

# Functional and Comparative Genomics

# Comparative and Functional Genomics

## Session 3

# **Phylogenetic analysis**

Gene family tree reconstruction.

Inference of gene duplication and other evolutionary events.

Detection of functional divergence.

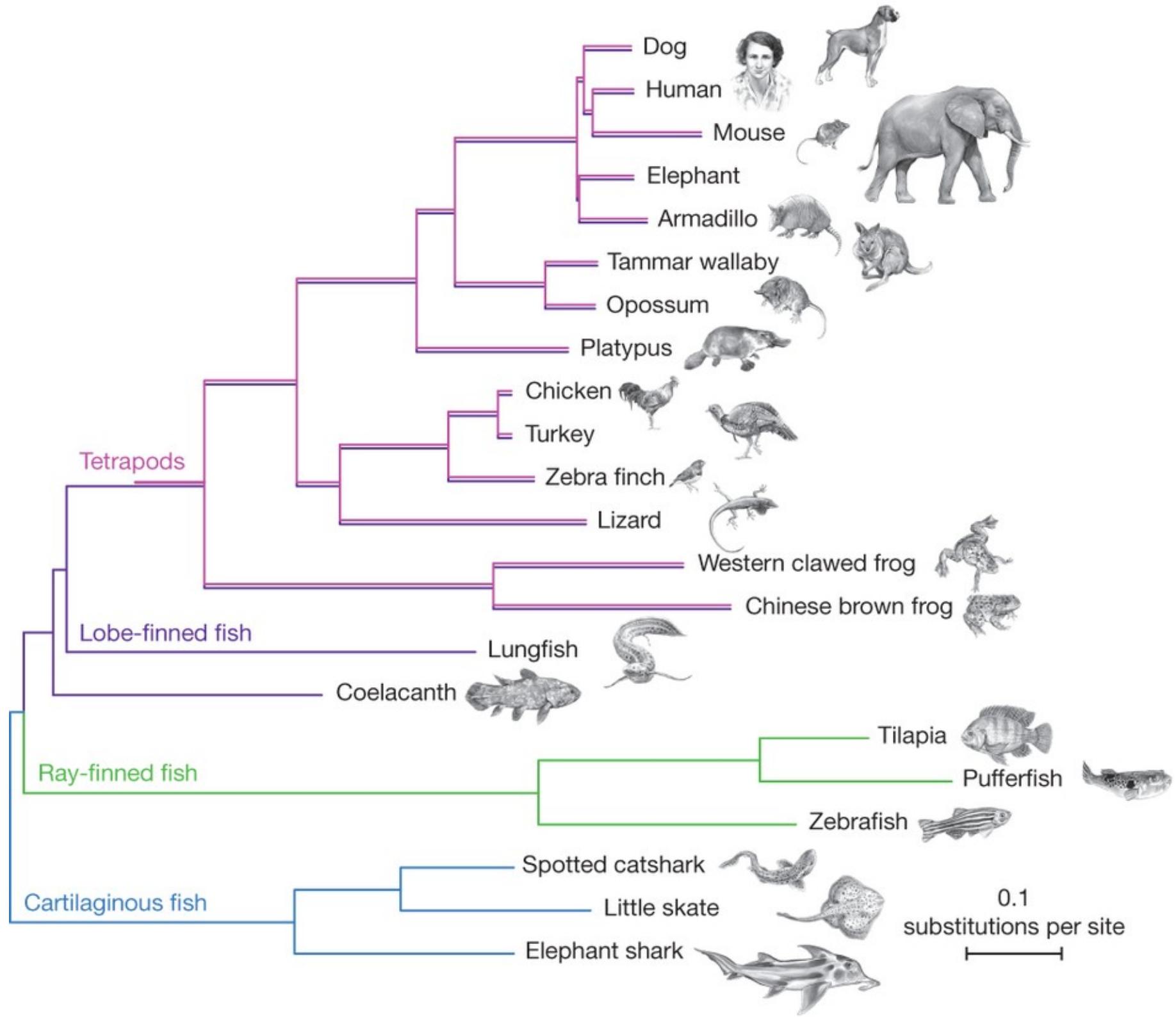
dn/ds analysis.

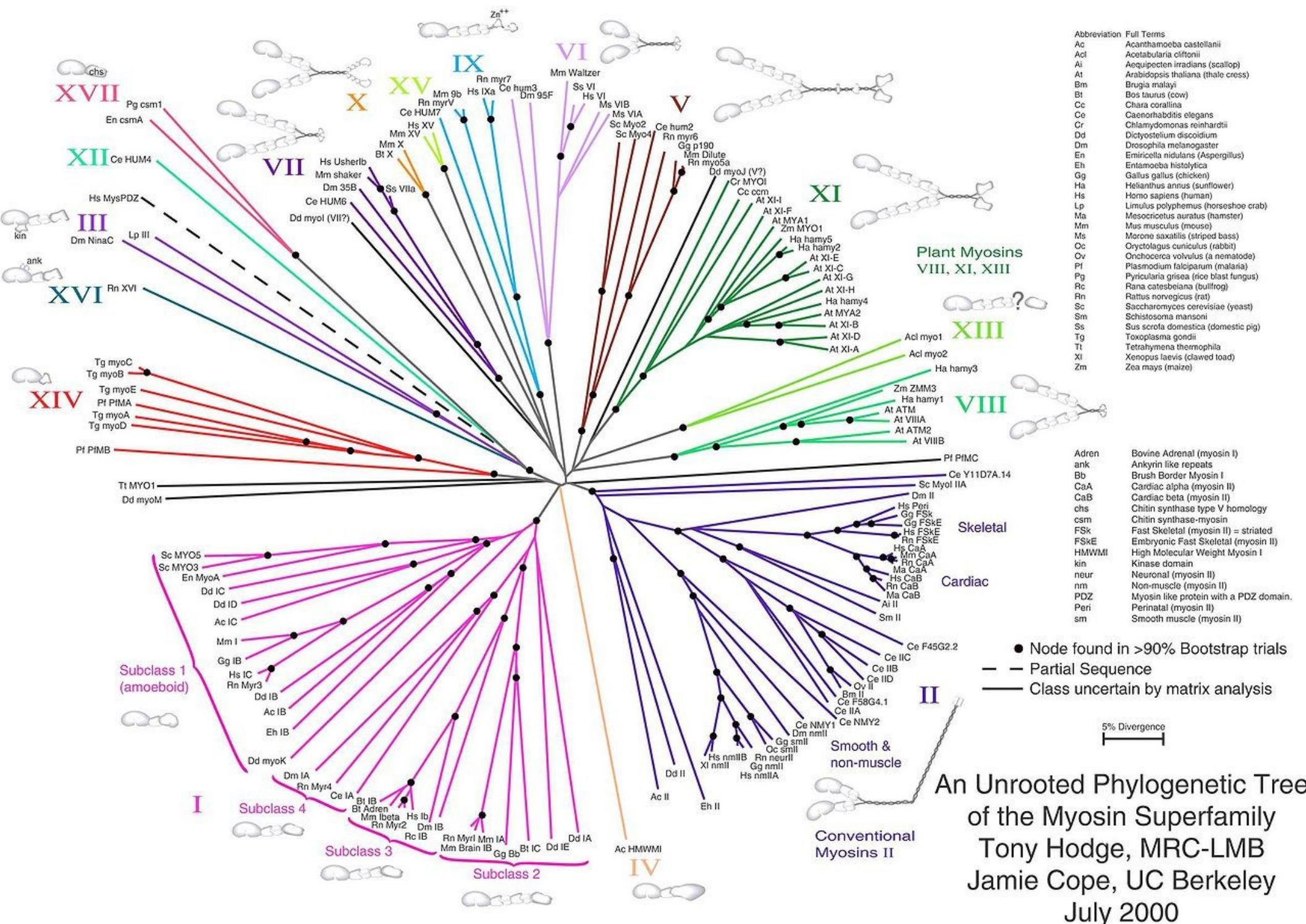
# A phylogenetic tree

A branching diagram (**bipartite graph**) showing the inferred **evolutionary relationships** among various **biological species or other entities** (e.g sequences) **based on similarities and differences** in their physical and/or genetic characteristics.

Species tree

Gene tree





# An Unrooted Phylogenetic Tree of the Myosin Superfamily

Tony Hodge, MRC-LMB  
Jamie Cope, UC Berkeley  
July 2000

# Why care about (phylogenetic) trees?



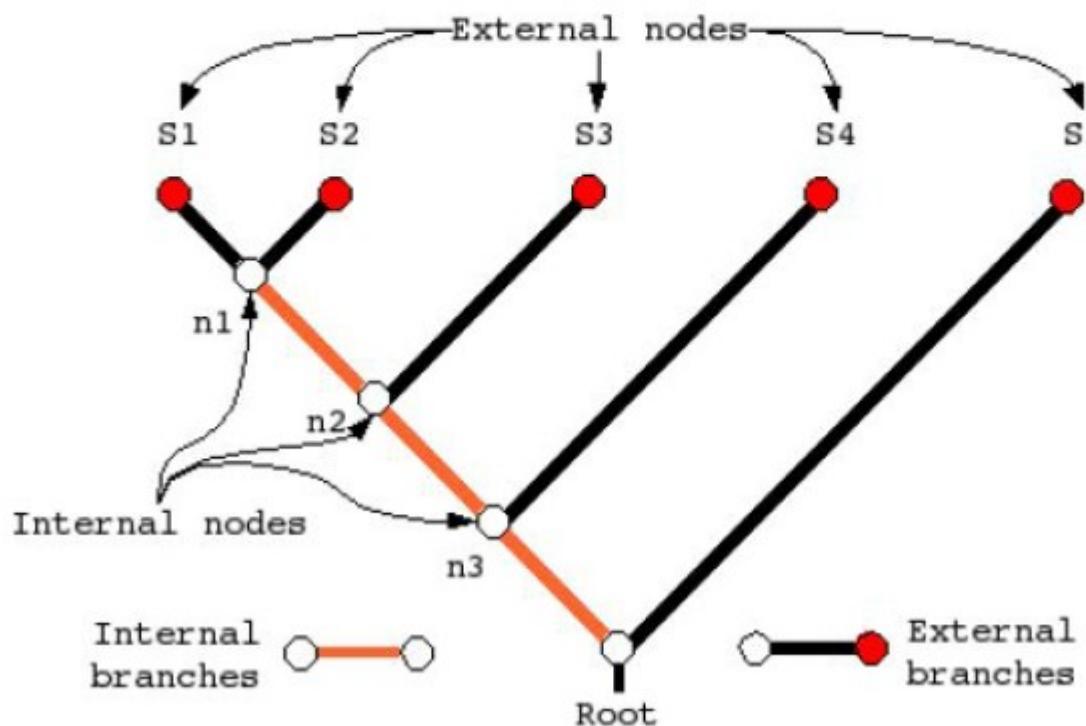
Biological systems are the result of the evolutionary process. Trees represent this process.

Assumptions on statistical tests usually assume “independence of data”, but no data in comparative biology is “independent”, as all organisms are related. Phylogenetic trees are the first step to detect and remove the effect of evolutionary relatedness.



“Nothing make sense in biology if not in the light of evolution”

- **Nodes & branches.** Trees contain internal and external nodes and branches. In molecular phylogenetics, **external nodes** are sequences representing **genes, populations or species!**. Sometimes, **internal nodes** contain the ancestral information of the clustered species. A **branch** defines the relationship between sequences in terms of descent and ancestry.

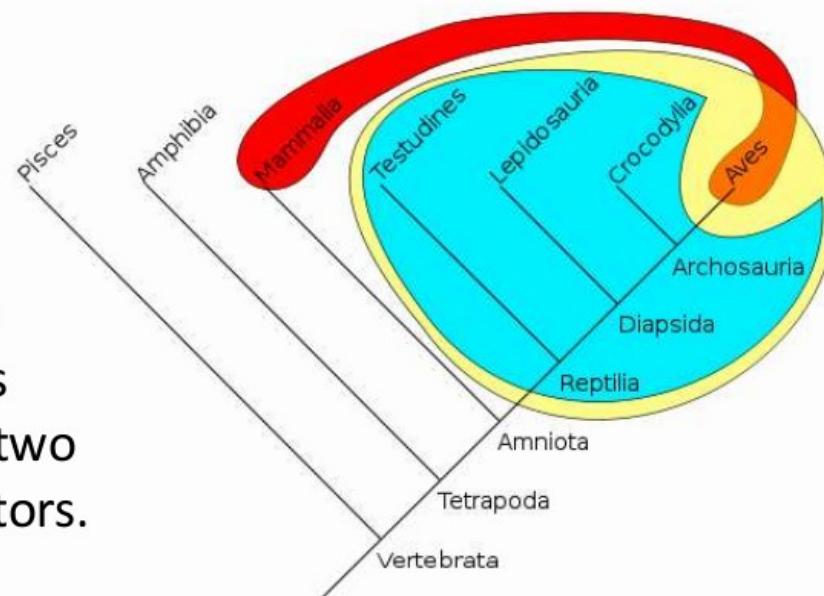


## Paraphyletic group:

modern reptiles (cyan)  
contain not all  
descendants from a  
common ancestor.

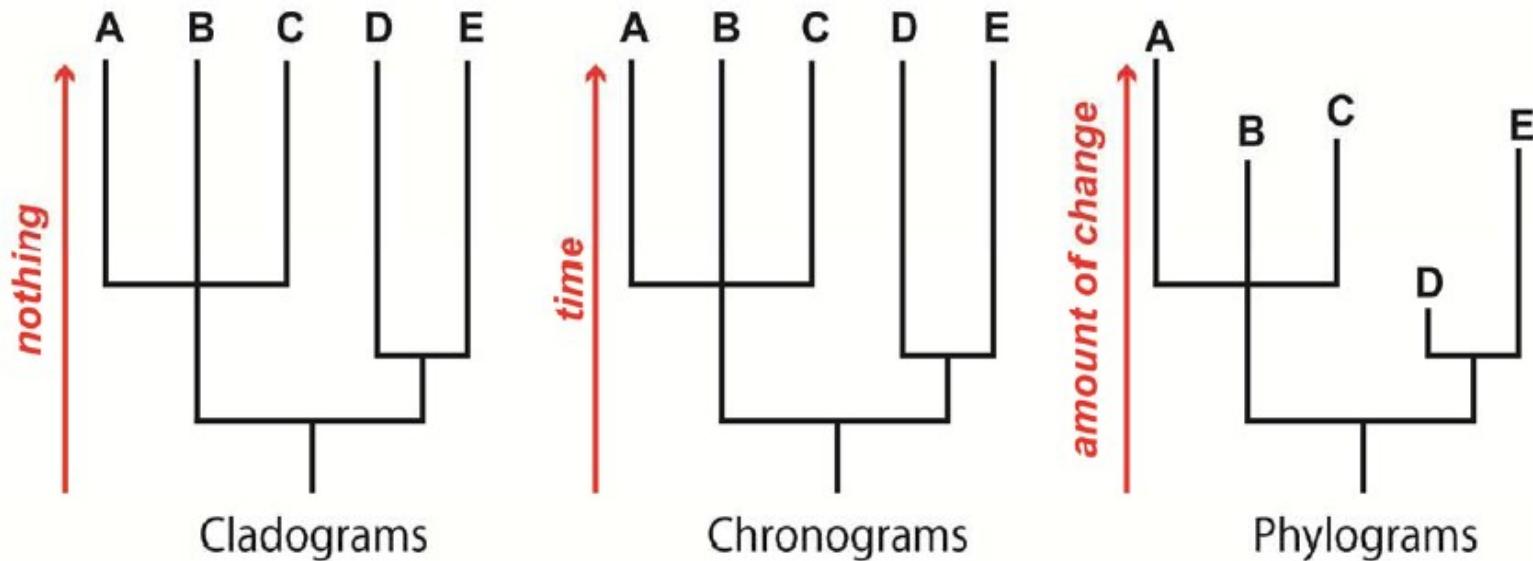
**Polyphyletic group:** warm-blood animals (mammals and birds; red) are from two different common ancestors.

