857143 1.652174 1.551724 3.122449 ...
#SAKLOA00198g 3.800000 NA 1.924779 1.913043 2.129032 4.136364 ...
#SAKLOA00220g 3.198529 1.666667 1.741573 1.756757 2.000000 ...
#SAKLOA00242g 4.500000 NA 2.095890 2.000000 1.408163 3.734043 ...
```

We can also compare the distribution of selection coefficients to the CAI values estimated from a reference set of genes. Figure [2.1](#), produced by the code below, shows that selection coefficients for the same codon can vary greatly between the genes.

```
selectionCoefficients <- getSelectionCoefficients(genome = genome,
                                                    parameter = parameter, samples = 1000)
s <- exp(selectionCoefficients)
```

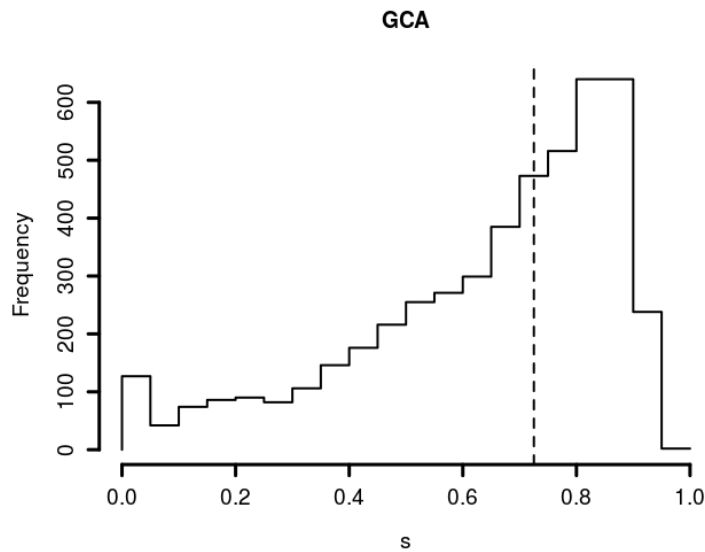


Figure 2.1: Distribution of  $s$  for codon GCA for amino acid alanine. Dashed line indicates the CAI weight for GCA. The comparison provides a more nuanced picture as we can see that the selection on GCA varies across the genome.

```
caiWeights <- getCAIweights(referenceGenome = ref.genome)
codonNames <- colnames(s)
h <- hist(s[, 1], plot = F)
plot(NULL, NULL, axes = F, xlim = c(0,1),
      ylim = range(c(0,h$counts)),
      xlab = "s", ylab = "Frequency",
      main = codonNames[1], cex.lab = 1.2)
lines(x = h$breaks, y = c(0,h$counts), type = "S", lwd=2)
abline(v = cai.weights[1], lwd=2, lty=2)
axis(1, lwd = 3, cex.axis = 1.2)
axis(2, lwd = 3, cex.axis = 1.2)
```

## Diagnostic Plots

A first step after every run should be to determine if the sampling routine has converged. To do that, AnaCoDa provides plotting routines to visualize all sampled parameter traces

from which the posterior sample is obtained (Figure 2.2). First we have to obtain the **trace** object stored within our **parameter** object. Now we can simply plot the **trace** object. The argument **what** specifies which type of parameter should be plotted. Here we plot the selection parameter  $\Delta\eta$  of the ROC model. These parameters are mixture specific and one can decide which mixture set to visualize using the argument **mixture**.

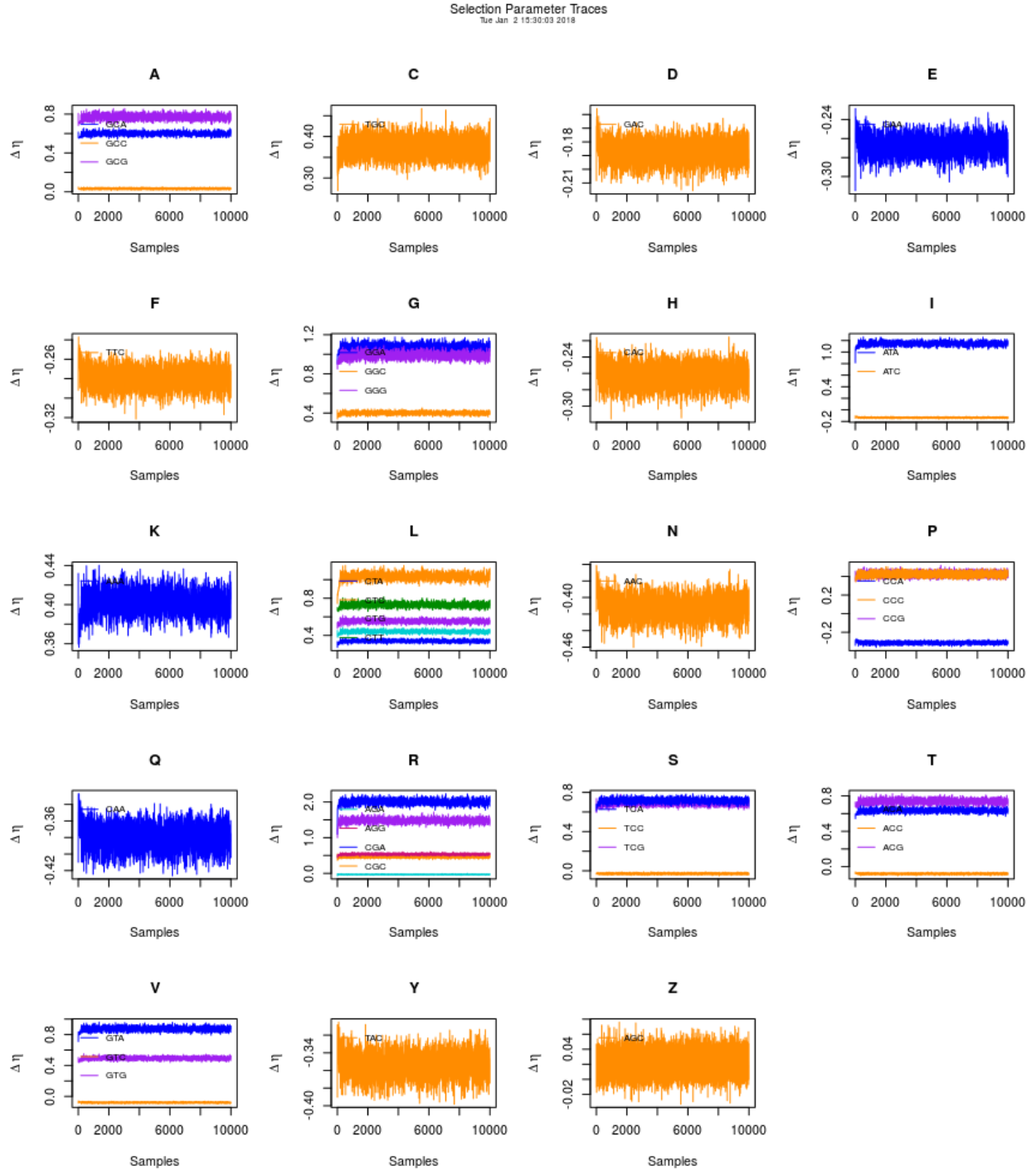


Figure 2.2: Trace plot showing the traces of all 40 codon specific selection parameters  $\Delta\eta$  organized by amino acid.

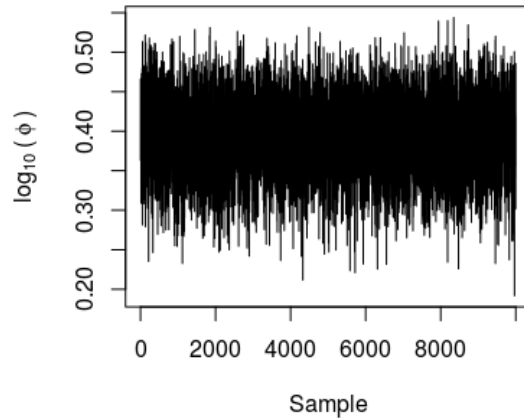


Figure 2.3: Trace plot showing the protein synthesis trace  $\phi$  for gene 669.

```
trace <- getTrace(parameter)
plot(x = trace, what = "Selection", mixture = 1)
```

A special case is the plotting of traces of the protein synthesis rate  $\phi$  (Figure 2.3). As the number of traces for the different  $\phi$  traces is usually in the thousands, a **geneIndex** has to be passed to determine for which gene the trace should be plotted. This allows to inspect the trace of every gene under every mixture assignment.

```
trace <- parameter$getTraceObject()
plot(x = trace, what = "Expression", mixture = 1, geneIndex = 669)
```

We find the likelihood and posterior trace of the model fit in the **mcmc object**. The trace can be plotted by just passing the **mcmc** object to the **plot** routine. Again we can switch between  $\log(\text{likelihood})$  and  $\log(\text{posterior})$  using the argument **what**. The argument **zoom.window** is used to inspect a specified window in more detail. It defaults to the last 10 % of the trace. The  $\log(\text{posterior})$  displayed in the figure title is estimated over the **zoom.window** (Figure 2.4).

```
plot(mcmc, what = "LogPosterior", zoom.window = c(9000, 10000))
```

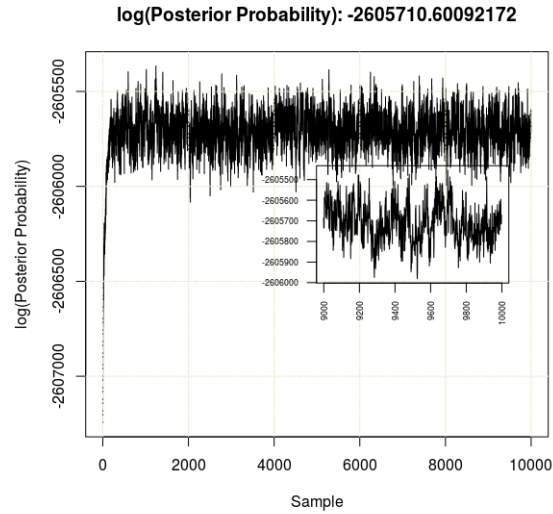


Figure 2.4: Trace plot showing the  $\log(\text{posterior})$  trace for the current model fit. Window inset shows the last 1.000 samples

## Model visualization

We can visualize the results of the model fit by plotting the **model** object (Figure 2.5). For this we require the model and the **genome** object. We can adjust which mixture set we would like to visualize and how many samples should be used to obtain the posterior estimate for each parameter. For more details see [GILCHRIST \*et al.\* \(2015\)](#).

```
# use the last 500 samples from mixture 1 for posterior estimate.
plot(x = model, genome = genome, samples = 500, mixture = 1)
```

As AnaCoDa is designed with the idea to allow gene-sets to have independent gene-set specific parameters, AnaCoDa also provides the option to compare different gene-sets by plotting the parameter object. Figure 2.6 allows us to compare the selection parameter estimated by ROC for seven yeast species. The code below illustrates how the figure is plotted.

```
# use the last 500 samples from mixture 1 for posterior estimate.
plot(parameter, what = "Selection", samples = 500)
```

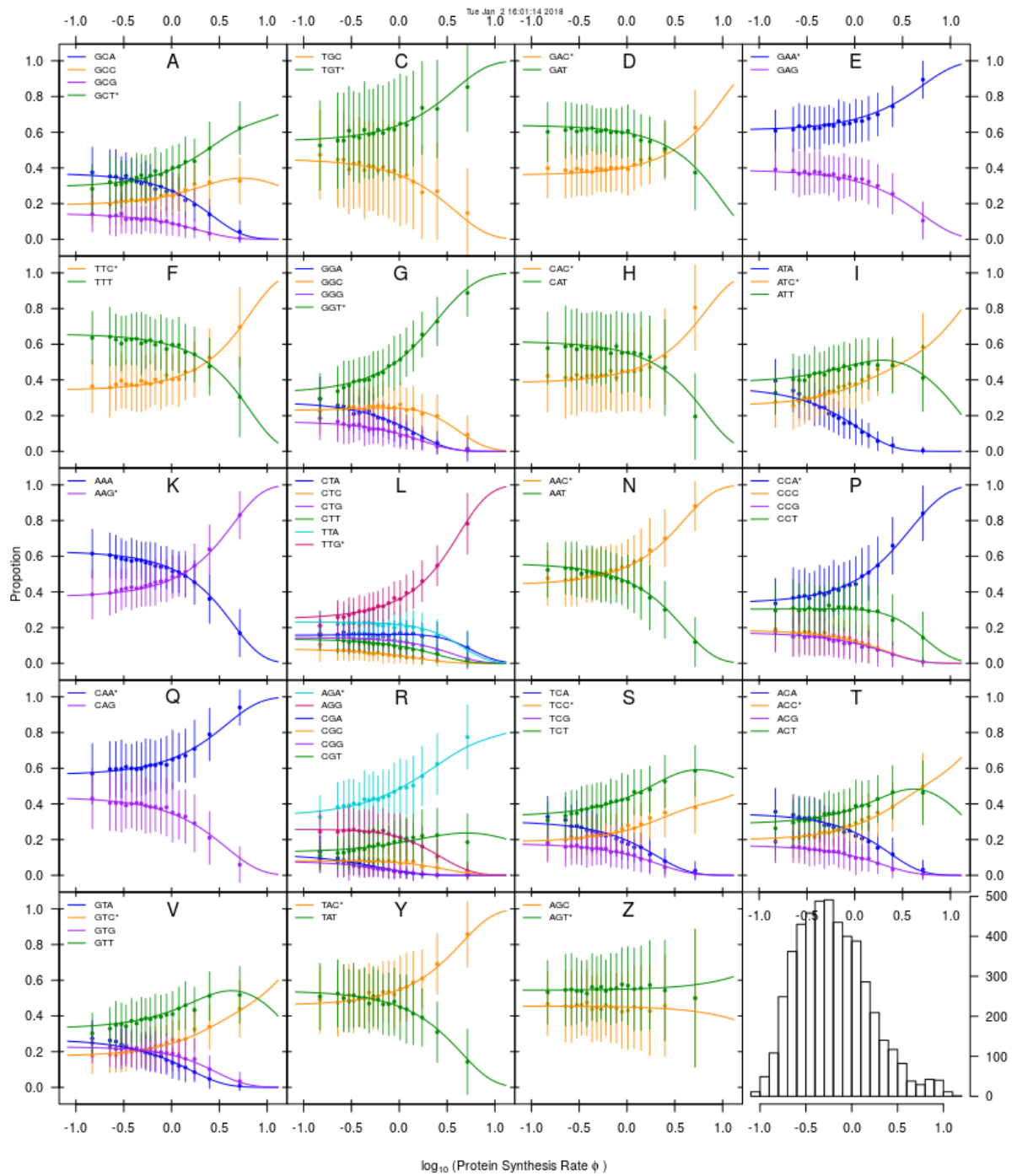


Figure 2.5: Fit of the ROC model for a random yeast. The solid line represent the model fit from the data, showing how synonymous codon frequencies change with gene expression. The points are the observed mean frequencies of a codon in that synthesis rate bin and the whisks indicate the standard deviation within the bin. The codon favored by selection is indicated by a ”\*”. The bottom right panel shows how many genes are contained in each bin



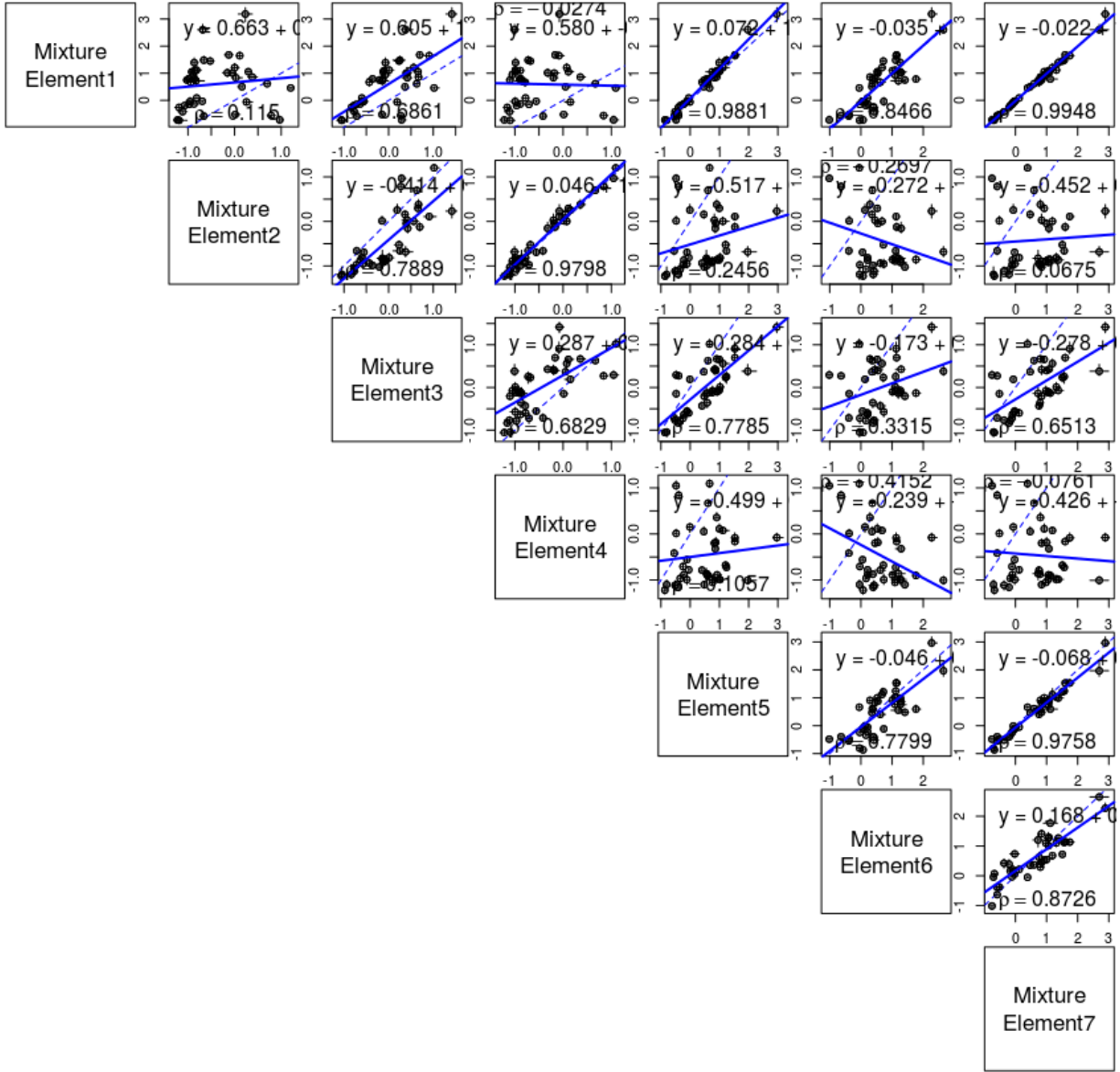


Figure 2.6: Comparison of the selection parameter of seven yeast species estimated with ROC-SEMPPR.

## Chapter 3

### Decomposing mutation and selection to identify mismatched codon usage

This chapter is a lightly revised version of a paper to be submitted to Genome Biology and Evolution and co-authored with Michael A. Gilchrist and Russel Zaretzki.

C. Landerer, R. Zaretzki, M.A. Gilchrist, Decomposing mutation and selection to identify mismatched codon usage

### 3.1 Abstract

For decades, codon usage has been used as a measure of adaptation for translational efficiency of a gene's coding sequence. These patterns of codon usage reflect both the selective and mutational environment in which the coding sequences evolved. Over this same period, gene transfer between lineages has become widely recognized as an important biological phenomenon. Nevertheless, most studies of codon usage implicitly assume that all genes within a genome evolved under the same selective and mutational environment, an assumption violated when introgression occurs. In order to better understand the effects of introgression on codon usage patterns and vice versa, we examine the patterns of codon usage in the yeast which has experienced a large introgression, *Lachancea kluyveri*. We quantify the effects of mutation bias and selection for translation efficiency on the codon usage pattern of the endogenous and introgressed exogenous genes using a Bayesian mixture model, ROC SEMPPR, which is built on mechanistic assumptions of protein synthesis and grounded in population genetics. We find substantial differences in codon usage between the endogenous and exogenous genes, and show that these differences can be largely attributed to a shift in mutation bias from A/T ending codons in the endogenous genes to C/G ending codons in the exogenous genes. Recognizing the two different signatures of mutation and selection bias improves our ability to predict protein synthesis rate by 17% and allowed us to accurately assess codon preferences. In addition, using our estimates of mutation and selection bias, we

to identify *Eremothecium gossypii* as the most likely source lineage, estimate the introgression occurred  $\sim 6 \times 10^8$  generation ago, and estimate its historic and current genetic load. Together, our work illustrates the advantage of mechanistic, population genetic models like ROC SEMPPR and the quantitative estimates they provide when analyzing sequence data.

## 3.2 Introduction

Synonymous codon usage patterns varies within a genome and between taxa, reflecting differences in mutation bias, selection, and genetic drift. The signature of mutation bias is largely determined by the organism’s internal or cellular environment, such as their DNA repair genes or UV exposure. While this mutation bias is an omnipresent evolutionary force, its impact can be obscured or even amplified by selection. The signature of selection on codon usage is also largely determined by an organism’s cellular environment, such as its tRNA species, their copy number, and post-transcriptional modifications. The strength of selection on the codon usage of an individual gene is largely determined by its expression level which, in turn, is also largely determined by the organism’s external environment. In general, the strength of selection on codon usage increases with its expression level (GOUY and GAUTIER, 1982; IKEMURA, 1985; BULMER, 1990), specifically its protein synthesis rate (GILCHRIST, 2007). Thus as gene expression increases, codon usage shifts from a process dominated by mutation to a process dominated by selection. The overall efficacy of selection on codon usage is a function of the organism’s effective population size  $N_e$  which, in turn, is largely determined by its external environment. By explicitly modeling the combined forces of mutation, selection, and drift, ROC SEMPPR allows us disentangle the evolutionary forces responsible for the patterns of codon usage bias (CUB) encoded in an species’ genome (GILCHRIST, 2007; SHAH and GILCHRIST, 2011a; WALLACE *et al.*, 2013; GILCHRIST *et al.*, 2015), should provide biologically meaningful information about the lineage’s historical cellular and external environment.

Most studies implicitly assume that the CUB of a genome is shaped by a single cellular environment. As genes are horizontally transferred, introgress, or combined to form novel hybrid species, one would expect to see the influence of multiple cellular environments on a genomes codon usage pattern (MDIGUE *et al.*, 1991; LAWRENCE and OCHMAN, 1997). Given that transferred genes are likely to be less adapted than endogenous genes to their new cellular environment, we expect a greater genetic load of transferred genes if donor and

recipient environment differ greatly in their selection bias, making such transfers less likely. More practically, if differences in codon usage of transferred genes are unaccounted for, they may distort parameter estimates. Such distortion could lead to the wrong codon preference for an amino acid, underestimate the variation in protein synthesis rate, or bias mutation estimates when analyzing a genome.

