

codon substitution models and the analysis of natural selection pressure

Joseph P. Bielawski

Department of Biology

Department of Mathematic and Statistics

Institute of Comparative genomics

Dalhousie University, Halifax, Nova Scotia, Canada

Details matter!

Codon models: waaaaay **too much** to cover in this talk

Macro-evolutionary inference of selection intensity:

- Very complex and diverse modelling strategies
- Deep statistical issues
- Model testing and interpretation
- Strong opinions about “the right thing to do”





Chapter 13

Looking for Darwin in Genomic Sequences: Validity and Success Depends on the Relationship Between Model and Data

Christopher T. Jones, Edward Susko, and Joseph P. Bielawski

Abstract

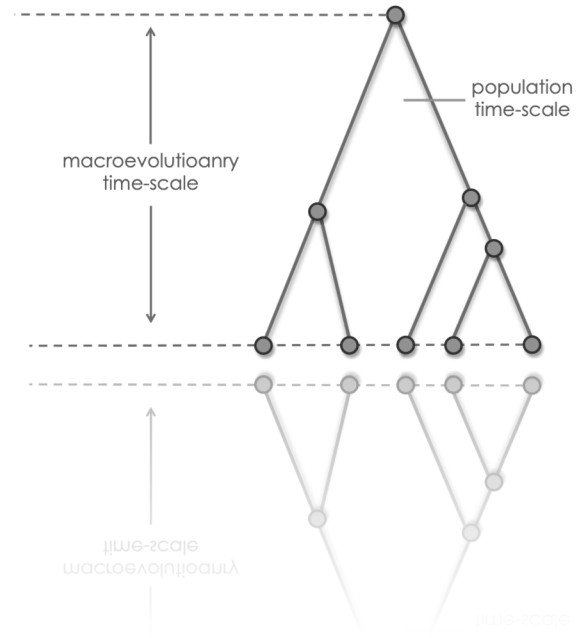
Codon substitution models (CSMs) are commonly used to infer the history of natural selection for a set of protein-coding sequences, often with the explicit goal of detecting the signature of positive Darwinian selection. However, the validity and success of CSMs used in conjunction with the maximum likelihood (ML) framework is sometimes challenged with claims that the approach might too often support false conclusions. In this chapter, we use a case study approach to identify four legitimate statistical difficulties associated with inference of evolutionary events using CSMs. These include: (1) model misspecification, (2) low information content, (3) the confounding of processes, and (4) phenomenological load, or PL. While past criticisms of CSMs can be connected to these issues, the historical critiques were often misdirected, or overstated, because they failed to recognize that the success of any model-based approach depends on the relationship between model and data. Here, we explore this relationship and provide a candid assessment of the limitations of CSMs to extract historical information from extant sequences. To aid in this assessment, we provide a brief overview of: (1) a more realistic way of thinking about the process of codon evolution framed in terms of population genetic parameters, and (2) a novel presentation of the ML statistical framework. We then divide the development of CSMs into two broad phases of scientific activity and show that the latter phase is characterized by increases in model complexity that can sometimes negatively impact inference of evolutionary mechanisms. Such problems are not yet widely appreciated by the users of CSMs. These problems can be avoided by using a model that is appropriate for the data; but, understanding the relationship between the data and a fitted model is a difficult task. We argue that the only way to properly understand that relationship is to perform *in silico* experiments using a generating process that can mimic the data as closely as possible. The mutation-selection modeling framework (MutSel) is presented as the basis of such a generating process. We contend that if complex CSMs continue to be developed for testing explicit mechanistic hypotheses, then additional analyses such as those described in here (e.g., penalized LRTs and estimation of PL) will need to be applied alongside the more traditional inferential methods.

Key words Codon substitution model, dN/dS, False positives, Maximum likelihood, Mechanistic model, Model misspecification, Mutation-selection model, Parameter confounding, Phenomenological load, Phenomenological model, Positive selection, Reliability, Statistical inference, Site-specific fitness landscape

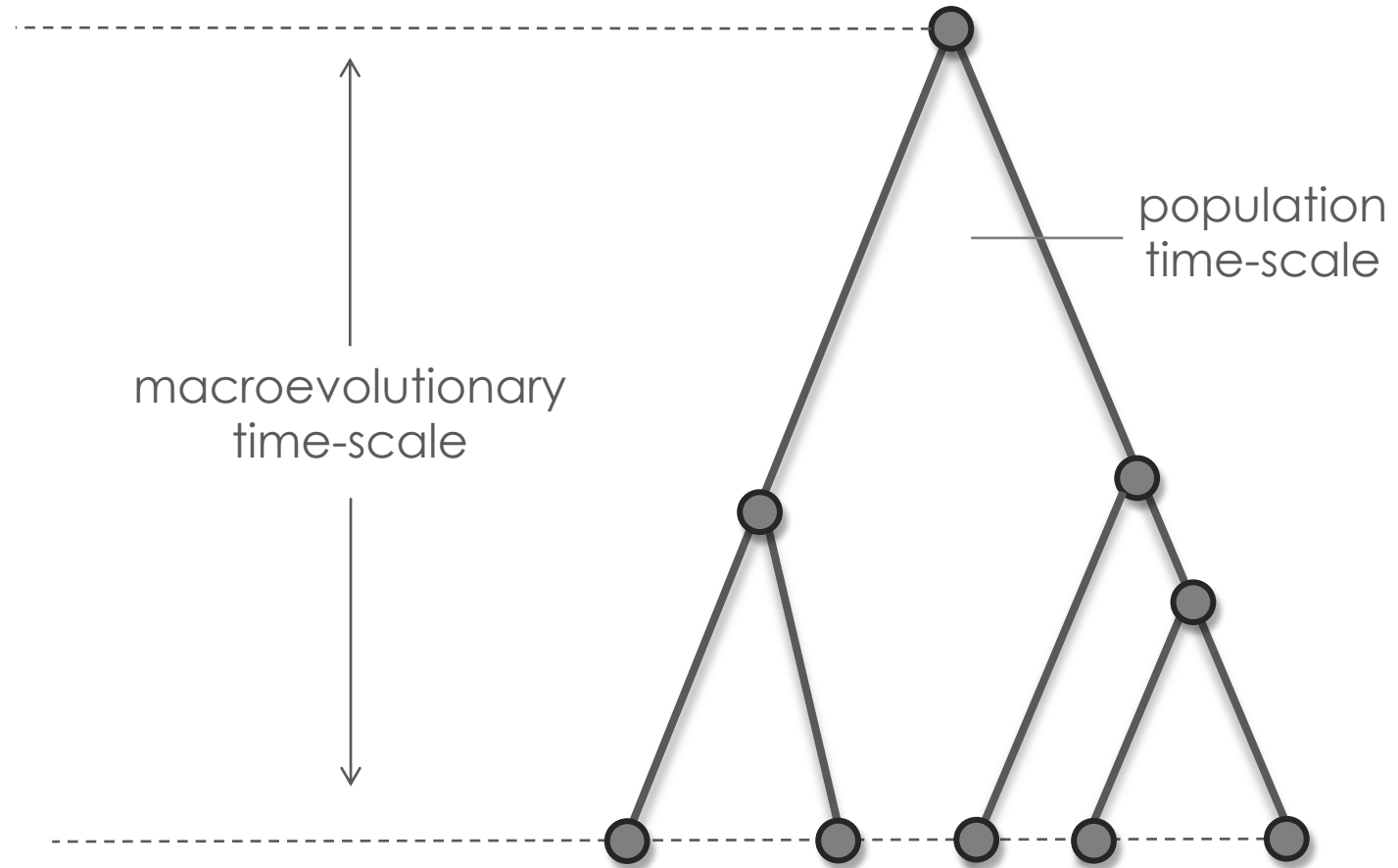
Details: **book chapter** (PDF) on course website

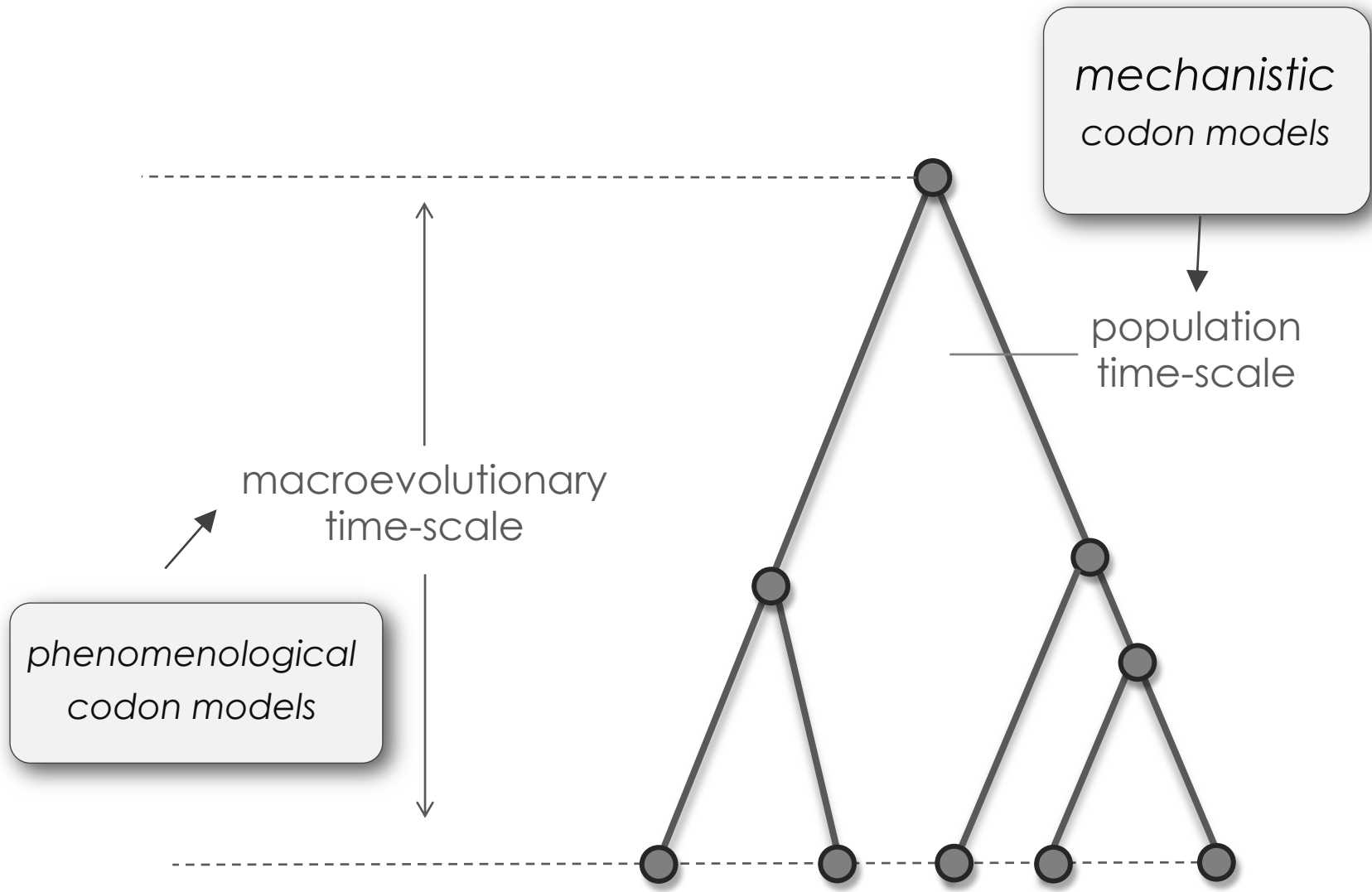
1. mechanistic codon models

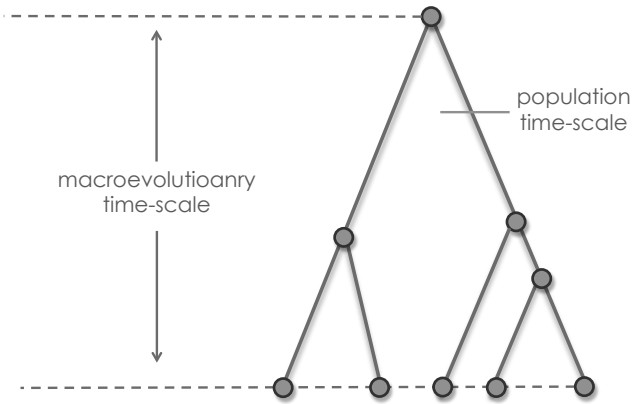
1. MECHANISTIC CODON MODELS



reconciling evolutionary time scales







mechanistic
codon models

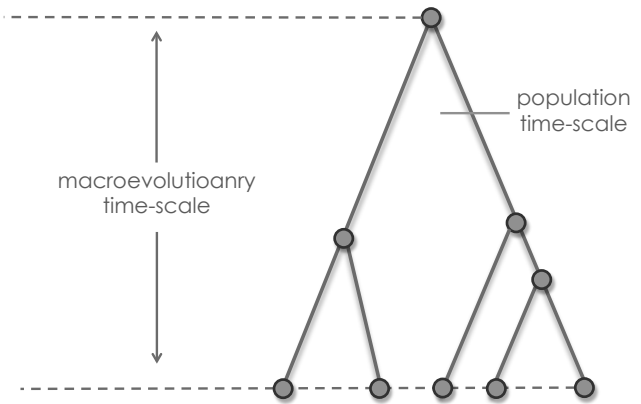
“MUTSEL MODELS”

$$\text{Pr} = \begin{cases} \mu_{ij} & \text{if neutral} \\ \mu_{ij} N \times \frac{2s_{ij}}{1 - e^{-2Ns_{ij}}} & \text{if selected} \end{cases}$$

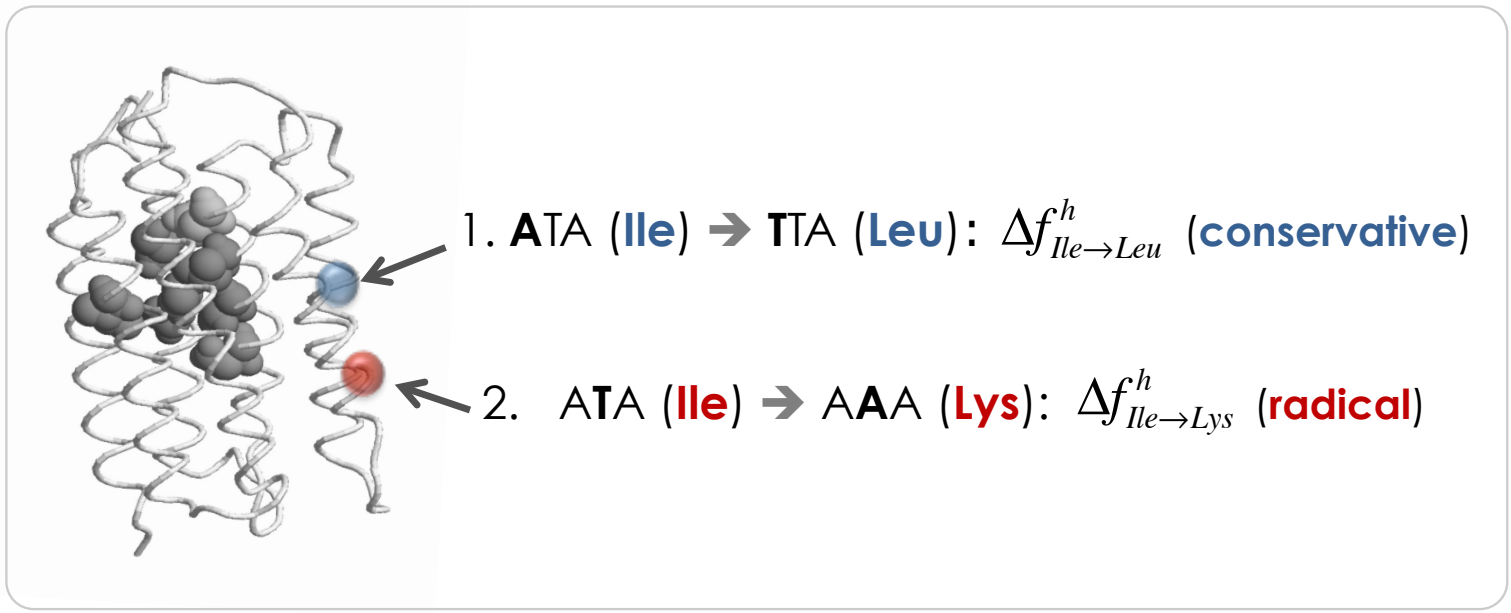
$$s_{ij} = \Delta f_{ij}$$

Halpern and Bruno (1998)

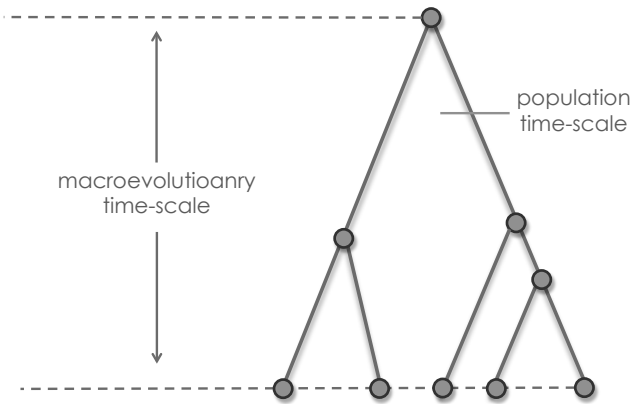
- Wright-Fisher population
- drift: N
- mutation: μ
- selection: s_{ij}
- s_{ij} vary among sites AND amino acids



mechanistic
codon models



- **realism:** fixation probability depends on fitness of ancestral and derived amino acids in the context of the protein.
- **the cost of realism: usually too complex to fit such a model to real data** (caveat: some versions will allow new ways to analyze big datasets)



*mechanistic
codon models*

population genetics at a single codon site (h)

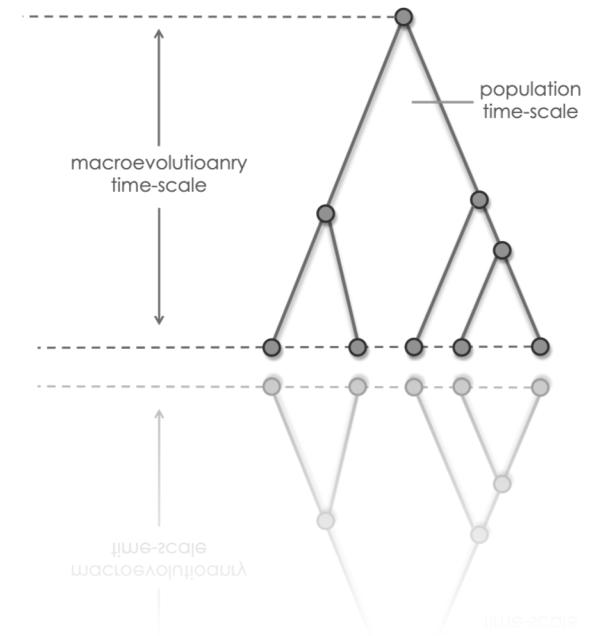
fitness coefficients $f^h = \langle f_1, \dots, f_{61} \rangle$