- Monophly
- Paraphly
- Polyphly

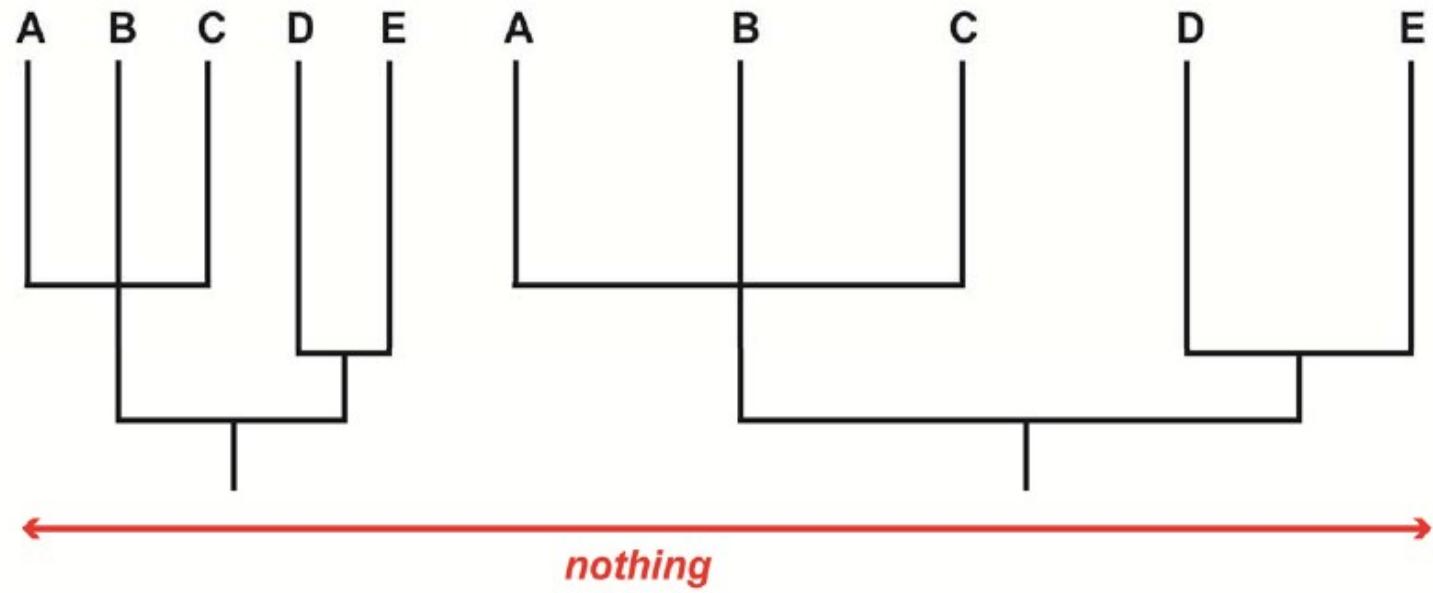


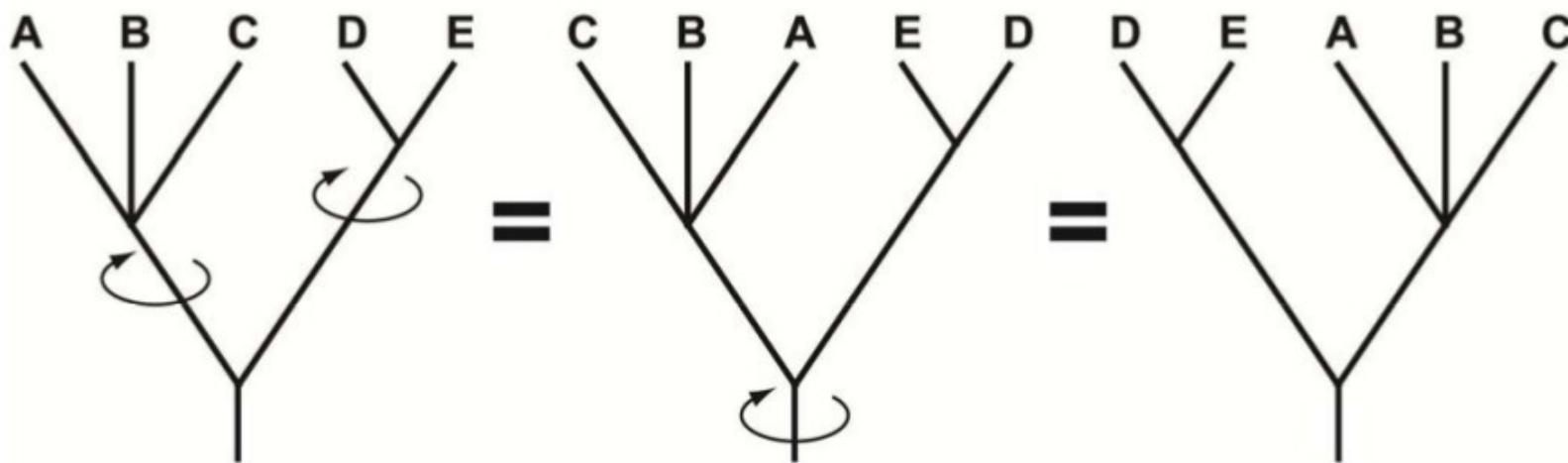
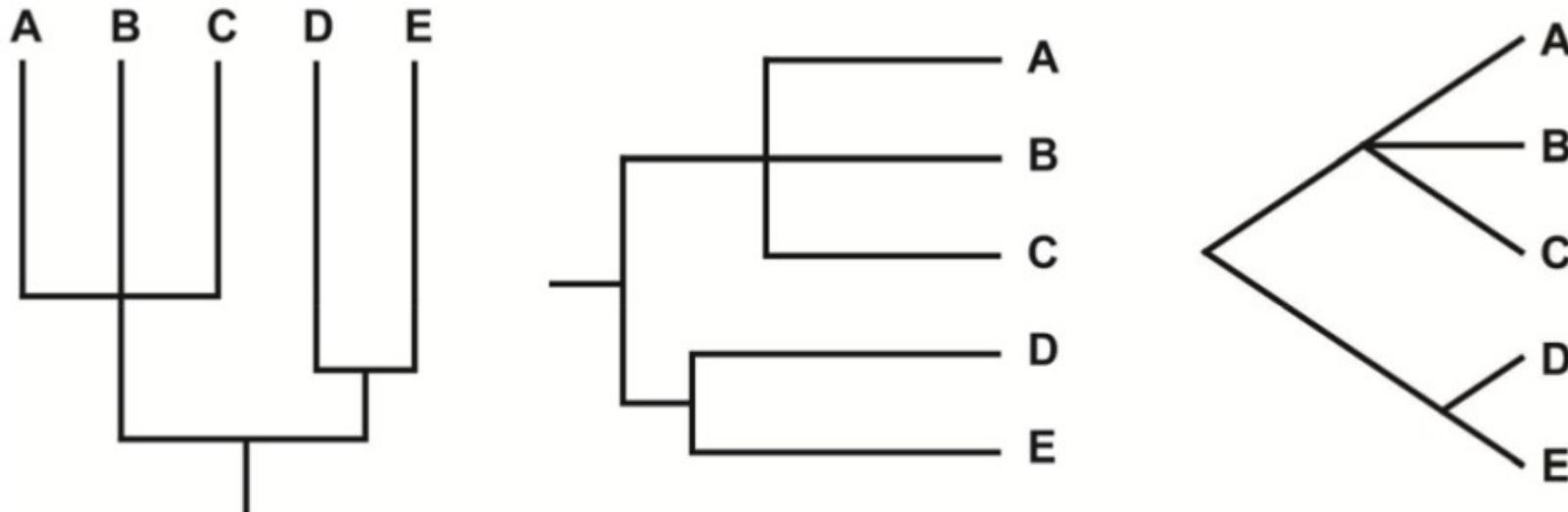
**Monophyletic group (clade):** birds and reptiles evolved from one common ancestor. (yellow)

*Vertical Axis:*



*Horizontal Axis:*



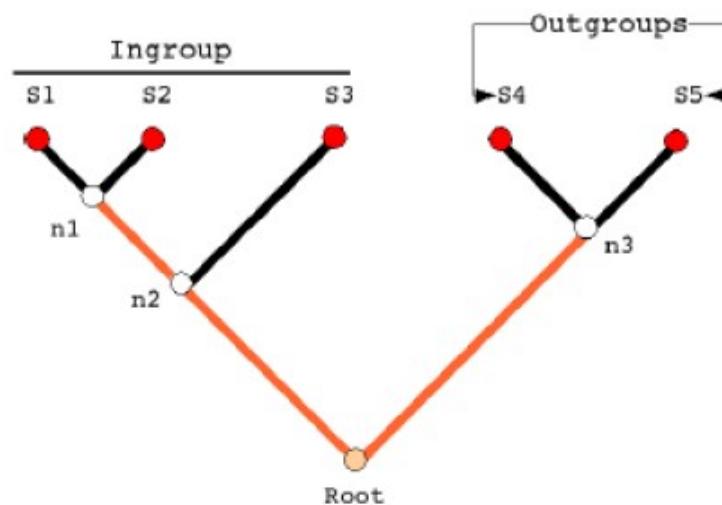
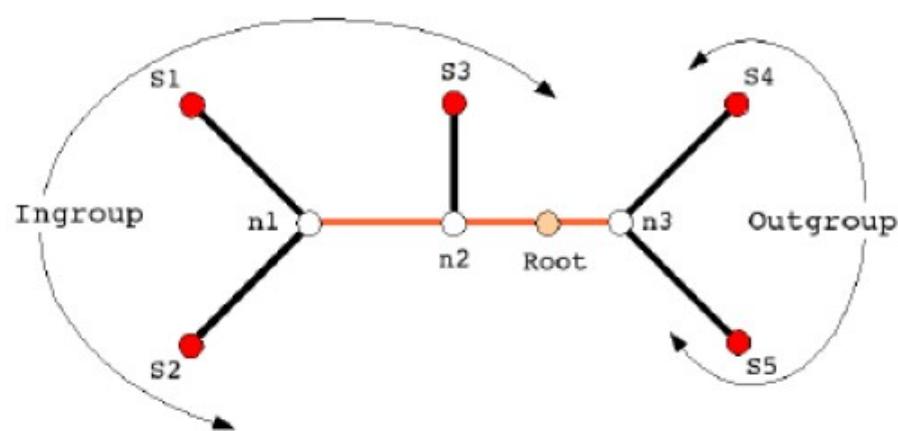


**Newick standard:**  $((A,B,C),(D,E));$

Trees can be rooted or un-rooted. Several strategies can be used for rooting (midpoint, outrgroup, etc)

Number of rooted trees with n OTUs:  $(2n - 3)!! = \frac{(2n - 3)!}{2^{n-2}(n - 2)!}$ , for  $n \geq 2$

Number of unrooted trees:  $(2n - 5)!! = \frac{(2n - 5)!}{2^{n-3}(n - 3)!}$ , for  $n \geq 3$



## How do we reconstruct phylogenetic trees?

Phylogenetic trees can be based on anything that can tell us on similarities and differences.  
(e.g RFLPs, phenotypic characters, sequences).

Molecular Phylogenetics is nowadays the most widely used method

Sequences are compared by means of alignments and this constitutes the information used to make the tree.

→ quality of the tree depends on the quality of the underlying alignment

So, how to find the best tree?

**Exhaustive** search: make ALL trees first, and then see which one best fits the data (you need an optimality criterion)

**Heuristic** search: Try to find a way to find an optimal tree (hopefully the best) without testing them all. You also need an optimality criterion and you are not guaranteed to find the best, but you save time.

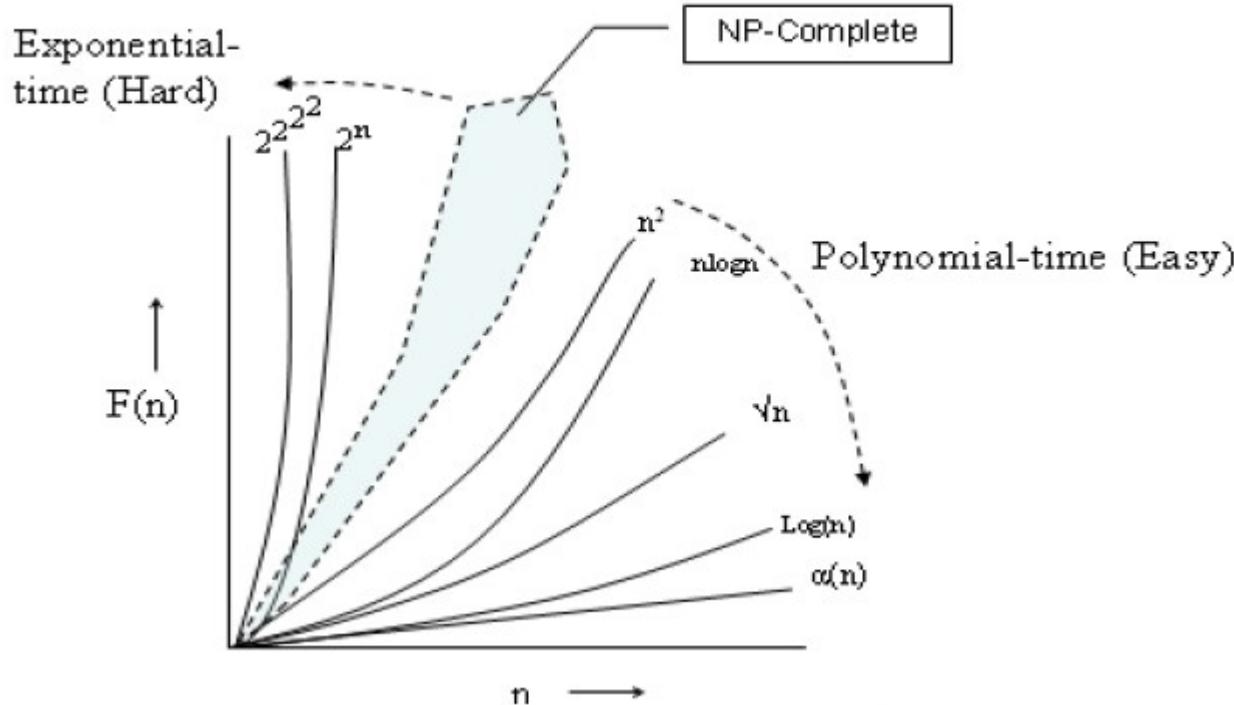
Number of taxa $T$	Number of unrooted bifurcating trees $B(T)$
3	1
4	3
5	15
6	105
7	945
8	10,395
9	135,135
10	2,027,025
22	$3 \times 10^{23}$
50	$3 \times 10^{74}$

Phylogenetic approaches:

Distance methods (NJ, UPGMA)

Maximum Parsimony

Probabilistic Methods (Maximum Likelihood and Bayesian Inference)



**Distance-based methods:** if there are no errors (in distances)then, the correct tree can be obtained in **polynomial time**. Otherwise, optimization problems are **NP-hard**.

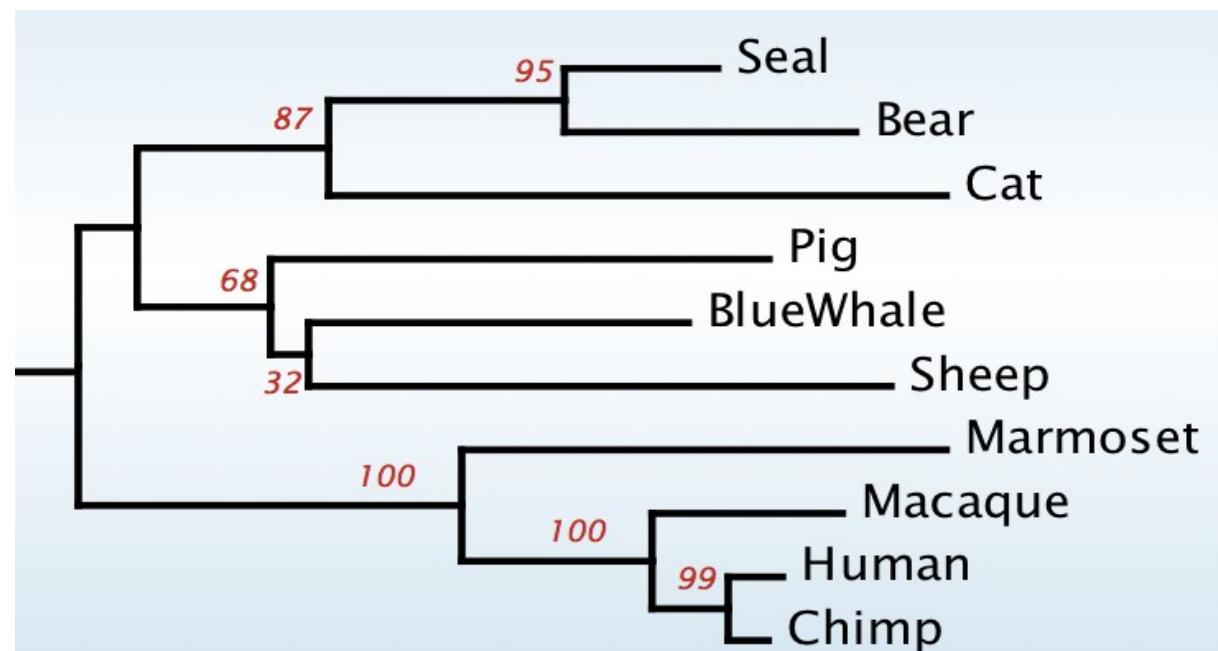
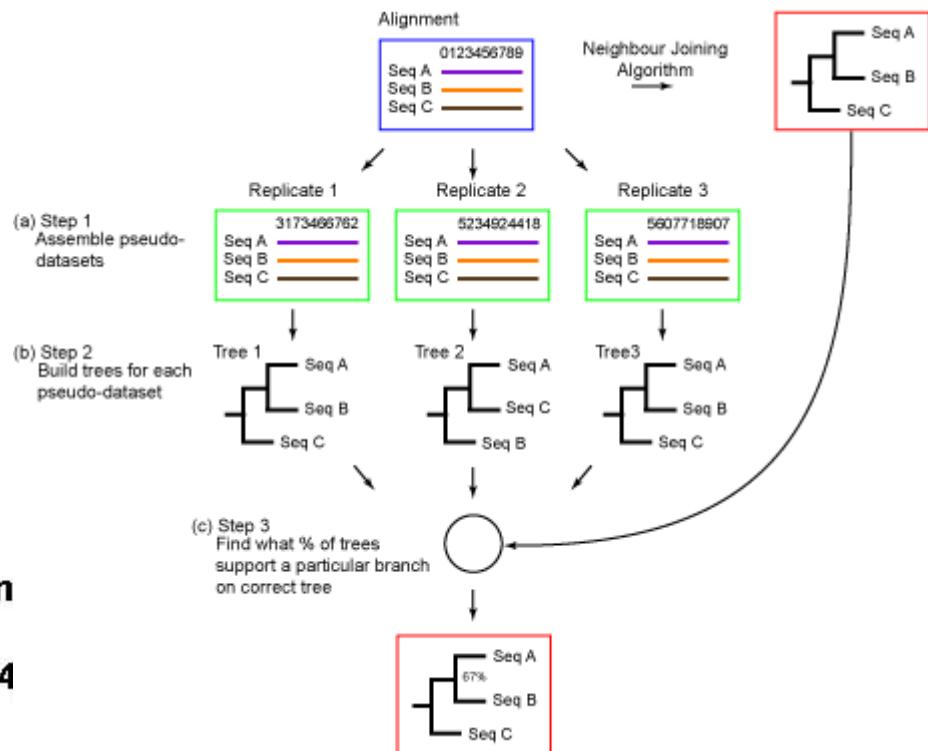
**Maximum Parsimony, Probabilistic Methods:** **NP-hard**

## Confidence values: bootstrapping

	Original sequence	Bootstrap Sequence
Human	A T G A C C	G T A A C A
Rat	A T A A C T	A T A A C A
Mouse	A T A A C T	A T A A C A
Chimp	A T G A C T	G T A A C A

Site 3 → is placed in first position

(Then the next five randomly chosen sites: 2, 1, 1, 5, 4 are placed in the next five positions.)

