

ON MODELS FOR DETECTING EVIDENCE OF
MOLECULAR ADAPTATION IN HOMOLOGOUS
SEQUENCES OF PROTEIN CODING GENES

by

Christopher T Jones

Submitted in partial fulfillment of the requirements
for the degree of Doctorate of Philosophy

at

Dalhousie University
Halifax, Nova Scotia
August 2019

© Copyright by Christopher T Jones, 2019

To new beginnings.

Table of Contents

List of Tables	vi
List of Figures	viii
Abstract	x
List of Abbreviations and Symbols Used	xii
Acknowledgements	xv
Chapter 1 Introduction	1
1.1 Thinking About Molecular Evolution	1
1.1.1 Biochemical Information Technology	1
1.1.2 Drift and Selection	3
1.1.3 The Diffusion Approximation	4
1.1.4 Site-Specific Substitution Rate Matrix	6
1.1.5 Evolution as a Markov Process	8
1.1.6 Stationary Frequencies and Time Reversibility	9
1.1.7 Substitution rates and the Canonical Signature of Positive Selection	11
1.2 Codon Substitution Models	14
1.2.1 The Data	14
1.2.2 A Simple CSM	17
1.2.3 The Pruning Algorithm	19
1.2.4 Inferring Positive Selection	21
1.2.5 <i>Post Hoc</i> Analysis	23
1.2.6 Heterotachy and the Covarion-like Model	23
1.2.7 Simplifying Assumptions	24
1.3 Thesis Outline	26
Chapter 2 Non-adaptive Shifting Balance on a Static Mutation-selection Landscape: A Novel Scenario of Positive Selection.	30
2.1 Introduction	30
2.1.1 Chapter Outline	32
2.2 Results	33
2.2.1 Defining dN/dS under the Mutation-Selection Framework	33
2.2.2 M_0 is Equivalent to a Model for Frequency-dependent Selection under the MS Framework	36
2.2.3 Shifting Balance on a static MS Landscape	37

2.2.4	Split MS Landscapes	40
2.2.5	A Mechanistic Model for Non-adaptive Shifting Balance	43
2.2.6	Detecting Transient Changes in ω Caused by Non-adaptive Shifting Balance	49
2.2.7	Detecting $\omega > 1$ Caused by Non-adaptive Shifting Balance	51
2.2.8	Changing Fitness Landscapes	53
2.3	Discussion	55
2.4	Methods	59
2.4.1	Alignment Generation	59
2.4.2	The M3($k = 2$) vs CLM3($k = 2$) Contrast	60
2.4.3	BUSTED	60
2.5	Appendix	61
2.5.1	A Demonstration of $p_+^h = p_-^h$	61
2.5.2	$dN^h/dS^h \leq 1$ When \mathbf{f}^h is Fixed	61
2.5.3	A Demonstration that A^h is Time-Reversible	62
2.5.4	Visualizing Substitution Dynamics	64
Chapter 3	Phenomenological Load on Model Parameters can lead to False Biological Conclusions.	68
3.1	Introduction	68
3.2	New Approaches	72
3.2.1	Modeling a Mixture of Static and Switching Sites: RaMoSS	72
3.2.2	Quantifying Phenomenological Load	73
3.2.3	Assessing the Realism of Alignments Simulated under MS	75
3.3	Results	77
3.3.1	Putative DT Mutations are Detectable in a Real Alignment	77
3.3.2	The Extent to which DT Parameters Carry PL is Related to Model Misspecification	81
3.3.3	Simulation Study 1	82
3.3.4	Simulation Study 2	83
3.3.5	Simulation Study 3	85
3.3.6	Alignments Generated under MSmmtDNA are Realistic by Several Measures of Comparison	88
3.3.7	Evidence of Confounding	90
3.3.8	Assessing PL in a Model for Detecting Relaxation of Selection Pressure.	92
3.3.9	Assessing PL in a Model for Detecting variations in dS .	92
3.4	Discussion	94
3.5	Methods	99
3.5.1	RaMoSS	99
3.5.2	Model Contrasts	99
3.5.3	Generating Alignments using MSmmtDNA	101

3.5.4	Constructing PDFs for Scaled Selection Coefficients . . .	104
3.6	Appendix	106
3.6.1	Tables of Median MLEs for Simulation Studies	106
3.6.2	Observed versus Simulated	107
Chapter 4	A Phenotype-Genotype Codon Substitution Model for Detecting Adaptive Evolution.	109
4.1	Introduction	109
4.2	Materials and Methods	112
4.2.1	Background	112
4.2.2	The PG-BSM	116
4.2.3	Rigorous Model Assessment Requires a Realistic Data Generating Process	124
4.3	Results	126
4.4	Simulations	126
4.4.1	Simulation 1: Generating under the Null PG-BSM . . .	127
4.4.2	Simulation 2: Generating under MSmmtDNA and MST- GdR	129
4.4.3	Simulation 3: A Scenario with Four Phenotypic States	133
4.5	Analysis with Real Data	136
4.5.1	mmtDNA	136
4.5.2	Phytochrome A&CF	143
4.5.3	Invertebrate Cytochrome B	147
4.6	Discussion	151
4.7	Methods	157
4.7.1	Data Generation using MSmmtDNA and MSTGdR . .	157
4.7.2	Generating Ancestral Phenotypes	160
4.7.3	Computing Scaling Constants for the PG-BSM	161
4.7.4	False Discovery Control	162
4.7.5	Dealing with Underflow	163
Chapter 5	Discussion	165
5.1	Historical Development of CSMs.	165
5.1.1	Phase I: Simple Models, Simple Problems	165
5.1.2	Phase II: The Rise in Complexity	170
5.1.3	An Argument for Phenomenological CSMs	174
5.2	Other Examples of Confounding and PL	178
5.2.1	Confounding in the Automated Detection of SST Fronts	180
5.2.2	PL and the Efficiency of Photosynthesis	182
5.3	Final Thoughts: What is the Canonical Signature of Molecular Adaptation?	184
Bibliography	187

List of Tables

Table 1.1	Amino acids and their codon aliases for the standard genetic code.	3
Table 2.1	The codon-specific rate ratios.	43
Table 2.2	Substitution Probabilities.	48
Table 2.3	Detecting Heterotachy.	50
Table 2.4	Detecting Positive Selection.	52
Table 3.1	Parameter estimated for the real data.	78
Table 3.2	Model contrasts for the real data.	79
Table 3.3	Site pattern analysis.	80
Table 3.4	Simulation 1 Results.	82
Table 3.5	Simulation 2 Results.	83
Table 3.6	Simulation 3 Results.	85
Table 3.7	Simulated versus real selection coefficient distributions.	89
Table 3.8	List of Critical Values.	100
Table 3.9	Simulation 1 Medians.	106
Table 3.10	Simulation 2 Medians.	106
Table 3.11	Simulation 3 Medians.	107
Table 4.1	Empirical vs Expected.	129
Table 4.2	Mean/median,(standard deviation) of select MLEs for Simulation 1.	129
Table 4.3	Simulation 2 Results.	132
Table 4.4	Simulation 2 <i>post hoc</i> analysis.	133
Table 4.5	Selected results for Simulation 3.	135
Table 4.6	Simulation 3 <i>post hoc</i> analysis.	136

Table 4.7	Results of the analysis of the mmtDNA.	140
Table 4.8	Results of the fit of alternate YN-BSM A and B to the mmtDNA data.	141
Table 4.9	Results of the analysis of the phyA&CF data.	145
Table 4.10	Results of the fit of alternate YN-BSM A and B to the phyA&CF data.	146
Table 4.11	Simulation 4 Results.	146
Table 4.12	Amino acid compositions, YN-BSM A.	150
Table 4.13	Amino acid compositions, PG-BSM.	150
Table 5.1	Simulated selection regimes.	168

List of Figures

Figure 1.1	An example of the object of analysis.	14
Figure 1.2	An arbitrary branching element in a binary tree.	19
Figure 2.1	A site-specific fitness landscape.	40
Figure 2.2	A MS and McCandlish landscape.	40
Figure 2.3	A landscape under relaxed selection pressure.	42
Figure 2.4	Distributions of site-specific rate ratios.	44
Figure 2.5	Distributions for the parameters of the mechanistic non-adaptive shifting balance model.	45
Figure 2.6	A MS landscape dominated by Methionine.	47
Figure 2.7	Alternate distributions for the parameters of the mechanistic non-adaptive shifting balance model.	49
Figure 2.8	Distributions of MLEs under CLM3 and BUSTED.	53
Figure 2.9	Investigation of the MSES model.	54
Figure 3.1	The phylogeny for the concatenation of 12 H-strand mitochondrial DNA sequences (3331 codon sites).	77
Figure 3.2	PL versus PRD.	88
Figure 3.3	Investigating confounding.	90
Figure 3.4	PL in other CSMs.	94
Figure 3.5	Distributions of scaled selection coefficients.	104
Figure 3.6	A comparison of the observed versus simulated amino acid frequencies.	107
Figure 3.7	A comparison of the observed versus simulated relative pairwise amino acid frequencies.	108
Figure 3.8	A comparison of the observed versus simulated distribution of the number of amino acids.	108

Figure 4.1	An illustration of the difference between the cladewise (CW and rCW) and branchwise (BW) evolutionary processes.	118
Figure 4.2	The clocked and unclocked trees used in Simulations 1, 2, and 3.	127
Figure 4.3	Branch lengths estimated by fitting the null and alternate PG-BSM to the mammalian mtDNA.	136
Figure 4.4	Patterns for sites in the mmtDNA alignment.	141
Figure 4.5	Patterns for sites in the mmtDNA alignment inferred to be in category 2a or 2b by the YN-BSM B.	142
Figure 4.6	Branch lengths estimated by fitting the alternate PG-BSM to the phytochrome A&CF alignment.	143
Figure 4.7	Patterns for sites in the phytochrome alignment inferred to be in category 2a or 2b by the YN-BSM A.	146
Figure 4.8	Branch lengths estimated by fitting the null PG-BSM to the cytochrome B alignment.	147
Figure 4.9	An arbitrary branching element in a binary tree.	160
Figure 5.1	A cartoon of the $(61^N - 1)$ -dimensional simplex containing all possible site-pattern distributions for an N-taxon alignment.	174
Figure 5.2	An Example of Confounding.	179
Figure 5.3	Satellite images.	181
Figure 5.4	Two typical $P(E)$ curves.	183

Abstract

Codon substitution models (CSMs) are commonly fitted to alignments of homologous protein-coding sequences with the objective of determining whether sites in the gene underwent positive selection. Under the standard paradigm such evidence is often assumed to be enough to conclude the gene evolved adaptively. CSMs are commonly validated using simulated alignments. A central theme of this dissertation is use of relatively realistic alignment-generating processes grounded in mutation-selection (MS) theory (Chapter 1). The MS framework permits sites to be evolved each on their own site-specific fitness landscape defined by a vector of fitness coefficients for the twenty amino acids. A novel MS alignment-generating process was used to show that evidence for variation in site-specific rate ratios (a.k.a. heterotachy) with episodic positive selection can be produced by episodic adaptive changes in site-specific fitness coefficients, consistent with the standard paradigm, but also by a second previously unrecognized process that I call non-adaptive shifting balance. This finding undermines sophisticated CSMs specifically designed to infer episodic adaptation by detecting heterotachy with episodic positive selection (Chapter 2). Processes that tend to generate similar patterns in data are said to be confounded. Confounding can lead to a novel statistical pathology that I call phenomenological load. A series of novel CSMs fitted to alignments generated under a version of MS uniquely formulated to mimic real data were used to demonstrate that phenomenological load can lead to false biological conclusions. These analyses were accompanied by a novel method to assess the potential impact of phenomenological load on any given model parameter (Chapter 3). Confounding of adaptive and non-adaptive processes that generate heterotachy can be avoided by abandoning positive selection as an indicator of adaptation and instead using evidence of changes in site-specific amino acid fitnesses. This approach was realized by constructing the phenotype-genotype branch-site model (PG-BSM), a descendant of traditional branch-site models

that combines alignment data with a discrete phenotype (i.e., contextual information) under a unified statistical framework. The PG-BSM was validated using extensive simulations and produced plausible results when applied to real data (Chapter 4). This dissertation ends with a discussion of implications of my findings (Chapter 5).

List of Abbreviations and Symbols Used

AIC	Akaike Information Criterion
ASM	Amino Acid Substitution Model
BG	Background Branches under the YN-BSM
BSREL	Branch-Site Random Effects Likelihood
BUSTED	Branch-Site Unrestricted Statistical Test for Episodic Diversification
BW	Branchwise
CL	Covarion-like
CLM3	Model M3 with Covarion-like Switching
CSM	Codon Substitution Model
CW	Cladewise
DT	Doublet and Triple Nucleotide Substitutions
FD	Functional Divergence
FDC	False Discovery Counts
FG	Foreground Branches under the YN-BSM
LL	Log-Likelihood
LLR	Log-Likelihood Ratio
ML	Maximum Likelihood
MLE	Maximum Likelihood Estimate
MS	Mutation-Selection

PCM	Phylogenetic Comparative Methods
PG	Phenotype-Genotype
PG-BSM	Phenotype-Genotype Branch-Site Model
rCW	Reverse Cladewise
RaMoSS	Random Mixture of Static and Switching Sites
SNS	Single Nucleotide Substitution
YN-BSM	Yang-Nielsen Branch-Site Model
A	A MS Substitution Rate Matrix
α	Double Nucleotide Substitution Rate
β	Triple Nucleotide Substitution Rate
δ	Covarion-like Switching Rate
dN	Nonsynonymous Substitution Rate
dS	Synonymous Substitution Rate
\mathbf{F}	Vector of Discrete Phenotypes
f_x	A Fitness Coefficient for Genotype x
h	Site Index
κ	Transition Bias
ℓ_x	Scalar Indicator for Condition x
$\mathbf{\ell}_x$	Matrix Indicator for Condition x
λ	Rate Parameter for a Model of Phenotype Evolution
π_x	Stationary Frequency of Codon x

π_{jk}^*	Empirical Frequency of Nucleotide j at Codon Position k
$\hat{\pi}_{\mathbf{z}}$	Proportion of Change Maps \mathbf{z} in the Sample
M	Mutation Rate Matrix
N_e	Population Size
$P^h(t)$	Substitution Probability Matrix
p_x	Probability of Condition x
p_+	Expected Proportion of Beneficial Substitutions
p_-	Expected Proportion of Detrimental Substitutions
$\boldsymbol{\pi}$	Stationary Codon Frequencies
Q	A Rate Matrix for a Standard CSM
s_{ij}	A Selection Coefficient
\mathbf{t}	Vector of Branch Lengths
τ	The Topology of a Phylogenetic Tree
$\boldsymbol{\theta}$	An arbitrary Vector of Model Parameters
$\hat{\boldsymbol{\theta}}$	The MLE of $\boldsymbol{\theta}$
ω	Nonsynonymous-to-Synonymous Rate Ratio
\mathbf{x}	A Site Pattern
X	An Alignment
x^h	Parameter x for the h^{th} Codon Site
\mathbf{z}	Change Map for a Discrete Phenotype

Acknowledgements

My sincerest thanks to Ed and Joe for their excellent supervision.

Chapter 1

Introduction

This first chapter covers background information in two sections. The first, entitled “Thinking About Molecular Evolution”, introduces the mechanistic mutation-selection (MS) framework (Halpern and Bruno, 1998) as a means to conceptualize evolutionary processes at individual codon sites in a protein-coding gene. The MS framework provided much insight and in particular led to the discovery of the significance of the site-specific dynamic I call “non-adaptive shifting balance” whereby a population first drifts to a suboptimal codon at a site and subsequently returns to the optimal codon for that site by a combination of positive selection and drift. The second section, “Codon Substitution Models”, introduces standard methods used to extract information from genetic data. Codon substitution models (CSMs) provide phenomenological approximations of mechanistic MS processes, and as such are founded on a number of simplifying assumptions. By presenting MS first, I aim to inculcate in the reader a mechanistic view of evolutionary processes before presenting CSMs with all of their necessary simplifications.

1.1 Thinking About Molecular Evolution

1.1.1 Biochemical Information Technology

We find ourselves in a universe in which matter spontaneously congregates to form organized low entropy structures that seem to defy explanation. The causal processes that give rise to organization apparently stem from the laws of thermodynamics, and in particular the laws that govern systems in disequilibrium (Schneider and Kay, 1994). For example, gradients of various kinds in the Earth’s atmosphere (temperature, pressure, humidity, etc.) prompt the formation of structures such as convection cells and hurricanes that persist for a time in a state of disequilibrium. These structures, we are told, serve

to reduce gradients and move the atmosphere toward equilibrium and a state of maximum entropy (Schneider and Kay, 1994). It is only the continuous influx of solar energy that prevents the global atmosphere from ever reaching this quiescent state. Some would characterize living systems in the same way, as functioning to reduce gradients (England, 2013). Photoautotrophic organisms, likely among the first life forms on Earth, make use of compounds such as H_2 , H_2S and H_2O to fix carbon (see Lenton and Watson, 2011, for a history of life on Earth) via processes that ultimately convert high energy photons into heat. Hence, one might say that such living systems serve to dissipate the solar influx and increase entropy, just as a hurricane does. However, there is a fundamental difference between purely physical structures and living systems, namely biochemical information technology, which on Earth consists of a double-strand of deoxyribose nucleic acid monomers called DNA (i.e., long sequences of the four nucleotides thymine, cytosine, guanine and adenine denoted T, C, G and A) in combination with an array of protein machines that maintain, transcribe, translate and replicate the stored information. This thesis is centered on models of the processes of change that can impact a protein-coding gene over macroevolutionary time scales.

1.1.2 Drift and Selection

Amino Acid	Codon Aliases
alanine	GCT, GCC, GCA, GCG
arginine	CGT, CGC, CGA, CGG, AGA, AGG
asparagine	AAT, AAC
aspartic acid	GAT, GAC
cysteine	TGT, TGC
glutamine	CAA, CAG
glutamic acid	GAA, GAG
glycine	GGT, GGC, GGA, GGG
histidine	CAT, CAC
isoleucine	ATT, ATC, ATA
methionine	ATG (start)
leucine	TTA, TTG, CTT, CTC, CTA, CTG
lysine	AAA, AAG
phenylalanine	TTT, TTC
proline	CCT, CCC, CCA, CCG
serine	TCT, TCC, TCA, TCG, AGT, AGC
threonine	ACT, ACC, ACA, ACG
tryptophan	TGG
tyrosine	TAT, TAC
valine	GTT, GTC, GTA, GTG
stop	TAA, TAG, TGA

Table 1.1: Amino acids and their codon aliases for the standard genetic code. The stop codons TAA, TAG and TGA indicate the end of a segment of protein-coding DNA. The codon ATG serves as both an indicator of the start of a segment of protein-coding DNA and also as the single codon alias for methionine.

The physical structure of DNA is a double helix (Watson and Crick, 1953), which in eukaryotes is mostly stored in the nucleus. When appropriately signaled, protein machines within the nucleus unwind a portion of the double helix to transcribe a segment of DNA (i.e., a gene¹) into a strand of ribonucleic acid (RNA). RNA is organized in nucleotide triplets called codons (Table 1.1), and serves as a template with which to build a protein molecule. Each of the $4^3 = 64$ possible codons indicates one of 20 possible amino acids. In eukaryotes, the codon template is translated into its corresponding sequence of amino acids by structures called ribosomes that reside outside the nucleus.

¹More broadly, a gene can be defined as a discrete locus of heritable DNA that can impact organismal phenotype when it is expressed either as a protein product or as a regulator of the expression of other genes.

The resulting protein molecule then folds either spontaneously or with the help of other proteins to take on its functional three-dimensional configuration.

This sequence of events, DNA transcribed to RNA translated to protein (i.e., the “central dogma” of molecular biology, Crick, 1970), is one of two fundamental processes. The other process occurs when strands of DNA are replicated². During replication it can happen that one nucleotide is mistakenly replaced by another, giving rise to a mutation. Although mutations are quite rare (e.g., the mutation rate was estimated to be 3×10^{-8} mutations/nucleotide/generation in the Human Y-chromosome, Xue *et al.*, 2009), they nevertheless provide variations essential for evolution by natural selection (Darwin, 1859). When a mutated gene is introduced into a population via a single individual it will either eventually vanish all together or spread to the entire population (i.e., be “fixed”). Its eventual fate depends in part on stochastic population processes (drift) and in part on its relative fitness compared to the wild type gene common to the rest of the population (selection).

1.1.3 The Diffusion Approximation

Imagine a diploid population of size N_e in which every individual has the same version A of a particular gene³. Over the generations there might on rare occasion occur a mutation in A to give a new variant B that exists in a single individual. There are many ways a gene can mutate during replication, from a change at a single nucleotide site up to large scale insertions and deletions of whole segments of DNA (Patthy, 2008). For the purpose of modeling codon evolution, it is standard practice to consider mutations that alter a single

²Additional processes include variations in gene expression caused by (i) changes in the conformation of a chromosome (Krogh *et al.*, 2018) and (ii) changes in the pattern of epigenetic biochemical markers (Jablonka and Lamb, 2006).

³The symbol N_e typically denotes the effective population size, the number of individuals an idealized population would have to exhibit the same dynamics as observed in a real population. In the present theoretical context N_e is equated to the number of individuals in the population from which genes that enter the next generation are randomly drawn. The symbol N is used throughout this thesis to indicate the number of sequences in an alignment.

codon only. Suppose the two variants A and B are as follows:

A : ATC ATA GTA CTC CAA GCC CTA **TTA** GCC ACC TAT

B : ATC ATA GTA CTC CAA GCC CTA **ATA** GCC ACC TAT

The two genes are identical everywhere except at one codon site, where a single nucleotide mutation changed the codon TTA, corresponding to the amino acid leucine, to the codon ATA, corresponding to the amino acid isoleucine. Each individual in a diploid population has two copies of each gene. Hence, the population initially consists of $N_e - 1$ individuals with two copies of variant (or allele) A and one individual with one copy of A and one of B. To approximate the probability that B will eventually be fixed, suppose pairs of alleles are passed to the next generation by random sampling, and that the selection coefficient is $s_{AB} = f_B - f_A$, where f_A is the fitness of A and f_B the fitness⁴ of B. Then the probability that B will be fixed is (Fisher, 1922, 1930; Haldane, 1927, 1932; Wright, 1931; Kimura, 1962):

$$P(\text{B is fixed}) \approx \frac{1 - \exp(-2s_{AB})}{1 - \exp(-4N_e s_{AB})} \quad (1.1)$$

A plot of s_{AB} versus $P(\text{B is fixed})$ is sigmoid in shape, and converges asymptotically to one as $s_{AB} \rightarrow +\infty$ and asymptotically to zero as $s_{AB} \rightarrow -\infty$. If A and B are equally fit (i.e., if $s_{AB} = 0$) then allele B will be fixed just by chance over the generations with probability $1/(2N_e)$. In this case the selection regime is said to be neutral and fixation is said to occur by genetic drift. The effect of selection is to cause the probability of fixation to deviate from the probability of fixation by drift alone. Hence, beneficial mutations (i.e., when $s_{AB} > 0$) will be fixed with probability $P(\text{B is fixed}) > 1/(2N_e)$, and deleterious mutations with probability $P(\text{B is fixed}) < 1/(2N_e)$. Equation (1.1) can be approximated by:

$$P(\text{B is fixed}) \approx \frac{2s_{AB}}{1 - \exp(-4N_e s_{AB})} \quad (1.2)$$

Equation (1.2) will be referred to as the diffusion approximation. Note that equation (1.1) and equation (1.2) both require that $|s_{AB}| \ll 1$ and also that the

⁴Following common practice, I use phrases like “the fitness of B” or “the fitness of a protein” as shorthand meaning the fitness of the organism in which the gene or protein resides subject to the usual *ceteris paribus* assumption.

mutation rate M per gene per generation is small enough that B will either be fixed or eliminated long before the next mutation occurs. A population is segregating for the longest time by a neutral mutation that is destined to be fixed, when the expected waiting time before fixation is on the order of $4N_e$ generations (Kimura and Ohta, 1969; McCandlish and Stoltzfus, 2014). Hence, the maximum expected number of mutations that arise before B is fixed or eliminated is on the order of (waiting time) \times (total mutation rate) = $4N_e \times 2N_e M$. The diffusion approximation therefore requires mutations to occur infrequently enough that $8N_e^2 M \ll 1$ or $M \ll 1/(8N_e^2)$. It commonly occurs that several variants of a given gene exist in a population at the same time. This is known as a polymorphism. There are four common blood groups present in human populations (designated A, B, AB and O), for example, all of which correspond to different versions of a single gene encoding the enzyme glycosyltransferase located on chromosome 9. This polymorphism is apparently stable (Ségurel *et al.*, 2012), meaning that it is unlikely that one of the four variants will eventually be fixed. Under equations (1.1) and (1.2) it is assumed that stable polymorphisms do not occur. An idealized Wright-Fisher population is also assumed, meaning that the population size N_e is fixed and that the pairs of alleles that are passed on to the next generation are selected randomly (Fisher, 1922; Wright, 1931).

1.1.4 Site-Specific Substitution Rate Matrix

To model⁵ the effect of any given mutation at a codon site in a gene, let $\mathbf{f}^h = \langle f_1^h, \dots, f_{64}^h \rangle$ be a 1×64 row vector that gives the fitness of the gene as a function of the codon occupying the h^{th} site holding everything else (e.g., the codons at other sites in the gene, other genes in the same genome, the environment, etc.) constant. The probability that a mutation from wild type A (for which codon i occupies the h^{th} site) to mutant B (for which codon

⁵The model presented in this section is based on the mutation-selection (MS) framework (Halpern and Bruno, 1998). Note that MS is seldom used analytically (i.e., it is not fitted to data), but is presented here because it provides a means to think about the way an individual codon site evolves in terms of the underlying population dynamics. MS will also play a key role in the generation of genetic data for the purpose of model testing in Chapters 2 to 4.

j occupies that site) occurs in an individual and is subsequently fixed in a diploid population of size N_e sometime during the time interval Δt can be approximated up to $o(\Delta t)$ as follows for all $j \neq i$:

$$P_{ij}^h(\Delta t) \approx \begin{cases} M_{ij}2N_e\Delta t \times 1/(2N_e) & \text{if } s_{ij}^h = 0 \\ M_{ij}2N_e\Delta t \times \frac{2s_{ij}^h/N_e}{1-\exp(-4s_{ij}^h)} & \text{otherwise} \end{cases} \quad (1.3)$$

The term M_{ij} represents the $i \rightarrow j$ mutation rate per individual per unit time. Since all versions of A are equally likely to mutate and each individual in the population contains two copies of A, the probability that the $i \rightarrow j$ mutation will arise sometime during the time interval Δt is $M_{ij}2N_e\Delta t$. Note that the selection coefficient $s_{ij}^h = N_e(f_j^h - f_i^h)$ in equation (1.3) is scaled by the population size.

The diffusion approximation gives the probability that a mutation will eventually be fixed, but does not give the time of fixation. However, the temporal scale of the population dynamic leading to fixation or elimination is on the order of thousands of generations (e.g., typically something less than $4N_e$ in a diploid population). This is virtually instantaneous compared to the macroevolutionary timescales represented by Δt (e.g., millions of years) meaning that the fate of a mutation that arose at $t = 0$ will be decided long before the end of the time interval $(0, \Delta t)$. Hence, the substitution probability per unit time $P_{ij}^h(\Delta t)/\Delta t$ approximates an instantaneous substitution rate. Such rates comprise the elements of a 64×64 site-specific substitution rate matrix A^h for all $j \neq i$:

$$A_{ij}^h \propto \begin{cases} M_{ij} & \text{if } s_{ij}^h = 0 \\ M_{ij} \frac{4s_{ij}^h}{1-\exp(-4s_{ij}^h)} & \text{otherwise} \end{cases} \quad (1.4)$$

Diagonal elements A_{ii}^h are specified to make rows sum to zero: $A_{ii}^h = -\sum_{j \neq i} A_{ij}^h$.

Of the 64 possible codons in the standard genetic code, three (TAA, TAG and TGA) are stop codons that mark the end of a protein-coding segment of DNA. It is standard practice to ignore the possibility of mutations to stop codons because such mutations would truncate the sequence of amino acids and result in a disfunctional protein. A^h can therefore be reduced to a 61×61 matrix. The proportionality constant implied by equation (1.4) can

vary depending on the context, but is set in such a way that time can be expressed as the expected number of single nucleotide substitutions per codon site. See Nielsen and Yang (2003); Yang and Nielsen (2008) for the first codon substitution models based on the mutation-selection framework.

1.1.5 Evolution as a Markov Process

The evolutionary model in equation (1.4) implies a stochastic process in the form of a chain of substitution events $i \rightarrow j$ over time. Explicit in its formulation is the defining characteristic of a Markov chain, namely that the probability distribution of the next event depends only on the current state of the chain. This is often referred to as the Markov property. If the evolutionary process is also assumed to be homogeneous, meaning that A^h does not change over time, then the probability that i is replaced by j at the end of the time interval $(s, s + t)$ is given by the $(i, j)^{th}$ entry of the matrix of substitution probabilities $P^h(t)$. This matrix is computed by solving the system of differential equations $dP^h(t)/dt = P^h(t)A^h$ subject to the constraint that $P^h(s) = I$ where I is the 61×61 identity matrix (Yang, 2006). The solution is obtained by matrix exponentiation $P^h(t) = \exp(tA^h)$ or equivalently:

$$P^h(t) = \sum_{k=0}^{\infty} \frac{(tA^h)^k}{k!} = I + tA^h + \sum_{k=2}^{\infty} \frac{(tA^h)^k}{k!} \quad (1.5)$$

The $(i, j)^{th}$ element of $P^h(t)$ is the conditional probability that the state will be j at time $s + t$ given that the state at any starting time s is i : $P_{ij}^h(t) = P(X(s + t) = j \mid X(s) = i)$. Time-homogeneity implies that the vector of fitness coefficients \mathbf{f}^h is constant. In a real protein \mathbf{f}^h (if it could be measured) can change over time via a number of processes. Detecting such changes is of great interest to biologists because they can be an indication of adaptive evolution ⁶. The majority of fitted codon substitution models nevertheless

⁶The terms “adaptive evolution” and “positive selection” are often considered to be synonymous with a selection coefficient $s_{ij} > 0$ and a rate-ratio $\omega > 1$ (e.g., see Nielsen and Yang, 2003). Here I use “adaptive molecular evolution” to refer to the specific case where the fitness coefficients for the twenty amino acids at a site in a protein-coding sequence change. This is justified by the fact that such changes, which suggest a change in the functional constraints acting on a protein molecule, are what biologists are most interested in detecting. The relevance of this statement will become clear in Chapter 2, where it is

assume time-homogeneity (e.g., the M-series models, Yang *et al.*, 2000a). The analytic model introduced in Chapter 4 relaxes this assumption by making use of contextual information in the form of a discrete phenotype.

1.1.6 Stationary Frequencies and Time Reversibility

Imagine a population in which all members are assigned the i^{th} codon at the h^{th} site in a particular gene at time $t = 0$. The codon occupying the site will change over time via a series of substitution events. Let $\mathbf{v}^h(s) = \langle v_1^h(s), \dots, v_{61}^h(s) \rangle$ be the 1×61 row vector indicating the proportion of time the population was fixed at each of the 61 codons by the end of the time interval $[0, s]$, so that $\sum_{i=1}^{61} v_i^h(s) = 1$. It is useful to construe the transition probability matrix $P^h(t) = \exp(tA^h)$ as an operator that projects $\mathbf{v}^h(s)$ forward in time. That is, $\mathbf{v}^h(s+t) = \mathbf{v}^h(s)P^h(t)$ gives the expected frequencies with which the 61 codons will have occupied the site by time $s+t$ given that the distribution was initially $\mathbf{v}^h(s)$. The operator $P^h(t) = \exp(tA^h)$ can be expressed in terms of the eigensystem (U, Λ) of A^h (i.e., $A^h = U\Lambda U^{-1}$), where the columns of U are orthonormal eigenvectors and Λ is a diagonal matrix of eigenvalues sorted in descending order:

$$P^h(t) = \exp(tA^h) = \sum_{k=0}^{\infty} \frac{(tA^h)^k}{k!} = \sum_{k=0}^{\infty} \frac{(tU\Lambda U^{-1})^k}{k!} \quad (1.6)$$

$$= \sum_{k=0}^{\infty} \frac{t^k (U\Lambda U^{-1})_1 (U\Lambda U^{-1})_2 \dots (U\Lambda U^{-1})_k}{k!} \quad (1.7)$$

$$= \sum_{k=0}^{\infty} \frac{t^k U \Lambda (U^{-1}U) \Lambda (U^{-1}U) \dots (U^{-1}U) \Lambda U^{-1}}{k!} \quad (1.8)$$

$$= \sum_{k=0}^{\infty} U \frac{(t\Lambda)^k}{k!} U^{-1} \text{ since } UU^{-1} = I \quad (1.9)$$

$$= U \left(\sum_{k=0}^{\infty} \frac{(t\Lambda)^k}{k!} \right) U^{-1} \text{ since } U \text{ and } U^{-1} \text{ are constant} \quad (1.10)$$

$$= U \exp(t\Lambda) U^{-1} \quad (1.11)$$

shown that $s_{ij} > 0$ and a rate-ratio $\omega > 1$ can both occur even when fitness coefficients are constant.

The substitution processes at a site can be characterized by the $\lim_{t \rightarrow \infty} \mathbf{v}^h(t) = \mathbf{v}^h(0) \lim_{t \rightarrow \infty} P^h(t)$. To calculate this limit, first note that the largest eigenvalue is always $\Lambda_{11} = 0$, and that the remaining eigenvalues are always less than zero, $\Lambda_{ii} < 0$ for all $i > 1$. Hence $\exp(t\Lambda)$ converges to a matrix with zeros everywhere but with a one in the top-left corner:

$$\lim_{t \rightarrow \infty} \exp(t\Lambda) = \lim_{t \rightarrow \infty} \begin{bmatrix} e^{\Lambda_{11}t} & \dots & 0 \\ & \ddots & \\ 0 & \dots & e^{\Lambda_{61,61}t} \end{bmatrix} = \begin{bmatrix} 1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & 0 \end{bmatrix} \equiv Z \quad (1.12)$$

The effect of Z in $\lim_{t \rightarrow \infty} P^h(t) = UZU^{-1} \equiv P^h(\infty)$ is to make the rows of $P^h(\infty)$ all the same. Let $\boldsymbol{\pi}^h = \langle \pi_1^h, \dots, \pi_{61}^h \rangle$ be the common row of $P^h(\infty)$. Since $\mathbf{v}^h(0)$ is a 1×61 row vector of zeros but with a one in the i^{th} position, the product $\mathbf{v}^h(0)P^h(\infty)$ is just the i^{th} row of $P^h(\infty)$ or in other words $\boldsymbol{\pi}^h$. The initial value i is therefore forgotten in the limit that $t \rightarrow \infty$, a consequence of the Markov property, and so $\lim_{t \rightarrow \infty} \mathbf{v}^h(t) = \boldsymbol{\pi}^h$ for any $i \in \{1, \dots, 61\}$.

An interesting property of $P^h(t)$ is that it projects $\boldsymbol{\pi}^h$ to itself. To see why, consider that $\mathbf{v}^h(s)P^h(t) = \mathbf{v}^h(s+t)$ for any starting time s . Taking the limit as $s \rightarrow \infty$ gives $\boldsymbol{\pi}^h P^h(t) = \boldsymbol{\pi}^h$. This means that the vector $\boldsymbol{\pi}^h$ is stationary with respect to the operator $P^h(t)$ (or equivalently, $\boldsymbol{\pi}^h$ is an eigenvector of $P^h(t)$ with a unit eigenvalue). $\boldsymbol{\pi}^h$ is therefore often referred to as the vector of stationary frequencies for the Markov chain. An alternative way to think about $\boldsymbol{\pi}^h$ is in terms of rates. From equation (1.5) we have that $P_{ij}^h(t) = tA_{ij}^h + o(t)$ for all $i \neq j$ and small t and that $P_{ii}^h(t) = 1 + tA_{ii}^h(t) + o(t)$. It follows that:

$$\begin{aligned} & \lim_{t \rightarrow 0} \frac{P_{ij}^h(s+t) - P_{ij}^h(s)}{t} \\ &= \lim_{t \rightarrow 0} \frac{(s+t)A_{ij}^h - sA_{ij}^h + o(t)}{t} = A_{ij}^h \end{aligned} \quad (1.13)$$

$$\begin{aligned} & \lim_{t \rightarrow 0} \frac{P_{ii}^h(s+t) - P_{ii}^h(s)}{t} \\ &= \lim_{t \rightarrow 0} \frac{1 + (s+t)A_{ii}^h - 1 - sA_{ii}^h + o(t)}{t} = A_{ii}^h \end{aligned} \quad (1.14)$$

Hence

$$A^h = \lim_{t \rightarrow 0} \frac{P^h(s+t) - P^h(s)}{t} \quad (1.15)$$

Multiplying both sides by $\boldsymbol{\pi}^h$:

$$\boldsymbol{\pi}^h A^h = \lim_{t \rightarrow 0} \frac{\boldsymbol{\pi}^h P^h(s+t) - \boldsymbol{\pi}^h P^h(s)}{t} = \lim_{t \rightarrow 0} \frac{\boldsymbol{\pi}^h - \boldsymbol{\pi}^h}{t} = 0 \quad (1.16)$$

The equality $\boldsymbol{\pi}^h A^h = 0$ suggests a substitution process in dynamic equilibrium (i.e., where the net rate of change is zero). For this reason $\boldsymbol{\pi}^h$ is also called the vector of equilibrium frequencies for the Markov chain.

A Markov chain is said to be time-reversible if the probability of consecutive events does not depend on the order of those events. If $X(t)$ is the state of the chain at time t , time-reversibility means that:

$$P(X(0) = i, X(t) = j) = P(X(0) = j, X(t) = i) \quad (1.17)$$

Equation (1.17), known as the detailed balance criterion, can be expressed in terms of substitution probabilities and stationary frequencies:

$$\begin{aligned} &P(X(0) = i, X(t) = j) \\ &= P(X(0) = i) P(X(t) = j \mid X(0) = i) = \pi_i^h P_{ij}^h(t) \end{aligned} \quad (1.18)$$

$$\begin{aligned} &P(X(0) = j, X(t) = i) \\ &= P(X(0) = j) P(X(t) = i \mid X(0) = j) = \pi_j^h P_{ji}^h(t) \end{aligned} \quad (1.19)$$

The necessary and sufficient condition for a Markov chain to be time-reversible is therefore $\pi_i^h P_{ij}^h(t) = \pi_j^h P_{ji}^h(t)$ (or equivalently $\pi_i^h A_{ij}^h = \pi_j^h A_{ji}^h$) for all (i, j) . This is known as Kolomogorov's criterion or more commonly the detailed balance equation. See the Appendix in Chapter 2 for a proof that the substitution processes defined by A^h is time-reversible.

1.1.7 Substitution rates and the Canonical Signature of Positive Selection

The 61 possible codons (excluding stop codons) map to only 20 amino acids. Hence, most amino acids have more than one codon alias. Arginine, leucine and serine, for example, have six aliases each in the standard genetic code, and most amino acids have two or four aliases (Table 1.1). It follows that there are two types of mutations. Synonymous mutations are those that do not change the amino acid (e.g., GCT \rightarrow GCC, both of which code for alanine). The

selection coefficient for any synonymous $i \rightarrow j$ mutation is usually assumed to be $s_{ij}^h = 0$. By this assumption synonymous mutations are fixed at the neutral rate $1/(2N_e)$ (cf. equation 1.1). Mutations that change the amino acid are called nonsynonymous (e.g., GCT \rightarrow ACT changes alanine to threonine). The selection coefficient for any nonsynonymous $i \rightarrow j$ mutation is usually either $s_{ij}^h < 0$ (with fixation rate $< 1/(2N_e)$), indicating that the mutation had a deleterious effect on the fitness of the protein, or perhaps $s_{ij}^h \approx 0$, indicating a neutral or nearly-neutral effect. It is expected that $s_{ij}^h > 0$ (with fixation rate $> 1/(2N_e)$), indicating that the random mutation markedly improved the fitness of the protein, will happen only rarely because most sites are likely to be occupied by an optimal or near-optimal amino acid most of the time. The overall rate of fixation of nonsynonymous mutations is therefore expected to be something less than the neutral rate $1/(2N_e)$ most of the time. The canonical exception to this arises in a gene that has undergone adaptive evolution, when an excess of nonsynonymous substitutions at some sites can be evident, resulting in a nonsynonymous substitution rate greater than $1/(2N_e)$. Hence, an excess of nonsynonymous substitutions can sometimes provide a means to infer the fixation of beneficial mutations, which is commonly interpreted as evidence of adaption (Kimura, 1983; Nei and Gojobori, 1986; Hughes and Nei, 1988).

The nonsynonymous-to-synonymous rate ratio, commonly denoted dN^h/dS^h for the h^{th} site, provides a measure that can be used to detect an excess in the nonsynonymous substitution rate above the neutral rate. The neutral rate dS^h is the ratio of the synonymous substitution rate (rS^h) to the rate at which synonymous mutations are expected to arise (rS_0^h). Working from equation (1.4):

$$dS^h = \frac{rS^h}{rS_0^h} = \frac{\sum_{(i,j)} \pi_i^h M_{ij} 2N_e \Delta t \ell_S \times \frac{1}{2N_e}}{\sum_{(i,j)} \pi_i^h M_{ij} 2N_e \Delta t \ell_S} = \frac{1}{2N_e} \quad (1.20)$$

where ℓ_S is an indicator for synonymous (i, j) pairs. The ratio of the nonsynonymous substitution rate (rN^h) to the rate at which nonsynonymous mutations arise (rN_0^h) is similarly expressed, the only difference being that the

fixation probabilities are different for each nonsynonymous (i, j) pair:

$$\begin{aligned} dN^h &= \frac{rN^h}{rN_0^h} = \frac{\sum_{(i,j)} \pi_i^h M_{ij} 2N_e \Delta t \ell_N \times \frac{2s_{ij}^h/N_e}{1-\exp(-4s_{ij}^h)}}{\sum_{(i,j)} \pi_i^h M_{ij} 2N_e \Delta t \ell_N} \\ &= \frac{\sum_{(i,j)} \pi_i^h A_{ij}^h \ell_N}{\sum_{(i,j)} \pi_i^h M_{ij} \ell_N} \times \frac{1}{2N_e} \end{aligned} \quad (1.21)$$

The expected site-specific rate ratio is therefore:

$$dN^h/dS^h = \frac{\sum_{(i,j)} \pi_i^h A_{ij}^h \ell_N}{\sum_{(i,j)} \pi_i^h M_{ij} \ell_N} \quad (1.22)$$

The fixation rate under a neutral regime, when $s_{ij}^h = 0$ for all nonsynonymous (i, j) pairs, is $dN^h = 1/(2N_e)$ and so corresponds to $dN^h/dS^h = 1$. Similarly, the fixation rate under a stringent selection regime, when $s_{ij}^h < 0$ for all nonsynonymous (i, j) pairs (i.e., at a site occupied by its fittest codon), is $dN^h < 1/(2N_e)$ and so corresponds to $dN^h/dS^h < 1$. And the fixation rate when $s_{ij}^h > 0$ for all nonsynonymous (i, j) pairs (i.e., at a site occupied by its least fit codon) is $dN^h > 1/(2N_e)$, and so corresponds to $dN^h/dS^h > 1$, the canonical signature of positive selection (Goldman and Yang, 1994).

Kimura's neutral theory of molecular evolution posits that, because the majority of mutations are selectively deleterious and rapidly eliminated from populations while very few are beneficial, most of molecular evolution is due to the fixation of selectively neutral alleles by drift (Kimura, 1968, 1983). By this theory, the probability that a neutral mutation arises in the next generation of a diploid population of size N_e is $\mu 2N_e$, where μ is the rate at which neutral mutations arise in a gene per individual per generation. The probability that the mutation is fixed is $1/(2N_e)$. It follows that the rate at which a given gene accumulates neutral differences is $\mu 2N_e \times 1/(2N_e) = \mu$. The rate at which neutral mutations are fixed over macroevolution time scales (millions of years) is therefore equal to the rate at which neutral mutations arise in a population over microevolutionary time scales (thousands of generations).

1.2 Codon Substitution Models

1.2.1 The Data

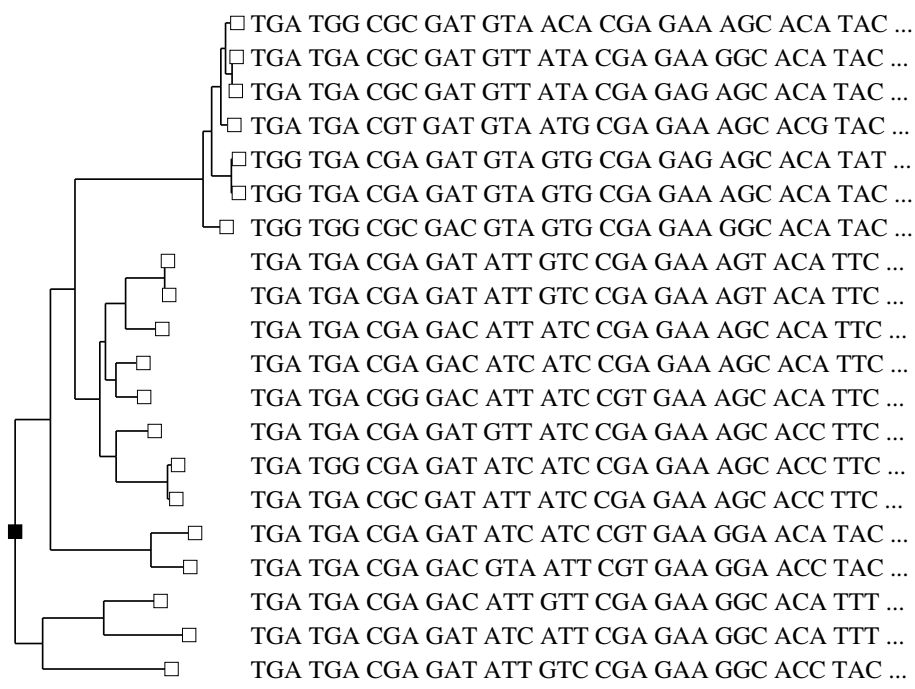


Figure 1.1: An example of the object of analysis. A phylogenetic tree $T = (\tau, \mathbf{t})$ (left) and an alignment of homologous protein-coding sequences X (right). The tree is composed of two components, a topology τ and a vector of branch lengths \mathbf{t} . Terminal nodes (white boxes) indicate extant representatives of different taxonomic lineages (e.g., species). Sequences are assumed to be homologous, meaning that they evolved from a common ancestor indicated by the root node of the tree (black box). Each column of the alignment represents a site pattern. A codon substitution model (CSM) can be fitted to X conditioned on τ using maximum likelihood. The fitted CSM assigns to each possible site pattern \mathbf{x} a probability $P(\mathbf{x}; \hat{\boldsymbol{\theta}}, \hat{\mathbf{t}})$ such that $\sum_{\mathbf{x}} P(\mathbf{x}; \hat{\boldsymbol{\theta}}, \hat{\mathbf{t}}) = 1$.

The MS framework introduced in the previous section provides a way to think about evolutionary processes at an individual codon site. Here attention is turned to the problem of inference. The objective is to extract meaningful summaries of the evolutionary processes that produced a given set of extant genetic data. To emphasize the magnitude of the problem, consider the ubiquitous protein cytochrome B. CytB is an essential component of the electron

transport chain that drives the synthesis of the energy-storage molecule adenosine triphosphate (ATP). CytB is consequently found in a wide range of life forms, meaning that its evolutionary history extends billions of years backward in time. Yet everything that can be said about the evolutionary history of cytB can only be inferred from extant variants of the gene. A rooted phylogenetic tree of N variants of cytB has $N - 1$ internal nodes each representing an unknown ancestral sequence. Given that there are 61^{N-1} possible ancestral histories for each codon site, and that the true history is unknown, one might ask whether it is possible to infer anything about the evolution of the gene. In fact it is possible: the strategy is to sum over all ancestral histories at a site (i.e., all possible combinations of codons at the ancestral nodes of the tree) each weighted by its probability under a given fitted CSM. Summation can be accomplished using the efficient pruning algorithm (Felsenstein, 1981) provided codon sites are assumed to have evolved independently of one another.

Prior to introducing the basic elements of a CSM and typical methods of inference, I first briefly outline the data and where it comes from. The data consists of a set of homologous sequences, each representing a different taxonomic grouping thought to have evolved from a common ancestor. The set is referred to as an alignment, and can be represented by an $N \times n$ matrix X where N is the number of sequences (rows) and n their common length in codons (columns, e.g., the $(r, c)^{th}$ element of X is some $x_{rc} \in \{1, \dots, 61\}$ that represents the codon in the r^{th} sequence at the c^{th} site). The alignment is usually accompanied by a topology τ that specifies the relationship between sequences via a series of ancestral bifurcations. The alignment X and topology τ can be assembled as follows. The first step is to obtain a sequence S_1 corresponding to one variant of the target gene. A good place to look is the UCSC Genome Browser (University of Santa Cruz Genomic Institute), which provides access to a large genetic data base along with a variety of search tools. Next, S_1 must be fed into an automated search algorithm to find homologous sequences likely to be related to S_1 by common ancestry. This can be done using the Basic Local Alignment Search Tool (BLAST)

hosted by the National Center for Biotechnology Information, which provides a ranked list of candidate homologues from which the user can choose based on various statistical measures of similarity. Homologous sequences can vary in composition. Eukaryotic genes, for example, often have introns or segments of non-coding DNA interspersed among exons, the segments of coding DNA. The arrangement of introns and exons can vary from one homologue to the next. The selected sequences S_1, \dots, S_N must therefore be aligned so that homologous sites appear in the same column of X using an algorithm such as Clustal Omega hosted by the European Molecular Biology Laboratory. The last step is to estimate the topology τ of the phylogenetic relationships between the aligned sequences. A variety of tree-estimation tools are available. A good place to start is RAxML (Randomized Axelerated Maximum Likelihood) because it is fast and easy to use. All of the above mentioned tools are free to use and available online.

Each step from S_1 to (X, τ) is the subject of various degrees of continuing research. I leave these efforts to experts in the relevant fields. For the purpose of this dissertation, the application of my expertise in CSMs to real data will usually start with previously assembled (X, τ) assumed to be error-free. A CSM can be fitted to X under the assumed topology using maximum likelihood (ML: here I focus on ML and for convenience use CSM to indicate a model that is used in conjunction with the ML approach; see Huelsenbeck and Dyer (2004) for an example of the Bayesian approach). The result is a vector $\hat{\boldsymbol{\theta}}$ of maximum-likelihood estimates (MLEs) of all the parameters included in the model and a vector of estimated branch lengths $\hat{\mathbf{t}}$. Together these define a probability distribution that assigns to each of the 61^N possible site patterns \mathbf{x} at the terminal nodes of the tree a probability $P(\mathbf{x}; \hat{\boldsymbol{\theta}}, \hat{\mathbf{t}})$ such that $\sum_{\mathbf{x}} P(\mathbf{x}; \hat{\boldsymbol{\theta}}, \hat{\mathbf{t}}) = 1$ (Figure 2.6). Under the assumption that sites evolved independently⁷, the likelihood of the alignment under the fitted model is a product across sites:

$$L(X; \boldsymbol{\theta}, \mathbf{t}) = \prod_{h=1}^n P(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) \quad (1.23)$$

⁷The validity of various model assumptions are discussed in the last subsection of this chapter.

where \mathbf{x}^h is the h^{th} column of X . The MLE $(\hat{\boldsymbol{\theta}}, \hat{\mathbf{t}})$ is the vector that maximizes equation (1.23).

1.2.2 A Simple CSM

The majority of CSMs are based on an underlying continuous-time homogeneous and time-reversible Markov process that describes the rate at which substitutions occur under a neutral regime (i.e., for which $dN/dS = 1$). This process can be specified by a 61×61 substitution rate matrix M . There are numerous ways to define M (e.g., Muse and Gaut, 1994; Goldman and Yang, 1994; Zaheri *et al.*, 2014). The following model is used throughout this thesis (Jones *et al.*, 2018, 2019b):

$$M_{ij} \propto \begin{cases} \kappa^{s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 1 \\ \alpha \kappa^{s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 2 \\ \beta \kappa^{s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 3 \end{cases} \quad (1.24)$$

Equation (1.24) applies to all pairs of codons (i, j) that differ by $s \in \{1, 2, 3\}$ nucleotides, s_t of which are transitions (substitutions of the form $T \leftrightarrow C$ between the pyrimidine nucleotides or $A \leftrightarrow G$ between the purine nucleotides) and $s - s_t$ of which are transversions (substitutions of the form $\{T, C\} \leftrightarrow \{A, G\}$ from a pyrimidine to a purine or vice versa). The $\pi_{j_k}^*$ represent position-specific nucleotide frequencies, $\kappa \geq 1$ the transition bias (i.e., transition substitutions tend to occur more frequently than transversion substitutions), and α and β the rate at which double and triple (DT) substitutions arise, respectively. Diagonal elements M_{ii} are adjusted to make rows sum to zero. To illustrate, the rate at which GCT (alanine) is replaced by TTT (phenylalanine) under this model is $\alpha \kappa \pi_{T_1}^* \pi_{T_2}^*$, where α accounts for the double substitution, κ the $C \rightarrow T$ transition, $\pi_{T_1}^*$ the frequency of T in the first codon position, and $\pi_{T_2}^*$ the same for T in the second codon position. The stationary frequency of any codon under equation (1.24) is proportional to the product of its position-specific nucleotide frequencies (e.g., $\pi_{ACT} \propto \pi_{A_1}^* \pi_{C_2}^* \pi_{T_3}^*$). The majority of CSMs allow single nucleotide substitutions only, which can be effected by setting $\alpha = \beta = 0$. Note that equation (1.24) is interpreted differently depending on whether it is used in a CSM or as part of the MS model introduced in the first section of this

chapter. It is not possible to separate the mutation and selection processes that produced an alignment because only those mutations that were fixed can be observed. In the context of a CSM, it is therefore better to think of M as the rate at which substitutions occur under a neutral selection regime. Under MS by contrast, which is presented in this thesis as a data-generating process (i.e., not a fitted model), M characterizes the mechanisms by which mutations arise, and is therefore correctly construed as a mutation rate matrix.

Selection effects can be introduced into the model via a parameter ω representing a nonsynonymous-to-synonymous rate ratio as follows (where \circ represents the element-wise matrix product):

$$Q(\omega) = M \circ (\boldsymbol{\ell}_S + \omega \boldsymbol{\ell}_N) / r, \text{ where } r = \sum_{j \neq i} \pi_i Q_{ij}(\omega) \{\ell_1 + 2\ell_2 + 3\ell_3\} \quad (1.25)$$

$\boldsymbol{\ell}_S$ is a 61×61 indicator matrix whose $(i, j)^{th}$ element is one if i and j are synonymous and zero otherwise, and $\boldsymbol{\ell}_N$ is a similar indicator matrix for non-synonymous codon pairs. Diagonal elements of $Q(\omega)$ are modified to make rows sum to zero. The indicator ℓ_k is one if i and j differ by $k \in \{1, 2, 3\}$ nucleotides and zero otherwise. The constant r scales $Q(\omega)$ so that branch lengths give the expected number of single nucleotide substitutions per codon site. Note that the stationary frequencies satisfy both $\boldsymbol{\pi}M = 0$ and $\boldsymbol{\pi}Q(\omega) = 0$, and in fact, any pair of matrices $Q(\omega_1)$ and $Q(\omega_2)$ defined by equation (1.25) share the same vector of stationary frequencies.

1.2.3 The Pruning Algorithm

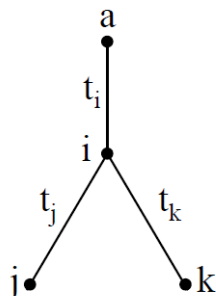


Figure 1.2: An arbitrary branching element in a binary tree. Node a is ancestral and i descendant.

Suppose the CSM represented by equation (1.25) is to be fitted to an $N \times n$ alignment X and topology τ . That is, suppose our objective is to maximize equation (1.23) with respect to parameters $\boldsymbol{\theta} = \langle \kappa, \omega \rangle$ and branch lengths \mathbf{t} using an optimization algorithm. Roughly speaking, optimization is achieved by stepping through parameter space $\langle \boldsymbol{\theta}, \mathbf{t} \rangle$ in “uphill” directions until a maximum is reached. Each step requires re-evaluation of equation (1.23), which entails computation of $P(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t})$ at each site in the alignment. An efficient method to compute $P(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t})$ is presented in this section.

Let $\mathbf{x}^h = \langle x_1^h, \dots, x_N^h \rangle$ represent the vector of codons at the h^{th} site at the N terminal nodes of a rooted tree, and let $\mathbf{c}^h = \langle c_{N+1}^h, \dots, c_{2N-1}^h \rangle$ represent one of the 61^{N-1} possible ancestral histories for the site, the codons that occupied the site at each of the $N-1$ internal nodes of the tree, where c_{2N-1}^h is presumed to be the codon at the root node. Let $a_b^h \in \{1, \dots, 61\}$ be the codon at the ancestral node and $d_b^h \in \{1, \dots, 61\}$ the codon at the daughter node of any given branch b . Given the ancestral state a_b^h , the substitution process at the site along that branch is assumed to be independent of the process at that site on other branches of the tree. Hence, if the site evolved via the Markov chain specified by $Q(\omega)$, then the probability of $\langle \mathbf{x}^h, \mathbf{c}^h \rangle$ is just a product over branches:

$$P(\langle \mathbf{x}^h, \mathbf{c}^h \rangle; \boldsymbol{\theta}, \mathbf{t}) = \pi_{c_{2N-1}^h} \prod_{b=1}^{2N-2} P(a_b^h \rightarrow d_b^h) \quad (1.26)$$

Note that the codon at the root of the tree is assigned the unconditional probability $\pi_{c_{2N-1}^h}$ equal to its stationary frequency. The probability $P(a_b^h \rightarrow d_b^h)$ is given by the $(a_b^h, d_b^h)^{th}$ element of the transition probability matrix $P(t_b) = \exp(t_b Q(\omega))$, where t_b is the length of branch b . The marginal probability of \mathbf{x}^h is the sum over all possible histories:

$$P(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) = \sum_{\mathbf{c}^h} P(\langle \mathbf{x}^h, \mathbf{c}^h \rangle; \boldsymbol{\theta}, \mathbf{t}) \quad (1.27)$$

This sum is over 61^{N-1} elements, and would therefore seem to be infeasible for all but the smallest number of sequences. However, the probabilities $P(\langle \mathbf{x}^h, \mathbf{c}^h \rangle; \boldsymbol{\theta}, \mathbf{t})$ in the summation includes factors $P(a_b^h \rightarrow d_b^h)$ that appear many times over. Efficiency can be gained by computing such factors only once. This is essentially what is achieved by the pruning algorithm (Felsenstein, 1981).