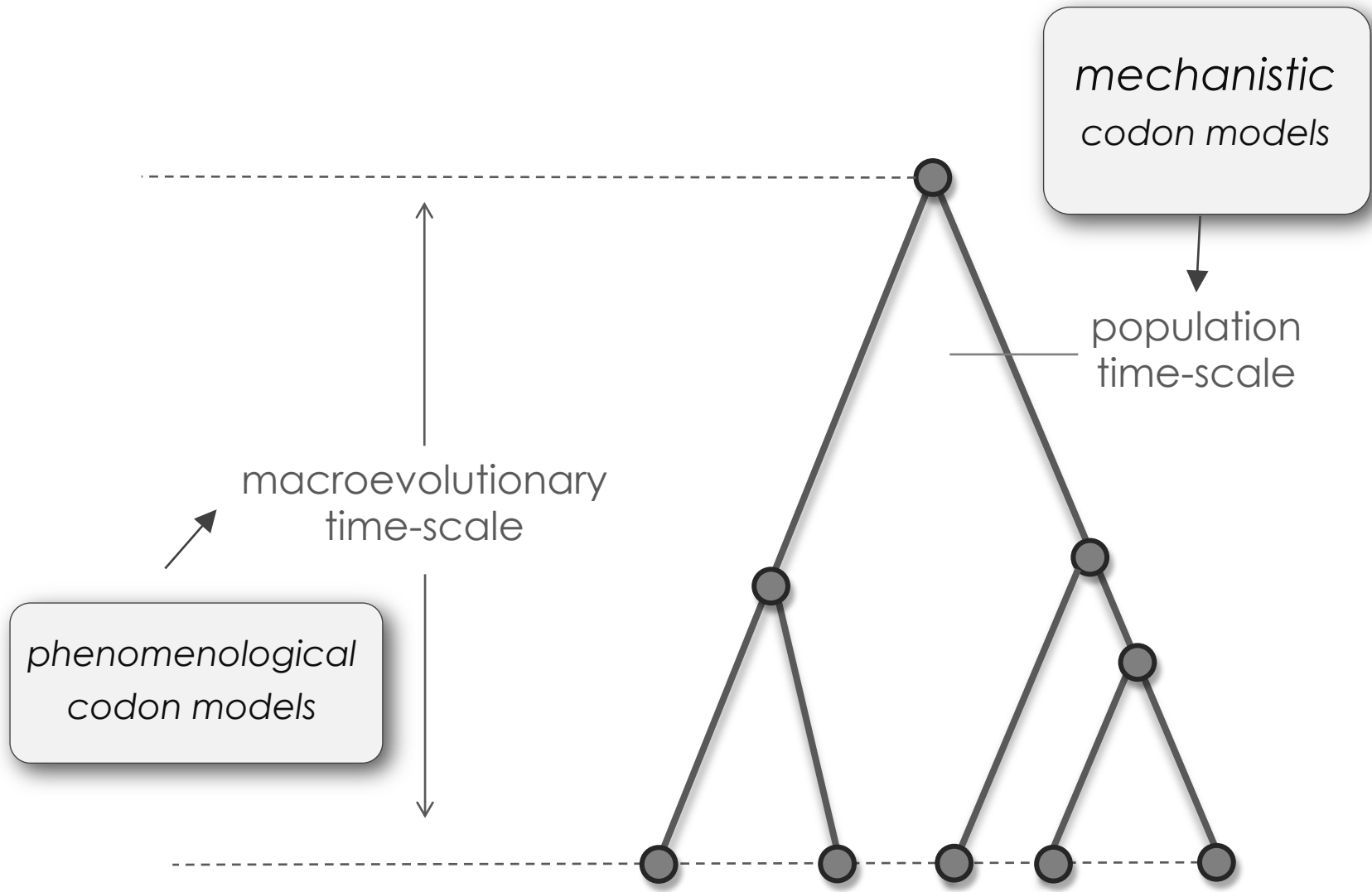
selection coefficients $s_{ij}^h = f_j^h - f_i^h$

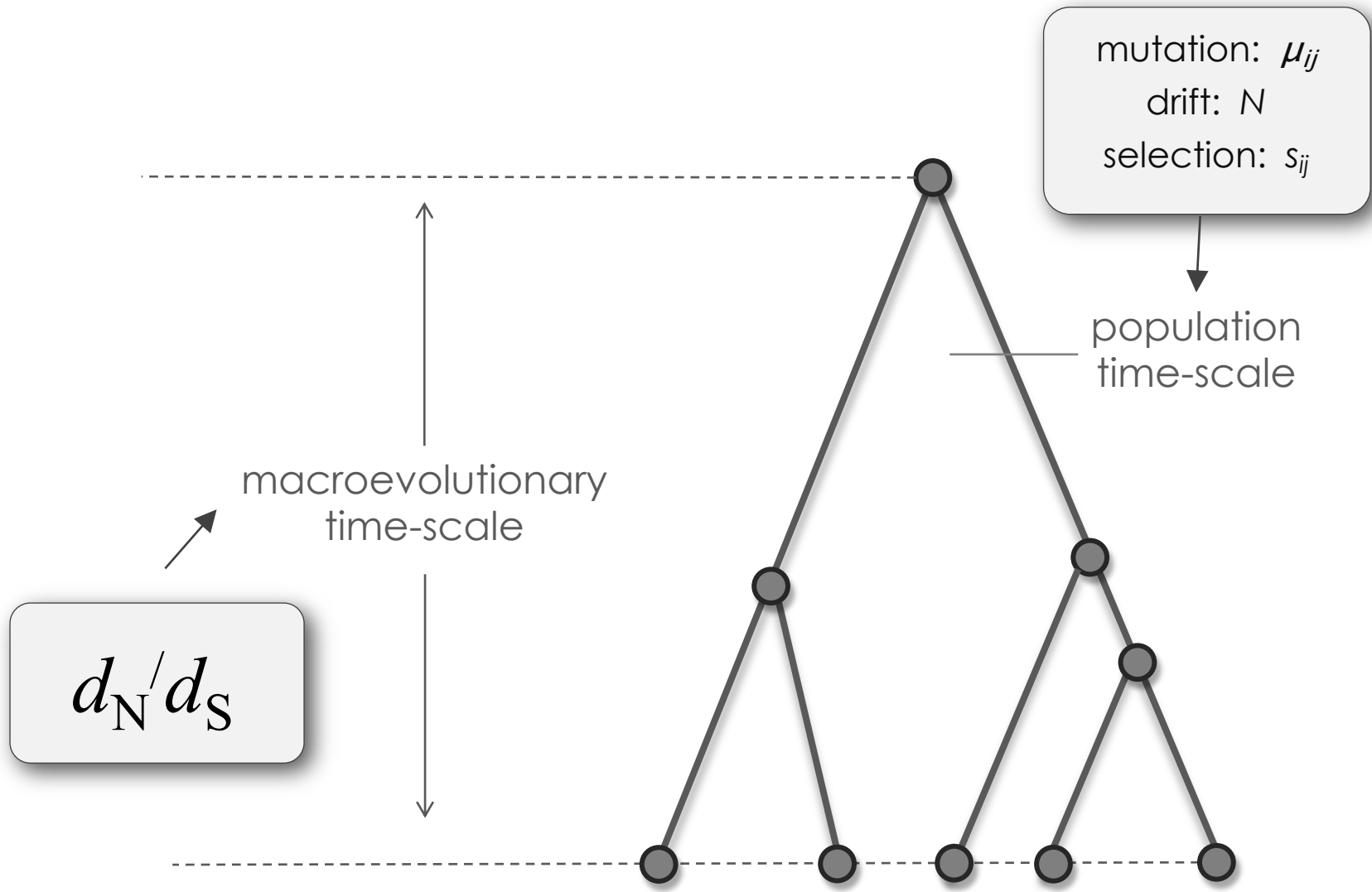
fixation probability (Kimura, 1962) $\Pr(s_{ij}^h) = \frac{2s_{ij}^h}{1 - e^{-2Ns_{ij}^h}}$

2. phenomenological codon models

2. phenomenological codon models

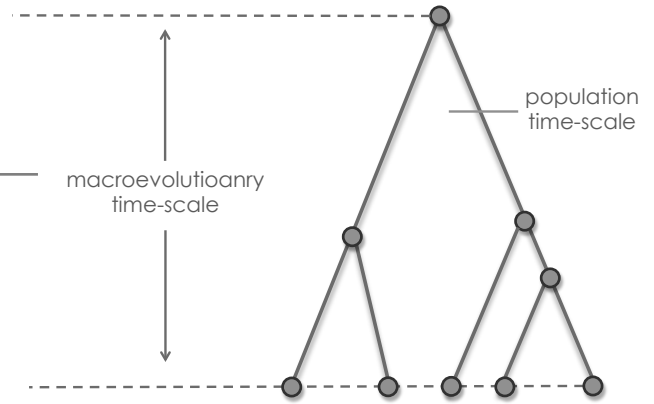






- phenomenological parameters
- ts/tv ratio: κ
- codon frequencies: π_j
- $\omega = dN/dS$
- parameter estimation via ML
- *stationary process*

phenomenological models



“OMEGA MODELS”

$$q_{ij} = \begin{cases} 0 & \text{if } i \text{ and } j \text{ differ by } > 1 \\ \pi_j & \text{for synonymous tv.} \\ \kappa\pi_j & \text{for synonymous ts.} \\ \omega\pi_j & \text{for non-synonymous tv.} \\ \omega\kappa\pi_j & \text{for non-synonymous ts.} \end{cases}$$

Goldman and Yang (1994)
Muse and Gaut (1994)

$$\omega = \frac{dN}{dS}$$

the instantaneous **rate matrix**, Q , is very big: 61×61

$$Q_{ij} = \begin{cases} 0 & \text{if } i \text{ and } j \text{ differ by } > 1 \\ \pi_j & \text{for synonymous tv.} \\ \kappa\pi_j & \text{for synonymous ts.} \\ \omega\pi_j & \text{for non-synonymous tv.} \\ \omega\kappa\pi_j & \text{for non-synonymous ts.} \end{cases}$$



Phenomenological codon models: just a few parameters can cover the 3721 changes between codons!



From codon below:	to codon below:							→	GGG (Gly)
	TTT (Phe)	TTC (Phe)	TTA (Leu)	TTG (Leu)	CTT (Leu)	CTC (Leu)	→		
TTT (Phe)	---	$\kappa\pi_{TTC}$	$\omega\pi_{TTA}$	$\omega\pi_{TTG}$	$\omega\kappa\pi_{TTT}$	0	→	0	
TTC (Phe)	$\kappa\pi_{TTT}$	---	$\omega\pi_{TTA}$	$\omega\pi_{TTG}$	0	$\omega\kappa\pi_{CTC}$	→	0	
TTA (Leu)	$\omega\pi_{TTT}$	$\omega\pi_{TTC}$	---		0	0	→	0	
TTG (Leu)	$\omega\pi_{TTT}$	$\omega\pi_{TTC}$	$\kappa\pi_{TTA}$	---	0	0	→	0	
CTT (Leu)	$\omega\kappa\pi_{TTT}$	0	0	0	---	$\kappa\pi_{CTC}$	→	0	
CTC (Leu)	0	$\omega\kappa\pi_{TTC}$	0	0	$\kappa\pi_{TTT}$	---	→	0	
↓	↓	↓	↓	↓	↓	↓	↘		
GGG (Gly)	0	0	0	0	0	0	0	---	

* This is equivalent to the codon model of Goldman and Yang (1994). Parameter ω is the ratio d_N/d_S , κ is the transition/transversion rate ratio, and π_i is the equilibrium frequency of the target codon (i).

the instantaneous rate matrix, Q , is very big: 61×61

$$Q_{ij} = \begin{cases} 0 & \text{if } i \text{ and } j \text{ differ by } > 1 \\ \pi_j & \text{for synonymous tv.} \\ \kappa\pi_j & \text{for synonymous ts.} \\ \omega\pi_j & \text{for non-synonymous tv.} \\ \omega\kappa\pi_j & \text{for non-synonymous ts.} \end{cases}$$



Phenomenological codon models: just a few parameters can cover the 3721 changes between codons!



context matters...

From codon below:	to codon below:							→	GGG (Gly)
	TTT (Phe)	TTC (Phe)	TTA (Leu)	TTG (Leu)	CTT (Leu)	CTC (Leu)			
TTT (Phe)	---	$\kappa\pi_{TTC}$	$\omega\pi_{TTA}$	$\omega\pi_{TTG}$	$\omega\kappa\pi_{TTT}$	0	→	0	
TTC (Phe)	$\kappa\pi_{TTT}$	---	$\omega\pi_{TTA}$	$\omega\pi_{TTG}$	0	$\omega\kappa\pi_{CTC}$	→	0	
TTA (Leu)	$\omega\pi_{TTT}$	$\omega\pi_{TTC}$	---		0	0	→	0	
TTG (Leu)	$\omega\pi_{TTT}$	$\omega\pi_{TTC}$	$\kappa\pi_{TTA}$	---	0	0	→	0	
CTT (Leu)	$\omega\kappa\pi_{TTT}$	0	0	0	---	$\kappa\pi_{CTC}$	→	0	
CTC (Leu)	0	$\omega\kappa\pi_{TTC}$	0	0	$\kappa\pi_{TTT}$	---	→	0	
↓	↓	↓	↓	↓	↓	↓	↘		
GGG (Gly)	0	0	0	0	0	0	0	---	

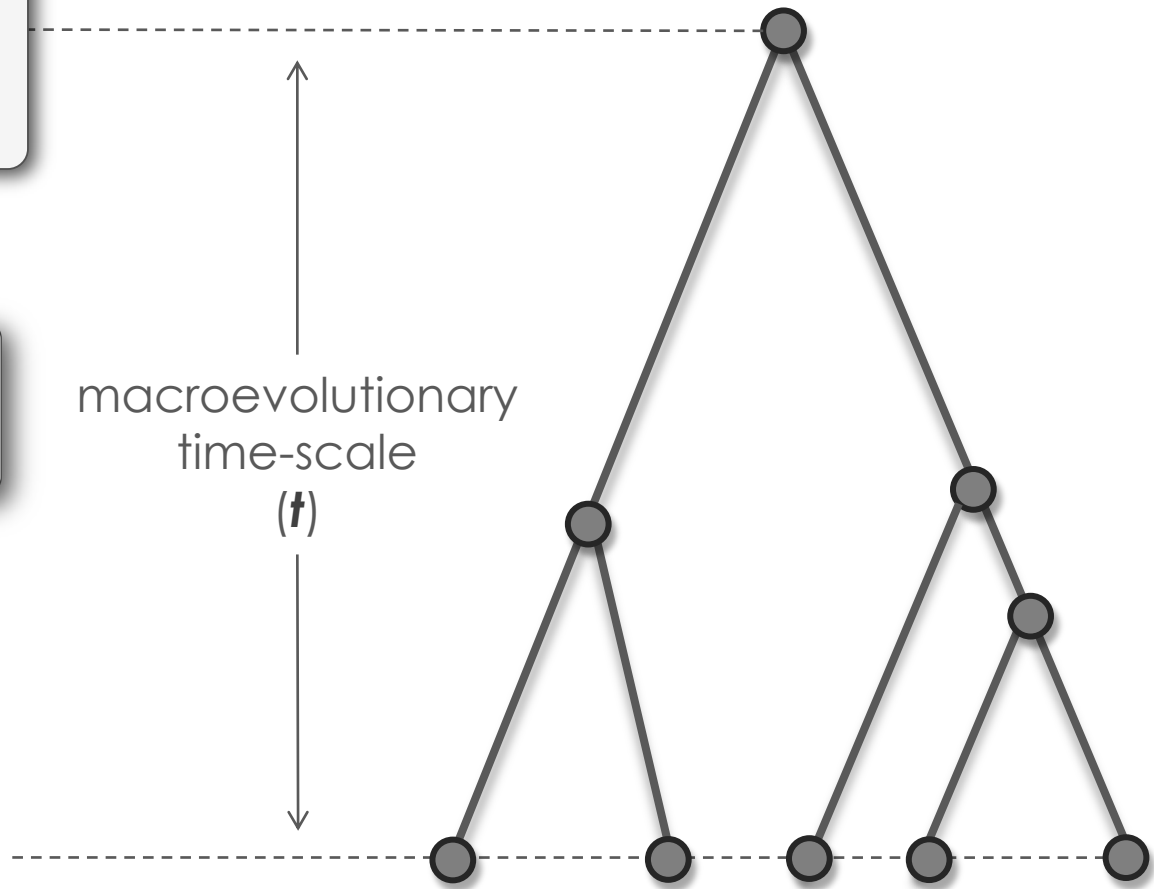
* This is equivalent to the codon model of Goldman and Yang (1994). Parameter ω is the ratio d_N/d_S , κ is the transition/transversion rate ratio, and π_i is the equilibrium frequency of the target codon (i).

probability of substitution between codons over time, $P(t)$

$$Q_{ij} = \begin{cases} 0 & \text{if } i \text{ and } j \text{ differ by } > 1 \\ \pi_j & \text{for synonymous tv.} \\ \kappa\pi_j & \text{for synonymous ts.} \\ \omega\pi_j & \text{for non-synonymous tv.} \\ \omega\kappa\pi_j & \text{for non-synonymous ts.} \end{cases}$$

$$P(\mathbf{t}) = \{p_{ij}(\mathbf{t})\} = e^{Q\mathbf{t}}$$

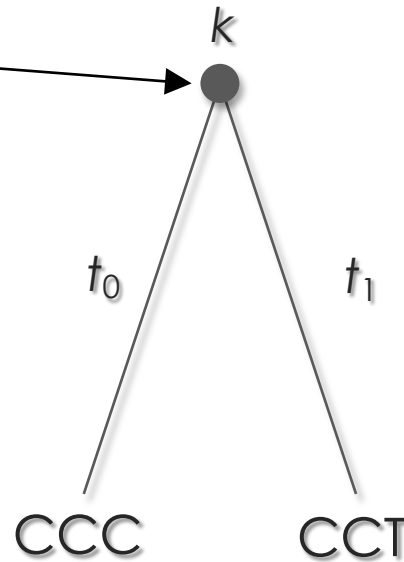
recall that **Paul Lewis** introduced **Q matrices** and how to obtain **transition probabilities**



likelihood of the data at a site

$$L_h(\text{CCC}, \text{CCT}) = \sum_k \pi_k p_{k\text{CCC}}(t_0) p_{k\text{CCT}}(t_1)$$

the likelihood is a **sum over all possible ancestral codon states** that could have been observed at node k



recall that **Paul Lewis** described how to compute the likelihood of the data at a site for a DNA model. The only difference here is that the states are codons rather than nucleotides

note: analysis is typically done by using an unrooted tree

likelihood of the data at all sites

The likelihood of observing the entire sequence alignment is the product of the probabilities at each site.