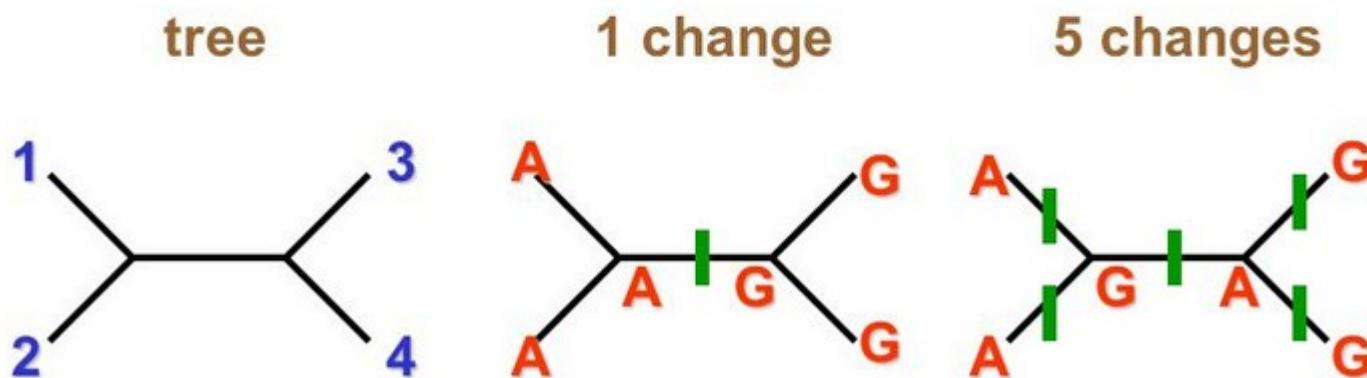


# Maximum Parsimony

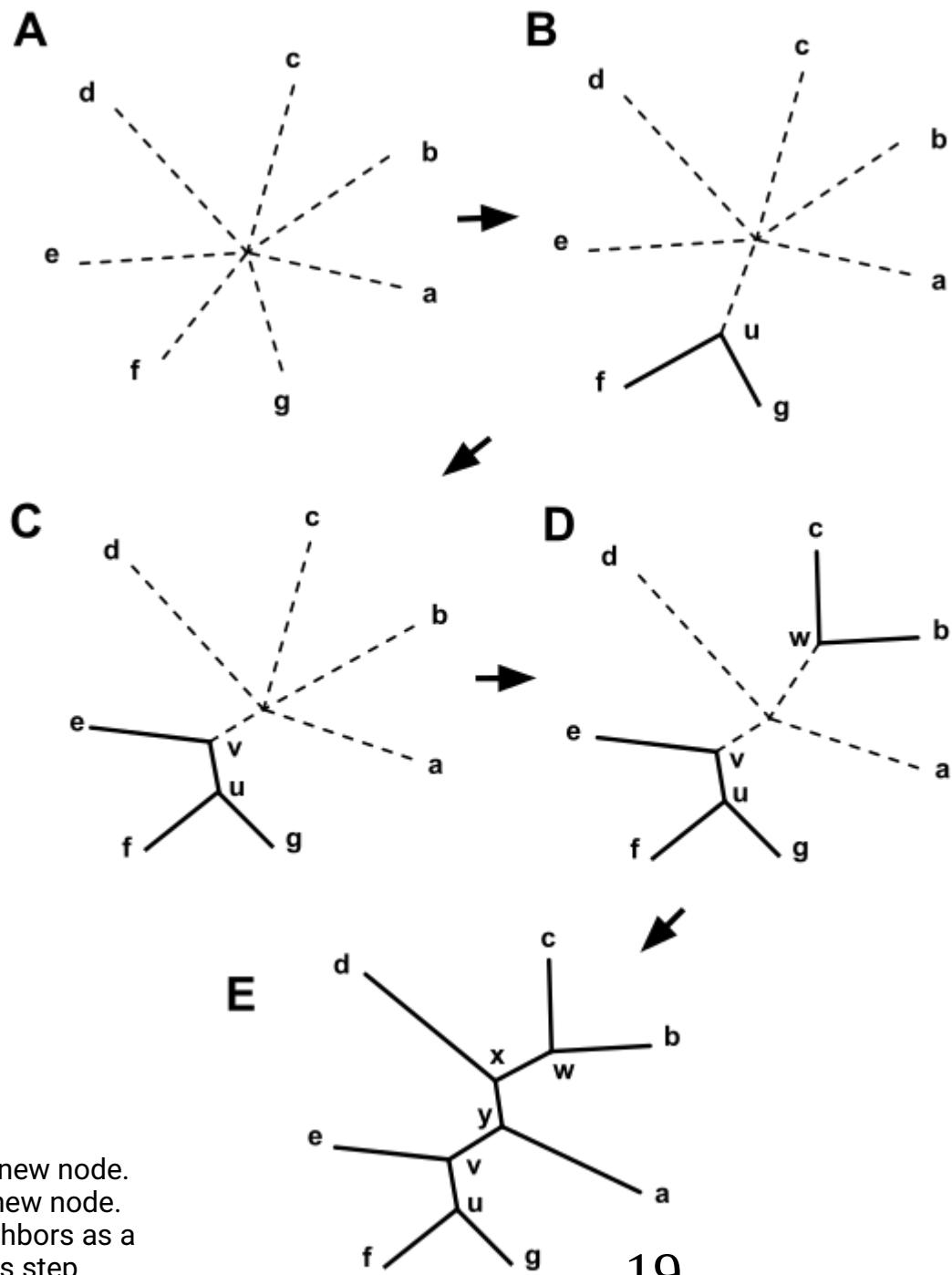
Finding the tree with that implies the minimal number of changes along its branches.

Taxon-1	ATATT
Taxon-2	ATCGT
Taxon-3	GCAGT
Taxon-4	GCCGT

For each site, the goal is to reconstruct the evolution of that site on a tree subject to the constraint of invoking the fewest possible evolutionary changes.



# Neighbor Joining



Based on the current distances matrix calculate the matrix Q

2. Find the pair of taxa in Q with the lowest value.

Create a node on the tree that joins these closest neighbors.

3. Calculate the distance of each of the taxa in the pair to this new node.

4. Calculate the distance of all taxa outside of this pair to the new node.

5. Start the algorithm again, considering the pair of joined neighbors as a single taxon and using the distances calculated in the previous step

## 10. Statistical Methods

### 10.1. Maximum Likelihood

- ♣ The phylogenetic methods described inferred the history (*or the set of histories*) that were most consistent with a set of observed data. All the methods explained used sequences as data and give one or more trees as phylogenetic hypotheses. Then, they use the logic of:

$$P(H/D)$$

- ♠ Maximum Likelihood (ML)<sup>28</sup> methods (*or maximum probability*) computes the probability of obtaining the data (*the observed aligned sequences*) given a defined hypothesis (*the tree and the model of evolution*). That is:

$$P(D/H)$$

#### A coin example

The ML estimation of the heads probabilities of a coin that is tossed  $n$  times.

---

<sup>28</sup>ML was invented by Ronald A. Fisher [27]. Likelihood methods for phylogenies were introduced by Edwards and Cavalli-Sforza for gene frequency data [9]. Felsenstein showed how to compute ML for DNA sequences [24].

# Likelihood

Given some data (**D**) a decision must be made about an adequate explanation (**H**, hypothesis)

**D**: alignment

**H**: Model of evolution, tree topology, branch lengths, parameters of the model

--> Each **H** will have a certain probability of producing the data  
 $P(D|H)$

The best **H** is that of the greatest **P**

## **Important remark!!**

The likelihood function **is not** the probability of a hypothesis being correct!!

The likelihood function is defined in terms of probability of producing the observed events not of the unknown parameters

**Thus:** the probability of observing the data has nothing to do with the probability that the underlying model is correct.



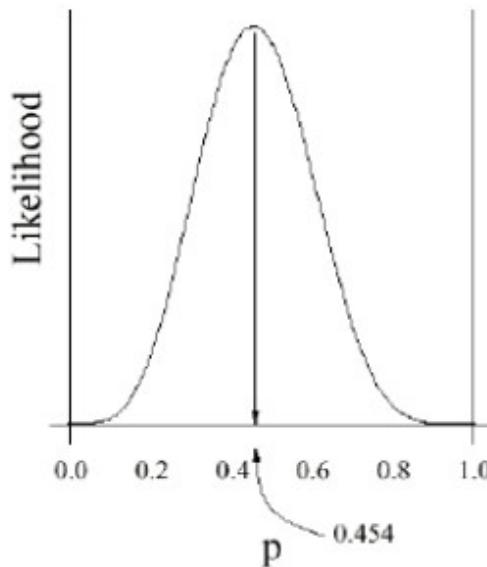
If tosses are all **independent**, and all have the same **unknown heads probability**  $p$ , then the observing sequence of tosses:

**HHTTHHTHHTTT**

we can calculate the ML of these data as:

$$L = \text{Prob}(D/p) = pp(1-p)(1-p)p(1-p)pp(1-p)(1-p)(1-p) = p^5(1-p)^6$$

Ploting  $L$  against  $p$ , we observe the probabilities of the same data ( $D$ ) for different values of  $p$ .



Thus the ML or the maximum probability to observe the above sequence of events is at  $p = 0.4545$ ,

Suppose we have:

- **Data:**

Sequence 1    **C C A T**

Sequence 2    **C C G T**

- **Model:**<sup>29</sup>

$$\pi = [0.1, 0.4, 0.2, 0.3]$$

$$P = \begin{bmatrix} 0.976 & 0.01 & 0.007 & 0.007 \\ 0.002 & 0.983 & 0.005 & 0.01 \\ 0.003 & 0.01 & 0.979 & 0.007 \\ 0.002 & 0.013 & 0.005 & 0.979 \end{bmatrix}$$

$$\begin{aligned} L_{(Seq_1 \rightarrow Seq_2)} &= \pi_C P_{C \rightarrow C} \pi_C P_{C \rightarrow C} \pi_A P_{A \rightarrow G} \pi_T P_{T \rightarrow T} \\ &= 0.4 \times 0.983 \times 0.4 \times 0.983 \times 0.1 \times 0.007 \times 0.3 \times 0.979 \\ &= 0.0000300 \end{aligned}$$

$$\ln L_{tree: Seq_1 \rightarrow Seq_2} = -10.414$$

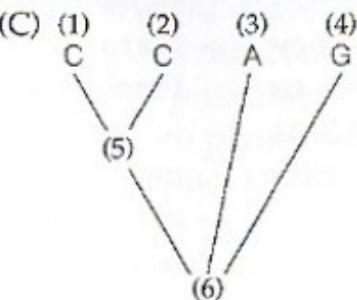
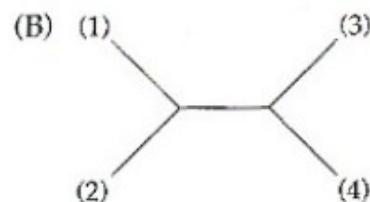
---

<sup>29</sup>Note that the base composition sum one, but indeed the the rows of substitution matrix sum one. Why?

## computation in a real problem

(A)      1                   $j$                    $N$

(1) C ... G G A C A C **G** T T T A ... C  
 (2) C ... A G A C A C C T C T C T A ... C  
 (3) C ... G G A T A A G T T T A A ... G  
 (4) C ... G G A T A G C C T A G ... C



(D)

$$L_{(j)} = \text{Prob} \left( \begin{array}{c} \text{C} & \text{C} & \text{A} & \text{G} \\ \backslash & / & & \\ \text{A} & & \text{A} & / \\ & & \backslash & \\ & & \text{A} & \end{array} \right) + \text{Prob} \left( \begin{array}{c} \text{C} & \text{C} & \text{A} & \text{G} \\ \backslash & / & & \\ \text{C} & & \text{A} & / \\ & & \backslash & \\ & & \text{A} & \end{array} \right)$$

$$+ \dots + \text{Prob} \left( \begin{array}{c} \text{C} & \text{C} & \text{A} & \text{G} \\ \backslash & / & & \\ \text{G} & & \text{A} & / \\ & & \backslash & \\ & & \text{C} & \end{array} \right)$$

$$+ \dots + \text{Prob} \left( \begin{array}{c} \text{C} & \text{C} & \text{A} & \text{G} \\ \backslash & / & & \\ \text{T} & & \text{A} & / \\ & & \backslash & \\ & & \text{T} & \end{array} \right)$$

- Tree after rooting in an arbitrary node (reversible model).
- The likelihood for a particular site is the sum of the probabilities of every possible reconstruction of ancestral states given some model of base substitution.
- The likelihood of the tree is the product of the likelihood at each site.

$$L = L_{(1)} \cdot L_{(2)} \cdot \dots \cdot L_{(N)} = \prod_{j=1}^N L_{(j)}$$

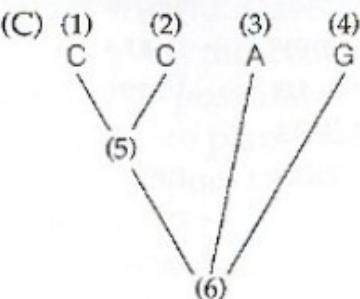
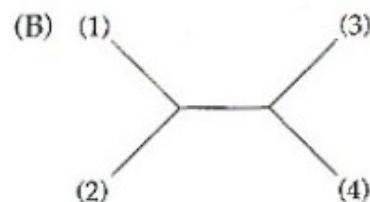
- The likelihood is reported as the sum of the log likelihood of the full tree.

$$\ln L = \ln L_{(1)} + \ln L_{(2)} + \dots + \ln L_{(N)} = \sum_{j=1}^N \ln L_{(j)}$$

## computation in a real problem

(A)      1                   $j$                    $N$

(1) C ... G G A C A C **G** T T T A ... C  
 (2) C ... A G A C A C C T C T C T A ... C  
 (3) C ... G G A T A A G T T T A A ... G  
 (4) C ... G G A T A G C C T A G ... C



(D)

$$L_{(j)} = \text{Prob} \left( \begin{array}{c} \text{C} & \text{C} & \text{A} & \text{G} \\ \backslash & / & & \\ \text{A} & & \text{A} & / \\ & & \backslash & \\ & & \text{A} & \end{array} \right) + \text{Prob} \left( \begin{array}{c} \text{C} & \text{C} & \text{A} & \text{G} \\ \backslash & / & & \\ \text{C} & & \text{A} & / \\ & & \backslash & \\ & & \text{A} & \end{array} \right)$$

$$+ \dots + \text{Prob} \left( \begin{array}{c} \text{C} & \text{C} & \text{A} & \text{G} \\ \backslash & / & & \\ \text{G} & & \text{A} & / \\ & & \backslash & \\ & & \text{C} & \end{array} \right)$$

$$+ \dots + \text{Prob} \left( \begin{array}{c} \text{C} & \text{C} & \text{A} & \text{G} \\ \backslash & / & & \\ \text{T} & & \text{A} & / \\ & & \backslash & \\ & & \text{T} & \end{array} \right)$$

- Tree after rooting in an arbitrary node (reversible model).
- The likelihood for a particular site is the sum of the probabilities of every possible reconstruction of ancestral states given some model of base substitution.
- The likelihood of the tree is the product of the likelihood at each site.

$$L = L_{(1)} \cdot L_{(2)} \cdot \dots \cdot L_{(N)} = \prod_{j=1}^N L_{(j)}$$

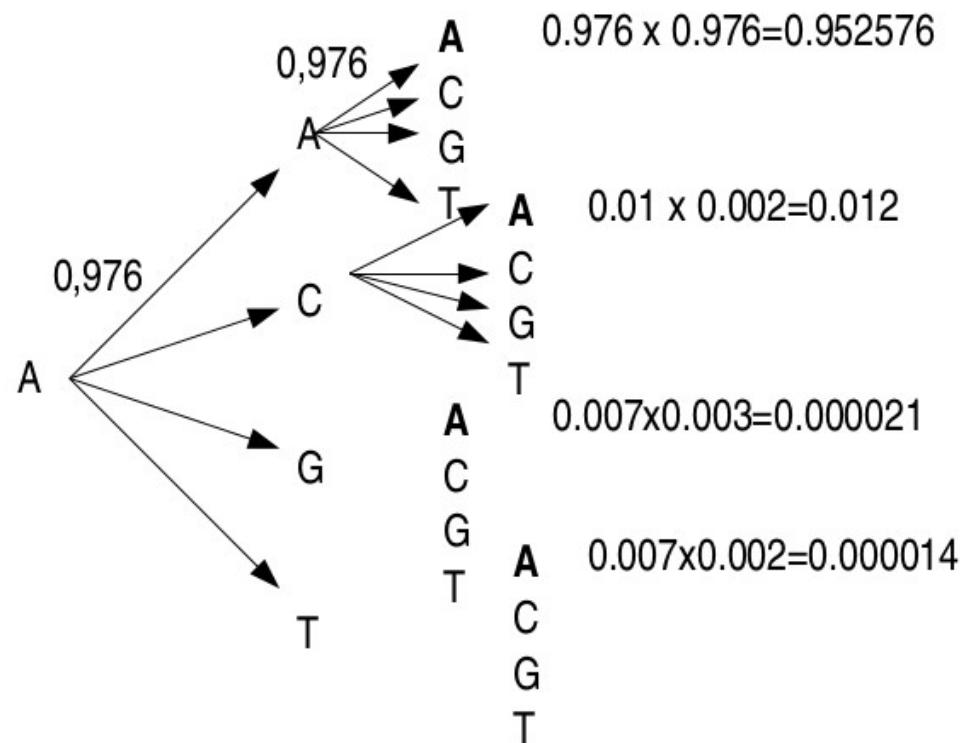
- The likelihood is reported as the sum of the log likelihood of the full tree.

$$\ln L = \ln L_{(1)} + \ln L_{(2)} + \dots + \ln L_{(N)} = \sum_{j=1}^N \ln L_{(j)}$$

Ancestral= A

Observed= A

Time = 2x CED



$$P = A \rightarrow A = 0.976 \text{ (CED=1)}, 0.9646 \text{ (CED=2)}, \dots$$

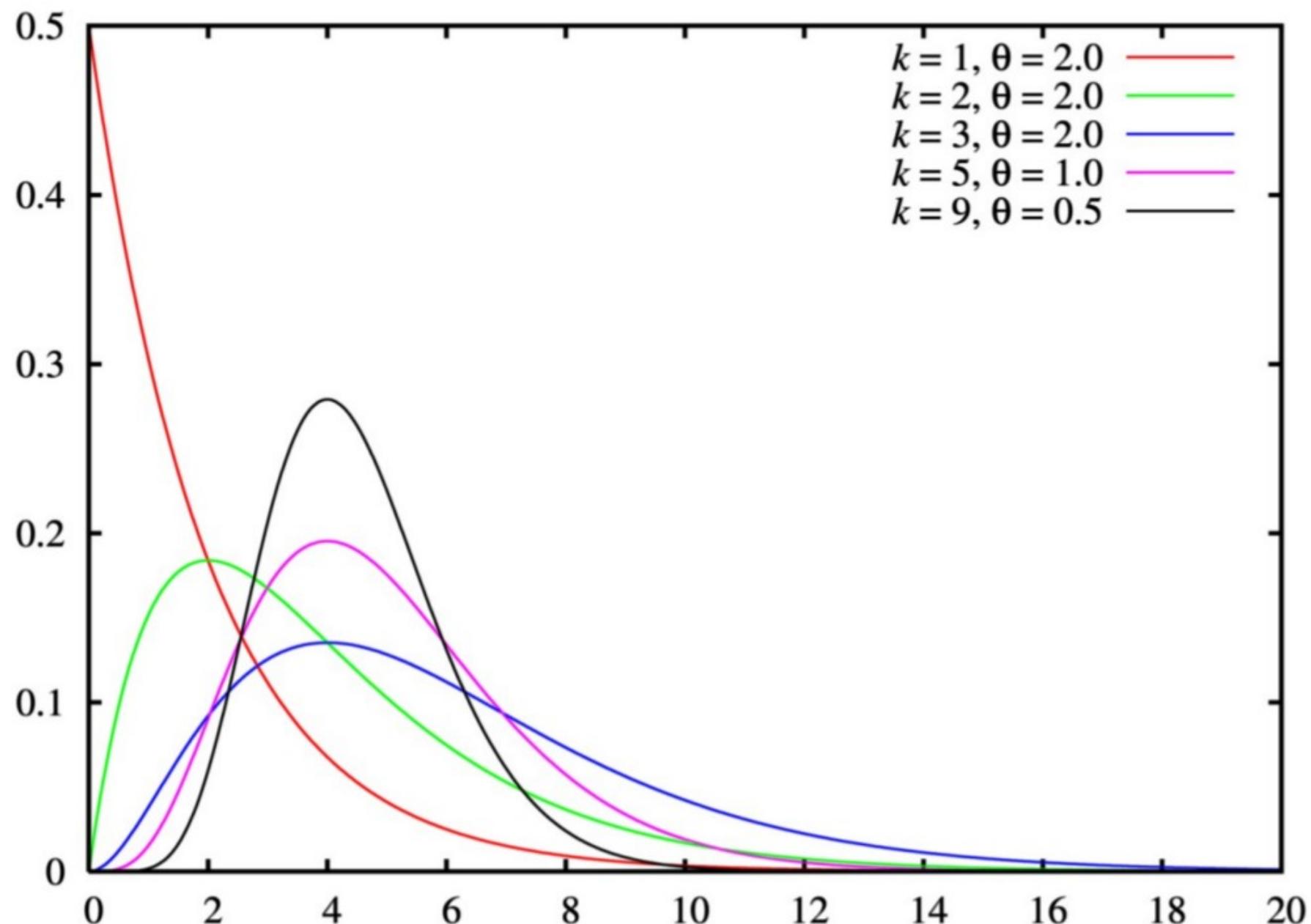
## Variation of evolutionary rates over sites

So far we have assumed 1 evolutionary rate for all the sites is that necessary so?

