

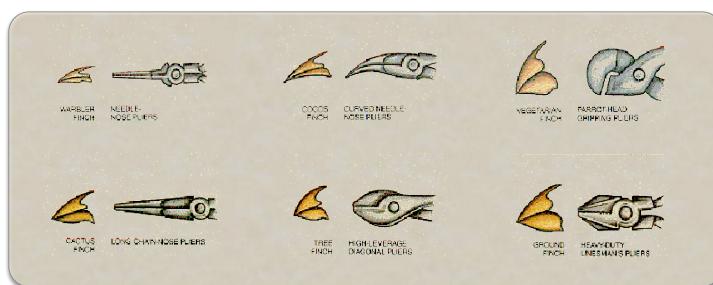
## codon substitution models and the analysis of natural selection pressure



Joseph P. Bielawski  
Department of Biology  
Department of Mathematics & Statistics  
Dalhousie University

### introduction

#### morphological adaptation

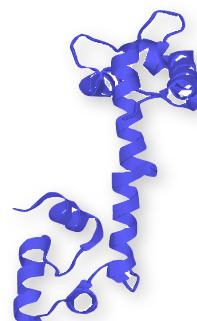


introduction

## protein structure



### Troponin C: fast skeletal



Troponin C: cardiac and slow skeletal

introduction

gene sequences

---

introduction

Powerful analytical tools:

1. Population genetic data
2. Comparative analysis of codon sequences
3. Comparative analysis of amino acid sequences

**“**there is no single statistic which is best for testing the most general departures from neutrality**”**  
(Watterson 1977)

---

introduction

---

overview

1. introduction to modeling codon evolution
2. model based inference
3. PAML introduction & real data exercises

part I

outline

1. introduction to the  $\omega$  ratio
2. markov model of codon evolution
3. codon models for  $\omega$  variation over branches & sites
4. model realism vs. model complexity

1. the  $\omega$  ratio

an index of natural selection pressure

	U	C	A	G	
U	UUU Phe UUC Phe UUA Leu UUG Leu	UCU Ser UCC Ser UCA Ser UCG Ser	UAU Tyr UAC Tyr UAA Stop UAG Stop	UGU Cys UGC Cys UGA Stop UGG Trp	U C A G
C	CUU Leu CUC Leu CUA Leu CUG Leu	CCU Pro CCC Pro CCA Pro CCG Pro	CAU His CAC His CAA Gln CAG Gln	CGU Arg CGC Arg CGA Arg CGG Arg	U C A G
A	AUU Ile AUC Ile AUA Ile AUG Met	ACU Thr ACC Thr ACA Thr ACG Thr	AAU Asn AAC Asn AAA Lys AAG Lys	AGU Ser AGC Ser AGA Arg AGG Arg	U C A G
G	GUU Val GUC Val GUA Val GUG Val	GCU Ala GCC Ala GCA Ala GCG Ala	GAU Asp GAC Asp GAA Glu GAG Glu	GGU Gly GGC Gly GGA Gly GGG Gly	U C A G

Kimura (1968)

 $d_S$ :number of synonymous substitutions per synonymous site ( $K_S$ ) $d_N$ :number of nonsynonymous substitutions per nonsynonymous site ( $K_A$ ) $\omega$ :the ratio  $d_N/d_S$ ; it measures selection at the protein level<http://www.langara.bc.ca/biology/mario/Assets/Geneticcode.jpg>

The genetic code determines how random changes to the gene brought about by the process of mutation will impact the function of the encoded protein.

### 1. the $\omega$ ratio

index of natural selection pressure:  $\omega$  ratio

rate ratio	mode	example
$\omega < 1$	purifying (negative) selection	histones
$\omega = 1$	Neutral Evolution	pseudogenes
$\omega > 1$	Diversifying (positive) selection	MHC, Lysin

### 1. the $\omega$ ratio

the basics

#### Why use $d_N$ and $d_S$ ? (Why not use raw counts?)

Example of counts:

- 300 codon gene from a pair of species
- 5 synonymous differences
- 5 nonsynonymous differences

$$5/5 = 1$$

Why don't we conclude that rates are equal (i.e.,  
**neutral evolution**)?

