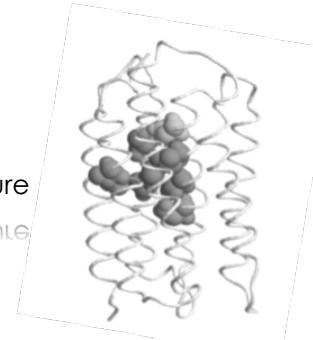


part 3: analysis of natural selection pressure

part 3: analysis of natural selection pressure



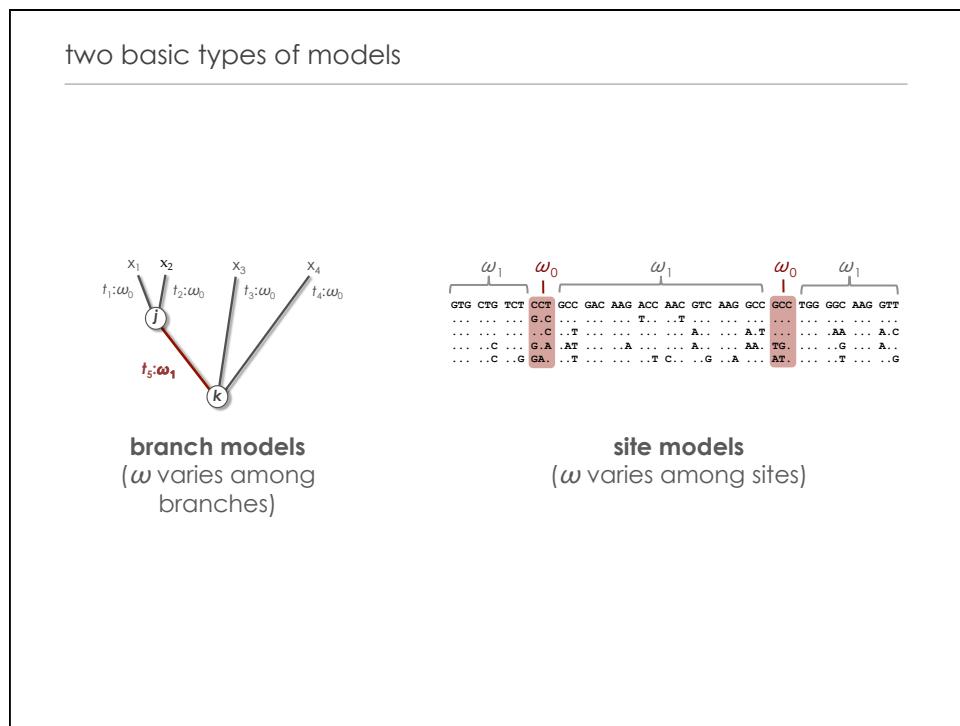
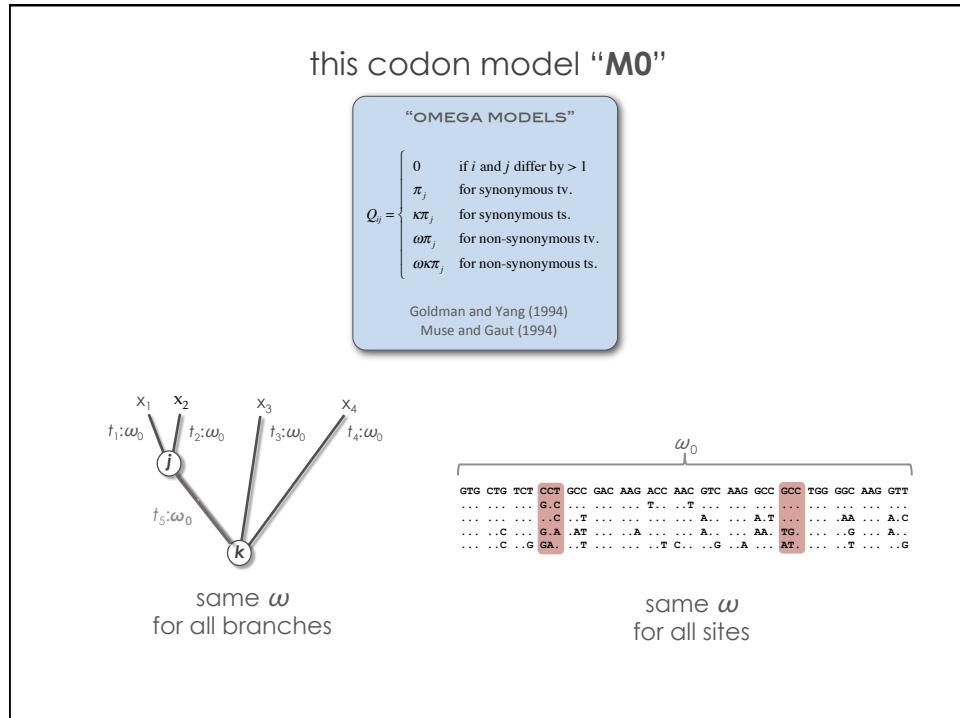
types of codon models

types of codon models

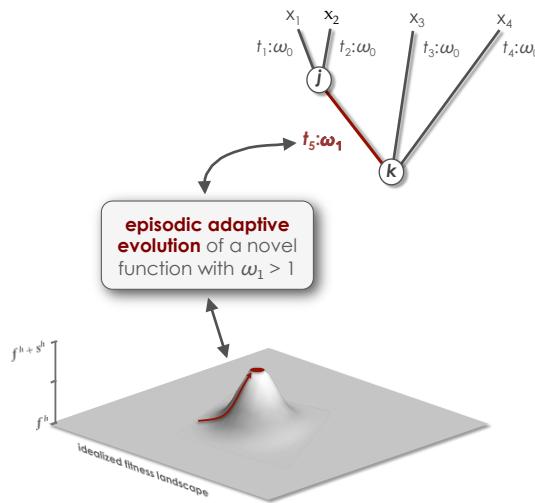
"OMEGA MODELS"

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

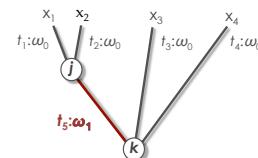
Goldman and Yang (1994)
Muse and Gaut (1994)



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**

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 .. A.. A.T .. AA .. A.C .. AGC ..
C .. G.A AT .. A .. AA TG .. G .. A .. T GC .. 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)

- useful when at some sites evolve under **diversifying selection** pressure over long periods of time
- this is not a comprehensive list