Paul Lewis covered this with the “**AND**” rule in his likelihood lecture

$$L = L_1 \times L_2 \times L_3 \times \dots \times L_N = \prod_{h=1}^N L_h$$

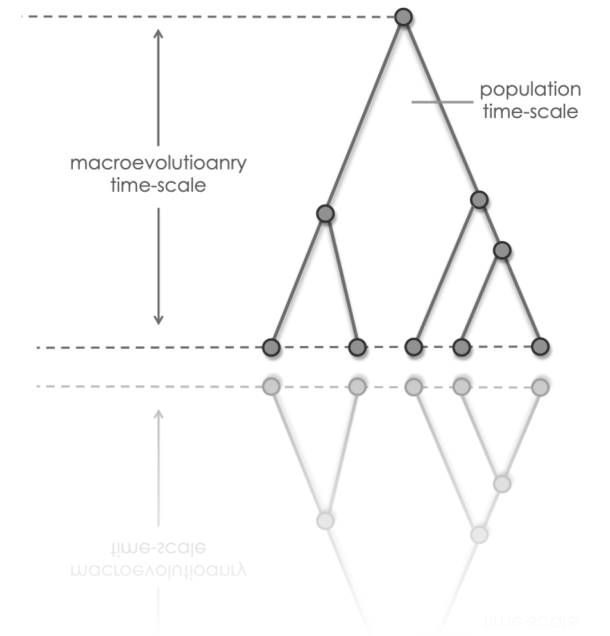
see **Paul Lewis's** lecture slides for more about likelihoods vs. log-likelihoods

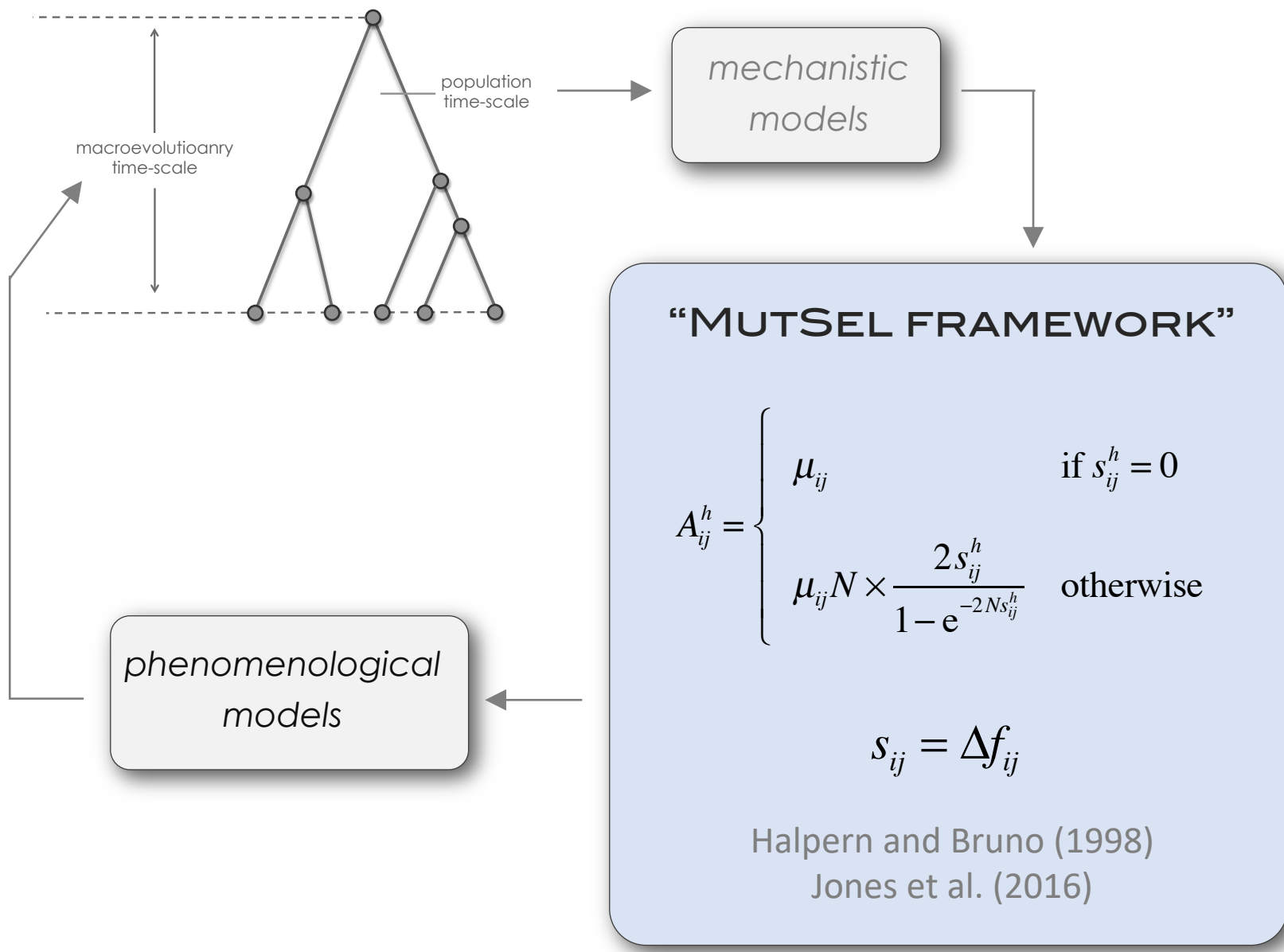
The log likelihood is a sum over all sites.

$$\ell = \ln\{L\} = \ln\{L_1\} + \ln\{L_2\} + \ln\{L_3\} + \dots + \ln\{L_N\} = \sum_{h=1}^N \ln\{L_h\}$$

3. bridging selection between time-scales

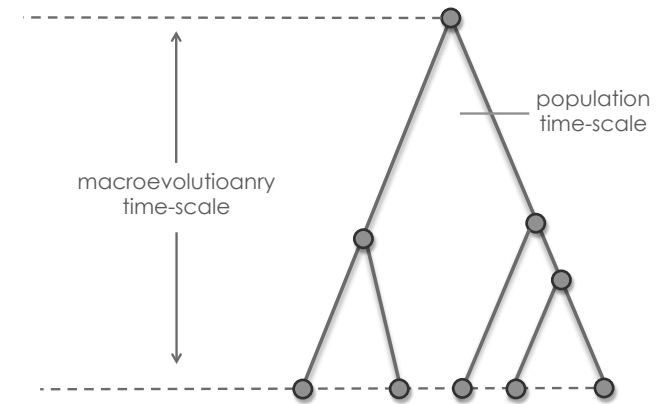
3. bridging selection between time-scales





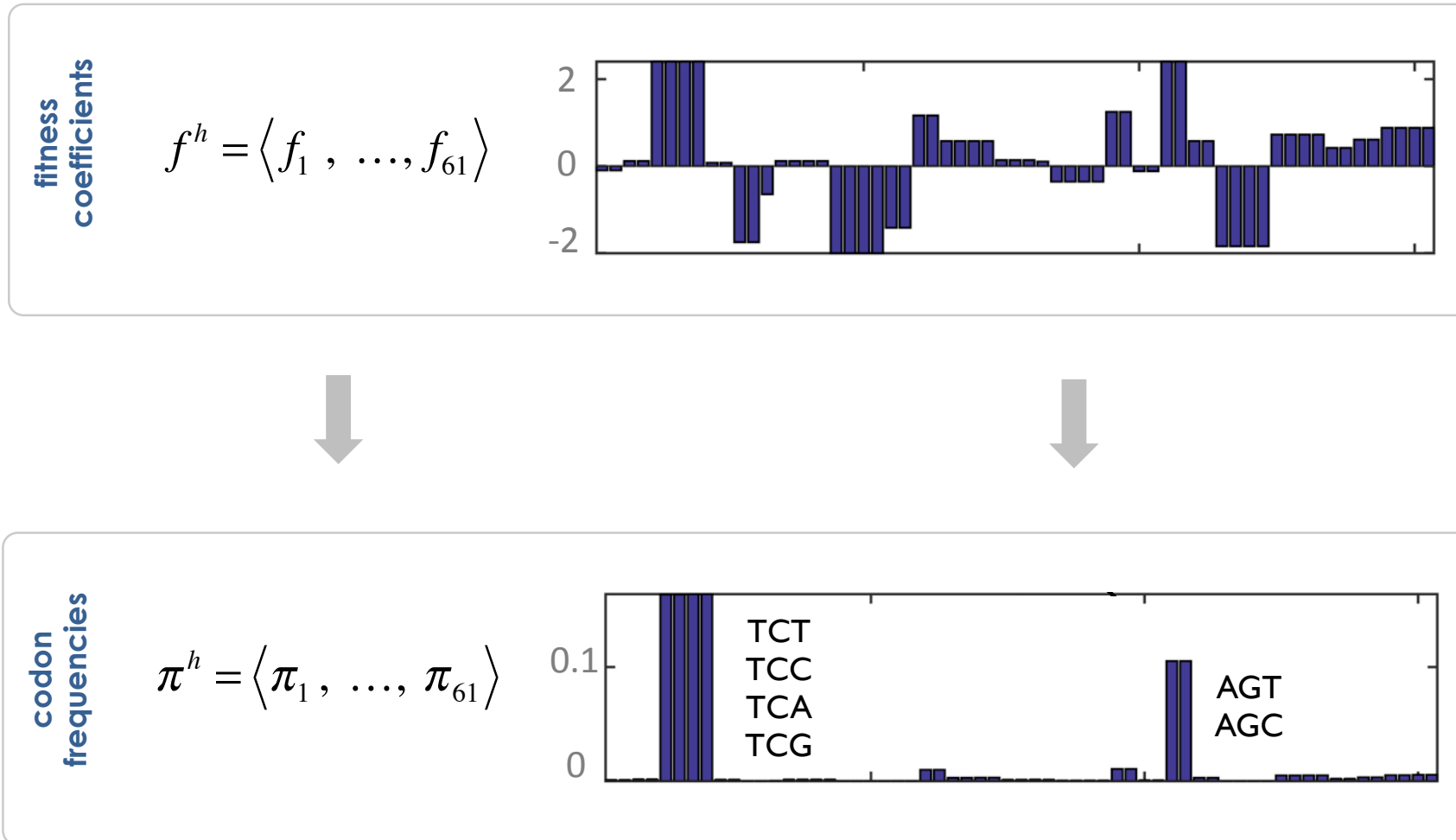
Two explicit ways to reconcile **population genetics** and **macroevolution**:

1. map fitness to equilibrium frequencies
2. expected index of selection intensity



(1) Sella and Hirsh 2005; (2) Jones et al. 2016

1. fitness coefficients map to stationary codon frequencies



(Sella and Hirsh 2005)

2. from fitness coefficients to expected dN/dS

MUTSEL RATE MATRIX

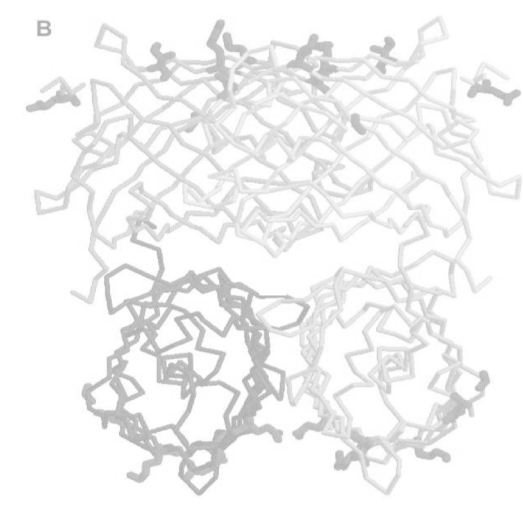
$$dN^h / dS^h = \frac{E[\text{evolution w/ selection}]}{E[\text{evolution by drift alone}]}$$

$$dN^h / dS^h = \frac{\sum_{i \neq j} \pi_i^h A_{ij}^h I_N}{\sum_{i \neq j} \pi_i^h \mu_{ij} I_N}$$

- $dN/dS = \omega$ when matrix A^h is replaced by matrix Q of model M0
- dN/dS is an analog of ω under MutSel

4. three positive selection scenarios

4. three positive selection scenarios

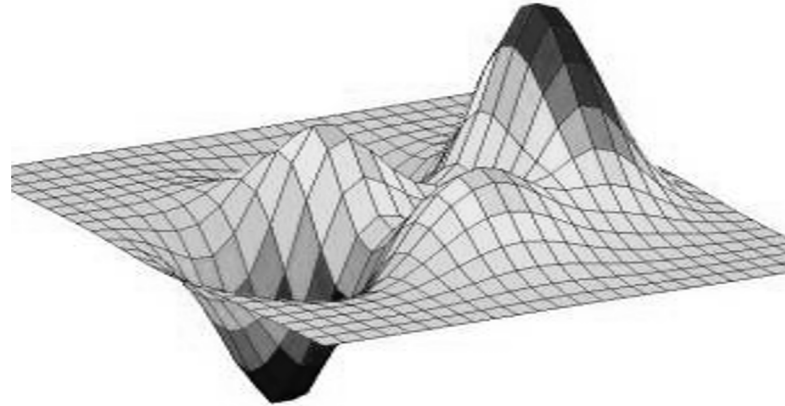


1932: adaptive landscapes and “shifting balance”



Sewall Wright

- introduces “ADAPTIVE LANDSCAPE” as a metaphor



- introduces “SHIFTING BALANCE” as a model
(SBT **more complex** than I will present)

positive selection: 3 evolutionary scenarios

1

frequency dependent
selection

2

episodic adaptation

3

non-adaptive shifting
balance

dynamic
fitness
landscape

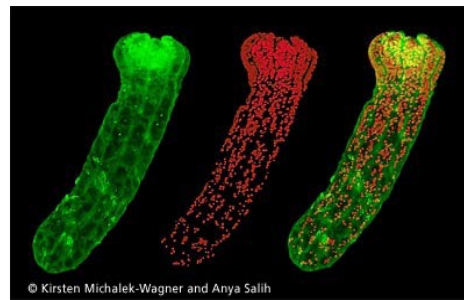
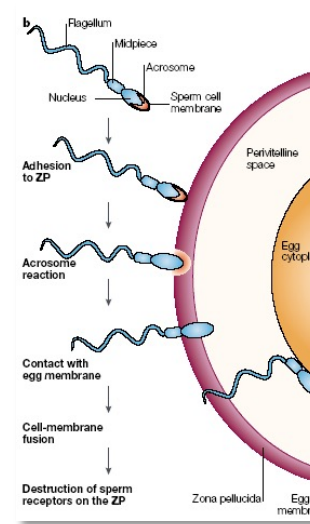
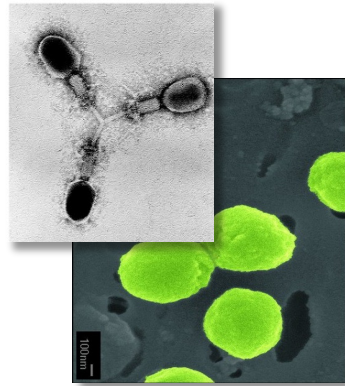
static
fitness
landscape

1. antagonistic evolutionary interaction

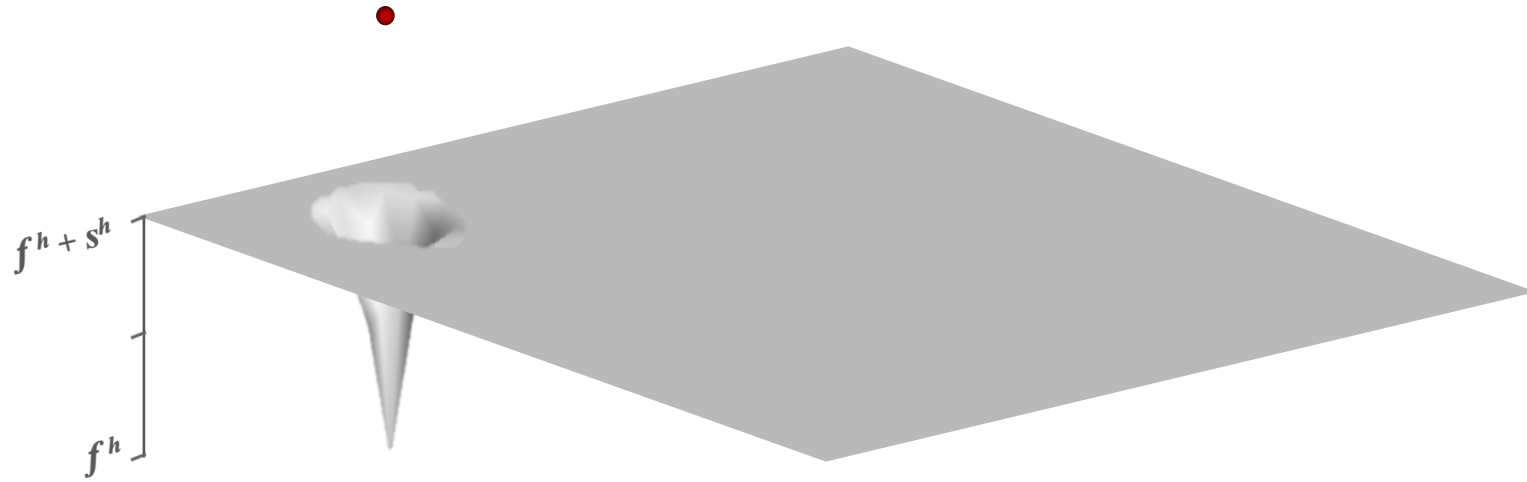
host-pathogen

sexual-conflict

molecular-interactions

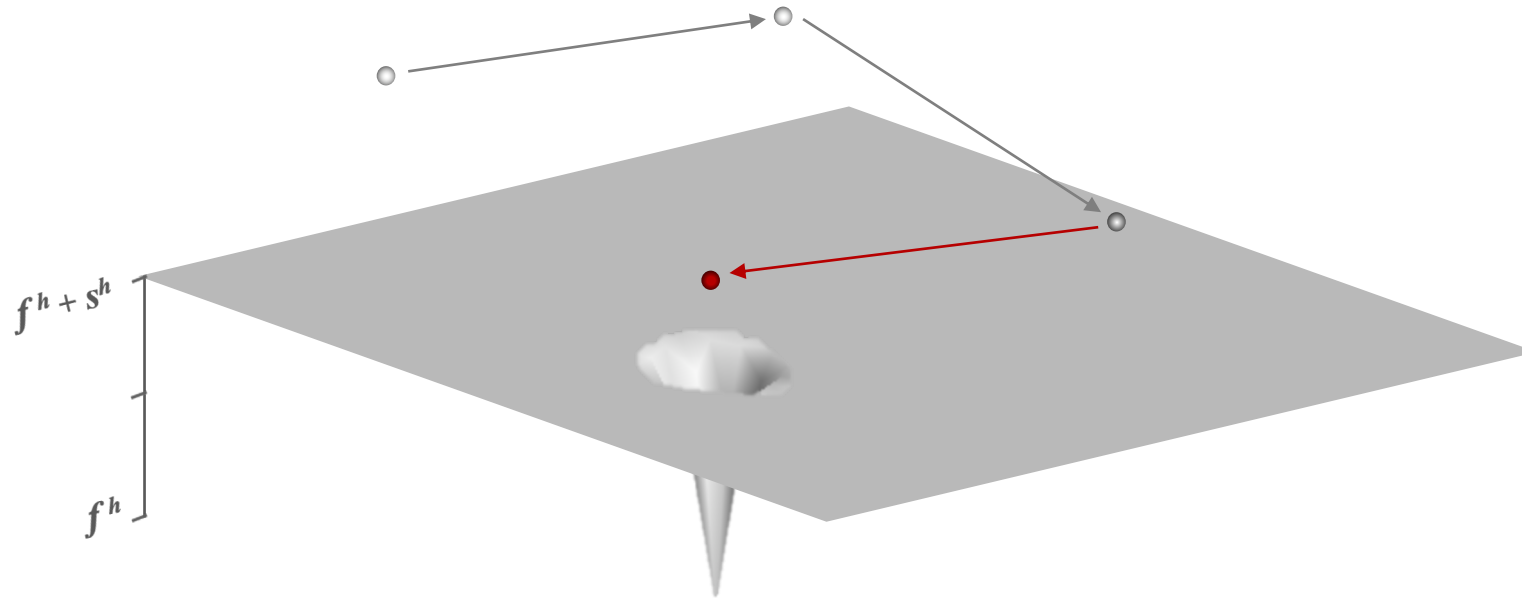


1 frequency-dependent adaptive landscape (weird)



1 frequency-dependent adaptive landscape (weird)

1. amino acid at a site has f^h ; all others have $f^h + s$
2. fitness values swap when a substitution occurs