To illustrate these ideas, we analyze the CUB of the genome of *Lachancea kluyveri*, the earliest diverging lineage of the Lachancea clade. The Lachancea clade diverged from the Saccharomyces clade, prior to its whole genome duplication  $\sim 100$  Mya ago (MARCET-HOUBEN and GABALDN, 2015; BEIMFORDE *et al.*, 2014). Since that time, *L. kluyveri* has experienced a large introgression of exogenous genes found in all populations (FRIEDRICH *et al.*, 2015). The introgression replaced the left arm of the C chromosome and displays a 13% higher GC content than the endogenous *L. kluyveri* genome (PAYEN *et al.*, 2009; FRIEDRICH *et al.*, 2015). These characteristics make *L. kluyveri* an ideal model to study the effects of an introgressed cellular environment and the resulting mismatch in codon usage.

Using ROC SEMPPR, a Bayesian population genetics model based on a mechanistic description of ribosome movement along an mRNA, allows us to quantify the cellular environment in which genes have evolved by separately estimating the effects of mutation bias and selection bias on codon usage. ROC SEMPPR’s resulting predictions of protein synthesis rates have been shown to be on par with laboratory measurements (SHAH and GILCHRIST, 2011a; GILCHRIST *et al.*, 2015). We use ROC SEMPPR to independently describe two cellular environments reflected in the *L. kluyveri* genome; the signature of the current environment in the endogenous genes and the decaying signature of the exogenous environment in the introgressed genes. Our results indicate that the difference in GC content between endogenous and exogenous genes is mostly due to the differences in mutation bias of their ancestral environments. Accounting for these different signatures of mutation bias and selection bias of the endogenous and exogenous sets of genes substantially improves our ability to predict present day protein synthesis rates. These endogenous and exogenous gene

set specific estimates of mutation bias and selection bias, in turn allow us to address more refined questions of biological importance. For example, it allows us to identify *E. gossypii* as the most likely source of the introgressed genes out of the 38 yeast lineages with sequenced genomes, estimate the age of the introgression to be on the order of 0.2-1 Mya, estimate the genetic load of these genes, both at the time of introgression and now, as well as make predictions about how the CUB of the introgressed genes will evolve in the future.

### 3.3 Results

#### 3.3.1 The Signatures of two Cellular Environments within *L. kluyveri*'s Genome

We used our software package AnaCoDa (LANDERER *et al.*, 2018) to compare model fits of ROC SEMPPR to the entire *L. kluyveri* genome and its genome partitioned into two sets of 4,864 endogenous and 497 exogenous genes. AIC values strongly support the hypothesis that the *L. kluyveri* genome consists of genes with two different and distinct patterns of codon usage bias ( $\Delta\text{AIC} = 75,462$ ; Table 4.1). We find additional support for this hypothesis when we compare our predictions of gene expression to empirically observed values. Specifically, the explanatory power between our predictions and observed values improved by  $\sim 42\%$ , from  $R^2 = 0.33$  to 0.46 (Figure 3.1).

Table 3.1: Model selection of the two competing hypothesis. Reported are the log-likelihood,  $\log(\mathcal{L})$ , the number of parameters estimated  $n$ , AIC, and  $\Delta\text{AIC}$  values.

Hypothesis	$\log(\mathcal{L})$	$n$	AIC	$\Delta\text{AIC}$
Separated	-2,612,397	5,402	5,235,598	0
Combined	-2,650,047	5,483	5,311,060	75,462

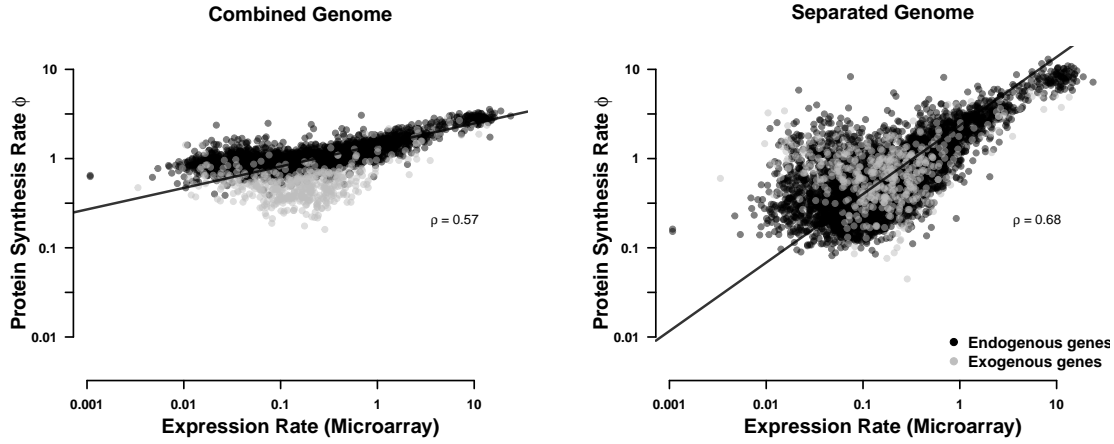


Figure 3.1: Comparison of predicted protein synthesis rate  $\phi$  to microarray data from [TSANKOV \*et al.\* \(2010\)](#) for (a) the combined genome and (b) the separated endogenous and exogenous genes. Endogenous genes are displayed in black and exogenous genes in red. Black line indicates type II regression line ([SOKAL and ROHLF, 1981](#)).

### 3.3.2 Comparing Differences in the Endogenous and Exogenous Codon Usage

To better understand the differences in the endogenous and exogenous cellular environments, we compared our parameter estimates of mutation bias  $\Delta M$  and selection  $\Delta\eta$  for the two sets of genes. Our estimates of  $\Delta M$  for the endogenous and exogenous genes were negatively correlated ( $\rho = -0.49$ ), indicating weak concordance of  $\sim 5\%$  between the two mutation environments (Figure 3.2). For example, the endogenous genes show a mutational preference for A and T ending codons in  $\sim 95\%$  of the codon families. In contrast, the exogenous genes display an equally consistent mutational preference towards C and G ending codons (Table 3.2). As a result, only the two codon amino acid Phenylalanine (Phe, F) shares the same rank order across the endogenous and exogenous  $\Delta M$  estimates.

In contrast, our estimates of  $\Delta\eta$  for the endogenous and exogenous genes were positively correlated ( $\rho = 0.69$ ) and showing concordance of  $\sim 53\%$  between the two selection environments (Figure 3.2). We find that the strength of selection within each codon family differs between sets of genes. Overall, the endogenous genes only show a selection preference





the complete *L. kluyveri* genome is estimated to share the selection preference with the exogenous genes in  $\sim 60\%$  of codon families that show discordance between endogenous and exogenous genes. These results clearly show that it is important to recognize the difference in endogenous and exogenous genes and treat these genes as separate sets to avoid the inference of incorrect synonymous codon preferences and better predict protein synthesis.

### 3.3.3 Determining Source of Exogenous Genes

We combined our estimates of mutation bias  $\Delta M$  and selection bias  $\Delta\eta$  with synteny information and searched for potential source lineages of the introgressed exogenous region. We examined 38 yeast lineages (Table 3.4) of which two (*Eremothecium gossypii* and *Candida dubliniensis*) showed a strong positive correlation in codon usage (Figure 3.3). The endogenous *L. kluyveri* genome exhibits codon usage very similar to most yeast lineages examined, indicating little variation in codon usage among the examined yeasts (Figure 3.5). Four lineages show a positive correlation for  $\Delta M$  and  $\Delta\eta$  with the exogenous genes and have a weak to moderate positive correlation in selection bias with the endogenous genes; but, like the exogenous genes, tend to have a negative correlation in  $\Delta M$  with the endogenous genes.

Comparing synteny between the exogenous genes, which are restricted to the left arm of chromosome C, and *E. gossypii* and *C. dubliniensis* as well as closely related yeast species we find that *E. gossypii* displays the highest synteny (Figures 3.7 & 3.8). *C. dubliniensis*, even though it displays similar codon usage does not show synteny with the exogenous region. Furthermore, the synteny relationship between the exogenous region and other yeasts appears to be limited to the Saccharomycetacease clade (Figure 3.8). Given these results, we conclude that of the 38 examined yeast lineages the *E. gossypii* lineage is the most likely source of the introgressed exogenous genes.

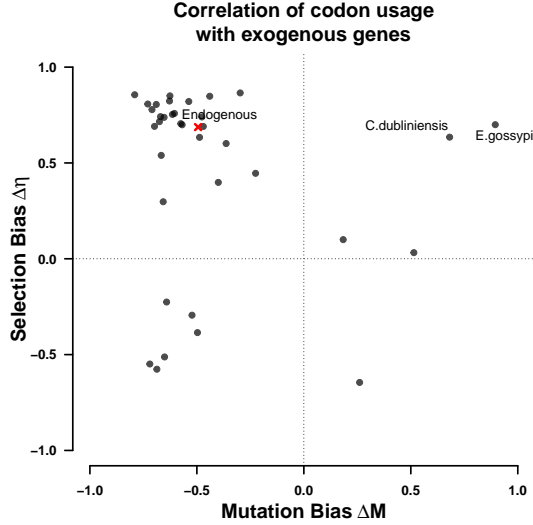


Figure 3.3: Correlation coefficients of  $\Delta M$  and  $\Delta\eta$  of the exogenous genes with 38 examined yeast lineages. Dots indicate the correlation of  $\Delta M$  and  $\Delta\eta$  of the lineages with the endogenous and exogenous parameter estimates. All regressions were performed using a type II regression (SOKAL and ROHLF, 1981).

### 3.3.4 Estimating Introgression Age

We modeled the change in codon frequency as a model of exponential decay, we estimated the age of the introgression assuming that *E. gossypii* still represents the mutation bias of its ancestral source lineage at the time of the introgression and a constant mutation rate. We infer the age of the introgression to be on the order of  $6.2 \pm 1.2 \times 10^8$  generations. Assuming *L. kluyveri* experiences between one and eight generations per day, we estimate the introgression to have occurred between 212,000 to 1,700,000 years ago. Our estimate places the time of the introgression earlier than previously assumed (FRIEDRICH *et al.*, 2015).

Using the same approach, we also estimated the persistence of the signal of the exogenous cellular environment. We assume that differences in mutation bias will decay more slowly than differences in selection bias to be able to utilize our bias free estimates of  $\Delta M$ . We predict that the  $\Delta M$  signal of the source cellular environment will have decayed to be within one percent of the *L. kluyveri* environment in  $\sim 5.4 \pm 0.2 \times 10^9$  generations, or

between 1,800,000 and 15,000,000 years. Together, these results indicate that the mutation signature of the exogenous genes will persist for a very long time.

### 3.3.5 Genetic Load due to Mismatching Codon Usage of the Exogenous Genes

We define genetic load as the difference between the fitness of an expected, replaced endogenous gene and the exogenous gene,  $s \propto \phi \Delta \eta$  due to the mismatch in codon usage parameters (See Methods for details). Estimates of selection bias for the exogenous genes show that, while well correlated with the endogenous genes, only nine amino acids share the same optimal codon. Exogenous genes are, therefore, expected to represent a significant reduction in fitness, or genetic load for *L. kluyveri* due to this mismatch in codon usage. As the introgression occurred before the diversification of *L. kluyveri* and has fixed throughout all populations (FRIEDRICH *et al.*, 2015), we can not observe the original endogenous sequences that have been replaced by the introgression. Using our estimates of  $\Delta M$  and  $\Delta \eta$  from the endogenous genes and assuming that the current exogenous amino acid composition of genes is representative of the replaced endogenous genes, we estimate the genetic load of the exogenous genes at the time of introgression (Figure 3.4a) and currently (Figure 3.4b). We find that the genetic load due to mismatched codon usage was -0.0008 at the time of the introgression and still represents a genetic load of -0.0003 today.

In order to account for differences in the efficacy of selection on codon usage between the donor lineage and *L. kluyveri* using a linear scaling factor  $\kappa$  (See Methods for details). We predict that a small number of low expression genes ( $\phi < 1$ ) were weakly exapted at the time of the introgression (Figure 3.4a). High expression genes ( $\phi > 1$ ) are predicted to have carried the largest genetic load in the novel cellular environment. These highly expressed genes are inferred to have the greatest degree of adaptation since the time of the introgression to the *L. kluyveri* cellular environment (Figures 3.4a & 3.10).

Explicitly

mention

value for

genetic

load

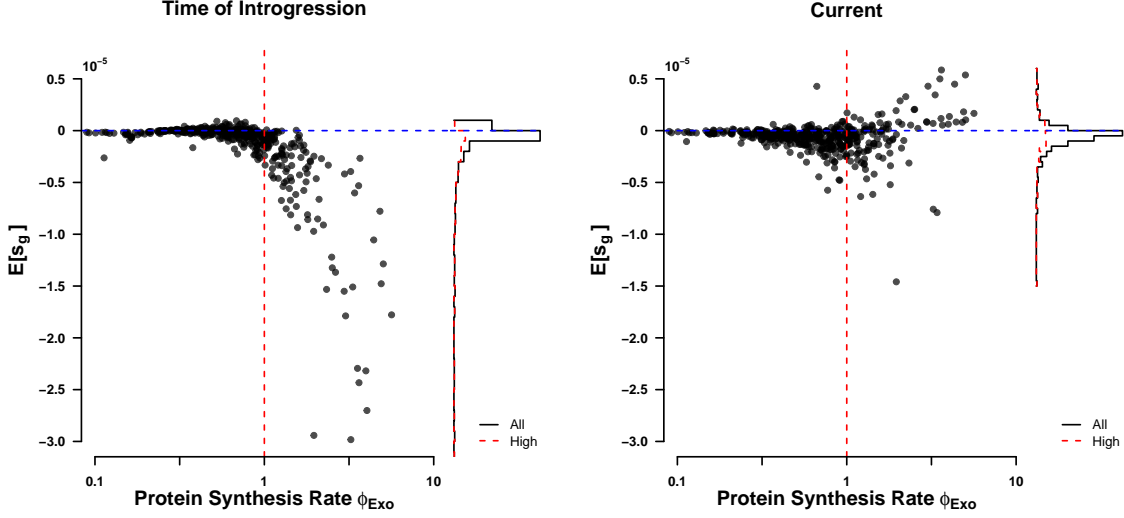


Figure 3.4: Genetic load  $s = \Delta\eta\phi$  (a) at the time of introgression ( $\kappa = 5$ ), and (b) currently ( $\kappa = 1$ ).

### 3.4 Discussion

In order to study the evolutionary effects of an introgression, we used ROC SEMPFR, a mechanistic model of ribosome movement along an mRNA. Our parameter estimates indicate that the *L. kluyveri* genome contains distinct signatures of mutation and selection bias from both an endogenous and exogenous cellular environment. By fitting ROC SEMPFR separately to *L. kluyveri*'s endogenous and exogenous sets of genes we generate a quantitative description of their signatures of mutation bias and natural selection for efficient protein translation. Our results indicate that the difference in GC content between endogenous and exogenous genes is mostly due to differences in mutation bias, but we also show that the strength and rank order of selection within a codon family differ between endogenous and exogenous cellular environments. Even though the exogenous genes make up only  $\sim 10\%$  of the *L. kluyveri* genome, when we fail to recognize these differences our estimates of  $\Delta M$  and  $\Delta\eta$  deviate substantial from their actual values (Figure 3.6). While this sensitivity of our parameters to a second cellular environment may be surprising, it highlights the importance of recognizing different cellular environments reflected by a genome. Furthermore, our results

indicate that we can attribute the increased GC content in the exogenous genes mostly to differences in mutation bias favoring G/C ending codons rather than selection.

The separation of the endogenous and exogenous genes improves our estimates of protein synthesis rate  $\phi$  by 42% relative to the full genome estimate ( $R^2 = 0.32$  vs.  $0.46$ , respectively). Furthermore, failing to separately analyze the endogenous and exogenous genes results in an unrealistically small amount of intergenic variation in  $\phi$  (compare Figure 3.1a & b). This behavior is due, in part, to constraining  $E[\phi] = 1$  which allows us to compare the efficacy of selection  $sN_e$  across genomes. Extremely small variances in the  $\phi$  values estimated by ROC SEMPPR could indicate that a genome contains the signature of multiple cellular environments.

The mutation and selection bias parameters  $\Delta M$  and  $\Delta \eta$  of the introgressed exogenous genes contain information, albeit decaying, about its previous cellular environment. We, therefore, utilize  $\Delta M$  and  $\Delta \eta$  to identify potential source lineages. The *E. gossypii* and *C. dubliniensis* lineages stand out from the other 36 yeast lineages in that the correlation coefficients between their  $\Delta M$  and  $\Delta \eta$  parameters and those of the exogenous genes are  $> 0.5$  (Figure 3.2). In terms of gene order, we found that synteny with the exogenous genes is limited to the Saccharomycetaceae clade, which *C. dubliniensis* is outside of. Overall, the synteny coverage extends along the whole exogenous regions with the exception of the 3' and 5' ends of the exogenous region (Figure 3.8b). Further, of the 38 species examined, *E. gossypii* is the only genome with a GC content  $> 50\%$ , making it most similar to the exogenous genes. Thus, only the *E. gossypii* genome displays strong correlations in  $\Delta M$  and  $\Delta \eta$ , synteny, and similar GC content with the exogenous genes.

With *E. gossypii* identified as potential source lineage of the introgressed region, we inferred the time since the introgression occurred using our estimates of mutation bias  $\Delta M$ . Our  $\Delta M$  estimates are well suited for this task as they are free of the influence of selection and unbiased by  $N_e$  and other scaling terms, which is in contrast to our estimates of  $\Delta \eta$  (GILCHRIST *et al.*, 2015). Our estimated age of the introgression of  $6.2 \pm 1.2 \times 10^8$

generations is  $\sim 10$  times longer time than a previous minimum estimate by [FRIEDRICH \*et al.\* \(2015\)](#) of  $5.6 \times 10^7$  generations. Our estimate assumes that the current *E. gossypii* and *L. kluyveri* cellular environment reflect their ancestral states at the time of the introgression. If the ancestral mutation environments were more similar (dissimilar) at the time of the introgression than now our result is an overestimate (underestimate).

In order to estimate the introgression’s genetic load due to codon mismatch, we had to make three key assumptions: 1) at the time of introgression the amino acid sequences of the endogenous genes and exogenous genes were highly similar, 2) the current *L. kluyveri* cellular environment is reflective of the cellular environment at the time of the introgression, and 3) the *E. gossypii* cellular environment reflects its ancestral environment at the time of the introgression. In general due to their very nature, low expression genes contribute little to the genetic load. Indeed,  $\sim 30\%$  of low expression exogenous genes ( $\phi < 1$ ) appeared to be exapted at the time of the introgression. These exapted genes are likely due to the mutation bias in the endogenous genes matching the selection bias in the exogenous genes for G/C ending codons. In contrast, highly expressed genes are predicted to have imposed a large genetic load. Many of these genes appear to still represent a significant genetic load. Overall, our estimates of codon mismatch genetic load, therefore, suggest strong selection against the introgression.

It is hard to contextualize the probability of this introgression being fixed as we are not aware of any estimates of the frequency at which such large scale introgressions of genes occur. A related example of a large scale merger of genomic material can be found in *S. bayanus*, which is currently believed to be a hybrid of *S. cerevisiae*, *S. eubayanus*, and *S. uvarum* lineages. Unlike with *L. kluyveri* and *E. gossypii*, the progenitor lineages of *S. bayanus* have similar codon usage parameters. For example, the correlation between  $\Delta M$  and  $\Delta \eta$  for these two lineages are  $\rho = 0.83$  and  $0.98$  (data not shown). These similarities in  $\Delta M$  and  $\Delta \eta$  parameters suggest that the genetic load for *S. bayanus* due to codon usage mismatch is small relative to the exogenous genes considered here. The large genetic load

Get  
uvarum  
genome  
and  
analyze!

of the exogenous genes due to codon mismatch at the time of the introgression would seem to indicate that the fixation of the introgression was either a fluke event or the codon mismatch genetic load was countered by one or more highly advantageous loci within the introgression.

Under the first scenario, our best estimate of the selection coefficient against the introgression based on expected codon mismatch at that time is  $s = -0.0008$  and an effective population size  $N_e$  on the order of  $10^8$  (WAGNER, 2005) yields an approximate fixation probability of  $(1 - \exp[-s])/(1 - \exp[2 - sN_e]) \approx 10^{-6950}$  (SELLA and HIRSH, 2005). Even though *L. kluyveri* diverged from the rest of the Lachancea clade around 85 Mya (KENSCHKE *et al.*, 2008; MARCET-HOUBEN and GABALDN, 2015), if we assume 1 to 8 generations/day, which implies  $10^{10}$  to  $10^{11}$  generations since the time of divergence, one round of meiosis for every 1000 rounds of mitosis based on *S. paradoxus* (TSAI *et al.*, 2008), and  $N_e \approx 10^8$  there were only  $10^{15}$  to  $10^{16}$  opportunities for such an introgression to have occurred and fixed. Clearly, unless there was a severe bottleneck with  $N_e < 1/|s| \approx 1,250$  around the time of introgression, which conceivably could have been triggered by a speciation event, this scenario seems very unlikely.

In the second scenario, where we assume the introgression contained advantageous loci, one may wonder why recombination events did not limit the introgression to only the adaptive loci. PAYEN *et al.* (2009) found that the exogenous region has a lower rate of recombination, presumably due to the dissimilarity in GC content and/or a lower than average sequence homology between the exogenous region and the one it replaced. Compatible with this explanation is the possibility of several highly advantageous loci distributed across the region which then drove a rapid selective sweep and/or the population through a bottleneck speciation process. A careful analysis of intra-specific genetic variation within the endogenous and exogenous regions could provide help us distinguish between these various scenarios.

Overall, our results show the usefulness of the separation of mutation bias and selection bias and the importance of recognizing the presence of multiple cellular environments in the study of codon usage. We also illustrate how a mechanistic model like ROC SEMPFR and



the quantitative estimates it provides can be used for more sophisticated hypothesis testing in the future. In contrast to other approaches used to study codon usage like CAI (SHARP, 1987) or tAI (DOS REIS *et al.*, 2004), ROC SEMPPR incorporates the effects of mutation bias and amino acid composition explicitly COPE *et al.* (2018). We highlight potential issues when estimating codon preferences, as estimates can be biased by the signature of a second, historical cellular environment. In addition, we show how quantitative estimates of mutation bias and selection relative to drift can be obtained from codon data and used to infer the fitness cost of an introgression as well as its history and potential future.

Discuss  
CAI and  
tAI in  
intro  
and  
start of  
discus-  
sion.

## 3.5 Materials and Methods

### 3.5.1 Separating Endogenous and Exogenous Genes

A GC-rich region was identified by PAYEN *et al.* (2009) in the *L. kluyveri* genome extending from position 1 to 989,693 of chromosome C. This region was later identified as an introgression by FRIEDRICH *et al.* (2015). We obtained the *L. kluyveri* genome from SGD Project <http://www.yeastgenome.org/download-data/> (on 09-27-2014) and the annotation for *L. kluyveri* NRRL Y-12651 (assembly ASM14922v1) from NCBI (on 12-09-2014). We assigned 457 genes located on chromosome C with a location within the  $\sim 1Mb$  window to the exogenous gene set. All other 4864 genes of the *L. kluyveri* genome were assigned to the endogenous genes. All genes could be uniquely assigned to one or the other gene set.

### 3.5.2 Model Fitting with ROC SEMPPR

ROC SEMPPR was fitted to each genome using AnaCoDa (0.1.1) (LANDERER *et al.*, 2018) and R (3.4.1) (R CORE TEAM, 2015). ROC SEMPPR was run from multiple starting values for at least 250,000 iterations, only every 50th step was collected as a sample to

reduce autocorrelation. After manual inspection to verify that the MCMC had converged, parameter posterior means were estimated from the last 500 samples.

