

ECCB'12

11th European Conference on

Computational Biology

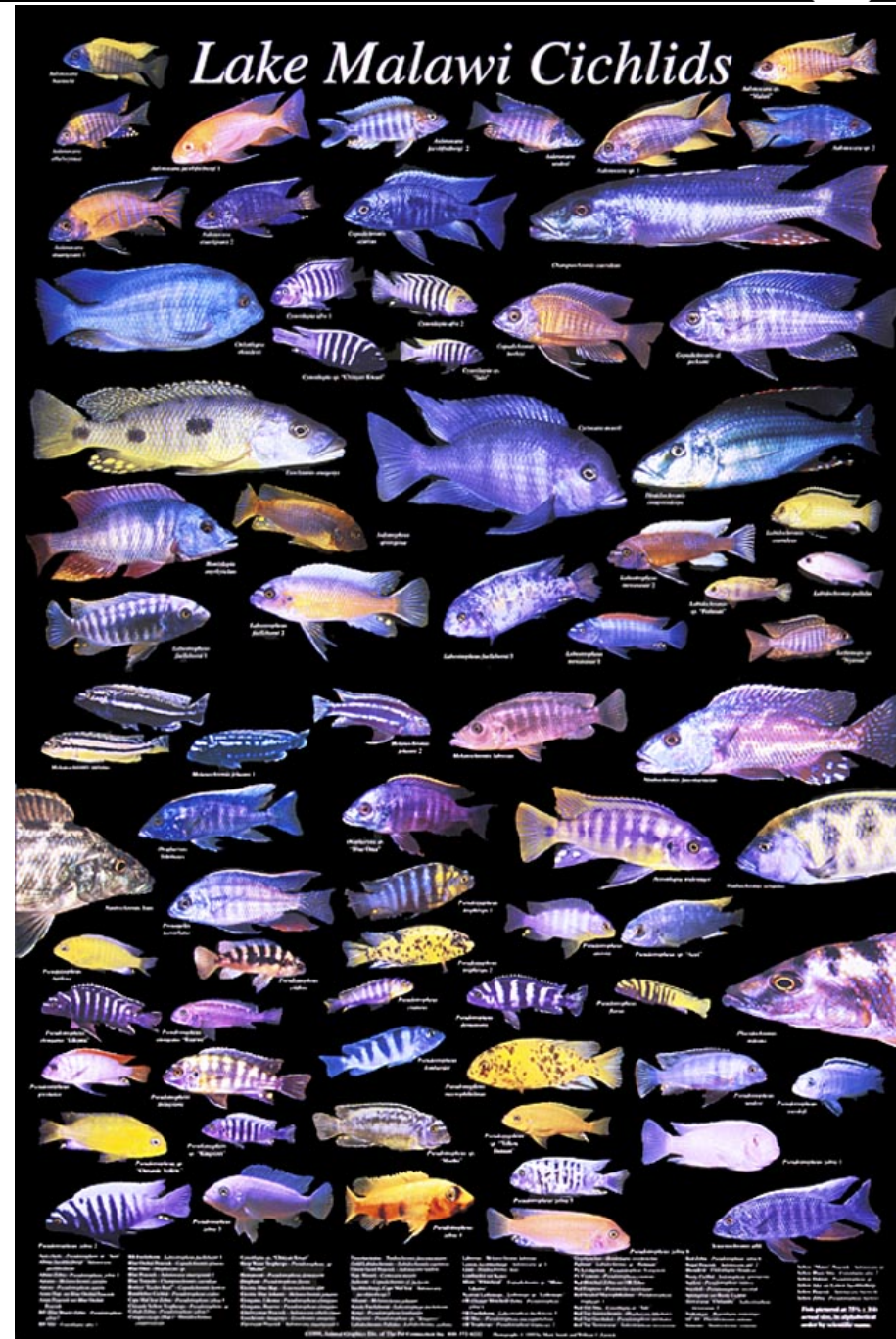
9-12 September 2012, Basel, Switzerland



Detection of evolutionary shifts in protein sequences

Evolution: morphological differences

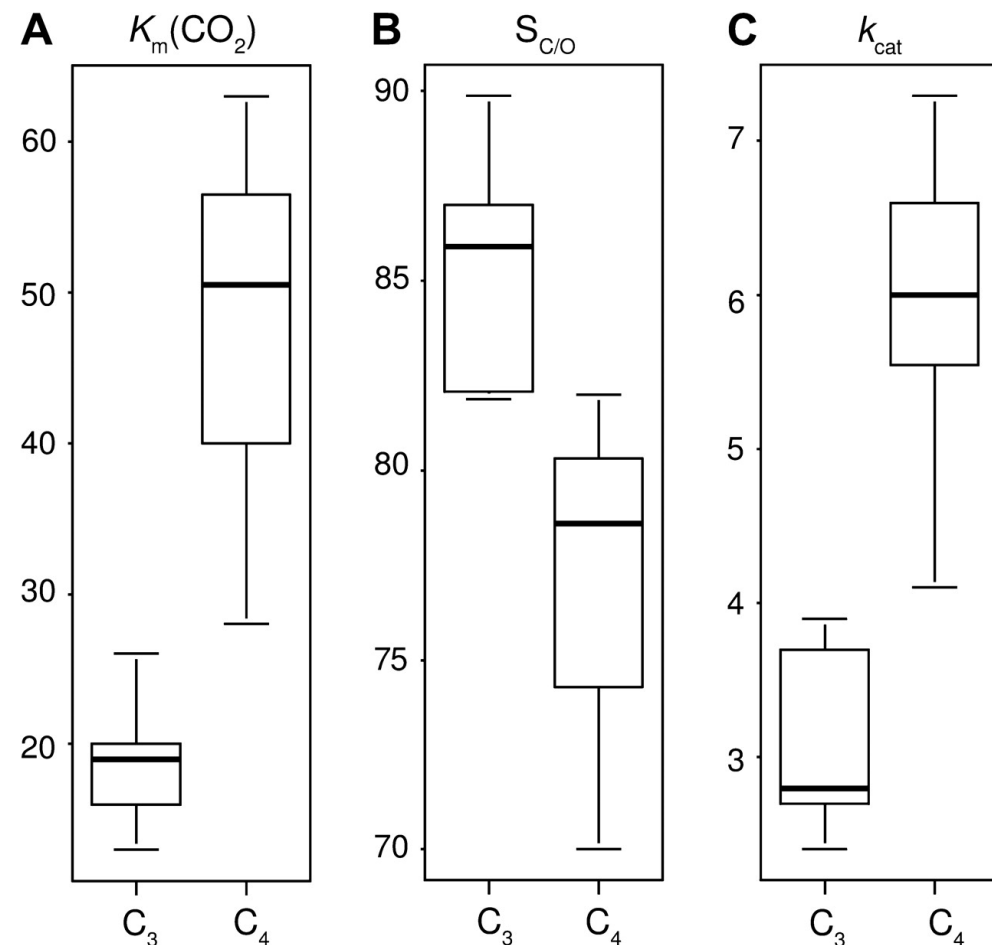
Can we explain the different
colour patterns?



Evolution: catalytic activities

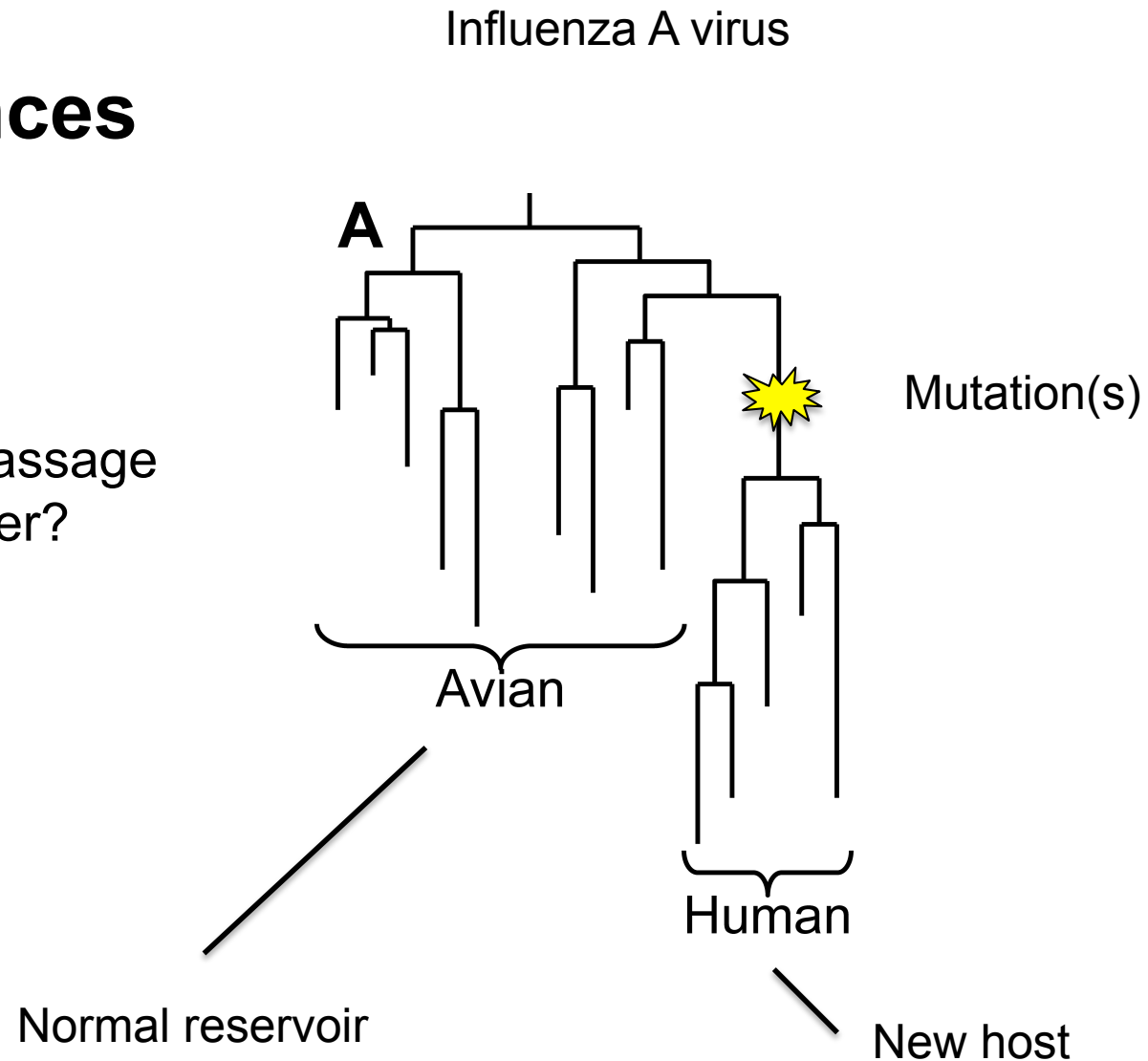
Can we explain the different
Kinetic parameters?

RubisCO enzyme from C₃ and C₄
forms in plants diverge by different
kinetic parameters (affinity, specificity
and activity)



Evolution: host preferences

Can we explain the passage
from an host to another?



Vectors of evolution

- Epigenetics (non genome modification, i.e. DNA methylation).
- Genetics (genome modification):
 - Changes in expression pattern or level.
 - Changes in non coding RNA (ncRNA).
 - Changes in proteins biochemistry:
 - => Change in the nucleotide sequence.
 - => Change in the amino acid sequence.

From genotype to phenotype

- Genotype (nucleotides, amino acids).



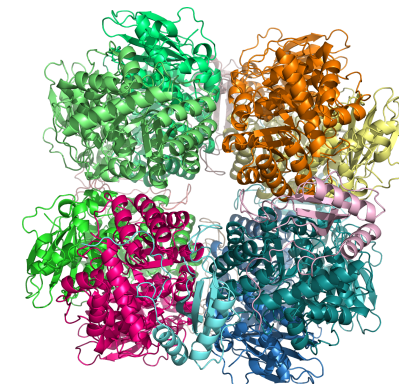
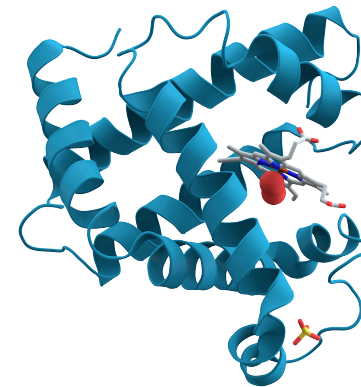
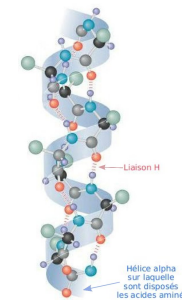
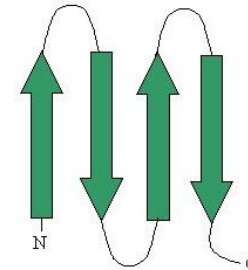
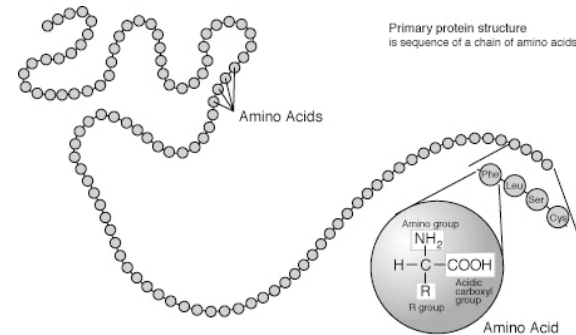
- Phenotype.



- Submitted to selective pressure.

Proteins: 4 levels of organisation

- Primary structure:
 - amino acids sequence.
- Secondary structure
 - simple elements.
- Tertiary structure:
 - folded in 3D.
- Quarternary structure:
 - assembly of multiple monomeres.



Conservation

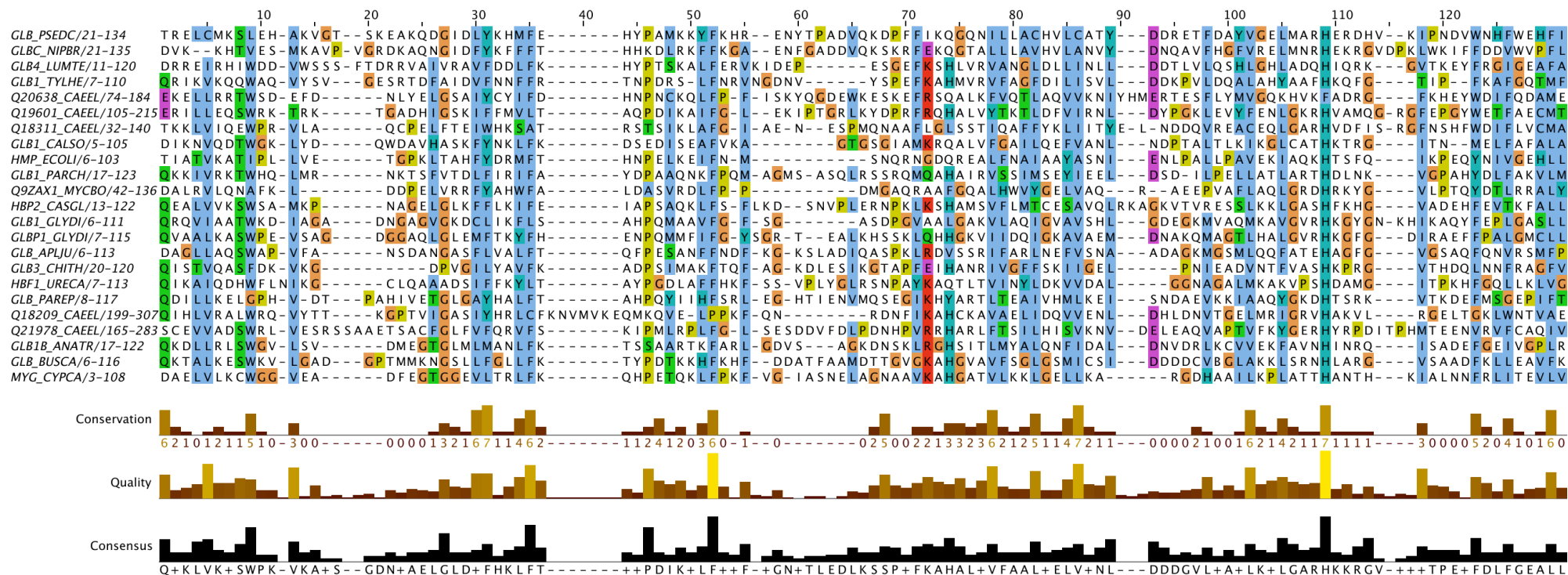
Quality

Consensus

Q+K+L+V++SWP+K+V+A+G+G+G+G+A+I+L+F+L+F+L+F+P+E+Q+K+A+F+F+F+K+G+E++G+L+A+A+S+K+++++A+V+V+A+L+L+L+V+L+L+D+K+D+K+A+V+L+A+G+H+G+H++F+G+R+G+F+++F+K+L+F+A+F+M+E--F-----