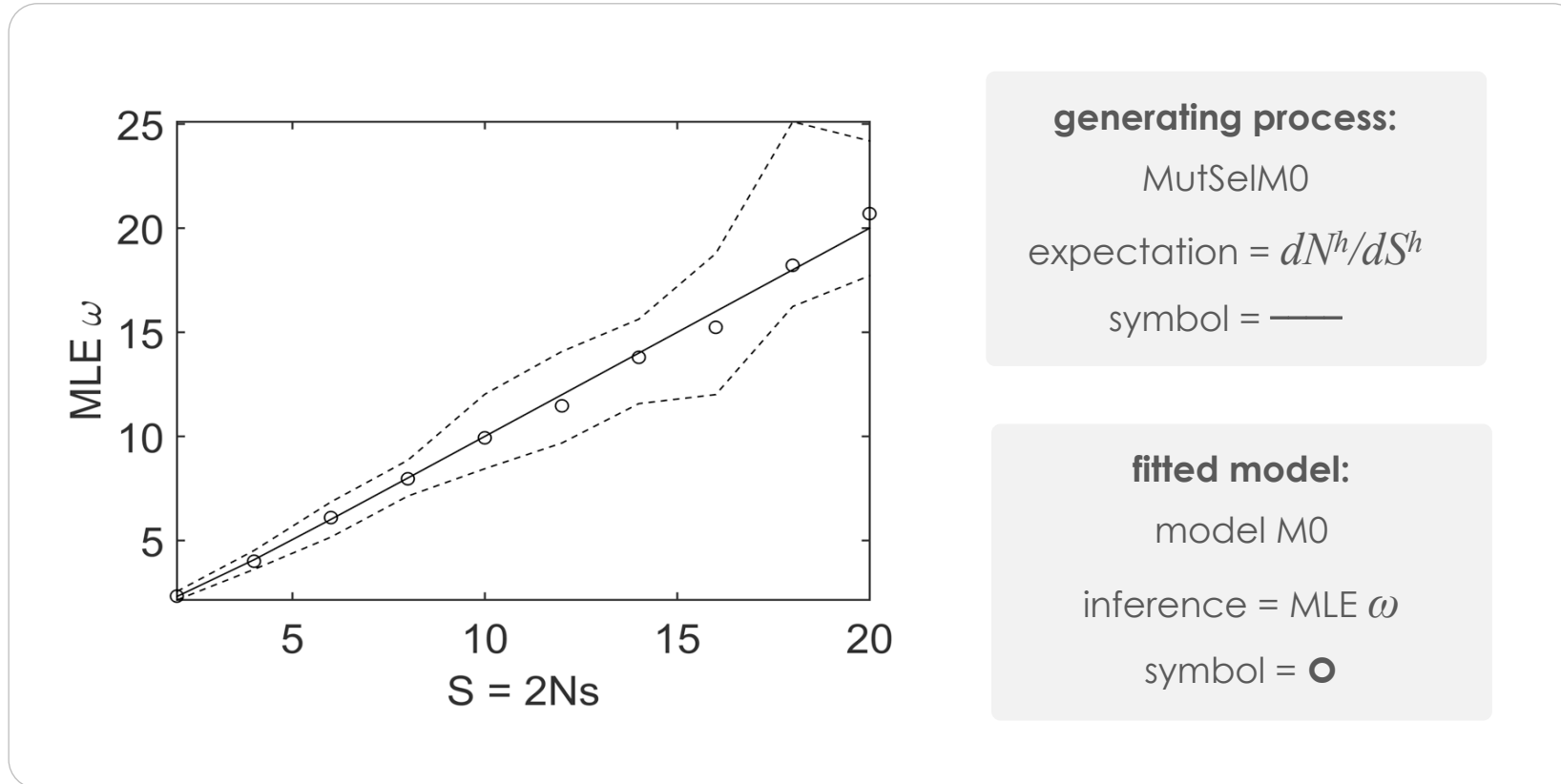
MutSelM0: (1) and (2) above imply Markov chain properties with the same rate matrix Q as **codon model M0**

“OMEGA MODELS”

$$q_{ij} = \begin{cases} 0 & \text{if } i \text{ and } j \text{ differ by } > 1 \\ \pi_j & \text{for synonymous tv.} \\ \kappa\pi_j & \text{for synonymous ts.} \\ \omega\pi_j & \text{for non-synonymous tv.} \\ \omega\kappa\pi_j & \text{for non-synonymous ts.} \end{cases}$$

Goldman and Yang (1994)
Muse and Gaut (1994)

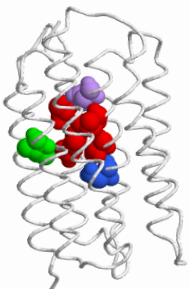
1 frequency-dependent adaptive landscape (weird)



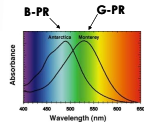
conclusion: phenomemolgical codon models
assume frequency-dependent selection

2. episodic Darwinian adaptation

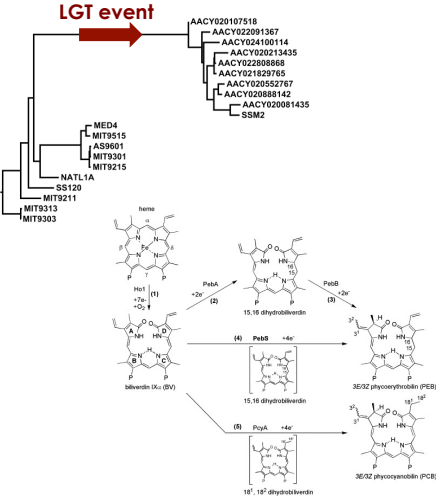
exploitation of a new niche



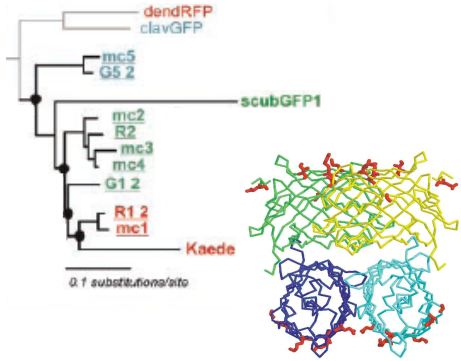
Spectral tuning switch (105)
Green (540) to Blue (490nm)



lateral gene transfer (LGT)

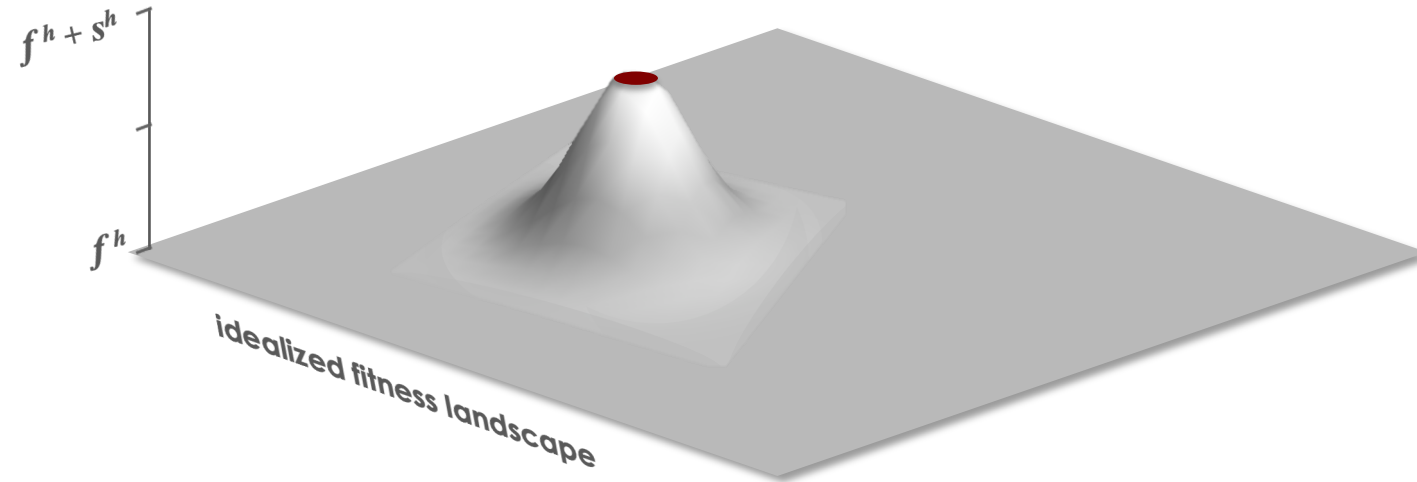


gene duplication



2 adaptive peak shift: evolution of novel function

optimal function in a stable environment



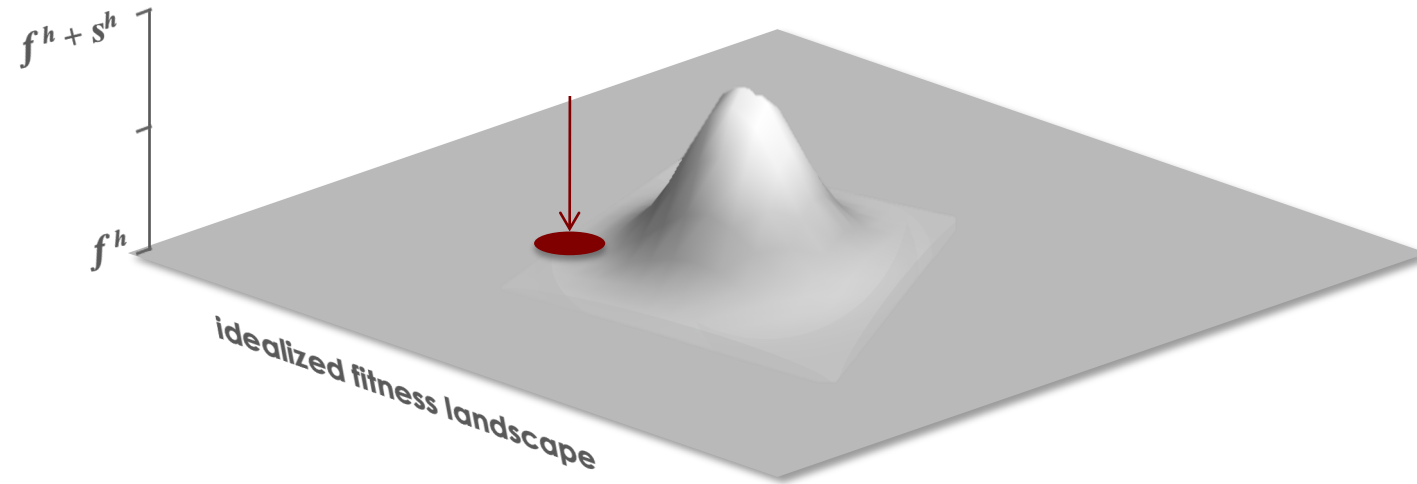
population: at fitness peak

fitness peak: stationary

FFTNS: keeps population at peak

2 adaptive peak shift: evolution of novel function

sub-optimal function in a novel environment



population: lower fitness

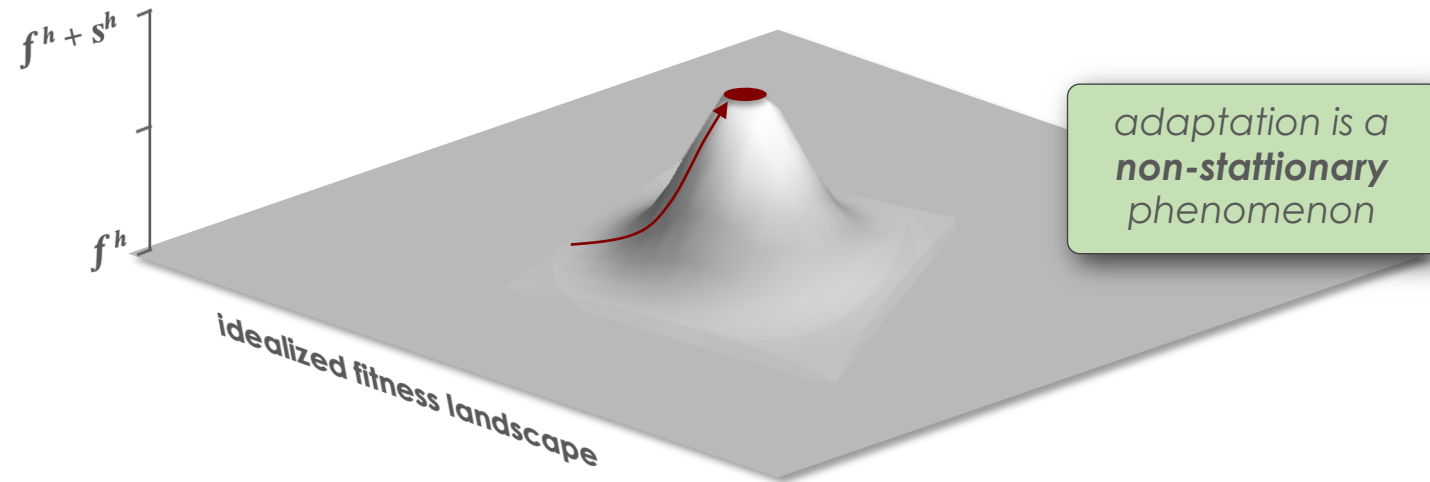
fitness peak: moving

FFTNS: increase population mean fitness

(non-stationary process)

2 adaptive peak shift: evolution of novel function

episodic adaptive evolution of a novel function



population: returns to peak

fitness peak: stabilized

FFTNS: increases population mean
fitness until at peak

BIOLOGY LETTERS

rsbl.royalsocietypublishing.org



CrossMark
click for updates

Research

Cite this article: dos Reis M. 2015 How to calculate the non-synonymous to synonymous rate ratio of protein-coding genes under the Fisher–Wright mutation–selection framework. *Biol. Lett.* **11**: 20141031.
<http://dx.doi.org/10.1098/rsbl.2014.1031>

Received: 8 December 2014

Accepted: 16 March 2015

Molecular evolution

How to calculate the non-synonymous to synonymous rate ratio of protein-coding genes under the Fisher–Wright mutation–selection framework

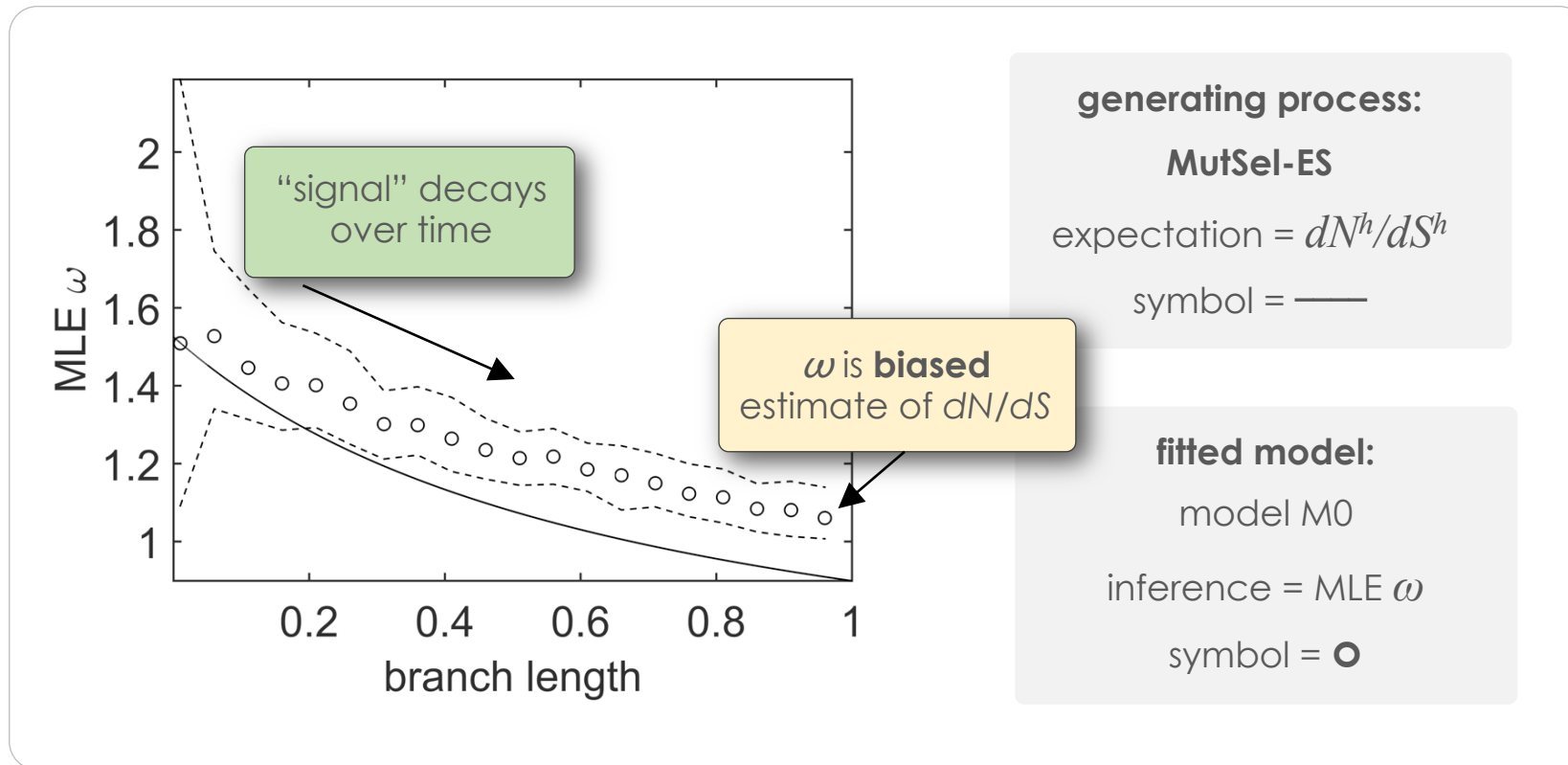
Mario dos Reis

Department of Genetics, Evolution and Environment, University College London, Gower Street, London WC1E 6BT, UK

First principles of population genetics are used to obtain formulae relating the non-synonymous to synonymous substitution rate ratio to the selection coefficients acting at codon sites in protein-coding genes. Two theoretical cases are discussed and two examples from real data (a chloroplast gene and a virus polymerase) are given. The formulae give much insight into the dynamics of non-synonymous substitutions and may inform the development of methods to detect adaptive evolution.