### 3.5.3 Comparing Codon Specific Parameter Estimates

Choice of reference codon does reorganize codon families coding for an amino acid relative to each other, therefore all parameter estimates are relative to the mean for each codon family.

$$\Delta M_{i,a}^c = \Delta M_{i,a} - \overline{\Delta M_a} \quad (3.1)$$

$$\Delta \eta_{i,a}^c = \Delta \eta_{i,a} - \overline{\Delta \eta_a} \quad (3.2)$$

Comparison of codon specific parameters ( $\Delta M$  and  $\Delta \eta$ ) was performed using the function `lmodel2` in the R package `lmodel2` (1.7.3) (LEGENDRE, 2018) and R version 3.4.1 (R CORE TEAM, 2015). Type II regression was performed with re-centered parameter estimates, accounting for noise in dependent and independent variable (SOKAL and ROHLF, 1981).

### 3.5.4 Synteny Comparison

We obtained complete genome sequences from NCBI (on: 02-05-2017). Genomes were aligned and checked for synteny using SyMAP (4.2) with default settings (SODERLUND *et al.*, 2006, 2011). We assess synteny as percentage coverage of the exogenous gene region (Figure 3.8b).

### 3.5.5 Estimating Age of Introgression

We modeled the change in codon frequency over time using an exponential model for all two codon amino acids, and describing the change in codon  $c_1$  as

$$\frac{dc_1}{dt} = -\mu_{1,2}c_1 - \mu_{2,1}(1 - c_1) \quad (3.3)$$

where  $\mu_{i,j}$  is the rate at which codon  $i$  mutates to codon  $j$  and  $c_1$  is the frequency of the reference codon. Our estimates of  $\Delta M_{\text{endo}}$  can be used to calculate the steady state of equation 3.3.

$$\frac{\mu_{2,1}}{\mu_{1,2} + \mu_{2,1}} = \frac{1}{1 + \exp[\Delta M_{\text{endo}}]} \quad (3.4)$$

Solving for  $\mu_{1,2}$  gives us  $\mu_{1,2} = \Delta M_{\text{endo}} \exp[\mu_{2,1}]$  which allows us to rewrite and solve equation 3.3 as

$$c_1(t) = \frac{\exp[-t(1 + \Delta M_{\text{endo}})\mu_{2,1}] \exp[t(1 + \Delta M_{\text{endo}})\mu_{2,1}] + (1 + \Delta M_{\text{endo}})K}{1 + \Delta M_{\text{endo}}} \quad (3.5)$$

where  $K$  is

$$K = c_1(0) - \frac{1}{1 + \Delta M_{\text{endo}}} \quad (3.6)$$

Equation 3.5 was solved with a mutation rate  $m_{2,1}$  of  $3.8 \times 10^{-10}$  per nucleotide per generation (LANG and MURRAY, 2008). Initial codon frequencies  $c_1(0)$  for each codon family were taken from our mutation parameter estimates for *E. gossypii*  $\Delta M_{\text{gos}}$ . Current codon frequencies for each codon family were taken from our estimates of  $\Delta M$  from the exogenous genes. Mathematica (11.3) (WOLFRAM RESEARCH INC., 2017) was used to calculate the time  $t_{\text{intro}}$  it takes for the initial codon frequencies  $c_1(0)$  for each codon family to equal the current exogenous codon frequencies. The same equation was used to determine the time  $t_{\text{decay}}$  at which the signal of the exogenous cellular environment has decayed to within 1% of the endogenous environment.

## Estimating Genetic Load

To estimate the genetic load due to mismatched codon usage, we made three key assumptions. First, we assumed that the current exogenous amino acid sequence of a gene is representative of its ancestral state and the replaced endogenous gene it replaced. Second, we assume that the currently observed cellular environment of *E. gossypii* reflects the cellular environment

that the exogenous genes experienced before transfer to *L. kluyveri*. Lastly, we assume that the difference in the efficacy of selection between the cellular environments due to differences in either effective population size  $N_e$  or the selective cost of an ATP  $q$  of the source lineage and *L. kluyveri* can be expressed as a scaling constant and that protein synthesis rate  $\phi$  has not changed between the replaced endogenous and the introgressed exogenous genes. Using estimates for  $N_e = 1.36 \times 10^7$  (WAGNER, 2005) for *Saccharomyces paradoxus* we scale our estimates of  $\Delta\eta$  and define  $\Delta\eta' = \frac{\Delta\eta}{N_e}$ .

We scale the difference in the efficacy of selection on codon usage between the donor lineage and *L. kluyveri* using a linear scaling factor  $\kappa$ . As  $\Delta\eta$  is defined as  $\Delta\eta = 2N_e q(\eta_i - \eta_j)$ , we can not distinguish if  $\kappa$  is a scaling on protein synthesis rate  $\phi$ , effective population size  $N_e$ , or the selective cost of an ATP  $q$  (GILCHRIST, 2007; GILCHRIST *et al.*, 2015). We calculated the genetic load each gene represents due to its mismatched codon usage assuming additive fitness effects as

$$s_g = \sum_{i=1}^{n_g} -\kappa \phi_g \Delta\eta'_i \quad (3.7)$$

where  $s_g$  is the overall strength of selection for translational efficiency on gene,  $g$  in the exogenous gene set,  $\kappa$  is a constant, scaling the efficacy of selection between the endogenous and exogenous cellular environments,  $n_g$  is length of the protein,  $\phi_g$  is the estimated protein synthesis rate of the gene in the endogenous environment, and  $\Delta\eta'_i$ , is the  $\Delta\eta'$  for the codon at position  $i$ . As stated previously, our  $\Delta\eta$  are relative to the mean of the codon family. We find that the genetic load of the introgressed genes is minimized at  $\kappa \sim 5$  (Figure 3.9b). Thus, we expect a five fold difference in the efficacy of selection between *L. kluyveri* and *E. gossypii*, either due to differences in either protein synthesis rate  $\phi$ , effective population size  $N_e$ , or the selective cost of an ATP  $q$ . Therefore, we set  $\kappa = 1$  if we calculate the  $s_g$  for the endogenous and the current exogenous genes, and  $\kappa = 5$  for  $s_g$  for the genetic load at the time of introgression.

Move  
 $\Delta\eta$  def-  
 inition  
 to first  
 usage of  
 $\Delta\eta$  in  
 Meth-  
 ods.

Since we are unable to observe codon counts for the replaced endogenous genes and for the exogenous genes at the time of introgression, we calculate expected codon counts

$$E[n_{g,i}] = \frac{\exp[-\Delta M_i - \Delta \eta_i \phi_g]}{\sum_j^C \exp[-\Delta M_j - \Delta \eta_j \phi_g]} \times m_{a_i} \quad (3.8)$$

$m_{a_i}$  is the number of occurrences of amino acid  $a$  that codon  $i$  codes for. We report the genetic load due to mismatched codon usage of the introgression as  $E[s_g] = s_{\text{intro},g} - s_{\text{endo},g}$  where  $s_{\text{intro},g}$  is the genetic load of an introgressed gene  $g$  either at the time of the introgression or presently.

## 3.6 Acknowledgments

This work was supported in part by NSF Awards MCB-1120370 (MAG and RZ) and DEB-1355033 (BCO, MAG, and RZ) with additional support from The University of Tennessee Knoxville. CL received support as a Graduate Student Fellow at the National Institute for Mathematical and Biological Synthesis, an Institute sponsored by the National Science Foundation through NSF Award DBI-1300426, with additional support from UTK. The authors would like to thank Brian C. O’Meara and Alexander Cope for their helpful criticisms and suggestions for this work.

### 3.7 Appendix: Supplementary Material

Table 3.2: Synonymous codon preference in the various data sets based on our estimates of  $\Delta M$

Amino Acid	<i>E. gossypii</i>	Endogenous	Exogenous	<i>L. kluyveri</i>
Ala A	GCG	GCA	GCG	GCG
Cys C	TGC	TGT	TGC	TGC
Asp D	GAC	GAT	GAC	GAC
Glu E	GAG	GAA	GAG	GAG
Phe F	TTC	TTT	TTT	TTT
Gly G	GGC	GGT	GGC	GGC
His H	CAC	CAT	CAC	CAC
Ile I	ATC	ATT	ATC	ATA
Lys K	AAG	AAA	AAG	AAA
Leu L	CTG	TTG	CTG	CTG
Asn N	AAC	AAT	AAC	AAT
Pro P	CCG	CCA	CCG	CCG
Gln Q	CAG	CAA	CAG	CAG
Arg R	CGC	AGA	AGG	CGG
Ser <sub>4</sub> S	TCG	TCT	TCG	TCG
Thr T	ACG	ACA	ACG	ACG
Val V	GTG	GTT	GTG	GTG
Tyr Y	TAC	TAT	TAC	TAC
Ser <sub>2</sub> Z	AGC	AGT	AGC	AGC

Table 3.3: Synonymous codon preference in the various data sets based on our estimates of  $\Delta\eta$

Amino Acid	<i>E. gossypii</i>	Endogenous	Exogenous	<i>L. kluyveri</i>
Ala A	GCT	GCT	GCT	GCT
Cys C	TGT	TGT	TGT	TGT
Asp D	GAT	GAC	GAT	GAT
Glu E	GAA	GAA	GAA	GAA
Phe F	TTT	TTC	TTC	TTC
Gly G	GGA	GGT	GGT	GGT
His H	CAT	CAC	CAT	CAT
Ile I	ATA	ATC	ATT	ATT
Lys K	AAA	AAG	AAA	AAG
Leu L	TTA	TTG	TTG	TTG
Asn N	AAT	AAC	AAT	AAC
Pro P	CCA	CCA	CCT	CCA
Gln Q	CAA	CAA	CAA	CAA
Arg R	AGA	AGA	AGA	AGA
Ser <sub>4</sub> S	TCA	TCC	TCT	TCT
Thr T	ACT	ACC	ACT	ACT
Val V	GTT	GTC	GTT	GTT
Tyr Y	TAT	TAC	TAT	TAC
Ser <sub>2</sub> Z	AGT	AGT	AGT	AGT

Taxon	Abbreviation	NCBI taxonomic ID	Codon Table	% GC	GC Source
<i>Candida albicans</i>	Calb	5476	12	34	NCBI Genom
<i>Saccharomyces bayanus</i>	Sbay	4931	1	40	NCBI Genom
<i>Trichophyton benhamiae</i>	Tben	63400	1	49	NCBI Genom
<i>Tetrapisispora blattae</i>	Tbla	1071379	1	32	NCBI Genom
<i>Saccharomyces castellii</i>	Scas	27288	1	37	NCBI Genom
<i>Saccharomyces cerevisiae</i>	Scer	4932	1	38	NCBI Genom
<i>Eremothecium cymbalariae</i>	Ecym	45285	1	40	NCBI Genom
<i>Torulaspora delbrueckii</i>	Tdel	4950	1	42	NCBI Genom
<i>Candida dubliniensis</i>	Cdub	42374	12	33	NCBI Genom
<i>Lodderomyces elongisporus</i>	Lelo	36914	1	37	NCBI Genom
<i>Saccharomyces eubayanus</i>	Seub	1080349	1	40	NCBI Genom
<i>Debaryomyces fabryi</i>	Dfab	58627	1	36	NCBI Genom
<i>Candida glabrata</i>	Cgla	5478	1	39	NCBI Genom
<i>Eremothecium gossypii</i>	Egos	33169	1	52	NCBI Genom
<i>Meyerozyma guilliermondii</i>	Mgui	4929	12	44	NCBI Genom
<i>Debaryomyces hansenii</i>	Dhan	4959	12	36	NCBI Genom
<i>Lachancea kluyveri</i>	Lku	4934	1	40/53	Payen et al.
<i>Saccharomyces kudriavzevii</i>	Skud	114524	1	41	NCBI Genom
<i>Kluyveromyces lactis</i>	Klac	28985	1	39	NCBI Genom
<i>Lachancea lanzarotensis</i>	Llan	1245769	1	44	NCBI Genom
<i>Yarrowia lipolytica</i>	Ylip	4952	1	49	NCBI Genom
<i>Clavispora lusitaniae</i>	Clus	36911	12	45	NCBI Genom
<i>Kluyveromyces marxianus</i>	Kmar	4911	1	40	NCBI Genom
<i>Saccharomyces mikatae</i>	Smik	114525	1	38	NCBI Genom
<i>Sphaerulina musiva</i>	Smus	85929	1	51	NCBI Genom
<i>Kazachstania naganishii</i>	Knag	588726	1	46	NCBI Genom
<i>Saccharomyces paradoxus</i>	Spar	27291	1	38	NCBI Genom
<i>Candida parapsilosis</i>	Cpar	5480	12	38	NCBI Genom
<i>Spathaspora passalidarum</i>	Spas	340170	12	38	NCBI Genom
<i>Tetrapisispora phaffii</i>	Tpha	113608	1	34	NCBI Genom
<i>Vanderwaltozyma polyspora</i>	Vpol	36033	1	33	NCBI Genom
<i>Lachancea quebecensis</i>	Lque	1654605	1	47	Freel et al. 2
<i>Zygosaccharomyces rouxii</i>	Zrou	4956	1	40	NCBI Genom
<i>Scheffersomyces stipitis</i>	Ssti	4924	12	41	NCBI Genom
<i>Lachancea thermotolerans</i>	Lthe	381046	1	47	NCBI Genom
<i>Candida tropicalis</i>	Ctro	5482	12	33	NCBI Genom
<i>Lachancea waltii</i>	Lwal	4914	1	44	NCBI Genom
<i>Cladophialophora yegresii</i>	Cyeg	470704	1	54	NCBI Genom

Table 3.4:



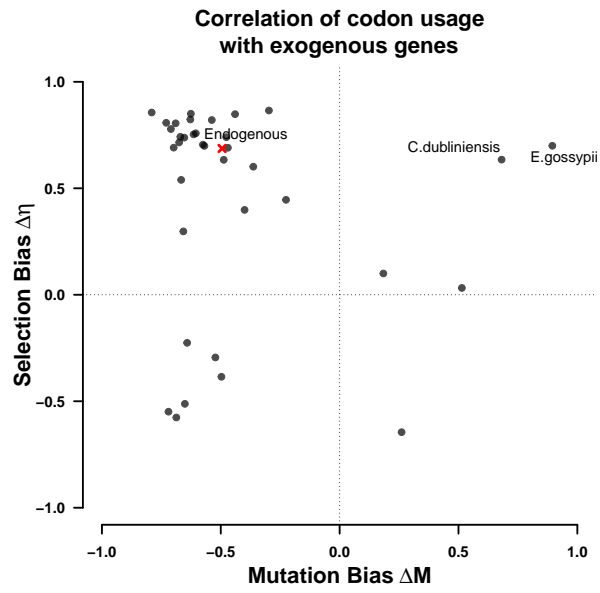


Figure 3.5: Correlation coefficient of  $\Delta M$  and  $\Delta\eta$  of the endogenous genes with 38 examined yeast lineages. Dots indicate the correlation of  $\Delta M$  and  $\Delta\eta$  of the lineages with the endogenous and exogenous parameter estimates. All regressions were performed using a type II regression line (SOKAL and ROHLF, 1981).



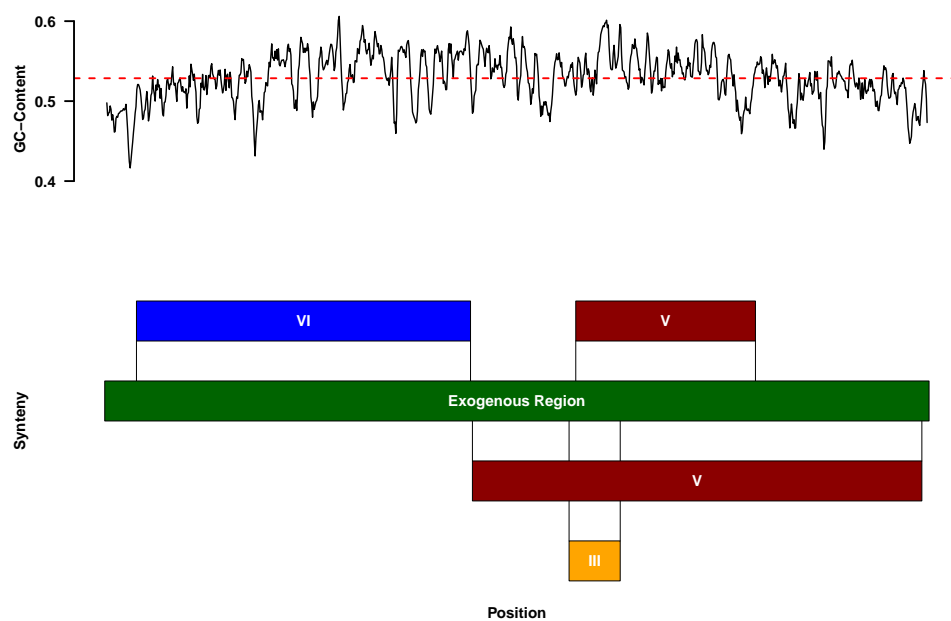


Figure 3.7: Synteny relationship of *E. gossypii* and the exogenous genes. Indicated is the GC content along the introgression.

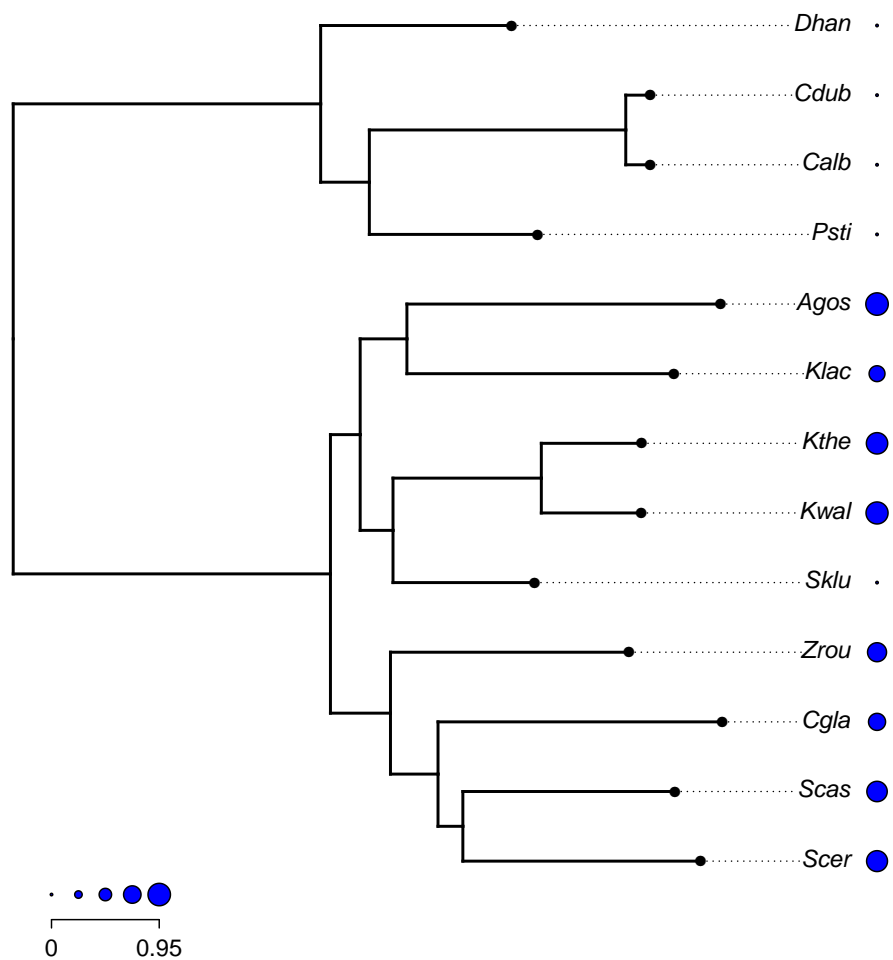


Figure 3.8: Amount of synteny for each species in units of standard deviations for selected species.

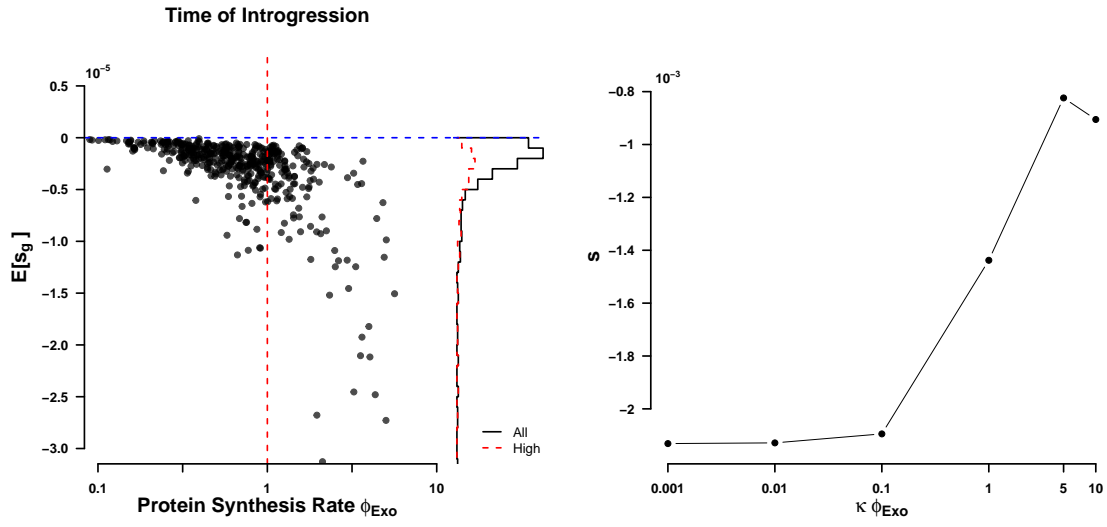


Figure 3.9: Genetic load (left) without scaling of  $\phi$  per gene, and change of total genetic load with scaling  $\kappa$  between *E. gossypii* and *L. kluyveri* (right)

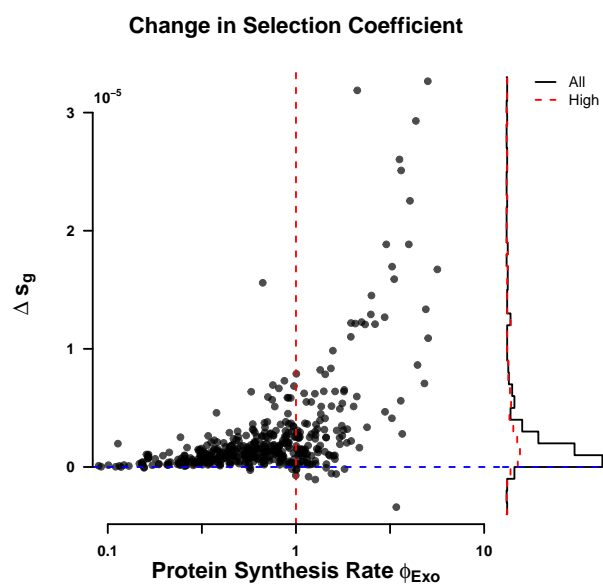


Figure 3.10: Total amount of adaptation estimated to have occurred between time of introgression and currently observed per gene.