### 1. the $\omega$ ratio

the basics

Relative proportion of different types of mutations in hypothetical protein coding sequence.				
Type	Expected number of changes (proportion)			
	All 3 Positions	1 <sup>st</sup> positions	2 <sup>nd</sup> positions	3 <sup>rd</sup> positions
Total mutations	549 (100)	183 (100)	183 (100)	183 (100)
Synonymous	134 (25)	8 (4)	0 (0)	126 (69)
Nonsynonymous	392 (71)	166 (91)	176 (96)	57 (27)
nonsense	23 (4)	9 (5)	7 (4)	7 (4)

Modified from Li and Graur (1991). Note that we assume a hypothetical model where all codons are used equally and that all types of point mutations are equally likely.

### 1. the $\omega$ ratio

the basics

#### Why use $d_N$ and $d_S$ ?

Same example, but using  $d_N$  and  $d_S$ :

Synonymous sites = 25.5%

$$S = 300 \times 3 \times 25.5\% = 229.5$$

Nonsynonymous sites = 74.5%

$$N = 300 \times 3 \times 74.5\% = 670.5$$

$$\text{So, } d_S = 5/229.5 = 0.0218$$

$$d_N = 5/670.5 = 0.0075$$

$$d_N/d_S (\omega) = 0.34, \text{ purifying selection !!!}$$

### 1. the $\omega$ ratio

mutational opportunity

**Relative proportion of different types of mutations in hypothetical protein coding sequence.**

Type	Expected number of changes (proportion)			
	All 3 Positions	1 <sup>st</sup> positions	2 <sup>nd</sup> positions	3 <sup>rd</sup> positions
Total mutations	549 (100)	183 (100)	183 (100)	183 (100)
Synonymous	134 (25)	8 (4)	0 (0)	126 (69)
Nonsynonymous	392 (71)	166 (91)	176 (96)	57 (27)
nonsense	23 (4)	9 (5)	7 (4)	7 (4)

Modified from Li and Graur (1991). Note that we assume a hypothetical model where all codons are used equally and that all types of point mutations are equally likely.

Note: by framing the counting of sites in this way we are using a "mutational opportunity" definition of the sites. Not everyone agrees that this is the best approach. For an alternative view see Bierne and Eyre-Walker 2003 Genetics 168:1587-1597.

### 1. the $\omega$ ratio

real data have biases (Drosophila *GstD1* gene)



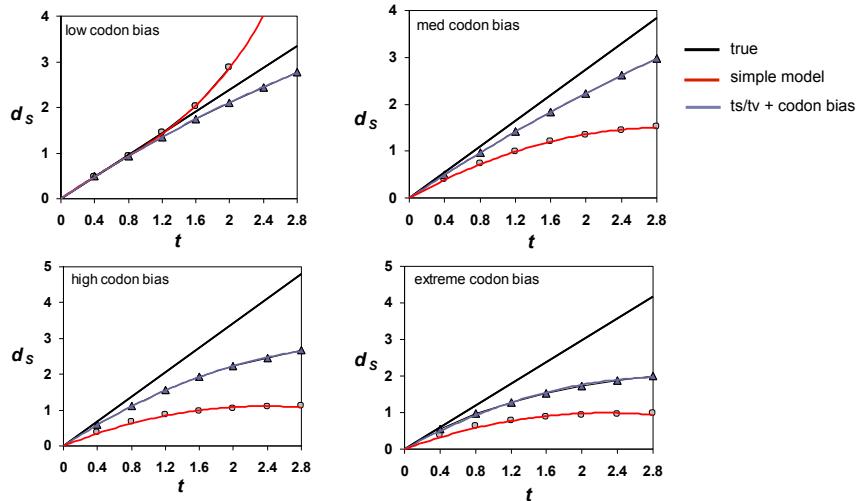
**preferred vs. un-preferred codons:**

Partial codon usage table for the *GstD* gene of *Drosophila*

Phe F	TTT	0   Ser S	TCT	0   Tyr Y	TAT	1   Cys C	TGT	0
	<b>TTC</b>	<b>27</b>	<b>TCC</b>	<b>15</b>	<b>TAC</b>	<b>22</b>	<b>TGC</b>	<b>6</b>
Leu L	TTA	0	TCA	0   *** * TAA		0   *** * TGA		0
	TTG	1	TCG	1   TAG		0   Trp W	<b>TGG</b>	<b>8</b>
Leu L	CTT	2	Pro P	CCT	1   His H	CAT	0   Arg R	CGT
	CTC	2		<b>CCC</b>	<b>15</b>	CAC	4	<b>CGC</b>
	CTA	0		CCA	3   Gln Q	CAA	0	CGA
	<b>CTG</b>	<b>29</b>		CCG	1	<b>CAG</b>	<b>14</b>	CGG

