

[bit.ly/selection-slides-2017](https://bit.ly/selection-slides-2017)

# Quantifying Natural Selection in Coding Sequences

**Sergei L Kosakovsky Pond**

*Professor, Department of Biology*

*Institute for Genomics and Evolutionary Medicine (iGEM)*

*Temple University*

spond@temple.edu

[www.hyphy.org/sergei](http://www.hyphy.org/sergei)

# Preliminaries

---

- Please confirm access to **HyPhy**: <http://hyphy.org/download/>
- General user questions and feedback: <https://github.com/veg/hyphy/issues>
- **Datamonkey** web-app:
  - <http://www.datamonkey.org>
  - <http://test.datamonkey.org>
- Test datasets and practical instructions: [bit.ly/hyphy-selection-tutorial](http://bit.ly/hyphy-selection-tutorial)

# Outline

---

- Brief background and examples of natural selection
- **dN/dS** as a tool to measure the action of natural selection, explained using the first counting method for estimating dN/dS (Nei-Gojobori, 1986) and its extensions.
- Codon substitution models – the basis of modern (1998-) dN/dS estimation approaches
- Different types of selection analyses enabled by **dN/dS**, told by examples from West Nile virus and HIV and analogies from image analysis
  - Gene-wide selection (BUSTED)
  - Lineage-specific selection (aBSREL)
  - Site-level **episodic** selection (MEME)
  - Site-level **pervasive** selection (FUBAR)
  - Relaxed or intensified selection (RELAX)
- Confounding processes (synonymous rate variation, recombination)
- On the suitability of dN/dS for within-species inference

# A bit of trivia

---

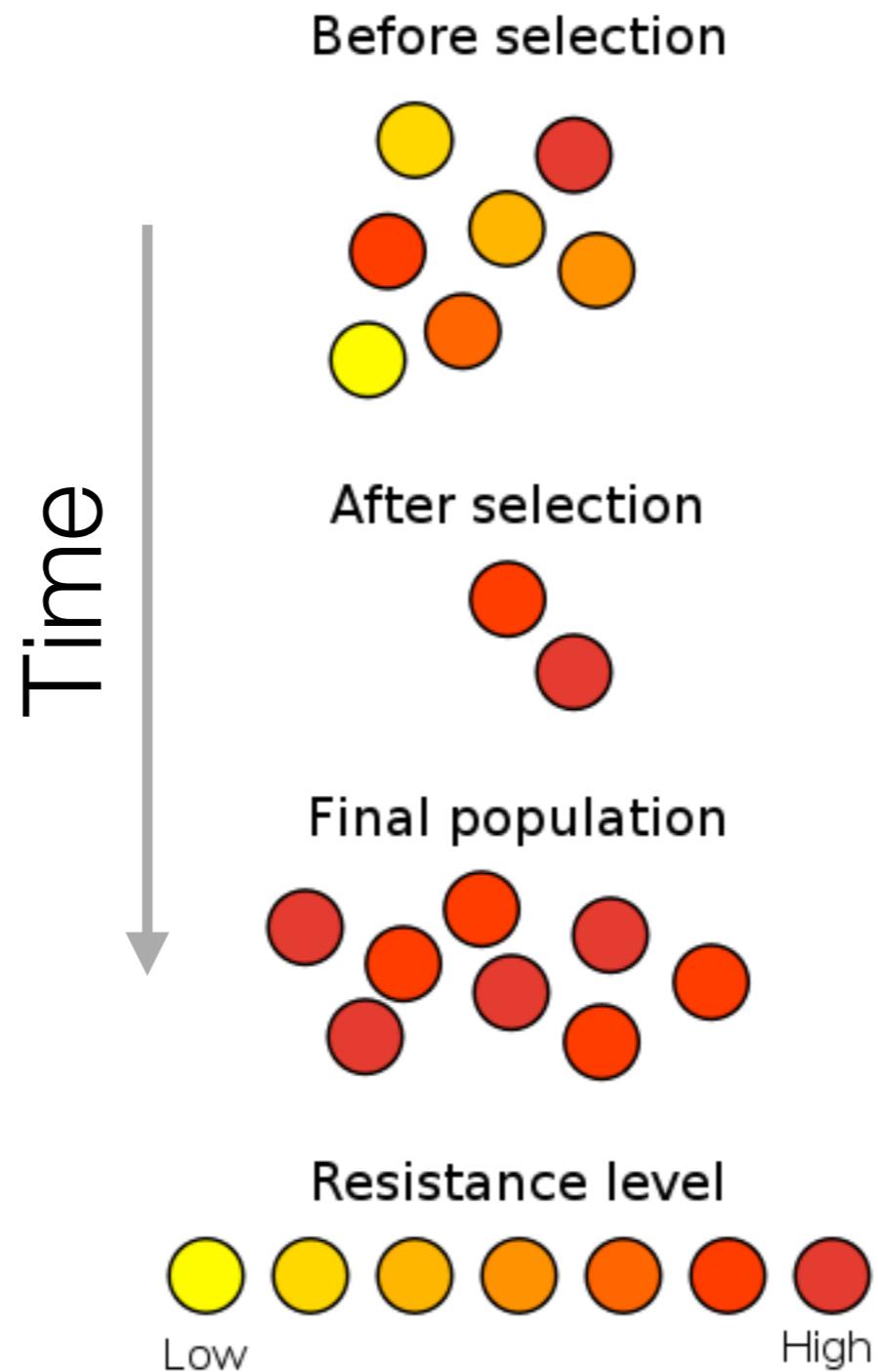
- The theory of natural selection was first proposed by ...*Patrick Matthew*
- Matthew seemed to regard the idea as more or less self-evident and not in need of further development.
- In a stunning example of how **not** to communicate science, he published his ideas in appendices B and F of his book “*On Naval Timber and Arboriculture*” (1831).
- Unsurprisingly, his peers failed to discover his ideas in such an obscure source, and his work had no impact on the subsequent, more developed, work of Darwin and Wallace (1859).
- Do **not** emulate Patrick Matthew.



# Natural Selection

---

- Mutation, recombination and other processes introduce variation into genomes of organisms
- The fitness of an organism describes how well it can survive/grow/function/replicate in a given environment, or how well it can pass on its genetic material to future generations
- Any particular mutation can be
  - Neutral: no or little change in fitness (the majority of genetic variation falls into this class according to the neutral theory)
  - **Deleterious**: reduced fitness
  - **Adaptive**: increased fitness
- The same mutation can have different fitness costs in different environments (fitness landscape), and different genetic backgrounds (epistasis)

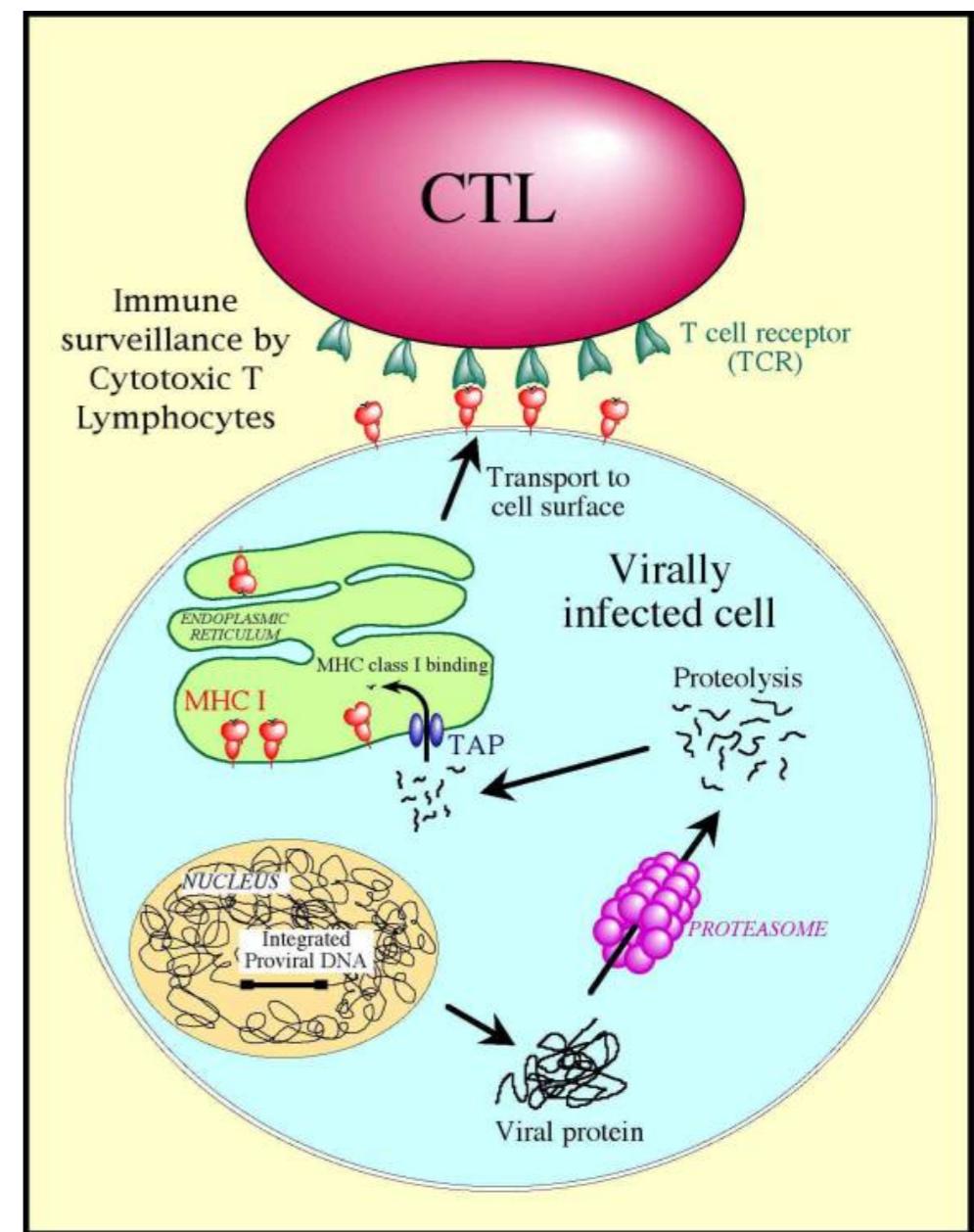






# Example: MHC-restricted CTL killing of infected cells

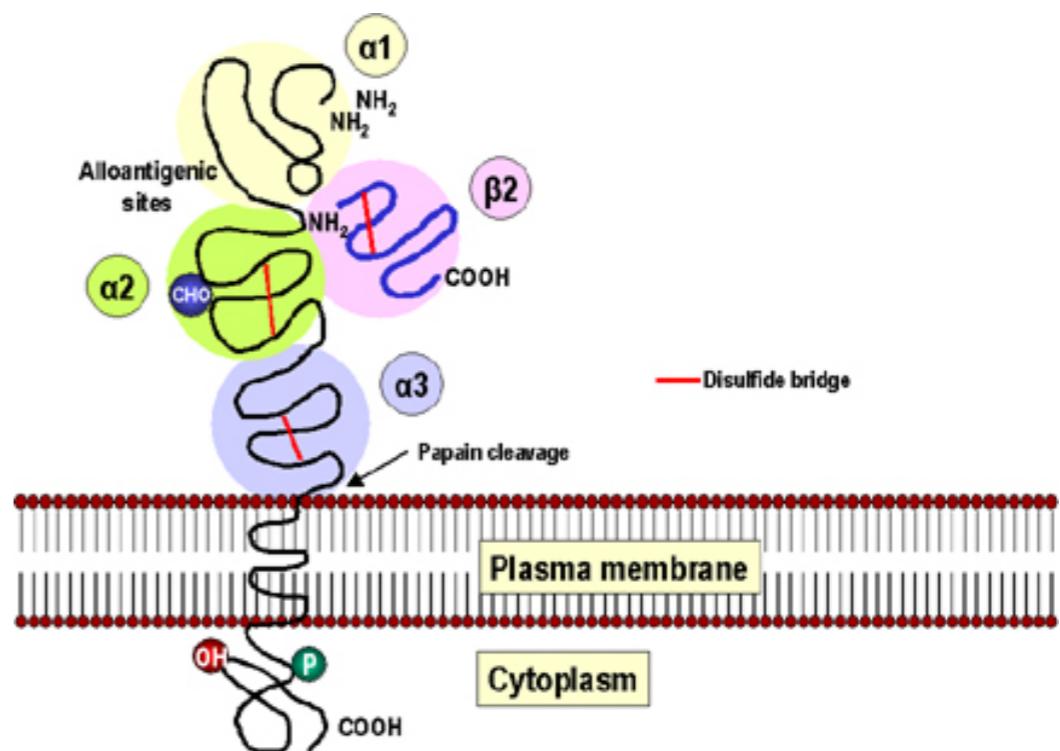
- Cytotoxic T-lymphocytes effect cell-mediated immune response
- Foreign (e.g., viral) proteins are cleaved by the proteasome, transported by TAP and loaded onto the MHC Class 1 molecule.
- MHC Class 1 presents a restricted polypeptide (epitope) on the surface of the cell.
- A CD8+ cell binds to presented foreign peptides via a T cell receptor (TCR) and initiates infected cell apoptosis.



# MHC Class 1 Molecules

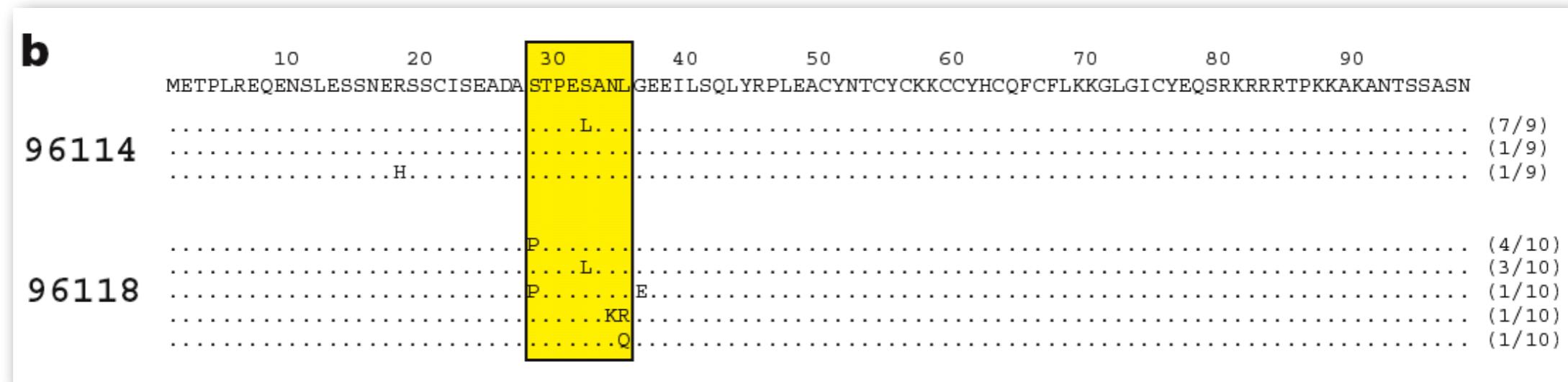
- Present **linear** foreign peptides which are most commonly 9 or 10 aminoacids long
- Anchor sites (2 and 9) are usually important for binding and recognition
- Mutations which alter the peptide can hinder or prevent CTL response activation

## Antigen Binding Site



# Rapid SIV sequence evolution in macaques in response to CTL-driven selection

- SIV: the only animal model of HIV (rhesus macaques)
  - Experimental infection with MHC-matched strain of SIV
  - Virus sequenced from a sample 2 weeks post infection
  - Only variation was in an epitope recognized by the MHC
    - CTL escape



# Key drivers of adaptation in pathogens

---

- Zoonoses and transmission to new hosts (both species and individuals)
- Immune selection (CTL, innate, antibody)
- Development of drug resistance
- Virulence/transmissibility
- Host/pathogen arms-races, e.g. host antiviral factors
- **Most of the time, most of the viral genome is conserved**

# Evolution of Coding Sequences

---



# Evolution of Coding Sequences

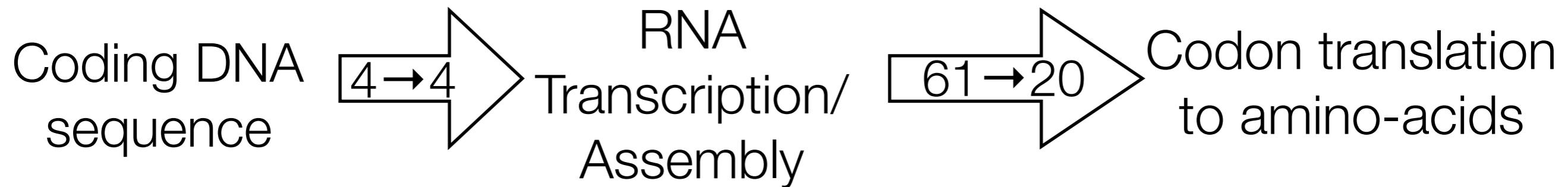
---



- Proper unit of evolution is a triplet of nucleotides — a **codon**

# Evolution of Coding Sequences

---



- Proper unit of evolution is a triplet of nucleotides — a **codon**
- **Mutation** happens at the **DNA level**

# Evolution of Coding Sequences

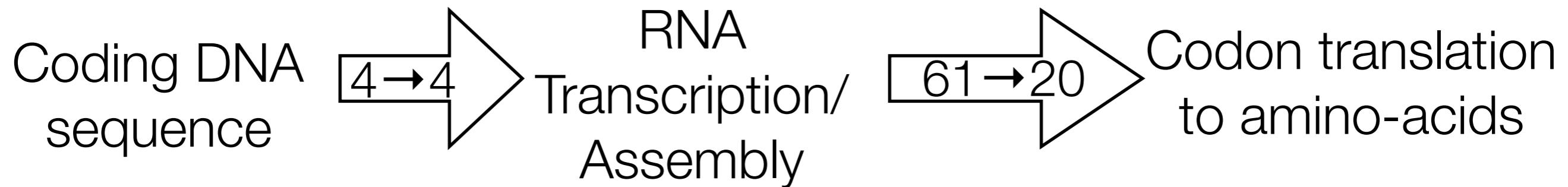
---



- Proper unit of evolution is a triplet of nucleotides — a **codon**
  - **Mutation** happens at the **DNA level**
  - **Selection** happens (by and large) at the **protein level**

# Evolution of Coding Sequences

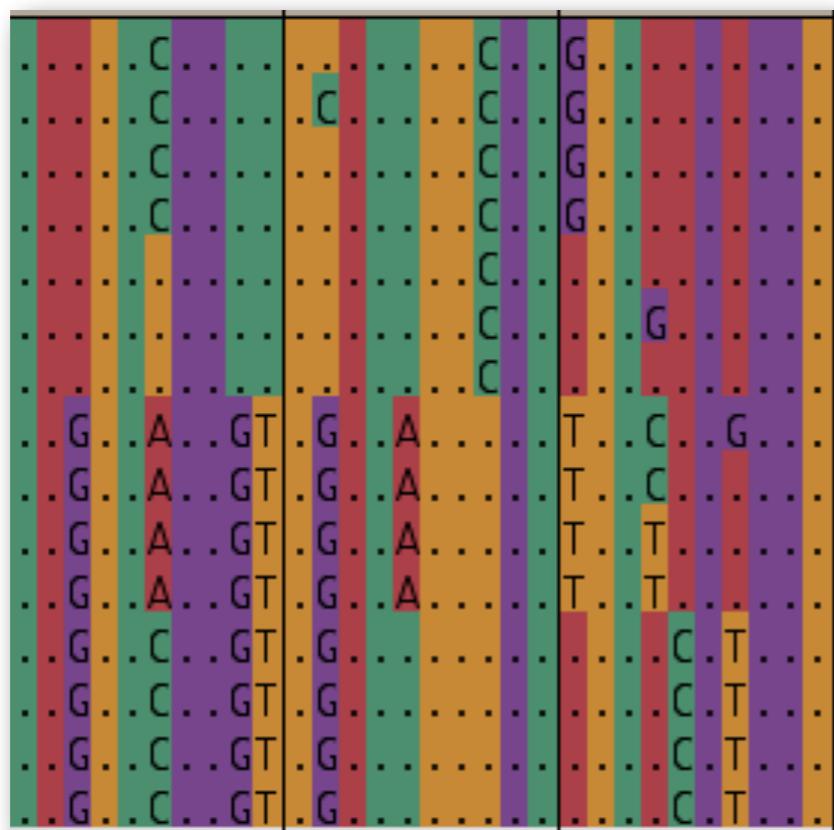
---



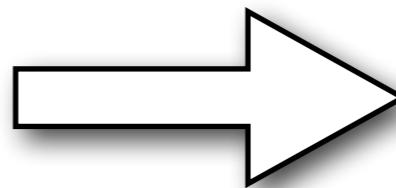
- Proper unit of evolution is a triplet of nucleotides — a **codon**
  - **Mutation** happens at the **DNA level**
  - **Selection** happens (by and large) at the **protein level**
- **Synonymous** (protein sequence **unchanged**) and **non-synonymous** (protein sequence **changed**) substitutions are fundamentally different

# Conservation

Measles, rinderpest, and *peste-de-petite* ruminant viruses nucleoprotein.



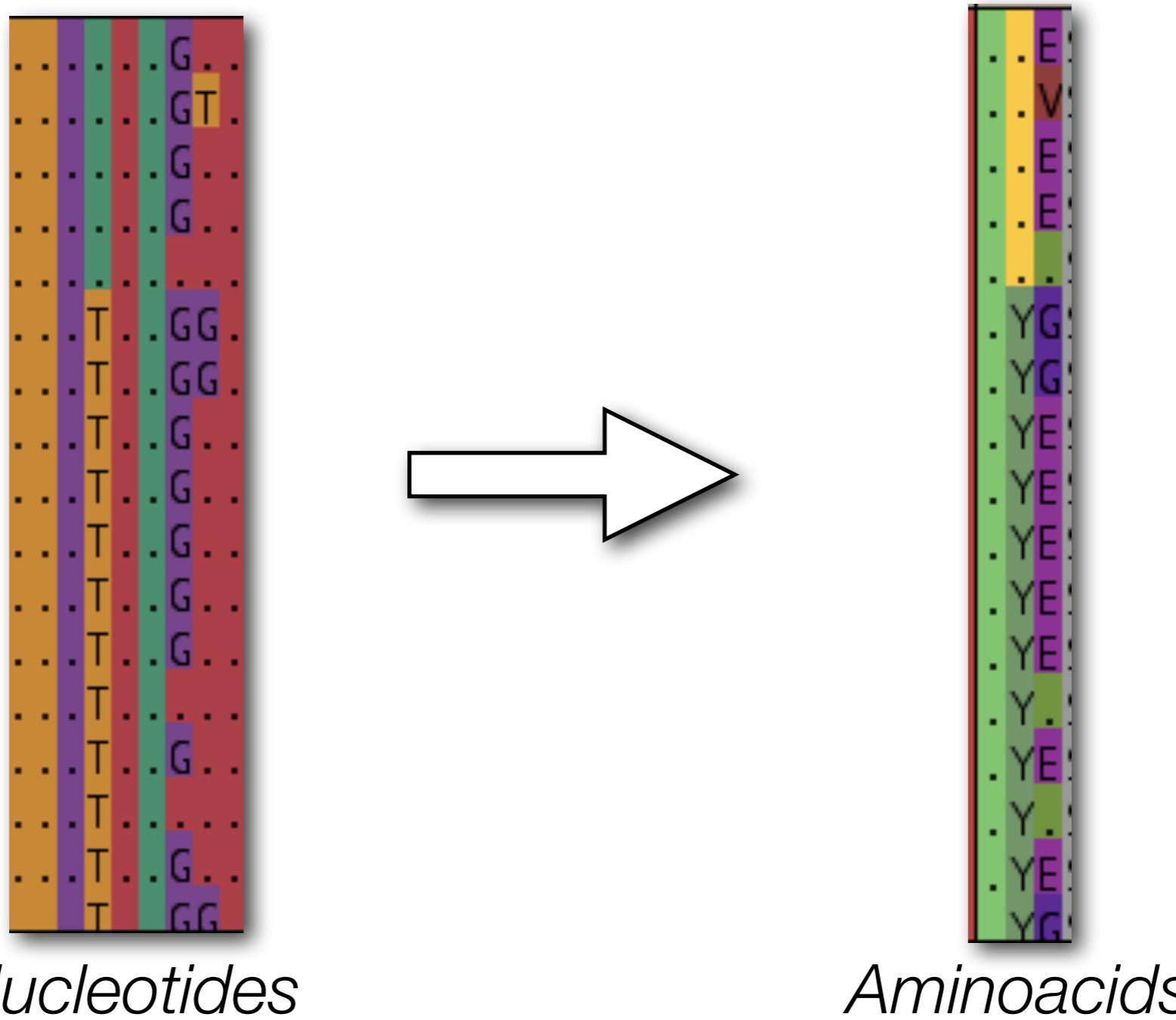
Nucleotides



Aminoacids

# Diversification

# An antigenic site in H3N2 IAV hemagglutinin



# Molecular signatures of selection

---

- Because synonymous substitutions do not alter the protein, we often posit that they are neutral
- The **rate** of accumulation of synonymous substitutions (**dS**) gives the neutral background
- We can compare the **rate** of accumulation of non-synonymous substitutions (**dN**), which alter the protein sequence, to classify the nature of the evolutionary process

$$dS \sim \frac{\text{number of fixed synonymous mutations}}{\text{proportion of random mutations that are synonymous}}$$

$$dN \sim \frac{\text{number of fixed non-synonymous mutations}}{\text{proportion of random mutations that are non-synonymous}}$$

# Evolutionary Modes

---

Positive Selection  
(Diversifying)

$dS < dN$  or  
 $\omega := dN/dS > 1$

Negative Selection

$dS > dN$  or  $\omega < 1$

Neutral Evolution

$dS \approx dN$  or  $\omega \approx 1$

# Estimating dS and dN

---

Consider two aligned homologous sequences

ACA	ATA	ATC	TTT	AAT	CAA
<i>T</i>	<i>I</i>	<i>I</i>	<i>F</i>	<i>N</i>	<i>Q</i>
<hr/>					
ACA	ATA	ACC	TTT	AAC	CAA
<i>T</i>	<i>I</i>	<b><i>T</i></b>	<i>F</i>	<i>N</i>	<i>Q</i>

# Estimating dS and dN

---

Consider two aligned homologous sequences

ACA	ATA	ATC	TTT	AAT	CAA
<i>T</i>	<i>I</i>	<i>I</i>	<i>F</i>	<i>N</i>	<i>Q</i>
<hr/>					
ACA	ATA	ACC	TTT	AAC	CAA
<i>T</i>	<i>I</i>	<b><i>T</i></b>	<i>F</i>	<i>N</i>	<i>Q</i>

Can one claim that **dN/dS = 1**, because there is **one** synonymous and **one** non-synonymous substitution?

# Universal genetic code

This genetic code has 61 sense (non-termination) codons

## Substitution types

	Synonymous			I	Non-synonymous			I	To a stop codon	Total
	Transitions	Transversions	Total		Transitions	Transversions	Total		Total	
1st position:	8	0	8		140		26	166	9	
2nd position:	0	0	0		148		28	176	7	
3rd position:	58	68	126		2		48	50	7	
<hr/>										
Total	66	68	134		290		102	392	23	

- Approximately 3:1 (392 N : 134 S) ratio when mutations are generated and **fixed** at random
- Non-random distribution over codon positions
  - All second position mutations are non-synonymous
  - Most synonymous mutations are confined to the third position

# Neutral expectation

---

- A random mutation is **~3 times more likely to be non-synonymous than synonymous**, depending on the variety of factors, such as codon composition, transition/transversion ratios, etc.
- We need to estimate the proportion of random mutations that are synonymous, and use it as a reference to compute **dS**.
- In early literature, these quantities were codified as synonymous and non-synonymous “sites” and/or mutational opportunity.
- As a very crude approximation (assuming that third positions ~ synonymous), each codon has 1 synonymous and 2 non-synonymous sites.

