



# Quantifying Natural Selection in Coding Sequences

**Sergei L Kosakovsky Pond**  
Professor of Biology  
Institute for Genomics and Evolutionary Medicine @ Temple University

✉ [spond@temple.edu](mailto:spond@temple.edu)  
🏡 <http://lab.hyphy.org>  
/github.com/spond  
🐦 [@sergeikp](https://twitter.com/sergeikp)

# Preliminaries

- Please confirm access to **HyPhy**: <http://hyphy.org/download/>
  - <https://youtu.be/fgNrPbOTpxE>
  - You can do a [datammonkey.org](http://datammonkey.org) based tutorial, but if you have Linux or OS X, you can also do a command line tutorial for more features.
- General user questions and feedback: <https://github.com/veg/hyphy/issues>
- **Datammonkey** web-app:
  - <http://www.datammonkey.org>
  - YouTube example videos (channel HyPhy vision)
  - <https://www.youtube.com/channel/UCIgRnbJjbOWhshe5ThhaWGw/videos>
- Test datasets and practical instructions: [www.hyphy.org](http://www.hyphy.org) (search for “Detect Selection”)
- Example datasets at <https://github.com/veg/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 (FEL, MEME)
- Site-level **pervasive** selection (SLAC, FEL, FUBAR)
- Relaxed or intensified selection (RELAX)
- Detecting **differences** in selective pressure (CONTRAST-FEL)
- Confounding processes (synonymous rate variation, recombination, multiple nucleotide substitutions)
- On the suitability of dN/dS for within-species inference

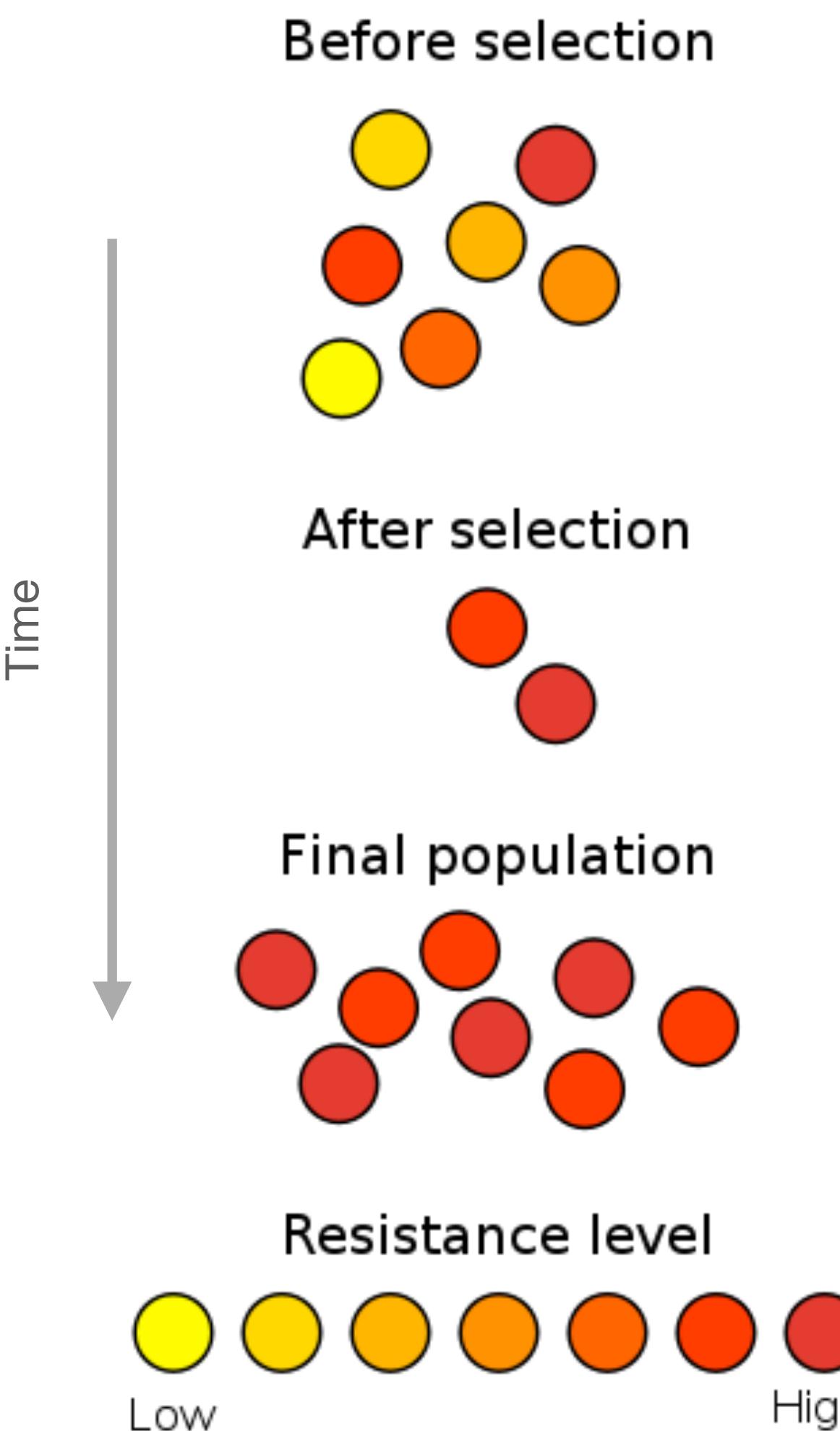
# A bit of trivia

- 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)





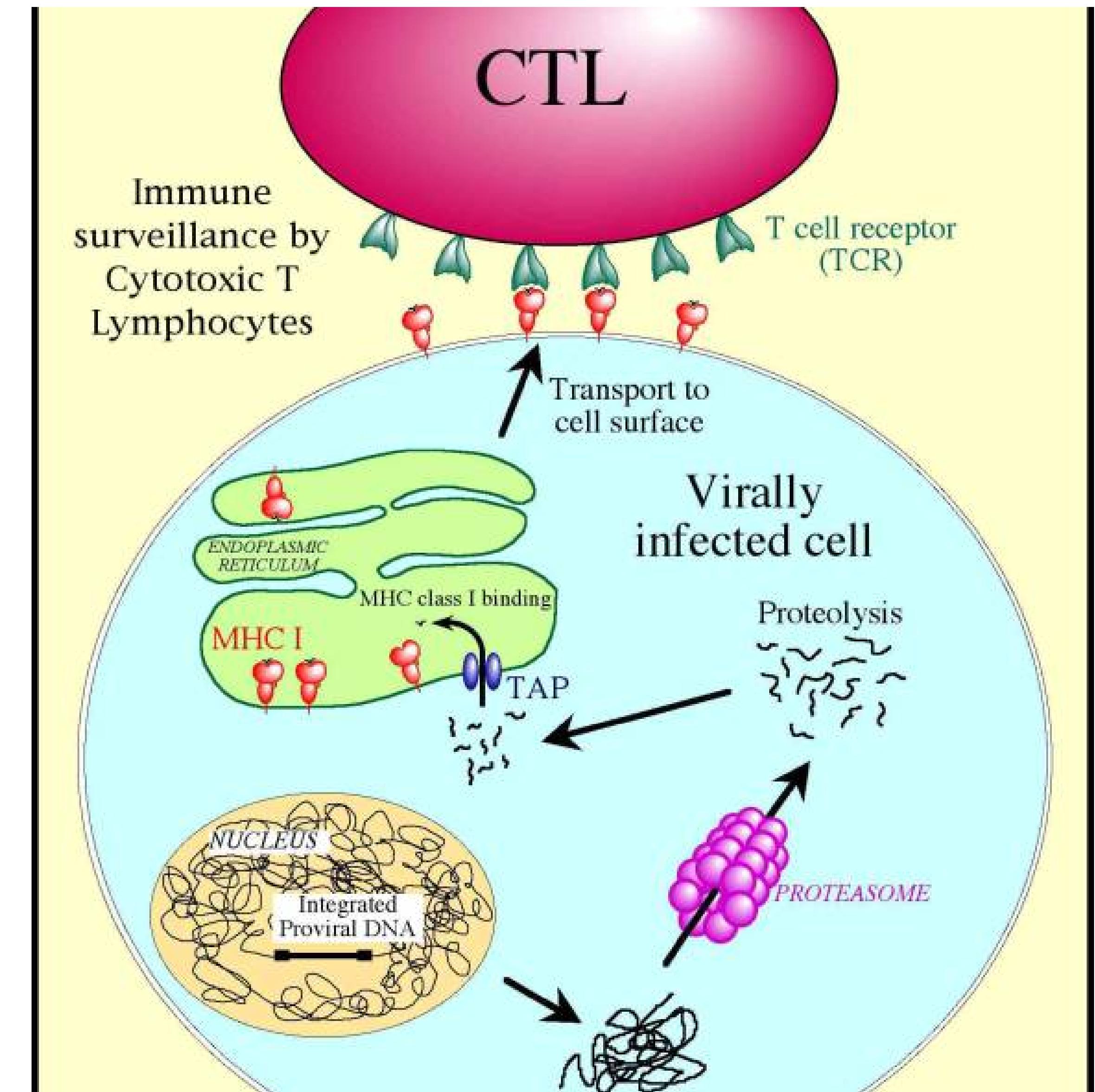
<https://www.youtube.com/watch?v=pIVk4NVIUh8>



<https://www.youtube.com/watch?v=pIVk4NVIUh8>

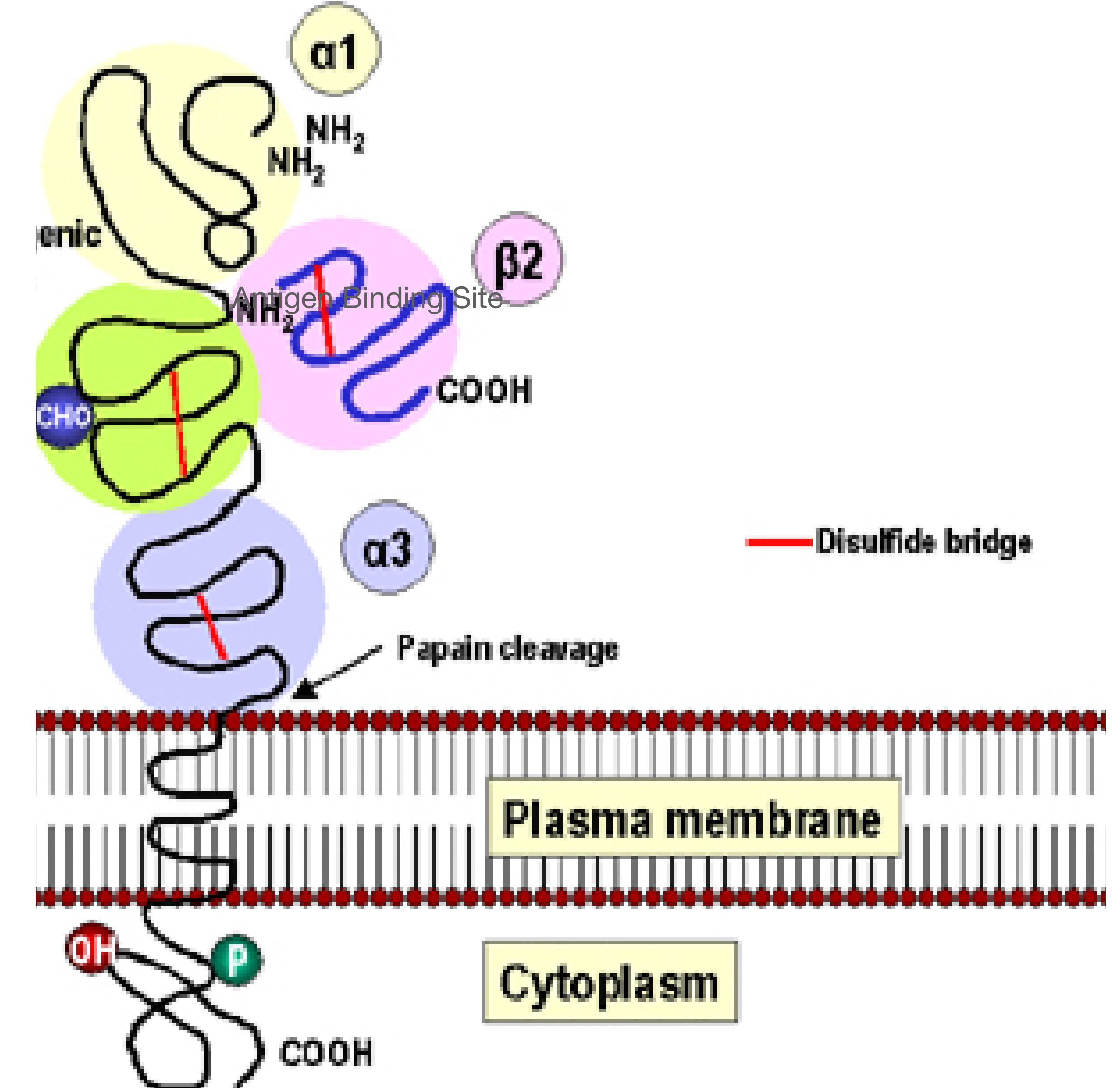
## Example: MHC-restricted CTL killing of infected cells

- 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.



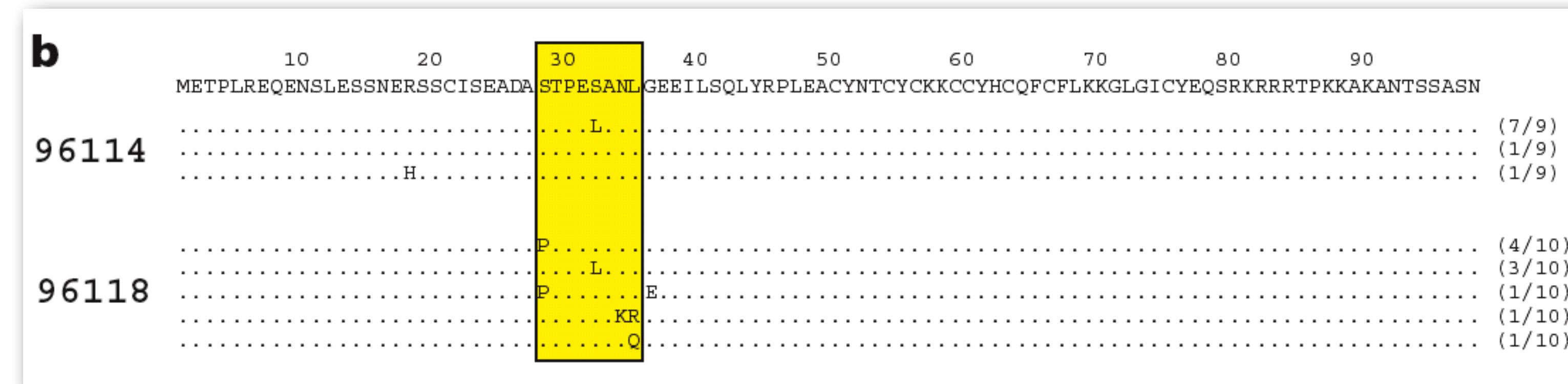
# MHC Class 1 Molecules

- 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



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

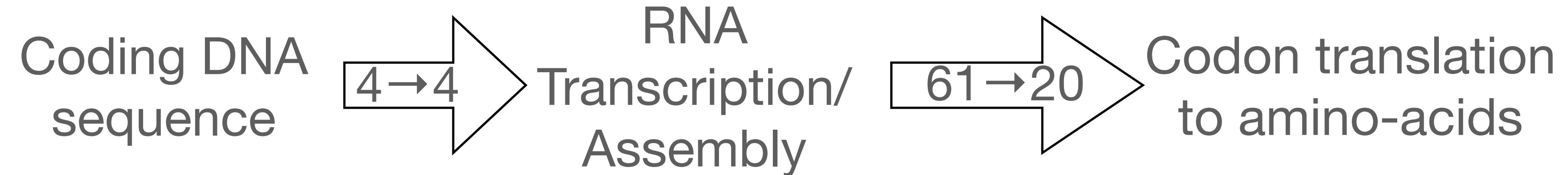
- SIV: the only animal model of HIV (rhesus macaques)
- Experimental infection with MHC-matched strain of SIV
- Virus sequenced from a sample 2 weeks post infection
- Only variation was in an epitope recognized by the MHC
  - CTL escape



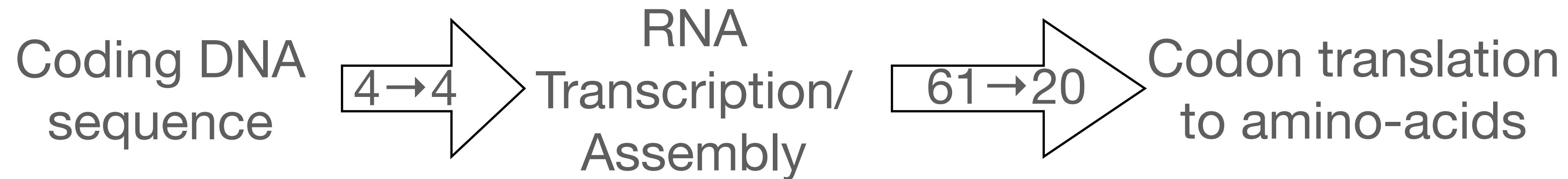
# Key drivers of adaptation in pathogens

- 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

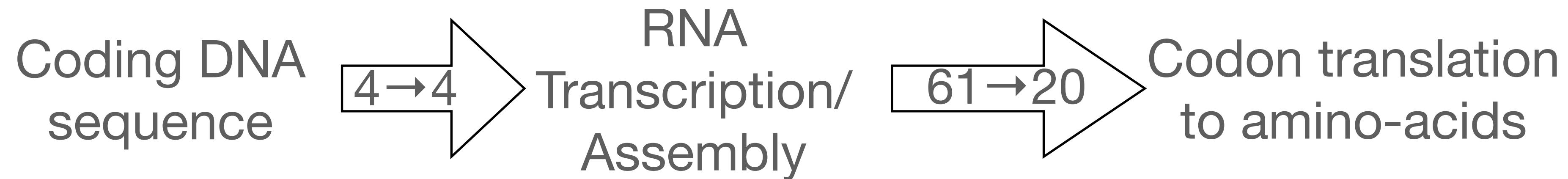


# Evolution of Coding Sequences



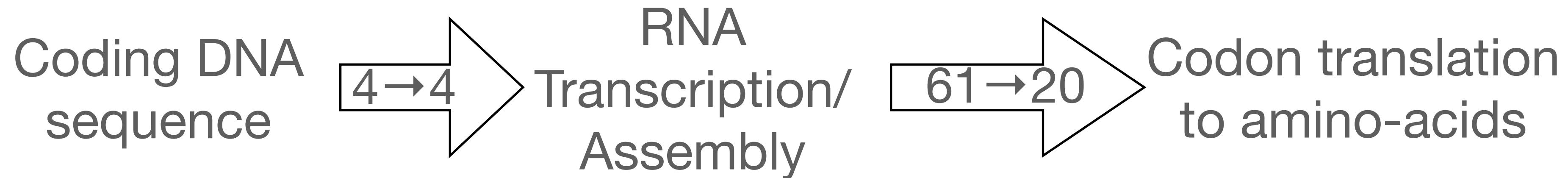
- Proper unit of evolution is a triplet of nucleotides — a **codon**

# Evolution of Coding Sequences



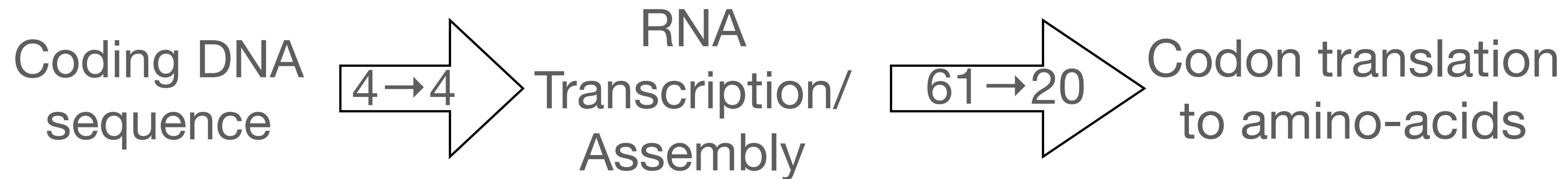
- Proper unit of evolution is a triplet of nucleotides — a **codon**
  - **Mutation** happens at the **DNA level**

# Evolution of Coding Sequences



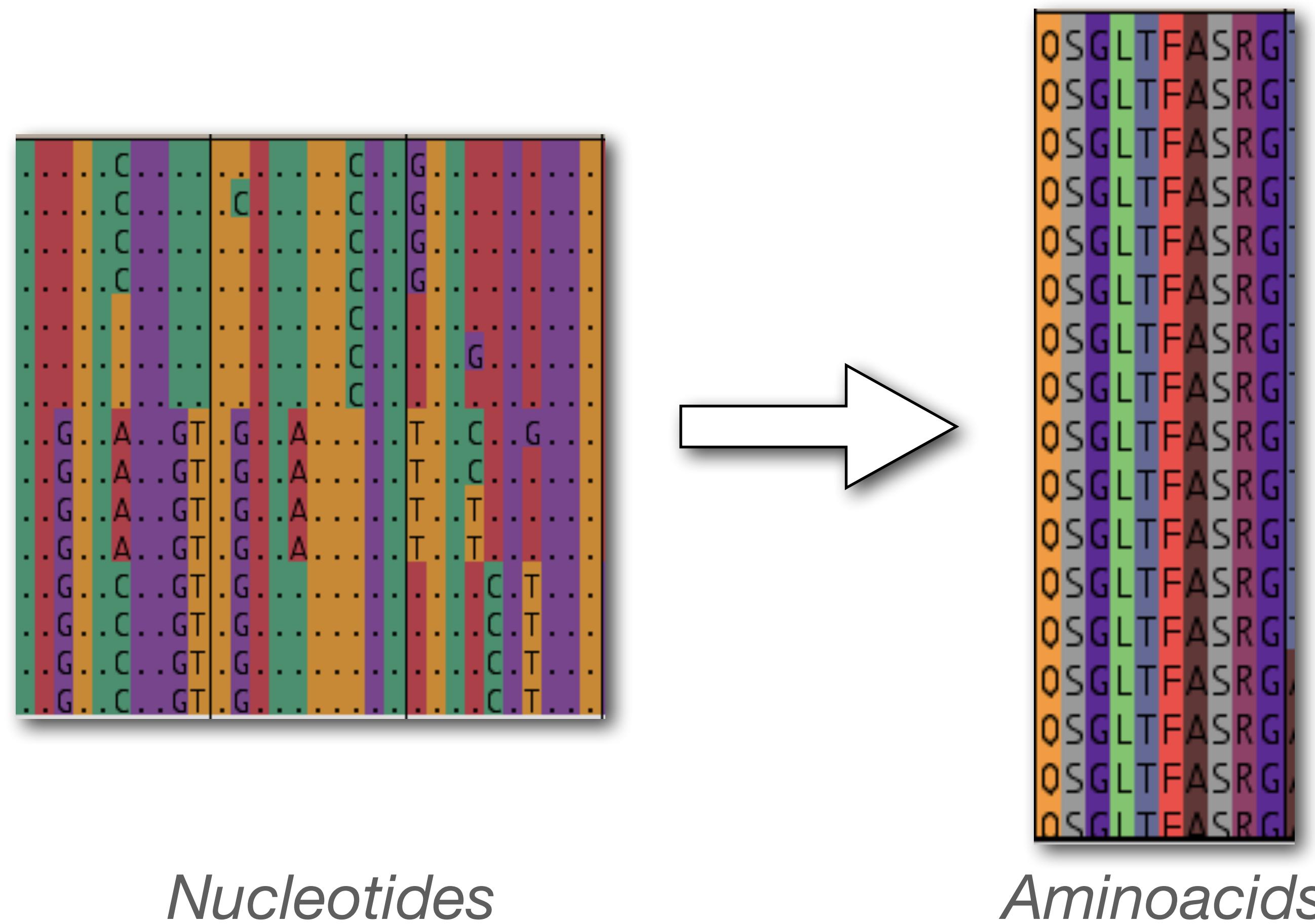
- Proper unit of evolution is a triplet of nucleotides — a **codon**
  - **Mutation** happens at the **DNA level**
  - **Selection** happens (by and large) at the **protein level**