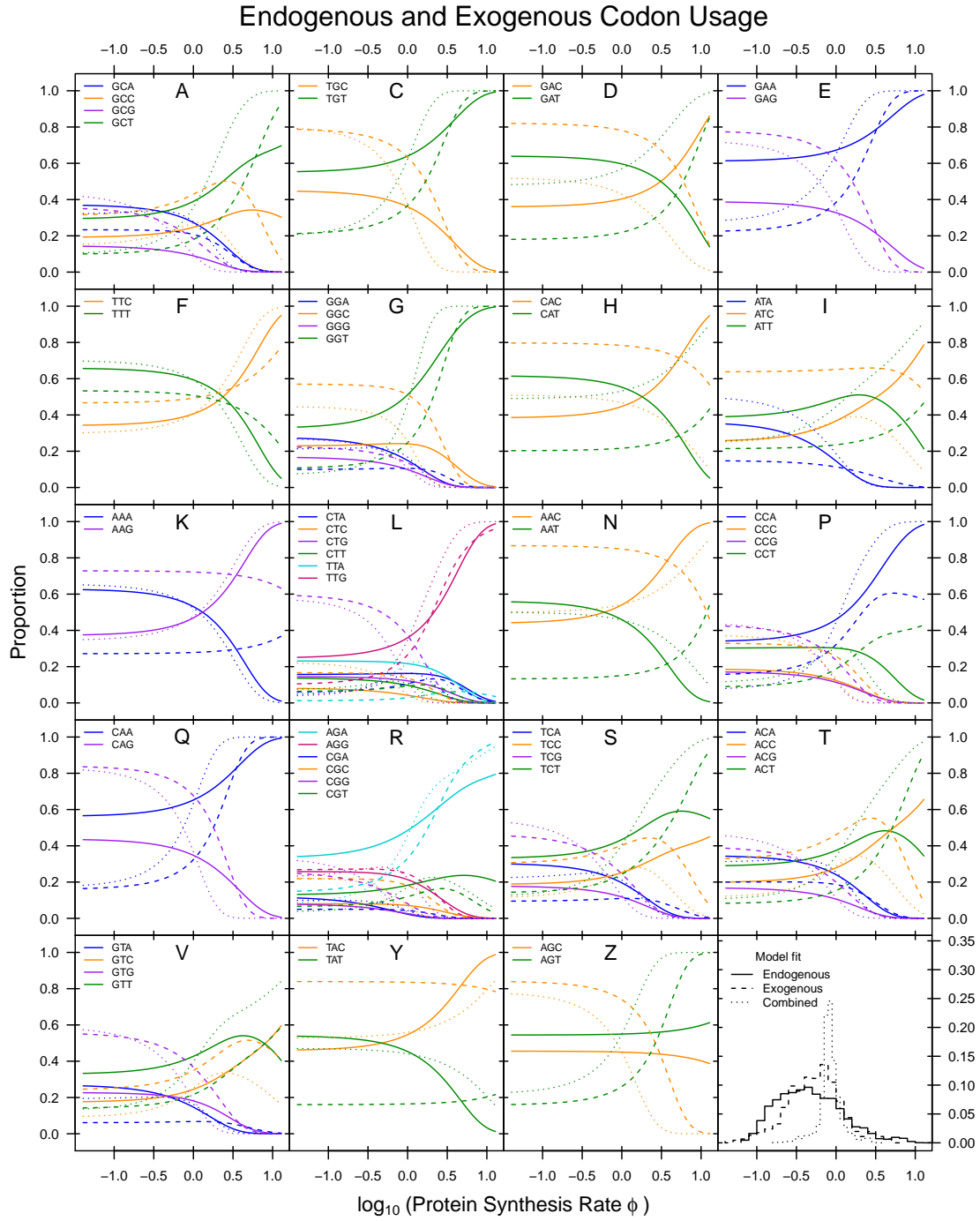


Figure 3.11: Codon usage patterns for 19 amino acids. Amino acids are indicated as one letter code. The amino acids Serine was split into two groups (S and Z) as Serine is coded for by two groups of codons that are separated by more than one mutation. Solid line indicates the endogenous codon usage, dashed line indicates the exogenous codon usage, dotted line indicates the combined codon usage.

## Chapter 4

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



This chapter is an early version of a paper to be submitted to Genome Biology and Evolution and co-authored with Michael A. Gilchrist and Brian C. O’Meara.

C. Landerer, B.C. O’Meara, M.A. Gilchrist, Phylogenetic model of stabilizing selection is more informative about site specific selection than extrapolation from laboratory estimates

## 4.1 Abstract

Here we examine the ability of experimentally inferred, site specific selection for amino acids to improve phylogenetic inferences of sequence evolution. Previous work has shown that laboratory estimates of selection can improve model fit, but did not assess their adequacy. We assess the ability of experimentally inferred site specific selection for antibiotic resistance from deep mutation scanning to inform phylogenetic models. In this study, we use the  $\beta$ -lactamase TEM for which empirical estimates of site specific selection on amino acids are readily available. TEM is an enzyme that catalyzes antibiotics with a  $\beta$ -lactam ring and is found in gram-negative bacteria like *Escherichia coli*. We compare the experimentally inferred site specific selection to our results obtained using *SelAC*, a new phylogenetic model of stabilizing selection. Using simulations we assess model adequacy, and find that experimentally inferred selection does not adequately reflect evolution in the wild. In contrast, *SelAC* fits to the data better over models informed by experimentally inferred selection and provides higher model adequacy. We demonstrate the capability of *SelAC* by estimating the site specific genetic load of the observed TEM variants.

## 4.2 Introduction

Numerous attempts to incorporate selection into phylogenetic models have been made. Early models only focused on the influence of selection on the substitution rate and fixation probability between a resident and a mutant introduced into a population (GOLDMAN and YANG, 1994; MUSE and GAUT, 1994; THORNE *et al.*, 1996). These models however, lack site specific equilibrium codon or amino acid frequencies. The importance of site specific equilibrium frequencies has long been noted (FELSENSTEIN, 1981; GOJOBORI, 1983). Individual amino acid sites along the protein show differences in evolutionary rates, and wide range of preferences for specific amino acids (HALPERN and BRUNO, 1998; ASHENBERG *et al.*, 2013; ECHAVE *et al.*, 2016). The usage of site specific selection acknowledges the heterogeneity in selection and amino acid preferences along the protein sequence (HILTON *et al.*, 2017).

HALPERN and BRUNO (1998) first introduced a model to incorporate site specific equilibrium frequencies of amino acids. However, they had to concede that their model was too parameter rich and therefore intractable for biological data sets without additional simplifying assumptions. More recent models incorporating site specific equilibrium frequencies still require a large number of parameters to be estimated from the sequence data (LARTILLOT and PHILIPPE, 2004; LE *et al.*, 2008; WANG *et al.*, 2008; HOLDER *et al.*, 2008; WU *et al.*, 2013; TAMURI *et al.*, 2014). Other approaches treat site specific selection as a random effect (RODRIGUE *et al.*, 2010; RODRIGUE, 2013; RODRIGUE and LARTILLOT, 2014). A full parameterization of site specific equilibrium frequencies for amino acids requires  $19 \times N$  parameters where  $N$  is the length of the sequence in amino acids. It is therefore an attractive option to utilize laboratory experiments to empirically estimate site specific strength of selection on amino acids and infer their equilibrium frequencies (BLOOM, 2014; THYAGARAJAN and BLOOM, 2014; BLOOM, 2017).

Empirical estimates of site specific selection can greatly reduce the number of parameters estimated from phylogenetic data, making it applicable for smaller data sets and allowing

for the fitting of more complex models. Deep mutation scanning (DMS) has recently been used to generate comprehensive site specific estimates of selection (FOWLER *et al.*, 2014). The ability to estimate site specific selection allows to estimate site specific amino acid preferences and the fitness consequences a mutation introduces at a particular site (BLOOM, 2014; FIRNBERG *et al.*, 2014; STIFFLER *et al.*, 2016). There are, however, also shortcomings. The quality of empirical estimates from DMS, however, depends on many factors including the initial library of mutants and the applied selection (FIRNBERG and OSTERMEIER, 2012). Mutation libraries have to be extensive and, therefore, produce a heterogeneous population of competing organisms not usually found in nature. In addition, estimates of selection can only be obtained for fast growing organisms that can be manipulated under laboratory conditions. As many organism can not be cultivated under laboratory conditions or have long generation times, this is a severe limitation to experimentally informed models.

Even in the cases where empirical estimates of site specific selection can be obtained, their utility for phylogenetic reconstruction is questionable. In this study, we assess the ability of experimentally inferred site specific selection to inform phylogenetic models and offer an alternative approach to determine site specific selection. We use site specific estimates of selection for the class A  $\beta$ -lactamase TEM from STIFFLER *et al.* (2016). TEM is an enzyme found in gram-negative bacteria like *Escherichia coli* that catalyzes antibiotics with a  $\beta$ -lactam ring and provides antibiotic resistance (NEU, 1969). The selection pressure imposed during the DMS experiment was limited to ampicillin and focused solely on TEM-1 (STIFFLER *et al.*, 2016). However, TEM variants can also confer resistance to a wide range of other antibiotics (SOUGAKOFF *et al.*, 1988, 1989; GOUSSARD *et al.*, 1991; MABILAT *et al.*, 1992; CHANAL *et al.*, 1992; BRUN *et al.*, 1994).

In order to do so, we fitted 227 nucleotide and codon models using IQTree and compared their model fits to site specific models of stabilizing selection with (*phydms*, *SelAC*+DMS) and without (*SelAC*) experimentally estimated site specific selection coefficients (NGUYEN *et al.*, 2015; HILTON *et al.*, 2017; BEAULIEU *et al.*, in review). We find that experimentally

inferred selection, while improving model fit, does not adequately reflect observed wild type sequences. In contrast, *SelAC* (BEAULIEU *et al.*, in review) a mechanistic phylogenetic model of stabilizing selection rooted in first principles with site specific equilibrium frequencies improves model fit, and better predicts sequences found in the wild. In addition, it was previously proposed to extrapolate to the fitness landscape of related proteins using experimentally inferred site specific selection (BLOOM, 2014, 2017). We utilize the TEM homologue SHV, another class A  $\beta$ -lactamase, to demonstrate the problematic with this approach and further highlight the generality of *SelAC*.

## 4.3 Results

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

We compared *phydms* (HILTON *et al.*, 2017) and *SelAC* (BEAULIEU *et al.*, in review), models of site specific stabilizing selection on amino acids, to 227 other codon and nucleotide models. We fitted all models to 49 observed sequences of the  $\beta$ -lactamase TEM (BLOOM, 2014). The *phydms* and *SelAC* models with site specific selection improved model fits by 366 and 934 AICc units, respectively, over the best performing codon or nucleotide models which lacks site specific selection (Table 4.1). In addition, *SelAC* outperformed the experimentally informed model *phydms* by 562 to 568 AICc units, depending whether site specific selection was inferred by *SelAC* or experimentally informed.

*SelAC* utilizes a hierarchical model and estimates 263 site specific parameters,  $\sim 5\%$  of the  $19 \times N = 4997$  parameters necessary to fully describe site specific selection. In contrast, *phydms* does not infer any site specific parameters from the phylogenetic data, but utilizes site specific selection estimated from deep mutation scanning experiments. In order to assess the quality of the *SelAC* fit with experimentally determined site specific selection, we fixed the optimal amino acid at each site to the experimentally determined one in *SelAC*

Table 4.1: 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 (GOLDMAN and YANG, 1994; ZHARKIKH, 1994). Reported are the log-likelihood  $\log(\mathcal{L})$ , the number of parameters estimated  $n$ , AIC,  $\Delta$ AIC, AICc, and  $\Delta$ AICc values. See Table 4.3 for results from all models we tested.

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

and refitted the model to the 49 TEM sequences (*SelAC*+DMS). Incorporating site specific selection estimated from deep mutation scanning experiments into *SelAC* (*SelAC*+DMS) yields a similar, but slightly better AICc value to *SelAC* without that information.

This improvement in AICc of *SelAC*+DMS over *SelAC* is solely due to a decrease in the number of parameters estimated. In contrast to AICc, the log-likelihood  $\log(\mathcal{L})$  of *SelAC*+DMS is 270  $\log(\mathcal{L})$  units worse than the *SelAC* one (Table 4.1). However, it is statistically speaking unclear if discrete parameters bias ones estimate the Kullback-Leibler divergence in the same way. This is important since 263 of the 374 parameters estimated by *SelAC* are the discrete optimal amino acid state at each site. Thus, it is possible that we are over penalizing. Therefore, the number of parameter for *SelAC* we use is conservative. For example, there are only 27 unique site patterns in the TEM alignment, which would yield a total of 138 parameters. This however would likely be an under estimate of the number of parameters estimated. The true number of parameters remains unclear at this point due to the inherent non-independence of the underlying data and the discrete nature of the optimized parameters.

We observe differences in the topology between model fits. The *SelAC* model is currently too slow to estimate the topology, therefore the topology was estimated using the codon model of KOSIOL *et al.* (2007). At this point, it is therefore unclear if the difference in topology can be attributed to the experimentally inferred selection. We find that the best

codon model (*GY94*) (GOLDMAN and YANG, 1994) is outperformed by several nucleotide model e.g. *SYM*+R2 (ZHARKIKH, 1994). This could be an indication that negative frequency dependent selection like it is modeled in *GY94* is not appropriate for TEM (GOLDMAN and YANG, 1994; BEAULIEU *et al.*, in review). Figure 4.1 shows that the estimated phylogenetic trees shift from long terminal branches (*SelAC*) to longer internal branches (*phydms*). While the *SelAC* model fit shows 84% of all evolution happening at the tips, this reduces to 79% in the *SelAC*+DMS model fit, and 77% in the *phydms* and *GY94* model fits. All models produce polytomies but their location differs between models. Surprisingly, the largest polytomies appear in the experimentally informed phylogeny of *phydms*. The position of the sequences with the longest branches also differ between *SelAC* and *phydms*.

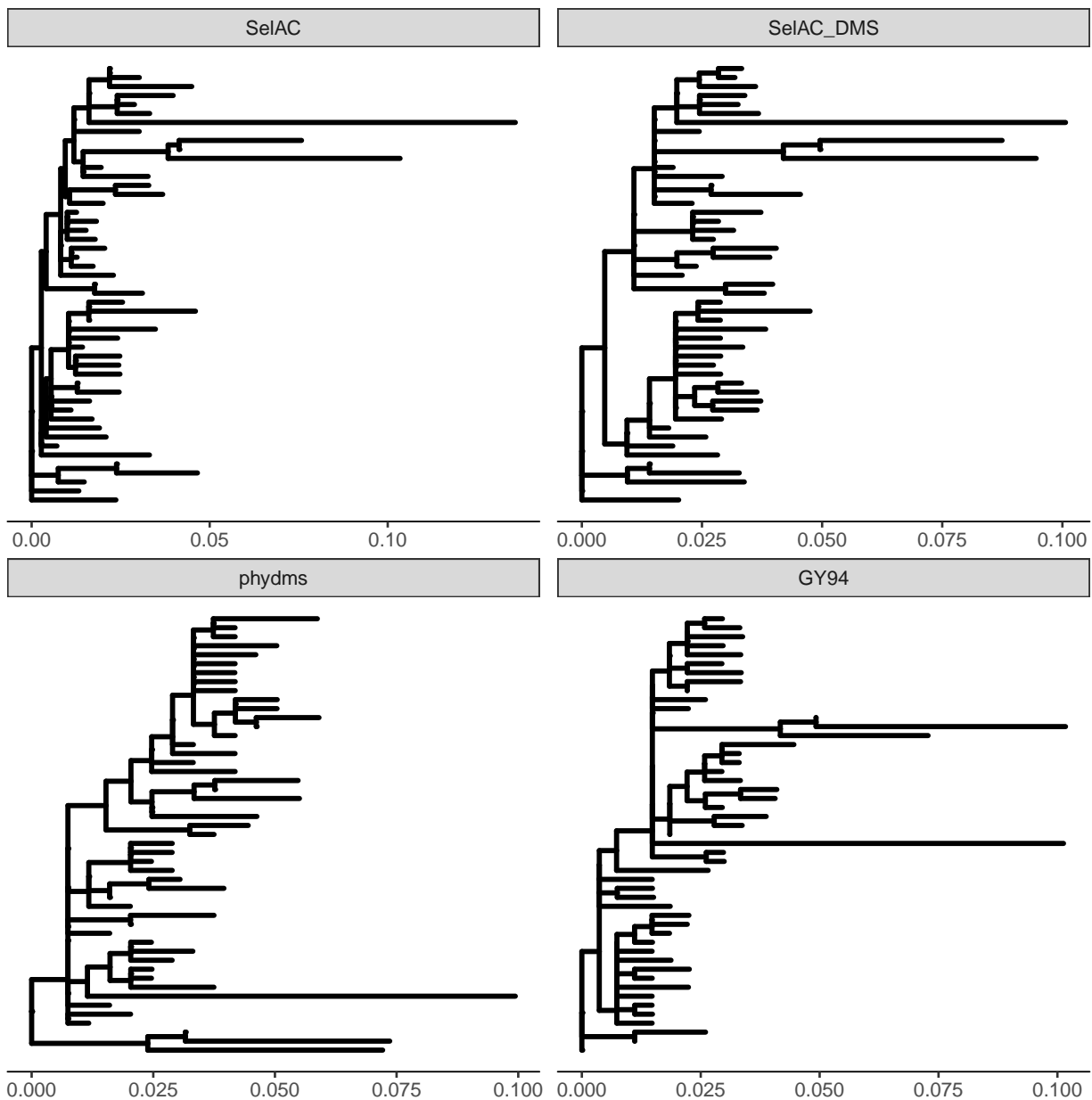


Figure 4.1: Phylogenies resulting from *SelAC*, *SelAC*+DMS, *phydms*, and *GY94*. As *SelAC* is currently too slow for the inference of topologies, the topology for the *SelAC* phylogenies was inferred using the codon model of KOSIOL *et al.* (2007).

TEM2016_SelAC/1-263	1	HPETLVKVKDAEDQLGARVGYIELDLNSGKILESFRPEERFPMMS	53
TEM2016_SelAC_Simulated/1-263	1	HPETRVKVKGAEECLGAGRGYIELDLNSGKILESFRPEERFPMRSTFKVLLCG	53
TEM2016_Consensus/1-263	1	HPETLVKVKDAEDQLGARVGYIELDLNSGKILESFRPEERFPMMS	53
TEM2016_DMS/1-263	1	SEKVKMAVQQMEWRMGHVGFIQIDIMDGDVLEAWRSKERFPMMS	53
TEM2016_DMS_Simulated/1-263	1	HEKTKTKVRDAERRMGRVGYLQIDIHDDVLESFRQKERFPMMS	53
TEM2016_SelAC/1-263	54	AVLSRVDAGQEQLGRRIHYSQNDLVEYSPVTEKHLTDGMTVRE	106
TEM2016_SelAC_Simulated/1-263	54	AELSRGDAGQEQLGRRIHYSQADEVEYSPVTEKHLTDGMTVRE	106
TEM2016_Consensus/1-263	54	AVLSRVDAGQEQLGRRIHYSQNDLVEYSPVTEKHLTDGMTVRE	106
TEM2016_DMS/1-263	54	CILERVNDNGFLKLRQKVKFQVNDLVAWSPITMMYIITGMTIQDL	106
TEM2016_DMS_Simulated/1-263	54	AILYRVDAETELGRRVHFVTVNDLVAYSPITSQYINDGMTIAD	106
TEM2016_SelAC/1-263	107	NTAANLLLTITIGGPKELTAF LHNMGDHSVTRLDRWEPELNEAIPN	159
TEM2016_SelAC_Simulated/1-263	107	NTAADLLLTTIGGRGELTAF LHNMTDHSVTRLARGAPELGEAIPG	159
TEM2016_Consensus/1-263	107	NTAANLLLTITIGGPKELTAF LHNMGDHSVTRLDRWEPELNEAIPN	159
TEM2016_DMS/1-263	107	NTAANI L LKELGGPIMLT MWMNMMDMYTRLDRWEPYLNMA	159
TEM2016_DMS_Simulated/1-263	107	NTAANI L LKSLGGPIELTEYMNMGDNVTRLDRWEPYLNAATP	159
TEM2016_SelAC/1-263	160	AMATT LRKLLTGELLTLASRQQLIDWMEADKVAGPLLRSLPAGWF	212
TEM2016_SelAC_Simulated/1-263	160	AMATTLRGLLTEELLTLASRARLIDWMEADKVAGPLLRSLPAGWF	212
TEM2016_Consensus/1-263	160	AMATT LRKLLTGELLTLASRQQLIDWMEADKVAGPLLRSLPAGWF	212
TEM2016_DMS/1-263	160	SMADTIKQMLKTHHSFNSSQILISWMYMDKVAGPLLRQKIPADWY	212
TEM2016_DMS_Simulated/1-263	160	VMAKTIHELKDHRLSKGSSQILIEWMKLDKVAGPLLRQAIPADWY	212
TEM2016_SelAC/1-263	213	GERGSRGIIAALGPDGKPSRI VVIYMTGSQATMDERNRQIAEIGASL	263
TEM2016_SelAC_Simulated/1-263	213	EVRGSGGIIAALGPDGKPSRI VVIYVTGRQATMDERSRQGEI	263
TEM2016_Consensus/1-263	213	GERGSRGIIAALGPDGKPSRI VVIYTTGSQATMDERNRQIAEIGASL	263
TEM2016_DMS/1-263	213	GDHGSRGIVALMGPNKHMERV I IYMTGSNANMIQRNQWFKEIGKN	263
TEM2016_DMS_Simulated/1-263	213	GKHGSRGIVAAIGPAGVASRVI IYLTGSNNNMDARNQWF AEIGKN	263

Figure 4.2: Alignment of TEM optimal and simulated sequences. Indicated is the percentage identity at each site.

### 4.3.2 Laboratory Inferences Inconsistent with Observed Sequences.

The improved model fits of *phydms* relative to more common nucleotide and codon models are, however, deceiving. The site specific selection inferred by DMS is inconsistent with the observed TEM sequences. We find that the sequence of selectively favored amino acids has only 52% sequence similarity with the observed consensus sequence (Figure 4.2). In contrast, the sequence of selectively favored amino acids estimated by *SelAC* shows 99% sequence similarity with the observed consensus sequence. In addition, 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.

In order to determine if we would expect the observed genetic load under the experimental selection estimates we reconstructed the ancestral TEM sequence and used it as initial conditions for our simulation studies. We simulate under a wide range of effective population



sizes  $N_e$ , and find that the experimentally inferred site specific selection is very strong. The estimated ancestral state is identical to the observed consensus sequence. Simulations of codon sequences under the experimentally inferred site specific selection for amino acids reveal that we would not expect to see the observed TEM sequences. With an effective population size  $N_e$  of  $10^7$ , we find that the simulated sequences show 62% sequence similarity to the observed consensus sequence (Figure 4.3a). Thus, the simulated sequences show a 10% higher similarity to the observed consensus sequence than the sequence of selectively favored amino acids estimated using deep mutation scanning.

In our simulations, only when  $N_e$  is reduced to one individual does drift overwhelm selection (Figure 4.3b). The genetic load of the simulated sequences decrease slowly with increasing  $N_e$ . After simulating until the sequences reached 0.1 expected mutation per site with an effective population size  $N_e = 10^7$  the simulated sequences showed an average sequence specific load of 6.68 and an average site specific genetic load of 0.025. This is less than half of the average sequence and site specific genetic load of the observed sequences. Thus it appears unlikely that the observed sequences have evolved under the DMS inferred site specific selection values.