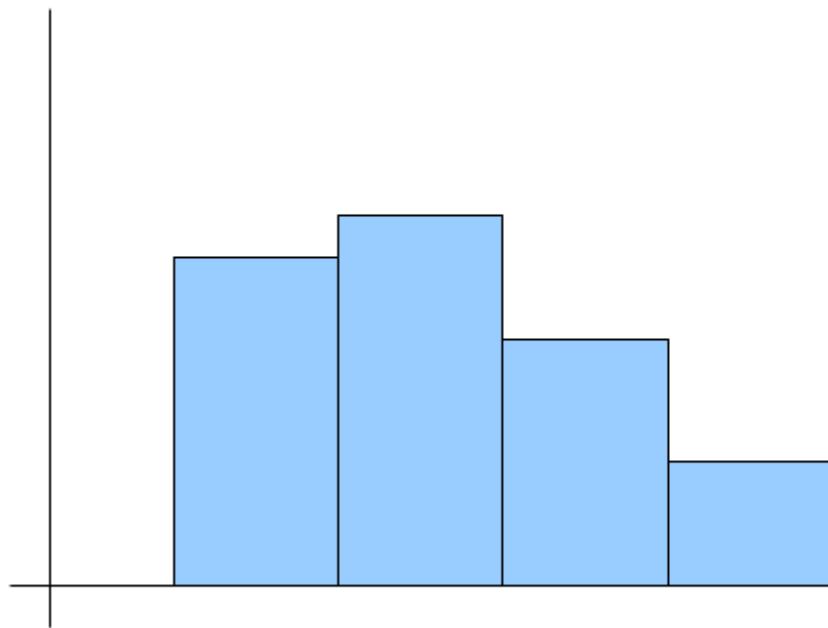
-> different sites, different evolutionary constraints

### **Approximation of rate-heterogeneity by a gamma-distribution**

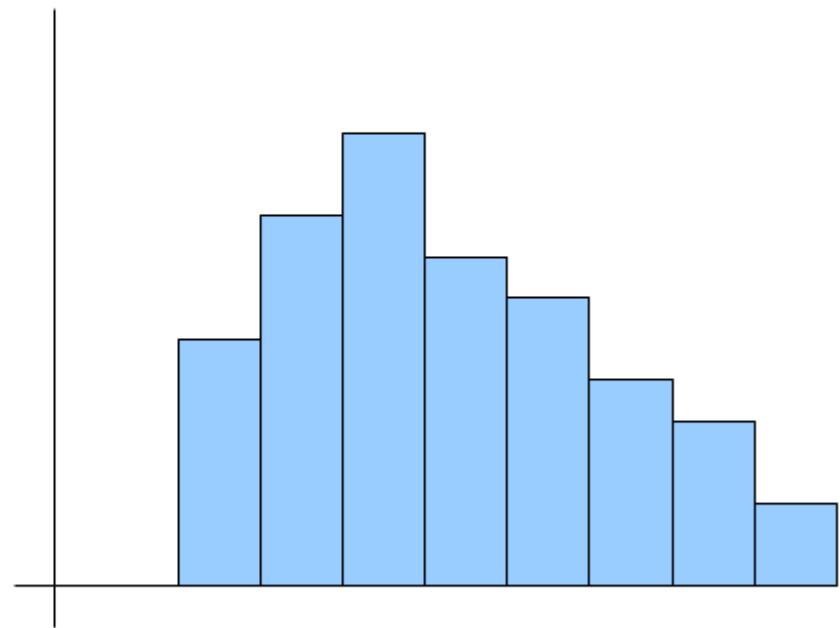
Parameters  $k$  (shape),  $\Theta$  (scale). mean =  $k\Theta$



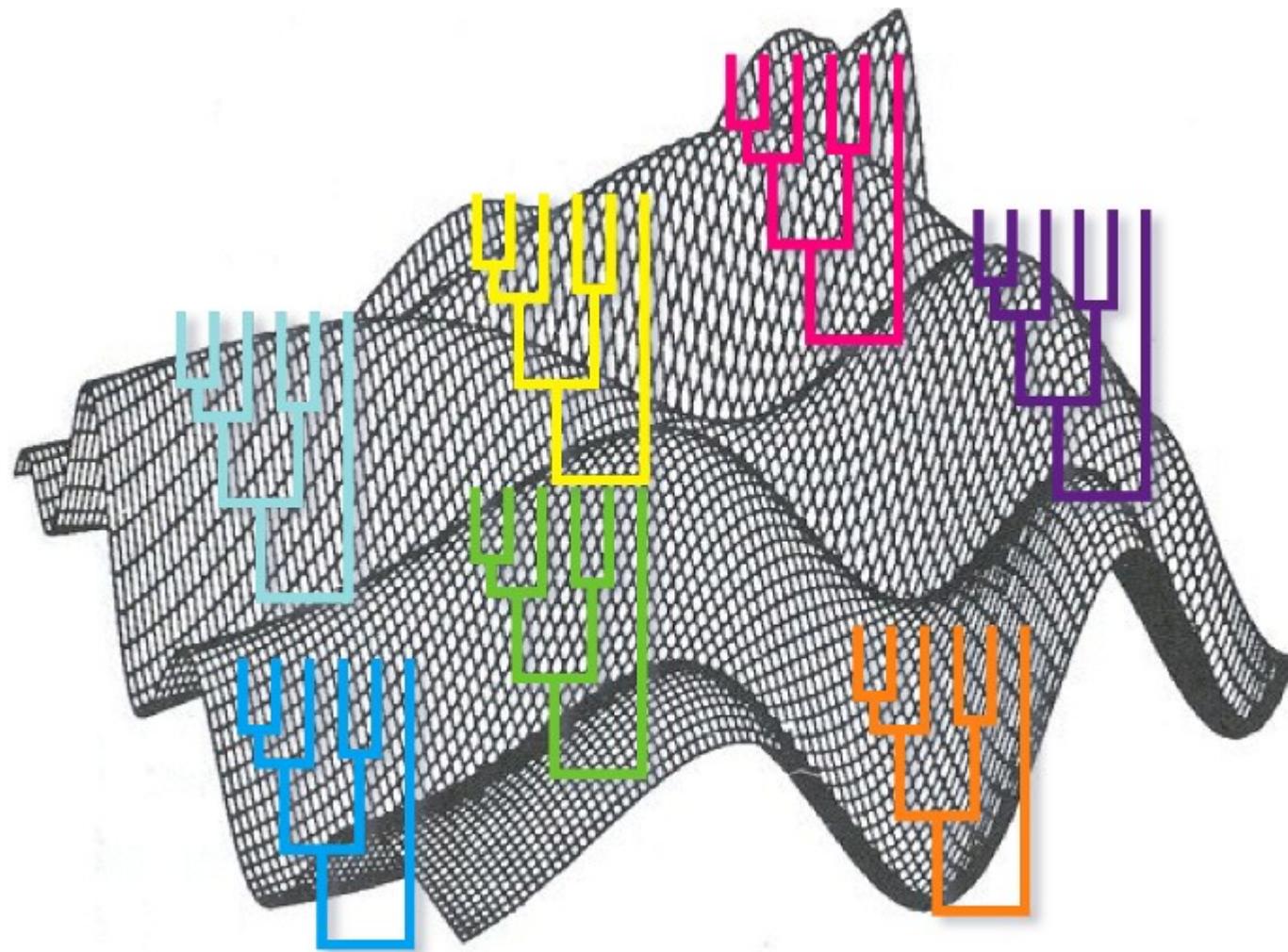
approximation by a discrete distribution



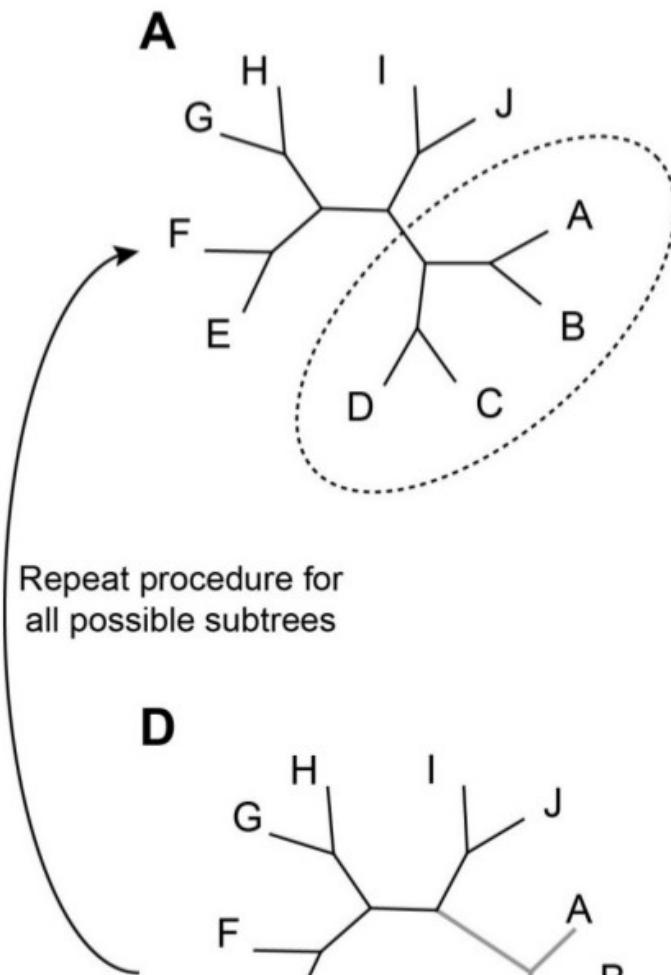
$n=4$



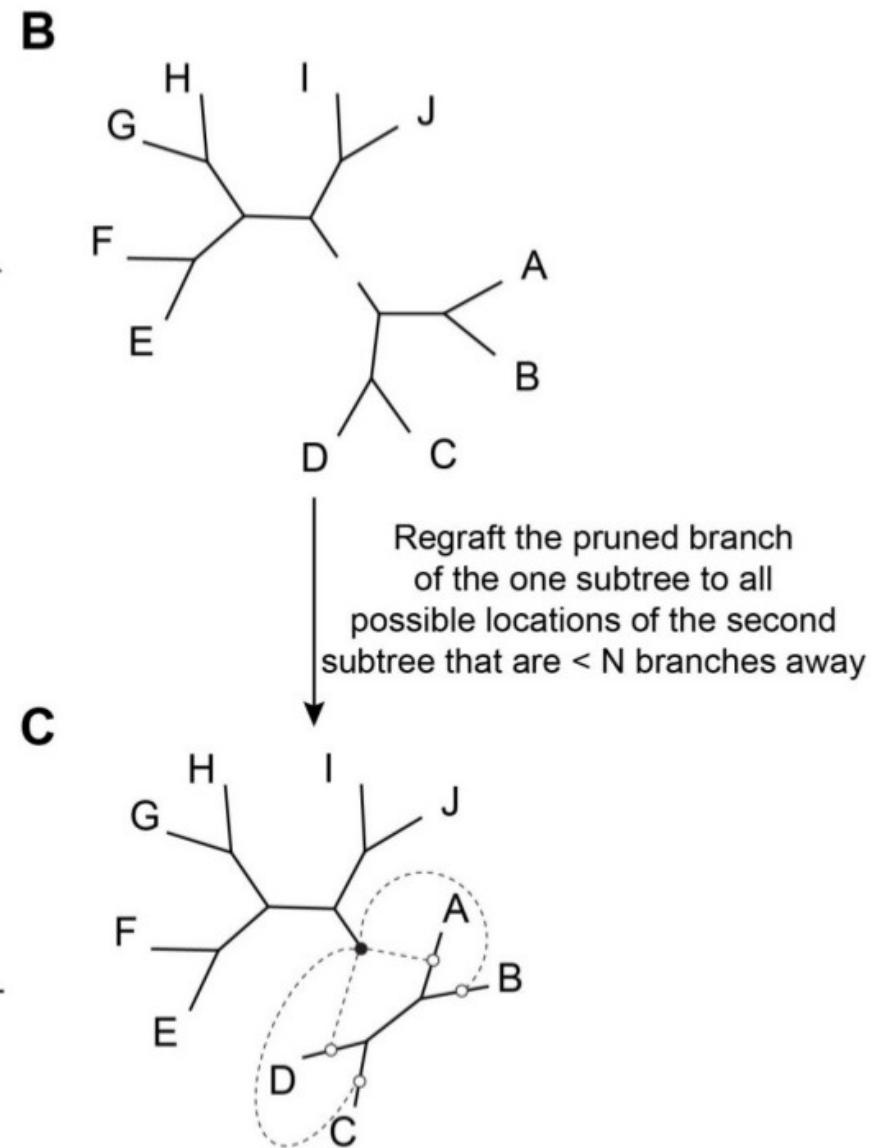
$n=8$



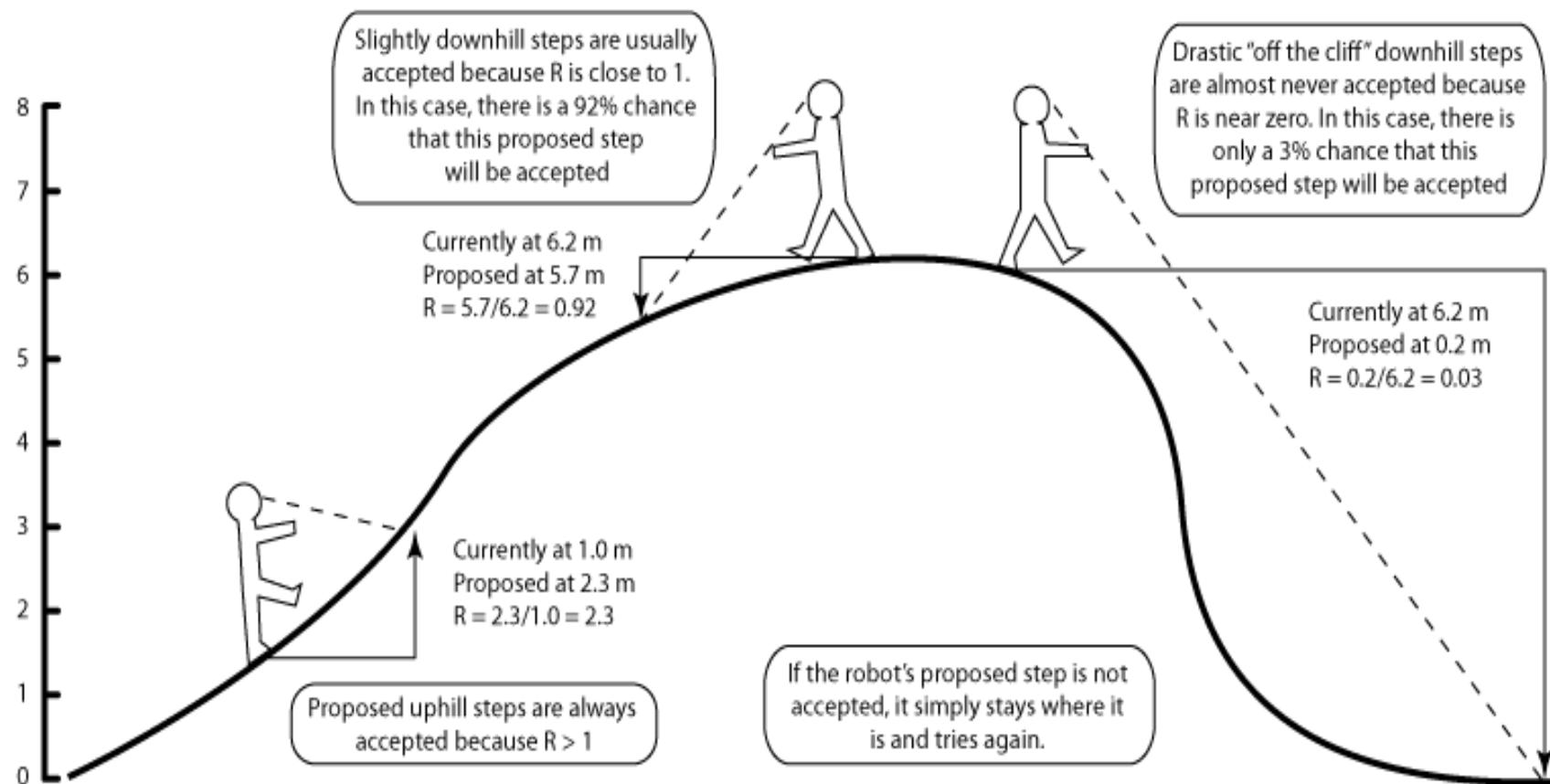
# Heuristic Search



Break a branch,  
separate the subtrees



Heuristic searches don't guarantee discovery of optimal topology



# Bayesian Inference

- ♣ **Maximum Likelihood** will find the tree that is most likely to have produced the observed sequences, or formally  $P(D/H)$  (the probability of seeing the data given the hypothesis).
- ♠ **A Bayesian approach** will give you the tree (or set of trees) that is most likely to be explained by the sequences, or formally  $P(H/D)$  (the probability of the hypothesis being correct given the data).
- ◊ **Bayes Theorem** provides a way to calculate the probability of a model (*tree topology and evolutionary model*) from the results it produces (*the aligned sequences we have*), what we call a **posterior probability**<sup>31</sup>.

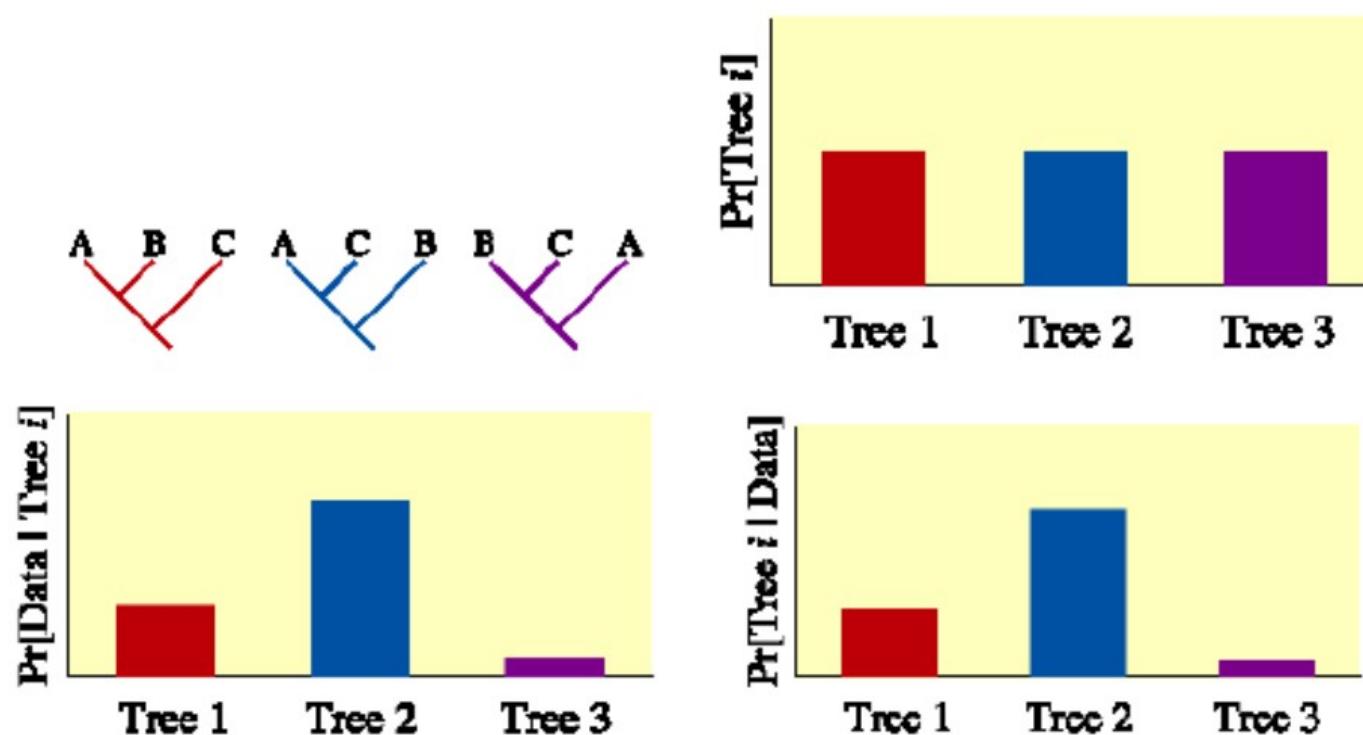
*Thomas Bayes (1702-1761)*



$$P(\theta/D) = \frac{P(\theta) \cdot P(D/\theta)}{P(D)}$$

## The main components of Bayes analysis

- $P(\theta)$  The **prior probability** of a tree represents the probability of the tree before the observations have been made. Typically, all trees are considered equally probable.