The pruning algorithm populates a $(2N - 1) \times 61$ matrix V with entries v_{ij} whose rows correspond to nodes $1, \dots, 2N - 1$ of the tree and columns to codons $1, \dots, 61$. For any $i \in \{1, \dots, N\}$ corresponding to the terminal nodes of the tree, the row v_i of V is an indicator vector constructed as follows:

$$v_{ij} = \begin{cases} 1 & \text{if the codon at terminal node } i \text{ is } j \\ 0 & \text{otherwise} \end{cases} \quad (1.28)$$

The remaining rows of V represent conditional probabilities. Consider the node $i \in \{N + 1, \dots, 2N - 2\}$ as depicted in Figure 1.2. The rows v_j and v_k will already have been constructed. They, along with the transition probability matrices $P(t_j) = \exp(t_j Q(\omega))$ and $P(t_k) = \exp(t_k Q(\omega))$, determine v_i :

$$v_i = (v_j P(t_j)^T) \circ (v_k P(t_k)^T) \quad (1.29)$$

where $P(t)^T$ is the transpose of $P(t)$. Let \mathbf{x}_i^h be the vector of codons at the terminal nodes that descend from the i^{th} internal node of the tree. Elements of v_i give the conditional probability of \mathbf{x}_i^h given the state c_i^h at the i^{th} node:

$$v_i = \langle P(\mathbf{x}_i^h | c_i^h = 1), \dots, P(\mathbf{x}_i^h | c_i^h = 61) \rangle \quad (1.30)$$

The probability of \mathbf{x}^h is computed from the last row of V (i.e., row $2N - 1$)

corresponding to the root of the tree, here expressed as a dot product:

$$P(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) = \langle \pi_1, \dots, \pi_{61} \rangle \cdot v_{2N-1} = \sum_{i=1}^{61} \pi_i P(\mathbf{x}^h \mid c_{2N-1}^h = i) \quad (1.31)$$

The efficiency of the pruning algorithm is evident in the fact that the number 61^{N-1} of elements in the summation in equation (1.27) increases exponentially with the number of taxa N whereas the number of rows $2N - 1$ that need to be computed to calculate equation (1.31) increases only linearly with N .

1.2.4 Inferring Positive Selection

CSMs were initially formulated to provide a means to test whether some of the differences between homologous protein-coding sequences might have been due to positive selection. Suppose the null model is $Q(\omega)$ with ω constrained to be equal to one, and the alternative $Q(\omega)$ with ω constrained to be greater than one. Further suppose that each model was fitted to an alignment X to produce a pair of likelihoods $L_{\text{nul}}(X; \hat{\boldsymbol{\theta}}_{\text{nul}}, \hat{\mathbf{t}}_{\text{nul}})$ and $L_{\text{alt}}(X; \hat{\boldsymbol{\theta}}_{\text{alt}}, \hat{\mathbf{t}}_{\text{alt}})$. Since the two models are nested, an omnibus test for the $\omega > 1$ signature of positive selection can be conducted by comparing the resulting log-likelihood ratio (LLR) to its theoretical limiting distribution. The LLR is computed as follows:

$$\text{LLR} = 2 \left(\ln \left\{ L_{\text{alt}}(X; \hat{\boldsymbol{\theta}}_{\text{alt}}, \hat{\mathbf{t}}_{\text{alt}}) \right\} - \ln \left\{ L_{\text{nul}}(X; \hat{\boldsymbol{\theta}}_{\text{nul}}, \hat{\mathbf{t}}_{\text{nul}}) \right\} \right) \quad (1.32)$$

Since the two models differ by a single parameter (i.e., $\omega = 1$ is fixed under the null and estimated under the alternative), the theoretical limiting distribution of the LLR is χ_1^2 . Hence, we could infer that the gene underwent positive selection at the 5% level of significance whenever $\text{LLR} > 3.84$. Note that the omnibus test is conducted with the assumption that the data was generated under the null model. This can be true in simulation, but will seldom if ever be the case for a real alignment. Furthermore, even if the data was generated under the null, the assumed distribution is still an approximation because it is an expectation that holds only in the limit that information (i.e., number of site patterns in the alignment) is infinite. Violation of regularity assumptions can also be an issue. There are cases in which the null model is derived from the alternate model by setting a parameter to a value on the boundary of

parameter space. Some of these cases have a known limiting distribution, while others do not (Self and Liang, 1987). And in some cases a parameter in the alternate model becomes unidentifiable under the null. The assumed distribution of the LLR is therefore only an approximation when the test is applied to real data.

Models similar to $Q(\omega)$ that estimate a single rate ratio for all sites and branches (e.g., M0, Nielsen and Yang, 1998) tend to have low power to detect sites that underwent positive selection because most sites evolve under stringent selection with $dN/dS \ll 1$ most of the time. Exceptions include genes under constant selection pressure (e.g., the class I major histocompatibility complex first analyzed by Hughes and Nei (1988); see Yang and Bielawski (2000) for other examples). The quest for greater power to detect positive selection has led to models that account for spatial and temporal variations in ω . There currently exists numerous CSMs known as branch-site models designed to detect evidence of positive selection at individual sites along individual branches (e.g., Yang and Nielsen, 2002; Zhang *et al.*, 2005; Kosakovsky Pond *et al.*, 2011; Murrell *et al.*, 2015; Smith *et al.*, 2015). Several models of this type are described in Chapters 2 to 4.

The CSM commonly designated M3(k) is one of a number of M-series models introduced by Yang *et al.* (2000a) to account for variation in dN/dS across sites. Here k designates the number of ω -categories included in the model. For example, the model M3($k = 2$) assumes that each site in an alignment evolved under one of two possible rate ratios $\omega_1 < \omega_2$ across the entire tree in proportions p_1 and $p_2 = 1 - p_1$. In this case, the likelihood is a weighted average:

$$L_{M3}(X; \boldsymbol{\theta}_{M3}, \mathbf{t}) = p_1 L_{\omega_1}(X; \kappa, \omega_1, \mathbf{t}) + p_2 L_{\omega_2}(X; \kappa, \omega_2, \mathbf{t}) \quad (1.33)$$

Here $\boldsymbol{\theta}_{M3} = \langle \kappa, \omega_1, \omega_2, p_1 \rangle$ and $L_{\omega_k}(X; \kappa, \omega_k, \mathbf{t})$ is the likelihood of the alignment under the model $Q(\omega_k)$ of equation (1.25) for $k \in \{1, 2\}$. A test for positive selection using M3($k = 2$) could be conducted by contrasting a null model for which $\omega_1 < 1$ and $\omega_2 = 1$ with an alternate model for which $\omega_1 < 1$ and $\omega_2 > 1$. Rejection of the null provides evidence for two categories of sites, those that evolved under a relatively stringent selection regime with $\hat{\omega}_1 < 1$

and those that evolved under positive selection with $\hat{\omega}_2 > 1$.

1.2.5 *Post Hoc Analysis*

Suppose the LLR for the contrast between M3($k = 2$) with $\omega_1 < 1$ and $\omega_2 = 1$ and M3($k = 2$) with $\omega_1 < 1$ and $\omega_2 > 1$ was large enough to reject the null hypothesis (i.e., $\text{LLR} > 3.84$). This would indicate that the gene underwent positive selection. The question remains Which sites underwent positive selection? The most commonly used method to answer this question is to compute naive empirical Bayes (NEB, Yang *et al.*, 2005) posteriors, which for M3($k = 2$) are as follows:

$$P(\hat{\omega}_2 | \mathbf{x}^h) = \frac{L_{\omega_2}(\mathbf{x}^h; \hat{\kappa}, \hat{\omega}_2, \hat{\mathbf{t}}) \hat{p}_2}{L_{\text{M3}}(\mathbf{x}^h; \hat{\boldsymbol{\theta}}_{\text{M3}}, \hat{\mathbf{t}})} \quad (1.34)$$

$P(\hat{\omega}_2 | \mathbf{x}^h)$ approximates the posterior probability that the h^{th} site evolved under $\hat{\omega}_2 > 1$ given the estimated prior \hat{p}_2 and other MLEs. This approach is called “naive” because it treats MLEs as if they were known without error. The approach can be problematic when the information content of the alignment (i.e. the number of substitutions) is low. Bayes empirical Bayes (BEB Yang *et al.*, 2005) and smoothed bootstrap aggregation (Mingrone *et al.*, 2016) are alternatives to NEB that mitigate problems associated with errors in MLEs. By any method, a site is inferred to have evolved under positive selection when the posterior is greater than some threshold (e.g., when $P(\hat{\omega}_2 | \mathbf{x}^h) > 0.95$).

1.2.6 *Heterotachy and the Covarion-like Model*

Covarion-like (CL) models (e.g., Galtier, 2001; Guindon *et al.*, 2004) have been proposed to account for intragenetic epistatic interactions thought to be the cause of temporal variations in site-specific evolutionary rates (i.e., the covarion phenomenon, Fitch and Markowitz, 1970; Fitch, 1971). We say “covarion-like” because, although the covarion phenomenon was thought to arise due to dependencies between interacting sites, CL models maintain the assumption of site independence and only mimic the phenomenological signature of epistasis. The simplest expression of this signature is captured by the covarion-like version of M3($k = 2$) called CLM3($k = 2$) (Jones *et al.*, 2017). I introduce CLM3($k = 2$) here because the model appears in all subsequent

chapters of this dissertation and plays a particularly significant role in the phenotype-genotype model presented in Chapter 4.

CLM3($k = 2$) assumes two ω -categories, $\omega_1 < \omega_2$, just as M3($k = 2$) does. But unlike M3($k = 2$), where each site is assumed to have evolved under the same rate ratio over the entire tree, a site under CLM3($k = 2$) is assumed to have switched between $\omega_1 < \omega_2$ randomly over time at a rate δ switches per unit branch length. The rate matrix for CLM3($k = 2$) can be constructed by expanding the state space from the 61 possible codons to the 122 possible (codon, ω) pairings where $\text{codon} \in \{1, \dots, 61\}$ and $\omega \in \{\omega_1, \omega_2\}$. The corresponding rate matrix is a concatenation of $Q(\omega_1)$ and $Q(\omega_2)$:

$$Q_{\text{CLM3}} = \frac{1}{c_1} \left[\begin{array}{c|c} Q(\omega_1) & 0 \\ \hline 0 & Q(\omega_2) \end{array} \right] + \frac{\delta}{c_2} \left[\begin{array}{c|c} -p_2 I & p_2 I \\ \hline p_1 I & -p_1 I \end{array} \right] \quad (1.35)$$

Here I is the identity matrix that is the same size as $Q(\omega_1)$ and $Q(\omega_2)$ (e.g., 61×61 for the standard genetic code), p_1 is the average proportion of time sites evolved under ω_1 , and $p_2 = 1 - p_1$ the average proportion of time sites evolved under ω_2 . The vector of stationary frequencies for all possible (codon, ω) pairs is the 1×122 vector $\langle p_1 \boldsymbol{\pi}, p_2 \boldsymbol{\pi} \rangle$, where $\boldsymbol{\pi} Q(\omega_1) = \boldsymbol{\pi} Q(\omega_2) = 0$. The scaling factor $1/c_1$ is set to make branch lengths equal to the expected number of single nucleotide substitutions per codon: $c_1 = p_1 r_1 + p_2 r_2$, where $r_k = \sum_{j \neq i} \pi_i Q_{ij}(\omega_k) \{\ell_1 + 2\ell_2 + 3\ell_3\}$ for $k \in \{1, 2\}$. Including the scaling factor $1/c_2 = 1/(2p_1 p_2)$ permits δ to be interpreted as the expected number of switches between ω_1 and ω_2 per unit branch length (Jones *et al.*, 2018).

1.2.7 Simplifying Assumptions

CSMs require a number of simplifying assumptions. For example, it is assumed that all sites share the same vector of stationary frequencies $\boldsymbol{\pi}$. This is unrealistic because the fittest codon at many sites depends on the physicochemical properties required to make the protein fold and function properly (Patthy, 2008), and will typically vary across sites. Yet this assumption is necessary in most cases due to sparse information. Even a large alignment consisting of 100 sequences yields only 100 samples from which to estimate a site-specific vector of 61 codon frequencies. Nevertheless, models have been

developed to account for differences in site-specific frequencies (e.g., Tamuri *et al.*, 2012, 2014), and appear to be reliable when fitted to large alignments (e.g., of 512 sequences, Spielman and Wilke, 2016). It is also assumed that all codon sites evolve independently. A sequence of amino acids bound together in a specific folded structure will often include networks of closely interacting residues. The fittest codon at a site can change in response to substitutions at sites it interacts with (i.e., via intraprotein epistasis, Pollock *et al.*, 2012; Starr and Thornton, 2016). It is therefore unclear what impact the independence assumption might have on the validity of inferences. What is clear is that the assumption allows the likelihood of an alignment to be handily expressed as a product of site-specific likelihoods (see equation 1.23) and thereby vastly simplifies what could otherwise be an unmanageable calculation.

More subtle is the fact that the basic model represented by $Q(\omega)$ in equation (1.25) assumes that the rate at which the codon i occupying a site is replaced by a nonsynonymous mutation j is the same for all nonsynonymous (i, j) pairs. Mechanistically speaking, it is as if fitness applies not to the codons themselves (as it does under MS), but to the codon site. Hence, a site evolving with $\omega > 1$ is perpetually unfit under $Q(\omega)$, no matter what codon occupies the site. This observation emphasizes the difference between the mechanistic construal of evolutionary processes at a codon site under the MS framework and the phenomenological summary of such processes offered by CSMs. Given all of their simplifications, it is somewhat surprising that CSMs can provide reliable inferences under a wide range of conditions, but that they do has been demonstrated by numerous simulation studies (e.g., Anisimova *et al.*, 2001, 2002; Wong *et al.*, 2004; Zhang, 2004; Kosakovsky Pond and Frost, 2005; Yang *et al.*, 2005; Zhang *et al.*, 2005; Yang and dos Reis, 2011; Kosakovsky Pond *et al.*, 2011; Lu and Guindon, 2013). Furthermore, CSMs have produced numerous biologically plausible results (e.g., Yang and Bielawski, 2000; Yang *et al.*, 2000b; Bielawski *et al.*, 2004; Yang, 2005; Field *et al.*, 2006; Anisimova and Kosiol, 2009; Zhai *et al.*, 2012; Romero *et al.*, 2016).

1.3 Thesis Outline

During the early phase or Phase I of CSM development, the main concern was the possibility that $\omega > 1$ might be incorrectly inferred due to low information content or the mismatch between the fitted model and the actual alignment-generating process. Concern was assuaged by numerous simulation studies in which CSMs were shown to be reliable when fitted to alignments generated *in silico* under a wide range of scenarios. A notable feature of Phase I, however, was the practice of testing CSMs using alignments generated from models based on the same CSM framework. It is becoming increasingly clear that, whereas CSMs are appropriate as tools to extract meaningful phenomenological summaries from genetic data, they are inappropriate as a means to generate data for model testing. Specifically, the research presented in this dissertation demonstrates that traditional simulation methods do not generate site-patterns consistent with variations in site-specific rate ratios (a.k.a. heterotachy) that match what is commonly observed in real data.

The mutation-selection (MS) framework introduced in Chapter 1 provides a way to generate heterotachy commensurate to that observed in real data. Each codon site can be assigned its own vector of fitness coefficients \mathbf{f}^h for the twenty amino acids to characterize a site-specific fitness landscape. Substitutions at the site can be thought of as movement over this landscape that generates temporal dynamics in rate ratio. When selection is not too stringent, a site can occasionally move away from its optimal amino acid A by drift to a suboptimal residue B and then back again to A . This process can be accompanied by evidence of a change in the rate ratio at the site from some $\omega_A < 1$ while the site was occupied by A to $\omega_B > \omega_A$ following substitution to B , and then back to ω_A once the site is reoccupied by A . The episodic occurrence of this non-adaptive shifting balance process across sites and over time can result in phenomenological shifts between rate ratios $\omega_1 < \omega_2$, often with $\omega_2 > 1$. The possibility of episodic positive selection via non-adaptive shifting balance undermines the standard paradigm that equates $\omega > 1$ to adaptation. This issue is explored in Chapter 2.

The majority of CSMs assume that sites evolve by single-nucleotide steps.

However, double and triple (DT) substitutions do sometime occur, and so it seems reasonable to introduce parameters to account for them. This was done with a variety of CSMs including the novel model RaMoSS (for Random Mixture of Static and Switching sites). RaMoSS accounts for a mixture of sites that evolved under a constant rate ratio and sites that evolved with heterotachy. A contrast based on RaMoSS fitted to a real alignment of mammalian mtDNA indicated that 9.7% of substitutions were DT. To test this result, alignments were generated without DT using either standard Phase I methods (Phase I alignments) or a novel generating model based on the MS framework and tuned to produce data similar to the real mtDNA (MS alignments). The false positive rate was less than 5% when the contrast was fitted to Phase I alignments. However, the same contrast produced a substantial number of false positives when fitted to the more realistic MS alignments. This demonstrates that Phase I data-generating methods can be inadequate for model testing.

The use of realistic data-generating methods for model testing marks the beginning of a new Phase II of CSM development. Phase II is also marked by the discovery of a novel statistical pathology called phenomenological load (PL). Suppose ψ represents some process P_1 that did not occur when a particular set of data was generated. Further, suppose process P_2 did occur when the data was generated, and that P_2 tends to produce patterns in the data similar to those produced by P_1 . Processes P_1 and P_2 are said to be confounded. Rejection of the null hypothesis in a contrast testing the significance of ψ is likely because ψ can account for variations in the data generated by P_2 . And although rejection of the null would be correct as an indication that the inclusion of ψ improved model fit, it would also lead to the false conclusion that process P_1 actually occurred. When this happens $\hat{\psi}$ is said to carry phenomenological load. Confounding and PL in the context of detecting fixation of DT mutations are covered in Chapter 3.

Heterotachy in the form of random shifts between $\omega_1 < \omega_2$ can be generated not only by non-adaptive shifting balance but also by episodic adaptive

changes in site-specific fitness coefficients (a.k.a. peak shifts). The two processes are therefore confounded in alignment data, and cannot be disentangled using the traditional $\omega > 1$ criterion alone. The novel Phenotype-Genotype Branch-Site Model (PG-BSM) was formulated to break confounding by including contextual information in the form of a discrete phenotype to discriminate heterotachy-by-any-cause (the null model) from specific patterns of change between $\omega_1 < \omega_2$ that occurred in association with changes in phenotype (the alternate model). Extensive simulation studies were used to demonstrate the reliability of the PG-BSM as a tool to detect sites associated with phenotype via specific mechanisms of adaptation. Analyses of real data sets were conducted to demonstrate the potential utility of the approach. Significantly, the model was shown to be capable of inferring adaptive evolution in association with phenotype at individual codons sites even when the estimate of ω_2 was < 1 . The PG-BSM is the topic of Chapter 4.

The first section of Chapter 5 consists of case studies that illustrate problems and solutions associated with Phase I and Phase II of CSM development. It is argued that the discovery of the problems associated with confounding and PL introduced in Chapters 2 to 4 mark the transition between the historical Phase I and the current Phase II. As to the way forward, some maintain that the next phase of CSM development should see an increase in the mechanistic content of fitted models. My findings suggest that this is likely to be infeasible due to limitations in the information contained in alignment data in which many processes are likely confounded. Instead, it would seem more appropriate to increase the mechanistic content of alignment-generating models to provide more realistic data to test phenomenological CSMs. The first section of Chapter 5 ends with an argument supporting this view. The second section of Chapter 5 provides other examples to illustrate that problems associated with confounding and phenomenological load are common to many areas of study. Confounding and confounding-breaking using contextual information is illustrated in the context of satellite remote sensing, where the objective was to discriminate image features associated with sea-surface temperature fronts (e.g., the north wall of the Gulf Stream) from those associated

with atmospheric processes. PL is shown to invalidate attempts to provide a mechanistic explanation for variations in the efficiency of photosynthesis. Given my rejection of $\omega > 1$ as proof of adaptive evolution, I end this dissertation with my thoughts on what the canonical signature of adaptation at the molecular level should in fact be.

Chapter 2

Non-adaptive Shifting Balance on a Static Mutation-selection Landscape: A Novel Scenario of Positive Selection.

2.1 Introduction

Codon substitution models (CSMs) have provided the basis for the most commonly used methods of inferring positive selection in protein-coding sequences since the pioneering efforts of Muse and Gaut (1994) and Goldman and Yang (1994). Such models produce estimates of the ratio of the nonsynonymous substitution rate (after adjusting for neutral opportunity, dN) to the synonymous substitution rate (likewise adjusted, dS). The rate ratio dN/dS is represented by the parameter ω in a 61×61 substitution rate matrix that is the building block for a variety of popular CSMs. The simplest CSM is M0 (Nielsen and Yang, 1998; Yang *et al.*, 2000a), which estimates a single ω for all sites and branches. The limited statistical power of M0 spurred the development of models that account for variation in ω across branches (Yang and Nielsen, 1998), across sites (e.g., the M-series models of Yang *et al.*, 2000a), and across both branches and sites (Guindon *et al.*, 2004; Yang *et al.*, 2005; Kosakovsky Pond *et al.*, 2011; Murrell *et al.*, 2015; Smith *et al.*, 2015). Positive selection is inferred when a model that permits ω to be greater than one fits the data significantly better than a nested version of the same model for which all ω are restrained to be one or less. Such inferences are characteristic of two positive selection scenarios: episodic changes in functional constraints causing transient increases in ω , and frequency-dependent selection causing sustained elevations in ω along an entire lineage. The signal for episodic selection is typically restricted to a few branches of a phylogeny, and can occur in association with events such as horizontal gene transfer (e.g., Yang *et al.*,

2013), gene duplication (e.g., Pegueroles *et al.*, 2013), or colonization of a new niche (e.g., Bielawski *et al.*, 2004). The signal of frequency-dependent selection, which has been consistently connected to immune surveillance (e.g., Hughes and Nei, 1988) and reproductive conflict (e.g., Swanson *et al.*, 2003) to name two scenarios (see Yang and Bielawski, 2000, for a more comprehensive list of examples of this type), differs in that ω is elevated at some sites over much longer periods of evolutionary history. Frequency-dependent selection is consequently easier to detect.

Analytic CSMs (i.e., those fitted to data) are phenomenological in the sense that they summarize the “net resultants of selection” (Rodrigue and Philippe, 2010), with only limited consideration of the actual generating mechanisms. The same CSMs are frequently used to simulate alignments for the purpose of model testing (e.g., Anisimova *et al.*, 2001, 2002; Wong *et al.*, 2004; Zhang, 2004; Kosakovsky Pond and Frost, 2005; Yang *et al.*, 2005; Zhang *et al.*, 2005; Yang and dos Reis, 2011; Kosakovsky Pond *et al.*, 2011; Lu and Guindon, 2013), despite their lack of realism as a generating process. More mechanistically realistic parameter rich models, such as the mutation-selection (MS) framework (Halpern and Bruno, 1998) introduced in Chapter 1, can easily be used to simulate alignments. Sites can each be assigned their own fitness landscape under MS, and this can result in realistic variations in site-specific dynamics. An estimate of ω obtained by fitting a CSM to data generated under MS with static site-specific landscapes is expected to be consistent with purifying selection (i.e., $\omega < 1$). Spielman and Wilke (2015b) investigated this scenario and showed that the expected dN/dS at any site will always be less than one when fitness coefficients are fixed provided synonymous codons have equal fitness (Spielman and Wilke, 2015b). However, characterizing evolution on a static landscape using the long-run average rate ratio ignores temporal dynamics that can arise from population level processes of mutation and drift. In this chapter the MS framework is employed to show that a site-specific temporally dynamic dN/dS can result from non-adaptive shifting balance, a process whereby drift causes substitution to a less-than-optimal amino acid at a site and a combination of positive selection and drift subsequently causes

a rapid series of nonsynonymous substitutions that end when the site is re-occupied by its optimal amino acid. This process can be detected by CSMs designed to identify site-specific episodic positive selection, and estimates of ω obtained by these models can be significantly greater than one under some conditions. This suggests that temporal variations in ω generated by a non-adaptive process on a static fitness landscape can be misinterpreted as evidence of positive selection due to episodic selection and lead to the false conclusion that adaptive evolution had occurred.

2.1.1 Chapter Outline

The MS framework of Halpern and Bruno (1998) introduced in Chapter 1 was used to derive an expression for the expected rate ratio at a site, dN^h/dS^h (equation 1.22). This chapter starts with a discussion of how dN/dS can be defined under the MS framework in comparison with the CSM framework. The interpretation of the M-series CSMs of Yang *et al.* (2000a) as being implicitly designed to detect frequency-dependent selection is then validated. Although hinted at by other authors (e.g., Nielsen and Yang, 2003; Kryazhimskiy and Plotkin, 2008; Mugal *et al.*, 2014) a demonstration has never been published (but see dos Reis, 2013, unpublished). The demonstration presented herein helps to elucidate differences between the mechanistic MS framework and the standard phenomenological CSM approach. The notion of a site-specific MS landscape is then introduced, and a theoretical explanation for non-adaptive shifting balance is presented. Two methods of representing a site-specific landscape illustrate an interesting ramification of MS that has not been fully appreciated, namely that a site can be occupied by a suboptimal amino acid for long periods when selection is stringent.

Moving to the main point, the next section provides a tentative mechanistic model for non-adaptive shifting balance. This is used to show that site-specific variations in dN/dS are expected to be most pronounced when the substitution process is not dominated by selection or drift but admits interplay between the two. It is then shown that the covarion-like model CLM3($k = 2$)

introduced in Chapter 1 can detect temporal variations in ω generated by non-adaptive shifting balance when this interplay exists. It is also shown that both CLM3($k = 2$) and the branch-site model known as BUSTED (Murrell *et al.*, 2015) can sometimes detect positive selection due to non-adaptive shifting balance. These results suggest that the two models cannot distinguish adaptive changes in function (i.e., where amino acid fitness have changed) from temporal dynamics on static site-specific fitness landscapes. This is followed by a minor investigation using pairs of sequences generated under an MS model with peak shifts (dos Reis, 2015) to show that standard CSMs that assume a stationary process might overestimate branch lengths when fitted to data generated with nonstationary changes in fitness landscapes.

2.2 Results

2.2.1 Defining dN/dS under the Mutation-Selection Framework

Consider the standard genetic code shown in Table 1.1. Of all 61^2 possible codon pairs only 526 differ by a single nucleotide. Of these codon pairs 392 are nonsynonymous and only 134 are synonymous. Suppose a codon site is evolved *in silico* under the rate matrix $Q(\omega) = M \circ (\ell_S + \omega \ell_N)/r$ defined in equation (1.25) using stationary frequencies $\boldsymbol{\pi} = \langle \pi_1, \dots, \pi_{61} \rangle$ (e.g., as estimated from a real alignment). The expected rate at which nonsynonymous and synonymous substitutions would occur at the site can be calculated as follows:

$$rN = \sum_{(i,j)} \Pi Q(\omega) \circ \ell_N = \frac{1}{r} \sum_{(i,j)} \pi_i \omega M_{ij} \ell_N \quad (2.1)$$

$$rS = \sum_{(i,j)} \Pi Q(\omega) \circ \ell_S = \frac{1}{r} \sum_{(i,j)} \pi_i M_{ij} \ell_S \quad (2.2)$$

where Π is the diagonal matrix with entries $\boldsymbol{\pi}$ and the summation is over all of the elements in the 61×61 matrix argument. If dN/dS was equated to rN/rS then its value would be approximately $(392/134)\omega = 2.93\omega$ depending on $\boldsymbol{\pi}$. To account for the excess of possible nonsynonymous single nucleotide mutations, dN is instead defined as the expected nonsynonymous

substitution rate under the current selection regime (rN) divided by the expected nonsynonymous mutation rate¹ given by $rN_0 = \frac{1}{r} \sum_{(i,j)} \pi_i M_{ij} \ell_N$. dS is similarly defined as the expected synonymous substitution rate under the current selection regime (rS) divided by the expected synonymous mutation rate $rS_0 = \frac{1}{r} \sum_{(i,j)} \pi_i M_{ij} \ell_S$. By this normalization, $dN = rN/rN_0 = \omega$, $dS = rS/rS_0 = 1$ and $dN/dS = (rN/rN_0)/(rS/rS_0) = \omega$.

Under the MS framework the amino acids each have their own site-specific fitnesses that determine the vector of site-specific stationary frequencies $\boldsymbol{\pi}^h$. These are not the same as the stationary frequencies $\boldsymbol{\pi}$ for the mutation process defined by M , although the two are related by $\pi_i^h \propto \pi_i \exp(4N_e f_i^h)$ assuming a ploidy of two. Taking the differences in frequencies into consideration, dos Reis (2015) defined site-specific rN^h and rS^h under the MS framework as follows:

$$rN^h = \sum_{(i,j)} \Pi^h A^h \circ \ell_N \quad (2.3)$$

$$rS^h = \sum_{(i,j)} \Pi^h A^h \circ \ell_S \quad (2.4)$$

where Π^h is the diagonal matrix with entries $\boldsymbol{\pi}^h$ that satisfy $\boldsymbol{\pi}^h A^h = 0$, and the normalizing factors rN_0 and rS_0 as:

$$rN_0 = \sum_{(i,j)} \Pi M \circ \ell_N \quad (2.5)$$

$$rS_0 = \sum_{(i,j)} \Pi M \circ \ell_S \quad (2.6)$$

Under these definitions $dN^h = rN^h/rN_0$ is the expected nonsynonymous substitution rate at the h^{th} site divided by the expected nonsynonymous mutation rate at an unrelated site evolving under a neutral selection regime (i.e., with $s_{ij} = 0$ for all $i \neq j$) with stationary frequencies $\boldsymbol{\pi}$. $dS^h = rS^h/rS_0$ is similarly defined. By this approach the normalizing factor rS_0 has the “desirable property of being constant over sites” (dos Reis, 2015) similar to rS_0 under the CSM framework.

¹The matrix M defined by equation (1.24) is usually interpreted under the CSM framework as giving substitution rates under a neutral selection regime, a macroevolutionary process. However, the rate at which neutral substitutions occur over macroevolution time scales is equal to the rate at which neutral mutations arise over microevolutionary time scales (Kimura, 1983). For the purpose of the theoretical discussion in this section, M will therefore be construed as a matrix of mutation rates.

In my view it is more appropriate to define *site-specific* mutation rates as follows:

$$rN_0^h = \sum_{(i,j)} \Pi^h M \circ \ell_N \quad (2.7)$$

$$rS_0^h = \sum_{(i,j)} \Pi^h M \circ \ell_S \quad (2.8)$$

By this approach, $dN^h = rN^h/rN_0^h$ is the expected nonsynonymous substitution rate at the h^{th} site divided by the expected nonsynonymous mutation rate *at that same site*, taking into account the site's stationary frequencies $\boldsymbol{\pi}^h$. $dS^h = rS^h/rS_0^h$ is similarly defined. This gives the following normalized rates:

$$dN^h = \frac{rN^h}{rN_0^h} = \frac{\sum_{(i,j)} \pi_i^h A_{ij}^h \ell_N}{\sum_{(i,j)} \pi_i^h M_{ij} \ell_N} \times \frac{1}{2N_e}, \quad dS^h = \frac{rS^h}{rS_0^h} = \frac{1}{2N_e} \quad (2.9)$$

The site-specific rate ratio is therefore:

$$dN^h/dS^h = \frac{\sum_{(i,j)} \Pi^h A^h \circ \ell_N}{\sum_{(i,j)} \Pi^h M \circ \ell_N} = \frac{\sum_{(i,j)} \pi_i^h A_{ij}^h \ell_N}{\sum_{(i,j)} \pi_i^h M_{ij} \ell_N} \quad (2.10)$$

Equation (2.10) is consistent with the approach taken by Spielman and Wilke (2015b), and can be used to compute the theoretical long-run average rate ratio at a site evolving under A^h with site-specific fitnesses $\mathbf{f}^h = \langle f_1^h, \dots, f_{61}^h \rangle$. It is also consistent with the CSM approach in the sense that the same vector of frequencies is used for both the rate terms (rN^h and rS^h) and the normalization factors (rN_0^h and rS_0^h). The only difference is that $\boldsymbol{\pi}^h$ can vary across sites under the MS framework whereas $\boldsymbol{\pi}$ is the same at all sites under the CSM framework. Note that both frameworks make the assumption that the normalized synonymous mutation rates $dS = 1$ and $dS^h = 1/(2N_e)$ are uniform across sites and over time. See Rubinstein and Pupko (2012) for a review of sources of variation in dS and Davydov *et al.* (2019) for a discussion of the potential impact of the assumption of a constant dS on the false detection of positive selection.

2.2.2 M0 is Equivalent to a Model for Frequency-dependent Selection under the MS Framework

The rate ratio at a site evolving under $Q(\omega)$ is constantly ω . When viewed under the MS framework, this implies a site-specific fitness landscape that changes with each substitution. Hence, $\omega > 1$ under any CSM for which sites evolve under the same rate ratio over the entire tree can be interpreted as an indication of adaptive evolution by something akin to frequency-dependent selection. Here I demonstrate the veracity of this statement. For an alternative demonstration see dos Reis (2013).

$Q(\omega)$ characterizes the substitution process either for all sites (e.g., as it would under M0) or for some subset of sites (e.g., under M3 where sites are apportioned between several ω -categories). $Q(\omega)$ is in some ways similar to the site-specific rate matrix A^h in equation (1.4). In both cases the substitution rate is usually assumed to be zero for codons that differ by more than a single nucleotide substitution, and is proportional to M_{ij} when i and j are synonymous. The two rate matrices differ only in their treatment of nonsynonymous substitutions: $Q_{ij}(\omega) = \omega M_{ij}/r$ for all pairs of nonsynonymous codons, whereas under A^h the substitution rate between nonsynonymous codons can be different for each (i, j) pair depending on the selection coefficients s_{ij}^h .

Consider a variation of MS where: (i) the incumbent amino acid at a site has one fitness coefficient f^h while all others have fitness $f^h + \Delta f^h$; and (ii) when a substitution occurs, the incumbent and incoming amino acids swap fitnesses so that condition (i) still holds (Nielsen and Yang, 2003; Mugal *et al.*, 2014). The vector \mathbf{f}^h of site-specific fitness coefficients under assumptions (i) and (ii) is a time-dependent random variable that has no analogue in the CSM framework. Nevertheless, because the parameters of the process at a site do not change until a substitution occurs, Markov chain properties imply:

- (1) The probability that the codon i occupying a site is substituted by a codon j is proportional to M_{ij} if (i, j) are synonymous and $\omega^h M_{ij}$ if (i, j) are nonsynonymous, where $\omega^h = 4N_e \Delta f^h / (1 - \exp(-4N_e \Delta f^h))$.
- (2) For any codon i that occupies the site, the time until a substitution

occurs is an exponential random variable with mean $r_i = -1/A_{ii}^h$.

Significantly, (1) and (2) define a Markov process with rate matrix $A^h = Q(\omega^h)$ (Ross, 1996, Chapter 5). Note that the vector of fitness coefficients \mathbf{f}^h at a site is dynamic since it depends on the codon currently occupying the site. It can therefore be different from one site to the next at any instant. All sites that share the same Δf^h nevertheless evolve under the same phenomenological rate matrix $Q(\omega^h)$.

The equivalence of the rate matrix A^h for a MS process under (i) and (ii) to the rate matrix $Q(\omega^h)$ suggests that M-series models can be interpreted as being designed to detect signatures of frequency-dependent selection where, for instance, antagonistic interactions between proteins cause the fitness of any given variant to be inversely proportional to its frequency in the population. This interpretation makes sense only when $\Delta f^h > 0$ (i.e., $\omega^h > 1$), however, as was pointed out by dos Reis (2013). It is more appropriate to think of $Q(\omega^h)$ as a model for purifying selection when $\Delta f^h < 0$ ($\omega^h < 1$), or neutral selection when $\Delta f^h \approx 0$ ($\omega^h \approx 1$). Furthermore, even when $\Delta f^h > 0$, $Q(\omega^h)$ only captures the phenomenological effect of frequency-dependent selection, the sustained elevation in rate ratio to a value greater than one over a branch or lineage.

2.2.3 Shifting Balance on a static MS Landscape

Returning now to the general MS model, consider equation (1.4), repeated here for convenience:

$$A_{ij}^h \propto \begin{cases} M_{ij} & \text{if } s_{ij}^h = 0 \\ M_{ij} \frac{4s_{ij}^h}{1 - \exp(-4s_{ij}^h)} & \text{otherwise} \end{cases} \quad (2.11)$$

Spielman and Wilke (2015b) proved that positive selection under A^h as indicated by a long-run average $dN^h/dS^h > 1$ is not possible when fitness coefficients are fixed provided synonymous codons have equal fitness and the mutation process is symmetrical (i.e., $M_{ij} = M_{ji}$, but see Appendix for my own proof that relaxes this assumption). Here I propose a different interpretation that takes into account the temporal dynamics of mutation and drift on a static fitness landscape. The amino acid occupying the h^{th} site will vary

over time as long as at least two amino acids have non-negligible equilibrium frequencies. The expected proportion p_+^h of substitutions $i \rightarrow j$ that are beneficial with $s_{ij}^h > 0$ is:

$$p_+^h = \frac{\sum_{j \neq i} \pi_i^h (A_{ij}^h - M_{ij}) \ell_+}{\sum_{j \neq i} \pi_i^h A_{ij}^h} \quad (2.12)$$

where ℓ_+ is an indicator for $s_{ij}^h > 0$. The summation $\sum_{j \neq i} \pi_i^h A_{ij}^h \ell_+$ accounts for the rate at which beneficial substitutions are expected to occur, whereas $\sum_{j \neq i} \pi_i^h M_{ij} \ell_+$ accounts for the rate at which the same substitutions would be expected to occur under neutral selection (i.e., by drift). The difference therefore quantifies the rate at which beneficial substitutions are expected to occur above what would be expected by drift alone. The denominator $\sum_{j \neq i} \pi_i^h A_{ij}^h$ is the rate at which all types of substitutions are expected to occur.

If the nonsynonymous substitution $i \rightarrow j$ occurs then so does its reverse. One of the two must have a positive selection coefficient provided the site is not evolving under a strictly neutral regime. And since $s_{ij}^h > 0 \rightarrow A_{ij}^h > M_{ij}$, it is evident that p_+^h must be greater than zero unless the site happens to be fixed at one amino acid (i.e., when all alternative amino acids are lethal). Equation (2.12) therefore demonstrates that positive selection can occur on a static fitness landscape. Let p_-^h represent the proportion of deleterious substitutions $i \rightarrow j$ for which $s_{ij}^h < 0$. It can be shown that $p_+^h = p_-^h$ (see Appendix). Hence, beneficial substitutions due to positive selection can be thought of as “repairing” previously deleterious substitutions caused by drift (Sella and Hirsh, 2005; Mustonen and Lässig, 2009).

The dynamic implied by the balance $p_+^h = p_-^h$ can be illustrated using a site-specific MS landscape (cf. Bazykin, 2015), an analogue of the traditional fitness landscape constructed by sorting site-specific stationary frequencies from largest to smallest as depicted in Figure 2.1. There the frequencies were derived from fitness coefficients drawn from a normal distribution with a standard deviation $\sigma = 0.001$. The population size was assumed to be $N_e = 1000$. Positive selection can be seen to occur on this landscape by considering how dN^h/dS^h varies over time. The dN^h/dS^h ratio for the site depicted in Figure 2.1 is 0.58.

But this is a long-run average. The rate ratio in fact varies depending on the codon currently occupying the site. The codon-specific rate ratio dN_i^h/dS_i^h (where i is the codon currently occupying the site) can be computed from the i^{th} row of A^h as follows:

$$dN_i^h/dS_i^h = \frac{\sum_{j \neq i} A_{ij}^h \ell_N}{\sum_{j \neq i} M_{ij} \ell_N} \quad (2.13)$$

The line plot in Figure 2.1 shows dN_i^h/dS_i^h (scaled on the right y-axis) for each codon. When a codon with low fitness (one far to the right or “tail” of the MS landscape) occupies the site, the majority of mutations are “up-slope” with $s_{ij}^h > 0$. The codon-specific rate ratio (2.13) is consequently greater than $dN^h/dS^h = 0.58$ (as large as 4.10 in this example). As the site moves in the up-slope direction, the proportion of mutations that are further up-slope diminishes and dN_i^h/dS_i^h decreases to a value below 0.58 (as small as 0.24 in this example). By this process, chance substitutions (i.e., drift) that move a site down-slope are balanced by a combination of drift and positive selection ($dN_i^h/dS_i^h > 1$) that move the site back toward its peak. I call this dynamic process “non-adaptive shifting balance” because it is evocative of Wright’s theory² of the same name (Wright, 1932, 1982). A rate ratio $dN/dS > 1$ indicates fixation by positive selection by definition. Hence, the possibility of positive selection on a static fitness landscape is verified by the fact that dN_i^h/dS_i^h can be greater than one.

²Wright introduced shifting-balance theory to explain how a sub-population might move from one fitness peak across a fitness valley to another higher peak on a fixed landscape and subsequently cause the entire population to move to that new peak. The process is therefore adaptive. Non-adaptive shifting balance refers to the movement of an entire population away from and back to the same peak on a fixed landscape.

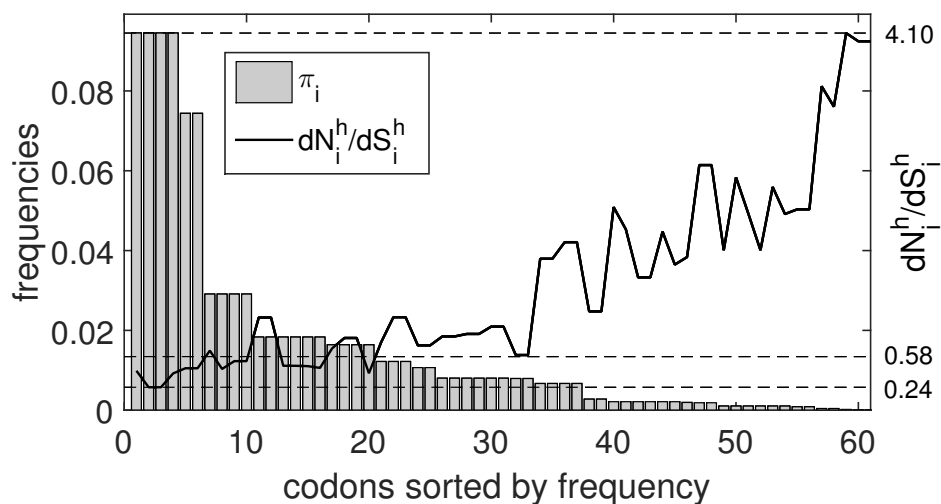


Figure 2.1: A site-specific fitness landscape. The bar plot depicts a MS landscape consisting of equilibrium frequencies sorted from largest to smallest. The line shows the codon-specific rate ratio dN_i^h/dS_i^h for the sorted codons. The rate ratio varies depending on the codon currently occupying the site, and can be greater than one following a chance substitution into the tail (to the right) of the landscape. In this case the codon specific rate ratio for the site ranges from 0.24 to 4.10 with a temporal average of $dN^h/dS^h = 0.58$.

2.2.4 Split MS Landscapes

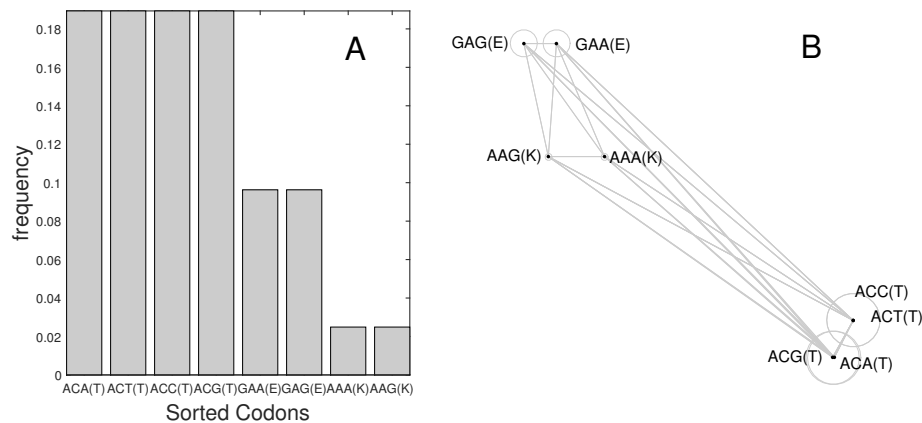


Figure 2.2: A MS and McCandlish landscape. The site depicted is under stringent selection pressure. A: The MS landscape shows that only three amino acids have non-negligible frequencies. B: The 2-dimensional McCandlish landscape provides information about the substitution dynamic. Vertices indicate codons; circle diameters are proportional to stationary frequencies; edge lengths are proportional to the expected number of single nucleotide substitutions required for the site to move from one vertex to the other and back again (McCandlish, 2011).

A landscape can be very sparse under a stringent selection regime, meaning that it is possible for the site-specific vector of stationary frequencies to be nearly zero for all but a few amino acids. A site can reside for non-negligible periods of time at suboptimal amino acids under this condition, especially when movement between viable amino acids via single-nucleotide steps requires substitution to one or more lethal residues. I call such cases “split landscapes”. A site can transition from a split landscape that limits substitutions between some pairs of codons to a broader landscape under which such substitutions can occur more readily following a decrease in population size.

Figure 2.1 is but one way to visualize site-specific dynamics. An alternative method makes use of the eigensystem of a transformation of the substitution probability matrix to produce a graphical representation of a site-specific landscape (following McCandlish, 2011, see Appendix). The vertices $i \in \{1, \dots, 61\}$ of the graph, which represent codons, are arranged in such a way that the length of any edge connecting i to j is approximately proportional to the expected time it will take for the site to move from i to j and back again to i . An example is shown in Figure 2.2. The bar plot in Figure 2.2 A shows the stationary frequencies for codons corresponding to amino acids threonine (T, with stationary frequency $\pi_T = 0.76$), glutamic acid (E, $\pi_E = 0.19$) and lysine (K, $\pi_K = 0.05$) for a site under stringent selection. All other amino acids have low fitness in this example, meaning that their stationary frequencies are essentially zero. The graph in Figure 2.2 B shows the relative location of each viable codon. Circles drawn with diameters proportional to the stationary frequency of the corresponding codon. The graph depicts a site that is occupied by T most of the time, with rapid substitutions between its four codon aliases ACA, ACC, ACG, ACT. If the site is currently occupied by ACG(T) it can move to AAG(K) via a single nucleotide substitution C→A in the second codon position. This would rarely occur, as indicated by the length of the edge connecting ACG(T) to AAG(K). But once it does occur, the system will tend to move between the more closely spaced E and K for some time before returning to T.

Figure 2.3 depicts the same site as in Figure 2.2 but with a 10-fold reduction

in N_e from 1000 to 100. Mutations with selection coefficients $s_{ij}^h = N_e(f_j^h - f_i^h) < 0$ that are rarely fixed when $N_e = 1000$ have a greater probability of being fixed by drift when $N_e = 100$. The reduction in N_e consequently allows the site to be occupied by a wider range of codons over time. This is evident in the resulting MS landscape shown in Figure 2.3 A, which is much broader than before reflecting an increase in frequencies that were previously negligible. Concordantly, the McCandlish landscape in Figure 2.3 B depicts a much larger network of connections between viable amino acids. Increasing the role of drift reduced the effect of selection at the site to the extent that it is now free to move rapidly between the two fittest amino acids T and E (third column, Table 2.1). Their stationary frequencies are now nearly the same ($\pi_T = 0.083, \pi_E = 0.078$). This demonstrates that site-specific codon frequencies, and the expected rate of nonsynonymous substitutions (Table 2.1), can change dramatically over time due to changes in population size even while the site's vector of fitness coefficients remains constant.

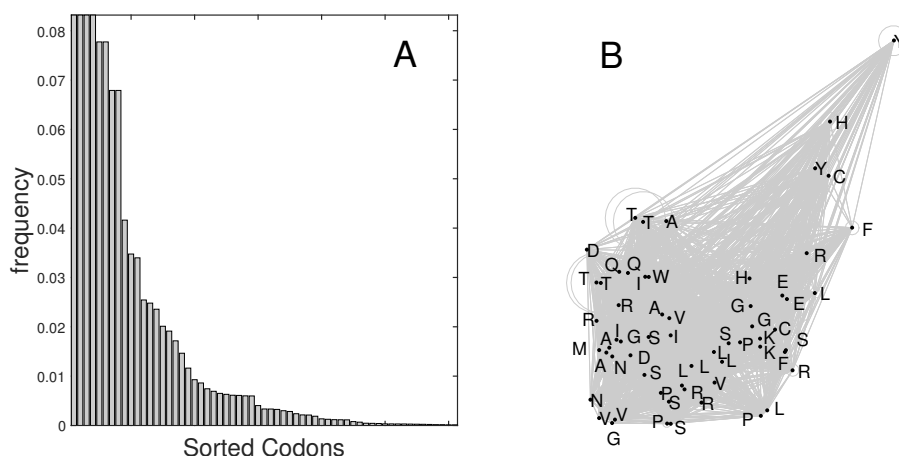


Figure 2.3: A landscape under relaxed selection pressure. The landscape is that depicted in Figure 2.2 after a 10-fold decrease in the population size. A: The MS landscape is now broader meaning that many codons have non-negligible equilibrium frequencies. B: The McCandlish landscape now depicts a site that is free to move between T and E via a large network of connections.

codon	$N_e = 1000$	$N_e = 100$
ACA(T)	0.026	0.20
ACT(T)	5.8×10^{-10}	0.072
ACC(T)	5.1×10^{-4}	0.35
ACG(T)	0.026	0.19
GAA(E)	0.15	0.38
GAG(E)	0.15	0.40
AAA(K)	0.74	0.55
AAG(K)	0.74	0.49

Table 2.1: The codon-specific rate ratios. The codon-specific rate ratios for each codon listed in the first column as a function of the population size N_e . The site-specific MS landscapes for $N_e = 1000$ and $N_e = 100$ are shown in Figures 2.2 and 2.3, respectively.

2.2.5 A Mechanistic Model for Non-adaptive Shifting Balance

In this section the MS framework is used to derive expressions for parameters that have meaningful interpretations in the context of CLM3($k = 2$). The purpose of this exercise is two-fold: first, to demonstrate that there is a mechanism by which a site can switch between two rate ratios on a static landscape; and second, to identify conditions under which such switches are expected to be most pronounced. Let ℓ_p^h be an indicator for codons i for which $dN_i^h/dS_i^h \leq 1$ (e.g., near the peak of the MS landscape), and let ℓ_t^h be the same for codons for which $dN_i^h/dS_i^h > 1$ (in the landscape's tail). A site will shift between its peak and tail, corresponding to switches between $\omega_1^h \leq 1$ and $\omega_2^h > 1$, with equilibrium proportions

$$p_1^h = \sum_i \pi_i^h \ell_p^h \quad (2.14)$$

$$p_2^h = 1 - p_1^h \quad (2.15)$$

The corresponding expected rate ratios can be computed using equation (2.10) by restricting the sum to either the peak or tail of the MS landscape:

$$\omega_1^h = \frac{\sum_{(i,j)} \frac{\pi_i^h}{p_1} A_{ij}^h \ell_N \ell_p^h}{\sum_{(i,j)} \frac{\pi_i^h}{p_1} M_{ij} \ell_N \ell_p^h}, \quad \omega_2^h = \frac{\sum_{(i,j)} \frac{\pi_i^h}{p_2} A_{ij}^h \ell_N \ell_t^h}{\sum_{(i,j)} \frac{\pi_i^h}{p_2} M_{ij} \ell_N \ell_t^h} \quad (2.16)$$

The expected number of switches between ω_1^h and ω_2^h per single nucleotide substitution is given by:

$$\delta^h = \frac{\sum_{(i,j)} \pi_i^h A_{ij}^h \ell_{\text{switch}}}{\sum_{j \neq i} \pi_i^h A_{ij}^h} \quad (2.17)$$

where ℓ_{switch} is an indicator for pairs of codons (i, j) for which one codon is in the peak and the other in the tail of the MS landscape. Since a switch can only occur upon a substitution, δ^h can be no greater than one.

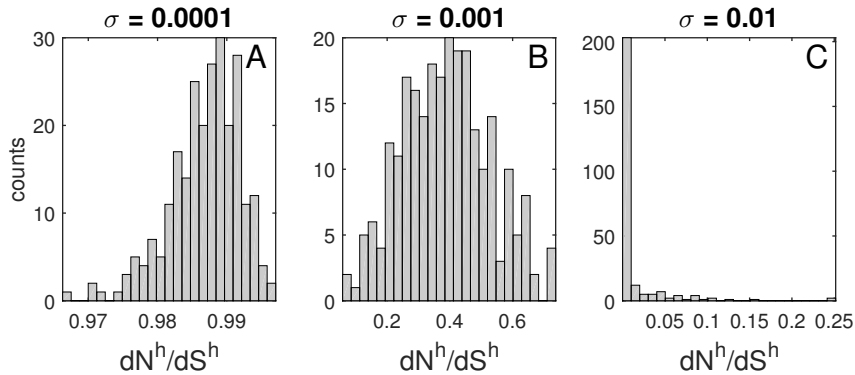


Figure 2.4: Distributions of site-specific rate ratios. The distributions of the site-specific rate ratio dN^h/dS^h for 250 sites under each of the three values of σ used in this study. A: $\sigma = 0.0001$ is consistent with nearly neutral evolution, as most sites evolve under a rate ratio close to one. B: the rate ratio varies over a broad range of values when $\sigma = 0.001$. C: $\sigma = 0.01$ is consistent with stringent selection for which most sites evolve under a rate ratio close to zero.

Sets of 250 vectors of site-specific fitness coefficients \mathbf{f}^h were drawn from a multivariate normal distribution centered at zero and with covariance $\sigma^2 I$ where I is the 61×61 identity matrix. Each was used to construct a site-specific rate matrix A^h after the fitnesses of each group of synonymous codons were adjusted to make them equal. The assumed mutation process was that described in equation (1.24) with $\kappa = 4$, and with $\alpha = \beta = 0$ to allow single nucleotide substitutions only. Uniform position-specific nucleotide frequencies were assumed. The site-specific rate ratio was calculated for each A^h using equation (2.10). Figure 2.4 shows how the distribution of dN^h/dS^h changes with σ . Figure 2.4 A demonstrates that $\sigma = 0.0001$ corresponds to a nearly neutral selection regime, as the site-specific rate ratio is very nearly one for most of the 250 sites. Figure 2.4 B shows that sites evolve over a wide range of site-specific rate ratios when $\sigma = 0.001$, where dN^h/dS^h ranges from 0.06 to 0.74 with a median of 0.39. And sites are mostly under stringent selection with $\omega^h \approx 0$ in Figure 2.4 C, where $\sigma = 0.01$. In the remainder of this section the differences between these three selection regimes are characterized in terms of the distributions of theoretical parameters $\langle \omega_1^h, \omega_2^h, p_2^h, \delta^h \rangle$.

Figure 2.5 shows box plots for parameter values computed from an additional draw of 250 vectors of fitness coefficients for each value of σ . First consider the case of nearly neutral evolution ($\sigma = 0.0001$). The median expected proportion of single nucleotide substitutions attributed to positive selection (p_+^h , Figure 2.5 A) is only 3.4%. The median probability that a site is in the tail of its MS landscape (p_2^h , Figure 2.5 B) is just under 43%. The median switching rate (δ^h , Figure 2.5 C) is 0.44, indicating approximately one switch for every two substitutions. The rate ratio from the tail (ω_2^h , Figure 2.5 D) is tightly distributed around a median of 1.1; ω_1^h (not shown) is similarly distributed but with a median of 0.91. In this scenario the population moves easily in and out of the tail, with only a small difference between ω_1^h and ω_2^h , because the landscape is nearly flat (i.e., with only slight variations in the π_i^h and dN_i^h/dS_i^h across codons).

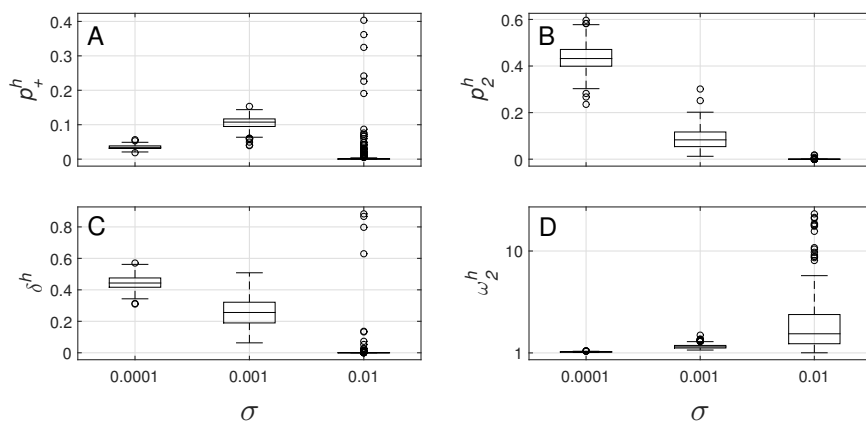


Figure 2.5: Distributions for the parameters of the mechanistic non-adaptive shifting balance model. Values were computed for 250 simulated sites. The span of each box represents the 50% of values that fall within the inter-quartile range; the midline shows the median value; whiskers show the range of the data excluding outliers, which are indicated by circles. A: p_+^h , the expected proportion of substitutions due to positive selection; B: p_2^h , the proportion of time a site is expected to be found in the tail of its MS landscape; C: δ^h , the expected number of switches between the peak and tail of its MS landscape per single nucleotide substitution; D: ω_2^h , the expected rate ratio for the site when in the tail of its MS landscape.

Next consider the case where a population is being tightly held to its fitness peak ($\sigma = 0.01$). The median value p_2^h is much less than one percent, reflecting a low probability of drift away from the peak. The rate ratio ω_2^h from the tail is

relatively large, with median 2.7. The median proportion of single nucleotide substitutions due to positive selection p_+^h and the median switching rate δ^h are both very low. This scenario is consistent with strong selective pressure that inhibits non-adaptive shifting balance by preventing movement away from the fittest amino acid.