### 4.3.3 Stabilizing Selection for Optimal Physicochemical Properties Improves Model Adequacy

Model adequacy of *SelAC* assessed based on sequence similarity and genetic load and shows that *SelAC* better explains the observed TEM sequences than the experimentally determined site specific selection on amino acids. The observed consensus sequence has 99% sequence similarity with the sequence of selectively favored amino acids estimated by *SelAC*, this is in contrast to the average sequence similarity of 98% among all 49 observed sequences. In addition, 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  $2.4 \times 10^{-7}$ .

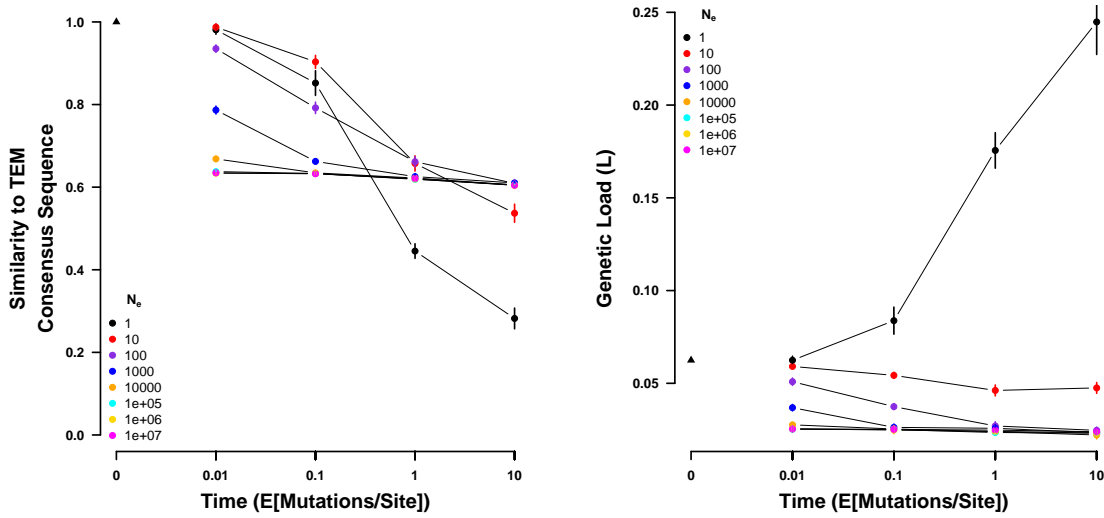


Figure 4.3: Sequences simulated from the ancestral state under the site specific selection on amino acids estimated using deep mutation scanning. (left) Sequence similarity to the observed consensus sequence at various times for a range of values of  $N_e$ . (right) Genetic load of the simulated sequences at various times for a range of values of  $N_e$ . Time is given in number of expected mutations per site, which equals the substitution rate of a neutral mutation. Points indicate sample means and vertical bars indicate standard deviations. Initial sequence is the inferred ancestral state of the TEM variants and indicated by a black triangle.

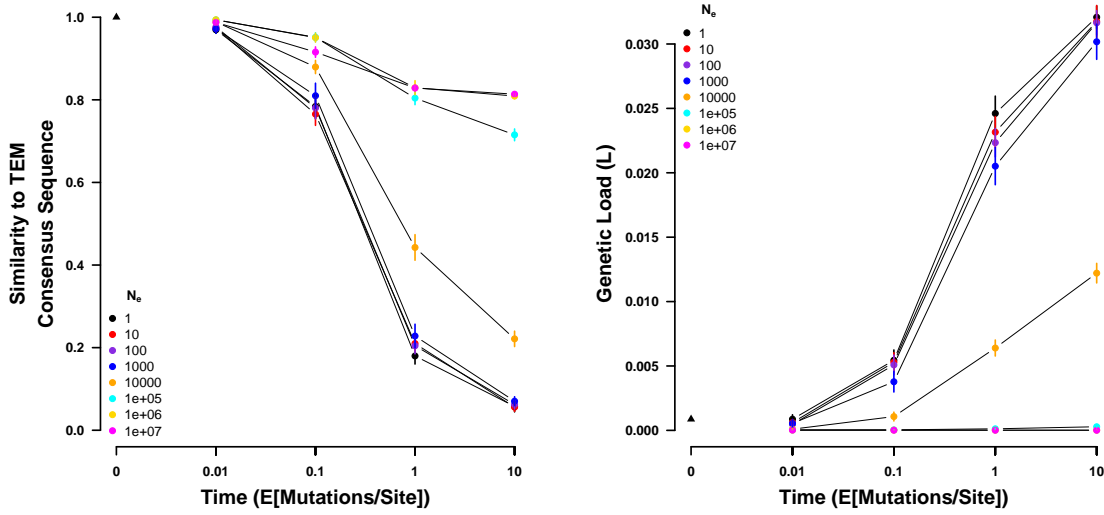


Figure 4.4: Sequences simulated from the ancestral state under the site specific selection on amino acids estimated using *SelAC*. (left) Sequence similarity to the observed consensus sequence at various times for a range of values of  $N_e$ . (right) Genetic load of the simulated sequences at various times for a range of values of  $N_e$ . Time is given in number of expected mutations per site, which equals the substitution rate of a neutral mutation. Points indicate sample means and vertical bars indicate standard deviations. Initial sequence is the inferred ancestral state of the TEM variants and indicated by a black triangle.

Using the *SelAC* inferred site specific selection and the reconstructed ancestral TEM sequence as initial conditions we simulated codon sequences forward in time for various time periods and  $N_e$  values. As expected, for small  $N_e$ , simulated sequences drift away from the observed consensus sequence (Figure 4.4a). Because of the high similarity between the optimal amino acid sequence estimated by *SelAC* and the observed consensus sequence, the genetic load increases drastically as a result. Increasing  $N_e$  to  $10^4$  or above the simulated sequences sequence similarity declines to 83%, indicating that *SelAC* underestimates selection. After simulating until the sequences reached 0.1 expected mutation per site with an effective population size  $N_e = 10^7$  the simulated sequences showed an average sequence specific load of  $1.3 \times 10^{-5}$  and an average site specific genetic load of  $4.8 \times 10^{-8}$  (Figure 4.4b). Thus, the simulated sequences show a lower genetic load despite the greater divergence from the observed consensus sequence. This be an indication that the selection differs between lineages.

To further demonstrate the consistency of *SelAC*, we simulated codon sequences over the same time periods using 10 sequences where codons were sampled uniformly. We find that the sequence similarity increases with effective population size  $N_e$ . The random sequences start off with a similarity of  $\sim 6\%$  and increase  $N_e$  to  $\sim 28\%$  (Figure 4.8a). The same initial sequences simulated under the site specific selection inferred by the deep mutation scanning experiment increase only to  $\sim 18\%$  in sequence similarity over the same period of time.

### 4.3.4 Estimating Site Specific Selection on Amino Acids

*SelAC* allows for the estimation of site specific selection on amino acids and the genetic load of an observed amino acid relative to the inferred optimal amino acid. Figures 4.5 and 4.6 illustrate how the genetic load varies along the TEM sequence. The region between residue 80 and 120, where three consecutive helices are located, consists only of selectively favored amino acids and does not show any genetic load. The highest genetic load is found in the unstructured regions and the lowest genetic load is found in  $\beta$ -sheets. However, this difference is not statistically significant ( $p = 0.17$ ). The largest increase in genetic load is located at the beginning of the last helix. This region strongly contributes to the estimate of similar genetic loads for helices and unstructured regions in the observed TEM sequences (Table 4.2). However, exclusion of this site as outlier does not yield significance ( $p = 0.09$ ).

*SelAC* assumes that the efficacy of selection  $G$  is  $\Gamma$ -distributed with a mean of 1. However, it is possible to estimate site specific values using the parameters estimated by *SelAC*. We constraint  $G$  to a maximum value of 300 in all cases. While this biases our estimate of  $G$ , the bias is consistent across all estimates and does not prohibit the comparison of  $G$  terms.

The highest efficacy of selection  $G$  is estimated in the  $\beta$ -sheet regions which is consistent with the lowest genetic load in these regions. Residues forming the substrate binding site appear to be under the strongest selection, with no accumulated genetic load. However, this is not the case for the two active sites. We find in one sequence (*Acinetobacter baumannii*, TEM-193) a Lysine, a proton donor, at the proton acceptor site 143 driving the reduced

Table 4.2: 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	$G$		Genetic Load $L_i$	
			Mean	SE	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	$1.5 \times 10^{-8}$	$1.5 \times 10^{-8}$

efficacy of selection  $G$ . This is in concordance with the experimental DMS estimates, where proton acceptors are selectively favored. Again, any differences between secondary structure elements are not statistically significant.

It was previously proposed that experimentally inferred site specific selection for amino acids can be used to extrapolate the fitness landscape of related proteins (BLOOM, 2014, 2017). We therefore compared the genetic load, the *SelAC* selection parameters of our *SelAC* TEM model fit to a *SelAC* model fit of SHV, and site specific efficacy of selection  $G$ . The genetic load observed in SHV sequences appears to be lower than in TEM with the exception of residues found in  $\beta$ -sheets and the active site (Table 4.2). This is consistent with the elevated efficacy of selection  $G$  in SHV. However, only differences in genetic load in the unstructured regions are significantly different between the TEM and SHV sequences, but only at the  $\alpha = 0.05$  significant level ( $p = 0.04$ ). While the average genetic load across secondary structures is not significantly different, the sites causing increases genetic load differ between SHV and TEM (Figure 4.7). In contrast to TEM, we find the highest genetic load among SHV secondary structure features in the helices (Table 4.2). The highest genetic load in SHV is observed at the end of the first helix. We do find a peak of similar magnitude

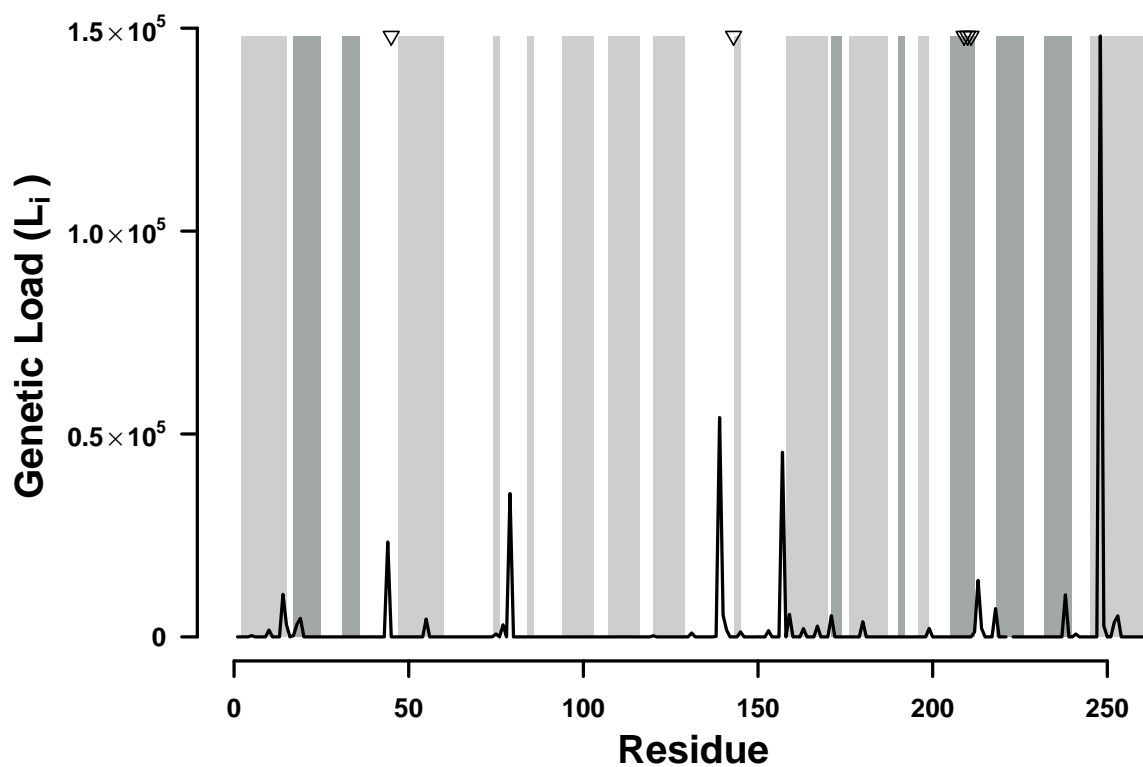


Figure 4.5: Distribution of average site specific genetic load in TEM over all observed TEM variants. Average site specific genetic load is indicated by the black line. Light gray bars indicate where helices are found, and dark gray bars indicate  $\beta$ -sheets. The three residues forming the binding site and the two residues forming the active are indicated by triangles at the top of the plot.

in the TEM sequence at the end of the first helix, but this peak is overshadowed by the increased genetic load at the beginning of the last helix.

We find that the site specific efficacy of selection  $G$  differs greatly between SHV and TEM ( $\rho = 0.12$ ), despite a similar estimate of  $\alpha_G$  describing the distribution of  $G$  values (Figure 4.10a). Most *SelAC* selection parameters are very similar between the TEM and the SHV model fit. An exception is the weight for the physicochemical composition property  $\alpha_c$  (Figure 4.10b). Furthermore, we find that the sequences of selectively favored amino acids estimated by *SelAC* for TEM and SHV only show 68% sequence similarity. These results indicate that the extrapolation from one proteins fitness landscape to another is problematic.

## 4.4 Discussion

Here we revisited how well site specific estimates of selection from deep mutation scanning experiments inform sequence evolution and compared it to *SelAC*, a novel phylogenetic model of stabilizing selection. Previous work has shown that laboratory estimates of selection can improve model fit over classical approaches like GY94 (BLOOM, 2014, 2017). While our study confirms this notion, we identify important shortcomings of these laboratory estimates for phylogenetic studies. In contrast, *SelAC* is a phylogenetic model of stabilizing selection based on physicochemical properties and does not require costly laboratory estimates of selection and is, nevertheless, favored by model selection (Table 4.1). More specifically, it estimates site specific selection on amino acids from the sequence data based on distances between amino acids in physicochemical space (GRANTHAM, 1974; BEAULIEU *et al.*, in review). This allows *SelAC* to be applied to any set of protein coding sequences, eliminating the need to extrapolate from one homologous gene family to the next (e.g. from TEM to SHV). In addition this generality allows for the comparison of model parameters for these proteins.

While previous work showed the advantages of experimentally informed phylogenetics, they did not assess how adequate the estimated selection reflects observed wild-type

sequences. The low sequence similarity between the observed consensus sequence and the sequence of selectively favored amino acids estimated by deep mutation scanning experiments is evidence for that. This begs the question how well the experimental selection coefficients represent selection on these sequences in nature. Deep mutation scanning experiments are performed using a comprehensive library of mutants and a strong artificial selection pressure (FIRNBERG and OSTERMEIER, 2012; JAIN and VARADARAJAN, 2014; FOWLER and FIELDS, 2014; FOWLER *et al.*, 2014). This results in very large selection coefficients  $s$  and a heterogeneous population of competing individuals unlikely to occur in nature.

The selection pressure imposed during the DMS experiment was limited to ampicillin and focused solely on TEM-1 (STIFFLER *et al.*, 2016). However, TEM variants can also confer resistance to a wide range of other antibiotics, including penicillins, cephalosporins, cefotaxime, ceftazidime, or aztreonam (SOUGAKOFF *et al.*, 1988, 1989; GOUSSARD *et al.*, 1991; MABILAT *et al.*, 1992; CHANAL *et al.*, 1992; BRUN *et al.*, 1994). Thus, the inferred selection is biased towards ampicillin and is inconsistent with the observed TEM sequences (Figure 4.3). This may very well be appropriate to explore the selection on TEM in a hospital environment but is unlikely to be representative of the selection faced by *E. coli* and other gram-negative bacteria in nature.

If we assume that the DMS selection coefficients underly the evolution of the observed TEM sequences we can think of two possible explanations for the observed sequences. First, the sequences are unable to reach a fitness peak, potentially due to a weak selection pressure or not enough time. Alternatively, the TEM sequences found in nature are highly maladapted, yet with very similar sequences. Both explanations seem unlikely. For example, *E. coli* has a large effective population size  $N_e$ , estimates are on the order of  $10^8$  to  $10^9$  (OCHMAN and WILSON, 1987; HARTL *et al.*, 1994). We, therefore, expect the observed sequence variants to be near mutation-selection-drift equilibrium. This expectation is supported by our simulations in which we observe a higher sequence similarity with the observed TEM consensus sequence and decreased genetic load even with much smaller  $N_e$ .



(Figure 4.3). Furthermore, previous work showed that the catalytic reaction performed by TEM of penicillin-class antibiotics is close the diffusion limit. As a result some researchers refer to TEM as a perfect enzyme (MATAGNE *et al.*, 1998; STIFFLER *et al.*, 2016). The very large effective population size, however, also raises a concern that the population mutation rate of *E. coli*  $\Theta = 4N_e\mu$  exceeds 0.1 and, thus, violated *SelAC*'s weak mutation assumption (DE KONING and DE SANCTIS, 2018). If the weak mutation assumption is violated evolution is no longer mutation limited and the time between fixation events increases. However, *phydms* operates under the same weak mutation assumption and the experimentally inferred selection clearly violates the weak mutation assumption.

As experimental selection estimates are not readily available for most organisms and proteins, a possible approach is the extrapolation of empirical estimates to homologous gene families (BLOOM, 2014, 2017). When extrapolating the selection estimates from the  $\beta$ -lactamase family TEM to the SHV family, the sequence similarity between the observed consensus sequence and the sequence of selectively favored amino acids estimated from deep mutation scanning experiments drops only slightly from 52% to 49%. This may have contributed to the notion that extrapolation to homologous gene families is possible. In contrast, estimates of site specific efficacy of selection  $G$  revealed large differences in the site specific selection on amino acids between TEM and SHV. The mismatched physicochemical weights further indicate differences in selection constraints. While the polarity of amino acids is of similar importance in TEM and SHV, amino acid composition appears to play a much greater role in SHV than in TEM. In contrast to the experimental selection estimates, extrapolated from TEM to SHV, the *SelAC* selection estimates are consistent with the observed sequences, e.g. the selectively favored amino acids estimated by *SelAC* shows a high sequence similarity with the observed TEM and SHV consensus sequence (99%). Furthermore, *SelAC* allows to compare parameters between fits to homologous proteins instead of relying on extrapolation.

While *SelAC* better explains the observed TEM sequences than the experimental estimates of site specific selection on amino acids, it is not without shortcomings itself. *SelAC* is currently too slow to be used in topology searches, therefore it is unclear if the differences in topology between *phyloms* and *SelAC* can be attributed to the same inadequacies of experimentally inferred selection. The formulation and implementation of *SelAC* can and should be improved upon as the simulation of TEM evolution from the ancestral state under the *SelAC* inferred site specific selection revealed. Starting from the ancestral sequence, the simulated sequences diverge despite stabilizing selection for the optimal amino acid, indicating that *SelAC* may underestimate selection. While *SelAC* allows for site heterogeneity in selection for amino acids, it still ignores epistasis. This however, is a shortcoming that is shared with experimental estimates as each mutant typically only carries one mutation (FIRNBERG and OSTERMEIER, 2012; JAIN and VARADARAJAN, 2014). Furthermore, not every protein is under stabilizing selection, however, *SelAC* is a model of stabilizing selection and may therefore not be adequate for every protein. TEM plays a role in chemical warfare with conspecifics and other microbes, therefore some sites may be under negative frequency dependent selection. This potential heterogeneity in selection highlights another shortcoming of *SelAC*. *SelAC* assumes the same distribution for the efficacy of selection  $G$  and physicochemical sensitivities across the whole protein. However, it is possible that residues in different secondary structures or at active sites do not share a common distribution.

As *SelAC* assumes that the fitness of an amino acid at a site declines with its distance in physicochemical space to the optimal amino acid, the choice of physicochemical properties becomes important. In this study, we used composition, polarity, and molecular volume (GRANTHAM, 1974) for all sites and only estimated their weighting. However, a wide range of additional physicochemical properties of amino acids have been described (KAWASHIMA *et al.*, 2008). A more optimal choice of physicochemical properties may be possible as well as relaxing the assumption that the same properties apply to all sites equally.

In conclusion, experimental estimates of site specific selection on amino acids have to be treated with skepticism and their adequacy should be assessed before using them to inform phylogenetic inferences. This study was initiated to assess the quality of *SelAC* with the expectation that *SelAC* could be a faster, cheaper, and more readily available alternative to experimentally inferred selection; specifically in organisms where these experiments are not feasible. Intuitively one would expect that selection coefficients of mutations estimated in living organisms would provide more information on the evolution of proteins than a model relying on many simplifying assumptions. As we show in this study, not only can *SelAC* estimate site specific selection on amino acids but our approach is a more adequate description of selection on amino acids in nature than experimental estimates.

## 4.5 Materials and Methods

### 4.5.1 Phylogenetic Inference and Model selection

TEM and SHV sequences were obtained from [BLOOM \(2017\)](#) already aligned. We separated the TEM and SHV sequences into individual alignments. Experimentally fitness values for TEM were taken from [STIFFLER \*et al.\* \(2016\)](#). We followed ([BLOOM, 2017](#)) to convert the experimental fitness values into site specific equilibrium frequencies for *phymds*. *phymds* (version 2.5.1) was fitted to the site specific selection from [STIFFLER \*et al.\* \(2016\)](#) using python (version 3.6). *SelAC* (version 1.6.1) was fitted to the TEM alignment using R (version 3.4.1) ([R CORE TEAM, 2015](#)) with and without experimental site specific selection. We assumed the physicochemical properties estimated by [GRANTHAM \(1974\)](#). We choose the constraint free general unrestricted model ([YANG, 1994](#)) as mutation model for *SelAC*. All other models were fitted using IQTree ([NGUYEN \*et al.\*, 2015](#)). We report each model's  $\log(\mathcal{L})$ , AIC, and AICc. Models were selected based on the AICc values.

### 4.5.2 Sequence Simulation

Sequences were simulated by stochastic simulations using a Gillespie algorithm (GILLESPIE, 1976) that was model independent. To calculate fixation probabilities during the simulation we followed SELLA and HIRSH (2005). The fitness values were estimated using *SelAC* or taken from STIFFLER *et al.* (2016). We choose the fitness values resulting from the highest concentration (2500  $\mu\text{g/mL}$ ) treatment of ampicillin for our comparison. We rescaled the experimental fitness such that the amino acid with the highest fitness at each site has a value of one. Mutation rates for the simulations were taken from the *SelAC* or *SelAC*+DMS fit, respectively. The initial sequences were either a random sequence sampled with uniform codon probabilities or the ancestral sequence reconstructed using FastML (ASHKENAZY *et al.*, 2012) (last accessed: 30.09.2018). Each sequence was simulated 10 times and we report average genetic load and sequence similarity and the standard error. The sequences were sampled at times 0.01, 0.1, 1, and 10 expected mutations per site.

### 4.5.3 Estimating site specific efficacy of selection $G$

*SelAC* does not by default estimate site specific values for  $G$  but assumes  $G$  values follow a  $\Gamma$ -distribution (FELSENSTEIN, 2001). Site specific values for  $G$  were optimized by fixing all estimated parameters and performing a maximum likelihood search without the integration over  $G$ . In contrast to *SelAC* that assumes  $G$  to be purely positive, we allowed negative values for  $G$  but constraint the search to values between  $-300$  and  $300$  to ensure numerical stability.