### 1. the $\omega$ ratio

"corrections" and estimation bias in  $d_s$

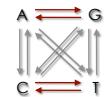


Data from: Dunn, Bielawski, and Yang (2001) Genetics, 157: 295-305

### 2. markovian codon models



#### Markov models of codon evolution



1. assumptions are explicit
2. "corrections" are not *ad hoc*
3. explicit use of a phylogeny improves power
4. principled framework for modelling and inference of the biology

Goldman & Yang 1994 MBE 11:725-736

Muse & Gaut 1994 MBE 11:715-724

2. markovian codon models

“GY-style” codon models (mechanistic)

some important parameters:

- transition/transversion rate ratio:  $\kappa$
- biased codon usage:  $\pi_j$  for codon  $j$
- nonsynonymous/synonymous rate ratio:  $\omega = d_N/d_S$

2. markovian codon models

“GY-style” codon models (mechanistic)

### let's model a change to CTG

Synonymous

**CTC** (Leu) → **CTG** (Leu):                    $\pi_{CTG}$

**TTC** (Leu) → **CTG** (Leu):                    $\kappa\pi_{CTG}$

Nonsynonymous

**GTC** (Val) → **CTG** (Leu):                    $\omega\pi_{CTG}$

**CCG** (Pro) → **CTG** (Leu):                    $\kappa\omega\pi_{CTG}$

## 2. markovian codon models

“GY-style” codon models (mechanistic)

From codon below:	to codon below:							
	TTT (Phe)	TTC (Phe)	TTA (Leu)	TTG (Leu)	CTT (Leu)	CTC (Leu)	...>	GGG (Gly)
<b>TTT (Phe)</b>	—	$\kappa\pi_{TTC}$	$\omega\pi_{TTA}$	$\omega\pi_{TTG}$	$\omega\kappa\pi_{TTT}$	0	...>	0
<b>TTC (Phe)</b>	$\kappa\pi_{TTT}$	—	$\omega\pi_{TTA}$	$\omega\pi_{TTG}$	0	$\omega\kappa\pi_{CTC}$	...>	0
<b>TTA (Leu)</b>	$\omega\pi_{TTT}$	$\omega\pi_{TTC}$	—	—	0	0	...>	0
<b>TTG (Leu)</b>	$\omega\pi_{TTT}$	$\omega\pi_{TTC}$	$\kappa\pi_{TTA}$	—	0	0	...>	0
<b>CTT (Leu)</b>	$\omega\kappa\pi_{TTT}$	0	0	0	—	$\kappa\pi_{CTC}$	...>	0
<b>CTC (Leu)</b>	0	$\omega\kappa\pi_{TTC}$	0	0	$\kappa\pi_{TTT}$	—	...>	0
⋮ ⋮	⋮ ⋮	⋮ ⋮	⋮ ⋮	⋮ ⋮	⋮ ⋮	⋮ ⋮	⋮ ⋮	⋮ ⋮
<b>GGG (Gly)</b>	0	0	0	0	0	0	0	—

\* This is equivalent to the codon model of Goldman and Yang (1994). Parameter  $\omega$  is the ratio  $d_N/d_S$ ,  $\kappa$  is the transition/transversion rate ratio, and  $\pi_i$  is the equilibrium frequency of the target codon ( $i$ ).

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

## 2. markovian codon models

“GY-style” codon models (mechanistic)

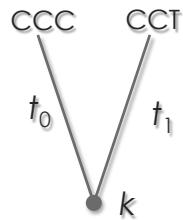
(Goldman & Yang 1994 MBE 11:725-736)

$$q_{ij} = \begin{cases} 0 & \text{if } i \text{ and } j \text{ differ at 2 or 3 positions} \\ \pi_j, & \text{for syn. transversion} \\ \kappa\pi_j, & \text{for syn. transition} \\ \omega\pi_j, & \text{for nonsyn. transversion} \\ \omega\kappa\pi_j, & \text{for nonsyn. transition} \end{cases}$$

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

## 2. markovian codon models

likelihood of the data at a site (only two codons)



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

Note: analysis is typically done by using an unrooted tree

## 2. markovian codon models

likelihood of the data at all sites

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

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

The log likelihood is a sum over all sites.

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

## 2. markovian codon models

we made some progress ...