# Computing synonymous and non-synonymous sites for GAA (Glutamic Acid)

---

Start codon:	G	A	A
Site/Change to	1	2	3
A	AAA <b>Lysine</b>	*	*
C	CAA <b>Glutamine</b>	GCA <b>Alanine</b>	GAC <b>Aspartic Acid</b>
G	*	GGA <b>Glycine</b>	GAG <b>Glutamic Acid</b>
T	TAA <b>Stop</b>	GTA <b>Valine</b>	GAT <b>Aspartic Acid</b>
Synonymous changes	0	0	1
Non-synonymous changes	3	3	2
Synonymous sites	0	0	<b>1/3</b>
Non-synonymous sites	<b>1</b>	<b>1</b>	<b>2/3</b>

# Computing synonymous and non-synonymous sites for GAA (Glutamic Acid)

Start codon:	G	A	A
Site/Change to	1	2	3
<b>A</b>	AAA <b>Lysine</b>	*	*
<b>C</b>	CAA <b>Glutamine</b>	GCA <b>Alanine</b>	GAC <b>Aspartic Acid</b>
<b>G</b>	*	GGA <b>Glycine</b>	GAG <b>Glutamic Acid</b>
<b>T</b>	TAA <b>Stop</b>	GTA <b>Valine</b>	GAT <b>Aspartic Acid</b>
Synonymous changes	0	0	1
Non-synonymous changes	3	3	2
Synonymous sites	0	0	1/3
Non-synonymous sites	1	1	2/3

AMINOACID	CODONS	REDUNDANCY
ALANINE	GC*	4
CYSTEINE	TGC,TGT	2
ASPARTIC ACID	GAC,GAT	2
GLUTAMIC ACID	GAA,GAG	2
PHENYLALANINE	TTC,TTT	2
GLYCINE	GG*	4
HISTIDINE	CAC,CAT	2
ISOLEUCINE	ATA,ATC,ATT	3
LYSINE	AAA,AAG	2
LEUCINE	CT*,TTA,TTG	6
METHIONINE	ATG	1
ASPARGINE	AAC,AAT	2
PROLINE	CC*	4
GLUTAMINE	CAA,CAG	2
ARGININE	AGA,AGG,CG*	6
SERINE	AGC,AGT,TC*	6
THREONINE	AC*	4
VALINE	GT*	4
TRYPTOPHAN	TGG	1
TYROSINE	TAC,TAT	2
STOP	TAATAG,TGA	3

# Computing synonymous and non-synonymous sites for GAA (Glutamic Acid)

Start codon:	G	A	A
Site/Change to	1	2	3
<b>A</b>	AAA <b>Lysine</b>	*	*
<b>C</b>	CAA <b>Glutamine</b>	GCA <b>Alanine</b>	GAC <b>Aspartic Acid</b>
<b>G</b>	*	GGA <b>Glycine</b>	GAG <b>Glutamic Acid</b>
<b>T</b>	TAA <b>Stop</b>	GTA <b>Valine</b>	GAT <b>Aspartic Acid</b>
Synonymous changes	0	0	1
Non-synonymous changes	3	3	2
Synonymous sites	0	0	1/3
Non-synonymous sites	1	1	2/3

AMINOACID	CODONS	REDUNDANCY
ALANINE	GC*	4
CYSTEINE	TGC,TGT	2
ASPARTIC ACID	GAC,GAT	2
GLUTAMIC ACID	GAA,GAG	2
PHENYLALANINE	TTC,TTT	2
GLYCINE	GG*	4
HISTIDINE	CAC,CAT	2
ISOLEUCINE	ATA,ATC,ATT	3
LYSINE	AAA,AAG	2
LEUCINE	CT*,TTA,TTG	6
METHIONINE	ATG	1
ASPARGINE	AAC,AAT	2
PROLINE	CC*	4
GLUTAMINE	CAA,CAG	2
ARGININE	AGA,AGG,CG*	6
SERINE	AGC,AGT,TC*	6
THREONINE	AC*	4
VALINE	GT*	4
TRPTOPHAN	TGG	1
TYROSINE	TAC,TAT	2
STOP	TAA,TAG,TGA	3

8/3 non-synonymous sites  
1/3 synonymous sites

# Nei-Gojobori dN/dS estimate (NG86)

~4,000 citations

Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions

M. Nei and T. Gojobori  
*Mol. Biol. Evol.* **3** 418–426 (1986)

- For each codon  $C$  we define  $ES(C)$  and  $EN(C)$  - the numbers of synonymous and non-synonymous **sites** of a codon
  - e.g.,  $ES(GAA) = 1/3$ ,  $EN(GAA) = 8/3$ .
- May also define them as fractions of substitutions that do not lead to stop codons,
  - e.g.,  $ES(GAA) = 1/3$ ,  $EN(GAA) = 7/3$ .
- The sum of  $ES$  and  $EN$  over all codons in a sequence gives an estimate of expected synonymous and non-synonymous **sites** in a sequence.
- For two sequences (the target of the original method), we average  $ES(C)$  and  $EN(C)$  at each site.
- $EN/ES$  is thus the ***expected ratio of non-synonymous to synonymous substitutions counts under neutral evolution***

# NG86 example

<b>Seq1</b>	<b>ACA</b>	<b>ATA</b>	<b>ATC</b>	<b>TTT</b>	<b>AAT</b>	<b>CAA</b>
Syn	1	2/3	2/3	1/3	1/3	1/3
NonSyn	2	7/3	7/3	8/3	8/3	7/3
<b>Seq2</b>	<b>ACA</b>	<b>ATA</b>	<b>ACC</b>	<b>TTT</b>	<b>AAC</b>	<b>CAA</b>
Syn	1	2/3	1	1/3	1/3	1/3
NonSyn	2	7/3	2	8/3	8/3	7/3
Mean						
Syn	1	2/3	5/6	1/3	1/3	1/3
NonSyn	2	7/3	13/6	8/3	8/3	7/3

**ES** =  $3\frac{1}{2}$ , **EN** =  $14\frac{1}{6}$ : under neutrality, would expect the ratio of non-synonymous to synonymous substitutions of **EN/ES** ~ 4

# NG86 example

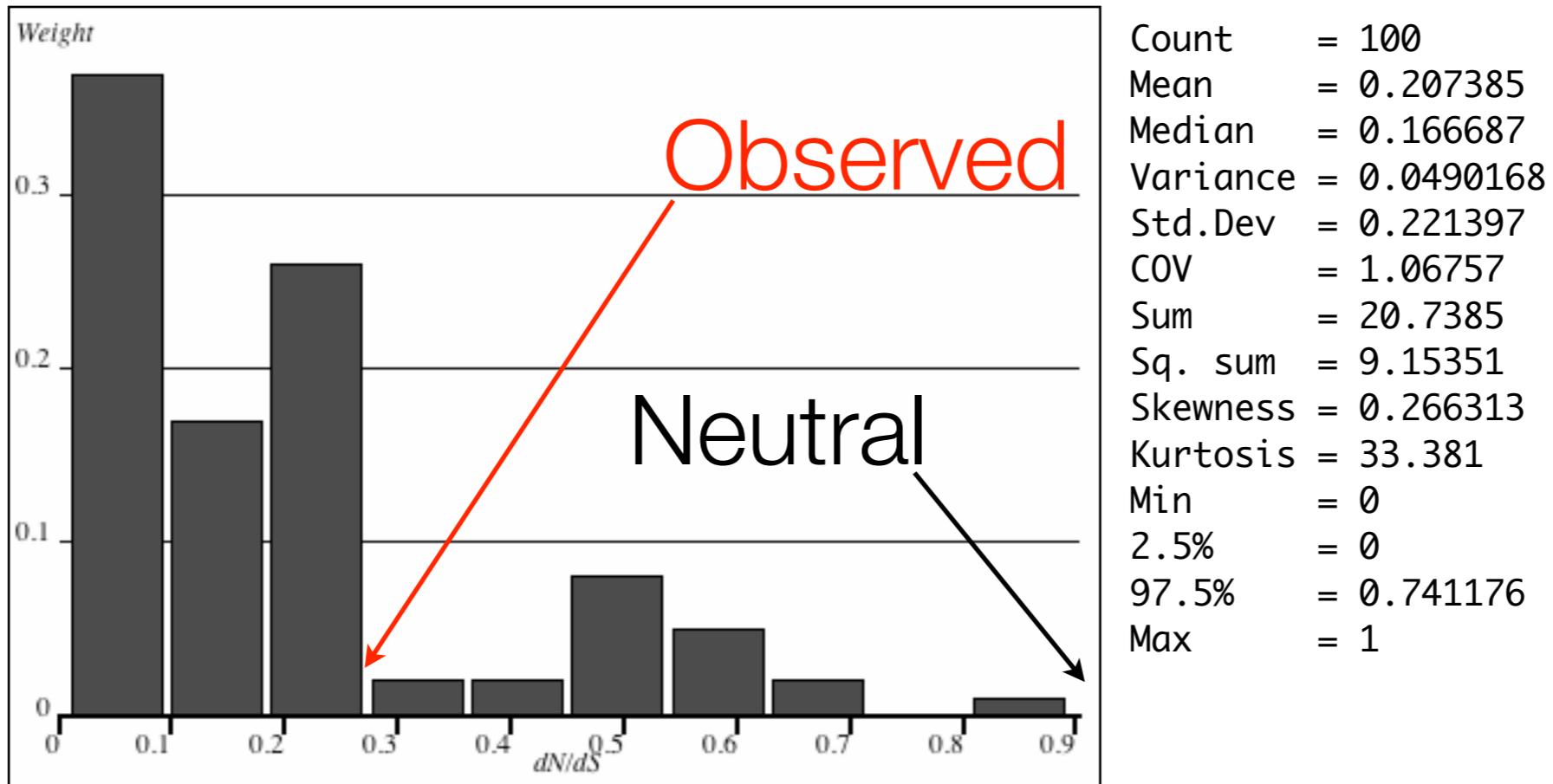
---

- The observed **N/S** ratio (1 . 0) is **lower** than the expected **EN/ES** ratio (4 . 05).
- The ratio of the ratios **(N:S) / (EN:ES)** yields  $dN/dS = 1/4.05 \sim 0.25$ .
- This ratio quantifies the excess or *paucity* of non-synonymous substitutions and is near  $dN/dS = 1$  for neutrally evolving sequences/sites.
- Because there are **fewer non-synonymous substitutions than expected**, we conclude that most non-synonymous mutations are removed by natural selection, i.e., the sequences are under **negative selection**.

# NG86 example

- How reliable is the inference based on only 6 codons?
- Obtain sampling variance via bootstrap (or by limiting approximations)
- In this case,  $dN/dS$  is **significantly** less than 1.0 ( $p \sim 0.01$ )

## Bootstrapped distribution of $dN/dS$

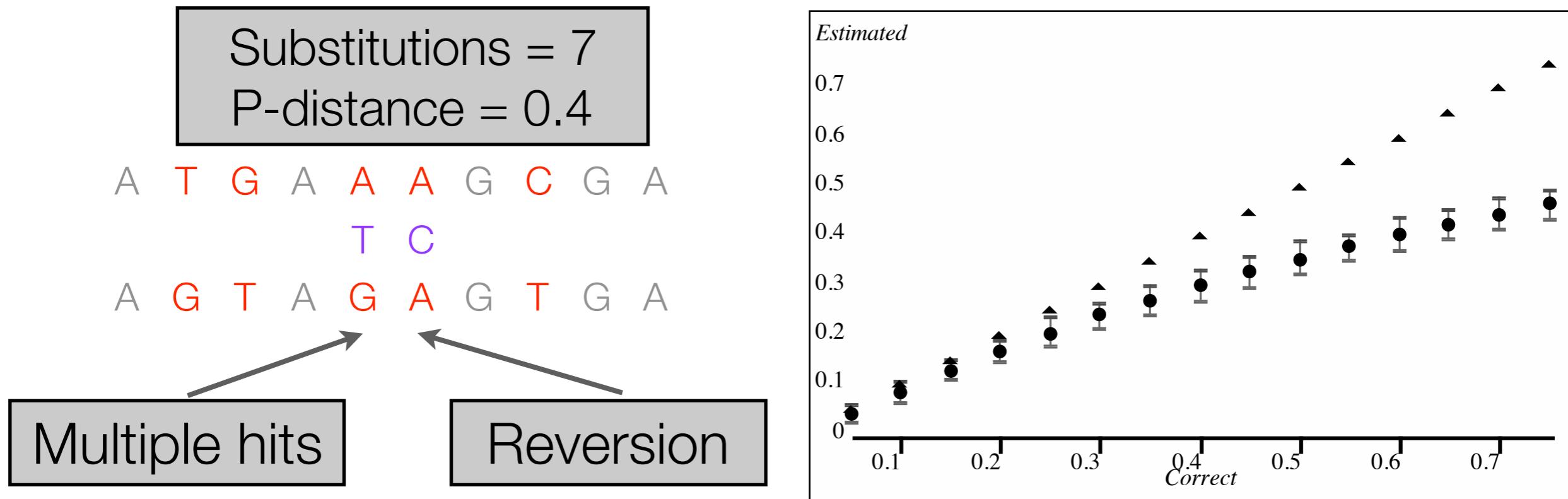


# NG86 limitations: multiple substitutions

---

- How many synonymous and how many non-synonymous substitutions does it take to replace **CCA** with **CAG**?
- **Assume** the shortest path (minimum of 2 substitutions)
  - CCA (Proline)  $\Rightarrow$  CAA (Histidine)  $\Rightarrow$  CAG (Glutamine)
  - CCA (Proline)  $\Rightarrow$  CCG (Proline)  $\Rightarrow$  CAG (Glutamine)
- Average over the two possible paths: **0.5** synonymous and **1.5** non-synonymous substitutions.
- Intuitively, paths should **not** be equiprobable, e.g., because it should be more expensive to route evolution through (presumably) suboptimal intermediate aminoacids.

# NG86 limitations: underestimation of substitution counts for higher divergence levels



- Simulated 100 replicates of 1000 nucleotide long sequences for various divergence levels (substitutions/site)
- Plotted simulated divergence vs that estimated by p-distance.
- Even for divergence of 0.25 (1/4 sites have mutation on average), p-distance already significantly underestimates the true level: 0.2125 (0.19–0.241 95% range)
- Underestimation becomes progressively worse for larger divergence levels

# NG86 limitations: ignoring phylogenies

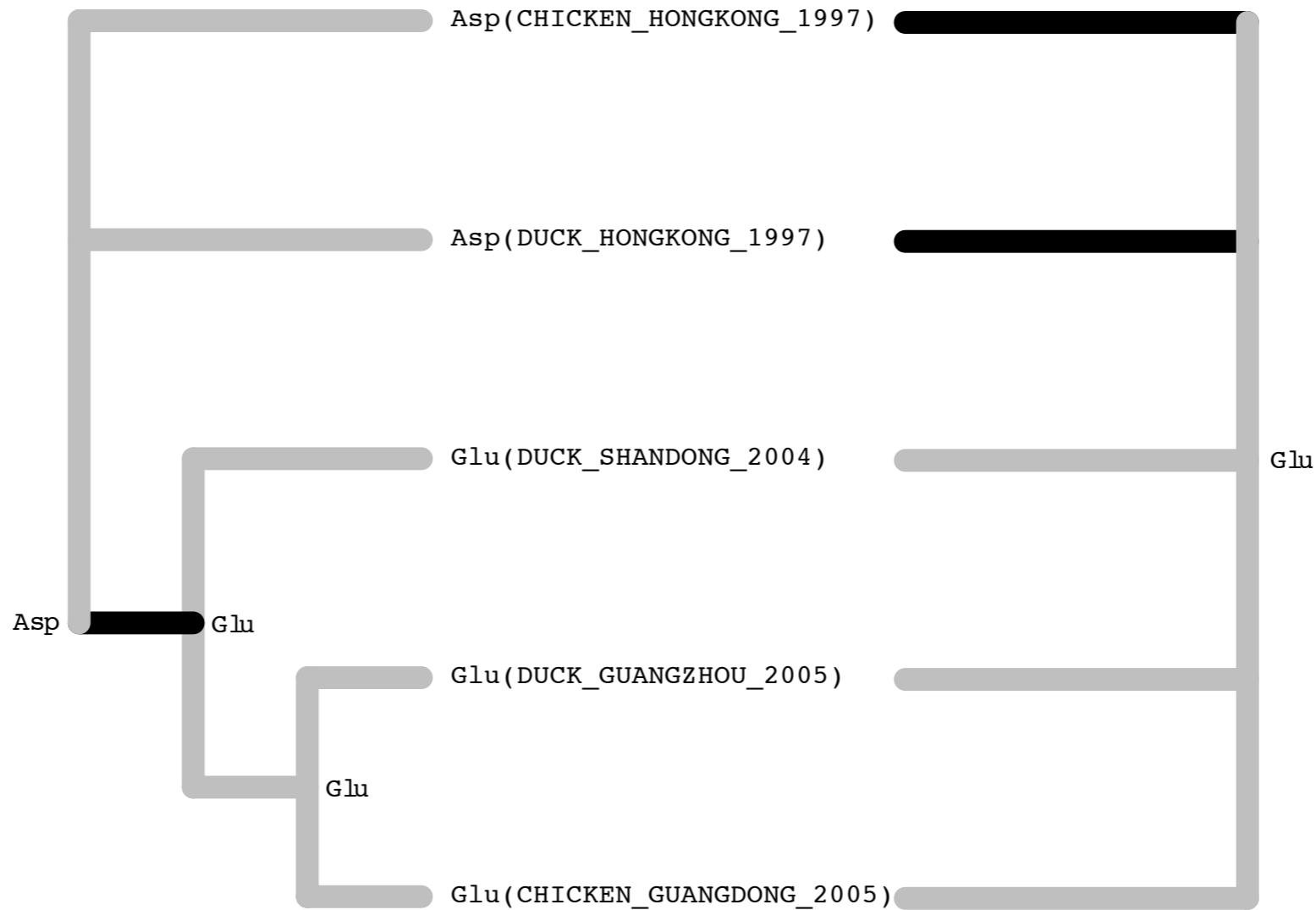


Fig. 1.1. Effect of phylogeny on estimating synonymous and nonsynonymous substitution counts in a dataset of Influenza A/H5N1 haemagglutinin sequences. Using the maximum likelihood tree on the left, the observed variation can be parsimoniously explained with one nonsynonymous substitution along the darker branch, whereas the star tree on the right involves at least two.

# NG86 limitations: averaging across all sites in a gene

---

- Different sites in a gene will be subject to different selective forces.
- A *gene-wide* measure of selection is going to average these effects.
- **Most** sites in **most** genes will be maintained by purifying selection.
- Positively selected sites are of great biological interest, because they point to how a particular gene can respond to selective pressures.
- Negatively selected sites are also of interest, because they point to functional constraint, and could be used to guide drug or vaccine design.
- Must develop methods that are able to disentangle the contributions of individual sites.

~450 citations

A method for detecting positive selection at single amino acid sites  
Y. Suzuki and T. Gojobori  
*Mol Biol Evol* **16** 1315-1328 (1999)

# Suzuki-Gojobori (SG99): the penultimate extension of NG86

---

## Uses a tree to compute dN/dS at a given site

1. Reconstruct ancestral sequences by nucleotide-level parsimony
2. Compute **EN** and **ES** using labeled branches; define  $p_e = ES/EN$
3. Compute **S** and **NS** for each site (minimum evolution)
4. Estimate the probability that the number of synonymous substitutions **S** is unusually low (positive selection) or unusually high (negative selection), using the binomial distribution given  $p_e$  from step 2.

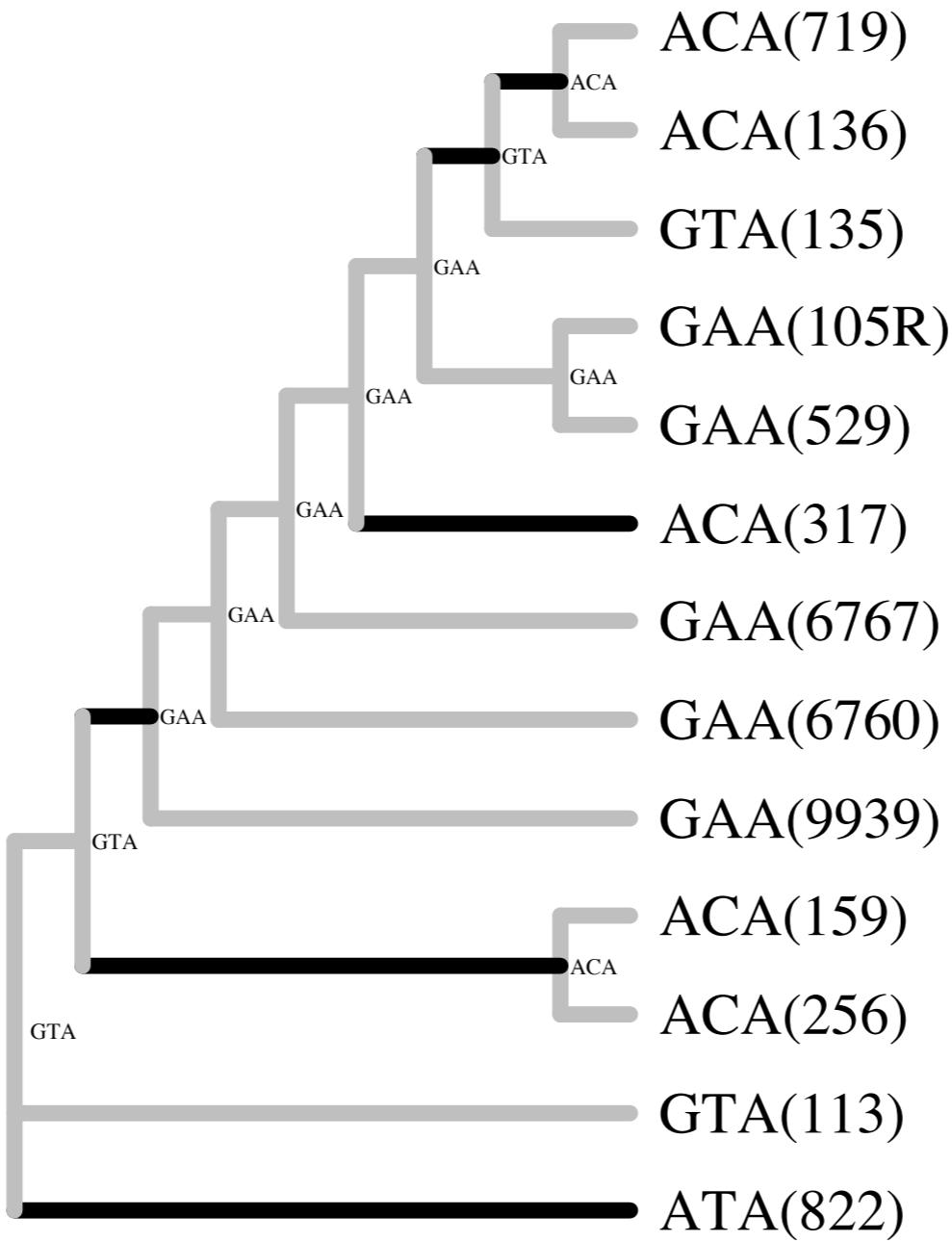


Fig. 1.6. An illustration of SLAC method, applied to a small HIV-1 envelope V3 loop alignment. Sequence names are shown in parentheses. Likelihood state ancestral reconstruction is shown at internal nodes. The parsimonious count yields 0 synonymous and 9 non-synonymous substitutions (highlighted with a dark shade) at that site. Based on the codon composition of the site and branch lengths (not shown), the expected proportion of synonymous substitutions is  $p_e = 0.25$ . An extended binomial distribution on 9 substitutions with the probability of success of 0.25, the probability of observing 0 synonymous substitutions is 0.07, hence the site is borderline significant for positive selection.

# Codon-substitution models

A codon-based model of nucleotide substitution for protein-coding DNA sequences.  
N. Goldman and Z. Yang  
*Mol Biol Evol* 11 725--736 (1994)

1620 citations

- In 1994, first tractable mechanistic evolutionary models for codon sequences were proposed by **Muse and Gaut** (MG94), and, independently, by **Goldman and Yang** (GY94) [in the same issue of MBE, back to back]
- Markov models of codon substitution provide a powerful framework for **estimating substitution rates** from coding sequence data, as they
  - *encode our mechanistic understanding of the evolutionary process,*
  - *enable one to compute phylogenetic likelihood,*
  - *permit hypothesis testing or Bayesian inference,*
  - *systematically account for confounding processes (unequal base frequencies, nucleotide substitution biases, etc.),*
  - *afford many opportunities for extension and refinement (still happening today).*

# Rate matrix for an MG-style codon model

$$(\text{Rate})_{X,Y}(dt) = \begin{cases} \alpha R_{xy} \pi_t dt & , \text{ one-step, synonymous substitution,} \\ \beta R_{xy} \pi_t dt & , \text{ one-step, non-synonymous substitution,} \\ 0 & , \text{ multi-step.} \end{cases}$$

X, Y = AAA...TTT (excluding stop codons),

$\pi_t$  - frequency of the target nucleotide.

Example substitutions:

AAC → AAT (one step, synonymous - Asparagine)

CAC → GAC (one step, non-synonymous - Histidine to Aspartic Acid)

AAC → GTC (multi-step).

$\alpha R_{CT}$

$\beta R_{CG}$

**α** (syn. rate) and **β** (non-syn. rate)  
are the key quantities for all selection analyses

# Computing the transition probabilities

---

- In order to recover transition probabilities  $T(t)$  from the rate matrix  $Q$ , one computes the matrix exponential  $T(t) = \exp(Qt)$ , same as with standard nucleotide models, e.g. HKY85 or GTR
- Because the computational complexity of matrix exponentiation scales as the cube of the matrix dimension, codon based models require roughly  $(61/4)^3 \approx 3500$  more operations than nucleotide models
- This explains why codon probabilistic models were not introduced until the 1990s, even though they are relatively straightforward extensions of 4x4 nucleotide models

# Multiple substitutions

---

- The model assumes that point mutations alter one nucleotide at a time, hence most of the instantaneous rates ( $3134/3761$  or  $84.2\%$  in the case of the universal genetic code) are 0.
- This restriction, however, does not mean that the model disallows any substitutions that involve multiple nucleotides (e.g.,  $\text{ACT} \Rightarrow \text{AGG}$ ).
- Such substitutions must simply be realized via several single nucleotide steps, e.g  $\text{ACT} \Rightarrow \text{AGT} \Rightarrow \text{AGG}$
- In fact the  $(i, j)$  element of  $T(t) = \exp(Qt)$  sums the probabilities of all such possible pathways of duration  $t$ , including reversions
- Compare this to the naive NG86 parsimony approach.