# Evolution of Coding Sequences

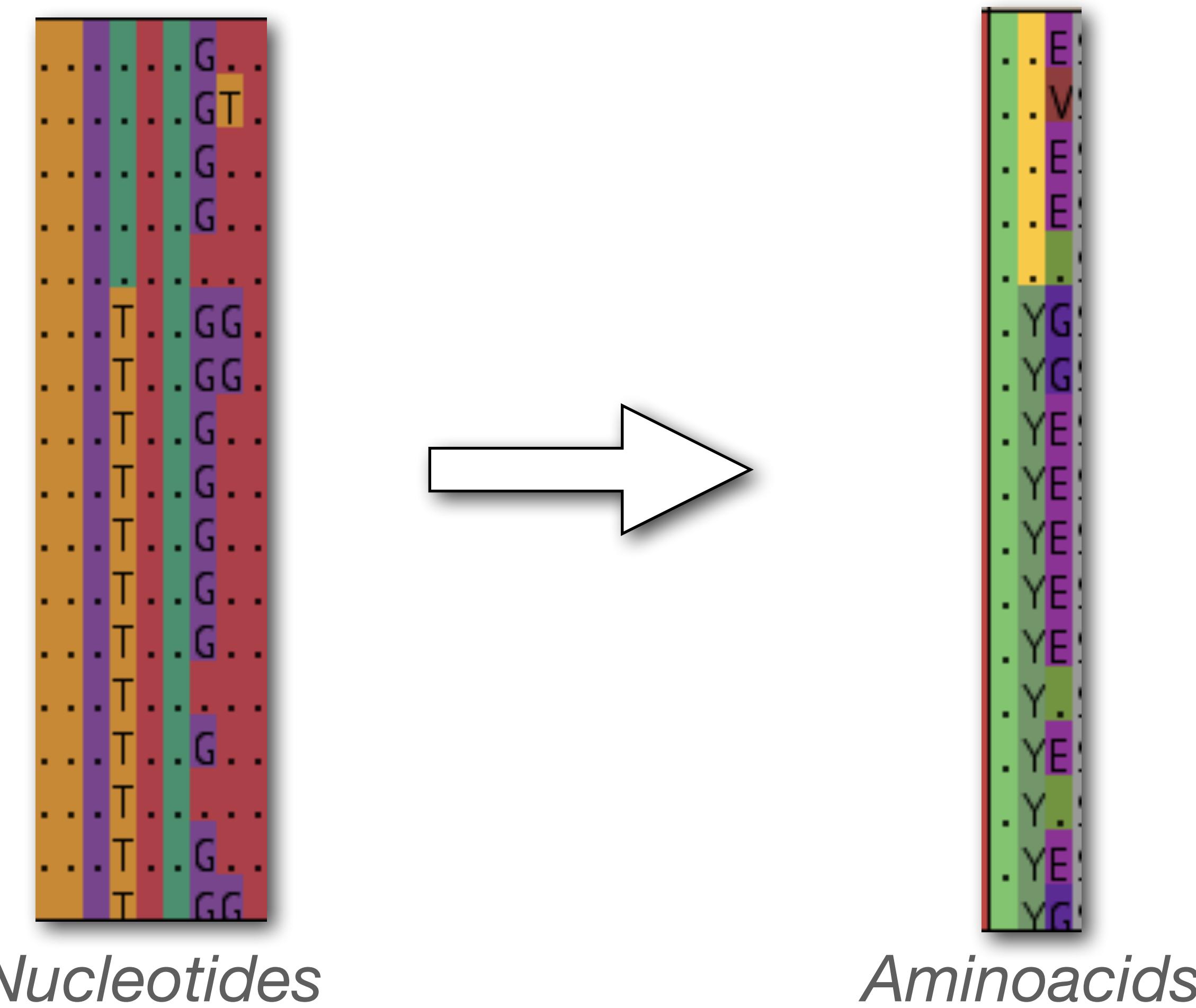


- Proper unit of evolution is a triplet of nucleotides — a **codon**
  - **Mutation** happens at the **DNA level**
  - **Selection** happens (by and large) at the **protein level**
- **Synonymous** (protein sequence **unchanged**) and **non-synonymous** (protein sequence **changed**) substitutions are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

# **Conservation: measles, rinderpest, and peste-de-petite ruminant viruses nucleoprotein.**



# 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**) can serve as the neutral background evolutionary rate
- We can compare the **rate** of accumulation of non-synonymous substitutions (**dN**), which alter the protein sequence, to **dS** and use their ratio to classify the nature of the evolutionary process

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

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

# 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**) can serve as the neutral background evolutionary rate
- We can compare the **rate** of accumulation of non-synonymous substitutions (**dN**), which alter the protein sequence, to **dS** and use their ratio to classify the nature of the evolutionary process

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

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

*What can the denominator proportions depend on?*

# 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	AT <b>C</b>	TTT	AAT <b>T</b>	CAA
T	I	<b>I</b>	F	N	Q
ACA	ATA	<b>ACC</b>	TTT	AAC <b>C</b>	CAA
T	I	<b>T</b>	F	N	Q

# Estimating dS and dN

Consider two **aligned homologous** sequences

ACA	ATA	ATC	TTT	AAT	CAA
T	I	I	F	N	Q
ACA	ATA	ACC	TTT	AAC	CAA
T	I	T	F	N	Q

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			Non-synonymous			To a stop codon		
	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 completely at random**
- Non-random distribution over codon positions
  - **All** second position mutations are non-synonymous
  - **Most** (but not all) 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
<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
<b>Synonymous sites</b>	0	0	<b>1/3</b>
<b>Non-synonymous sites</b>	<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
<b>Synonymous sites</b>	<b>0</b>	<b>0</b>	<b>1/3</b>
<b>Non-synonymous sites</b>	<b>1</b>	<b>1</b>	<b>2/3</b>

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
Asparagine	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	TAA,TAG,TGA	3

# 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	1	1	<b>2/3</b>

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
Asparagine	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	TAA,TAG,TGA	3

8/3 non-synonymous sites (or 7/3 + 1/3 “stop” site)  
1/3 synonymous sites

~4,500 citations

# Nei-Gojobori dN/dS estimate (NG86)

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

Seq1	<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
Seq2	<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
Syn	1	2/3	5/6	1/3	1/3	1/3
NonSyn	2	7/3	13/6	8/3	8/3	7/3

Mean

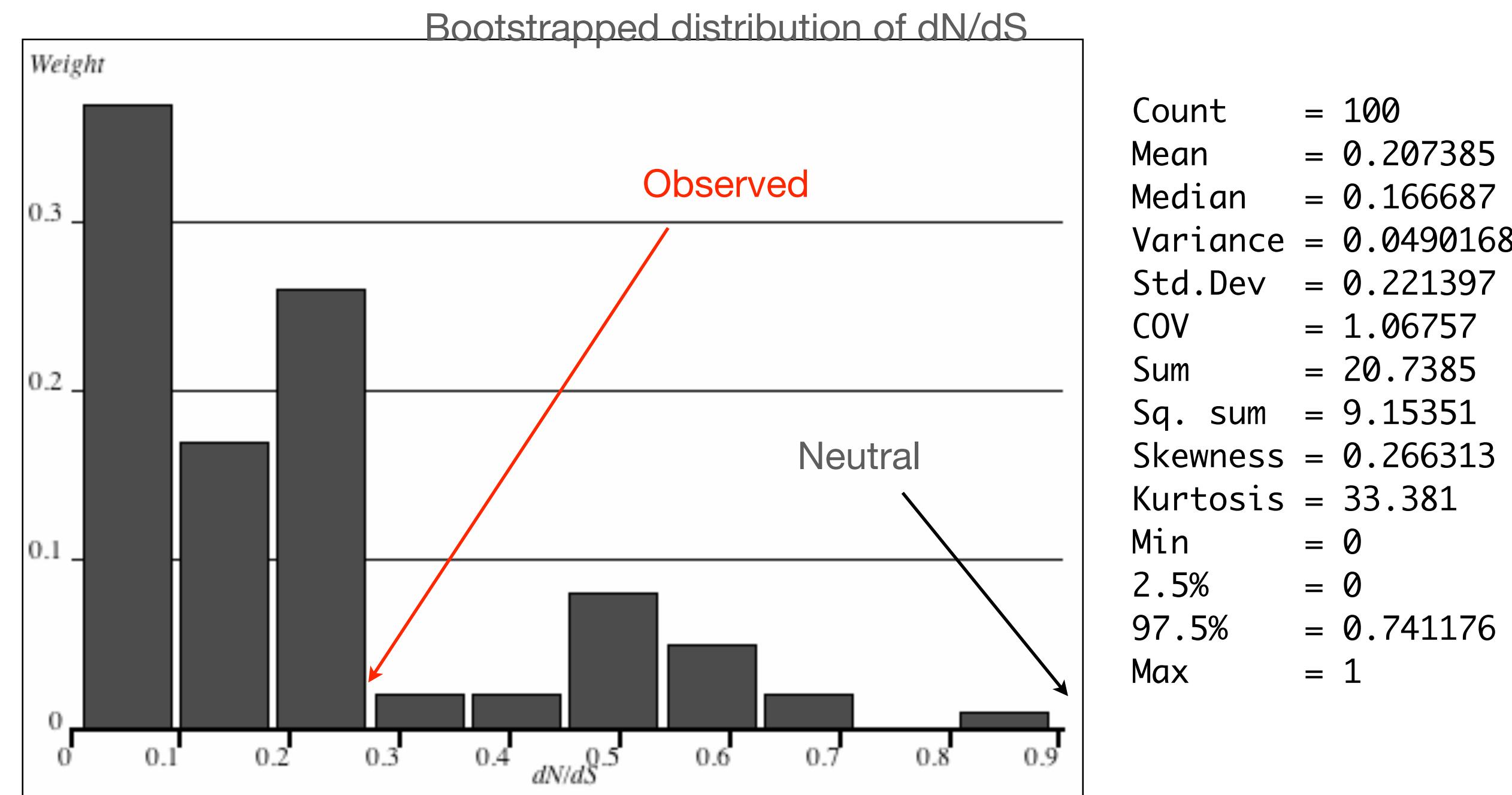
**ES** =  $3\frac{1}{2}$ , **EN** =  $14\frac{1}{6}$ : under neutrality, we expect the ratio of non-synonymous to synonymous substitutions of **EN/ES** ~ 4.05

# 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 under neutrality, we conclude that most non-synonymous mutations are **removed by natural selection**, i.e., the sequences are under **negative selection**
- **If there were more** non-synonymous substitutions than expected, we would conclude that many non-synonymous mutations are **fixed due to natural selection**, i.e., the sequences are under **positive selection**

# NG86 example

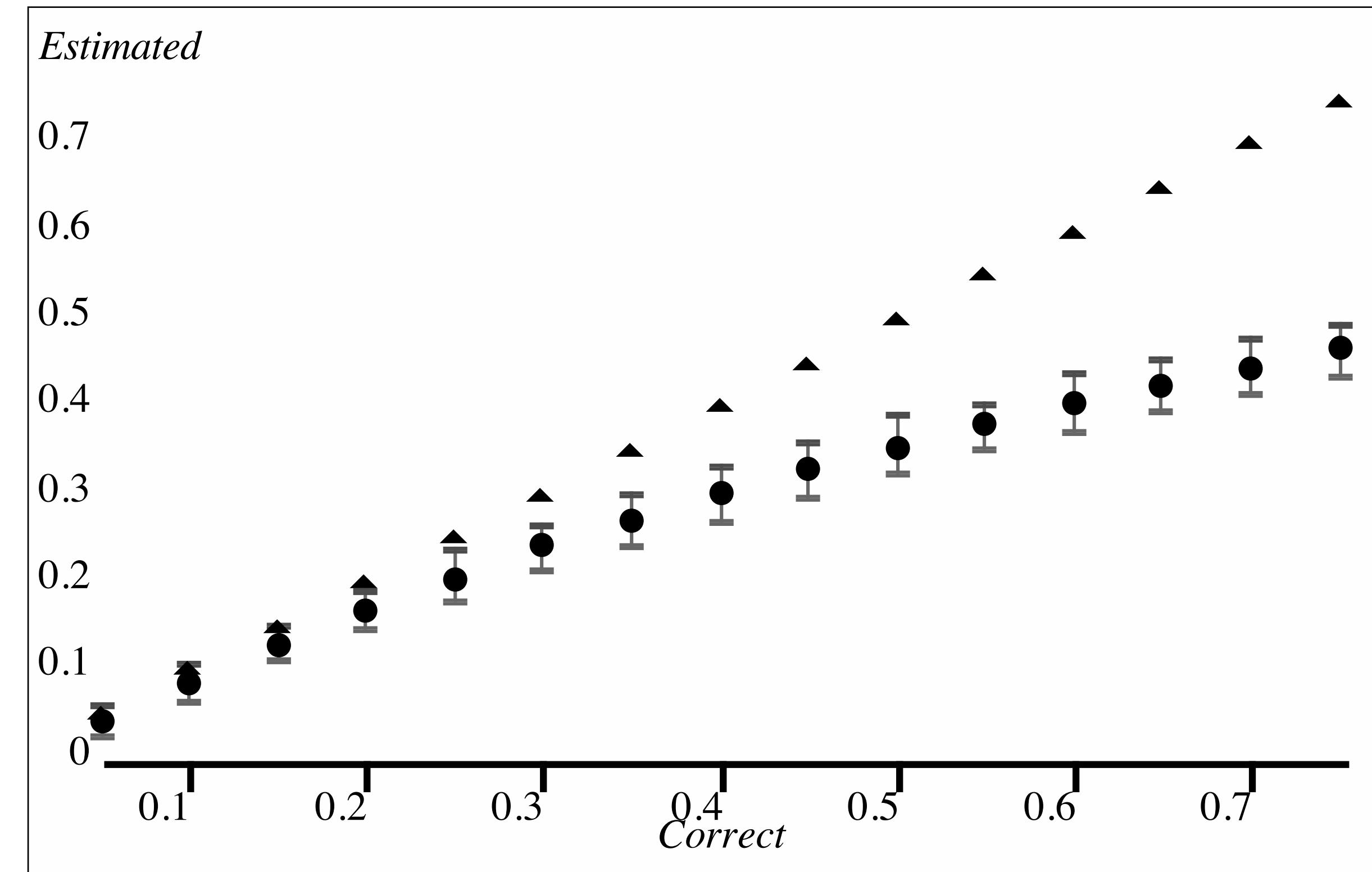
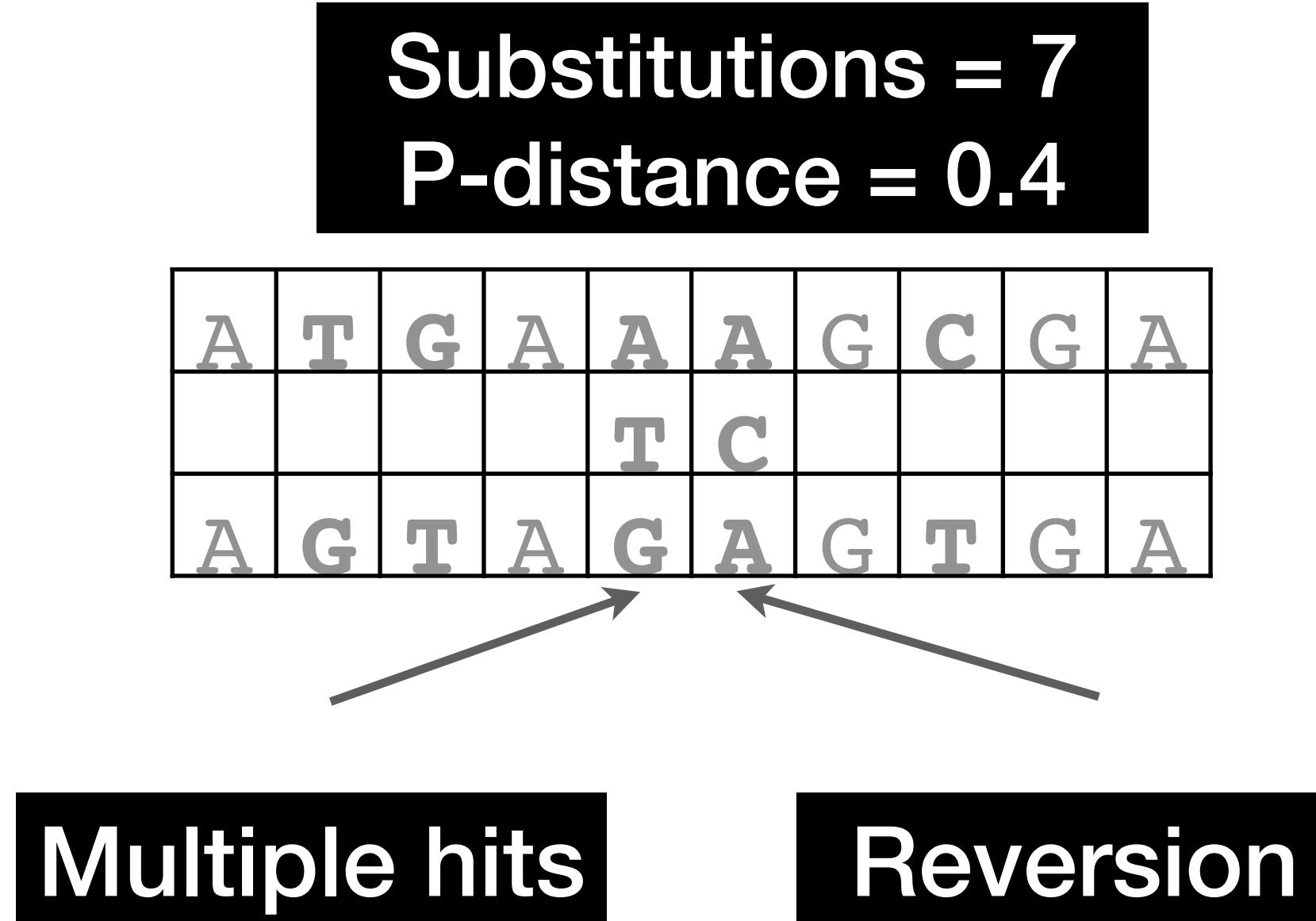
- How reliable is the inference based on only **6 codons**?
- Obtain sampling variance via bootstrap (or by limiting approximations)
- In this case,  $dN/dS$  is **significantly** less than 1.0 ( $p \sim 0.01$ )



# 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 amino-acids.

# 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.

~500 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.

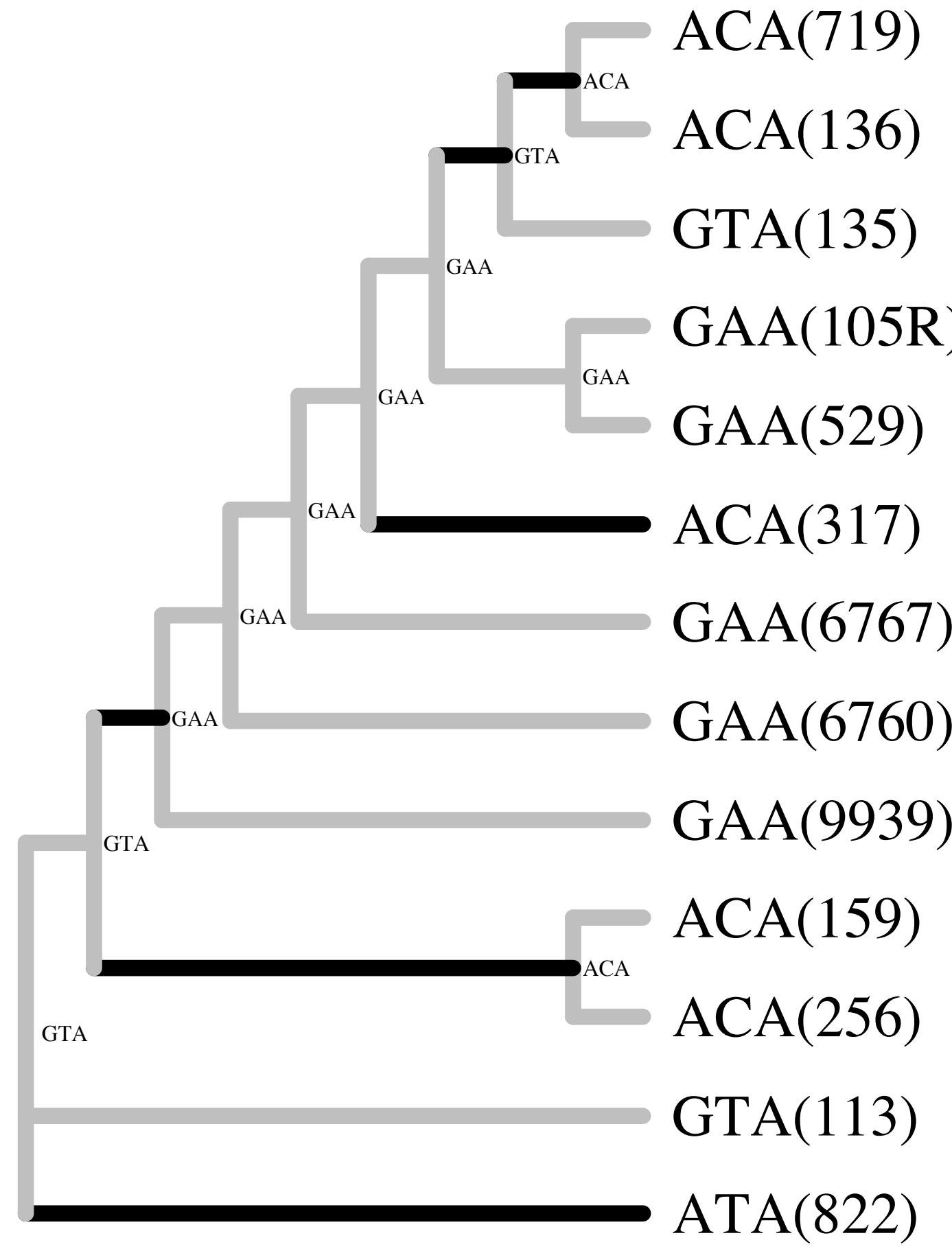


Fig. 1.6. An illustration of SLAC method, applied to a small HIV-1 envelope V3 loop alignment. Sequence names are shown in parentheses. Likelihood state ancestral reconstruction is shown at internal nodes. The parsimonious count yields 0 synonymous and 9 non-synonymous substitutions (highlighted with 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)

~1820 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 the 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).