- $P(D/\theta)$  The **likelihood** is proportional to the probability of the observations (data sets) conditional on the tree.

# Posterior probability $\sim$ likelihood $\times$ prior probability



Usually flat distributions  
uninformative  
our beliefs before “seeing” the data

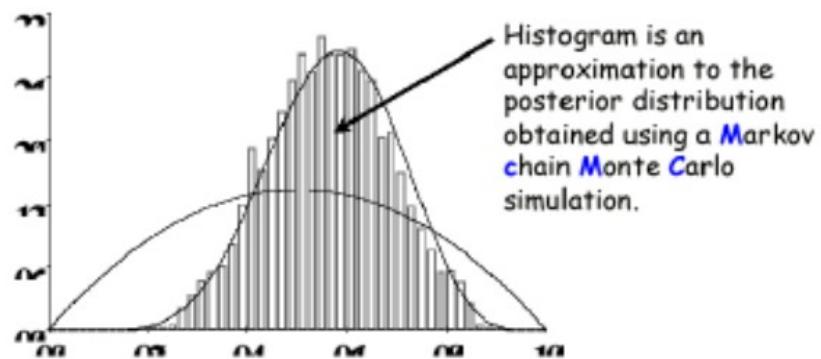


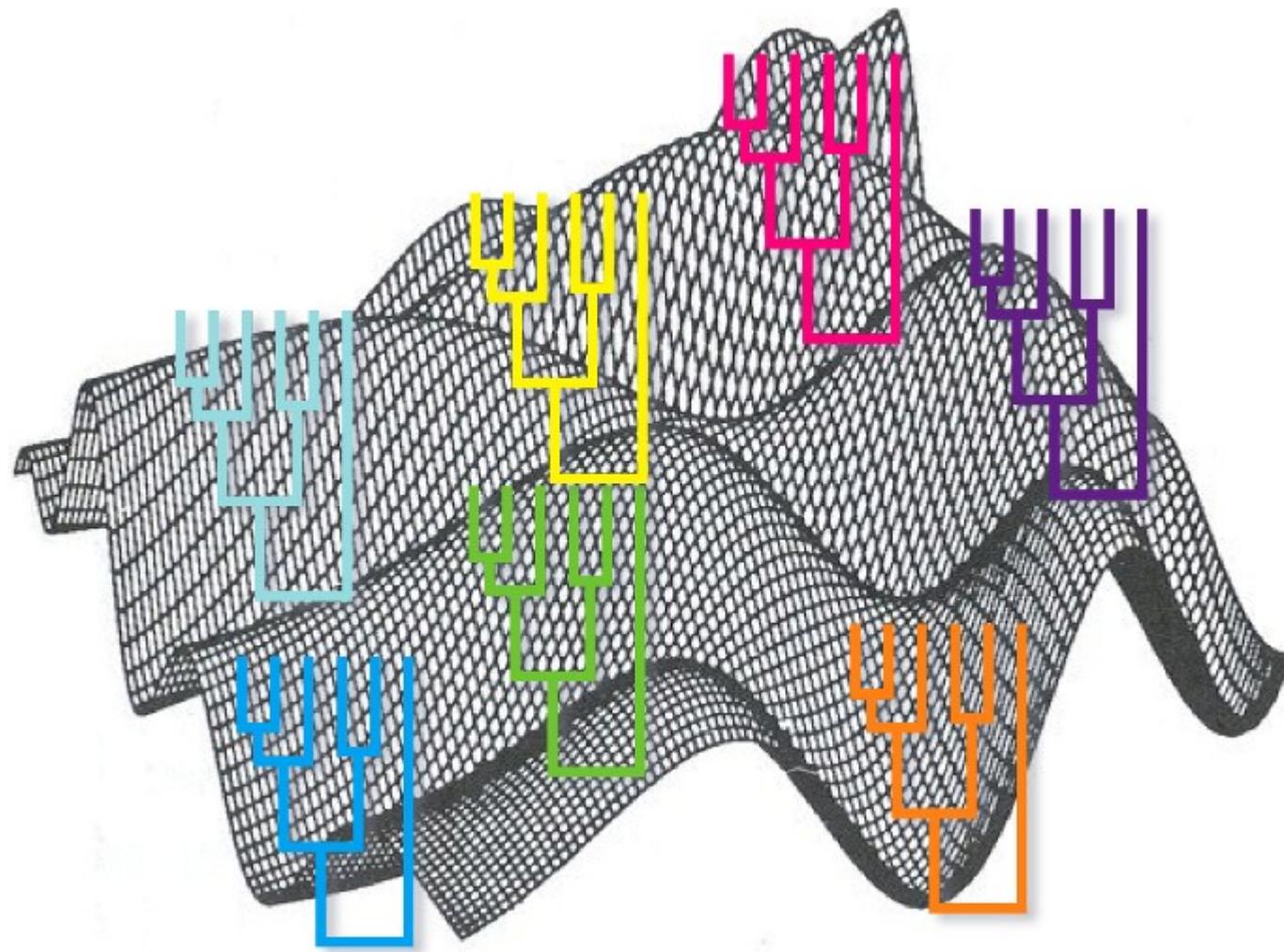
Therefore most of the differences in the posterior probability are due to differences in likelihood (-> strong connection to ML)

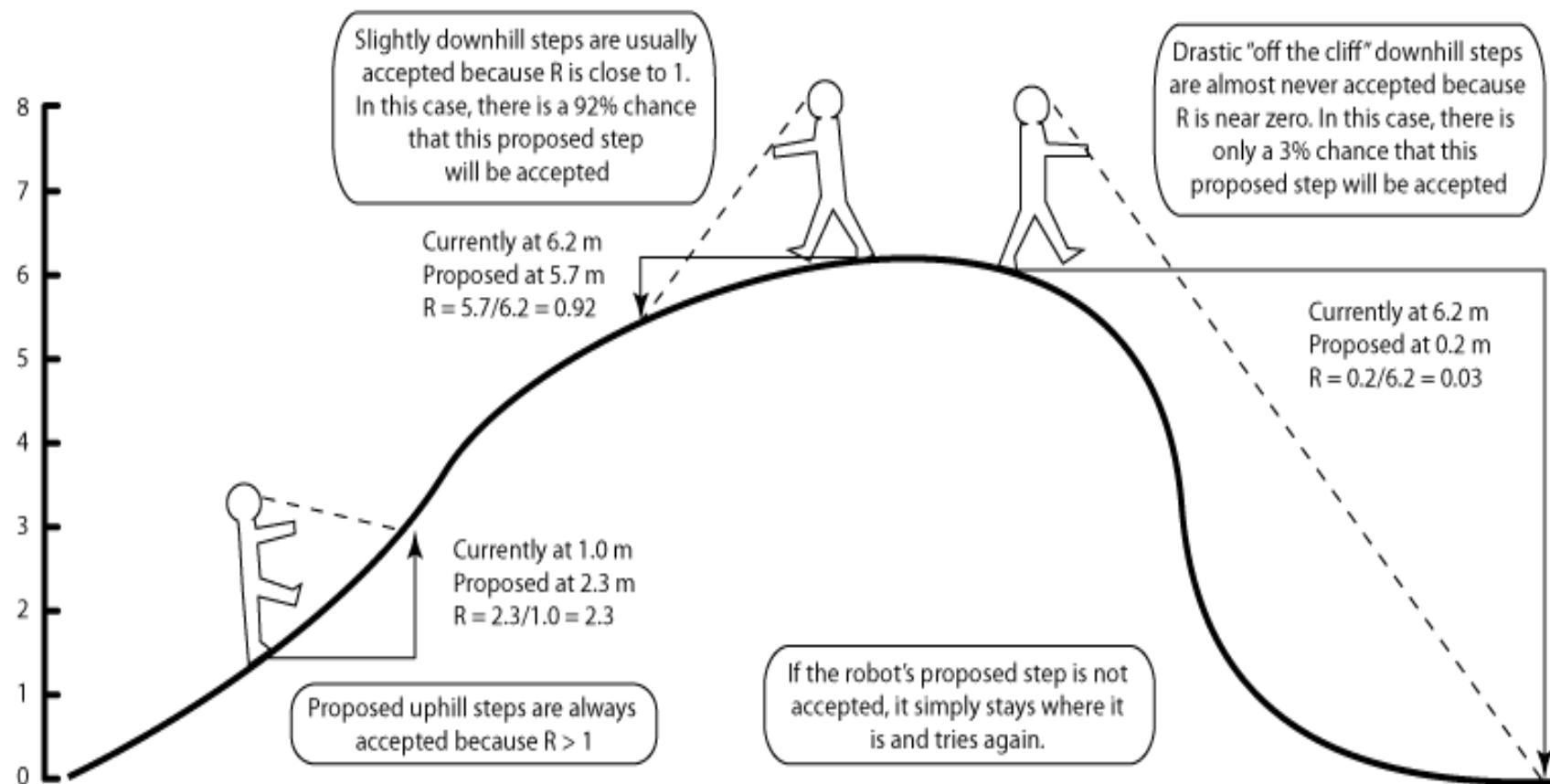
## How to find the solution

There's no analytical solution for a Bayesian system. However, giving:

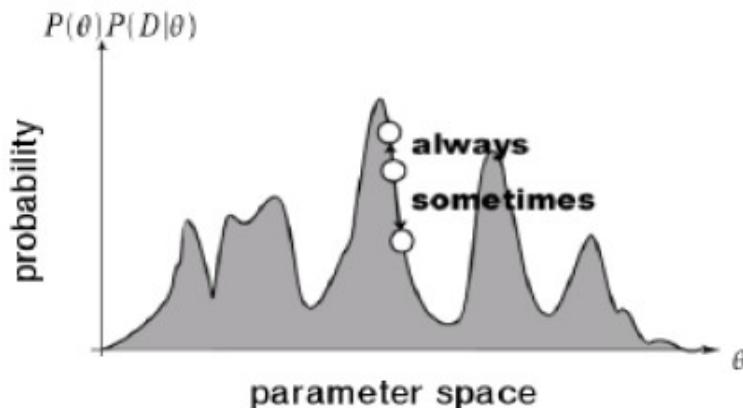
- **Data:** Sequence data,
- **Model:** The evolutionary model, base frequencies, among site rate variation parameters, a tree topology, branch lengths
- **Priors** distribution on the model parameters, and
- **A method** for calculating posterior distribution from prior distribution and data: **MCMC** technique<sup>32</sup>







- Each step in a Markov chain **a random modification** of the tree topology, a branch length or a parameter in the substitution model (e.g. substitution rate ratio) is assayed.
- If the **posterior computed is larger** than that of the current tree topology and parameter values, the proposed step **is taken**.
- Steps downhill are not automatic accepted, depending on the magnitude of the decrease.



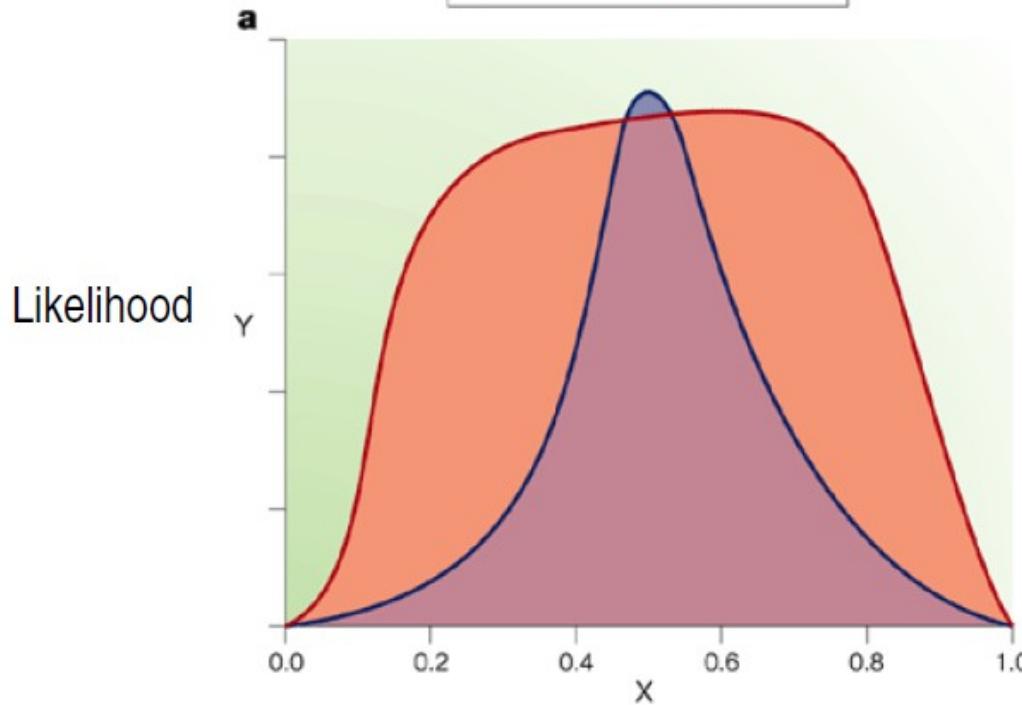
- Using these rules, the **Markov chain visits regions of the tree space in proportion of their posterior**.
- Suppose you sample 100,000 trees and a particular clade appears in 74,695 of the sampled trees. The probability (giving the observed data) that the group is monophyletic is 0.746, because **MC visits trees in proportion to their posterior probabilities**.

A fundamental difference:

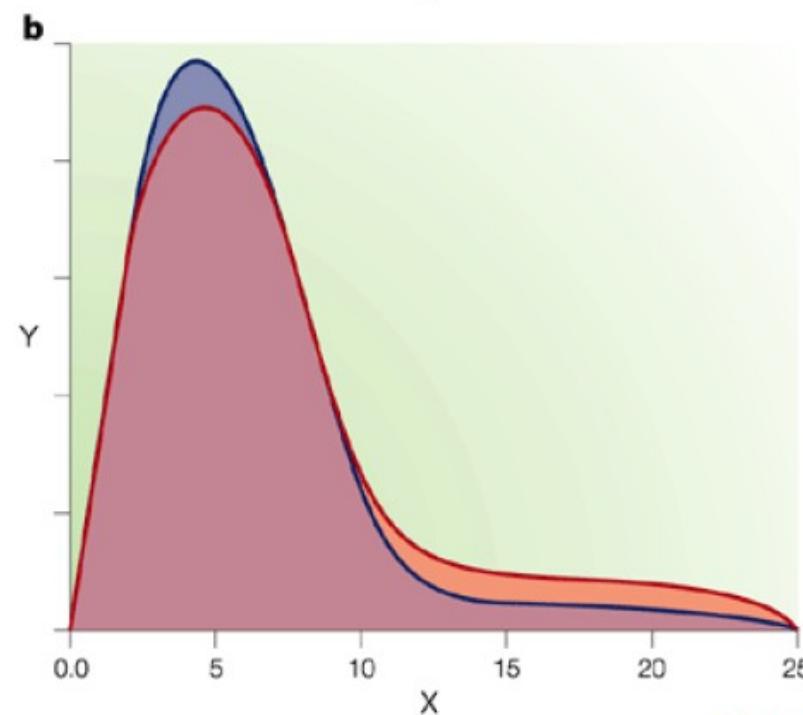
**ML** commonly uses **Joint Estimation**: finding the highest point in the parameter landscape

**Bayesian analyses** measure the volume under the posterior probability surface, the parameters are integrated (marginalized) to obtain the marginal posterior probability of a tree (**Marginal estimation**)

## Joint versus Marginal estimation



Every part of the surface affects the results, so the prior distributions may be seriously considered.



1) When there are few parameters and plenty of data, the likelihood and posterior landscapes become thin spires with high peaks:

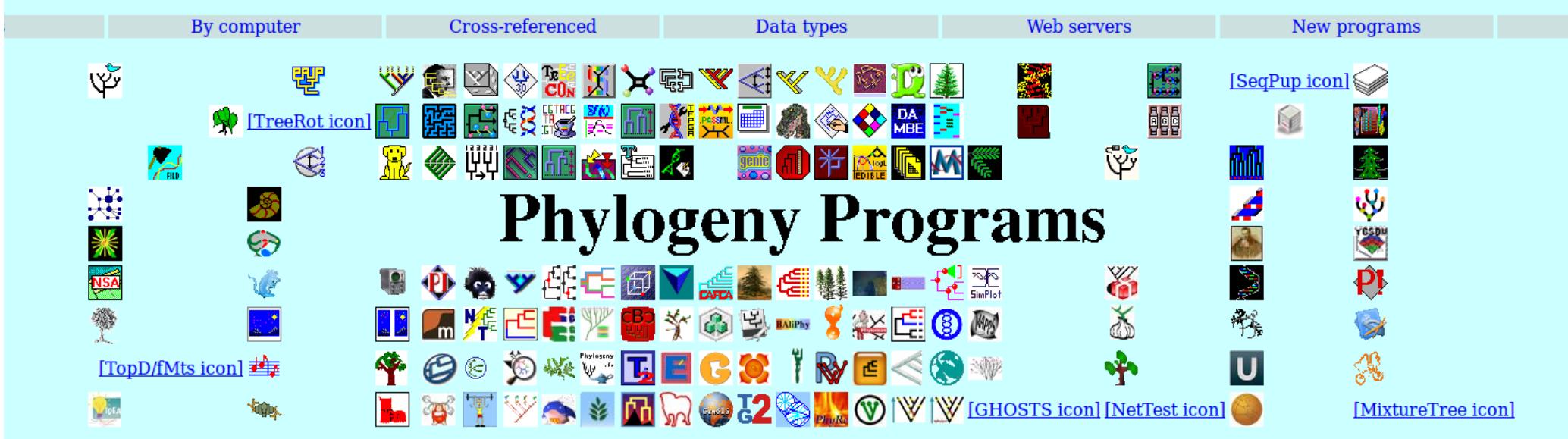
---> The height of the peak is a good predictor of the integral over the whole surface.....  $ML \sim BI$

2) As the amount of data decreases relative to the number of parameters (complex models), the landscapes become soft hills, and consideration on the uncertainty of the parameters is necessary.