# Alignment-wide estimates

---

- Using standard MLE approaches it is straightforward to obtain point estimates of  $dN/dS := \beta/\alpha$
- Can also easily test whether or not  $dN/dS > 1$ , or  $< 1$  using the likelihood ratio test (LRT)
- Codon models also support the concepts of synonymous and non-synonymous distances between sequences using standard properties of Markov processes (exponentially distributed waiting times)

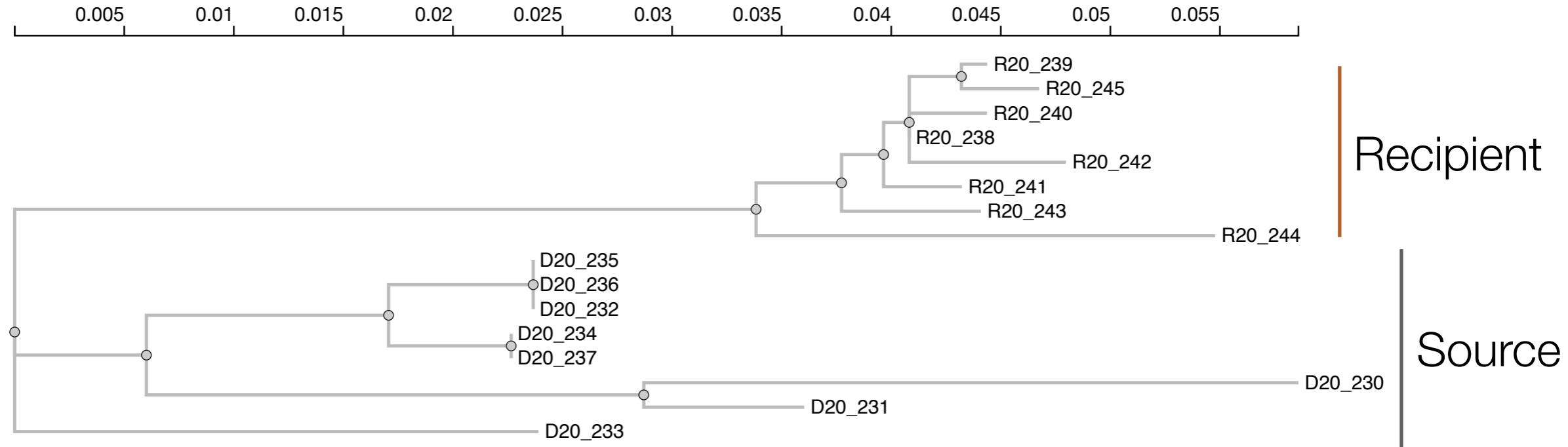
$$E[\text{subs}] = - \sum_i \pi_i \hat{q}_{ii}, \quad E[\text{subs}] = E[\text{syn}] + E[\text{nonsyn}] = - \sum_i \pi_i \hat{q}_{ii}^s - \sum_i \pi_i \hat{q}_{ii}^{ns}.$$

# Two example datasets

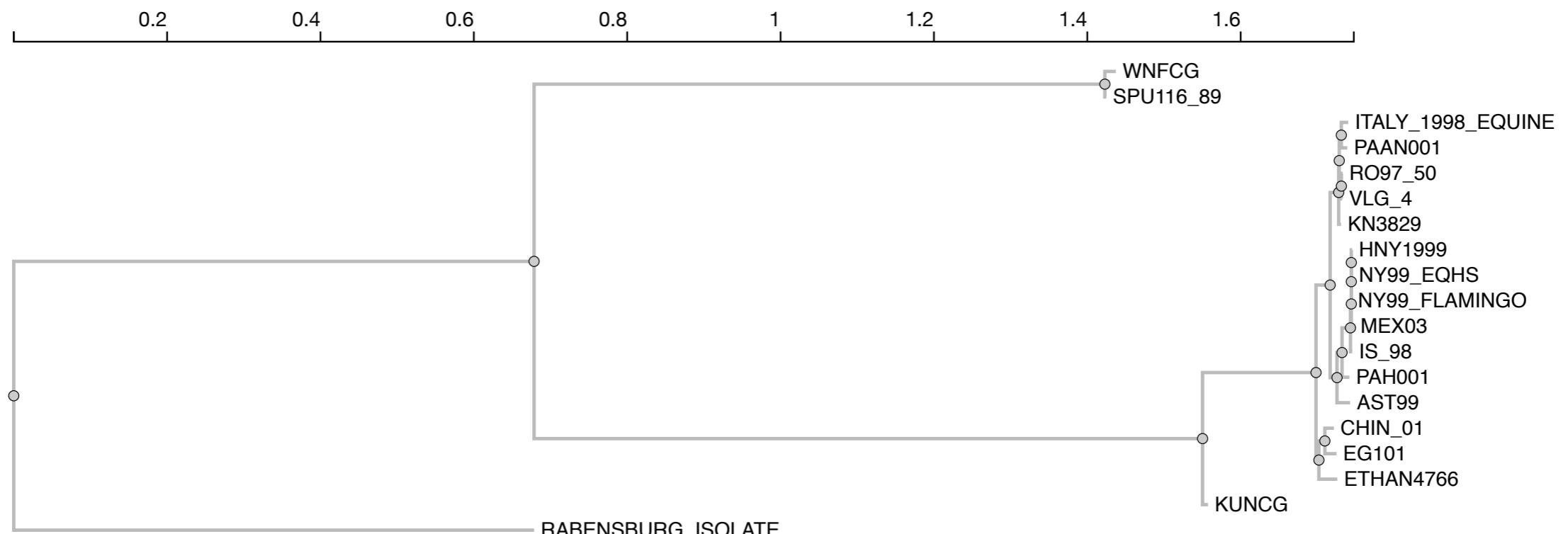
---

- **West Nile Virus NS3 protein**
  - An interesting case study of how positive selection detection methods lead to testable hypotheses for function discovery
  - Brault et al 2007, *A single positively selected West Nile viral mutation confers increased virogenesis in American crows*
- **HIV-1 transmission pair**
  - Partial *env* sequences from two epidemiologically linked individuals
  - An example of multiple selective environments (source, recipient, transmission)

# HIV-1 env



# WN NS3



# Information content of the alignments

	WNV NS3	HIV-1 <i>env</i>
Sequences	19	16
Codons	619	288
Tree Length <i>MG94 model, subs/site</i>	3.32	0.20

How do you expect these measures to correlate with the ability to detect selection?

## WNV NS3

Model	Log L	# p	dN/dS	LRT	p-value
Null	-7668.7	49	1		
Alternative	-6413.5	50	0.009	2510.4	~0

*Very strongly conserved*

## HIV-1 env

Model	Log L	# p	dN/dS	LRT	p-value
Null	-2078.3	40	1		
Alternative	-2078.2	41	1.128	0.2	~0.6

*Not significantly different from neutral*

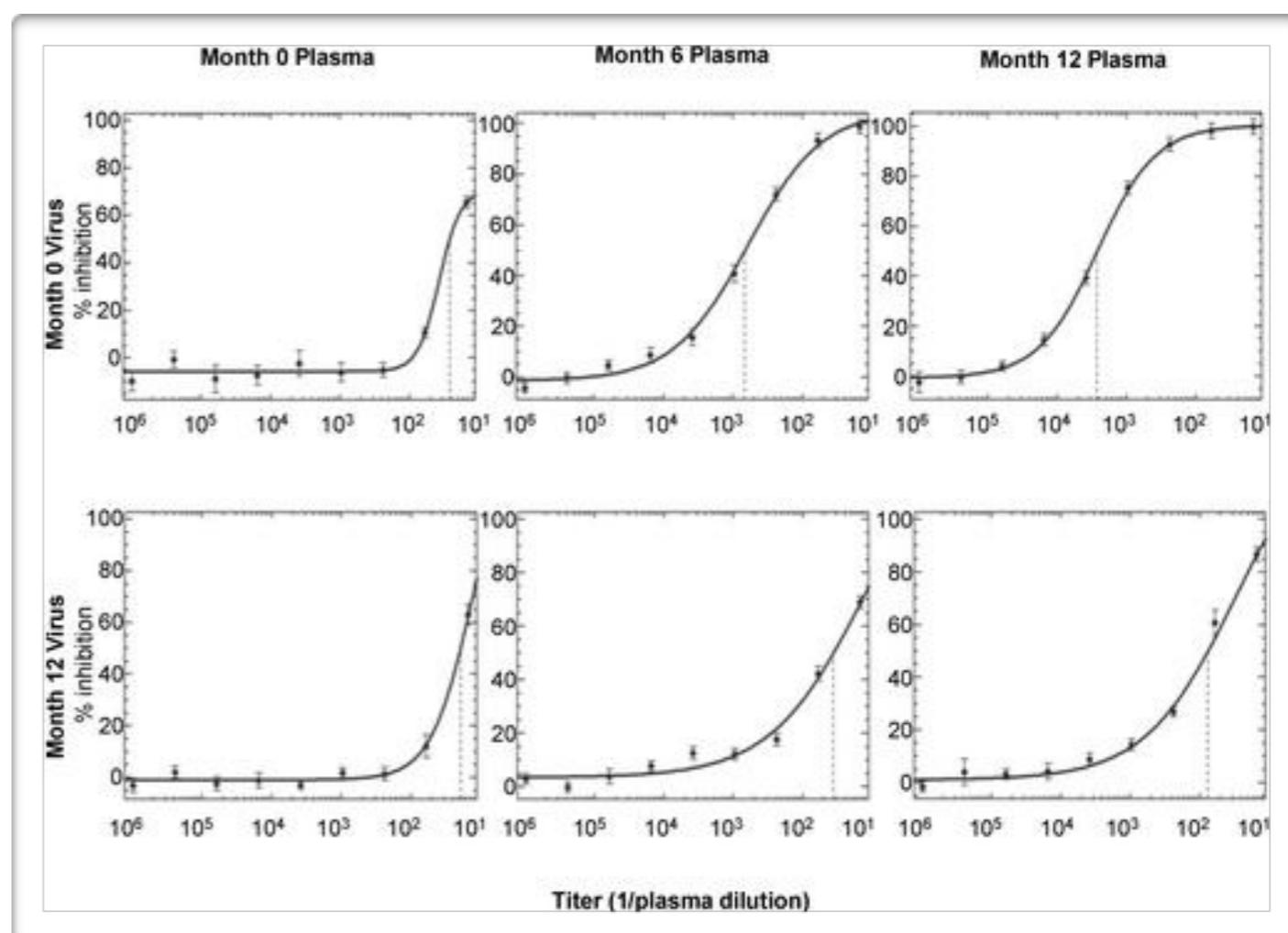
# Mean gene-wide dN/dS estimates

---

- Are not the way to go, **except** when you have very small (2-3 sequence) datasets
- For example:
  - The humoral arm of the immune system mounts a potent defense against viral infections
  - Existing successful vaccines are based on raising a neutralizing antibody (nAb) response to the pathogen
  - No simple host genetic basis (epitopes) of the specificity of neutralizing antibody responses is known
  - Need to measure these responses

# Neutralization curves from an individual with early HIV infection

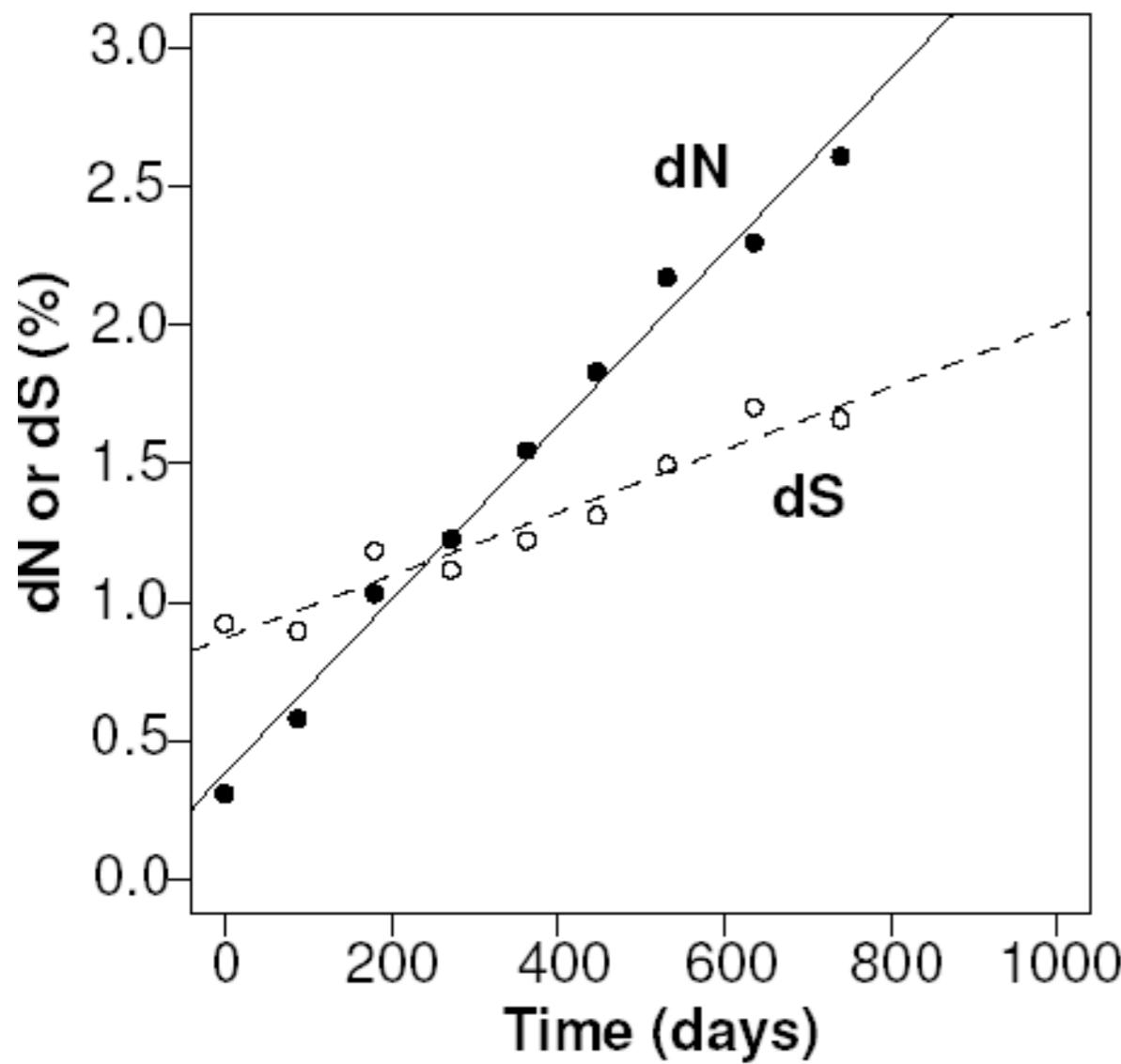
- Neutralization can be measured by the serum dilution needed to reduce viral replication by 50% (typically presented as the inverse of the titer)
- Although variable between individuals, the rate of escape from neutralizing antibodies can be very high during acute/early HIV infection
- Sera are effective at neutralizing earlier viruses, but significantly less effective at neutralizing contemporaneous viruses
- The immune system loses the arms race



# Amino acid substitutions in HIV-1 env accumulate faster during rapid escape

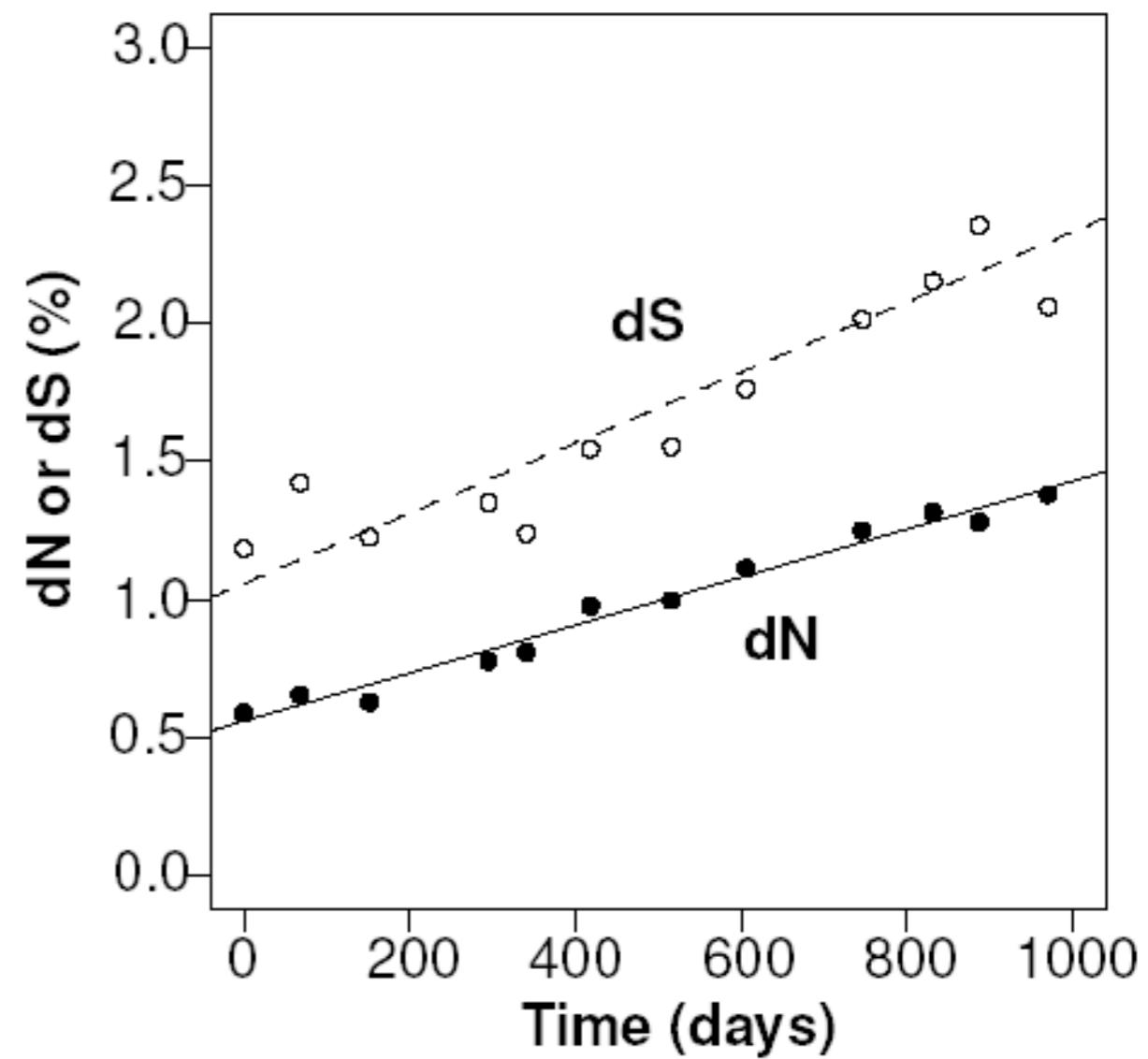
Patient 01–0127, rapid escape

(a)

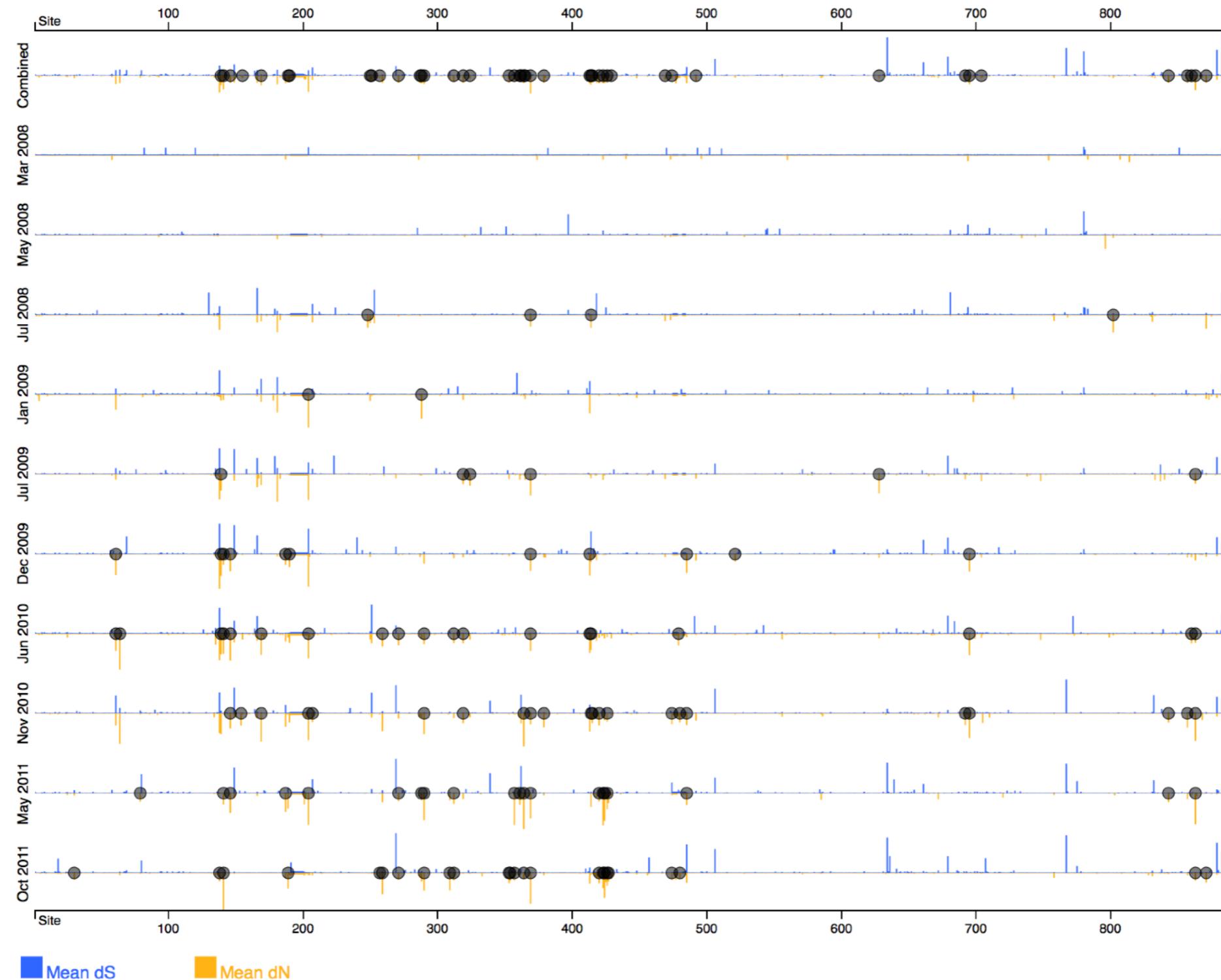


Patient 01–0083, slow escape

(b)



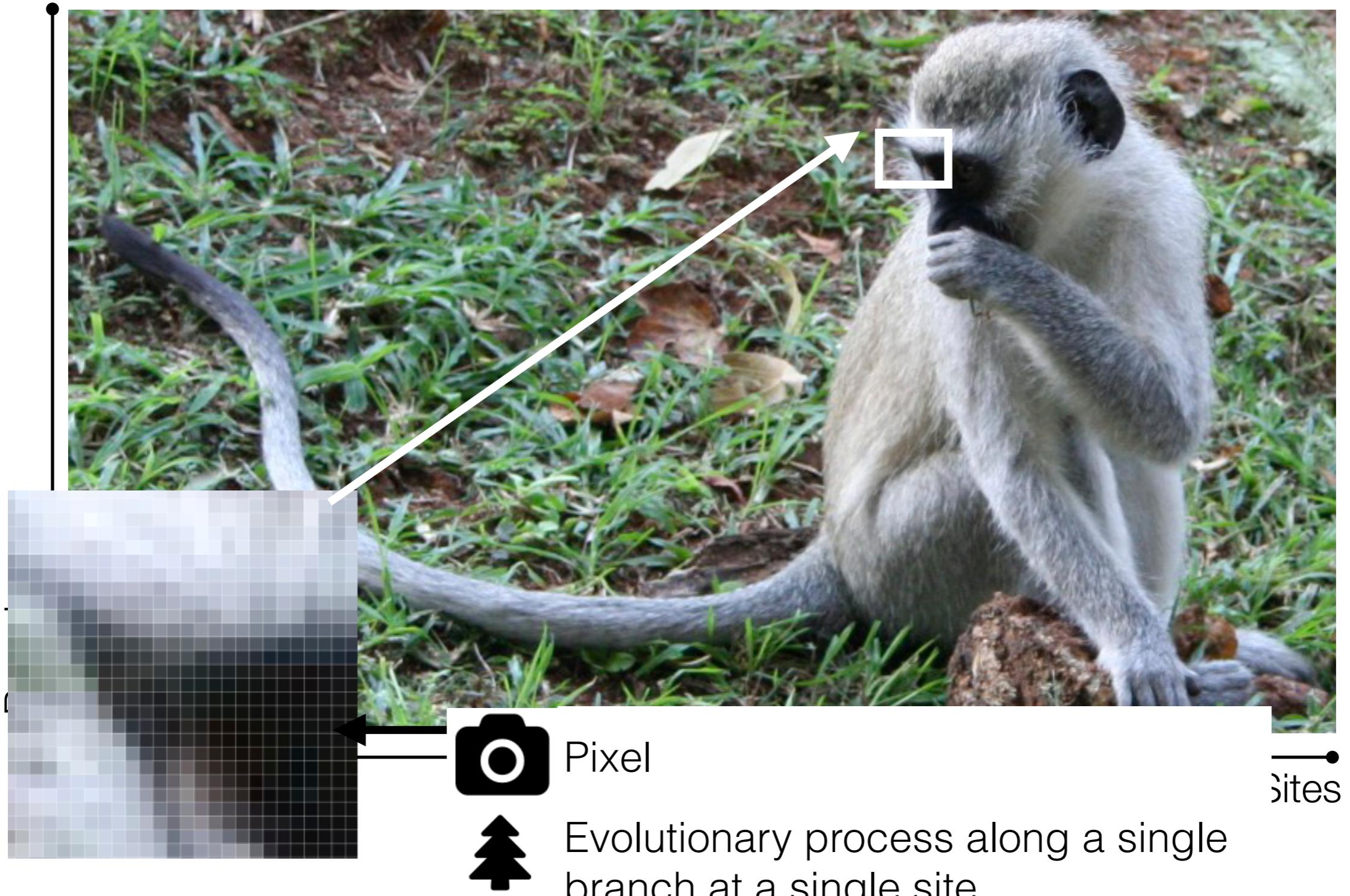
But upon closer look, this pattern is highly variable both across a gene and through time.



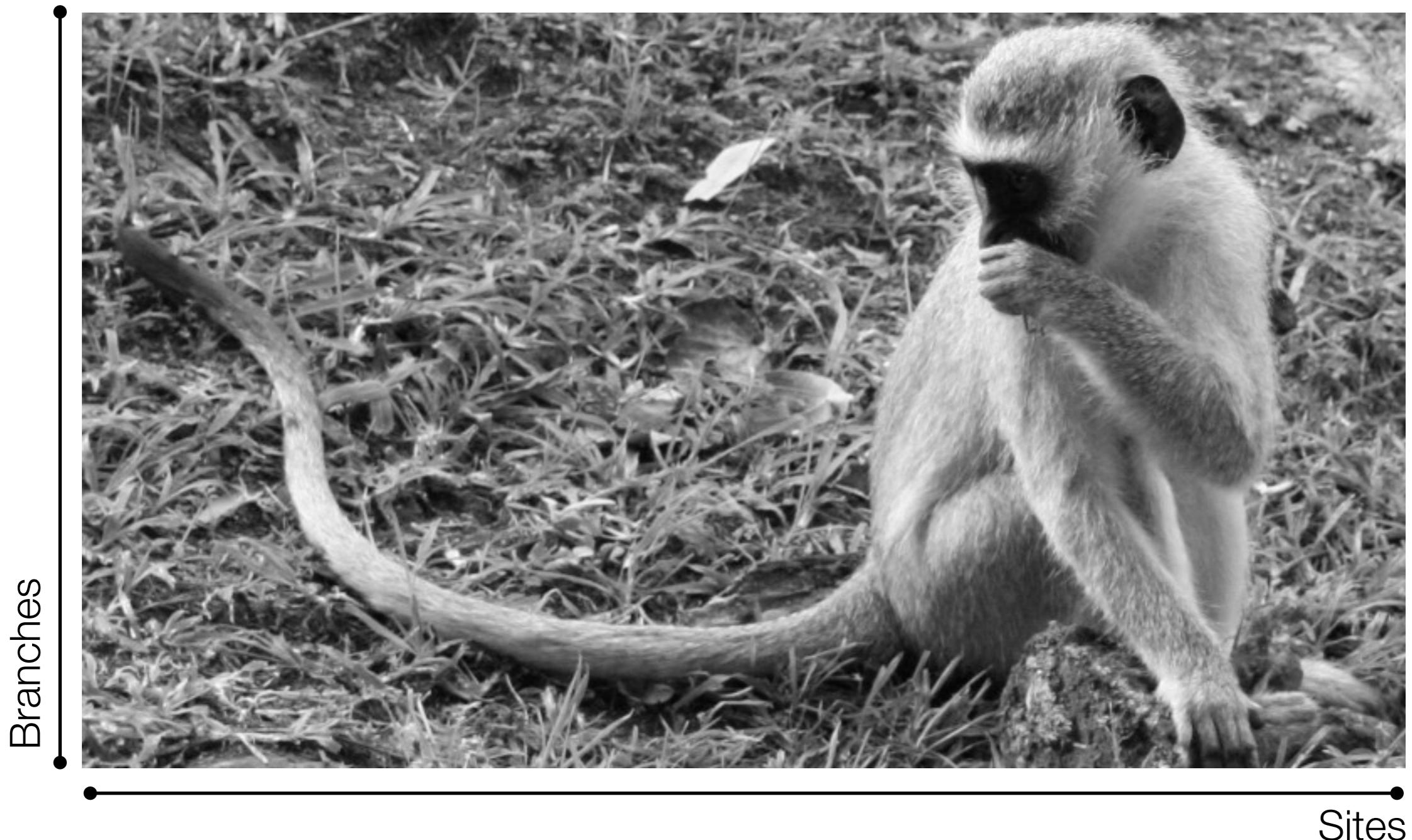
# Selection inference as image processing



# Selection inference as image processing



# Forget about the color



Intensity/brightness



Evolutionary rate ( $dN/dS$ )