$\alpha R_{CT}$

$\beta R_{CG}$

$\alpha$  (syn. rate) and  $\beta$  (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.,  $\text{ACT} \Rightarrow \text{AGG}$ ).
  - This can be further relaxed with models supporting multiple nucleotide changes.
- 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

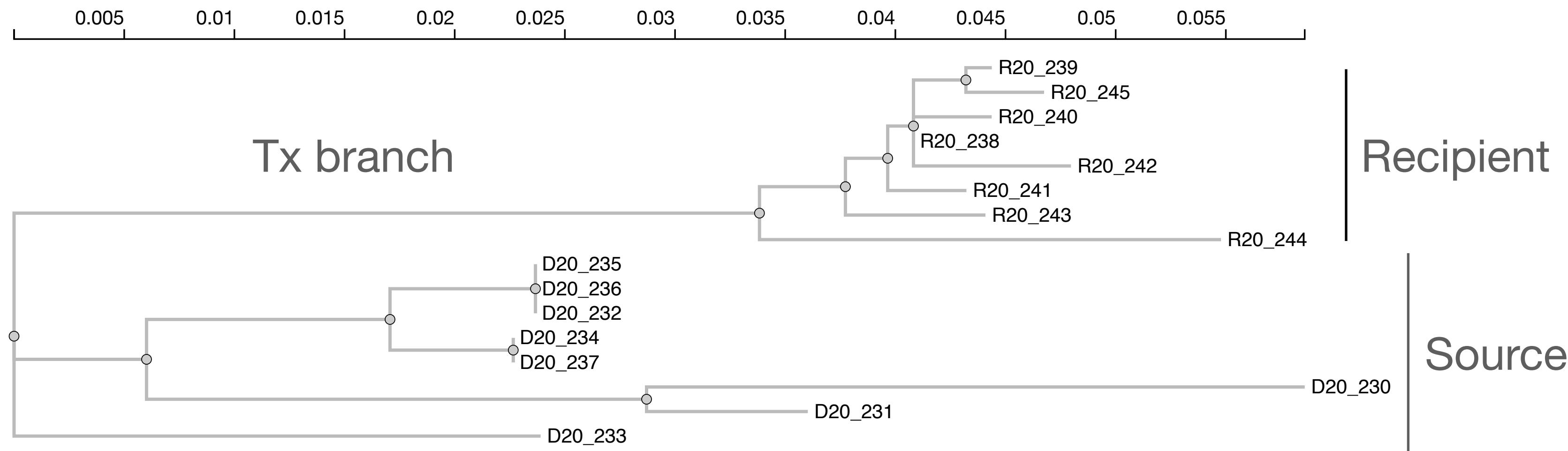
- Using standard MLE approaches it is straightforward to obtain point estimates of  $dN/dS := \beta/\alpha$
- Can also easily test whether or not  $dN/dS > 1$ , or  $< 1$  using the likelihood ratio test (LRT)
- Codon models also support the concepts of synonymous and non-synonymous distances between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

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

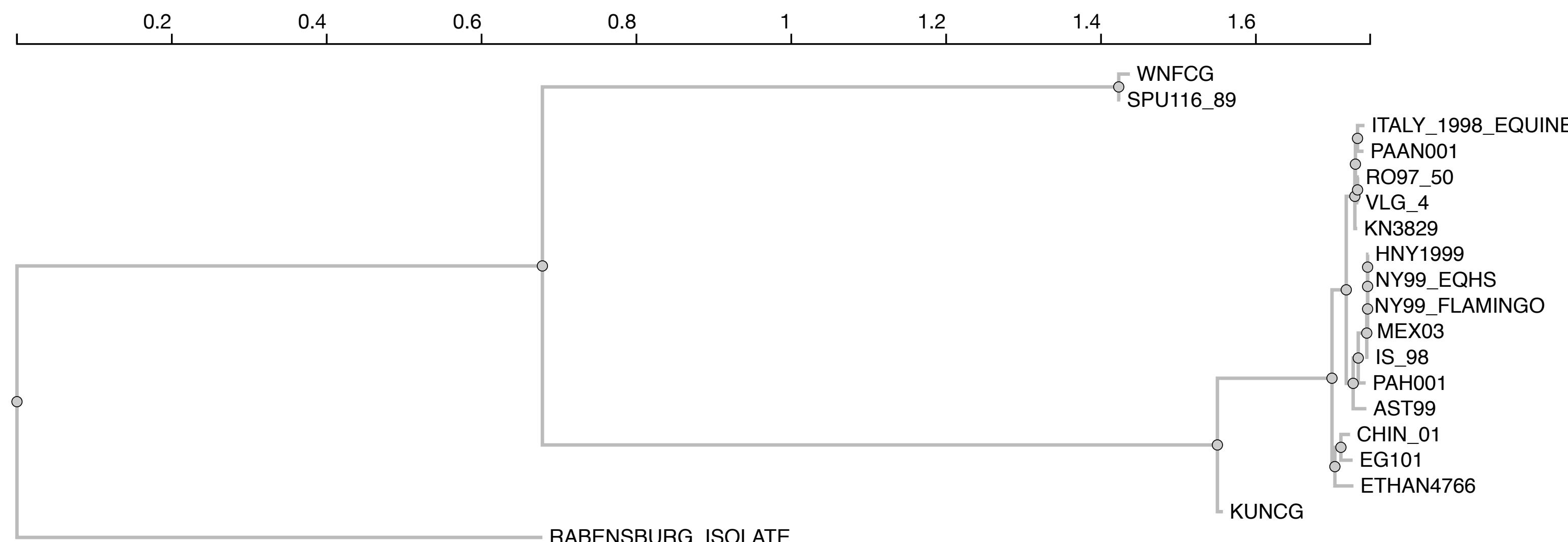
# Three example datasets

- **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)
- **SARS-CoV-2 Spike**
  - Full length spike sequences chosen to represent viral diversity
  - Good example for analyzing selection in population samples with many “dead-end” intra-host variants

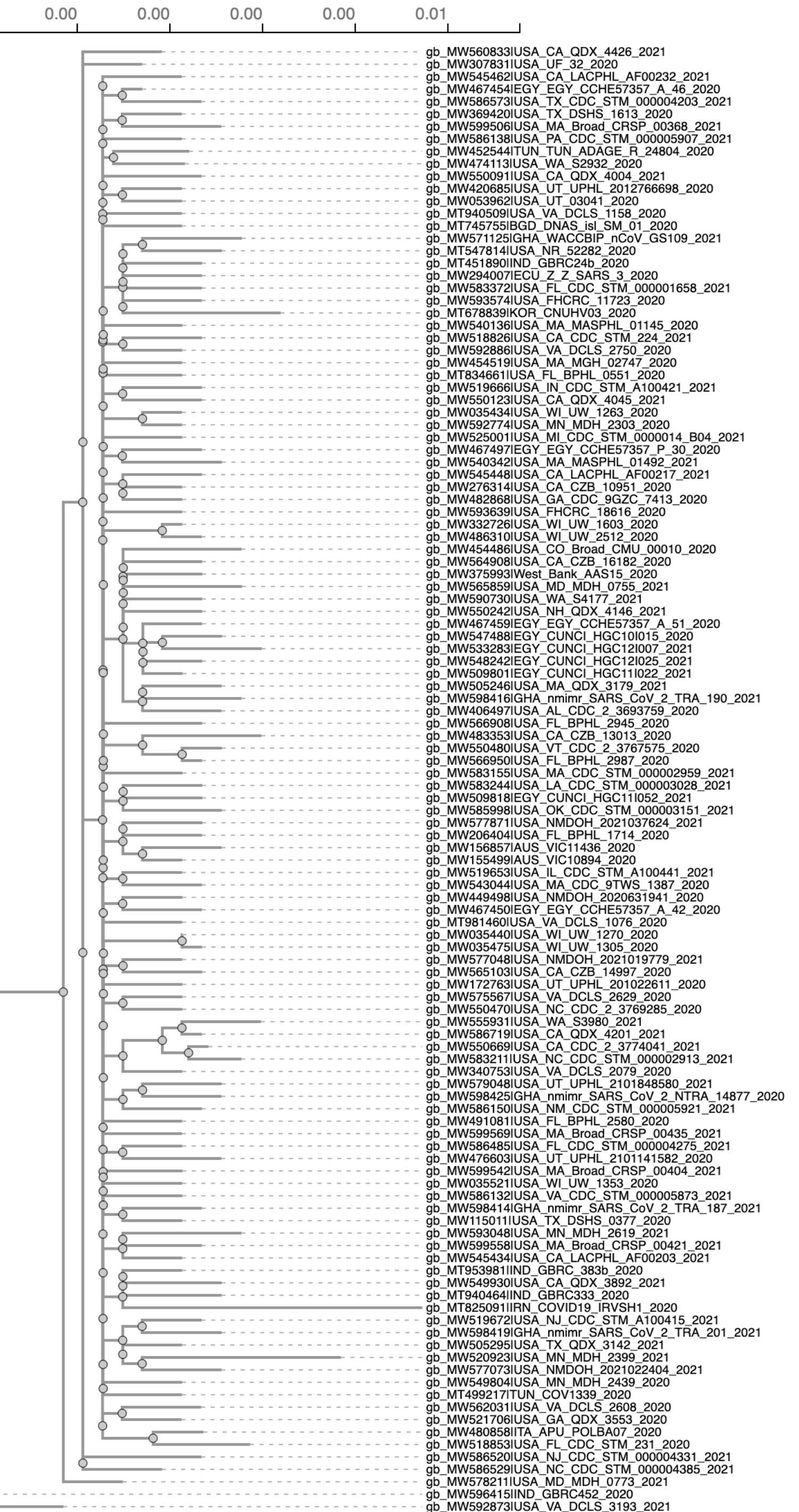
# HIV-1 env



# WN NS3



# SARS-CoV-2 spike



# Information content of the alignments

	WNV NS3	HIV-1 env	SARS-CoV-2 spike
Sequences	19	16	118
Codons	619	288	1273
Tree Length <i>MG94 model, subs/ site</i>	0.67	0.20	0.134

# Information content of the alignments

	WNV NS3	HIV-1 env	SARS-CoV-2 spike
Sequences	19	16	118
Codons	619	288	1273
Tree Length <i>MG94 model, subs/site</i>	0.67	0.20	0.134

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

MSA

Tree

Settings

## HyPhy analysis

MarkDown screen output

JSON file with analysis results



[vision.hyphy.org](http://vision.hyphy.org)



```
$hyphy ~/Development/hyphy-analyses/FitMG94/FitMG94.bf --help
```

Available analysis command line options

Use --option VALUE syntax to invoke

If a [required] option is not provided on the command line, the analysis will prompt for its value  
[conditionally required] options may or not be required based on the values of other options

rooted

Accept rooted trees  
default value: No

code

Which genetic code should be used  
default value: Universal

alignment [required]

An in-frame codon alignment in one of the formats supported by HyPhy

tree [conditionally required]

A phylogenetic tree  
applies to: Please select a tree file for the data:

type

Model type: global (single dN/dS for all branches) or local (separate dN/dS)  
default value: terms.global [computed at run time]  
applies to: Model Type

frequencies

Equilibrium frequency estimator  
default value: CF3x4

lrt

Perform LRT to test which for dN/dS == 1 (global model only)  
default value: No

output

Write the resulting JSON to this file (default is to save to the same path as the alignment file + 'MG94.json')  
default value: fitter.codon\_data\_info[terms.json.json] [computed at run time]

save-fit

Save MG94 model fit to this file (default is not to save)  
default value: /dev/null

```
$hyphy ~/Development/hyphy-analyses/FitMG94/FitMG94.bf --lrt Yes --alignment WestNileVirus_NS3.fas
```

#### Analysis Description

---

Fit an MG94xREV model with several selectable options frequency estimator and report the fit results including dN/dS ratios, and synonymous and non-synonymous branch lengths. v0.2 adds LRT test for dN/dS != 1

- Requirements: in-frame codon alignment and a phylogenetic tree
- Written by: Sergei L Kosakovsky Pond
- Contact Information: spond@temple.edu
- Analysis Version: 0.2

rooted: No

>code -> Universal  
>Loaded a multiple sequence alignment with \*\*19\*\* sequences, \*\*619\*\* codons, and \*\*1\*\* partitions from `/Users/sergei/Dropbox/Talks/VEME-2021/data/WestNileVirus\_NS3.fas`

>type -> global

>frequencies -> CF3x4

>lrt -> Yes

### Obtaining branch lengths and nucleotide substitution biases under the nucleotide GTR model

>kill-zero-lengths -> Yes

```

### Deleted 2 zero-length internal branches: `Node1, Node2`
* Log(L) = -7745.48, AIC-c = 15577.06 (43 estimated parameters)
* 1 partition. Total tree length by partition (subs/site) 0.672

### Fitting Standard MG94
* Log(L) = -6413.46, AIC-c = 12923.32 (48 estimated parameters)
* non-synonymous/synonymous rate ratio = 0.0086 (95% profile CI 0.0069– 0.0106)

### Running the likelihood ratio tests for dN/dS=1

>Testing _non-synonymous/synonymous rate ratio_ == 1

Likelihood ratio test for _non-synonymous/synonymous rate ratio == 1_, **p = 0.0000**.

### **Synonymous tree**
(HNY1999:0.001081533713434183, NY99_EQHS:0.001066739556141331, NY99_FLAMINGO:0, (((((RABENSBURG_ISOLATE:1.02609485578009,
(WNFCG:0.009747728897547157, SPU116_89:0.006173701209230282)Node11:0.4965243777785416)Node9:0.5660609164185494, KUNCG:0.085800899625552
03)Node8:0.06825862226912914, (ETHAN4766:0.02340495761620409,
(CHIN_01:0.01183597486736843, EG101:0.0147869746725866)Node17:0.007539849787066069)Node15:0.003381601852478645)Node7:0.018166959127327
35, (((ITALY_1998_EQUIINE:0.008869163186759085, PAAN001:0.007726247531366057)Node22:0.002644807163484107,
(R097_50:0.001612646775301341, VLG_4:0.001063048479450396)Node25:0.002742256244585967)Node21:0.000710558183635453, KN3829:0.0030011303
30081092)Node20:0.01077404359987302)Node6:0.009119438361664152, AST99:0.01648564314579341)Node5:0.006362164798999309, PAH001:0.00975889
8063241322)Node4:0.01061022608779064, IS_98:0.002196841638364542)Node3:0.001024862946203058, MEX03:0.003213545672362693)

### **Non-synonymous tree**
(HNY1999:2.027602010061738e-05, NY99_EQHS:1.99986670907985e-05, NY99_FLAMINGO:0, (((((RABENSBURG_ISOLATE:0.0192366818181515,
(WNFCG:0.0001827452483514863, SPU116_89:0.0001157412739507511)Node11:0.009308575534198456)Node9:0.01061220965829361, KUNCG:0.0016085497
32522618)Node8:0.001279676426151033, (ETHAN4766:0.00043878372461769,
(CHIN_01:0.0002218945755829819, EG101:0.0002772183538658889)Node17:0.0001413531024869812)Node15:6.339647694885701e-05)Node7:0.00034058
45087002921, (((ITALY_1998_EQUIINE:0.0001662743646514413, PAAN001:0.0001448475884777466)Node22:4.958344112897044e-05,
(R097_50:3.023304592824938e-05, VLG_4:1.992946874381341e-05)Node25:5.141036478624048e-05)Node21:1.332116774085881e-05, KN3829:5.6263598
75928158e-05)Node20:0.0002019860517360077)Node6:0.0001709663908120937, AST99:0.0003090640889356378)Node5:0.0001192744893160473, PAH001:
0.0001829546419425591)Node4:0.0001989148879557604, IS_98:4.11852211947375e-05)Node3:1.921358663116306e-05, MEX03:6.024584886973394e-05)

### Writing detailed analysis report to `/Users/sergei/Dropbox/Talks/VEME-2021/data/WestNileVirus_NS3.fas.FITTER.json'
```

```

355 "Standard MG94":{
356   "AIC-c":12923.31524090407,
357   "Confidence Intervals":{
358     "non-synonymous/synonymous rate ratio":{
359       "LB":0.006858130361216103,
360       "UB":0.01057945031810191
361     }
362   },
363   "Log Likelihood":-6413.456800779904,
364   "Rate Distributions":{
365     "Substitution rate from nucleotide A to nucleotide C":0
366     .2435700972629257,
367     "Substitution rate from nucleotide A to nucleotide G":1,
368     "Substitution rate from nucleotide A to nucleotide T":0
369     .3060821438629754,
370     "Substitution rate from nucleotide C to nucleotide G":0
371     .02086612098116157,
372     "Substitution rate from nucleotide C to nucleotide T":1
373     .979041445845451,
374     "Substitution rate from nucleotide G to nucleotide T":0
375     .230594853422287,
376     "non-synonymous/synonymous rate ratio":0.008593681136758865
377   },
378   "display order":1,
379   "estimated parameters":48
380 }
381 },
382 "input":{
383   "file name":"/Users/sergei/Dropbox/Talks/VEME-2021/data
384   /WestNileVirus_NS3.fas",
385   "number of sequences":19,
386   "number of sites":619,
387   "partition count":1,
388   "trees":{
389     "0":"(HNY1999,NY99_EQHS,NY99_FLAMINGO,((((((RABENSBURG_ISOLATE
390     ,(WNFCG,SPU116_89)Node11)Node9,KUNCG)Node8,(ETHAN4766,(CHIN_01
391     ,EG101)Node17)Node15)Node7,(((ITALY_1998_EQUIINE,PAAN001)Node22
392     ,(R097_50,VLG_4)Node25)Node21,KN3829)Node20)Node6,AST99)Node5
393     ,PAH001)Node4,IS_98)Node3,MEX03)"
394   }
395 },
396   "test results":{
397     "non-synonymous/synonymous rate ratio":{
398       "LRT":2512.583763184408,
399       "p-value":0
400     }
401 }

```

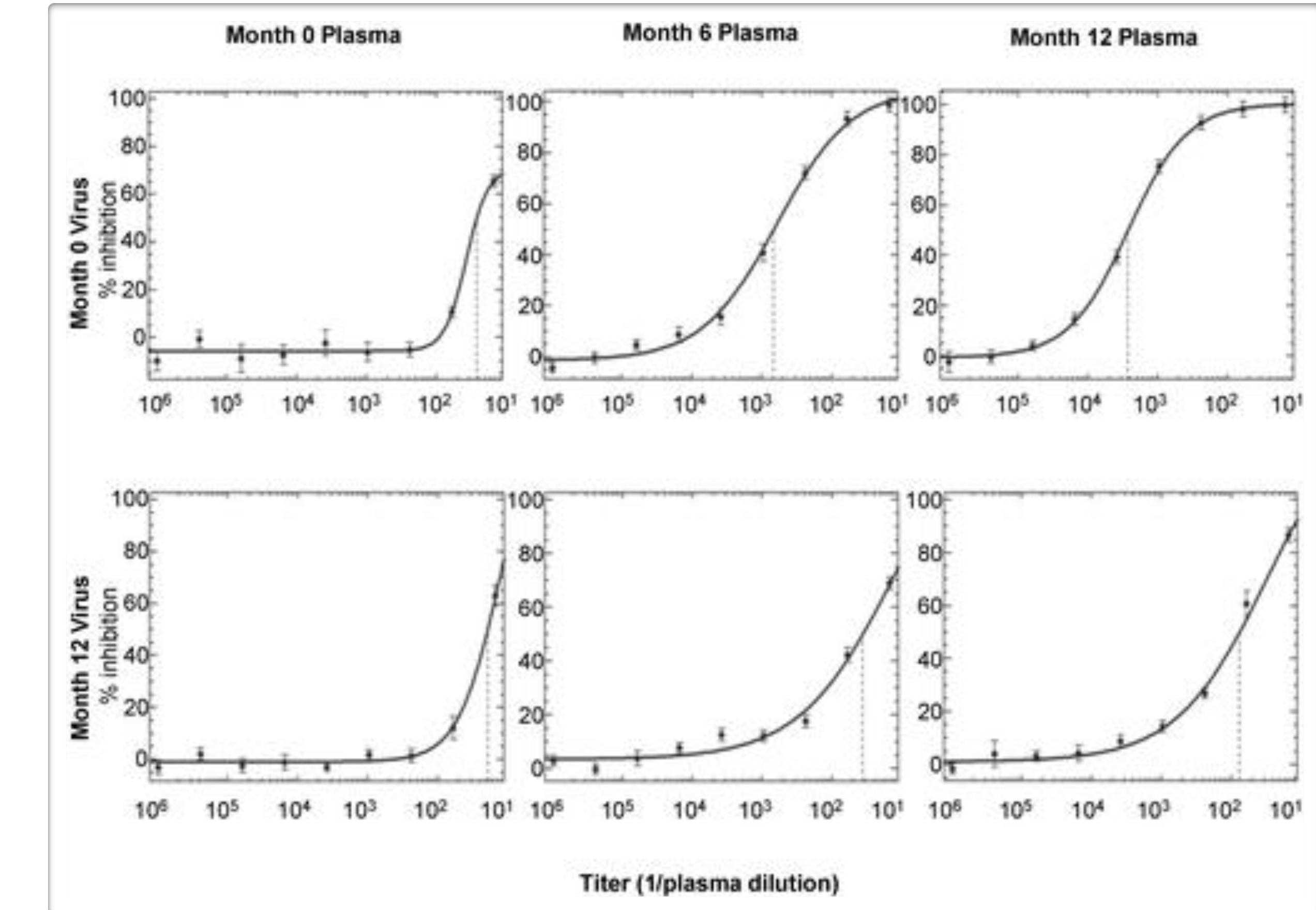
WNV NS3	Model	Log L	# p	dN/dS	LRT	p-value
	<i>Null</i>	-7668.7	49	1		
	Alternative	-6413.5	50	0.009 [0.007-0.011]	2512.6	~0
	<i>Very strongly conserved</i>					
HIV-1 env	Model	Log L	# p	dN/dS	LRT	p-value
	<i>Null</i>	-2078.3	40	1		
	Alternative	-2078.2	41	1.122 [0.94-1.33]	0.33	~0.6
	<i>Not significantly different from neutral</i>					
SARS-CoV-2 spike	Model	Log L	# p	dN/dS	LRT	p-value
	<i>Null</i>	-9311.0	176	1		
	Alternative	-9292.0	177	0.54 [0.48-0.61]	37.94	~0
	<i>Very strongly conserved</i>					

# Mean gene-wide dN/dS estimates

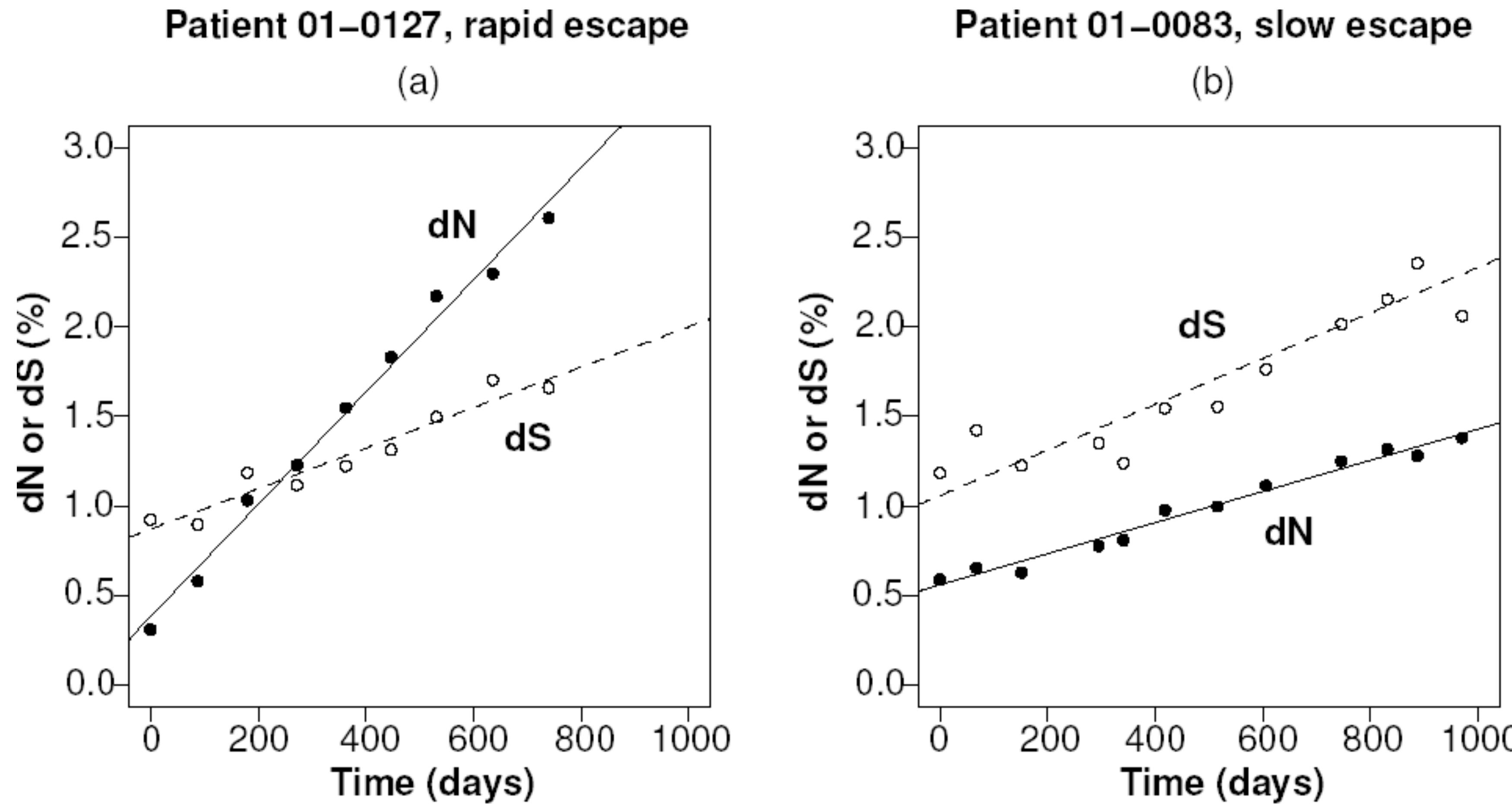
- Are not the way to go, **except** when you have very small (2-3 sequence) datasets
- For example:
  - The humoral arm of the immune system mounts a potent defense against viral infections
  - Existing successful vaccines are based on raising a neutralizing antibody (nAb) response to the pathogen
  - 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

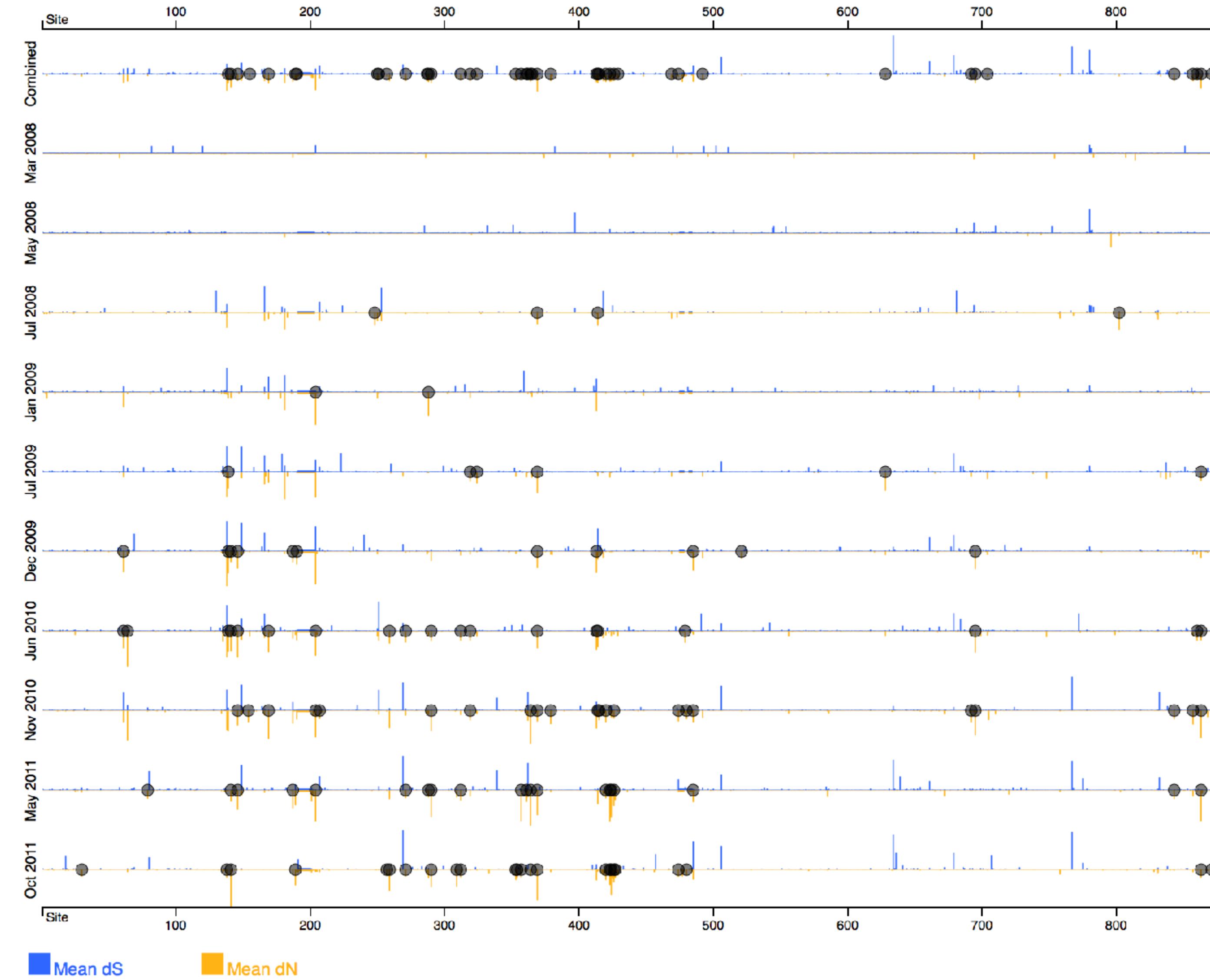
- Neutralization can be measured by the serum dilution needed to reduce viral replication by 50% (typically presented as the inverse of the titer)
- Although variable between individuals, the rate of escape from neutralizing antibodies can be very high during acute/early HIV infection
  - Sera are effective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