Some parameters have outliers under the $\sigma = 0.01$ scenario. Outliers in ω_2^h tend to correspond to cases where p_2^h is very small. Among the 53 trials for which ω_2^h was greater than 5, for example, the median value of p_2^h was less than 2×10^{-10} . Such values indicate very strong selection pressure that prevents movement into the tail. Outliers in δ^h and p_+^h can be attributed to chance relationships in a reduced space of viable codons (i.e., a sparse landscape). An example is depicted in Figure 2.6, where codons were sorted by dN_i^h/dS_i^h rather than frequency so that the point of separation between peak and tail could be represented (i.e., where $dN_i^h/dS_i^h = 1$, vertical dashed line). Two amino acids M and I dominate the landscape; almost all substitutions are between them and consist of single nucleotide substitutions in the third position. Since M has only one codon alias, any substitution from the peak can only be nonsynonymous. And since a substitution from M across the $dN_i^h/dS_i^h = 1$ boundary to ATA(I) is about 25% more likely than a substitution to either ATT(I) or ATC(I) due to transition bias $\kappa = 4$ (Table 2.2), δ^h and p_+^h are both unusually large for the $\sigma = 0.01$ scenario ($\delta^h = 0.58$ and $p_+^h = 0.28$).

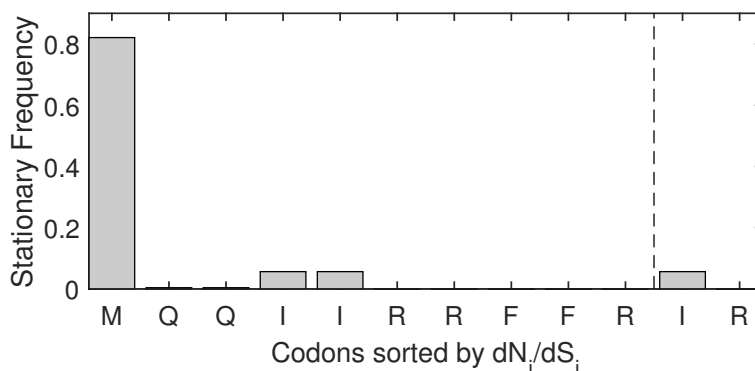


Figure 2.6: A MS landscape dominated by Methionine. The MS landscape depicted here was constructed by sorting codons by dN_i^h/dS_i^h . The vertical dashed line shows the point of separation between peak (where $dN_i^h/dS_i^h < 1$) and tail (where $dN_i^h/dS_i^h > 1$). The site is dominated by substitutions between M and the three codons for I. Although substitutions away from M are rare, when they do occur they are nonsynonymous since M has only one alias (ATG), and they are almost always to I. Codon aliases for I appear in order: ATT, ATC, and ATA. Substitutions to the right-most alias ATA(I) in the tail of the landscape are favored due to transition bias ($\kappa = 4$). As a result, both $\delta^h = 0.58$ and $p_+^h = 0.28$ are unusually large for this scenario.

Now consider the case between the nearly neutral and stringent selection scenarios, where $\sigma = 0.001$. Population dynamics on this landscape lead to a relatively large median value of p_+^h , which indicates that about 10% of single nucleotide substitutions are due to positive selection. The median switching rate is 0.26, close to one switch for every four substitutions. The median rate ratio from the tail is 1.5. Whereas the previous two scenarios represent extreme cases where one process strongly dominates (i.e., drift dominates when $\sigma = 0.0001$, and selection dominates when $\sigma = 0.01$), this scenario reflects an interplay between both processes. Here, the population occasionally moves away from its peak. But such events are quickly corrected because selection remains an effective force for moving the population back. This is the scenario that produces the strongest transient signature of positive selection on a fixed landscape.

	ATG(M)	ATT(I)	ATC(I)	ATA(I)
ATG(M)	0	0.30	0.30	0.39
ATT(I)	0.76	0	0.17	0.07
ATC(I)	0.76	0.17	0	0.07
ATA(I)	0.87	0.06	0.06	0

Table 2.2: Substitution Probabilities. Numbers give the probabilities that the incumbent codon in a row is next substituted by the codon in a column for the MS landscape depicted in Figure 2.6.

The way fitness coefficients are selected introduces a phenomenological component to the mutation-selection framework, since an assumed distribution is used in lieu of actual values. Although a few investigations suggest that the s_{ij}^h (and therefore the f_i^h) are sometimes consistent with a normal distribution (Nielsen and Yang, 2003; Tamuri *et al.*, 2012), there is no reason not to try alternatives. Figure 2.7 shows the result of a repeat of the experiment that generated the data in Figure 2.5 but with fitnesses drawn from exponential distributions with variance $\mu^2 \in \{1 \times 10^{-8}, 1 \times 10^{-6}, 1 \times 10^{-4}\}$. The patterns are similar: a lower variance corresponds to the nearly neutral scenario dominated by drift, a higher variance to the stringent scenario where selection dominates, and something in between to a balance between selection and drift under which shifting balance is strongest.

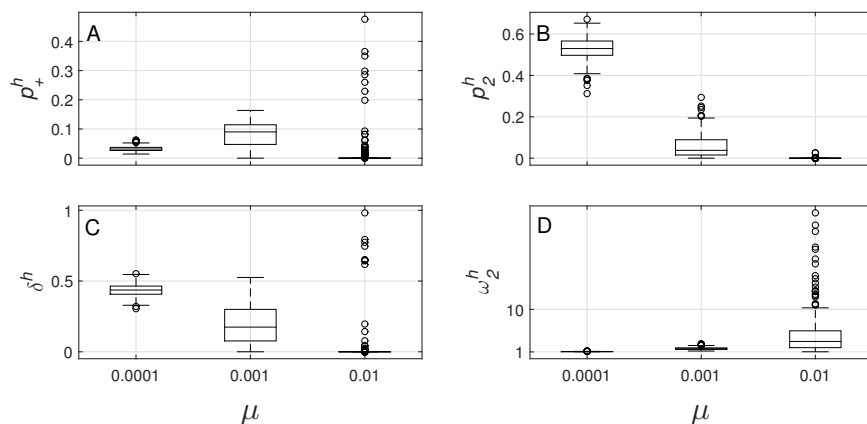


Figure 2.7: Alternate distributions for the parameters of the mechanistic non-adaptive shifting balance model. Fitness coefficients were drawn from exponential distributions. The mean μ (or variance, μ^2) of the exponential distribution plays a similar role here as σ did in Figure 2.5. With smaller values of μ , most fitness coefficients are similar to one another, all being close to zero. As μ becomes larger it becomes more likely that a few or one amino acid will draw a fitness coefficient much larger than the rest. As in Figure 2.5, three scenarios are indicated. When $\mu = 0.0001$ the rate ratio is always something very close to one, a site spends about half the time in the tail of its distribution, and switches about once every two substitutions. When $\mu = 0.01$ a site is typically held to its peak indicating stringent selection with $p_+^h \approx 0$ but with the exception of a small number of outliers. And when $\mu = 0.001$ the non-adaptive shifting balance phenomenon is more pronounced, with the median $p_+^h \approx 0.10$ as before.

2.2.6 Detecting Transient Changes in ω Caused by Non-adaptive Shifting Balance

In the previous section a mechanistic process by which a site can theoretically switch between two rate ratios as it moves over its fixed MS landscape was investigated. The objective in this section is to demonstrate that CLM3($k = 2$) can detect site-specific variations in rate ratio under certain conditions. To this end, alignments were generated on an 8-taxon symmetrical tree with branch lengths $b \in \{0.25, 0.5, 1\}$ using fitness coefficients with $\sigma \in \{0.0001, 0.001, 0.01\}$ and $N_e = 1000$ as described in Methods. For each scenario defined by (σ, b) the same set of 500 vectors of fitness coefficients was used to generate 50 unique alignments. CLM3($k = 2$) and M3($k = 2$) were fitted to each to provide a test for the significance of the switching rate δ . Table 2.3 shows the number of trials out of 50 for which the M3-CLM3 contrast rejected the null hypothesis and therefore detected switching. The test,

conducted at the 5% level of significance, seldom detected evidence for switching under the nearly neutral ($\sigma = 0.0001$) and stringent selection ($\sigma = 0.01$) scenarios, the exception being the $(\sigma, b) = (0.01, 1)$ scenario where the test was significant in 15/50 trials. Shifting was detected in all trials when $\sigma = 0.001$ and $b \in \{0.5, 1\}$, and in most trials when $\sigma = 0.001$ and $b = 0.25$. These results are in agreement with the mechanistic model that predicted that the scenario where neither drift nor selection dominate ($\sigma = 0.001$) would produce the strongest covarion-like signal due to non-adaptive shifting balance.

b/σ	0.0001	0.001	0.01
1.00	0 (0.51,1.2)	50 (0.00,0.77)	15 (0.00,0.06)
0.50	1 (0.75,1.2)	50 (0.12,0.68)	1 (0.00,0.08)
0.25	2 (0.69,1.1)	39 (0.05,0.79)	1 (0.00,0.08)

Table 2.3: Detecting Heterotachy. The left-most column gives the branch length and the top-most row the value of σ used to generate 50 alignments for the nine (b, σ) scenarios. Each cell shows the number of cases out of 50 for which the M3-CLM3 contrast detected site-specific switches between ω_1 and ω_2 . The numbers inside the brackets give the median MLEs for ω_1 and ω_2 .

Previous investigations indicated that a covarion-like model can detect switching even when data is generated without switching, and that this may occur when the number of ω -categories used to generate the data is greater than the number assumed by the fitted model (Lu and Guindon, 2013). The site specific rate ratio can vary greatly under the generating scenario with $\sigma = 0.001$, with values as small as $dN^h/dS^h = 0.06$ and as large as $dN^h/dS^h = 0.74$ for the 250 trials depicted in Figure 2.4 B. To rule out the possibility that this variation might produce false signatures of switching, an additional set of vectors of fitness coefficient was drawn with $\sigma = 0.001$. The rate ratio $\omega^h = dN^h/dS^h$ was computed for each vector. Each rate ratio was used to construct a site-specific phenomenological substitution rate matrix $Q(\omega^h)$. The resulting generating model was thus similar to an M-series model but with a different ω^h for each site. This model was used to generate fifty 500-codon alignments on a symmetrical 8-taxon rooted tree with all branch lengths $b = 1$. Since each site was evolved under its own rate matrix with $\omega^h = dN^h/dS^h$, the alignments had a similar distribution of rate ratios across sites as data generated under MS but without the covarion-like rate shifts that can occur

under MS. The M3-CLM3 contrast failed to reject the null at the 5% level of significance in all 50 trials, indicating no detectable switching. Hence, the rejection of the null under the $\sigma = 0.001$ scenario in Table 2.3 can be attributed to covarion-like rate shifts caused by non-adaptive shifting balance.

2.2.7 Detecting $\omega > 1$ Caused by Non-adaptive Shifting Balance

The results in Table 2.3 show that shifting balance can manifest as heterotachy under some circumstances. In this section I show that non-adaptive shifting balance can also manifest as episodic positive selection with $\hat{\omega}_2 > 1$. For this purpose, alignments simulated under MS were fitted to a new test for positive selection based on CLM3($k = 2$), as well as a popular analytical framework called the branch-site unrestricted statistical test for episodic diversification or BUSTED (Murrell *et al.*, 2015). In its original form, CLM3($k = 2$) allows sites to switch between $\omega_1 < \omega_2$ at a rate of δ switches per unit branch length. This framework can be used to construct a test for positive selection by placing restrictions on ω_1 and ω_2 . To that end, I define the null model CLM3a as CLM3($k = 2$) but restricted so that $\omega_1 < 1$ and $\omega_2 = 1$. Positive selection is therefore not permitted under CLM3a. Under the alternative model CLM3b the larger rate ratio ω_2 is estimated with the restriction that $\omega_2 \geq 1$. Hence, the CLM3a versus CLM3b contrast provides a likelihood ratio test for episodic positive selection. Under BUSTED, sites are assumed to switch randomly between three rate ratios over time (see Methods). Unlike CLM3, under which a site can switch between ω_1 and ω_2 multiple times along a branch, it is assumed under BUSTED that the rate ratio at a site is constant along any given branch, but can change from one branch to the next. The null hypothesis under BUSTED is that $\omega_0 \leq \omega_1 \leq \omega_2 = 1$ in contrast to the alternative for which $\omega_0 \leq \omega_1 \leq 1 \leq \omega_2$.

b/σ	0.0001	0.001	0.01
1.00	(1, 3)	(20, 11)	(10, 3)
0.50	(1, 2)	(20, 1)	(3, 1)
0.25	(2, 3)	(5, 0)	(0, 0)

Table 2.4: Detecting Positive Selection. The left-most column gives the branch length and the top-most row the value of σ used to generate 50 alignments for the nine (b, σ) scenarios. Each cell shows the number of cases (x, y) out of 50 for which positive selection was detected by BUSTED (x) and the CLM3a-CLM3b (y) contrast.

The null and alternative model for CLM3a vs CLM3b and BUSTED were fitted to the same alignments that were used to generate the data in Table 2.3. In each case, the test for positive selection was conducted only if $\hat{\omega}_2 > 1$ under the alternative. All tests were conducted at the 5% level of significance. Table 2.4 shows the number of trials in each scenario for which BUSTED and the CLM3a-CLM3b contrast found evidence of positive selection. Both models inferred positive selection in substantially more than 5% of the trials with $(\sigma, b) = (0.001, 1)$ (40% of trials under BUSTED and 22% under CLM3a vs CLM3b). BUSTED also detected positive selection in 40% of trials with $(\sigma, b) = (0.001, 0.5)$. Both models detected positive selection in data generated under the $\sigma = 0.0001$ and $\sigma = 0.01$ scenarios at a rate consistent with what would be expected by chance under each of their null models (i.e., close to 5% or 2 to 3 trials out of 50), except that BUSTED found signal in 20% of trials in the $(\sigma, b) = (0.01, 1)$ scenario. The distribution of the MLEs for ω_2 and p_2 among the trials for which $\hat{\omega}_2 > 1$ and the null was rejected are shown in Figure 2.8, with the exception of 13 trials for which the MLE for ω_2 under BUSTED was greater than 50. These results demonstrate that both models can detect the phenomenological signature of positive selection due to non-adaptive shifting balance under the $\sigma = 0.001$ scenario where neither selection nor drift dominate.

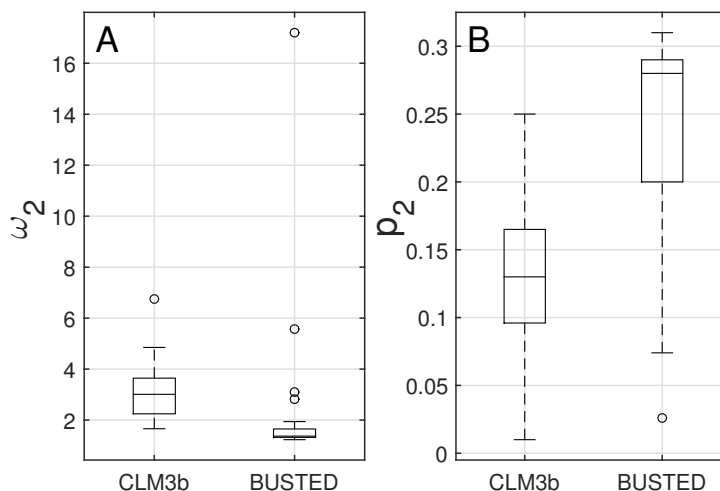


Figure 2.8: Distributions of MLEs estimated under CLM3 and BUSTED. Box plots show the distribution for the MLEs of A: ω_2 and B: p_2 estimated under the alternative model for BUSTED and under CLM3b for the trials where the null hypothesis of no positive selection was rejected under each model respectively. The plots for BUSTED do not show 13 trials for which $\hat{\omega}_2 > 50$.

2.2.8 Changing Fitness Landscapes

Up to this point it was shown that the substitution process at a site under the MS framework can be dynamic even if site-specific landscapes are fixed. An episodic shift in a fitness landscape can produce a similar dynamic that can be detected if such were to occur at a number of sites at the same time. The dynamic following a change in a fitness landscape was recently illustrated by dos Reis (2015) (also see Mustonen and Lässig, 2009). Under his environmental shift (MSES) model, A^h is the rate matrix defining the selection regime for the h^{th} site of an ancestral sequence. At $t = 0$ the regime switches to a different matrix B^h . This change initiates a non-stationary substitution process characterized by an elevated rate ratio $dN^h(t)/dS^h(t)$ that can be quantified in the following way (cf. equation 2.10):

$$dN^h(t)/dS^h(t) = \frac{\sum_{(i,j)} \Pi^h(t) B^h \circ \ell_N}{\sum_{(i,j)} \Pi^h(t) M \circ \ell_N} \quad (2.18)$$

where $\Pi^h(0)$ is the diagonal matrix with entries $\boldsymbol{\pi}^h(0) = \langle \pi_1^h(0), \dots, \pi_{61}^h(0) \rangle$ that give the stationary frequencies for the site consistent with A^h , and $\Pi^h(t)$ is the diagonal matrix with entries $\boldsymbol{\pi}^h(t) = \boldsymbol{\pi}^h(0) \exp(tB^h)$ that converges to the

new set of stationary frequencies consistent with B^h as $t \rightarrow \infty$. Consider what happens when M0 is fitted to pairs of sequences (S_1, S_2) generated by the non-stationary process that follows simultaneous changes in fitness coefficients at n codon sites. Modeling a non-stationary process as stationary can be thought of as a way of estimating an average effect. It is therefore reasonable to interpret estimates of ω under M0 as a mean taken across all codon sites (n) over the branch length (b)³. One possible way to formulate this under MSES is:

$$\bar{\omega}(b) = \frac{1}{b} \int_0^b \frac{1}{n} \sum_{h=1}^n dN^h(t)/dS^h(t) dt \quad (2.19)$$

This provides a means to predict the estimate of ω under M0.

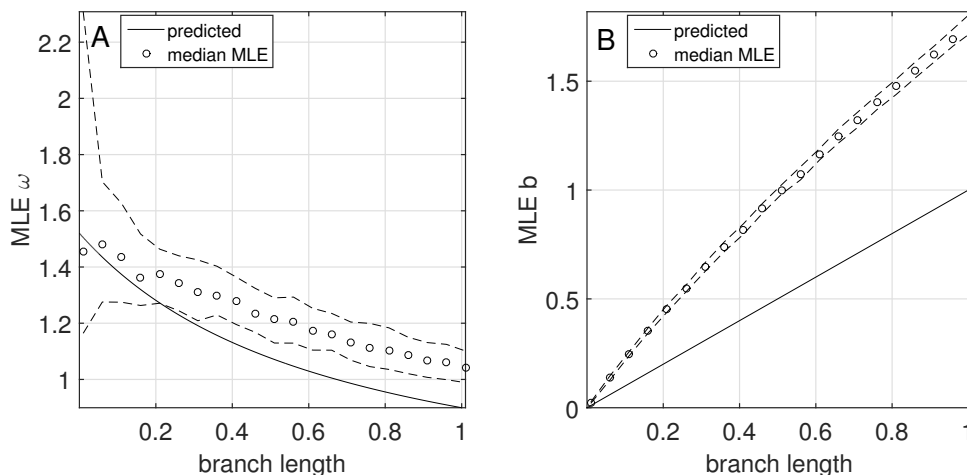


Figure 2.9: Investigation of the MSES model. Circles are median values of maximum likelihood estimates (MLE ω) produced by fitting data generated under MSES to M0. Dashed lines indicate the inter-quartile range. A: A comparison of the predicted versus estimated rate ratios computed from pairs of sequences generated under the MSES model. B: A comparison of generating and estimated branch lengths.

To compare predictions with MLEs, 200 pairs of sequences (S_1, S_2) with $n = 1000$ codon sites were generated under MSES with $\sigma = 0.001$ and $N_e = 1000$, and with branch lengths ranging between 0 and 1. Each site had its own pair of rate matrices A^h and B^h . The first sequence S_1 was generated by evolving each site of a random starting sequence under its assigned rate matrix A^h long enough to reach a codon near or at the peak of its site-specific fitness

³Analyses in Chapter 3 suggest that the complexity of the MSES generating model compared to the simplicity of M0 may well obfuscate this desirable interpretation.

landscape. The second sequence S_2 was then generated by evolving each site under its new rate matrix B^h along a branch of length b . M0 was subsequently fitted to each pair of sequences. Figure 2.9 A shows that the median MLE for ω estimated by M0 is highly correlated with the prediction $\bar{\omega}(b)$ ($\rho = 0.98$, p-value $\ll 0.0001$), decreases as b gets larger, but with a positive bias compared to predicted values, especially for longer branch lengths. This effect might account in part for the observed negative correlation between $\hat{\omega}$ and \hat{b} estimated from real pairs of sequences (dos Reis and Yang, 2013). Figure 2.9 B shows that branch lengths \hat{b} are consistently overestimated compared to the generating branch length b . Figure 2.9 A and B suggest that (i) elevations in site-specific rate ratios following peak shifts are less likely to be detected on longer branches; and (ii) failing to adjust for non-stationary processes following peak shifts can result in overestimation of branch lengths. These issues will be revisited in Chapter 4 where my phenotype-genotype branch-site model will be shown to mitigate (ii) by accounting for non-stationary processes in the form of site-specific peak shifts.

2.3 Discussion

The mutation-selection (MS) framework of Halpern and Bruno (1998) described in Chapter 1 provides a mechanistic description of the codon substitution process that is more realistic than that implied by commonly used phenomenological CSMs. CSMs implicitly assume that all amino acids have the same fitness save the one currently occupying the site, for example, as was demonstrated in Section 2.2.2. MS, by contrast, permits amino acids to have different fitnesses, and for these to vary across sites. The difference between the two approaches has many implications, several of which were explored in this chapter via theoretical arguments and computer simulations.

A site-specific MS landscape can in theory be split between two amino acids. Such landscapes can lead to patterns consistent with what is commonly called Type II functional divergence (FD) (Gu, 1999, 2001, 2006). Suppose a site were to evolve in two segregated populations under the site-specific landscape depicted in Figure 2.3 long enough for each population to be fixed

at a different amino acid. Further, suppose that each population were to subsequently undergo a 10-fold increase in population size, so that the site now evolves under the MS landscape depicted in Figure 2.2. It could happen that one population becomes fixed at T, and the other at E, depending in part on the starting codon for each population. Subsequent changes at other sites and/or in other genes might then canalize this difference as the populations diverge over macroevolutionary time scales (e.g., Pollock *et al.*, 2012). By this process, a site might eventually exhibit the constant-but-different pattern of Type II FD without any change in its fitness coefficients. To cite a real case, Gu (2006) identified sites in the COX gene that exhibited Type II FD in the form of physicochemical properties among the amino acids at a site that were similar within two clades (labeled COX1 and COX2) but radically different between the clades. This included a site for which T (categorized as polar) and E (charge-negative) dominated the COX1 and COX2 clades respectively. Although the apparent change in the physicochemical requirement of the site from polar to charge-negative might suggest an adaptive change in fitness coefficients, the observed pattern could also have arisen without adaptation under a split static MS landscape.

The mechanistic model for non-adaptive shifting balance in Section 2.2.5 indicated that sites are relatively free to move across their MS landscapes under the scenario where neither selection nor drift dominate. Box plots for p_2^h in Figure 2.5 B show that approximately 10% of sites can be expected to be in the tails of their respective landscapes at any instant when $\sigma = 0.001$ and $N_e = 1000$. These sites would be forced toward their fitness peaks all at the same time if a population evolving under this scenario were to undergo a rapid increase in effective size. This could result in a transient $\omega > 1$ signature of positive selection very similar to that following simultaneous peak shifts at a number of sites. The two processes might therefore be indistinguishable without accounting for changes in effective population size as part of the fitted CSM.

Setting aside the effect of changes in N_e , it has been commonly assumed that statistical evidence for $\omega > 1$ at some sites and/or branches in a tree is

indicative of positive selection due to episodic or continuous (e.g., frequency-dependent) changes in fitness coefficients. However, it was demonstrated in Section 2.2.5 that there are conditions under which a site evolving with fixed coefficients can undergo episodic positive selection by non-adaptive shifting balance. Furthermore, it was demonstrated in Section 2.2.7 that non-adaptive shifting balance can manifest as phenomenological switches between $\omega_1 \leq 1$ and $\omega_2 > 1$ that are detectable by commonly used branch-site models. Hence, it might not be possible to determine whether a site inferred to have undergone positive selection did so as a result of changes in fitness coefficients or non-adaptive shifting balance based on estimates of ω alone. Additional information about the role of the protein or the history of the organism would no doubt decide the issue in many cases. A protein implicated in an arms race, such as an immune surveillance protein in conflict with a pathogenic immune evasion protein, is very likely to have undergone changes in fitness coefficients at some sites (e.g., Hughes and Nei, 1988). So too for a protein that has been linked to variations in phenotype correlated with changes in habitat (e.g., Yokoyama *et al.*, 2008). In the absence of corroborating evidence of some kind however, positive selection by non-adaptive shifting balance might be the appropriate null hypothesis.

Although branch-site models commonly used to infer site-specific changes in the rate ratio can detect non-adaptive shifting balance under some conditions, M-series models that do not allow site-specific variations in ω are insensitive to this process. Spielman and Wilke (2015b) proved that the theoretical site-specific rate ratio dN^h/dS^h cannot exceed one when synonymous codons are equally fit. Whereas this proof applies to a single site, empirical results suggest that an equivalent statistical statement holds for the estimate of a single ω from an alignment with variations in dN^h/dS^h across sites. The likelihood ratio test for the contrast of M0 with $\omega = 1$ versus M0 with $\omega > 1$ was applied to alignments generated under each of the nine (σ, b) scenarios of Tables 2.3 and 2.4. The test never rejected the null at the 5% level of significance, and so ω was never inferred to be greater than one. It was shown in Section 2.2.2 that an M-series model is consistent with frequency-dependent

selection among sites evolving under $\omega > 1$. Taken together, these results suggests that it might be more appropriate to use M-series models when analyzing a gene suspected to have undergone frequency-dependent selection, if only to remove the possibility of detecting non-adaptive shifting balance and confusing it for episodic positive selection. However, evidence of a sustained elevation in ω over only a portion of a tree might easily be overlooked, resulting in reduced power.

The analysis of the MSES model in Section 2.2.8 underlines a potential difference between non-adaptive shifting balance and episodic changes in fitness landscapes. Non-adaptive shifting balance is a site-wise process, meaning that sites might be expected to undergo positive selection due to non-adaptive shifting balance randomly across sites and over time. Under the MSES model, by contrast, it is possible for a change in environment to impact a collection of sites all at the same time (e.g., sites that correspond to a functional domain or epitope). Hence, whereas M0 never detected $\omega > 1$ in the alignments generated with fixed fitness coefficients when signatures of positive selection were generated randomly across sites and over time, it was able to detect evidence of positive selection in alignments generated under MSES where positive selection occurred at all sites at the same time (Figure 2.9). By contrast, CLM3($k = 2$) and BUSTED detected $\omega > 1$ in a sizable proportion the alignments generated with fixed fitness coefficients. The apparent inability of models such as CLM3($k = 2$) and BUSTED to discriminate between episodic positive selection caused by adaptive versus non-adaptive processes might be addressed by introducing dependencies between sites to account for simultaneous changes in landscapes. This approach is investigated in detail in Chapter 4.

Non-adaptive shifting balance has implications beyond codon substitution models. For example, consider the method of estimating the proportion α of amino acid substitutions attributed to positive selection. Estimation of α is based on a comparison of the fixation ratio D_n/D_s , the number of observed nonsynonymous (D_n) and synonymous (D_s) differences between two closely related species, with the mutation ratio P_n/P_s , the total number of nonsynonymous (P_n) and synonymous (P_s) polymorphisms within each species (Smith

and Eyre-Walker, 2002). Neutrality is implied when $P_n/P_s = D_n/D_s$ making $\alpha = 1 - (P_n/P_s)/(D_n/D_s) = 0$, whereas an excess of nonsynonymous substitutions D_n above that expected under neutrality implies that some nonsynonymous mutations were fixed by positive selection, making $\alpha > 0$. Estimates of α range between about 10% for humans to more than 50% for *Drosophila* (Grossmann *et al.*, 2014, and references therein). While substitution by positive selection is often taken as an indication of adaptive evolution, the analysis in Section 2.2.5 suggest that as much as 10% of substitutions by positive selection can be attributed to non-adaptive shifting balance on a static fitness landscape. Thus, it seems possible that the human genome might not be evolving in response to changes in selection pressure, but merely experiencing non-adaptive shifting balance. Likewise, some fraction of positive selection in *Drosophila* might be due to the same process. The key question is “How prevalent is non-adaptive shifting balance in real data?”.

2.4 Methods

2.4.1 Alignment Generation

Alignments were generated using the MS framework as follows. Mutations were modeled using equation (1.24) with $\kappa = 4$, uniform nucleotide frequencies, and $\alpha = \beta = 0$ to allow single nucleotide substitutions only. Vectors of site-specific fitness coefficients \mathbf{f}^h were drawn from a zero-mean multivariate normal distribution with covariance matrix $\sigma^2 I$, where I is the 61×61 identity matrix. Each vector was modified to make synonymous substitutions equally fit before using equation (2.11) to construct a site-specific rate matrix A^h . All A^h were divided by the mean rate $\bar{r} = (1/n) \sum_{h=1}^n \sum_{j \neq i} \pi_i^h A_{ij}^h$ to make branch lengths interpretable as the expected number of single nucleotide substitutions per codon. The population size was set to $N_e = 1000$ for all simulations. All alignments were generated on a symmetrical eight-taxa rooted tree with uniform branch lengths $b \in \{0.25, 0.50, 1.00\}$ except where otherwise indicated. Variations in the strength of the non-adaptive shifting balance phenomenon were effected by using values of $\sigma \in \{0.0001, 0.001, 0.01\}$.

2.4.2 The M3($k = 2$) vs CLM3($k = 2$) Contrast

CLM3($k = 2$) is equivalent to M3($k = 2$) when $\delta = 0$, on the boundary of its parameter space, $\delta \in [0, 1]$. The M3($k = 2$) vs CLM3($k = 2$) contrast therefore provides a likelihood ratio test for heterotachy ($\delta > 0$). The limiting distribution for the LLR for the M3($k = 2$) vs CLM3($k = 2$) contrast is an equal mixture of a point-mass at zero and the χ_1^2 distribution (e.g., Case 5 in Self and Liang, 1987), with a critical value of 2.71 at the 5% level of significance. The variant CLM2b with $\omega_2 > 1$ is equivalent to CLM3a when ω_2 is fixed at one. The CLM3a vs CLM3b contrast therefore provides a test for episodic positive selection ($\omega_2 > 1$). The theoretical limiting distribution for the LLR for this contrast is χ_1^2 , with a critical value of 3.84 at the 5% level of significance.

2.4.3 BUSTED

The branch-site unrestricted statistical test for episodic diversification or BUSTED (Murrell *et al.*, 2015) assumes that each site evolved under one of three rate ratios $\{\omega_0, \omega_1, \omega_2\}$ along each branch of the tree selected randomly in proportions $\{p_0, p_1, 1 - p_0 - p_1\}$. The unrestrained model allows each of the three rate ratios to take on any non-negative value. If the largest estimated rate ratio is greater than one then a test for positive selection contrasting the model with $\omega_0 \leq \omega_1 \leq \omega_2 = 1$ versus the model with $\omega_0 \leq \omega_1 \leq 1 \leq \omega_2$ is conducted. The distribution of the LLR for this test is an unknown mixture of χ_0^2 , χ_1^2 and χ_2^2 . To be conservative, BUSTED uses χ_2^2 to compute p-values. BUSTED differs from CLM3 in that sites switch between three rate ratios instead of two, and switching is from branch-to-branch instead of at any point along any branch.

2.5 Appendix

2.5.1 A Demonstration of $p_+^h = p_-^h$

Here I establish the expected balance between beneficial and detrimental substitutions on a static site-specific MS fitness landscape. Consider the expected proportion p_+^h of substitutions that are beneficial with $s_{ij}^h > 0$:

$$p_+^h = \frac{\sum_{(i,j)} \pi_i^h A_{ij}^h \ell_+}{\sum_{j \neq i} \pi_i^h A_{ij}^h} - \frac{\sum_{(i,j)} \pi_i^h M_{ij}^h \ell_+}{\sum_{j \neq i} \pi_i^h A_{ij}^h} \quad (2.20)$$

where ℓ_+ is an indicator for $s_{ij}^h > 0$. The first addend of p_+^h accounts for the rate at which beneficial substitutions occur. The second addend accounts for the rate at which the same substitutions would be expected to occur if they were neutral (i.e., if $s_{ij}^h = 0$). Since $s_{ij}^h > 0 \leftrightarrow s_{ji}^h < 0$, it follows that the counterpart expression for the proportion of substitutions that are deleterious is:

$$p_-^h = \frac{\sum_{(i,j)} \pi_j^h A_{ji}^h \ell_+}{\sum_{j \neq i} \pi_j^h A_{ji}^h} - \frac{\sum_{(i,j)} \pi_j^h M_{ji}^h \ell_+}{\sum_{j \neq i} \pi_j^h A_{ji}^h} \quad (2.21)$$

Since every term in p_+^h has a counterpart in p_-^h , the difference between the two can be written as:

$$p_+^h - p_-^h = \frac{\sum_{(i,j)} (\pi_i^h A_{ij}^h - \pi_j^h A_{ji}^h) \ell_+}{\sum_{j \neq i} \pi_i^h A_{ij}^h} + \frac{\sum_{(i,j)} (\pi_i^h M_{ij}^h - \pi_j^h M_{ji}^h) \ell_+}{\sum_{j \neq i} \pi_i^h A_{ij}^h} \quad (2.22)$$

Since both M and A^h are time reversible, it follows from the detailed balance equation that $\pi_i^h M_{ij}^h - \pi_j^h M_{ji}^h = 0$ and $\pi_i^h A_{ij}^h - \pi_j^h A_{ji}^h = 0$. Hence, $p_+^h = p_-^h$.

2.5.2 $dN^h/dS^h \leq 1$ When \mathbf{f}^h is Fixed

Equation (2.10) gives the expected or long-run average rate ratio at a site. An important property of dN^h/dS^h is that it is bounded above by one. Following Spielman and Wilke (2015b), this property can be demonstrated by proving that the substitution rate between any nonsynonymous (i, j) pair is bounded by the mutation rate:

$$\pi_i^h A_{ij}^h + \pi_j^h A_{ji}^h \leq \pi_i^h M_{ij}^h + \pi_j^h M_{ji}^h \quad (2.23)$$

Simplifying:

$$2\pi_i^h A_{ij}^h \leq \pi_i^h M_{ij} + \pi_j^h M_{ji}$$

since the process is time-reversible: $\pi_i^h A_{ij}^h = \pi_j^h A_{ji}^h$

$$2A_{ij}^h \leq M_{ij} + \frac{\pi_j^h}{\pi_i^h} M_{ji}$$

$$2A_{ij}^h \leq M_{ij} + M_{ij} \exp(4s_{ij}^h)$$

since $\pi_j^h/\pi_i^h = (M_{ij} \exp(4N_e f_j^h))/(M_{ji} \exp(4N_e f_i^h))$

$$M_{ij} \frac{4s_{ij}^h}{1 - \exp(-4s_{ij}^h)} \leq M_{ij} (1 + \exp(4s_{ij}^h))$$

Notice that the term M_{ij} can be divided out, meaning that the proof does not require that $M_{ij} = M_{ji}$ (as was assumed by Spielman and Wilke, 2015b). If we let $x = 4s_{ij}^h$, the proof is reduced to showing that:

$$\frac{x}{1 - e^{-x}} \leq 1 + e^x \quad (2.24)$$

To prove (2.24), first note that $e^x - e^{-x} = \sum_{n=0}^{\infty} 2x^{2n+1}/(2n+1)!$. When $x > 0$, multiplying both sides of (2.24) by $1 - e^{-x}$ gives $x \leq e^x - e^{-x}$. This inequality is true because:

$$e^x - e^{-x} = 2x + \sum_{n=1}^{\infty} 2x^{2n+1}/(2n+1)! > x \quad (2.25)$$

When $x < 0$, multiplying both sides of (2.24) by $1 - e^{-x}$ gives $x \geq e^x - e^{-x}$ or $-x \leq -(e^x - e^{-x})$. Since both x and $e^x - e^{-x}$ are odd functions (i.e., $-f(x) = f(-x)$) and since $x < 0$, the inequality is equivalent to $x \leq (e^x - e^{-x})$ with $x > 0$, which was already shown to be true. Hence the inequality (2.24) holds for all $x = 2s_{ij}^h$. It follows that $dN^h/dS^h \leq 1$ since equation (2.23) applies to all nonsynonymous (i, j) pairs.

2.5.3 A Demonstration that A^h is Time-Reversible

The matrix M was defined in Chapter 1 as follows (equation 1.24):

$$M_{ij} \propto \begin{cases} \kappa^{st} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 1 \\ \alpha \kappa^{st} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 2 \\ \beta \kappa^{st} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 3 \end{cases} \quad (2.26)$$

Recall that equation (2.26) applies to any pair of codons (i, j) that differ by s nucleotides s_t of which are transversions, and that the π_{jk}^* are position-specific nucleotide frequencies, κ the transition bias, and α and β rate parameters that account for pairs of codons that differ by $s = 2$ or $s = 3$ nucleotides. When M is used to construct a phenomenological rate matrix $Q(\omega) \propto M \circ (\ell_S + \omega \ell_N)$ it is interpreted as quantifying the rate at which substitutions occur under a neutral selection regime. When M is used to construct a site-specific rate matrix A^h (as in equation 2.11) it is interpreted as quantifying the rate at which mutations arise. Only single nucleotide differences were permitted under both usages for all analyses presented in this chapter (i.e., α and β were always assumed to be zero). Nevertheless, the following demonstrates that M as defined in equation (2.26) is time-reversible.

Let $i = n_{i_1} n_{i_2} n_{i_3}$ and $j = n_{j_1} n_{j_2} n_{j_3}$ represent an arbitrary pair of nucleotide triplets each of which occurs with stationary frequency $\pi_i^0 = \pi_{i_1}^* \pi_{i_2}^* \pi_{i_3}^* / r$ and $\pi_j^0 = \pi_{j_1}^* \pi_{j_2}^* \pi_{j_3}^* / r$, where $r = \sum_{i=1}^{61} \pi_{i_1}^* \pi_{i_2}^* \pi_{i_3}^*$ is a common normalization constant. M is time-reversible if the following holds:

$$\pi_i^0 M_{ij} = \pi_{i_1}^* \pi_{i_2}^* \pi_{i_3}^* (c_{ij}/r) \prod_{i_k \neq j_k} \pi_{j_k}^* = \pi_{j_1}^* \pi_{j_2}^* \pi_{j_3}^* (c_{ji}/r) \prod_{i_k \neq j_k} \pi_{i_k}^* = \pi_j^0 M_{ji} \quad (2.27)$$

where $c_{ij} = c_{ji}$ is a coefficient composed of the appropriate combination of κ , α and β . The proof therefore requires a demonstration that:

$$\pi_{i_1}^* \pi_{i_2}^* \pi_{i_3}^* \prod_{i_k \neq j_k} \pi_{j_k}^* = \pi_{j_1}^* \pi_{j_2}^* \pi_{j_3}^* \prod_{i_k \neq j_k} \pi_{i_k}^* \quad (2.28)$$

The products $\prod_{i_k \neq j_k} \pi_{i_k}^*$ and $\prod_{i_k \neq j_k} \pi_{j_k}^*$ are composed of frequencies at positions where i and j differ, whereas the products $\pi_{i_1}^* \pi_{i_2}^* \pi_{i_3}^*$ and $\pi_{j_1}^* \pi_{j_2}^* \pi_{j_3}^*$ may contain common factors. The removal of common factors from $\pi_{i_1}^* \pi_{i_2}^* \pi_{i_3}^*$ leaves only $\prod_{i_k \neq j_k} \pi_{i_k}^*$. Likewise, the removal of the same factors from $\pi_{j_1}^* \pi_{j_2}^* \pi_{j_3}^*$ leaves only $\prod_{i_k \neq j_k} \pi_{j_k}^*$. Equation (2.28) can therefore be reduced to:

$$\prod_{i_k \neq j_k} \pi_{i_k}^* \prod_{i_k \neq j_k} \pi_{j_k}^* = \prod_{i_k \neq j_k} \pi_{j_k}^* \prod_{i_k \neq j_k} \pi_{i_k}^* \quad (2.29)$$

which is evidently true.

Next I demonstrate that $\pi_i^h \propto \pi_i^0 e^{4N_e f_i^h}$. Following Wang *et al.* (2014), let p_{ij}^h be the probability that the $i \rightarrow j$ mutation is fixed, and note that:

$$\frac{p_{ij}^h}{p_{ji}^h} = \frac{4s_{ij}^h}{4s_{ji}^h} \frac{1 - e^{-4s_{ji}^h}}{1 - e^{-4s_{ij}^h}} = \frac{1 - e^{4s_{ij}^h}}{1 - \frac{1}{e^{4s_{ij}^h}}} = e^{4s_{ij}^h} \quad (2.30)$$

Given equation 2.30 and the fact that M is time reversible, it follows that:

$$\frac{A_{ij}^h}{A_{ji}^h} = \frac{\pi_j^0 \pi_i^0 M_{ij} p_{ij}^h}{\pi_i^0 \pi_j^0 M_{ji} p_{ji}^h} = \frac{\pi_j^0}{\pi_i^0} e^{4s_{ij}^h} = \frac{\pi_j^0 e^{4N_e f_j^h}}{\pi_i^0 e^{4N_e f_i^h}} \quad (2.31)$$

Now recall that $\pi^h A^h = 0$ and that the rows of A^h sum to 0:

$$\begin{aligned} \sum_{k=1}^{61} \pi_k^h A_{kj}^h &= \pi_i^h A_{ij}^h + \sum_{k \neq i} \pi_k^h A_{kj}^h = 0 \\ \sum_{k=1}^{61} A_{jk}^h &= A_{ji}^h + \sum_{k \neq i} A_{jk}^h = 0 \end{aligned} \quad (2.32)$$

Combining the A_{ij}^h/A_{ji}^h ratio in (2.31) with a similar ratio constructed from the two sums in (2.32) yields:

$$\pi_i^h \frac{A_{ij}^h}{A_{ji}^h} = \pi_i^h \frac{\pi_j^0 e^{4N_e f_j^h}}{\pi_i^0 e^{4N_e f_i^h}} = \frac{-\sum_{k \neq i} \pi_k^h A_{kj}^h}{-\sum_{k \neq i} A_{jk}^h} \quad (2.33)$$

Solving for π_i^h yields the required relation:

$$\pi_i^h = \pi_i^0 \left(\frac{1}{\pi_j^0 e^{4N_e f_j^h}} \frac{\sum_{k \neq i} \pi_k^h A_{kj}^h}{\sum_{k \neq i} A_{jk}^h} \right) e^{4N_e f_i^h} \propto \pi_i^0 e^{4N_e f_i^h} \quad (2.34)$$

The constraint that $\sum_{i=1}^{61} \pi_i^h = 1$ implies the proportionality constant $c = \sum_{i=1}^{61} \pi_i^0 e^{4N_e f_i^h}$. Returning to equation (2.31):

$$\left(\frac{\pi_i^h}{\pi_j^h} \right) \frac{A_{ij}^h}{A_{ji}^h} = \left(\frac{\pi_i^h}{\pi_j^h} \right) \frac{\pi_j^0 e^{4N_e f_j^h}}{\pi_i^0 e^{4N_e f_i^h}} = \left(\frac{\pi_i^0 e^{4N_e f_i^h} / c}{\pi_j^0 e^{4N_e f_j^h} / c} \right) \frac{\pi_j^0 e^{4N_e f_j^h}}{\pi_i^0 e^{4N_e f_i^h}} = 1 \quad (2.35)$$

It follows that A^h is time-reversible.

2.5.4 Visualizing Substitution Dynamics

In section 1.1.6 it was shown that $P^h(t) = \exp(tA^h)$ can be expressed in terms of the eigensystem (U, Λ) of A^h :

$$P^h(t) = U \exp(t\Lambda) U^{-1} \quad (2.36)$$

It follows that $P^h(t)$ has the same eigenvectors as A^h but with eigenvalues $0, e^{t\lambda_2}, \dots, e^{t\lambda_{61}}$. Let $0 = \lambda_1 > \lambda_2 \geq \dots \geq \lambda_{61}$ be the eigenvalues of A^h sorted in descending order and let $\mathbf{1}/\sqrt{61}, \mathbf{u}_2, \dots, \mathbf{u}_{61}$ be the corresponding eigenvectors. It is convenient to write $P^h(t)$ as a sum of matrices:

$$P^h(t) = \sum_{k=1}^{61} e^{t\lambda_k} \mathbf{u}_k \mathbf{u}_k^T = \mathbf{1}\mathbf{1}^T \Pi^h + \sum_{k=2}^{61} e^{t\lambda_k} \mathbf{u}_k \mathbf{u}_k^T \quad (2.37)$$

where $\mathbf{1}\mathbf{1}^T\Pi^h$ is a matrix whose rows are all equal to $\boldsymbol{\pi}^h$. Suppose i is the codon occupying the h^{th} site of a gene in all members of a population as represented by a 1×61 row vector $\mathbf{v}^h(0)$ consisting of zeros but with one in the i^{th} position. The probability distribution for the codons that will be fixed at the site during the interval $(0, t)$ is given by $\mathbf{v}^h(t) = \mathbf{v}^h(0)P^h(t)$:

$$\mathbf{v}^h(t) = \mathbf{v}^h(0)\mathbf{1}\mathbf{1}^T\Pi^h + \mathbf{v}^h(0) \sum_{k=2}^{61} e^{t\lambda_k} \mathbf{u}_k \mathbf{u}_k^T = \boldsymbol{\pi}^h + \sum_{k=2}^{61} e^{t\lambda_k} \mathbf{u}_{ik} \mathbf{u}_k^T \quad (2.38)$$

Since $\lambda_k < 0$ for $k \in \{2, \dots, 61\}$ it follows that $\lim_{t \rightarrow \infty} \mathbf{v}^h(t) = \boldsymbol{\pi}^h$ (cf. section 1.1.6). Equation (2.38) characterizes the transient dynamic of the evolution of the h^{th} site as the site moves back to its optimal amino acid following fixation to a suboptimal codon i in terms of the row vectors \mathbf{u}_k^T and weights $e^{t\lambda_k} \mathbf{u}_{ik}$ that diminish over time. This is the essential idea behind the McCandlish landscape but for details surrounding the way the landscape is actually constructed. Those details start with a Laplacian matrix as described in the next paragraph.

Following Koren (2005) let $G = (V, E, W)$ be any weighted graph, where V is a set of n vertices, E a set of edges and W a set of weights. Each $w_{ij} \in W$ measures the similarity of the two vertices connected by edge $(i, j) \in E$, meaning that $w_{ij} = w_{ji}$. Each vertex $i \in V$ has a set of neighbourhoods defined as $N(i) = \{j \mid (i, j) \in E\}$. The degree of each vertex is given by $\text{deg}(i) = \sum_{j \in N(i)} w_{ij}$. Assuming that every pair of distinct vertices is connected by an edge, the elements of the Laplacian L of G are defined as follows:

$$L_{ij} = \begin{cases} \text{deg}(i) & \text{if } i = j \\ -w_{ij} & \text{if } i \neq j \end{cases} \quad (2.39)$$

Let $\mathbf{x}^T = \langle x_1, \dots, x_n \rangle$ and $\mathbf{y}^T = \langle y_1, \dots, y_n \rangle$ be any pair of $n \times 1$ vectors whose elements (x_i, y_i) give the location of the i^{th} vertex in 2-dimensions. Suppose the objective is to represent G in 2-dimensions in such a way that pairs of vertices that are more similar (with larger w_{ij}) tend to be closer together than pairs of vertices that are less similar (with smaller w_{ij}). Following the criterion specified by Koren (2005), the desired two-dimensional representation of G is given by the pair (\mathbf{x}, \mathbf{y}) that minimizes the weighted sum of squared edge

lengths:

$$\mathbf{x}^T L \mathbf{x} + \mathbf{y}^T L \mathbf{y} = \sum_{(i,j) \in E} w_{ij} [(x_i - x_j)^2 + (y_i - y_j)^2] \quad (2.40)$$

subject to the constraint that (\mathbf{x}, \mathbf{y}) are orthonormal. The Laplacian is symmetric and so has n real eigenvalues $0 = \lambda_1 < \lambda_2 \leq \dots \leq \lambda_n$ with corresponding eigenvectors $\mathbf{1}/\sqrt{n}, \mathbf{q}_2, \dots, \mathbf{q}_n$. It can be shown that the second and third eigenvectors (corresponding to smallest two non-zero eigenvalues) satisfy the minimization problem (Koren, 2005):

$$(\mathbf{q}_2, \mathbf{q}_3) = \arg \min \{ \mathbf{x} L \mathbf{x}^T + \mathbf{y} L \mathbf{y}^T \} \quad (2.41)$$

The graph G can therefore be depicted in 2-dimensions by placing vertices at points $\{(\mathbf{q}_{i2}, \mathbf{q}_{i3}), i = 1, \dots, n\}$.

This result was used by McCandlish (2011) to construct his site-specific landscapes as follows. Let $\mathbf{q}_1, \dots, \mathbf{q}_{61}$ be the 61×1 orthonormal eigenvectors for the transformation $D^{1/2} P^h(t) D^{-1/2}$ of $P^h(t)$ with corresponding eigenvalues $1 = \lambda_1 > \lambda_2 \geq \dots \geq \lambda_{61}$, where $D^{1/2} = (\Pi^h)^{1/2}$. $P^h(t)$ can be expressed in terms of the \mathbf{q}_k as follows:

$$\begin{aligned} P^h(t) &= D^{-1/2} (D^{1/2} P^h(t) D^{-1/2}) D^{1/2} \\ &= D^{-1/2} \left(\sum_{k=1}^{61} \lambda_k^t \mathbf{q}_k \mathbf{q}_k^T \right) D^{1/2} \end{aligned} \quad (2.42)$$

Equation (2.38) can therefore be re-expressed in terms of the left $\mathbf{q}_k^T D^{1/2}$ and right $D^{-1/2} \mathbf{q}_k$ eigenvectors of $P^h(t)$, where $\mathbf{q}_1^T D^{1/2} = \boldsymbol{\pi}^h$ and $D^{-1/2} \mathbf{q}_1 = \mathbf{1}$:

$$\begin{aligned} \mathbf{v}^h(t) &= \boldsymbol{\pi}^h + \mathbf{v}^h(0) \sum_{k=2}^{61} \lambda_k^t (D^{-1/2} \mathbf{q}_k \mathbf{q}_k^T D^{1/2}) \\ &= \boldsymbol{\pi}^h + \sum_{k=2}^{61} \frac{\lambda_k^t \mathbf{q}_{ik}}{\sqrt{\pi_i^h}} \mathbf{q}_k^T D^{1/2} \end{aligned} \quad (2.43)$$

The dynamic following fixation to the suboptimal codon i is therefore characterized by the left eigenvectors $\mathbf{q}_k^T D^{1/2}$ with weights $\lambda_k^t \mathbf{q}_{ik} / \sqrt{\pi_i^h}$. In terms of graphical representation, McCandlish (2011) defined the alternative Laplacian $L = D(I - P^h(t = 1)) = D(I - P)$ for which $w_{ij} = \pi_i^h P_{ij}$ for $i \neq j$ and

$\deg(i) = \pi_i^h(1 - P_{ii}) = \sum_{j \neq i} \pi_i^h P_{ij}$. It can be shown that the right eigenvectors of $P^h(t)$ are the generalized right eigenvectors of the $L = D(I - P)$ but with eigenvalues $0 = 1 - \lambda_1 < 1 - \lambda_2 \leq \dots, \leq 1 - \lambda_{61}$ (McCandlish, 2011):

$$LD^{-1/2}\mathbf{q}_k = D(I - P)D^{-1/2}\mathbf{q}_k \quad (2.44)$$

$$= DD^{-1/2}\mathbf{q}_k - DPD^{-1/2}\mathbf{q}_k \quad (2.45)$$

$$= D^{1/2}\mathbf{q}_k - \lambda_k D^{1/2}\mathbf{q}_k \quad (2.46)$$

$$= (1 - \lambda_k)D(D^{-1/2}\mathbf{q}_k) \quad (2.47)$$

The criterion in equation (2.41) is therefore satisfied by

$$(\mathbf{x}, \mathbf{y}) = (D^{-1/2}\mathbf{q}_2, D^{-1/2}\mathbf{q}_3).$$

An additional innovation proposed by McCandlish (2011) was to include scaling factors to make $(\mathbf{x}, \mathbf{y}) = (D^{-1/2}\mathbf{q}_2/\sqrt{1 - \lambda_2}, D^{-1/2}\mathbf{q}_3/\sqrt{1 - \lambda_3})$ to make the length of an edge $(i, j) \in E$ approximately proportional to the expected time it will take for the population to move from i to j and back again to i (McCandlish, 2011).

Chapter 3

Phenomenological Load on Model Parameters can lead to False Biological Conclusions.

3.1 Introduction

There are two ways to quantitatively describe a natural process. The phenomenological approach is to summarize relationships between variables with little or no reference to causation. The alternative is to specify a model based on known or hypothetical mechanistic links between variables that explain their relationships. For example, although Newton's law of universal gravitation provides a highly accurate description of the apparent force of attraction between objects, it does so without explaining the cause of that attraction. Newton's law is therefore phenomenological. Einstein, by contrast, described gravitation as the result of the causal process of mass generating curvature in space-time. Einstein's general theory of relativity is therefore comparatively mechanistic. Highly complex biological processes pose a particular challenge to modellers. On the one hand, there is the natural desire to build mechanistic models that capture as much of the complexity and richness of a process as possible. On the other hand, limitations in information and computational resources often make simplifying assumptions unavoidable, thereby forcing a more phenomenological approach. This tension often results in models that include both phenomenological and mechanistic components together under the same framework (Rodrigue and Philippe, 2010).

The defining feature of a parameter characterized as mechanistic is that it is interpretable with respect to some real data generating process (Liberles *et al.*, 2013). This suggests that a parameter can be thought of as being mechanistic only to the extent that its interpretation is valid, meaning that it will account for variations in the data generated only by the processes it was intended to represent. Consider the two components of a typical codon

substitution model (CSM), one for neutral substitution processes (the DNA submodel) and the other for the effects of selection at the amino acid level (the selection submodel). The DNA submodel used throughout this thesis is:

$$M_{ij} \propto \begin{cases} \kappa^{\sum s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 1 \\ \alpha \kappa^{\sum s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 2 \\ \beta \kappa^{\sum s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 3 \end{cases} \quad (3.1)$$

To review, equation (3.1) applies to any pair of codons (i, j) that differ by s nucleotides s_t of which are transversions, where the $\pi_{j_k}^*$ are position-specific nucleotide frequencies, κ the transition bias, and α and β rate parameters that account for pairs of codons that differ by $s = 2$ or $s = 3$ nucleotides. These parameters are often thought of as mechanistic (e.g., Miyazawa, 2011; Zaheri *et al.*, 2014) despite the lack of explicit connections to mechanism. The parameter for transition bias κ , for example, is not linked to any specific mechanism that promotes transitions over transversions. Instead, it is interpreted as accounting for all mechanisms that do so. It would seem that this is valid in part because the neutral substitution process is considered to be constant across sites, making the maximum likelihood estimate $\hat{\kappa}$ unambiguously interpretable at any particular site in the alignment¹. By contrast, consider a selection submodel that includes a single rate ratio ω . This parameter might be construed as mechanistic since it is correctly interpreted as accounting for all mechanisms that might impact the dN/dS rate ratio. However, $\hat{\omega}$ has little meaning when applied to any individual site because it is an average of a process that is in fact heterogeneous across sites and over time. It is therefore incorrect to interpret $\hat{\omega}$ in the context of the mechanisms behind the evolution of any particular site. The rate ratio ω in a CSM is therefore seldom if ever considered to be mechanistic.

Under the maximum likelihood (ML) framework the likelihood of a set of model parameters such as ω and the vector of branch lengths \mathbf{t} is expressed in the form of a likelihood function $L(\omega, \mathbf{t} \mid X, \tau)$, where X represents the

¹Selection effects in fact undermine the intended mechanistic interpretation of κ because its estimate is influenced not only by synonymous substitutions that are assumed to be selectively neutral but also by nonsynonymous substitutions some of which may have been fixed by positive selection.

alignment and τ the assumed topology of the tree. The maximum likelihood estimate (MLE) for (ω, \mathbf{t}) is the vector $(\hat{\omega}, \hat{\mathbf{t}})$ that maximizes the likelihood $L(\hat{\omega}, \hat{\mathbf{t}} \mid X, \tau)$ of the data. A key feature of the ML framework is that the likelihood of the data always increases when a new parameter is added to the fitted model. For example, $L(\hat{\omega}, \hat{\mathbf{t}}, \hat{\psi} \mid X, \tau)$ must theoretically be greater than $L(\hat{\omega}, \hat{\mathbf{t}} \mid X, \tau)$. The new parameter ψ is said to have improved the fit of the model in proportion to the size of the increase in likelihood it engendered, and is said to be statistically significant if the increase is larger than some prespecified threshold. Under this framework it is possible to find that ψ is statistically significant even if the process it represents did not actually play a role in the generation of the data. Any mechanistic interpretation assigned to ψ can thereby be invalidated. The risk of this is in part a function of differences between the fitted model and the actual data-generating process.

All CSMs are misspecified, meaning that they do not match the true generating process. The substitution rate matrix $Q(\omega)$, for example, characterizes the selection process with one rate ratio ω for all sites and branches. It is underspecified because it does not account for the type of heterogeneity in the selection process typically observed in real data. In general, if the selection submodel of a CSM fails to absorb a substantial proportion of the variation in site patterns due to selection effects, some of this variation might be inappropriately absorbed by parameters of the DNA submodel. This is especially likely to occur when a parameter of the DNA submodel represents a process that is confounded with selection effects. Two processes are confounded if they produce similar patterns or “signatures” in the data. It was shown in Chapter 2 that heterotachy can arise by episodic movement away from and back to the optimal amino acid for a given site provided neither selection nor drift dominates. This process of non-adaptive shifting balance would be confounded with any other process that produces similar variations in rate ratio over time, such as episodic adaptive changes in site-specific landscapes over time.