Color

Type of evolutionary/  
function/property change

Evolution is largely unobserved and noisy

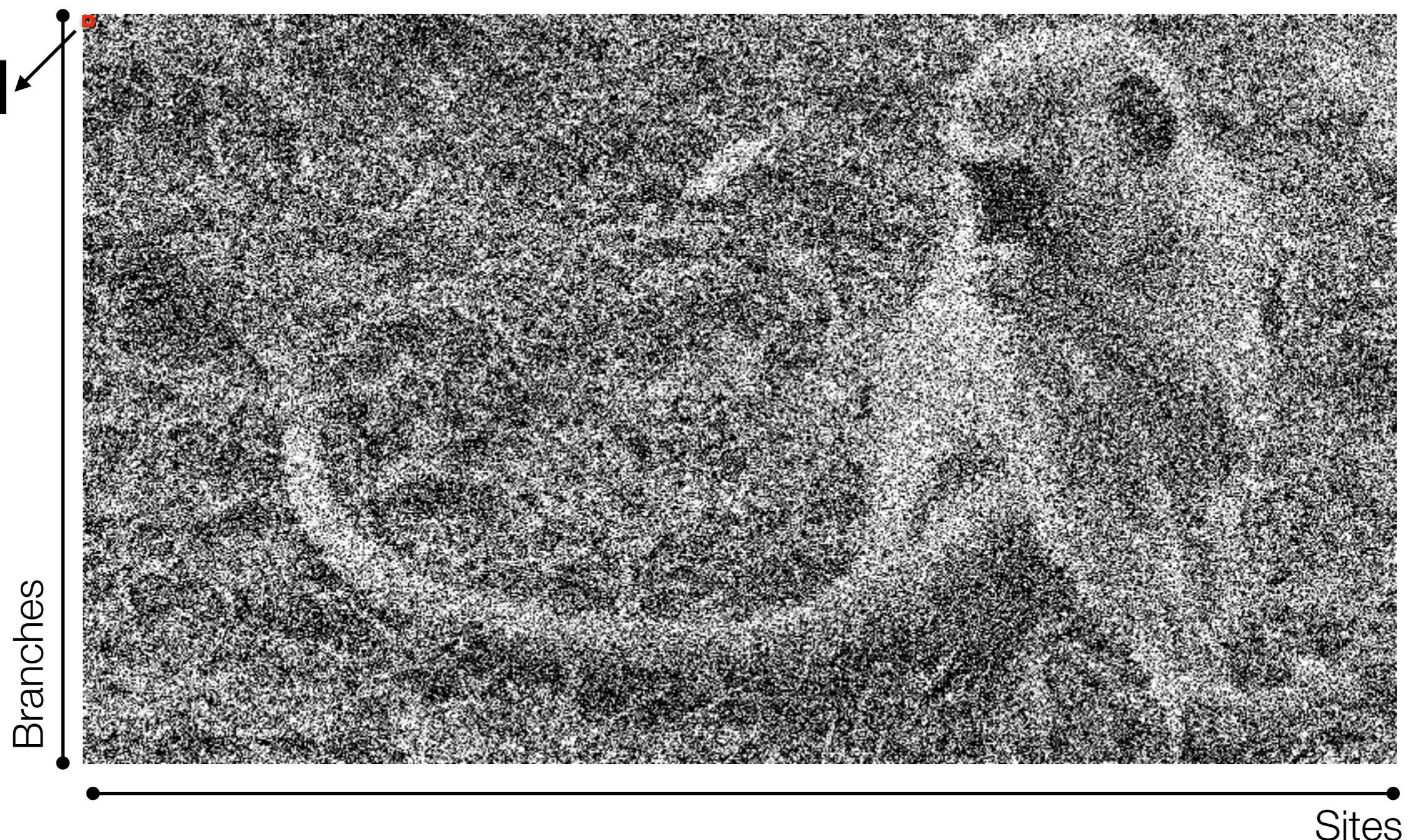


Visual noise



Saturation, missing data, model misspecification,  
sampling variation

Evolution is largely unobserved and noisy (another replicate)



Visual noise



Saturation, missing data, model misspecification,  
sampling variation

Evolution is largely unobserved and noisy (another replicate)



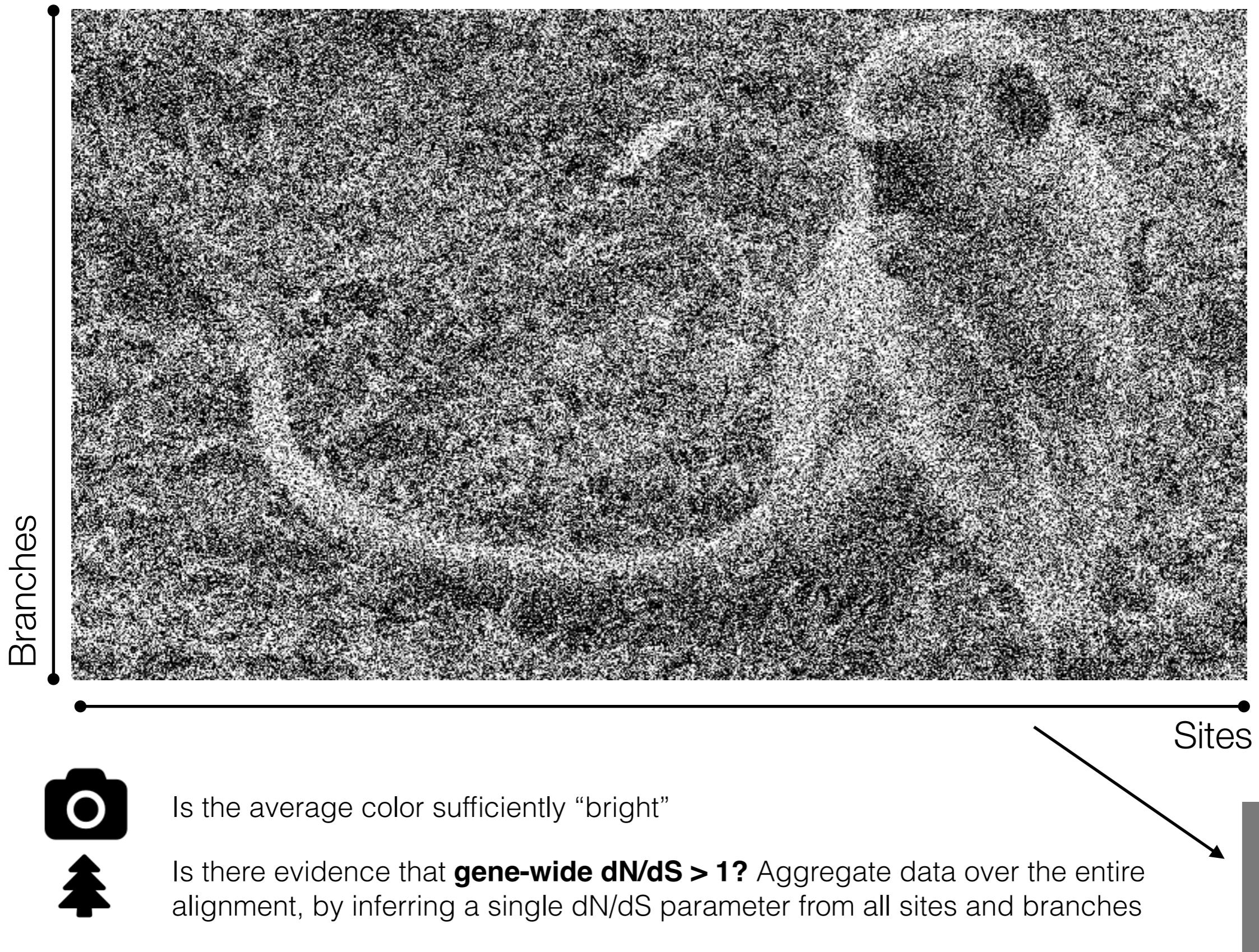
Visual noise



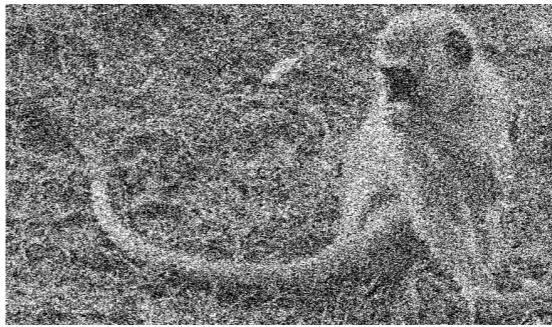
Saturation, missing data, model misspecification,  
sampling variation

- ➊ High local variability
  - ➋ Stable global (monkey) and local (head, tail) patterns, easily discernible
- 
- ☰ Desired resolution (branch-site) is not attainable
  - ☰ Global (and some local) patterns should be inferable and testable
  - ☰ Statistical inference draws power from sample (and effect) size, need to aggregate data to gain power

## Gene-wide selection (mean dN/dS)

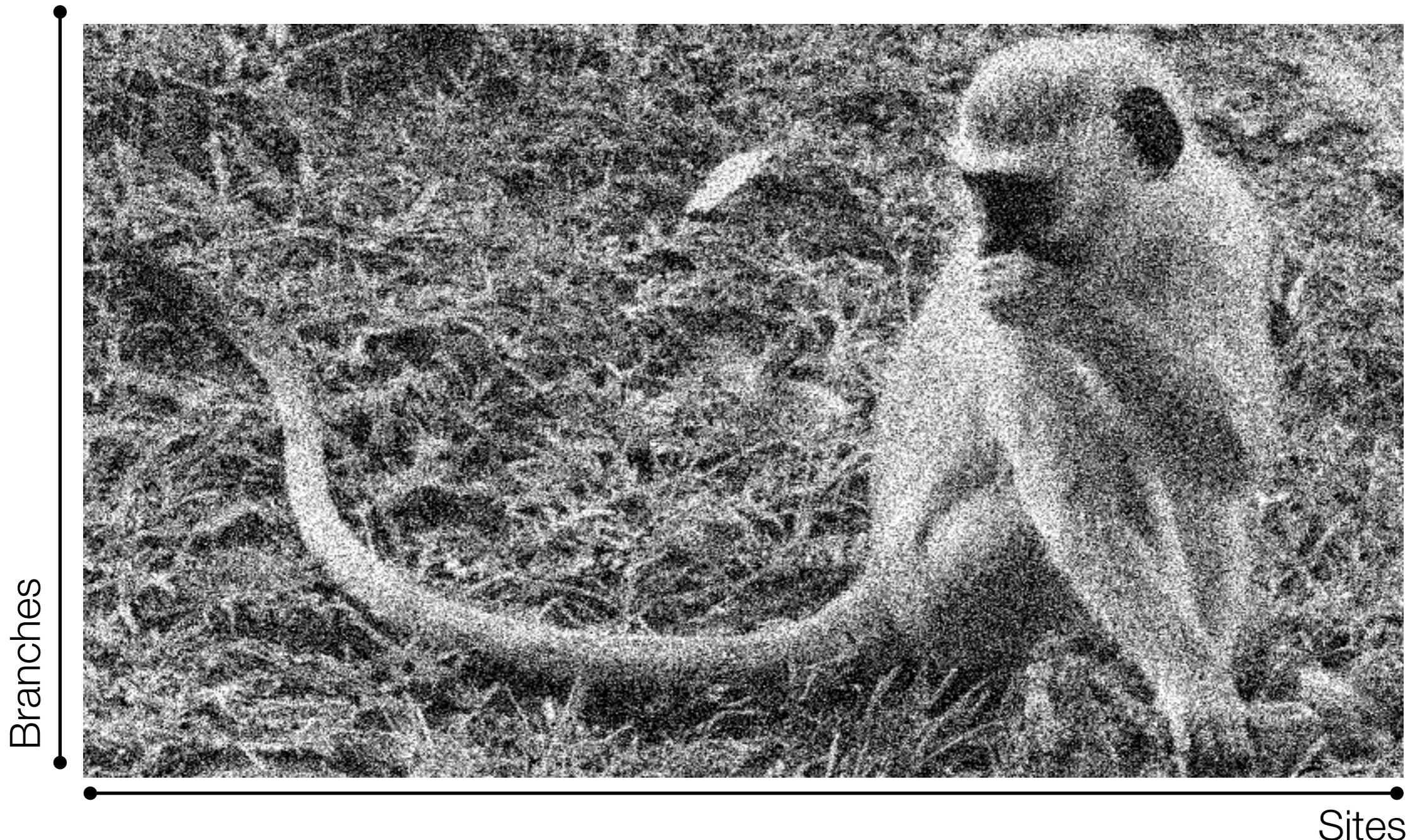


- Simple
  - single rate parameter
  - relatively compute-light
- Very robust to local variation
- Sample size  $\sim$  sites  $\times$  branches
- Very low power
  - most genes are **on average** conserved
- No resolution
  - if selection occurred, how much of the gene was involved, and when did it happen
- Rate variation model is definitely misspecified



# Gene-wide selection

random effects over sites and branches [BUSTED]



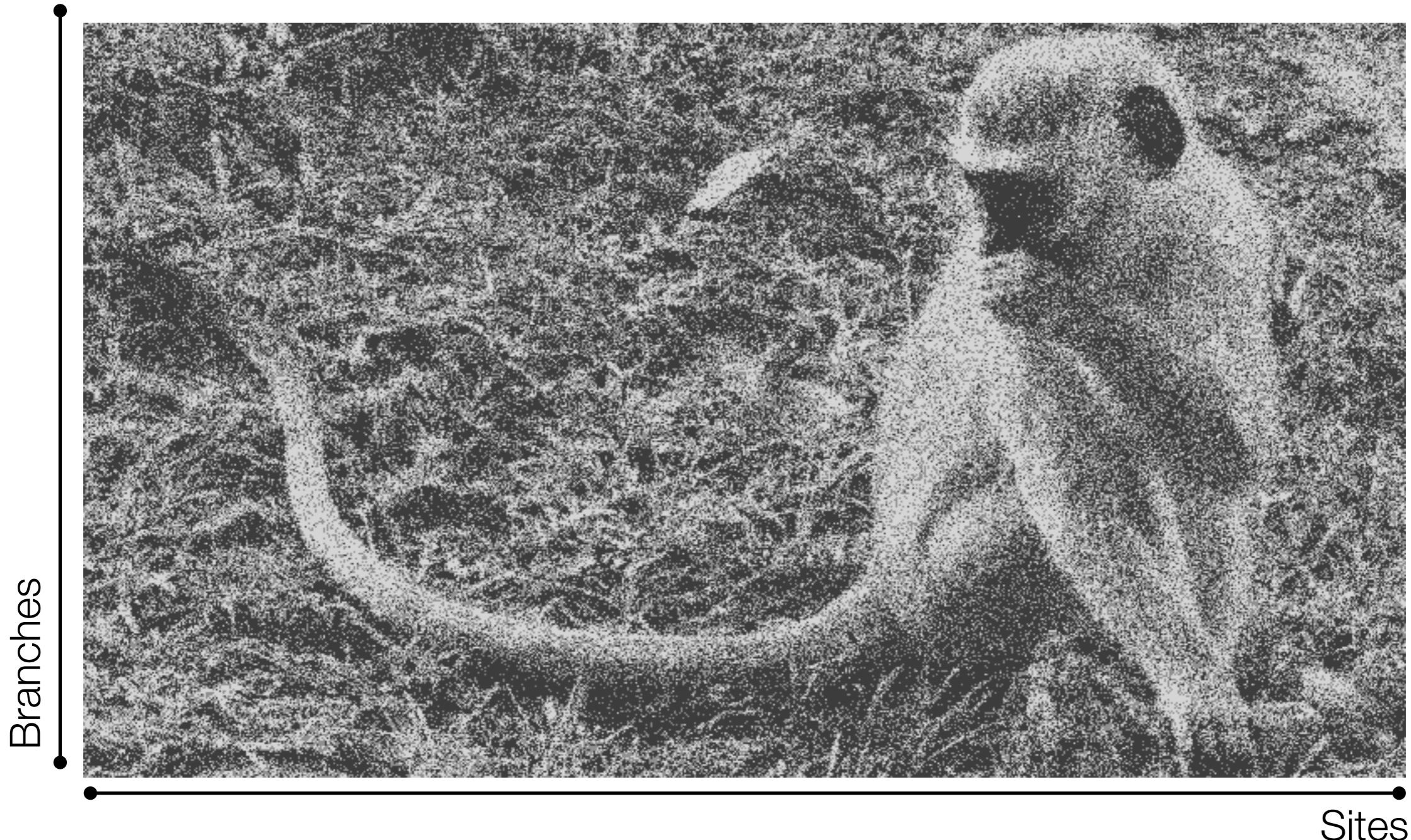
Is there enough **image area** that is sufficiently bright; allow each pixel to be one of 3 colors, chosen adaptively, e.g. to minimize perceptual differences



[BUSTED]: each branch-site combination is drawn from a 3-bin ( $dS, dN$ ) distribution. The distribution is estimated from the entire alignment. Tests if  $dN/dS > 1$  for some branch/site pairs in the alignment

# Gene-wide selection

random effects over sites and branches [BUSTED]



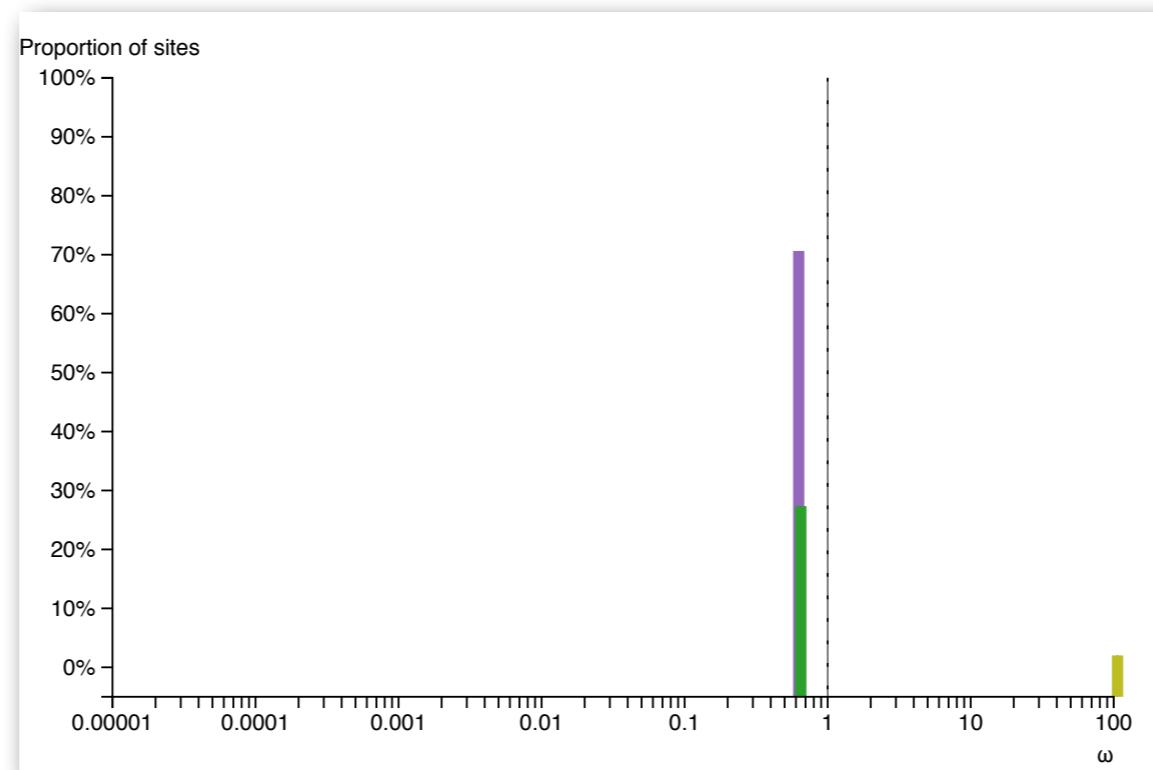
Is there enough **image area** that is sufficiently bright; allow each pixel to be one of 3 colors, chosen adaptively, e.g. to minimize perceptual differences



[BUSTED]: each branch-site combination is drawn from a 3-bin ( $dS, dN$ ) distribution. The distribution is estimated from the entire alignment. Tests if  $dN/dS > 1$  for some branch/site pairs in the alignment

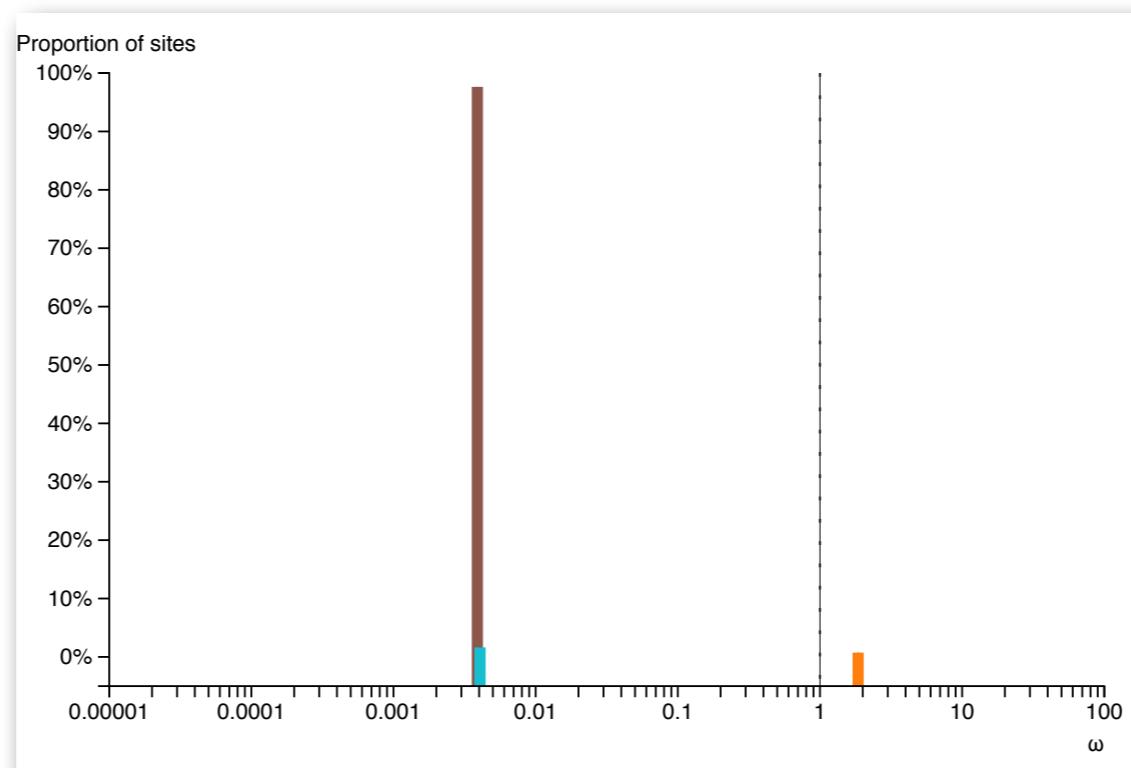
# Gene-wide selection analysis using a branch-site method (BUSTED), HIV-1 *env*

Gene-wide <b>dN/dS</b> distribution	$\omega_1 = 0.627$ (71%) $\omega_2 = 0.649$ (27%) $\omega_3 = 106$ (2%)
<b>p-value</b> for selection ( $H_0 : \omega_3 = 1$ )	$<10^{-15}$
<b>Log L</b> (no variation)	-2078.20
<b>Log L</b> (branch-site; 4 addt'l parameters)	-2039.99



# Gene-wide selection analysis using a branch-site method (BUSTED), WN NS3

Gene-wide $dN/dS$ distribution	$\omega_1 = 0.004$ (99.3%) $\omega_2 = (\text{n/a})$ $\omega_3 = 1.86$ (0.73%)
p-value for selection ( $H_0 : \omega_3 = 1$ )	0.54
Log L (no variation)	-6413.50
Log L (branch-site; 4 addt'l parameters)	-6396.18



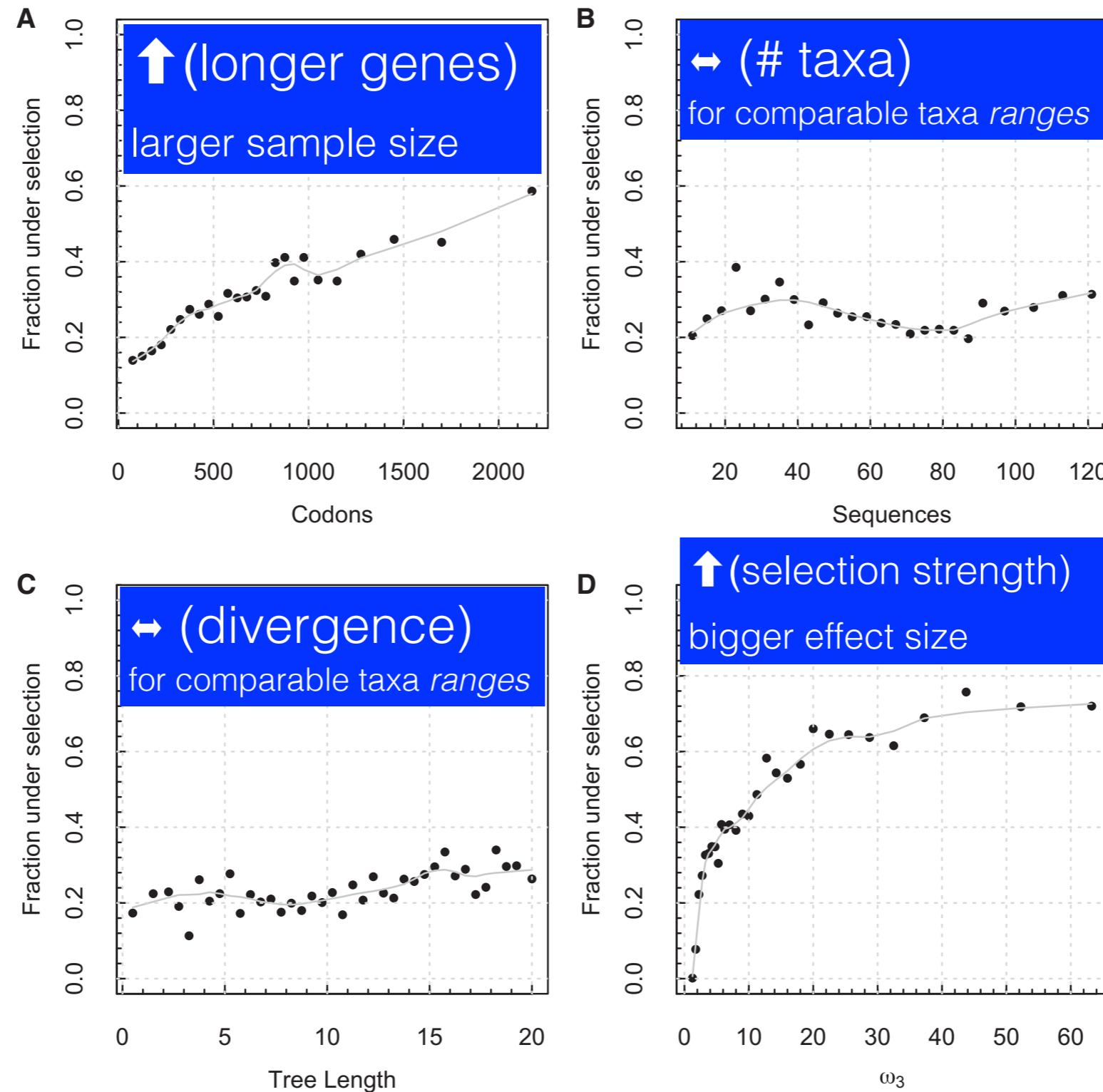
# BUSTED analysis

---

- West Nile Virus NS3 protein
  - No statistical support for selection; ML point estimate allocates a small proportion of sites (~1%) to the selected group ( $dN/dS \sim 2$ )
  - The rest of the gene is very strongly conserved ( $dN/dS = 0.004$ )
- HIV-1 transmission pair
  - Very strong evidence of strong episodic diversification ( $dN/dS \sim 100$ ) on a small proportion of sites (2%)
  - The rest of the gene evolves with weak purifying selection ( $dN/dS = 0.6-0.7$ )

# Where does the power come from for BUSTED?

An analysis of ~9,000 curated gene alignments from [selectome.unil.ch](http://selectome.unil.ch)