4. The non-synonymous rate during adaptive evolution

2 adaptive peak shift: MutSel-ES



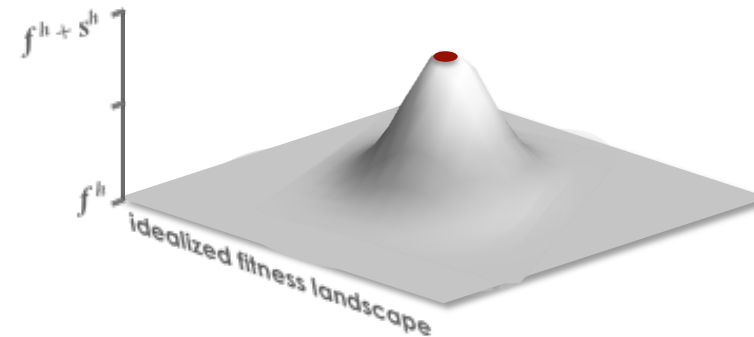
conclusion : episodic models “work” because $w > 1$ is a consequence of a system moving towards a new fitness peak.

conclusion : episodic models “work” because they are sensitive to non-stationary behavior

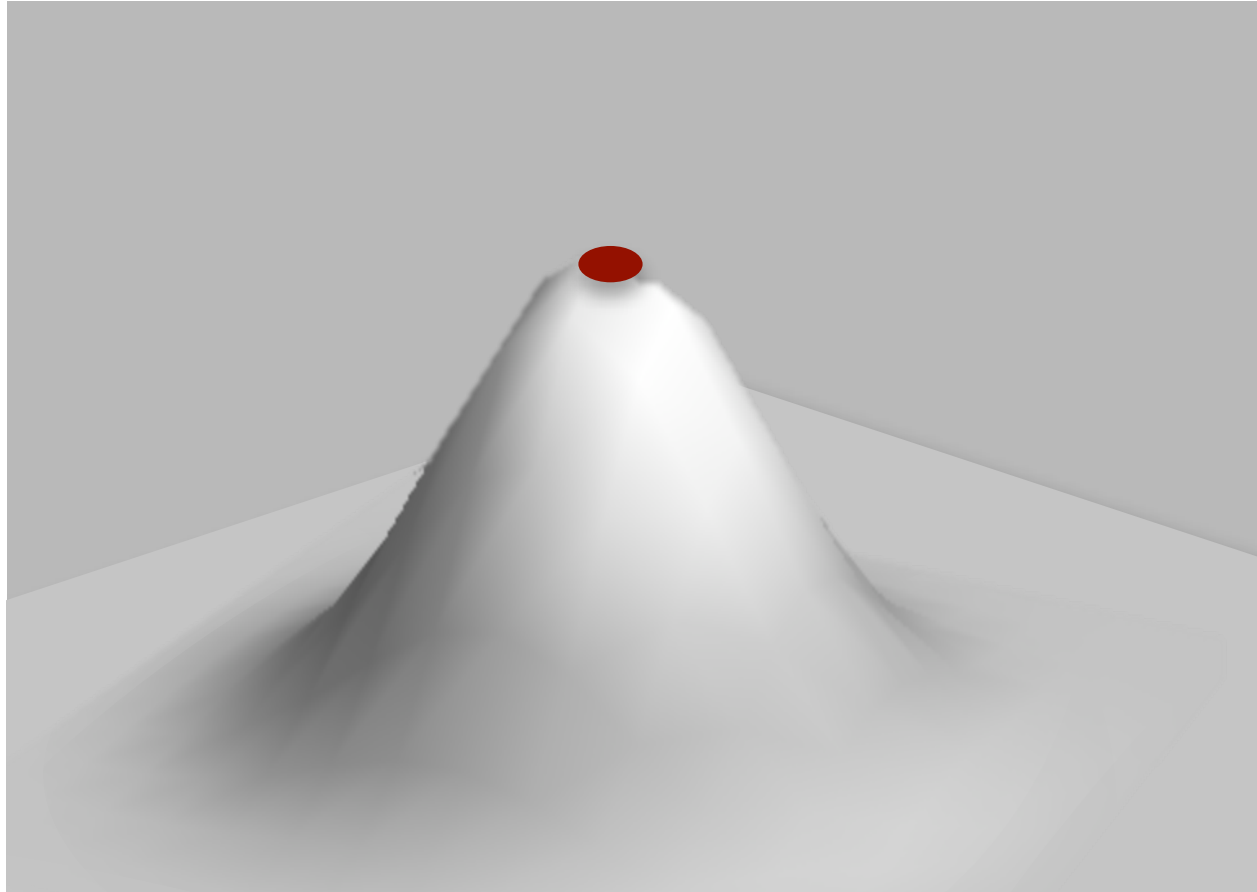
**3. fitness
coefficients are
constant
(fixed-peak)**

Spielman and Wilke (2015)

- dN/dS must be ≤ 1 when fitness coefficients are fixed.
- positive selection is not possible on a stationary fitness peak

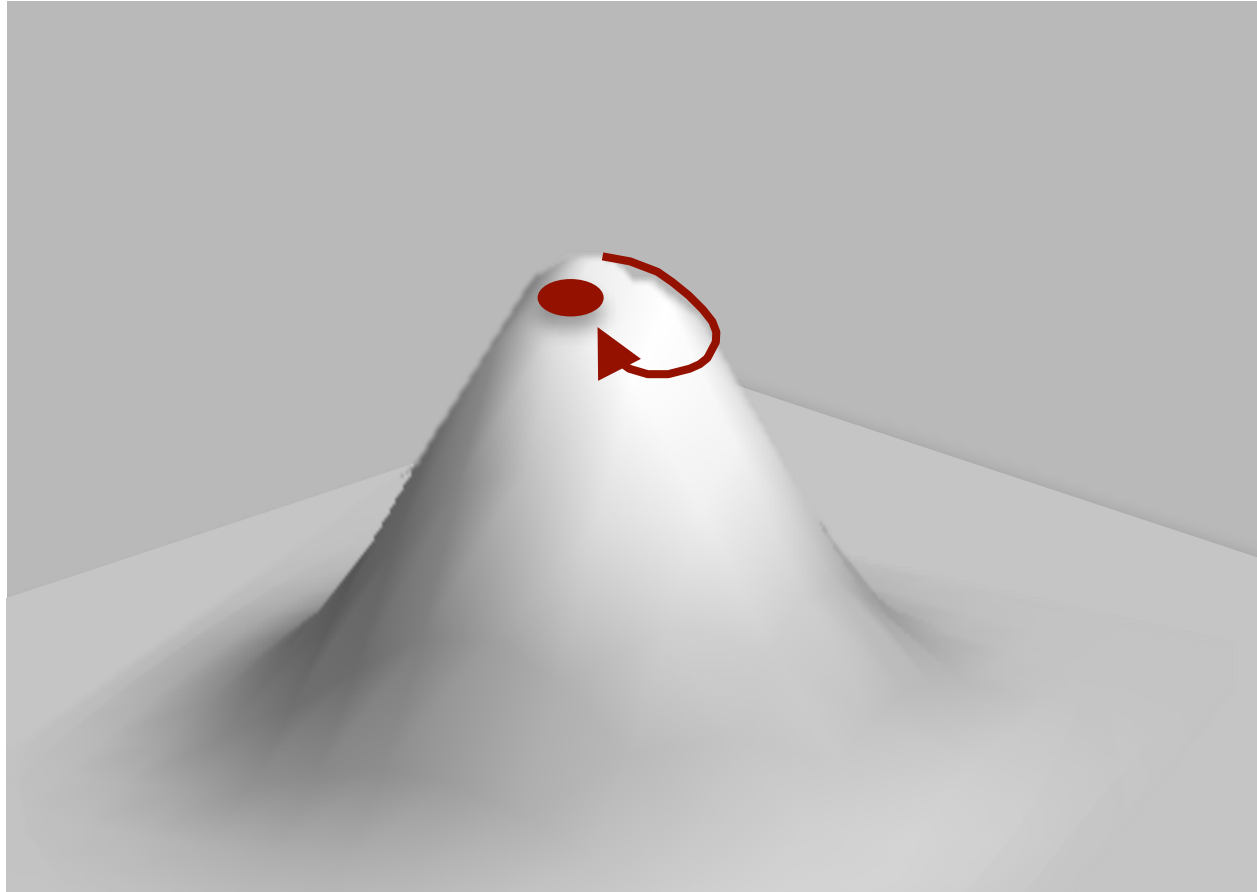


3 shifting balance: movement around stationary peak (non-adaptive)



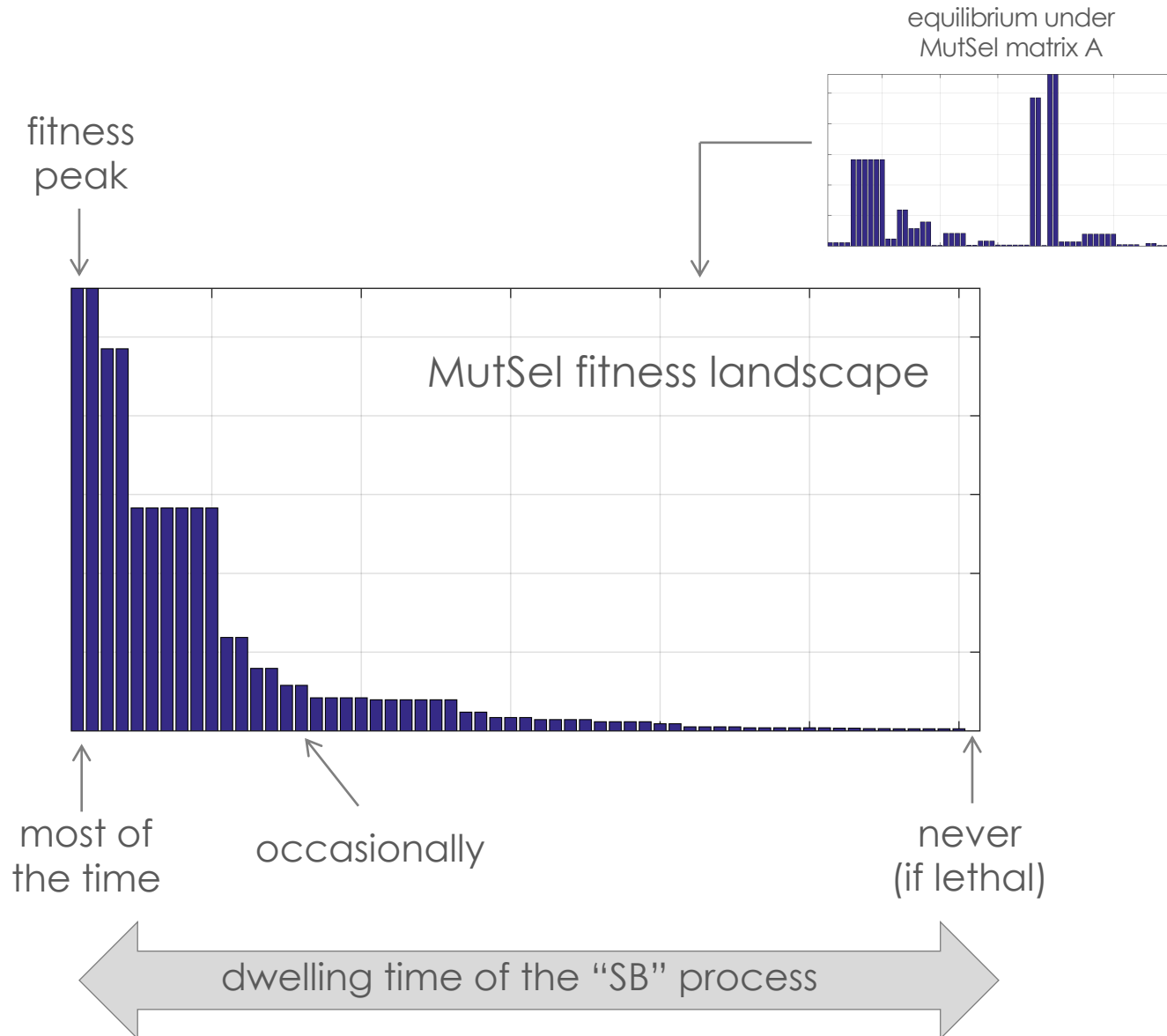
mutation and **drift** can move a pop. off a fitness peak

3 shifting balance: movement around stationary peak (non-adaptive)

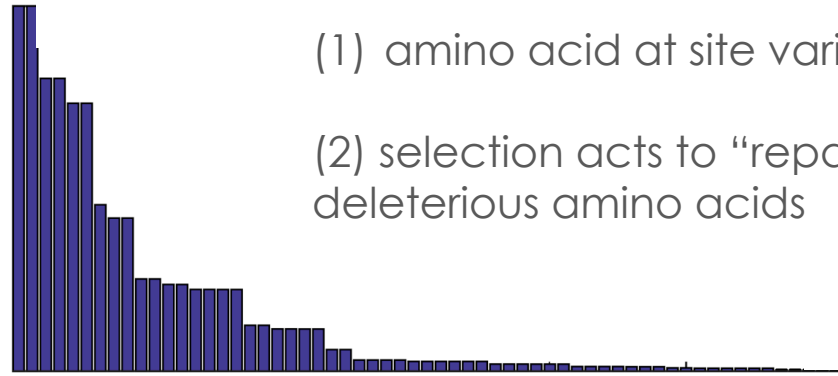


mutation and **drift** can move a pop. off a fitness peak

3 shifting balance: the MutSel landscape (Jones et al. 2016)



3 shifting balance: positive selection on a MutSel landscape



(1) amino acid at site varies over time

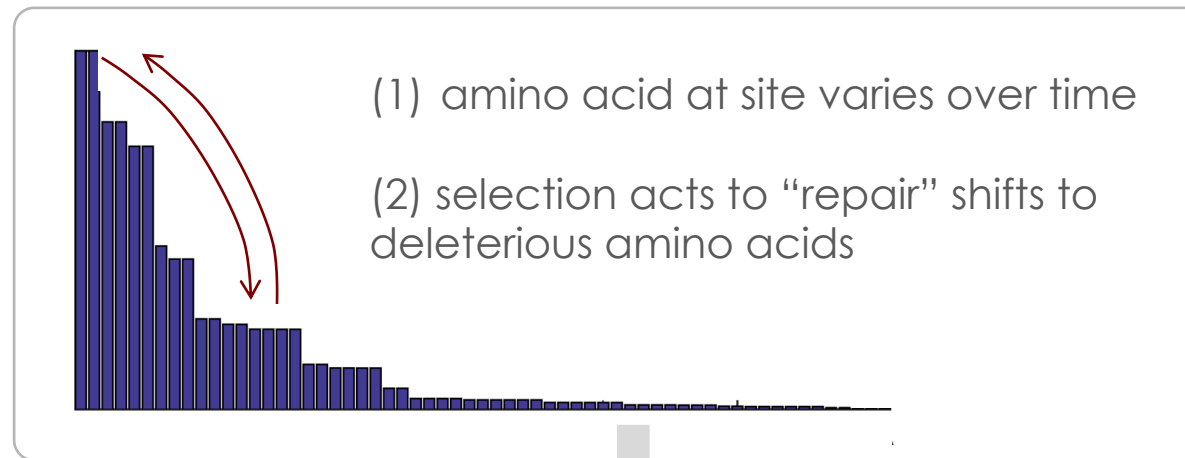
(2) selection acts to “repair” shifts to deleterious amino acids



EXPECTED PROPORTION OF
MUTATIONS FIXED BY SELECTION

$$p_+^h = \frac{\sum_{(i,j)} \pi_i^h (A_{ij}^h - \mu_i) I_+}{\sum_{i \neq j} \pi_i^h A_{ij}^h}$$

conclusion: $p_+ > 0$ as long as number of viable amino acids > 1 at a site



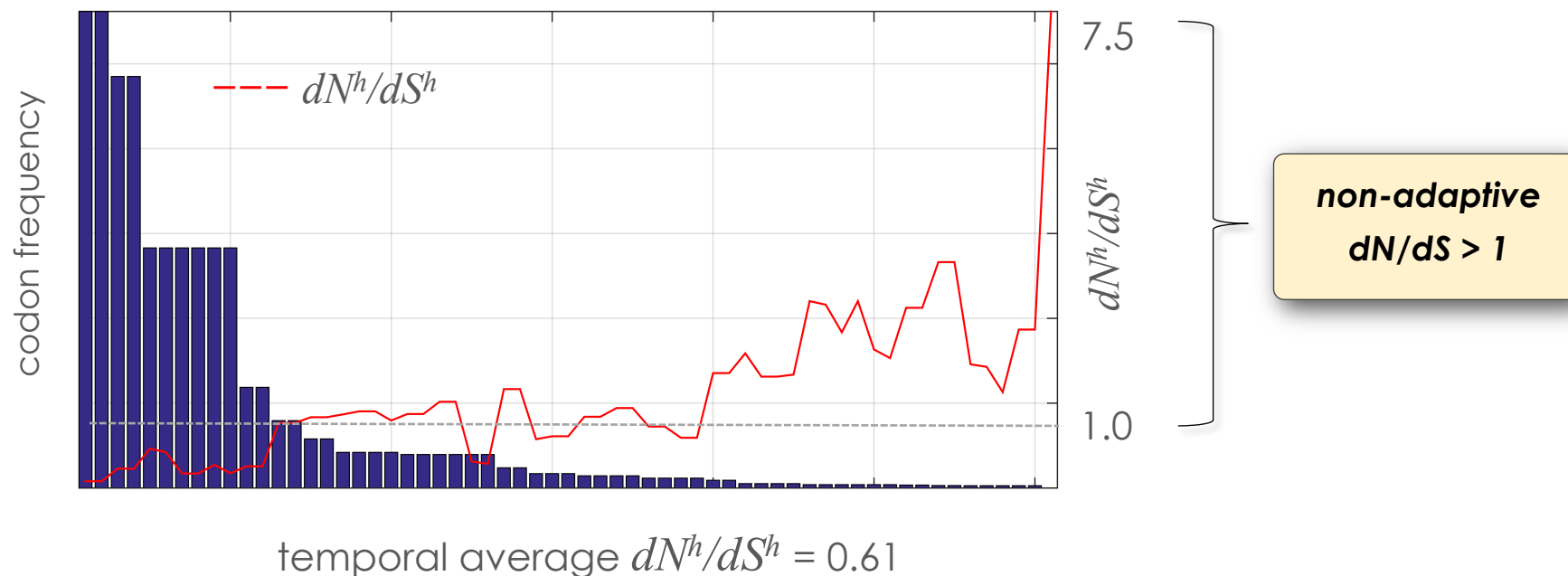
key result:
purifying selection: $p_+ = p_-$
(static landscape)

p_+ = positive selection without adaptation (***maintenance!***)

p_- = related to “fixed drift load”

3 shifting balance: the MutSel landscape

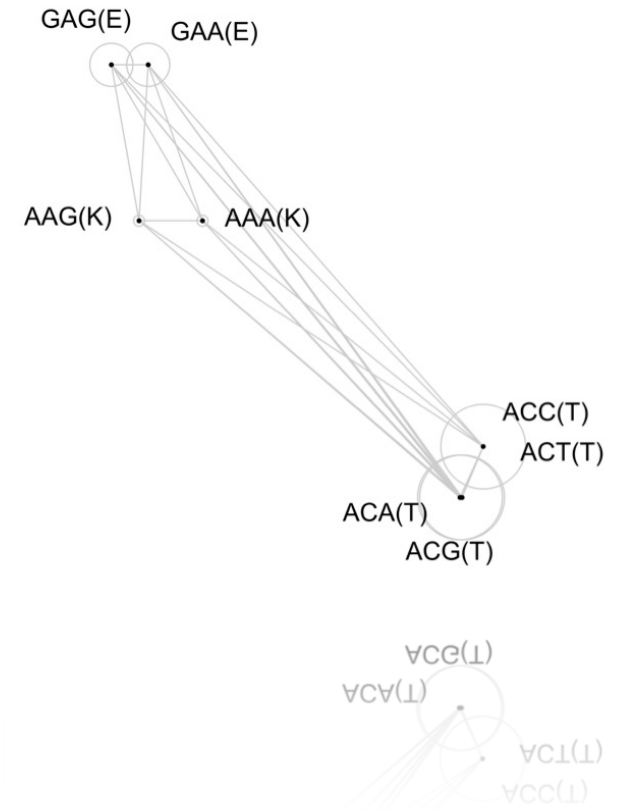
dN^h/dS^h depends on the current amino acid



conclusion: positive selection operates on a stationary fitness peak in the same way as when there is an adaptive peak shift

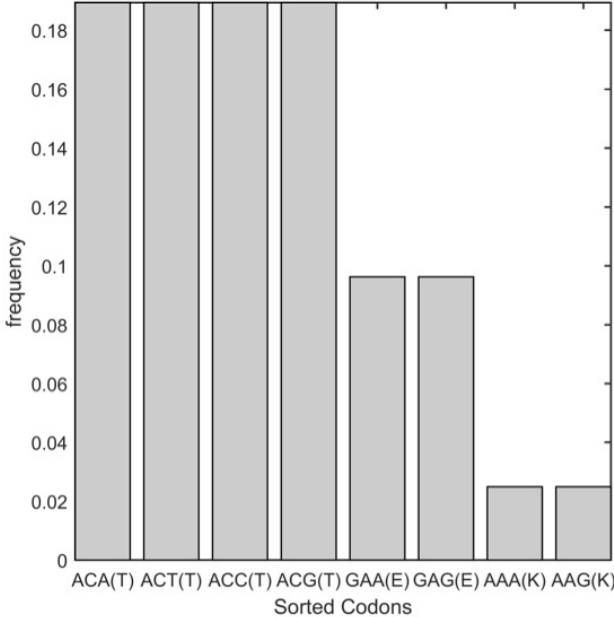
5. nearly-neutral theory and “heterotachy”

ρ: μεγάλη-μικρή μεταβολή ανα „μειοταχυλ”

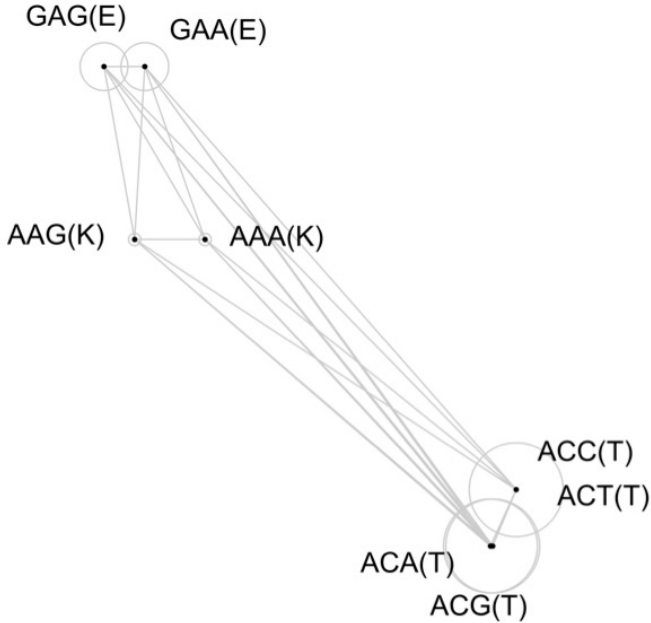


landscapes have unique structures

MutSel landscape



McCandlish landscape

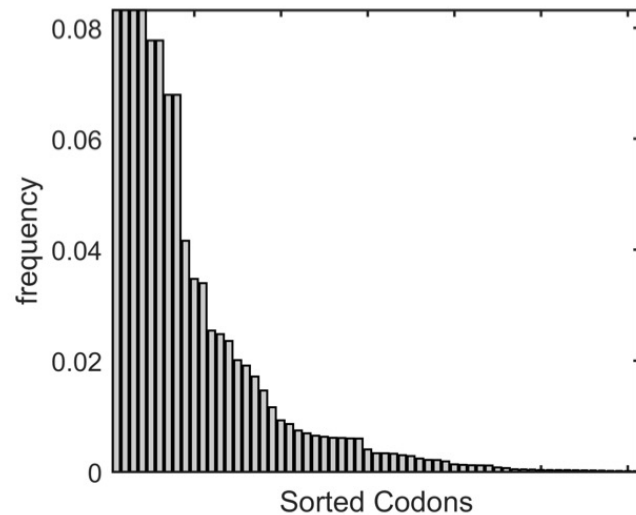


conclusion: A population can get to a sub-optimal codon (E) by drift and reside there for some time (b/c moving between T and E requires changes ≥ 2 codons).