The concepts of percent reduction in deviance (PRD) and phenomenological load (PL) are introduced in this chapter to provide a means to assess the impact of confounding on the MLE of a model parameter. Deviance is the

difference between the maximum log-likelihood (LL) of a given CSM and the maximum log-likelihood of the saturated model (Ms) when both are fitted to the same alignment. The saturated model, analogous to a regression model in which there are as many predictor variables as observations, will always provide the largest LL of any CSM. The difference between this and the LL of M0 provides a baseline deviance score for comparison with differences between other pairs of models. The deviance under a model M can be reduced by the addition of a new parameter ψ . The PRD is the decrease in deviance the inclusion of ψ engenders normalized by the baseline score. A large PRD is generally considered to indicate that the new parameter improved model fit. However, better fit does not imply a better model. If ψ has a mechanistic interpretation, and if the process it represents either did not actually occur when the data was generated or is confounded with other processes, then its mechanistic interpretation is invalid and its PRD is equated to PL. A large PL is a concern because it not only invalidates the mechanistic interpretation of $\hat{\psi}$, but also increases the likelihood that ψ will found to be statistically significant. Under this scenario, the model M with ψ will provide a better fit, but might also lead to false conclusions about the true data generating process if its mechanistic interpretation is taken at face value.

The analyses in this chapter are focused on detection of the fixation of simultaneous double and triple (DT) mutations. The majority of CSMs assume that sites evolve by a series of single nucleotide substitutions, despite evidence for fixation of DT mutations (Whelan and Goldman, 2004; Kosiol *et al.*, 2007; Tamuri *et al.*, 2012). Several authors have argued that it would be beneficial to add a few extra parameters to the DNA submodel of any standard CSM to account for DT mutations (e.g., Miyazawa, 2011; Zaheri *et al.*, 2014). To investigate the utility of this recommendation, DT parameters (α and β in equation 1.24) were added to a variety of CSMs. The main study in this chapter illustrates the propensity of these models to detect DT substitutions in alignments generated with single nucleotide substitutions only. Two smaller studies were also conducted, including an analysis of the impact of PL on a

parameter intended to measure changes in the intensity of selection (the RELAX model of Wertheim *et al.*, 2014), and an analysis of the impact of PL on parameters intended to account for variations in the synonymous substitution rate (dS) across sites (Kosakovsky Pond and Muse, 2005). The results of all three analyses provide evidence for the universal applicability of the PL concept.

3.2 New Approaches

3.2.1 Modeling a Mixture of Static and Switching Sites: RaMoSS

Many commonly used CSMs assume either that rate ratios vary across sites but not time (e.g., the M-series models of Yang *et al.*, 2000a), or that temporal variations occur at all sites (e.g., the branch-site models of Guindon *et al.*, 2004; Kosakovsky Pond *et al.*, 2011; Murrell *et al.*, 2015). One exception is the branch-site model of Yang and Nielsen (2002), which allows some sites to evolve under one rate ratio across the whole tree, and others to switch from a stringent or neutral selection regime to positive selection at a specific location in the tree. The location of the switch, based on prior information, is treated as a fixed effect. Although this approach is well suited for identifying episodic directional selection on a specific branch, it is inappropriate for detecting random site-specific variations in rate ratio. Since real alignments might include both static and switching sites, I propose a new model, RaMoSS (for **R**andom **M**ixture **o**f **S**tatic and **S**witching sites), that combines the standard M-series model M3($k = 2$) with the covarion-like model CLM3($k = 2$) (cf., Galtier, 2001; Guindon *et al.*, 2004). Specifically, RaMoSS mixes with proportion p_{M3} one selection submodel with two rate-ratio categories $\omega_1 < \omega_2$ that are constant over the entire tree with a second selection submodel with proportion $p_{CLM3} = 1 - p_{M3}$ under which sites switch randomly in time between $\omega'_1 < \omega'_2$ at an average rate of δ switches per unit branch length.

3.2.2 Quantifying Phenomenological Load

Phenomenological load can be quantified using the concept of statistical deviance. Consider the likelihood of a vector of parameters $\boldsymbol{\theta}_M$ for model M given an alignment X and topology τ under the usual assumption of site independence:

$$L_M(\boldsymbol{\theta}_M | X, \tau) = \prod_{h=1}^n L_M(\boldsymbol{\theta}_M | \mathbf{x}^h, \tau) \quad (3.2)$$

It is standard practice to apply a natural-log transform to (3.2) to obtain the log-likelihood (LL):

$$LL_M(\boldsymbol{\theta}_M | X, \tau) = \ln\{L_M(\boldsymbol{\theta}_M | X, \tau)\} = \sum_{i=1}^k y_i \ln\{L_M(\boldsymbol{\theta}_M | \mathbf{x}_i, \tau)\} \quad (3.3)$$

Here $\mathbf{x}_i \in \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ represent the unique site patterns in the alignment, each of which occurs y_i times. The objective of the ML approach is to find the vector $\hat{\boldsymbol{\theta}}_M$ that maximizes LL_M or equivalently minimizes the deviance of the fitted model. Deviance is defined as the difference between the LL of the fitted model compared to the LL of the most general iid model (a.k.a., the saturated model, M_s). The MLE for the probability of a site pattern \mathbf{x}_i under the saturated model can be shown to be its observed relative frequency y_i/n . Hence, the LL for the saturated model is:

$$LL_{M_s}(X) = \sum_{i=1}^k y_i \ln(y_i/n) \quad (3.4)$$

Any CSM fitted to an N -taxon alignment can be thought of as a multinomial distribution for the 60^N possible site patterns (60 for mammalian mitochondrial DNA, 61 for the standard genetic code). M_s is the unique multinomial distribution defined by the vector of observed relative frequencies $(y_1/n, \dots, y_k/n)$. In other words, the saturated model is specified by the empirical site-pattern distribution of X . Because it takes none of the mechanisms of mutation or selection into account, ignores the phylogenetic relationships between sequences (i.e., is independent of τ), and excludes the possibility of site patterns that were not actually observed (i.e., the probability of a site pattern that was

not observed is assumed to be zero), Ms can be construed as the maximally phenomenological explanation of X (e.g., none of its parameters can be interpreted in terms of real data-generating processes). The salient feature of (3.4) is that it is always larger than the LL for any CSM. It is in this sense that Ms is saturated, akin to a regression model with the same number of predictor variables as observations. Ms therefore provides a natural benchmark for model comparisons.

The selection submodel under M0 consists of one rate ratio for all sites, similar to a regression model that fits an intercept only. The deviance under M0 is defined as:

$$D_{M_0} = 2\{\text{LL}_{M_s}(X) - \text{LL}_{M_0}(\hat{\theta}_{M_0} | X, \tau)\} \quad (3.5)$$

Equation (3.5) provides a baseline with which to compare changes in deviance for other model contrasts. For example, suppose M_ψ is the same model as M but with one extra parameter ψ . The statistical significance of this new parameter can be assessed by conducting a hypothesis test based on the log-likelihood ratio (LLR) for the M vs M_ψ contrast:

$$\text{LLR} = D_M - D_{M_\psi} = 2\{\text{LL}_{M_\psi}(\hat{\theta}_{M_\psi} | X, \tau) - \text{LL}_M(\hat{\theta}_M | X, \tau)\} \quad (3.6)$$

Equation (3.6) is an absolute measure of the decrease in deviance caused by the addition of ψ to M. An alternative relative measure is what I call the percent reduction in deviance (PRD):

$$\text{PRD}(\hat{\psi}) = \frac{D_M - D_{M_\psi}}{D_{M_0}} \times 100\% \quad (3.7)$$

This quantity can be construed as reflecting the strength of the signature for the process represented by ψ combined with random error and possibly PL. However, if an alignment is generated without the process represented by ψ , then $\text{PRD}(\hat{\psi})$ is due to phenomenological load and random error only. The notation $\text{PRD}(\hat{\psi})$ is replaced by $\text{PL}(\hat{\psi})$ when this is the case. Hence, the notation $\text{PL}(\hat{\psi})$ is used when ψ is estimated from alignments generated *in silico* without the process ψ represents, whereas $\text{PRD}(\hat{\psi})$ is used when ψ is estimated from a real alignment for which the true generating process is unknown.

3.2.3 Assessing the Realism of Alignments Simulated under MS

The standard way to assess the sensitivity of a CSM to misspecification has been to fit the model to alignments simulated using another, perhaps more complex, CSM (Anisimova *et al.*, 2001, 2002; Wong *et al.*, 2004; Zhang, 2004; Kosakovsky Pond and Frost, 2005; Yang *et al.*, 2005; Zhang *et al.*, 2005; Yang and dos Reis, 2011; Kosakovsky Pond *et al.*, 2011; Lu and Guindon, 2013). The problem with this approach is that alignments generated under even a relatively complex CSM are not misspecified in the same way as real data. Off-the-shelf CSMs make the unrealistic assumption that all sites evolve under the same vector of stationary frequencies, and assume that all nonsynonymous substitutions have the same probability of fixation for a given rate ratio. These assumptions preclude the generation of realistic levels of variation in rate ratio across sites and over time, and have until now prevented recognition of the problem of PL. The MS framework of Halpern and Bruno (1998) introduced in Chapter 1 provides a way to evolve a codon site over a tree that is consistent with the dynamics of an ideal Wright-Fisher population on a static fitness landscape. Under this framework, each site can be assigned its own vector of fitness coefficients. Amino acid proclivities and the stringency of selection reflected by the average rate ratio at a site can therefore be made to vary across sites in a way that is consistent with a real alignment. Alignments generated under MS can also exhibit heterotachy, which may comprise a significant proportion of the total variation in a real alignment (Lopez *et al.*, 2002; Jones *et al.*, 2017). MS therefore seems to be the ideal framework for generating realistic data with which to assess the reliability of a CSM.

The degree to which an alignment generated under MS mimics real data is in large part dependent on how site-specific fitness coefficients are specified. The most direct approach is to make use of site-specific amino acid frequencies derived from real data. For example, Spielman and Wilke (2016) estimated vectors of site-specific fitness coefficients from frequencies observed in structurally curated alignments of at least 150 taxa. These were then fed into Pyvolve (Spielman and Wilke, 2015a), among the first open-source software packages with the option to evolve sites using the MS framework, to produce

simulations of the original alignments. To explore how model misspecification might have influenced the analysis of the 20-taxon alignment of mammalian mtDNA presented in this chapter it was necessary to simulate alignments consistent with those data. Unfortunately, the methods presented in Spielman and Wilke (2016) are inappropriate for such a limited number of taxa. A new method, called MSm(ammalian)mtDNA, was therefore devised to generate plausible vectors of site-specific fitness coefficients. See Methods for details.

More important than the new method of simulation is the way it was assessed for realism. MSmmtDNA was validated by comparing distributions of summary statistics from simulated alignments to those of the real mtDNA alignment. The summary statistics considered were (i) the distribution of the number of amino acids per codon site, (ii) the overall amino acid and codon frequencies, and (iii) the frequency with which each pair of amino acids occur in the same site pattern. In addition, the expected distribution of simulated scaled selection coefficients for all mutations, all substitutions, all nonsynonymous mutations and all nonsynonymous substitutions generated under MSmmtDNA were compared to their empirical counterparts reported by Tamuri *et al.* (2012). The choice to use an alignment of mammalian mtDNA was largely motivated by the availability of these empirical distributions, which were derived from a concatenated alignment of 12 genes (3598 codon sites) from 244 mammal species.

3.3 Results

3.3.1 Putative DT Mutations are Detectable in a Real Alignment

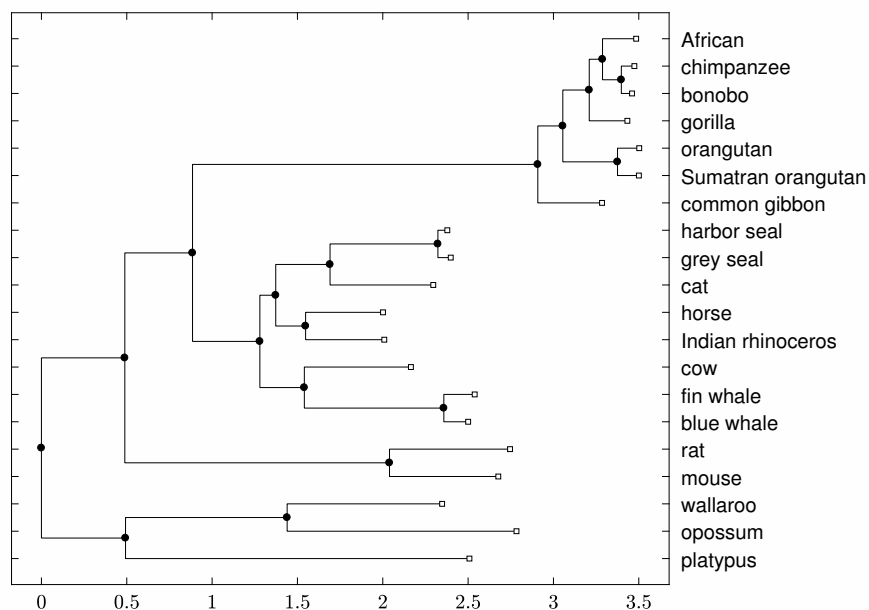


Figure 3.1: The phylogeny for the concatenation of 12 H-strand mitochondrial DNA sequences (3331 codon sites). The data is from 20 mammalian species as distributed by the PAML software package Yang (2007). The topology is that reported in Cao *et al.* (1998). Branch lengths (expected number of single nucleotide substitutions per codon) were estimated using RaMoSSwDT (the best fitting of the models used in this study). The scale on the horizontal axis is the expected number of single nucleotide substitutions per codon.

My objective was to use DT as a case study to test the hypothesis that an underspecified selection submodel combined with confounding can lead to false biological conclusions due to PL. The first step was to identify DT in a real alignment. To that purpose, models M0, M3($k = 2$), CLM3($k = 2$), and the new model RaMoSS, as well as their counterpart models that allow fixation of DT mutations, were fitted to an alignment of 20 mammalian mtDNA sequences. The tree with branch lengths estimated under the best fitting model (RaMoSSwDT) is depicted in Figure 3.1. Table 3.1 lists the LL and parameter estimates for each model. Table 3.2 shows the results for various model contrasts.

Model	LL	rate ratios	proportions	δ	% S,D,T
Ms	-26,752				
M0	-92,006	0.04			
M3	-89,162	0.01,0.15	$\hat{p}_1 = 0.71$		
CLM3	-88,880	0.00,0.21	$\hat{p}_1 = 0.77$	0.06	
RaMoSS	-88,677	0.00,0.08 0.01,0.44	$\hat{p}_{M3} = 0.73$ $\hat{p}_1 = 0.80$ $\hat{p}'_1 = 0.66$	0.21	
M0wDT	-91,280	0.03			76.4,21.5,2.1
M3wDT	-88,930	0.01,0.12	$\hat{p}_1 = 0.71$		83.0,16.7,0.3
CLM3wDT	-88,786	0.00,0.16	$\hat{p}_1 = 0.75$	0.06	86.5,13.5,0.0
RaMoSSwDT	-88,635	0.00,0.08 0.02,0.34	$\hat{p}_{M3} = 0.68$ $\hat{p}_1 = 0.80$ $\hat{p}'_1 = 0.73$	0.12	90.3,9.7,0.0

Table 3.1: Parameter estimated for the real data. Log-likelihood (LL) and parameter estimates for each model fitted to the mammalian mtDNA alignment shown in Figure 3.1.

The LLR was statistically significant for M0 vs M3($k = 2$), M3($k = 2$) vs CLM3($k = 2$), and CLM3($k = 2$) vs RaMoSS (Table 3.2). Collectively, these contrasts provide evidence for variation in rate ratio across sites and branches, and support the existence of both static and temporally dynamic sites within the alignment. The four contrasts of the form model-M vs model-MwDT were also statistically significant, and therefore apparently detected DT substitutions. However, the signal for DT became weaker with each increment in the complexity (i.e., number of parameters) of the selection submodel. The proportion of DT substitutions was inferred to be 23.6% under the simplest model contrast M0 vs M0wDT. Accounting for variations in rate ratio across sites (M3 vs M3wDT) reduced this to 17.0%. By allowing sites to switch rate ratio (CLM3 vs CLM3wDT), and allowing a mixture of static and switching sites (RaMoSS vs RaMoSSwDT), the DT proportion was further reduced to 13.5% and 9.7%, respectively. Similarly, the PRD($\hat{\alpha}, \hat{\beta}$) decreased from 1.11% for the M0 vs M0wDT contrast to 0.36%, 0.14%, and 0.06% under M3($k = 2$) vs M3wDT, CLM3($k = 2$) vs CLM3wDT and RaMoSS vs RaMoSSwDT, respectively.

Contrast	LLR	crit.val.	Detected	PRD
M0 vs M3($k = 2$)	5668	5.99	yes	4.36%
M3 vs CLM3($k = 2$)	564	2.71	yes	0.43%
CLM3 vs RaMoSS	406	9.49	yes	0.31%
M0 vs M0wDT	1452	5.99	yes	1.11%
M3 vs M3wDT	464	5.99	yes	0.36%
CLM3 vs CLM3wDT	188	5.99	yes	0.14%
RaMoSS vs RaMoSSwDT	84	5.99	yes	0.06%

Table 3.2: Model contrasts for the real data. Results for model contrasts applied to the mammalian mtDNA alignment shown in Figure 3.1.

An investigation was conducted to determine which site patterns contributed the most to the 42-point difference in LL for RaMoSS (LL = -88,677) compared to RaMoSSwDT (LL = -88,635). Of the 3331 sites patterns, 83 were invariant, 25 had nonsynonymous differences only, 1730 had synonymous differences only, and 1493 were mixed with both synonymous and nonsynonymous differences. The contribution of each of these site-pattern categories to the total LL under RaMoSS and RaMoSSwDT is listed in Table 3.3. RaMoSSwDT provided a slightly better fit to the 2.49% of site patterns that were invariant. This is because invariant sites become less likely as branch lengths increase, and RaMoSS produced larger branch lengths than RaMoSSwDT. These sites accounted for only 3 the total difference of 42 LL points. Less than 1% of all site patterns had nonsynonymous differences only. RaMoSSwDT fitted these sites slightly better as expected given that allowing DT substitutions increases the probability that a nonsynonymous substitution will occur. But since there were so few site patterns in this category, the total contribution was only 1 out of 42 LL points. Approximately 52% of site patterns had synonymous differences only. RaMoSSwDT provided a slightly worse fit to these sites. Most synonymous differences can be explained by a single nucleotide substitution at the third codon position. Allowing DT substitutions, most of which are nonsynonymous, apparently reduces the probability of a site pattern with synonymous differences only. This effect was very small however, contributing a difference of only -1 LL points, despite the large number of site patterns in this category. Approximately 45% of sites had mixed site patterns, and these accounted for 39 out of the 42 LL points difference between RaMoSS

and RaMoSSwDT. This demonstrates that mixed site patterns are more likely when the model permits DT. Critically, mixed site patterns are also more likely to exhibit heterotachy. Of the 297 site patterns with a posterior probability of switching > 0.80 (computed using equation 3.11 in Methods), 289 had mixed site patterns. The remaining 8 were among the site patterns with nonsynonymous differences only. This suggests that episodic but rare DT substitutions can be confounded with non-adaptive shifting balance since both processes can produce signatures consistent with heterotachy.

site-pattern cat.	number (%)	ΔLL	Post. > 0.80
invariant	83 (2.49)	3	0
nonsynonymous	25 (0.75)	1	8
synonymous	1730 (51.94)	-1	0
mixed	1493 (44.82)	39	289
TOTAL	3331	42	297

Table 3.3: Site pattern analysis. Each row reports the number (and %) of sites in the corresponding site-pattern category, the total change in LL associated with each category, and the number of sites for which the posterior probability of switching was > 0.80 .

A heuristic for inferring DT substitutions is to examine sites occupied by serine only (Averof *et al.*, 2000). Codon aliases for serine include TCN, where N is any nucleotide, and AGY, where Y is a pyrimidine. Minimum paths between TCN and AGY by single nucleotide steps require substitution to cystine or threonine. But these amino acids are physicochemically different than serine, and can be assumed to be less fit than serine at a site observed to be occupied by serine only. The existence of serine sites with a mix of TCN and AGY would therefore suggest that some double mutations of the form $TC \leftrightarrow AG$ were fixed. Of the 112 serine sites in the real mtDNA alignment, one site was occupied by a single alias for serine, 19 were occupied by a combination of AGT and AGC, and 92 were occupied by a combination of TCC, TCT, TCA and TCG. Aliases from the AGY and TCN groups did not appear together in any of the sites. This result, combined with the observed decrease in the strength of the evidence for DT with each incremental increase in the complexity of the selection submodel, casts doubt on the veracity of the detection of DT substitutions in the real mtDNA under RaMoSS vs RaMoSSwDT. Simulation studies were therefore conducted to investigate the possibility of false detection

of DT, as reported in the next section.

3.3.2 The Extent to which DT Parameters Carry PL is Related to Model Misspecification

There is substantial heterogeneity in selection pressure across sites within the mammalian mitochondrial genome (Garvin *et al.*, 2015). It is therefore likely that the single rate ratio of M0 provides a highly inadequate summary of variations due to selection effects in the real mtDNA alignment. The analysis reported in the previous section revealed a substantial PRD for the M0 vs M0wDT contrast (1.11% PRD, corresponding to a highly significant LLR of 1452) and a relatively large estimated proportion of DT substitutions (23.6%). One would expect a reduction in both PRD and %DT with an increase in the complexity of the selection submodel if the analysis was influenced by PL. This is exactly what was observed. RaMoSS vs RaMoSSwDT resulted in a much smaller PRD (only 0.06%) and indicated a smaller proportion of DT substitutions (9.7%). However, even the selection submodel under RaMoSS is likely to be underspecified compared to the actual data-generating process.

Three simulation studies were conducted to assess the relationship between model misspecification and the PL detected by the four M vs MwDT contrasts. Each study was conducted using a different alignment generating method, none of which included DT substitutions. Alignments were generated under RaMoSS in the first simulation. Hence, this study covered the scenario under which the selection submodel of the RaMoSS vs RaMoSSwDT contrast was not misspecified. Alignments were generated under a substantially more complex CSM in the second simulation in which each site was assigned an independent rate ratio ω^h and evolved under $Q(\omega^h)$. Alignments were therefore generated with more variation in rate ratio across sites than accounted for by any of the M vs MwDT contrasts. Alignments were generated under MSmmtDNA in the third simulation study to mimic variations in rate ratio across sites and over time comparable (as will be shown) to the real mtDNA alignment. The analyses of these simulations revealed that, although PL was readily identified

in all simulation studies under the contrast with the simplest selection submodel (M0 vs M0wDT), it was only detected under the most complex contrast (RaMoSS vs RaMoSSwDT) in alignments generated using the more realistic MSmmtDNA generating model.

3.3.3 Simulation Study 1: MLEs for the DT Process Carry Substantial PL when the Selection Submodel is Underspecified, but False Conclusions are Avoided when the Selection Submodel is Correctly Specified

Contrast	Testing For ...	median LLR	Detected
M0 vs M3	var. across sites	171	98
M3 vs CLM3	var. across time	34.0	99
CLM3 vs RaMoSS	static and switching sites	19.6	95
M0 vs M0wDT	DT mutations	29.7	99
M3 vs M3wDT	DT mutations	5.69	49
CLM3 vs CLM3wDT	DT mutations	0.01	3
RaMoSS vs RaMoSSwDT	DT mutations	0.00	0

Table 3.4: Simulation 1 Results. Median values for log-likelihood ratios (LLR) and the number of times DT was detected from 100 alignments generated under RaMoSS with $\alpha = \beta = 0$.

In the first simulation study, one-hundred 300-codon alignments were generated on the tree depicted in Figure 3.1 using RaMoSS as the generating model. A starting sequence was constructed by selecting codons in proportion to their empirical frequencies estimated from the real mtDNA. All alignments were generated starting with this same sequence. Parameters for the selection submodel (including $\omega_1, \omega_2, p_1, \omega'_1, \omega'_2, p'_1, p_{M3}$ and δ) were set to values estimated from the real mtDNA alignment using RaMoSSwDT (i.e., the best fitting model; see Table 3.1 for parameter values). Similarly, parameters for the DNA submodel, including κ and position-specific nucleotide frequencies, were set to values estimated from the real alignment, except that α and β were set to zero to exclude DT substitutions. Table 3.4 shows median results for the various likelihood ratios tests (see Appendix for median parameter estimates).

The contrast with the simplest selection submodel (M0 vs M0wDT) incorrectly rejected the null hypothesis and inferred DT mutations in almost all

trials (99/100). Improving the selection submodel by accounting for variations in rate ratio across sites (M3 vs M3wDT) yielded a substantial reduction in the false positive rate (49/100). Accounting for heterotachy (CLM3 vs CLM3wDT and RaMoSS vs RaMoSSwDT) further reduced the number of false positives to 3/100 and 0/100 respectively. RaMoSS provided the best fit in all trials, and produced median parameter estimates similar to their generating values: $\hat{\omega}_1 = 0.00$ (generating $\omega_1 = 0.00$), $\hat{\omega}_2 = 0.03$ (0.08), $\hat{p}_1 = 0.86$ (0.80), $\hat{\omega}'_1 = 0.01$ (0.01), $\hat{\omega}'_2 = 0.44$ (0.44), $\hat{p}'_1 = 0.88$ (0.66), $\hat{p}_{M3} = 0.72$ (0.73) and $\hat{\delta} = 0.17$ (0.21). It was no surprise to find that RaMoSS produced reliable parameter estimates, and that RaMoSS vs RaMoSSwDT did not falsely detect the fixation of DT mutations, since RaMoSS was an exact match to the generating process. However, it was interesting to find that the MLEs $\hat{\alpha}$ and $\hat{\beta}$ for the DT process carry substantial PL when the selection submodel was underspecified, as indicated by the high false detection rate, and that DT was only detected by the two models that do not account for heterotachy (M0 and M3). This demonstrates that random variations in site-specific rate ratios, produced in this simulation study by the CLM3($k = 2$) component of the generating model, can create false signal for DT substitutions when heterotachy is not accounted for by the selection submodel.

3.3.4 Simulation Study 2: Adding Complexity to the Selection Submodel Reduces PL even when the Submodel is Substantially Underspecified

Contrast	Testing For ...	median LLR	Detected
M0 vs M3	var. across sites	867	100
M3 vs CLM3	var. across time	5.88	75
CLM3 vs RaMoSS	static and switching sites	62.9	98
M0 vs M0wDT	DT mutations	63.3	100
M3 vs M3wDT	DT mutations	1.03	10
CLM3 vs CLM3wDT	DT mutations	0.20	3
RaMoSS vs RaMoSSwDT	DT mutations	0.01	5

Table 3.5: Simulation 2 Results. Median values for log-likelihood ratios (LLR) and the number of times DT was detected from 100 alignments generated under M3($k = n$) with $\alpha = \beta = 0$.

In the second simulation study, one-hundred 300-codon alignments were generated on the tree depicted in Figure 3.1 using $M3(k = n)$ as the generating model (where n is the number of codon sites). The objective was to produce the same level of variation in rate ratio across sites as in the real mtDNA but without heterotachy. First, a vector of codon fitness coefficients \mathbf{f}^h was drawn for each site using the MSmmtDNA model. The MS rate matrix A^h was then constructed with $\alpha = \beta = 0$ and used to determine the expected rate ratio for the site (from equation 2.10):

$$\omega^h = \frac{\sum_{(i,j)} \pi_i^h A_{ij}^h \ell_N}{\sum_{(i,j)} \pi_i^h M_{ij}^h \ell_N} \quad (3.8)$$

The rate matrix $Q(\omega^h)$ for the site-specific generating model was then constructed. Note that the rate ratio at a site evolving under $Q(\omega^h)$ is always ω^h regardless of the incumbent codon. An alignment generated using the set of $Q(\omega^1), \dots, Q(\omega^n)$ will therefore have the same level of variation in the expected rate ratio across sites as an alignment generated using the A^h (e.g., using MSmmtDNA), but without heterotachy.

Table 3.5 shows median results for the various likelihood ratios tests (see Appendix for median parameter estimates). As expected, the M0 vs M3($k = 2$) contrast detected substantial signal for variations in rate ratio across sites in all trials. Quite unexpected was the result that the M3($k = 2$) vs CLM3($k = 2$) contrast implied signal for heterotachy in 75/100 trials. This is in apparent contradiction to the design of the generating process, which precluded heterotachy. However, the signal for changes in rate ratio over time was relatively weak: the median switching rate was only $\hat{\delta} = 0.02$ or one switch per 50 single nucleotide substitutions. Furthermore, the median LLR for M3($k = 2$) vs CLM3($k = 2$) was only 5.88 (compared to the critical value of 2.71 for a 5% test) with a corresponding p-value of 0.008. Given that CLM3($k = 2$) is equivalent to M3($k = 2$) when $\delta = 0$, these results are not entirely inconsistent with sites evolving under fixed rate ratios. Nevertheless, they seem to indicate that $\hat{\delta}$ carried some PL in three-quarters of the trials. The CLM3($k = 2$) vs RaMoSS contrast similarly implied a fraction of sites with signal for heterotachy. The LLR for this contrast was significant in 98/100 trials (median LLR

= 62.9), but with a very small switching rate ($\hat{\delta} = 0.00$). RaMoSS is the same as M3($k = 4$) when $\delta = 0$, so in this case it seems that RaMoSS provided the better fit not because of PL on $\hat{\delta}$, but because four ω -categories provided a significantly better fit than two, reflecting the level of variation in ω across sites under the M3($k=n$) generating process.

Turning to the tests for fixation of DT mutations, the contrast involving the simplest selection submodel (M0 vs M0wDT) incorrectly inferred DT mutations in all 100 trials (Table 3.5). Again, improving the selection submodel substantially reduced the false positive rate. Even limited accommodation of variations in rate ratio across sites using M3($k = 2$) (e.g., with only two rate-ratio categories) reduced the false positive rate to 10/100. This was reduced further to only 3/100 and 5/100 by CLM3 vs CLM3wDT and RaMoSS vs RaMoSSwDT. These are consistent with the 5% level of significance of the likelihood ratio test, and seem to imply that both CLM3 vs CLM3wDT and RaMoSS vs RaMoSSwDT might reliably fail to detect the fixation of DT mutations when they do not occur. It must be remembered however that the generating model M3($k = n$) is unrealistic, and in particular does not simulate heterotachy. A more rigorous test of the reliability of the RaMoSS vs RaMoSSwDT contrast requires use of a more realistic alignment-generating process.

3.3.5 Simulation Study 3: RaMoSS vs RaMoSSwDT is Unreliable when Fitted to Data Generated using MSmmtDNA

Contrast	Testing For ...	median LLR	Detected
M0 vs M3	var. across sites	767	100
M3 vs CLM3	var. across time	30.6	100
CLM3 vs RaMoSS	static and switching sites	57.1	100
M0 vs M0wDT	DT mutations	147	100
M3 vs M3wDT	DT mutations	25.7	97
CLM3 vs CLM3wDT	DT mutations	12.3	76
RaMoSS vs RaMoSSwDT	DT mutations	4.34	41

Table 3.6: Simulation 3 Results. Median values for log-likelihood ratios (LLR) and the number of times DT was detected from 100 alignments generated under MSmmtDNA with $\alpha = \beta = 0$.

The M3($k = n$) generating model reflects the traditional approach of testing the impact of model misspecification by simulating alignments using a more complex CSM. However, the absence of heterotachy means that the simulated distribution of site patterns can only be unrealistic compared to the real mtDNA alignment. In the third simulation study, one-hundred 300-codon alignments were generated on the tree depicted in Figure 3.1 using the generating process called MSmmtDNA, which was formulated to produce alignments that match the real mtDNA alignment as closely as possible. The median value of parameters estimated by fitting RaMoSS to the simulated alignments were: $\hat{\omega}_1 = 0.00$ (compared to $\hat{\omega}_1 = 0.00$ for the real mtDNA) $\hat{\omega}_2 = 0.12$ (0.08), $\hat{p}_1 = 0.82$ (0.80), $\hat{\omega}'_1 = 0.00$ (0.01), $\hat{\omega}'_2 = 0.56$ (0.44), $\hat{p}'_1 = 0.60$ (0.66), $\hat{p}_{M3} = 0.80$ (0.73) and $\hat{\delta} = 0.20$ (0.21). These results suggest a substantial degree of “phenomenological similarity” between the real and simulated alignments. Note that this was not by design, since the MLEs derived from the real mtDNA alignment were not used in the formulation of MSmmtDNA. Instead, the similarity was a consequence of the method used to generate site-specific fitness coefficients (see Methods). Further comparisons between the real alignment and those simulated using MSmmtDNA appear in the next section. The remainder of this section reports the results of model fits.

The impact of PL when the models were fitted to alignments generated under MSmmtDNA is apparent in Table 3.6. The contrast involving the simplest selection submodel (M0 vs M0wDT) incorrectly inferred DT in all 100 trials, as might be expected given previous results. However, unlike the previous two simulation studies, accounting for variations in rate ratio across sites (M3 vs M3wDT) had negligible impact on the false positive rate (97/100). Although accounting for heterotachy (CLM3 vs CLM3wDT and RaMoSS vs RaMoSSwDT) reduced the number of false positives (to 76/100 and 41/100, respectively), the lowest rate was still too large given the 5% level of significance of the test. It seems that the selection submodel for RaMoSS is underspecified with respect to the MSmmtDNA generating process, with the result that substantial PL was conferred to $\hat{\alpha}$ and $\hat{\beta}$ in a large number of trials.

It now seems plausible that the detection of DT in the real mtDNA was a

false positive due to PL. If it is true that MSmmtDNA produces alignments consistent with the real data, then it can be used to estimate the distribution of $PL(\hat{\alpha}, \hat{\beta})$ for each of the M vs MwDT model contrasts for comparison with results from the real alignment. To this end, MSmmtDNA was used to generate 50 full-scale alignments, each with 3331 codon sites without DT. Each model contrast was fitted to these alignments to produce distributions of $PRD(\hat{\alpha}, \hat{\beta})$. Because α and β were set to zero in the generating process, PRD can be equated to PL. The resulting distributions are shown as boxplots in Figure 3.2. The previously described decline in the $PRD(\hat{\alpha}, \hat{\beta})$ obtained by fitting the contrasts to the real mtDNA (last column of Table 3.2) is reflected by a similar decline in the median $PL(\hat{\alpha}, \hat{\beta})$ with each incremental increase in the complexity of the selection submodel.

The diamond in each boxplot of Figure 3.2 marks the $PRD(\hat{\alpha}, \hat{\beta})$ for the corresponding contrast fitted to the real mtDNA alignment. This value falls just within the upper tail of the estimated distribution of $PL(\hat{\alpha}, \hat{\beta})$ for the RaMoSS vs RaMoSSwDT contrast (also see first boxplot in Figure 3.4). For comparison, a single full-sized alignment was generated using MSmmtDNA with α and β set to the values estimated by RaMoSSwDT fitted to the real mtDNA (e.g., with 9.7% double and 0.0% triple nucleotide mutations under MSmmtDNA, see Table 3.1). The small square in each boxplot marks the $PRD(\hat{\alpha}, \hat{\beta})$ obtained by fitting each contrast to this alignment. As the signal for DT was real in this case, $PRD(\hat{\alpha}, \hat{\beta})$ cannot be equated to $PL(\hat{\alpha}, \hat{\beta})$, but can be interpreted as an indication of real signatures for DT. The decrease in PRD with each increase in the complexity of the selection submodel is still evident however, and suggests that $\hat{\alpha}$ and $\hat{\beta}$ carry some PL. These comparisons, combined with the large number of false detections reported in Table 3.6, suggest that either the detection of DT substitutions in the real mtDNA was false or it was a true but that the rate of DT was overestimated. PL played a role in the analysis in either case.

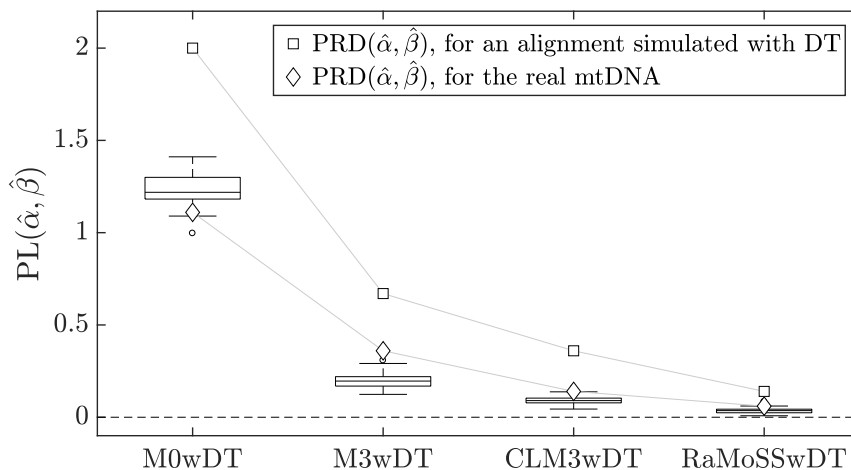


Figure 3.2: PL versus PRD. Boxplots show the distribution of $PL(\hat{\alpha}, \hat{\beta})$ for each of the M vs MwDT model contrasts fitted to 50 full-scale alignments (20 taxa, 3331 codon sites) generated under MSmmtDNA with $\alpha = \beta = 0$. Diamonds show $PRD(\hat{\alpha}, \hat{\beta})$ for each contrast fitted to the real mtDNA. Squares show the $PRD(\hat{\alpha}, \hat{\beta})$ for each contrast fitted to a full-scale alignment generated under MSmmtDNA with α and β set to values estimated from the real mtDNA using RaMoSSwDT. Circles indicate outliers in $PL(\hat{\alpha}, \hat{\beta})$ for the corresponding boxplot.

3.3.6 Alignments Generated under MSmmtDNA are Realistic by Several Measures of Comparison

The design of the third simulation study represents a substantial departure from the first two, and demonstrates that the role PL might have played in the analysis of the real mtDNA can be assessed only insofar as simulated alignments match real data. Hence, rather than using an M-series CSM as the generating process, the more realistic MS framework was used. Under MS, each site can be assigned its own vector of fitness coefficients \mathbf{f}^h . This determines the stringency of selection (the average rate ratio) and temporal dynamics (heterotachy) at the site. The purpose of using MS was to simulate alignments with heterogeneity in rate ratio across sites and over time and with heterogeneity in site patterns consistent with the real mammalian mtDNA. This section reports results that show that MSmmtDNA can produce alignments similar to the real data by several measures of comparison.

	$p(s_{ij} < -2)$	$p(-2 < s_{ij} < 2)$	$p(s_{ij} > 2)$
all mutations	0.61 (0.65)	0.39 (0.34)	0.00 (0.01)
nonsyn. mutations	0.90 (0.89)	0.09 (0.10)	0.01 (0.01)
all substitutions	0.03 (0.03)	0.94 (0.94)	0.03 (0.03)
nonsyn. substitutions	0.18 (0.14)	0.64 (0.72)	0.18 (0.14)

Table 3.7: Simulated versus real selection coefficient distributions. Comparison of interval probabilities for scaled selection coefficients s_{ij} under the generating model MSmmtDNA versus those derived empirically by Tamuri *et al.* (2012) (shown in parentheses).

Empirical distributions of scaled selection coefficients for all mutations, all substitutions, all nonsynonymous mutations and all nonsynonymous substitutions derived from mammalian mtDNA have already been published (Tamuri *et al.*, 2012). MSmmtDNA was therefore adjusted to make the estimated probability density functions of generated scaled selection coefficients match those published as closely as possible. The predicted distributions derived from 10^5 sites simulated under the resulting MSmmtDNA model were similar in shape to their empirical counterparts (cf. Figure 3.5 in Appendix vs Figure 2 in Tamuri *et al.*, 2012) and had similar probabilities $p(s_{ij} < -2)$, $p(-2 < s_{ij} < 2)$ and $p(s_{ij} > 2)$ (Table 3.7). Further comparisons between MSmmtDNA and the real mtDNA were based on a full-sized simulated alignment of 3331 codon sites. Amino acid frequencies for the simulated alignment were highly correlated with those in the real data (correlation = 0.91, p-value $\ll 0.001$, Appendix Figure 3.6), as were the codon frequencies (correlation = 0.83, p-value $\ll 0.001$). The frequencies with which each pair of amino acids was observed within a given site pattern were found to be strongly concordant (correlation = 0.91, p-value $\ll 0.001$, Appendix Figure 3.7). The distributions of the number of amino acids realized at each site were also very similar (Appendix Figure 3.8). And the simulated alignment had a similar number of invariant, nonsynonymous, synonymous, and mixed site patterns compared to the real data: (87, 18, 2052, 1174) in the simulated alignment versus (83, 25, 1730, 1493) as reported in Table 3.3.

3.3.7 Evidence of Confounding

Simulation studies demonstrate that $PL(\hat{\alpha}, \hat{\beta})$ is related to the degree to which the selection submodel is underspecified with respect to the data-generating process. Misspecification alone is likely insufficient to produce PL however. There must also be some measure of confounding between the processes governed by the mechanistic parameters in the DNA submodel with processes that generate variations in selection effects. To further illustrate this issue, the effects of changes in κ and α (both of which are parameters of the DNA submodel) on the expected number of nonsynonymous substitutions per unit branch length rN and the predicted switching rate δ (measures that reflect variations in selection effects) were examined. Specifically, the changes in rN and δ when κ was increased from 1 to 10 with $\alpha = \beta = 0$, and the changes in the same when α was increased from 0.015 (corresponding to 2.5% double mutations) to 0.075 (11% double mutations) with κ fixed at 4 and β fixed at zero were assessed.

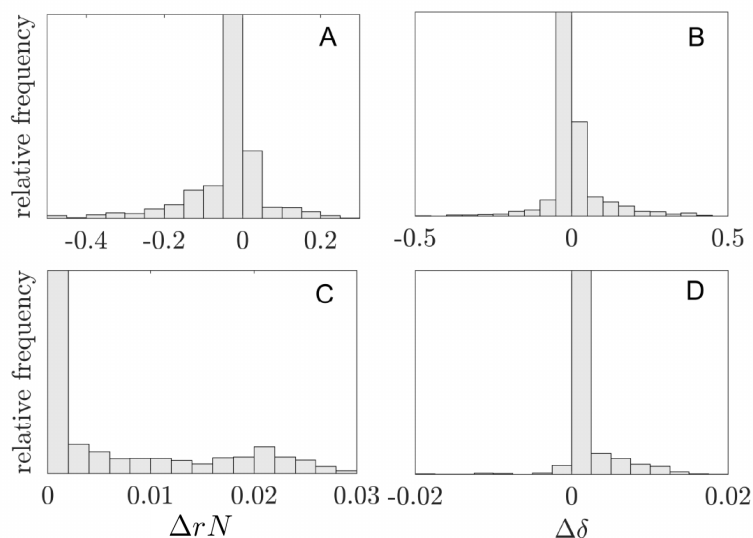


Figure 3.3: Investigating confounding. Distributions of the change in the expected number of nonsynonymous substitutions per unit branch length (ΔrN) and the expected switching rate ($\Delta\delta$) for 1000 sites with fitness coefficients generated using MSmmtDNA. A: ΔrN when κ is increased from 1 to 10 with $\alpha = \beta = 0$, B: $\Delta\delta$ when κ is increased from 1 to 10 with $\alpha = \beta = 0$, C: ΔrN when α is increased from 0.015 to 0.075 with $\kappa = 4$ and $\beta = 0$, D: $\Delta\delta$ when α is increased from 0.015 to 0.075 with $\kappa = 4$ and $\beta = 0$.

Vectors of site-specific fitness coefficients were first generated using the MSmmtDNA generating model with $\alpha = \beta = 0$. Each was used to compute theoretical predictions of site-specific values for rN and δ , once with $\kappa = 1$ and again with $\kappa = 10$ using the following equations:

$$rN^h = \sum_{(i,j)} \pi_i^h A_{ij}^h \ell_N, \delta^h = \frac{\sum_{(i,j)} \pi_i^h A_{ij}^h \ell_{\text{switch}}}{\sum_{j \neq i} \pi_i^h A_{ij}^h} \quad (3.9)$$

The resulting distributions for the change in rN and δ (ΔrN and $\Delta \delta$ respectively) with κ were both roughly symmetric and centered at zero (Figure 3.3 A and B). Hence, the same change in κ sometimes increased and sometimes decreased both rN^h and δ^h . The net effect of changes in κ on these two quantities is therefore negligible when averaged across sites. In a subsequent simulation study, 100 alignments were generated using MSmmtDNA with $\kappa = 1$. Both M0($\kappa = 1$) (i.e., M0 with κ fixed at 1) and M0 (under which κ is estimated) were fitted to these alignments. The M0($\kappa = 1$) vs M0 contrast reliably failed to reject the null hypothesis of no transition bias in all trials (i.e., estimates of κ carried no PL), despite the fact that the submodel for selection under M0 is highly underspecified with respect to MSmmtDNA.

In the second analysis, vectors of fitness coefficients were generated under MSmmtDNA with $\kappa = 4$ and $\beta = 0$ (to prohibit fixation of triple mutations). Each vector was used to compute site-specific values for rN and δ , with α set to either 0.015 (corresponding to 2.5% double mutations) or 0.075 (11% double mutations). The distributions for ΔrN and $\Delta \delta$ show that these values are almost always non-negative (Figure 3.3 C and D). Hence, an increase in α can result in an increase in both the expected nonsynonymous substitution rate and expected degree of heterotachy measured by δ when effects are averaged across sites. These simulations support the view that the process under which rare DT substitutions occur can be confounded with selection effects, and illustrate a method by which the potential for a parameter to take on PL might be assessed.

3.3.8 Assessing PL in a Model for Detecting Relaxation of Selection Pressure.

The utility of the PL framework for assessing the validity of the interpretation of model parameters in other CSMs is illustrated in this and in the next section by applying the methods to two other inferential scenarios. The first is a test for changes in selection intensity in one clade compared to the remainder of the tree (RELAX, Wertheim *et al.*, 2014). Under the RELAX model, it is assumed that each site evolved with a rate ratio randomly drawn from $\omega_R = \{\omega_1, \dots, \omega_k\}$ on a set of pre-specified reference branches, and from a modified set of rate ratios $\omega_T = \{\omega_1^m, \dots, \omega_k^m\}$ on test branches, where m is an exponent. A value $0 < m < 1$ moves the rate ratios in ω_T closer to one compared to their corresponding values in ω_R , consistent with relaxation of selection pressure at all sites on the test branches. Relaxation is indicated when the contrast of the null hypothesis that $m = 1$ versus the alternative that $m < 1$ is statistically significant. RELAX was fitted to the real mtDNA with three ω -categories using the HyPhy software package (Kosakovsky Pond *et al.*, 2004). Test branches were set to all of those in the primate clade, including the long branch leading to that clade (see Figure 3.1). The test revealed significant evidence for relaxation of selection pressure among the branches in the primate clade ($m = 0.81$, LLR = 18, p-value = 2.2×10^{-5} , $\text{PRD}(\hat{m}) = 0.015\%$). The model was also fitted to the 50 full-scale alignments generated using MSmmtDNA, under which no relaxation occurred. The null was falsely rejected in 31/50 trials. Furthermore, $\text{PRD}(\hat{m})$ estimated from the real alignment fell well within the distribution of $\text{PL}(\hat{m})$ from the 50 simulated alignments (Figure 3.4, middle boxplot). These results suggest that PL provides a plausible explanation for the detection of relaxation in selection pressure in the primate clade of the real mtDNA.

3.3.9 Assessing PL in a Model for Detecting variations in dS .

The vast majority of CSMs assume that the synonymous substitution rate is constant across sites, despite evidence that dS can vary (particularly in mitochondrial DNA, e.g., Bielawski and Gold, 2002). The second scenario is a

test for variation in dS across sites (Kosakovsky Pond and Muse, 2005). This test has no moniker that I am aware of, so it will be designated here as M3wdS (M3 with changes in dS) due to its similarity to the M-series model M3 (Yang *et al.*, 2000a). Under M3wdS, it is assumed that there are k dS categories and k dN categories that combine to produce k^2 ω -categories. M3wdS(k) is contrasted with the null model M3(k) that assumes dS constant across sites. Rejection of the null is interpreted as evidence for variations in dS across sites. M3(k) was first fitted to the real mtDNA using HyPhy with $k \in \{3, 4, 5\}$. Four categories were sufficient to account for all of the variation in rate ratio across sites (i.e., four categories fit the alignment better than three and just as well as five). M3($k = 4$) was then contrasted with M3wdS($k = 4$) using HyPhy. The contrast was found to be significant (LLR = 252, PRD(dS) = 0.19%). The M3($k = 4$) vs M3wdS($k = 4$) contrast was then fitted to the 50 full-scale alignments generated using MSmmtDNA under which dS is constant. The null was rejected in 0/50 trials. Furthermore, the phenomenological load PL(dS) associated with the parameters for dS variation was very close to zero in all fifty trials (Figure 3.4, right-most boxplot). These results support the interpretation of M3wdS as detecting genuine variations in dS in the real mtDNA alignment. Interestingly, there is biochemical support for the notion of spatial variation in dS within the mitochondrial genome: due to the different amount of time that mtDNA spends in the single-strand state during its replication process (Clayton, 1982), it will be subject to different probabilities of spontaneous mutational damage (Tanaka and Ozawa, 1994), which is expected to lead to different synonymous substitution rates (Reyes *et al.*, 1998; Bielawski and Gold, 2002; Raina *et al.*, 2005).

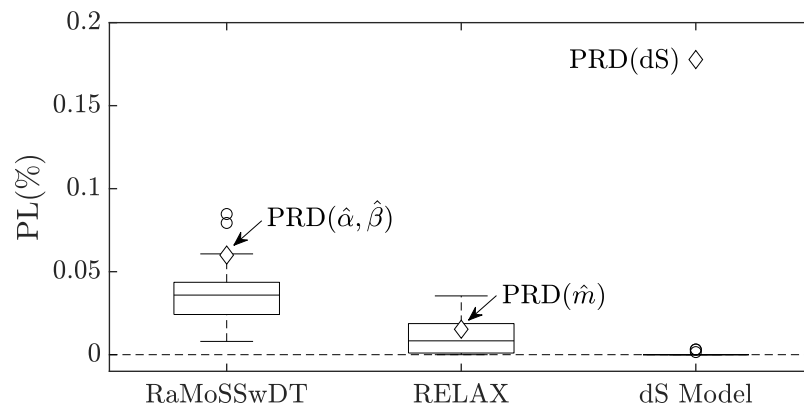


Figure 3.4: PL in other CSMs. Boxplots show the distributions of PL for parameters in models fitted to the same 50 full-scale alignments generated under MSmmtDNA (20 taxa, 3331 codon sites). Circles indicate outliers. Diamonds show PRD for each contrast fitted to the real mmtDNA. The left-most boxplot is the same as that shown in Figure 3.2 for the assessment of PL in RaMoSSwDT. The middle boxplot is for the assessment of PL in the RELAX model described in section 3.3.8. The right-most boxplot is for the assessment of PL in the test for variations in dS described in section 3.3.9. The PRD was statistically significant in 48/50 trials under RaMoSS vs RaMoSSwDT contrast, in 31/50 trials under RELAX, and in 0/50 trials under M3($k = 4$) vs M3wdS($k = 4$).

3.4 Discussion

Codon substitution models have evolved toward ever increasing complexity since their introduction by Muse and Gaut (1994) and Goldman and Yang (1994), motivated in part by the rapid increase in the quantity of information available. With greater information comes greater opportunity to tease out the effects of subtle processes. This can be achieved by adding parameters for such processes to a standard CSM. Or so it would seem. Sites for which mutation and selection are in balance can exhibit signatures consistent with random changes in site-specific rate ratios or heterotachy (e.g., mixed site patterns, Table 3.3) caused by non-adaptive shifting balance. But signatures of heterotachy can also be produced by the occasional fixation of double or triple mutations. Hence, non-adaptive shifting balance and DT are confounded processes. Consequently, if a CSM accounts for DT in its DNA submodel but fails to account for non-adaptive shifting balance in its selection submodel, the rate parameters (α, β) will be forced to account for signatures of heterotachy alone. The analyses in this chapter demonstrate that confounding can result in false inference for DT. They also provide an example of a general

principle, namely that it is possible for a parameter meant to support a specific interpretation (a so-called mechanistic parameter) to be inferred to be statistically significant even if the process it represents did not occur. When this happens the parameter’s MLE is said to carry phenomenological load and its intended interpretation is lost.

The basal cause of PL was shown to be confounding. It was stated in the introduction that two processes are confounded if they produce a common signature in the data. This implies that the generating process is the ultimate source of confounding. But whether or not confounding manifests depends on the relationship between the fitted model and the data. This is in part because different CSMs are sensitive to different signatures. The existence of multiple processes that generate heterotachy has little or no impact on the estimation of the parameters in M0 because this model ignores temporal dynamics, for example. By contrast, the above analyses suggest that RaMoSSwDT overestimated DT in the real mmtDNA ($\approx 10\%$ DT substitutions), and that this occurred because the parameter for site-specific shifts in rate ratio (δ), and the parameters for the rates of double (α) and triple (β) mutation, are all sensitive to the same signatures in the alignment, those consistent with heterotachy. An analysis of a similar although larger set of mammalian mtDNA (244 taxa with 3598 codon sites) using the site-wise mutation-selection model (swMutSel, Tamuri *et al.*, 2012) produced an estimate of DT an order of magnitude smaller (approximately 1% DT substitutions). swMutSel utilizes signatures in the alignment that RaMoSSwDT is insensitive to in the form of empirical site-specific codon frequencies. The temporal dynamics at a site is to a large degree characterized by its site-specific frequencies, as illustrated in Chapter 2, so their inclusion in swMutSel likely captured some variation in selection effects due to non-adaptive shifting balance. This apparently facilitated the detection of distinct signatures for DT (if DT was real), or else reduced the PL carried by DT parameters (if it was not). Hence, the degree to which confounding impacts inference is dependent on the signatures present in the data that the model is sensitive to, or in other words on the relationship between model and data (Jones *et al.*, 2019a).

It can happen that two processes produce signatures that differ only slightly and in such a way that they are confounded under a given CSM when information is sparse, but readily disentangled when information is rich. Such a scenario might not be uncommon, particularly among mixture models (Mingrone *et al.*, 2018), but is an issue only if the amount of data required to ameliorate associated pathologies (e.g., false positives due to PL) is prohibitively large. Under this scenario it might be said that the processes are only nearly confounded. In the analyses presented in this chapter, by contrast, the reduction in deviance engendered by the inclusion of the parameters for DT was associated with mixed site patterns in the real mtDNA alignment (Table 3.3). A larger taxonomic sample or the addition of more genes to the concatenation would result in more mixed site patterns, and would presumably increase the probability of falsely inferring DT. This is supported by the observation that the false positive rate for DT under RaMoSS vs RaMoSSwDT increased from 41% (41/100) to 96% (48/50) among alignments generated using MSmmtDNA when the number of sites was increased from 300 to 3331. It therefore seems that false detection of DT substitutions by RaMoSS vs RaMoSSwDT was not driven by lack of information, but by an abundance of information (cf. Kumar *et al.*, 2011). It is under the scenario where more site patterns (or more taxa) only worsen PL that the two processes are said to be perfectly confounded. To be clear, the introduction of information of a different type into the analysis, such as site-specific codon frequencies in the case of swMutSel (or a discrete phenotype as shown in Chapter 4), can potentially allay pathologies associated with perfect confounding.

It would be helpful to have a means to assess in advance whether confounding might impact inference under a given CSM. This was attempted in the section 3.3.7, where it was shown that an increase in κ (a property of the mutation process) is not correlated with an increase in the rate of fixation of nonsynonymous mutations (a property of the substitution process), but that an increase in the double mutation rate is. Such a result is intuitive, since the fixation of a double mutation at a site along a branch (e.g., TTA(L) \rightarrow GCA(A)) can be consistent with the fixation of two single mutations in rapid

succession (e.g., TTA(L) \rightarrow GTA(V) \rightarrow GCA(A)), and therefore manifest as a transient elevation in the nonsynonymous to synonymous rate ratio at that site under a model that does not allow DT. Indeed, such intuition might have been used to predict the possibility of confounding between episodic elevations in dN/dS and episodic fixation of DT mutations. The analysis in section 3.3.7 was based on predictions derived from a mechanistic model however, not by fitting a CSM to data. Given the supposition that the impact of confounding on inference depends on the relationship between a CSM and the actual data it is to be fitted to, it would seem that the only currently available method to identify PL is a case-by-case application of the approach illustrated in Figure 3.4. To reiterate: suppose a mechanistic parameter ψ were introduced into a substitution model M to give the model M_ψ . Further suppose that the M vs M_ψ contrast indicated a significant $\text{PRD}(\hat{\psi})$ when fitted to a real alignment. To determine whether the cause of the balance of this reduction was real signal or PL, one can first generate alignments in such a way as to resemble the real alignment as closely as possible but without the mechanistic process represented by ψ . These would be fitted to M vs M_ψ to produce a null distribution for $\text{PL}(\hat{\psi})$. Confounding would be inferred to have influenced the analysis when the $\text{PRD}(\hat{\psi})$ computed from the real alignment is found to be no greater than the 95% percentile of the $\text{PL}(\hat{\psi})$ distribution. This approach requires a method to mimic the real alignment. The Pyvolve software package (Spielman and Wilke, 2015a) provides a way to generating alignments consistent with a large real alignment. The methods used in this chapter (described in Methods) provide an alternative approach for smaller alignments.

The results presented in this chapter have implications about how the performance of a CSM should be assessed. Early efforts to test the reliability of CSMs made use of the comparatively simplistic generating models available at the time under the assumption that the findings of such analyses would be applicable to real alignments (Anisimova *et al.*, 2001, 2002; Wong *et al.*, 2004; Zhang, 2004; Yang *et al.*, 2005; Zhang *et al.*, 2005; Yang and dos Reis, 2011; Lu and Guindon, 2013). Implicit in this methodology is the presupposition that the reliability of a CSM has little to do with the data. In its original

instantiation, for example, the Yang-Nielsen Branch-Site Model (YN-BSM, Yang and Nielsen, 2002) was evaluated using real data only. It was later shown via simulation that the original YN-BSM is prone to falsely infer positive selection under certain testing scenarios (Zhang, 2004). A modified version of the model was subsequently shown to be reliable under the same testing scenarios (Zhang *et al.*, 2005). Hence the problem was implicitly assumed to be with the model, with little consideration of the role the data might have played in the observed pathology. The problem with this approach is that it leaves open the possibility that the modified YN-BSM might still be unreliable when fitted to alignments simulated using an alternative, more realistic, generating scenario. Indeed, some of the results presented in Chapter 4 support this conjecture.