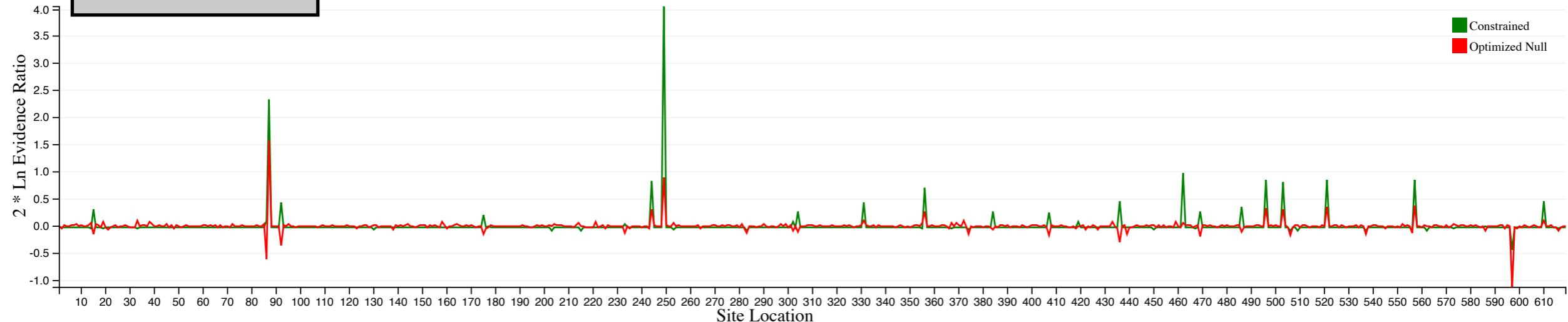


# BUSTED site-level inference

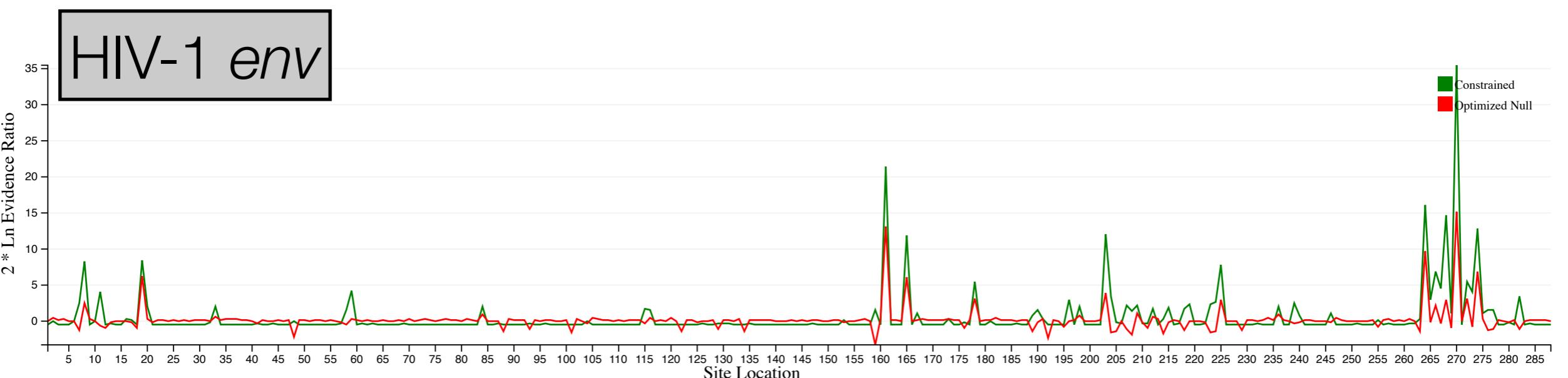
---

- Because BUSTED is a random-effects method, it **pools** information across multiple sites and branches to gain power
- The cost to this pooling is lack of site-level **resolution**, i.e., it is not immediately obvious which sites and/or branches drive the signal
- Standard ways to extract individual site contributions to the overall signal is to perform a post-hoc analysis, such as empirical Bayes, or “category loading”
- For BUSTED, “category loading” is faster and experimentally better

# WN NS3

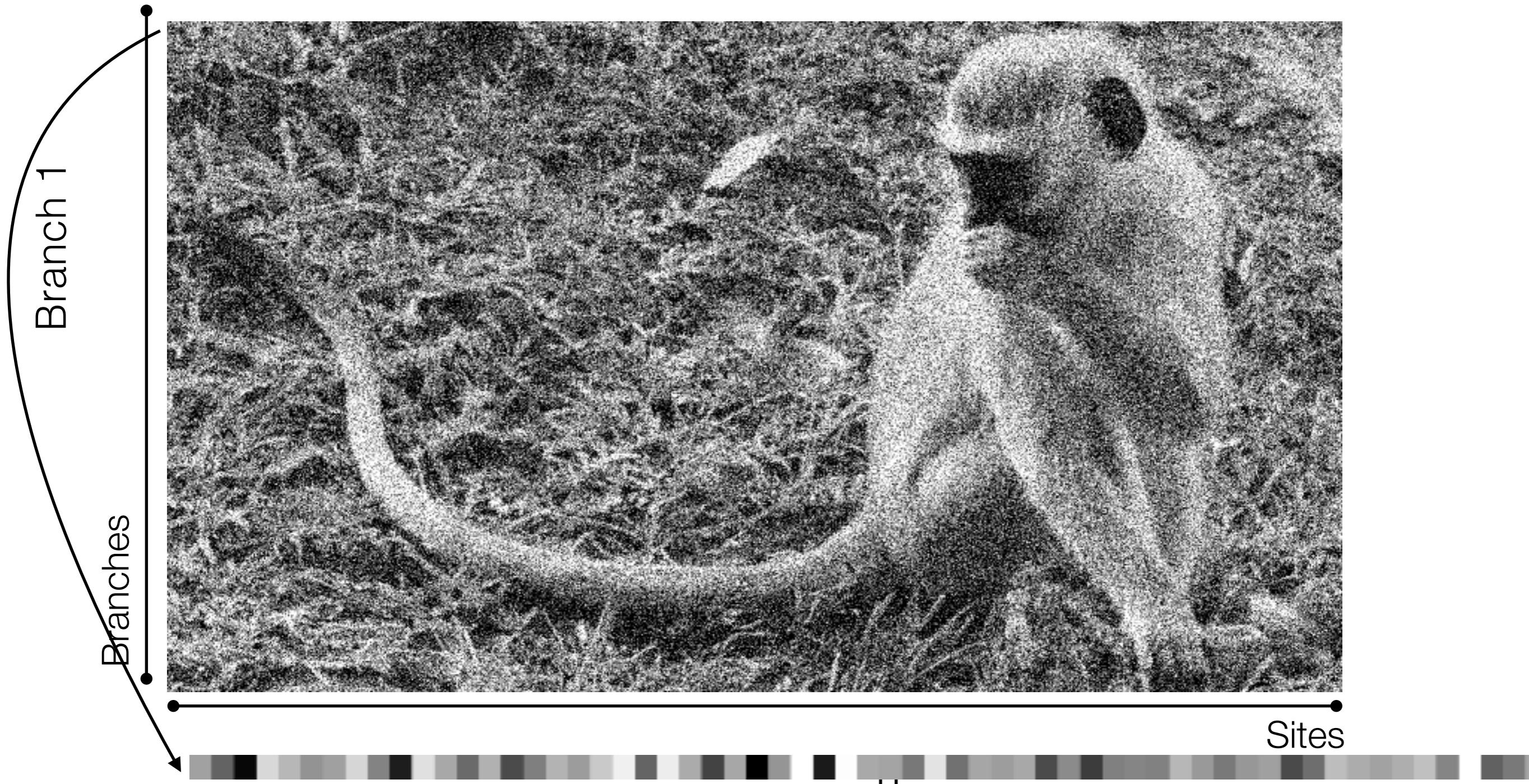


Site Index	Unconstrained Likelihood	Constrained Likelihood	Optimized Null Likelihood	Constrained Evidence Ratio	Optimized Null Evidence Ratio
249 - 249					
249	-65.6860	-67.7124	-66.1326	2.0264	0.4467



Site Index	Unconstrained Likelihood	Constrained Likelihood	Optimized Null Likelihood	Constrained Evidence Ratio	Optimized Null Evidence Ratio
161 - 274					
161	-17.4988	-28.2260	-24.1050	10.7272	6.6062
165	-17.4502	-23.4271	-20.5521	5.9769	3.1019
203	-26.0973	-32.1071	-28.0348	6.0098	1.9375
264	-11.6149	-19.6889	-16.4445	8.0741	4.8297
268	-21.5339	-28.9201	-23.0637	7.3862	1.5298

# Which branches are under selection?



For each image **row**, is there a significant proportion of bright pixels, once the column has been reduced to **N** colors only?

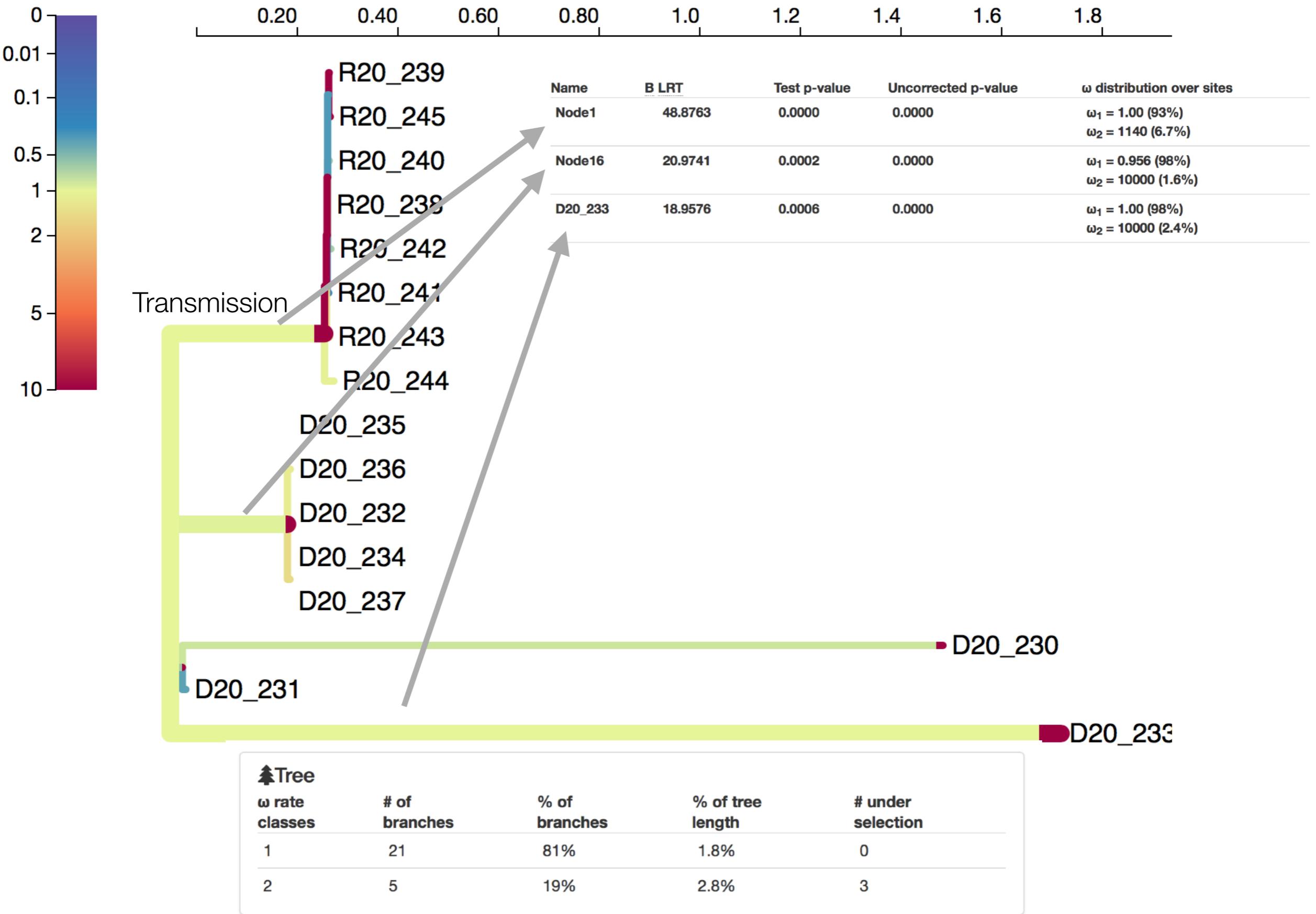


[aBSREL]: at a given branch, each site is a draw from an N-bin ( $dN/dS$ ) distribution, which is inferred from all data for the branch. Test if there is a proportion of sites with  $dN/dS > 1$  (LRT). **N** is derived adaptively from the data.

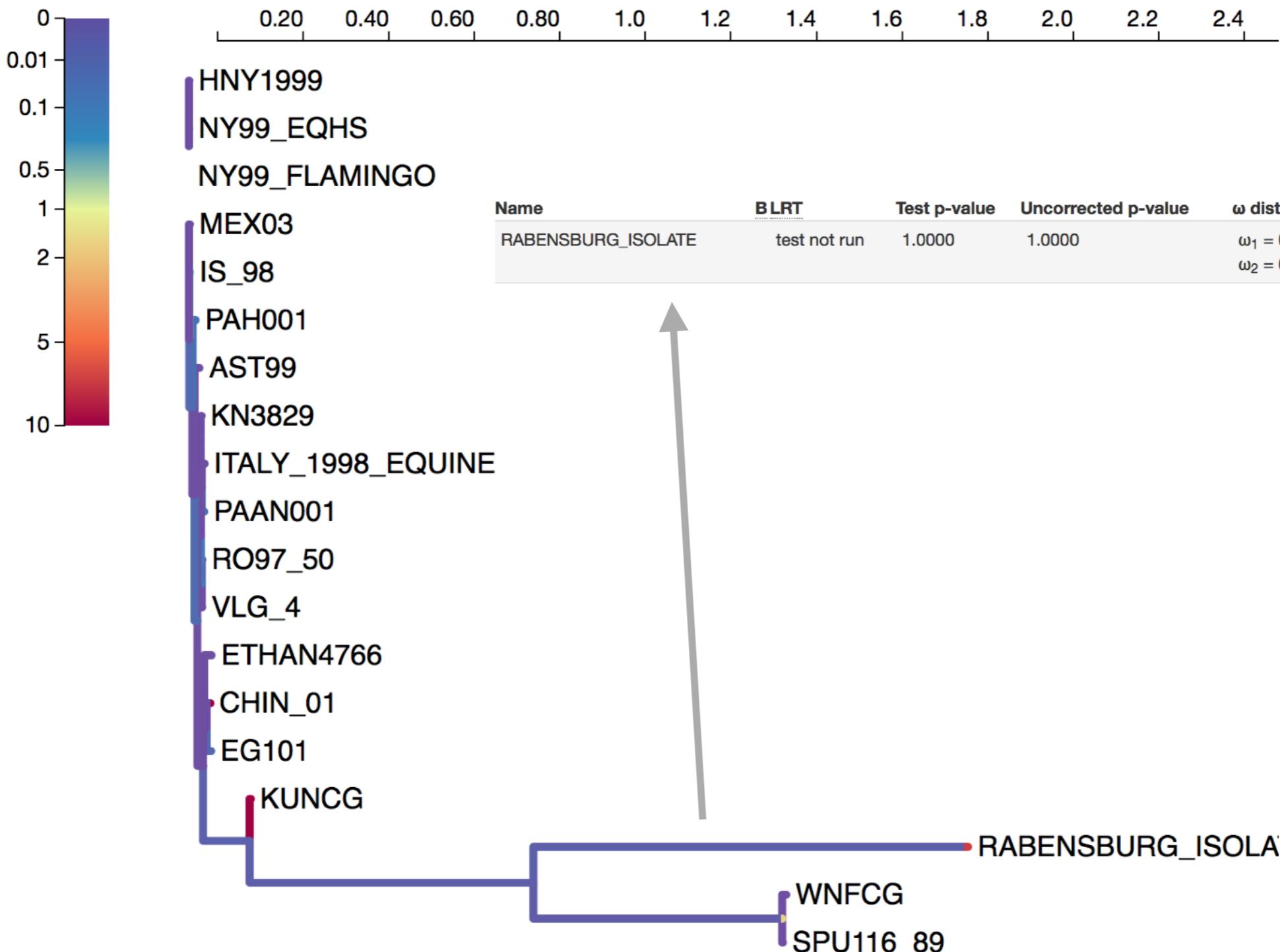
- Best-in-class power
- Able to detect episodes of selection, not just selection on average at a branch
- Does not make unrealistic assumptions for tractability, improves statistical behavior
- Sample size is ~sites, branch level rate estimates could be imprecise
- Cannot reliably estimate which individual sites are subject to selection
- Exploratory testing of all branches leads to loss of power for large data sets (multiple test correction)

- Uses a computationally simple trick to compute the likelihood of data, efficiently summing over all possible assignments of rate classes to branches
  - These cannot be factored into products, unlike sites, because evolution across tree branches is correlated, i.e. a change in the process along one branch affects many others.
- Uses a greedy (but well-performing) step-up procedure to decide how many rate classes to allocate to each branch, prior to testing for selection
  - Perform an evolutionary complexity analysis first (the *adaptive* part), then run selection tests.

# HIV-1 env



# WN NS3



## Tree

$\omega$ rate classes	# of branches	% of branches	% of tree length	# under selection
1	32	91%	35%	0
2	3	8.6%	60%	0

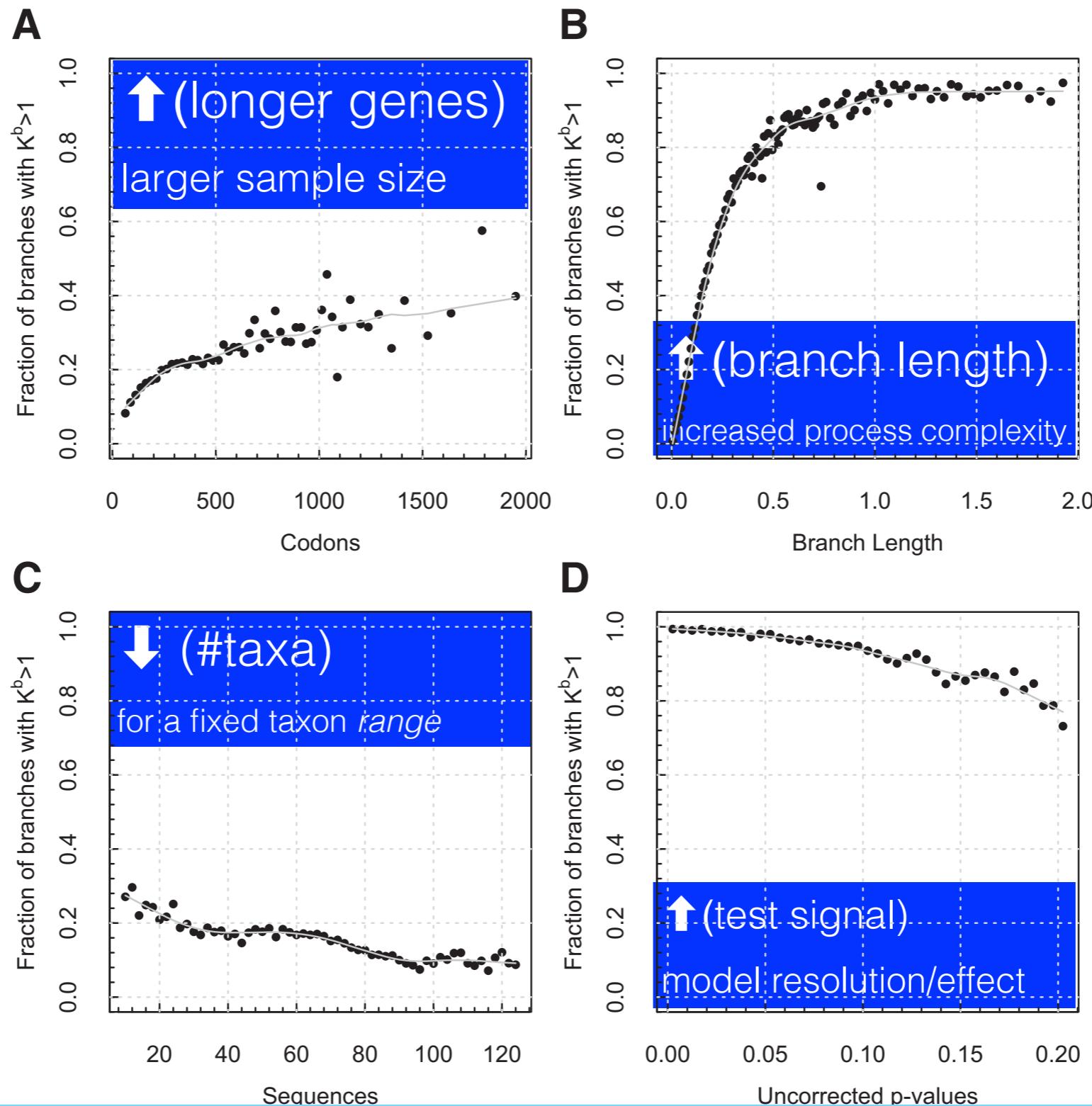
# aBSREL analysis

---

- West Nile Virus NS3 protein
  - 91% branches can be explained with simple (single  $dN/dS$ ) models
  - 3 branches (9%, 60% of tree length) have evidence of multiple  $dN/dS$  rate classes over sites, but **none** with significant proportions of sites with  $dN/dS > 1$
- HIV-1 transmission pair
  - 81% branches can be explained with simple (single  $dN/dS$ ) models
  - 5 branches (19%, 90+% of tree length) have evidence of multiple  $dN/dS$  rate classes over sites
  - 3 branches have small (1–7%), but statistically significant ( $p < 0.05$ , multiple testing corrected) proportions of sites with  $dN/dS > 1$ , including the **transmission** branch

# Correlates of evolutionary complexity

An analysis of ~9,000 curated gene alignments from [selectome.unil.ch](http://selectome.unil.ch)



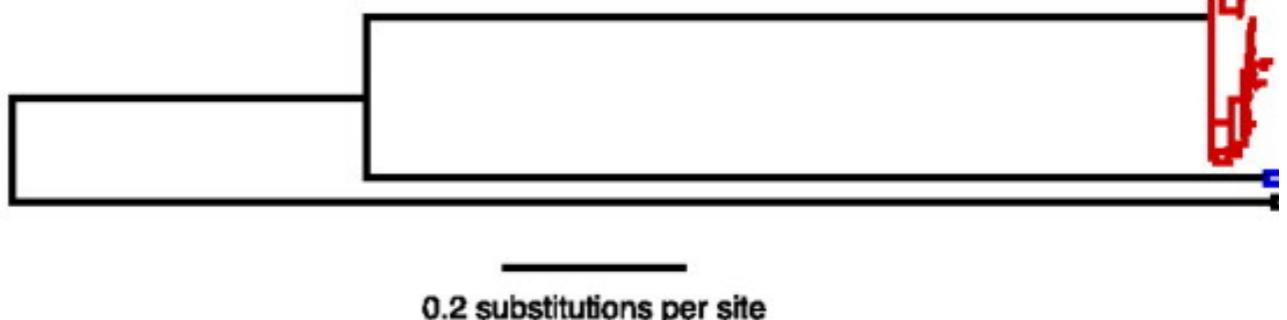
# Unanticipated effects of bad modeling assumptions

---

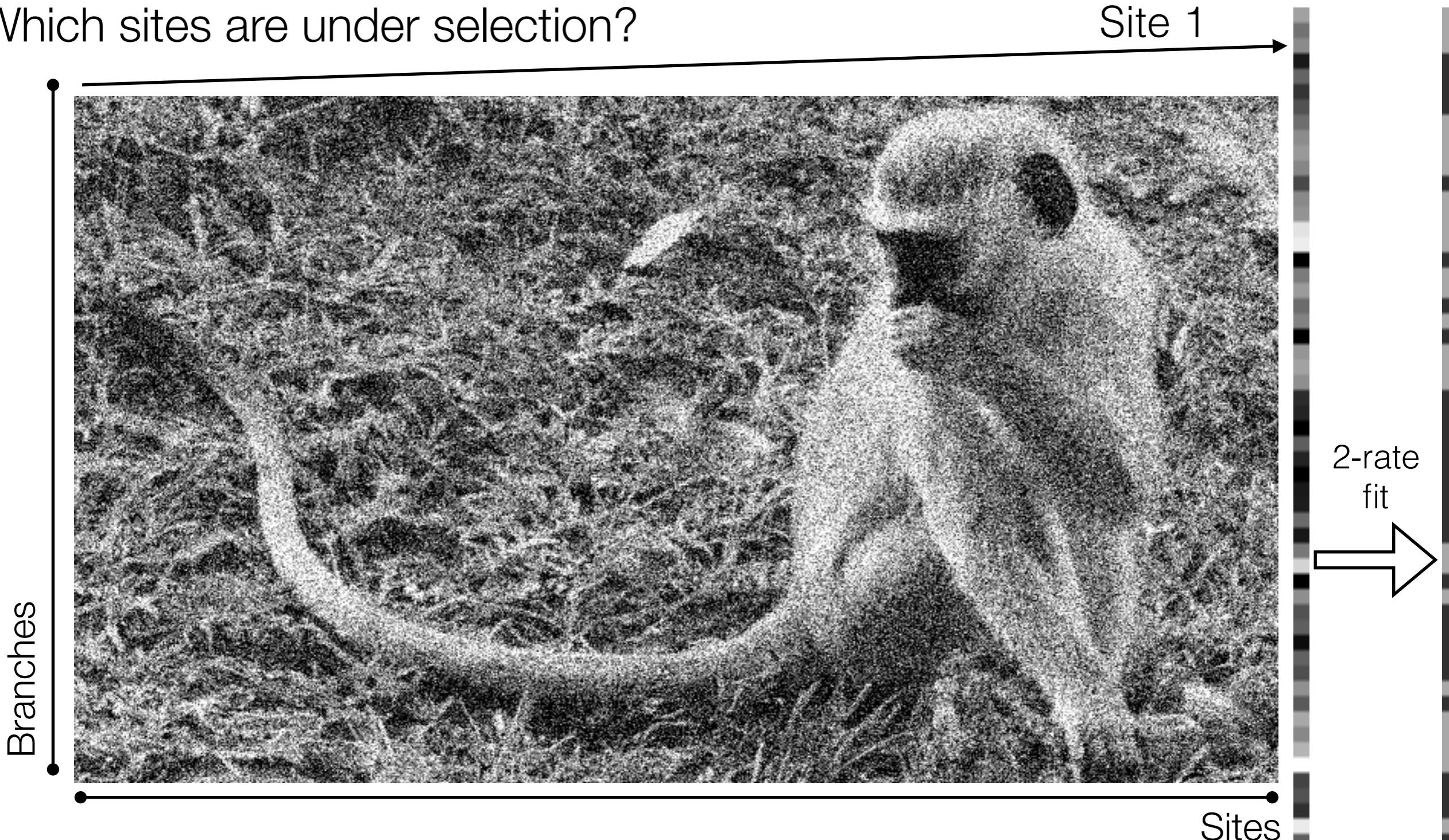
- Models that fail to account for significant shifts in selective pressures through lineages also significantly underestimate branch lengths
- An instructive example is long-range molecular dating of pathogens, where recent isolates (e.g., 30-50 years of sampling) are used to extrapolate the date when a particular pathogen had emerged
- This creates the situation when terminal branches in the tree have relatively high dN/dS (within-host level evolution), which deep interior branches have very low dN/dS (long term conservation)

- Using models that do not vary selection pressure across lineages **A GTR +  $\Gamma_4$**  yields a patently false “*too young*” estimate for the origin of **measles** (about 600 years ago)
- This estimate is refuted by clear historical records which suggest that measles is at least 1,500-5,000 years old
  - *This includes a treatise by a Persian physician Rhazes about **differential diagnosis of measles and smallpox** published circa 600 AD.*
  - Same patterns found for corona-viruses, ebola, avian influenza and herpesvirus

**B Lineage+Dual (two rate)**



Which sites are under selection?