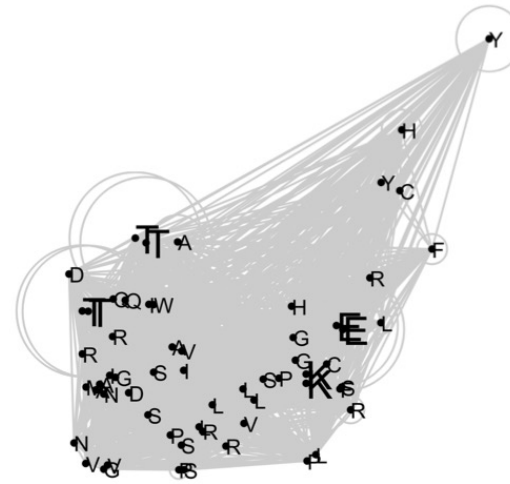
landscapes have unique structures

same site... 10x decrease in N (f^h have not changed!)

MutSel landscape



McCandlish landscape



Nearly-neutral theory:
**selection and drift
interact!**

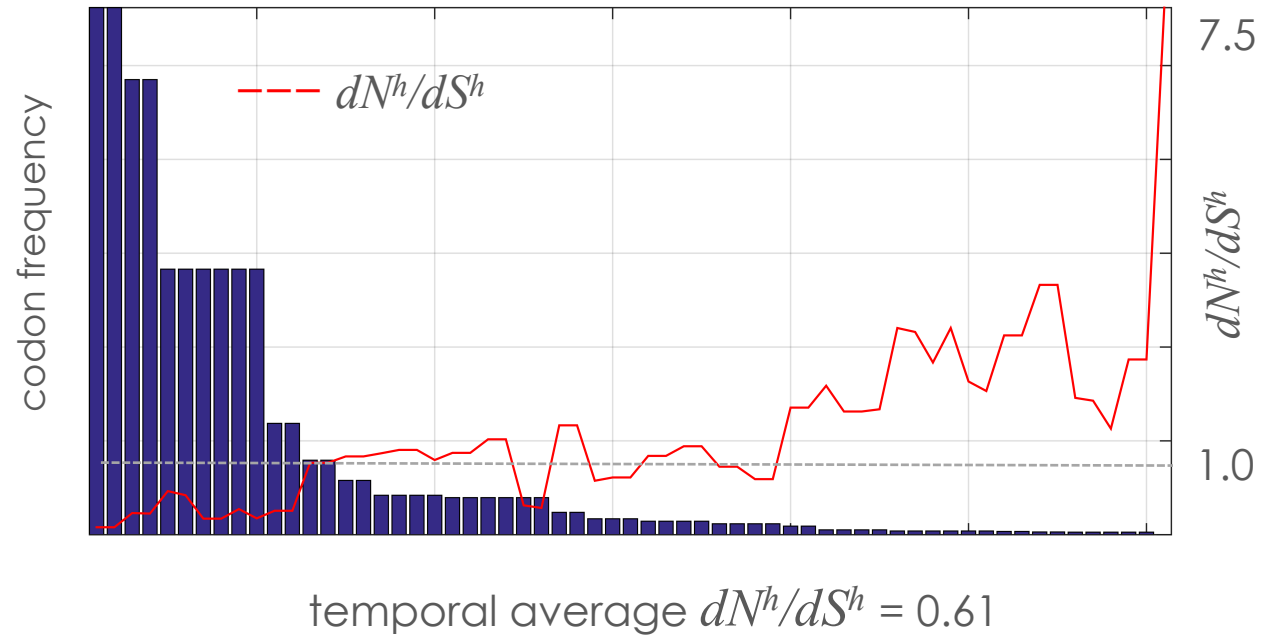
Rate of evolution
depends on
population size

conclusion: decreasing N changes:

- i. the “space” for shifting balance
- ii. mean dN/dS
- iii. equilibrium frequencies

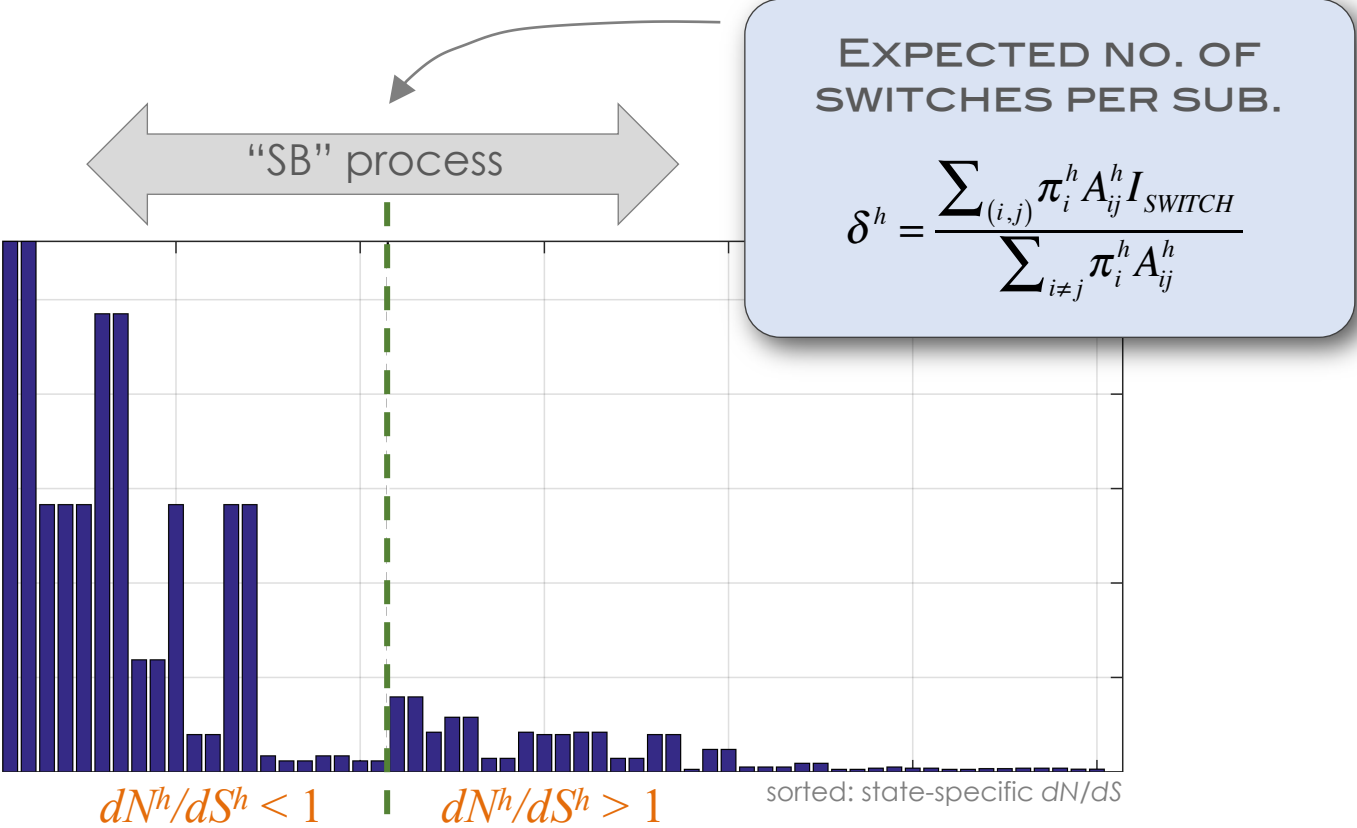
shifting balance: the MutSel landscape

dN^h/dS^h depends on the current amino acid



shifting balance: a mechanistic model

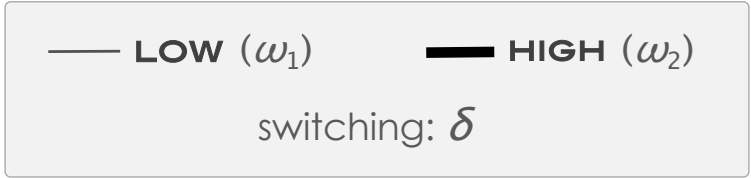
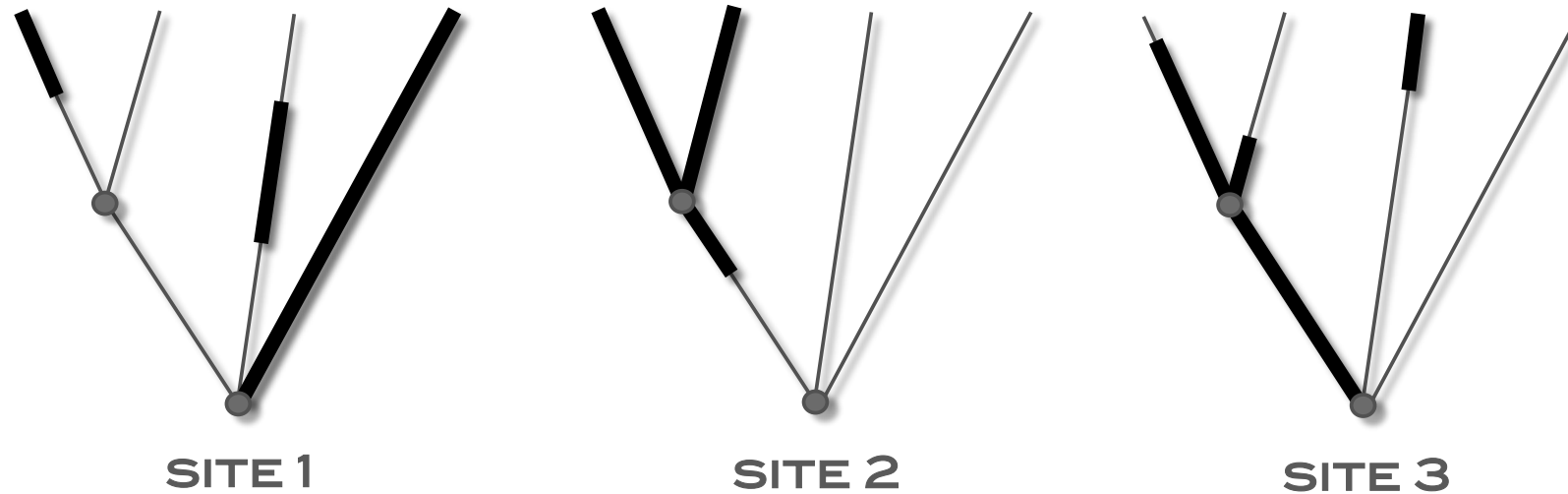
NOTE: Rate switching like this is called **"Heterotachy"**



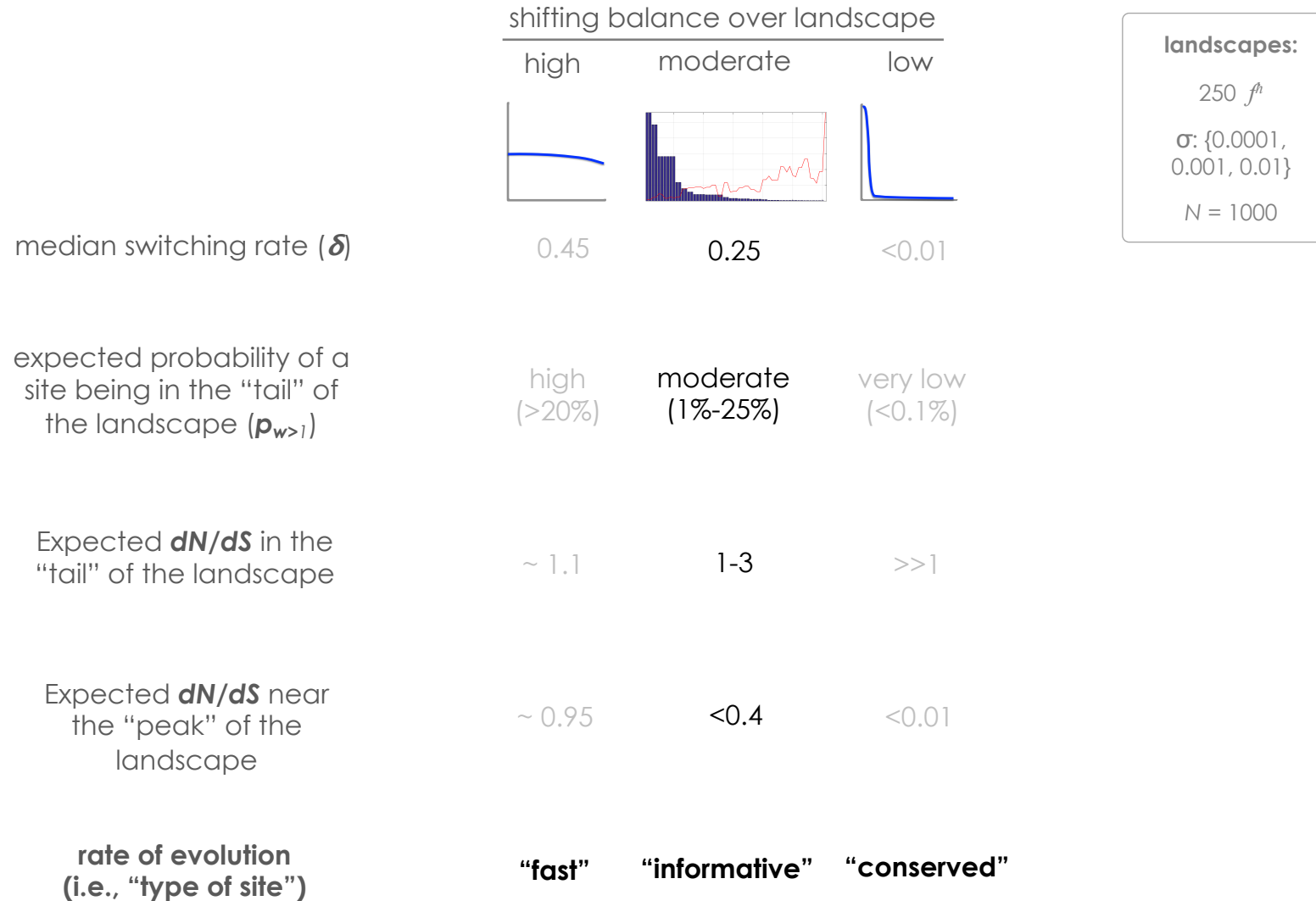
$$\omega^h < 1 = \frac{\sum_{i \in I_p^h} \frac{\pi_i^h}{p_1} A_{ij}^h I_N}{\sum_{i \in I_p^h} \frac{\pi_i^h}{p_1} \mu_{ij} I_N}$$

$$\omega^h > 1 = \frac{\sum_{i \in I_t^h} \frac{\pi_i^h}{p_2} A_{ij}^h I_N}{\sum_{i \in I_t^h} \frac{\pi_i^h}{p_2} \mu_{ij} I_N}$$

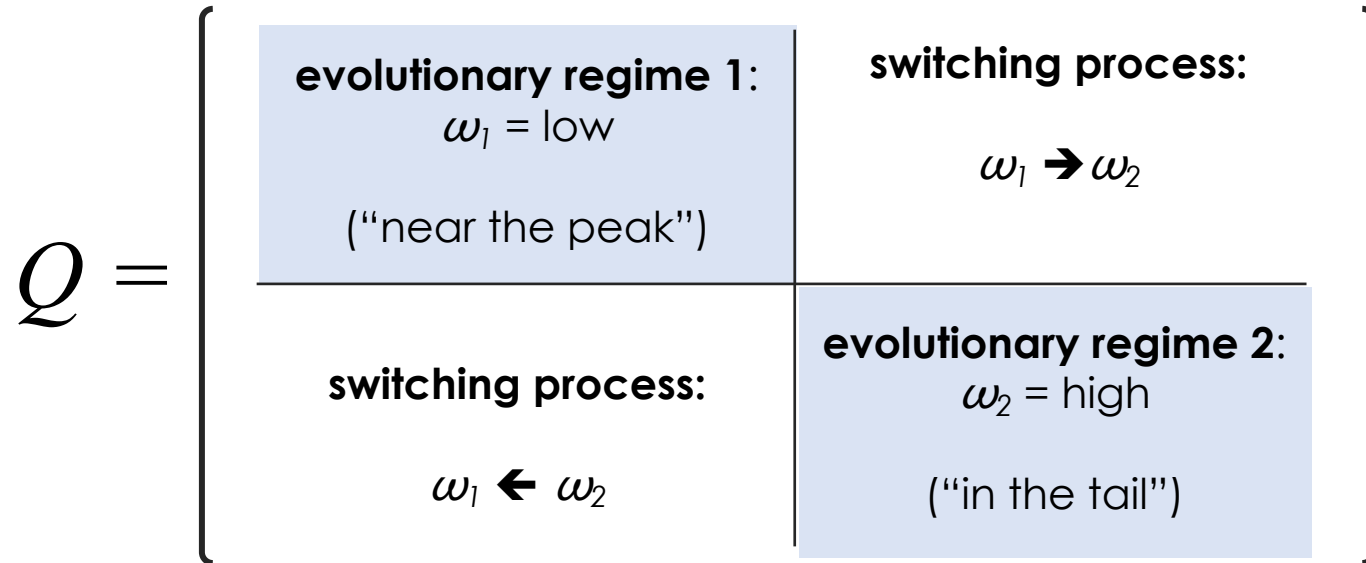
What does heterotachy “look” like on a tree?