**the good:** we now have a framework for ...

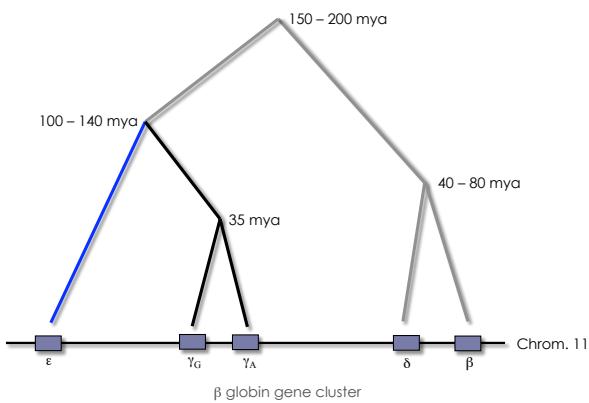
- o avoiding ad hoc corrections of "counting" methods
- o computation of transition probabilities \*
- o principled framework for statistical inference

**a new issue:** averaging  $\omega$  over a pair of sequences has very low power to detect positive selection if the question is about "**when**" or "**where**"  $\omega > 1$ !

\* Computation of transition probabilities accomplishes, in just one step, (1) a proper correction for multiple substitutions, (2) weighting for alternative pathways between codons and (3) is the basis for estimating the values of the model parameters from the data in hand.

## 2. markovian codon models

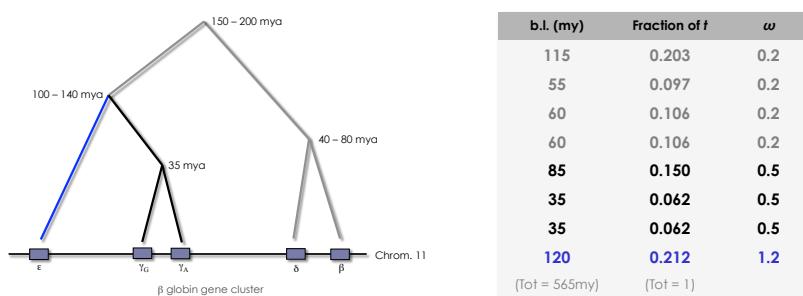
selection pressure ( $\omega$ ) varies in time



Our question: **When?**

## 2. markovian codon models

selection pressure ( $\omega$ ) varies in time



Grey branches:  $\omega = 0.2$

Black branches:  $\omega = 0.5$

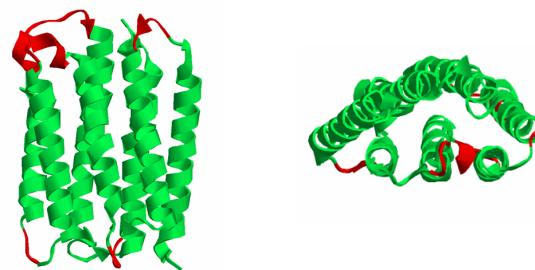
Blue branches:  $\omega = 1.2$

If we average over the tree,  
we do NOT detect positive  
selection;

$$\omega = 0.49.$$

## 2. markovian codon models

selection pressure ( $\omega$ ) varies among sites



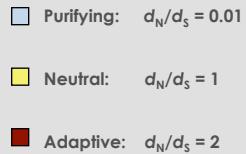
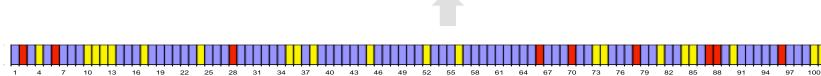
Our question: **Where?**

## 2. markovian codon models

selection pressure ( $\omega$ ) varies among sites

If we average over sites, we do NOT detect positive selection;

$$\omega = 0.31$$

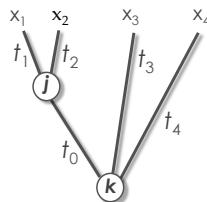


## 2. markovian codon models

likelihood of a phylogeny

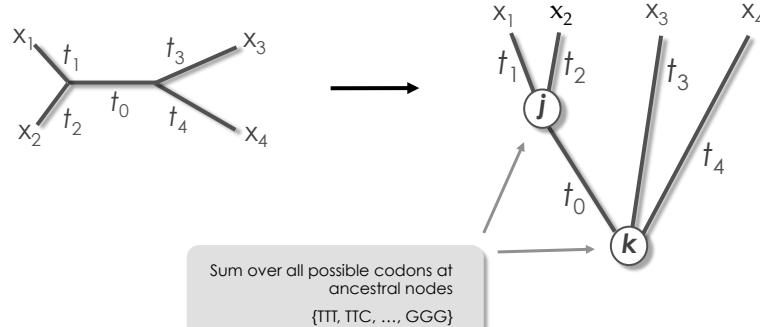
**problem:** averaging  $\omega$  over a pair has very low power if the questions are about “**when**” or “**where**”!

