The neutral and nearly neutral theories of molecular evolution

Joseph P. Bielawski

Department of Biology Department of Mathematic and Statistics Institute of Comparative genomics Dalhousie University, Halifax, Nova Scotia, Canada Models of molecular evolution seek to explain the *origin* and *maintenance* of genetic variation



~ 5 million nucleotide differences



~ 35 million nucleotide differences

brief review of the fundamental evolutionary "forces"...



1. Neo-Darwinism



2. Neutral Theory



3. Nearly Neutral



natural selection dominates

genetic drift dominates

Key population concepts:

1. mutation

- 2. fixation
- 3. substitution



Key concepts:

1. mutation

2. fixation

3. substitution



Key concepts:

- 1. mutation
- 2. fixation
- 3. substitution



Mutations can be **fixed** or **lost**



1. Neo-Darwinism



2. Neutral Theory



3. Nearly Neutral



natural selection dominates

genetic drift dominates

3. Nearly Neutral



2. Neutral Theory

natural selection dominates

1. Neo-Darwinism

genetic drift dominates

1. Neo-Darwinism



2. Neutral Theory



3. Nearly Neutral



natural selection dominates

genetic drift dominates

Neutral theory of molecular evolution (Kimura 1968)

k = rate of nucleotide substitution [at mutation-drift equilibrium]
k = new mutations × probability of fixation







N = large



Why is the equilibrium substitution rate $k = \mu$?

neutral theory:
$$k=\mu$$

In words ...

Small populations: lower number of new mutants each generation, but each has a higher probability of fixation

Large populations: higher number of new mutants each generation but probability of fixation is lower

Neutral Theory: precise expectations when mutation & drift are at equilibrium



- 1. evolutionary rate is independent of population size
- 2. constant neutral rate : "molecular clock"
- 3. evolutionary rate is inverse of functional constraint



Distribution of fitness effects (DFE) of MUTATIONS according to Kimura's Neutral Theory



selection coefficients (s) for new mutations





selection coefficients (s) for new mutations

Neutral Theory: precise expectations when mutation & drift are at equilibrium



- 1. evolutionary rate is independent of population size
- 2. constant neutral rate : "molecular clock"
- 3. evolutionary rate is inverse of functional constraint

Distribution of fitness effects (DFE) according to Kimura's Neutral Theory



selection coefficients (s) for new mutations

Distribution of fitness effects (DFE) according to Kimura's Neutral Theory

lethal & strongly deleterious mutations:

- rapidly removed by natural selection
- **never observed** in natural populations
- Kimura ignores them

The **ratio** of these determines the rate of evolution

neutral mutations = vast majority of:

- polymorphism
- species divergence



neutral theory: $k=\mu$

Lower-rate gene:

- more sites functional
- more mutations removed by selection



Distribution of fitness effects (DFE) according to Kimura's Neutral Theory

lethal & strongly deleterious mutations:

- rapidly removed by natural selection
- **never observed** in natural populations
- Kimura ignores them

The **ratio** of these determines the rate of evolution

neutral mutations = vast majority of:

- polymorphism
- species divergence



neutral theory: $k=\mu$

High-rate gene:

- more neutral sites
- more mutations fixed by drift



neutral theory predicts:

The evolutionary rate is inverse of functional constraint.

1. Neo-Darwinism

2. Neutral Theory



3. Nearly Neutral



natural selection dominates

genetic drift dominates



evolution is modelled as an "all or nothing affair"

DFE for Ohta's Nearly-Neutral Theory



evolution is modelled as an **interaction** between genetic drift and natural selection

selection coefficients (s) for new mutations

Distribution of fitness effects (DFE): broad importance to evolutionary biology

Evolutionary Importance:

- mutations are ultimate source of variation
- rate of evolution
- species adaptation
- mutational load (genome decay & reduced survival)

Inference DFEs:

- longstanding goal
- hard to estimate
- much variation in observed DFEs

Distribution of fitness effects: non-synonymous mutations in viral PB2 gene

example: Tamuri et al. (2012) Genetics. 190:1101-1115.