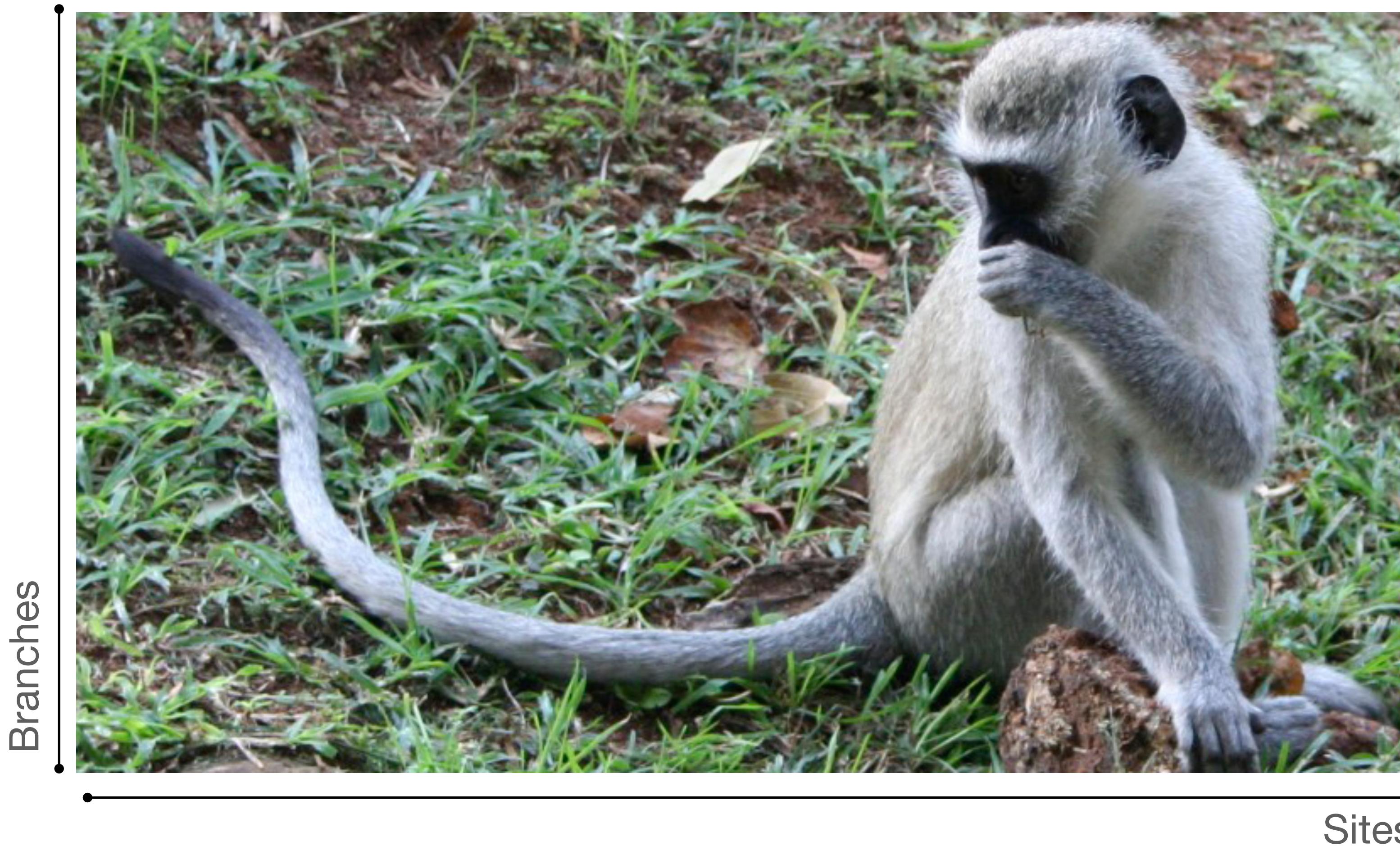
# Amino acid substitutions in HIV-1 *env* accumulate faster during rapid escape



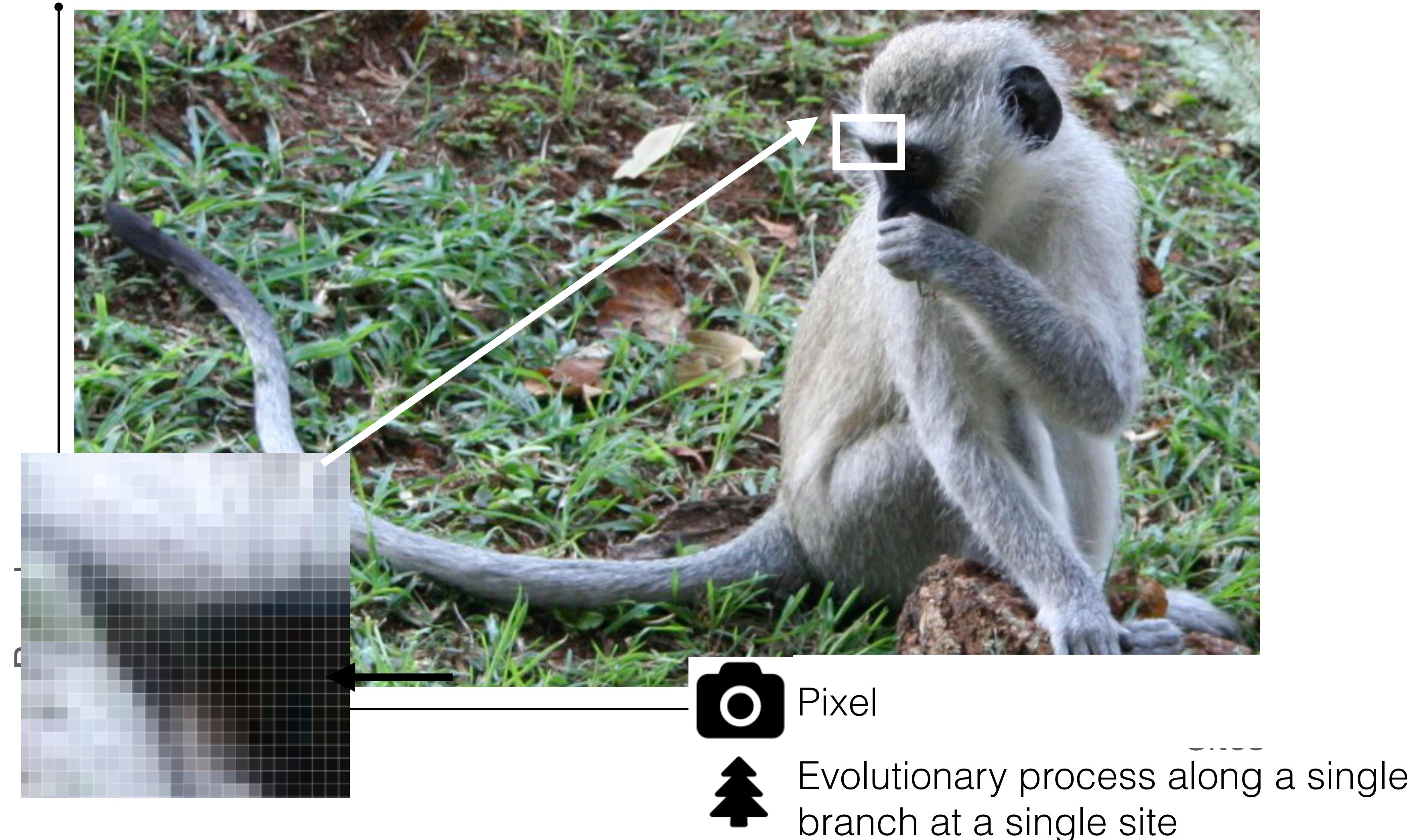
**But upon closer look, this pattern is highly variable both across a gene and through time.**



# Selection inference as image processing



# Selection inference as image processing



# Forget about the color



Intensity/brightness

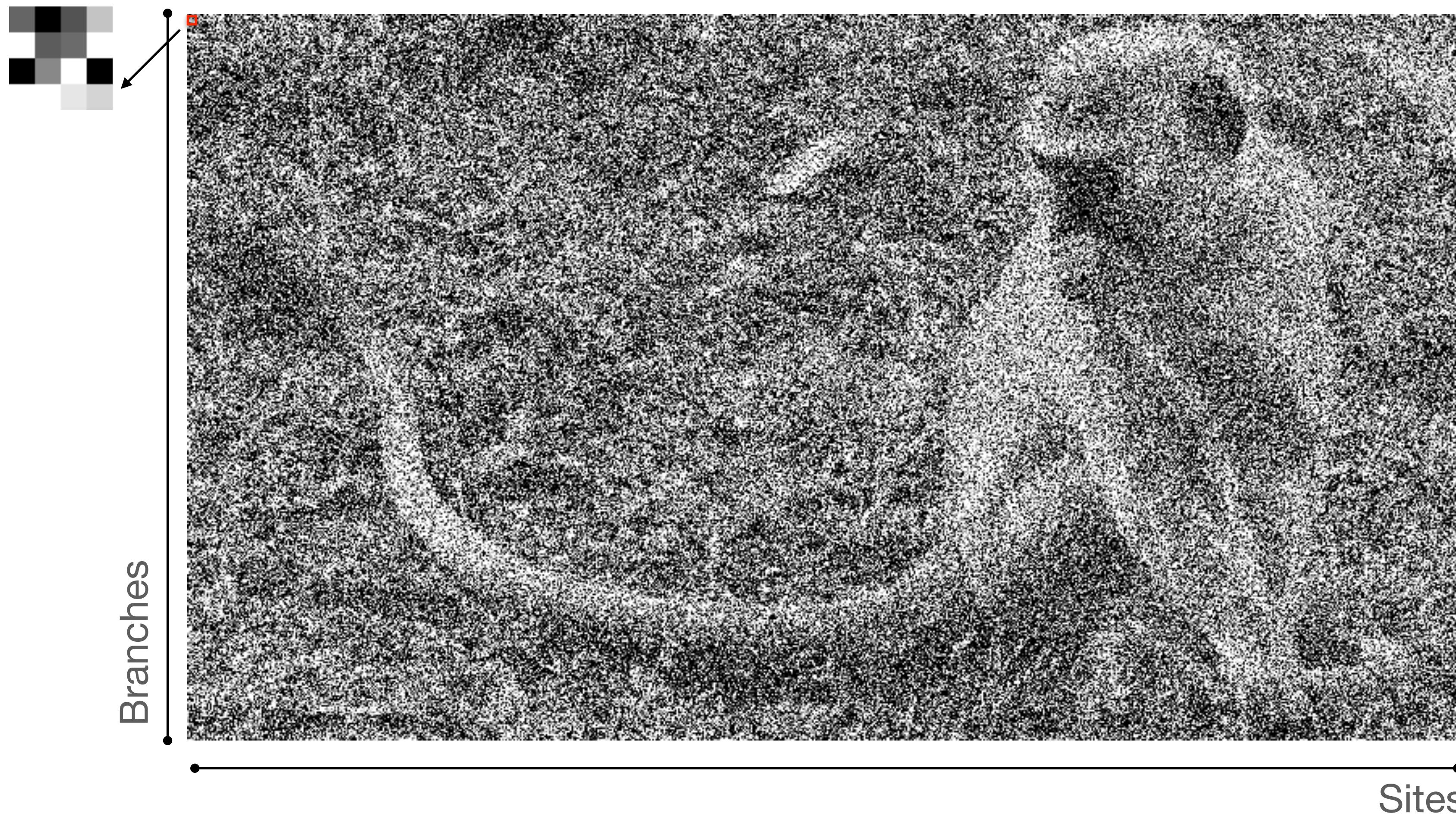
Color



Evolutionary rate ( $dN/dS$ )

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)

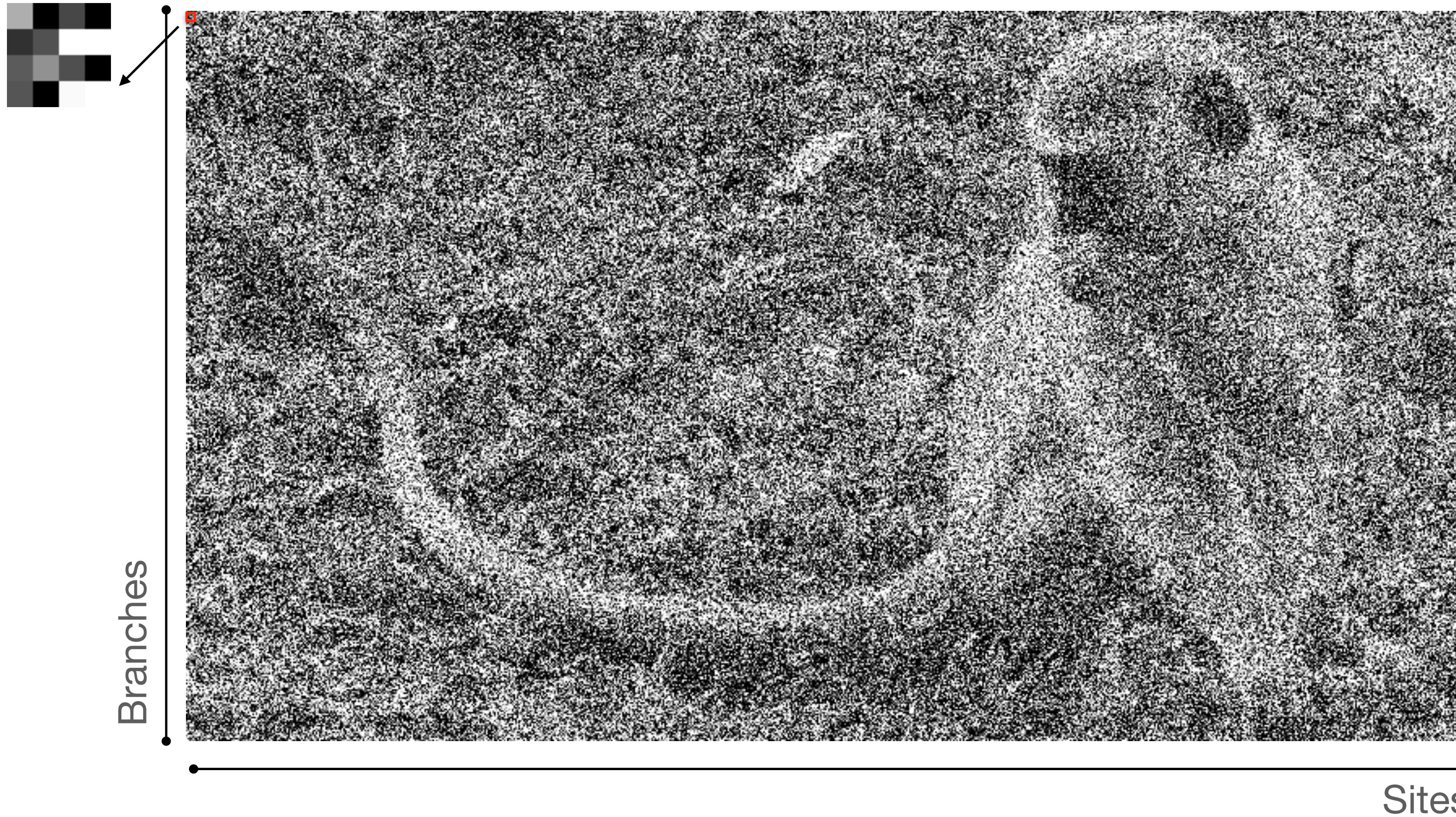


Visual noise



Saturation, missing data, model misspecification,  
sampling variation

# Evolution is largely unobserved and noisy (another replicate)



Visual noise



Saturation, missing data, model misspecification,  
sampling variation

- ⌚ High local variability
  - ⌚ 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)



Is the average color sufficiently “bright”

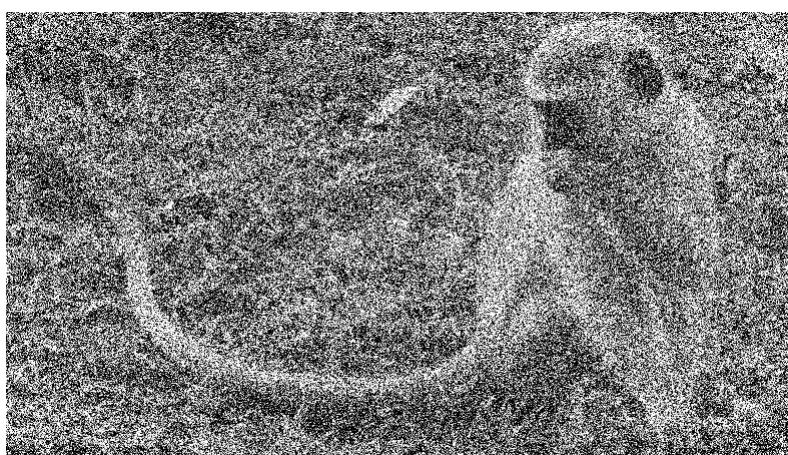
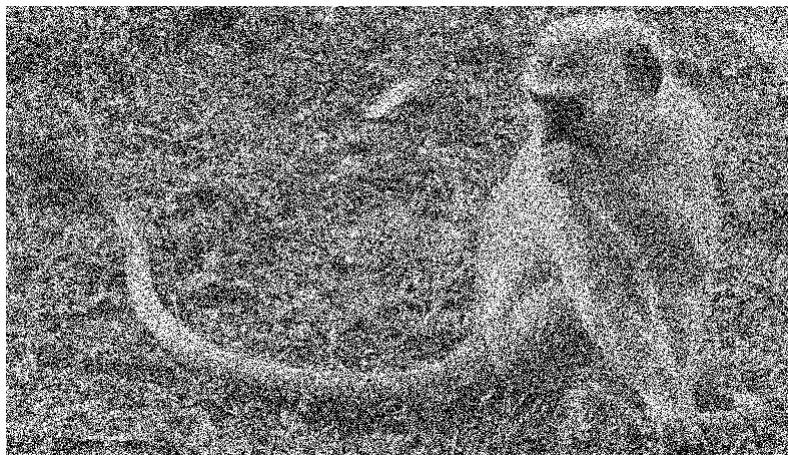


Is there evidence that **gene-wide dN/dS > 1?** Aggregate data over the entire alignment, by inferring a single dN/dS parameter from all sites and branches

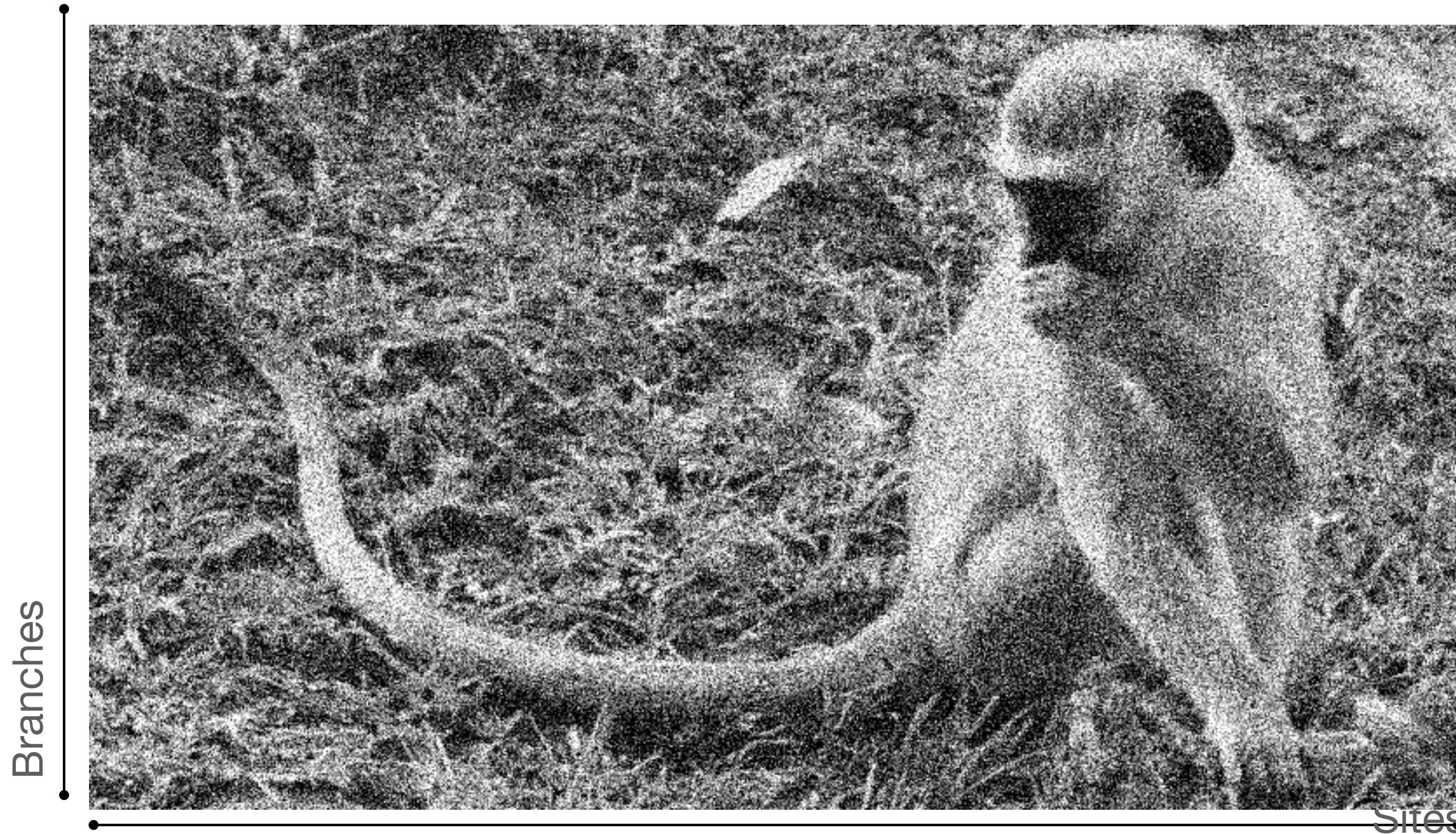
Sites



- 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]

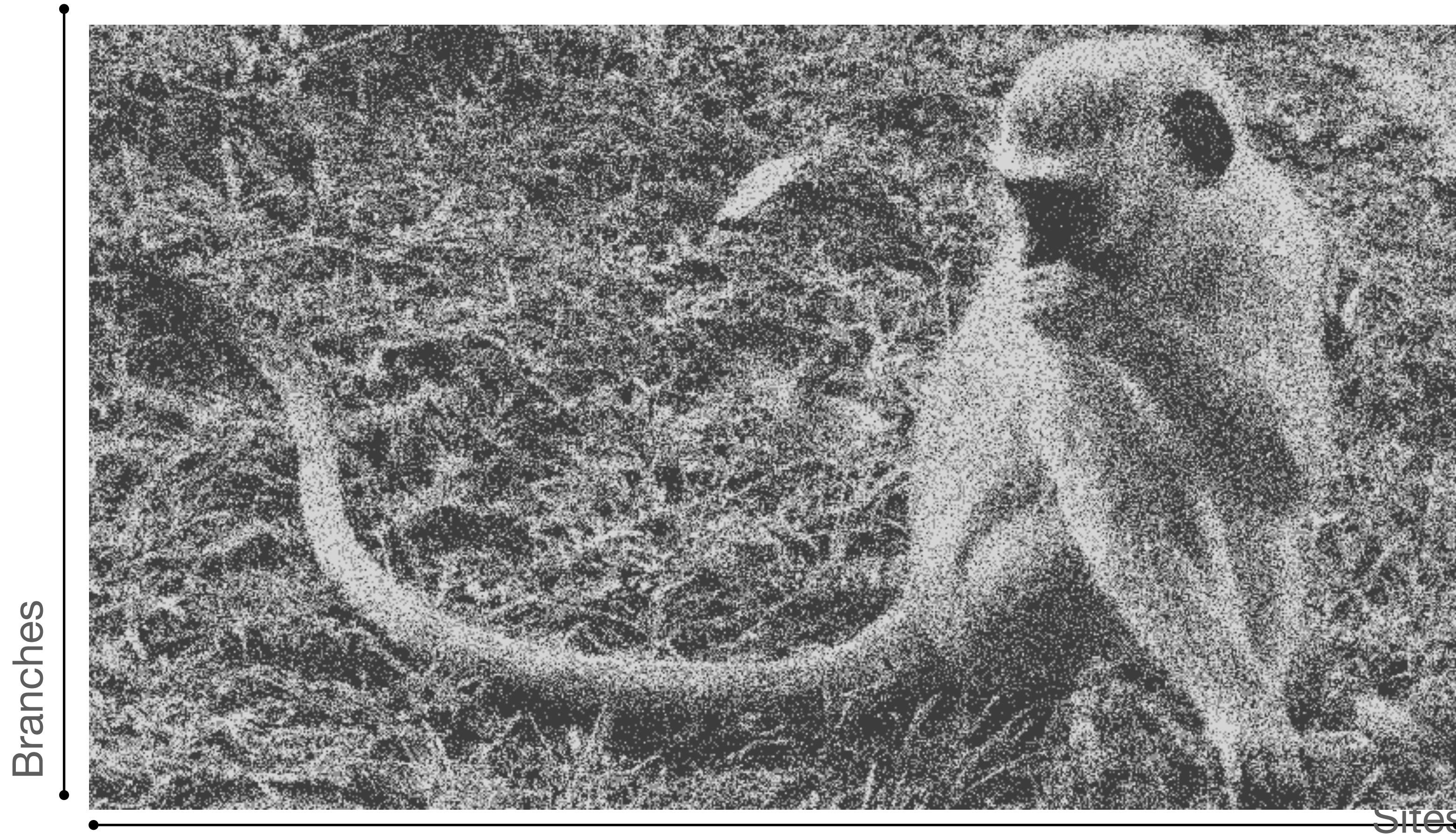


Is there enough **image area** that is sufficiently bright; allow each pixel to be one of K (=3) colors, chosen adaptively, e.g. to minimize perceptual differences



[BUSTED]: each branch-site combination is drawn from a K-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 K (=3) colors, chosen adaptively, e.g. to minimize perceptual differences



[BUSTED]: each branch-site combination is drawn from a K-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

## Model fits



Model	log L	#. params	AIC <sub>c</sub>	Branch set	ω <sub>1</sub>	ω <sub>2</sub>	ω <sub>3</sub>	
Unconstrained model	-2040.0	45	4170.9	Test	0.58 (85.37%)	0.73 (12.50%)	93.41 (2.13%)	
Constrained model	-2076.6	44	4242.1	Test	0.00 (29.28%)	1.00 (54.27%)	1.00 (16.45%)	

This table reports a statistical summary of the models fit to the data. Here, **Unconstrained model** refers to the BUSTED alternative model for selection, and **Constrained model** refers to the BUSTED null model for selection.