The importance of using a realistic data-generating process was illustrated by three simulation studies. In the first study, alignments were generated using RaMoSS to assess the reliability of the RaMoSS vs RaMoSSwDT contrast in the absence of any model misspecification. The fixation of double and triple mutations was not inferred in any of the 100 simulated alignments. To assess performance in the presence of some misspecification, alignments in the second simulation study were generated using $M3(k = n)$, a model that typifies traditional methods to assess model reliability. The false positive rate for DT under the RaMoSS vs RaMoSSwDT contrast was only 5/100. In the past this result would have been sufficient to conclude that the contrast is a reliable instrument with which to detect signatures of DT in real data, similar to the conclusion implicit in Zhang *et al.* (2005) about the modified YN-BSM. However, alignments in the third simulation study were generated to have variations in selection effects across sites and over time that mimic the real mtDNA alignment. Under this generating scenario, RaMoSS vs RaMoSSwDT falsely detected DT in 41/100 of the 300-codon alignments and 48/50 of the full-scale alignments. These results illustrate that pathologies associated with confounding might only be realized by fitting a contrast to be applied to a real data set to alignments that are comparable with that data. It was shown that the generating model MSmmtDNA can produce alignments that

are similar in many respects to the real mtDNA alignment used in this study. However, MSmmtDNA neglects many important aspects of molecular evolution that might further impact inference. For example, MSmmtDNA does not include changes in site-specific fitness coefficients that initiate site-specific dynamics consistent with adaptive evolution (dos Reis, 2015), and does not take into account epistatic interactions that preserve thermodynamic stability (Pollock *et al.*, 2012). It might therefore be necessary to continue to work toward the formulation of data-generating methods that include these and other such processes.

3.5 Methods

3.5.1 RaMoSS

RaMoSS is a mixture of two standard CSMs: M3($k = 2$) to account for static sites (those evolving under one of two rate ratios ω_1 or ω_2 across the tree) and CLM3($k = 2$) to account for switching sites (those that change between ω'_1 and ω'_2 randomly over time). The likelihood for RaMoSS is therefore a weighted average of the likelihoods for the M3($k = 2$) and CLM3($k = 2$) components:

$$L_{\text{RaMoSS}}(\boldsymbol{\theta}_{\text{RaMoSS}} | X) = p_{\text{M3}} L_{\text{M3}}(\boldsymbol{\theta}_{\text{M3}} | X) + (1 - p_{\text{M3}}) L_{\text{CLM3}}(\boldsymbol{\theta}_{\text{CLM3}} | X) \quad (3.10)$$

where X represents the alignment, and $\boldsymbol{\theta}_{\text{RaMoSS}} = \langle p_{\text{M3}}, \boldsymbol{\theta}_{\text{M3}}, \boldsymbol{\theta}_{\text{CLM3}} \rangle$ a vector that includes the model parameters for both M3($k = 2$) and CLM3($k = 2$), as well as the additional parameter p_{M3} for the proportion of sites evolving under M3($k = 2$). The posterior probability that the rate ratio at the h^{th} site switched episodically between ω'_1 and ω'_2 can be estimated from the MLE for $\boldsymbol{\theta}_{\text{RaMoSS}}$ using the standard naive empirical Bayesian approach:

$$P(\text{switching} | \mathbf{x}^h, \hat{\boldsymbol{\theta}}_{\text{RaMoSS}}) = \frac{L_{\text{CLM3}}(\mathbf{x}^h | \hat{\boldsymbol{\theta}}_{\text{CLM3}})(1 - \hat{p}_{\text{M3}})}{L_{\text{RaMoSS}}(\mathbf{x}^h | \hat{\boldsymbol{\theta}}_{\text{RaMoSS}})} \quad (3.11)$$

3.5.2 Model Contrasts

Nested models (a null model versus an alternative, e.g., M0 vs M0wDT) can be compared using a log-likelihood ratio test. The null hypothesis is that

the data was generated under the simpler of the two models (e.g., M0). This is rejected if the log-likelihood ratio (LLR) for the test is larger than a critical value determined by the limiting distribution of the LLR statistic and the level of significance of the test. In this chapter the models M0, M3($k = 2$), CLM3($k = 2$) and RaMoSS were fitted to real and simulated alignments. Each allows single nucleotide substitutions only (e.g., $\alpha = \beta = 0$). The four models have counterparts that allow double and triple substitutions (e.g., α and β are estimated): M0wDT, M3wDT, CLM3wDT, and RaMoSSwDT. The contrast between M and MwDT provides a test for DT mutations, where $M \in \{M0, M3(k = 2), CLM3(k = 2), RaMoSS\}$. In a similar fashion, the M0 vs M3 contrast provides a tests for variation in the rate ratio across sites; M3($k = 2$) vs CLM3($k = 2$) a test for variations in the rate ratio over time; and CLM3($k = 2$) vs RaMoSS a test for a combination of static and switching sites in the same alignment compared to switching sites only. The limiting distribution of the LLR statistic is often unknown when parameters are on the boundary under the null hypothesis. In such cases, it is standard practice to use a distribution that is thought to be more conservative (i.e., less likely to reject the null hypothesis) than the unknown true distribution.

contrast	d.f.	distribution	implemented	crit. val.
M vs MwDT	2	n/a	χ_2^2	5.99
M0 vs M3($k = 2$)	2	n/a	χ_2^2	5.99
M3($k = 2$) vs CLM3($k = 2$)	1	$0.5\chi_0^2 + 0.5\chi_1^2$	$0.5\chi_0^2 + 0.5\chi_1^2$	2.71
CLM3($k = 2$) vs RaMoSS	4	n/a	χ_4^2	9.49

Table 3.8: List of Critical Values. Critical values used for the log-likelihood ratios tests in this article. d.f. is the number of extra parameters in the larger model compared to its nested counterpart.

The distributions used for the tests in this study are listed in Table 3.8, along with the corresponding critical values for 5% level of significance. The null hypothesis for all of the M vs MwDT contrasts places both $\alpha = 0$ and $\beta = 0$ on the boundary of the parameter space. The theoretical limiting distribution is therefore a mixture of the χ_0^2 , χ_1^2 and χ_2^2 distributions (Self and Liang, 1987). The mixing weights are unknown however, so the χ_2^2 distribution was used to be conservative. The proportion of sites in the ω_2 category under

the null for the M0 vs M3($k = 2$) contrast is $p_2 = 0$, making ω_2 unidentifiable. Similarly, the proportion of sites evolving with constant rate ratio under the null for the CLM3($k = 2$) vs RaMoSS contrast is $p_{\text{M3}} = 0$, making ω'_1 , ω'_2 and p'_1 unidentifiable. The theoretical limiting distributions for these contrasts are not available from Self and Liang (1987). The conventional χ^2_{df} distribution with degrees of freedom (df) equal to the difference in the number of parameters (Table 3.8) was used in these cases. The theoretical limiting distribution for the M3($k = 2$) vs CLM3($k = 2$) contrast is known to be an equal mixture of a χ^2_0 and a χ^2_1 (Self and Liang, 1987).

3.5.3 Generating Alignments using MSmmtDNA

The most direct way to simulate alignments consistent with real data is to estimate site-specific fitness coefficients $\{\mathbf{f}^1, \dots, \mathbf{f}^n\}$ from that data (Tamuri *et al.*, 2012, 2014). The MS framework can then be used to construct site-specific substitution rate matrices $\{A^1, \dots, A^n\}$ from which to generate data. The Pyvolve software package (Spielman and Wilke, 2015a) includes modules for this purpose. In the course of the study presented in this chapter it was necessary to generate data consistent with the 20-taxon concatenated alignment of H-strand mammalian mitochondrial DNA sequences provided by the PAML software package (Yang, 2007). This alignment is too small for the direct approach. It was therefore necessary to devise a more approximate generating procedure, which is described in this section.

The degree to which an alignment generated under MS can be said to be realistic is in large part dependent on how site-specific amino acid fitnesses are constructed. One method is to draw vectors of fitness coefficients from a normal distribution as described in Chapter 2. But since fitnesses are random, it is possible to draw a vector that assigns nearly the same fitness to a pair of amino acids with very different physicochemical properties. Hence, a site pattern might easily contain both isoleucine (hydrophobic) and serine (polar), for example. This is unlikely to occur at a real site evolving under stringent selection without a drastic change in the physicochemical requirements for that site. It would be more realistic to see isoleucine occur together with the

similarly hydrophobic aliphatic amino acids leucine and valine. Furthermore, the stringency of selection, determined by the standard deviation of the normal distribution, must have a realistic level of variance across sites.

Taking these requirements into account, the following MSmmtDNA generating model was used to produce a vector of fitness coefficients for a codon site:

1. A codon for the h^{th} site was randomly drawn using a multinomial distribution with probabilities equal to the empirical codon frequencies for the real mtDNA.
2. The amino acid X corresponding to the chosen codon was assigned a provisional fitness of 0.25.
3. A provisional vector of fitnesses for the remaining amino acids was constructed by dividing $\langle v_{Y_1}, \dots, v_{Y_{19}} \rangle$ by its largest element, where v_Y is the number of site patterns in the real mtDNA that included both amino acids X and Y. This gave the amino acid that paired most frequently with X (call it Z) a fitness of one and all other amino acids $Y \neq X$ a fitness less than one.
4. Each element of $\langle v_{Y_1}, \dots, v_{Y_{19}} \rangle$ was then reduced by a random draw from a half-normal distribution with mean zero and standard deviation one. The expected value of the half-normal distribution is $\sqrt{2/\pi} \approx 0.80$. The expected fitness of Z was therefore 0.20, slightly less than the fitness of X. Other amino acids tended to have lower fitness.
5. A scaling factor $\sigma^h \sim 0.001 + (0.01 - 0.001) \times B$ was drawn to determine the stringency of selection at the site, where $B \in [0, 1]$ is a beta random variable with shape parameters $u, v > 0$. Values of $\sigma^h \in [0.001, 0.01]$ closer to the upper bound correspond to greater stringency, whereas values closer to the lower bounded correspond to a balance between selection and drift that typically results in heterotachy (Jones *et al.*, 2017). Parameters u and v for the beta distribution were chosen to

make the distributions of scaled selection coefficients s_{ij} match those reported by Tamuri *et al.* (2012) as closely as possible.

6. A vector \mathbf{f}^h of fitness coefficients for the 60 codons (i.e., for mammalian mitochondrial DNA) was then constructed from the amino acid fitnesses assuming synonymous codons to be equally fit. This vector was scaled to make its standard deviation equal to σ^h .
7. \mathbf{f}^h was then used to construct the site-specific rate matrix A^h .

The following describes how the parameters (u, v) for the beta distribution were determined. For a given (u, v) , 1000 draws of \mathbf{f}^h were used to approximate the PDFs of the s_{ij} for all mutations, nonsynonymous mutations, all substitution and nonsynonymous substitutions (as detailed in the next section). Probabilities $p(s_{ij} < -2)$, $p(-2 < s_{ij} < 2)$ and $p(s_{ij} > 2)$ were calculated and compared with empirical values reported by Tamuri *et al.* (2012). This process was repeated over a grid of (u, v) coordinate pairs. The coordinate corresponding to the smallest sum of squared differences between simulated and empirical probabilities was found to be $(u, v) = (0.08, 0.02)$. These values give σ^h a U-shaped density function with most of its mass near the upper and lower bounds of its domain $[0.001, 0.01]$.

Site-specific fitness coefficients $\mathbf{f}^h = \langle f_1^h, \dots, f_{60}^h \rangle$ for the 60 codons were converted into scaled selection coefficients $s_{ij}^h = N_e(f_j^h - f_i^h)$ assuming an effective population size of $N_e = 1000$ and a ploidy of one for mtDNA. These were used to construct a site-specific substitution rate matrix A^h as follows:

$$A_{ij}^h \propto \begin{cases} M_{ij} & \text{if } s_{ij}^h = 0 \\ M_{ij} \frac{2s_{ij}^h}{1 - \exp(-2s_{ij}^h)} & \text{otherwise} \end{cases} \quad (3.12)$$

Diagonal elements A_{ii}^h were specified to make rows sum to zero. The transition bias κ and position-specific nucleotide frequencies $\{\pi_{i_k}^* \mid i \in \{T, C, G, A\}, k \in \{1, 2, 3\}\}$ for the mutation rate matrix M were set to values estimated from the real mtDNA alignment. The double and triple nucleotide mutation rate α and β were both set to zero for all simulations unless otherwise indicated (e.g., Figure 3.2). Each resulting A^h has its own vector of stationary frequencies

$\boldsymbol{\pi}^h = \langle \pi_1^h, \dots, \pi_{60}^h \rangle$ and its own expected rate r^h :

$$r^h = \sum_{j \neq i} \pi_i^h A_{ij}^h \{\ell_1 + \ell_2 + \ell_3\} \quad (3.13)$$

The indicator ℓ_k is one if i and j differ by $k \in \{1, 2, 3\}$ nucleotides and zero otherwise. All A^h were divided by the mean $r = (1/n) \sum_{h=1}^h r^h$ so that branch length could be interpreted as the expected number of single nucleotide substitutions per codon.

3.5.4 Constructing PDFs for Scaled Selection Coefficients

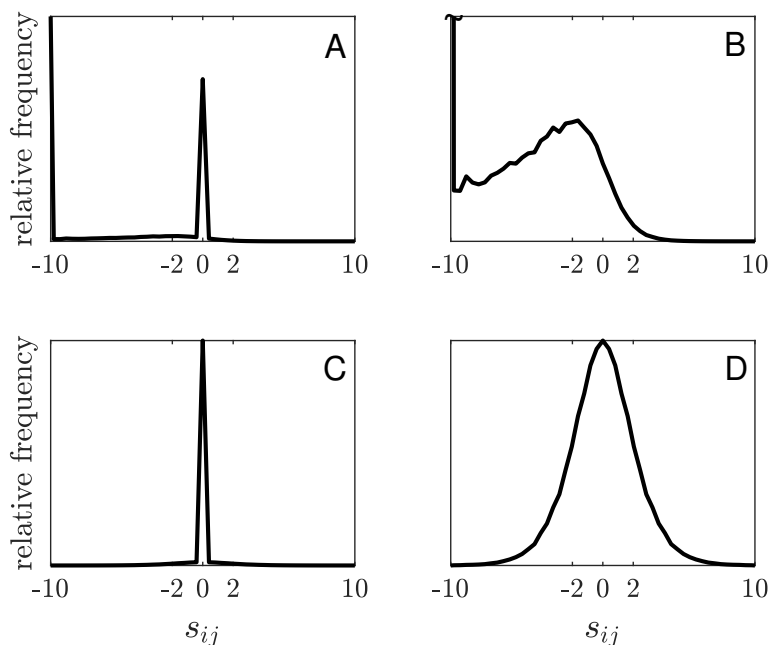


Figure 3.5: Distributions of scaled selection coefficients. Each distribution resulted from the proposed method of simulating vectors of site-specific fitness coefficients for mammalian mtDNA. A: all mutations; B: nonsynonymous mutations; C: all substitutions; and D: nonsynonymous substitutions. These are very similar to empirical distributions obtained in an analysis of 12 mitochondrial genes from 244 placental mammal species (Tamuri *et al.*, 2012).

The probability density functions for the scaled selection coefficients depicted in Figure 3.5 were approximated by discrete probability mass functions (PMFs). This section explains how the PMFs were constructed. It started with a fixed set of $n = 10^5$ vectors of site-specific fitness coefficients from which scaled selection coefficients were produced. The PMF for all mutations was then constructed as follows:

1. $p_{ij}^h = \pi_i^h M_{ij}$ was computed for each s_{ij}^h ; p_{ij}^h is proportional to the long-run probability that a mutation will occur at site h and correspond to $i \rightarrow j$ with associated scaled selection coefficient s_{ij}^h .
2. The elements of S (the set of all s_{ij}^h) were then partitioned into 50 bins. The left-most bin was the interval $(-\infty, -10)$ and the right-most bin was $(10, +\infty)$. The remaining bins between ± 10 were constructed with bin width ≈ 0.4 .
3. Each bin was assigned a sum $c_b = \sum_{i,j,h} p_{ij}^h \ell(s_{ij}^h \in \text{the } b^{\text{th}} \text{ bin})$ where $\ell(s_{ij}^h \in \text{the } b^{\text{th}} \text{ bin})$ is one if s_{ij}^h is in the b^{th} bin and zero otherwise.
4. Each c_b was then divided by $\sum_{b=1}^{50} c_b$.
5. The resulting values were plotted against the bin centers, except for the end points c_1 and c_{50} , for which the abscissa was -10 and +10, respectively.

The PMF for all substitutions was constructed by first setting $p_{ij}^h = \pi_i^h A_{ij}^h$, where A_{ij}^h is the site-specific substitution rate matrix, followed by the same steps 2 to 5. The PMFs for nonsynonymous mutations and nonsynonymous substitutions were similarly constructed using s_{ij}^h and p_{ij}^h corresponding to nonsynonymous pairs of codons i and j . The resulting PMFs approximate continuous distributions of scaled selection coefficients s_{ij} , and can be used to approximate integrals.

3.6 Appendix

3.6.1 Tables of Median MLEs for Simulation Studies

Model	LL	rate ratios	proportions	δ	% S,D,T
M0	-6972	0.02			
M3	-6890	0.00,0.08	$\hat{p}_1 = 0.84$		
CLM3	-6866	0.00,0.21	$\hat{p}_1 = 0.93$	0.06	
RaMoSS	-6859	0.00,0.03 0.01,0.44	$\hat{p}_{M3} = 0.72$ $\hat{p}_1 = 0.86$ $\hat{p}'_1 = 0.88$	0.17	
M0wDT	-6960	0.01			90.5,7.1,2.4
M3wDT	-6888	0.00,0.07	$\hat{p}_1 = 0.85$		94.1,4.6,1.3
CLM3wDT	-6866	0.00,0.19	$\hat{p}_1 = 0.92$	0.06	98.6,1.1,0.3
RaMoSSwDT	-6859	0.00,0.03 0.01,0.43	$\hat{p}_{M3} = 0.73$ $\hat{p}_1 = 0.85$ $\hat{p}'_1 = 0.88$	0.18	99.6,0.4,0.0

Table 3.9: Simulation 1 Medians. Median values for parameter estimates derived from 100 alignments generated under RaMoSS with $\alpha = \beta = 0$.

Model	LL	rate ratios	proportions	δ	% S,D,T
M0	-10,023	0.13			
M3	-9589	0.01,0.42	$\hat{p}_1 = 0.70$		
CLM3	-9585	0.01,0.45	$\hat{p}_1 = 0.71$	0.02	
RaMoSS	-9551	0.00,0.31 0.05,0.69	$\hat{p}_{M3} = 0.69$ $\hat{p}_1 = 0.74$ $\hat{p}'_1 = 0.65$	0.00	
M0wDT	-9985	0.11			88.7,4.5,6.8
M3wDT	-9588	0.01,0.41	$\hat{p}_1 = 0.70$		98.5,0.7,0.7
CLM3wDT	-9583	0.01,0.44	$\hat{p}_1 = 0.71$	0.02	99.1,0.4,0.5
RaMoSSwDT	-9550	0.00,0.31 0.05,0.70	$\hat{p}_{M3} = 0.69$ $\hat{p}_1 = 0.74$ $\hat{p}'_1 = 0.65$	0.00	99.5,0.2,0.5

Table 3.10: Simulation 2 Medians. Median values for parameter estimates derived from 100 alignments generated under M3(k=n) with $\alpha = \beta = 0$.

Model	LL	rate ratios	proportions	δ	% S,D,T
M0	-8056	0.05			
M3	-7713	0.01,0.25	$\hat{p}_1 = 0.76$		
CLM3	-7698	0.00,0.30	$\hat{p}_1 = 0.78$	0.05	
RaMoSS	-7670	0.00,0.12 0.00,0.62	$\hat{p}_{M3} = 0.79$ $\hat{p}_1 = 0.82$ $\hat{p}'_1 = 0.57$	0.22	
M0wDT	-8018	0.04			80.9,13.7,5.4
M3wDT	-7702	0.00,0.21	$\hat{p}_1 = 0.76$		90.3,7.6,2.1
CLM3wDT	-7691	0.00,0.25	$\hat{p}_1 = 0.78$	0.05	92.8,5.7,1.5
RaMoSSwDT	-7666	0.00,0.12 0.00,0.56	$\hat{p}_{M3} = 0.80$ $\hat{p}_1 = 0.82$ $\hat{p}'_1 = 0.60$	0.20	96.3,2.9,0.8

Table 3.11: Simulation 3 Medians. Median values for parameter estimates derived from 100 alignments generated under MutSel-mmtDNA with $\alpha = \beta = 0$.

3.6.2 Observed versus Simulated

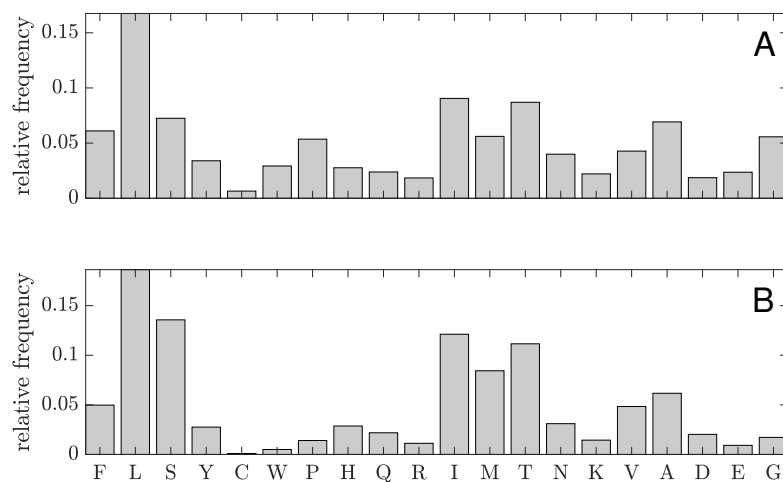


Figure 3.6: A comparison of the observed versus simulated amino acid frequencies. A: Frequencies obtained from the real data; B: the same for the simulated alignment (20 taxon, 3331 sites).

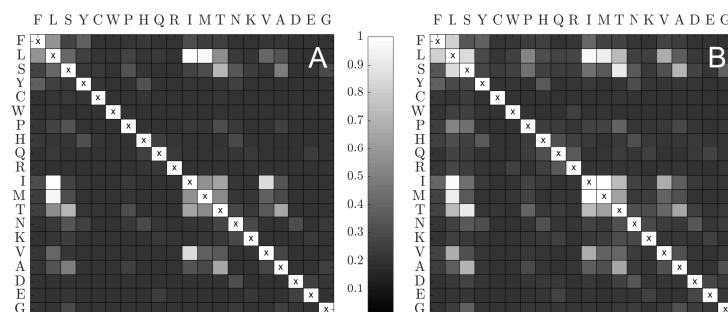


Figure 3.7: A comparison of the observed versus simulated relative pairwise amino acid frequencies. For any cell, the value is the proportion of sites where the amino acids indicated were both present. A: Values obtained from the real mtDNA alignment; B: values obtained from a simulated alignment (20 taxon, 3331 sites). The same grayscale applies to both panels.

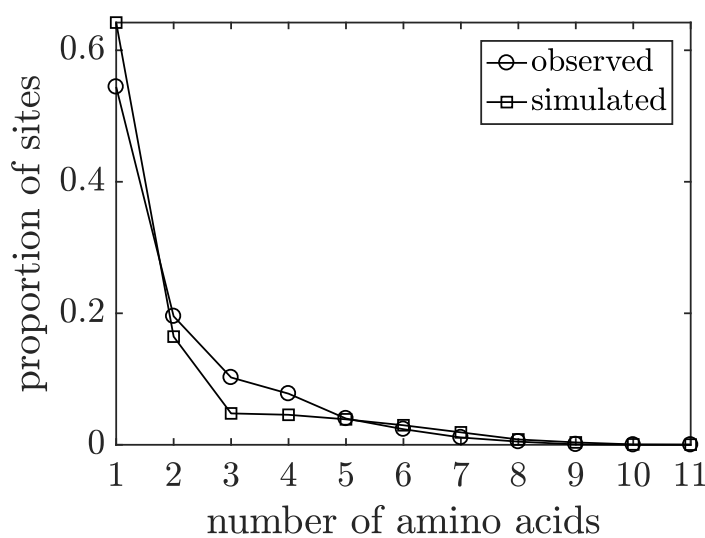


Figure 3.8: A comparison of the observed versus simulated distribution of the number of amino acids realized at a site for mammal mtDNA (20 taxon, 3331 sites).

Chapter 4

A Phenotype-Genotype Codon Substitution Model for Detecting Adaptive Evolution.

4.1 Introduction

Statistical models for the evolution of phenotypes have traditionally been formulated independently of models for the evolution of gene sequences. Yet the two approaches share a common motivation, namely to provide a means to test various evolutionary hypotheses regarding apparent structural and/or functional novelties that might have occurred as a result of adaptation. Analyzing the two data types separately neglects any possible advantage of combining information and belies the fundamental objective of identifying individual genes whose evolution can be mechanistically linked to adaptive changes in phenotype. The centrality of this objective underlines the need for models that combine the two types of data under a common statistical framework. This chapter presents such a model.

Among the first models for the evolution of phenotype were those developed to infer the rate and mode (e.g., gradual or punctuated) of phenotypic evolution, or to infer correlations between two phenotypic measures or between a phenotype and a contextual variable (for a brief review see Cornwell and Nakagawa, 2017). Such models typically assume either a continuous phenotype that evolved via Brownian Motion (Felsenstein, 1985) or a discrete phenotype that evolved via a Markov process (Pagel, 1994). These provide the basis for a wide variety of phylogenetic comparative methods or PCMs. Sophisticated PCMs include models that assume an Ornstein-Uhlenbeck “mean-reverting” evolutionary process (Hansen, 1997), models that account for temporal dynamics in the form of changes in the rate of change in a phenotype over the tree (Butler and King, 2004; O’Meara *et al.*, 2006; Eastman *et al.*, 2011), and

models that can be used to infer a relationship between phenotype and diversification (e.g., lineage diversification PCMs such as the binary state speciation-extinction model, Maddison *et al.*, 2007). More recently, several models for the analysis of multivariate data have been proposed (for a critical assessment of such methods see Adams and Collyer, 2018). The relevant point here is that the majority of PCMs use alignments of homologous protein-coding genes to estimate phylogenetic relationships that are treated as fixed and known for the remainder of an analysis based on the phenotype data alone.

Codon substitution models (CSMs) were developed to detect evidence of adaptation at the molecular level. Under the current paradigm, the canonical signature of positive selection in the form of a nonsynonymous-to-synonymous rate ratio (typically denoted ω) greater than its neutral expectation (i.e., $\omega > 1$) is considered evidence of adaptation (e.g., Yang *et al.*, 2000a). Among the more sophisticated CSMs in common use today are the branch-site models (BSMs) designed to detect evidence of adaptation at some sites along particular branches of the tree (Yang and Nielsen, 2002; Yang *et al.*, 2005; Zhang *et al.*, 2005). An alternative approach, based on amino acid substitution models (ASMs), was formulated to detect clade-specific changes in the replacement rate (Type I functional divergence or FD) or the preferred amino acid at a site (Type II FD) (Gu, 1999, 2001, 2006; Gaston *et al.*, 2011). Both approaches (CSMs and ASMs) require *a priori* specification of the branches over which changes in the substitution process are thought to have occurred. This is often realized via informal use of external information such as phenotype.

Models that account for molecular and phenotypic evolution under a unified statistical framework have been proposed (Mayrose and Otto, 2011; Lartillot and Poujol, 2011; O'Connor and Mundy, 2013; Karin *et al.*, 2017). In CoEvolve (Lartillot and Poujol, 2011), for example, $\log(\omega)$ is assumed to have evolved continuously over the tree via Brownian motion and the model objective is to estimate correlations between it and other continuous variables, such as body size, longevity, and metabolic rate. Similarly, in TrateRateProp (Karin *et al.*, 2017) the objective is to determine whether a subset of nucleotide sites evolved under one of two substitution rates depending on the state of a binary

phenotype. Neither model appeals to mechanisms by which evolution of the phenotype might be linked to evolution of the gene. A novel approach is proposed in this chapter, the phenotype-genotype branch-site model (PG-BSM), the objective of which is to link phenomenological signatures of site-specific variations in ω (a.k.a. heterotachy, Lopez *et al.*, 2002) to specific mechanistic processes, including those that occurred in association with changes in a discrete character state (e.g., a phenotype).

The mutation-selection (MS) framework of Halpern and Bruno (1998) presented in Chapter 1 together with the notion of a site-specific fitness landscape presented in Chapter 2 (McCandlish, 2011; Jones *et al.*, 2017) provides a means to think about mechanistic processes that can give rise to heterotachy in real alignments. Under MS, each site is assumed to evolve independently with its own vector of fitness coefficients for the twenty amino acids (i.e., a site-specific fitness landscape). A site evolving on a static landscape can undergo chance fixation to a suboptimal amino acid followed by a period of positive selection that restores the site to its optimal state. This results in heterotachy via non-adaptive shifting balance as was demonstrated in Chapter 2. Heterotachy can also be caused by episodic changes in site-specific landscapes congruent with molecular adaptation, such as a change in the optimal amino acid (i.e., a peak shift) or a change in the stringency of selection at a site. Non-adaptive shifting balance and episodic changes in site-specific landscapes can both be represented phenomenologically as random switches between two rate ratios $\omega_1 < \omega_2$. It was shown in Chapter 2 that non-adaptive shifting balance on static fitness landscapes and episodic adaptive changes in landscapes can both manifest as transient elevations to $\omega_2 > 1$. It follows that the canonical $\omega > 1$ signature of positive selection does not necessarily provide unequivocal evidence of adaptation.

The PG-BSM was formulated to identify sites that likely underwent adaptive events without appealing to evidence of positive selection. It does so by making formal use of a discrete phenotype assigned to the terminal nodes of the tree. Branches over which the phenotype might have changed are determined by a distribution of ancestral phenotypic states at the internal nodes of

the tree derived from a model for phenotype evolution (cf. Karin *et al.*, 2017). It is assumed under the null hypothesis that all heterotachous sites evolved on static fitness landscapes independently of the phenotype and that their observed site patterns are consistent with the phenomenological CL process of random shifts between $\omega_1 < \omega_2$. The alternative model permits specific modes of switching between $\omega_1 < \omega_2$ that occurred in coordination with changes in the discrete phenotype. The modes are specified to be consistent with either a change in the stringency of selection or a change in the optimal amino acid at a site. The PG-BSM therefore represents a paradigm shift in both the information utilized (genotype and phenotype) and the evidence used to infer molecular adaptation (specific modes of heterotachy).

4.2 Materials and Methods

4.2.1 Background

The traditional way of characterizing codon evolution is to estimate the ratio of the nonsynonymous substitution rate dN to the synonymous substitution rate dS , accounting for differences in the rate at which nonsynonymous and synonymous mutations arise. Selection regimes are categorized according to $\omega = dN/dS$, such that $\omega < 1$ indicates a conservative regime, $\omega \approx 1$ a neutral regime, and $\omega > 1$ the canonical positive selection regime. CSMs can be used to infer $\omega > 1$ by contrasting a null model that allows sites to evolve under a set of ω -categories all constrained to be ≤ 1 with an alternate model that includes an additional category for sites with $\omega > 1$. Rejection of the null is interpreted as evidence that positive selection occurred somewhere in the gene. Subsequent analysis can be conducted to identify sites at which positive selection is most likely to have occurred (e.g., Yang and Nielsen, 1998).

The majority of CSMs are based on a continuous-time homogeneous and time-reversible Markov process that describes the rate at which nucleotide substitutions occur under a neutral regime (i.e., for which $dN/dS = 1$). As

we've seen in previous chapters, this can be represented by the following:

$$M_{ij} \propto \begin{cases} \kappa^{s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 1 \\ \alpha \kappa^{s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 2 \\ \beta \kappa^{s_t} \prod_{i_k \neq j_k} \pi_{j_k}^* & \text{if } s = 3 \end{cases} \quad (4.1)$$

Equation (4.1) applies to all pairs of codons (i, j) that differ by $s \in \{1, 2, 3\}$ nucleotides, s_t of which are transitions (substitutions of the form $T \leftrightarrow C$ or $A \leftrightarrow G$) and $s - s_t$ of which are transversions (substitutions of the form $\{T, C\} \leftrightarrow \{A, G\}$). $\pi_{j_k}^*$ is the frequency of the j^{th} nucleotide in the $k^{th} \in \{1, 2, 3\}$ codon position, κ the transition/transversion rate ratio, and α and β the rate at which double and triple (DT) substitutions arise, respectively. Diagonal elements M_{ii} are adjusted to make rows sum to zero. The selection process is parameterized by ω , which can be introduced into the model via an element-wise matrix product:

$$Q(\omega) = M \circ (\ell_S + \omega \ell_N) / r, \text{ where } r = \sum_{j \neq i} \pi_i Q_{ij}(\omega) \quad (4.2)$$

Equation 4.2 represents the simplest possible CSM and provides the basis for more sophisticated models such as those that take into account variations in ω across sites and/or over time. It was emphasized in Chapter 1 that $Q(\omega)$ is unsuitable as a means of thinking about the substitution process at a site. For example, the rate ratio ω , a proxy for the strength of selection for ($\omega > 1$) or against ($\omega < 1$) the i to j substitution, is assumed to be the same for all nonsynonymous (i, j) pairs. This is conceptually misleading for the majority of proteins because it implies that the fitness of an amino acid at a site is independent of its physicochemical properties. It is more useful to think of the evolutionary process at a codon site in terms of the dynamic on its site-specific fitness landscape as described in Chapter 2. If codon sites are assumed to evolve independently, a site-specific fitness landscape can be defined for the h^{th} site by a vector of fitness coefficients \mathbf{f}^h or its implied vector of equilibrium codon frequencies $\boldsymbol{\pi}^h$ (Sella and Hirsh, 2005). These determine the way it moves across its landscape over macroevolutionary time scales. Possible dynamic regimes include: non-adaptive shifting balance, under which the site moves episodically away from the peak of its static fitness landscape

(i.e., the fittest amino acid) via drift and back again by positive selection; adaptive evolution, under which a change in the landscape in the form of a peak shift is followed by movement of the site toward its new fitness peak; and neutral or nearly-neutral evolution, under which drift dominates and the site is free to move over a relatively flat landscape constrained primarily by biases in the mutation process.

The objective of many CSMs is to identify sites that likely underwent positive selection at some point in time. Among the most sophisticated CSMs developed for this purpose are the branch-site models. Two approaches have been considered. Under the fixed-effects approach (e.g., the branch-site models of Yang and Nielsen, 2002; Yang *et al.*, 2005; Zhang *et al.*, 2005, herein referred to as YN-BSM) branches over which positive selection is thought to have occurred (a.k.a. the foreground or FG) are specified *a priori*. The contrast between the null and alternative models provides a means to detect sites that switched from $\omega \leq 1$ on background (BG) branches to $\omega > 1$ on FG branches. The fixed-effect approach has good power when the FG is correctly specified (Yang and dos Reis, 2011). However, there is often no means to identify the correct FG, in which case the random-effect approach (Kosakovsky Pond *et al.*, 2011; Smith *et al.*, 2015; Murrell *et al.*, 2015) is applicable. Under the adaptive branch-site random effects likelihood model (aBSREL Smith *et al.*, 2015), a different ω -distribution is estimated independently for each branch of the tree. This is meant to account for variations in site-specific rate ratios from one branch to the next and can include episodic switches to $\omega > 1$. Rejection of the null under either the fixed-effects or random-effects BSM has traditionally been interpreted as evidence of adaptive evolution.

The amino acid model of Gu (1999, 2001, 2006) offers an alternative approach for detecting adaptive evolution that does not rely on evidence of positive selection. The model can be used to infer whether a protein had undergone a change in function following a duplication event, or what the author calls functional divergence (FD). A distinction is made between changes in substitution rate (Type I FD) and changes in the preferred amino acid at a site (Type II FD). Type I FD can either entail a change in the substitution

rate (e.g., an increase due to relaxation of selection pressure on one copy of the gene) without a peak shift, or a change in both rate and the preferred amino acid. Type II FD is associated with site patterns that are constant within clades but different between clades consistent with a peak shift at a site that otherwise evolved under stringent selection. Gu (2006) also makes a distinction between radical and conserved substitutions by partitioning the amino acids into four groups: positive charge (K, R, H); negative charge (D, E); hydrophilic (S, T, N, Q, C, G, P); and hydrophobic (A, I, L, M, F, W, V, Y). Substitutions between members of the same group are considered to be physicochemically conservative, while those between members of different groups are labeled radical. The underlying assumption is that a site that is conserved within each clade but radically different between clades is consistent with a change in function, and that the detection of such sites can be adduced as evidence of adaptive evolution.

The fixed-effects BSM for codon sequences requires specification of the FG branches. Similarly, ASMs for detecting FD require specification of the clades to be compared. In both cases this can be effected by informal use of information gleaned from tree topology and branch lengths combined with other sources of information such as phenotype. The random-effects BSM avoids this prerequisite but by doing so can be sensitive to positive selection by non-adaptive shifting balance as was shown in Chapter 2. The challenges posed by previous approaches can be minimized by combining alignment data with contextual information under a unified statistical framework. The potential of this approach is demonstrated in this chapter using the novel PG-BSM, the formulation of which was motivated by two considerations. First was the realization that heterotachy with episodic shifts to $\omega > 1$ is possible at a site without adaptation via non-adaptive shifting balance on a static fitness landscape (Jones *et al.*, 2017). Second was the idea that specific types of site patterns might be more indicative of adaptation (Gu, 1999, 2001; Pupko and Galtier, 2002; Philippe *et al.*, 2003; Gu, 2006). With these considerations in mind, the PG-BSM was formulated to distinguish site patterns consistent with non-adaptive shifting balance from those less likely to occur without changes

in their site-specific fitness landscapes.

4.2.2 The PG-BSM

The PG-BSM consists of three components : (1) a model for the evolution of the codon sequence; (2) a model for the evolution of a discrete phenotype; and (3) a model that accounts for the mechanism(s) by which (1) and (2) are associated. The model chosen for the evolution of the phenotype is analogous to F81 for DNA (Felsenstein, 1981) as characterized by the following rate matrix (in this case assuming $k = 3$ distinct phenotypes):

$$Q_F = \frac{\lambda}{r_F} \begin{bmatrix} * & \pi_F^2 & \pi_F^3 \\ \pi_F^1 & * & \pi_F^3 \\ \pi_F^1 & \pi_F^2 & * \end{bmatrix} \quad \text{where } r_F = \sum_{i \neq j} \pi_F^i \pi_F^j \quad (4.3)$$

Each $*$ in (4.3) indicates a value that makes the corresponding row sum to zero. The π_F^i represent the stationary frequencies for the three phenotypic states $i \in \{1, 2, 3\}$. The scaling factor r_F is included so that the rate constant λ gives the expected number of changes in the discrete phenotypic state per unit branch length. Equation (4.3) is used not only to compute the probability of the vector of observed phenotypes, but also to generate samples from the distribution of ancestral phenotypes that are required to inform the mechanism(s) of PG association described below.

For sequence evolution it was decided to use a model that accounts for both invariant sites (those with no observed replacement substitutions) and sites with some degree of heterotachy (covarion-like sites that exhibit at least one replacement substitution). The null PG-BSM assumes that some proportion π_0 of sites evolved under $\omega_0 = 0$ over the tree while the remaining sites switched randomly between $\omega_1 < \omega_2$ over time under the covarion-like model CLM3($k = 2$). Alignments typically exhibit variations in rate ratio across sites (RAS) in addition to site-specific variations over time. It is possible to account for RAS using an M-series model such as M3($k = 2$) (Yang *et al.*, 2000a), which assumes sites evolved under either ω_1 or ω_2 over the entire tree without heterotachy. However, by accounting for variable amino acid site patterns via the process of random switching between ω_1 and ω_2 over time, the covarion-like model

implicitly accounts for variations in site-specific time-averaged rate ratios (cf. Wu and Susko, 2009). Hence, with only one extra parameter (the switching rate δ), CLM3($k = 2$) captures both heterotachy and RAS, and consequently often provides a better fit to real alignments compared to M3($k = 2$) (e.g., see Table 3.1). Furthermore, the CLM3($k = 2$) component of the PG-BSM provides a means to account for non-adaptive shifting balance, and in doing so reduces the probability of falsely rejecting the null hypothesis of no PG association (as shown in the analysis of the phytochrome A&CF gene).

Genetic information is assumed to consist of an alignment X of N homologous protein-coding sequences of length n with a known rooted topology τ . The phenotype, encoded by a vector \mathbf{F} , can be any discrete character state, such as a property of the gene's protein product (i.e., a molecular phenotype), some characteristic of the organism, or an environmental variable. Given τ , the likelihood function under the null hypothesis that the phenotype and genotype evolved independently is computed as follows:

$$L_{\text{null}}(X, \mathbf{F}; \lambda, \boldsymbol{\theta}, \mathbf{t}) = P(\mathbf{F}; \lambda, \mathbf{t}) \prod_{h=1}^n (\pi_0 P_0(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) + (1 - \pi_0) P_{\text{CL}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t})) \quad (4.4)$$

$P(\mathbf{F}; \lambda, \mathbf{t})$ is the probability of the vector of phenotypes given the rate parameter λ and branch lengths \mathbf{t} , $P_0(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t})$ is the probability of the site pattern \mathbf{x}^h assuming the site evolved under $\omega_0 = 0$, $P_{\text{CL}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t})$ is the probability of the site pattern assuming it evolved under CLM3($k = 2$), and $\boldsymbol{\theta} = (\omega_1, \omega_2, p_1, \delta, \kappa)$ is a vector of parameters for sequence evolution. All probabilities are computed using the pruning algorithm (Felsenstein, 1981).

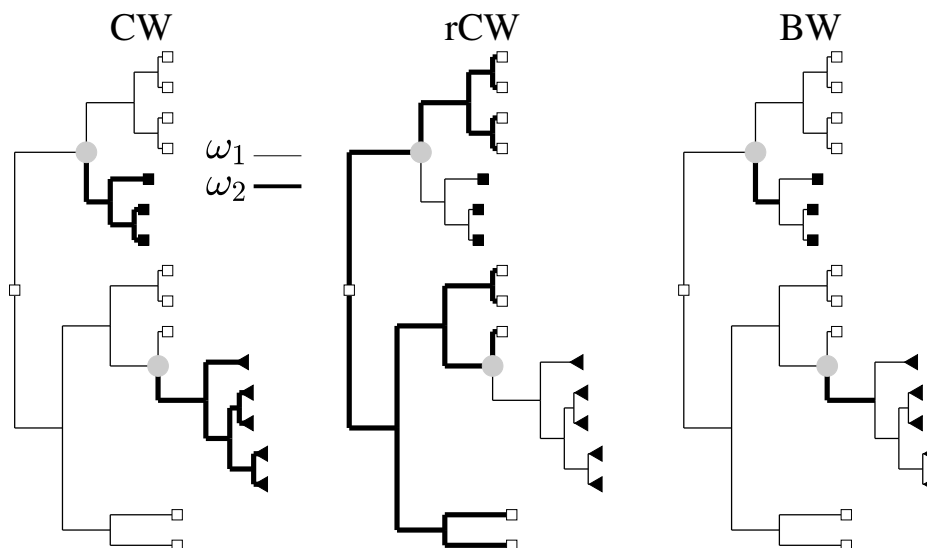


Figure 4.1: An illustration of the difference between the cladewise (CW and rCW) and branchwise (BW) evolutionary processes. Each process accounts for a specific form of heterotachy associated with changes in phenotype. The empty and filled markers at the terminal nodes indicate three phenotypic states. $\omega_1 < \omega_2$ are dN/dS rate ratios. The grey disks indicate the nodes at which a change in phenotype occurred. CW sites are assumed to have evolved under ω_1 prior to a change in phenotype and under ω_2 after a change. rCW sites are assumed to have evolved under ω_2 prior to a change in phenotype and under ω_1 after a change. BW sites are assumed to have evolved under ω_2 over branches along which a change in phenotype occurred and under ω_1 everywhere else in the tree. The model assumes a rooted tree because the interpretation of the CW and rCW processes require a particular order of change in rate ratio.

The alternative PG-BSM enforces dependencies between phenotype and genotype evolution at some fraction of sites. Various mechanisms of dependency are amenable to phenomenological representation as distinct modes of heterotachy (Figure 4.1). Here I consider three. First, a change in phenotype along a branch can coincide with a reduction in the stringency of selection at a site due to loss of functional importance (e.g., the site's role in the maintenance of the protein's tertiary structure, Pupko and Galtier, 2002) in the descendant clade. This mechanism is expressed phenomenologically in the PG-BSM as the cladewise (CW) process under which a proportion π_{CW} of sites are assumed to have evolved under the smaller ω_1 prior to a change in phenotype and under the larger ω_2 over the entire clade descending from the branch over which a change in phenotype occurred (CW tree in Figure 4.1).

Second, a change in phenotype along a branch can coincide with an increase in the stringency of selection at a site reflecting an increase in its functional importance. This mechanism is represented in the PG-BSM by the reverse cladewise (rCW) process under which a proportion π_{rCW} of sites are assumed to have evolved under the larger ω_2 prior to a change in phenotype and under the smaller ω_1 over the entire clade descending from the branch over which a change in phenotype occurred (rCW tree in Figure 4.1). Third, a change in phenotype can coincide with changes in site-specific fitness landscapes in the form of peak shifts. This mechanism is represented in the PG-BSM by the branchwise (BW) process under which a proportion π_{BW} of sites are assumed to have evolved under the larger ω_2 over branches along which the phenotype changed and under the smaller ω_1 everywhere else in the tree (BW tree in Figure 4.1).

Sites consistent with the phenomenological CW, rCW or BW process (herein referred to as CW, rCW or BW sites) represent a subset of those assumed by the null PG-BSM to have evolved under the CL process (i.e., sites deemed to have undergone at least one replacement substitution). It is therefore assumed that all four processes (CW, rCW, BW and CL) share the same ω_1 and ω_2 . This is in contrast to the standard approach exemplified by the YN-BSM A. That model partitions sites into four categories according to the way they are assumed to have evolved along background (BG) and foreground (FG) branches. Category 0 sites evolved under $\omega_0 < 1$ over the entire tree; category 1 sites evolved neutrally with $\omega_1 = 1$ over the entire tree; category 2a sites evolved under $\omega_0 < 1$ on BG branches and under $\omega_2 \geq 1$ on the FG; and category 2b sites evolved under $\omega_1 = 1$ on the BG and under $\omega_2 \geq 1$ on the FG. Hence, ω_2 applies to category 2 sites only. This approach gives the YN-BSM A the power to detect evidence of positive selection (i.e., $\omega_2 > 1$) at a small number of sites along the FG, but also introduces the risk of issues related to irregularity (e.g., Baker *et al.*, 2016; Mingrone *et al.*, 2018). Specifically, ω_2 becomes nearly unidentifiable when the proportion p_2 of category 2 sites is small. Its maximum likelihood estimate (MLE) can consequently be very large and potentially misleading. Using ω_1 and ω_2 for all non-invariant sites

under the PG-BSM makes it much less likely that these parameters will be unidentifiable, but undoubtedly reduces the statistical power to detect a small number of sites that evolved under an exceptionally large rate ratio. The potential impact of this loss is mitigated by the fact that the PG-BSM does not rely on evidence of $\omega > 1$ to reject the null.

To infer PG associations of any kind requires knowledge of the ancestral branches over which the phenotype changed or remained constant. This information is provided by realizations of ancestral phenotypes at the internal nodes of the tree generated using equation (4.3), each of which is converted to a change map $\mathbf{z} = (z_1, \dots, z_{2N-2})$ as follows (see Figure 4.1):

$$z_b = \begin{cases} 0 & \text{if the phenotype is the same at the two ends of branch } b \\ 1 & \text{if the phenotype is different at the two ends of branch } b \end{cases} \quad (4.5)$$

The likelihood function under the alternative hypothesis can in principle be computed by summing over all possible change maps:

$$L_{\text{alt}}(X, \mathbf{F}; \lambda, \boldsymbol{\theta}, \mathbf{t}) = P(\mathbf{F}; \lambda, \mathbf{t}) \sum_{\mathbf{z}} \left\{ P(\mathbf{z} | \mathbf{F}, \lambda, \mathbf{t}) \prod_{h=1}^n g(\mathbf{x}^h; \mathbf{z}) \right\} \quad (4.6)$$

where $P(\mathbf{z} | \mathbf{F}, \lambda, \mathbf{t})$ is the probability of the change map \mathbf{z} given the terminal states \mathbf{F} , the rate constant λ and a vector of branch lengths \mathbf{t} . The mixture probability $g(\mathbf{x}^h; \mathbf{z})$ depends on the particular combination of processes included in the alternate PG-BSM. If the objective was to detect sites consistent with either the CW or BW process, for example, the mixture probability would be as follows:

$$\begin{aligned} g(\mathbf{x}^h; \mathbf{z}) &= \pi_0 P_0(x^h; \boldsymbol{\theta}, \mathbf{t}) + \pi_{\text{CL}} P_{\text{CL}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) \\ &+ \pi_{\text{CW}} P_{\text{CW}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}, \mathbf{z}) + \pi_{\text{BW}} P_{\text{BW}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}, \mathbf{z}) \end{aligned} \quad (4.7)$$

Here $P_{\text{CW}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}, \mathbf{z})$ and $P_{\text{BW}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}, \mathbf{z})$ give the probabilities of the site pattern \mathbf{x}^h assuming the site evolved under the CW and BW process, respectively, and $\pi_0 + \pi_{\text{CL}} + \pi_{\text{CW}} + \pi_{\text{BW}} = 1$.

The number of possible change maps can be very large depending on the number of taxa and phenotypic states. It is therefore often infeasible to compute equation (4.6) exactly by summing over all \mathbf{z} . Instead, $P(\mathbf{z} | \mathbf{F}, \lambda, \mathbf{t})$

was approximated by the relative frequency $\hat{\pi}_{\mathbf{z}}$ of \mathbf{z} in a sample of 10^5 realizations of ancestral phenotypes. These were generated using the MLEs for λ and \mathbf{t} as described in Methods. The summation in equation (4.6) was then over all unique change maps that appeared in the sample using $\hat{\pi}_{\mathbf{z}}$ in place of $P(\mathbf{z} \mid \mathbf{F}, \lambda, \mathbf{t})$. Computation time increases with the number of unique change maps. To reduce computational load, change maps \mathbf{z} that occurred with a frequency $< 10^{-3}$ were excluded and the probabilities $\hat{\pi}_{\mathbf{z}}$ were renormalized to sum to one. Note that estimates of λ and \mathbf{t} are required to generate ancestral phenotypes. An exact but costly approach would be to resample from $P(\mathbf{z} \mid \mathbf{F}, \lambda, \mathbf{t})$ with each iterative update of λ and \mathbf{t} inside the optimization function. A less costly approximation was implemented instead as follows. The MLEs $(\hat{\lambda}, \hat{\mathbf{t}})$ obtained by fitting the null PG-BSM to the data were first used to generate a preliminary sample. The alternate PG-BSM was then fitted to the data using this sample to produce new MLEs $(\hat{\lambda}, \hat{\mathbf{t}})$. To account for any resulting changes in $(\hat{\lambda}, \hat{\mathbf{t}})$, a second sample was generated using the new MLEs and the alternate PG-BSM was fitted once more to produce the final results. The alternate PG-BSM characterized by equation (4.7) was fitted to all simulated alignments and to the real data. In some analyses of real data the alternate model was modified to detect the phenomenological signature of a single mechanism of PG association alone (e.g., either the CW, rCW or BW process by itself).

An omnibus log-likelihood ratio (LLR) test is conducted to contrast the null and alternative components of the PG-BSM. The components can differ by $m \in \{1, 2, 3\}$ parameters among the proportions $\{\pi_{\text{CW}}, \pi_{\text{rCW}}, \pi_{\text{BW}}\}$ depending on which mechanistic processes of PG association are included in the alternate model. In all cases the null PG-BSM is the same as the alternative when the proportions are on the boundary of the parameter space (i.e., when $\pi_{\text{CW}} = \pi_{\text{rCW}} = \pi_{\text{BW}} = 0$). The theoretical limiting distribution of the LLR is therefore a 50:50 mixture of the χ_0^2 and χ_1^2 distributions when $m = 1$ and an unknown mixture of the $\chi_0^2, \chi_1^2, \dots, \chi_m^2$ distributions when $m \in \{2, 3\}$ (Self and Liang, 1987). Note that Self and Liang (1987) does not apply to mixture models in general, but can be shown to apply in this case because ω_1 and ω_2 can

always be estimated under the CL component of the null model. Although the mixture weights for the LLR distribution are unknown when $m \in \{2, 3\}$, the 95th percentile of a mixture of the $\chi_0^2, \chi_1^2, \dots, \chi_m^2$ distributions is always at most that for the χ_m^2 distribution (i.e., 3.84, 5.99, and 7.81 for $m = 1, m = 2$, and $m = 3$, respectively). Using these as critical values for the omnibus test should therefore be conservative and produce less than 5% type I error rate.

Rejection of the null hypothesis of the omnibus test provides evidence for an association between the gene and the phenotype. Naive empirical Bayes (NEB) analysis is then used *post hoc* to identify the most likely category for each site. Let $c \in \{\text{CW}, \text{rCW}, \text{BW}\}$ index the three categories of PG-associated sites possibly included in the alternative model. The posterior probability that a site evolved under the process indicated by category c is evaluated at the MLE for the alternate PG-BSM as follows:

$$P(c | \mathbf{x}^h) = \frac{L_{\text{alt}}(\mathbf{x}^h | c) \hat{\pi}_c}{L_{\text{alt}}(\mathbf{x}^h)} \quad (4.8)$$

A false positive occurs when a site is incorrectly assigned to one of the three categories. The false positive rate (the proportion of all sites falsely inferred to be associated with the phenotype) is usually controlled by assigning sites to category $c \in \{\text{CW}, \text{rCW}, \text{BW}\}$ only when their posteriors $P(c | \mathbf{x}^h)$ are greater than some threshold such as 0.95 (e.g., Yang *et al.*, 2000a). An alternative approach, also based on the posteriors, is to aim to control the proportion of sites assigned to category c that in fact did not evolve under category c (i.e., the false discovery rate or FDR, Newton *et al.*, 2004; Guindon *et al.*, 2006). To see the difference between the two approaches, consider an analysis of an alignment with 1000 codon sites. Suppose 10 sites were inferred to have evolved under the CW process (i.e., there were 10 “discoveries”) and that 5 of these were incorrect. Then the false positive rate would only be 0.5% (5 sites out of 1000), whereas the FDR would be 50% (5 sites out of 10). Hence, a low false positive rate does not necessarily imply a low FDR, particularly when the number of discoveries is small.

The FDR approach was used in all analyses but with one modification. Rather than controlling the *proportion* of false discoveries of a given category

$c \in \{\text{CW}, \text{rCW}, \text{BW}\}$, it was decided to control the *number* of false discoveries or the “false discovery counts” (FDC) for each category. For example, to assess the results of the *post hoc* analyses for our simulation studies, the mean observed FDC (and the mean power) were computed for each set of S alignments as follows:

$$\text{FDC}(c) = \frac{1}{S} \sum_{i=1}^S F_i(c), \text{ Power} = \frac{1}{S} \sum_{i=1}^S \frac{D_i(c) - F_i(c)}{n(c)} \quad (4.9)$$

Here $F_i(c)$ is the number of false discoveries of category c sites and $D_i(c)$ the total number of discoveries of category c sites, both for the i^{th} alignment, and $n(c)$ is the number of sites in the alignment that were evolved under the process indicated by category c . The expected number of false discoveries $\text{E}\{\text{FDC}\} = \text{E}\{F_i(c)\}$ can be controlled by setting a posterior threshold that is specific to the alignment under consideration (see Methods). This threshold can change from one data set to another. For the analyses presented in this chapter I used $\text{E}\{\text{FDC}\} \in \{1, 2\}$ for each category $c \in \{\text{CW}, \text{rCW}, \text{BW}\}$ of sites included in the alternate PG-BSM.

To quantify the evidential support for branches over which the phenotype is thought to have changed, the probability of the most frequently sampled change map \mathbf{z}^* conditioned on the combined data is estimated as follows:

$$\text{P}(\mathbf{z}^* \mid X, \mathbf{F}, \lambda, \mathbf{t}) = \frac{\text{L}_{\text{alt}}(X, \mathbf{F} \mid \mathbf{z}^*) \hat{\pi}_{\mathbf{z}^*}}{\text{L}_{\text{alt}}(X, \mathbf{F})} \quad (4.10)$$

where $\hat{\pi}_{\mathbf{z}^*}$ is the frequency of the most commonly sampled change map. Equation (4.10) is evaluated at the MLE for the alternate PG-BSM. The algorithm that generates realizations of ancestral phenotypes makes use of estimates of λ and \mathbf{t} . But λ is independent of the alignment. The observed frequency $\hat{\pi}_{\mathbf{z}^*}$ therefore depends on X only through branch length estimates. The likelihood $\text{L}_{\text{alt}}(X, \mathbf{F} \mid \mathbf{z}^*)$, by contrast, also depends on the existence of individual site patterns that match to greater or lesser degree the patterns of PG association indicated by \mathbf{z}^* . An alignment generated with no PG association will tend to result in $\text{L}_{\text{alt}}(X, \mathbf{F} \mid \mathbf{z}^*)/\text{L}_{\text{alt}}(X, \mathbf{F}) \approx 1$ making $\text{P}(\mathbf{z}^* \mid X, \mathbf{F}, \lambda, \mathbf{t}) \approx \hat{\pi}_{\mathbf{z}^*}$, in concordance with an alignment that contains no information about ancestral phenotypic states. When there is PG association at some sites the ratio of

likelihoods will be greater than one depending on the strength of the signal for PG association (e.g., the proportion of sites generated under the CW, rCW and BW processes) and on the extent to which \mathbf{z}^* matches the change map that generated the data. A good match combined with strong signal will tend to result in $L_{\text{alt}}(X, \mathbf{F} \mid \mathbf{z}^*)/L_{\text{alt}}(X, \mathbf{F}) > 1$ making $P(\mathbf{z}^* \mid X, \mathbf{F}, \lambda, \mathbf{t}) > \hat{\pi}_{\mathbf{z}^*}$. Hence, confidence in the most probable ancestral history can be increased by accounting for PG associations in the combined data when such associations exist.