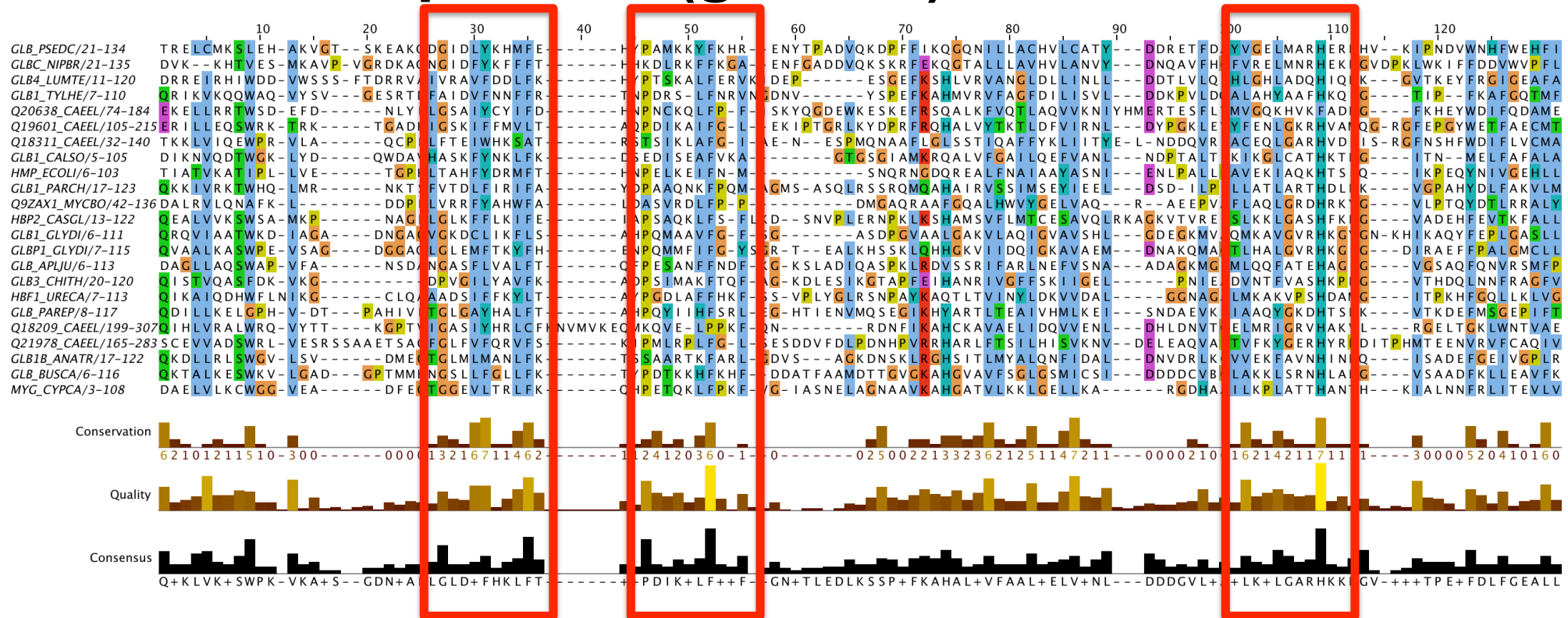
Group of highly divergent sequences

Protein sequences (globins)



Multiple aligner: ClustalW, MUSCLE, MAFFT, PROBCONS, PRANK

Protein sequences (globins)

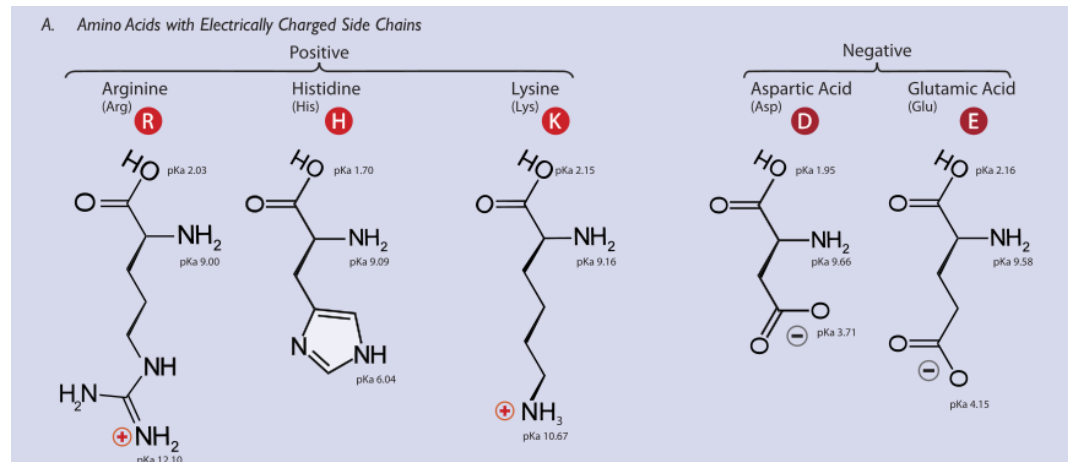


Tools to select blocks: Gblocks, Guidance, M-coffee, Pagan

Amino acids have different physicochemical properties

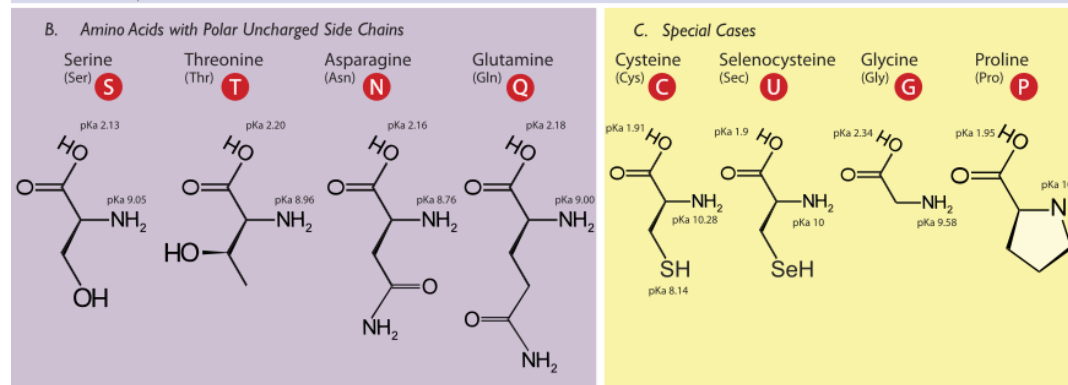
Charged:
Acid (-)

Charged:
Basic (+)



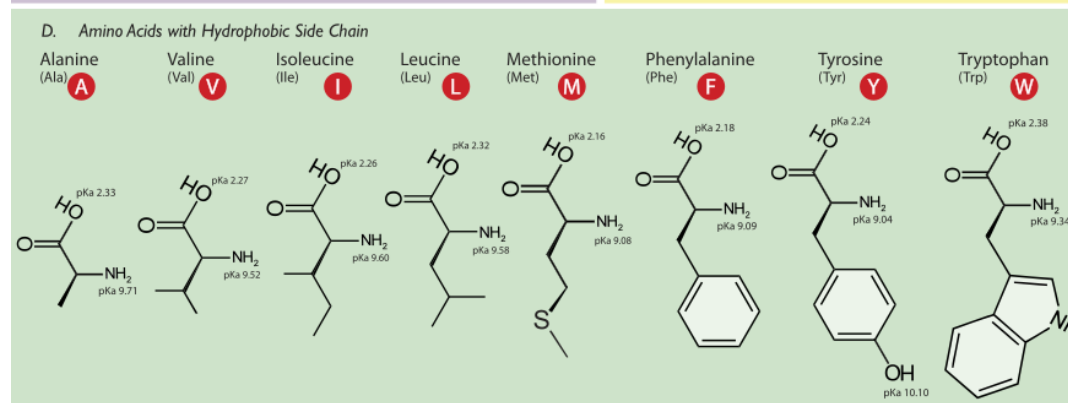
Polar

Special cases

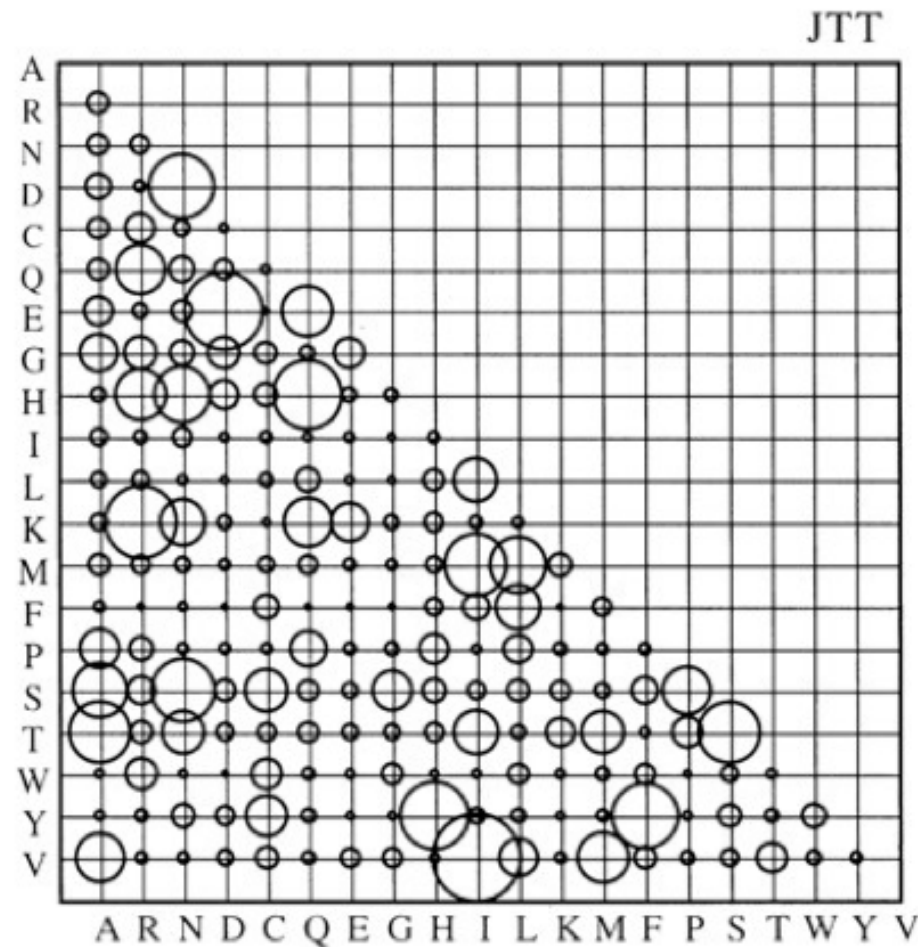


Hydrophobic
small

Hydrophobic
big

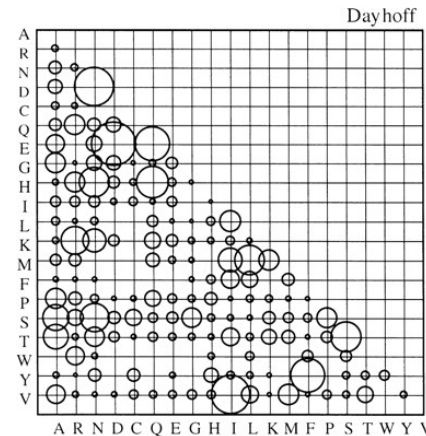


**Amino acids have different physicochemical properties
=> Replacements of amino acid are not equivalent**

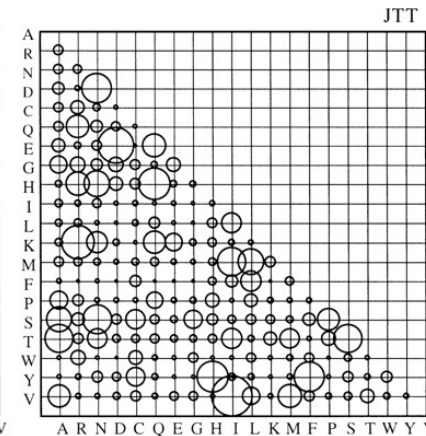


Schematic representations of amino acid replacement matrices.

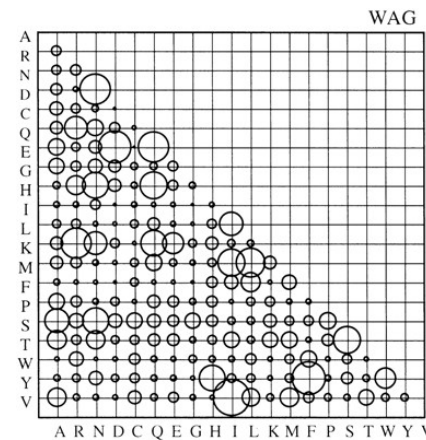
Dayhoff
(Dayhoff,
Schwartz,
Orcutt,
1978)



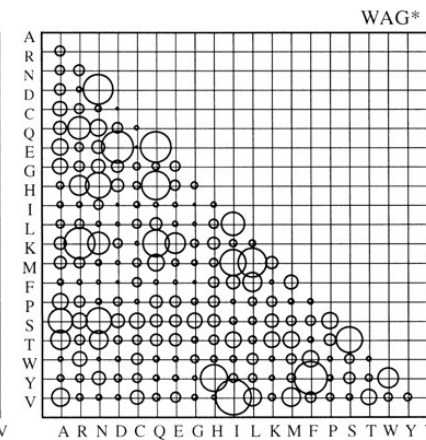
JTT
(Jones,
Taylor WR and
Thornton,
1992)



WAG
(Whelan and
Goldman,
2001)



WAG
(Whelan and
Goldman,
2001)



○ Average rate

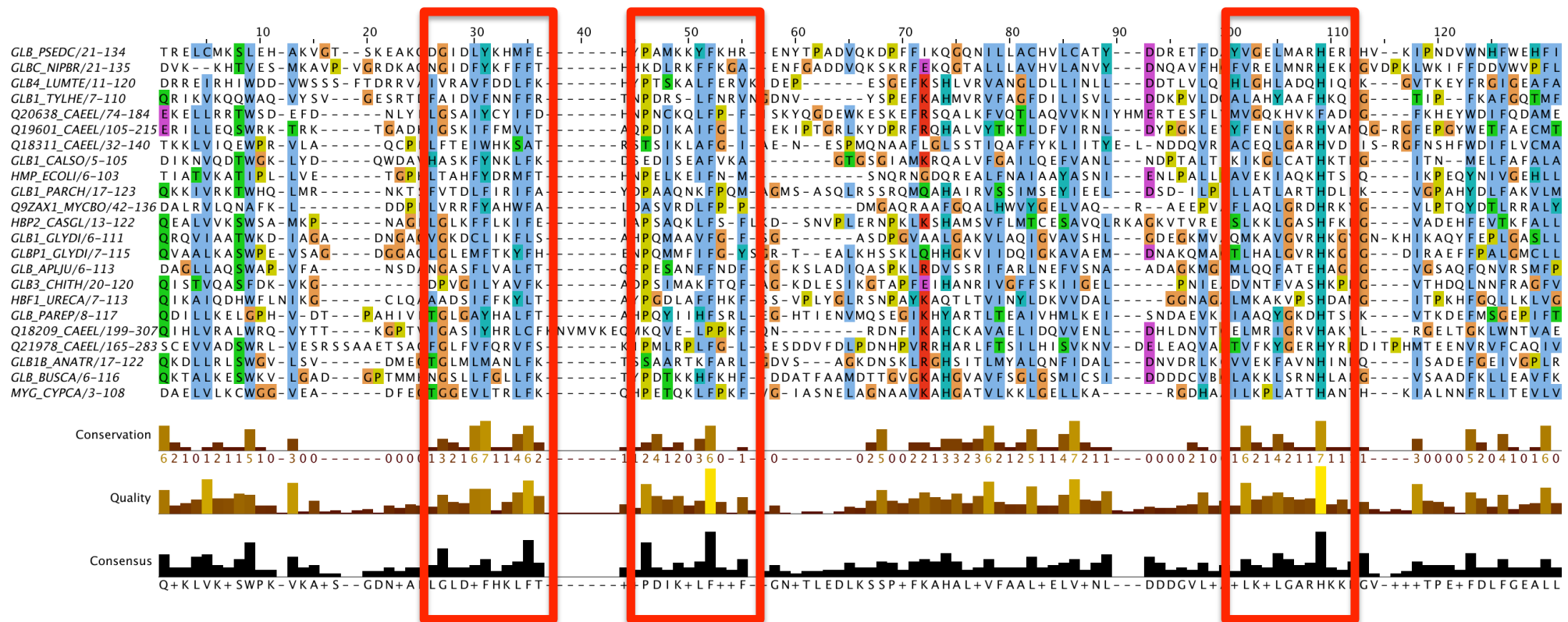
LG (Le and Gascuel, Mol Biol Evol 2008)

LG4M and LG4X (Le, Dang, Gascuel, Mol Biol Evol 2012)

How to choose the best matrix?