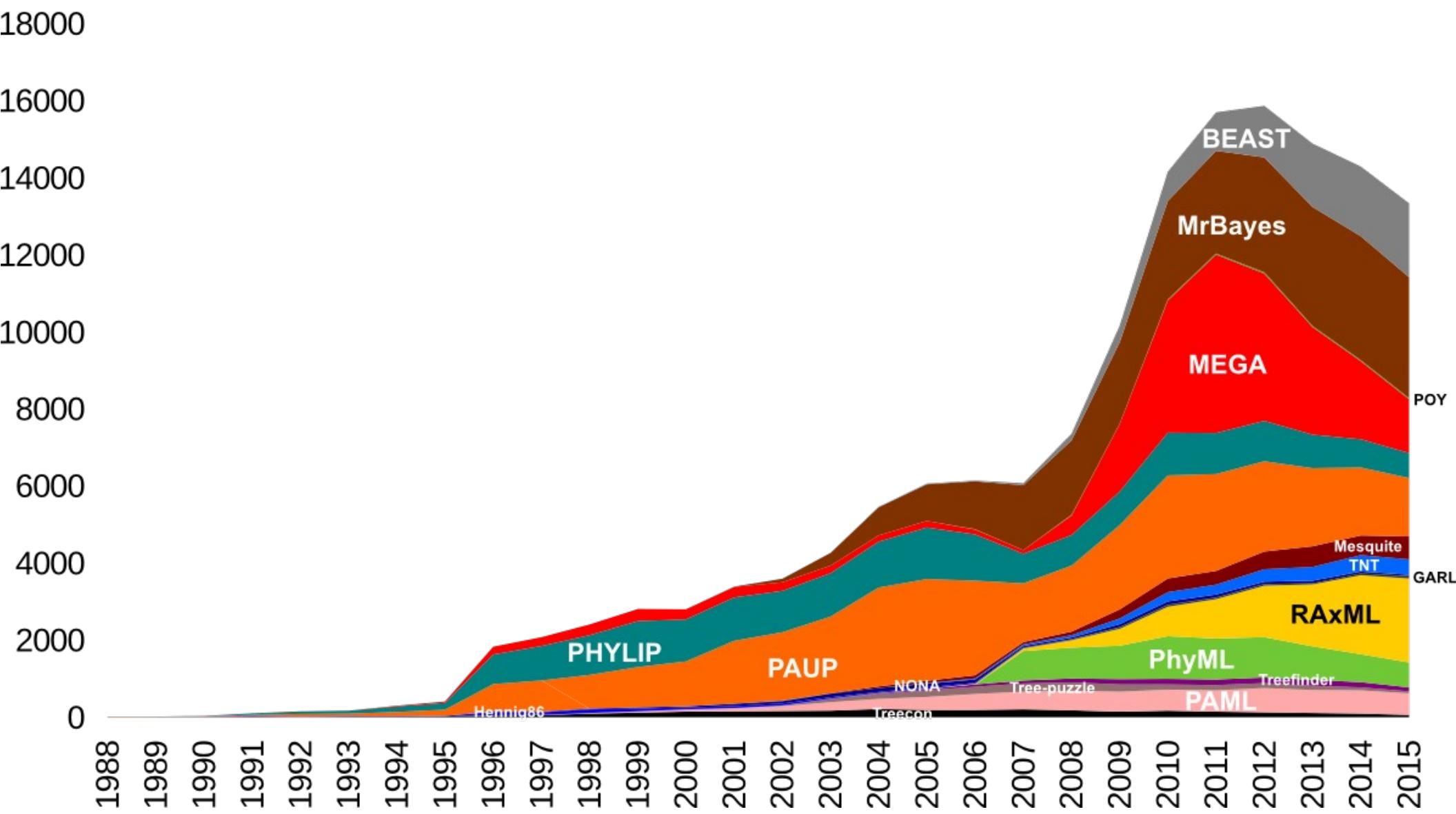
--> marginal estimation (BI) is more reliable

<http://evolution.genetics.washington.edu/phylip/software.html>

Owing to other pressures on my time, I cannot devote time to searching for new programs, so **their authors are begged to (please!) use the submission form instead**. That form will be found at the "Submitting" link below. If you are upset that your program is not included, but it's too much trouble for you to fill out the submission form, then I will not listen to you. This list of software is now aging and its links are becoming more and more outdated. I will make attempts to fix them when I can. If anyone else wants to help with this, let me know.



[https://en.wikipedia.org/wiki/List\\_of\\_phylogenetics\\_software](https://en.wikipedia.org/wiki/List_of_phylogenetics_software)



# ngphylogeny.fr

ATAAAAGAACTCAAGCTCTGGGAAGTCAGTAG  
ATAAAGAACACAGAGTTGGGAAGTCCATT

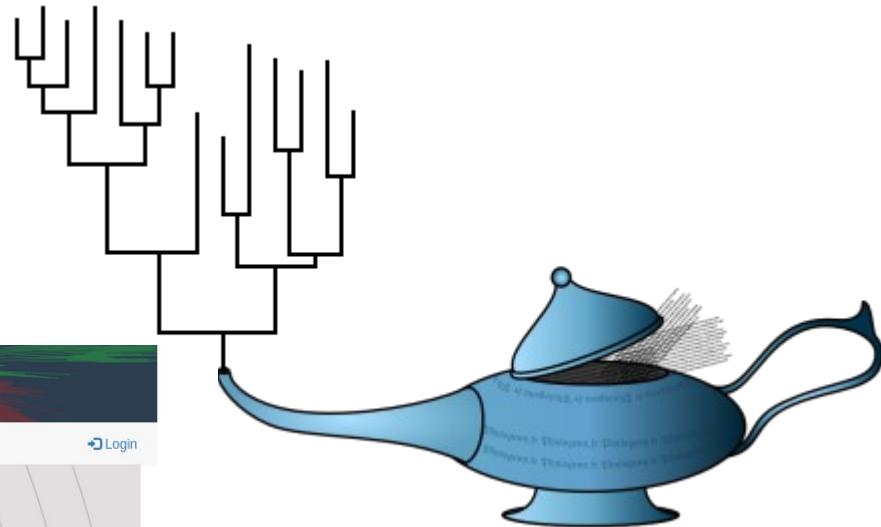
ngPhylogeny.fr

Home Phylogeny Analysis Tools Workspace (0) Documentation About Login

Robust phylogenetic analysis for everyone.

Free, simple to use web service dedicated to reconstructing and analysing phylogenetic relationships between molecular sequences.

Let's GO ! with One Click Workflow



## » One Click

Fully automatic workflow  
Default tools + default parameters.



## » Advanced

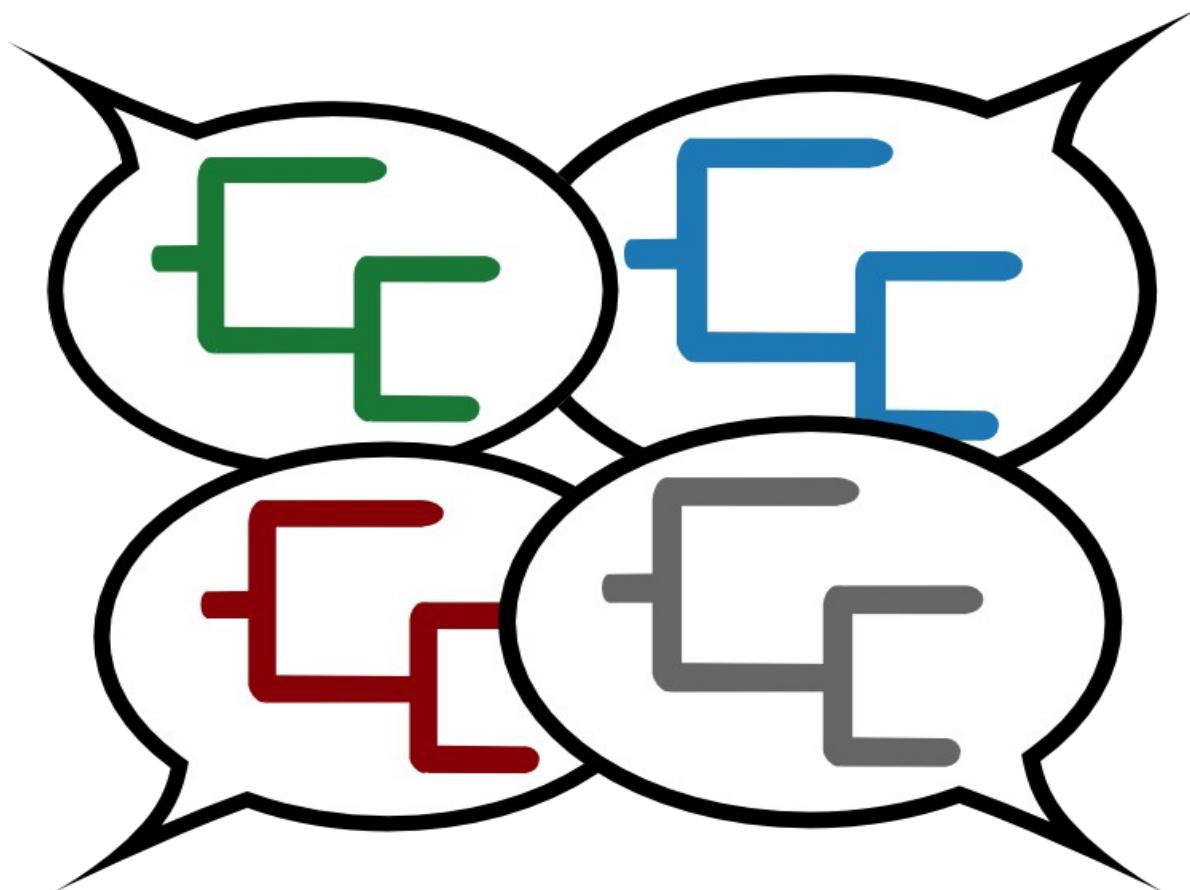
Semi automatic workflow  
Default tools + custom parameters.



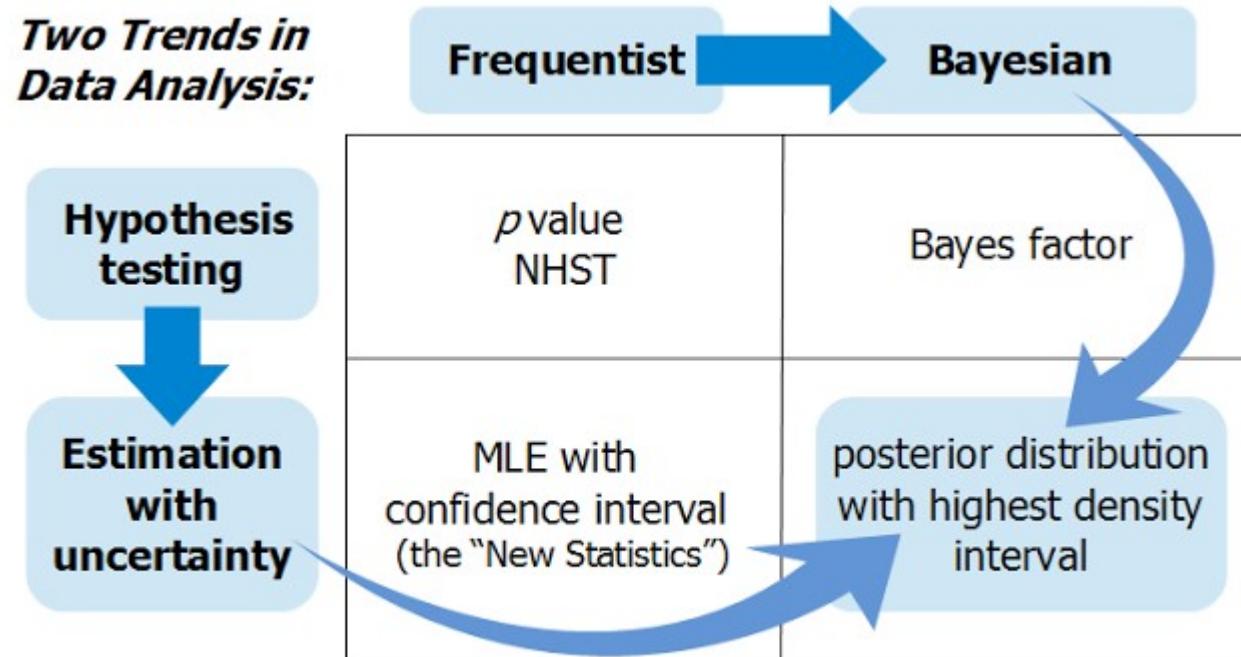
## » A la Carte

Custom workflow  
Custom tools + Custom parameters.

## Hypothesis testing on gene family trees

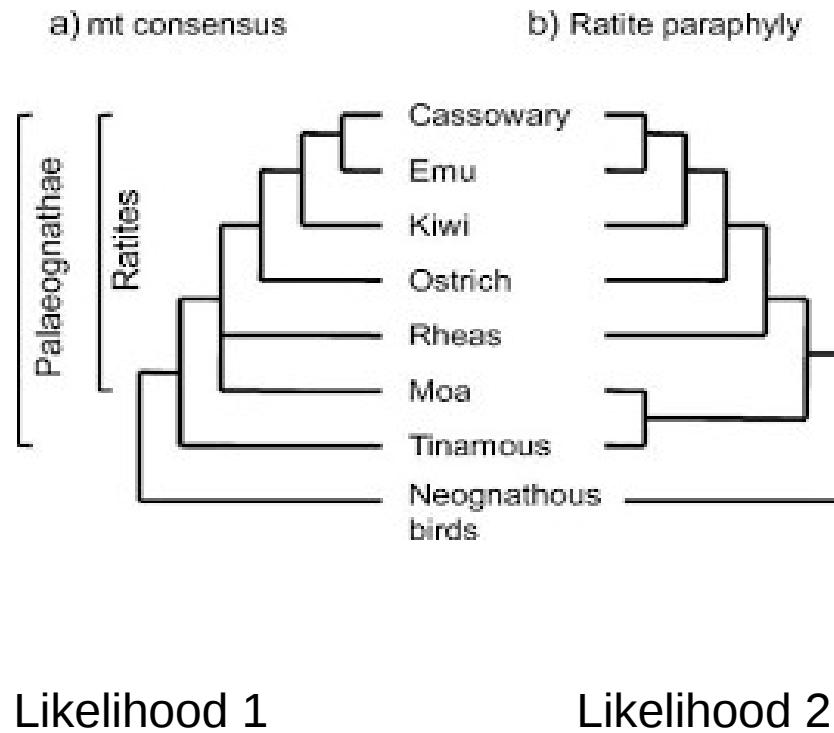


## Probabilistic methods render themselves for testing



Copyright © 2015 John K. Kruschke

# Comparing (testing) alternative topologies



$$AIC = 2k - 2 \ln(\hat{L})$$

Number of model parameters  
Maximum likelihood value

The Akaike information criterion (AIC) is an [estimator](#) of the relative quality of [statistical models](#) for a given set of data. Given a collection of models for the data, AIC estimates the quality of each model, relative to each of the other models. Thus, AIC provides a means for [model selection](#).

Given a set of model, we can use AIC to choose the one with the minimum value, achieves Better likelihood with the least number of parameters