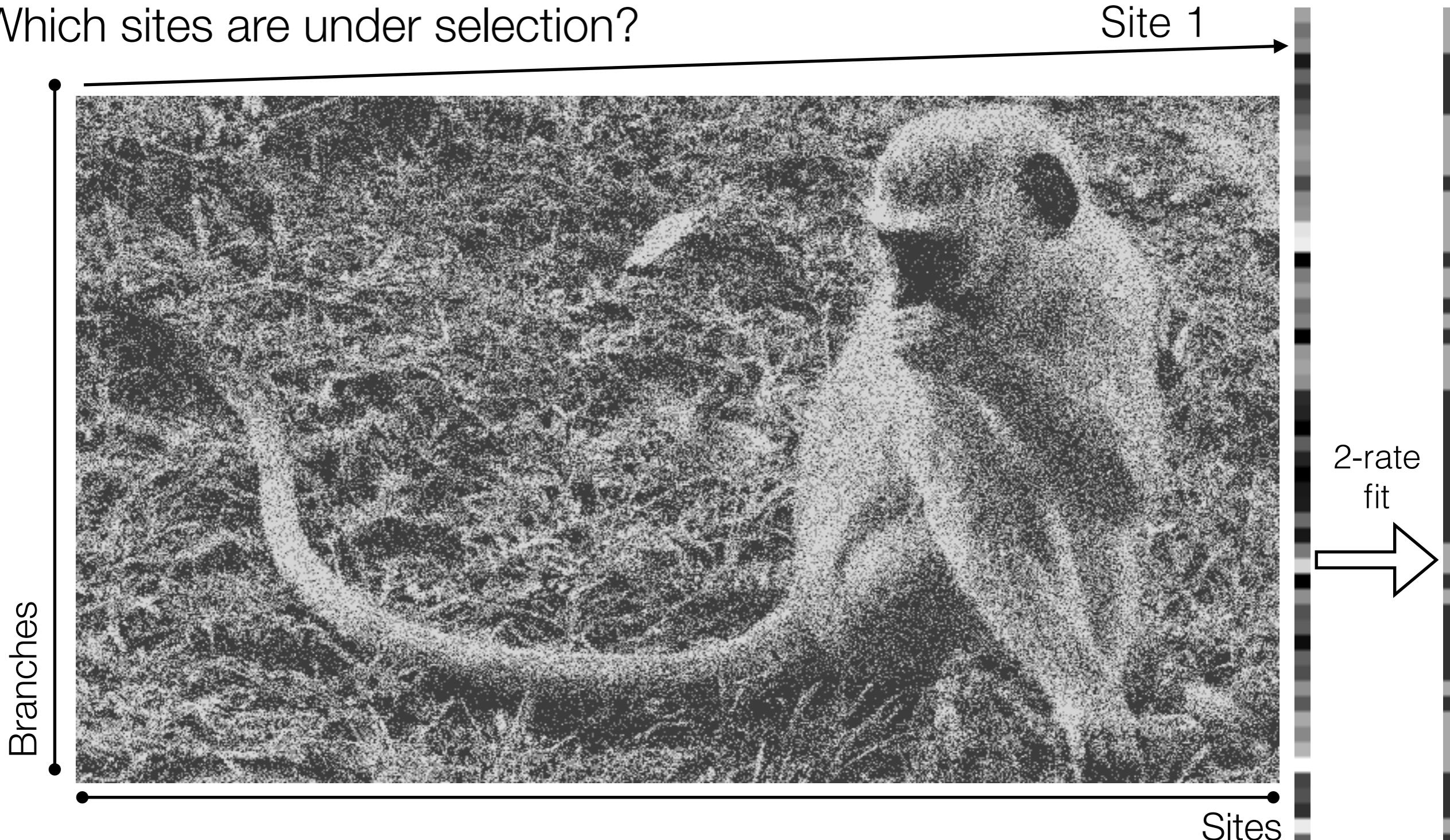


For each image column, is there a significant proportion of bright pixels, once the column has been reduced to 2 colors only?



[MEME]: at a given **site**, each branch is a draw from a 2-bin ( $dS$ ,  $dN$ ) distribution, which is inferred from that site only. Test if there is a proportion of branches with  $dN > dS$  (LRT)

Which sites are under selection?



For each image column, is there a significant proportion of bright pixels, once the column has been reduced to 2 colors only?



[MEME]: at a given **site**, each branch is a draw from a 2-bin ( $dS$ ,  $dN$ ) distribution, which is inferred from that site only. Test if there is a proportion of branches with  $dN > dS$  (LRT)

# Detecting Individual Sites Subject to Episodic Diversifying Selection

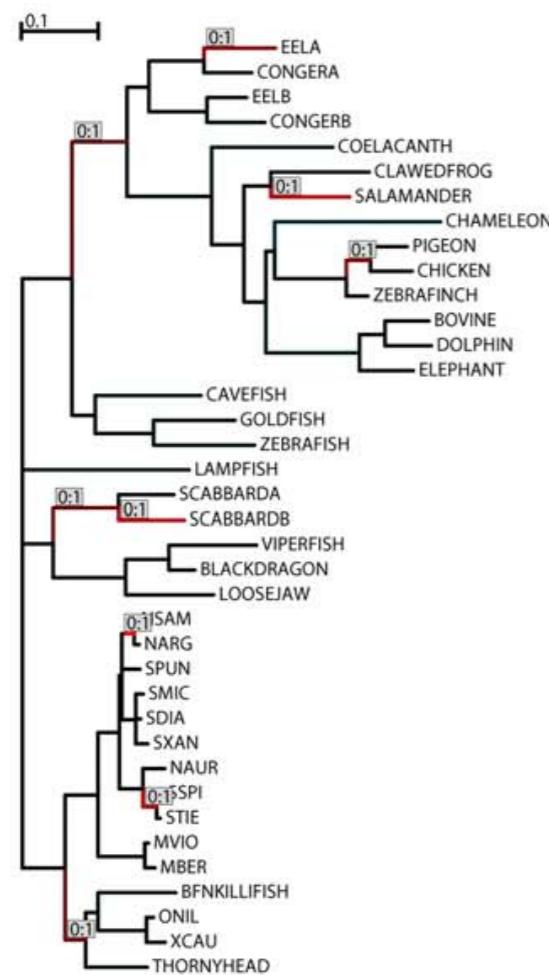
Ben Murrell<sup>1,2</sup>, Joel O. Wertheim<sup>3</sup>, Sasha Moola<sup>2</sup>, Thomas Weighill<sup>2</sup>, Konrad Scheffler<sup>2,4</sup>,  
Sergei L. Kosakovsky Pond<sup>4\*</sup>

PLOS Genetics | www.plosgenetics.org

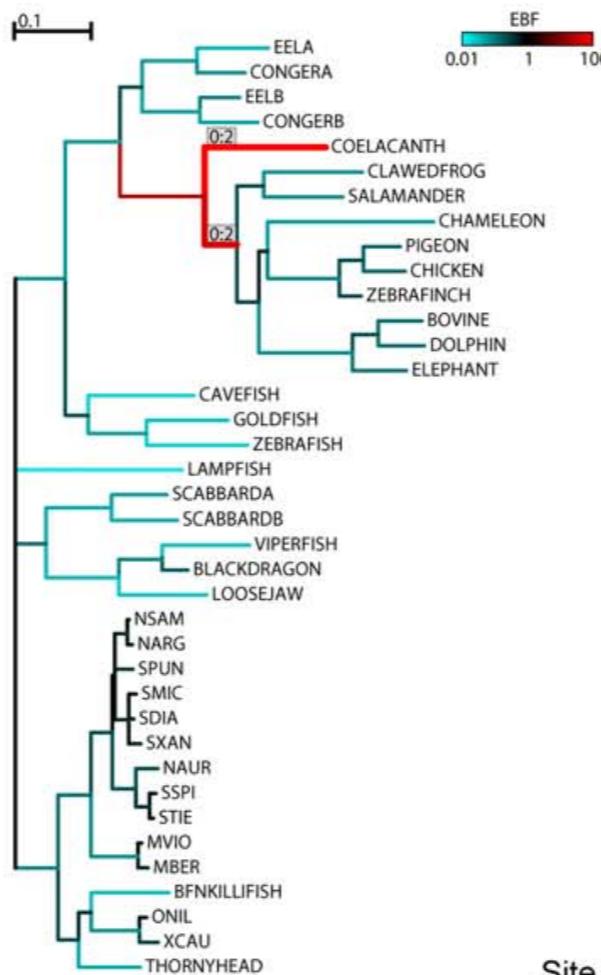
1

July 2012 | Volume 8 | Issue 7 | e1002764

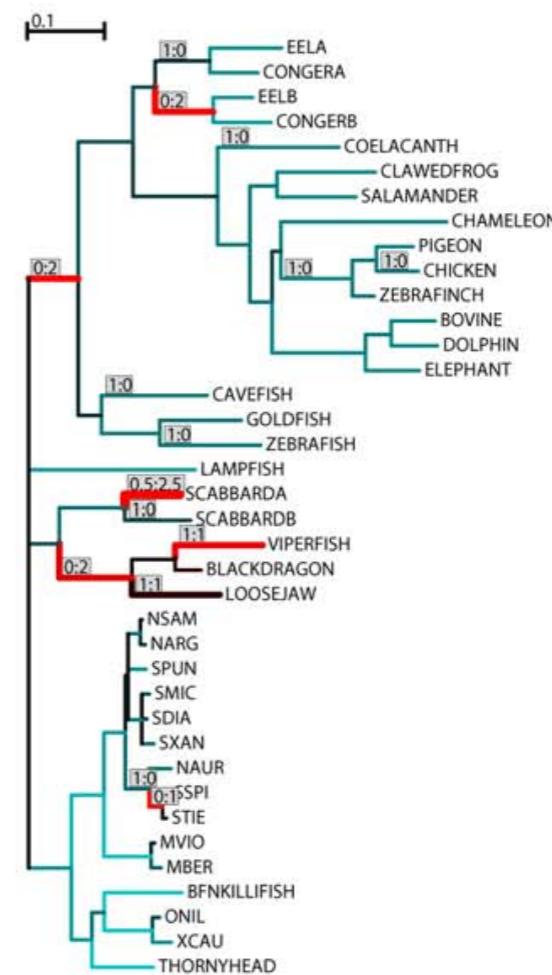
- Best-in-class power
- Able to detect episodes of selection, not just selection on average at a site
- Embarrassingly parallel (farm out each site), so runs reasonably fast
- Sample size is ~sequences, site level rate estimates imprecise
- Cannot estimate which individual branches are subject to selection
- Does not scale especially well with the number of sequences



Site 54



Site 273



Site 210

Pervasive selection,  
also picked up by  
older methods

Episodic selection,  
missed by old  
methods

Episodic selection,  
followed by  
conservation.  
Miscalled by old  
methods as purifying  
selection only

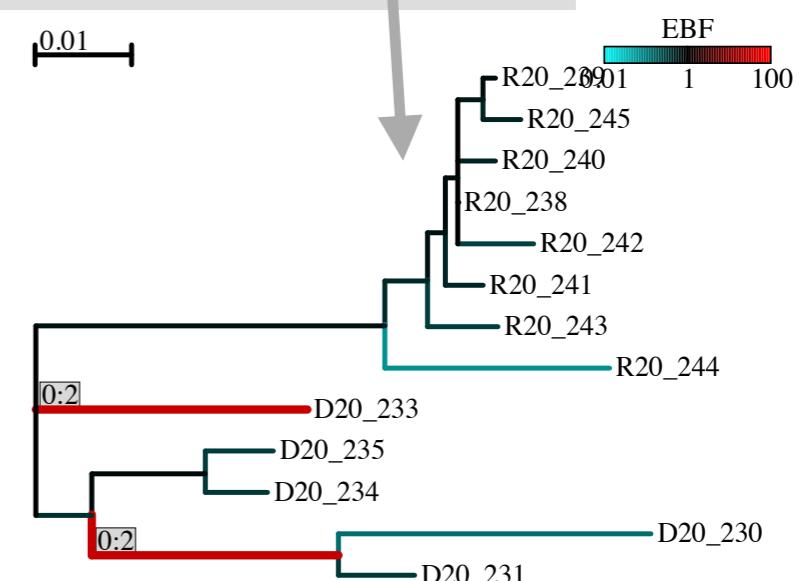
# HIV-1 env

Found 11 sites with evidence of episodic diversifying selection ( 0.1 significance level Retabulate )

This summary table reports the distribution of synonymous ( $\alpha$ ) and non-synonymous ( $\beta$ ) substitution rates over sites inferred by the MEME model, where the proportion of branches with  $\beta > \alpha$  is significantly greater than 0. p-value is derived using a mixture of  $\chi^2$  distributions, and q-values are obtained using Simes' procedure, which controls the false discovery rate under the strict neutral null (likely to be conservative).

Codon	$\alpha$	$\beta^-$	Pr[ $\beta = \beta^-$ ]	$\beta^+$	Pr[ $\beta = \beta^+$ ]	p-value	q-value	Branch-site information
19	0	0	0.943201	773.518	0.0567992	0.0317833	1	[Display]
161	0	0	0.829382	115.825	0.170618	0.00844037	1	[Display]
165	0	0	0.779355	60.0702	0.220645	0.0688644	1	[Display]
178	0	0	0.902811	67.331	0.0971885	0.0357418	1	[Display]
225	0	0	0.78381	72.3626	0.21619	0.0458125	1	[Display]
261	0.000305384	0.00027966	0.913635	881.437	0.0863648	0.0895985	1	[Display]
268	0	0	0.895791	10000	0.104209	0.0350278	1	[Display]
270	0	0	0.760392	226.778	0.239608	0.057291	1	[Display]
272	0	0	0.874576	59.5745	0.125424	0.0525324	1	[Display]
274	4.99073	0	0.764401	260.143	0.235599	0.062812	1	[Display]
282	0	0	1e-09	10.8217	1	0.0881541	1	[Display]

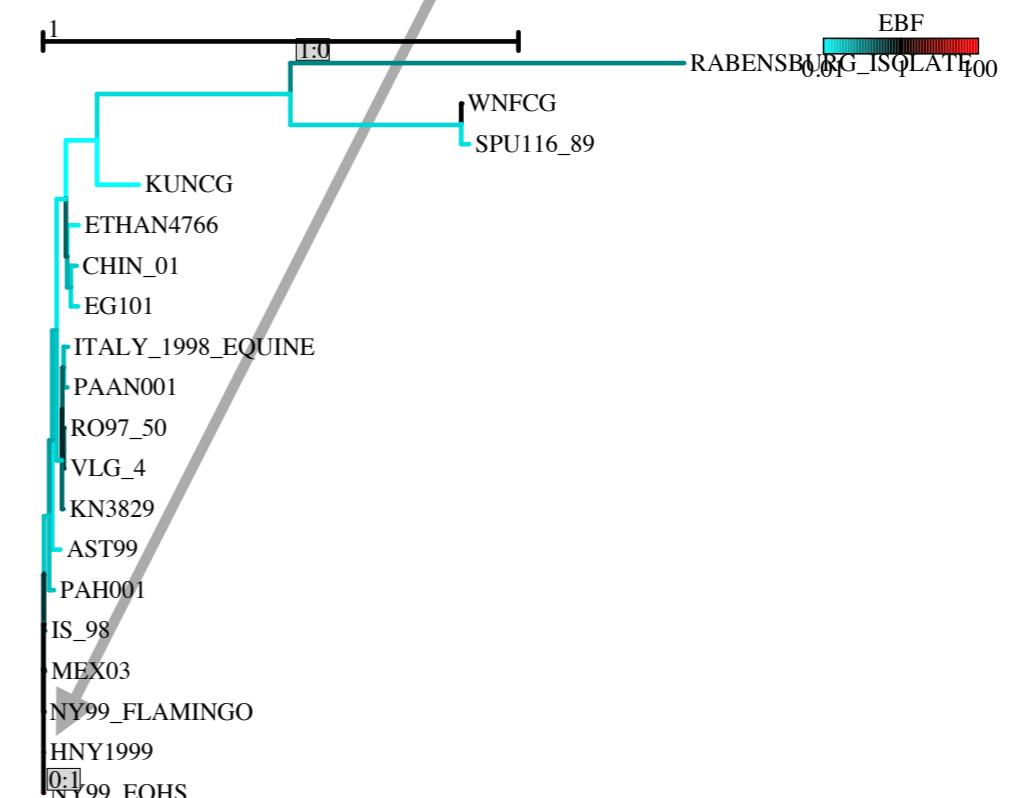
Site 161  
82% of branches with  $\alpha=\beta=0$   
18% of branches with  $\alpha=0, \beta=116$



Found 3 sites with evidence of episodic diversifying selection ( 0.1 significance level Retabulate )

This summary table reports the distribution of synonymous ( $\alpha$ ) and non-synonymous ( $\beta$ ) substitution rates over sites inferred by the MEME model, where the proportion of branches with  $\beta > \alpha$  is significantly greater than 0. p-value is derived using a mixture of  $\chi^2$  distributions, and q-values are obtained using Simes' procedure, which controls the false discovery rate under the strict neutral null (likely to be conservative).

Codon	$\alpha$	$\beta^-$	Pr[ $\beta = \beta^-$ ]	$\beta^+$	Pr[ $\beta = \beta^+$ ]	p-value	q-value	Branch-site information
249	0	0	1.00003e-09	2.44107	1	0.0166364	1	[Display]
496	0.947581	0	0.955206	74.0758	0.0447943	0.0838408	1	[Display]
557	0.275201	0	0.963363	171.171	0.036637	0.0261761	1	[Display]



Site 557  
96% of branches with  $\alpha=0.28$ ,  $\beta=0$   
4% of branches with  $\alpha=0.28$ ,  $\beta=171$

# MEME results

---

- West Nile Virus NS3 protein
  - Three sites, (including 249) with significant evidence of episodic (or pervasive) diversifying selection.
- HIV-1 transmission pair
  - Eleven sites with significant evidence of episodic (or pervasive) diversifying selection.

# Why MEME?

---

- Affords a much greater power to detect selection
- Mitigates the pathological effect when adding sequences to a sample can reduce, or remove, signal of selection

"The greater power of MEME indicates that selection acting at individual sites is considerably more widespread than constant  $\omega$  models would suggest. It also suggests that natural selection **is predominantly episodic**, with transient periods of adaptive evolution masked by the prevalence of purifying or neutral selection on other branches. We emphasize that MEME is not just a quantitative improvement over existing models: for 56 sites in our empirical analyses, we obtain qualitatively different conclusions. FEL asserts that these sites evolved under significant **purifying** selection, but MEME is able to identify the signature of **positive selection on some branches**"

# Why MEME?

---

- Affords a much greater power to detect selection
- Mitigates the pathological effect when adding sequences to a sample can reduce, or remove, signal of selection

"Although a previous analysis of 38 vertebrate rhodopsin sequences found no sites under selection at posterior probability >95%, the same authors found 7 selected sites in the subset of 11 squirrelfish sequences, and 2 selected sites when the subset of 28 fish sequences was analyzed. These results run counter to the expectation that more data should provide greater power to detect selection. MEME, on the other hand, [typically] detects more selected sites when more sequences are included."

# Analysis summary

---

	WNV NS3	HIV-1 <i>env</i>
Gene-wide episodic selection (BUSTED)	No	Yes
Branch-level selection (aBSREL)	No	Yes, three branches, including transmission
Site-level episodic selection (MEME)	Yes, 3 sites	Yes, 11 sites

It is **not** unexpected that site-level positive results can occur when a gene-level test does not yield a positive result

---

- **Lack of power for the global test:** if the proportion of sites under selection is very small, a mixture-model test, like BUSTED, will miss it.
- **Model violations:** MEME supplies much more flexible distributions of  $dN/dS$  over sites; compared to alignment-wide 3-bit BUSTED distribution.
- **False positives at site-level:** our site-level tests have good statistical properties, but each positive site result could be a false positive; FWER correction would make site-level tests too conservative.
- **Summary:** gene-level selection tests need a minimal proportion of sites to be under selection to be powered; site-level tests should not be used to make inferences about gene-level selection.

However, we caution that despite obvious interest in identifying specific branch-site combinations subject to diversifying selection, such inference is based on very limited data (the evolution of one codon along one branch), and cannot be recommended for purposes other than data exploration and result visualization. This observation could be codified as the “***selection inference uncertainty principle***” — one cannot simultaneously infer both the site and the branch subject to diversifying selection. In this manuscript [MEME], we describe how to infer the location of sites, pooling information over branches; previously [aBSREL] we have outlined a complementary approach to find selected branches by pooling information over sites.

*Murrell et al 2012*

# Purpose-built models

---

- It is tempting to “hack” existing tools to answer questions that they are not designed to answer
- A recent example we tackled is a rigorous test for relaxation of selection (or more generally a difference in selective regimes) in a part of the tree, relative to the rest of the tree
- Typical approaches have been to estimate dN/dS ratios from two sets of branches, and interpret an *elevation* in dN/dS as evidence of selective constraint relaxation
- Two problems with this approach
  - An increase in mean dN/dS could also be caused by an **intensification** of selective forces.
  - *Post-hoc* analyses (e.g., estimate branch-level dN/dS and then compare [t-test, etc] them as if they were observed quantities) discard a lot of information (e.g., variance of individual estimates), and make obviously wrong assumptions (e.g., estimates are uncorrelated).

# Testing for selective relaxation



Partition the image into horizontal bands (a priori); compare whether or not there is visual benefit to using separate 3-color palettes in two sets of bands instead of a single 3-color palette



[RELAX]: Compare whether or not the set of branches of interest (test set) has a significantly different dN/dS distribution than the rest of the tree (background), fitted jointly to the entire alignment. For relaxation testing, the two dN/dS distributions are related via a power transformation.

# Testing for selective relaxation



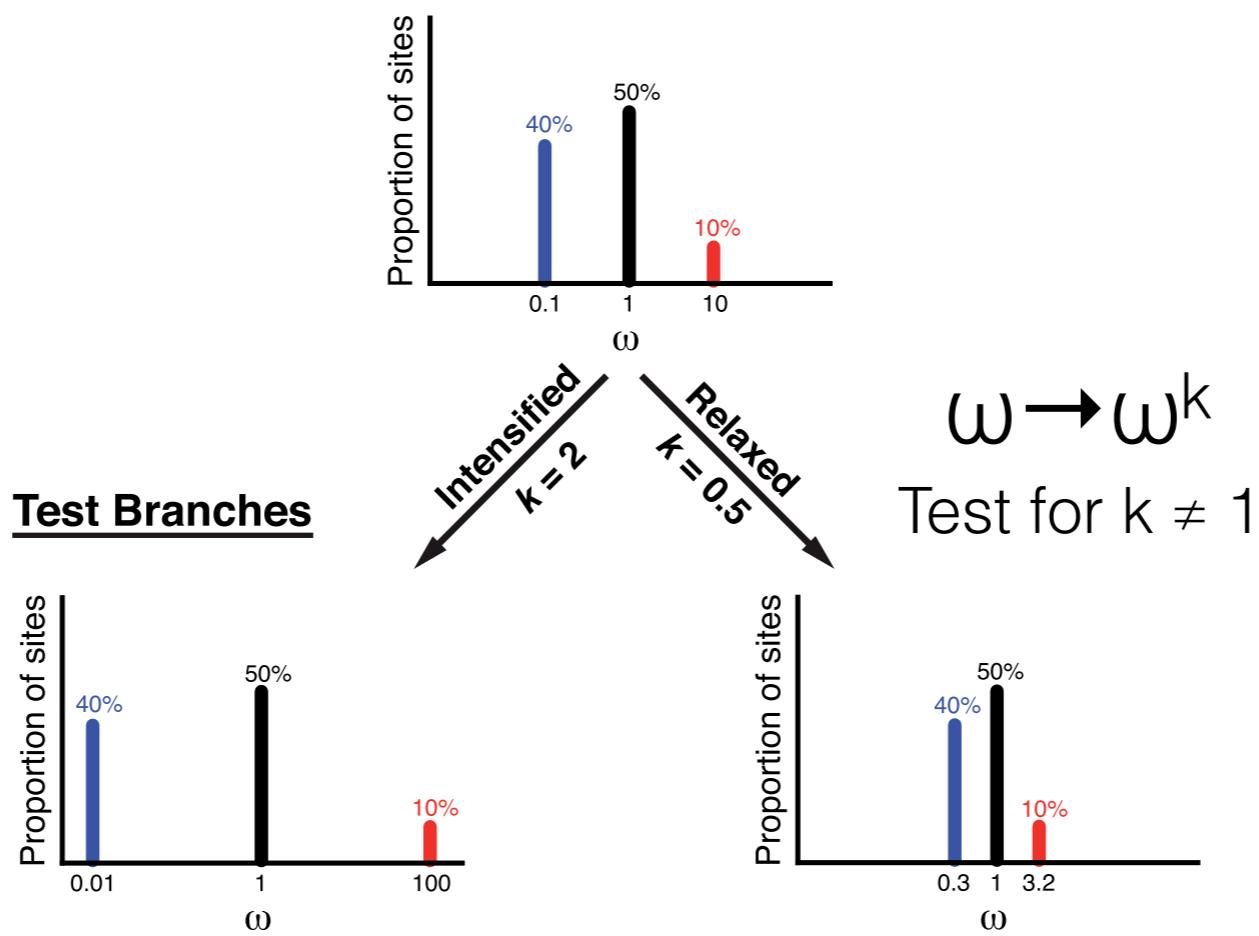
Partition the image into horizontal bands (a priori); compare whether or not there is visual benefit to using separate 3-color palettes in two sets of bands instead of a single 3-color palette



[RELAX]: Compare whether or not the set of branches of interest (test set) has a significantly different dN/dS distribution than the rest of the tree (background), fitted jointly to the entire alignment. For relaxation testing, the two dN/dS distributions are related via a power transformation.

## Reference Branches

Mol. Biol. Evol. 32(3):820–832



**Table 1.** Test for Relaxed Selection Using RELAX in Various Taxonomic Groups.

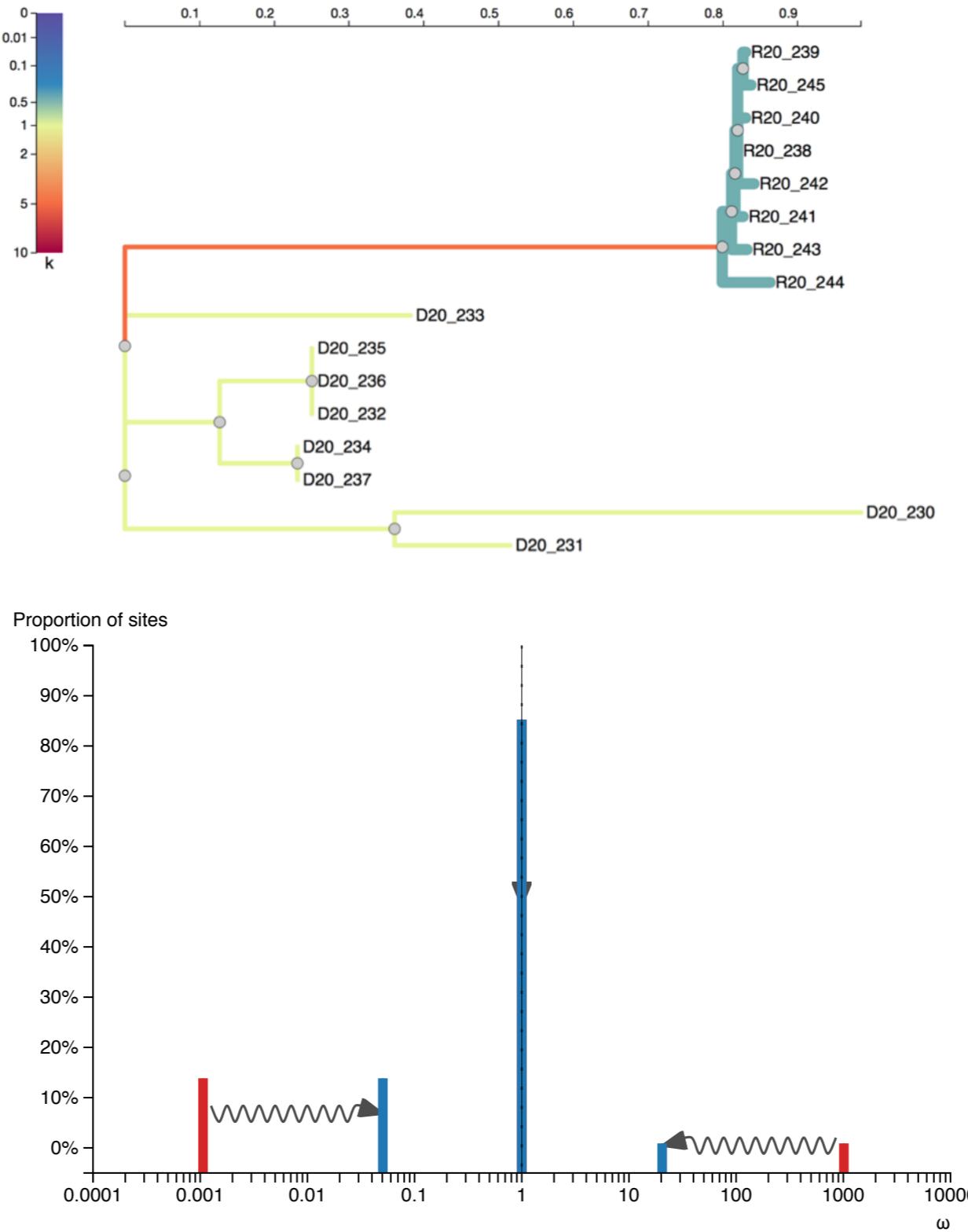
Taxa	Gene/Genes	Test Branches	Reference Branches	$k^a$	P-Value
$\gamma$ -proteobacteria	Single-copy orthologs	Primary/secondary endosymbionts	Free-living $\gamma$ -proteobacteria	0.30	< 0.0001
		Primary endosymbionts	Free-living $\gamma$ -proteobacteria	0.28	< 0.0001
		Secondary endosymbionts	Free-living $\gamma$ -proteobacteria	0.61	< 0.0001
		Primary endosymbionts	Secondary endosymbionts	0.56	< 0.0001
Bats	SWS1	HDC echolocating and cave roosting (pseudogenes)	LDC echolocating and tree roosting (functional genes)	0.16	< 0.0001
		LDC echolocating	Tree roosting	1.07	0.577
	M/LWS1	HDC echolocating and cave roosting	LDC echolocating and tree roosting	0.70	0.495
		Echolocating species	Tree- and cave-roosting species	0.21	0.0005
Bornavirus	Nucleoprotein	HDC echolocating	LDC echolocating	0.84	0.427
	Mitochondrial protein-coding genes	Endogenous viral elements	Exogenous virus	0.02	< 0.0001
<i>Daphnia pulex</i>	Mitochondrial protein-coding genes	Asexual	Sexual	0.63	< 0.0001

<sup>a</sup>Estimated selection intensity.