- By experience (i.e. JTT/WAG for Vertebrates, mtREV24 for mitochondrial gene).
- Using tools: ProtTest, ReplacementMatrix

Positions don't evolve at the same rate



Positions don't evolve at the same rate

Evolutionary rate
(number of mutation
per site)

Slow:

- Presence of structural/functional constraints



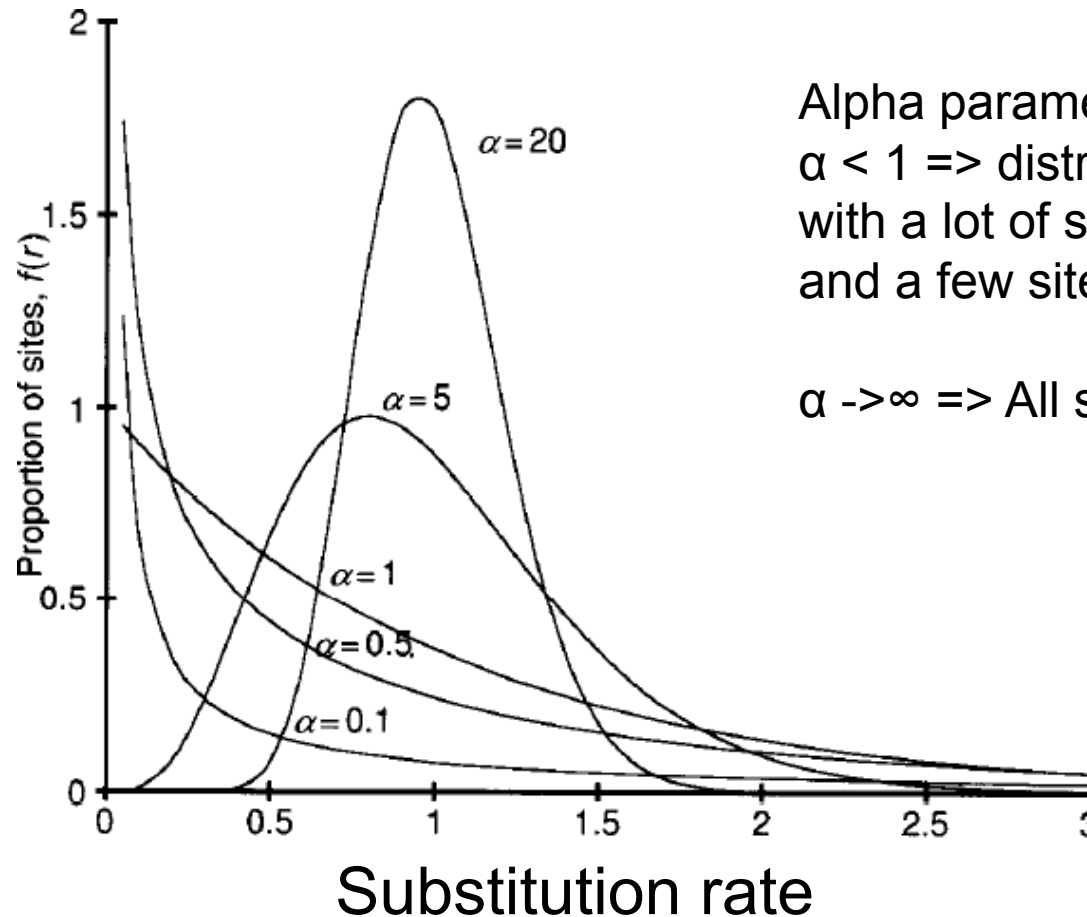
Fast:

- Absence of structural/functional constraints
- Change in structural/functional constraints

- Can be described by a **gamma distribution**.
 - Need to cluster sites into different categories of evolutionary rates (generally, 4 or 8 categories).

Positions don't evolve at the same rate

% of sites



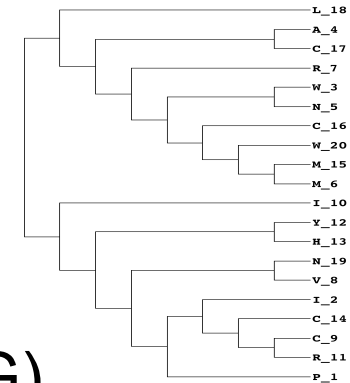
Alpha parameter:

$\alpha < 1 \Rightarrow$ distribution shaped in **L form**, with a lot of sites evolving slowly and a few sites evolving fast.

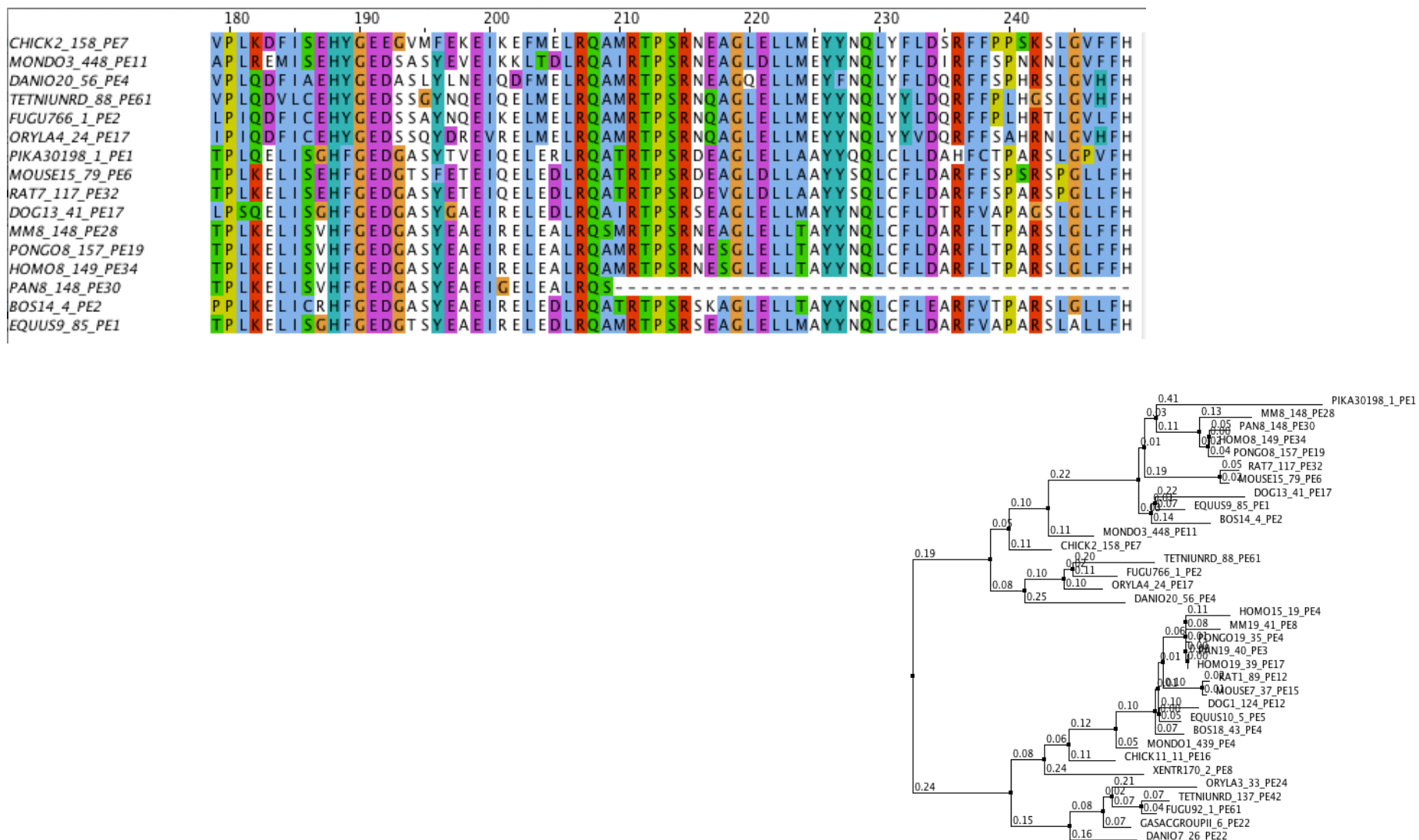
$\alpha \rightarrow \infty \Rightarrow$ All sites evolve at the same rate

Phylogenetic tree building

- Data: Multiple alignment of sequence.
- Matrix of replacement rate (i.e. JTT, WAG, LG).
- Use of gamma distribution. => Estimate alpha parameter and categorise sites.
- Method to build the tree:
 - Unweighted Pair Group Method with Arithmetic Mean (UPGMA)
 - Neighbour-Joining (BioNJ)*
 - Maximum parsimony **
 - Minimum Evolution (ME) (FastME) **
 - Maximum likelihood (ML) (PhyML, RAxML) ***
 - Bayesian (MrBayes) ***



Results: multiple alignment and tree



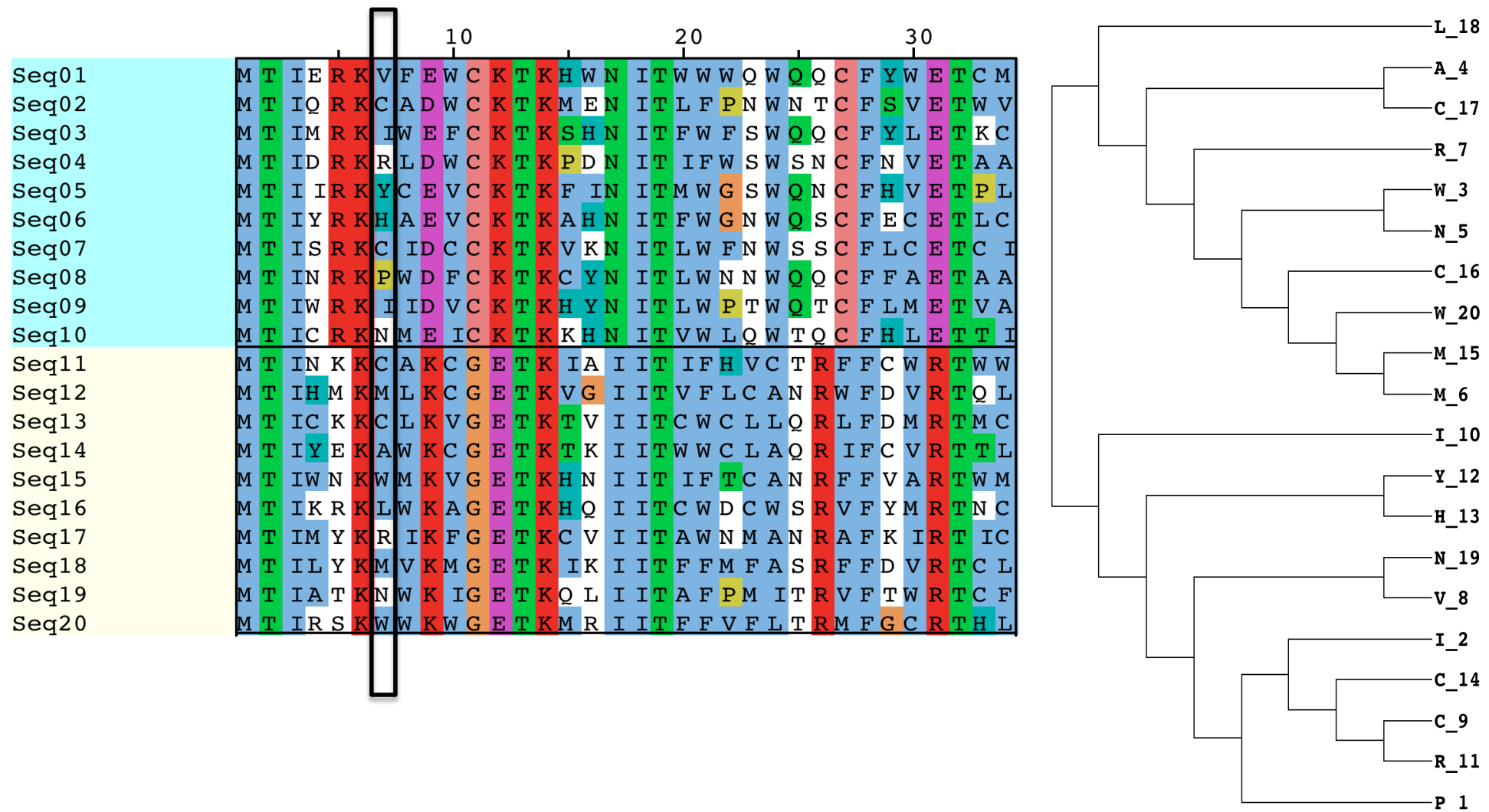
Ancestral Sequence Reconstruction

- Multiple Sequence Alignment
- Phylogenetic tree
- Substitution parameters (Matrix, Gamma distribution)

=> Infer each amino acid at each node.

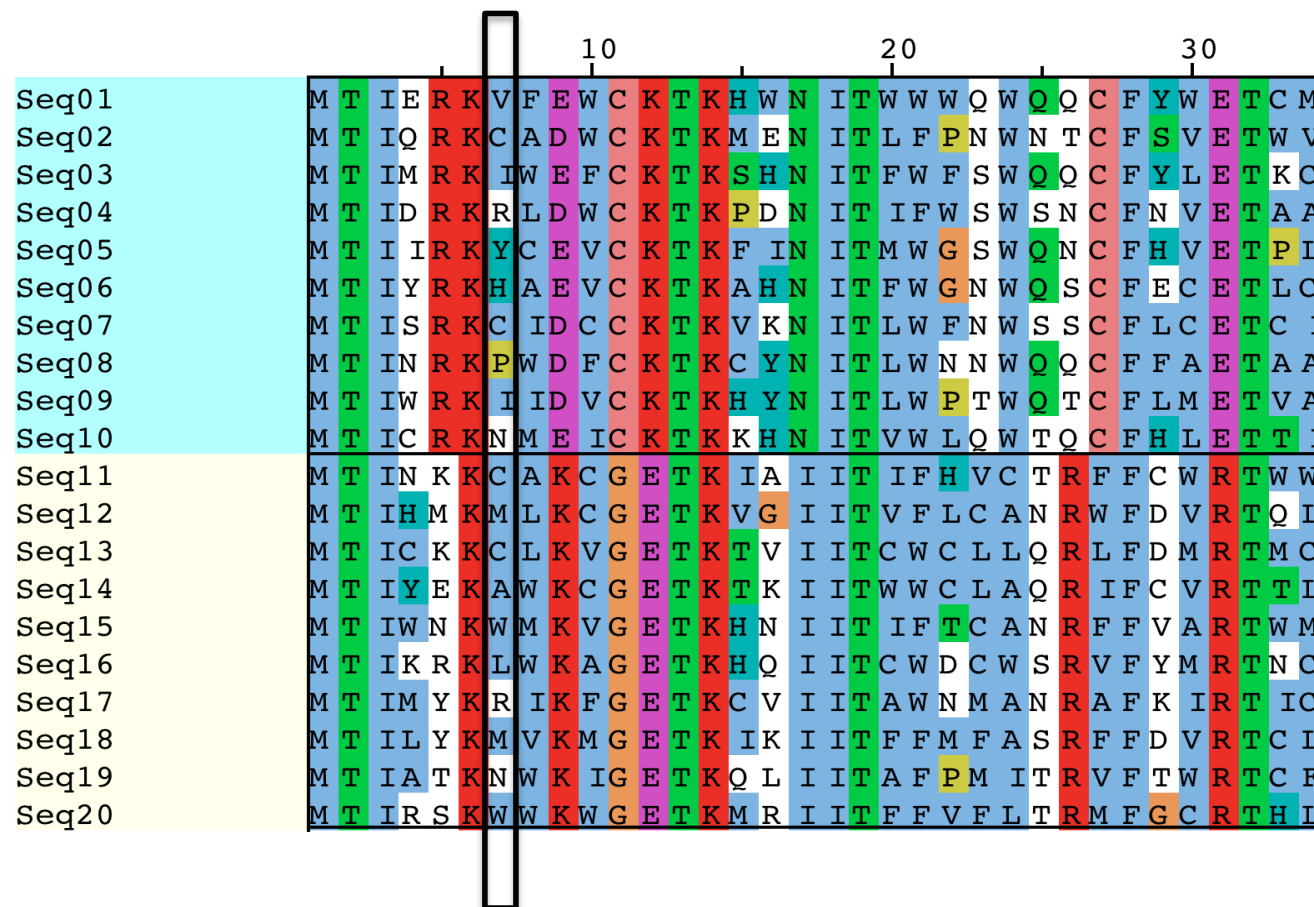
=> Reconstruct ancestral sequences.