Distribution of fitness effects: generalized compilation of inferences



3. Nearly Neutral





3. Nearly Neutral



drift and selection interact

OHTA EXTENDS THE SELECTION AND NEUTRAL MODELS

Ohta adds evolution of "nearly neutral" mutations:

- small populations:
 - larger "neutral space"
 - more mutations are "effectively neutral"
 - more mutations evolve by genetic drift
- large populations:
 - smaller "neutral space"
 - deleterious: mostly eliminated by natural selection
 - beneficial: fixed more frequently than by drift (but fixation is not certain)

now... molecular clock unlikely (the rate of evolution is affected by changes in N)

Selective implications of near neutrality

- 1. rate slows as population becomes adapted
- 2. population approaches an equilibrium
- 3. population reaches a state of "detailed balance"



phenotypic trait

concave (saturating) fitness curve



Implications of nearly neutral theory



Implications of nearly neutral theory



An now, an actual molecular evolutionary process...

example: epistasis and protein stability



Implications of nearly neutral theory



This is why some forms of Nearly-Neutral models are sometimes called "**steady state models**" or "**balance mutation models**"

- Hartl DL, Taubes CH. Compensatory nearly neutral mutations: selection without adaptation. J Theor Biol. **1996**. 182(3):303-309.
- Sella G, Hirsh AE. The application of statistical physics to evolutionary biology. Proceedings of the National Academy of Sciences. **2005**. 102(27):9541-9546.
- Razeto-Barry P, Díaz J, Vásquez RA. The nearly neutral and selection theories of molecular evolution under the fisher geometrical framework: substitution rate, population size, and complexity. Genetics. **2012**. 191(2):523-534.
- Jones CT, Youssef N, Susko E, Bielawski JP. Shifting balance on a static mutation-selection landscape: a novel scenario of positive selection. Molecular biology and evolution. 2016. 34(2):391-407.

The equilibrium phenotype is NOT the most fit type. (**adaptation** ≠ optimal "engineering" state)



Further reading on stability-mediated epistasis and protein evolution...

• Sella G, Hirsh AE. The application of statistical physics to evolutionary biology. Proceedings of the National Academy of Sciences. **2005**. 102(27):9541-9546.

• Goldstein RA. The evolution and evolutionary consequences of marginal thermostability in proteins. Proteins: Structure, Function, and Bioinformatics. **2011**. 79(5):1396-1407.

• Pollock DD, Thiltgen G, Goldstein RA. Amino acid coevolution induces an evolutionary Stokes shift. Proceedings of the National Academy of Sciences. **2012.** 109(21):E1352-9.

• Youssef N, Susko E, Roger AJ, Bielawski JP. Evolution of amino acid propensities under stability-mediated epistasis. Molecular Biology and Evolution. **2022**. 39(3):msac030.

Summary



1. Neo-Darwinism:

- almost everything is adaptive (too strong)
- evolution "seeks optimality" (promotes an engineering perspective)
- remains THE evolutionary mechanism for origin of adaptations
- 2. Neutral theory:
 - elegant simplicity
 - assumes simplistic DFE
 - predictions are correct for effectively neutral mutations
 - many predictions serve as "principles of evolution"
- 3. Nearly-neutral theory:
 - predicts more complex evolutionary dynamics
 - depends on populations size (unlike neutral theory)
 - some predictions closer to natural populations that neutral theory
 - predicts equilibrium where phenotype is not necessarily optimal
 - natural selection acts to balance "mutational load" on fitness (maintenance)





An index of the intensity of natural selection for proteins

Neutral and Nearly-neutral theory

Evolutionary rate depends on intensity of selection

Neutral theory: independent of N



Nearly-neutral theory: **depends on N**



Let's apply these ideas to individual sites...