### 4.5.4 Estimating site specific fitness values $w_i$

Following BEAULIEU *et al.* (in review)  $w_i$  is proportional to

$$w_i \propto \exp(-A_0\eta\psi) \tag{4.1}$$

where  $A_0$  describes the decline in fitness with each high energy phosphate bond wasted per unit time, and  $\psi$  is the protein's production rate.  $\eta$  is the cost/benefit ratio of a protein (see (BEAULIEU *et al.*, in review) for details). However, *SelAC* only estimates a composition parameter  $\psi' = A_0\psi N_e$  thus

$$\psi = \frac{\psi'}{A_0 N_e q} \quad (4.2)$$

*SelAC* assumes that the effective population size  $N_e = 5 \times 10^6$  and that  $A_0 = 4 \times 10^{-7}$  (GILCHRIST, 2007).

#### 4.5.5 Model Adequacy

Model adequacy was assessed by comparing the observed sequences and simulations under the site specific selection inferred by the deep mutation scanning experiment or *SelAC*. First, similarity between the sequence of selectively favored amino acids and the observed TEM sequences was assessed. Sequence similarity was measured as the number of differences in the aligned amino acid sequences. Second, the genetic load of the observed and the simulated sequences was calculated using either the site specific selection inferred by the deep mutation scanning experiment or *SelAC*. The average genetic load for site  $i$  in the alignment was calculated as

$$L_i = \frac{w_{max,i} - \overline{w}_i}{w_{max,i}} \quad (4.3)$$

where  $w_{max,i}$  is the fitness of the selectively favored amino acids at position  $i$ , either estimated using the site specific selection inferred by DMS or *SelAC*.  $\overline{w}_i$  represents the average fitness of the residues observed at position  $i$ . The average sequence specific genetic load  $L$  was calculated as the sum of the site specific genetic loads  $L = \frac{1}{n} \sum_{i=1}^n L_i$  where  $n$  is the number of amino acid sites.

## 4.6 Acknowledgments

This work was supported in part by NSF Award and DEB-1355033 (BCO, MAG, and RZ) with additional support from The University of Tennessee Knoxville. CL received support as a Graduate Student Fellow at the National Institute for Mathematical and Biological Synthesis, an Institute sponsored by the National Science Foundation through NSF Award DBI-1300426, with additional support from UTK. The authors would like to thank Russel Zaretzki, Jeremy Beaulieu and Alexander Cope for their helpful criticisms and suggestions for this work.

## 4.7 Appendix: Supplementary Material

Table 4.3: Model selection of 230 models of nucleotide and codon evolution.

No.	Model	LnL	n	AIC	$\Delta$ AIC	AICc	$\Delta$ AICc
1	<i>SelAC</i> +DMS +G4	-1768	111	3758	14	3760	0
2	<i>SelAC</i> +G4	-1498	374	3744	0	3766	6
3	<i>phydms</i>	-2060.85	102	4326	582	4328	568
4	SYM+R2	-2229.616	102	4663.232	919.232	4693.862	933.862
5	TIMe+R2	-2232.406	100	4664.811	920.811	4694.172	934.172
6	TVMe+R2	-2232.838	101	4667.677	923.677	4697.668	937.668
7	TIM3e+R2	-2234.332	100	4668.664	924.664	4698.024	938.024
8	TIM2e+R2	-2234.381	100	4668.763	924.763	4698.123	938.123
9	K3P+R2	-2235.777	99	4669.553	925.553	4698.291	938.291
10	TNe+R2	-2236.078	99	4670.155	926.155	4698.892	938.892
11	SYM+R3	-2229.616	104	4667.232	923.232	4699.162	939.162
12	TIM+F+R2	-2230.958	103	4667.915	923.915	4699.191	939.191
13	TIMe+R3	-2232.404	102	4668.808	924.808	4699.437	939.437
14	GTR+F+R2	-2228.537	105	4667.073	923.073	4699.665	939.665
15	K3Pu+F+R2	-2232.617	102	4669.234	925.234	4699.864	939.864
16	TVM+F+R2	-2230.105	104	4668.21	924.21	4700.14	940.14
17	TVMe+R3	-2232.838	103	4671.676	927.676	4702.952	942.952
18	K2P+R2	-2239.424	98	4674.847	930.847	4702.969	942.969
19	TIM3e+R3	-2234.332	102	4672.664	928.664	4703.293	943.293
20	TIM2e+R3	-2234.381	102	4672.762	928.762	4703.391	943.391
21	TIM3+F+R2	-2233.064	103	4672.127	928.127	4703.403	943.403
22	TIM2+F+R2	-2233.114	103	4672.227	928.227	4703.503	943.503
23	K3P+R3	-2235.777	101	4673.553	929.553	4703.545	943.545
24	TN+F+R2	-2234.624	102	4673.249	929.249	4703.878	943.878
25	TPM3u+F+R2	-2234.673	102	4673.347	929.347	4703.977	943.977
26	TPM3+F+R2	-2234.674	102	4673.348	929.348	4703.978	943.978
27	TPM2u+F+R2	-2234.681	102	4673.363	929.363	4703.993	943.993
28	TPM2+F+R2	-2234.683	102	4673.365	929.365	4703.995	943.995
29	TNe+R3	-2236.077	101	4674.155	930.155	4704.146	944.146
30	TIM+F+R3	-2230.958	105	4671.915	927.915	4704.507	944.507
31	HKY+F+R2	-2236.266	101	4674.531	930.531	4704.522	944.522
32	GTR+F+R3	-2228.536	107	4671.073	927.073	4705.011	945.011
33	K3Pu+F+R3	-2232.617	104	4673.234	929.234	4705.163	945.163
34	TVM+F+R3	-2230.105	106	4672.21	928.21	4705.471	945.471
35	K2P+R3	-2239.192	100	4678.384	934.384	4707.745	947.745
36	TIM3+F+R3	-2233.063	105	4676.127	932.127	4708.718	948.718
37	TIM2+F+R3	-2233.113	105	4676.227	932.227	4708.818	948.818

Table 4.3 Continued

No.	Model	LnL	n	AIC	$\Delta$ AIC	AICc	$\Delta$ AICc
38	TN+F+R3	-2234.624	104	4677.249	933.249	4709.178	949.178
39	TPM3u+F+R3	-2234.673	104	4677.347	933.347	4709.277	949.277
40	TPM3+F+R3	-2234.674	104	4677.348	933.348	4709.277	949.277
41	TPM2u+F+R3	-2234.681	104	4677.363	933.363	4709.293	949.293
42	TPM2+F+R3	-2234.682	104	4677.364	933.364	4709.294	949.294
43	HKY+F+R3	-2236.074	103	4678.148	934.148	4709.424	949.424
44	SYM+I+G4	-2243.212	102	4690.424	946.424	4721.054	961.054
45	TVMe+I+G4	-2244.533	101	4691.066	947.066	4721.057	961.057
46	TIMe+I+G4	-2246.457	100	4692.914	948.914	4722.275	962.275
47	K3P+I+G4	-2248.166	99	4694.332	950.332	4723.069	963.069
48	TVM+F+I+G4	-2241.853	104	4691.707	947.707	4723.636	963.636
49	TIM3e+I+G4	-2247.379	100	4694.758	950.758	4724.119	964.119
50	K3Pu+F+I+G4	-2245.156	102	4694.311	950.311	4724.941	964.941
51	GTR+F+I+G4	-2241.484	105	4692.968	948.968	4725.559	965.559
52	TIM+F+I+G4	-2244.418	103	4694.836	950.836	4726.112	966.112
53	TPM3u+F+I+G4	-2246.03	102	4696.06	952.06	4726.69	966.69
54	TPM3+F+I+G4	-2246.069	102	4696.138	952.138	4726.768	966.768
55	TIM2e+I+G4	-2248.934	100	4697.868	953.868	4727.228	967.228
56	TNe+I+G4	-2250.587	99	4699.174	955.174	4727.911	967.911
57	TIM3+F+I+G4	-2245.534	103	4697.068	953.068	4728.344	968.344
58	K2P+I+G4	-2252.181	98	4700.362	956.362	4728.484	968.484
59	TPM2u+F+I+G4	-2247.579	102	4699.158	955.158	4729.788	969.788
60	TPM2+F+I+G4	-2247.685	102	4699.371	955.371	4730	970
61	HKY+F+I+G4	-2249.065	101	4700.13	956.13	4730.121	970.121
62	TIM2+F+I+G4	-2247.009	103	4700.018	956.018	4731.294	971.294
63	TN+F+I+G4	-2248.511	102	4701.023	957.023	4731.652	971.652
64	TVMe+I	-2254.804	100	4709.608	965.608	4738.968	978.968
65	K3P+I	-2257.72	98	4711.439	967.439	4739.561	979.561
66	SYM+I	-2254.11	101	4710.221	966.220	4740.212	980.212
67	TIMe+I	-2257.074	99	4712.149	968.149	4740.886	980.886
68	TVM+F+I	-2252.157	103	4710.315	966.315	4741.591	981.591
69	K3Pu+F+I	-2254.856	101	4711.712	967.712	4741.704	981.704
70	TIM3e+I	-2257.796	99	4713.592	969.592	4742.33	982.33
71	TPM3+F+I	-2255.771	101	4713.543	969.543	4743.534	983.534
72	TPM3u+F+I	-2255.771	101	4713.543	969.543	4743.534	983.534
73	K2P+I	-2261.218	97	4716.436	972.436	4743.949	983.949
74	GTR+F+I	-2252.067	104	4712.133	968.133	4744.063	984.063
75	TIM+F+I	-2254.783	102	4713.566	969.566	4744.195	984.195
76	TNe+I	-2260.579	98	4717.158	973.158	4745.28	985.28
77	TIM3+F+I	-2255.684	102	4715.368	971.368	4745.998	985.998
78	HKY+F+I	-2258.352	100	4716.703	972.703	4746.064	986.064
79	TIM2e+I	-2259.878	99	4717.757	973.757	4746.494	986.494



Table 4.3 Continued

No.	Model	LnL	n	AIC	$\Delta$ AIC	AICc	$\Delta$ AICc
80	TVMe+G4	-2258.853	100	4717.705	973.705	4747.066	987.066
81	SYM+G4	-2257.573	101	4717.146	973.146	4747.137	987.137
82	TPM2+F+I	-2257.712	101	4717.423	973.423	4747.415	987.415
83	TPM2u+F+I	-2257.712	101	4717.423	973.423	4747.415	987.415
84	K3P+G4	-2261.922	98	4719.844	975.844	4747.966	987.966
85	TIMe+G4	-2260.683	99	4719.365	975.365	4748.103	988.103
86	TN+F+I	-2258.28	101	4718.561	974.561	4748.552	988.552
87	TIM3e+G4	-2261.255	99	4720.51	976.51	4749.247	989.247
88	TVM+F+G4	-2256.108	103	4718.216	974.216	4749.492	989.492
89	TIM2+F+I	-2257.643	102	4719.286	975.286	4749.915	989.915
90	K3Pu+F+G4	-2258.971	101	4719.941	975.941	4749.933	989.933
91	TPM3u+F+G4	-2259.716	101	4721.433	977.433	4751.424	991.424
92	TPM3+F+G4	-2259.717	101	4721.434	977.434	4751.425	991.425
93	GTR+F+G4	-2255.75	104	4719.5	975.5	4751.43	991.43
94	TIM+F+G4	-2258.638	102	4721.276	977.276	4751.906	991.906
95	K2P+G4	-2265.454	97	4724.907	980.907	4752.421	992.421
96	TNe+G4	-2264.219	98	4724.437	980.437	4752.559	992.559
97	TIM3+F+G4	-2259.366	102	4722.732	978.732	4753.361	993.361
98	TIM2e+G4	-2263.57	99	4725.141	981.141	4753.878	993.878
99	JC+R2	-2266.233	97	4726.466	982.466	4753.98	993.98
100	F81+F+R2	-2262.327	100	4724.654	980.654	4754.015	994.015
101	HKY+F+G4	-2262.499	100	4724.999	980.999	4754.359	994.359
102	TPM2+F+G4	-2261.915	101	4725.829	981.829	4755.82	995.82
103	TPM2u+F+G4	-2261.915	101	4725.829	981.829	4755.82	995.82
104	TN+F+G4	-2262.169	101	4726.338	982.338	4756.329	996.329
105	TIM2+F+G4	-2261.585	102	4727.17	983.17	4757.8	997.8
106	F81+F+R3	-2262.028	102	4728.056	984.056	4758.685	998.685
107	JC+R3	-2265.997	99	4729.994	985.994	4758.731	998.731
108	F81+F+I+G4	-2274.845	100	4749.69	1005.69	4779.05	1019.05
109	JC+I+G4	-2279.318	97	4752.636	1008.636	4780.149	1020.149
110	F81+F+I	-2283.56	99	4765.119	1021.119	4793.857	1033.857
111	JC+I	-2287.984	96	4767.968	1023.968	4794.881	1034.881
112	F81+F+G4	-2287.834	99	4773.669	1029.669	4802.406	1042.406
113	JC+G4	-2292.095	96	4776.19	1032.19	4803.103	1043.103
114	$GY94 + F1X4 + R2$	-2242.963	102	4689.926	945.926	4821.251	1061.251
115	MGK+F1X4+R2	-2243.111	102	4690.221	946.221	4821.546	1061.546
116	$GY94 + F1X4 + R3$	-2238.022	104	4684.043	940.043	4822.271	1062.271
117	MGK+F3X4+R2	-2229.923	108	4675.846	931.846	4828.729	1068.729
118	$GY94 + F1X4 + I + G4$	-2247.179	102	4698.359	954.359	4829.684	1069.684
119	MGK+F1X4+I+G4	-2247.292	102	4698.583	954.583	4829.908	1069.908
120	MGK+F1X4+R3	-2241.989	104	4691.978	947.978	4830.206	1070.206
121	MGK+F3X4+R3	-2224.78	110	4669.559	925.559	4830.217	1070.217

Table 4.3 Continued

No.	Model	LnL	n	AIC	$\Delta$ AIC	AICc	$\Delta$ AICc
122	<i>GY94</i> +F1X4+G4	-2251.144	101	4704.287	960.287	4832.263	1072.263
123	MGK+F1X4+G4	-2251.472	101	4704.944	960.944	4832.919	1072.919
124	<i>GY94</i> +F3X4+R3	-2227.048	110	4674.096	930.096	4834.754	1074.754
125	<i>GY94</i> +F3X4+R2	-2233.068	108	4682.136	938.136	4835.019	1075.019
126	MGK+F3X4+I+G4	-2233.539	108	4683.078	939.0781	4835.962	1075.962
127	MGK+F3X4+G4	-2237.512	107	4689.024	945.024	4838.134	1078.134
128	<i>GY94</i> +F3X4+I+G4	-2238.243	108	4692.485	948.485	4845.368	1085.368
129	<i>GY94</i> +F3X4+R4	-2227.106	112	4678.213	934.213	4846.96	1086.96
130	<i>GY94</i> +F3X4+G4	-2242.394	107	4698.789	954.789	4847.899	1087.899
131	<i>GY94</i> +F1X4+I	-2260.085	101	4722.169	978.169	4850.144	1090.144
132	MGK+F1X4+I	-2260.345	101	4722.69	978.69	4850.665	1090.665
133	MGK+F3X4+I	-2246.112	107	4706.225	962.225	4855.335	1095.335
134	MG+F1X4+R2	-2268.482	101	4738.963	994.963	4866.938	1106.938
135	<i>GY94</i> +F3X4+I	-2252.532	107	4719.064	975.064	4868.174	1108.174
136	MG+F3X4+R2	-2254.453	107	4722.906	978.906	4872.015	1112.015
137	MG+F1X4+I+G4	-2272.057	101	4746.113	1002.113	4874.089	1114.089
138	MG+F1X4+R3	-2267.523	103	4741.047	997.047	4875.789	1115.789
139	MG+F1X4+G4	-2276.171	100	4752.342	1008.342	4877.033	1117.033
140	MG+F3X4+I+G4	-2257.945	107	4729.891	985.891	4879.001	1119.001
141	MG+F3X4+G4	-2261.949	106	4735.898	991.898	4881.309	1121.309
142	MG+F3X4+R3	-2253.514	109	4725.027	981.027	4881.759	1121.759
143	SYM	-2329.878	100	4859.756	1115.756	4889.116	1129.116
144	TIMe	-2333.105	98	4862.21	1118.21	4890.332	1130.332
145	TIM3e	-2333.481	98	4862.961	1118.961	4891.083	1131.083
146	TVMe	-2333.164	99	4864.328	1120.328	4893.065	1133.065
147	GTR+F	-2328.404	103	4862.809	1118.809	4894.085	1134.085
148	K3P	-2336.391	97	4866.783	1122.783	4894.297	1134.297
149	MG+F1X4+I	-2284.946	100	4769.892	1025.892	4894.583	1134.583
150	TVM+F	-2330.086	102	4864.172	1120.172	4894.802	1134.802
151	TIM+F	-2331.48	101	4864.96	1120.96	4894.952	1134.952
152	TNe	-2336.729	97	4867.458	1123.458	4894.972	1134.972
153	K3Pu+F	-2333.162	100	4866.323	1122.323	4895.684	1135.684
154	TIM3+F	-2331.971	101	4865.942	1121.942	4895.934	1135.934
155	TPM3+F	-2333.648	100	4867.297	1123.297	4896.657	1136.657
156	TPM3u+F	-2333.648	100	4867.297	1123.297	4896.657	1136.657
157	TIM2e	-2336.292	98	4868.584	1124.584	4896.706	1136.706
158	MG+F3X4+I	-2270.442	106	4752.885	1008.885	4898.295	1138.295
159	K2P	-2340.015	96	4872.03	1128.03	4898.943	1138.943
160	TN+F	-2335.102	100	4870.204	1126.204	4899.565	1139.565
161	HKY+F	-2336.783	99	4871.566	1127.566	4900.303	1140.303
162	TIM2+F	-2334.7	101	4871.401	1127.401	4901.392	1141.392
163	TPM2u+F	-2336.381	100	4872.761	1128.761	4902.122	1142.122

Table 4.3 Continued

No.	Model	LnL	n	AIC	$\Delta$ AIC	AICc	$\Delta$ AICc
164	TPM2+F	-2336.381	100	4872.762	1128.762	4902.123	1142.123
165	JC	-2366.286	95	4922.571	1178.571	4948.892	1188.892
166	F81+F	-2362.554	98	4921.108	1177.108	4949.229	1189.229
167	<i>GY94</i> +F1X4	-2315.788	100	4831.575	1087.575	4956.267	1196.267
168	KOSI07+FU+R2	-2325.725	97	4845.45	1101.45	4960.675	1200.675
169	MGK+F1X4	-2318.048	100	4836.095	1092.095	4960.787	1200.787
170	KOSI07+FU+R3	-2323.063	99	4844.126	1100.126	4965.599	1205.599
171	MGK+F3X4	-2304.357	106	4820.713	1076.713	4966.124	1206.124
172	<i>GY94</i> +F3X4	-2306.17	106	4824.339	1080.339	4969.749	1209.749
173	KOSI07+FU+I+G4	-2335.554	97	4865.108	1121.108	4980.332	1220.332
174	KOSI07+FU+G4	-2339.513	96	4871.026	1127.026	4983.218	1223.218
175	KOSI07+F3X4+R2	-2315.814	106	4843.627	1099.627	4989.038	1229.038
176	KOSI07+F3X4+R3	-2310.509	108	4837.018	1093.018	4989.901	1229.901
177	KOSI07+F1X4+R2	-2333.491	100	4866.983	1122.983	4991.674	1231.674
178	KOSI07+F1X4+R3	-2328.692	102	4861.383	1117.383	4992.708	1232.708
179	SCHN05+FU+R2	-2344.705	97	4883.411	1139.411	4998.635	1238.635
180	KOSI07+F1X4+I+G4	-2337.965	100	4875.93	1131.93	5000.621	1240.621
181	KOSI07+F1X4+G4	-2341.156	99	4880.312	1136.312	5001.784	1241.784
182	SCHN05+FU+R3	-2341.179	99	4880.358	1136.358	5001.831	1241.831
183	KOSI07+FU+I	-2349.617	96	4891.233	1147.233	5003.426	1243.426
184	KOSI07+F3X4+I+G4	-2323.767	106	4859.534	1115.534	5004.944	1244.944
185	MG+F1X4	-2342.797	99	4883.593	1139.593	5005.065	1245.065
186	KOSI07+F3X4+G4	-2327.376	105	4864.751	1120.751	5006.534	1246.534
187	MG+F3X4	-2328.539	105	4867.078	1123.078	5008.861	1248.861
188	SCHN05+F1X4+R3	-2340.927	102	4885.854	1141.854	5017.179	1257.179
189	KOSI07+F1X4+I	-2349.1	99	4896.2	1152.2	5017.672	1257.672
190	SCHN05+F3X4+R3	-2324.472	108	4864.944	1120.944	5017.827	1257.827
191	SCHN05+FU+I+G4	-2354.523	97	4903.046	1159.046	5018.27	1258.27
192	SCHN05+F1X4+R2	-2348.226	100	4896.452	1152.452	5021.143	1261.143
193	SCHN05+F3X4+R2	-2331.916	106	4875.833	1131.833	5021.243	1261.243
194	SCHN05+FU+G4	-2358.682	96	4909.365	1165.365	5021.558	1261.558
195	KOSI07+F3X4+I	-2336.826	105	4883.653	1139.653	5025.436	1265.436
196	SCHN05+F1X4+I+G4	-2351.096	100	4902.192	1158.192	5026.883	1266.883
197	SCHN05+F1X4+G4	-2353.895	99	4905.79	1161.79	5027.263	1267.263
198	SCHN05+F1X4+R4	-2340.593	104	4889.187	1145.187	5027.414	1267.414
199	SCHN05+F3X4+R4	-2324.102	110	4868.203	1124.203	5028.861	1268.861
200	SCHN05+F3X4+I+G4	-2338.345	106	4888.69	1144.69	5034.101	1274.101
201	SCHN05+F3X4+G4	-2341.811	105	4893.621	1149.621	5035.404	1275.404
202	SCHN05+FU+I	-2370.471	96	4932.943	1188.943	5045.135	1285.135
203	SCHN05+F1X4+I	-2363.696	99	4925.391	1181.391	5046.864	1286.864
204	SCHN05+F3X4+I	-2352.81	105	4915.621	1171.621	5057.404	1297.404
205	KOSI07+FU	-2394.782	95	4979.563	1235.563	5088.785	1328.785

Table 4.3 Continued

No.	Model	LnL	n	AIC	$\Delta$ AIC	AICc	$\Delta$ AICc
206	KOSI07+F1X4	-2398.44	98	4992.88	1248.88	5111.197	1351.197
207	KOSI07+F3X4	-2383.159	104	4974.318	1230.318	5112.546	1352.546
208	SCHN05+FU	-2419.333	95	5028.665	1284.665	5137.887	1377.887
209	SCHN05+F1X4	-2416.544	98	5029.088	1285.088	5147.405	1387.405
210	SCHN05+F3X4	-2402.838	104	5013.675	1269.675	5151.903	1391.903
211	<i>GY94</i> +F+R2	-2208.59	159	4735.181	991.181	5229.161	1469.161
212	<i>GY94</i> +F+G4	-2217.694	158	4751.388	1007.388	5234.504	1474.504
213	<i>GY94</i> +F+I+G4	-2213.659	159	4745.319	1001.319	5239.299	1479.299
214	<i>GY94</i> +F+R3	-2202.599	161	4727.198	983.198	5243.673	1483.673
215	<i>GY94</i> +F+I	-2228.346	158	4772.691	1028.691	5255.807	1495.807
216	<i>GY94</i> +F+R4	-2202.61	163	4731.219	987.219	5271.26	1511.26
217	<i>GY94</i> +F	-2282.254	157	4878.509	1134.509	5351.004	1591.004
218	KOSI07+F+R2	-2291.643	157	4897.286	1153.286	5369.781	1609.781
219	KOSI07+F+G4	-2301.662	156	4915.325	1171.325	5377.438	1617.438
220	KOSI07+F+I+G4	-2298.418	157	4910.835	1166.835	5383.33	1623.33
221	KOSI07+F+R3	-2286.723	159	4891.446	1147.446	5385.426	1625.426
222	KOSI07+F+I	-2311.78	156	4935.559	1191.559	5397.672	1637.672
223	SCHN05+F+R2	-2310.015	157	4934.03	1190.03	5406.525	1646.525
224	SCHN05+F+G4	-2316.684	156	4945.369	1201.369	5407.482	1647.482
225	SCHN05+F+I+G4	-2313.733	157	4941.467	1197.467	5413.962	1653.962
226	SCHN05+F+R3	-2303.732	159	4925.463	1181.463	5419.444	1659.444
227	SCHN05+F+I	-2327.127	156	4966.254	1222.254	5428.367	1668.367
228	SCHN05+F+R4	-2303.45	161	4928.9	1184.9	5445.375	1685.375
229	KOSI07+F	-2357.579	155	5025.157	1281.157	5477.12	1717.12
230	SCHN05+F	-2379.264	155	5068.528	1324.528	5520.491	1760.491

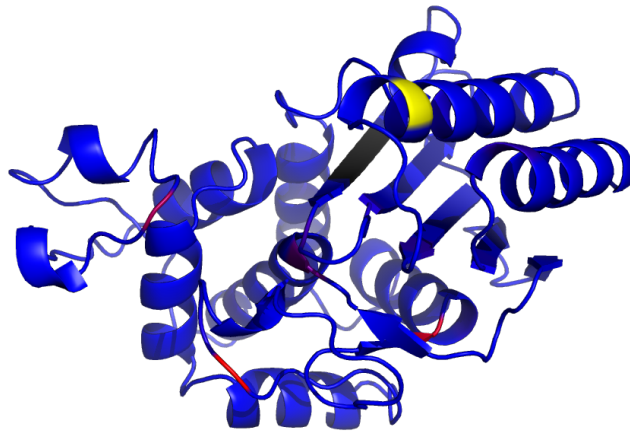


Figure 4.6: Distribution of genetic load in TEM mapped on its structure (1xpb). Average genetic load over all observed TEM variants is indicated by the color, blue low, red medium, yellow high genetic load. Active site is indicated in black.