shifting balance: a mechanistic model



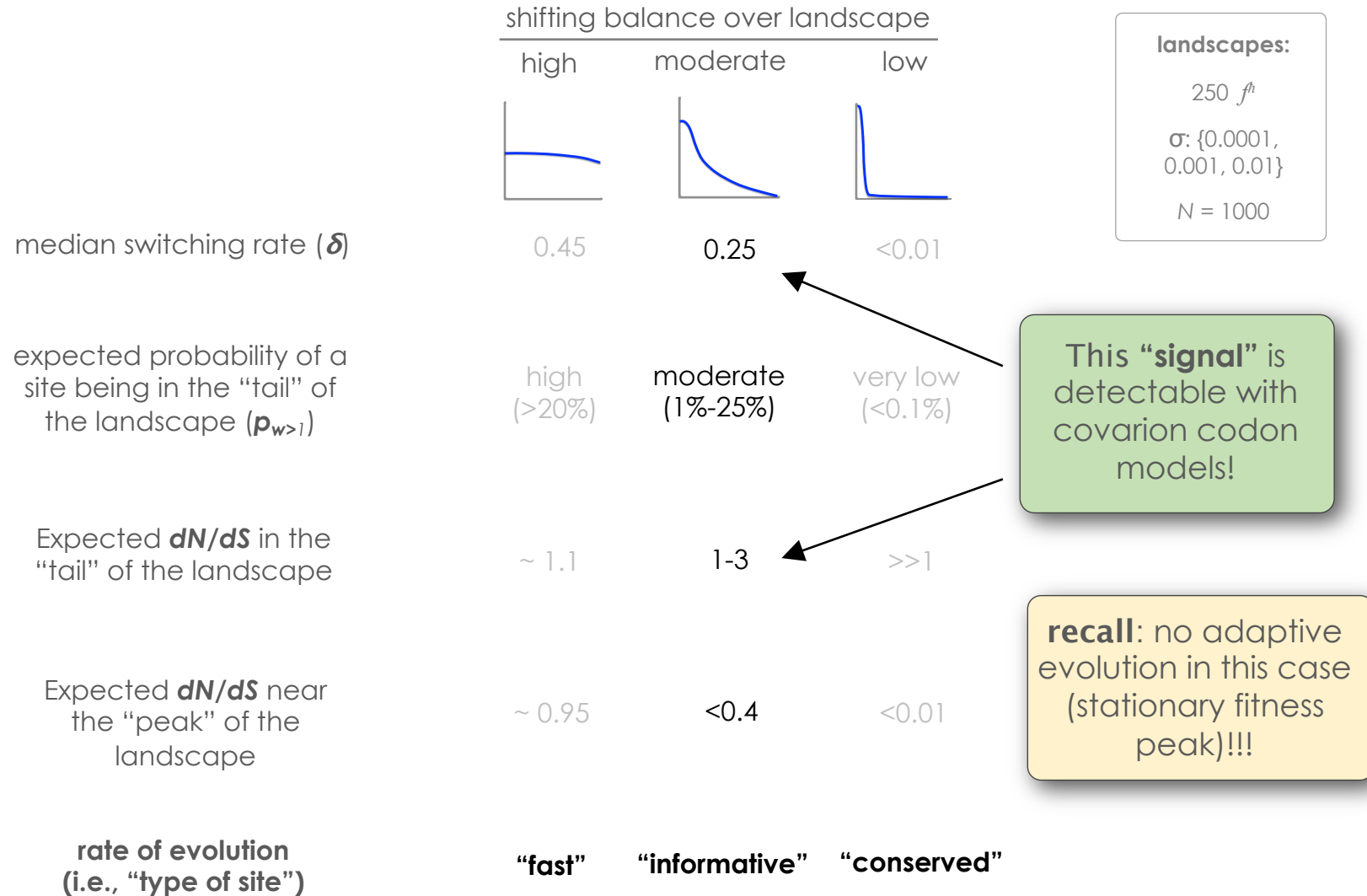
We can model heterotachy with a covarion-like model



[Guindon et al., (2004); Jones et al. (2016); Jones et al. (2018); Jones et al. 2019]

the covarion-like codon model can be **fit to real data**

shifting balance: a mechanistic model



6. some common types of codon models

6. some common types of codon models

“OMEGA MODELS”

$$Q_{ij} = \begin{cases} 0 & \text{if } i \text{ and } j \text{ differ by } > 1 \\ \pi_j & \text{for synonymous tv.} \\ \kappa\pi_j & \text{for synonymous ts.} \\ \omega\pi_j & \text{for non-synonymous tv.} \\ \omega\kappa\pi_j & \text{for non-synonymous ts.} \end{cases}$$

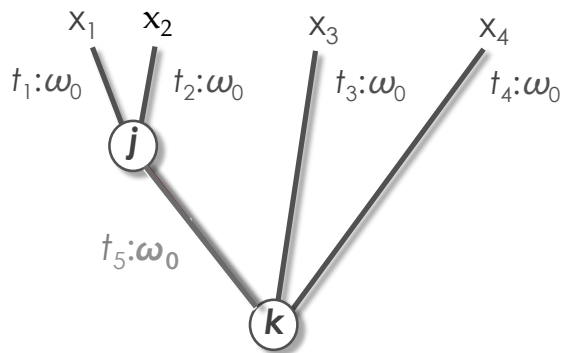
Goldman and Yang (1994)
Muse and Gaut (1994)

this codon model “**MO**”

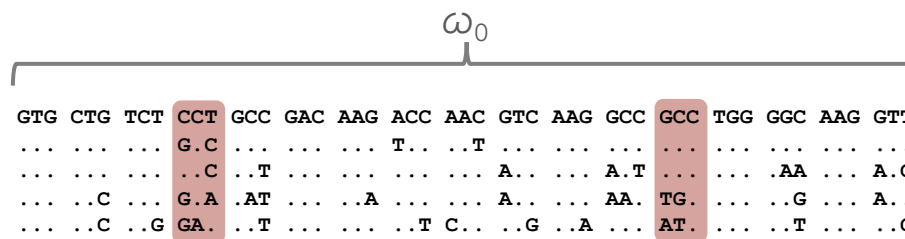
“OMEGA MODELS”

$$Q_{ij} = \begin{cases} 0 & \text{if } i \text{ and } j \text{ differ by } > 1 \\ \pi_j & \text{for synonymous tv.} \\ \kappa\pi_j & \text{for synonymous ts.} \\ \omega\pi_j & \text{for non-synonymous tv.} \\ \omega\kappa\pi_j & \text{for non-synonymous ts.} \end{cases}$$

Goldman and Yang (1994)
Muse and Gaut (1994)

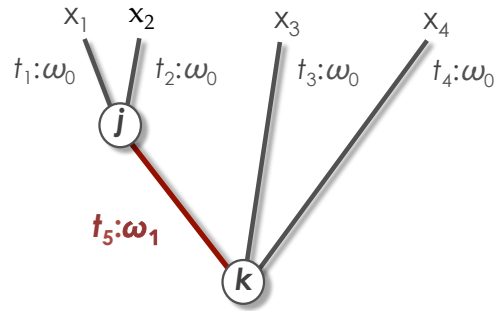


same ω
for all branches



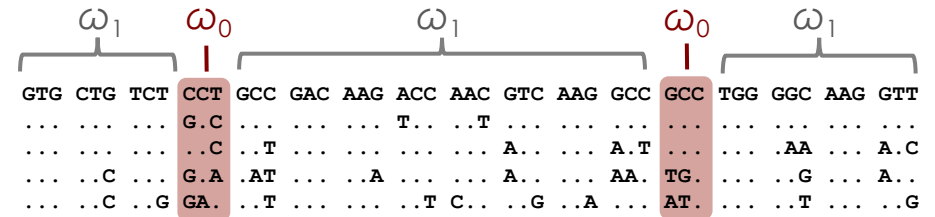
same ω
for all sites

two basic types of models...



branch models

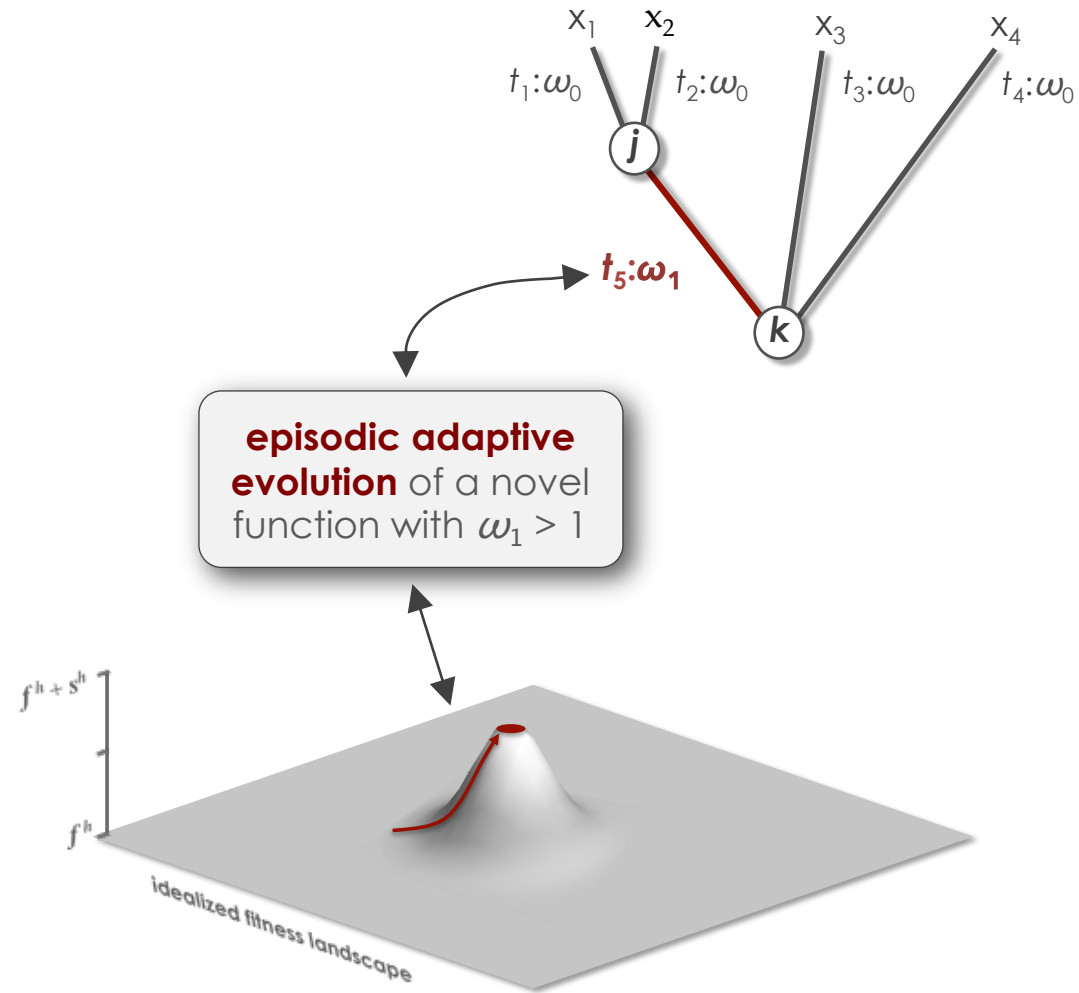
(ω varies among branches)



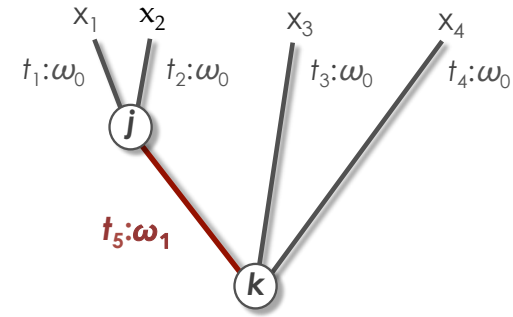
site models

(ω varies among sites)

interpretation of a branch model



branch models*



variation (ω) among branches:	approach
Yang, 1998	fixed effects
Bielawski and Yang, 2003	fixed effects
Seo et al. 2004	auto-correlated rates
Kosakovsky Pond and Frost, 2005	genetic algorithm
Dutheil et al. 2012	clustering algorithm

* these methods can be useful when selection pressure is strongly **episodic and functional change is substantial**

site models*

```

GTG CTG TCT CCT GCC GAC AAG ACC AAC GTC AAG GCC GCC TGG GGC AAG GTT GGC GCG CAC
... .. G.C ... .. T.. ..T ... .. ..GC A..
... ..C ..T ... .. ..A.. ..A.T ... ..AA ... A.C ... AGC ...
... ..C ... G.A .AT ... ..A ... ..A.. ..AA. TG. ... ..G ... A.. ..T .GC ..T
... ..C ..G GA. ..T ... ..T C.. ..G ..A ... AT. ... ..T ... ..G ..A .GC ...

```

variation (ω) among sites:	approach
Yang and Swanson, 2002	fixed effects (ML)
Bao, Gu and Bielawski, 2006	fixed effects (ML)
Massingham and Goldman, 2005	site wise (LRT)
Kosakovsky Pond and Frost, 2005	site wise (LRT)
Nielsen and Yang, 1998	mixture model (ML)
Kosakovsky Pond, Frost and Muse, 2005	mixture model (ML)
Huelsenbeck and Dyer, 2004; Huelsenbeck et al. 2006	mixture (Bayesian)
Rubenstein et al. 2011	mixture model (ML)
Bao, Gu, Dunn and Bielawski 2008 & 2011	mixture (LiBaC/MBC)
Murell et al. 2013	mixture (Bayesian)

This is **NOT** a comprehensive list!

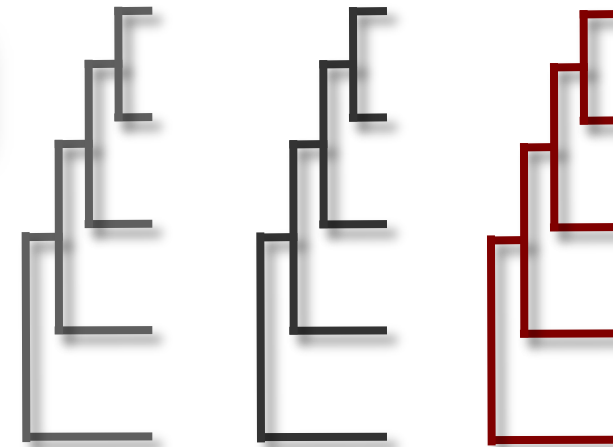
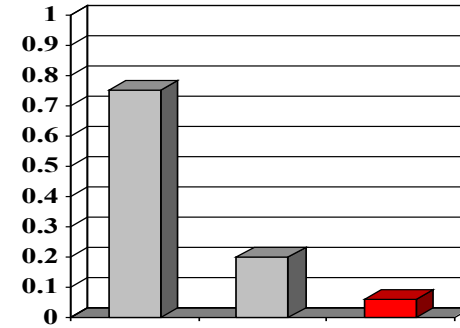
- useful when at some sites evolve under **diversifying selection** pressure over long periods of time

site models: discrete mixture model (**M3**)

MIXTURE-MODEL LIKELIHOOD

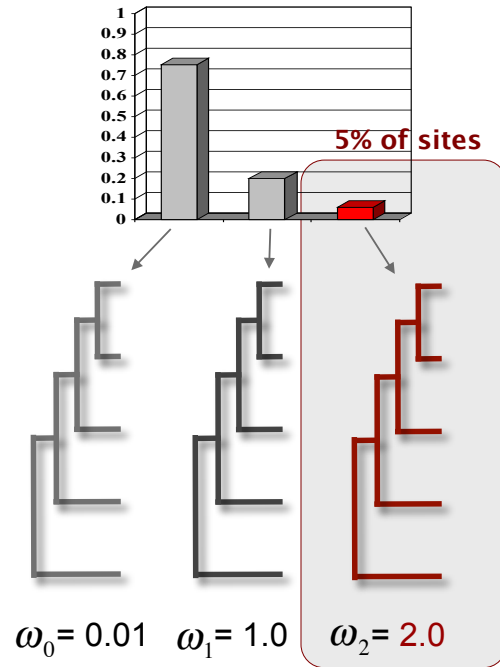
$$P(\mathbf{x}_h) = \sum_{i=0}^{K-1} p_i P(\mathbf{x}_h | \omega_i)$$

conditional likelihood
calculation (see part 1)

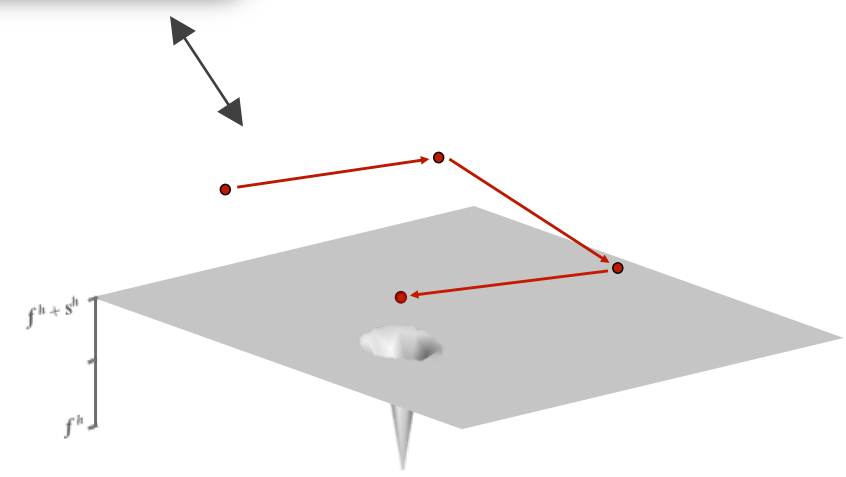


$\omega_0 = 0.01$ $\omega_1 = 1.0$ $\omega_2 = 2.0$

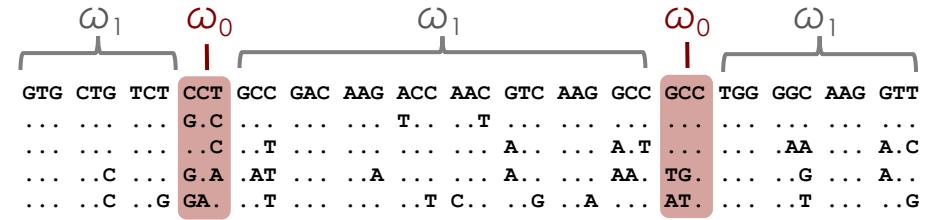
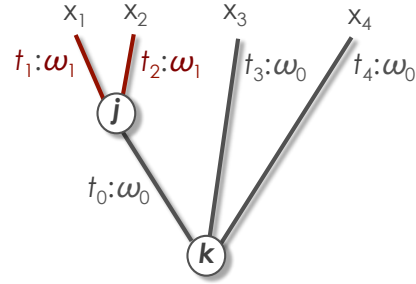
interpretation of a sites-model



**diversifying selection
(frequency dependent)**
at 5% of sites with
 $\omega_2 = 2$



models for variation among branches & sites



branch models
(ω varies among
branches)

site models
(ω varies among sites)

branch-site models
(combines the features of above models)

models for variation among branches & sites

This is **NOT** a comprehensive list!

variation (ω) among branches & sites:	approach
Yang and Nielsen, 2002	fixed+mixture (ML)
Forsberg and Christiansen, 2003	fixed+mixture (ML)
Bielawski and Yang, 2004	fixed+mixture (ML)
Giundon et al., 2004	covarion-like (ML)
Zhang et al. 2005	fixed+mixture (ML)
Kosakovsky Pond et al. 2011, 2012	full mixture (ML)
Jones et al., 2016, 2018, 2020	mix-covarion-like (ML)