Ancestral Sequence Reconstruction

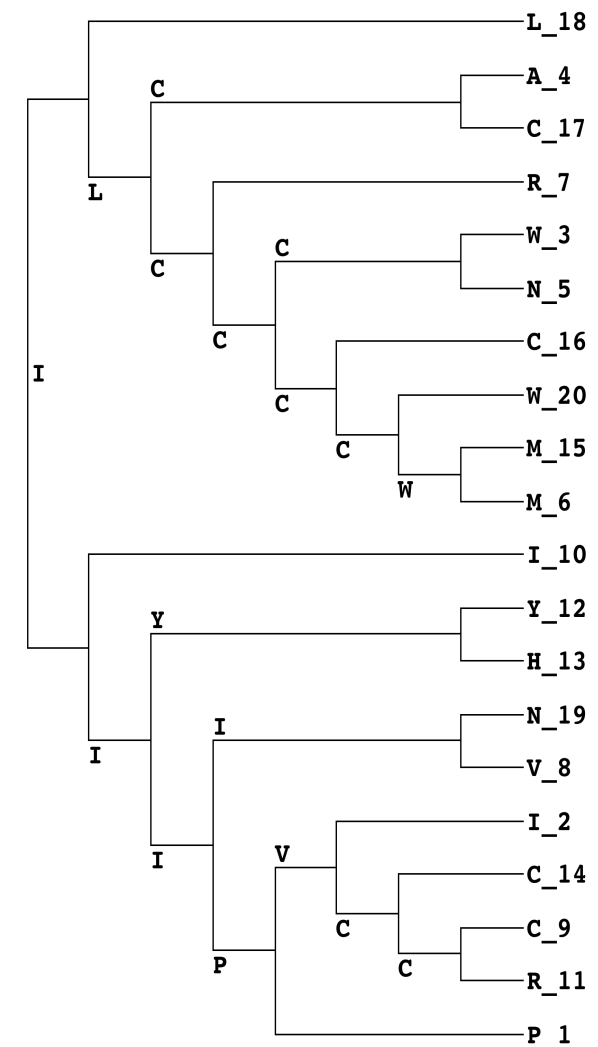


(Simulated alignment)

Ancestral Sequence Reconstruction



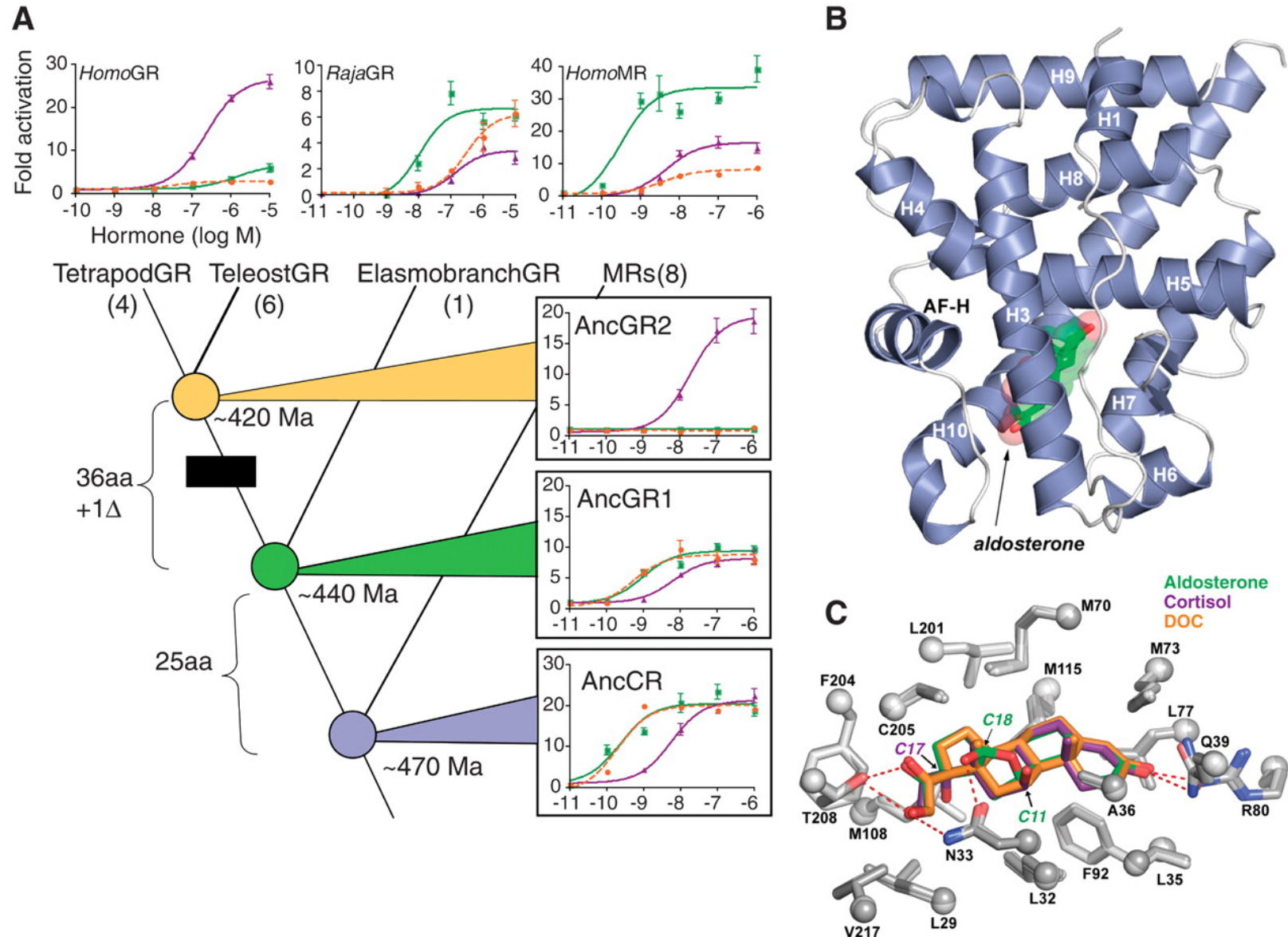
Amino acid at site 7



Control quality with Posterior Probabilities

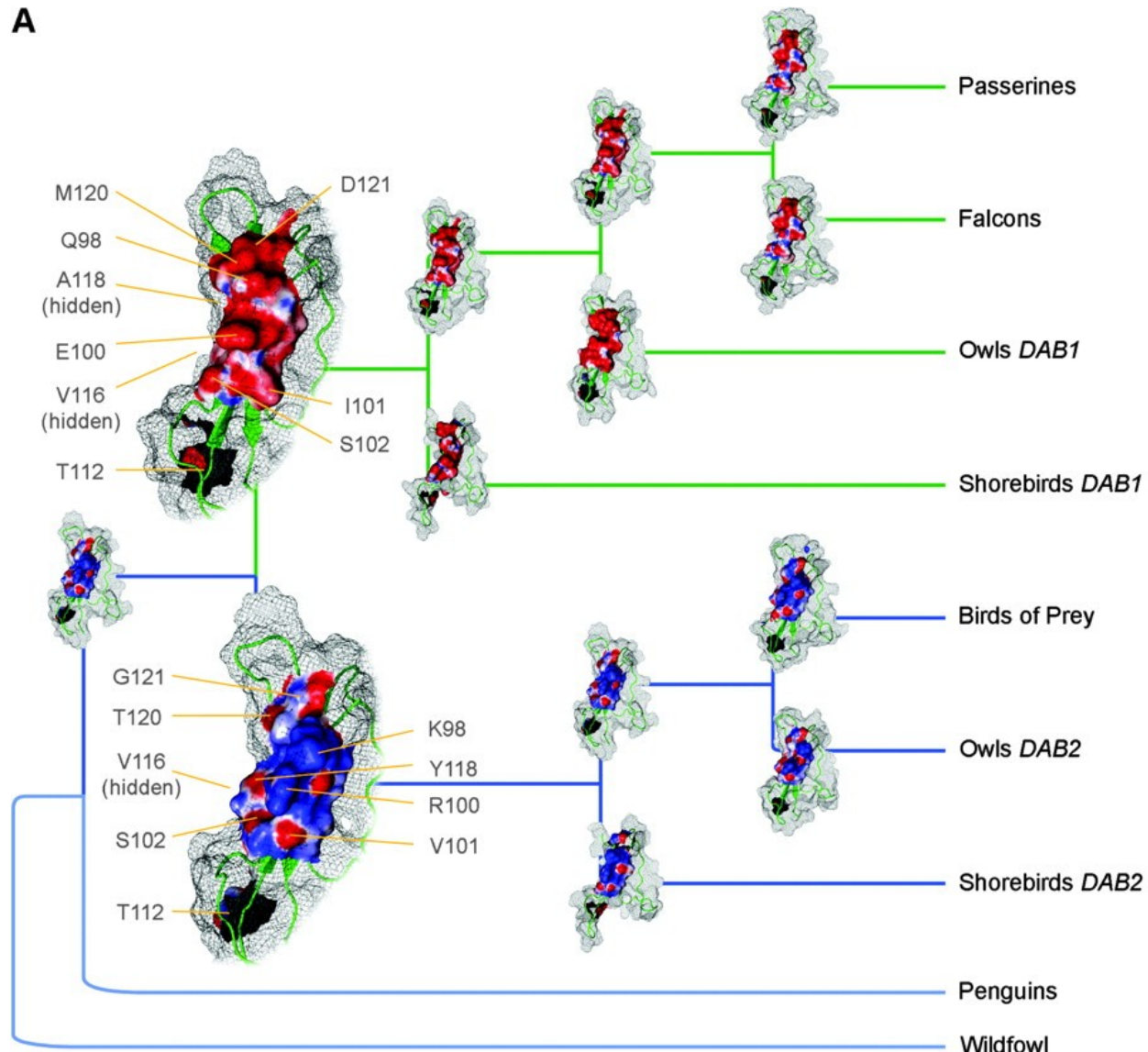
(Simulated alignment)

Ancestral Sequence Reconstruction



Ancestral Sequence Reconstruction

A

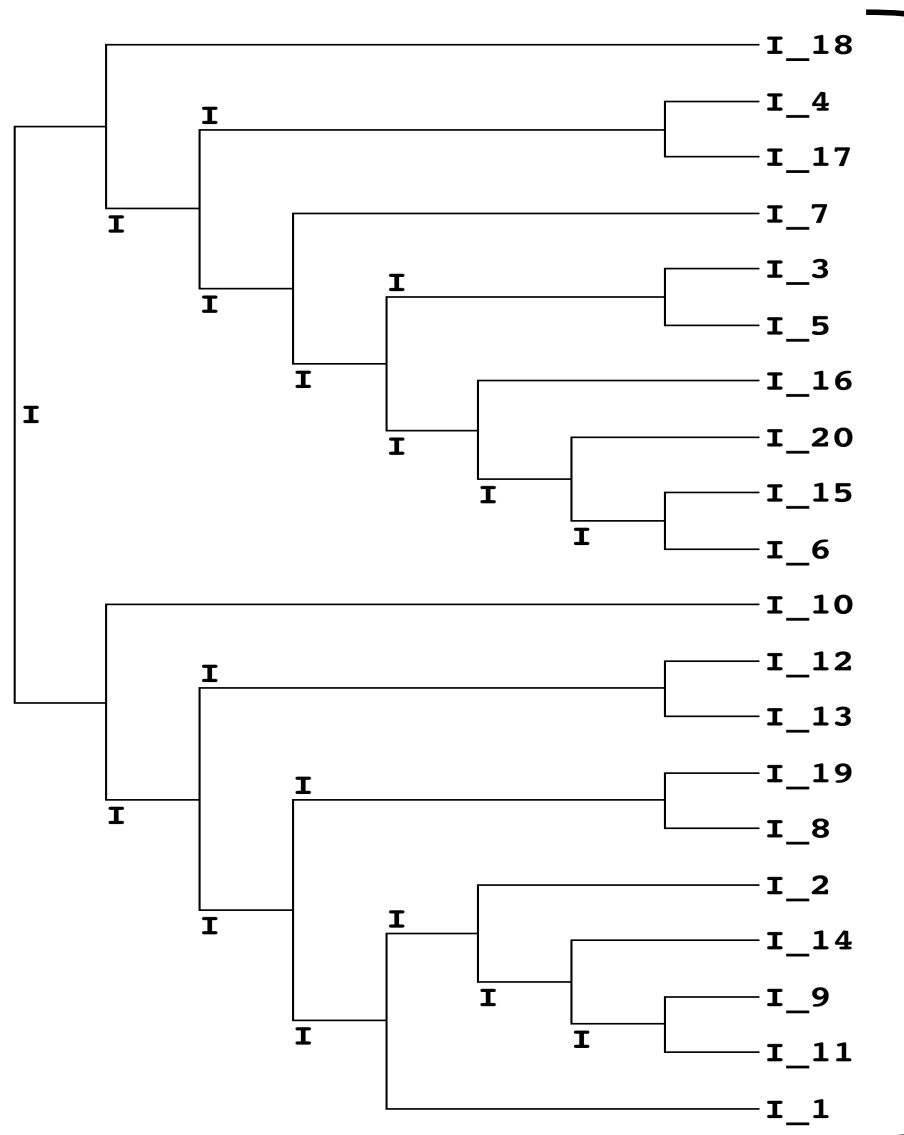


Why detecting events of adaptation in protein?

- Lot of mutations in proteins. Some are important, other are not.
- => How to find those that are important?

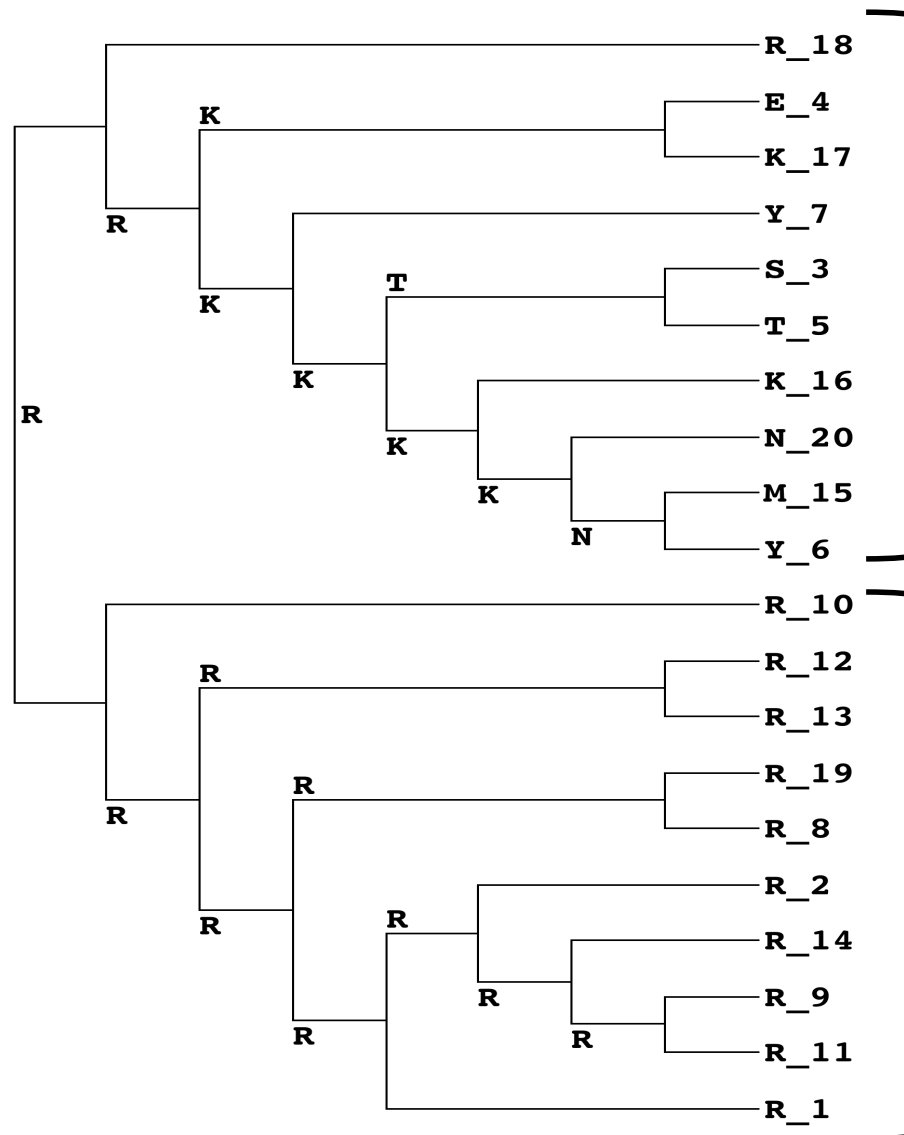
[illegible]