**solution:** phylogenetic estimation of selection pressure



## 2. markovian codon models

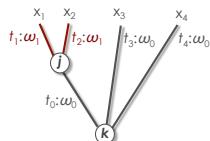
likelihood of a phylogeny



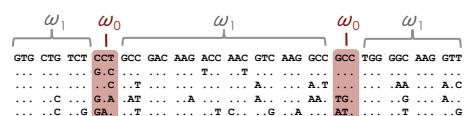
$$L(x_h) = \sum_k \sum_j [\pi_k p_{kx_4}(t_4) p_{kx_3}(t_3) p_{kj}(t_0) p_{jx_2}(t_2) p_{jx_1}(t_1)]$$

## 3. $\omega$ variation

improvements ....



**3.1. branch models**  
( $\omega$  varies among branches)



**3.2. site models**  
( $\omega$  varies among sites)

3.  $\omega$  variation

3.1. branch models \*

Variation ( $\omega$ ) 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**

3.  $\omega$  variation

3.1. branch models

- species colonization of a new niche
- altered context for gene expression
- gene duplication event(s)
- lateral gene transfers (LGTs)
- cross-species virus transmission & host switching
- organismal adaptive radiations

### 3. $\omega$ variation

### 3.2. site models \*

```

GTC CTG TCT CCT GCC GAC AAG ACC AAC GTC AAG GCC GCC TGG GCC AAG GTT GGC CGG CAC
..... .G.C..... .T...T..... .A..... A.T..... .AA..... A.C ..RGC...
..... .C..... .T..... .A..... A..... AA.TG..... .G..... A..T.G.C..T.
..... .C..... G.A..... AT..... A..... A..... AA.TG..... .G..... A..T.G.C..T.
..... .C..... G.A..... TA..... T.C..... G.A..... AT..... T..... G.A..G.C..

```

Variation ( $\omega$ ) 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.

### 3. $\omega$ variation

### 3.2. site models: “M-series”

Model	Code	NP	Parameters
One-ratio	M0	1	$\omega$
Neutral	M1a	2	$P_0, \omega_0$
Selection	M2a	4	$P_0, P_1, \omega_0, \omega_2$
Discrete	M3	2K-1	$P_0, P_1, \dots, P_{K-2}$ $\omega_0, \omega_1, \dots, \omega_{K-2}$
Frequency	M4	5	$P_0, P_1, \dots, P_4$
Gamma	M5	2	$\alpha, \beta$
2Gamma	M6	4	$P_0, \alpha_0, \beta_0, \alpha_1$
Beta	M7	2	$p, q$
Beta& $\omega$	M8	4	$P_0, p, q, \omega$
Beta&gamma	M9	5	$P_0, p, q, \alpha, \beta$
Beta&normal+1	M10	5	$P_0, p, q, \alpha, \beta$
Beta&normal>1	M11	5	$P_0, p, q, \mu, \alpha$
0&2normal>1	M12	5	$P_0, P_1, \mu_2, \alpha_1, \alpha_2$
3normal>0	M13	6	$P_0, P_1, \mu_2, \alpha_0, \alpha_1, \alpha_2$

### 3. $\omega$ variation

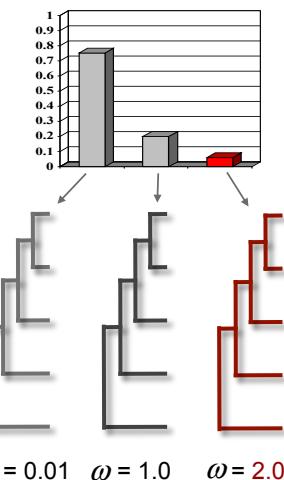
### 3.2. discrete model

```

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

```

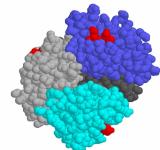
$$P(\mathbf{x}_h) = \sum_{i=0}^{K-1} p_i P(\mathbf{x}_h \mid \omega_i)$$



### 3.2. site models

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

“globin model”



directional  
selection

## “MHC model”



diversifying selection  site models

Site models detect diversifying positive selection; **absence of “signal” under a sites model does not mean changes in a protein did not have fitness consequences** (Bielawski et al. 2004 PNAS 101: 14824–29 )