1.) selectively constrained:

- neutral space < 100%
- rate < strictly neutral

2.) strictly neutral:

- 100% neutral space
- rate = neutral rate

3.) adaptive evolution:

- large adaptive space
- rate > neutral rate (?)



Prote/ins have a "built in ruler" for their own neutral rate of molecular evolution!



Genetic code

.

Second letter							
		U	С	А	G		
	U	UUU UUC UUA UUG Leu	UCU UCC UCA UCG	UAU UAC UAA Stop UAG Stop	UGU UGC UGA UGG Trp	U C A G	
	с	CUU CUC CUA CUG	CCU CCC CCA CCG	$ \begin{array}{c} CAU \\ CAC \end{array} \\ \begin{array}{c} His \\ CAA \\ CAG \end{array} \\ \begin{array}{c} GIn \end{array} \end{array} $	CGU CGC CGA CGG	UCAG	letter
	A	AUU AUC AUA AUG Met	ACU ACC ACA ACG	AAU AAC AAA AAG Lys	AGU AGC AGA AGG AGG	UCAG	Third
	G	GUU GUC GUA GUG	GCU GCC GCA GCG	GAU GAC GAA GAG GIu	GGU GGC GGA GGG	UCAG	

all possible mutations → **two types**

synonymous (S): no change to protein

- no effect on protein
- selectively neutral
- rate = neutral rate (w.r.t. protein evolution)

non-synonymous (**N**): changes the amino acid composition of protein

- changes AA of protein
- **deleterious**, or neutral, or **positive**
- rate depends on intensity of selection

First letter

"built in ruler" = synonymous substitution rate

The rate at which that proteins would have evolved if it had been 100% free from selection (*at the protein level*).



Genetic code

Cocond latter

Second letter							
	U	С	А	G			
U	UUU }Phe UUC }Phe UUA }Leu UUG }Leu	UCU UCC UCA UCG	UAU UAC UAA Stop UAG Stop	UGU UGC UGA UGG Trp	UCAG		
с	CUU CUC CUA CUG	CCU CCC CCA CCG	$\begin{array}{c} CAU \\ CAC \end{array} \\ \begin{array}{c} His \\ CAA \\ CAG \end{array} \\ \begin{array}{c} GIn \end{array} \end{array}$	CGU CGC CGA CGG	UCAG	letter	
A	AUU AUC AUA AUG Met	ACU ACC ACA ACG	AAU AAC AAA AAG Lys	AGU AGC AGA AGG AGG	UCAG	Third	
G	GUU GUC GUA GUG	GCU GCC GCA GCG	GAU GAC GAA GAG GIu	GGU GGC GGA GGG	UCAG		

all possible mutations → **two types**

synonymous (**S**): no change to protein

- no effect on protein
- selectively neutral
- rate = neutral rate
- neutral rate = $K_{\rm S}$

non-synonymous (**N**): changes the amino acid composition of protein

- changes AA of protein
- deleterious, or neutral, or positive
- rate depends on intensity of selection
- purifying selection: $K_N < K_S$

First letter

Kimura (1983) : the **rate ratio** is an index of selection intensity



Note: not calling this "positive" selection (yet)

1.) selectively constrained:	2.) strictly neutral:	3.) adaptive evolution:
$d_{\rm N} / d_{\rm S} < 1$	$d_{\rm N} / d_{\rm S} = 1$	$d_{\rm N} / d_{\rm S} > 1$



a useful perspective...

Consequences of Stability-Induced Epistasis for Substitution Rates

Noor Youssef,*^{1,2} Edward Susko,^{2,3} and Joseph P. Bielawski^{1,2,3}

¹Department of Biology, Dalhousie University, Halifax, Nova Scotia, Canada ²Centre for Genomics and Evolutionary Bioinformatics, Dalhousie University, Halifax, Nova Scotia, Canada ³Department of Mathematics and Statistics, Dalhousie University, Halifax, Nova Scotia, Canada ***Corresponding author**: E-mail: n.voussef@dal.ca.