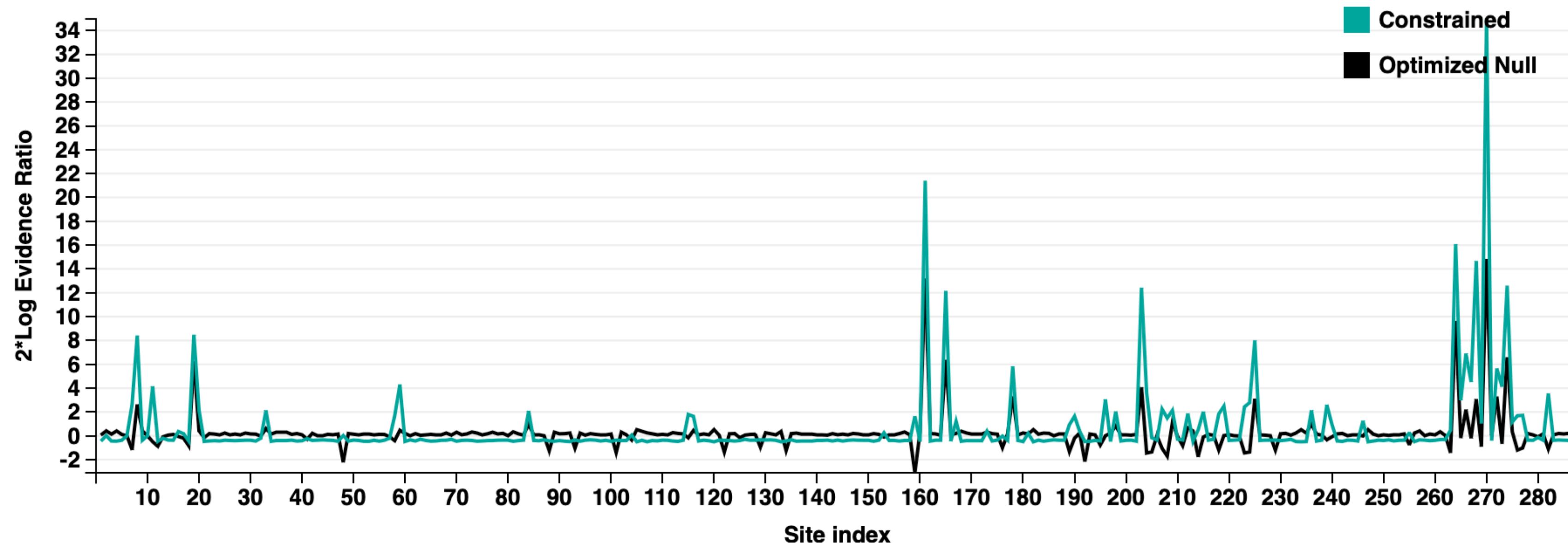
hyphy busted --srv No --alignment HIV-sets.fas --starting-points 5

Produces *HIV-sets.fas.BUSTED.json* file  
View in <http://vision.hyphy.org/BUSTED>

# 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

# Model Test Statistics Per Site


[Export Chart to PNG](#)
[Export Chart to SVG](#)


Showing entries 1 through 13 out of 13.

[Export Table to CSV](#)

« < > »

Site index	Unconstrained likelihood	Constrained likelihood	Optimized Null Likelihood	Constrained Statistic	Optimized Null Statistic
◆	◆	◆			
8	-25.12	-29.31	-26.42	8.38	2.6
19	-13.12	-17.34	-16.24	8.45	6.23

...  
...

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

BUSTED without synonymous rate variation **found no evidence** (LRT, p-value = 0.262  $\geq .05$ ) of gene-wide episodic diversifying selection in the selected test branches of your phylogeny. Therefore, there is no evidence that any sites have experienced diversifying selection along the test branch(es). 

See [here](#) for more information about this method.

Please cite [PMID 25701167](#) if you use this result in a publication, presentation, or other scientific work.

## Model fits



Model	log L	#. params	AIC <sub>c</sub>	Branch set	ω <sub>1</sub>	ω <sub>2</sub>	ω <sub>3</sub>	
Unconstrained model	-6396.1	52	12896.7	Test	0.00 (75.70%)	0.00 (23.57%)	1.89 (0.73%)	
Constrained model	-6396.7	51	12895.9	Test	0.00 (25.78%)	0.00 (73.15%)	1.00 (1.07%)	

This table reports a statistical summary of the models fit to the data. Here, **Unconstrained model** refers to the BUSTED alternative model for selection, and **Constrained model** refers to the BUSTED null model for selection.

hyphy busted --srv No --alignment WestNileVirus\_NS3.fas --starting-points 5

# Gene-wide selection analysis using a branch-site method (BUSTED), SARS-CoV-2 spike

---

BUSTED without synonymous rate variation **found evidence** (LRT, p-value = 0.019 ≤ .05) of gene-wide episodic diversifying selection in the selected test branches of your phylogeny. Therefore, there is evidence that at least one site on at least one test branch has experienced diversifying selection.



See [here](#) for more information about this method.

Please cite [PMID 25701167](#) if you use this result in a publication, presentation, or other scientific work.

---

## Model fits



Model	log L	#. params	AIC <sub>c</sub>	Branch set	$\omega_1$	$\omega_2$	$\omega_3$	
Unconstrained model	-9287.4	181	18937.2	Test	0.07 (94.22%)	1.00 (0.58%)	9.07 (5.20%)	
Constrained model	-9290.6	180	18941.7	Test	0.00 (46.91%)	1.00 (4.87%)	1.00 (48.23%)	

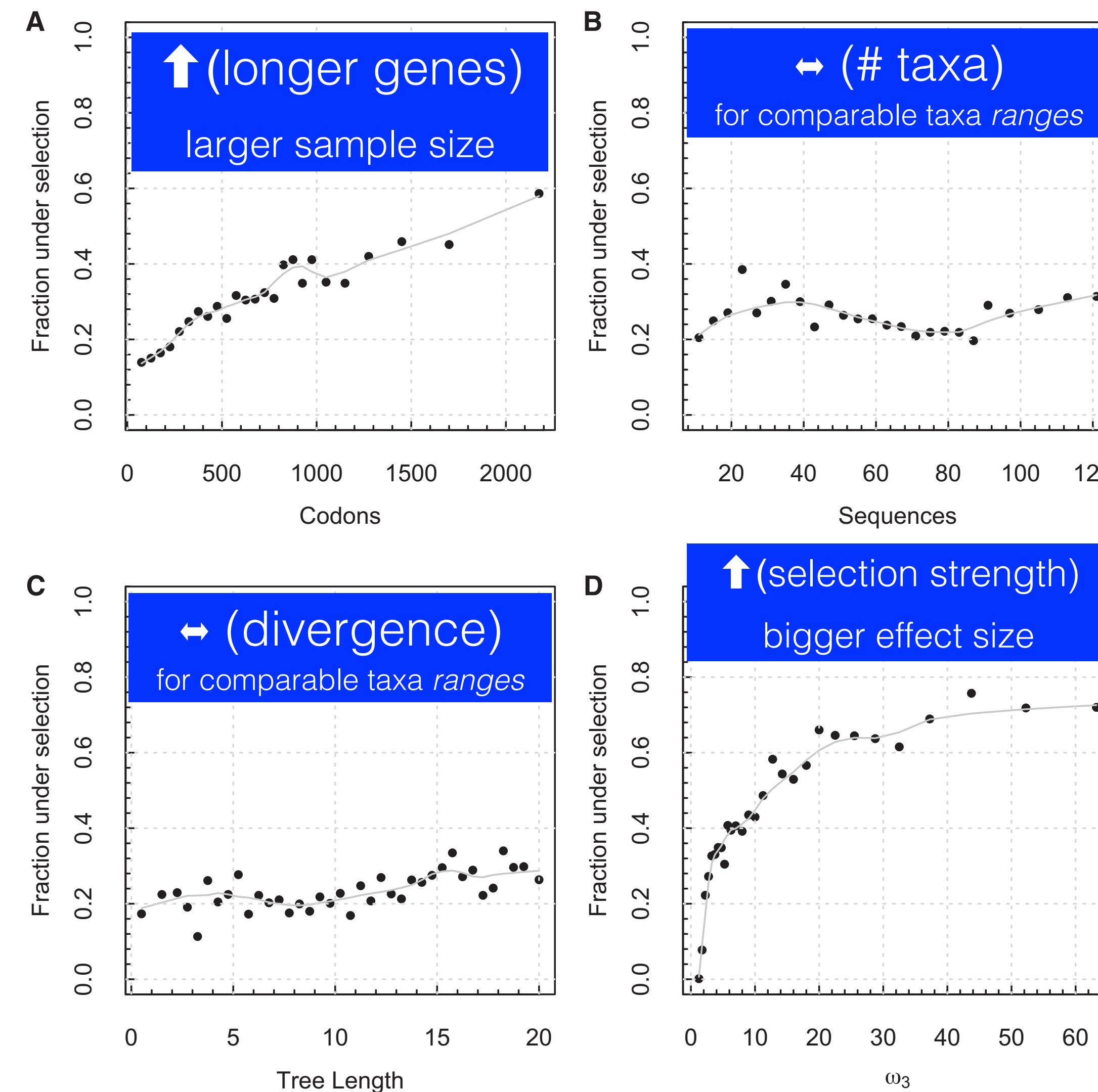
hyphy busted --srv No --alignment spike.fas --tree spike.tree --starting-points 5

# BUSTED analysis

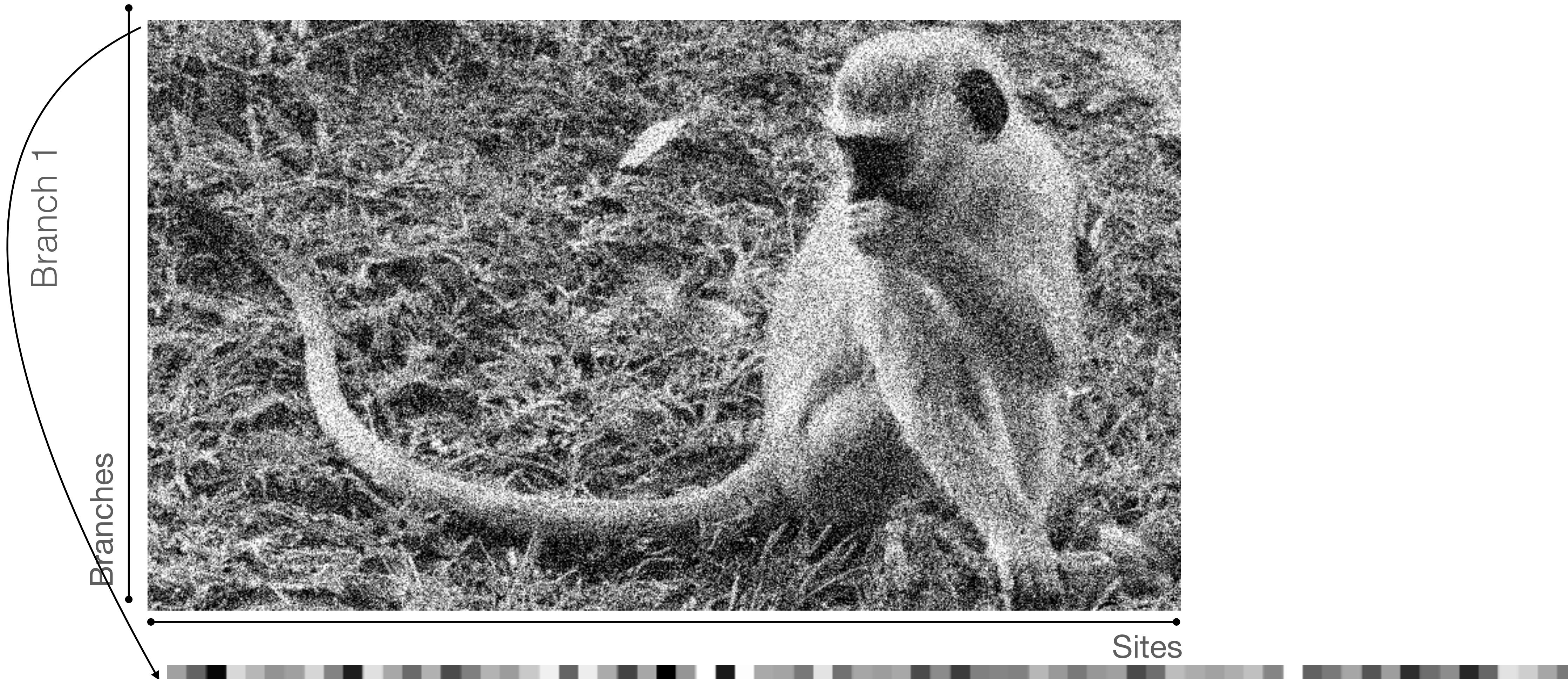
- **West Nile Virus NS3 protein**
  - No statistical support for selection; ML point estimate allocates a small proportion of sites (~1%) to the selected group ( $dN/dS \sim 2$ )
  - The rest of the gene is very strongly conserved ( $dN/dS = 0.004$ )
- **HIV-1 transmission pair**
  - Very strong evidence 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 ( $dN/dS = 0.6-0.7$ )
- **SARS-CoV-2 spike**
  - Evidence of episodic diversification ( $dN/dS \sim 9$ ) on a small proportion of sites (5.2%)
  - Most of the rest of the gene evolves with very strong 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)



# 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.

Less Is More: An Adaptive Branch-Site Random Effects Model  
for Efficient Detection of Episodic Diversifying Selection

Martin D. Smith,<sup>1</sup> Joel O. Wertheim,<sup>2</sup> Steven Weaver,<sup>2</sup> Ben Murrell,<sup>2</sup> Konrad Scheffler,<sup>2,3</sup> and  
Sergei L. Kosakovsky Pond<sup>\*2</sup>

*Mol. Biol. Evol.* 32(5):1342–1353

- 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)

# Less Is More: An Adaptive Branch-Site Random Effects Model for Efficient Detection of Episodic Diversifying Selection

Martin D. Smith,<sup>1</sup> Joel O. Wertheim,<sup>2</sup> Steven Weaver,<sup>2</sup> Ben Murrell,<sup>2</sup> Konrad Scheffler,<sup>2,3</sup> and Sergei L. Kosakovsky Pond<sup>\*2</sup>

*Mol. Biol. Evol.* 32(5):1342–1353

- Fix the tree; estimate and fix some of the nuisance model parameters that are shared by all branches (GTR biases, frequency counts)
- Fit a simple baseline model (one  $\omega$  per branch); use this model to get initial guesses for all other parameters
- Perform a greedy step-up procedure (like forward variable selection in regression models, but not as statistically bad)
- For each branch (longest first) try two  $\omega$  rate classes, then three  $\omega$  rate classes etc, until no more goodness-of-fit improvement (AIC-c)
  - Fix the number of rates and move on to the next longest branch
  - Perform selection testing on the overall model (different number of  $\omega$  classes on branches), using the likelihood ratio test
  - Each branch specified a priori (could be all branches)
  - Appropriate multiple testing correction

# adaptive Branch Site REL results summary

INPUT DATA | HIV-sets.fas | 16 sequences | 288 sites

 Export ▾

aBSREL found evidence of episodic diversifying selection on 3 out of 26 branches in your phylogeny. 

A total of 26 branches were formally tested for diversifying selection. Significance was assessed using the Likelihood Ratio Test at a threshold of  $p \leq 0.05$ , after correcting for multiple testing. Significance and number of rate categories inferred at each branch are provided in the [detailed results](#) table.

See [here](#) for more information about this method.

Please cite [PMID 25697341](#) if you use this result in a publication, presentation, or other scientific work.

## Tree summary

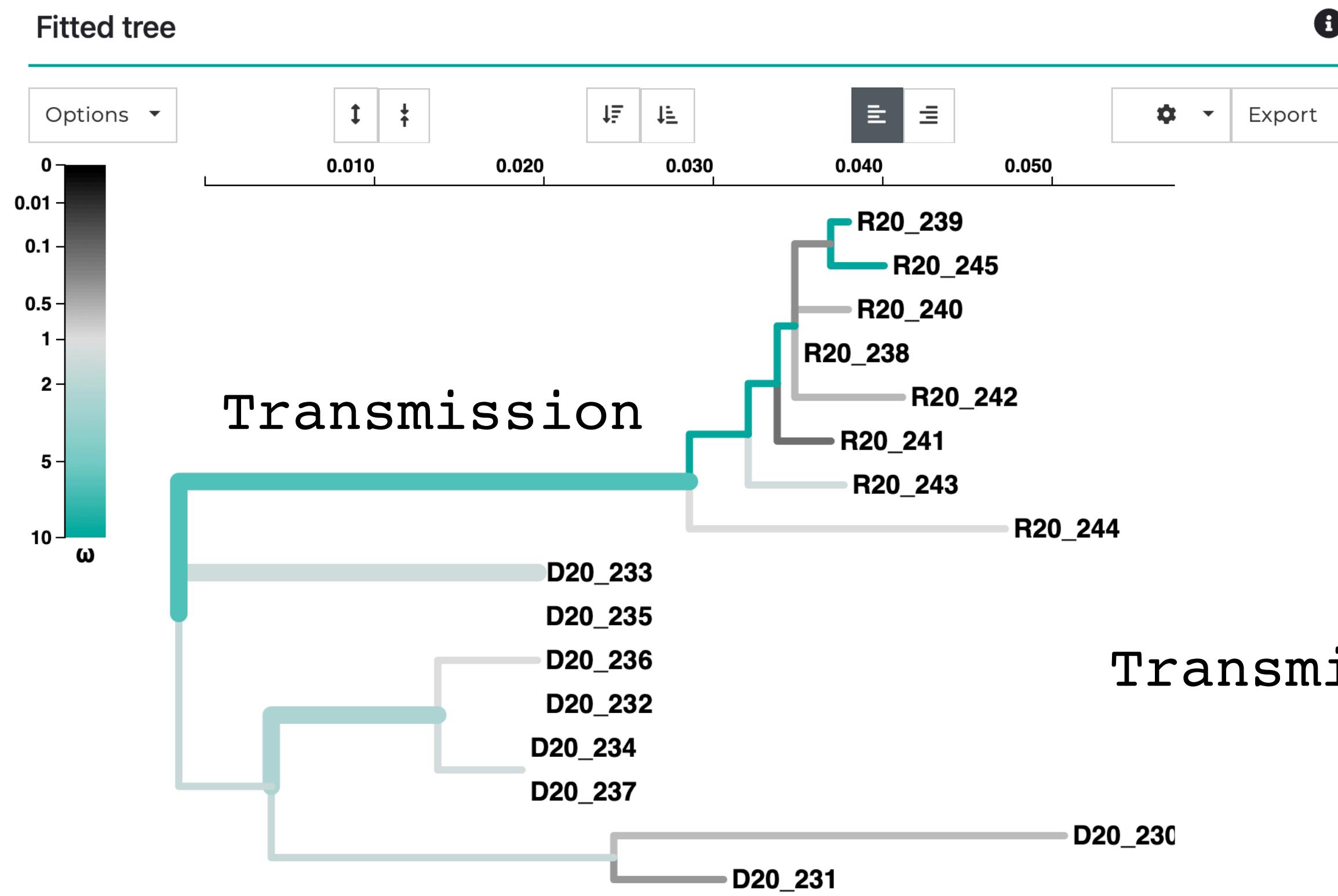
ω rate classes	# of branches	% of branches	% of tree length	# under selection
1	21	81%	0.49%	0
2	5	19%	100%	3

This table contains a summary of the inferred aBSREL model complexity. Each row provides information about the branches that were best described by the given number of ω rate categories.

hyphy absrel --alignment HIV-sets.fas

# HIV-1 env

Fitted tree



One dN/dS per branch

Fitted tree



Adaptive dN/dS per branch

# adaptive Branch Site REL results summary

INPUT DATA | WestNileVirus\_NS3.fas | 19 sequences | 619 sites

 Export ▾

aBSREL found no evidence of episodic diversifying selection in your phylogeny.



A total of 33 branches were formally tested for diversifying selection. Significance was assessed using the Likelihood Ratio Test at a threshold of  $p \leq 0.05$ , after correcting for multiple testing. Significance and number of rate categories inferred at each branch are provided in the [detailed results table](#).

See [here](#) for more information about this method.

Please cite [PMID 25697341](#) if you use this result in a publication, presentation, or other scientific work.

## Tree summary

ω rate classes	# of branches	% of branches	% of tree length	# under selection
1	30	91%	37%	0
2	3	9.1%	63%	0

This table contains a summary of the inferred aBSREL model complexity. Each row provides information about the branches that were best described by the given number of ω rate categories.

hyphy absrel --alignment WestNileVirus\_NS3.fas

# adaptive Branch Site REL results summary

INPUT DATA | spike.fas | 118 sequences | 1273 sites

 Export ▾

aBSREL found no evidence of episodic diversifying selection in your phylogeny. 

A total of **44** branches were formally tested for diversifying selection. Significance was assessed using the Likelihood Ratio Test at a threshold of  $p \leq 0.05$ , after correcting for multiple testing. Significance and number of rate categories inferred at each branch are provided in the [detailed results](#) table.

See [here](#) for more information about this method.

Please cite [PMID 25697341](#) if you use this result in a publication, presentation, or other scientific work.