Amino acids: detecting sites with different evolutionary pattern



Strictly conserved
for Isoleucine

Amino acids: detecting sites with different evolutionary pattern

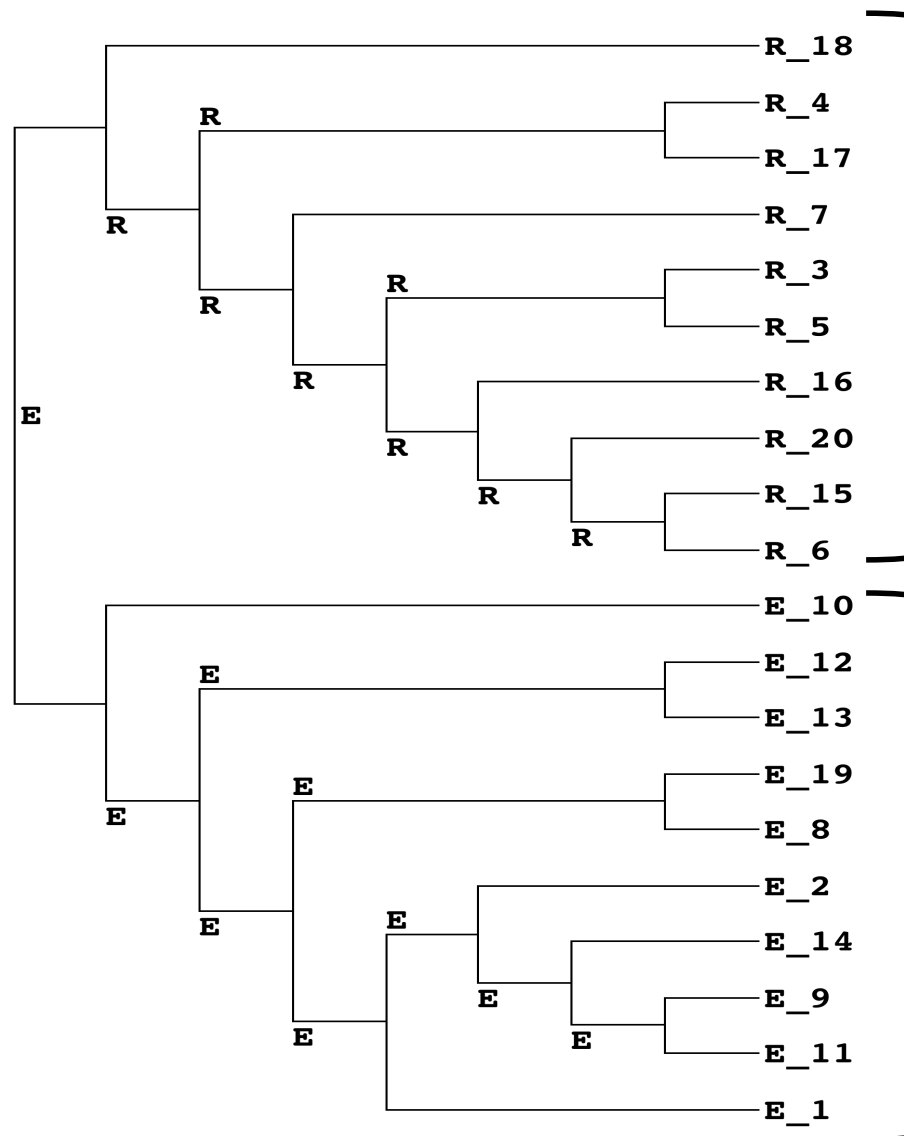


Highly divergent

- Type I of functional divergence
- Heterotachy
- Covarion-like

Highly conserved

Amino acids: detecting sites with different evolutionary pattern

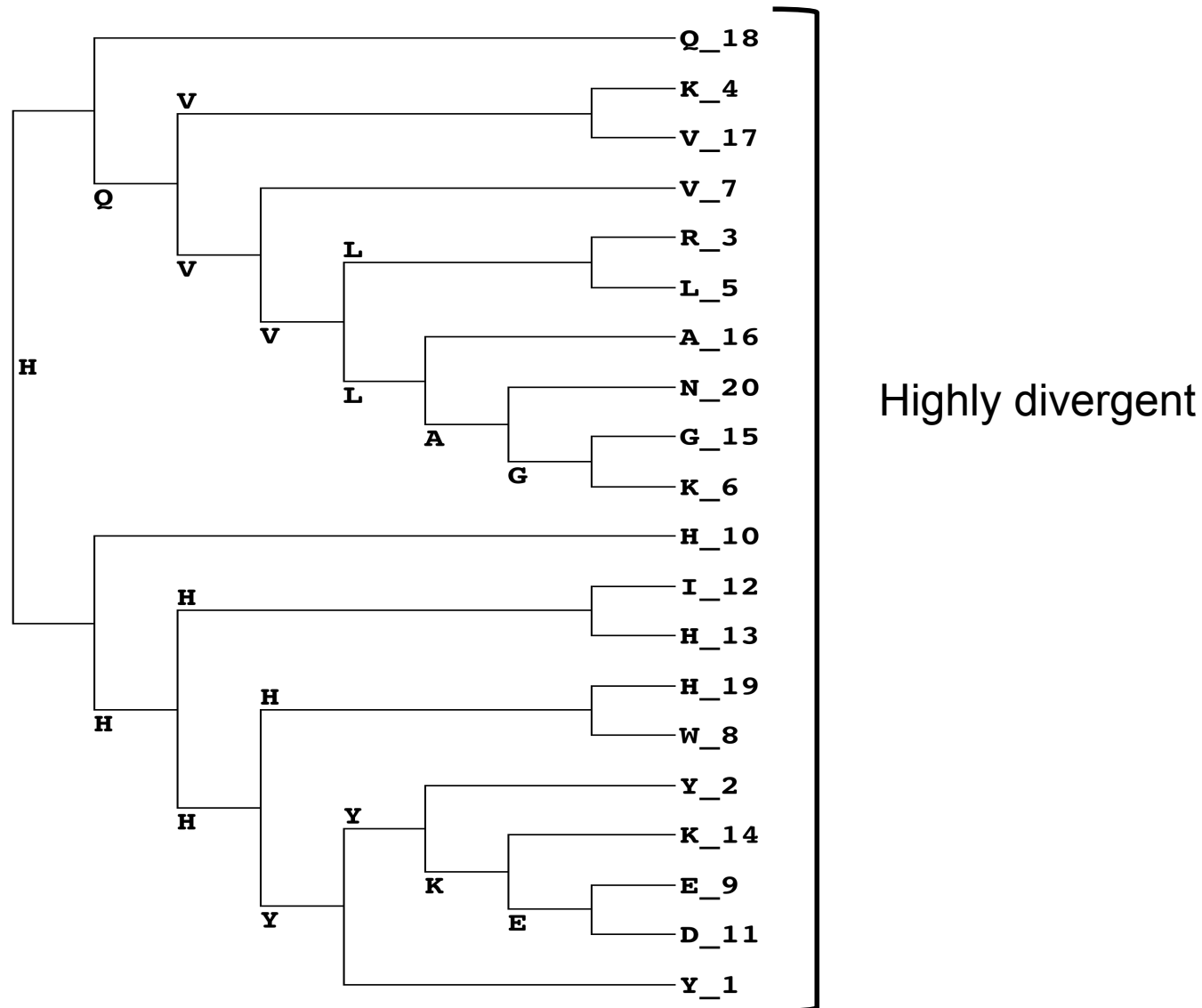


Conserved for basic residues

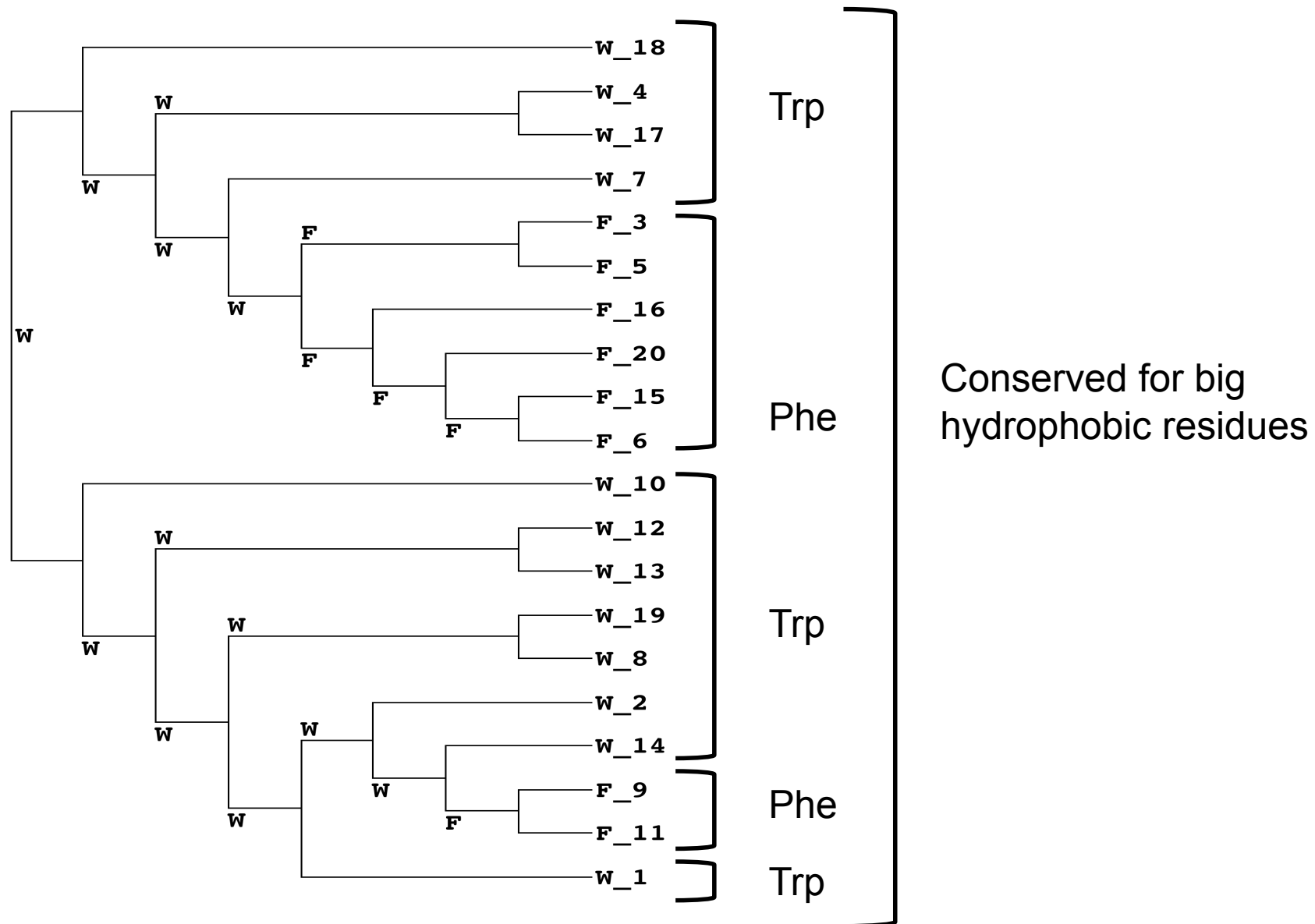
- Type II of functional divergence
- Constant-but-different

Conserved for acidic residues

Amino acids: detecting sites with different evolutionary pattern



Amino acids: detecting sites with different evolutionary pattern



Summary of changes

Group 1	Group 2	Type	Comment
Slow evolution (highly conserved)	Slow evolution (highly conserved)	Type 0	Conserved for structural and/or biochemical properties
Fast evolution (highly divergent)	Fast evolution (highly divergent)	Type 0	Either not important, or adapted in each species.
Slow evolution (highly conserved)	Fast evolution (highly divergent)	Type I	Change in evolutionary rate. Either a residue is recruited for a new function, or a position presents relaxation in selection pressures.
Fast evolution (highly divergent)	Slow evolution (highly conserved)	Type I	
Slow evolution (highly conserved)	Slow evolution, but another type (highly conserved)	Type II	A residue is recruited for a new function.

Fixation of mutations in the genome

- Have a deleterious effect on the fitness:
 - Not likely to be fixed.
 - => **Negative (purifying) selection** (Darwin theory)
- Have a positive effect on the fitness:
 - Likely to be fixed.
 - => **Positive selection** (Darwin theory)
- Have no effect on the fitness:
 - May or may not be fixed (random process)
 - => **Neutral evolution** (Kimura theory, early 70s)

Fixation of mutation in the genome

- Difficulty:

Differentiate between

NEUTRAL EVOLUTION (Most cases)

and

POSITIVE SELECTION (A few cases)

Statistical measure of selection / adaptation

- Null hypothesis:
 - All amino acids / nucleotides are fixed under neutral evolution.
- Alternative hypothesis:
 - A subset of amino acids / nucleotides are fixed under positive selection.
- Test: Do we have more amino acids under positive selection than expected?

Amino acids: detecting sites with different evolutionary pattern

- Various tools exist:
 - **DIVERGE** (Gu *et al.*, Bioinformatics 2002): Graphical User Interface. Works only in Windows.
 - **BADASP** (Edwards and Shields, Bioinformatics 2005): Easy to use. Doesn't provide statistical measure.
 - **TDG09** (Tamuri *et al.* Plos Comp Biol 2009): Provide likelihood inference per site. Can be used on convergent data.
 - **FUNDI** (Gaston *et al.*, Bioinformatics. 2011): Mix Type I and Type II. Use phylogenetic tool to estimate parameters.

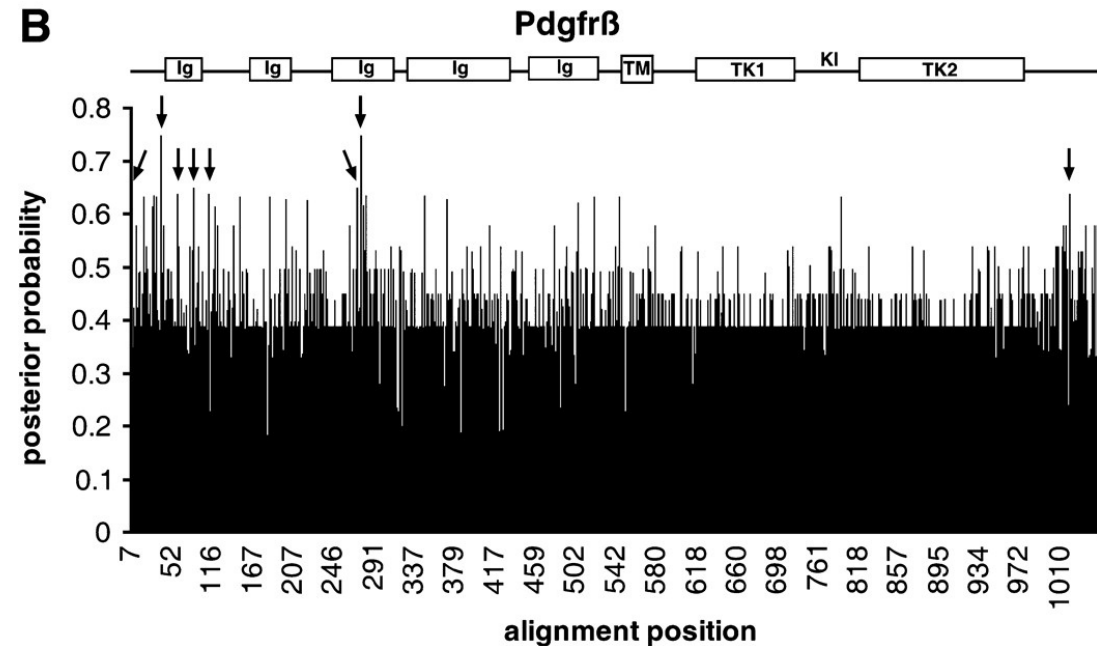
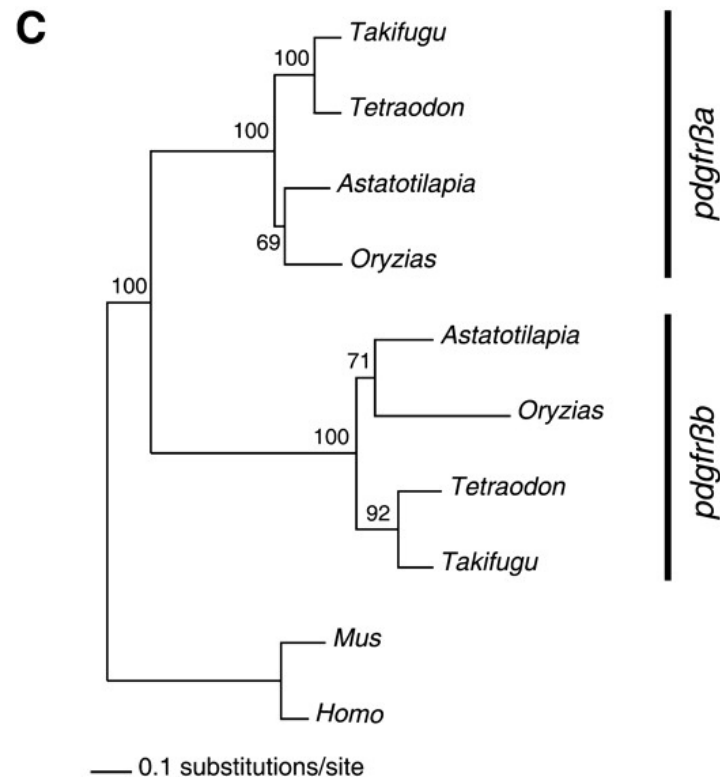
Other tools:

- SPEER (Chakrabarti *et al.*, Mol Biol. 2007)
- CHECKCOV (Pupko and Galtier, Proc Biol Sci. 2002)
- SHIFT-FINDER (Pontarotti's lab, unpublished)

Parameters to be estimated

- Alpha shape parameter of the gamma distribution
- Amino acid frequencies
- Branch size

Asymmetric Evolution in Two Fish-Specifically Duplicated Receptor Tyrosine Kinase Paralogons Involved in Teleost Coloration



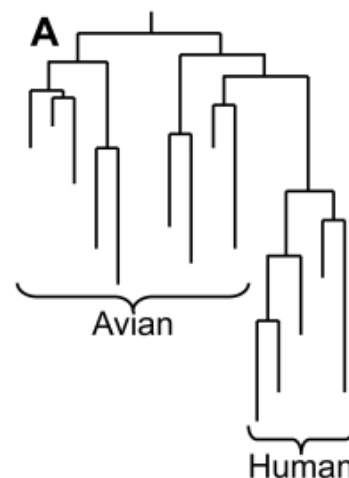
DIVERGE, used to detect Type I of functional divergence

R
R
R
R
E
F
W
C

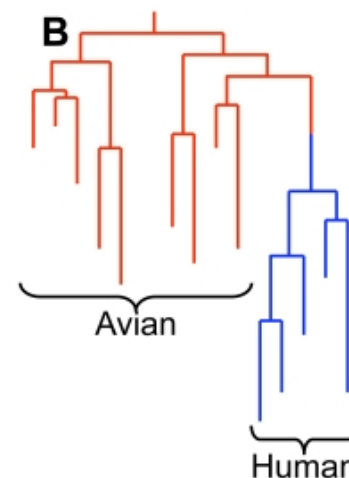
TDG09 Algorithm: Amino Acids

- Identifying changes in selective constraints: host shifts in influenza.
Tamuri AU, Dos Reis M, Hay AJ, Goldstein RA.
PLoS Comput Biol. 2009 Nov;5(11):e1000564.
- Sitewise non-homogeneous phylogenetic model that explicitly takes into account **differences in the equilibrium frequencies of amino acids** in different hosts and locations.