4.2.3 Rigorous Model Assessment Requires a Realistic Data Generating Process

The accuracy and power of a new CSM is usually assessed by fitting the model to alignments generated under a similar model (Anisimova *et al.*, 2001, 2002; Wong *et al.*, 2004; Zhang, 2004; Kosakovsky Pond and Frost, 2005; Yang *et al.*, 2005; Zhang *et al.*, 2005; Yang and dos Reis, 2011; Kosakovsky Pond *et al.*, 2011; Lu and Guindon, 2013). One drawback to this approach is that standard CSMs constructed from rate matrices of the form $Q(\omega) = M \circ (\ell_S + \omega \ell_N)$ cannot mimic site-specific variations in ω caused by processes such as adaptation following episodic peak shifts (dos Reis, 2015) or non-adaptive shifting balance (Jones *et al.*, 2017). This is an issue because heterotachy may well be pervasive in real alignments (e.g., Fitch and Markowitz, 1970; Fitch, 1971; Lopez *et al.*, 2002; Philippe *et al.*, 2003; Wang *et al.*, 2007; Whelan *et al.*, 2011) and can engender novel statistical pathologies in models fitted to such data (e.g., phenomenological load, Jones *et al.*, 2018). It follows that rigorous assessment of a CSM requires fitting the model to alignments generated with realistic levels of heterotachy.

A direct way to produce such alignments is to base the generating model on the MS framework (Halpern and Bruno, 1998; Spielman and Wilke, 2015a,b; Jones *et al.*, 2017). Two such generating models were used in this study. The first, MSm(ammalian)mtDNA, was developed in Chapter 3 to mimic 12 concatenated H-strand mitochondrial DNA sequences (3331 codon sites) from 20 mammalian species as distributed in the PAML software package (Yang,

2007). MSmmtDNA was shown to produce data similar to the real alignment by several measures of comparison (Jones *et al.*, 2018). In particular, it was shown to produce alignments with similar levels of heterotachy. However, while MSmmtDNA is more realistic as a generating model, it represents only a small proportion of the space of all distributions of vectors of site-specific fitness coefficients that might arise in reality. A second generating model that samples with replacement from a set of 3598 such vectors estimated from an alignment of 12 mitochondrial genes taken from 244 mammalian species (Tamuri *et al.*, 2014) was therefore used for some of simulations. This generating process will be referred to as MSTGdR (Tamuri, Goldman and dos Reis) after the authors of that study. Substitutions between codons that differ by two or three nucleotides can only occur in single nucleotide steps under the PG-BSM, consistent with the majority of CSMs in common use today. It was demonstrated in Chapter 3 that the occasional fixation of double or triple mutations can manifest as an additional source of heterotachy. It was not clear what effect DT might have on power and accuracy when left unaccounted for by the PG-BSM. Alignments were therefore generated using both MSmmtDNA with 0% DT and MSmmtDNA with 6% DT (recent studies suggest that DT mutations comprise between 1% and 3% of all mutations, Keightley *et al.*, 2009; Schrider *et al.*, 2014; De Maio *et al.*, 2013; Harris and Nielsen, 2014). To assess the statistical properties of the PG-BSM without misspecification the null PG-BSM was also used to generate alignments.

The mutation-selection framework has been used in several recent studies to simulate alignments (Spielman and Wilke, 2015b, 2016; Jones *et al.*, 2017, 2018). In all cases the substitution process was stationary, meaning that fitness coefficients and the stringency of selection were made to be constant at each site over the entire tree. In this study, MSmmtDNA and MSTGdR were formulated to include a subset of sites evolved under non-stationary processes in the form of changes in the stringency of selection and/or fitness coefficients at specific nodes of the tree (Figure 4.1). Changes in the stringency of selection starting along a single branch leading to a clade can manifest as a

cladewise difference in the mean rate ratio ω to produce site patterns phenomenologically consistent with the CW or rCW processes. Similarly, changes in fitness coefficients (a peak shift) along a single branch can result in site patterns phenomenologically consistent with the BW process, particularly if they occur at sites otherwise evolved under stringent selection. In this way, the MS framework was used to produce alignments with realistic levels of heterotachy due to multiple processes. The purpose was to assess the ability of the PG-BSM to distinguish heterotachy due to non-adaptive processes from that due to changes in site-specific landscapes. The proposed approach represents a significant improvement over traditional methods of model testing based on data generated using rate matrices of the form $Q(\omega) = M \circ (\ell_S + \omega \ell_N)$. Details of all generating processes are provided in Methods.

4.3 Results

4.4 Simulations

The results of three broad simulation studies encompassing a wide variety of evolutionary scenarios are reported in this section. Simulation 1 was designed to test the statistical properties of the PG-BSM by fitting the model to alignments generated under the null PG-BSM. Simulation 2 was conducted to assess the impact of differences between the process as assumed under the fitted model and the process used to generate the data (i.e., misspecification). For this purpose, alignments were generated using MSmmtDNA with 0% DT or 6% DT and MSTGdR with 0% DT. Furthermore, in some cases the alternate PG-BSM was fitted with the phenotype designated incorrectly. Simulation 3 was designed to assess the performance of the model under a scenario with four phenotypes. Increasing the number of phenotypes introduces greater uncertainty in the distribution of change maps (i.e., $P(\mathbf{z} \mid \mathbf{F}, \lambda, \mathbf{t})$ becomes more dispersed) and therefore represents a greater challenge to the model. The analyses of real alignments that follow this section includes further simulation studies designed to mimic specific aspects of those data sets.

4.4.1 Simulation 1: Generating under the Null PG-BSM

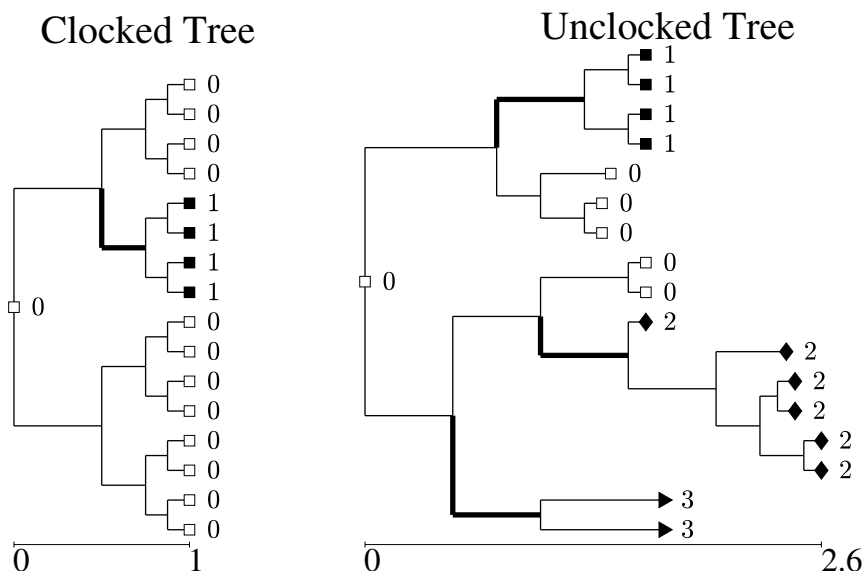


Figure 4.2: The clocked and unlocked trees used in Simulations 1, 2, and 3. Tree depths give the expected number of single nucleotide substitutions per codon. Symbols at the nodes indicate different phenotypes, with $k = 2$ phenotypes on the clocked tree and $k = 4$ on the unlocked tree.

According to ML theory, when the PG-BSM is fitted to data generated under the null PG-BSM, and as information (e.g., the number of codon sites) increases without bound, (i) the distribution of the LLR for the contrast between the null (equation 4.4) and alternate PG-BSM (equations 4.6 and 4.7) will converge to some unknown mixture of the χ_0^2 , χ_1^2 and χ_2^2 distributions (Self and Liang, 1987), and (ii) the distribution of the MLE for each model parameter will converge to a normal centered on the parameter's generating value. The objective of the first simulation was to assess how well these expectations hold. To that end, 100 alignments 300 codons in length and 100 alignments 1000 codons in length were generated on the clocked and unlocked trees shown in Figure 4.2. The generating model was the null PG-BSM with the following parameters: $\pi_0 = 0.65$, $\omega_1 = 0.10$, $\omega_2 = 1.50$, $p_1 = 0.80$, $\delta = 0.20$, $\pi_{CW} = 0$, $\pi_{BW} = 0$. The parameters for the mutation process, including position-specific nucleotide frequencies and the transition/transversion rate ratio, were set to

values estimated from the alignment of 12 concatenated H-strand mitochondrial DNA sequences from 20 mammalian species. The phenotypes assumed under the alternate PG-BSM were those indicated by the different symbols at the terminal nodes in Figure 4.2 (e.g., $\mathbf{F} = (0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0)$ for the clocked tree and $\mathbf{F} = (1, 1, 1, 1, 0, 0, 0, 0, 2, 2, 2, 2, 2, 2, 3, 3)$ for the unclocked tree).

As the limiting distribution of the LLR is unknown, the test was conducted as if it was χ_2^2 to be conservative, using the critical value 5.99 for a 5% test. Table 4.1 compares the relative frequency of the empirical LLR in each of the three intervals $[0, 0.50)$, $[0.50, 5.99)$ and $(5.99, +\infty)$ for all four simulation scenarios to that expected under the χ_2^2 distribution. Using the 5.99 cut-off gave a false positive rate of at most 2/100 among Simulation 1 scenarios. Furthermore, the relative frequencies in the $[0, 0.50)$ interval fell between 0.65 and 0.78 compared to the expected probability 0.22 for the χ_2^2 distribution. This result is not inconsistent with the fact that the actual LLR distribution places a substantial weight of 1/2 on the χ_0^2 distribution (i.e., a point mass at zero) and $1/2-p$ on the χ_1^2 distribution for some unknown $p \in [0, 1/2)$. Nevertheless, the χ_2^2 distribution was used for the remainder of analyses as a buffer against inevitable misspecifications and/or issues associated with low information content, as is standard practice when the exact distribution of the LLR is unknown (e.g., Wong *et al.*, 2004; Zhang *et al.*, 2005; Yang, 2007, 2017).

The mean, median and standard deviation of the MLEs of select model parameter for each generating scenario are shown in Table 4.2. In each case the mean and median were either the same or nearly so, indicating symmetrical distributions. A one-sample Kolomogorov-Smirnov test for normality applied to each set of 100 MLEs failed to reject the null hypothesis of a normal distribution in all cases (p-value ≥ 0.16). Furthermore, the mean MLE for each parameter was either the same or very close to its generating value in all four scenarios. And in each case the standard deviation was smaller for the 1000 codon scenario compared to its counterpart 300 codon scenario. These results suggest that the PG-BSM is statistically well behaved when fitted to

alignments generated under the scenarios considered.

scenario	0 to 0.50	0.50 to 5.99	5.99 to $+\infty$	false +ve
PG-BSM C 300 sites	0.66	0.32	0.02	2/100
PG-BSM UC 300 sites	0.68	0.32	0.00	0/100
PG-BSM C 1000 sites	0.65	0.34	0.01	1/100
PG-BSM UC 1000 sites	0.78	0.22	0.00	0/100
expectation under χ_2^2	0.22	0.73	0.05	5/100

Table 4.1: Empirical vs Expected. Comparison of the empirical LLR with χ_2^2 for Simulation 1 scenarios on the clocked (C) and unclocked (UC) trees shown in Figure 4.2. The last column shows the number of times the omnibus test incorrectly rejected the null to give a false positive (false +ve).

parameter	generating	C 300	C 1000
π_0	0.65	0.65/0.65,(0.03)	0.65/0.65,(0.01)
ω_1	0.10	0.10/0.10,(0.03)	0.10/0.10,(0.01)
ω_2	1.50	1.51/1.50,(0.25)	1.57/1.54,(0.18)
p_1	0.80	0.78/0.78,(0.03)	0.83/0.83,(0.02)
δ	0.20	0.19/0.18,(0.07)	0.21/0.20,(0.04)
κ	4.61	4.71/4.66,(0.42)	4.53/4.53,(0.20)
parameter	generating	UC 300	UC 1000
π_0	0.65	0.65/0.65,(0.01)	0.65/0.65,(0.01)
ω_1	0.10	0.10/0.10,(0.02)	0.10/0.10,(0.01)
ω_2	1.50	1.56/1.55,(0.26)	1.51/1.52,(0.14)
p_1	0.80	0.80/0.80,(0.03)	0.79/0.79,(0.02)
δ	0.20	0.20/0.19,(0.06)	0.20/0.19,(0.03)
κ	4.61	4.60/4.63,(0.37)	4.54/4.54,(0.19)

Table 4.2: Mean/median,(standard deviation) of select MLEs for Simulation 1.

4.4.2 Simulation 2: Generating under MSmmtDNA and MSTGdR

The second simulation study was conducted to assess the statistical accuracy and power of the PG-BSM when fitted to alignments simulated using a more complex generating model compared to the null PG-BSM. In particular, the aim was to generate alignments using the MS framework in such a way as to mimic realistic levels of heterotachy caused by non-adaptive shifting balance and episodic changes in site-specific fitness landscapes. The simulation is comprised of five scenarios, each of which was tested under three different sequence generating processes, yielding 15 cases in total (Table 4.3). In the

first scenario (denoted 2a $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 0\%)$) alignments were generated with no PG association, but with substantial heterotachy due to non-adaptive shifting balance, and therefore contained signal for the CL process that could potentially be misconstrued as signal for the CW and BW processes under the alternate PG-BSM. The second scenario (denoted 2b $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 0\%)$) included signal in the form of a small fraction ($\pi_{\text{CW}} = 5\%$ of 300 sites) of sites generated with a reduction in the stringency of selection. The third scenario (denoted 2c $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 5\%)$) included sites generated with peak shifts ($\pi_{\text{BW}} = 5\%$). In the fourth scenario (denoted 2d $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 0\%)$) the effect of phenotype error was investigated by using the data generated for 2c but with a misspecified vector of phenotypes. Signal for PG association was increased in the final scenario (denoted 2e $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 5\%)$) by including both sites generated with a reduction in the stringency of selection ($\pi_{\text{CW}} = 5\%$) and sites with peak shifts ($\pi_{\text{BW}} = 5\%$). The three generating processes used were: MSmmtDNA with 0% DT, MSmmtDNA with 6% DT, and MSTGdR with 0% DT. In each case, 50 alignments 300-codons in length were generated on the clocked tree in Figure 4.2. Changes in the stringency of selection and/or peak shifts were effected along the branch marked in bold. The correct phenotype designation was $\mathbf{F} = (0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0)$, while the incorrect designation used in 2d was $\mathbf{F} = (0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0)$.

The omnibus test correctly failed to reject the null in all Scenario 2a $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 0\%)$ trials (Table 4.3). In Scenario 2b $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 0\%)$, the null was correctly rejected in 47/50, 46/50 and 42/50 trials, indicating good power, and the CW and BW processes were inferred at an average of (7%,1%), (7%,1%) and (4%,0%) of sites, in approximate agreement with their generating values. The agreement was also good in Scenario 2c $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 5\%)$ where the null was rejected in 46/50, 38/50 and 50/50 trials, and the CW and BW processes were inferred at an average of (1%,9%), (2%,7%) and (1%,7%). The null was rejected in only 1/50, 1/50 and 2/50 trials, in Scenario 2d $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 0\%)$, well below the expected 5% error rate. And in Scenario 2e $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 5\%)$ the null was rejected in all trials and the CW and BW processes were inferred at an average of (6%,9%), (6%,8%) and

(6%,6%) of sites.

Results of the *post hoc* analysis applied to alignments with signal for PG association (Scenarios 2b, 2c, and 2e) are summarized in Table 4.4. Analyses were conducted with $E\{\text{FDC}\} = 1$ for each category $c \in \{\text{CW}, \text{BW}\}$. For Scenario 2b ($\pi_{\text{CW}}, \pi_{\text{BW}} = (5\%, 0\%)$) MSmmtDNA 0% DT (first row of Table 4.4), for example, the average FDC was 0.80 CW sites and 0.86 BW sites per alignment. The average power to detect CW sites was 0.31 corresponding to an average of 4.72 correctly identified sites per alignment out of the 15 sites generated with a reduction in the stringency of selection. Among all scenarios included in Table 4.4, the average FDC ranged between 0.62 and 2.12 with mean 1.32 and standard deviation 0.44, and were approximately normal in distribution. The FDC was therefore slightly biased toward values greater than the nominal $E\{\text{FDC}\} = 1$. Note that $\text{FDC} = 1.32$ corresponds to a mean false positive rate of $1.32/300 \text{ sites} \times 100\% = 0.44\%$.

To determine the impact of the generating scenario on the number of false discoveries, a generalized linear model was fitted to the counts of each type for each alignment using as predictor variables: the generating MS model ($x_{\text{MS}} = 1, 2$ or 3 for MSmmtDNA 0% DT, MSmmtDNA 6% DT, and MST-GdR 0% DT); the presence or absence of sites evolved under the CW process ($x_{\text{CW}} = 1$ or 0); and the presence or absence of sites evolved under the BW process ($x_{\text{BW}} = 1$ or 0). The generating MS model was found to have no significant effect on either the expected $\text{FDC}(\text{CW})$ (p-value = 0.08) or the expected $\text{FDC}(\text{BW})$ (p-value = 0.12). After removing x_{MS} from the analysis, the following significant relationships were found:

$$\log \{\text{FDC}(\text{CW})\} = -0.37 + 0.80x_{\text{BW}} \quad (4.11)$$

$$\log \{\text{FDC}(\text{BW})\} = -0.27 + 0.34x_{\text{CW}} + 0.51x_{\text{BW}} \quad (4.12)$$

For alignments generated with no PG association (i.e., with $x_{\text{CW}} = x_{\text{BW}} = 0$) the first model predicts an average $\text{FDC}(\text{CW}) = 0.69$ per alignment and the second model an average $\text{FDC}(\text{BW}) = 0.76$ per alignment, both slightly under the target $E\{\text{FDC}\} = 1$. The first model predicts an average $\text{FDC}(\text{CW}) = 1.54$ for alignments generated with BW sites whether or not CW sites were present.

The second model predicts $\text{FDC}(\text{BW}) = 1.07$ for alignments generated with CW sites only and $\text{FDC}(\text{BW}) = 1.27$ for alignments generated with BW sites only. The predicted rate is nearly two false detections per alignment when the generating process includes both CW and BW sites, with $\text{FDC}(\text{BW}) = 1.78$ per alignment, corresponding to a false positive rate of $1.78/300 \text{ sites} \times 100\% = 0.59\%$.

Turning to the last two columns of Table 4.4, the model performed well with respect to identifying the correct evolutionary history of the phenotype. The change map \mathbf{z}^* corresponding to the most frequently sampled history matched that used to generate the alignment in no less than 44/50 trials and often in 50/50 trials. Furthermore, $P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t})$ was always greater than $\hat{\pi}_{\mathbf{z}^*}$ with average differences ranging between $0.22 < P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t}) - \hat{\pi}_{\mathbf{z}^*} < 0.27$. This result illustrates how accounting for PG associations can substantially reduce uncertainty in the inferred history of the phenotype.

generating model	$\hat{\pi}_0$	$\hat{\omega}_2$	$\hat{\pi}_{\text{CW}}$	$\hat{\pi}_{\text{BW}}$	rejections
Scenario 2a ($\pi_{\text{CW}}, \pi_{\text{BW}} = (0\%, 0\%)$)					
MSmmtDNA 0% DT	0.59	1.10	0.01	0.00	0/50 false
MSmmtDNA 6% DT	0.65	1.40	0.01	0.01	0/50 false
MSTGdR 0% DT	0.76	2.31	0.00	0.00	0/50 false
Scenario 2b ($\pi_{\text{CW}}, \pi_{\text{BW}} = (5\%, 0\%)$)					
MSmmtDNA 0% DT	0.55	1.11	0.07	0.01	47/50 true
MSmmtDNA 6% DT	0.58	1.25	0.07	0.01	46/50 true
MSTGdR 0% DT	0.71	1.38	0.04	0.00	42/50 true
Scenario 2c ($\pi_{\text{CW}}, \pi_{\text{BW}} = (0\%, 5\%)$)					
MSmmtDNA 0% DT	0.59	1.10	0.01	0.09	46/50 true
MSmmtDNA 6% DT	0.57	2.20	0.02	0.07	38/50 true
MSTGdR 0% DT	0.73	1.35	0.01	0.07	50/50 true
Scenario 2d ($\pi_{\text{CW}}, \pi_{\text{BW}} = (0\%, 0\%)$)					
MSmmtDNA 0% DT	0.60	1.27	0.01	0.01	1/50 false
MSmmtDNA 6% DT	0.56	2.30	0.01	0.00	1/50 false
MSTGdR 0% DT	0.72	1.37	0.00	0.00	2/50 false
Scenario 2e ($\pi_{\text{CW}}, \pi_{\text{BW}} = (5\%, 5\%)$)					
MSmmtDNA 0% DT	0.55	1.10	0.06	0.09	50/50 true
MSmmtDNA 6% DT	0.57	1.42	0.06	0.08	50/50 true
MSTGdR 0% DT	0.70	1.36	0.06	0.06	50/50 true

Table 4.3: Simulation 2 Results. Select mean MLEs and omnibus test results for Simulation 2. Note that Scenario 2d used the same alignments as Scenario 2c but with a misspecified vector of phenotypes.

model	CW FDC	CW Power	BW FDC	BW Power	(prior, post)	=
Scenario 2b (π_{CW}, π_{BW}) = (5%, 0%)						
1	0.80	0.31(4.72)	0.86	-	(0.73,0.96)	48
2	0.66	0.34(5.12)	1.16	-	(0.70,0.96)	47
3	0.62	0.36(5.44)	1.20	-	(0.62,0.88)	44
Scenario 2c (π_{CW}, π_{BW}) = (0%, 5%)						
1	1.30	-	1.54	0.37(5.60)	(0.76,1.00)	50
2	1.86	-	1.08	0.24(3.54)	(0.72,0.98)	49
3	1.76	-	1.20	0.59(8.88)	(0.73,1.00)	50
Scenario 2e (π_{CW}, π_{BW}) = (5%, 5%)						
1	1.04	0.29(4.38)	2.12	0.28(4.14)	(0.75,1.00)	50
2	1.90	0.26(3.90)	1.76	0.34(5.10)	(0.77,1.00)	50
3	1.44	0.46(6.86)	1.48	0.42(6.18)	(0.78,1.00)	50

Table 4.4: Simulation 2 *post hoc* analysis. Results of the *post hoc* analysis of Simulation 2 datasets. Column headings: CW FDC - the average number of CW sites discovered per alignment that were false; CW Power - average proportion of sites generated under the CW processes that were correctly identified (in brackets is the average number of true discoveries per alignment); BW FDC and BW power are similarly defined; prior = $\hat{\pi}_{\mathbf{z}^*}$ - the frequency of the most frequently sampled change map \mathbf{z}^* averaged over trials; post = $P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t})$ - the probability of \mathbf{z}^* conditioned on all of the data averaged over trials; “=” represents matches - the number of trials for which \mathbf{z}^* matched the generating change map. Models: 1 = MSmmtDNA 0% DT; 2 = MSmmtDNA 6% DT, 3 = MSTGdR 0% DT.

4.4.3 Simulation 3: A Scenario with Four Phenotypic States

The third simulation study was conducted to assess the statistical accuracy and power of the PG-BSM under scenarios with more than two phenotypic states. In this case only MSmmtDNA was used to generate data because the results of Simulation 2 indicated no substantial differences between the three MS generating processes used there. Fifty 300-codon alignments were generated on the unlocked tree in Figure 4.2 under four scenarios. In the first scenario (denoted 3a (π_{CW}, π_{BW}) = (0%, 0%)) alignments were generated with no PG association. In the second scenario (denoted 3b (π_{CW}, π_{BW}) = (5%, 0%)) alignments were generated with a reduction in the stringency of selection at 5% of sites. In the third scenario (denoted 3c (π_{CW}, π_{BW}) = (0%, 5%)) alignments were generated with peak shifts at 5% of sites. And in the last scenario (denoted 3d (π_{CW}, π_{BW}) = (5%, 5%)) alignments were generated with both a reduction in the stringency of selection at 5% of sites and peak shifts at 5% of sites. In all cases the stringency of selection and/or peak shifts were effected along the branches marked in bold. The phenotype assumed by the alternate

PG-BSM was $\mathbf{F} = (1, 1, 1, 1, 0, 0, 0, 0, 0, 2, 2, 2, 2, 2, 3, 3)$ in all scenarios. The unclocked tree is arguably more consistent with real data in both its irregular topology and depth compared to the clocked tree in Figure 4.2, and was chosen, in combination with the increase in the number of phenotypes, to provide a more challenging test of model performance.

The omnibus test correctly failed to reject the null hypothesis in all Scenario 3a $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 0\%)$ trials under which alignments were generated with no PG association (Table 4.5). The null was correctly rejected in all Scenario 3c $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 5\%)$ and Scenario 3d $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 5\%)$ trials. However, in Scenario 3b $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 0\%)$ the null was correctly rejected in only 15/50 trials. The PG-BSM apparently had difficulty identifying sites generated with a reduction in the stringency of selection. Concordantly, the average power to detect CW sites was 0.16 for Scenario 3b $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 0\%)$ alignments and 0.36 for Scenario 3d $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 5\%)$ alignments, as reported in Table 4.6, much less than the average power to detect BW sites, which was 0.67 for Scenarios 3c $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (0\%, 5\%)$ and 3d $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 5\%)$. The power to detect BW sites was substantially better overall in Simulation 3 (0.67) compared to Simulation 2 ($0.24 \leq \text{power} \leq 0.59$), and the FDCs for both CW and BW sites were significantly lower (no more than 0.78 false discoveries per alignment compared to as much as 2.12 among Simulation 2 scenarios). There was also a marked increase in uncertainty in the ancestral phenotypes in Simulation 3, with an average $0.38 \leq \hat{\pi}_{\mathbf{z}^*} \leq 0.56$ (Table 4.6) compared to $0.62 \leq \hat{\pi}_{\mathbf{z}^*} \leq 0.78$ (Table 4.4) for Simulation 2. However, the combined use of information in both the alignment and the vector of phenotypes made up the difference, as $0.80 \leq P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t}) \leq 1.00$ for Simulation 3 scenarios compared to $0.88 \leq P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t}) \leq 1.00$ for Simulation 2. More telling is the fact that $P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t})$ was at least twice as large as $\hat{\pi}_{\mathbf{z}^*}$ in all Simulation 3 scenarios with PG association, underlining one benefit of combining phenotype and sequence data when such associations exist.

A possible explanation for the low power of the omnibus test in Scenario 3b $(\pi_{\text{CW}}, \pi_{\text{BW}}) = (5\%, 0\%)$ is confounding due to “branch-length effects”. As

was stated in Chapter 3, two alignment-generating processes (real or simulated) are nearly confounded if the site-pattern distributions they produce are approximately the same (Jones *et al.*, 2018). Sites evolved under MSmmtDNA on the unlocked tree in Figure 4.2 with relaxation of selection pressure along the three branches marked in bold tend to produce site patterns consistent with the phenomenological CW process, (i.e., with greater diversity among amino acids at the terminal nodes indicated by the filled markers and less diversity among terminal nodes indicated by the open marker). But similar patterns can arise on that tree at sites evolved on static fitness landscapes due to the fact that the distances from the root to the terminal nodes indicated by the filled markers are relatively long (increasing the probability of replacement substitutions) whereas the tip-to-root distances for terminal nodes indicated by the open marker are relatively short (decreasing the probability of replacement substitutions). Heterotachous site patterns x^h generated with relaxation of selection pressure in Scenario 3b $(\pi_{CW}, \pi_{BW}) = (5\%, 0\%)$ therefore tended to be approximately as likely under the CL process as they were under the CW process: $P_{CL}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) \approx \sum_{\mathbf{z}} P_{CW}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}, \mathbf{z}) \hat{\pi}_{\mathbf{z}}$. The LLR for the contrast between the null and alternate PG-BSM consequently tended to be small, and often something less than the critical value 5.99 (assuming the χ_2^2 distribution for the LLR and a 5% significance test). Note that modifying the alternate PG-BSM to test for the CW process only, which permits the use of the χ_1^2 distribution for the LLR and a critical value of 3.84 for a 5% test, increased the power of the omnibus test only slightly (19/50 rejections instead of 15/50).

Scenario	$\hat{\pi}_0$	$\hat{\omega}_2$	$\hat{\pi}_{CW}$	$\hat{\pi}_{BW}$	rejections
3a, $(\pi_{CW}, \pi_{BW}) = (0\%, 0\%)$	0.62	1.00	0.00	0.00	0/50 false
3b, $(\pi_{CW}, \pi_{BW}) = (5\%, 0\%)$	0.59	1.05	0.02	0.00	15/50 true
3c, $(\pi_{CW}, \pi_{BW}) = (0\%, 5\%)$	0.60	2.21	0.00	0.06	50/50 true
3d, $(\pi_{CW}, \pi_{BW}) = (5\%, 5\%)$	0.57	1.52	0.04	0.05	50/50 true

Table 4.5: Selected results for Simulation 3 where there were four phenotypes evolved from three discrete changes over the tree.

Scenario	CW FDC	CW Power	BW FDC	BW Power	(prior, post)	=
3b	0.08	0.16(2.46)	0.14	-	(0.38,0.80)	49
3c	0.44	-	0.32	0.67(10.18)	(0.50,1.00)	50
3d	0.68	0.36(5.38)	0.78	0.67(10.12)	(0.56,1.00)	49

Table 4.6: Simulation 3 *post hoc* analysis. Results of the *post hoc* analysis of Simulation 3 datasets. See caption for Table 4.4 for definitions of each column heading.

4.5 Analysis with Real Data

4.5.1 mmtDNA

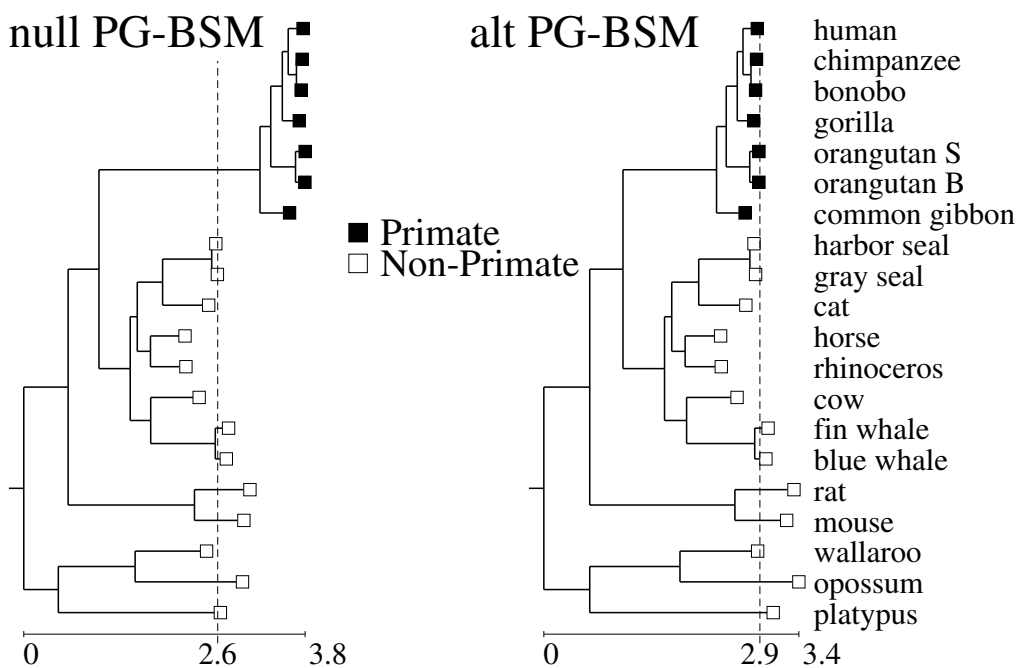


Figure 4.3: Branch lengths (the expected number of single nucleotide substitutions per codon) estimated by fitting the null and alternate PG-BSM to the alignment of 12 concatenated H-strand mitochondrial DNA sequences (3331 codon sites) from 20 mammalian species distributed in alignment form by the PAML software package (Yang, 2007). Trees are drawn to the same scale. Vertical dashed lines indicate the average tip-to-root distance (d_i , $i = 1, \dots, 20$) across terminal nodes. Accounting for non-stationary processes under the alternate model caused the tree to be more consistent with the molecular clock hypothesis, with less variation in the d_i .

Mitochondrial genes encode proteins involved in basic metabolic processes and are therefore thought to be functionally constrained most of the time. Signatures of adaptation in mammalian mmtDNA have nevertheless been detected. For example, Pupko and Galtier (2002) detected a significant difference

in the replacement rates at some sites in mitochondrial genes within a clade of 7 simian primates compared to a clade of 27 other mammalian lineages. The authors interpreted this as evidence of adaptation and speculated that it might have been due to changes in metabolic requirements related to an increase in the size of the neocortex among primates compared to the other mammals included in their analysis (Pupko and Galtier, 2002). In this section I present an analysis of a similar alignment of 12 mmtDNA genes 3331 codons in length taken from 20 mammalian species (Yang, 2007). Primate/non-primate was used as a binary phenotype to facilitate comparison with Pupko and Galtier (2002), but consider this to be a proxy for an unobserved and unknown phenotype and view our analysis as exploratory.

Four models of increasing complexity were fitted to the alignment: (i) the model M3($k = 2$) (Yang *et al.*, 2000a), which assumes that some proportion p_1 of sites evolved under a smaller ω_1 and the remaining sites under a larger ω_2 across the tree without heterotachy; (ii) the covarion-like model CLM3($k = 2$) (Jones *et al.*, 2017), which assumes that sites evolved under ω_1 a proportion p_1 of the time but switched randomly between ω_1 and ω_2 at a rate of δ switches per single nucleotide substitution; (iii) the null PG-BSM, which combines CLM3($k = 2$) with a category of sites that evolved under $\omega_0 = 0$; and (iv) the alternate PG-BSM, which adds to the null model two categories of sites for the CW and BW processes. Note that CLM3($k = 2$) is fitted to the alignment X alone, whereas the PG-BSM is fitted to X together with the vector of phenotypes \mathbf{F} . The log-likelihood of the PG-BSM was therefore reported as the sum of two values. From equation (4.4), where $\boldsymbol{\theta} = \langle \omega_1, \omega_2, p_1, \delta, \kappa, \rangle$:

$$\begin{aligned} & \ln\{L_{\text{mul}}(X, \mathbf{F}; \lambda, \boldsymbol{\theta}, \mathbf{t})\} \\ = & \ln\{P(\mathbf{F}; \lambda, \mathbf{t})\} + \ln \left\{ \prod_{h=1}^n (\pi_0 P_0(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) + (1 - \pi_0) P_{\text{CL}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t})) \right\} \\ = & \text{LL}_{\mathbf{F}}(\mathbf{F}; \lambda, \mathbf{t}) + \text{LL}_X(X; \boldsymbol{\theta}, \mathbf{t}) \end{aligned}$$

Similarly:

$$\text{LL}_{\text{CL}}(X; \boldsymbol{\theta}, \mathbf{t}) = \ln \left\{ \prod_{h=1}^n P_{\text{CL}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) \right\}$$

The $\text{LLR} = 2 \{ \text{LL}_X(X; \boldsymbol{\theta}, \mathbf{t}) - \text{LL}_{\text{CL}}(X; \boldsymbol{\theta}, \mathbf{t}) \}$ evaluated at the MLEs for the two models provides a means to test the significance of accounting for sites evolving under $\omega_0 = 0$. Note however that, whereas \mathbf{t} is estimated from (X, \mathbf{F}) under the null PG-BSM, it is estimated from X alone under $\text{CLM3}(k = 2)$. Hence the LLR is only approximate. It is likely to be a close approximation however, since the information in \mathbf{F} generally has a very small impact on the MLE for \mathbf{t} , which is determined primarily by the information in X .

Accounting for heterotachy resulted in a large improvement in fit, as the M3 vs CLM3 contrast produced $\text{LLR} = 2(89,162 - 88,880) = 564$ on one parameter (δ) (Table 4.7). Including a category for sites that evolved under $\omega_0 = 0$ also resulted in a large improvement in fit, with $\text{LLR} = 2(88,880 - 88,719) = 322$ on one parameter (π_0) for the CLM3 vs null PG-BSM contrast. Accounting for PG association resulted in a smaller but also highly significant improvement, with $\text{LLR} = 2(88,719 - 88,681) = 76$ on two parameters (π_{CW} and π_{BW}). Trees with branch lengths estimated under the null and alternate PG-BSM are shown in Figure 4.3. It is interesting that accounting for non-stationary CW and BW processes under the alternate PG-BSM reduced the length of the branch leading to the primate clade from 2.2 single nucleotide substitutions per codon under the null model to 1.3 under the alternate model. This was accompanied by an increase in the average of the distances d_1, \dots, d_{20} from each terminal node to the root from 2.6 to 2.9. The combined effect was a reduction in the variance of the d_i . Accounting for the non-stationary CW and BW processes therefore caused the tree to be more consistent with the molecular clock hypothesis.

The alternate PG-BSM inferred a fraction of sites ($\hat{\pi}_{\text{CW}} = 0.06, \hat{\pi}_{\text{BW}} = 0.04$) to be associated with the change from non-primate to primate. The *post hoc* analysis identified a total of 17 sites (11 CW and 6 BW) that likely underwent changes in their site-specific fitness landscapes along the branch leading to the primate clade. Substitution patterns for these sites are shown in Figure 4.4. CW sites with posteriors $0.77 \leq P(\text{CW}) \leq 0.87$ (where the upper bound 0.87 corresponds to the maximum value of $P(\text{CW})$ and the lower

bound was determined by setting $E\{\text{FDC}(\text{CW})\} = 2$) exhibit patterns consistent with Type I FD, consisting of two or more amino acids among primates and typically a single amino acid among non-primates. BW sites with posteriors $0.67 \leq P(\text{BW}) \leq 0.75$ (where the upper bound 0.75 corresponds to the maximum value of $P(\text{BW})$ and the lower bound was determined by setting $E\{\text{FDC}(\text{BW})\} = 2$) exhibit patterns consistent with Type II FD, with a single amino acid among primates and a single different amino acid among non-primates.

YN-BSM A and B were also fitted to the alignment using the long branch leading to the primate clade as the FG. The YN-BSM B is similar to Model A but estimates the three rate ratios ω_0 , ω_1 , and ω_2 freely and uses M3($k = 2$) as the null model (Yang and Nielsen, 2002). Note that the M3 vs Model B contrast is not a test for $\omega_2 > 1$, but only for a change to ω_2 at some sites along the FG branch. Model A detected no evidence of positive selection along the branch leading to the primate clade (LLR = 2.62 compared to a critical value of 3.84 for a 5% test assuming $\text{LLR} \sim \chi_1^2$, Table 4.8). Model B detected strong evidence (LLR = 54.44 compared to a critical value of 5.99 for a 5% test assuming $\text{LLR} \sim \chi_2^2$) for an elevation in the rate ratio at some sites along the same branch ($\hat{\omega}_2 = 0.98$, $\hat{p}_2 = 0.15$). Akaike's Information Criterion ($\text{AIC} = 2m - 2\text{LL}$, where m is the number of estimated model parameters and LL is the log-likelihood of the data under the model) is frequently used to compare non-nested models under the maximum likelihood framework. The better of any two models is the one with the smaller AIC. By this criterion the PG-BSM provided the best fit, as the AIC was 177,456 for the PG-BSM, 178,358 for Model B and 181,777 for Model A.

Figure 4.5 shows 13 sites inferred by YN-BSM B to have undergone an elevation in rate ratio to $\hat{\omega}_2 = 0.98$ on the FG branch after FDC control was applied to the BEB (or Bayes-Empirical-Bayes, Yang *et al.*, 2005) posteriors with $E\{\text{FDC}\} = 2$. Three of these sites have amino acid patterns that match the primate/non-primate phenotype (marked by filled circles in Figure 4.5). All three are among the 6 sites inferred by the PG-BSM to have undergone the BW process. None of the sites are consistent with the CW process. However,

seven of them are consistent with the rCW process, exhibiting one amino acid among primates and several among non-primates (marked by filled triangles in Figure 4.5). Similar site patterns were identified by Pupko and Galtier (2002) in their data, and were interpreted as evidence of functional shifts. Motivated by the presence of site patterns of this kind, the alternate PG-BSM accounting for the rCW process alone was fitted to the alignment. No evidence for the rCW process was detected ($\hat{\pi}_{\text{rCW}} = 0.00$) despite the fact that the alignment includes a fair number of site patterns that are apparently consistent with this process (e.g., there are 828 sites patterns with one amino acid among primates and 2 or more among the 13 non-primates).

The fact that the PG-BSM did not detect evidence for the rCW process suggests that the CL component of the null model provides a measure of protection against branch-length effects (although possibly at the expense of statistical power, cf. Scenario 3b ($\pi_{\text{CW}}, \pi_{\text{BW}} = (5\%, 0\%)$). A site that evolved on a static fitness landscape over the tree in Figure 4.3 is likely to exhibit a single amino acid among the primates because the branches within the primate clade are very short. The same site is likely to exhibit considerable amino acid diversity among the remaining terminal nodes because the non-primates consist of three clades each with relatively long terminal branches. It follows that the site-pattern distribution implied by evolution on static fitness landscapes is similar to that implied by the mechanism whereby sites undergo an increase in the stringency of selection along the branch leading to the primate clade. Heterotachous site patterns \mathbf{x}^h in the mmtDNA alignment consistent with the rCW process therefore tended to be approximately equally likely under the CL process: $P_{\text{CL}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}) \approx \sum_{\mathbf{z}} P_{\text{rCW}}(\mathbf{x}^h; \boldsymbol{\theta}, \mathbf{t}, \mathbf{z}) \hat{\pi}_{\mathbf{z}}$. Inclusion of the CL process in the null model consequently resulted in no evidence for the rCW process, possibly preventing a type I error.

Model	LL_X	LL_F	$\hat{\pi}_0$	$\hat{\omega}_1$	$\hat{\omega}_2$	\hat{p}_1	$\hat{\delta}$	$\hat{\pi}_{\text{CW}}$	$\hat{\pi}_{\text{BW}}$
M3($k = 2$)	-89,162	-	-	0.01	0.15	0.71	-	-	-
CLM3($k = 2$)	-88,880	-	-	0.00	0.21	0.77	0.06	-	-
null PG-BSM	-88,719	-4	0.47	0.02	0.35	0.76	0.12	-	-
alt PG-BSM	-88,681	-4	0.48	0.02	0.31	0.70	0.13	0.06	0.04

Table 4.7: Results of the analysis of the mmtDNA.

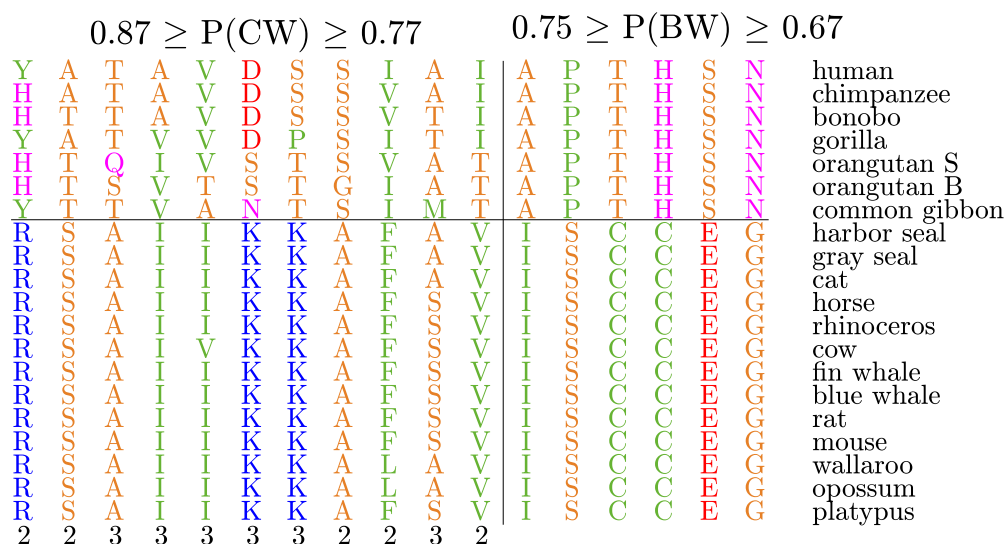


Figure 4.4: Patterns for sites in the mmtDNA alignment. Sites shown are those inferred to have evolved in association with a change in the primate/non-primate binary character state that remained after application of FDC control with $E\{\text{FDC}\} = 2$. Rows designate taxa in the same order as they appear in the tree shown in Figure 4.3. Columns designate sites. The horizontal line separates primates from non-primates. The vertical line separates CW from BW sites. Sites of each type are arranged in order of decreasing posterior probability estimated using equation (4.8). CW sites (left) exhibit patterns of Type I FD consistent with an elevated replacement rate along branches within the primate clade. BW sites (right) exhibit patterns of Type II FD with conserved but different amino acids within the two groups of taxa. The value below each CW site pattern indicates the number of different amino acids among primates at that site.

	Model A LLR = 2.62 < 3.84			Model B LLR = 54.44 > 5.99		
category	proportion	BG ω	FG ω	proportion	BG ω	FG ω
0	0.89	0.03	0.03	0.70	0.01	0.01
1	0.07	1.00	1.00	0.28	0.15	0.15
2a	0.04	0.03	1.79	0.01	0.01	0.98
2b	0.00	1.00	1.79	0.01	0.15	0.98

Table 4.8: Results of the fit of alternate YN-BSM A and B to the mmtDNA data.

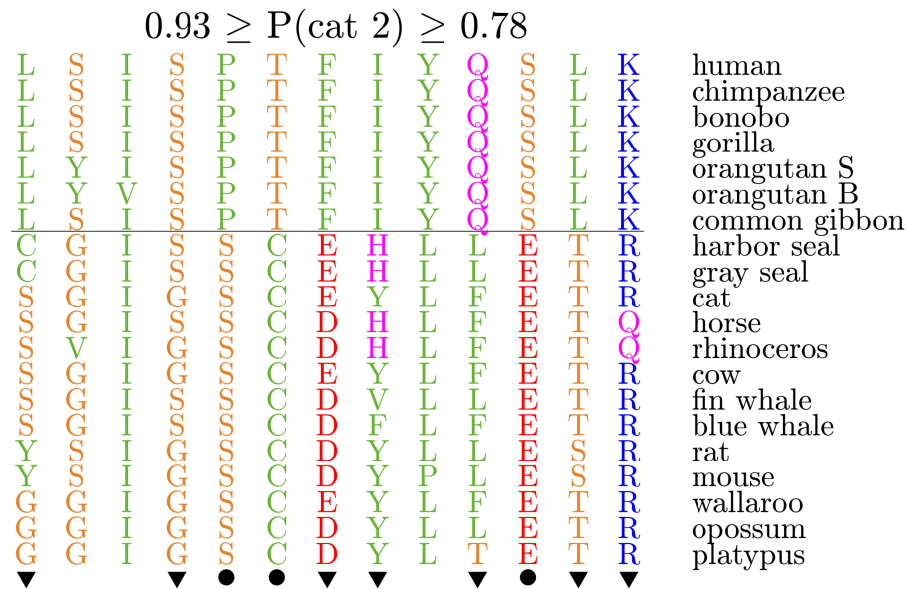


Figure 4.5: Patterns for sites in the mmtDNA alignment inferred to be in category 2a or 2b by the YN-BSM B. Sites shown remained after application of FDC control with $E\{\text{FDC}\} = 2$. Rows designate taxa in the same order as they appear in the tree shown in Figure 4.3. Columns designate sites. The horizontal line separates primates from non-primates. Sites are arranged in order of decreasing posterior probability estimated using YB-BSM B. Filled triangles and circles mark site patterns most consistent with the rCW and BW process, respectively.

4.5.2 Phytochrome A&CF

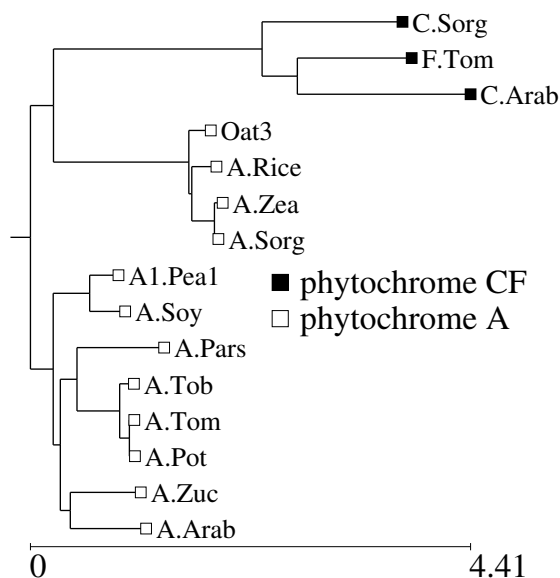


Figure 4.6: Branch lengths (the expected number of single nucleotide substitutions per codon) estimated by fitting the alternate PG-BSM to the phytochrome A&CF alignment.

Phytochrome is a plant photo-receptor associated with the regulation of developmental processes. Most seed plants contain variants phyA, phyB and phyC, but some also include variants phyD, phyE and phyF, which arose following duplication of phyB (Mathews, 2010). An analysis of an alignment of 15 angiosperm phytochrome sequences 1072 codons in length (Figure 4.6) was conducted. Previous analyses of the same data performed as a test case for the YN-BSM partitioned the sequences as phyCF versus phyA (Yang and Nielsen, 2002; Zhang *et al.*, 2005). The same partition was used here as a binary phenotype.

Models $M3(k = 2)$, $CLM3(k = 2)$, and the null and alternate PG-BSM were fitted to the data. Accounting for heterotachy resulted in a large improvement in fit, as the LLR for the $M3$ vs $CLM3$ contrast was $2(28,818 - 28,739) = 158$ on one parameter (δ) (Table 4.9). Including a category for sites that evolved under $\omega_0 = 0$ also improved the fit, with $LLR = 2(28,739 - 28,708) = 62$ on one parameter (π_0) for the $CLM3$ vs null PG-BSM contrast. Accounting for PG associations engendered no improvement ($LLR = 0$), meaning that PG

associations were not detected. Discordantly, both YN-BSM A and B rejected their respective nulls (Table 4.10). Model A detected evidence of positive selection at some sites along the branch leading to the phyCF clade (i.e., category 2 sites with $\hat{\omega}_2 = 30.17, \hat{p}_2 = 0.12$) and identified 29 sites with BEB posteriors $0.95 < P(\text{cat } 2) < 1.00$, all of which remained after FDC control was applied with $E\{\text{FDC}\} = 1$ (Figure 4.7). Model B detected evidence for an elevation in the rate ratio at some sites along the same branch ($\hat{\omega}_2 = 5.56, \hat{p}_2 = 0.05$).

The AIC was 57,510 for the null PG-BSM, 57,700 for Model B and 58,255 for Model A. Thus the null PG-BSM provided the best fit of the three models. The alignment nevertheless contains site patterns consistent with the CW and BW processes (e.g., those marked by filled triangles and circles in Figure 4.7). The data might contain true signal for PG association that went undetected due to the unusually large proportion (70%) of variable sites in the alignment. To test this hypothesis, a fourth simulation was conducted under which MSmmtDNA was used to generate sets of 50 alignments 1072 codons in length on the phytochrome tree (Table 4.11). The proportion of variable sites can be controlled under MSmmtDNA by changing the proportion of sites with landscapes that admit non-adaptive shifting balance (i.e., landscapes with a selection regime somewhere between stringent and neutral, Jones *et al.*, 2017). Alignments under scenarios 4a and 4b were generated with either $\approx 40\%$ or $\approx 70\%$ variable sites, including 5% CW and 5% BW sites. These were used to assess and compare the power of the PG-BSM and YN-BSM A omnibus tests. Alignments under scenarios 4c and 4d were generated with either $\approx 40\%$ or $\approx 70\%$ variable sites but with no PG association. These were used to assess and compare the accuracy of the omnibus tests.

The PG-BSM correctly detected PG association in 50/50 alignments generated with 40% variable sites (scenario 4a ($\pi_{\text{CW}}, \pi_{\text{BW}} = (5\%, 5\%)$)), but in only 42/50 alignments generated with 70% variable sites (scenario 4b ($\pi_{\text{CW}}, \pi_{\text{BW}} = (5\%, 5\%)$), see Table 4.11). Although 42/50 indicates substantial statistical power, the reduction in the number of detections in comparison with 40% variable sites is consistent with our hypothesis that the power of the PG-BSM can be reduced when sequences are highly divergent. The PG-BSM

produced 0/50 false positives when there was no PG association regardless of the proportion of variable sites. The YN-BSM A, by contrast, inferred positive selection (i.e., $\omega_2 > 1$) at some sites along the branch leading to the phyCF clade in 11/50 alignments generated with 40% variable sites (scenario 4c ($\pi_{\text{CW}}, \pi_{\text{BW}} = (0\%, 0\%)$)) and in 31/50 alignments generated with 70% variable sites (scenario 4d ($\pi_{\text{CW}}, \pi_{\text{BW}} = (0\%, 0\%)$)). Some of these might be true evidence of $\omega > 1$ on the FG since positive selection due to non-adaptive shifting balance is expected to occur some of the time (Jones *et al.*, 2017). However, they are all false positives when interpreted as evidence of adaptive evolution.

The PG-BSM was specifically designed to account for non-adaptive shifting balance with the inclusion of CLM3($k = 2$) as a component of the mixture in equation (4.7). The importance of this component is illustrated by fixing $\delta = 0$ and fitting the resulting modified PG-BSM to scenario 4d alignments (70% variable sites, no PG association). Setting the switching rate to zero has the effect of making CLM3($k = 2$) equivalent to M3($k = 2$), since $\omega_1 < \omega_2$ are still estimated but sites can no longer switch between them. Sites most consistent with M3($k = 2$) are those that evolved at a constant rate over the tree. The modified version of the alternate PG-BSM can therefore accommodate heterotachous sites only by appealing to the CW and BW processes. The modified PG-BSM incorrectly inferred PG association in 33/50 trials when fitted to scenario 4d alignments compared to 0/50 for the regular PG-BSM.

Model	LL_X	LL_F	$\hat{\pi}_0$	$\hat{\omega}_1$	$\hat{\omega}_2$	\hat{p}_1	$\hat{\delta}$	$\hat{\pi}_{\text{CW}}$	$\hat{\pi}_{\text{BW}}$
M3($k = 2$)	-28,818	-	-	0.03	0.22	0.55	-	-	-
CLM3($k = 2$)	-28,739	-	-	0.01	0.30	0.62	0.13	-	-
null PG-BSM	-28,708	-4	0.22	0.04	0.40	0.68	0.20	-	-
alt PG-BSM	-28,708	-4	0.22	0.04	0.40	0.68	0.20	0.00	0.00

Table 4.9: Results of the analysis of the phyA&CF data.

	Model A LLR = 22.67 > 3.84			Model B LLR = 23.27 > 5.99		
category	proportion	BG ω	FG ω	proportion	BG ω	FG ω
0	0.81	0.09	0.09	0.55	0.03	0.03
1	0.07	1.00	1.00	0.40	0.23	0.23
2a	0.11	0.09	30.17	0.03	0.03	5.56
2b	0.01	1.00	30.17	0.02	0.23	5.56

Table 4.10: Results of the fit of alternate YN-BSM A and B to the phyA&CF data.

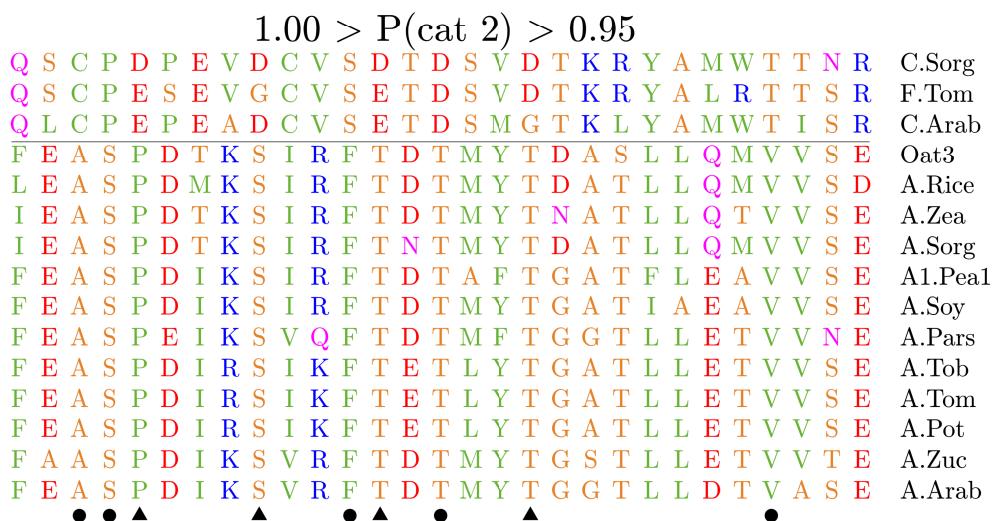


Figure 4.7: Patterns for sites in the phytochrome alignment inferred to be in category 2a or 2b by the YN-BSM A. Sites shown remained after application of FDC control with $E\{\text{FDC}\} = 1$. Rows designate taxa in the same order as they appear in the tree shown in Figure 4.6. Columns designate sites. The horizontal line separates phyA from phyCF. Sites are arranged in order of decreasing posterior probabilities estimated using YB-BSM A. Filled triangles and circles mark site patterns most consistent with the CW and BW processes, respectively.

scenario	% variant	% CW	% BW	PG-BSM	YN-BSM A
4a	40	5	5	50/50	50/50
4b	70	5	5	42/50	50/50
4c	40	0	0	0/50	11/50
4d	70	0	0	0/50	31/50

Table 4.11: Simulation 4 Results. Counts of the number of times the null was rejected for omnibus tests applied to Simulation 4 alignments. % variant gives the approximate proportion of site patterns in the simulated alignments that exhibited some degree of heterotachy. The columns PG-BSM and YN-BSM A show the number of trials out of 50 for which the relevant model rejected the null hypothesis. Rejections are true for scenarios 4a and 4b and false for scenarios 4c and 4d.