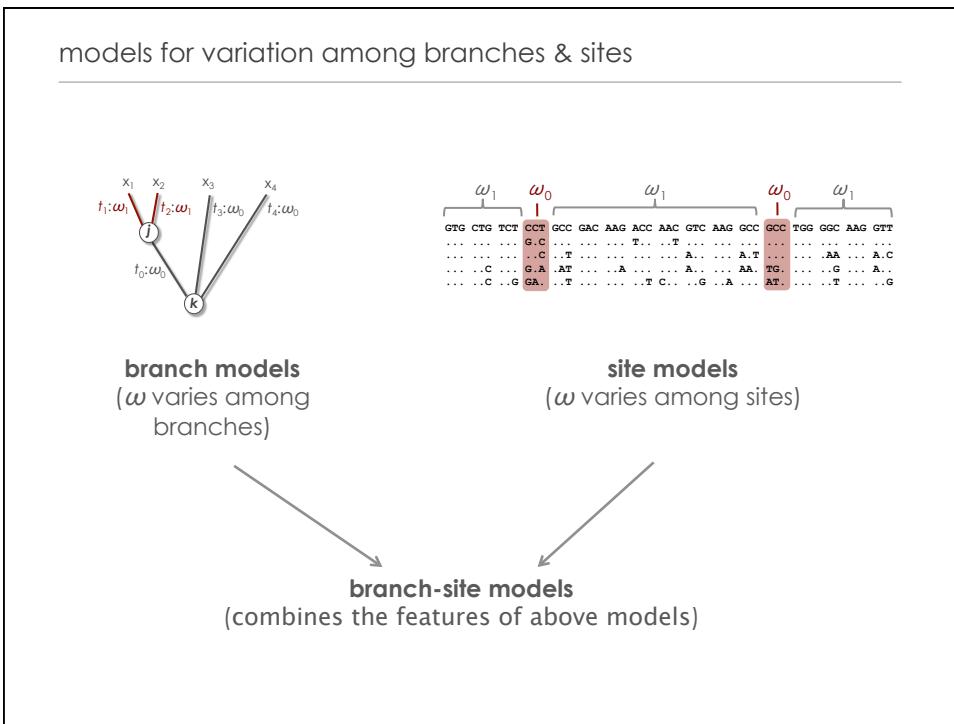
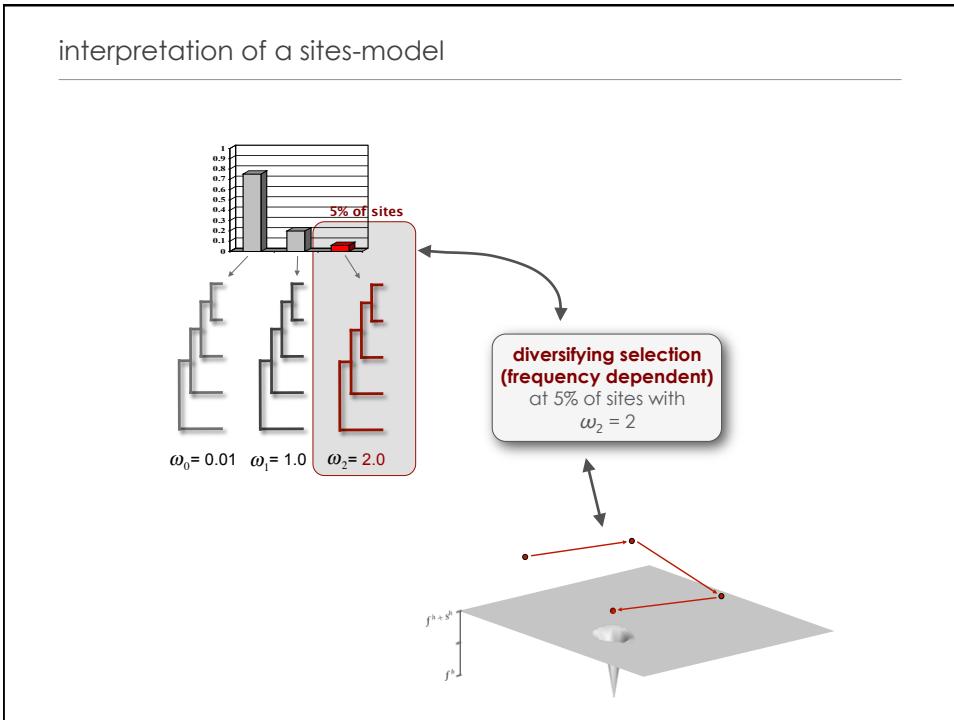
site models: discrete model (**M3**)

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)

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



models for variation among branches & sites

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	covarion-like (ML)

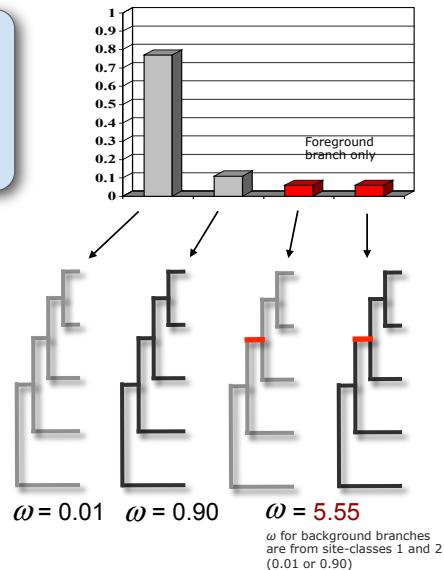
* these methods can be useful when selection **pressures change over time at just a fraction of sites**

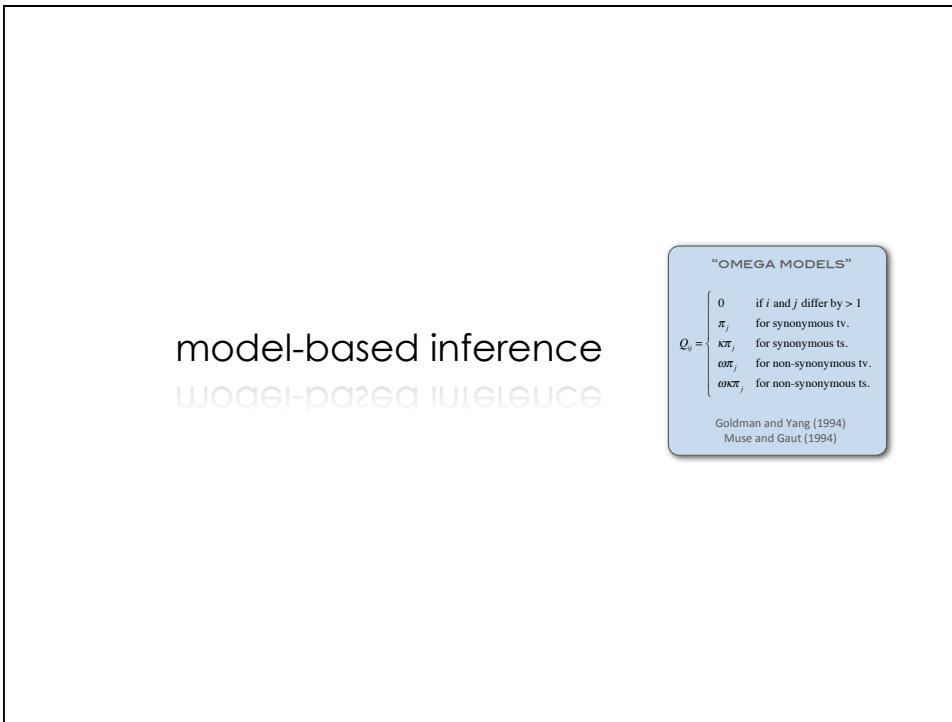
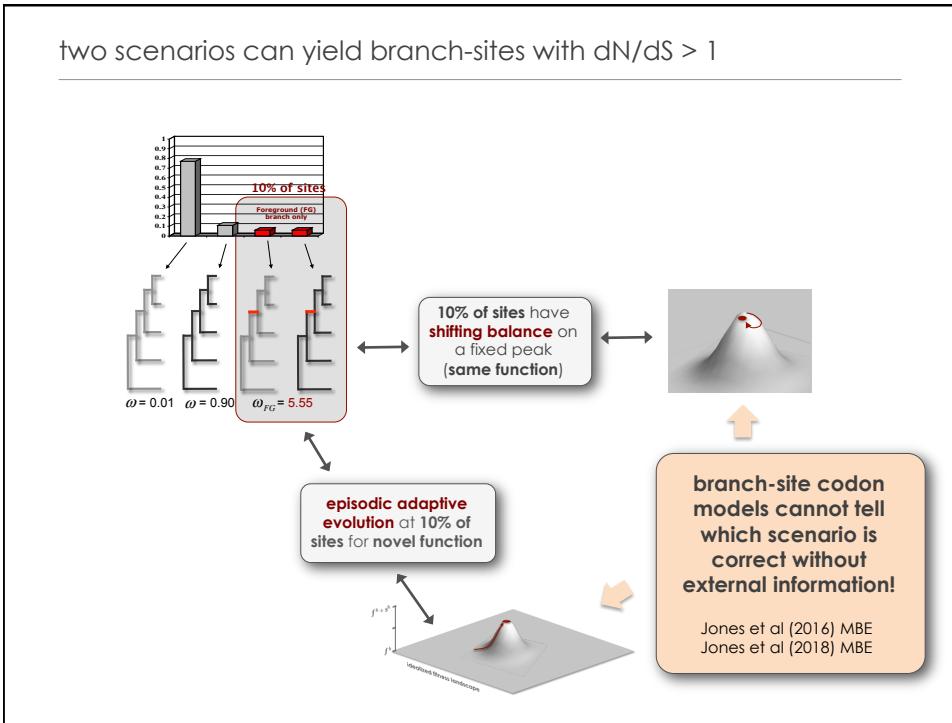
* it can be a challenge to apply these methods properly (**more about this later**)

branch-site “Model B”