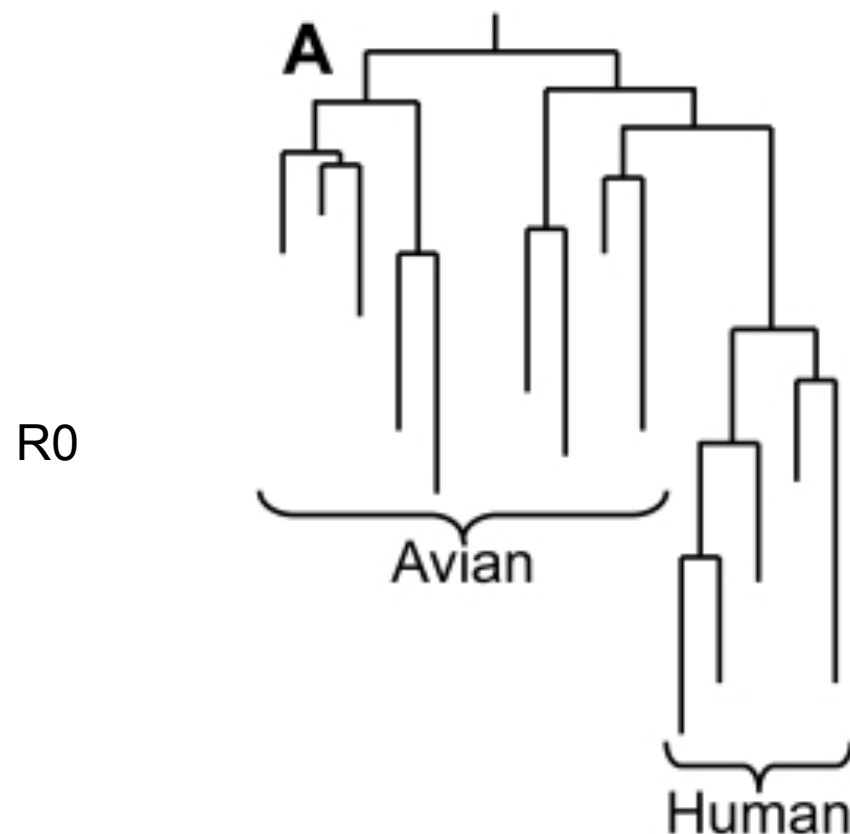
Homogeneous assumption



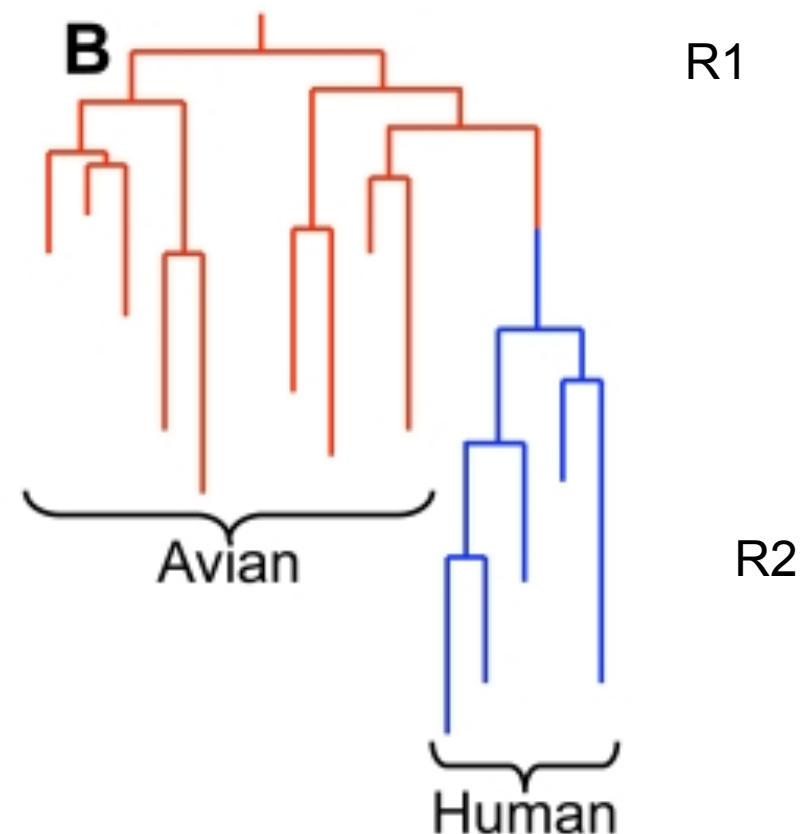
Non-homogeneous assumption



TDG09 Algorithm: Amino Acids



Substitution rate
=
identical throughout the tree



Substitution rate
=
different throughout the tree

Practical:

<http://beta.cathdb.info>

⇒ About

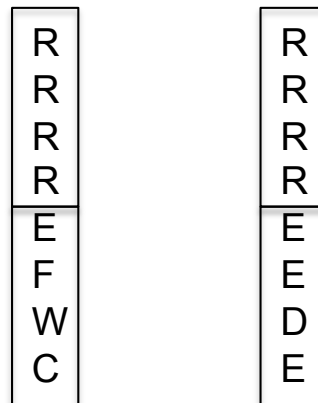
⇒ Tutorials

⇒ Detecting sites under functional divergence

⇒ Part I on amino acids dataset

- Protein: nucleotide to amino acid
- Measuring positive selection using the dN/dS ratio.

- Type I
- Type II



- dN/dS helps to identify the constraints.

Protein: nucleotide to amino acid

Inverse table (compressed using IUPAC notation)

AA	Codons
Ala/A	GCU, GCC, GCA, GCG
Arg/R	CGU, CGC, CGA, CGG, AGA, AGG
Asn/N	AAU, AAC
Asp/D	GAU, GAC
Cys/C	UGU, UGC
Gln/Q	CAA, CAG
Glu/E	GAA, GAG
Gly/G	GGU, GGC, GGA, GGG
His/H	CAU, CAC
Ile/I	AUU, AUC, AUA
START	AUG

AA	Codons
Leu/L	UUA, UUG, CUU, CUC, CUA, CUG
Lys/K	AAA, AAG
Met/M	AUG
Phe/F	UUU, UUC
Pro/P	CCU, CCC, CCA, CCG
Ser/S	UCU, UCC, UCA, UCG, AGU, AGC
Thr/T	ACU, ACC, ACA, ACG
Trp/W	UGG
Tyr/Y	UAU, UAC
Val/V	GUU, GUC, GUA, GUG
STOP	UAA, UGA, UAG

How to detect Positive Selection?

- Make use of the degenerate genetic code.
- Assume that dS (synonymous) substitutions are neutral.
- Assume that dN (non-synonymous) substitution are either neutral (fixed under random drift) or fixed under positive selection.

How to detect Positive Selection?

Species 1	A	L	P	H	Y
	GCC	C <u>T</u>	CCT	CAT	TAT <u>T</u>
Species 2	A	R	P	H	Y
	GCC	C <u>G</u> T	CCT	CAT	TAC <u>C</u>

How to detect Positive Selection?

Species 1	A	L	P	H	Y
	GCC	CTT	CCT	CAT	TAT
Species 2	A	R	P	H	Y
	GCC	CGT	CCT	CAT	TAC

15 nucleotides: 11 non-synonymous sites and 4 synonymous sites
 1 synonymous substitutions (S) and 1 non-synonymous substitution (N)

$$dN = \frac{\text{number of non-synonymous substitutions}}{\text{non-synonymous sites}} = 1 / 11 = 0.09$$

$$dS = \frac{\text{number of synonymous substitutions}}{\text{synonymous sites}} = 1 / 4 = 0.25$$

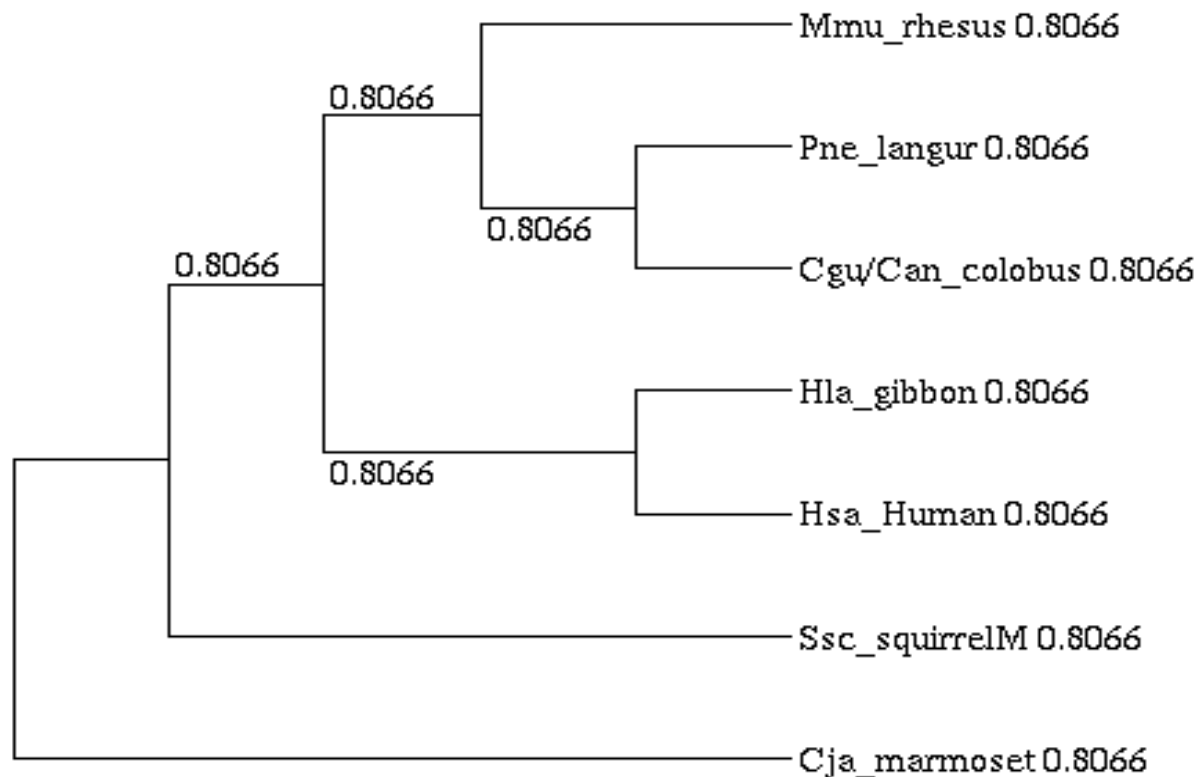
$$dN/dS (\omega) = 0.09 / 0.25 = 0.36$$

$dN/dS < 1$ -> purifying selection
 $dN/dS = 1$ -> neutral evolution
 $dN/dS > 1$ -> positive selection !!!!!!!!

Branch models

(Yang 1998; Yang and Nielsen 1998)

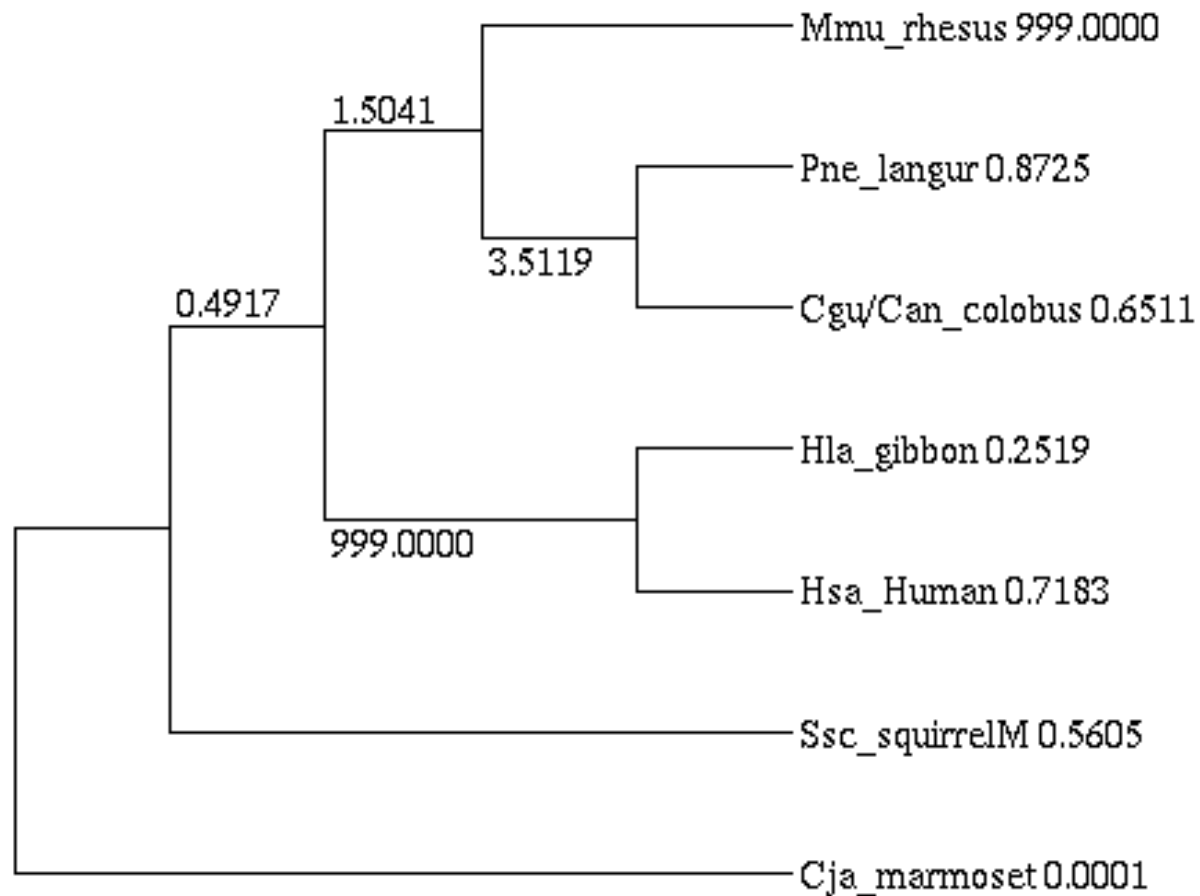
- One ratio model: estimate different dN and different dS, but keep the same ratio on all branches: null model
- All branches on the tree are under the same degree of selective pressure (i.e. same constraints).



Branch models

(Yang 1998; Yang and Nielsen 1998)

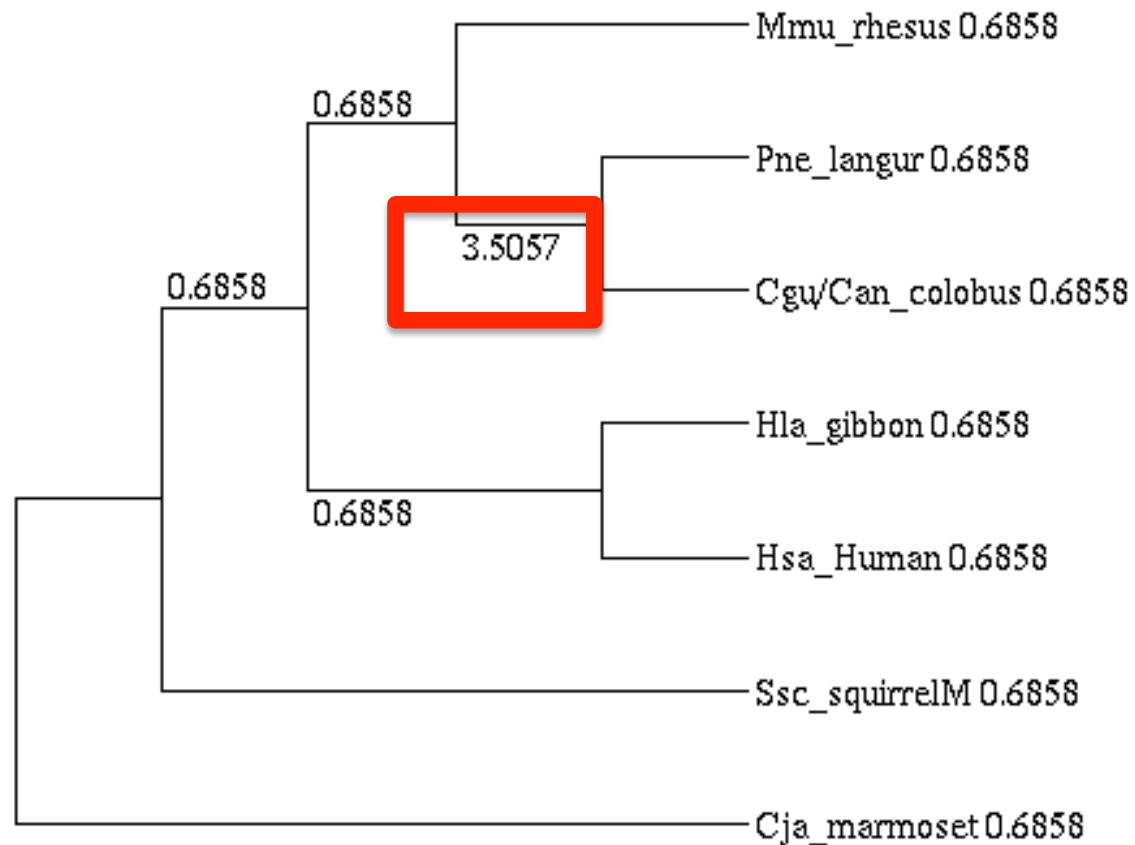
- Estimate all dN/dS ratio for each branch on the phylogeny