4.5.3 Invertebrate Cytochrome B

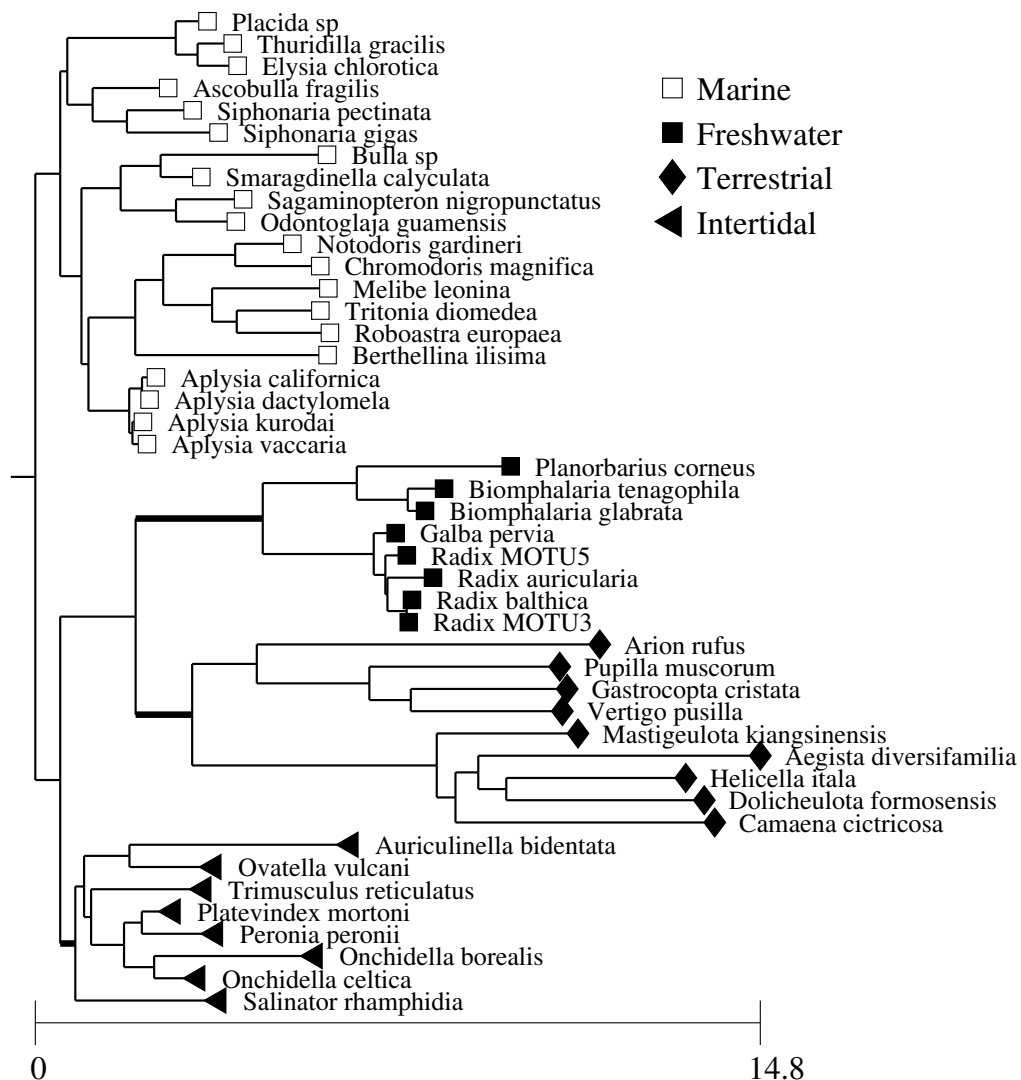


Figure 4.8: Branch lengths (the expected number of single nucleotide substitutions per codon) estimated by fitting the null PG-BSM to the cytochrome B alignment.

Euthyneura (snails and slugs) have adapted to diverse habitats, including marine, intertidal, terrestrial and freshwater. Their mitochondrial genome includes cytochrome B, an essential component of the electron transport chain

common to most life forms on Earth. Given its crucial role, one would expect *cytB* to be highly conserved. It is nevertheless reasonable to suspect that transition from the marine to the other three environments might have required some adaptations (e.g., for differences in osmotic pressure or the risk of desiccation). To test this hypothesis, the PG-BSM and YN-BSM were fitted to an alignment consisting of 45 *cytB* sequences 341 codons in length. The sequences were selected from a larger published data set (Romero *et al.*, 2016) to produce four homogeneous clades. Tree topology was estimated from the DNA sequences using RAxMLv0.6.0 with default settings, and the tree was rooted to produce that shown in Figure 4.8.

The PG-BSM was initially fitted to the alignment using the different environments to define a phenotype with four states, but no signal for PG association was found. The data were then re-analyzed using three different binary phenotypes: terrestrial vs non-terrestrial, freshwater vs non-freshwater, and intertidal vs non-intertidal. Furthermore, the PG-BSM was modified to detect either CW, rCW or BW sites alone (i.e., using three versions of the alternate model). The YN-BSM A was also fitted to the alignment using the branch leading to the terrestrial, freshwater, and intertidal clade each in turn as the foreground (marked in bold in Figure 4.8). Signal was detected for BW sites by the PG-BSM when phenotype was set to freshwater vs non-freshwater, with $LLR = 2(24,537 - 24,527) = 20$ compared to a critical value of 5.73 (assuming $LLR \sim \chi_1^2$ and using a level of significance $\alpha = 0.05/3$ to adjust for the fact that three tests were conducted on the alignment with freshwater vs non-freshwater as the phenotype). The MLEs for the analysis were $\hat{\omega}_1 = 0.00$, $\hat{\omega}_2 = 0.08$, $\hat{p}_1 = 0.58$, $\hat{\delta} = 0.06$ and $\hat{\pi}_{BW} = 0.06$. The YN-BSM A detected evidence of positive selection in two cases, once along the branch leading to the terrestrial clade ($LLR = 2(25,819 - 25,812) = 14$, $\hat{\omega}_2 = 999$, $\hat{p}_2 = 0.04$) and again along the branch leading to the freshwater clade ($LLR = 2(25,811 - 25,807) = 8$, $\hat{\omega}_2 = 999$, $\hat{p}_2 = 0.12$). As in the previous two analyses with real data, here the PG-BSM provided the better fit: AIC was 51,700 for YN-BSM A but 49,146 for the alternate PG-BSM with BW alone using the freshwater vs non-freshwater phenotype.

It is instructive to compare the distribution of amino acids among each clade at sites identified by the YN-BSM and PG-BSM via *post hoc* analysis. Following application of FDC control with $E\{\text{FDC}\} = 1$, the YN-BSM A identified 11 sites with strong evidence of having undergone episodic positive selection on the branch leading to the freshwater clade. Table 4.12 shows the distribution of the amino acids at six of those sites for which $P(\text{cat } 2) \geq 0.95$. Site 329, for example, is occupied by four amino acids among the 20 taxa in the marine clade, 12 by L (Leucine), 5 by M (Methionine), 2 by I (Isolucine) and 1 by V (Valine). A comparison of distributions across clades gives some clue as to the process that might have generated the data. Sites 153 and 239, for example, exhibit one amino acid among the freshwater clade (T at site 153 and K at site 239) but are dominated by a different amino acid among the other three clades (N at site 153 and M at site 239). These patterns are consistent with peak shifts at these sites along the branch leading to the freshwater clade (i.e., the BW process). Site 329 and 127 show one amino acid among the freshwater clade and two or more different amino acids among each of the remaining clades. These sites are consistent with intensification of selective constraint in the freshwater clade (i.e., the rCW process) possibly following a peak shift (since the amino acid in the freshwater clade, C at site 329 and M at site 127, does not occur in any of the other three clades).

The PG-BSM fitted using freshwater vs non-freshwater as the phenotype detected seven sites with $0.52 \leq P(\text{BW}) \leq 0.99$ after *post hoc* analysis was conducted with $E\{\text{FDC}\} = 1$. The first six of these are shown in Table 4.13. The first four sites (153, 144, 25, and 239) are highly consistent with the BW process, being dominated by one amino acid among the freshwater clade and a different amino acid among the non-freshwater clades. Sites 182 and 64 are both occupied by serine only. There are eight codon aliases for serine in the invertebrate mtDNA, including TCN and AGN where N is any nucleotide. Paths between TCN and AGN by single nucleotide substitutions require a minimum of one nonsynonymous change to either tryptophan, cysteine or threonine. The existence of serine sites with a mix of TCN and AGN would therefore suggest one or more replacement substitutions. The fact that sites

64 and 182 were identified in the *post hoc* analysis is explained by the codons that appear within each clade: both sites are occupied by AGN everywhere in the freshwater clade but are dominated by TCN among all remaining taxa data not shown. This suggests that substitutions to intermediate amino acids occurred along the branch leading to the freshwater clade. Note that the YN BSM A assigned the largest posterior to sites 64 and 182, and also assigned them equal probability $P(\text{cat } 2) = 0.9970$ despite the fact that the two sites have different codon substitution patterns (data not shown). The PG-BSM, by comparison, placed less weight on these sites and was apparently sensitive to their differences, since $P(\text{BW}) = 0.9221$ for site 182 but only $P(\text{BW}) = 0.7509$ for site 64.

Clade	site 64	site 182	site 153	site 329	site 239	site 127
Marine	S ₂₀	S ₂₀	N ₂₀	L ₁₂ M ₅ I ₂ V ₁	M ₂₀	L ₁₃ V ₅ A ₁ I ₁
Freshwater	S ₈	S ₈	T ₈	C ₈	K ₈	M ₈
Terrestrial	S ₉	S ₉	N ₉	L ₇ F ₁ A ₁	M ₈ L ₁	L ₇ F ₂
Intertidal	S ₈	S ₈	N ₈	L ₃ A ₃ S ₂	M ₈	L ₆ V ₂
P(cat 2)	0.9970	0.9970	0.9610	0.9590	0.9530	0.9500

Table 4.12: Amino acid compositions, YN-BSM A. Composition for sites with $P(\text{cat } 2) \geq 0.95$ are shown, as determined by the YN-BSM A using the branch leading to the Freshwater clade as the FG. Sites are shown in order of descending BEB posteriors. Letters represent amino acids and subscripts the number of taxa with that amino acid among the corresponding clade.

Clade	site 153	site 144	site 25	site 239	site 182	site 64
Marine	N ₂₀	G ₂₀	L ₂₀	M ₂₀	S ₂₀	S ₂₀
Freshwater	T ₈	S ₈	F ₈	K ₈	S ₈	S ₈
Terrestrial	N ₉	G ₉	I ₅ L ₄	M ₈ L ₁	S ₉	S ₉
Intertidal	N ₈	G ₈	L ₈	M ₈	S ₈	S ₈
P(BW)	0.9859	0.9799	0.9458	0.9433	0.9221	0.7509

Table 4.13: Amino acid compositions, PG-BSM. Composition for the first six sites identified by the PG-BSM to be associated with the Freshwater vs Other phenotype are shown. Sites are in order of descending P(BW). Letters represent amino acids and subscripts the number of taxa with that amino acid among the corresponding clade.

4.6 Discussion

Branch-site CSMs provide a means to detect evidence that a codon site underwent positive selection along a specified foreground branch of a phylogeny. Such evidence, in the form of an estimated rate ratio $\omega > 1$, is widely considered to be sufficient to infer adaptive evolution at a codon site (i.e., a site-specific peak shift). However, $\omega > 1$ does not necessarily imply adaptation. It is true that the dynamic at a codon site following a peak shift is characterized by a transient increase in the expected rate ratio, and that the increase can sometimes be to $\omega > 1$ (dos Reis, 2015). But the same can also occur on a static fitness landscape following chance fixation to a less-than-optimal amino acid (i.e., by non-adaptive shifting balance, Jones *et al.*, 2017, 2019a). It is therefore not possible to distinguish an episodic change in a site-specific landscape from non-adaptive shifting balance on a static landscape using estimates of ω alone. Furthermore, adaptation does not necessarily imply $\omega > 1$. The increase in the rate ratio following a peak shift rapidly diminishes as the site moves toward its new fitness peak (dos Reis, 2015). This suggests that the initial elevation in rate ratio can be more difficult to detect as sequences become more divergent. The analysis at the end of Chapter 2 supports this intuition. Peak shifts were implemented using the MS framework by simultaneously changing the fitness coefficients at all 1000 codon sites in an initial sequence S_1 that was subsequently evolved over a branch of length b to obtain a second sequence S_2 . M0 was then fitted to (S_1, S_2) to obtain $\hat{\omega}$. The median estimate across 200 trials was $\hat{\omega} \approx 1.4$ when $b = 0.2$, but $\hat{\omega} \approx 1.0$ when $b = 1.0$ (see Figure 2.9). It follows from all the above that adaptive evolution and the inference that $\omega > 1$ are not only not equivalent, but neither one necessarily implies the other.

The PG-BSM provides an approach for inferring adaptation that does not rely on evidence of positive selection. The method is based on the supposition that mechanisms of adaptation at the molecular level consist of changes in site-specific fitness landscapes. The mechanisms considered in this study consisted of either a persistent change in the stringency of selection at a site or a peak shift at a site along a particular branch of the tree. Changes in stringency

are represented by the CW and rCW processes as cladewise changes in rate ratio, whereas a peak shift is represented by the BW process as a transient elevation in rate ratio along specific branches of the tree. The locations of branches over which these processes may have occurred are informed by a discrete character state (e.g., a phenotype) via a model for the evolution of that character state. This constraint provides additional information that makes it possible to identify among all variant sites those with replacement patterns that imply PG association. The model also includes a covarion-like component to account for variant site patterns inconsistent with PG association. The CL component therefore provides the null hypothesis, which is rejected by the presence of site patterns that are more likely to have occurred under one of the CW, rCW or BW processes. Because these represent phenomenological outcomes consistent with changes in site-specific fitness landscapes, rejection of the null can be adduced as evidence of a change in functional constraint.

The rCW process implies an increase in functional constraint and the BW process a peak shift. Evidence of either therefore suggests molecular adaptation. The CW process, by contrast, was intended to detect sites that might have undergone a reduction in functional importance (expressed in the data as a reduction in the stringency of selection) along branches over which the phenotype changed, a process not necessarily associated with adaptation. However, many of the CW site patterns identified in the real mmtDNA suggest not only a reduction in the stringency of selection among primates but also a peak shift along the branch leading to the primate clade. The first CW site pattern in Figure 4.4, for example, consists of one amino acid (arginine) among non-primates and two amino acids among primates. A peak shift is suggested by the fact that the amino acids among primates are both different than arginine (histidine and tyrosine). Most of the CW sites in Figure 4.4 have similar patterns. It might therefore be the case that sites that underwent a reduction in the stringency of selection in combination with a peak shift are more readily detected by the CW component of the alternate PG-BSM than sites that underwent a reduction in stringency alone. This provides an alternative explanation for the low power to detect CW sites in Simulation Scenario

3b ($\pi_{\text{CW}}, \pi_{\text{BW}}$) = (5%, 0%), where CW sites were generated by reducing the stringency of selection alone (i.e., by rescaling the vector \mathbf{f}^h of site-specific fitness coefficients without changing the relative order of amino acid fitnesses as described in Methods). The alternate PG-BSM does not provide a formal test for the co-occurrence of a reduction in the stringency of selection with a peak shift however because both processes are expressed in the model as a phenomenological switch to the larger rate ratio ω_2 . That co-occurrence can be identified informally by appealing to contextual information (e.g., Figure 4.4) demonstrates that it might be formally detected by the PG-BSM via inclusion of measures to discriminate between amino acids (c.f., Gu, 2006). I leave this task for future efforts.

The PG-BSM framework offers several advantages over the YN-BSM. First, it includes a model for the evolution of a discrete phenotype that not only frees the analyst from the task of specifying the FG, but also automatically takes into account less likely but nevertheless possible evolutionary histories of the phenotype. Second, it includes a CL component to account for random shifts between $\omega_1 < \omega_2$ consistent with all processes that can potentially result in heterotachy. Covarion-like models (e.g., Galtier, 2001; Guindon *et al.*, 2004) were originally intended to account for epistatic interactions between codons sites thought to be the cause of the covarion (i.e., concomitantly variable codons, Fitch and Markowitz, 1970; Fitch, 1971) phenomenon. Potential sources of heterotachy include non-adaptive shifting balance and the fixation of DT mutations in addition to episodic changes in site-specific fitness landscapes, as demonstrated in Chapters 2 and 3. The utility of using the CL model as the null hypothesis was illustrated by simulations of the phytochrome A&CF alignment, where inclusion of the CL component of the PG-BSM was instrumental in reducing the false positive rate of the omnibus test. Third, pathologies such as false positives that can sometimes arise under the YN-BSM due to statistical irregularities (e.g., Baker *et al.*, 2016; Mingrone *et al.*, 2018) are avoided under the PG-BSM. The YN-BSM assumes that category 2 sites evolved under a separate rate ratio ω_2 on the FG. The rate ratio ω_2 is consequently nearly

unidentifiable when p_2 is small. Under this irregular condition, the maximum-likelihood estimate $\hat{\omega}_2$ is sometimes very large and potentially misleading (e.g., in our analysis of *cytB*, the YN-BSM A yielded $\hat{\omega}_2 = 999$ with $\hat{p}_2 = 0.04$ or $\hat{p}_2 = 0.12$). This issue is avoided under the PG-BSM because estimates of ω_1 and ω_2 make use of information contained in all variant sites. Fourth, the PG-BSM can identify sites consistent with specific mechanisms of adaptation even without evidence of positive selection. This key feature was empirically validated using simulation studies in which the null hypothesis was correctly rejected for the majority of alignments generated with changes in site-specific landscapes. Moreover, a fair proportion of sites generated under specific mechanisms (relaxation or intensification in the stringency of selection, a peak shift) were correctly identified via *post hoc* analysis.

The chance fixation into the tail of a static site-specific landscape and an adaptive change in a site-specific landscape both cause a site to be temporarily occupied by a less-than-optimal amino acid, say B. In either case the result is a transient increase in rate ratio to some value ω_B that decays exponentially while positive selection drives the site from B to the fittest amino acid A. Once A is fixed the rate ratio stabilizes to some value $\omega_A < \omega_B$. These processes manifest across the sites in an alignment as covarion-like switching between $\hat{\omega}_1 < \hat{\omega}_2$. The magnitude of $\hat{\omega}_2$ depends on the distribution of the ω_B , which in turn depends on the magnitude of the selection coefficients $s_{AB} = f_B - f_A < 0$. In Simulations 2 and 3, where sites were evolved using models based on the MS framework, the mean value of $\hat{\omega}_2$ was never less than one. This indicates that the magnitude of the s_{AB} tended to be large enough to make the $\omega_B > 1$. In the real data the rate ratio was always less than one, with $\hat{\omega}_2 = 0.31$ for the *mmtDNA*, $\hat{\omega}_2 = 0.40$ for phytochrome AC&F, and only $\hat{\omega}_2 = 0.08$ for cytochrome B. Sites in real proteins are undoubtedly subject to both intragenetic (e.g., Pollock *et al.*, 2012; Starr and Thornton, 2016) and intergenetic (e.g., Phillips, 2008) epistatic constraints. These can be difficult to model because they depend on unique aspects of the structure and function of a given protein as well as the nature of its interactions with other proteins. These and other potential sources of constraint are therefore absent in the majority

of generating models used in simulation studies (e.g., Anisimova *et al.*, 2001, 2002; Wong *et al.*, 2004; Zhang, 2004; Kosakovsky Pond and Frost, 2005; Yang *et al.*, 2005; Zhang *et al.*, 2005; Yang and dos Reis, 2011; Kosakovsky Pond *et al.*, 2011; Lu and Guindon, 2013), including those used in this chapter. Such constraints might have the effect pushing the s_{AB} closer to zero. For example, there is evidence that epistasis can cause the magnitude of s_{AB} at a site to diminish over time due to compensating substitutions at other sites (e.g., via an evolutionary Stokes shift, Pollock *et al.*, 2012). This can have the overall effect of reducing the ω_B . Differences in the depth of the tree might also have played a role in the lower rate ratios among the real alignments, since the real data was considerably more divergent than the simulated alignments (especially cytochrome B), and estimates of ω tend to diminish with larger divergences (dos Reis and Yang, 2013; Jones *et al.*, 2017). The PG-BSM managed to detect evidence of adaptive evolution in the mmtDNA and cytochrome B alignments in the form of site patterns consistent with either a reduction of the stringency of selection (the CW sites) or a change in the optimal amino acid (the BW sites) despite the small estimate of ω_2 . This was possible only because the model was designed to identify patterns of change in ω consistent with specific mechanisms of adaptation without imposing bounds on the magnitude of ω_2 . The YN-BSM A, by contrast, did not infer adaptive evolution in the mmtDNA precisely because it can do so only when there is evidence for $\omega_2 > 1$.

Like the vast majority of CSMs, the YN-BSM framework assumes evolution occurs via a series of single nucleotide substitutions (SNS). Consequently, whether or not a site is inferred to have undergone positive selection depends in part on the codon distribution implicitly inferred by the pruning algorithm at the two nodes of the FG branch. Positive selection is more often inferred when the codons that most likely occupied those two nodes differ by more than one SNS. Indeed, it was recently shown that the majority of support for positive selection in real data under the YN-BSM A consists of sites patterns that suggest multiple SNS along the FG (Venkat *et al.*, 2018). Yet instantaneous double and triple (DT) mutations can occur, with recent estimates indicating roughly 1% and 3% of all mutations being DT (Keightley *et al.*, 2009;

Schrider *et al.*, 2014; De Maio *et al.*, 2013; Harris and Nielsen, 2014). The chance fixation of a DT mutation along the FG can only be misconstrued by the YN-BSM A as evidence of multiple SNS. Hence, positive selection was often falsely inferred by the YN-BSM A in simulated alignments generated with rare fixation of DT mutations (Venkat *et al.*, 2018). It follows that positive selection due to genuine episodic peak shifts can be confounded not only by non-adaptive shifting balance (Jones *et al.*, 2017), but also by the fixation of DT mutations (Venkat *et al.*, 2018). The PG-BSM was specifically formulated with the understanding that evidence of positive selection in the form of $\omega > 1$ can result from multiple processes, some of which are non-adaptive. This was the point of the move away from the standard $\omega > 1$ paradigm. The PG-BSM is apparently robust to DT, since, although the inclusion of 6% DT mutations resulted in larger $\hat{\omega}_2$ compared to simulations with 0% DT, the omnibus test never incorrectly rejected the null.

The current trend in CSM development is toward greater realism via the addition of parameters that represent specific mechanistic processes (e.g., Liberles *et al.*, 2013; Zaheri *et al.*, 2014; Pollock *et al.*, 2017; Venkat *et al.*, 2018). The study presented in Chapter 3 suggests that this approach is not guaranteed to give better models (Jones *et al.*, 2018, 2019a). Under the ML framework, the addition of any parameter ψ to a null model M will always result in a better fit (i.e., larger likelihood). To guard against a spurious increase in likelihood, the null is rejected only if $LLR = 2 \left(\ln \left\{ L_{\text{alt}}(\hat{\theta}_M, \hat{\psi}) \right\} - \ln \left\{ L_{\text{null}}(\hat{\theta}_M) \right\} \right)$ is greater than some prespecified threshold chosen to limit the false positive rate to some maximum upper bound (e.g., 5%). The trend toward realism implicitly assumes that rejection of the null can only occur if the model with ψ provides a better representation of the actual data-generating process than the model without ψ . It is becoming increasingly clear that this assumption is incorrect due to the problem of confounding and phenomenological load, as was illustrated in Chapter 3.

The covarion-like component of the PG-BSM confers some robustness against PL. Under the null PG-BSM, the CL component accounts for all mechanisms that might generate heterotachy, whether adaptive (i.e., episodic peak shifts)

or non-adaptive (shifting balance, DT substitutions, epistasis). Hence, the parameters $\langle \omega_1, \omega_2, \delta \rangle$ for the CL process account for multiple mechanisms. By contrast, the parameters α and β in equation (4.1) have specific mechanistic interpretations as the rate at which instantaneous double and triple nucleotide substitution arise. It was recently suggested that existing CSMs should be modified to account for the possibility of DT substitutions (Venkat *et al.*, 2018). However, CSMs that include α and β as estimated parameters can result in false detection of DT substitutions due to PL (Jones *et al.*, 2018). PL is avoided under the PG-BSM by allowing that the MLEs for $\langle \omega_1, \omega_2, \delta \rangle$ result from an unknown combination of mechanisms, including DT substitutions. Hence, for example, finding that $\hat{\delta}$ is significantly > 0 in a contrast between the null PG-BSM with $\delta = 0$ versus the null PG-BSM with δ estimated need not be interpreted as evidence for any particular mechanism of heterotachy, but only for “heterotachy-by-any-cause”. In this way, the possibility of DT substitutions is subsumed in the parameters for the CL process, and false conclusions due to the confounding of processes that generate heterotachous site patterns are avoided. The PG-BSM framework therefore not only provides a means to identify site patterns consistent with specific adaptive mechanisms, but through the addition of external phenotypic information also offers a solution to several recently discovered problems associated with confounding and PL (Jones *et al.*, 2017, 2018; Venkat *et al.*, 2018).

4.7 Methods

4.7.1 Data Generation using MSmmtDNA and MSTGdR

Simulations under the MS framework were conducted as follows. First, parameters of the mutation process, including position-specific nucleotide frequencies and the transition/transversion rate ratio, were estimated from a real alignment consisting of 12 concatenated H-strand mitochondrial DNA sequences (3331 codon sites) from 20 mammalian species as distributed in alignment form by the PAML software package (Yang, 2007). Most alignments were simulated without fixation of DT mutations by setting $(\alpha, \beta) = (0, 0)$. In cases where alignments were generated with fixation of DT mutations,

$(\alpha, \beta) = (0.0371, 0.0030)$ so that $\approx 6\%$ of all mutations would be double ($\approx 5.8\%$) or triple ($\approx 0.2\%$). Recent studies suggest that DT mutations comprise between 1% and 3% of all mutations (Keightley *et al.*, 2009; Schrider *et al.*, 2014; De Maio *et al.*, 2013; Harris and Nielsen, 2014). A larger value (i.e., 6%) was used to investigate the impact of DT substitutions on model power and accuracy when unaccounted for. Next, vectors of amino acid fitness coefficients were drawn for each site using either MSmmtDNA (as described in Jones *et al.*, 2018) or MSTGdR (i.e., drawn with replacement from the set of 3598 vectors estimated from real mmtDNA using swMutSel with a Dirichlet-based penalty and with $\sigma = 0.1$ as described in Tamuri *et al.*, 2014). These were scaled and converted to site-specific substitution rate matrices as follows:

1. A scaling factor $\sigma^h \sim 0.001 + (0.01 - 0.001) \times B$ was drawn to determine the stringency of selection at the site, where $B \in [0, 1]$ is a beta random variable with shape parameters $u, v > 0$. Values of $\sigma^h \in [0.001, 0.01]$ closer to the upper bound correspond to greater selection stringency, whereas values closer to the lower bound correspond to a balance between selection and drift that typically results in non-adaptive shifting balance (Jones *et al.*, 2017). Parameters u and v for the beta distribution were chosen to make the distributions of scaled selection coefficients s_{ij} drawn under MSmmtDNA match those reported by Tamuri *et al.* (2012) as closely as possible (Jones *et al.*, 2018) (i.e., $u = 0.08$ and $v = 0.02$).
2. A vector \mathbf{f}^h of fitness coefficients for the 60 codons for the mammalian mitochondrial genetic code was then constructed from the amino acid fitnesses with the assumption that synonymous codons are equally fit. This vector was scaled to make its standard deviation equal to σ^h .
3. $\mathbf{f}^h = \langle f_1^h, \dots, f_{60}^h \rangle$ for the 60 codons was converted into a matrix W^h of fixation probabilities computed from scaled selection coefficients $s_{ij}^h = N_e(f_j^h - f_i^h)$ assuming an effective population size of $N_e = 1000$ and a ploidy of one for mtDNA:

$$W_{ij}^h \propto \begin{cases} 1 & \text{if } s_{ij}^h = 0 \\ \frac{2s_{ij}^h}{1 - \exp(-2s_{ij}^h)} & \text{otherwise} \end{cases} \quad (4.13)$$

The corresponding site-specific rate matrix $A^h = M \circ W^h$ was then constructed by taking the element-wise product of the matrix of mutation rates M and the matrix of fixation probabilities W^h .

The expected number of single nucleotide substitutions per codon per unit branch length at each site was then computed as follows:

$$r^h = \sum_{i \neq j} \pi_i A^h(i, j) \{\ell_1 + 2\ell_2 + 3\ell_3\} \quad (4.14)$$

where the indicator ℓ_s is one if i and j differ by $s \in \{1, 2, 3\}$ nucleotides and zero otherwise. When generating an alignment with no PG association, all rate matrices were divided by $\bar{r} = (1/n) \sum_{h=1}^n r^h$ to make branch lengths equal to the expected number of single nucleotide substitutions per codon.

A CW shift was implemented by reducing the stringency of selection to $\sigma^h = 0.0001$ at a subset of sites. Such shifts were made to occur at the ancestral node of the branch over which the phenotype changed and were made to persist along all descendant branches. This was intended to mimic an increase in the replacement rate among a subset of sites over a clade; rCW shifts were similarly implemented but with the stringency of selection increased to $\sigma^h = 0.01$. A BW shift was implemented by drawing new vectors of fitness coefficients for a subset of sites to mimic peak shifts. These new vectors were scaled to increase the stringency of selection to $\sigma^h = 0.01$ and were applied starting at the ancestral node of the branch over which the phenotype changed and along all descendant branches. When CW, rCW and/or BW shifts were implemented, all rate matrices were scaled by dividing by $\bar{r}_a = (1/n) \sum_{h=1}^n r^h$ with site-specific rates averaged over branches:

$$r^h = \frac{\sum_{b=1}^{2N-2} r^h(b) \mathbf{t}(b)}{\sum_{b=1}^{2N-2} \mathbf{t}(b)} \quad (4.15)$$

Here $r^h(b)$ is the rate for the site along branch b computed using (4.14) and $\mathbf{t}(b)$ is the length of that branch.

4.7.2 Generating Ancestral Phenotypes

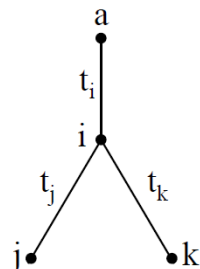


Figure 4.9: An arbitrary branching element in a binary tree. Node a is ancestral and i descendant.

Ancestral phenotypic states were sampled using the marginal approach described in Yang (2006) pp 121. Consider the root node of the tree where the state x_r is unknown. The probability of x_r conditioned on the vector of phenotypes \mathbf{F} at the terminal nodes of the tree can be expressed using Bayes' theorem as follows (omitting the rate parameter λ and the vector of branch lengths \mathbf{t} for brevity):

$$P(x_r = x \mid \mathbf{F}) = \frac{P(\mathbf{F} \mid x_r = x)P(x_r = x)}{P(\mathbf{F})} = \frac{P(\mathbf{F} \mid x_r = x)\pi_F^x}{\sum_{x=1}^3 P(\mathbf{F} \mid x_r = x)\pi_F^x} \quad (4.16)$$

Here we assume three discrete phenotypes. Note that the conditional probabilities $P(\mathbf{F} \mid x_r = x)$ are readily computed using the pruning algorithm (Felsenstein, 1981). The vector

$$\langle P(x_r = 1 \mid \mathbf{F}), P(x_r = 2 \mid \mathbf{F}), P(x_r = 3 \mid \mathbf{F}) \rangle$$

can be used to draw a realization of the state at the root node of the tree. This is used in turn to compute a vector of marginal probabilities for the two nodes that descend from the root node. Thus the algorithm moves inductively from the root to the terminal nodes of the tree.

Consider the i^{th} internal node of the tree. A realization x_a for the parent of the i^{th} node will already have been drawn. The conditional probability of the state at the i^{th} node can be computed as follows:

$$\begin{aligned} P(x_i = x \mid x_a, \mathbf{F}_i) &= \frac{P(\mathbf{F}_i \mid x_i = x)P(x_i = x \mid x_a)}{P(\mathbf{F}_i)} \\ &= \frac{P(\mathbf{F}_i \mid x_i = x)P_i(x_a, x)}{\sum_{x=1}^3 P(\mathbf{F}_i \mid x_i = x)P_i(x_a, x)} \end{aligned} \quad (4.17)$$

Here \mathbf{F}_i denotes the vector of phenotypes at terminal nodes that descend from the i^{th} node. $P_i(x_a, x)$ is the element of the transition probability matrix $P_i = \exp(t_i Q_F)$ corresponding to the $x_a \rightarrow x$ change of state, where t_i is the length of the branch connecting the i^{th} node to its parent. The vector

$$\langle P(x_i = 1 \mid x_a, \mathbf{F}_i), P(x_i = 2 \mid x_a, \mathbf{F}_i), x_a, P(x_i = 3 \mid \mathbf{F}_i) \rangle$$

can be used to draw a realization of the state at the i^{th} node. The algorithm continues in this way until all internal nodes have been assigned a phenotypic state. The resulting vector of states then constitutes one realization of the evolution of the phenotype conditioned on the values \mathbf{F} at the terminal nodes of the tree.

4.7.3 Computing Scaling Constants for the PG-BSM

All rate matrices included in the PG-BSM were constructed as follows for $k \in \{0, 1, 2\}$:

$$Q(\omega_k) = M \circ (\ell_S + \omega_k \ell_N) / r \quad (4.18)$$

The value of the common scaling constant r was specified so that estimated branch lengths would give the expected number of single nucleotide substitutions per codon. The scaling constant for any individual rate matrix depends only on ω_k and can be computed as follows:

$$r_k = \sum_{i \neq j} \pi_i Q_{ij}(\omega_k) \{\ell_1 + 2\ell_2 + 3\ell_3\} \quad (4.19)$$

If π_0 is the proportion of sites that evolved under $Q(\omega_0)$, $\pi_{\text{CL}} = 1 - \pi_0$ the proportion of sites that evolved covarion-like, and p_1 the proportion of time a CL site is expected to spend evolving under ω_1 , then the scaling constant for the null PG-BSM is:

$$r = \pi_0 r_0 + \pi_{\text{CL}} r_{\text{CL}} = \pi_0 r_0 + (1 - \pi_0)(p_1 r_1 + (1 - p_1) r_2) \quad (4.20)$$

Under the alternate PG-BSM accounting for CW and BW sites, the scaling constant is:

$$r = \pi_0 r_0 + (1 - \pi_0 - \pi_{\text{CW}} - \pi_{\text{BW}}) r_{\text{CL}} + \pi_{\text{CW}} r_{\text{CW}} + \pi_{\text{BW}} r_{\text{BW}} \quad (4.21)$$

Values for r_{CW} and r_{BW} can be calculated using an average that accounts for the rate ratio at a site on any particular branch weighted by the length of that branch, and the distribution of ancestral histories (i.e., via change maps). For example, the scaling factor for the BW process is:

$$r_{\text{BW}} = \sum_{\mathbf{z}} \hat{\pi}_{\mathbf{z}} \frac{\sum_{b=1}^{2N-2} r(z_b) \mathbf{t}(b)}{\sum_{b=1}^{2N-2} \mathbf{t}(b)} \quad (4.22)$$

where $r(z_b) = r_1$ if $z_b = 0$ and r_2 if $z_b = 1$. Scaling factors r_{CW} and r_{rCW} are similarly computed by taking into account branches over which the site was evolved under ω_1 or ω_2 .

4.7.4 False Discovery Control

The algorithm described by Newton *et al.* (2004) was applied to NEB posteriors to control the false discovery count (FDC) for our *post hoc* analyses. The objective was to identify as many sites with evidence of PG association as possible while controlling the number of false discoveries to some specified value. The procedure was as follows, here described for BW sites:

1. Suppose $p^h = 1 - \text{P}(\text{BW} \mid \mathbf{x}^h)$ is the conditional probability that assigning the h^{th} site to the BW category is a type I error.
2. Let p^1, \dots, p^n be a list of these probabilities for all n sites.
3. For any specified bound κ , assign a site to the BW category if $p^h \leq \kappa$.
4. By this rule, the expected number of false discoveries given the data is:

$$\text{E}\{\text{FDC}\} = \sum_{h=1}^n p^h \ell(p^h \leq \kappa) \quad (4.23)$$

where $\ell(p^h \leq \kappa)$ is 1 if $p^h \leq \kappa$ and zero otherwise (Newton *et al.*, 2004).

5. To control the FDC to be no more than $k \in \{1, 2\}$, κ can be set to the largest value for which $\text{E}\{\text{FDC}\} \leq k$. Note that this expectation is across data sets, and is only approximate because it depends on how well the fitted model matches the data-generating process (Newton *et al.*, 2004).

In practice p^h was approximated using the NEB approach in equation 4.8.

4.7.5 Dealing with Underflow

Likelihood functions are typically optimized in log-space. For example, the log-likelihood for the alternate component of PG-BSM is:

$$\ln \{L_{\text{alt}}(X, \mathbf{F}; \lambda, \boldsymbol{\theta}, \mathbf{t})\} = \ln \{P(\mathbf{F}; \lambda, \mathbf{t})\} + \ln \left\{ \sum_{\mathbf{z}} \hat{\pi}_{\mathbf{z}} \prod_{h=1}^n g(x^h; \mathbf{z}) \right\} \quad (4.24)$$

It is convenient to express the second addend of (4.24) as follows:

$$\ln \left\{ \sum_{\mathbf{z}} \hat{\pi}_{\mathbf{z}} \prod_{h=1}^n g(x^h; \mathbf{z}) \right\} = \ln \left\{ \sum_{\mathbf{z}} \hat{\pi}_{\mathbf{z}} \exp(\ell_{\mathbf{z}}) \right\} \quad (4.25)$$

$$\text{where } \ell_{\mathbf{z}} = \ln \left\{ \prod_{h=1}^n g(x^h; \mathbf{z}) \right\} \quad (4.26)$$

A problem arises when the probability $\exp(\ell_{\mathbf{z}})$ is too small to be represented digitally, an issue commonly referred to as underflow. This can be mitigated by a simple transformation:

$$\ln \left\{ \sum_{\mathbf{z}} \hat{\pi}_{\mathbf{z}} \prod_{h=1}^n g(x^h; \mathbf{z}) \right\} = \max\{\ell_{\mathbf{z}}\} + \ln \left\{ \sum_{\mathbf{z}} \hat{\pi}_{\mathbf{z}} \exp(\ell_{\mathbf{z}} - \max\{\ell_{\mathbf{z}}\}) \right\} \quad (4.27)$$

Since $\forall \mathbf{z} \ell_{\mathbf{z}} < 0$, the transformation $\ell_{\mathbf{z}} - \max\{\ell_{\mathbf{z}}\}$ moves the log-probabilities to the right toward zero (i.e., toward larger values), making underflow less likely.

The transformation method can also be used to avoid underflow in equation (4.10). Let \mathbf{z}^* represent the change map that maximizes $\ell_{\mathbf{z}}$ (i.e. maximizes the likelihood of the alignment). The natural log of (4.10) at \mathbf{z}^* is:

$$\begin{aligned} & \ln \{P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t})\} \\ &= \ln \{L_{\text{alt}}(X, \mathbf{F} | \mathbf{z}^*)\} + \ln \{\hat{\pi}_{\mathbf{z}^*}\} - \ln \{L_{\text{alt}}(X, \mathbf{F})\} \end{aligned} \quad (4.28)$$

Applying the transformation:

$$\begin{aligned} & \ln \{L_{\text{alt}}(X, \mathbf{F})\} \\ &= \ln \{P(\mathbf{F}; \lambda, \mathbf{t})\} + \ell_{\mathbf{z}^*} + \ln \left\{ \sum_{\mathbf{z}} \hat{\pi}_{\mathbf{z}} \exp(\ell_{\mathbf{z}} - \ell_{\mathbf{z}^*}) \right\} \end{aligned} \quad (4.29)$$

$$= \ln \{L_{\text{alt}}(X, \mathbf{F} | \mathbf{z}^*)\} + \ln \left\{ \sum_{\mathbf{z}} \hat{\pi}_{\mathbf{z}} \exp(\ell_{\mathbf{z}} - \ell_{\mathbf{z}^*}) \right\} \quad (4.30)$$

Substituting into (4.28) gives:

$$P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t}) = \frac{\hat{\pi}_{\mathbf{z}^*}}{\hat{\pi}_{\mathbf{z}^*} + \sum_{\mathbf{z} \neq \mathbf{z}^*} \hat{\pi}_{\mathbf{z}} \exp(\ell_{\mathbf{z}} - \ell_{\mathbf{z}^*})} \quad (4.31)$$

Equation (4.31) demonstrates that the strength of the evidence for the ancestral reconstruction corresponding to \mathbf{z}^* is a function of the relative frequency of the most frequently sampled change map $\hat{\pi}_{\mathbf{z}^*}$ and the differences $\ell_{\mathbf{z}} - \ell_{\mathbf{z}^*}$ in log-likelihoods of the alignment under the various other \mathbf{z} . When evidence for the PG processes dictated by \mathbf{z}^* is strong, it will be the case $\forall \mathbf{z} \neq \mathbf{z}^*$ that $\exp\{\ell_{\mathbf{z}} - \ell_{\mathbf{z}^*}\} \approx 0$ making $P(\mathbf{z}^* | X, \mathbf{F}, \lambda, \mathbf{t}) \approx 1$.

Chapter 5

Discussion

5.1 Historical Development of CSMs.

The evolution of CSMs can be divided in two phases. Phase I began with the pioneering efforts of Muse and Gaut (1994) and Goldman and Yang (1994) when they proposed CSMs to estimate a single rate ratio ω from an alignment. Subsequent models developed during this phase account for variations in ω across branches (Nielsen and Yang, 1998) across sites (Yang *et al.*, 2000a), and across both branches and sites (Yang and Nielsen, 2002; Zhang *et al.*, 2005), each by including multiple ω -categories. The use of multiple categories permits *post hoc* identification of sites and/or branches where $\omega > 1$ was likely to have occurred. The characteristic feature of Phase I was the development of simple CSMs (e.g., with few parameters) tested using data-generating algorithms based on the same CSM framework. The main concern during Phase I was the possibility of a false detection of $\omega > 1$ due either to model misspecification (the mis-match between the model and the data-generating process) or low information content. Phase II is marked by a leap in model complexity with the introduction of several parameter-rich CSMs (Kosakovsky Pond *et al.*, 2011; Zaheri *et al.*, 2014; Murrell *et al.*, 2015; Smith *et al.*, 2015). The main concern in the ongoing Phase II is the growing realization of problems caused by confounding and PL. Problems and solutions associated with Phase I and Phase II concerns are illustrated in this section via a series of case studies.

5.1.1 Phase I: Simple Models, Simple Problems

The first effort to detect positive selection at the molecular level (Hughes and Nei, 1988) relied on heuristic counting methods (Nei and Gojobori, 1986). Phase I of CSM development followed with the introduction of formal statistical approaches based on ML (Muse and Gaut, 1994; Goldman and Yang,

1994). The first CSMs were used to infer whether the estimate $\hat{\omega}$ of a single nonsynonymous-to-synonymous substitution rate ratio averaged over all sites and branches was significantly greater than one. This approach was found to have low power due to the pervasiveness of synonymous substitutions at most sites within a typical gene (Yang and Bielawski, 2000). An early attempt to increase the statistical power to infer positive selection was the CSM designed to detect $\hat{\omega} > 1$ on specific branches (Yang and Nielsen, 1998). Models accounting for variations in ω across sites were subsequently developed, the most prominent of which are the M-series CSMs (Yang *et al.*, 2000a). These were accompanied by methods to identify individual sites under positive selection. The quest for power culminated in the development of models that can be used to identify branches over which specific sites evolved under positive selection (e.g., Yang and Nielsen, 2002; Forsberg and Christiansen, 2003; Bielawski and Yang, 2004; Zhang *et al.*, 2005). Two case studies are employed to illustrate some of the inferential challenges associated with Phase I models. Case Study A exemplifies the problem of low information content. The subject of the study is the M1a vs M2a model contrast applied to the *tax* gene of the human T-cell lymphotropic virus type I (HTLV-I Suzuki and Nei, 2004; Yang *et al.*, 2005). Case Study B illustrates how model misspecification can lead to false inferences. The subject of this study is the Yang-Nielsen Branch-Site Model (YN-BSM, Yang and Nielsen, 2002) applied to simulated data.

Case Study A: Low Information Content

An example setting with low information content arises when there are a substantial number of invariant sites, since these provide little information about the substitution process. Consider the pair of nested M-series models known as M1a and M2a (Wong *et al.*, 2004; Yang *et al.*, 2005). Under M1a sites are partitioned into two categories, $0 < \omega_0 < 1$ and $\omega_1 = 1$ in proportions p_0 and $p_1 = 1 - p_0$. M2a includes an additional category for a proportion of sites $p_2 = 1 - p_0 - p_1$ that evolved under positive selection with $\omega_2 > 1$. The issue of low information content is well illustrated by the application of the M1a vs M2a contrast to the *tax* gene, HTLV-I (Suzuki and Nei, 2004). The alignment

consists of 20 sequences with 181 codon sites, 158 of which are invariant. The 23 variable sites appear to have undergone only one substitution each: 2 are synonymous and 21 are nonsynonymous. The high nonsynonymous-to-synonymous substitution ratio suggests that the gene underwent positive selection. This hypothesis was supported by analytic results: the LLR for the M1a vs M2a contrast was 6.96 corresponding to a p-value of approximately 0.03 (Yang *et al.*, 2005). However, the MLE for p_2 under M2a was $\hat{p}_2 = 1$. Using this in the *post hoc* analysis (i.e., using the method modeled by equation 1.34) gives a posterior probability of $\Pr(\omega > 1 \mid \mathbf{x}, \hat{\boldsymbol{\theta}}_{\text{M2a}}) = 1$ for all sites, including the 158 invariable sites. Such an unreasonable result can occur under NEB because $\hat{\boldsymbol{\theta}}_{\text{M2a}}$ is treated as known despite the possibility of large sampling errors in the MLEs.

The analysis of the *tax* gene led to the development of methods such as Bayes empirical Bayes (BEB, Yang *et al.*, 2005), smoothed bootstrap aggregation (SBA, Mingrone *et al.*, 2016), and the penalized likelihood ratio test (Mingrone *et al.*, 2018) to reduce errors in estimates of posterior probabilities associated with low information content. Using BEB in the analysis of the *tax* gene, for example, resulted in posterior probabilities $0.91 < \Pr(\omega > 1 \mid \mathbf{x}, \hat{\boldsymbol{\theta}}_{\text{M2a}}) < 0.93$ for the 21 sites with a single nonsynonymous change and $0.55 < \Pr(\omega > 1 \mid \mathbf{x}, \hat{\boldsymbol{\theta}}_{\text{M2a}}) < 0.61$ for the remaining sites (Yang *et al.*, 2005). Results such as these undoubtedly promoted confidence in the ability of CSMs to deal with low information content, which at the time was largely associated with the obvious case of low divergence. Interestingly, problems associated with low information content can also occur in alignments exhibiting an abundance of variant sites. The problem of confounding between non-adaptive shifting balance and episodic adaptation (Chapter 2), for example, is not due to a lack of variant site patterns, but to a lack of contextual information. The addition of such information in the form of a discrete phenotype was shown to break confounding and permit detection of changes in site-specific landscapes that co-occurred with changes in phenotype (Chapter 4). The discovery of this novel low-information scenario would have to wait until the increase in the complexity of analytic CSMs that marked the beginning of Phase II was

matched by the use of more realistic data-generating methods that reveal the problem of confounding (i.e., Chapter 2).

Case Study B: Model Misspecification

The mechanisms that give rise to the diversity of site patterns in a set of homologous genes are highly complex and not fully understood. CSMs are therefore necessarily simplified representations of the true generating process, and are in this sense misspecified. The extent to which misspecification might cause an omnibus test to falsely detect positive selection was of primary concern during Phase I. The first branch-site model (the YN-BSM, Yang and Nielsen, 2002) illustrates this issue. The YN-BSM in its original form assumes a null under which a proportion p_0 of sites evolved under stringent selection with $\omega_0 = 0$ and the remaining sites evolved under a neutral regime with $\omega_1 = 1$ on all branches of the tree (i.e., model M1 in Nielsen and Yang, 1998). This is contrasted with Model A, which is the same as M1 except that it assumes that some stringent sites and some neutral sites evolved under positive selection with $\omega_2 > 1$ on a pre-specified branch called the foreground (FG) branch. The original omnibus test contrasting M1 with Model A was therefore designed to detect a subset of sites that evolved adaptively on the same FG branch.

ω regime X	1.00	1.00	0.80	0.80	0.50	0.50	0.20	0.20	0.00	0.00
ω regime Z	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00

Table 5.1: Simulated selection regimes. Rate ratios (ω) for regimes X and Z taken from Zhang (2004). Each rate ratio applies to 20 codon sites. Alignments were generated with 100 sites in total.

The standard method used to test the impact of misspecification on the reliability of the LRT during Phase I was to generate alignments using a more complex version of the CSM to be tested. This usually involved adding more variability in ω across sites and/or branches than assumed by the fitted CSM while leaving all other aspects of the generating model the same. In Zhang (2004), for example, alignments were generated using site-specific rate matrices $Q(\omega)$ with rate ratios ω specified by pre-determined selection regimes, two of which are shown in Table 5.1. In one simulation, 200 alignments were

generated using regime Z on a single foreground branch and regime X on all of the remaining branches of a 10 or 16 taxon tree. The gene therefore underwent a mixture of stringent and neutral evolution over most of the tree (regime X), but with complete relaxation of selection pressure on the foreground branch (regime Z). Positive selection did not occur at any of the sites. Nevertheless, the M1 vs Model A contrast inferred positive selection in 20% to 55% of the alignments, depending on the location of the foreground branch. Such a high rate of false positives was attributed to the mismatch between the process used to generate the data compared to the process assumed by the null model M1 (Zhang, 2004).

The branch-site model was subsequently modified to allow $0 < \omega_0 < 1$ instead of $\omega_0 = 0$ (Modified Model A in Zhang *et al.*, 2005). Furthermore, the new null model assumes that some proportion p_0 of sites (the stringent sites) evolved with $0 < \omega_0 < 1$ everywhere in the tree except on the foreground branch, where those same sites evolved neutrally with $\omega_2 = 1$. All other sites in the alignment (the neutral sites) are assumed to have evolved neutrally with $\omega_1 = 1$ everywhere in the tree. This is contrasted with the Modified Model A, which assumes that some of the stringent sites and some of the neutral sites evolved under positive selection with $\omega_2 > 1$ on the FG. Hence, unlike the original omnibus test that contrasts M1 with Model A, the new test contrasts Modified Model A with $\omega_2 = 1$ against Modified Model A with $\omega_2 > 1$. These changes to the YN-BSM were shown to mitigate the problem of false inference. For example, using the same generating model with regimes X and Z, the modified omnibus test falsely inferred positive selection in only 1% to 7.5% of the alignments, consistent with the 5% level of significance of the test (Zhang *et al.*, 2005).

This case study demonstrates how problems associated with model misspecification were traditionally identified, and how they could sometimes be corrected through relatively minor changes to the fitted model. However, the generating methods employed by studies such as Zhang (2004) and Zhang *et al.* (2005), although sophisticated for their time, can only produce alignments that are highly unrealistic compared to real data. In particular, such

methods fail to generate heterotachy beyond simple branch-wise changes similar to those implied by the selection regimes in Table 5.1. While the mitigation of statistical pathologies due to low information content (e.g., using BEB) was a critical advancement during Phase I of CSM development, other statistical pathologies went unrecognized due to reliance on such unrealistic simulation methods.

5.1.2 Phase II: The Rise in Complexity

A typical protein-coding gene evolves adaptively only episodically (Struder and Robinson-Rechavi, 2009). The evidence of adaptive evolution of this type can be very difficult to detect. For example, it is assumed under the YN-BSM that an unknown subset of sites switched from a stringent or neutral selection regime to positive selection together on the same set of foreground branches. The power to detect a signal of this kind can be very low when the proportion of sites that switched together is small (Yang and dos Reis, 2011). Perhaps encouraged by the reliability of Phase I models demonstrated by extensive simulation studies (Anisimova *et al.*, 2001, 2002; Wong *et al.*, 2004; Zhang, 2004; Kosakovsky Pond and Frost, 2005; Yang *et al.*, 2005; Zhang *et al.*, 2005; Yang and dos Reis, 2011; Kosakovsky Pond *et al.*, 2011; Lu and Guindon, 2013), combined with experimental validation of results obtained from their application to real data (Yang and Bielawski, 2000; Yang, 2005; Anisimova and Kosiol, 2009), investigators began to formulate increasingly complex and parameter-rich CSMs (Rodrigue *et al.*, 2010; Kosakovsky Pond *et al.*, 2011; Tamuri *et al.*, 2012, 2014; Rodrigue and Lartillot, 2014; Murrell *et al.*, 2015; Smith *et al.*, 2015). The hope was that carefully selected increases in model complexity would yield greater power to detect subtle signatures of positive selection overlooked by Phase I models. Phase II models fall into three broad categories:

1. Phase I CSMs modified to account for more variability in selection effects across sites and branches than previously assumed with the aim of increasing the power to detect subtle signatures of positive selection (e.g., the Branch-Site Random Effects Likelihood model, Kosakovsky Pond

et al., 2011).

2. Phase I CSMs modified to contain parameters for mechanistic processes not directly associated with selection effects. Many such models have been motivated either by a particular interest in the added mechanism (e.g., the fixation of double and triple mutations, Miyazawa, 2011; Zaheri *et al.*, 2014; Jones *et al.*, 2018), or by the notion that increasing the mechanistic content of a CSM can only improve inferences about selection effects (e.g., by accounting for variations in the synonymous substitution rate, Kosakovsky Pond and Muse, 2005; Rubinstein *et al.*, 2011).
3. Models that abandon the traditional CSM approach in favor of a substitution process expressed in terms of explicit population genetic parameters, such as population size and selection coefficients (Nielsen and Yang, 2003; Rodrigue *et al.*, 2010; Tamuri *et al.*, 2012, 2014; Rodrigue and Lartillot, 2014, 2016).

An example of the first category of models is BUSTED (for Branch-site Unrestricted Statistical Test for Episodic Diversification, Murrell *et al.*, 2015), which is used to illustrate the problem of confounding in Case Study C. The second category of models includes those formulated by adding parameters for the rate of double and triple mutations to a traditional CSM, an example of which is RaMoSSwDT (for Random Mixture of Static and Switching sites with Double and Triple substitutions, Jones *et al.*, 2018). This model is used in Case Study D to illustrate the problem of phenomenological load. Models in the third category are the most ambitious CSMs currently in use, and are far more challenging to fit to real alignments than traditional models. One of the most impressive examples is the site-wise mutation-selection model (swMutSel: Tamuri *et al.*, 2012, 2014). Based on the mutation-selection framework of Halpern and Bruno (1998), swMutSel estimates a vector of selection coefficients for each site in an alignment. This and similar models (e.g., Rodrigue *et al.*, 2010; Rodrigue and Lartillot, 2014, 2016) appear to be reliable (Spielman and Wilke, 2016), but require a very large number of taxa (e.g., several

hundred). Phase II models of this category are therefore impractical for the majority of empirical datasets.

Case Study C: Confounding

The MS framework facilitates investigation of complex evolutionary dynamics, such as non-adaptive shifting balance on a fixed fitness landscape or adaptation to a change in selective constraints (i.e., a peak shift, dos Reis, 2013; Jones *et al.*, 2017). The phenomenological outcomes of these processes are difficult to mimic in alignments generated using traditional methods. MS can therefore be used to generate more variation in rate ratio across sites and over time than has been realized in past simulation studies (e.g., Table 5.1) and can be adjusted to produce alignments that closely mimic real data, as was shown in Chapter 3. It was only when the MS framework was used to generate data for the purpose of model testing that the problem of confounding between adaptive and non-adaptive processes was revealed. BUSTED (Murrell *et al.*, 2015), for example, is a model that was intended to detect episodic adaptive evolution by accounting for random variations in the intensity of selection over sites and branches. The rate ratio at each site/branch combination is assumed to be an independent draw from the distribution $\{(\omega_0, p_0), (\omega_1, p_1), (\omega_2, p_2)\}$. The model contrasts the null hypothesis that $\omega_0 \leq \omega_1 \leq \omega_2 = 1$ with the alternative that $\omega_0 \leq \omega_1 \leq 1 \leq \omega_2$. As was reported in Chapter 2, BUSTED inferred episodic positive selection due to non-adaptive shifting balance in as many as 40% of MS-generated alignments. Whereas positive selection did likely occur in those alignments, it could only have been due to non-adaptive shifting balance. Hence, the intended interpretation of BUSTED as a test for episodic adaptive evolution is negated by the possibility of confounding. Confounding between non-adaptive shifting balance and episodic peak shifts can be broken by the addition of contextual information (e.g., phenotype), as was demonstrated in Chapter 4.

Case Study D: Phenomenological Load

PL was illustrated by the RaMoSS vs RaMoSSwDT contrast in Chapter 3. Recall that RaMoSS is a mixture of the standard M-series model $M3(k = 2)$

with the covarion-like model CLM3($k = 2$), and that RaMoSSwDT includes two extra parameters (α, β) to account for the fixation of double and triple mutations. RaMoSS and RaMoSSwDT were fitted to fifty alignments simulated to mimic a real alignment of mammalian mtDNA using MSmmtDNA as the generating model. Since DT substitutions are not permitted under MSmmtDNA, any reduction in deviance (expressed as the percentage PRD, equation 3.7) caused by $(\hat{\alpha}, \hat{\beta})$ in each trial could only be attributed to PL plus noise. Non-adaptive shifting balance can produce site patterns similar to those produced by a process that includes DT substitutions¹. DT was consequently falsely inferred in 48 of 50 trials at the 5% level of significance. PL was only identified as an issue when model contrasts were fitted to data generated with realistic evolutionary dynamics using the MS framework.

¹It has previously been noted that the rapid fixation of compensatory mutations following substitution to an unstable base pair (e.g., AT→GT→GC) can also produce site patterns that suggest fixation of DT mutations (Yang, 2014, page 46).