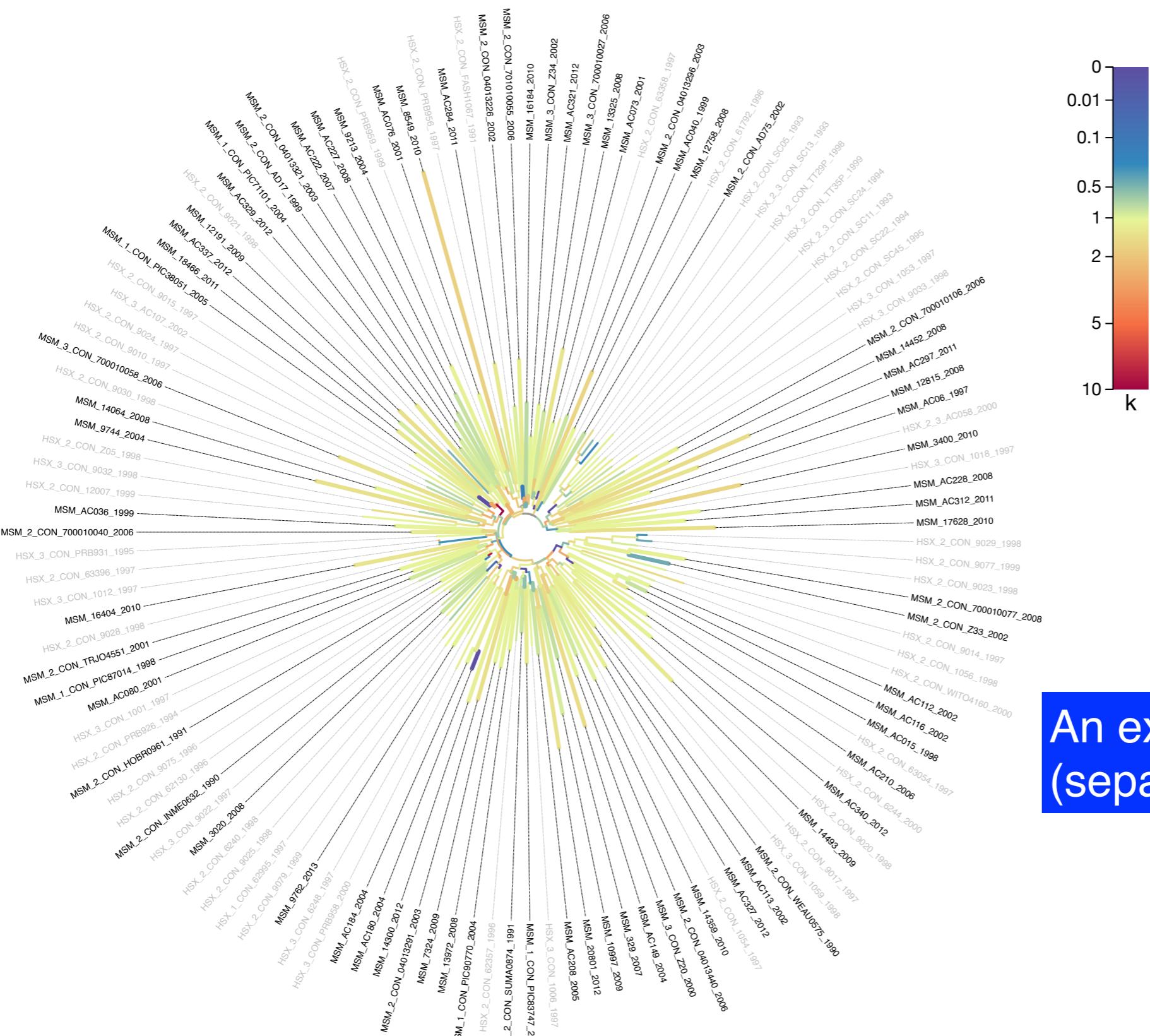
## RELAX(ed selection test) summary

Test for selection **relaxation** ( $K = 0.44$ ) was **significant** ( $p = 0.0002$ ,  $LR = 14.29$ )

Please cite PubMed ID undefined if you use this result in a publication, presentation, or other scientific work.



Another use of RELAX: test for difference of selective pressures between HSX and MSM HIV-1 isolates



## An exploratory model fit (separate k for each branch)

[RELAX] assigned **fewer codon sites in the MSM lineages to the positively selected category** (2.6% [2.3-2.9%] in MSM vs 5.4% [5.0-6.4%] in HSX, all confidence intervals are 95% profile likelihood approximations), and **inferred that selection on these sites was stronger in MSM** ( $\omega = 15.8$  [14.4-17.5] in MSM vs  $\omega = 9.2$  [8.2-9.6] in HSX).

## Different distributions fitted to sets of branches

### Model fits

Model	<i>log L</i>	# par.	AIC <sub>c</sub>	Time to fit	L <sub>tree</sub>	Branch set	$\omega_1$	$\omega_2$	$\omega_3$
Partitioned MG94xREV	-83169.97	277	166895.40	3 min. 39 sec.	7.49	Reference	0.634 (100%)		
						Test	0.558 (100%)		
General Descriptive	-81843.18	538	164767.88	40 min. 10 sec.	28.38	All	0.0839 (63%)	1.00 (33%)	11.9 (3.1%)
Null	-81960.41	358	164639.26	30 min. 11 sec.	26.79	Reference	0.00 (58%)	1.00 (38%)	13.0 (3.4%)
						Test	0.00 (58%)	1.00 (38%)	13.0 (3.4%)
						Unclassified	0.0000000750 (63%)	0.974 (36%)	15.4 (2.0%)
Alternative	-81959.64	359	164639.74	9 min. 58 sec.	27.03	Reference	0.00235 (58%)	1.00 (39%)	13.3 (3.2%)
						Test	0.00189 (58%)	1.00 (39%)	14.6 (3.2%)
						Unclassified	0.0000000750 (62%)	0.974 (36%)	15.5 (1.9%)
Partitioned Exploratory	-81952.79	363	164634.09	1 hrs. 39 min.	27.15	Reference	0.00 (62%)	1.00 (32%)	8.88 (5.7%)
						Test	0.00 (57%)	1.00 (40%)	17.5 (2.6%)
						Unclassified	0.0000000750 (60%)	0.969 (39%)	15.5 (1.8%)

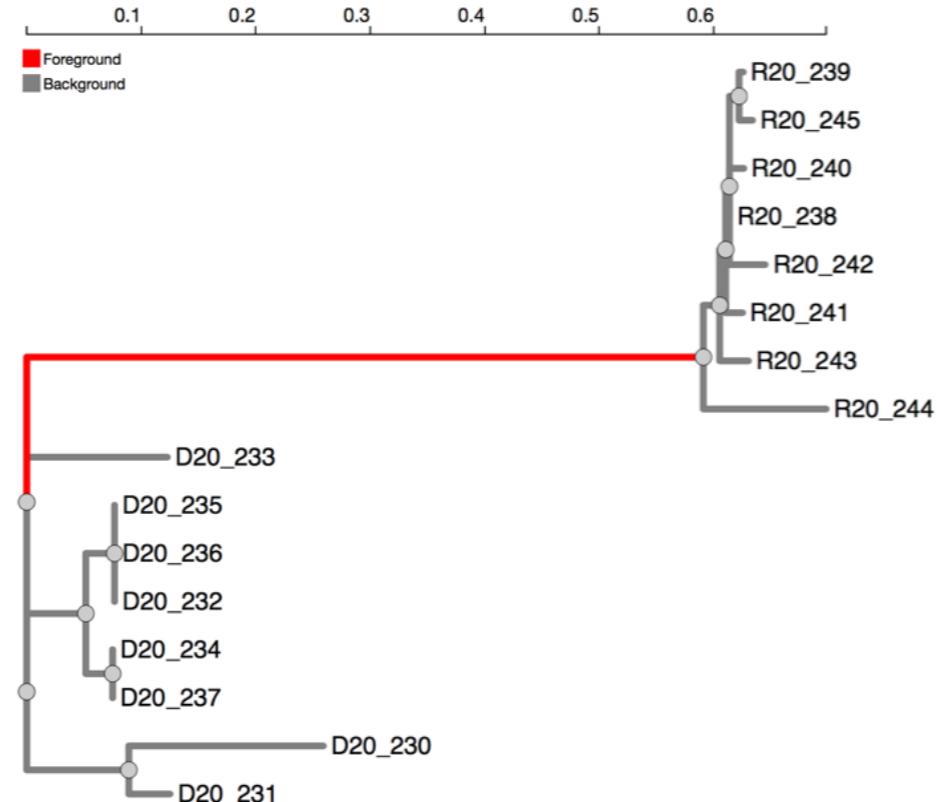
Models compared by AIC<sub>c</sub>  
(or LRT)

Nuisance branches  
explicitly modeled



# Branch testing; exploratory vs *a priori*

- aBSREL and BUSTED can test all branches for selection (exploratory), or apply the test to a set of branches defined *a priori* (e.g. defining a particular biological hypothesis).
- For BUSTED, *a priori* partitioning of branches can increase power, especially if selective regimes are markedly different on different parts of the tree.
- For example, BUSTED applied to the HIV dataset where the transmission branch is designated as foreground, found a greater proportion sites under stronger selection on this branch than the rest of the tree (8% vs 1%), and a lower **p-value**.



	Background	Foreground
Class 1	$\omega = 0.51$ $p = 0.08$	$\omega = 0.00$ $p = 0.92$
Class 2	$\omega = 0.72$ $p = 0.91$	
Class 3	$\omega = 116$ $p = 0.01$	$\omega = 510$ $p = 0.08$

Task	Test	Site strategy	Branch strategy	Complexity	Effective sample size	Parallelization	Practical # sequences limit
<b>Gene-wide selection</b>	BUSTED	Random Effects	Random Effects	Fixed	$\sim$ sites x taxa	SMP	$\sim$ 1,000
<b>Site-level selection</b>	MEME	Fixed Effects	Random Effects	Fixed	$\sim$ taxa	MPI	$\sim$ 5000 (cluster)
<b>Branch-level selection</b>	aBSREL	Random Effects	Fixed Effects	Adaptive	$\sim$ sites	SMP/MPI	$\sim$ 1,000
<b>Compare selective regimes between sets of branches</b>	RELAX	Random Effects	Mixed Effects	Fixed	$\sim$ sites x (branch set size)	SMP	$\sim$ 1,000

# FUBAR: selection testing done fast

Branches



Average colors over sites; use a relatively large but fixed palette to approximate the image



[FUBAR]: Fix a grid of  $dS$  and  $dN$  values, use the data to sample (Bayesian MCMC) weights to individual grid points; this forms the prior distribution on rates; use empirical Bayes to obtain site-level estimates of posterior probability that  $dN > dS$

# FUBAR: selection testing done fast

5 (best) color adaptive palette



Branches



Average colors over sites; use a relatively large but fixed palette to approximate the image



[FUBAR]: Fix a grid of  $dS$  and  $dN$  values, use the data to sample (Bayesian MCMC) weights to individual grid points; this forms the prior distribution on rates; use empirical Bayes to obtain site-level estimates of posterior probability that  $dN > dS$

# FUBAR: selection testing done fast

Fixed web palette (216 colors)

Branches



Average colors over sites; use a relatively large but fixed palette to approximate the image



[FUBAR]: Fix a grid of  $dS$  and  $dN$  values, use the data to sample (Bayesian MCMC) weights to individual grid points; this forms the prior distribution on rates; use empirical Bayes to obtain site-level estimates of posterior probability that  $dN > dS$

# FUBAR: selection testing done fast

Fixed web palette (216 colors)

Branches



Average colors over sites; use a relatively large but fixed palette to approximate the image



[FUBAR]: Fix a grid of dS and dN values, use the data to sample (Bayesian MCMC) weights to individual grid points; this forms the prior distribution on rates; use empirical Bayes to obtain site-level estimates of posterior probability that  $dN > dS$

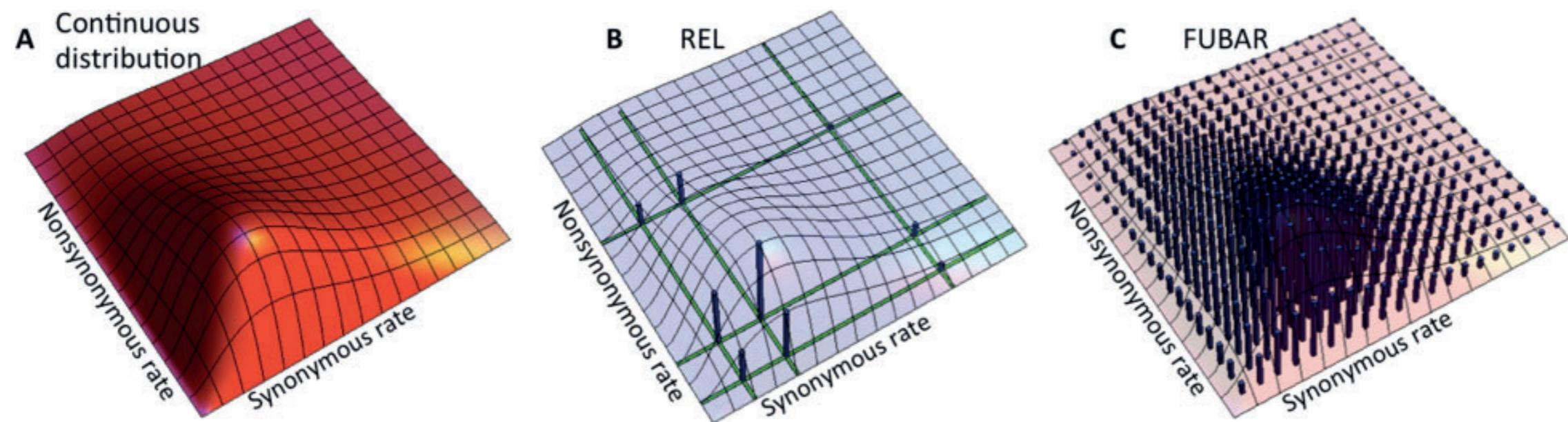
- The time consuming part of traditional random-effects models is the estimation of the alimenter-wide dN/dS distribution
- Each hyper-parameter adjustment entails an expensive phylogenetic likelihood calculation
- Larger data sets —> more complex mixtures needed to avoid smoothing, i.e., more parameters, more evaluations, and a non-linear dependence on data-set sizes

- With FUBAR we make the following approximations:
  - Branch lengths, GTR biases etc, are estimated using simple (nucleotide models) and held fixed
  - We fix a 15x15 or 20x20 grid of (dS,dN) values *a priori*; the data only inform how much weight will be allocated to each point
- Only need to evaluate the expensive codon-based phylogenetic likelihood once for each grid point: **complexity only increases linearly with the size of the data.** This step is also embarrassingly parallel.
- Allocating weights to individual points is done using MCMC (or Gibbs sampling, or variational Bayes); this step does not require ANY further evaluations of the phylogenetic likelihood, i.e., **its cost does not depend on the size of the alignment**

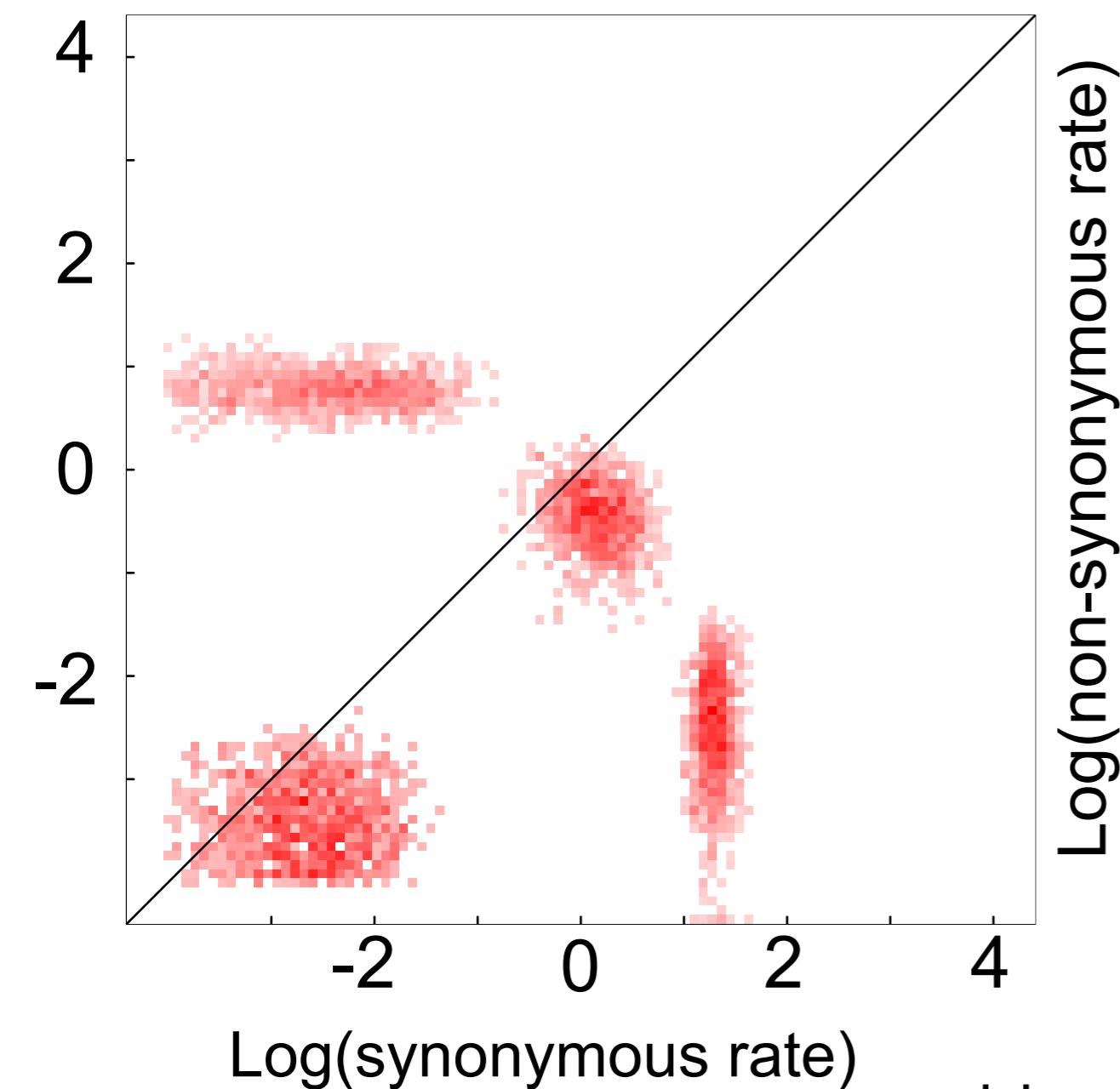
# FUBAR: A Fast, Unconstrained Bayesian AppRoximation for Inferring Selection

Ben Murrell,<sup>1,2,3</sup> Sasha Moola,<sup>1,3</sup> Amandla Mabona,<sup>1,4</sup> Thomas Weighill,<sup>1</sup> Daniel Sheward,<sup>5</sup> Sergei L. Kosakovsky Pond,<sup>6</sup> and Konrad Scheffler<sup>\*,1,6</sup>

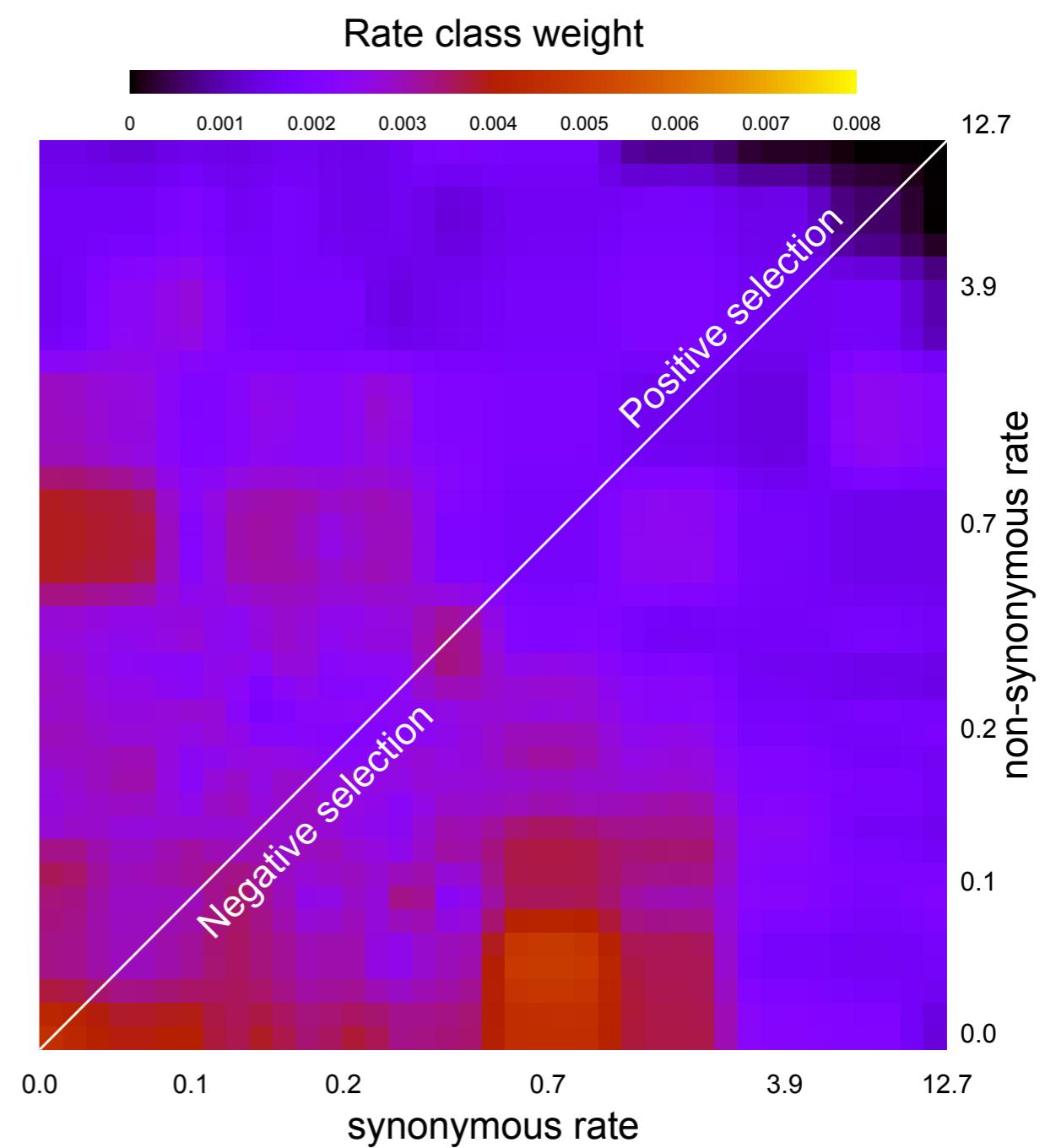
*Mol. Biol. Evol.* 30(5):1196–1205



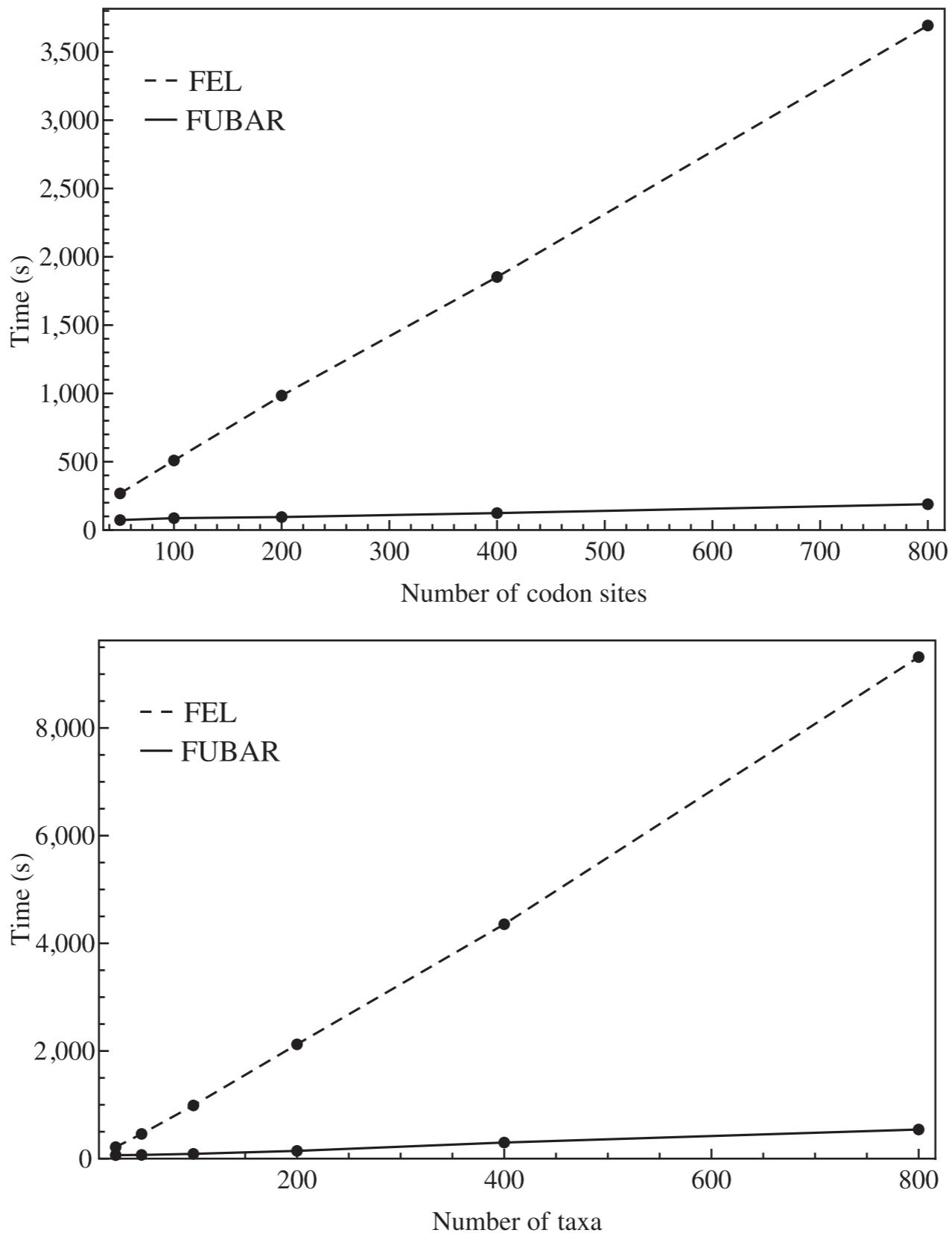
Fitting a small number (4) of  
dN and dS values directly  
*with post-hoc error estimates*



Using a FUBAR grid



Hepatitis E Virus Genotype 4 ORF3



**Fig. 2.** Execution times for FEL and FUBAR as a function of the number of codon sites (top) and number of taxa (bottom).

*FUBAR is dramatically faster (and as good or better)*

**Table 2.** Run Time Comparisons between Different Selection Detection Methods on 16 Empirical Data Sets, Sorted on the Duration of the FUBAR Run.