### 3. $\omega$ variation

### 3.2. site models

```

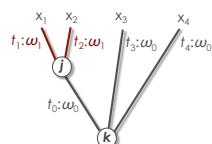
GTG CTG TCT CCT GCC GAC AAG ACC AAC GTC TGG GCC AAG GTT GGC GCG CAC
..... .C..... .T... .T..... .A..... .A.T..... .AA..... .TG..... .G... .A...
..... .C..... .T... .A..... .A..... .AA..... .TG..... .G... .A... .T... .G... .A...
..... .C..... .G.A... .T..... .T.C... .G... .A..... .AT..... .T... .G...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

```

- genetic incompatibilities in human infertility
  - non-hormonal contraception drugs
  - identify pathogenicity genes
  - venom-anti-venom co-evolution
  - identify candidate genes for drug therapies
  - identify immune and defense system genes
  - vaccine design
  - aid functional classification of unknown genes
  - incorporate in models of protein 2D and 3D structure

### 3. $\omega$ variation

we made some more progress ....



### 3.1. branch models ( $\omega$ varies among branches)

## 3.2. site models ( $\omega$ varies among sites)

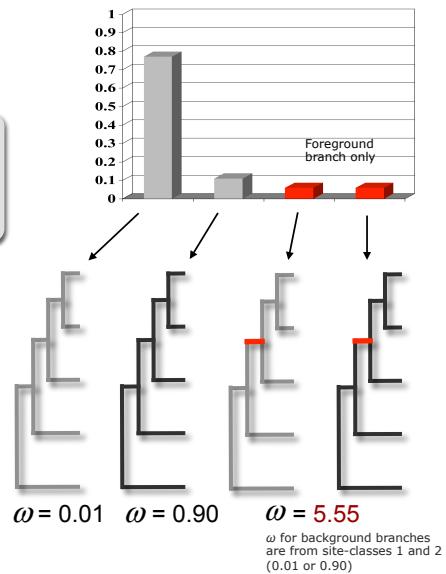
### **3.3. branch-site models**

(combines the features of above models)

### 3. $\omega$ variation

#### 3.3. branch-site model: “model B”

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



### 3. $\omega$ variation

#### 3.3. models for variation in branches & sites

Variation ( $\omega$ ) among branches & sites:	Approach
Yang and Nielsen, 2002	fixed+mixtures (ML)
Forsberg and Christiansen, 2003	fixed+mixtures (ML)
Bielawski and Yang, 2004	fixed+mixtures (ML)
Giundon et al., 2004	switching (ML)
Zhang et al. 2005	fixed+mixtures (ML)
Kosakovsky Pond et al. 2011, 2012	full mixture (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**).

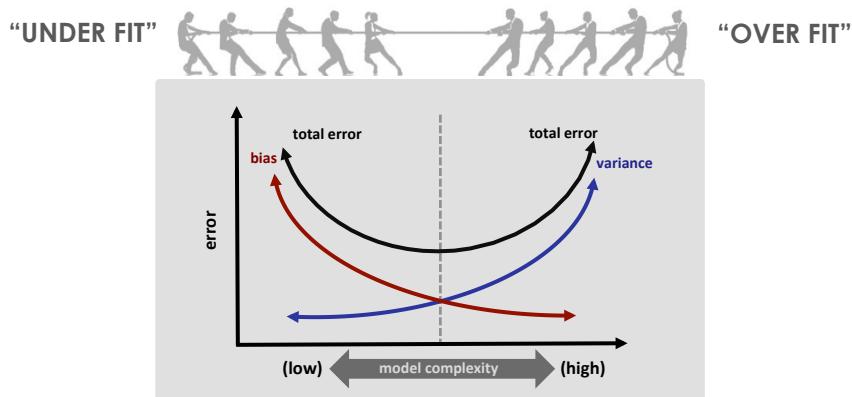
**a MODEL** is an intentional simplification of a complex situation designed to eliminate extraneous detail in order to focus attention on the essentials of the situation

(Daniel L. Hartl)



(images from Paul Lewis)