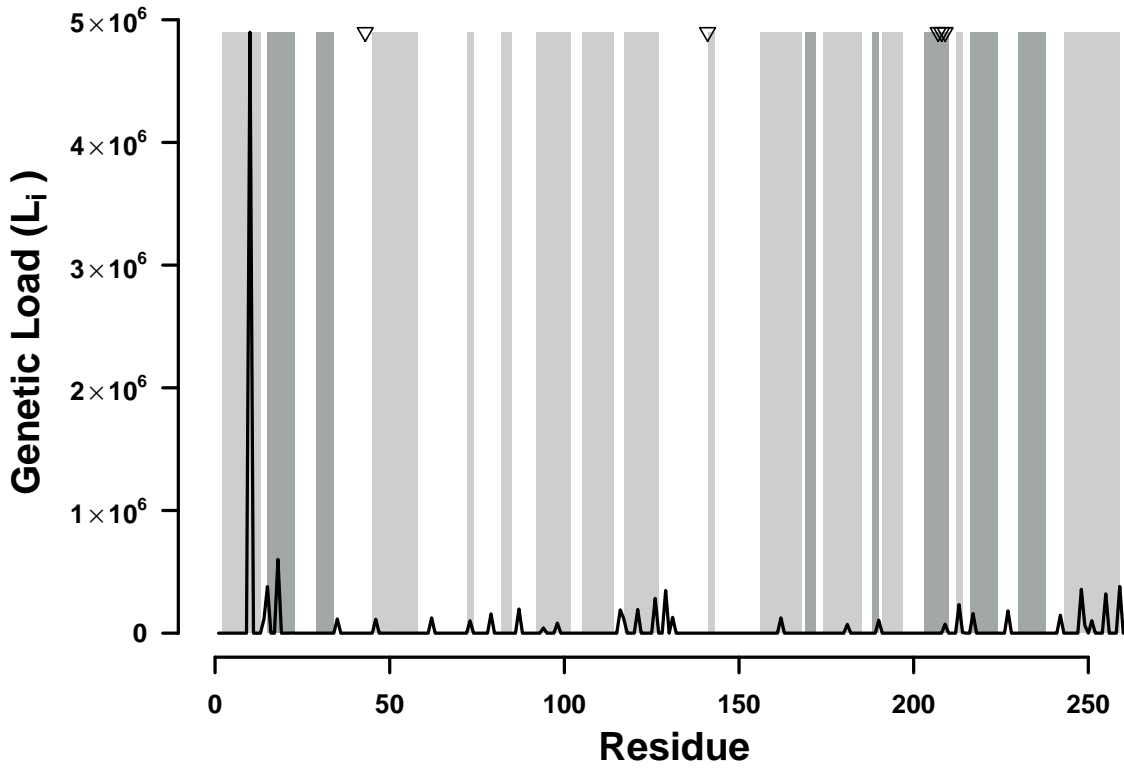


Figure 4.7: Distribution of genetic load in SHV. Average genetic load over all observed SHV variants is indicated by the black line. Light gray bars indicate where helices are found, and dark gray bars indicate  $\beta$ -sheets. The three residues forming the binding site and the two residues forming the active are indicated by triangles at the top of the plot.

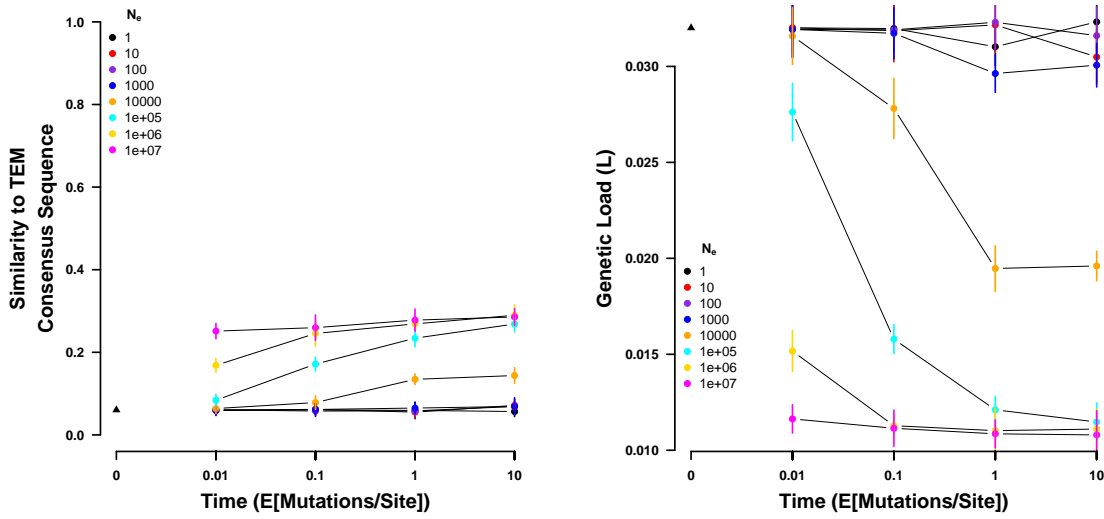


Figure 4.8: Sequences simulated from a random codon sequence under the site specific selection on amino acids estimated using *SelAC*. (left) Sequence similarity to the observed consensus sequence at various times for a range of values of  $N_e$ . (right) Genetic load of the simulated sequences at various times for a range of values of  $N_e$ . Time is given in number of expected mutations per site, which equals the substitution rate of a neutral mutation. Points indicate sample means and vertical bars indicate standard deviations. Initial sequence is the inferred ancestral state of the TEM variants and indicated by a black triangle.

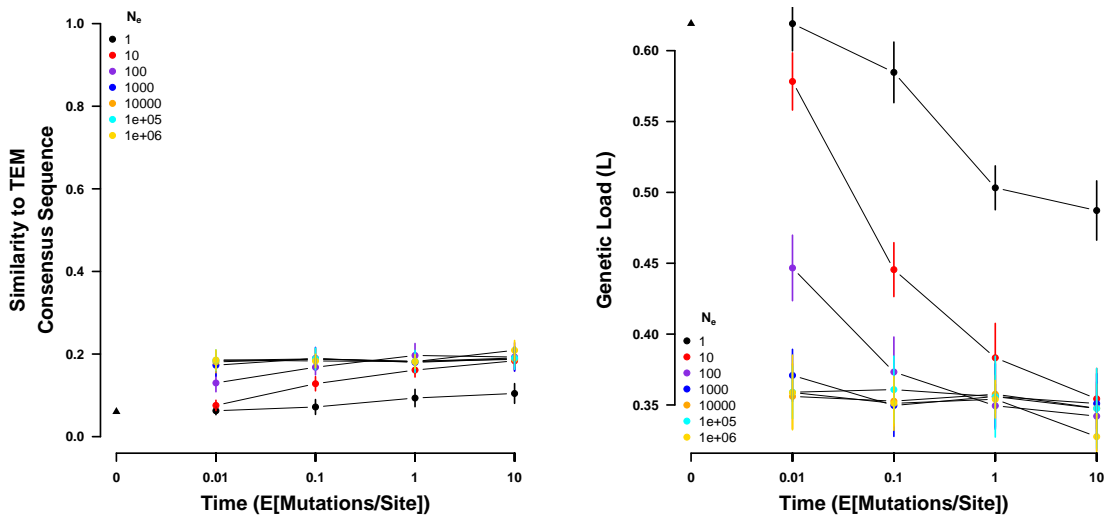


Figure 4.9: Sequences simulated from a random codon sequence under the site specific selection on amino acids estimated using deep mutation scanning. (left) Sequence similarity to the observed consensus sequence at various times for a range of values of  $N_e$ . (right) Genetic load of the simulated sequences at various times for a range of values of  $N_e$ . Time is given in number of expected mutations per site, which equals the substitution rate of a neutral mutation. Points indicate sample means and vertical bars indicate standard deviations. Initial sequence is the inferred ancestral state of the TEM variants and indicated by a black triangle.



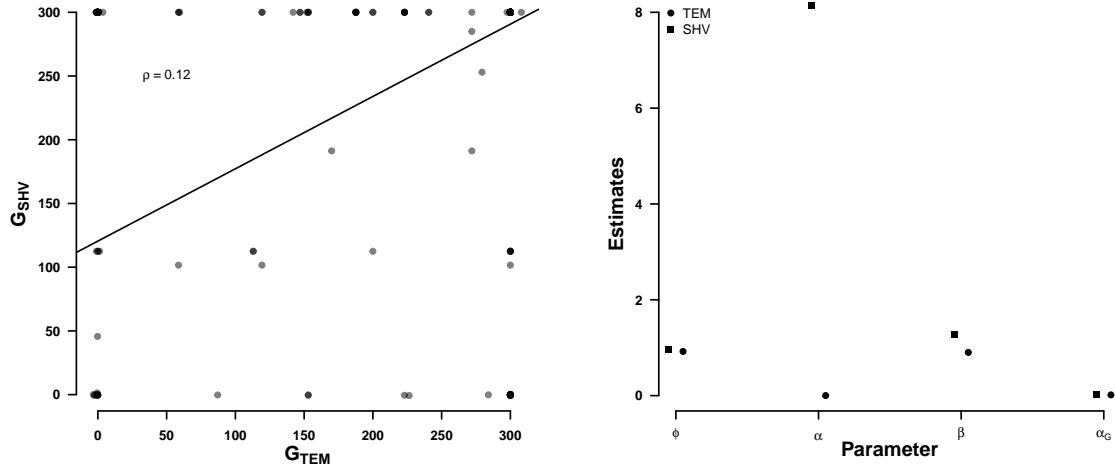


Figure 4.10: Comparison of selection related parameters between TEM and SHV. (left) Estimated site specific efficacy of selection  $G$ . (right) Selection related parameter estimates. Protein functionality production rate  $\psi$ , physicochemical weight for amino acid composition  $\alpha_c$ , physicochemical weight for amino acid polarity  $\alpha_p$ , and the parameter describing the distribution of  $G$ ,  $\alpha_G$  estimated by *SelAC*.

## Chapter 5

### Conclusion

#### 5.1 Summary

Protein synthesis from mRNA is the metabolically most expensive process a cell performs with about 20% of the cells total energy budget (REEDS *et al.*, 1985; WATERLOW and MILLWARD, 1989). The direct cost for the translation of a protein of length  $L$  requires  $4L+4$  high energy phosphate bonds provided by ATP and GTP molecules. Protein synthesis is the results of a complex interplay of many different metabolic and regulatory pathways. Each step of protein synthesis is under selection and prone to errors with consequences for downstream processes. This enormous energy expenditure for the translation of a protein from mRNA leads to strong selection for efficient translation (GILCHRIST, 2007; DRUMMOND and WILKE, 2008; GILCHRIST *et al.*, 2009; SHAH and GILCHRIST, 2011a; GILCHRIST *et al.*, 2015). However, the efficacy of selection varies with the effective population size  $N_e$  between organisms, the rate of protein synthesis, and absolute difference in metabolic expenditure with changes in amino acid and codon usage.

On the other hand, proteins are involved in almost all processes a cell performs. From communication between cells, over the processing of metabolites, to the transport of nutrients. This ratio of cost to benefit is the fundamental concept I applied to understand and separate the effects mutation, selection, and genetic drift have on protein sequence evolution. I approached cost and benefit by applying mechanistic models rooted in first

principles to protein coding sequences. In chapter 3, I focused on the cost of protein synthesis and explored the effects of mismatched codon usage. In chapter 4, I focused on the benefit of protein synthesis and estimated site specific selection on amino acids and assessed their adequacy.

### 5.1.1 The Value of Mechanistic Models

Mathematical and statistical models exist on a spectrum from descriptive over phenomenological to mechanistic with increasing power to extract information from data. Models allow us to summarize data and identify patterns. They are an essential tool to formalize verbal theory and allow for hypothesis testing. Well formulated models grounded in first principles can provide insights into underlying biological processes. Yet, we still have blackboxes in our models and many phenomena could lead to the model when approximated.

While descriptive and phenomenological models are important contributions to summarize processes, these models lack explanatory power. In contrast, mechanistic models allow researchers to extract information about the processes underlying the data. Mechanistic models, however, require an understanding about the underlying process which may not always be available. Even when this information is available, transition towards mechanistic models can be slow. For example, the most popular models used today to analyze codon usage are still phenomenological (IKEMURA, 1981; BENNETZEN and HALL, 1982; SHARP, 1987; WRIGHT, 1990; DOS REIS *et al.*, 2003, 2004). While these phenomenological models provide good heuristics to explore differences in codon usage or other phenomena, they do not directly account for the evolutionary forces shaping the observed patterns such as selection, mutation, or genetic drift. Accounting for these forces allows for the proposal and testing of more sophisticated hypothesis as I demonstrate in chapter 3 and chapter 4.

### 5.1.2 Mechanistic Models Supplement Experiments

In addition to extracting information about biological processes from data, mechanistic models can help supplement experimental procedures. Empirical estimates of site specific selection are a valuable resource to e.g. identify sites conferring antibiotic resistance (FIRNBERG *et al.*, 2014; STIFFLER *et al.*, 2016). While the unit that selection can act on is the amino acid, amino acids are a complex collection of physicochemical properties. It is, therefore, unclear for which properties amino acids are actually selected and when. Mechanistic models could be used to explore differences in the selection for physicochemical properties within and across proteins. Furthermore hypothesis could be formulated about the differing importance of physicochemical properties between e.g. sites or secondary structure elements

## 5.2 Estimating Protein Functional and Fitness Landscape

The selection on a protein sequence is highly complex. A protein of length  $L$  has  $20^L$  possible states it can occupy in a  $L$  dimensional fitness landscape. This enormous complexity makes it prohibitively expensive to study protein fitness landscapes without simplifying assumptions. It is therefore important to be aware of potential impacts such assumptions have on the obtained results and how models can be further improved. However, despite such simplifying assumptions, valuable information has been extracted from protein coding sequences.

### 5.2.1 The Importance of Translation Errors

We often think of genes evolving with natural selection favoring proteins that encode their function optimally, with mutations and genetic drift reducing protein functionality. The error rate of protein synthesis is five to six orders of magnitude higher than mutations,

causing between 10% and 20% of average length proteins to contain errors (GOLDSMITH and TAWFIK, 2009; DRUMMOND and WILKE, 2009), creating more erroneous high expression proteins such as ribosomal proteins than error free low expression proteins. Selection on a gene is, therefore, not based solely on the error free protein sequence, but on the average fitness of the population of proteins resulting from a gene by means of error prone protein synthesis. Previous work showed that proteins with functionality essential to an organism can adapt to increased error rates by increasing gene expression and showed increased selection for more stable proteins (GOLDSMITH and TAWFIK, 2009).

Organisms can take two routes to minimize the synthesis of proteins with altered functionality (DRUMMOND and WILKE, 2009). First, organisms can evolve to minimize the rate at which errors during protein synthesis occur, e.g. selecting for codons that minimize translation error rates (AKASHI, 2003; GILCHRIST and WAGNER, 2006). Second, selection could favor proteins with increased robustness to transcriptional and translational errors, e.g. increase protein stability or increase protein synthesis to compensate for non-functional proteins (GOLDSMITH and TAWFIK, 2009).

In chapter 3, I assumed that the translation process is error free, and that each produced protein functions optimal. Thus, I explicitly ignore any selection on the reduction of translation error rates. While selection for the reduction of translation error rates and selection on ribosome overhead cost do not have to be counteracting forces, they could be for some synonymous codon families. The employed ROC SEMPFR framework (GILCHRIST *et al.*, 2015) yields 100% usage of the most efficient codon if proteins synthesis rate is high enough. While individual genes may reach a 100% codon usage of the most efficient codon, we do not observe populations of high expression genes like that in nature. It is therefore unclear if selection for ribosome overhead cost can overpower counteracting selective forces if protein synthesis rate is high enough.

### 5.2.2 Homogeneous Selection

In ROC SEMPPR and *SelAC* functionality of a protein refers to the ability of a protein to perform its function and the overall need of an organism for the function. The functionality of a protein depends on many factors (DRUMMOND and WILKE, 2009). As a result, we can approximate the functionality of the protein sequence  $\vec{a}$  in a multitude of ways (GIBBS, 1873; GRANTHAM, 1974; COHEN *et al.*, 2009). However, none can capture the full complexity of a folded protein. It is easy to imagine how the strength, or direction of selection can vary between amino acid site, secondary structures and protein domains within the same protein and certainly between proteins. For example, the functionality of a well-adapted protein is unlikely to be increased by an amino acid substitution. However, the effect on the functionality and in turn fitness of a substitution may drastically differ between active sites and structural sites. Similarly, the exchange of an hydrophilic amino acid with a hydrophobic amino acid is likely to have different effects on the surface of a protein than at the core.

The *SelAC* framework (BEAULIEU *et al.*, in review) employed in chapter 4 assumes that the efficacy of selection follows a gamma-distribution. This distribution is applied to all sites. I, therefore, explicitly do not account for potential differences in the distribution of selection between e.g. secondary structure elements. Similarly, selection for physicochemical properties may differ between sites in e.g. the core or at the surface of a protein. Improvements to *SelAC* with regards to these shortcomings would allow for new hypothesis to be tested and novel information to be gained from the same data.

# Bibliography

## Bibliography

- AKASHI, H., 2003 Translational selection and yeast proteome evolution. *Genetics* **164**: 1291–1303. [101](#)
- AKASHI, H., and T. GOJOBORI, 2002 Metabolic efficiency and amino acid composition in the proteoms of *Escherichia coli* and *Bacillus subtilis*. *Proceedings of the National Academy of Sciences U.S.A* **99**: 3695–3670. [1](#)
- ALTSCHUL, S., 1991 Amino acid substitution matrices from an information theoretic perspective. *Journal of Molecular Biology* **219**: 555–565. [2](#)
- ANFINSEN, C., 1973 Principles that govern the folding of protein chains. *Science* **181**: 223–230. [2](#)
- ASHENBERG, O., L. GONG, and J. BLOOM, 2013 Mutational effects on stability are largely conserved during protein evolution. *Proceedings of the National Academy of Sciences U.S.A* **110**: 21071–21076. [66](#)
- ASHKENAZY, H., O. PENN, A. DORON-FAIGENBOIM, O. COHEN, G. CANNAROZZI, *et al.*, 2012 Fastml: a web server for probabilistic reconstruction of ancestral sequences. *Nucleic Acids Research* **40**: W580–4. [84](#)
- BEAULIEU, J., B. O’MEARA, R. ZARETZKI, C. LANDERER, J. CHAI, *et al.*, in review Population genetics based phylogenetics under stabilizing selection for an optimal amino acid sequence: A nested modeling approach. *Molecular Biology and Evolution* **X**: NA. [1](#), [6](#), [67](#), [68](#), [70](#), [79](#), [84](#), [85](#), [102](#)



- BEIMFORDE, C., K. FELDBERG, S. NYLINDER, J. RIKKINEN, H. TUOVILA, *et al.*, 2014 Estimating the phanerozoic history of the ascomycota lineages: combining fossil and molecular data. *Mol. Phylogenet. Evol.* **78**: 386–398. [38](#)
- BENNETZEN, J., and B. HALL, 1982 Codon selection in yeast. *J. Biol. Chem.* **257**: 3026–3031. [2](#), [99](#)
- BLOOM, J., 2014 An experimentally informed evolutionary model improves phylogenetic fit to divergent lactamase homologs. *Molecular Biology and Evolution* **31**: 2753–2769. [6](#), [66](#), [67](#), [68](#), [77](#), [79](#), [81](#)
- BLOOM, J., 2017 Identification of positive selection in genes is greatly improved by using experimentally informed site-specific models. *Biology Direct* **12**: 1. [6](#), [66](#), [68](#), [77](#), [79](#), [81](#), [83](#)
- BOOCH, G., 1993 *Object-oriented analysis and design with applications*. Benjamin-Cummings Publishing Co, Redwood City. [12](#)
- BRUN, T., J. PEDUZZI, M. CANICA, G. PAUL, P. NEVOT, *et al.*, 1994 Characterization and amino acid sequence of irt-4, a novel tem-type enzyme with a decreased susceptibility to beta-lactamase inhibitors. *FEMS Microbiology Letters* **120**: 111–117. [67](#), [80](#)
- BULMER, M., 1990 The selection-mutation-drift theory of synonymous codon usage. *Genetics* **129**: 897–907. [37](#)
- BUTTGEREIT, F., and M. BRAND, 1995 A hierarchy of atp-consuming processes in mammalian cells. *Biochemical Journal* **312**: 163–167. [1](#)
- CHANAL, C., M. POUPART, D. SIROT, R. LABIA, J. SIROT, *et al.*, 1992 Nucleotide sequences of caz-2, caz-6, and caz-7 beta-lactamase genes. *Antimicrob. Agents Chemother.* **36**: 1817–1820. [67](#), [80](#)

- COHEN, M., V. POTAPOV, and G. SCHREIBER, 2009 Four distances between pairs of amino acids provide a precise description of their interaction. *PLoS Comp. Biol.* **5**: e1000470. [102](#)
- COPE, A., R. HETTICH, and M. GILCHRIST, 2018 Quantifying codon usage in signal peptides: Gene expression and amino acid usage explain apparent selection for inefficient codons. *Biochimica et Biophysica Acta (BBA) - Biomembranes* **1860**: 2479–2485.