## Tree summary

ω rate classes	# of branches	% of branches	% of tree length	# under selection
1	161	99%	66%	0
2	1	0.62%	34%	0

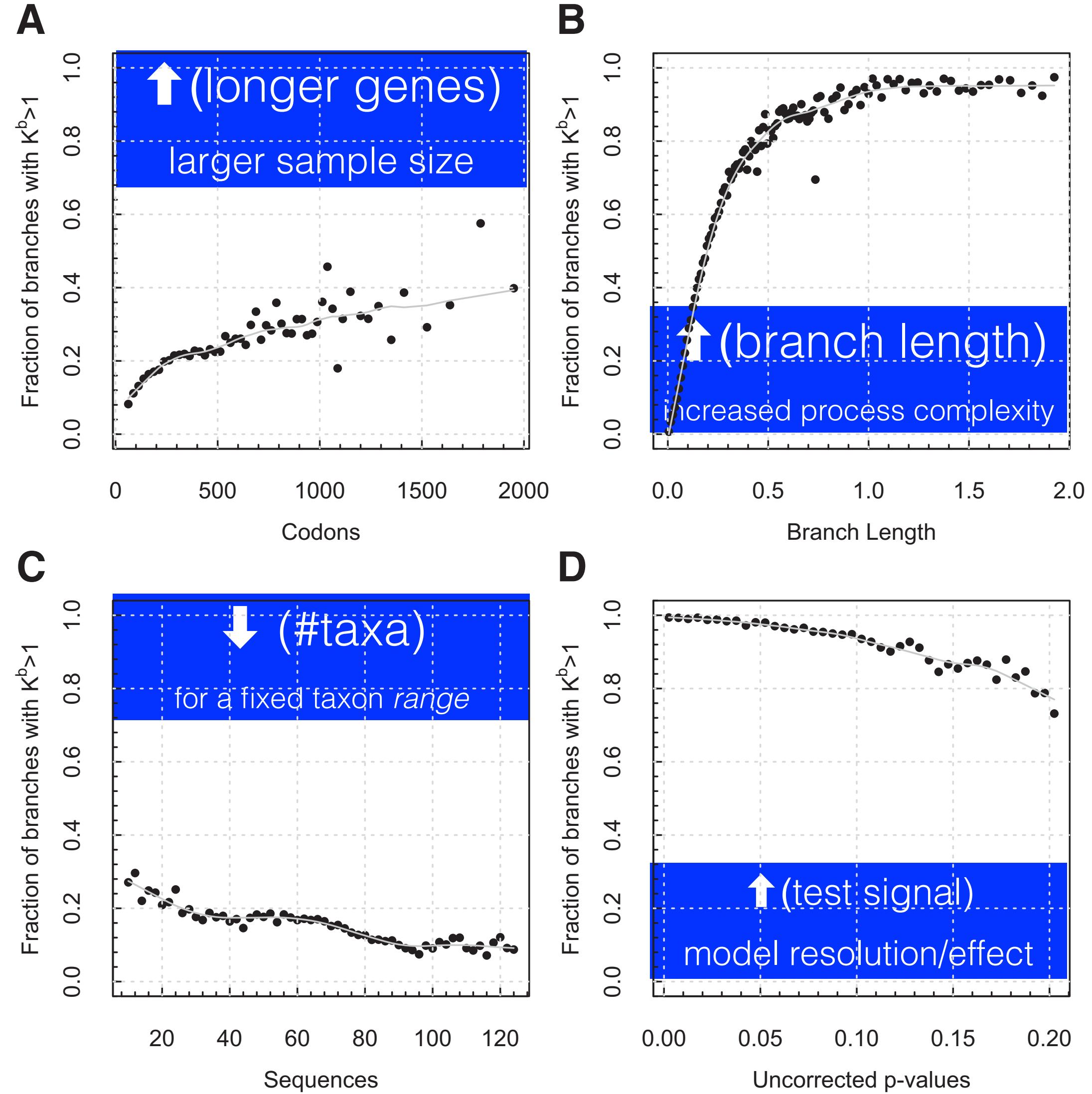
hyphy absrel --alignment spike.fas --tree spike.tree --branches Internal

# 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**
  - 76% branches can be explained with simple (single  $dN/dS$ ) models
  - 5 branches (24%, ~100% 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
- **SARS-CoV-2 spike**
  - All but one branches can be explained with simple single  $dN/dS$  models
  - 1 long terminal branches (~34% of tree length) has evidence of multiple  $dN/dS$  rate classes over sites
  - No evidence of branch level selection on internal branches.

# 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)

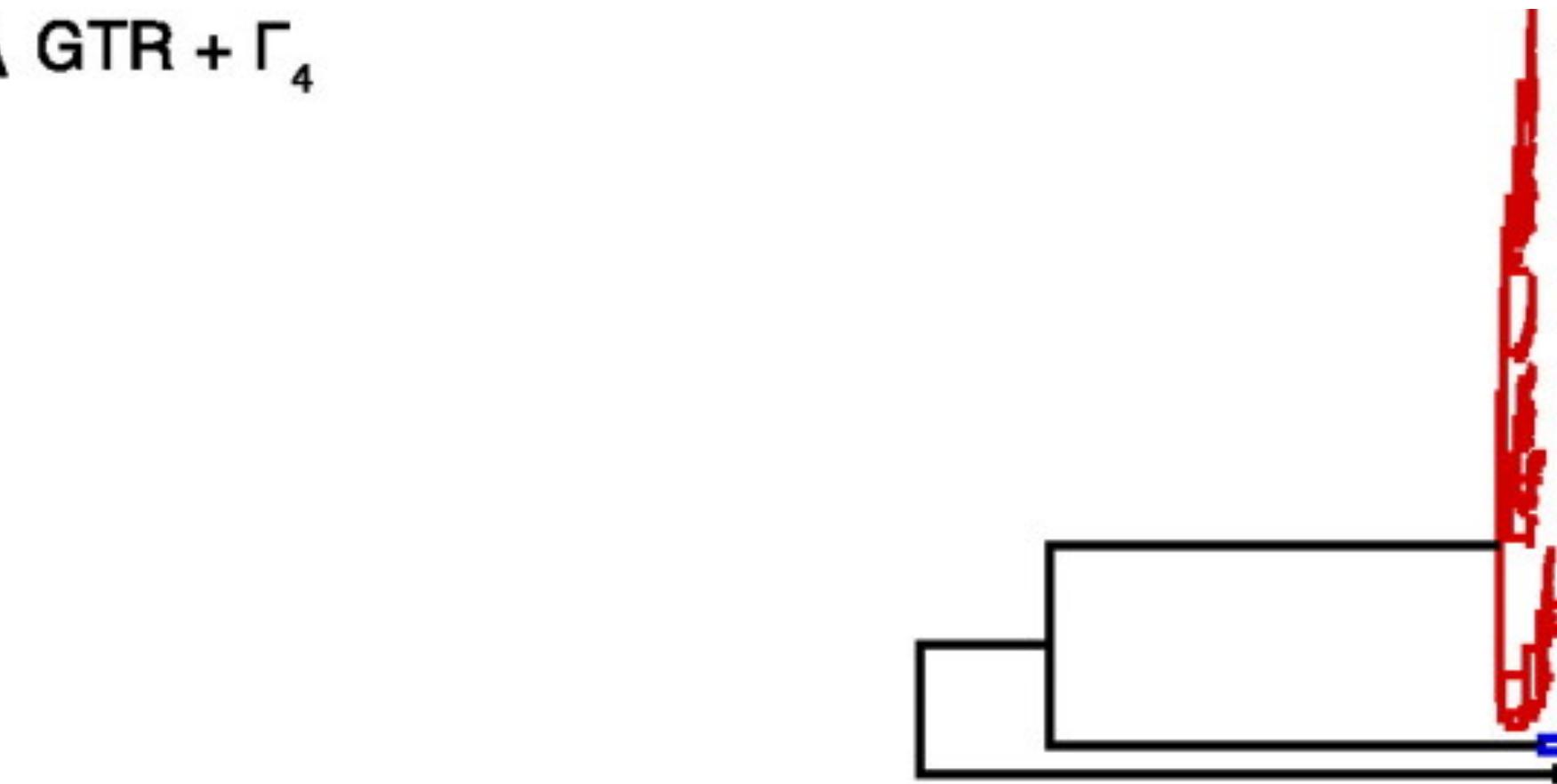
- Using models that do not vary selection pressure across lineages yields a patently false “*too young*” estimate for the origin of **measles** (about 600 years ago)
- This estimate is refuted by clear historical records which suggest

that measles is at least 1,500-5,000 years old

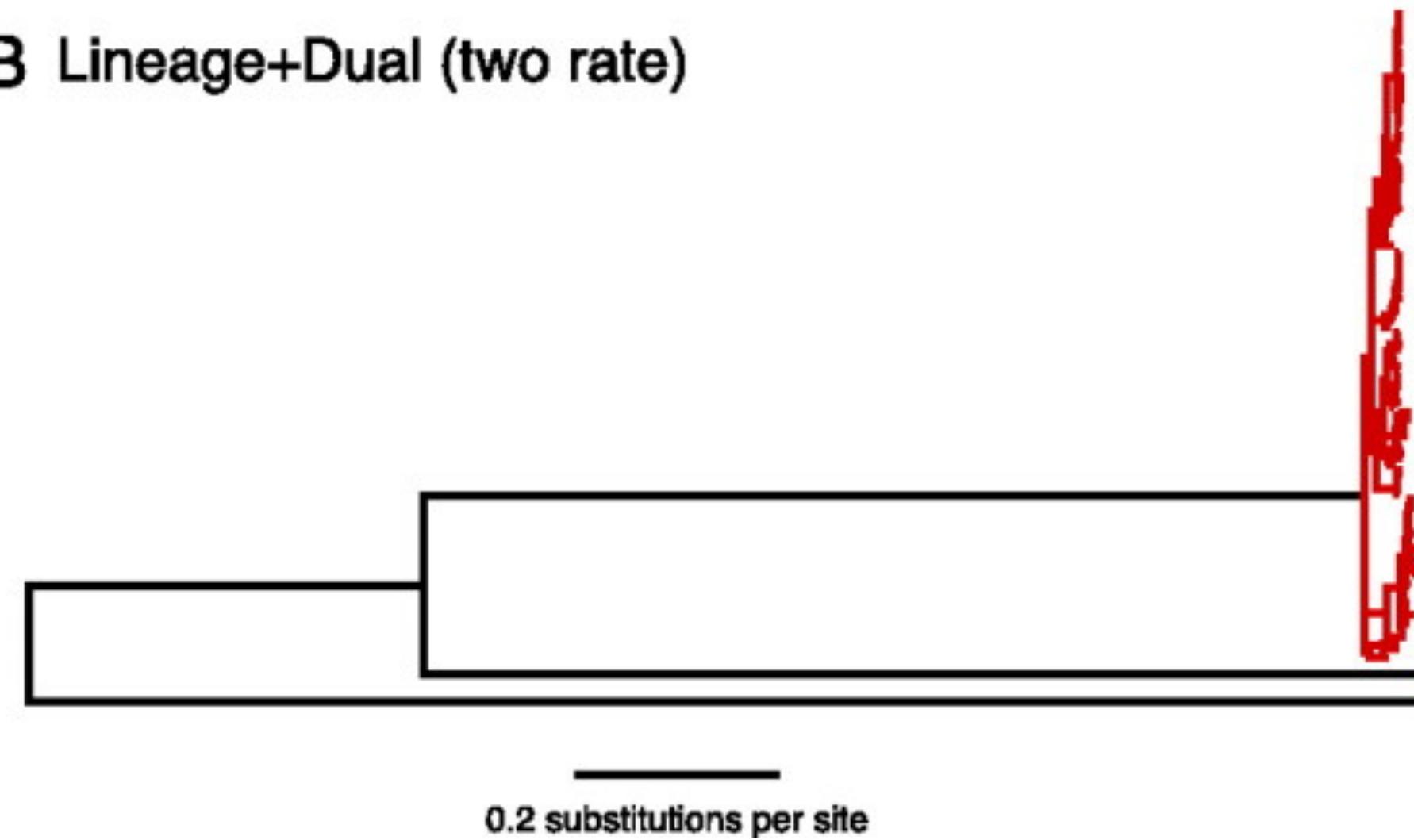
- *This includes a treatise by a Persian physician Rhazes about differential diagnosis of measles and smallpox published circa 600 AD.*

- Same patterns found for corona-viruses, ebola, avian influenza and herpesvirus

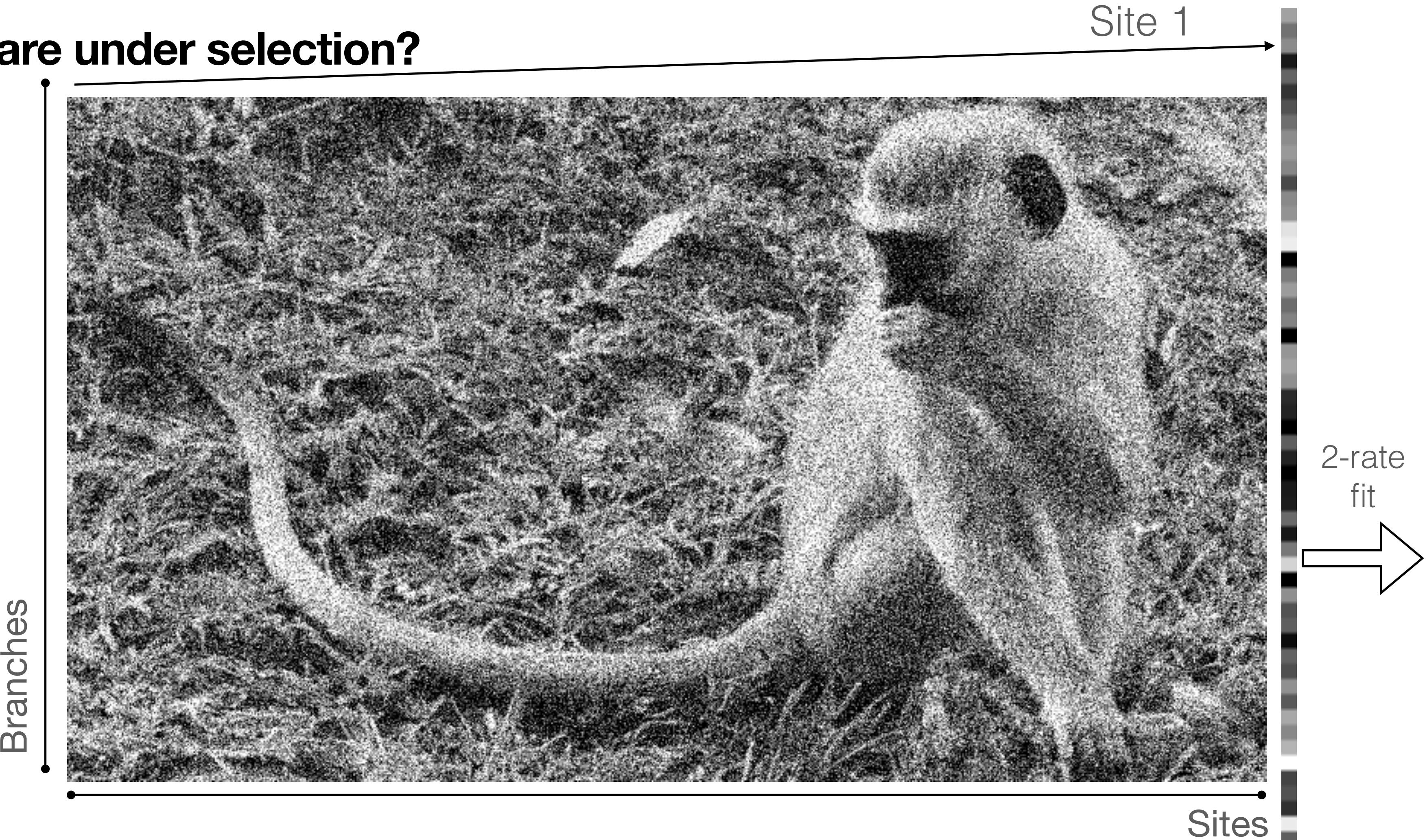
A GTR +  $\Gamma_4$



B Lineage+Dual (two rate)



# Which sites are under selection?

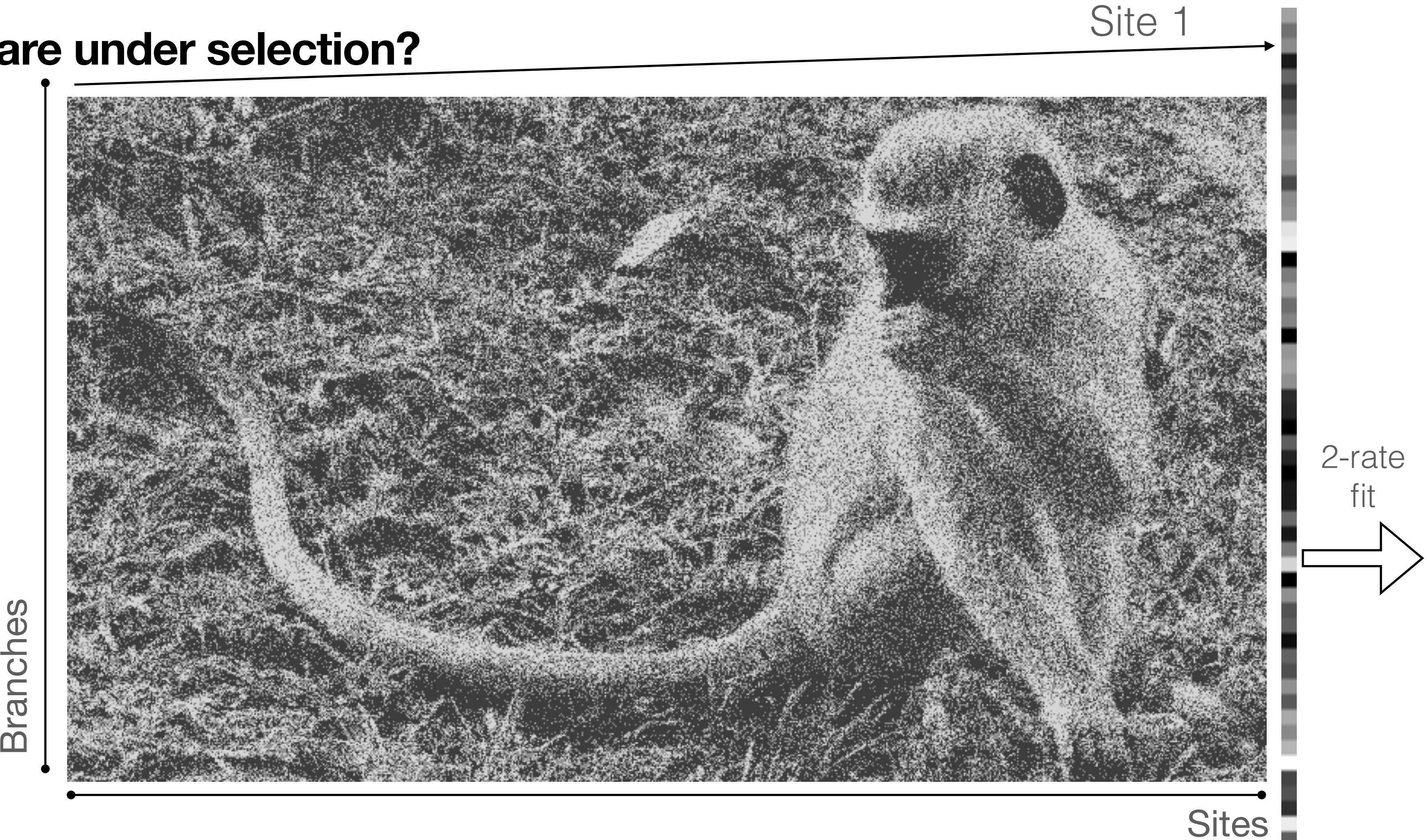


For each image column, is there a significant proportion of bright pixels, once the column has been reduced to 2 colors only?



[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 with any precision
- Does not scale especially well with the number of sequences

MEME found evidence of

+ episodic positive/diversifying selection at 9 sites

with p-value threshold of 0.1 .See [here](#) for more information about the MEME method.Please cite [PMID 22807683](#) if you use this result in a publication, presentation, or other scientific work.

## MEME Table



Sites that yielded a statistically significant result are highlighted in green.

Showing entries 1 through 20 out of 288.

[Export Table to CSV](#)   

Site	Partition	$\alpha$	$\beta^-$	$p^-$	$\beta^+$	$p^+$	LRT	p-value	# branches under selection	Total branch length	MEME LogL	FEL LogL
161	1	0.00	0.00	0.82	114.39	0.18	7.58	0.01	0.00	0.00	-16.57	-14.14
19	1	0.00	0.00	0.95	2106.40	0.05	6.61	0.02	0.00	0.00	-14.85	-11.57
274	1	2.77	2.77	0.95	10000.00	0.05	4.98	0.04	1.00	0.00	-20.16	-17.67
165	1	0.00	0.00	0.78	52.90	0.22	4.25	0.06	0.00	0.00	-15.51	-14.44
225	1	0.00	0.00	0.74	43.93	0.26	3.71	0.07	2.00	0.00	-13.87	-13.06
264	1	0.00	0.00	0.90	176.76	0.10	3.64	0.08	0.00	0.00	-11.73	-10.17
282	1	0.00	0.00	0.00	8.19	1.00	3.65	0.08	0.00	0.00	-19.33	-19.32
272	1	0.00	0.00	0.85	38.49	0.15	3.33	0.09	1.00	0.00	-10.45	-9.37

## Fitted tree

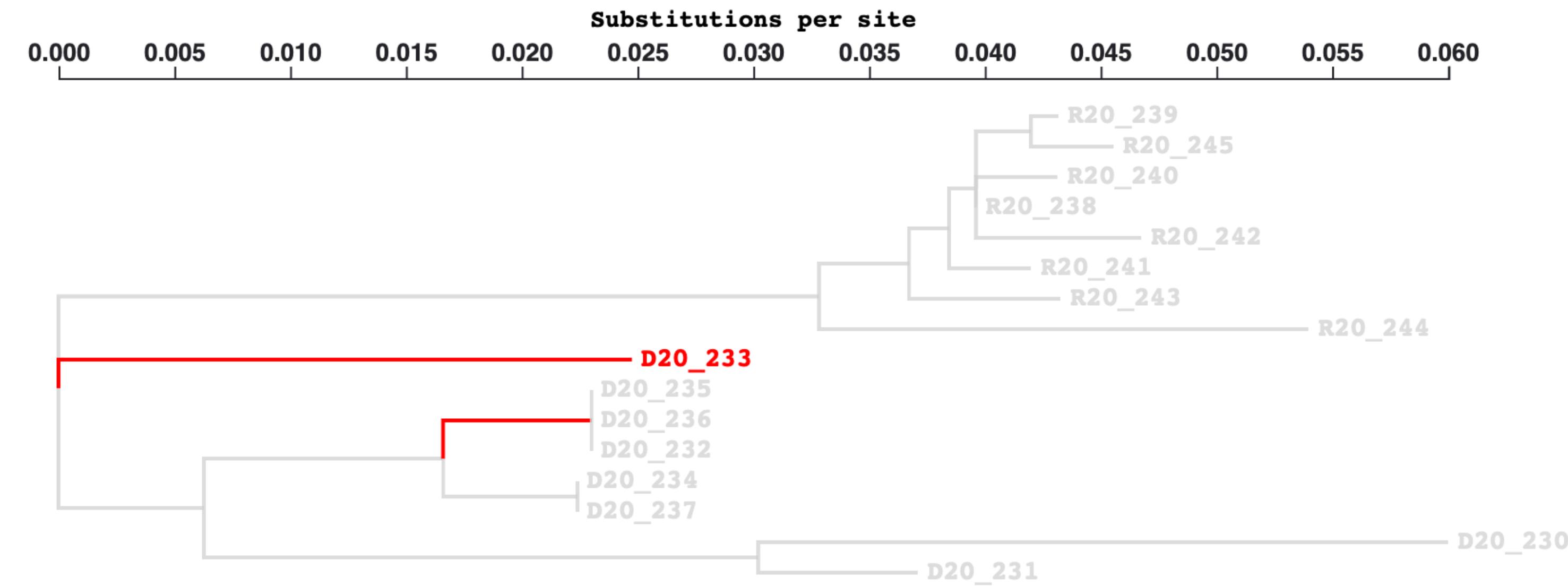


Options ▾   

Site to display: 225

EBF threshold:

Export ▾



- Where in the tree is there evidence for selection?
- Not a “strict” statistical test!
- Can use exploratory Empirical Bayes Factor analysis to find “hotspots”

# Mixed Effects Model of Evolution results summary

INPUT DATA | WestNileVirus\_NS3.fas | 19 sequences | 619 sites Export

**MEME found evidence of**

- + episodic positive/diversifying selection at 4 sites

with p-value threshold of 0.1.

See [here](#) for more information about the MEME method.

Please cite [PMID 22807683](#) if you use this result in a publication, presentation, or other scientific work.

---

**MEME Table** i

Sites that yielded a statistically significant result are highlighted in green.

Showing entries 1 through 20 out of 619. Export Table to CSV

Site	Partition	$\alpha$	$\beta^-$	$p^-$	$\beta^+$	$p^+$	LRT	$p^-$ value	# branches under selection	Total branch length	MEME LogL	FEL LogL
249	1	0.00	0.00	0.00	2.71	1.00	7.88	0.01	0.00	0.00	-34.22	-34.22
557	1	0.27	0.00	0.97	149.17	0.03	5.52	0.03	1.00	0.00	-17.69	-14.17
87	1	1.79	0.00	0.95	31.42	0.05	3.47	0.08	1.00	0.00	-23.52	-16.73
521	1	0.96	0.00	0.96	97.61	0.04	3.60	0.08	1.00	0.00	-17.26	-14.31



hyphy meme --alignment WestNileVirus\_NS3.fas

# Mixed Effects Model of Evolution results summary

INPUT DATA | spike.fas | 118 sequences | 1273 sites

 Export ▾

MEME found evidence of

+ episodic positive/diversifying selection at 6 sites

with p-value threshold of 0.1 .

See [here](#) for more information about the MEME method.

Please cite [PMID 22807683](#) if you use this result in a publication, presentation, or other scientific work.

## MEME Table



Sites that yielded a statistically significant result are highlighted in green.