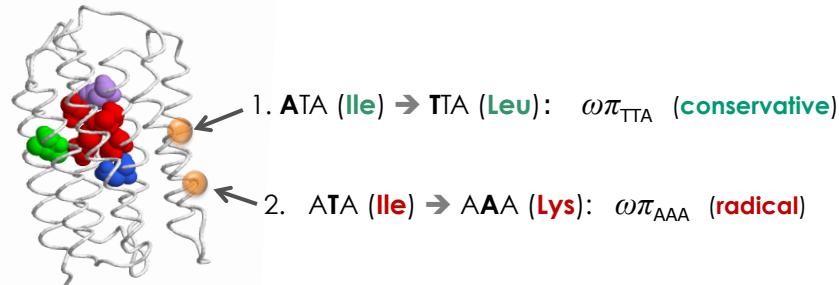
there is a kind of “tension” between model realism and model complexity



## 4. realism vs. complexity

## AA exchangeabilities

**intentional simplification:** all of the models listed above treat these two amino acid substitutions in the same way!



## 4. realism vs. complexity

## Q matrix for M0 codon models

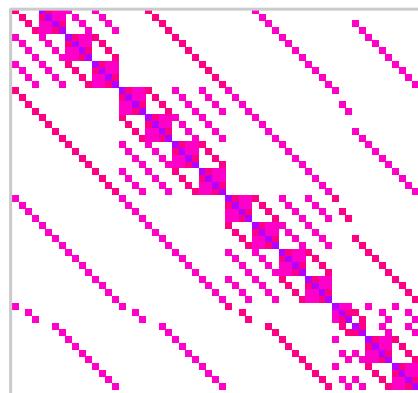
From codon below:	to codon below:							
	TTT (Phe)	TTC (Phe)	TTA (Leu)	TTG (Leu)	CTT (Leu)	CTC (Leu)	...>	GGG (Gly)
TTT (Phe)	—	$\kappa\pi_{TTC}$	$\omega\pi_{TTA}$	$\omega\pi_{TTG}$	$\omega\kappa\pi_{TTT}$	0	...>	0
TTC (Phe)	$\kappa\pi_{TTT}$	—	$\omega\pi_{TTA}$	$\omega\pi_{TTG}$	0	$\omega\kappa\pi_{CTC}$	...>	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  $\omega$  is the ratio  $d_N/d_S$ ,  $\kappa$  is the transition/transversion rate ratio, and  $\pi_i$  is the equilibrium frequency of the target codon ( $i$ ).

---

4. realism vs. complexity

structure of rate matrix: GY-M0

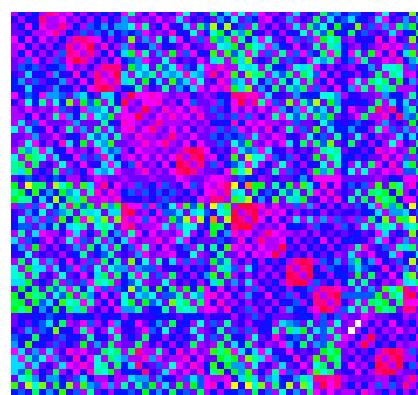
2 parameters ( $\kappa$  and  $\omega$ )too  
sparse?Structure of a Q matrix (log scaled and binned)  
derived from M0 for Abalone sperm lysin