Branch models

(Yang 1998; Yang and Nielsen 1998)

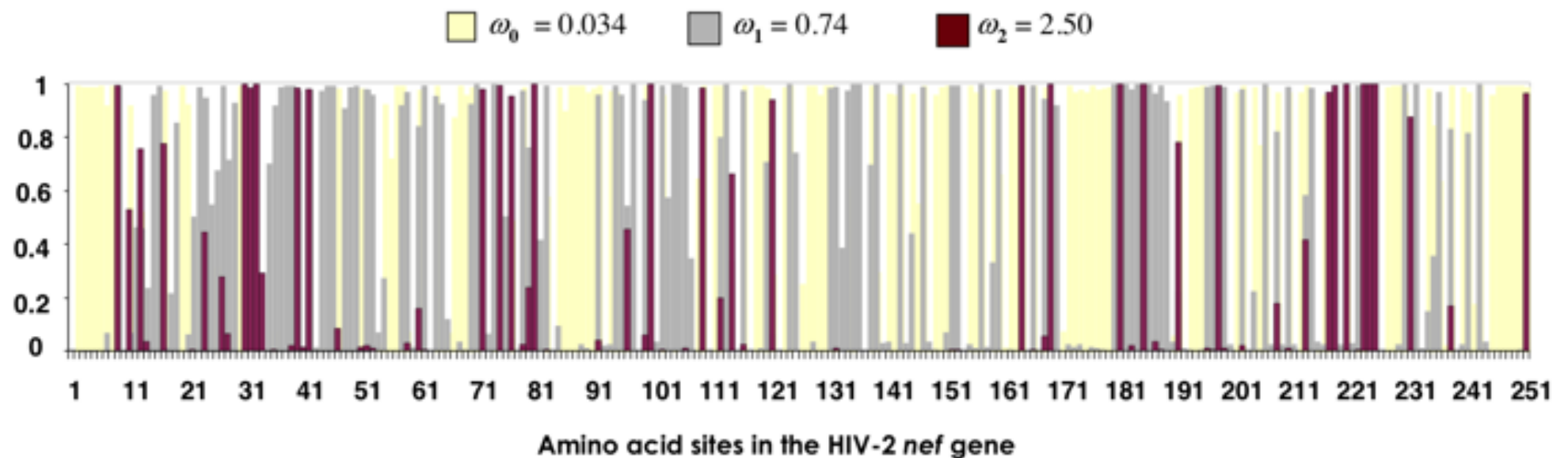
- Estimate different dN/dS ratio on one or more branches



Site models

(Nielsen and Yang 1998; Yang et al. 2000, Yang et al. 2005)

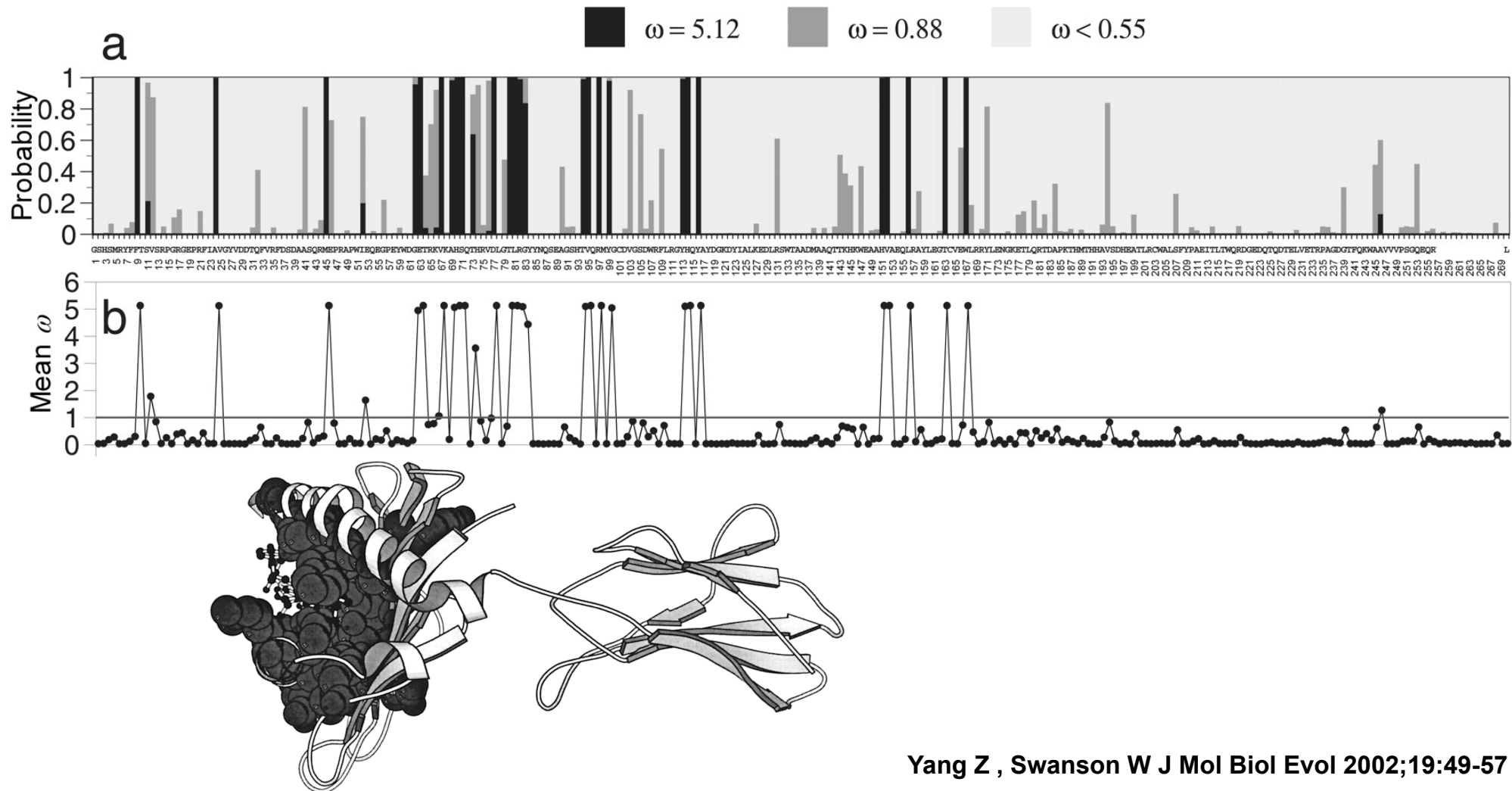
- Estimate all dN/dS (ω) at each site (codons) into three categories (or more):
 ω_0 (yellow) = negative selection
 ω_1 (grey) = neutral evolution (or nearly-neutral)
 ω_2 (rouge) = positive selection



Site models

(Nielsen and Yang 1998; Yang et al. 2000, Yang et al. 2005)

Posterior probabilities of site classes for sites along the MHC class I gene under the random-sites model M8

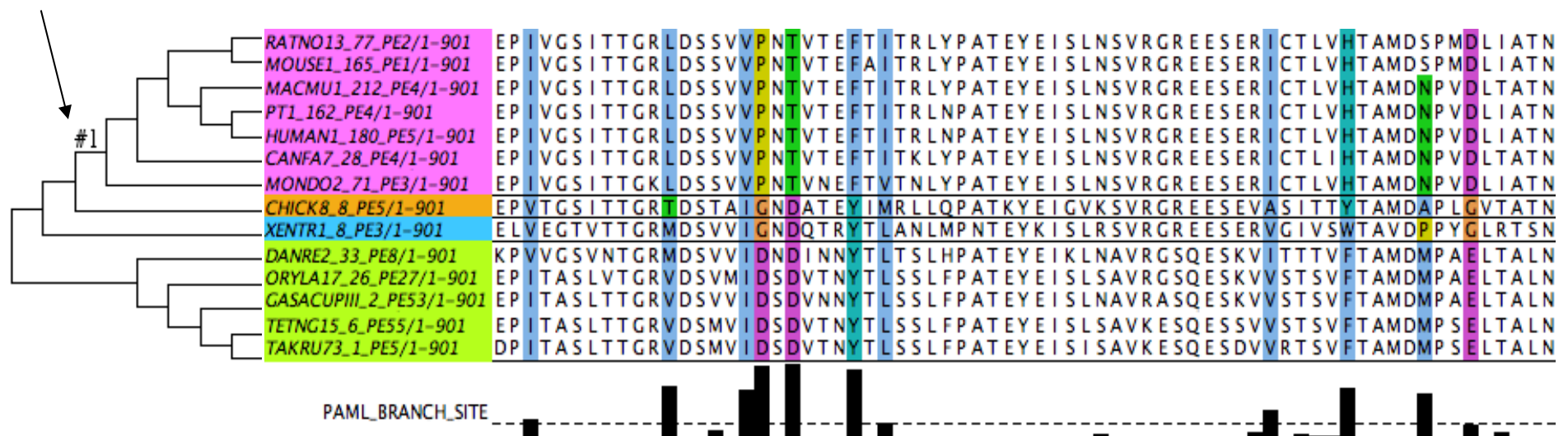


Branch-site models

(Yang and Nielsen 2002, Zhang *et al.* 2005)

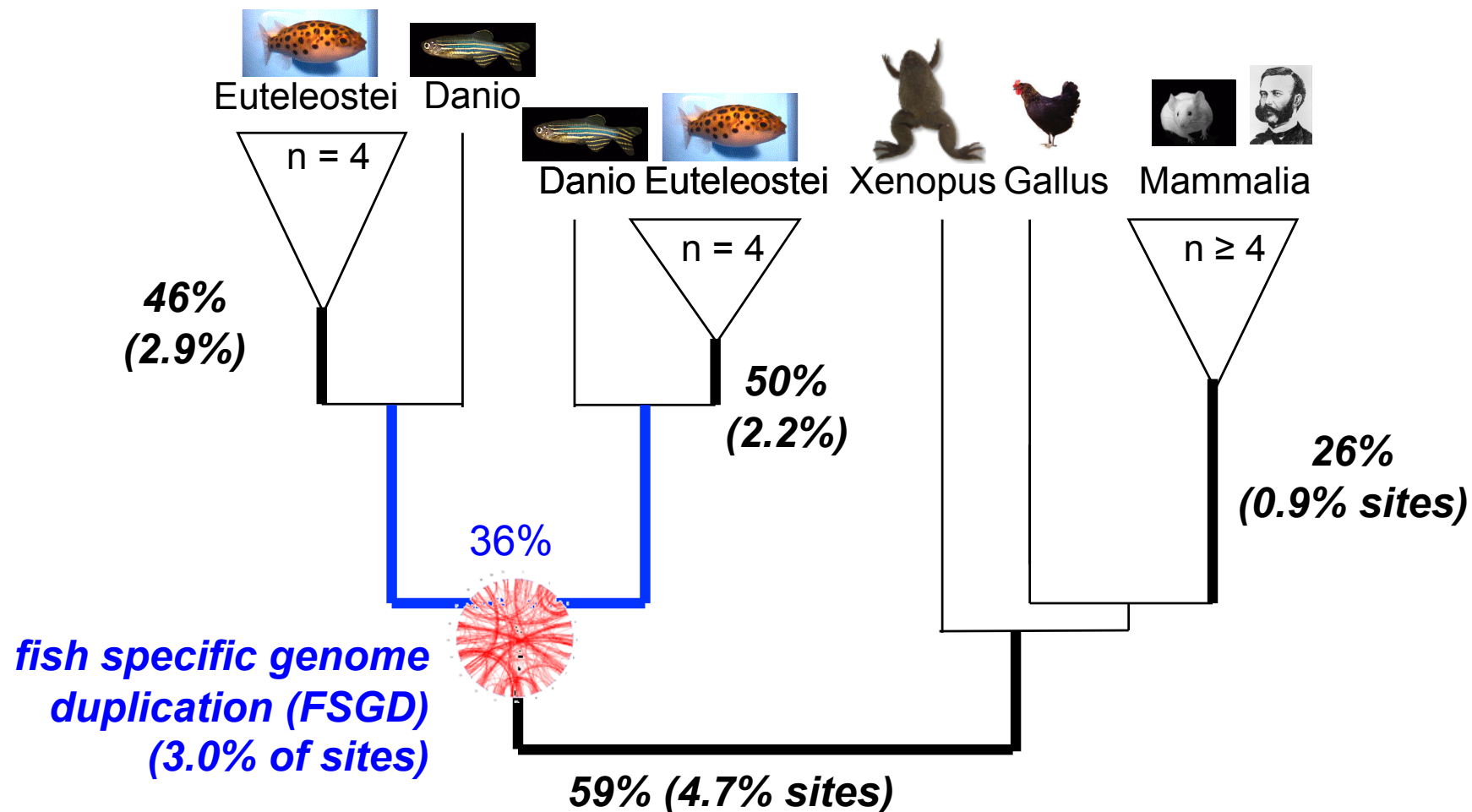
- Detect positive selection that affects only a **few sites** on pre-specified lineages.
- **foreground branches** = branches under test for positive selection.
- **background branches** = all other branches.
- LRT: branch-site model A is the alternative model ($\omega_2 > 1$), while the simpler null model is **model A but with $\omega_2 = 1$ fixed**.

Branch leading
to Mammals



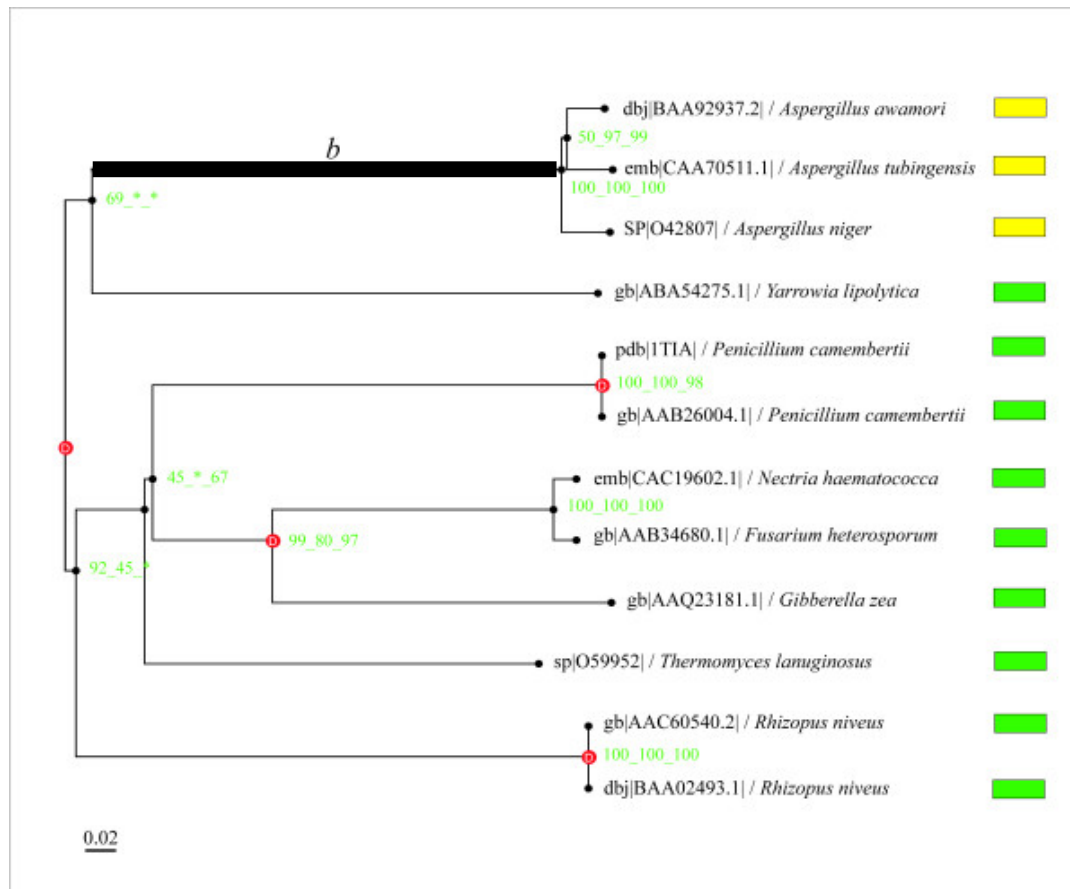
We can estimate ancient dN/dS.