5.1.3 An Argument for Phenomenological CSMs

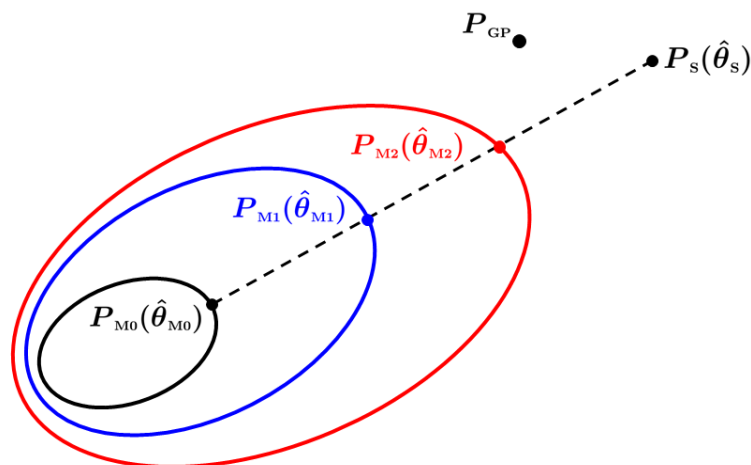


Figure 5.1: A cartoon of the $(61^N - 1)$ -dimensional simplex containing all possible site-pattern distributions for an N -taxon alignment. The inner-most ellipse represents the subspace $\{P_{M_0}(\theta_{M_0}) \mid \theta_{M_0} \in \Omega_{M_0}\}$ that is the family of distributions that can be specified using M_0 , the simplest of CSMs. This is nested in the family of distributions that can be specified using M_1 (blue ellipse), a hypothetical model that has the same parameters as M_0 plus some extra parameters. Similarly, M_1 is nested in M_2 (red ellipse). Whereas models are represented by subspaces of distributions, the true generating process is represented by a single point P_{GP} , the location of which is unknown. The empirical site-pattern distribution $P_S(\hat{\theta}_S)$ corresponds to the fitted saturated model; with large samples $P_S(\hat{\theta}_S) \approx P_{GP}$. For any other model M , the member $P_M(\hat{\theta}_M) \in \{P_M(\theta_M) \mid \theta_M \in \Omega_M\}$ most consistent with X is the one that minimizes deviance, which is twice the difference between the maximum log-likelihood of the data under the saturated model and the maximum log-likelihood of the data under M .

CSMs have become increasingly complex with the addition of more free parameters since the introduction of the M -series models in Yang *et al.* (2000a). The *prima facie* objective of this trend is to produce models that provide better mechanistic explanations of the data. It is assumed that this will lead to more accurate inferences about evolutionary processes, particularly as the volume of genetic data increases (Liberles *et al.*, 2013; Pollock *et al.*, 2017). However, the significance of a new model parameter is assessed under the ML framework by a comparison of site-pattern distributions without reference to mechanism. Combined with the possibility of confounding, this feature of maximum likelihood means that the objective of improving model fit does not necessarily coincide with the objective of providing a better mechanistic representation of the true generating process.

Given any CSM with parameters θ it is possible to compute a vector that assigns to each of the 61^N possible site patterns \mathbf{x} for an N -taxon alignment a probability $P(\mathbf{x}; \theta)$ such that $\sum_{\mathbf{x}} P(\mathbf{x}; \theta) = 1$ (i.e., a multinomial distribution for 61^N categories assuming the standard genetic code). Figure 5.1 depicts the space of all possible site-pattern distributions for an N -taxon alignment. Each ellipse represents the family of distributions $\{P_M(\theta_M) \mid \theta_M \in \Omega_M\}$, where $P_M(\theta_M)$ is the site-pattern distribution for model M given θ_M and Ω_M is the vector-space of all possible values of θ_M . For example, $\{P_{M_0}(\theta_{M_0}) \mid \theta_{M_0} \in \Omega_{M_0}\}$ is the family of distributions that can be specified using M_0 , the simplest CSM that assumes a common substitution rate matrix $Q(\omega)$ for all sites and branches. This is nested inside $\{P_{M_1}(\theta_{M_1}) \mid \theta_{M_1} \in \Omega_{M_1}\}$, where M_1 is a hypothetical model that is the same as M_0 but for a few extra parameters. Likewise, M_1 is nested in M_2 . The location of the site-pattern distribution for the true generating process is represented by P_{PG} in Figure 5.1. Its location is fixed but unknown. It is therefore not possible to assess the distance between it and any other distribution. Instead, comparisons are made using the site-pattern distribution inferred under the saturated model.

Whereas $\{P_M(\theta_M) \mid \theta_M \in \Omega_M\}$ represents a family of multinomial distributions, the fitted saturated model $P_s(\hat{\theta}_s)$ is the unique multinomial distribution defined by the MLE $\hat{\theta}_s = (y_1/n, \dots, y_m/n)^T$, where $y_i > 0$ is the observed frequency of the i^{th} site pattern, m is the number of unique site patterns, and n is the number of codon sites. In other words, the fitted saturated model is the empirical site-pattern distribution for a given alignment. Because it takes none of the mechanisms of mutation or selection into account, ignores the phylogenetic relationships between sequences, and excludes the possibility of site patterns that were not actually observed (i.e., $y_i/n = 0$ for site patterns i not observed in X), $P_s(\hat{\theta}_s)$ can be construed as the maximally phenomenological explanation of the observed alignment. An alignment is always more likely under the saturated model than it is under any other CSM. $P_s(\hat{\theta}_s)$ therefore provides a natural benchmark for model improvement.

The MLE over the family of distributions $\{P_M(\theta_M) \mid \theta_M \in \Omega_M\}$ is represented by a fixed point $P_M(\hat{\theta}_M)$ in Figure 5.1. $P_M(\hat{\theta}_M)$ is the distribution that

minimizes the statistical deviance between $P_M(\boldsymbol{\theta}_M)$ and $P_S(\hat{\boldsymbol{\theta}}_S)$ for any given data set. Deviance is defined as twice the difference between the maximum log-likelihood (LL) of the data under the saturated model and the maximum log-likelihood of the data under M:

$$D(\hat{\boldsymbol{\theta}}_M, \hat{\boldsymbol{\theta}}_S) = 2 \left\{ \text{LL}(\hat{\boldsymbol{\theta}}_S | X) - \text{LL}(\hat{\boldsymbol{\theta}}_M | X) \right\} \quad (5.1)$$

This is represented by the distance between $P_M(\hat{\boldsymbol{\theta}}_M)$ and $P_S(\hat{\boldsymbol{\theta}}_S)$ in Figure 5.1. A key feature of deviance is that it always decreases as more parameters are added to the model, corresponding to an increase in the probability of the data under that model. For example, suppose $\{P_{M_2}(\boldsymbol{\theta}_{M_2}) \mid \boldsymbol{\theta}_{M_2} \in \Omega_{M_2}\}$ is the same family of distributions as $\{P_{M_1}(\boldsymbol{\theta}_{M_1}) \mid \boldsymbol{\theta}_{M_1} \in \Omega_{M_1}\}$ but for the inclusion of one additional parameter ψ , so that $\boldsymbol{\theta}_{M_2} = (\boldsymbol{\theta}_{M_1}, \psi)$. The improvement in the probability of the data under $P_{M_2}(\hat{\boldsymbol{\theta}}_{M_2})$ over its probability under $P_{M_1}(\hat{\boldsymbol{\theta}}_{M_1})$ is assessed by the size of the reduction in deviance induced by ψ :

$$\begin{aligned} \Delta D(\hat{\boldsymbol{\theta}}_{M_1}, \hat{\boldsymbol{\theta}}_{M_2}) &= D(\hat{\boldsymbol{\theta}}_{M_1}, \hat{\boldsymbol{\theta}}_S) - D(\hat{\boldsymbol{\theta}}_{M_2}, \hat{\boldsymbol{\theta}}_S) \\ &= 2 \left\{ \text{LL}(\hat{\boldsymbol{\theta}}_{M_2} | X) - \text{LL}(\hat{\boldsymbol{\theta}}_{M_1} | X) \right\} \end{aligned} \quad (5.2)$$

Equation (5.2) is just the familiar log-likelihood ratio (LLR) used to compare nested models under the maximum likelihood framework.

Given this measure of model improvement, the *de facto* objective of model building is not to move closer to the true generating process (e.g., by adding parameters for mechanisms thought to have occurred when the data was generated), but only to move closer to the site-pattern distribution corresponding to the fitted saturated model, (which is maximally phenomenological and therefore minimally explanatory). Real alignments are limited in size, so there will always be some distance between $P_S(\hat{\boldsymbol{\theta}}_S)$ and P_{GP} due to sampling error (as represented in Figure 5.1). But even with an infinite number of codon sites, when $P_S(\hat{\boldsymbol{\theta}}_S)$ converges to P_{GP} , the criterion of minimizing deviance does not inevitably lead to a better explanation of the data because of the possibility of confounding. As stated in Chapter 3, two process are said to be confounded if they can produce similar patterns in the data. Hence, if ψ represents a process that did not actually occur when the data was generated, and if that process

is confounded with another process that did occur, the LLR in equation (5.2) can still be significant. Under this scenario, the addition of ψ to M1 would engender movement toward $P_s(\hat{\theta}_s)$ and P_{GP} , but the new model M2 would also provide a *worse* mechanistic explanation of the data because it would falsely indicate that the process represented by ψ actually occurred.

Models of molecular evolution require validation. The data-generating processes that gave rise to any particular real alignment are largely unknown, and so the most expedient method of validation is to generate alignments *in silico*. Sophisticated generating methods include those based on the MS framework represented in this thesis by MSmmtDNA and MSTGdR (Chapter 4), and various models based on physical considerations such as thermodynamic stability Pollock *et al.* (e.g., 2012). The possibility of confounding and PL suggests that it is unlikely that a CSM can ever produce MLEs that match all of the parameters used in such models. It is more realistic to expect a CSM to provide phenomenological summaries of major sources of variation that can be meaningfully interpreted in terms of mutation and selection processes. The null PG-BSM, for example, captures a major component of variation with a simple phenomenological model for heterotachy (i.e., the CLM3 component). This null provides an effective contrast for the alternate PG-BSM that accounts for specific modes of heterotachy consistent with the phenomenological outcomes expected to arise from changes in site-specific amino acid fitness coefficients that co-occurred with a change in phenotype. In this way, the relatively simple PG-BSM provides a meaningful phenomenological summary of the most biologically salient result (i.e., adaptation) of a complex generating process. As the volume of genetic data increases there is a natural tendency to want to make CSMs more complex with the addition of a greater number of ostensibly mechanistic parameters. I maintain that this is a mistake, and that the appropriate place for greater mechanistic realism is in alignment-generating processes used to test what should remain largely phenomenological CSMs.

5.2 Other Examples of Confounding and PL

Confounding was a common theme in the material presented in Chapters 2, 3 and 4. It first appeared in Chapter 3. There it was argued that positive selection caused by non-adaptive shifting balance can be difficult to distinguish from positive selection caused by adaptive changes in site-specific landscapes based on an analysis of alignment data alone. This is illustrated in Figure 5.2. The red dashed line in Figures 5.2A and 5.2B mark the point where the codon-specific rate ratio changes from a value $\omega < 1$ (to the left of the mark) to a value $\omega > 1$ (to the right of the mark). Phenylalanine (F, TTT) is the fittest amino acid in the landscape depicted in Figure 5.2A. During the course of evolution it can sometimes happen that the population becomes fixed at an amino acid such as valine (V, GTT) in the tail of that landscape. This will be followed by a temporary elevation in the rate ratio at the site to a value greater than one as positive selection moves the site back to F by a series of replacement substitutions e.g., V (GTT) \rightarrow G (GGT) \rightarrow C (TGT) \rightarrow F (TTT). Similarly, a change in one or more external factors that impact the functional significance of the site can change the landscape from that depicted in Figure 5.2A to that depicted in Figure 5.2B where glutamine (Q) is fittest. If at the time of the change the site is occupied by F, then the site-specific rate ratio would similarly be elevated to a value $\omega > 1$ as positive selection moves the site toward its new peak at Q e.g., F (TTT) \rightarrow Y (TAT) \rightarrow H (CAT) \rightarrow Q (CAA). Both processes can manifest a detectable elevation in rate ratio to $\omega > 1$.

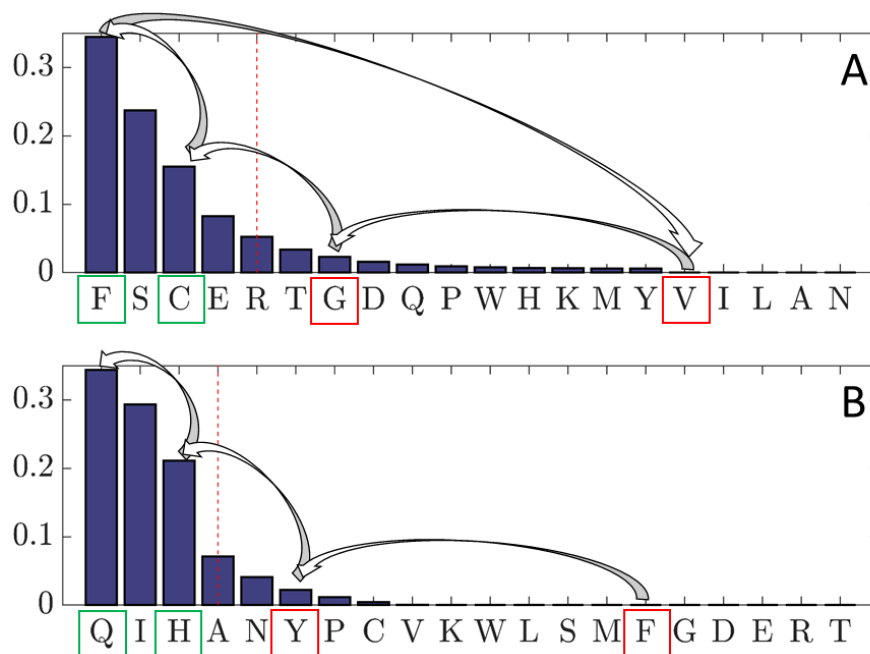


Figure 5.2: An Example of Confounding: Non-adaptive drift into the tail of a static site-specific fitness landscape (A) and an adaptive shift at a site from the landscape in (A) to the landscape in (B) can both be followed by a rapid series of substitutions that take the site to its optimal amino acid. It follows that both processes can generate evidence for a transient increase in a site-specific rate ratio sometimes to a value $dN/dS > 1$.

In this section I present brief summaries of two studies drawn from my previous work to further illustrate confounding and PL. The objective in the first study was to devise automated methods to identify signatures of sea-surface temperature (SST) fronts in radar images of the ocean surface. Confounding was a problem in that data because atmospheric and oceanographic processes tend to generate image features that look very similar. Confounding was overcome by the inclusion of additional contextual information in the form of concurrent wind vectors. The objective in the second study was to account for variations in the way the efficiency of photosynthesis declines with increasing irradiance via a single parameter p representing a specific mechanistic process. In the course of that study it was realized that confounding between a variety of processes that impact efficiency meant that p could only be interpreted phenomenologically. In the parlance of Chapter 3, I would now say that estimates of p carry phenomenological load. The SST study provides an excellent example of confounding because it can be seen directly in the imagery. The

photosynthesis study supports my view that a phenomenological approach can often be more appropriate than the mechanistic approach advocated by some (in the context of molecular evolution, e.g., Miyazawa, 2011; Liberles *et al.*, 2013; Zaheri *et al.*, 2014; Pollock *et al.*, 2017).

5.2.1 Confounding in the Automated Detection of SST Fronts

The location of the North Wall of the Gulf Stream (NWGS) is tactically important in submarine warfare, and also of interest to certain fisheries. Sea-surface temperature (SST) fronts are readily identified by spaceborne passive radiometry (e.g., the Moderate Resolution Imaging Spectroradiometer or MODIS instrument). However, MODIS is not a reliable source of information because clouds typically obscure its view of the western Atlantic between about 45% and 90% of the time depending on the season (Jones *et al.*, 2013). An alternative data source comes from Synthetic Aperture Radar (SAR), an active sensor that can penetrate cloud cover to measure variations in cm-scale surface waves (i.e., the roughness of the ocean-surface). Surface roughness is modulated by large-scale atmospheric and oceanographic processes to produce regions of high (bright) and low (dark) backscatter. SST fronts often appear as edges between dark and bright regions (a.k.a. brightness fronts) due to several mechanisms that increase surface roughness where the water is warmer compared to adjacent cooler water bodies.

Figure 5.3 A shows two contiguous RADARSAT-2 images acquired on 7 March 2009. The corresponding MODIS SST image acquired the same day is shown in Figure 5.3 B. A number of interesting features are evident. The most obvious is the large meander in the Gulf Stream at the bottom of both images. Other features in the RADARSAT-2 image include (1) evidence of horizontal mixing of cooler shelf water and warmer water just north of the Gulf Stream meander between 39°N and 40°N; (2) an intrusion of cooler water along the eastern side of the Gulf Stream meander; (3) the imprint of atmospheric gravity waves that appears as alternating dark and bright bands just south of Cape Cod; (4) evidence of surfactant accumulation (oils or biomass that smooth the ocean surface) in the lee of the Gulf Stream meander; (5) bright filaments

consistent with convergence zones commonly associated with strong currents; and finally (6) brightness fronts caused by large-scale atmospheric processes marked by WIN (for wind) in the lower portion of the image. All of these features can be correctly identified from the RADARSAT-2 image alone by a human operator after a little training. But they all look more-or-less the same from the point of view of an automated edge detector. Their generating processes are therefore confounded from the perspective of that edge detector.

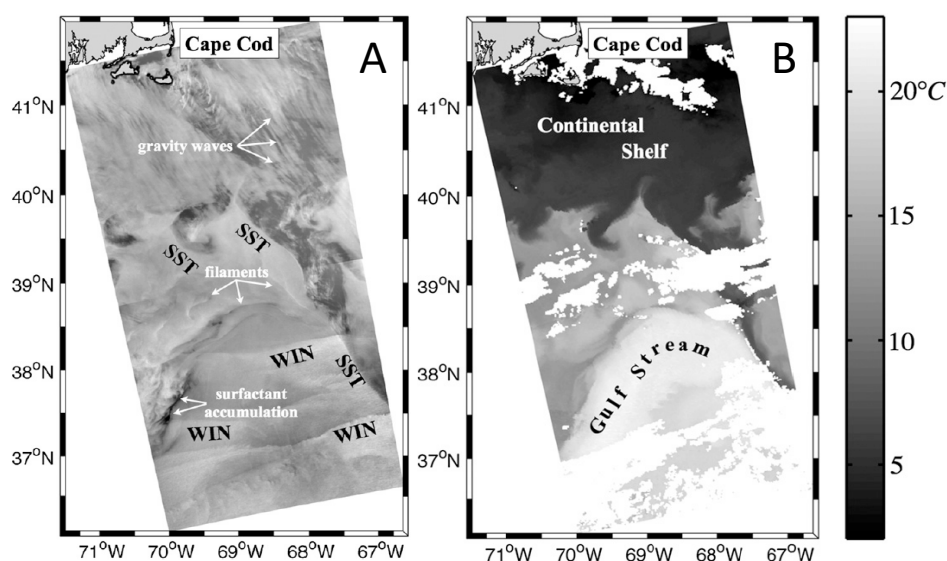


Figure 5.3: Satellite images. (A): Two contiguous RADARSAT-2 SCNA VV frames acquired on 7 March 2009 at approximately 2230 UTC. SST front signatures appear as well-defined brightness fronts. Signatures of horizontal wind shear can be seen in the lower portion of the image, identifiable as brightness fronts misaligned with SST fronts in the MODIS image on the right. (B): Composite MODIS SST image acquired between 1300 and 1900 UTC on the same day. (Images from Jones *et al.* (2013), RADARSAT-2 data and products ©2009 MacDonald, Dettwiler and Associates Ltd. - all rights reserved.)

Standard methods of automated feature classification are based on textual measures extracted from pixels surrounding features of interest. This approach proved to be completely ineffective as a means to discriminate SST fronts from WIN fronts. Instead, external contextual information in the form of wind vectors measured no more than a few hour before or after acquisition of the RADARSAT-2 image was found to be helpful. Specifically, it was found that any given brightness front identified by automated edge-detection software was statistically more likely to be SST when the wind was blowing

across the brightness front and more likely to be WIN when the wind was blowing along the brightness front. A decision rule based on the mean angle between a brightness front and its near-concurrent wind vectors was found to discriminate between the two signatures with an accuracy of 0.80 to 0.85 using cross-validation (Jones *et al.*, 2013).

Contextual information made it possible to discriminate SAR features generated by SST fronts from those generated from purely atmospheric processes. Similarly, I showed in Chapter 4 that the introduction of contextual information in the form of a discrete phenotype makes it possible to discriminate heterotachy caused by changes in site-specific fitness landscapes from heterotachy-by-any-means. The inclusion of other types of information might therefore provide a general solution to the problem of confounding. This idea is not new to the study of molecular evolution. The probability of the $i \rightarrow j$ substitution at a codon site in gene is a function of both time Δt and the substitution rate R_{ij} : $P_{ij}(\Delta t) = \Delta t R_{ij} + o(\Delta t)$. Rates are therefore confounded with chronological time, and furthermore can change over time as lineages diverge. Rate and time can be disentangled by explicitly accounting for variations in evolutionary rates (Thorne *et al.*, 1998) and additionally by applying calibrations using contextual information gleaned from the fossil record (Kishino *et al.*, 2001). The principle of using other sources of information seems to have been limited to such calibrations up to the introduction of the phenotype-genotype models cited in Chapter 4 (i.e., Mayrose and Otto, 2011; Lartillot and Poujol, 2011; O'Connor and Mundy, 2013; Karin *et al.*, 2017; Jones *et al.*, 2019b).

5.2.2 PL and the Efficiency of Photosynthesis

Most of carbon fixation by oxygenic photosynthesis that occurs in the world's oceans is performed by autotrophic phytoplankton. Estimations of the potential of local assemblages of phytoplankton to fix carbon are based on the relationship between the rate of photosynthesis P and the intensity of irradiance E . Measurements of $P(E)$ can be collected by exposing water samples to various levels of irradiance for a time and measuring the amount

of carbon fixed per unit biomass. Such data can be summarized by a simple $P(E)$ relationship.

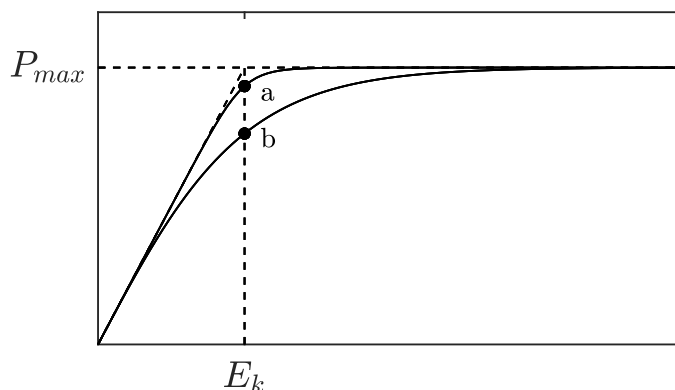


Figure 5.4: Two typical $P(E)$ curves. The rate of transition from a linear relationship at low levels of irradiance to a constant value at high levels of irradiance is a function of the efficiency with which photons are absorbed and used to fix carbon. Curve “a” represents a system that is more efficient than curve “b”.

A number of $P(E)$ models have been proposed (e.g., Blackman, 1905; Baly, 1935; Webb *et al.*, 1974; Jassby and Platt, 1976; Bannister, 1979), most of which contain two adjustable parameters, the saturated rate of photosynthesis P_{\max} and the saturating irradiance E_k . Each two-parameter model assumes an intrinsic level of the efficiency with which photons are converted to fixed carbon, expressed as different values for $P(E_k)$, where lower $P(E_k)$ correspond to lower efficiency (Figure 5.4). Many $P(E)$ data sets exhibit substantial variations in efficiency and are therefore not easily summarized by a single model. Motivated by this, Bannister (1979) introduced an additional parameter b to account for variations in efficiency. The Bannister equation often fits data better than any two-parameter model (e.g., Jones *et al.*, 2014), but is nevertheless seldom used, perhaps because b is not easily assigned a mechanistic interpretation. Two attempts were made to account for variations in efficiency via a mechanistic process, one based on variations in the availability of resources in the electron transport chain (the queue model, Honti, 2007) and another based on properties of the light-harvesting apparatus (the connectivity model, Jones *et al.*, 2014). In the latter case, the parameter $p \in [0, 1]$ was used to represent variations in efficiency via a specific mechanism. Variations in efficiency can

only manifest as variations in $P(E_k)$ however, so mechanisms by which efficiency might vary are confounded under normal laboratory conditions. Such mechanisms include not only those proposed by Honti (2007) and Jones *et al.* (2014), but also processes that can occur at the level of the cell (e.g., variations in the level of intracellular self-shading of closely packed chloroplasts) and at the level of the culture (variations in intercellular self-shading in dense cultures). Using the terminology introduced in Chapter 3, an efficiency parameter meant to account for a specific mechanistic cause is likely to carry substantial PL. Indeed, Jones *et al.* (2014) concluded that the parameter p can only be interpreted as accounting for variations in efficiency by any cause.

5.3 Final Thoughts: What is the Canonical Signature of Molecular Adaptation?

In Chapter 2 it was argued that the traditional signature of positive selection in the form of a rate ratio $\omega > 1$ is in general insufficient to infer adaptation at the level of an individual codon site. It can be sufficient in cases where an elevation in the nonsynonymous substitution rate is sustained due to intergenetic interactions (e.g., an immune surveillance/evasion conflict Hughes and Nei, 1988). It is not sufficient in general because episodic elevations to $\omega > 1$ can be generated by both adaptive and non-adaptive processes. So what should we consider to be the canonical signature of molecular adaptation? The material in Chapter 4 suggests that it might be defined as evidence for a change in a site-specific fitness landscape as detected by the PG-BSM. This was based on the results of simulation studies however, which necessarily used alignment generating processes that are simple in comparison to the actual processes of nature. Unrecognized sources of confounding in real data might therefore undermine the PG-BSM. Indeed, as was noted at the end of Chapter 4, an evolving protein is subject to both intragenetic and intergenetic epistatic constraints (Pollock *et al.*, 2012; Starr and Thornton, 2016; Phillips, 2008). These make it possible for a site to manifest a change in its landscape that only serves to maintain the current function of the protein.

Such function-preserving change is arguably non-adaptive. Changes in site-specific landscapes due to epistasis might therefore be confounded with those due to genuine adaptation where the protein undergoes some form of change in function.

Thermodynamic stability models have been used to investigate evolutionary processes associated with intragenetic epistasis. Such models mimic evolution over a sequence-to-sequence landscape on which the fitness of each possible sequence is fixed. It has been shown that the modeled process can include episodes of contingency-and-entrenchment at a site (Pollock *et al.*, 2012, 2017), in which (i) the stationary frequency of the resident amino acid at a site decreases over time until it becomes replaced by drift (contingency), and (ii) the stationary frequency of the newly fixed amino acid increases over time as other sites undergo substitutions that adjust for its fixation (entrenchment). Recall that, whereas the PG-BSM can be used to infer changes in site-specific landscapes, it does so via specific patterns of heterotachy informed by a discrete phenotype. Changes in site-specific landscapes due to processes arising from intragenetic epistasis can also manifest as heterotachy. Contingency, for example, might sometimes lead to a transient burst of amino acid substitutions (Pollock *et al.*, 2017), and these can be followed by a decline in the substitution rate due to entrenchment. But if the resulting changes in site-specific substitution rates occur independently of changes in phenotype, they will contribute to the signature of heterotachy-by-any-cause accounted for by the covarion-like component of the PG-BSM, and will therefore be unlikely to lead to false inference of adaptive peak shifts.

Nevertheless, as was stated more than once in this thesis, I believe that the reliability of inferences made by fitting a CSM to a real alignment can be assessed only insofar as the simulated data matches the real data. Alignment-generating methods based on the MS framework represent a big step forward, and arguably a new phase in CSM development. But it is unclear whether such methods (e.g., MSmmtDNA, MSTGdR) can mimic all of the variation found in a real alignment that might impact inference. Some argue for an increase in the mechanistic content of fitted models with the apparent objective of

extracting increasingly subtle signatures from alignment data (e.g., Liberles *et al.*, 2013; Pollock *et al.*, 2017). By contrast, I argue for an increase in the mechanistic realism of the processes used to generate data with the objective of identifying and accounting for all potential sources of confounding. I leave this challenging task for future research efforts. In the mean time I conclude that, whatever the signature for molecular evolution is taken to be, it should be based on statistical summaries extracted from alignment data combined with whatever contextual information (e.g., phenotype, protein structure) is required to break confounding.

Bibliography

- Adams, D. C. and Collyer, M. L. 2018. Multivariate phylogenetic comparative methods: evaluations, comparisons, and recommendations. *Sys. Biol.*, 67: 14–31.
- Anisimova, M. and Kosiol, C. 2009. Investigating protein-coding sequence evolution with probabilistic codon substitution models. *Mol. Biol. Evol.*, 26: 255–271.
- Anisimova, M., Bielawski, J. P., and Yang, Z. H. 2001. Accuracy and power of the likelihood ratio test in detecting adaptive molecular evolution. *Mol. Biol. Evol.*, 18: 1585–1592.
- Anisimova, M., Bielawski, J. P., and Yang, Z. H. 2002. Accuracy and power of Bayes prediction of amino acid sites under positive selection. *Mol. Biol. Evol.*, 19: 950–958.
- Averof, M., Rokas, A., Wolfe, K. H., and Sharp, P. M. 2000. Evidence for a high frequency of simultaneous double-nucleotide substitutions. *Science*, 287: 1283–1286.
- Baker, J. E., Dunn, K. A., Mingrone, J. M., Wood, B. A., Karpinski, B. A., Sherwood, C. C., Wildman, D. E., Maynard, T. M., and Bielawski, J. P. 2016. Functional divergence of the nuclear receptor nr2c1 as a modulator of pluripotentiality during hominid evolution. *Genetics*, 203: 905–922.
- Baly, E. C. C. 1935. The kinetics of photosynthesis. *Proc. R. Soc. B. Biol. Sci.*, 117: 218–239.
- Bannister, T. T. 1979. Quantitative description of the steady state, nutrient-saturated algal growth, including adaptation. *Limnol. Oceanogr.*, 24: 76–96.
- Bazykin, G. A. 2015. Changing preferences: deformation of single position amino acid fitness landscapes and evolution of proteins. *Biol. Lett.*, 11: 1–7.
- Bielawski, J. P. and Gold, J. R. 2002. Mutation patterns of mitochondrial H- and L-strand dna in closely related cyprinid fishes. *Genetics*, 161: 1589–1597.
- Bielawski, J. P. and Yang, Z. H. 2004. A maximum likelihood method for detecting functional divergence at individual codon sites, with application to gene family evolution. *J. Mol. Evol.*, 59: 121–132.

- Bielawski, J. P., Dunn, K. A., Sabehi, G., and B ej a, O. 2004. Darwinian adaptation of proteorhodopsin to different light intensities in the marine environment. *PNAS*, 101: 14824–14829.
- Blackman, F. F. 1905. Optima and limiting factors. *Ann. Bot.*, 19: 281–295.
- Bustamante, C. D. 2005. Population genetics of molecular evolution, in *Statistical methods in molecular evolution*. Springer, New York.
- Butler, M. A. and King, A. A. 2004. Phylogenetic comparative analysis: a modeling approach for adaptive evolution. *Am. Nat.*, 164: 683–695.
- Cao, Y., Janke, A., Waddell, P. J., Westerman, M., Takenaka, O., Murata, S., Okada, N., Paabo, S., and Hasegawa, M. 1998. Conflict among individual mitochondrial proteins in resolving the phylogeny of eutherian orders. *J. Mol. Evol.*, 47: 307–322.
- Clayton, D. A. 1982. Replication of animal mitochondrial DNA. *Cell*, 28: 693–705.
- Cornwell, W. and Nakagawa, S. 2017. Phylogenetic comparative methods. *Curr. Biol.*, 27: 327–338.
- Crick, F. 1970. Central dogma of molecular biology. *Nature*, 227: 561–563.
- Darwin, C. 1859. *On the origin of species by means of natural selection, of preservation of favored races in the struggle for life*. John Murray, London.
- Davydov, I. I., Salamin, N., and Robinson-Rechavi, M. 2019. Large-scale comparative analysis of codon models accounting for protein and nucleotide selection. *Mol. Biol. Evol.*, 36: 1316–1332.
- De Maio, N., Holmes, I., Schl otterer, C., and Kosiol, C. 2013. Estimating empirical codon hidden Markov models. *Mol. Biol. Evol.*, 30: 725–736.
- dos Reis, M. 2013. <http://arxiv:1311.6682v1>. last accessed November 26 2013.
- dos Reis, M. 2015. How to calculate the non-synonymous to synonymous rate ratio protein-coding genes under the Fisher-Wright mutation-selection framework. *Biology Letters*, 11: 1–4.
- dos Reis, M. and Yang, Z. H. 2013. Why do more divergent sequences produce smaller nonsynonymous/synonymous rate ratios in pairwise sequence comparisons. *Genetics*, 195: 195–204.
- Eastman, J. M., Alfaro, M. E., and *et. al.*, P. J. 2011. A novel comparative method for identifying shifts in the rate of character evolution on trees. *Evolution*, 65: 3578–3589.

- England, J. L. 2013. Statistical physics of self-replication. *J. Chem. Phys.*, 139: 1–8.
- Felsenstein, J. J. 1981. Evolutionary trees from dna sequences: a maximum likelihood approach. *J. Mol. Evol.*, 17: 368–376.
- Felsenstein, J. J. 1985. Phylogenies and the comparative method. *Am. Nat.*, 125: 1–15.
- Field, S. F., Bulina, M. Y., Kelmanson, I. V., Bielawski, J. P., and Matz, M. V. 2006. Adaptive evolution of multicolored fluorescent proteins in reef-building corals. *J. Mol. Evol.*, 62: 332–339.
- Fisher, R. A. 1922. On the dominance ratio. *Proc. R. Soc. Edinb.*, 42: 321–341.
- Fisher, R. A. 1930. The evolution of dominance in certain polymorphic species. *Am. Nat.*, 64: 385–406.
- Fitch, W. 1971. Rate of change of concomitantly variable codons. *J. Molec. Evolution*, 2: 84–96.
- Fitch, W. and Markowitz, E. 1970. An improved method for determining codon variability in a gene and its application to the rate of fixation of mutations in evolution. *Biochem. Genet.*, 4: 579–593.
- Forsberg, R. and Christiansen, F. B. 2003. A codon-based model of host-specific selection in parasites, with an application to the influenza a virus. *Mol. Biol. Evol.*, 20: 1252–1259.
- Galtier, N. 2001. Maximum-likelihood phylogenetic analysis under a covarion-like model. *Mol. Biol. Evol.*, 18: 866–873.
- Garvin, M. R., Bielawski, J. P., Sazanov, L. A., and Gharrett, A. J. 2015. Review and meta-analysis of natural selection in mitochondrial complex I in metazoans. *J. Zoolog. Syst. Evol. Res.*, 53: 1–17.
- Gaston, D., Roger, A. J., and Susko, E. 2011. A phylogenetic mixture model for the identification of functionally divergent protein residues. *Bioinformatics*, 27: 2655–2663.
- Goldman, N. and Yang, Z. H. 1994. Codon-based model of nucleotide substitution for protein-coding dna-sequences. *Mol. Biol. Evol.*, 11: 725–736.
- Grossmann, T. I., Waxman, D., and Eyre-Walker, A. 2014. Fluctuating selection models and mcDonald-kreitman type analyses. *PLOS ONE*, 9.
- Gu, X. 1999. Statistical methods for testing functional divergence after gene duplication. *Mol. Biol. Evol.*, 16: 1664–1674.

- Gu, X. 2001. Maximum-likelihood approach for gene family evolution under functional divergence. *Mol. Biol. Evol.*, 18: 453–464.
- Gu, X. 2006. A simple statistical model for estimating type-ii (cluster-specific) functional divergence of protein sequences. *Mol. Biol. Evol.*, 23: 1937–1945.
- Guindon, S., Rodrigo, A. G., Dyer, K. A., and Huelsenbeck, J. P. 2004. Modeling the site-specific variation of selection patterns along lineages. *PNAS*, 101: 12957–12962.
- Guindon, S., Black, M., and Rodrigo, A. G. 2006. Control of the false discovery rate applied to the detection of positively selected amino acid sites. *Mol. Biol. Evol.*, 23: 919–926.
- Haldane, J. B. S. 1927. The mathematical theory of natural and artificial selection. *Proc. Camb. Philos. Soc.*, 23: 838–844.
- Haldane, J. B. S. 1932. *The causes of evolution*. Harper and Brothers, New York.
- Halpern, A. L. and Bruno, W. J. 1998. Evolutionary distances for protein-coding sequences: modeling site-specific residue frequencies. *Mol. Biol. Evol.*, 15: 910–917.
- Hansen, T. F. 1997. Stabilizing selection and the comparative analysis of adaptation. *Evolution*, 51: 1341–1351.
- Harris, K. and Nielsen, R. 2014. Error-prone polymerase activity causes multinucleotide mutations in humans. *Genome Research*, 9: 1445–1554.
- Honti, M. 2007. Stochastic parallel processing can shape photosynthesis-irradiance curves in phytoplankton - the Q-model. *Hydrobiologia*, 592: 315–328.
- Huelsenbeck, J. P. and Dyer, K. A. 2004. Bayesian estimation of positively selected sites. *J. Mol. Evol.*, 58: 661–672.
- Hughes, A. L. and Nei, M. 1988. Pattern of nucleotide substitution at major histocompatibility complex class-1 loci reveals overdominant selection. *Nature*, 335: 167–170.
- Jablonka, E. and Lamb, M. J. 2006. *Evolution in four dimensions*. MIT Press, Massachusetts.
- Jassby, A. D. and Platt, T. 1976. Mathematical formulation of the relationship between photosynthesis and light for phytoplankton. *Limnol. Oceanogr.*, 21: 540–547.

- Jones, C. T., Sikora, T. D., Vachon, P. W., Wolfe, J., and DeTracey, B. 2013. Automated discrimination of certain brightness fronts in radarsat-2 images of the ocean surface. *JTECH*, 30: 2203–2215.
- Jones, C. T., Barnett, A. B., MacIntyre, H. L., and Cullen, J. J. 2014. Curvature in models of the photosynthesis-irradiance response. *Journal of Phycology*, 50: 341–355.
- Jones, C. T., Youssef, N., Susko, E., and Bielawski, J. P. 2017. Shifting balance on a static mutation-selection landscape: a novel scenario of positive selection. *Mol. Biol. Evol.*, 34: 391–407.
- Jones, C. T., Youssef, N., Susko, E., and Bielawski, J. P. 2018. Phenomenological load on model parameters can lead to false biological conclusions. *Mol. Biol. Evol.*, 35: 1473–1488.
- Jones, C. T., Susko, E., and Bielawski, J. P. 2019a. Looking for Darwin in genomic sequences: validity and success depends on the relationship between model and data, in *Evolutionary Genomics: Statistical and Computational Methods*. Humana Press, 2nd edition.
- Jones, C. T., Youssef, N., Susko, E., and Bielawski, J. P. 2019b. A phenotype-genotype codon model for detecting adaptive evolution. *submitted to Sys. Bio. May 2019*, 0: 1–50.
- Karin, E. L., Wicke, S., Pupko, T., and Mayrose, I. 2017. An integrated model of phenotypic trait changes and site-specific sequence evolution. *Syst. Biol.*, 66: 917–933.
- Keightley, P., Trivedi, U., Thomson, M., Oliver, F., Kumar, S., and Blaxter, M. 2009. Analysis of the genome sequences of three *Drosophila melanogaster* spontaneous mutation accumulation lines. *Genet. Res.*, 19: 1195–1201.
- Kimura, M. 1962. On the probability of fixation of mutant genes in a population. *Genetics*, 47: 713–719.
- Kimura, M. 1968. Evolutionary rate at the molecular level. *Nature*, 217: 624–626.
- Kimura, M. 1983. *The neutral theory of molecular evolution*. Cambridge University Press, Cambridge.
- Kimura, M. and Ohta, T. 1969. The average number of generations until fixation of a mutant gene in a finite population. *Genetics*, 61: 763–771.
- Kishino, H., Thorne, J. L., and Bruno, W. J. 2001. Performance of a divergence time estimation method under a probabilistic model of rate evolution. *Mol. Biol. Evol.*, 18: 352–361.

- Koren, Y. 2005. Drawing graphs by eigenvectors: theory and practice. *Comput. Math. Appl.*, 49: 1867–1888.
- Kosakovsky Pond, S. L. and Frost, S. D. W. 2005. Not so different after all: a comparison of methods for detecting amino acid sites under selection. *Mol. Biol. Evol.*, 22: 1208–1222.
- Kosakovsky Pond, S. L. and Muse, S. V. 2005. Site-to-site variations of synonymous substitution rates. *Mol. Biol. Evol.*, 22: 2375–2385.
- Kosakovsky Pond, S. L., Frost, S. D. W., and Muse, S. V. 2004. Hyphy: hypothesis testing using phylogenies. *Bioinformatics*, 21: 676–679.
- Kosakovsky Pond, S. L., Murrell, B., Fourment, M., Frost, S. D. W., Delport, W., and Scheffler, K. 2011. A random effects branch-site model for detecting episodic diversifying selection. *Mol. Biol. Evol.*, 28: 3033–3043.
- Kosiol, C., Holmes, I., and Goldman, N. 2007. An empirical codon model for protein sequence evolution. *Mol. Biol. Evol.*, 24: 1464–1479.
- Krogh, T. J., Møller-Jensen, J., and Kaleta, C. 2018. Impact of chromosomal architecture on the function and evolution of bacterial genomes. *Frontiers in Microbiology*, 9: 1–15.
- Kryazhimskiy, A. and Plotkin, J. B. 2008. The population genetics of dn/ds. *PLoS Genetics*, 4: e1000304.
- Kumar, S., Filipinski, A. J., Battistuzzi, F. U., Kosakovsky Pond, S. L., and Tamura, K. 2011. Statistics and truth in phylogenomics. *Mol. Biol. Evol.*, 29: 457–472.
- Lartillot, N. and Poujol, R. 2011. A phylogenetic model for investigating correlated evolution of substitution rates and continuous phenotypic characters. *Mol. Biol. Evol.*, 28: 729–744.
- Lenton, T. and Watson, A. 2011. *Revolutions that made the Earth*. Oxford University Press, Oxford.
- Liberles, D. A., Teufel, A. I., Liu, L., and Stadler, T. 2013. On the need for mechanistic models in computational genomics and metagenomics. *Genome Biol. Evol.*, 5: 2008–2018.
- Lopez, P., Casane, D., and Phillipe, H. 2002. Heterotachy, and important process of protein evolution. *Mol. Biol. Evol.*, 19: 1–7.
- Lu, A. and Guindon, S. 2013. Performance of standard and stochastic branch-site models for detecting positive selection among coding sequences. *Mol. Biol. Evol.*, 31: 484–495.

- Maddison, W. P., Midford, P. E., and Otto, S. P. 2007. Estimating a binary character's effect on speciation and extinction. *Syst. Biol.*, 56: 701–710.
- Mathews, S. 2010. Evolutionary studies illuminate the structural-functional model of plant phytochromes. *The Plant Cell*, 22: 4–16.
- Mayrose, I. and Otto, S. P. 2011. A likelihood method for detecting trait-dependent shifts in the rate of molecular evolution. *Mol. Biol. Evol.*, 28: 759–770.
- McCandlish, D. M. 2011. Visualizing fitness landscapes. *Evolution*, doi:10.1111/j.1558-5646.2011.01236.x: 1544–1558.
- McCandlish, D. M. and Stoltzfus, A. 2014. Modeling evolution using the probability of fixation: history and implications. *The Quarterly Review of Biology*, 89: 225–252.
- Mingrone, J., Susko, E., and Bielwaski, J. P. 2016. Smoothed bootstrap aggregation for assessing selection pressure at amino acid sites. *Mol. Biol. Evol.*, 33: 2976–2989.
- Mingrone, J., Susko, E., and Bielwaski, J. P. 2018. Modified likelihood ratio tests for positive selection. *submitted*.
- Miyazawa, S. 2011. Advantages of a mechanistic codon substitution model for evolutionary analysis of protein-coding sequences. *PLoS ONE*, 6: 20pp.
- Mugal, C. F., Wolf, J. B. W., and Kaj, I. 2014. Why time matters: codon evolution and the temporal dynamics of dn/ds . *Mol. Biol. Evol.*, 31: 212–231.
- Murrell, B., Weaver, S., Smith, M. D., Wertheim, J. O., Murrell, S., Aylward, A., Eren, K., Pollner, T., Martin, D. P., Smith, D. M., Scheffler, K., and Pond, S. L. K. 2015. Gene-wide identification of episodic selection. *Mol. Biol. Evol.*, 32: 1365–1371.
- Muse, S. V. and Gaut, B. S. 1994. A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with applications to the chloroplast genome. *Mol. Biol. Evol.*, 11: 715–724.
- Mustonen, V. and Lässig, M. 2009. From fitness landscapes to seascapes: non-equilibrium dynamics of selection and adaptation. *Trends in Genetics*, 25: 111–119.
- Nei, M. and Gojobori, T. 1986. Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Mol. Biol. Evol.*, 3: 418–426.

- Newton, M. A., Noueir, A., Sarkar, D., and Ahlquist, P. 2004. Detecting differential gene expression with a semiparametric hierarchical mixture model. *Biostatistics*, 5: 155–176.
- Nielsen, R. and Yang, Z. 2003. Estimating the distribution of selection coefficients from phylogenetic data with applications to mitochondrial and viral dna. *Mol. Biol. Evol.*, 20: 1231–1239.
- Nielsen, R. and Yang, Z. H. 1998. Likelihood models for detecting positively selected amino acid sites and applications to the HIV-1 envelope gene. *Genetics*, 148: 929–936.
- O'Connor, T. D. and Mundy, N. I. 2013. Evolutionary modeling of genotype-phenotype association and application to the primate coding and non-coding mtDNA rate variation. *Evolutionary Bioinformatics*, 9: 301–316.
- O'Meara, B. C., Ané, C., Sanderson, M. J., and Wainwright, P. C. 2006. Testing for different rates of continuous trait evolution using likelihood. *Evolution*, 60: 922–933.
- Pagel, M. 1994. Detecting correlated evolution on phylogenies: a general method for the comparative analysis of discrete characters. *Proc. R. Soc. Lond. B*, 255: 37–45.
- Patthy, L. 2008. *Protein Evolution 2nd Ed.* Blackwell Publishing, Victoria, Australia.
- Pegueroles, C., Laurie, S., and Albà, M. M. 2013. Accelerated evolution after gene duplication: a time-dependent process affecting just one copy. *Mol. Biol. Evol.*, 30: 1830–1842.
- Philippe, H., Casane, D., Gribaldo, S., Lopez, P., and Meunier, J. 2003. Heterotachy and functional shift in protein evolution. *IUBMB Life*, 55: 257–265.
- Phillips, P. C. 2008. Epistasis - the essential role of gene interactions in the structure and evolution of genetic systems. *Nature Reviews: Genetics*, 9: 855–867.
- Pollock, D. D., Thiltgen, G., and Goldstein, R. A. 2012. Amino acid coevolution induces an evolutionary Stokes shift. *PNAS*, 109: E1352–E1359.
- Pollock, D. D., Pollard, S. T., Shortt, J. A., and Goldstein, R. A. 2017. Mechanistic models of protein evolution, in *Evolutionary biology: Self/nonsel evolution, species and complex traits evolution, methods and concepts*. Springer.
- Pupko, T. and Galtier, N. 2002. A covarion-based method for detecting molecular adaptation: application to the evolution of primate mitochondrial genomes. *Proc. R. Soc. Lond.*, 269: 1313–1316.

- Raina, S. Z., Faith, J. J., Disotell, T. R., Seligmann, H., Stewart, C. B., and Pollock, D. D. 2005. Evolution of base-substitution gradients in primate mitochondrial genomes. *Genomes*, 15: 665–673.
- Reyes, A., Gissi, C., and Saccone, C. 1998. Asymmetric directional mutation pressure in the mitochondrial genome of mammals. *Mol. Biol. Evol.*, 15: 957–966.
- Rodrigue, N. and Lartillot, N. 2014. Site-heterogeneous mutation-selection models with the PhyloBayes-MPI package. *Bioinformatics*, 30: 1020–1021.
- Rodrigue, N. and Lartillot, N. 2016. Detection of adaptation in protein-coding genes using a bayesian site-heterogeneous mutation-selection codon substitution model. *Mol. Biol. Evol.*, 34: 204–214.
- Rodrigue, N. and Philippe, H. 2010. Mechanistic revisions of phenomenological modeling strategies in molecular evolution. *Trends in Genetics*, 26: 248–252.
- Rodrigue, N., Philippe, H., and Lartillot, N. 2010. Mutation-selection models of coding sequence evolution with site-heterogeneous amino acid fitness profiles. *PNAS*, 107: 4629–4634.
- Romero, P. E., Weigand, A. M., and Pfenninger, M. 2016. Positive selection on panpulmonate mitogenomes provide new clues on adaptations to terrestrial life. *BMC Evolutionary Biology*, 16: 1–13.
- Ross, S. M. 1996. *Stochastic Processes*. Wiley, New York.
- Rubinstein, N. D. and Pupko, T. 2012. Detection and analysis of conservation at synonymous sites, in *Codon evolution: mechanisms and models*. Oxford University Press Inc.
- Rubinstein, N. D., Doron-Faigenboim, A., Mayrose, I., and Pupko, T. 2011. Evolutionary model accounting for layers of selection in protein-coding genes and their impact on the inference of positive selection. *Mol. Biol. Evol.*, 28: 3297–3308.
- Schneider, E. D. and Kay, J. J. 1994. Life as a manifestation of the second law of thermodynamics. *Mathl. Comput. Modelling*, 19: 25–48.
- Schrider, D., Hourmozdi, J., and Hahn, M. 2014. Pervasive multinucleotide mutational events in eukaryotes. *Curr. Biol.*, 21: 1051–1054.
- Ségurel, L., Thompson, E. E., Flutre, T., Lovstad, J., Venkat, A., Margulis, S. W., Moyse, J., Ross, S., Gamble, K., Sella, G., Ober, C., and Przeworski, M. 2012. The abo blood group is a trans-species polymorphism in primates. *PNAS*, 109: 18493–18498.

- Self, S. G. and Liang, K. Y. 1987. Asymptotic properties of maximum likelihood estimators and likelihood ratio test under nonstandard conditions. *JASA*, 82: 605–610.
- Sella, G. and Hirsh, A. E. 2005. The application of statistical physics to evolutionary biology. *PNAS*, 102: 9541–9546.
- Smith, M. D., Wertheim, J. O., Weaver, S., Murrell, B., Scheffler, K., and Pond, S. L. K. 2015. Less is more: an adaptive branch-site random effects model for efficient detection of episodic diversifying selection. *Mol. Biol. Evol.*, 32: 1342–1353.
- Smith, N. G. C. and Eyre-Walker, A. 2002. Adaptive protein evolution in drosophila. *Nature*, 415: 1022–1024.
- Spielman, S. and Wilke, C. O. 2015a. Pyvolve: A flexible Python module for simulating sequences along phylogenies. *PLoS ONE*, 10: 1–7.
- Spielman, S. and Wilke, C. O. 2015b. The relationship between dN/dS and scaled selection coefficients. *Mol. Biol. Evol.*, 34: 1097–1108.
- Spielman, S. and Wilke, C. O. 2016. Extensively parameterized mutation-selection models reliably capture site-specific selective constraints. *Mol. Biol. Evol.*, 33: 2990–3001.
- Starr, T. N. and Thornton, J. W. 2016. Epistasis in protein evolution. *Protein Science*, 25: 1204–1218.
- Struder, R. A. and Robinson-Rechavi, M. 2009. Evidence for an episodic model of protein sequence evolution. *Biochem. Soc. Trans.*, 37: 783–786.
- Suzuki, Y. and Nei, M. 2004. False-positive selection identified by ML-based methods: examples from the Sig1 gene of the diatom *Thalassiosira weissflogii* and the tax gene of the human T-cell lymphotropic virus. *Mol. Biol. Evol.*, 21: 914–921.
- Swanson, W. J., Nielsen, R., and Yang, Q. F. 2003. Pervasive adaptive evolution in mammalian fertilization proteins. *Mol. Biol. Evol.*, 20: 18–20.
- Tamuri, A. U., dos Reis, M., and Goldstein, R. A. 2012. Estimating the distribution of selection coefficients from phylogenetic data using sitewise mutation-selection models. *Genetics*, 190: 1101–1115.
- Tamuri, A. U., Goldman, N., and dos Reis, M. 2014. A penalized-likelihood method to estimate the distribution of selection coefficients from phylogenetic data. *Genetics*, 197: 257–271.
- Tanaka, M. and Ozawa, T. 1994. Strand asymmetry in human mitochondrial mutations. *Genomics*, 22: 327–335.

- Thorne, J. L., Kishino, H., and Painter, I. S. 1998. Estimating the rate of evolution of the rate of molecular evolution. *Mol. Biol. Evol.*, 15: 1647–1657.
- Venkat, A., Hahn, M. W., and Thornton, J. W. 2018. Multinucleotide mutations cause false inferences of lineage-specific positive selection. *Nature Ecology & Evolution*, 2: 1280–1288.
- Wang, H., Spencer, M., Susko, E., and Rodger, A. J. 2007. Testing for covarion-like evolution in protein sequences. *Mol. Biol. Evol.*, 24: 294–305.
- Wang, H., Li, K., Susko, E., and Rodger, A. J. 2014. An amino acid substitution-selection model adjusts residue fitness to improve phylogenetic estimation. *Mol. Biol. Evol.*, 31: 779–792.
- Watson, J. D. and Crick, F. H. C. 1953. Molecular structure of nucleic acids: a structure for deoxyribose nucleic acid. *Nature*, 171: 737–738.
- Webb, W. L., Newton, M., and Starr, D. 1974. Carbon dioxide exchange of *alnus rubra*. a mathematical model. *Oecologia(Berl)*, 17: 281–291.
- Wertheim, J. O., Murrell, B., Smith, M. D., Pond, S. L. K., and Scheffler, K. 2014. Relax: Detecting relaxed selection in a phylogenetic framework. *Mol. Biol. Evol.*, 32: 820–832.
- Whelan, S. and Goldman, N. 2004. Estimating the frequency of events that cause multiple-nucleotide changes. *Genetics*, 167: 2027–2043.
- Whelan, S., Blackburne, B. P., and Spencer, M. 2011. Phylogenetic substitution models for detecting heterotachy during plastid evolution. *Mol. Biol. Evol.*, 28: 449–458.
- Wong, W. S. W., Yang, Z. H., Goldman, N., and Nielsen, R. 2004. Accuracy and power of statistical methods for detecting adaptive evolution in protein coding sequences and for identifying positively selected sites. *Genetics*, 168: 1041–1051.
- Wright, S. 1931. Evolution in Mendelian populations. *Genetics*, 16: 97–159.
- Wright, S. 1932. The roles of mutation, inbreeding, crossbreeding, and selection in evolution. *Proceeding of the Sixth International Congress on Genetics*, 1: 355–366.
- Wright, S. 1982. The shifting balance theory and macroevolution. *Annual Review of Genetics*, 16: 1–19.
- Wu, J. and Susko, E. 2009. General heterotachy and distance method adjustments. *Mol. Biol. Evol.*, 26: 2689–2697.

- Xue, Y., Wang, Q., Long, Q., Ng, B. L., Swerdlow, H., Burton, J., Skuce, C., Taylor, R., Abdellah, Z., Zhao, Y., Asan, MacArthur, D. G., Quail, M. A., Carter, N. P., Yang, H., and Tyler-Smith, C. 2009. Human y chromosome base-substitution mutation rate measured by direct sequencing in a deep-rooted pedigree. *Curr Biol.*, 19: 1453–1457.
- Yang, Z. F., Wang, Y. F., Zhou, Y., Gao, Q. S., Zhang, E. Y., Zhu, L., Hu, Y. Y., and Xu, C. W. 2013. Evolution of land plant genes encoding L-Ala-D/L-Glu epimerases (AEEs) via horizontal gene transfer and positive selection. *BCM Plant Biology*, 13: 34–43.
- Yang, Z. H. 2005. The power of phylogenetic comparison in revealing protein function. *PNAS*, 102: 3179–3180.
- Yang, Z. H. 2006. *Computational Molecular Evolution*. Oxford University Press, Oxford.
- Yang, Z. H. 2007. PAML4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.*, 24: 1586–1591.
- Yang, Z. H. 2014. *Molecular Evolution: A Statistical Approach*. Oxford University Press, Oxford.
- Yang, Z. H. 2017. *PAML: Phylogenetic Analysis by Maximum Likelihood*. <http://abacus.gene.ucl.ac.uk/software/pamlDOC.pdf>.
- Yang, Z. H. and Bielawski, J. P. 2000. Statistical methods for detecting molecular adaptation. *Trends in Ecology & Evolution*, 15: 496–503.
- Yang, Z. H. and dos Reis, M. 2011. Statistical properties of the branch-site test of positive selection. *Mol. Biol. Evol.*, 28: 1217–1228.
- Yang, Z. H. and Nielsen, R. 1998. Synonymous and nonsynonymous rate variation in nuclear genes of mammals. *J. Mol. Evol.*, 46: 409–418.
- Yang, Z. H. and Nielsen, R. 2002. Codon-substitution models for detecting molecular adaptation at individual sites along specific lineages. *Mol. Biol. Evol.*, 19: 908–917.
- Yang, Z. H. and Nielsen, R. 2008. Mutation-selection models of codon substitution and their use to estimate selective strengths on codon usage. *Mol. Biol. Evol.*, 25: 568–579.
- Yang, Z. H., Nielsen, R., Goldman, N., and Pedersen, A. M. K. 2000a. Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics*, 155: 431–449.

- Yang, Z. H., Swanson, W. J., and Vacquier, V. D. 2000b. Maximum-likelihood analysis of molecular adaptation in abalone sperm lysin reveals selective pressures along lineages and sites. *Mol. Biol. Evol.*, 17: 1446–1455.
- Yang, Z. H., Wong, S. W. S., and Nielsen, R. 2005. Bayes empirical bayes inference of amino acid sites under positive selection. *Mol. Biol. Evol.*, 22: 1107–1118.
- Yokoyama, S., Tada, T., Zhang, H., and Britt, L. 2008. Elucidation of phenotypic adaptations: Molecular analysis of dim-light vision proteins in vertebrates. *PNAS*, 105: 13480–13485.
- Zaheri, M., Dib, L., and Salamin, N. 2014. A generalized mechanistic codon model. *Mol. Biol. Evol.*, 31: 2528–2541.
- Zhai, W., Nielsen, R., Goldman, N., and Yang, Z. H. 2012. Looking for Darwin in genomic sequences - validity and success of statistical methods. *Mol. Biol. Evol.*, 20: 2889–2893.
- Zhang, J. 2004. Frequent false detection of positive selection by the likelihood method with branch-site models. *Mol. Biol. Evol.*, 21: 1332–1339.
- Zhang, J., Nielsen, R., and Yang, Z. H. 2005. Evaluation of an improved branch-site likelihood method for detecting positive selection at the molecular level. *Mol. Biol. Evol.*, 22: 2472–2479.