Showing entries 1 through 20 out of 1273.

 Export Table to CSV

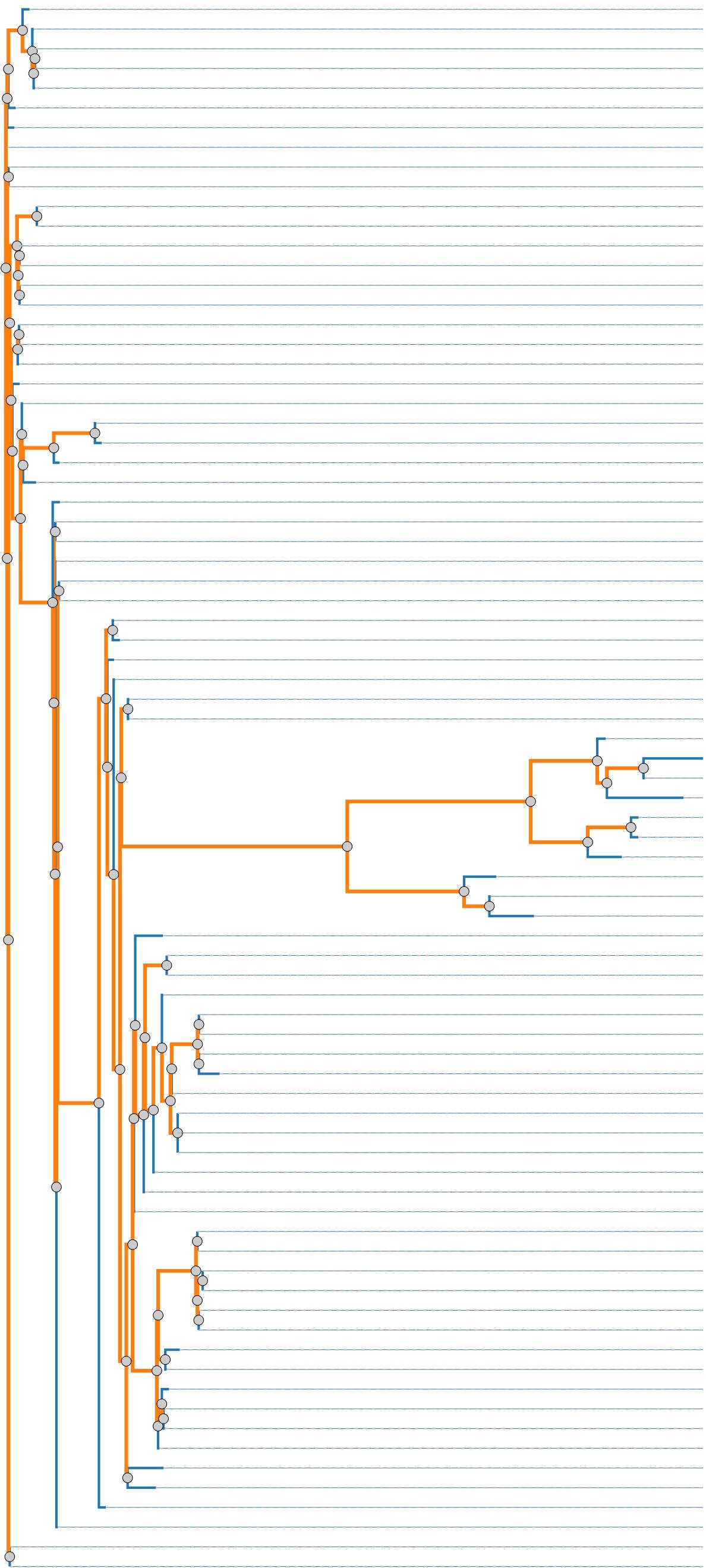


Site	Partition	$\alpha$	$\beta^-$	$p^-$	$\beta^+$	$p^+$	LRT	p-value	# branches under selection	Total branch length	MEME LogL	FEL LogL
1	1	1.22	0.55	0.99	2557.85	0.01	9.76	0.00	1.00	0.00	-19.14	-14.26
470	1	4.21	1.65	0.99	10000.00	0.01	7.46	0.01	1.00	0.00	-24.32	-20.11
1243	1	0.00	0.00	0.98	608.47	0.02	8.85	0.01	2.00	0.00	-23.28	-19.43
452	1	0.00	0.00	0.01	14.48	0.99	8.77	0.01	5.00	0.00	-36.49	-36.49
501	1	0.00	0.00	0.92	328.53	0.08	5.02	0.04	5.00	0.00	-37.07	-36.20
157	1	0.00	0.00	0.01	7.32	0.99	3.71	0.07	4.00	0.00	-29.55	-29.55

hyphy meme --alignment spike.fas --tree spike.tree

# 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

SARS CoV-2 Spike  
Internal Branches Only

Codon	Partition	alpha	beta+	p+	LRT	Episodic selection detected?	# branches	Most common codon substitutions at this site
367	1	0.000	97.046	0.481	9.056	Yes, p = 0.0047	2	[ 2 ]GTC>TTC
439	1	0.000	35.187	1.000	4.987	Yes, p = 0.0380	1	[ 1 ]AAC>AAA
452	1	0.000	30.833	1.000	5.519	Yes, p = 0.0289	1	[ 4 ]CTG>CGG   [ 1 ]CTG>ATG
477	1	0.000	51.490	0.460	4.330	Yes, p = 0.0533	1	[ 1 ]AGC>AAC
501	1	0.000	271.405	0.145	3.456	Yes, p = 0.0840	1	[ 3 ]AAT>TAT   [ 1 ]AAT>ACT, TAT>AAT

hyphy meme --alignment spike.fas --tree spike.tree --branches Internal

# MEME results

- **West Nile Virus NS3 protein**
  - **Four** sites (incl. 249, **previously reported**) with significant evidence of **episodic** (or pervasive) diversifying selection.
- **HIV-1 transmission pair**
  - **Nine** sites with significant evidence of **episodic** (or pervasive) diversifying selection. HIV-1 transmission pair
- SARS-CoV-2 spike (all)
  - **Six** sites with significant evidence of **episodic** (or pervasive) diversifying selection.
- SARS-CoV-2 spike (internal)
  - **Five** sites with significant evidence of **episodic** (or pervasive) diversifying selection.

# More on site-level selection

- Three more methods in HyPhy
- Fixed Effects Likelihood (**FEL**)
  - A simpler alternative to MEME (looks for pervasive selection)
  - May be more suited for smaller datasets or datasets of low divergence
  - Single Likelihood Ancestor Counting (**SLAC**)
- A counting-based approach
  - Good for data exploration and visualization
- Fast Unrestricted Bayesian AppRoximation (**FUBAR**)
  - A novel statistical approach for detecting pervasive adaptive evolution on large datasets (scales to 10000s of sequences)

# FEL on internal branches of Spike finds most selected sites, including many known to be of functional significance

Codon	Partition	alpha	beta	LRT	Selection detected?
5	1	0.000	19.127	2.890	Pos. p = 0.0891
12	1	0.000	20.331	2.989	Pos. p = 0.0838
18	1	0.000	19.110	2.885	Pos. p = 0.0894
138	1	0.000	26.771	2.736	Pos. p = 0.0981
367	1	0.000	44.309	9.045	Pos. p = 0.0026
439	1	0.000	34.548	4.987	Pos. p = 0.0255
452	1	0.000	30.618	5.518	Pos. p = 0.0188
477	1	0.000	23.671	4.325	Pos. p = 0.0376
501	1	0.000	38.285	3.317	Pos. p = 0.0686
570	1	0.000	21.073	3.047	Pos. p = 0.0809
614	1	0.000	22.073	3.099	Pos. p = 0.0784
681	1	0.000	18.366	2.818	Pos. p = 0.0932
1176	1	0.000	21.955	3.039	Pos. p = 0.0813

```
hyphy fel --alignment spike.fas --tree spike.tree --branches Internal
```

# More accurate testing via parametric bootstrap

- P-values for MEME/FEL etc are derived from asymptotic approximations (large N)
- Not clear how well these hold for smaller and low-divergence datasets
- Can use a **much slower** simulation based method to derive more accurate p-values at each site
- Can result both in improved power and lower rates of false positives

# FEL on internal branches of Spike finds most selected sites

**CAUTION: A VERY TIME CONSUMING ANALYSIS (SEVERAL HOURS)**

```
hyphy fel --alignment spike.fas --tree spike.tree --branches Internal --output  
Spike-pbs.FEL.json --resample 100
```

# Obtaining site-level dN/dS estimates with FEL

- dN/dS estimates at individual sites are not **precise**
- They are estimated from relatively small samples
- Precision improves with the number of sequences and divergence levels
- One approach to correct for this is to compute approximate site-level confidence intervals.

Codon	Partition	alpha	beta	LRT	Selection detected?	dN/dS with confidence intervals
2	1	1.816	0.000	7.506	Neg. p = 0.0062	0.000( 0.00- 0.09)
247	1	1.353	0.000	8.080	Neg. p = 0.0045	0.000( 0.00- 0.10)
248	1	0.451	0.000	3.493	Neg. p = 0.0616	0.000( 0.00- 0.29)
249	1	0.000	2.700	7.884	Pos. p = 0.0050	10000.000(7599.84-10000.00)
250	1	0.224	0.000	2.797	Neg. p = 0.0945	0.000( 0.00- 0.61)

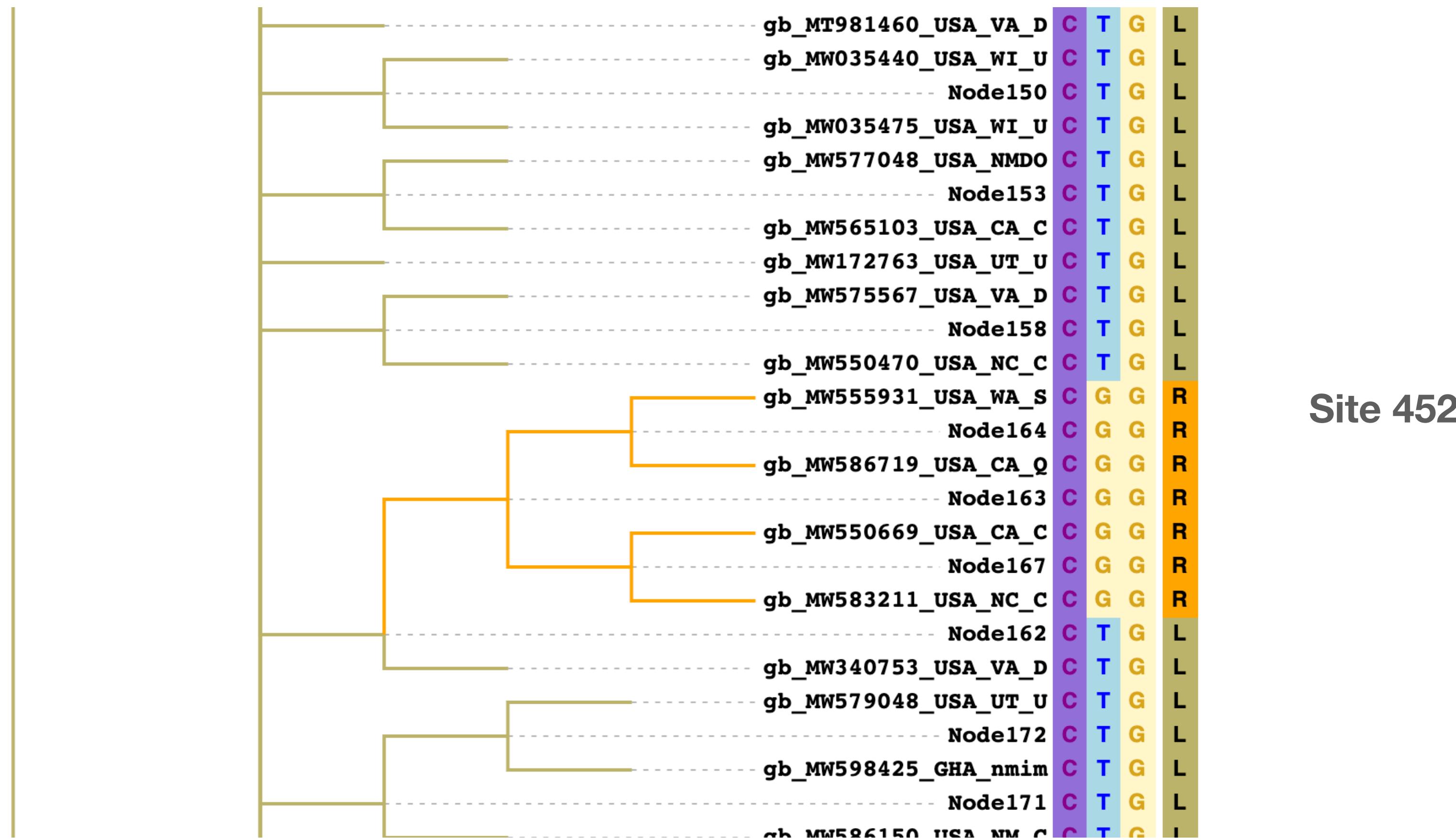
hyphy fel --alignment WestNileVirus\_NS3.fas --ci Yes

# Mapping substitutions with SLAC

- SLAC capable of detecting selection, is fast, but generally lacks power
- It provides a number of intuitive metrics for interpreting selection results
- SLAC recovers ancestral states and allows one to “map” evolutionary history onto a tree.

```
hyphy slac --alignment spike.fas --tree spike.tree --branches Internal
```

Partition	Site	ES	EN	S	N	P[S]	dS	dN	dN-dS	P [dN/dS > 1]	P [dN/dS < 1]	Total branch length
◆	◆	◆	◆	◆	◆	◆	◆	◆	◆	◆	◆	◆
1	452	1.75	1.25	0.00	1.00	0.584	0.00	0.801	49.6	0.416	1.00	0.0162
1	367	0.996	2.00	0.00	2.00	0.332	0.00	0.998	61.8	0.446	1.00	0.0162



# Analysis summary

	WNV NS3	HIV-1 env	SARS-CoV-2 spike
Gene-wide episodic selection (BUSTED)	No	Yes	Yes
Branch-level selection (aBSREL)	No	Yes, three branches, including transmission	No
Site-level episodic selection (MEME)	Yes, 1 site	Yes, 8 sites	Yes, sites found depend on which branches are included

## **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.

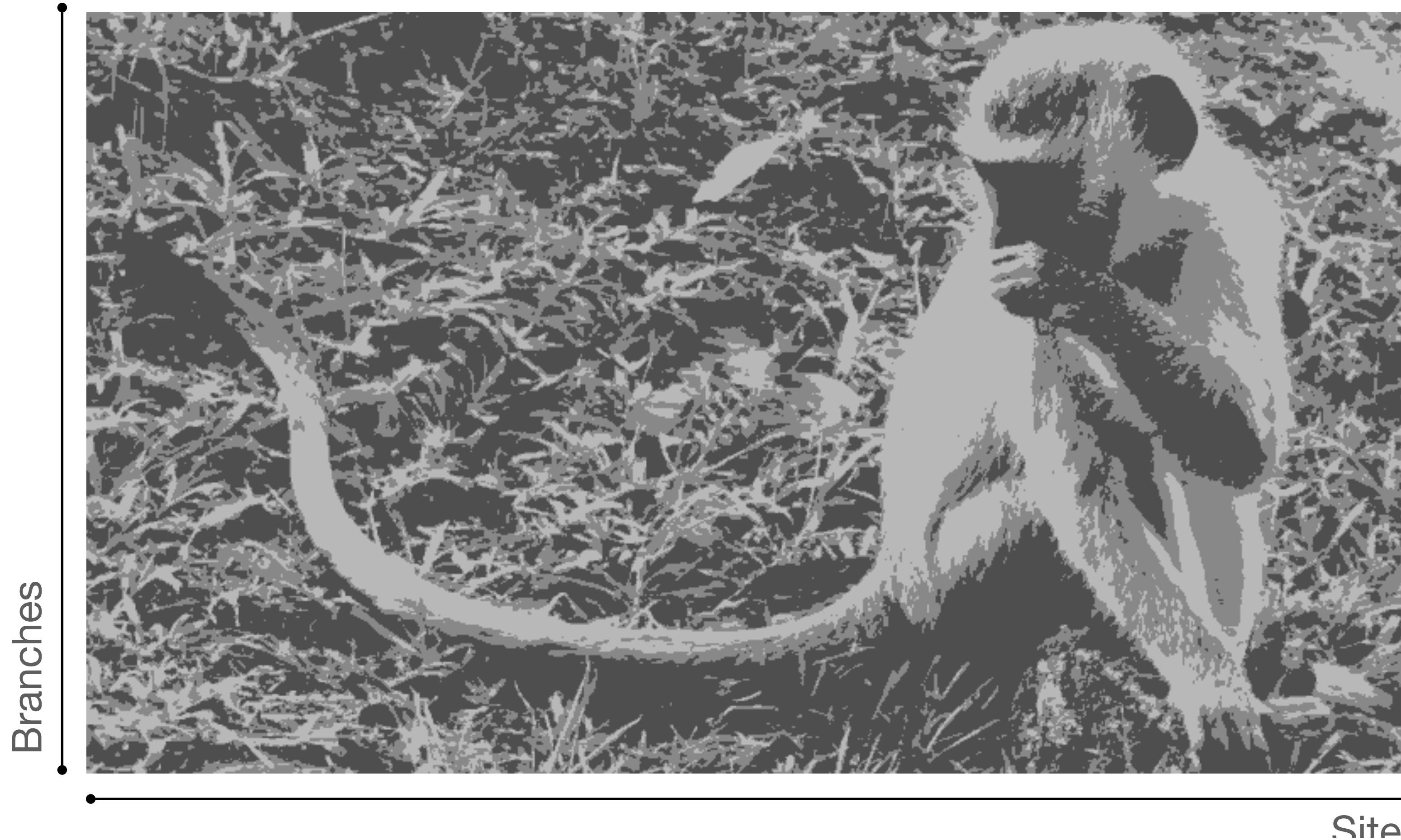
However, we caution that despite obvious interest in identifying specific branch-site combinations subject to diversifying selection, such inference is based on very limited data (the evolution of one codon along one branch), and cannot be recommended for purposes other than data exploration and result 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



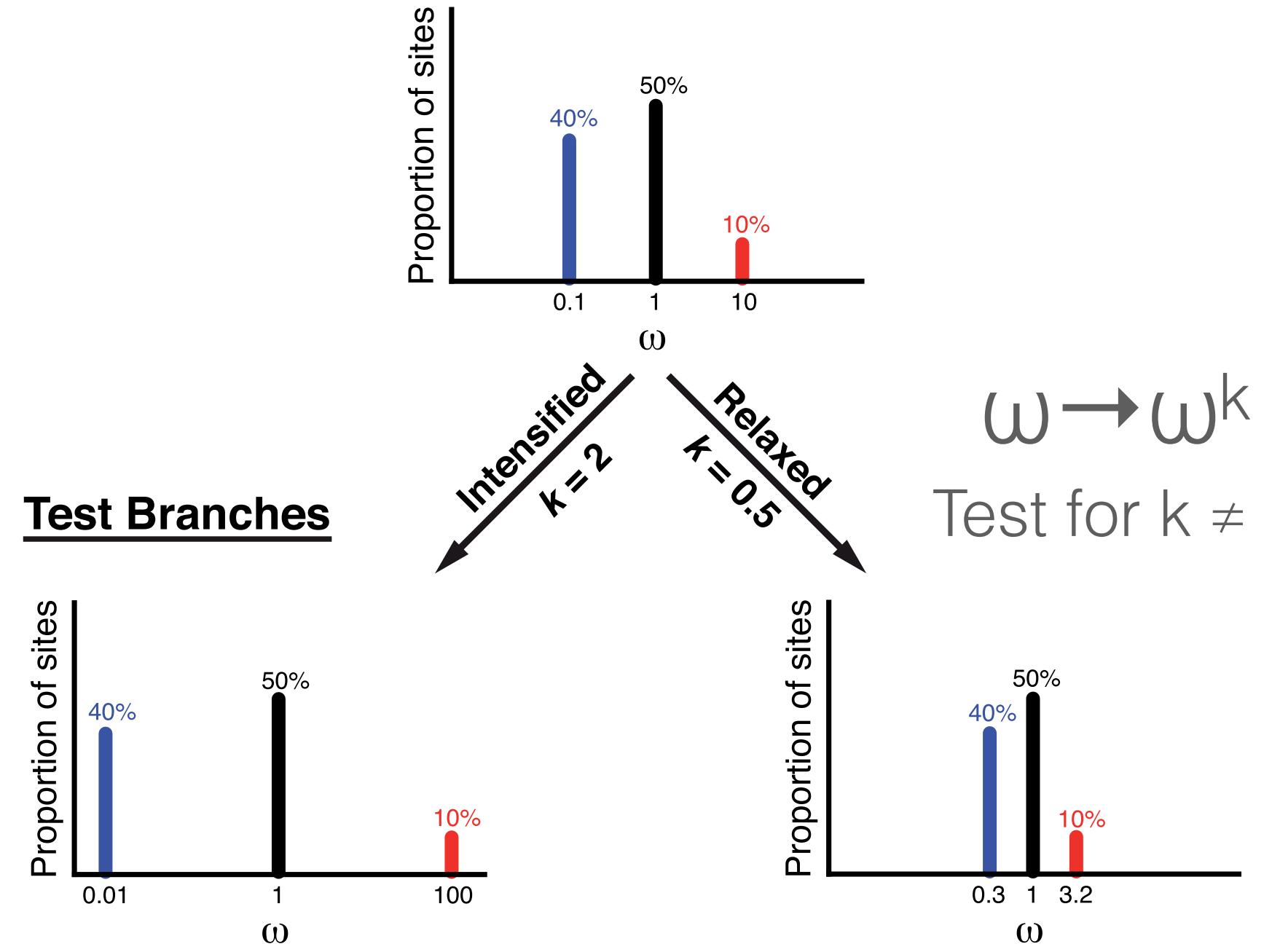
Partition the image into horizontal bands (a priori); compare whether or not there is visual benefit to using separate 3-color palettes in two sets of bands instead of a single 3-color palette



[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



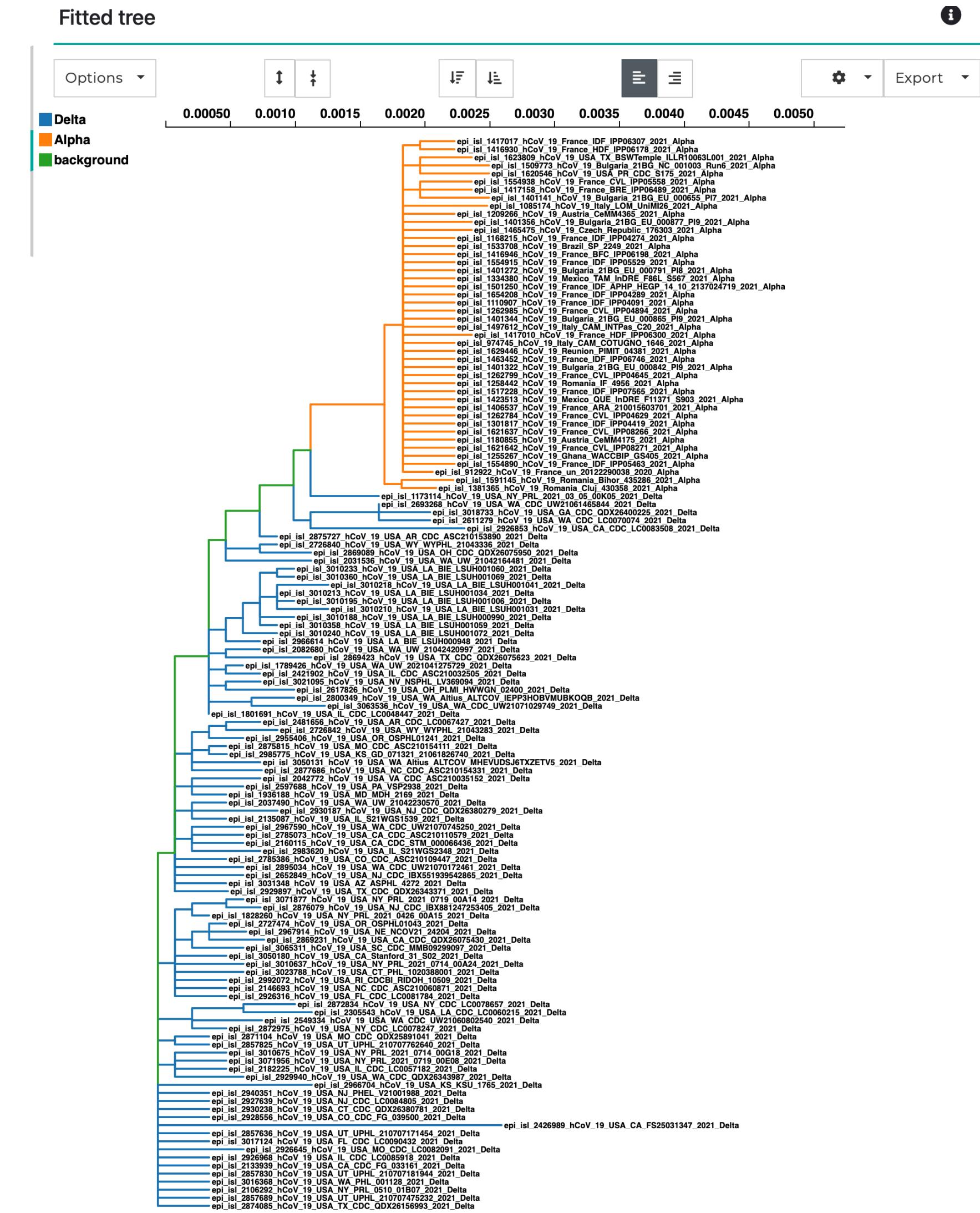
**Table 1.** Test for Relaxed Selection Using RELAX in Various Taxonomic Groups.