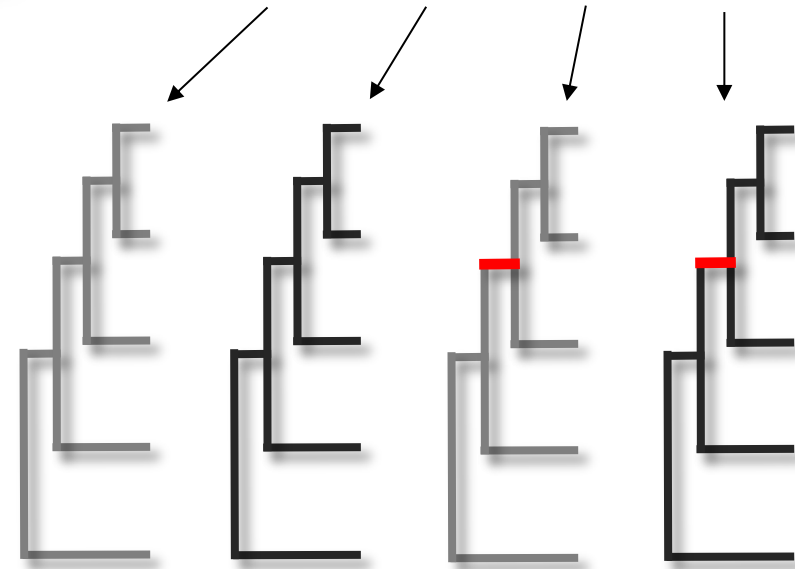
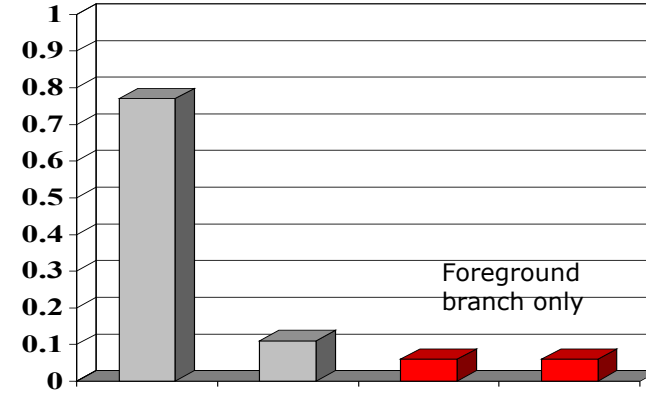
* *these methods can be useful when selection **pressures change over time at just a fraction of sites***

* *it can be a challenge to apply these methods properly*

branch-site "Model B"

MIXTURE-MODEL LIKELIHOOD

$$P(\mathbf{x}_h) = \sum_{i=0}^{K-1} p_i P(\mathbf{x}_h | \omega_i)$$



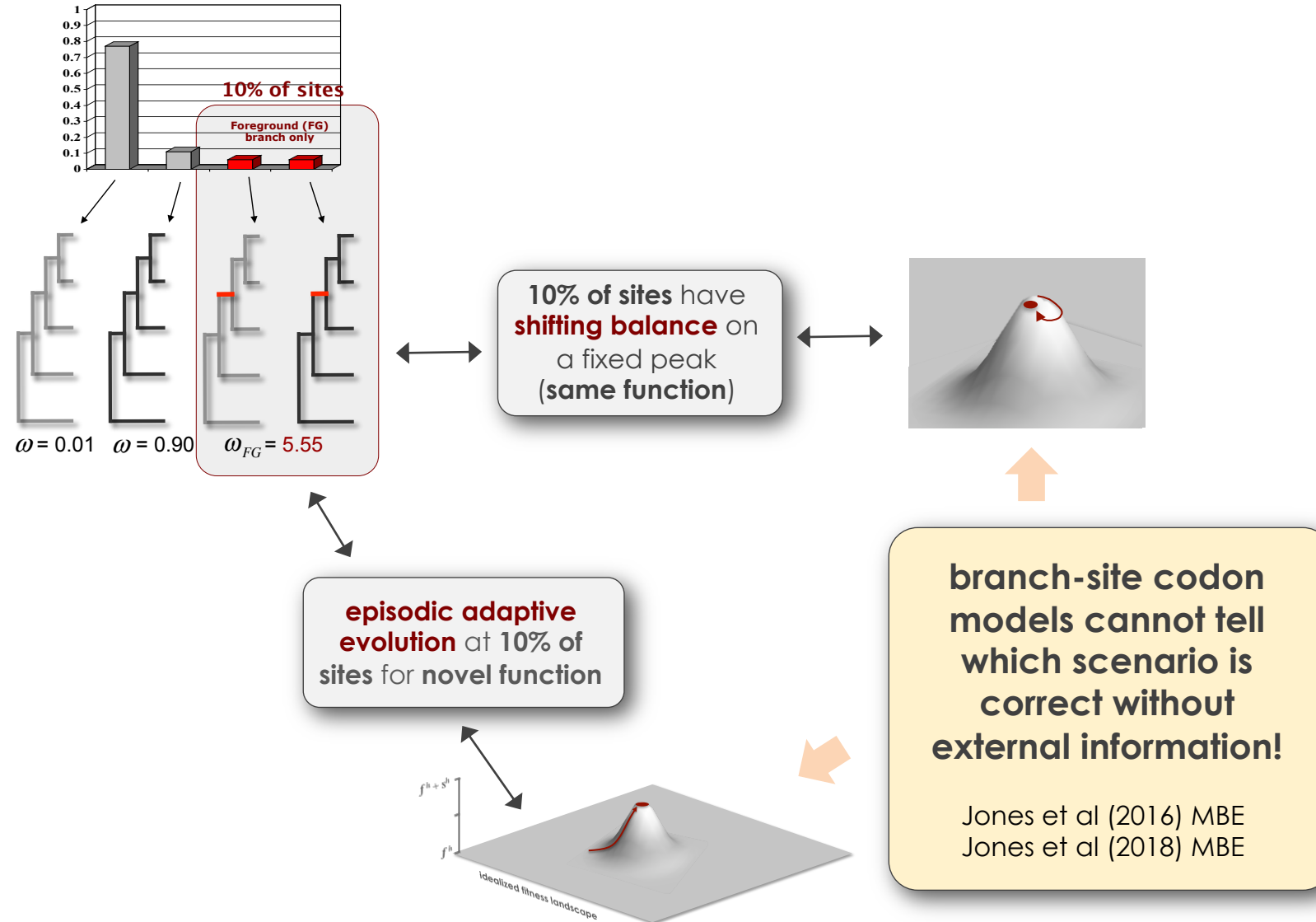
$\omega = 0.01$

$\omega = 0.90$

$\omega = 5.55$

ω for background branches are from site-classes 1 and 2 (0.01 or 0.90)

two scenarios can yield branch-sites with $d_N/d_S > 1$



7. “bells –n– whistles” ...

codon models + “other processes”

codon models + “other processes”

From codon below:	to codon below:						
	TTT (Phe)	TTC (Phe)	TTA (Leu)	TTG (Leu)	CTT (Leu)	CTC (Leu)	GGG (Gly)
TTT (Phe)	—	K π TTC	ω π TTA	ω π TTG	ω K π TTT	0	0
TTC (Phe)	K π TTT	—	ω π TTA	ω π TTG	0	ω K π CTC	0
TTA (Leu)	ω π TTT	ω π TTC	—	—	0	0	0
TTG (Leu)	ω π TTT	ω π TTC	K π TTA	—	0	0	0
CTT (Leu)	ω K π TTT	0	0	0	—	K π CTC	0
CTC (Leu)	0	ω K π TTC	0	0	K π TTT	—	0
GGG (Gly)	0	0	0	0	0	0	—

ccc (eA)	0	0	0	0	0	0	0	—
cic (ren)	0	ω π LLC	0	0	π K π LLL	—	—	0
cII (ren)	ω π LLL	0	0	0	—	π K π LLC	—	0
llc (ren)	ω π LLL	ω π LLC	π K π LLV	—	—	—	—	0
llv (ren)	ω π LLL	ω π LLC	—	—	0	0	—	0
llc (h μ e)	π K π LLL	—	ω π LLV	ω π LLC	0	ω π K π LLC	—	0
lll (h μ e)	—	π K π LLC	ω π LLV	ω π LLC	ω π LLL	—	—	0

“bells –n– whistles”... some general categories

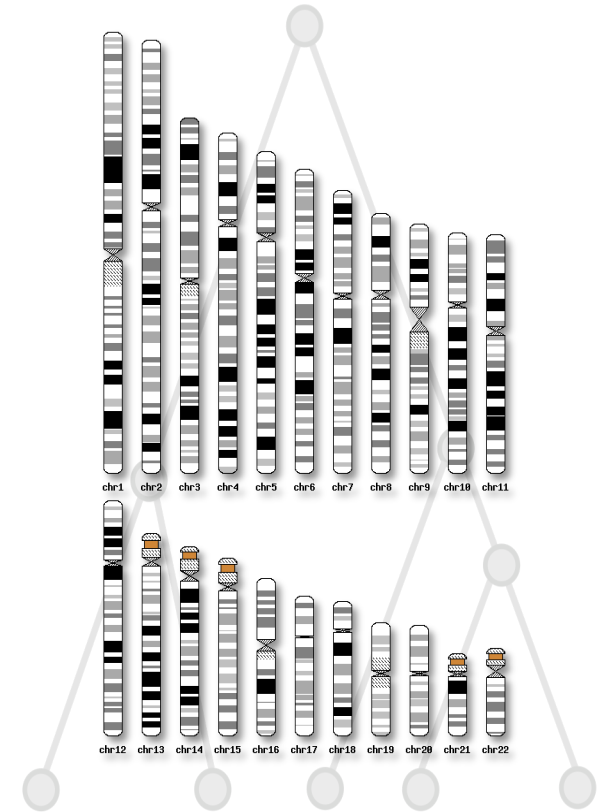
1. alternative models of codon frequencies *(the π 's parameters are important)*
2. GTR process at DNA-level *(this is NOT a mutational process)*
3. among-site synonymous rate (d_s) variation *(phenomenologically important)*
4. double & triple nucleotide changes *(confounded with heterotachy)*
5. amino acid exchangeabilities *(confounded with codon frequencies via fitness)*
6. multi-process variation among sites *(do we really want this much complexity?)*
7. multi-pattern (tree) variation among sites *(this can be important)*

From codon below:	to codon below:						GGG (Gly)
	TTT (Phe)	TTC (Phe)	TTA (Leu)	TTG (Leu)	CTT (Leu)	CTC (Leu)	
TTT (Phe)	—	κπTTC	ωπTTA	ωπTTG	ωκπTTT	0	0
TTC (Phe)	κπTTT	—	ωπTTA	ωπTTG	0	ωκπCTC	0
TTA (Leu)	ωπTTT	ωπTTC	—	—	0	0	0
TTG (Leu)	ωπTTT	ωπTTC	κπTTA	—	0	0	0
CTT (Leu)	ωκπTTT	0	0	0	—	κπCTC	0
CTC (Leu)	0	ωκπTTC	0	0	κπTTT	—	0
GGG (Gly)	0	0	0	0	0	0	—

Is there some way that we can become less-dependent on all this phenomenological complexity which can obscure evolutionary implications?

8. Phenotype-Genotype codon models

8. Phenotype-Genotype codon models



Opinion

Phylogenetics is the New Genetics
(for Most of Biodiversity)Stacey D. Smith,^{1,6,*} Matthew W. Pennell,² Casey W. Dunn,³ and Scott V. Edwards^{4,5}

Despite substantial progress in understanding the genetic basis for differences in morphology, physiology, and behavior, many phenotypes of interest are difficult to study with traditional genetic approaches because their origin traces to deep nodes in the tree of life. Moreover, many species are not amenable to either large-scale sampling or laboratory crosses. We argue that phylogenetic methods and theory provide tremendous power to identify the functional genetic variation underlying trait evolution. We anticipate that existing statistical comparative approaches will be more commonly applied to studying the genetic basis for phenotypic evolution as whole genomes continue to populate the tree of life. Nevertheless, new methods and approaches will be needed to fully capitalize on the power of clade-scale genomic datasets.

Highlights

Genome sequencing is rapidly spreading beyond model organisms, opening the door to comparative studies that can reveal the genetic basis for phenotypic variation across species. Nevertheless, statistical comparative methods have not been frequently applied to these data.

New phylogenetic methods have been developed with the explicit goal of linking genes and even specific mutations to species differences (PhyloG2P[®]). Applications of these methods show great promise for uncovering new sources of functional variation and tackling traits beyond the reach of traditional genetic approaches.

Parallel advances in statistical comparative methods present new avenues for expanding the phylogenetic toolkit and creating tailored approaches for mapping genotype to phenotype.

Most of Biodiversity Is Beyond the Reach of Classical Genetics

One of the fundamental goals of biology is to connect variation across genomes to differences in phenotypes. With advances in sequencing and molecular genetic techniques, this area of biology has blossomed in recent years, revealing the genetic basis for traits ranging from floral scent [1] to sociality [2] to herbivory [3]. At the same time, statistical methods for analyzing these data have also proliferated [4–6]. At their core, however, all classical and population genetic methods for **genotype-to-phenotype mapping** (see [Glossary](#)) work by associating genetic variation with differences in the trait of interest. Thus, they require a population with segregating phenotypic variation, which could be produced artificially through crosses or mutagenesis or could occur naturally, such as in polymorphic species or hybrid zones between species. As with any statistical approach, association methods [e.g., **genome-wide association studies (GWASs)**] have significant challenges and pitfalls [6,7]. Still, the loci uncovered by association mapping and similar methods have often been validated in subsequent functional studies [8,9], confirming their ability to identify regions of the genome that contribute to phenotypic differences.

Despite the success of this population genetic program for genotype–phenotype mapping, it presents significant limitations for understanding the genetic basis of phenotypes for most of biodiversity. First, many species cannot be propagated artificially or sampled in the wild at the scale needed for association mapping (usually hundreds of individuals, depending on the trait of interest). Second, and more importantly, many traits of interest are not found segregating in nature nor can different species with contrasting phenotypes be crossed. For example, mammals with and without pouches cannot be crossed, precluding the creation of a mapping population segregating for pouches. As a consequence, our understanding of the genetic basis for phenotypic diversity is concentrated around a narrow range of species and traits – often those that vary in model organisms amenable to genetic studies. Although loci discovered through genetic studies of model species often later help to explain variation at deeper phylogenetic levels (i.e., across species [10,11]), we wonder what we might discover if this research program were inverted (Figure 1). We suggest, and recent studies confirm, that beginning from a phylogenetic

¹Department of Ecology and Evolutionary Biology, University of Colorado, Boulder, CO 80309, USA
²Department of Zoology and Biodiversity Research Centre, University of British Columbia, Vancouver, BC V6T 1Z4, Canada
³Department of Ecology and Evolutionary Biology, Yale University, New Haven, CT 06520, USA
⁴Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA 02138, USA
⁵Museum of Comparative Zoology, Harvard University, Cambridge, MA 02138, USA
⁶<http://www.colorado.edu/smithlab>

*Correspondence: Stacey.D.Smith@colorado.edu (S.D. Smith).
[®]Twitter: @iochromaland (S.D. Smith).



Phenotype-Genotype codon models?

Phenotype only models:

Cornwell, W. and Nakagawa, S. 2017. Phylogenetic comparative methods. *Curr.Biol.*, 27: 327-338.

phenotype models

Phenotype + Genotype models:

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.

+ DNA model

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.

+ codon model

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.

+ DNA model

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.

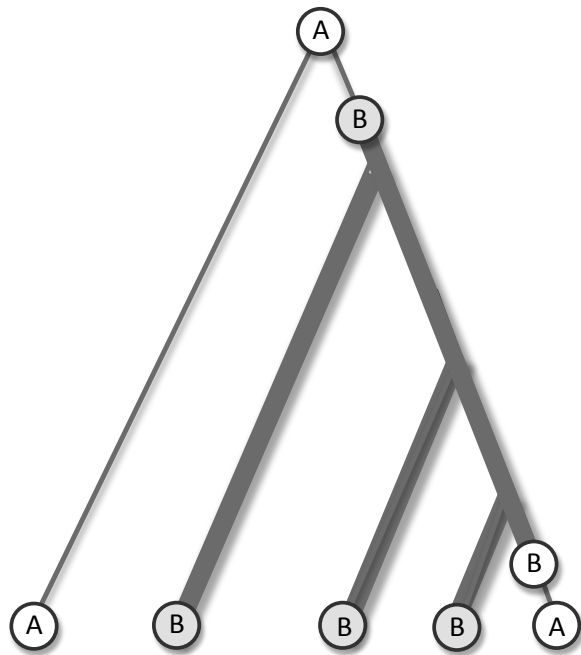
+ DNA model

Jones, C. T., Youssef, N., Susko, E., & Bielawski, J. P. (2020). *A Phenotype-Genotype Codon Model for Detecting Adaptive Evolution*. *Systematic Biology*, 69(4), 722-738.

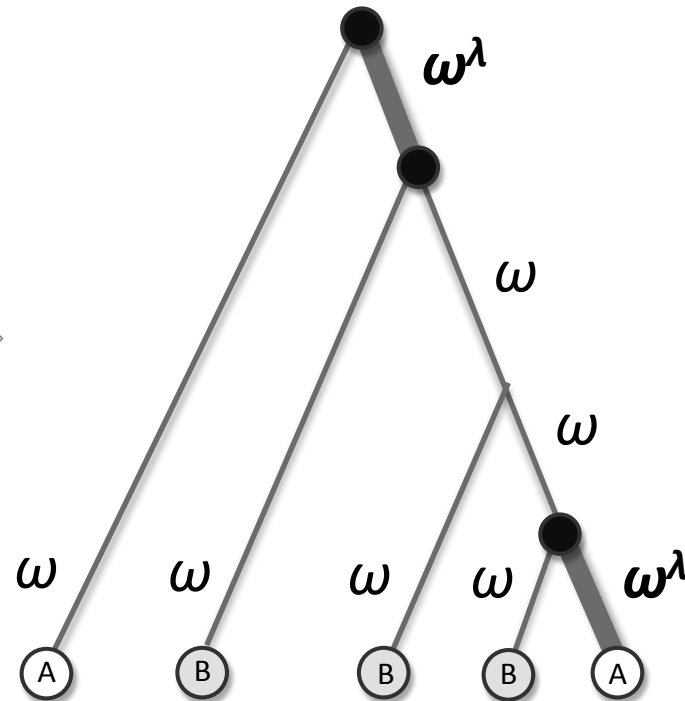
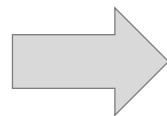
+ codon model

Halabi, K., Karin, E. L., Guéguen, L., & Mayrose, I. (2021). *A codon model for associating phenotypic traits with altered selective patterns of sequence evolution*. *Systematic Biology*, 70(3), 608-622.

+ codon model



phenotype mapping 1 (of many)



gene evolution

Phenotype only models:

Cornwell, W. and Nakagawa, S. 2017. Phylogenetic comparative methods. *Curr.Biol.*, 27: 327-338.

phenotype models

Phenotype + Genotype models:

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.

+ DNA model

Lartillot, N. and
substitution ra

Detect adaptive molecular evolution

+ codon model

O'Connor, T. D
association an
Evolutionary B

(possibly without $d_N/d_S > 1$)

+ DNA model

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.

+ DNA model

Jones, C. T., Youssef, N., Susko, E., & Bielawski, J. P. (2020). *A Phenotype-Genotype Codon Model for Detecting Adaptive Evolution*. *Systematic biology*, 69(4), 722-738.

+ codon model

Halabi, K., Karin, E. L., Guéguen, L., & Mayrose, I. (2021). *A codon model for associating phenotypic traits with altered selective patterns of sequence evolution*. *Systematic Biology*, 70(3), 608-622.

+ codon model

9. model-based inference

9. model-based inference



You will get “*the basics*”
in the evening
PAML-Lab