MIXTURE-MODEL LIKELIHOOD

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





model based inference

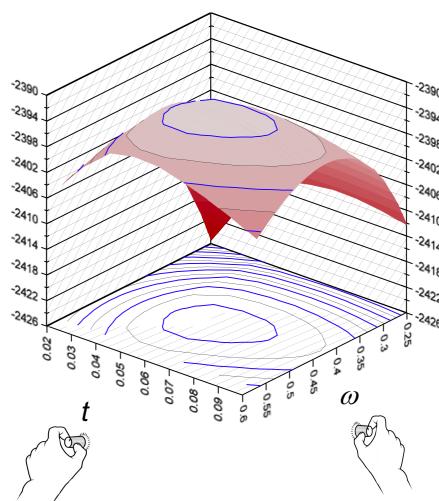
3 analytical tasks

task 1. parameter estimation (e.g., ω) ←

task 2. hypothesis testing

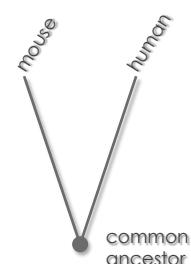
task 3. make predictions (e.g., sites having $\omega > 1$)

task 1: parameter estimation



Parameters: t and ω

Gene: acetylcholine α receptor



$\ln L = -2399$

task 2: statistical significance

task 1. parameter estimation (e.g., ω) ✓

task 2. hypothesis testing ← LRT

task 3. prediction / site identification

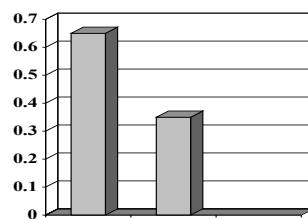
task 2: likelihood ratio test for positive selection

H_0 : variable selective pressure but NO positive selection (M1)

H_1 : variable selective pressure with positive selection (M2)

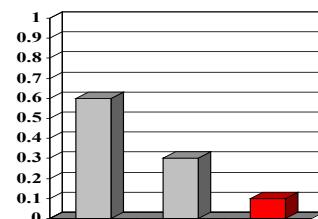
Compare $2\Delta I = 2(I_1 - I_0)$ with a χ^2 distribution

Model 1a (**M1a**)



$$\hat{\omega} = 0.5 \quad (\omega = 1)$$

Model 2a (**M2a**)



$$\hat{\omega} = 0.5 \quad (\omega = 1) \quad \hat{\omega} = 3.25$$

task 3: identify the selected sites

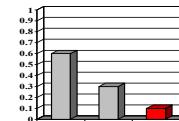
task 1. parameter estimation (e.g., ω) ✓

task 2. hypothesis testing ✓

task 3. prediction / site identification ← **Bayes' rule**

task 3: which sites have $dN/dS > 1$

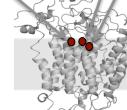
model:
10% have $\omega > 1$



Bayes' rule:
site 4, 12 & 13

GTC CTG TCT CCT GCC GAC RAG ACC AAC GTC AAG
...G.CTA ... A.T ... AA TG ...
... .C ... G.A AT ... A ... A ... AA AT ...
... C ... G AT ... TT CG ... A

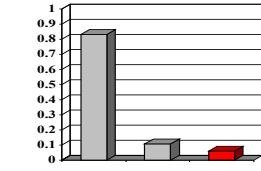
structure:
sites are in contact



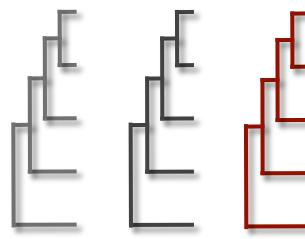
review the mixture likelihood (model **M3**)

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

↑ ↑ ↑
Total probability Prior Likelihood



$$\begin{aligned}\omega_0 &= 0.03 & \omega_1 &= 0.40 & \omega_2 &= 14.1 \\ p_0 &= 0.85 & p_1 &= 0.10 & p_2 &= 0.05\end{aligned}$$



Bayes' rule for identifying selected sites

- Site class 0: $\omega_0 = .03$, 85% of codon sites
- Site class 1: $\omega_1 = .40$, 10% of codon sites
- Site class 2: $\omega_2 = 14$, 05% of codon sites

Prior probability of hypothesis (ω_2)

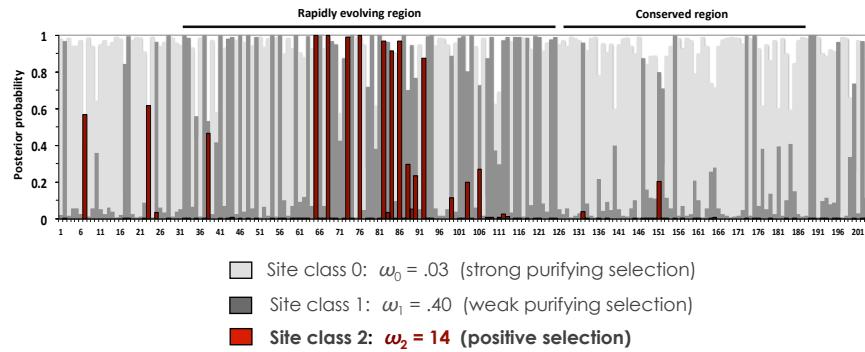
Likelihood of hypothesis (ω_2)

$$P(\omega_2 | x_h) = \frac{P(\omega_2) P(x_h | \omega_2)}{\sum_{i=0}^{K-1} P(\omega_i) P(x_h | \omega_i)}$$

Posterior probability of hypothesis (ω_2)

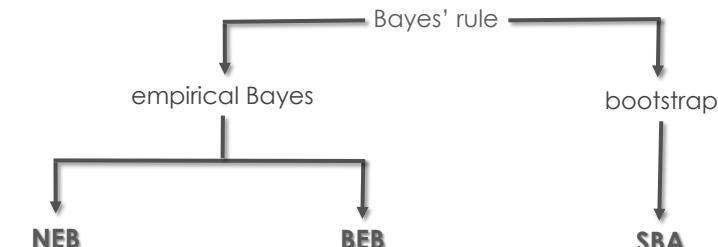
Marginal probability (Total probability) of the data

task 3: Bayes rule for which sites have $dN/dS > 1$



NOTE: The posterior probability should NOT be interpreted as a "P-value"; it can be interpreted as a measure of relative support, although there is rarely any attempt at "calibration".

task 3: Bayes rule for which sites have $dN/dS > 1$



Naive Empirical Bayes

- Nielsen and Yang, 1998
- assumes no MLE errors

Bayes Empirical Bayes

- Yang et al., 2005
- accommodate **MLE errors** for some model parameters via uniform priors

Smoothed bootstrap aggregation

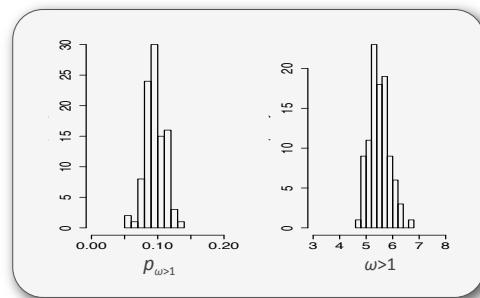
- Mingrone et al., MBE, 33:2976-2989
- accommodate **MLE errors** via bootstrapping
- ameliorates **biases** and **MLE instabilities** with kernel smoothing and aggregation

critical question:

Have the requirements for maximum likelihood inference been met?