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
	Daphnia pulex	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.

# Comparing alpha vs delta clades in SARS-CoV-2

- Are selective pressures on the Delta SARS-CoV-2 clade relaxed or intensified compared to the Alpha clade.
- Partition the tree into corresponding clades.
- See <http://www.hyphy.org/tutorials/CL-prompt-tutorial/#preparing-labeled-phylogenies> for how to label phylogenies



# RELAX(ed selection test) results summary

INPUT DATA | AlphaDeltaSpike.fas | 133 sequences | 1273 sites

 Export ▾

Test for selection **intensification** ( $K = 1.31$ ) was **not significant** ( $p = 0.558$ ,  $LR = 0.34$ ).

See [here](#) for more information about this method.

Please cite [PMID 123456789](#) if you use this result in a publication, presentation, or other scientific work.

## Model fits



Model	log L	#. params	AIC <sub>c</sub>	Branch set	$\omega_1$	$\omega_2$	$\omega_3$
General descriptive	-8790.1	367	18315.9	Shared	0.00 (11.36%)	0.86 (88.62%)	1288.13 (0.02%)
RELAX alternative	-8876.3	199	18151.1	Reference	1.00 (97.76%)	1.00 (2.24%)	1450.96 (0.00%)
				Test	1.00 (97.76%)	1.00 (2.24%)	13744.44 (0.00%)
RELAX null	-8876.5	198	18149.4	Reference	1.00 (98.05%)	1.00 (1.95%)	11625.16 (0.00%)
				Test	1.00 (98.05%)	1.00 (1.95%)	11625.16 (0.00%)

# Which sites are evolving differentially?

- We have established that in the HIV example, donor, recipient, and transmission branches evolve differently.
- Can we identify specific sites where this may be occurring?
  - Why is this of interest?
- More generally, given a tree with  $N$  sets of branches, we fish to find sites where evolution is different between these  $N$  sets, with a degree of statistical significance.
- Solution: use a fixed effects method (Contrast-FEL)
  - For each branch set  $i$ , estimate a  $dN/dS$  ratio ( $N$  total ratios)
  - Test whether or not any of the ratios are different (group test)
- For each pair of ratios, test if they are different [ up to  $N(N-1)/2$  tests ]
- Can identify subtle differences among selective pressures.

# Contrast-FEL results summary

**INPUT DATA** | **AlphaDeltaSpike.fas** | **133** sequences | **1273** sites

 Export

ContrastFEL found evidence of

1

Found 5 sites with different Overall dN/dS

with p-value threshold of 0.01

See [here](#) for more information about this method.

Please cite [PMID 15703242](#) if you use this result in a publication, presentation, or other scientific work.

## ContrastFEL Table

1

Showing entries 1 through 20 out of 1273

 Export Table to CSV

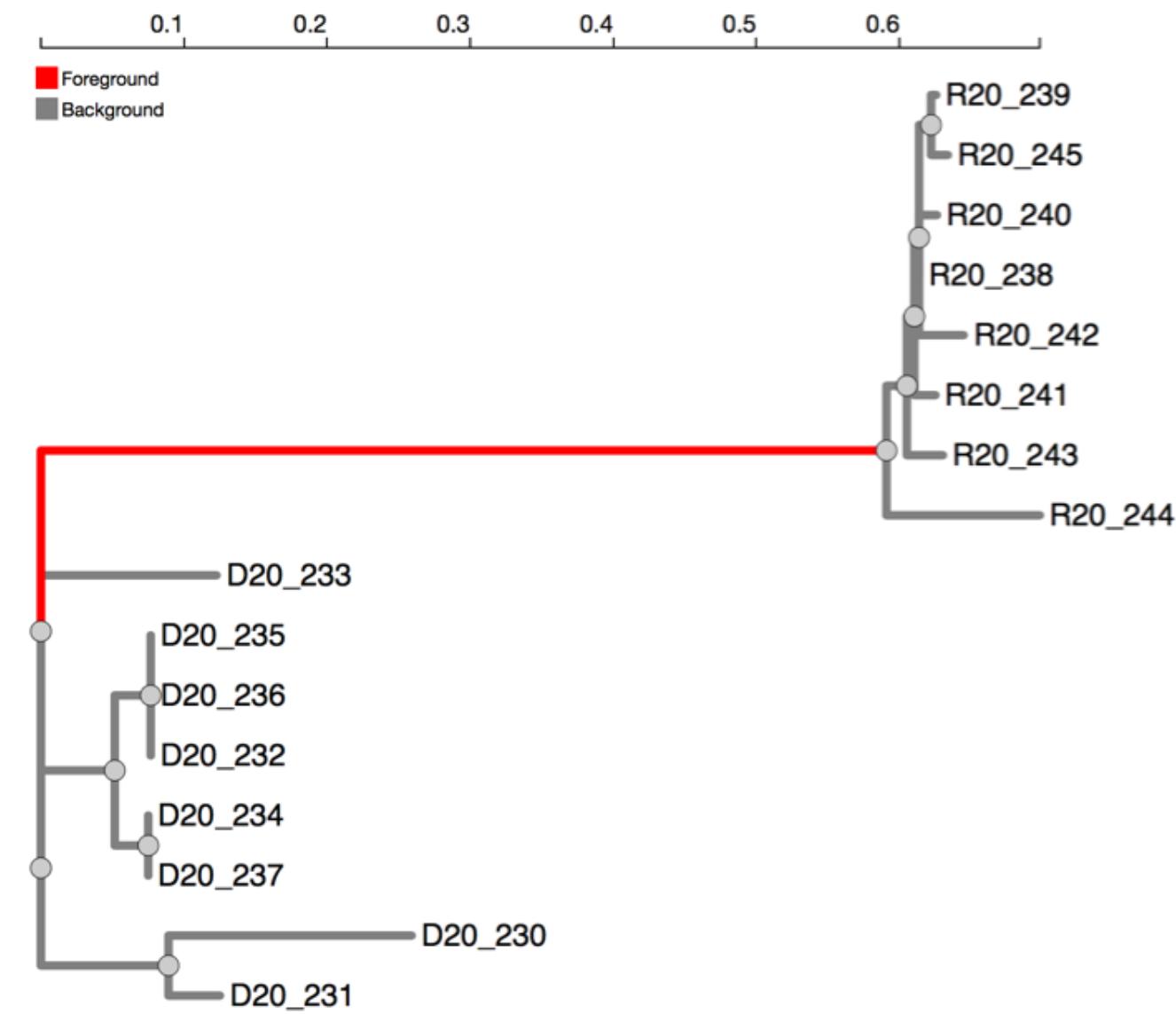
<< < > >>

Site	Partition	alpha	beta (Delta)	beta (Alpha)	beta (background)	subs (Delta)	subs (Alpha)	P-value (overall)	Q-value (overall)	Permutation p-value	b1
			◆	◆	◆	◆	◆	↓	◆	◆	◆
70	1	6.145	11.989	710.847	2.892	3.000	2.000	0.000	0.325	-1.000	3
157	1	7.355	73.889	0.000	0.445	5.000	0.000	0.001	0.623	-1.000	3
950	1	0.037	54.713	0.000	740.108	6.000	0.000	0.003	1.000	-1.000	3
142	1	0.000	40.094	7.705	460.016	6.000	1.000	0.004	1.000	-1.000	3
77	1	0.000	41.101	0.000	0.835	6.000	0.000	0.009	1.000	-1.000	3

```
hyphy contrast-fel --alignment AlphaDeltaSpike.fas --tree AlphaDeltaSpike.nwk --branch-set Alpha --branch-set Delta
```

# 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	~sites x taxa	SMP	~1,000
<b>Site-level selection / episodic</b>	MEME	Fixed Effects	Random Effects	Fixed	~ taxa	SMP/MPI	~5000 (cluster)
<b>Site-level selection / pervasive</b>	FEL	Fixed Effects	Fixed Effects	Fixed	~ taxa	SMP/MPI	~20000 (cluster)
<b>Branch-level selection</b>	aBSREL	Random Effects	Fixed Effects	Adaptive	~ sites	SMP/MPI	~ 1,000
<b>Compare selective regimes between sets of branches</b>	RELAX	Random Effects	Mixed Effects	Fixed	~sites x (branch set size)	SMP	~ 1,000
<b>Compare selective pressure between sets of branches on individual sites</b>	Contrast-FEL	Fixed Effects	Fixed Effects	Fixed	~ branch set	SMP/MPI	~5000 (cluster)

# Current suggested best practices.

There are lots of methods you could use to study positive selection, including >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>
Intra-species viral analyses for sites under selection	MEME/FEL internal branches	<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.

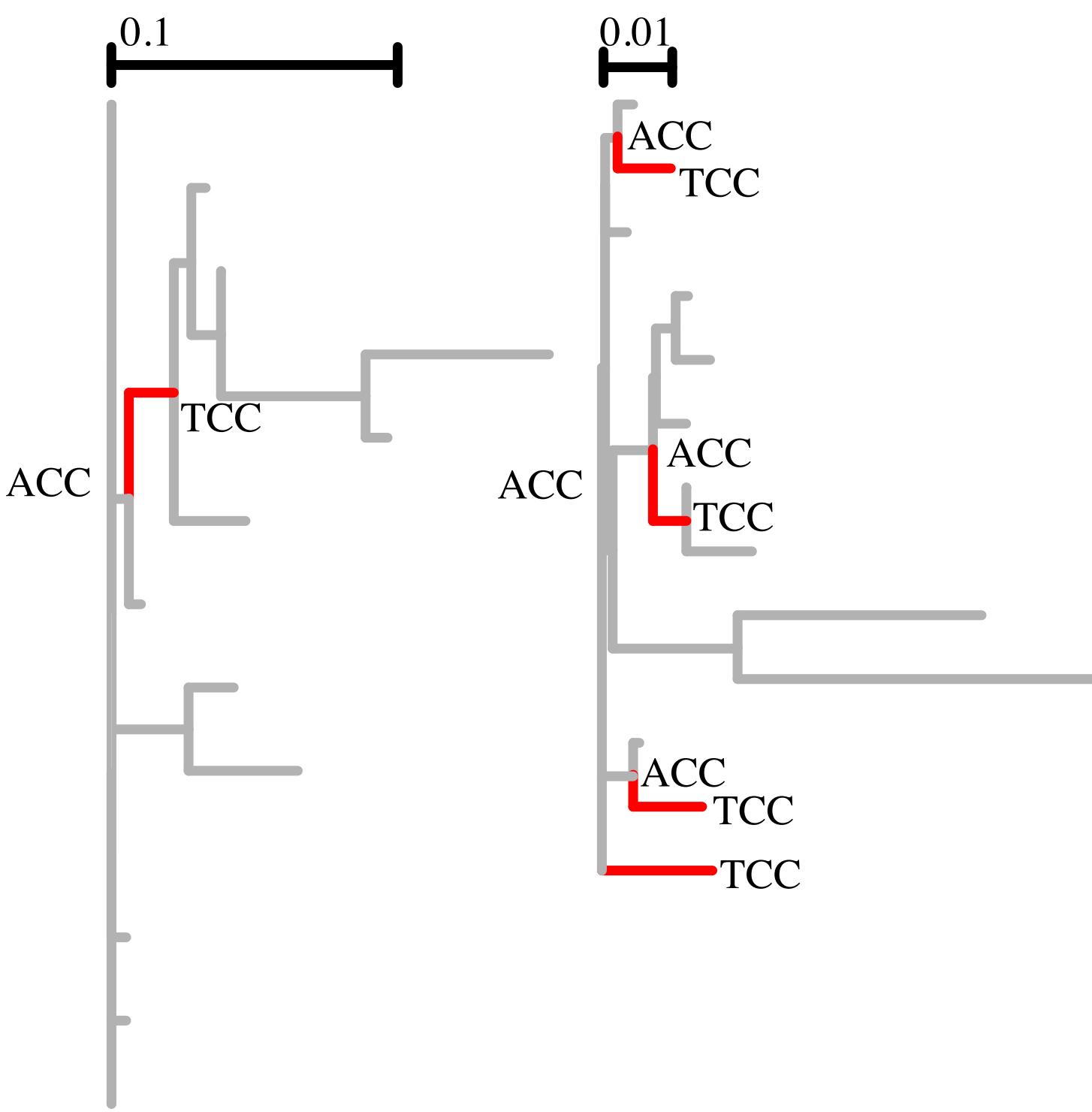


Figure 4.2: The effect of recombination on inferring diversifying selection. Reconstructed evolutionary history of codon 516 of the Cache Valley Fever virus glycoprotein alignment is shown according to GARD inferred segment phylogeny (left) or a single phylogeny inferred from the entire alignment (right). 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

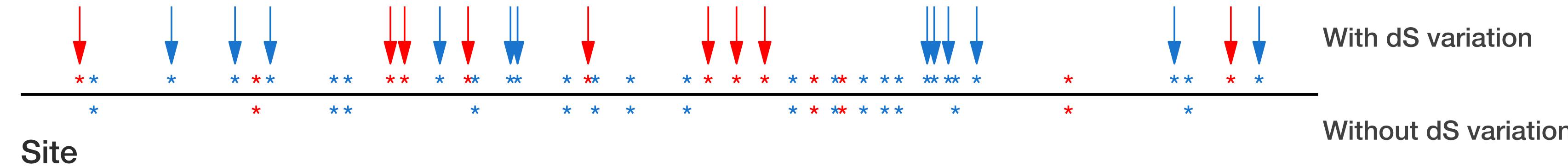
- First screen the alignment to find putative non-recombinant fragments (e.g. using GARD)
- Apply a model-based test (MEME, FUBAR) using multiple phylogenies (one per fragment), but inferring other parameters (e.g. nucleotide substitution biases and base frequencies) from the entire alignment
- This 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

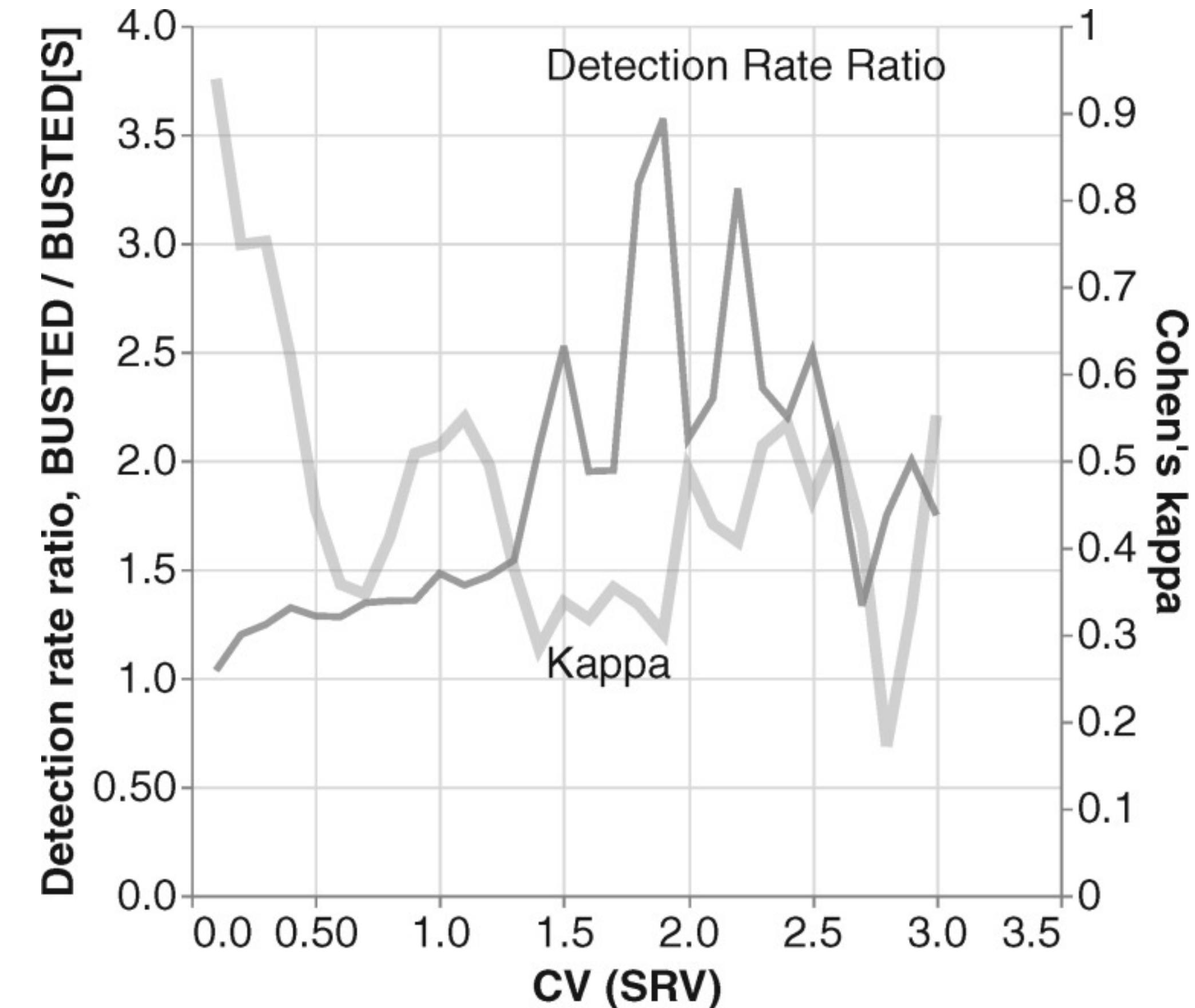
Test  $p < 0.1$  was used to classify sites as selected. Codon sites found under selection by both methods are shown in bold.

### Sites detected by FEL with and without dS variation



# 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 and high false positive rates
- Most of HyPhy methods provide support for including  $dS$



**Synonymous Site-to-Site Substitution Rate Variation Dramatically Inflates False Positive Rates of Selection Analyses: Ignore at Your Own Peril**

Sadie R Wisotsky <sup>1 2</sup>, Sergei L Kosakovsky Pond <sup>2</sup>, Stephen D Shank <sup>2</sup>, Spencer V Muse <sup>1 3</sup>

Affiliations + expand

PMID: 32068869 PMCID: PMC7403620 DOI: 10.1093/molbev/msaa037

Free PMC article

# Allowing multi-nucleotide substitutions

- Some of the methods (e.g. BUSTED, aBSREL, RELAX) can extend substitution models to allow instantaneous double- and triple-“hits” (e.g. ACC to AGG)
- Sometimes multi-nucleotide changes along short branches at a single site can drive selection signal (possible false positives?)
- HyPhy includes a simple standard analysis for estimating alignment-wide multiple-hit rates.

**Extra base hits: Widespread empirical support for instantaneous multiple-nucleotide changes**

Alexander G. Lucaci , Sadie R. Wisotsky , Stephen D. Shank, Steven Weaver, Sergei L. Kosakovsky Pond 

Published: March 12, 2021 • <https://doi.org/10.1371/journal.pone.0248337>

[See the preprint](#)

# Multi-Hit results summary

INPUT DATA | spike.fas | 118 sequences | 1273 sites

 Export ▾

## Likelihood Test Results

3H vs  1H

LRT 2.307

p-value 0.511

3H vs.  2H

LRT 0.000

p-value 1.000

3H vs  3HSI

LRT -0.000

p-value 1.000

2H vs  1H

LRT 2.307

p-value 0.129

3HSI vs  2H

LRT 0.000

p-value 0.988

See [here](#) for more information about this method.

Please cite [biorxiv](#) if you use this result in a publication, presentation, or other scientific work.

---

```
hyphy fmm --alignment spike.fas --tree spike.tree --triple-islands Yes
```

# results summary

INPUT DATA | WestNileVirus\_NS3.fas | 19 sequences | 619 sites

 Export ▾

## Likelihood Test Results

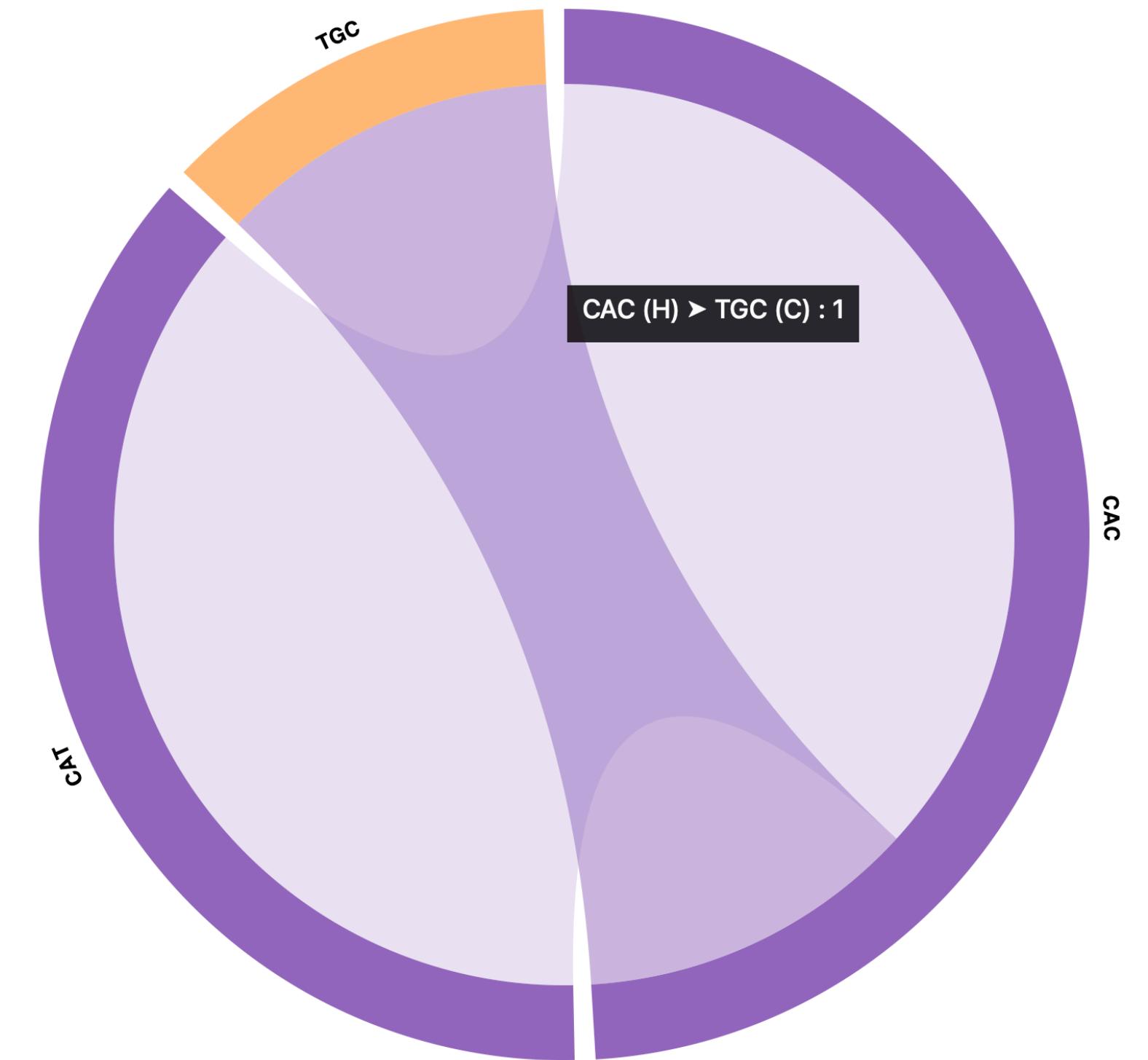
<input checked="" type="checkbox"/> 3H vs <input type="checkbox"/> 1H
LRT 7.406
p-value 0.060

<input checked="" type="checkbox"/> 3H vs. <input type="checkbox"/> 2H
LRT 0.006
p-value 0.997

<input checked="" type="checkbox"/> 3H vs <input checked="" type="checkbox"/> 3HSI
LRT -0.000
p-value 1.000

<input checked="" type="checkbox"/> 2H vs <input type="checkbox"/> 1H
LRT 7.400
p-value 0.007

<input checked="" type="checkbox"/> 3HSI vs <input type="checkbox"/> 2H
LRT 0.006
p-value 0.936



See [here](#) for more information about this method.

Please cite [bioRxiv](#) if you use this result in a publication, presentation, or other scientific work.

### 1 individual site which showed sufficiently strong preference for multiple-hit models

Site	ER (2 vs 1)	ER (3 vs 2)	ER (3-island vs 2)	ER (3-island vs 3)	Substitutions
87	54.0092	1.0065	1.0020	1.0045	CAC->CAT(3)TGC(1)

hyphy fmm --alignment WestNileVirus\_NS3.fas --triple-islands Yes