# **Inferring evolutionary events**

## **Speciations/Duplications.**

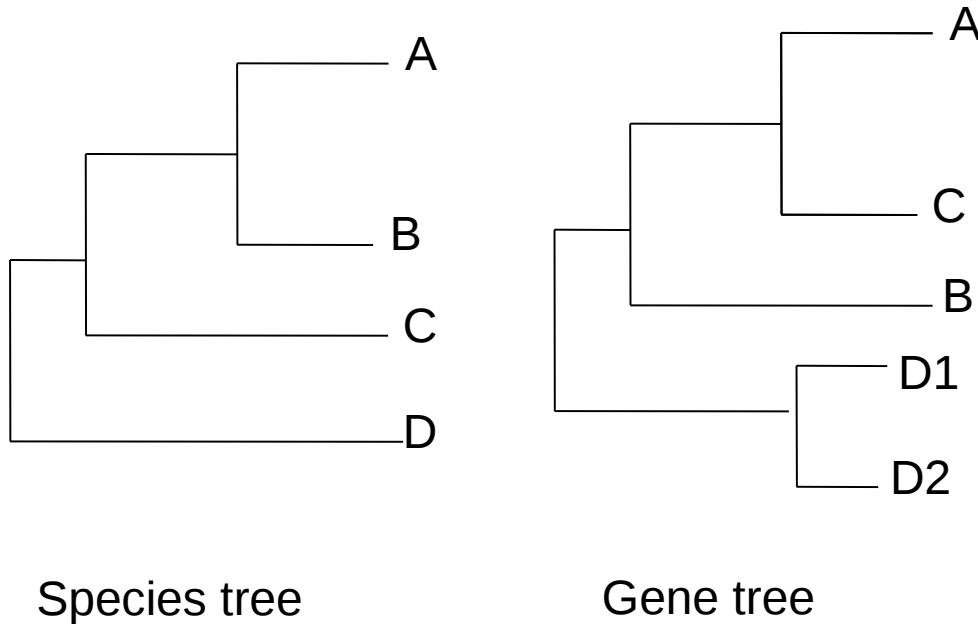
**1- Reconciliation algorithms**

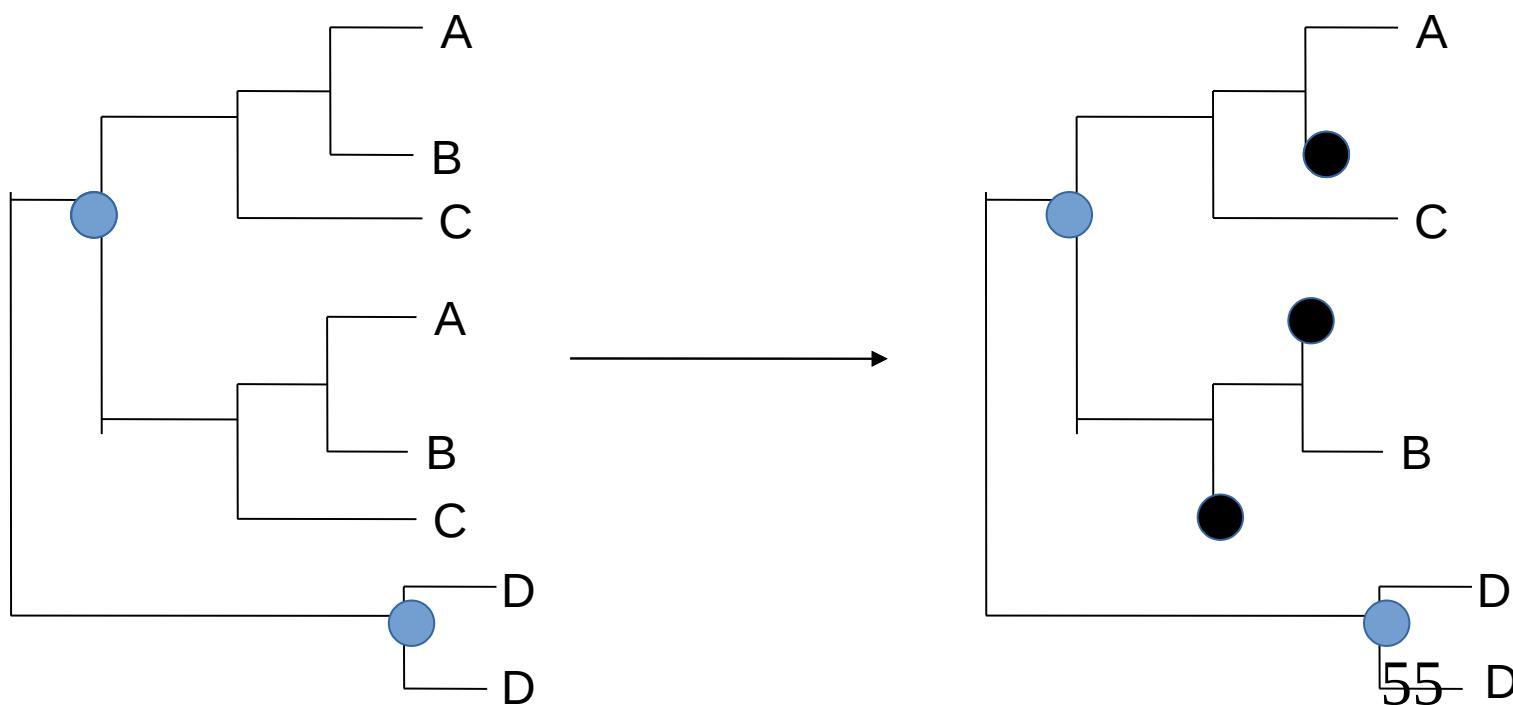
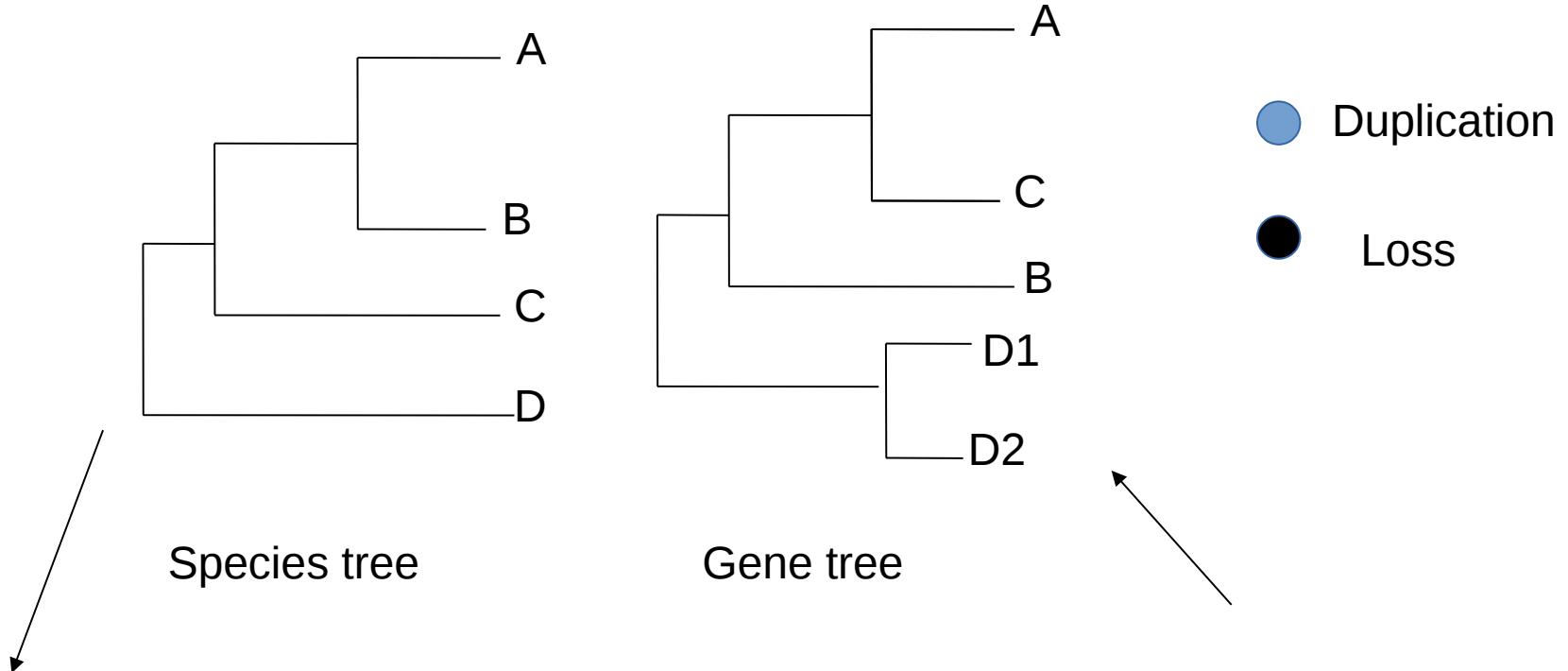
**2- Species overlap**

## Reconciliation algorithm.

**(Hard reconciliation)** Resolve any incongruence between gene tree and species tree by introducing the minimal number of gene duplications and losses.

**(Soft reconciliation)** Allow incongruences below a given support value

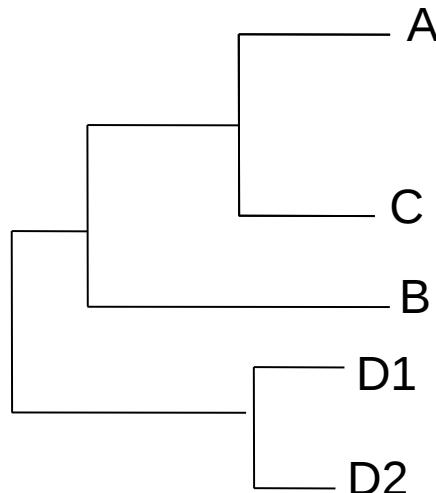




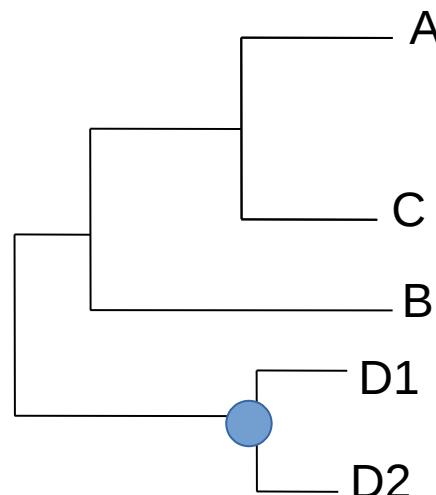
## Species overlap algorithm.

It does not require a species-tree but needs to know the species to which  
The genes belong  
In essence can be seen as a reconciliation with an unresolved species tree

For every node in the gene tree evaluate whether the daughter partitions  
share any species. If the overlap (number of species shared over total  
number of species ) is higher than the given threshold. Inpute a  
duplication at that node.



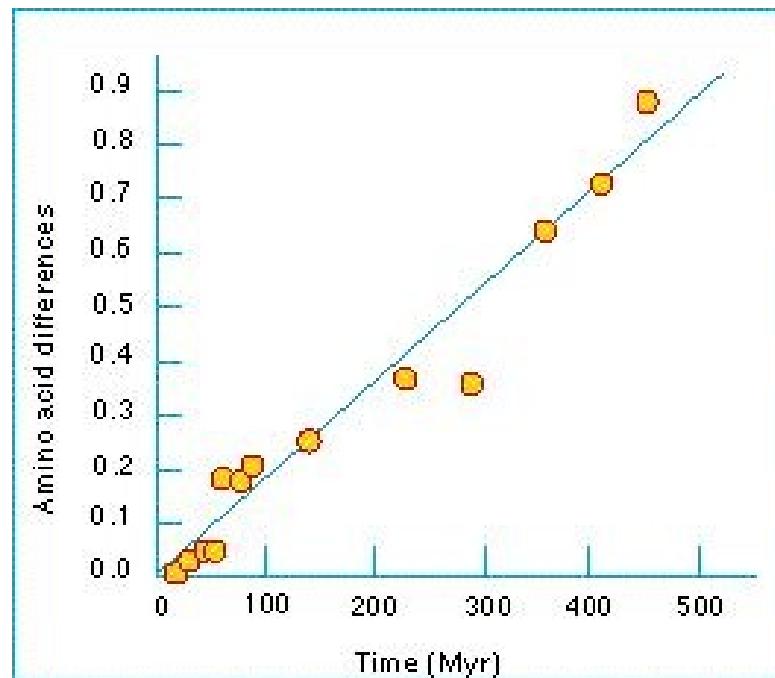
Gene tree



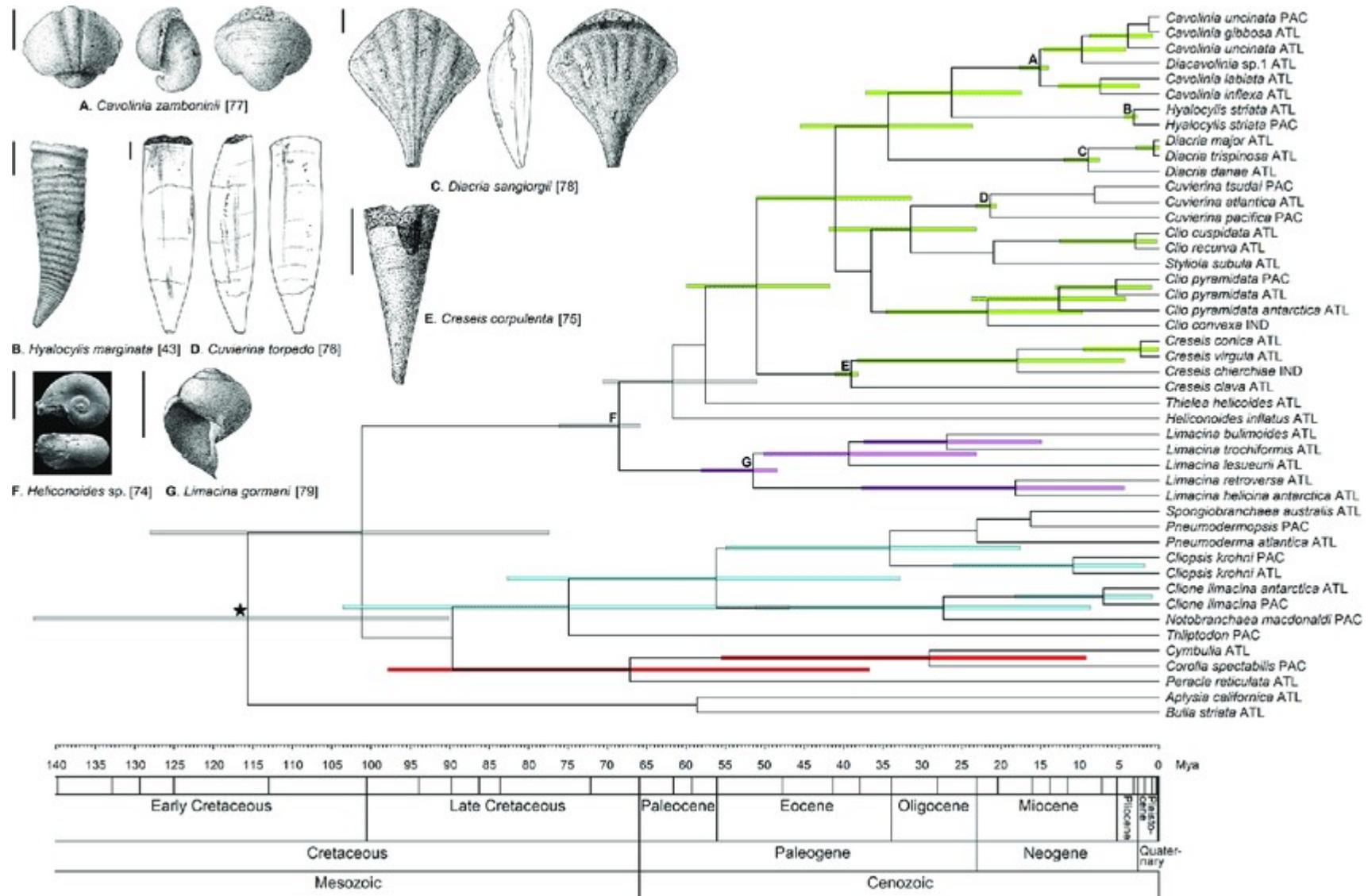
## Inferring relative timing of speciations and duplications

A **molecular clock** is said to exist when substitutions accumulate linearly with time. This assumption is violated most of the time, particularly at long evolutionary distances.

For a constant mutation rate, neutral substitutions are expected to behave more clock-like than non-neutral substitutions, that is why they are generally used as a **proxy** for time



Fossils can be used to calibrate timing of speciation nodes in a phylogeny,  
But fossil dating is also subject to errors



What about speciations and duplication nodes in a gene tree?

Genes do not have equal evolutionary rates across a genome

Synonymous substitutions still used as a proxy. However variations in GC content, Local variations in mutation rates, or repair can introduce noise

## A Nonsynonymous / Synonymous substitution

TCC	GAT	ATATGG	CAA	CCC	GAC	AAA	
S	D	I	W	Q	P	D	K

TCA	GAT	CTATGG	CAG	CCC	CAC	AAA	
S	D	L	W	Q	P	R	K

## B Radical / Conservative substitution

ATT	GACTATTCC	TGT	TGGTTT	GAA	CCAGGC	AGA				
I	D <sup>-</sup>	Y	S	C <sup>N</sup>	W	F	E <sup>-</sup>	P	G	R <sup>+</sup>

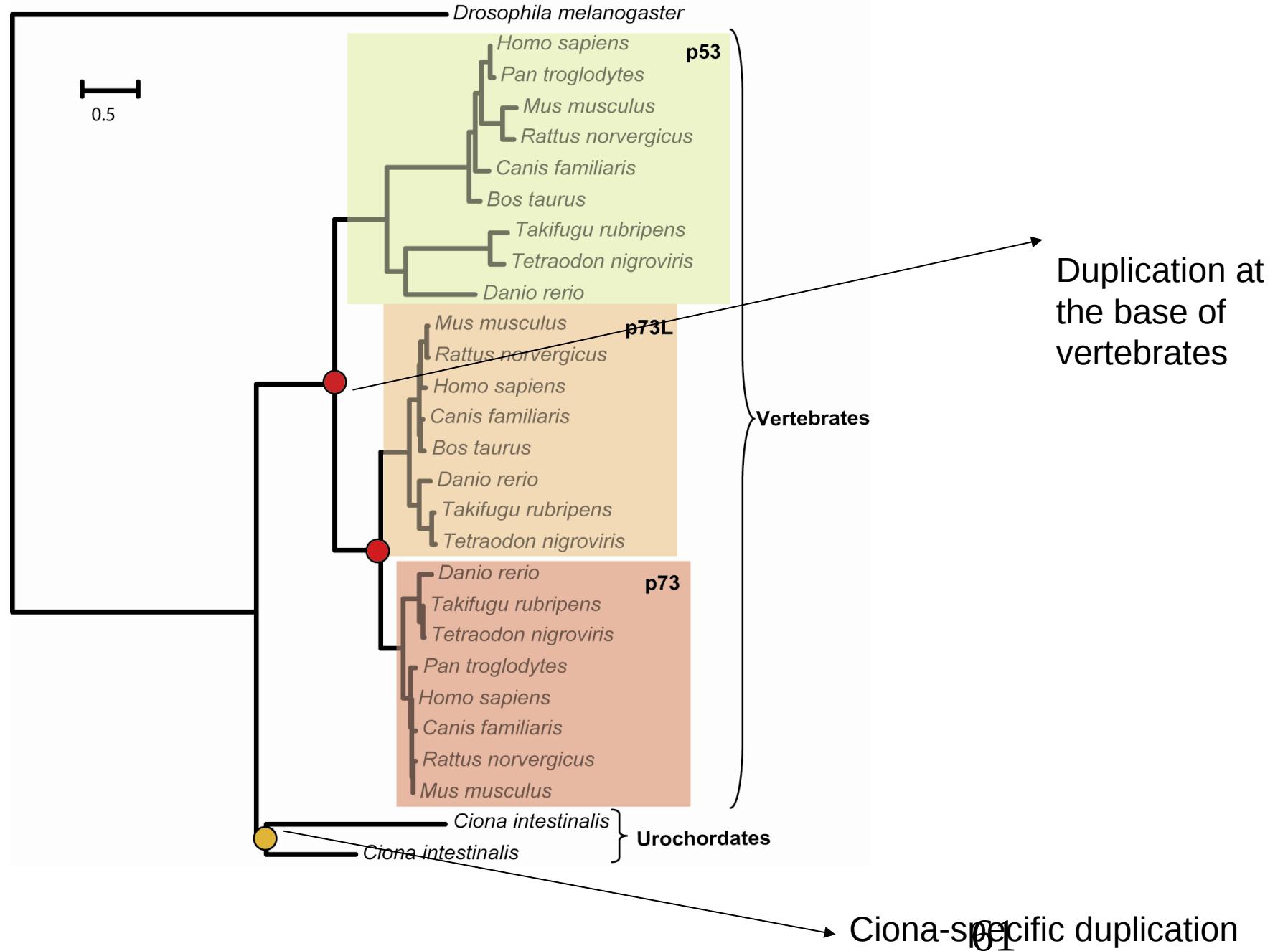
ATT	CACTACTCC	GGT	TGGTTTC	GCAC	CCAGGA	AAA				
I	R <sup>+</sup>	Y	S	G <sup>N</sup>	W	F	A <sup>N</sup>	P	G	K <sup>+</sup>