(rarely addressed in real data analyses)

regularity conditions have been met

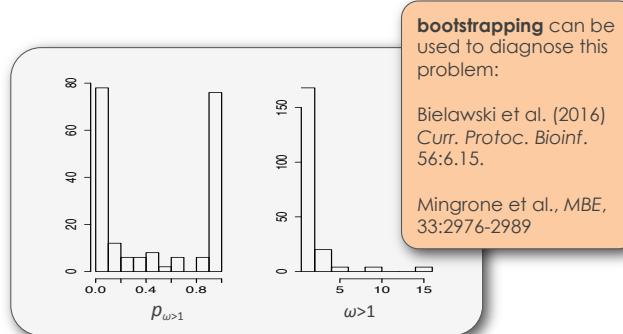


Normal MLE uncertainty (M2a)

- large sample size with regularity conditions
- MLEs approximately unbiased and minimum variance

$$\hat{\theta} \sim N\left(\theta, I(\hat{\theta})^{-1}\right)$$

regularity conditions have **NOT** been met



MLE instabilities (M2a)

- small sample sizes and θ on boundary
- continuous θ has been discretized (e.g., M2a)
- non-Gaussian, over-dispersed, divergence among datasets

software for codon models in the ML framework

PAML: a package of programs for process modeling

HyPhy: comparative sequence analysis using stochastic evolutionary models;
<http://www.hyphy.org/>

DataMonkey: a server that supports a variety of HYPHY tools at no cost;
<http://www.datammonkey.org/>

COLD: a program that implements a general-purpose parametric (GPP) codon model. Most codon models are special cases of the GPP codon model. <https://github.com/tjk23/COLD>

codeml_SBA: a program that implements smoothed Bootstrap Aggregation for Assessing Selection Pressure at Amino Acid Sites.
https://github.com/Jehops/codeml_sba.

ModL: a program for restoring regularity when testing for positive selection using codon models https://github.com/jehops/codeml_modl

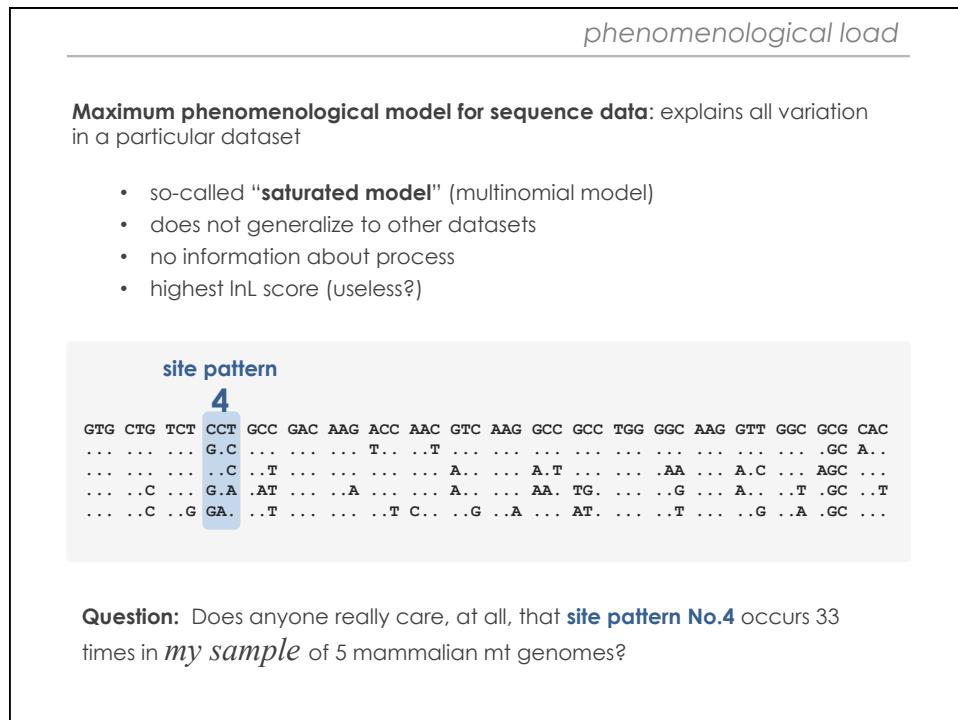
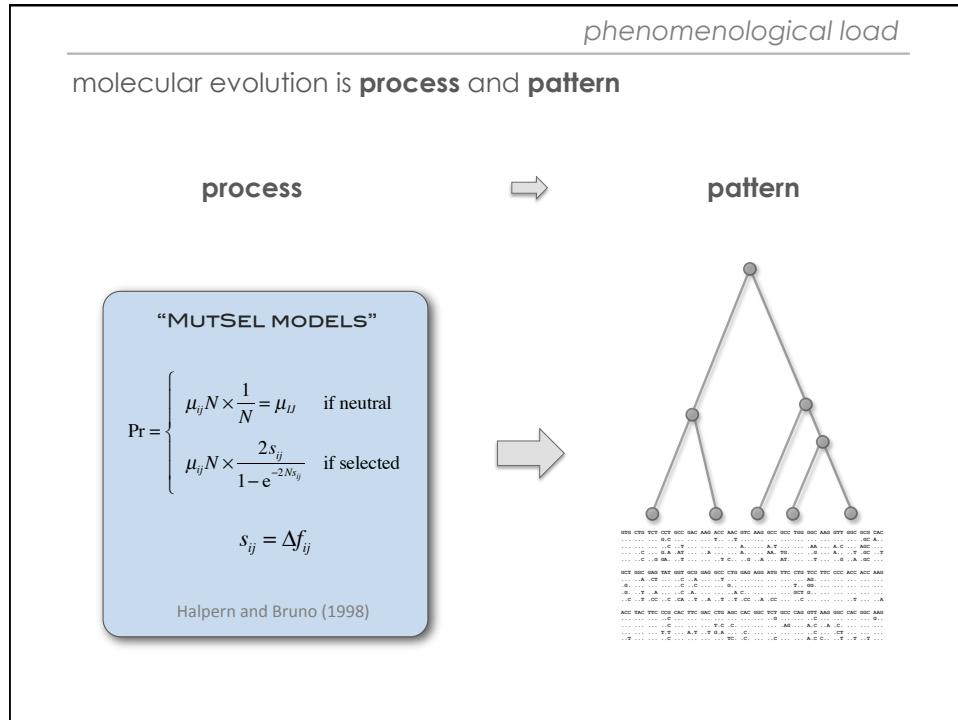
part 4: phenomenological load and biological inference

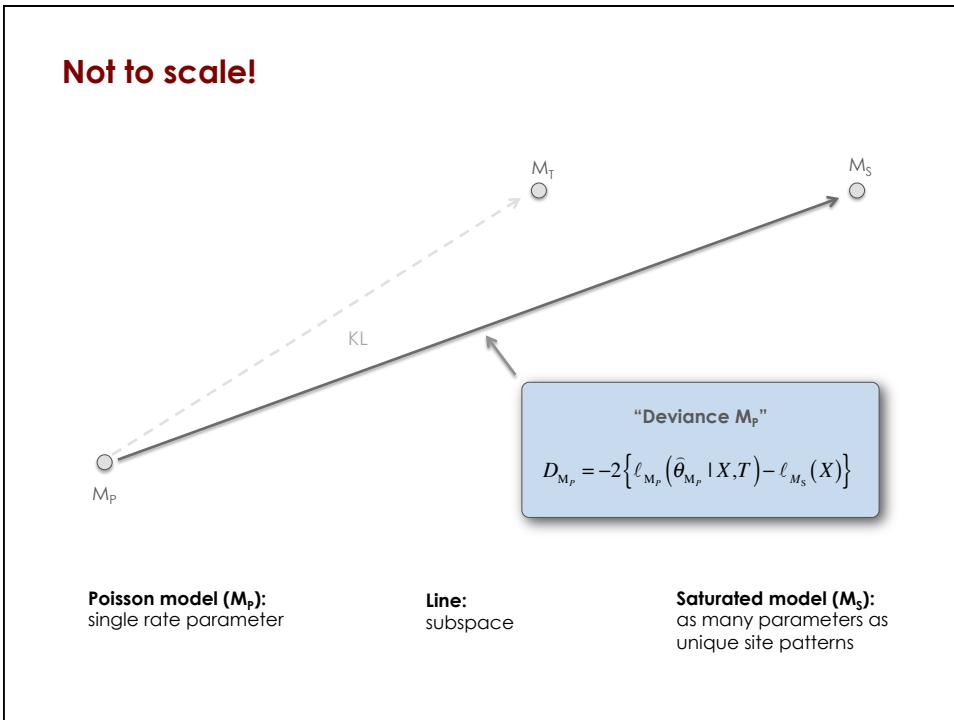
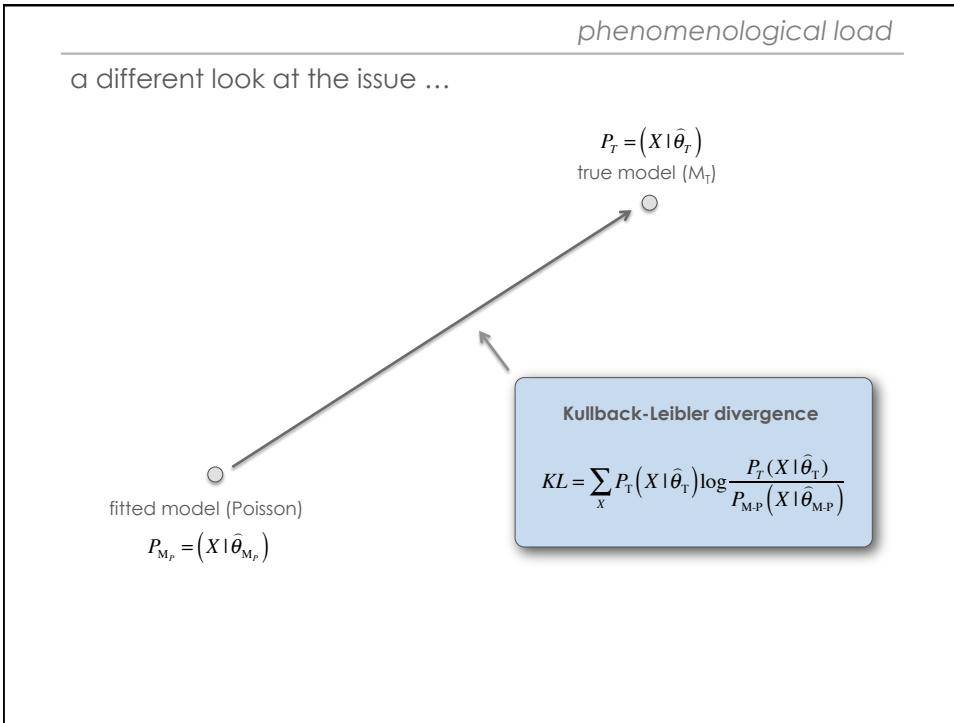
phenomenological load and biological inference part 4: that phenomenon load and biological inference