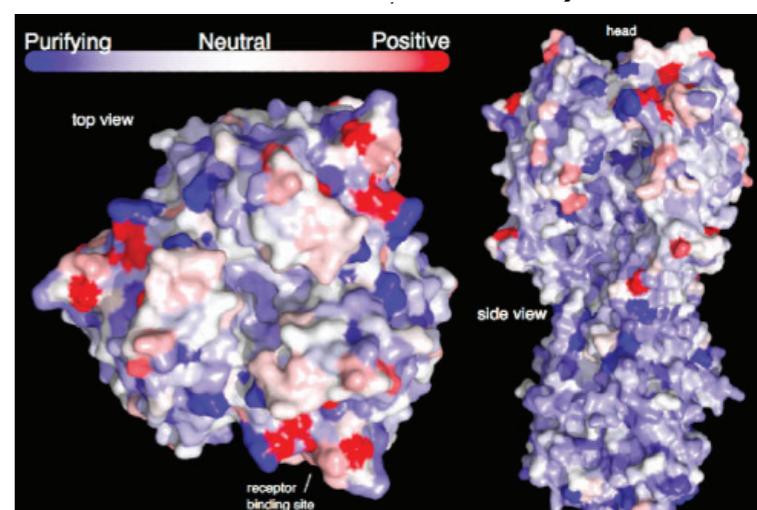
Data Set	Taxa	Codons	Mean Divergence Subs/Site	FUBAR Run Times (s)	Run Times (Times Slower than FUBAR)			
					FEL	REL	PAML M2a	PAML M8
Echinoderm H3	37	111	0.33	40	5.1	12.0	7.1	46.1
Flavivirus NS5	18	342	0.48	45	8.6	4.5	9.3	25.5
<i>Drosophila</i> adh	23	254	0.26	53	3.4	4.0	2.7	4.3
West Nile virus NS3	19	619	0.13	58	6.1	5.9	37.2	<u>105.5</u>
Hepatitis D virus Ag	33	196	0.29	59	4.0	3.3	10.1	22.4
Primate lysozyme	19	130	0.08	62	0.5	3.0	0.7	1.8
Vertebrate rhodopsin	38	330	0.34	62	12.0	4.9	8.4	18.2
Japanese encephalitis virus env	23	500	0.13	68	4.8	8.8	1.6	4.0
Mammalian $\beta$ -globin	17	144	0.38	74	1.5	8.4	2.3	5.6
Abalone sperm lysin	25	134	0.43	78	1.9	3.9	3.7	9.3
HIV-1 vif	29	192	0.08	84	2.6	3.8	2.3	4.5
<i>Salmonella</i> recA	42	353	0.04	102	2.1	2.9	2.6	12.3
Camelid VHH	212	96	0.27	120	6.3	17.2	<u>141.0</u>	<u>311.1</u>
Diatom SIT	97	300	0.54	136	10.2	5.1	21.5	19.3
Influenza A virus H3N2 HA	349	329	0.04	210	15.0	14.4	<u>221.1</u>	<u>616.4</u>
HIV-1 rt	476	335	0.08	278	15.2	14.4	$\emptyset^a$	$\emptyset^a$

NOTE.—Run times that are at least 10 times greater than those of FUBAR are italicized, and those at least 100 times greater are underlined.

<sup>a</sup>PAML reported an error regarding too many ambiguities in the data set.

*FUBAR is dramatically faster (and as good or better)*

We reconstructed the phylogeny for **3,142** complete H3 nucleotide sequences isolated from humans using FastTree 2. The FUBAR selection analysis (which we restricted to 10 CPUs, just as for the timing comparisons) took **one and a half hours**.



# Fast site-level analysis (FUBAR): no branch to branch variation; pervasive diversifying selection; random effects

## WNV NS3

THE EXPECTED NUMBER OF FALSE POSITIVES IS 0.01 (95% CI: [0-0]).

Codon	$\alpha$	$\beta$	$\beta-\alpha$	Posterior Prob $\beta>\alpha$	Emp. Bayes Factor	PSRF	N <sub>eff</sub>	3D rate plot?
249	0.138179	1.51208	1.3739	0.988732	622.977	1.03135	144.342	<a href="#">[SVG]</a> <a href="#">[PNG]</a>

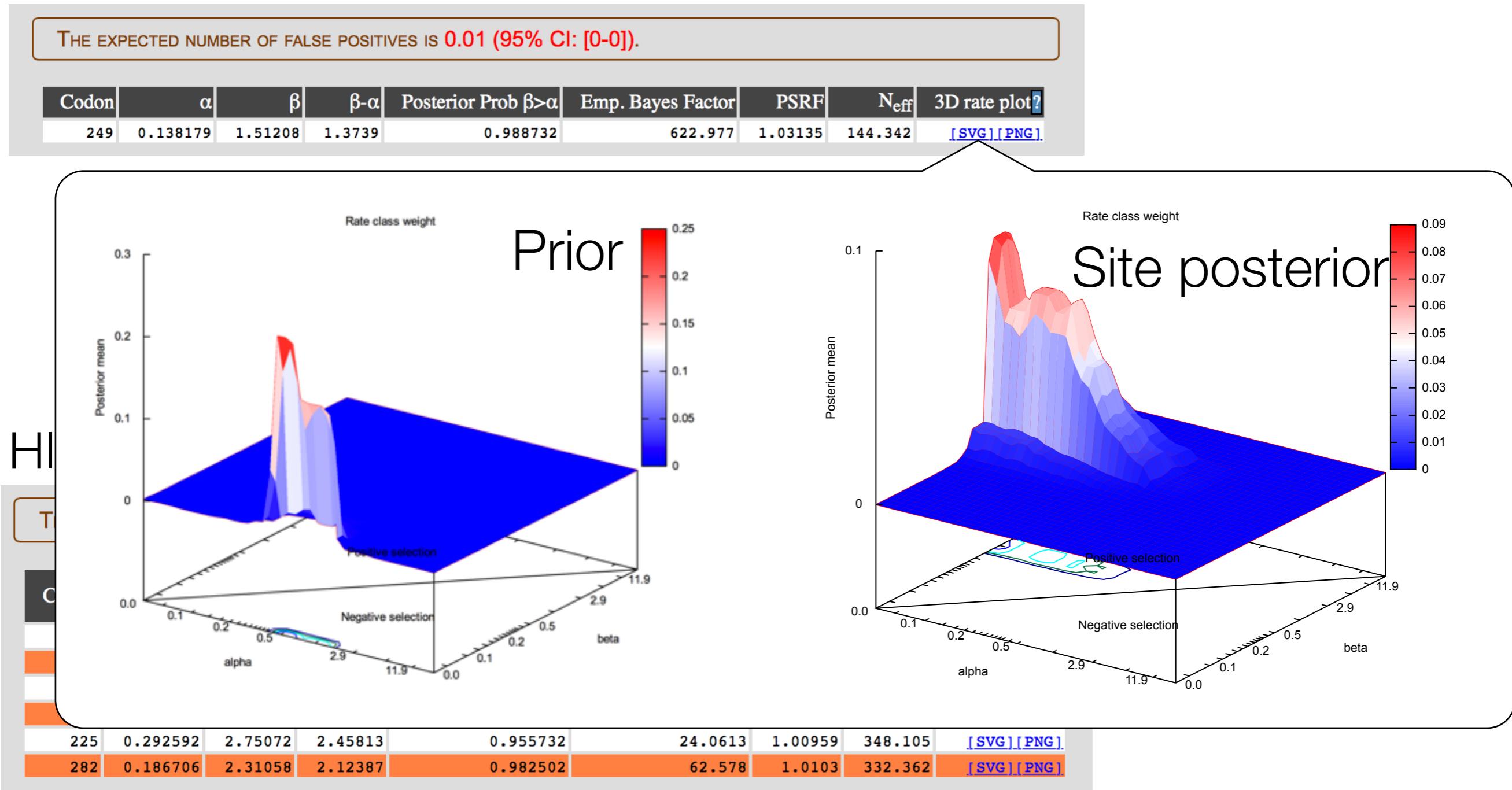
## HIV-1 env

THE EXPECTED NUMBER OF FALSE POSITIVES IS 0.20 (95% CI: [0-1]).

Codon	$\alpha$	$\beta$	$\beta-\alpha$	Posterior Prob $\beta>\alpha$	Emp. Bayes Factor	PSRF	N <sub>eff</sub>	3D rate plot?
161	0.401387	5.64609	5.2447	0.974565	42.7026	1.00373	576.851	<a href="#">[SVG]</a> <a href="#">[PNG]</a>
165	0.311605	2.85981	2.54821	0.963007	29.0123	1.00897	363.228	<a href="#">[SVG]</a> <a href="#">[PNG]</a>
203	0.399165	4.27264	3.87348	0.968713	34.5068	1.00481	514.136	<a href="#">[SVG]</a> <a href="#">[PNG]</a>
204	0.31539	2.80044	2.48505	0.951136	21.6933	1.00947	350.933	<a href="#">[SVG]</a> <a href="#">[PNG]</a>
225	0.292592	2.75072	2.45813	0.955732	24.0613	1.00959	348.105	<a href="#">[SVG]</a> <a href="#">[PNG]</a>
282	0.186706	2.31058	2.12387	0.982502	62.578	1.0103	332.362	<a href="#">[SVG]</a> <a href="#">[PNG]</a>

# Fast site-level analysis (FUBAR): no branch to branch variation; pervasive diversifying selection; random effects

## WNV NS3



# FUBAR results

---

- **West Nile Virus NS3 protein**
  - A single site (**249**, same as in Brault *et al*) with significant evidence of pervasive diversifying selection.
- **HIV-1 transmission pair**
  - 6 sites with significant evidence of **pervasive** diversifying selection.

# Current suggested best practices.

There are lots of methods you could use to study positive selection, including about 10 developed by our group. The field is still evolving, and this is our current suggestions of what to do with your data, depending on the question you want to answer.

Question	Method	Output
Is there episodic selection anywhere in my gene (or along a set of branches known a priori)?	Branch-site unrestricted statistical test of episodic diversification (BUSTED).	<ul style="list-style-type: none"><li>• p-value for gene-wide selection</li><li>• inferred dN/dS distributions</li><li>• a “quick and dirty” scan of sites where selection could have operated.</li></ul>
Are there branches in the tree where some sites have been subject to diversifying selection? <b>Also:</b> inferring ancient divergence times.	Adaptive branch site random effects likelihood (aBSREL)	<ul style="list-style-type: none"><li>• p-values for each branch</li><li>• dN/dS distributions for each branch</li><li>• evolutionary process complexity</li></ul>
Are there sites in the alignment where some of the branches have experienced diversifying selection?	Mixed effects model of evolution (MEME)	<ul style="list-style-type: none"><li>• p-values for each site</li><li>• dN/dS distributions for each site</li></ul>
Are there sites which have experienced diversifying selection <b>and</b> my alignment is large?	Fast unconstrained bayesian analysis of selection (FUBAR)	<ul style="list-style-type: none"><li>• Posterior probabilities of selection at each site</li><li>• An estimate of the the gene-wide dN/dS distribution</li></ul>
Are parts of the tree evolving with different selective pressures relative to other parts of the tree?	RELAX (a test for relaxed selection)	<ul style="list-style-type: none"><li>• p-value for whether or not there is relaxed or intensified selection</li><li>• inferred dN/dS distributions for different branch sets</li><li>• more flexible distribution companions possible</li></ul>

# Recombination

---

- Affects a large variety of organisms, from viruses to mammals (e.g. gene family evolution)
- Manifests itself by incongruent phylogenetic signal
- This can be exploited to detect which sequence regions recombined and which sequences were involved
- Recombination can influence or even mislead selection detection methods.
- Using an incorrect tree to analyze a segment of a recombinant analysis can bias **dS** and **dN** estimation
- The basic intuition is that an incorrect tree will generally break up identity by descent and hence make it appear as if more substitutions took place than did in reality.

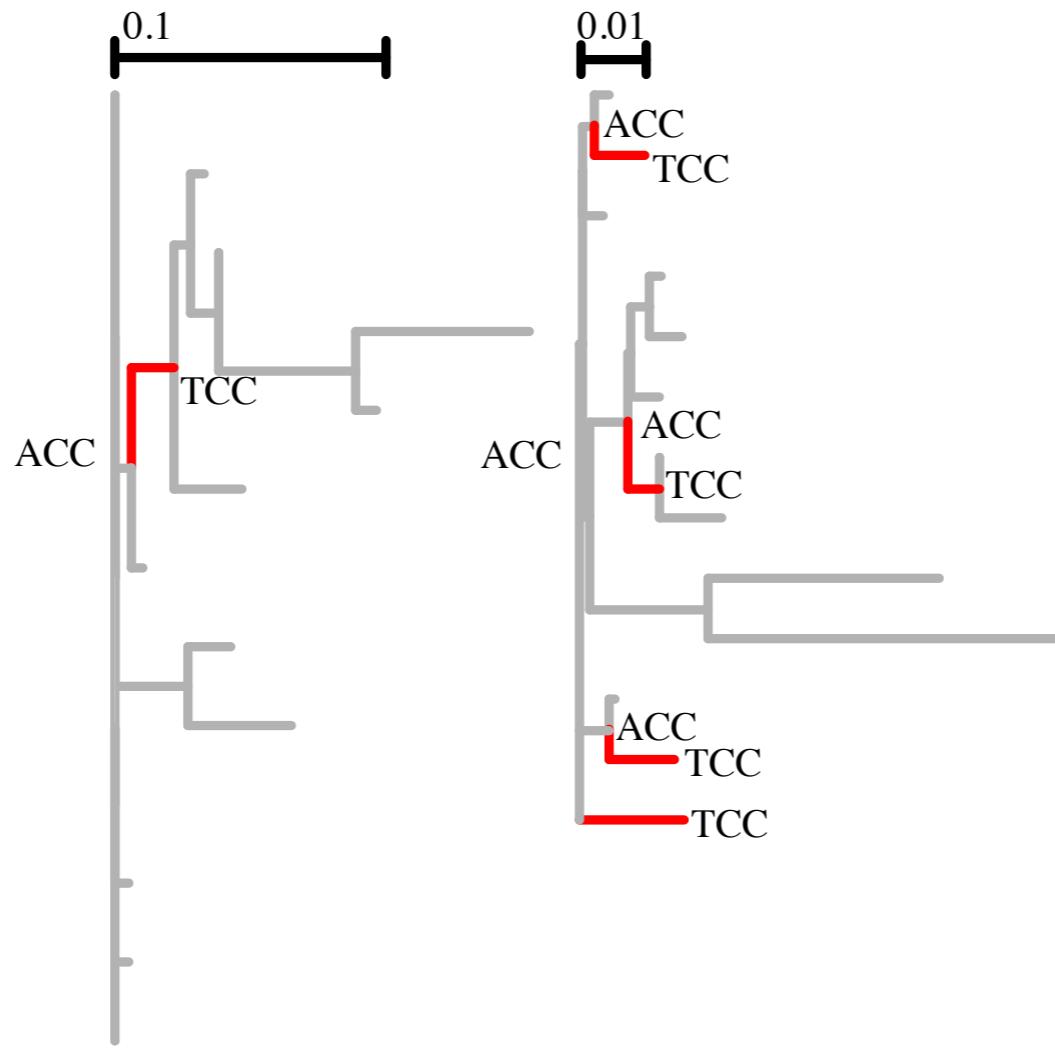


Figure 4.2: The effect of recombination on inferring diversifying selection. Reconstructed evolutionary history of codon 516 of the Cache Valley Fever virus glycoprotein alignment is shown according to GARD inferred segment phylogeny (left) or a single phylogeny inferred from the entire alignment (right). Ignoring the confounding effect of recombination causes the number of nonsynonymous substitutions to be overestimated. A fixed effects likelihood (FEL, Kosakovsky Pond and Frost (2005)) analysis infers codon 516 to be under diversifying selection when recombination is ignored ( $p = 0.02$ ), but not when it is corrected for using a partitioning approach ( $p = 0.28$ ).

# Accounting for recombination

---

- First screen the alignment to find putative non-recombinant fragments (e.g. using GARD)
- Apply a model-based test (MEME, FUBAR) using multiple phylogenies (one per fragment), but inferring other parameters (e.g. nucleotide substitution biases and base frequencies) from the entire alignment
- This has been shown to work very well on simulated and empirical data
- This approach does not work for analyses assuming a single tree (BUSTED, aBSREL).

**Table 4.** Effect of correcting for recombination when using fixed effects likelihood to detect positively selected sites.

Virus and gene	Positively Selected Codons	
	Uncorrected FEL	Corrected FEL
Cache Valley G	212,516,546,551	None
Canine Distemper H	<b>158, 179, 264, 444</b>	<b>179, 264, 444, 548</b>
Crimean Congo hemm. fever NP	<b>195</b>	<b>9,195</b>
Hantaan G2	None	None
Human Parainfluenza (1) HN	<b>37,91, 358, 556</b>	<b>91, 358</b>
Influenza A (human H2N2) HA	<b>87, 166, 252, 358</b>	<b>87, 147,252, 358</b>
Influenza B NA	<b>42,106,345,436</b>	<b>42,106,345,436</b>
Mumps F	<b>57, 480</b>	<b>57, 480</b>
Mumps HN	399	None
Newcastle disease F	<b>1,4,5,7,16,18,108,516</b>	<b>1,5,7,16,108,493,505</b>
Newcastle disease HN	<b>2,54,58,228,262,284,306,471</b>	<b>2,58,228,262,284,306,471</b>
Newcastle disease N	<b>425, 430, 466</b>	<b>425, 430, 462, 466</b>
Newcastle disease P	<b>12,56,65,174,179,188,189, 204, 208, 213,217,218,239,306,332</b>	<b>56, 65, 146, 153, 174, 179, 189, 193, 204,208, 213, 218, 261,306,332</b>
Puumala NP	79	None

Test  $p < 0.1$  was used to classify sites as selected. Codon sites found under selection by both methods are shown in bold.

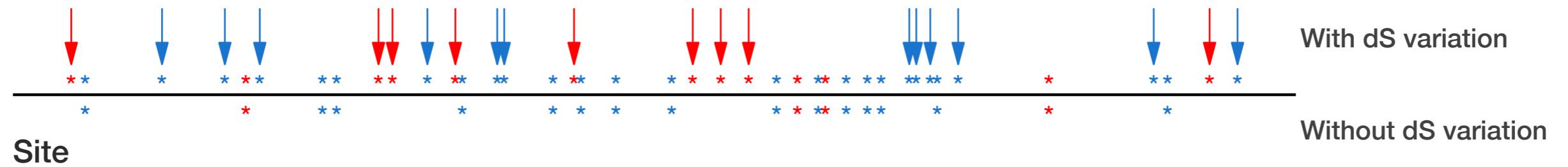
# Synonymous rate variation

- **dS** = constant for all sites (assumed by many models); this assumption appears to be nearly universally violated in biological data, due to e.g. secondary structure, localized codon usage bias, overlapping reading frames, etc.
- This can lead to, e.g. incorrect identification of relaxed constraint as selection
- FUBAR and MEME fully account for **dS** variation; BUSTED and aBSREL provide experimental support.

**Table 1**  
**Data Sets Analyzed for Presence of Synonymous Rate Variation**

Data	Reference	Sequences	Codons	MG94 × REV Nonsynonymous GDD 3		MG94 × REV Dual GDD 3 × 3		<i>P</i> Value	ΔAIC
				log <i>L</i>	Tree Length	log <i>L</i>	Tree Length		
Sperm lysin	(Yang and Swanson 2002)	25	135	-4,409	2.85 (0.06)	-4,397.3	2.93 (0.06)	0.0001	15.36
Primate COXI	(Seo, Kishino, and Thorne 2004)	21	506	-12,013.3	8.5 (0.22)	-11,976.6	5.8 (0.15)	<0.0001	65.27
Drosophila <i>adh</i>	(Yang et al. 2000)	23	254	-4,586.2	1.41 (0.03)	-4,583.4	1.47 (0.03)	0.23	-2.35
HIV-1 <i>vif</i>	(Yang et al. 2000)	29	192	-3,347.2	0.97 (0.02)	-3,334.4	0.99 (0.02)	<0.0001	17.63
β-globin	(Yang et al. 2000)	17	144	-3,659.3	2.6 (0.08)	-3,649.1	3.3 (0.1)	0.0004	12.43
Influenza A*	(Yang 2000)	349	329	-10,916.5	1.42 (0.002)	-10,860.7	1.42 (0.002)	<0.0001	103.7
Camelid VHH*	(Harmsen et al. 2000)	212	96	-16,540.8	14.9 (0.04)	-16,391.2	14.9 (0.04)	<0.0001	291.24
Encephalitis <i>env</i>	(Yang et al. 2000)	23	500	-6,774.4	0.85 (0.02)	-6,752.8	0.89 (0.02)	<0.0001	35.15
Flavivirus NS5	(Yang et al. 2000)	18	183	-9,137.8	6.3 (0.19)	-9,110.2	7.8 (0.24)	<0.0001	47.25
Hepatitis D antigen	(Anisimova and Yang 2004)	33	196	-5,137.7	1.9 (0.03)	-5,074.2	2.02 (0.03)	<0.0001	118.98

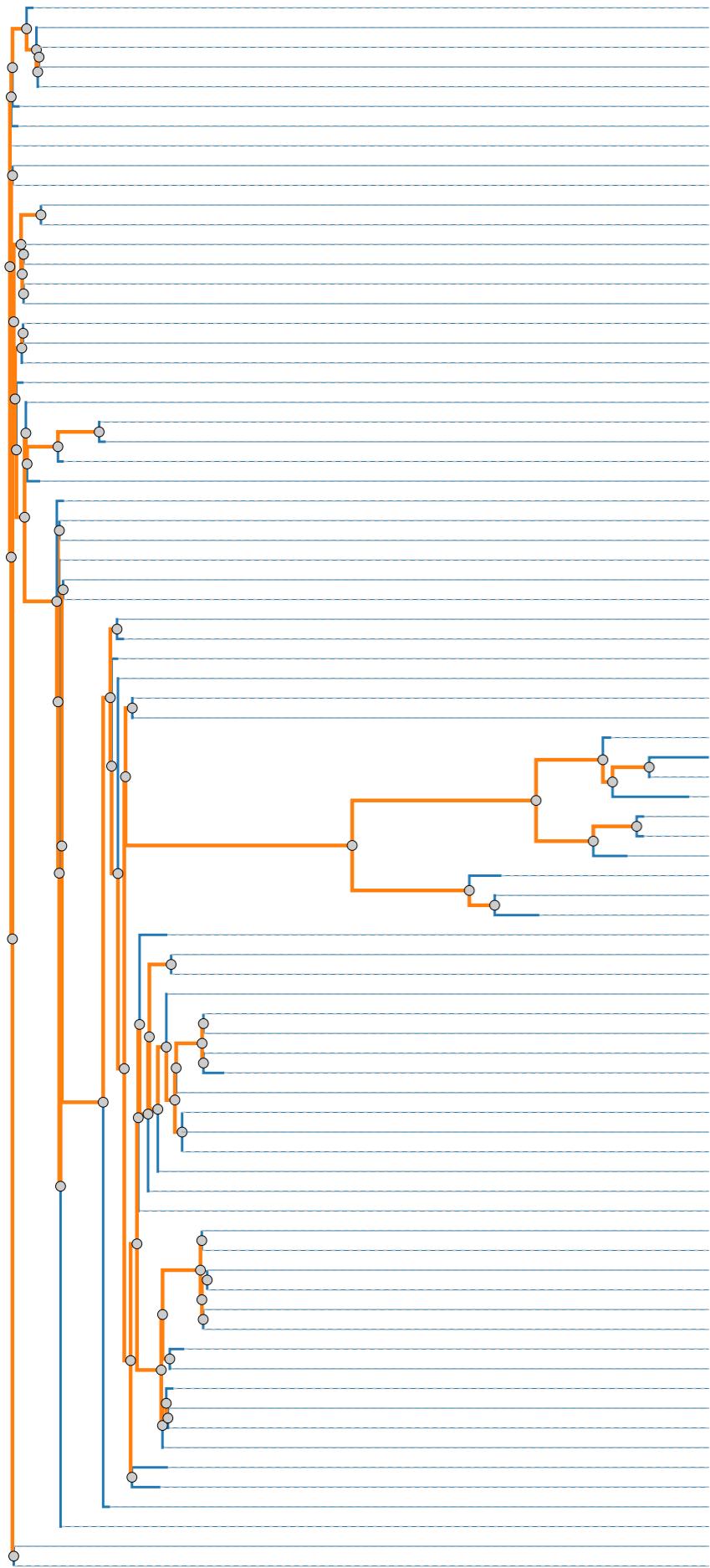
## Sites detected by FEL with and without dS variation



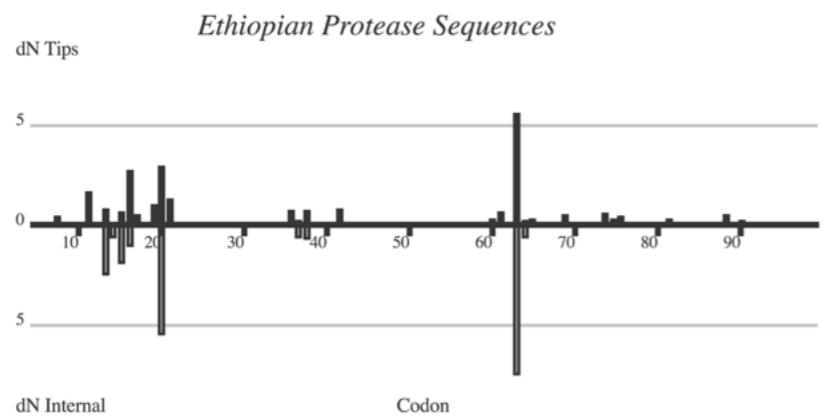
# Interpreting dN/dS for intra-host and intra-species pathogen

---

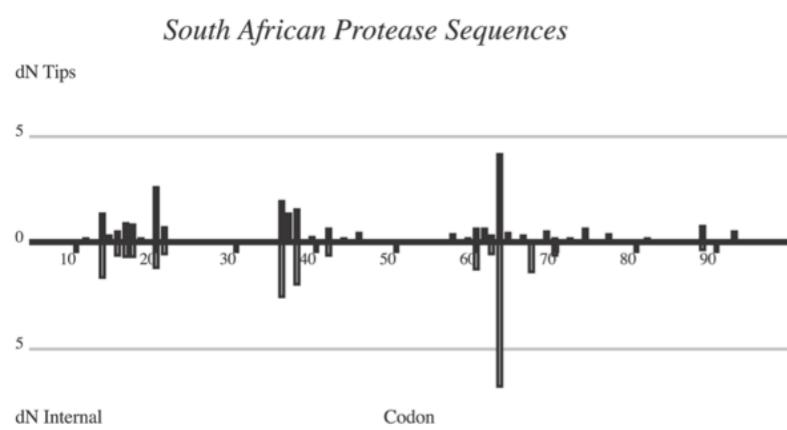
- **dN/dS** can be estimated for all sorts of sequence data (e.g., it has been done for cancer SNP data)
- Traditional interpretation of dN/dS is based on the assumption that **substitution ~ fixation**
- Not the same for intra-species / intra-host pathogens
  - Much of variation is due to polymorphism, or even dead-end mutations
  - This is because selection has not had a chance to “filter” mutations (except for patently deleterious ones)
  - This often manifests as differences in selective “regimes” between tips and internal branches



- Partition a pathogen tree into terminal and internal branches
- Terminal branches potentially include “dead-end” lineages, i.e. those which are maladaptive
- Internal branches include at least one “*transmission*” (intra-species) or “*replication*” (intra-host) events: stronger action of selection
- Focusing on a subset of branches can allow one to interpret dN/dS more precisely



Site Class	Codon Count	dN/dS along Terminal Branches	Is dN/dS elevated when dN internal > 0
Only dN Tips >0	17	0.49	LRT = 5.55
dN Internal >0	9	1.42	$p_A = 0.018, p_B = 0.02$



Site Class	Codon Count	dN/dS along Terminal Branches	Is dN/dS elevated when dN internal > 0
Only dN Tips >0	18	0.42	LRT = 10.46
dN Internal >0	15	1.19	$p_A = 0.001, p_B = <0.01$

- “... at least half of the amino acid sites selected within individuals are not selected at a population level”
- “... Based on the elevated rate of adaptation within individuals detected at codons subject to population-level selection, relative to the codons where only recent substitutions have been inferred, we conclude that recent substitutions are, on average, maladaptive at the level of the human population”