- Analyses done on 884 genes families (2'673 branches).
- Percentages of branches with sites under positive selection:

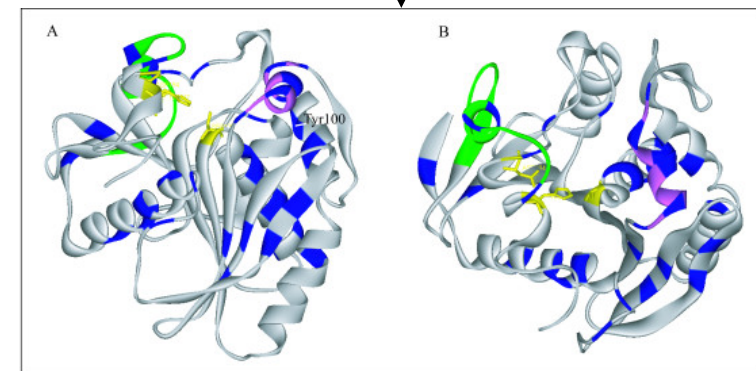


Example of positive selection in enzyme

Lipase  / feruloyl esterase A 

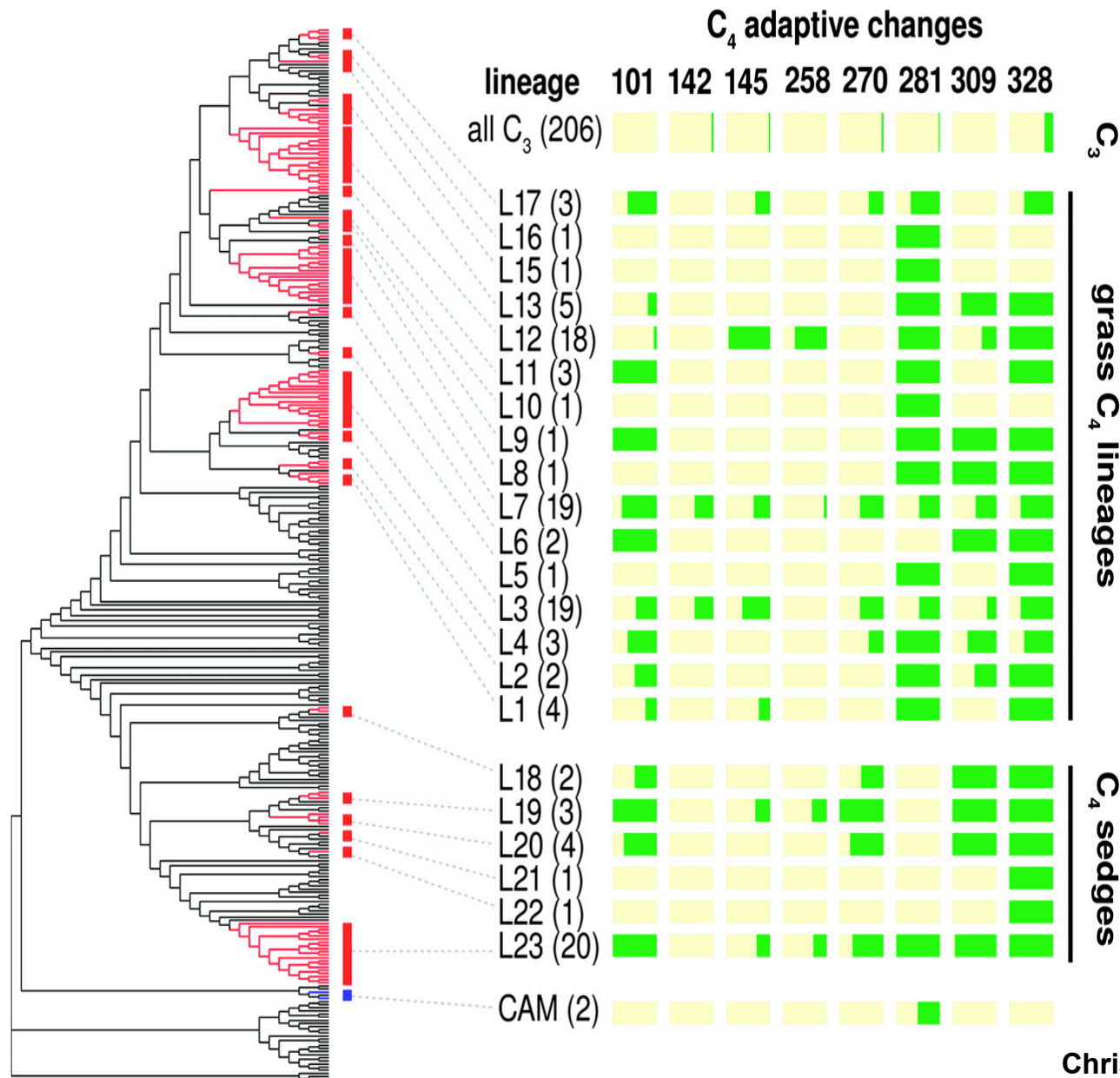


SITES IN 3D???



Positively selected sites in blue

Positive selection detected in RubisCO



C4 HAS A HIGHER
TURN-OVER
THAN C3.

Databases of positive selection

Selectome: database of positive selection

<http://selectome.unil.ch/>

NVTTEITA
NVTNEVTA
Selectome
NADTEMIA
NQDTRLTA

BETA Selectome: a Database of Positive Selection


UNIL | Université de Lausanne

Selectome © 2008/2012

[selectome \[AT\] unil.ch](http://selectome.unil.ch)

Search

[example](#)

Basic search:

Go!

Clear

Advanced search

About

[Quick guide](#)

[Methods](#)

[Downloads](#)

[How to cite ?](#)

[Contact](#)

Welcome on BETA Selectome home page

Selectome is a database of positive selection based on a rigorous branch-site specific likelihood test. Positive selection is detected using [CODEML](#) on all branches of animal gene trees. The web interface of Selectome enables queries according both to the results of positive selection tests, and to gene related criteria. Test results including positively selected sites can be visualized on the tree, and on the protein sequence alignment.

We are currently improving our source data selection as well as the filtering steps on multiple sequence alignments in order to get more reliable results.

The interface uses the multiple alignment viewer applet [Jalview](#).

The current data are built from the [Ensembl database](#) version 61.


Last News

2012-02-15	Add extra annotations in trees; Add new xrefs; various bug fixes
2011-10-24	Selectome release 05, based on Ensembl 61, <i>Primates</i> taxa
2010-08-18	Full Selectome release, based on TreeFam7 A & B

[See all news](#)

Selectome: database of positive selection

<http://selectome.unil.ch/>



[selectome \[AT\] unil.ch](#)

Search

[example](#)

Basic search:

Go! Clear

Advanced search

About

[Quick guide](#)
[Methods](#)
[Downloads](#)
[How to cite ?](#)
[Contact](#)

Advanced search

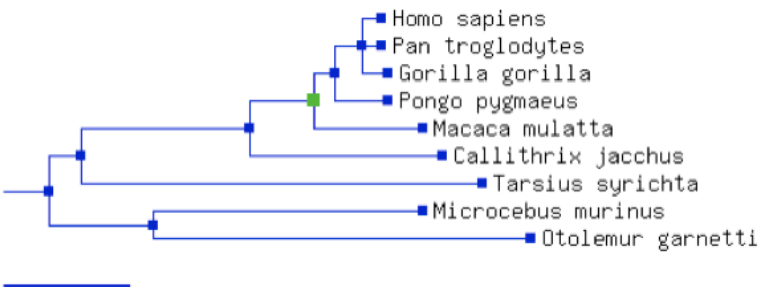
Search by terms
Search by branch
Tips

Search for a branch

Choose a branch

Taxon Primates

Selected taxon : Catarrhini



0.04303

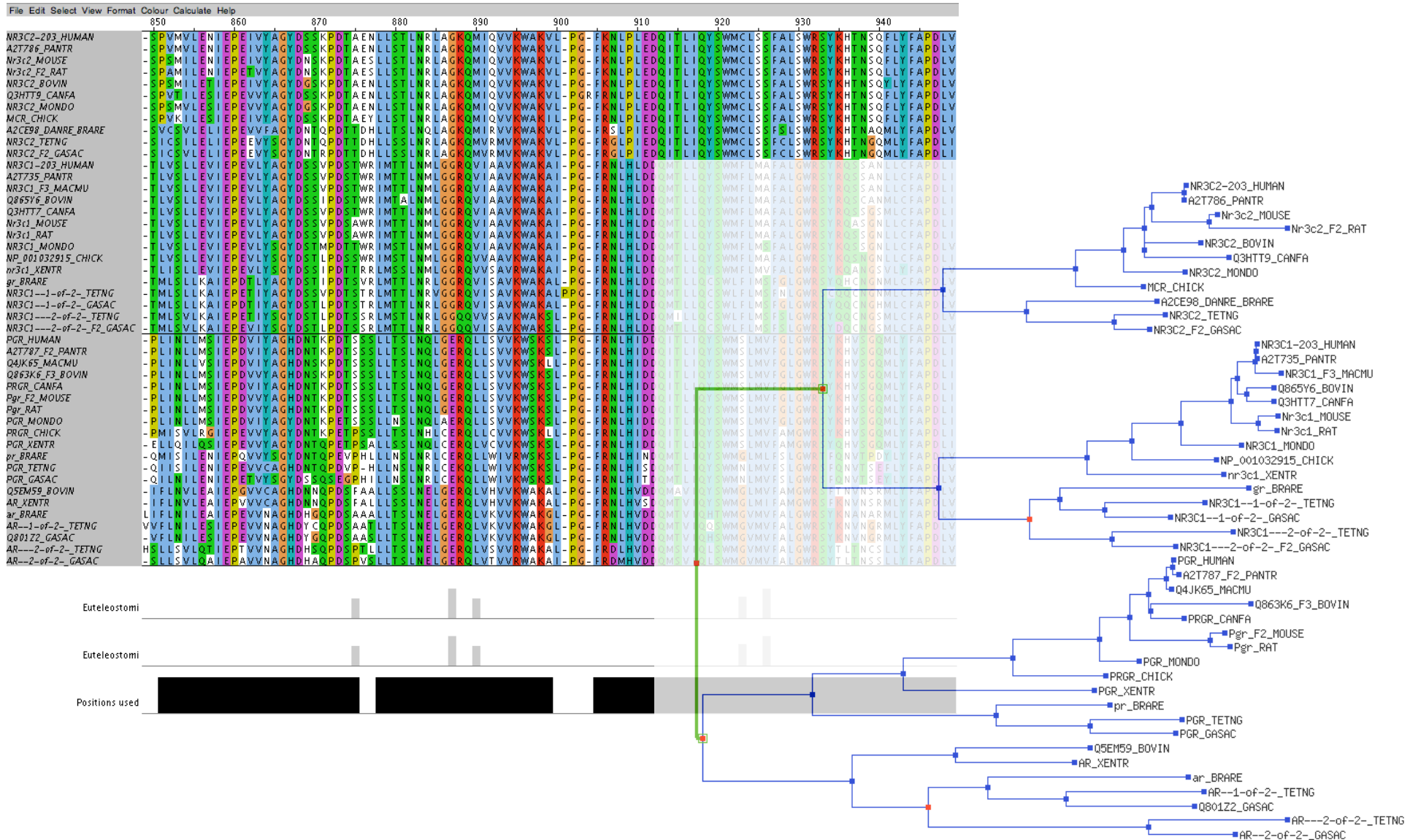
Positive selection : ☒ With ☐ Without ☐ I don't care

Submit

Not sure about the taxonomy? Visit the [NCBI Entrez Taxonomy](#) website.

Selectome: database of positive selection

<http://selectome.unil.ch/>



The Adaptive Evolution Database (TAED)

<http://www.wyomingbioinformatics.org/TAED/>

 UNIVERSITY OF WYOMING

[UW Home](#) | [About UW](#) | [Apply](#) | [A-Z Directory](#) | [Phone/E-mail](#) | [Search UW](#)

www.wyomingbioinformatics.org



[TAED Introduction](#)

[Tree Thrasher](#)

[TAED Search](#)

[Tree of Life](#)

[Alphabetical gene Search](#)

The Adaptive Evolution Database

TAED is a database of phylogenetically indexed gene families. It contains multiple sequence alignments from MAFFT¹, maximum likelihood phylogenetic trees from PhyML², bootstrap values for each node, dN/dS ratios for each lineage from the free ratios model in PAML³, and labels for each node of speciation or duplication from gene tree/species tree reconciliation using SoftParsMap⁴. The phylogenetic indexing enables simultaneous viewing of lineages with high dN/dS that occurred along the same species tree branches.

The current status of the database includes 8,060 gene families. In addition to making the complete resource available, future updates will also include cross referencing with PDB and with KEGG⁵.

The database can be entered through the species tree, where the species tree branch links to underlying gene families, either all gene families where that species tree lineage is represented, or the subset with dN/dS > 1. It can also be entered by searching for genes of interest, or through an alphabetical list of gene family annotations. It is ultimately useful in identifying candidates to answer the question, "What makes this species unique?" or which genes show signals for diversification under which lineages?

The Adaptive Evolution Database (TAED)

<http://www.wyomingbioinformatics.org/TAED/>



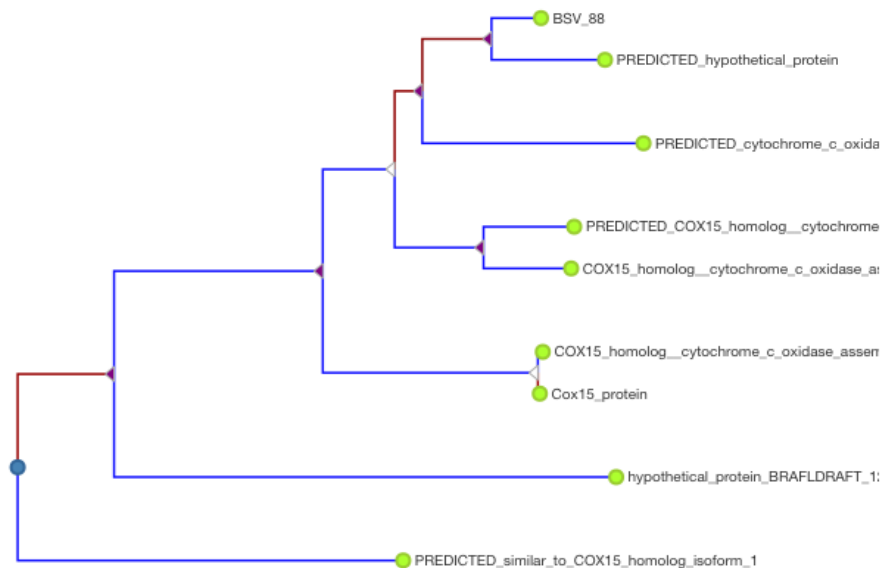
Branches
show
DN/DS:

- 0
- 0.25
- 0.5
- 0.75
- 1
- 1.5
- 2
- 4
- 8
- 16
- 25
- 50
- 100
- 200
- 300
- 400
- 500
- 600
- 700
- 800
- 900
- 999

COX15 homolog, cytochrome c oxidase assembly protein : 8259.nhx_rooted.isoforms

☐ branch lengths ☐ DN/DS ☐ species collapse beyond depth

☒ theRoot ☐ duplication ☐ speciation ☐ collapsed ☐ leaf



Click on the inner nodes to collapse or expand the tree one level at a time, double click (fast) to expand entire branches. Clicking on a leaf node will load the NCBI record (if any). Trees will not be visible in Internet Explorer 8 and earlier, if you do not see a tree in IE9 - try refreshing the page. Internet Explorer 9, Firefox, Safari, and Chrome should all render trees correctly.

This page developed by Benjamin Oswald at the University of Wyoming, using the D3 javascript package, and extending code by Jason Davies and kueda.

Some advices

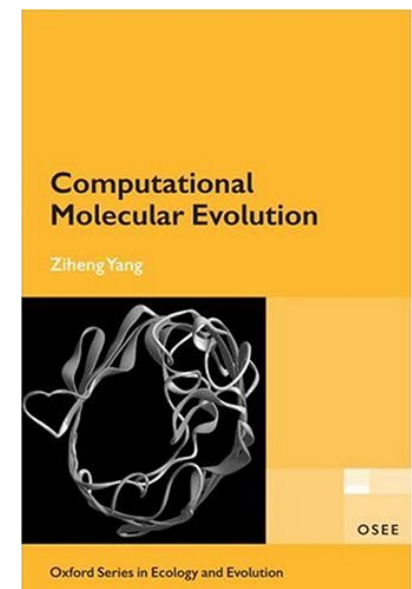
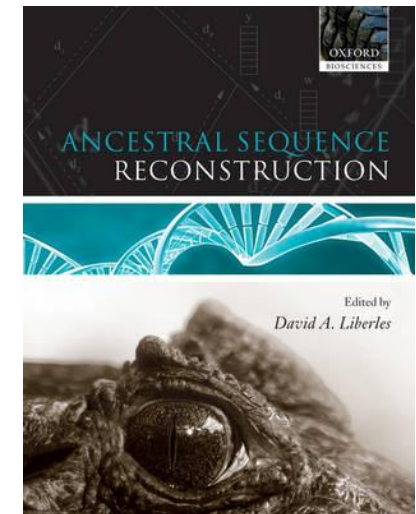
- Garbage in / garbage out
=> Sequences of good quality !
- Multiple alignment with best methods
- Remove badly aligned blocks
- Select most appropriate model of substitution (ie WAG or LG).
- Build trees with Maximum Likelihood / Bayesian.

Some advices

- Prior biological hypothesis:
 - ⇒ Apply test on the branch of interest.
- No prior hypothesis:
 - ⇒ Apply tests on different branches.
 - ⇒ Apply a False Discovery Rate (FDR) *a posteriori*.
- Always visualise your result!

Recommended books

- Ancestral Sequence Reconstruction,
by David A. Liberles
- Computational Molecular Evolution,
by Ziheng Yang



Practical:

<http://beta.cathdb.info>

⇒ About

⇒ Tutorials

⇒ Detecting sites under functional divergence

⇒ Part II on nucleotides dataset