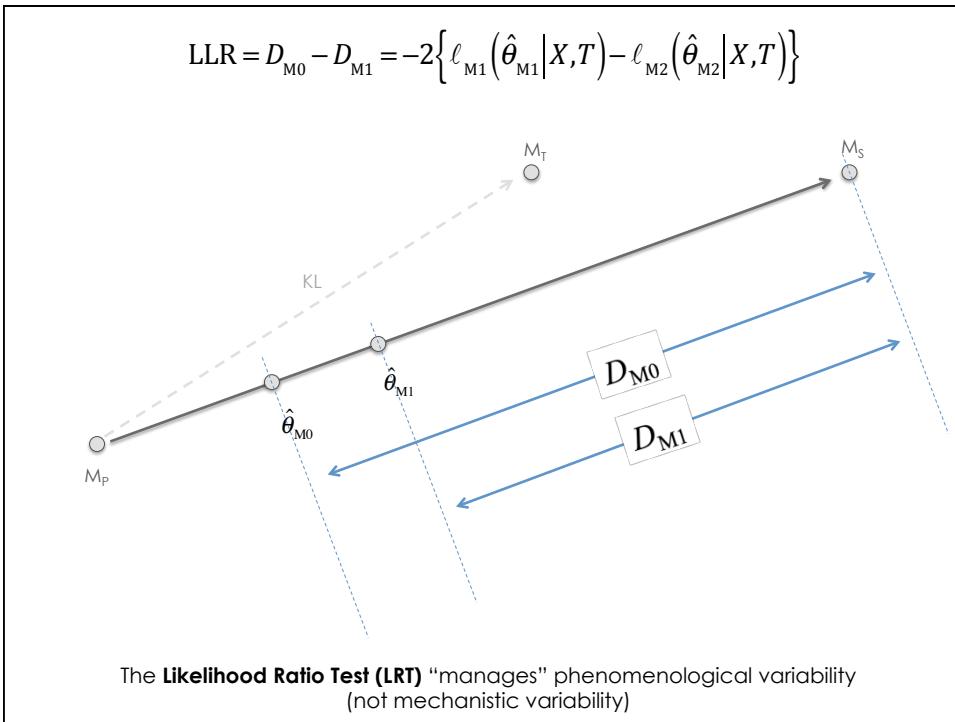
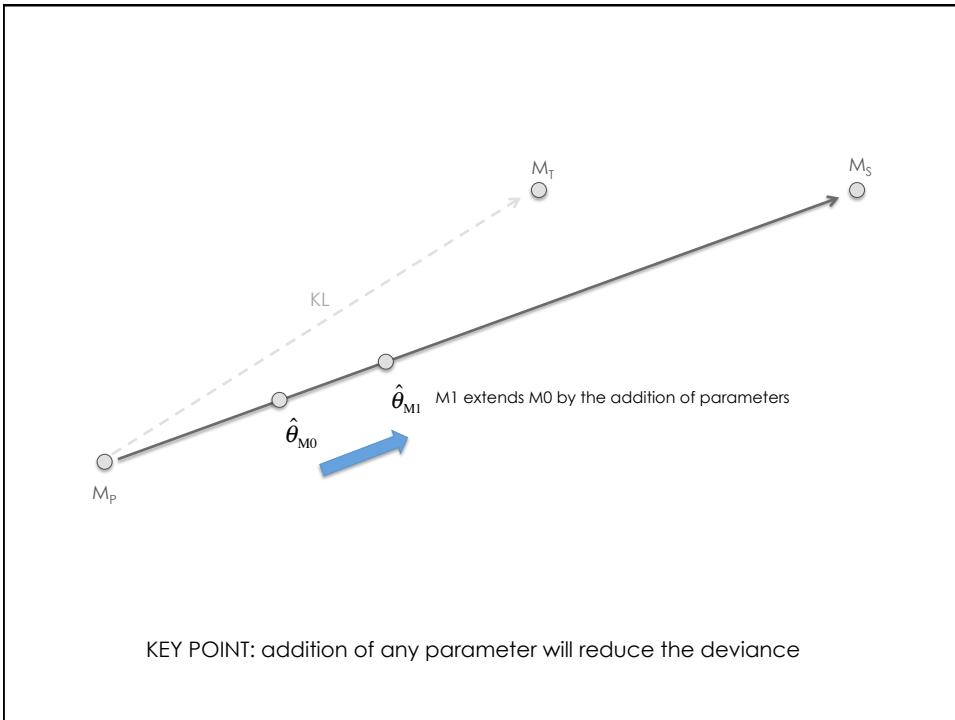
phenomenological load

review types of models

<p>phenomenological</p> <p>Newton</p> $F = -\frac{Gm_1m_2}{r^2}$	<p>mechanistic</p> <p>Einstein</p> $G_{\alpha\beta} = 8\pi T_{\alpha\beta}$
---	--







let's do a simulation study

and

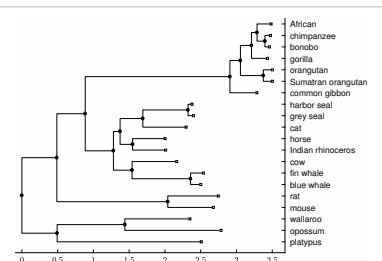
let's use "double mutations" and "triple mutations" as an example

example double (D): ATG (Met) → AAA (Lys)

example triple (T): AAA (Lys) → GGG (GLY)

the simulation and the outcomes...

process (M_T):

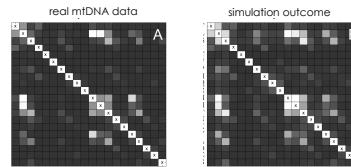
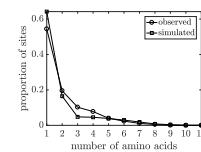


simulation

- MutSel
- μ^b differ for each site
- NO DT-mutations
- 12 mt proteins (3331 codons)
- 20 mammals

outcome (X):

we need outcomes to match up



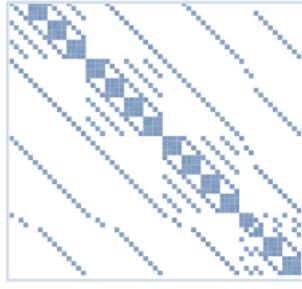
Our simulated data LOOKS LIKE the REAL DATA!