Associate editor: Jeffrey Thorne

Abstract

Do interactions between residues in a protein (i.e., epistasis) significantly alter evolutionary dynamics? If so, what consequences might they have on inference from traditional codon substitution models which assume site-independence for the sake of computational tractability? To investigate the effects of epistasis on substitution rates, we employed a mechanistic mutation-selection model in conjunction with a fitness framework derived from protein stability. We refer to this as the stability-informed site-dependent (S-SD) model and developed a new stability-informed site-independent (S-SI) model that captures the average effect of stability constraints on individual sites of a protein. Comparison of S-SI and S-SD offers a novel and direct method for investigating the consequences of stability-informed epistasis on protein evolution. We developed S-SI and S-SD models for three natural proteins and showed that they generate sequences consistent with real alignments. Our analyses revealed that epistasis tends to increase substitution rates compared with the rates under site-independent evolution. We then assessed the epistatic sensitivity of individual site and discovered a counterintuitive effect: Highly connected sites were less influenced by epistasis relative to exposed sites. Lastly, we show that, despite the unrealistic assumptions, traditional models perform comparably well in the presence and absence of epistasi and provide reasonable summaries of average selection intensities. We conclude that epistatic models are critical to understanding protein evolutionary dynamics, but epistasis might not be required for reasonable inference of selection pressure when averaging over time and sites.

Key words: epistasis, dN/dS, protein stability, substitution rates, mutation-selection model, protein evolution.

Introduction

Most proteins must fold into a native structure in which they are moderately stable before they are able to perform their biological function. Protein stability depends on the sequence of amino acids and their interactions in the folded threedimensional structures. Because of these interactions, evolutionary selective constraints to maintain adequate stability result in epistatic dependencies between residues. Specifically, epistasis manifests as a dependency in the fitness effect of a mutation on the background protein sequence in which it arose. For example, let $f_a^h(S)$ be the fitness of the protein provided amino acid a is occupying site h in the context of background sequence S. Then, $F^h(S) = \langle f_1^h(S), \ldots, \rangle$ $f_{20}^{h}(S)\rangle$ is the site-specific vector of amino acid fitness values specifying the fitness landscape at site h. Following a substitution at another position in the protein, so that the background sequence changes from S to X, the fitness of the same amino acid will subsequently change, $f_a^h(S) \neq f_a^h(X)$. Therefore, in the presence of epistatic dependencies, the fitness landscape at a site is subject to fluctuations as substitutions occur at other sites (fig. 1A). Stability constraints typically result in global epistasis, meaning that a change in the incumbent amino acid at one site induces shifts in the fitness landscapes at many, often all, other sites in the protein (Starr and Thomton 2016). Although such interdependencies inevitably occur, the magnitude and frequency of these shifts, and their impact on protein evolution, remain controversial.

Anticle

Using extensive computational experiments, Pollock et al. (2012) found that stability-induced epistasis results in frequent and substantial shifts in amino acid fitness landscapes. To the contrary, Ashenberg et al. (2013) used computational and experimental approaches and reported that although stability-induced fluctuations in site-specific amino acid fitness landscapes do occur, they are relatively minor in magnitude and are therefore inconsequential with regards to longterm evolutionary dynamics. This controversy has spurred multiple follow-up experiments, finding support for both claims and little consensus (Risso et al. 2015; Shah et al. 2015; Starr et al. 2018; Ferrada 2019). It remains unclear if and how stability-induced epistasis influences protein evolution.

Models used to infer evolutionary parameters from natural protein alignments commonly assume site-independence and other simplifying assumptions (e.g., time-stationary substitution rates, and low levels of among-site rate

© The Author(s) 2020. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution. All rights reserved. For permissions, please e-mail: journals.permissions@oup.com





"...despite the unrealistic assumptions, traditional models perform comparably well in the presence and absence of epistasis and provide reasonable summaries of average selection intensities."

> This does **NOT** mean that you can just plow ahead and *ignore the assumptions* of your models and your tests!!!