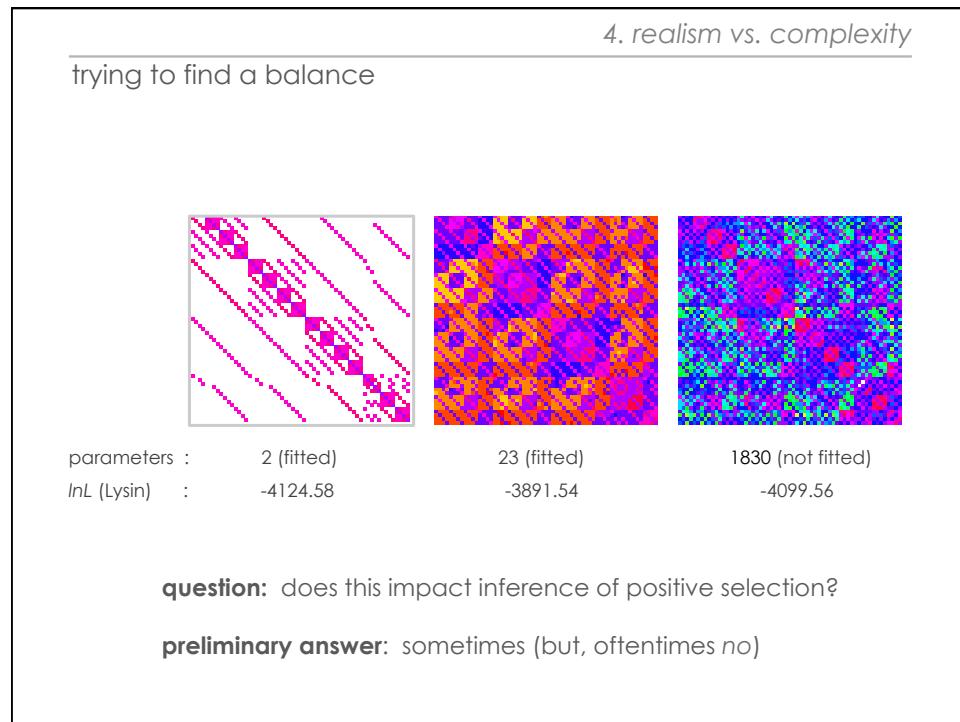
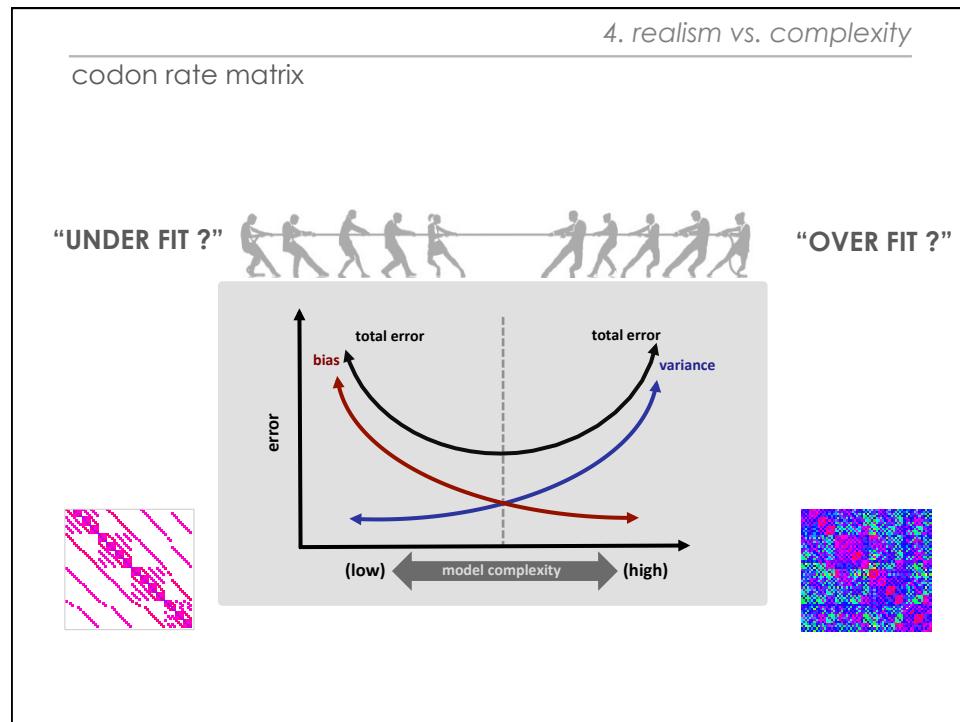
---

4. realism vs. complexity

structure of rate matrix: ECM

1830 parameters

way  
too  
many!Structure of a Q matrix (log scaled and binned) derived  
ECM (Kosiol et al. 2007) for all sequences in the Pandit  
database



#### 4. realism vs. complexity

modeling amino acid exchangeabilities

multiple amino acid exchangeabilities:	model type
Yang and Goldman 1994; Yang et al. 1998;	PCP
Sanudiin et al. 2005; Wong et al. 2006	PCP
Conant and Stadler (2009)	PCP
Kosiol et al. 2007; De Maio et al. 2012	ECM → MEP *
Doron-Faigenboim & Pupko 2007	MEP *
Miyazawa (2011)	MEP (LCEP) *
Zoller and Schneider (2012)	MEP (LCEP) *
Delpont et al. 2010	GPP (GA)
Zaheri, Dib and Salamin (2014)	GPP *
Bielawski et al. (in prep.)	GPP *

PCP: physiochemically constrained parametric model

MEP: mixed empirical and parametric model

ECM: empirical codon model

GPP: general purpose parametric model

\* models permitting double and triple changes among codons

#### codon models

lots of options ...

- methods have different pros and cons (that's OK)
- researchers have different preferences (that's OK)
- I use a wide variety of different methods; you might have different preferences (that's OK)
- most important: KNOW YOUR DATA and MAKE INFORMED DECISIONS.

**no single method will be suitable for all situations**