DT: Double and Triple mutations

Example double: ATG (Met) \rightarrow AAA (Lys) [α parameter]
 Example triple: AAA (Lys) \rightarrow GGG (GLY) [β parameter]

M₀ Q matrix

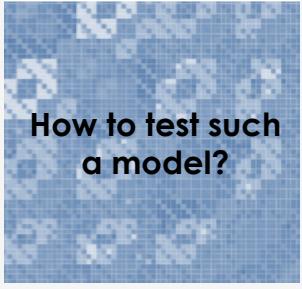
- 2 parameters (κ and ω)
- DT not allowed



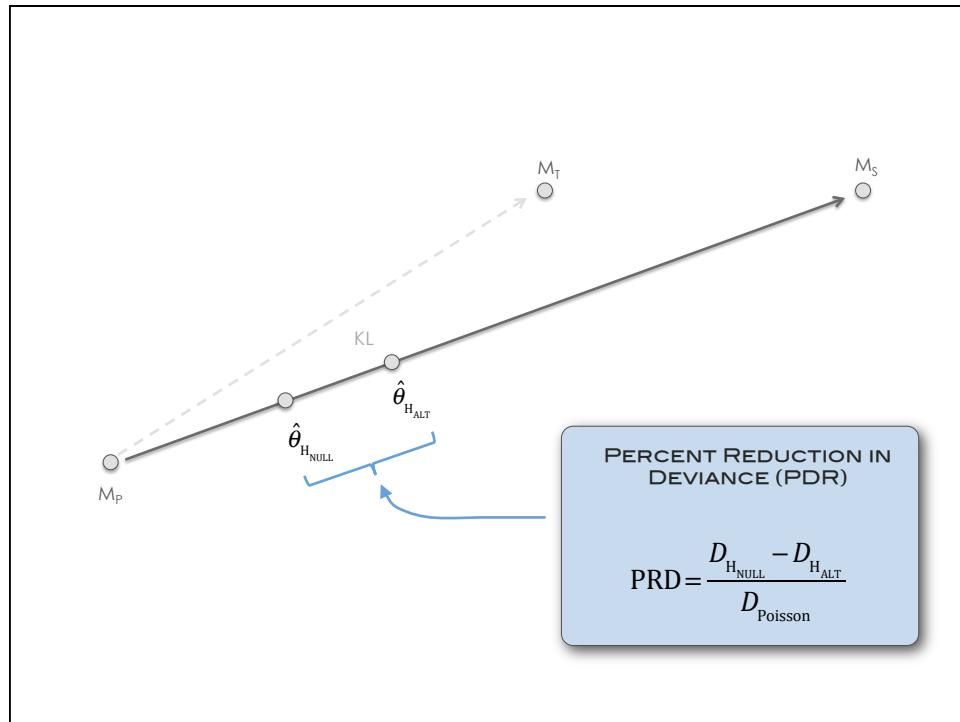
white: probability = 0

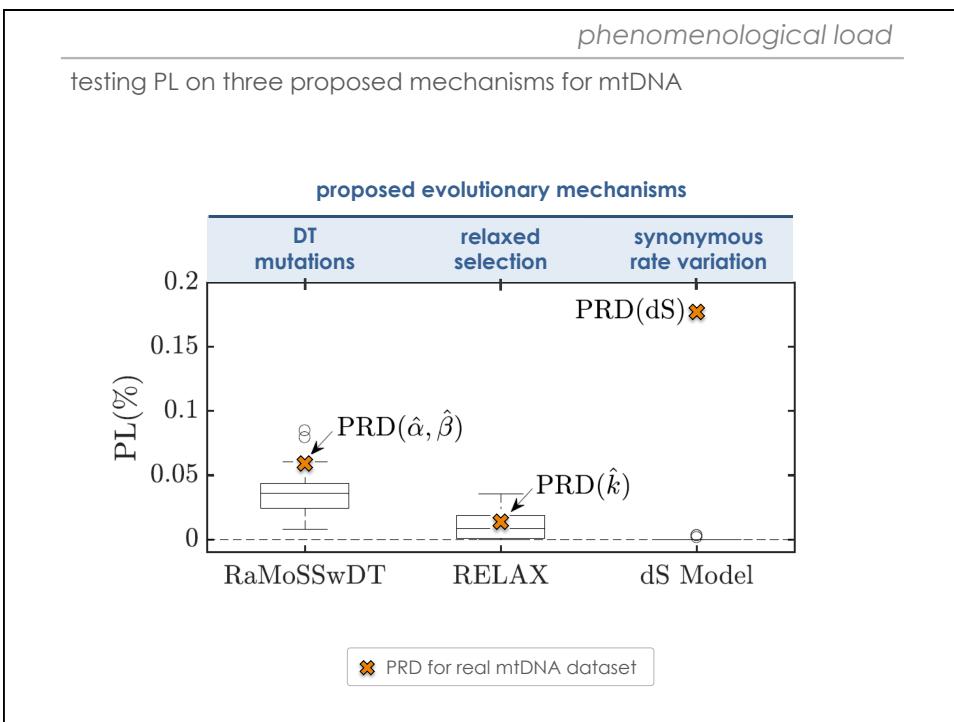
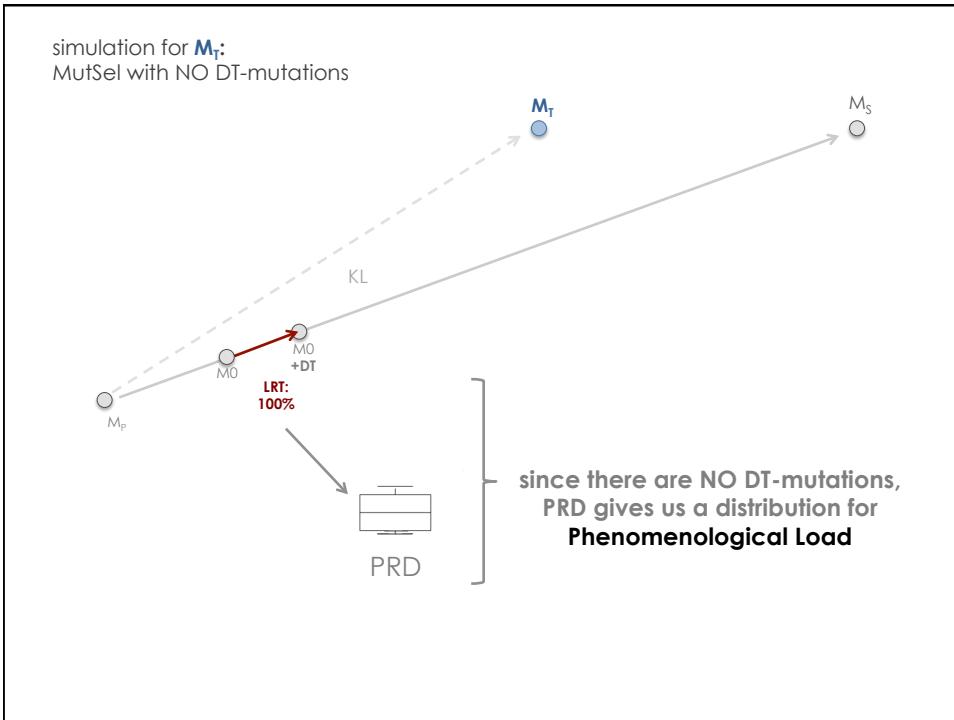
New Q matrix: M₀ + DT

- 4 parameters (κ , ω , α , β)
- DT allowed (via α and β)



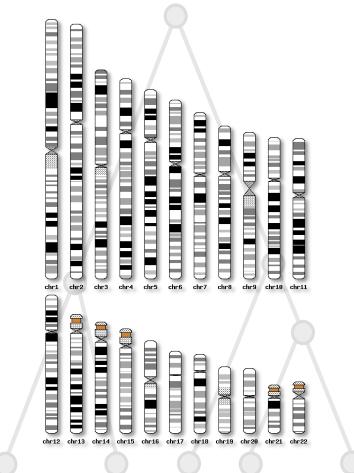
How to test such a model?