- DAVIS, M., and M. PELSOR, 2001 Experimental support for a resourcebased mechanistic model of invasibility. *Ecology Letters* **4**: 421–428. [2](#)
- DAYHOFF, M., R. SCHWARTZ, and B. ORCUTT, 1978 A model of evolutionary change in proteins. *Atlas of Protein Sequence and Structure* **5**: 345–352. [2](#)
- DE KONING, A., and B. DE SANCTIS, 2018 The rate of molecular evolution when mutation may not be weak. *bioRxiv* . [81](#)
- DORON-FAIGENBOIM, A., and T. PUPKO, 2007 A combined empirical and mechanistic codon model. *Molecular Biology and Evolution* **24**: 388–397. [2](#)
- DOS REIS, M., R. SAVVA, and L. WERNISCH, 2004 Solving the riddle of codon usage preferences: a test for translational selection. *Nucleic Acids Research* **32**: 5036–5044. [49](#), [99](#)
- DOS REIS, M., L. WERNISCH, and R. SAVVA, 2003 Unexpected correlations between gene expression and codon usage bias from microarray data for the whole escherichia coli k-12 genome. *Nucleic Acids Research* **31**: 6976–6985. [99](#)
- DRUMMOND, D., and C. WILKE, 2008 Mistranslation-induced protein misfolding as a dominant constraint on coding-sequence evolution. *Cell* **134**: 341–352. [98](#)
- DRUMMOND, D., and C. WILKE, 2009 The evolutionary consequences of erroneous protein synthesis. *Nature Reviews* **10**: 715–724. [101](#), [102](#)

- DUNN, C., F. ZAPATA, C. MUNRO, S. SIEBERT, and A. HEJNOL, 2018 Pairwise comparisons across species are problematic when analyzing functional genomic data. *Proc Natl Acad Sci USA* **115**: E409–E417. [13](#)
- ECHAVE, J., S. SPIELMAN, and C. WILKE, 2016 Causes of evolutionary rate variation among protein sites. *Nature Reviews Genetics* **17**: 109–121. [66](#)
- EDELBUETTEL, D., and R. FRANCOIS, 2011 Rcpp: Seamless r and c++ integration. *Journal of Statistical Software* **40**: 1–18. [12](#), [22](#)
- FELSENSTEIN, J., 1981 Evolutionary trees from dna sequences: a maximum likelihood approach. *Journal of Molecular Evolution* **17**: 368–376. [2](#), [66](#)
- FELSENSTEIN, J., 2001 Taking variation of evolutionary rates between sites into account in inferring phylogenies. *Journal of Molecular Evolution* **53**: 447–455. [84](#)
- FIRNBERG, E., J. LABONTE, J. GRAY, and M. OSTERMEIER, 2014 A comprehensive, high-resolution map of a gene’s fitness landscape. *Molecular Biology and Evolution* **31**: 1581–1592. [67](#), [100](#)
- FIRNBERG, E., and M. OSTERMEIER, 2012 Pfunkel: Efficient, expansive, user-defined mutagenesis. *PLOS ONE* **7**: e52031. [67](#), [80](#), [82](#)
- FITCH, W., 1976 Is there selection against wobble in codon-anticodon pairing? *Science* **194**: 1173–1174. [2](#)
- FOWLER, D., and S. FIELDS, 2014 Deep mutational scanning: a new style of protein science. *Nature Methods* **11**: 801–807. [80](#)
- FOWLER, D., J. STEPHANY, and S. FIELDS, 2014 Measuring the activity of protein variants on a large scale using deep mutational scanning. *Nature Protocols* **9**: 2267–2284. [67](#), [80](#)

- FRIEDRICH, A., C. REISER, G. FISCHER, and J. SCHACHERER, 2015 Population genomics reveals chromosome-scale heterogeneous evolution in a protoploid yeast. *Molecular Biology and Evolution* **32**: 184 – 192. [4](#), [38](#), [43](#), [44](#), [47](#), [49](#)
- GIBBS, J., 1873 A method of geometrical representation of the thermodynamic properties of substances by means of surfaces. *Transactions of the Connecticut Academy of Arts and Sciences* **2**: 382–404. [102](#)
- GILCHRIST, M., 2007 Combining models of protein translation and population genetics to predict protein production rates from codon usage patterns. *Molecular Biology and Evolution* **24**: 2362–2372. [2](#), [3](#), [37](#), [52](#), [85](#), [98](#)
- GILCHRIST, M., W. CHEN, P. SHAH, C. LANDERER, and R. ZARETZKI, 2015 Estimating gene expression and codon-specific translational efficiencies, mutation biases, and selection coefficients from genomic data alone. *Genome Biology and Evolution* **7**: 1559–1579. [1](#), [2](#), [3](#), [4](#), [10](#), [13](#), [15](#), [31](#), [37](#), [38](#), [46](#), [52](#), [98](#), [101](#)
- GILCHRIST, M., P. SHAH, and R. ZARETZKI, 2009 Measuring and detecting molecular adaptation in codon usage against nonsense errors during protein translation. *Genetics* **183**: 1493–1505. [1](#), [98](#)
- GILCHRIST, M., and A. WAGNER, 2006 A model of protein translation including codon bias, nonsense errors, and ribosome recycling. *J. of Theo. Biol.* **239**: 417–434. [101](#)
- GILLESPIE, D., 1976 A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of Computational Physics* **22**: 403–434. [84](#)
- GOJOBORI, T., 1983 Codon substitution in evolution and the "saturation" of synonymous changes. *Genetics* **105**: 1011–1027. [66](#)

- GOLDMAN, N., and Z. H. YANG, 1994 Codon-based model of nucleotide substitution for protein-coding DNA-sequences. *Molecular Biology and Evolution* **11**: 725–736. [x](#), [2](#), [6](#), [66](#), [69](#), [70](#)
- GOLDSMITH, M., and D. TAWFIK, 2009 Potential role of phenotypic mutations in the evolution of protein expression and stability. *Proceedings of the National Academy of Sciences U.S.A* **106**: 6197–6202. [101](#)
- GOUSSARD, S., W. SOUGAKOFF, C. MABILAT, A. BAUERNFEIND, and P. COURVALIN, 1991 An *isI*-like element is responsible for high-level synthesis of extended-spectrum beta-lactamase *tem-6* in enterobacteriaceae. *J. Gen. Microbiol.* **137**: 2681–2687. [67](#), [80](#)
- GOUY, M., and C. GAUTIER, 1982 Codon usage in bacteria: correlation with gene expressivity. *Nucleic Acids Research* **10**: 7055–7074. [37](#)
- GRANTHAM, R., 1974 Amino acid differences formula to help explain protein evolution. *Science* **185**: 862–864. [xi](#), [5](#), [79](#), [82](#), [83](#), [102](#)
- GRANTHAM, R., C. GAUTIER, and M. GOUY, 1980 Codon frequencies in 119 individual genes confirms consistent choices of degenerate bases according to genome type. *Nucleotide Acid Research* **8**: 1893–1912. [2](#)
- GRANTHAM, R., C. GAUTIER, M. GOUY, M. JACOBZONE, and R. MERCIER, 1981 Codon catalog usage is a genome strategy modulated for gene expressivity. *Nucleic Acids Research* **9**: 43–74. [2](#)
- HALPERN, A., and W. BRUNO, 1998 Evolutionary distances for protein-coding sequences: Modeling site-specific residue frequencies. *Molecular Biology and Evolution* **15**: 910–917. [66](#)
- HARTL, D., E. MORIYAMA, and S. SAWYER, 1994 Selection intensity for codon bias. *Genetics* **138**: 227–234. [80](#)

- HILTON, S., M. DOUD, and J. BLOOM, 2017 phydms: software for phylogenetic analyses informed by deep mutation scanning. *PeerJ* **5**: e3657. [66](#), [67](#), [68](#)
- HOLDER, M., D. ZWICKL, and C. DESSIMOZ, 2008 Evaluating the robustness of phylogenetic methods to among-site variability in substitution processes. *Philos Trans R Soc Lond B* **363**: 4013–4021. [5](#), [66](#)
- IKEMURA, T., 1981 Correlation between the abundance of *Escherichia coli* transfer rnas and the occurrence of the respective codons in its protein genes: a proposal for a synonymous codon choice that is optimal for the *E. coli* translational system. *Journal of Molecular Biology* **151**: 389–409. [2](#), [99](#)
- IKEMURA, T., 1985 Codon usage and trna content in unicellular and multicellular organisms. *Molecular Biology and Evolution* **2**: 13–34. [37](#)
- JAIN, P., and R. VARADARAJAN, 2014 A rapid, efficient, and economical inverse polymerase chain reaction-based method for generating a site saturation mutant library. *Analytical Biochemistry* **449**: 90–981. [80](#), [82](#)
- JUKES, T., and C. CANTOR, 1969 *Evolution of Protein Molecules*. Academic Press, 21–132. [2](#)
- KAWASHIMA, S., P. POKAROWSKI, M. POKAROWSKA, A. KOLINSKI, T. KATAYAMA, *et al.*, 2008 Aaindex: amino acid index database, progress report 2008. *Nucleic Acids Research* **36**: D202–D205. [82](#)
- KENSCHKE, P., M. OTI, B. DUTILH, and M. HUYNEN, 2008 Conservation of divergent transcription in fungi. *Trends Genet.* **5**: 207–211. [48](#)
- KIMURA, M., 1980 A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *Journal of Molecular Evolution* **16**: 111–120. [2](#)

- KOSIOL, C., I. HOLMES, and N. GOLDMAN, 2007 An empirical codon model for protein sequence evolution. *Molecular Biology and Evolution* **24**: 1464–1479. [xiv](#), [69](#), [71](#)
- LAFAY, B., P. SHARP, A. LLOYD, M. MCLEAN, K. DEVINE, *et al.*, 1999 Proteome composition and codon usage in spirochaetes: Species-specific and dna strand-specific mutational biases. *Nucleic Acids Research* **27**: 1642–1649. [4](#)
- LANDERER, C., A. COPE, R. ZARETZKI, and M. A. GILCHRIST, 2018 Anacoda: analyzing codon data with bayesian mixture models. *Bioinformatics* **34**: 2496–2498. [4](#), [39](#), [49](#)
- LANG, G. I., and A. W. MURRAY, 2008 Estimating the per-base-pair mutation rate in the yeast *saccharomyces cerevisiae*. *Genetics* **178**: 67 – 82. [51](#)
- LARTILLOT, N., and H. PHILIPPE, 2004 A bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Molecular Biology and Evolution* **21**: 1095–1109. [5](#), [66](#)
- LAUREAU, M., 1998 Biodiversity and ecosystem functioning: A mechanistic model. *Proceedings of the National Academy of Sciences U.S.A* **95**: 5632–5636. [2](#)
- LAWRENCE, J., and H. OCHMAN, 1997 Amelioration of bacterial genomes: Rates of change and exchange. *Journal of Molecular Biology* **44**: 383–397. [4](#), [37](#)
- LE, S., N. LARTILLOT, and G. O, 2008 Phylogenetic mixture models for proteins. *Philos Trans R Soc Lond B Biol Sci* **363**: 3965–3976. [5](#), [66](#)
- LEDER, P., and M. NIERENBERG, 1964 Rna codewords and protein synthesis, iii. on the nucleotide sequence of a cysteine and leucine rna codeword. *Proceedings of the National Academy of Sciences U.S.A* **52**: 1521–1529. [2](#)
- LEGENDRE, P., 2018 *lmodel2: Model II Regression*. R package version 1.7-3. [50](#)

- LIBERLES, D., A. TEUFEL, L. LIU, and T. STADLER, 2013 On the need for mechanistic models in computational genomics and metagenomics. *Genome Biology and Evolution* **5**: 2008–2018. [2](#)
- LINDQVIST, L., K. TANDOC, I. TOPISIROVIC, and L. FURIC, 2018 Cross-talk between protein synthesis, energy metabolism and autophagy in cancer. *Current Opinion in Genetics and Development* **48**: 104–111. [1](#)
- MABILAT, C., J. LOURENCAO-VITAL, S. GOUSSARD, and P. COURVALIN, 1992 A new example of physical linkage between tn1 and tn21: the antibiotic multiple-resistance region of plasmid pcff04 encoding extended-spectrum beta-lactamase tem-3. *Mol Gen Genet* **235**: 113–121. [67](#), [80](#)
- MARCEY-HOUBEN, M., and T. GABALDN, 2015 Beyond the whole-genome duplication: Phylogenetic evidence for an ancient interspecies hybridization in the baker’s yeast lineage. *PLoS Biology* **13**: e1002220. [38](#), [48](#)
- MATAGNE, A., J. LAMOTTE-BRASSEUR, and J. FRERE, 1998 Catalytic properties of class a beta-lactamases: efficiency and diversity. *Biochemistry Journal* **300**: 581–598. [81](#)
- MATTHAEI, J., and M. NIERENBERG, 1961 Characteristics and stabilization of dnaase-sensitive protein synthesis in *E. coli* extracts. *Proceedings of the National Academy of Sciences U.S.A* **47**: 1580–1588. [1](#)
- MAXWELL, E., 1962 Stimulation of amino acid incorporation into protein by natural and synthetic polyribonucleotides in a mammalian cell-free system. *Proceedings of the National Academy of Sciences U.S.A* **48**: 1639–1643. [2](#)
- MCGILL, B., R. ETIENNE, J. GRAY, D. ALONSO, M. ANDERSON, *et al.*, 2007 Species abundance distributions: moving beyond single prediction theories to integration within an ecological framework. *Ecology Letters* **10**: 995–1015. [2](#)



- MI, G., Y. DI, and D. SCHAFER, 2015 Goodness-of-fit tests and model diagnostics for negative binomial regression of rna sequenceing data. *PLOS ONE* **10**: e0119254. [12](#)
- MUSE, S., and B. GAUT, 1994 A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. *Molecular Biology and Evolution* **11**: 715–724. [6](#), [66](#)
- MDIGUE, C., T. ROUXEL, P. VIGIER, A. HNAUT, and A. DANCHIN, 1991 Evidence for horizontal gene transfer in escherichia coli speciation. *Journal of Molecular Miology* **222**: 851–856. [4](#), [37](#)
- NEU, H., 1969 Effect of beta-lactamase location in escherichia coli on penicillin synergy. *Appl Microbiol* **17**: 783–786. [67](#)
- NGUYEN, L., H. SCHMIDT, A. VON HAESELER, and B. MINH, 2015 Iq-tree: A fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution* **32**: 268–274. [67](#), [83](#)
- NIERENBERG, M., and J. MATTHAEI, 1961 The dependence of cell-free protein synthesis in *E. coli* upon naturally occurring or synthetic polyribonucleotides. *Proccedings of the National Academy of Sciences U.S.A* **47**: 1588–1602. [2](#)
- OCHMAN, H., and A. WILSON, 1987 *Evolutionary history of enteric bacterian*. ASM Press, 1649–1654. [80](#)
- PAYEN, C., G. FISCHER, C. MARCK, C. PROUX, D. J. SHERMAN, *et al.*, 2009 Unusual composition of a yeast chromosome arm is associated with its delayed replication. *Genome Research* **19**: 1710–1721. [38](#), [48](#), [49](#)
- R CORE TEAM, 2015 *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. [10](#), [49](#), [50](#), [83](#)

- REEDS, P., M. FULLER, and N. BA, 1985 *Metabolic basis of energy expenditure with particular reference to protein*. John Libby, 46–47. [1](#), [98](#)
- RODRIGUE, N., 2013 On the statistical interpretation of site-specific variables in phylogeny-based substitution models. *Genetics* **193**: 557–564. [66](#)
- RODRIGUE, N., and N. LARTILLOT, 2014 Site-heterogeneous mutation-selection models within the phylobayes-mpi package. *Bioinformatics* **30**: 1020–1021. [66](#)
- RODRIGUE, N., H. PHILIPPE, and N. LARTILLOT, 2010 Mutation-selection models of coding sequence evolution with site-heterogeneous amino acid fitness profiles. *Proceedings of the National Academy of Sciences U.S.A* **107**: 4629–4634. [66](#)
- ROMERO, H., A. ZAVALA, and H. MUSTO, 2000 Codon usage in chlamydia trachomatis is the result of strand-specific mutational biases and a complex pattern of selective forces. *Nucleic Acids Research* **28**: 2084–2090. [4](#)
- SAGI, D., R. RAK, H. GINGOLD, I. ADIR, G. MAAYAN, *et al.*, 2016 Tissue- and time-specific expression of otherwise identical trna genes. *PLOS Genetics* **12**: 1–27. [4](#)
- SELLA, G., and A. HIRSH, 2005 The application of statistical physics to evolutionary biology. *Proceedings of the National Academy of Sciences of the United States of America* **102**: 9541–9546. [48](#), [84](#)
- SHAH, P., and M. GILCHRIST, 2011a Explaining complex codon usage patterns with selection for translational efficiency, mutation bias, and genetic drift. *Proceedings of the National Academy of Sciences U.S.A* **108**: 10231–10236. [1](#), [2](#), [3](#), [37](#), [38](#), [98](#)
- SHAH, P., and M. GILCHRIST, 2011b Explaining complex codon usage patterns with selection for translational efficiency, mutation bias, and genetic drift. *Proc Natl Acad Sci USA* **108**: 10231–6. [10](#), [13](#)

- SHARP, P., 1987 The codon adaptatoin index - a meassure of directional synonymous codon usage bias, and its potential applications. *Nucleic Acids Research* **15**: 1281–1295. [2](#), [10](#), [49](#), [99](#)
- SHARP, P., E. COWE, D. HIGGINS, D. SHIELDS, K. WOLFE, *et al.*, 1988 Codon usage patterns in escherichia coli, bacillus subtilis, saccharomyces cerevisiae, schizosaccharomyces pombe, drosophila melanogaster and homo sapiens; a review of the considerable within species diversity. *Nucleic Acids Research* **16**: 8207–8211. [2](#)
- SODERLUND, C., M. BOMHOFF, and W. NELSON, 2011 Symap v3.4: a turnkey synteny system with application to plant genomes. *Nucleic Acids Research* **39**: e68. [50](#)
- SODERLUND, C., W. NELSON, A. SHOEMAKER, and A. PATERSON, 2006 Symap A system for discovering and viewing syntenic regions of fpc maps. *Genome Research* **16**: 1159 – 1168. [50](#)
- SOKAL, R., and F. ROHLF, 1981 *Biometry - The principles and practice of statistics in biological*. W. H. Freeman, 547–555. [xii](#), [xiii](#), [40](#), [41](#), [43](#), [50](#), [57](#), [58](#)
- SOUGAKOFF, W., S. GOUSSARD, and P. COURVALIN, 1988 The tem-3 beta-lactamase, which hydrolyzes broad-spectrum cephalosporins, is derived from the tem-2 penicillinase by two amino acid substitutions. *FEMS Microbiology Letters* **56**: 343–348. [67](#), [80](#)
- SOUGAKOFF, W., A. PETIT, S. GOUSSARD, D. SIROT, A. BURE, *et al.*, 1989 Characterization of the plasmid genes blat-4 and blat-5 which encode the broad-spectrum beta-lactamases tem-4 and tem-5 in enterobacteriaceae. *Gene* **78**: 339–348. [67](#), [80](#)
- STIFFLER, M., D. HEKSTRA, and R. R., 2016 Evolvability as a function of purifying selection in tem-1  $\beta$ -lactamase. *Cell* **160**: 882–892. [6](#), [67](#), [80](#), [81](#), [83](#), [84](#), [100](#)

- TAMURI, A., N. GOLDMAN, and M. DOS REIS, 2014 A penalized likelihood method for estimating the distribution of selection coefficients from phylogenetic data. *Genetics* **197**: 257–271. [6](#), [66](#)
- THORNE, J., N. GOLDMAN, and D. JONES, 1996 Combinng protein evolution and secondary structure. *Molecular Biology and Evolution* **13**: 666–673. [6](#), [66](#)
- THYAGARAJAN, B., and J. BLOOM, 2014 The inherent mutational tolerance and antigenic evolvability of influenza hemagglutinin. *eLife* **3**: e03300. [6](#), [66](#)
- TSAI, I., D. BENSASSON, A. BURT, and V. KOUFOPANOU, 2008 Population genomics of the wild yeast *saccharomyces paradoxus*: quantifying the life cycle. *Proc Natl Acad Sci U.S.A.* **105**: 4957–4962. [48](#)
- TSANKOV, A., D. THOMPSON, A. SOCHA, A. REGEV, and O. RANDO, 2010 The role of nucleosome positioning in the evolution of gene regulation. *PLoS Biol* **8**: e1000414. [xii](#), [40](#)
- WAGNER, A., 2005 Energy constraints on the evolution of gene expression. *Molecular Biology and Evolution* **22**: 1365–1374. [48](#), [52](#)
- WALLACE, E., E. AIROLDI, and D. DRUMMOND, 2013 Estimating selection on synonymous codon usage from noisy experimental data. *Molecular Biology and Evolution* **30**: 1438–1453. [10](#), [13](#), [37](#)
- WANG, H., K. LI, E. SUSKO, and A. ROGER, 2008 A class frequency mixture model that adjusts for site-specific amino acid frequencies and improves inference of protein phylogeny. *BMC Evolutionary Biology* **8**: 331. [5](#), [66](#)
- WARNER, J., 1999 The economics of ribosome biosynthesis in yeast. *Trends Biochem Sci* **24**: 437–440. [1](#)

- WATERLOW, J., and D. MILLWARD, 1989 *Energy cost of turnover of protein and other cellular constituents*. Georg Thieme Verlag, 277–282. [1](#), [98](#)
- WOLFRAM RESEARCH INC., 2017 *Mathematica 11*. [51](#)
- WRIGHT, F., 1990 The 'effective number of codons' used in a gene. *Gene* **87**: 23–29. [2](#), [10](#), [26](#), [99](#)
- WU, C., M. SUCHARD, and A. DRUMMOND, 2013 Bayesian selection of nucleotide substitution models and their site assignments. *Molecular Biology and Evolution* **30**: 669–688. [6](#), [66](#)
- YANG, Z., 1994 Maximum-likelihood phylogenetic estimation from DNA-sequences with variable rates over sites - approximate methods. *Journal of Molecular Evolution* **39**: 306–314. [83](#)
- ZHARKIKH, A., 1994 Estimation of evolutionary distances between nucleotide sequences. *Journal of Molecular Evolution* **39**: 315–329. [x](#), [69](#), [70](#)
- ZUCKERKANDL, E., and L. PAULING, 1962 *Molecular disease, evolution, and genic heterogeneity*. Academic Press, 189–225. [2](#)

## **Vita**

Cedric Landerer was born in Floersheim am Main, Germany on December 22, 1986 and raised in Frankfurt am Main, Germany. He graduated from Heinrich-Kleyer Highschool in Frankfurt, Germany in 2006. After that he moved to Munich, Germany to study Bioinformatics in a joint major at the University of Munich and Technical University, Munich. He received his Bachelor of Science in Bioinformatics in 2011 and his Masters of Science in Bioinformatics in 2013. He joined the Department of Ecology and Evolutionary Biology at the University of Tennessee, Knoxville in 2013 to pursue his Ph. D. He received his Ph. D. in Ecology and Evolutionary Biology in December 2018 and will start a postdoctoral position at the Max Planck Institute for Molecular Cell Biology and Genetics in Dresden, Germany in February 2019.