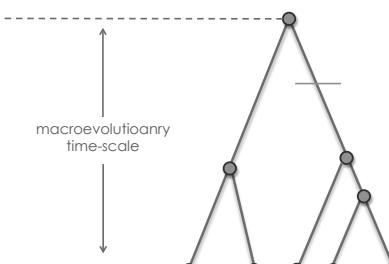
part 5: re-assessing long-held paradigms
for evidence of adaptive evolution

for evidence of adaptive evolution
simulations: re-assessing long-held paradigms



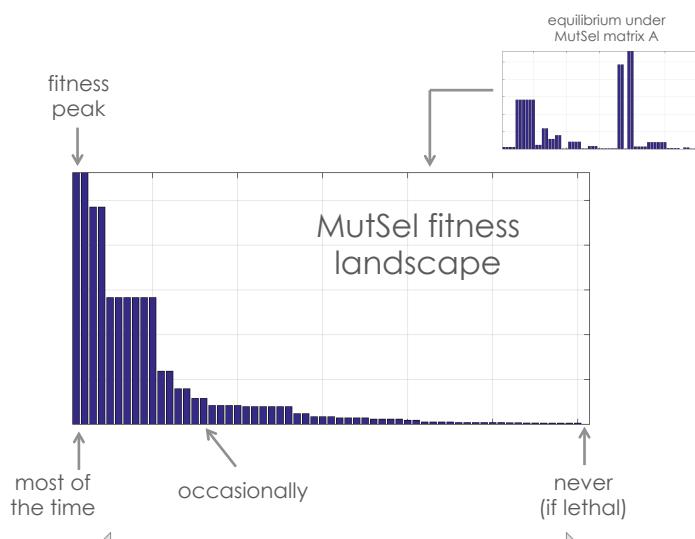
Three paradigms for "this" side:

1. codon substitution model: $d_N/d_S > 1$ is evidence of adaptive evolution
2. "mechanistic" substitution models are better
3. It's easy to test and predict model performance via simulation

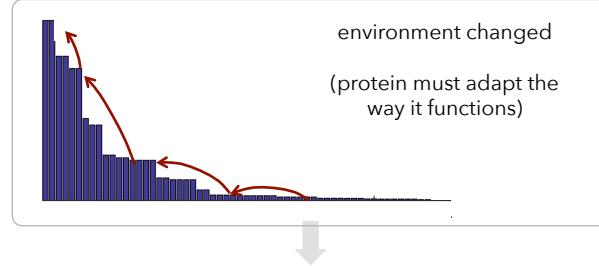


Paradigm 1: $d_N/d_S > 1$ is evidence of adaptive evolution of function

the MutSel fitness landscape



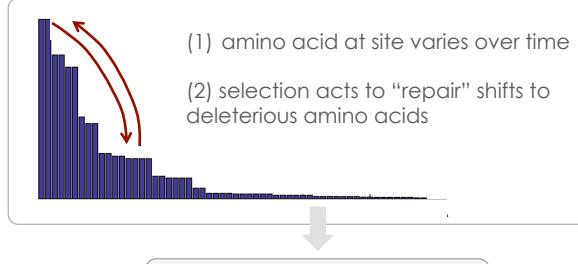
the MutSel fitness landscape: **adaptive evolution**



key result 1:
adaptive evolution: $p_+ > p_-$
("peak shift")

$$d_N/d_S > 1 \text{ (transient)}$$

the MutSel fitness landscape: **non-adaptive shifting balance**



key result 2:
purifying selection: $p_+ = p_-$
(static landscape)

$$d_N/d_S > 1 \text{ (transient)}$$

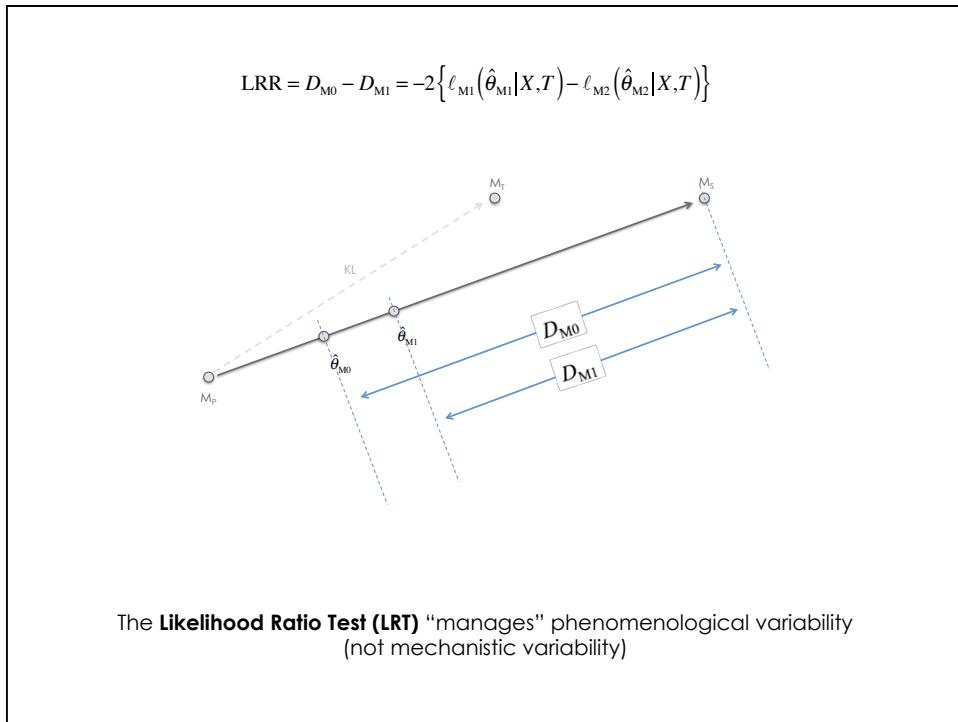
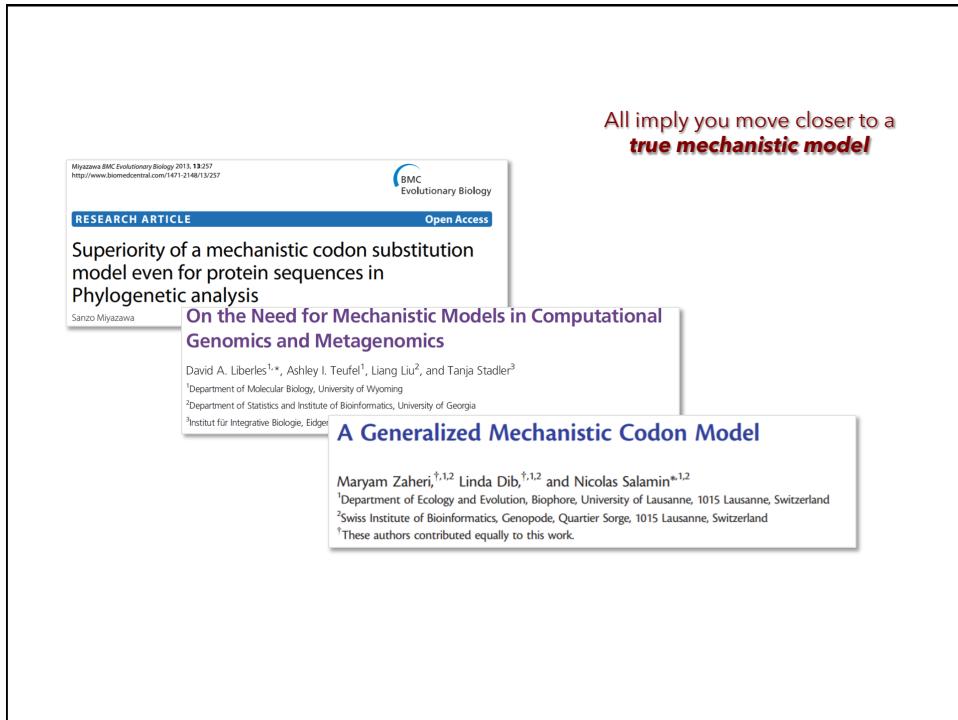
Paradigm 1:



Reality: $d_N/d_S > 1$ on a fixed landscape with **no change in function**

Proposal: develop new frameworks that do NOT depend on the $d_N/d_S > 1$ paradigm

Paradigm 2: "mechanistic" substitution models should be better



Paradigm 2:

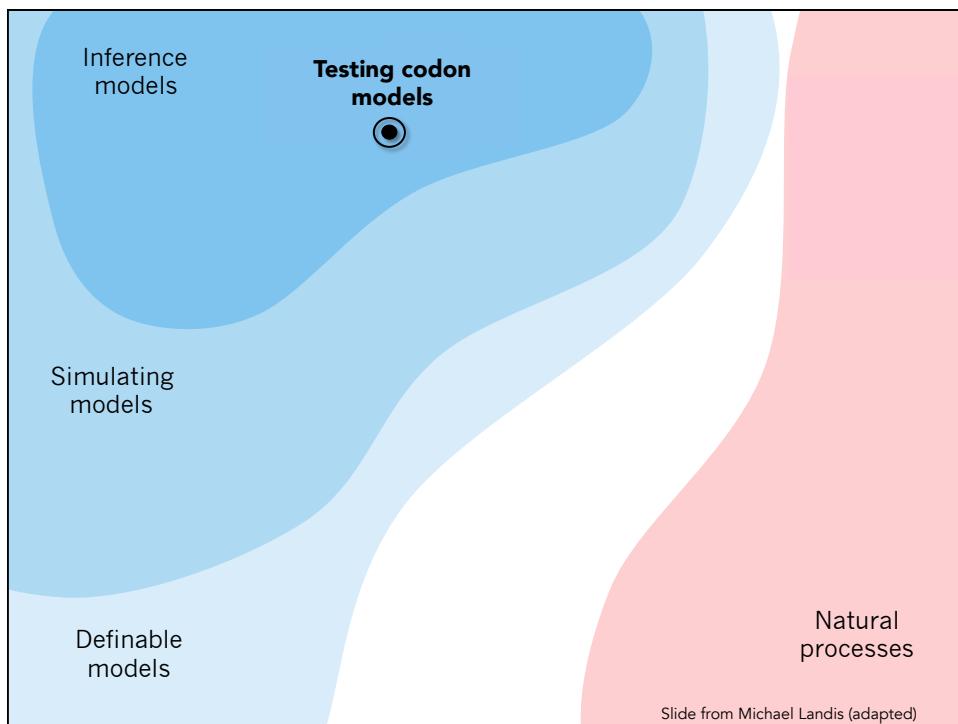
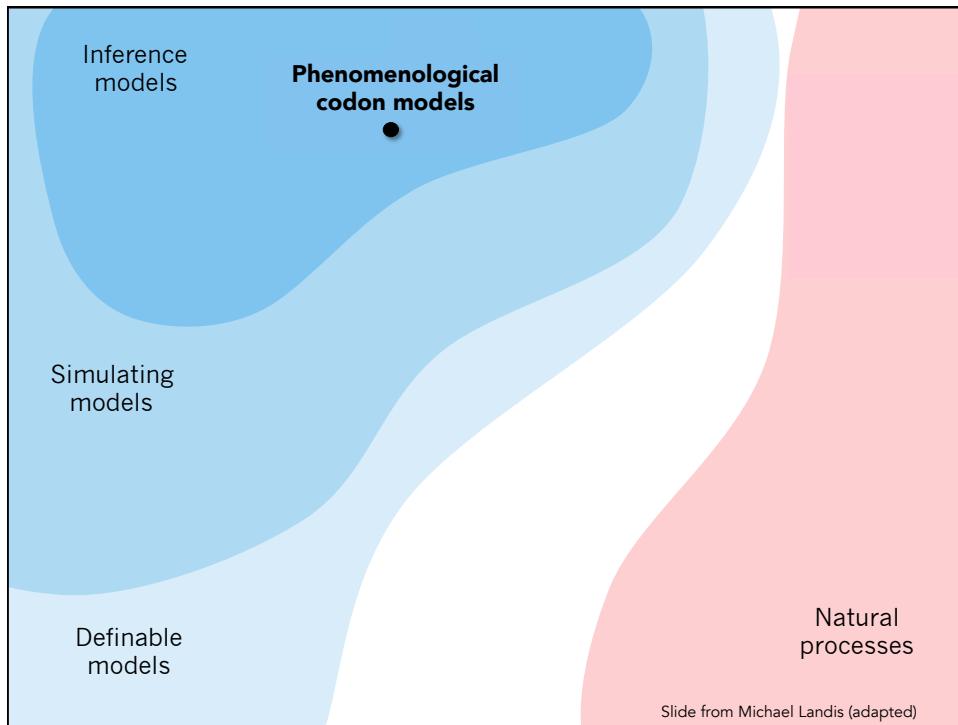


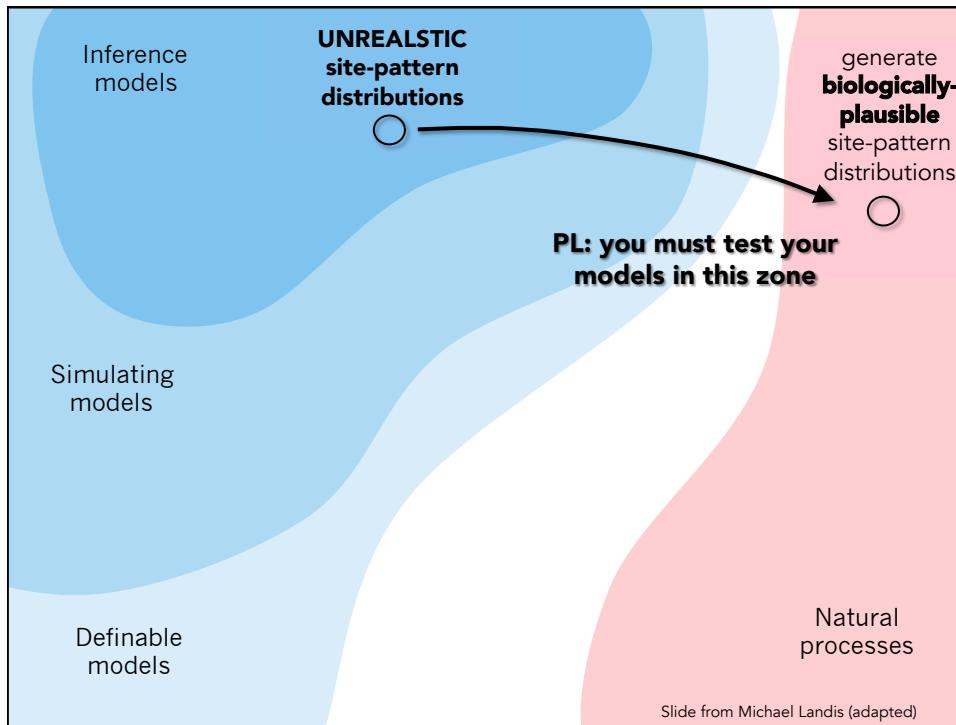
Statistical substitution models
should be better

For real data: mechanistic parameters within
models are expected to carry some
Phenomenological Load

Proposal: intentionally add phenomenological
parameters that improve inferences (e.g., covarion δ)

Paradigm 3: It's easy to test and predict model
performance via simulation





Paradigm 3: It's hard to test model performance against realistic site pattern distributions

MYTH BUSTED!

In reality it's hard to (1) compare complex site pattern distributions, and (2) identify models that produce biologically plausible distributions

Proposal: we need to do more work on how to generate "realistic" site pattern distributions and change the way we think about testing model performance

How can you really tell if you have learned anything relevant to the function of your protein?

- combine computational and **experimental approaches** (B. Chang, next lecture; “**Gold Standard**”)
- informal cross-validation via comparison with **external phenotypic information** (B. Chang, next lecture)
- formally include phenotypic information within the **likelihood inference framework** (we have this working; the paper is in revision... “stay tuned”)