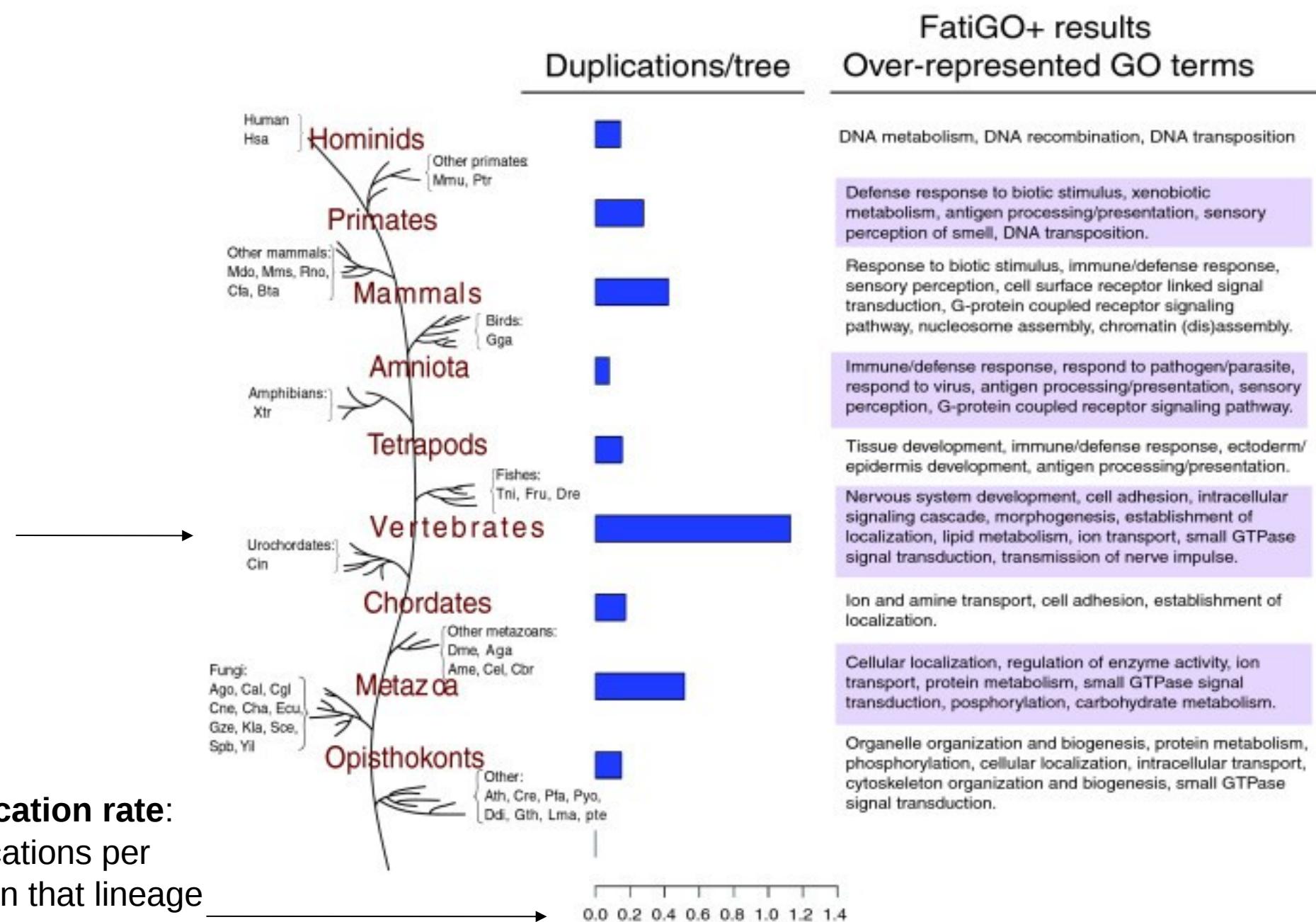
+ positive

- negative

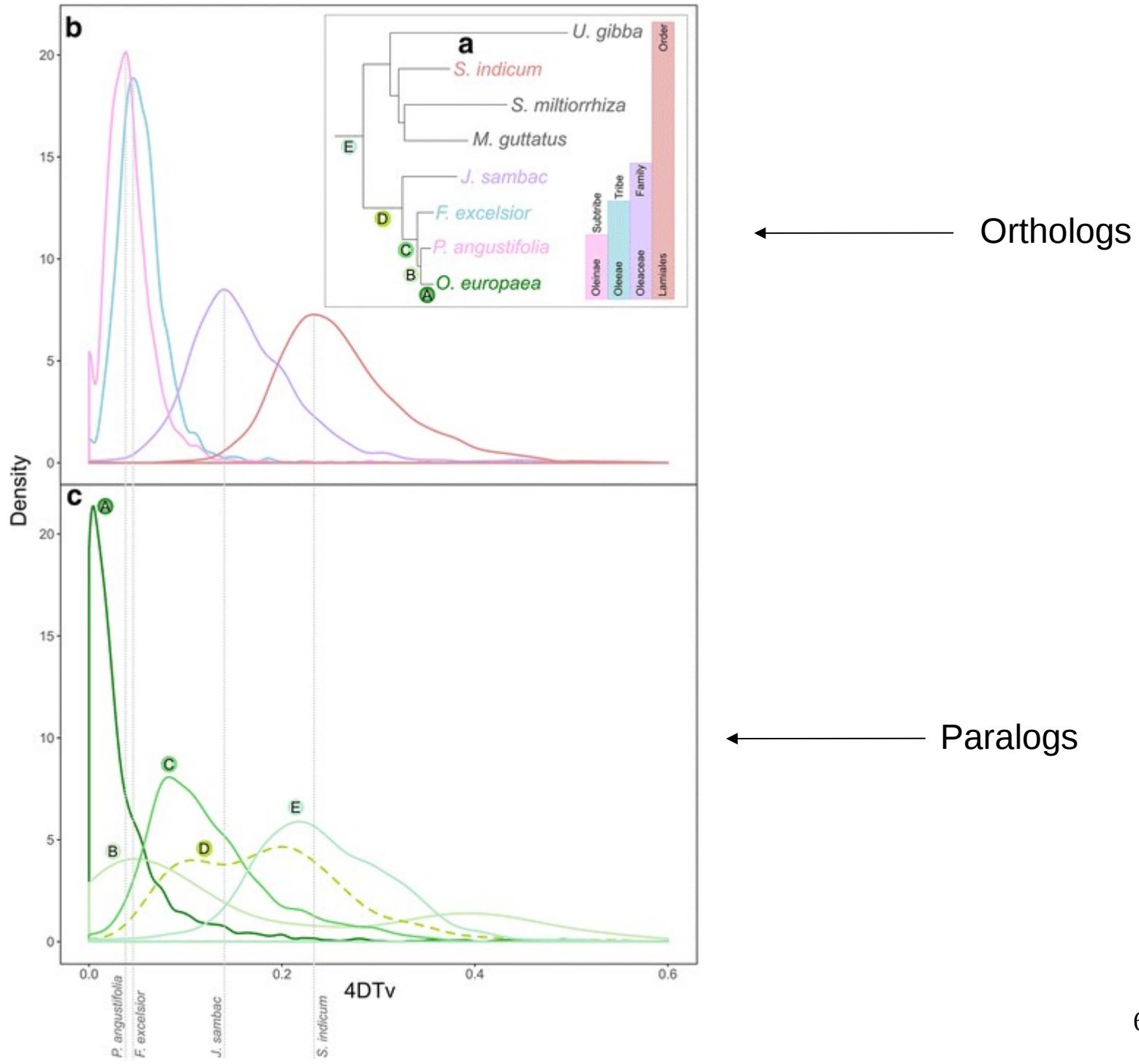
N neutral

Ks = rate of synonymous substitutions, 4Dtv = rate at **four fold degenerate sites**





Huerta-Cepas et. al. (2007)



If non-synonymous mutations are the result of selection, can we use them to detect positive selection, functional shifts, etc?

### A Nonsynonymous / Synonymous substitution

TCC	GAT	ATATGG	CAA	CCC	GAC	AAA	
S	D	I	W	Q	P	D	K

TCAG	ATCT	ATGG	CAG	CCC	CAC	AAA	
S	D	L	W	Q	P	R	K

### B Radical / Conservative substitution

ATT	<u>GAC</u>	TATTCC	<u>TGT</u>	TGGTTT	<u>GAA</u>	CCAGGC	<u>AGA</u>			
I	D-	Y	S	<u>C<sup>N</sup></u>	W	F	<u>E-</u>	P	G	<u>R<sup>+</sup></u>

ATT	<u>CAC</u>	TACTCC	<u>GGT</u>	TGGTTG	<u>GCA</u>	CCAGGA	<u>AAA</u>			
I	<u>R<sup>+</sup></u>	Y	S	<u>G<sup>N</sup></u>	W	F	<u>A<sup>N</sup></u>	P	G	<u>K<sup>+</sup></u>

+ positive

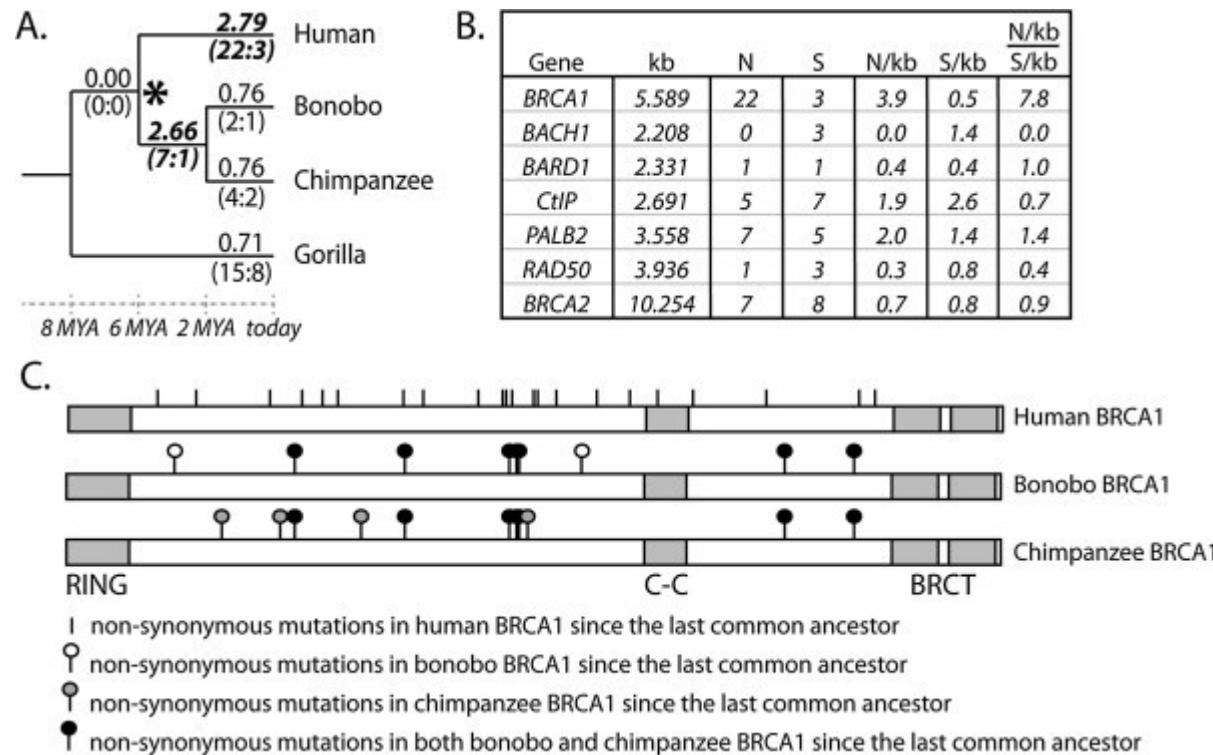
- negative

N neutral

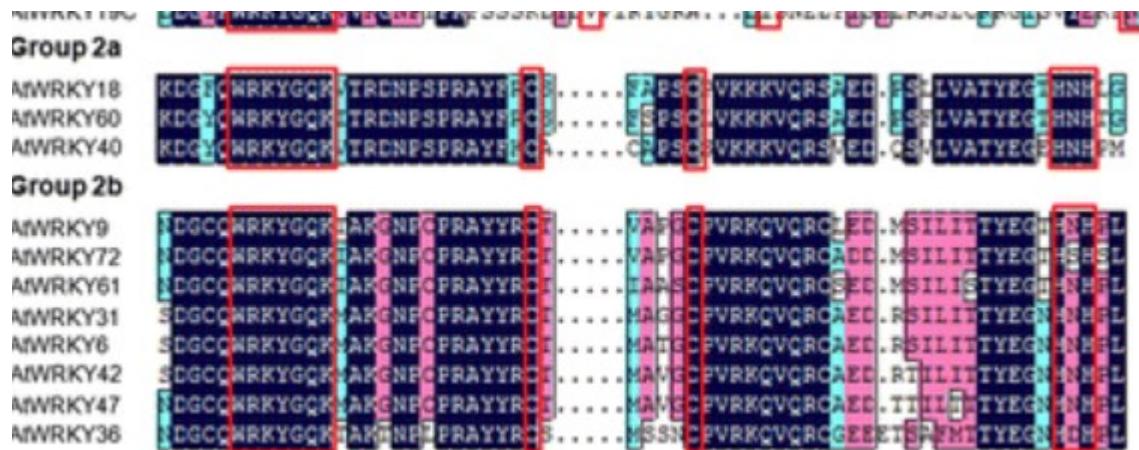
dN/dS ( )

The ratio of non-synonymous to synonymous substitutions ( $dN/dS$ ) is a useful measure of the strength and mode of natural selection acting on protein-coding genes.

Popular ML programs (i.e. PAML) can compute dN/dS ratios per branch, and per site



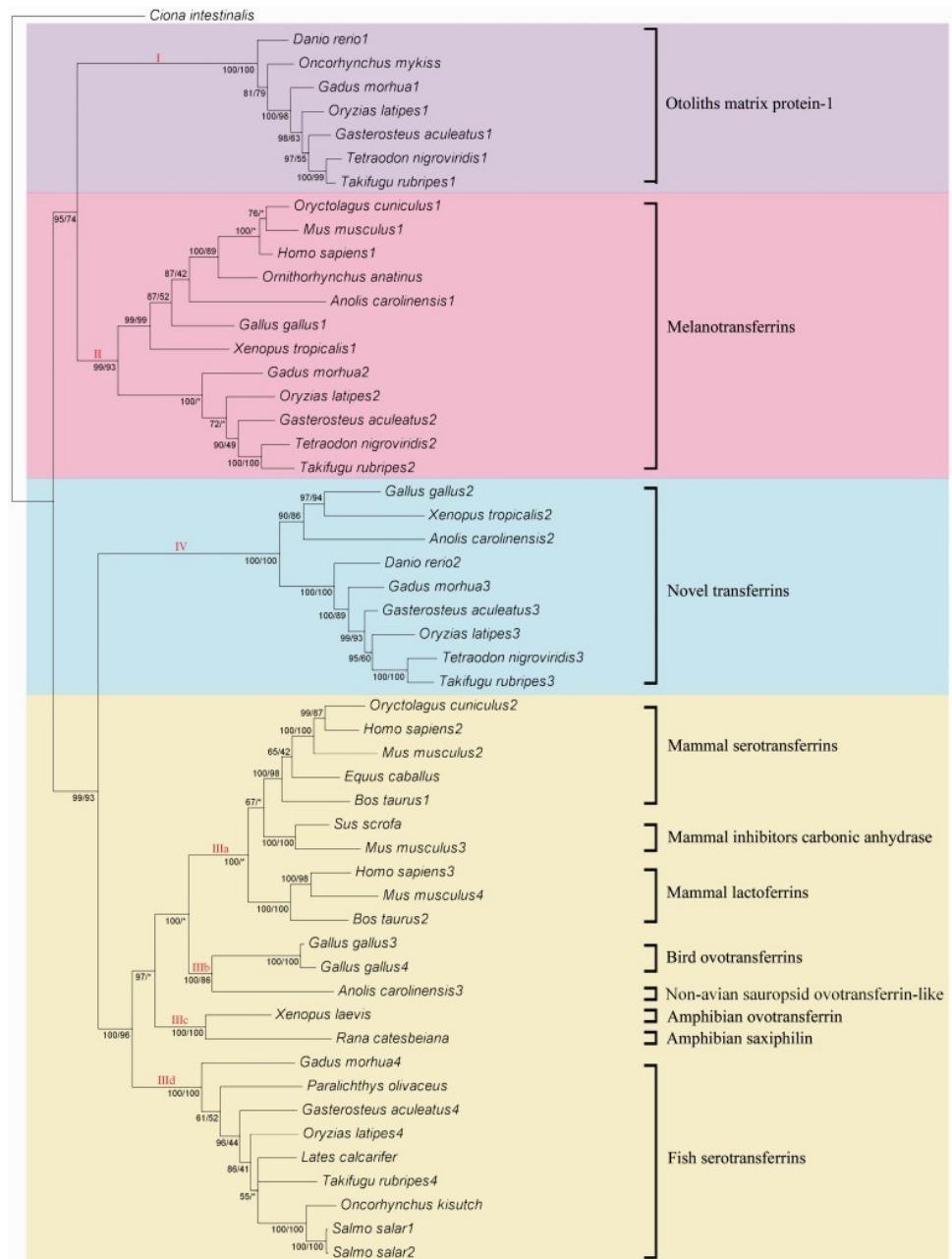
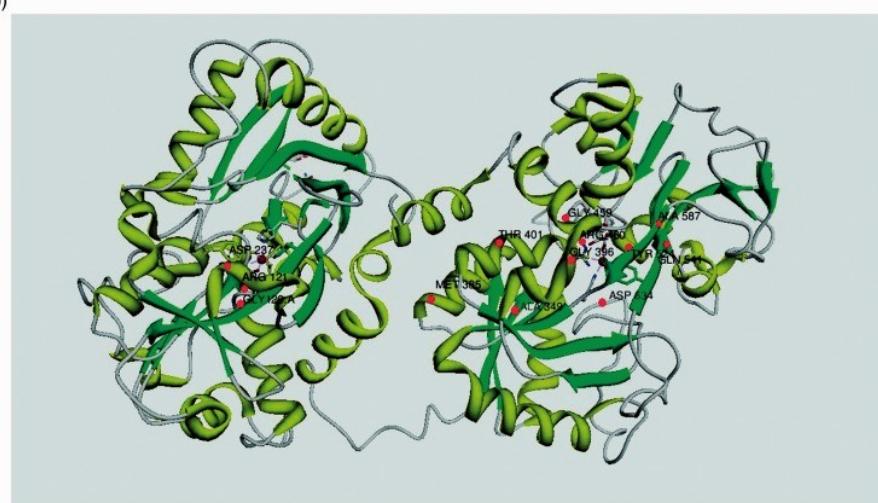
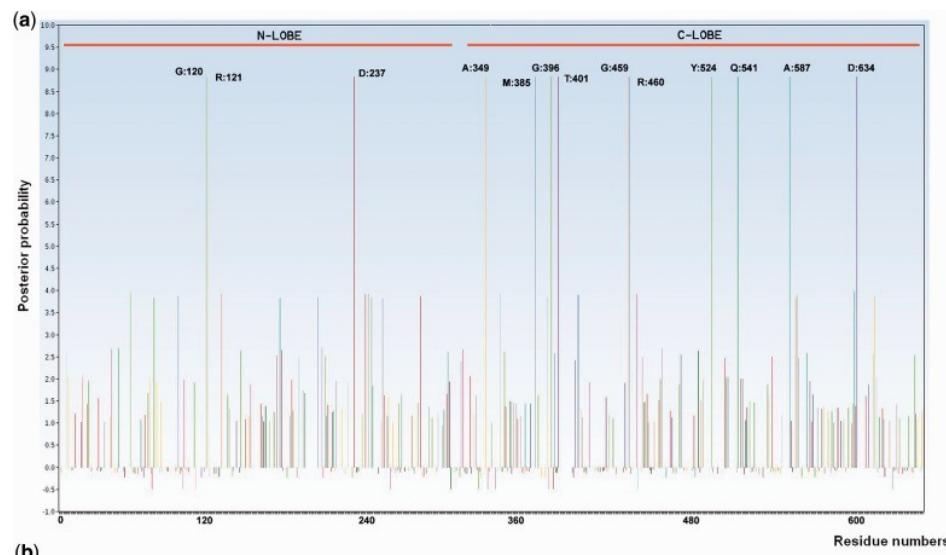
Detecting radical amino acid changes in a phylogeny can be indicative of functional shift (Program DIVERGE2)



# Phylogenetic Analyses Uncover a Novel Clade of Transferrin in Nonmammalian Vertebrates

Hirzahida Mohd-Padil,<sup>✉1</sup> Adura Mohd-Adnan,<sup>\*,1</sup> and Toni Gabaldón<sup>\*,2,3</sup>

[Author information](#) ► [Copyright and License information](#) ► [Disclaimer](#)



Ancestral Sequence Ressurection: directly testing for ancestral protein functions

Published online 5 September 2002 | Nature | doi:10.1038/news020902-7

News

## Triassic reptile saw red

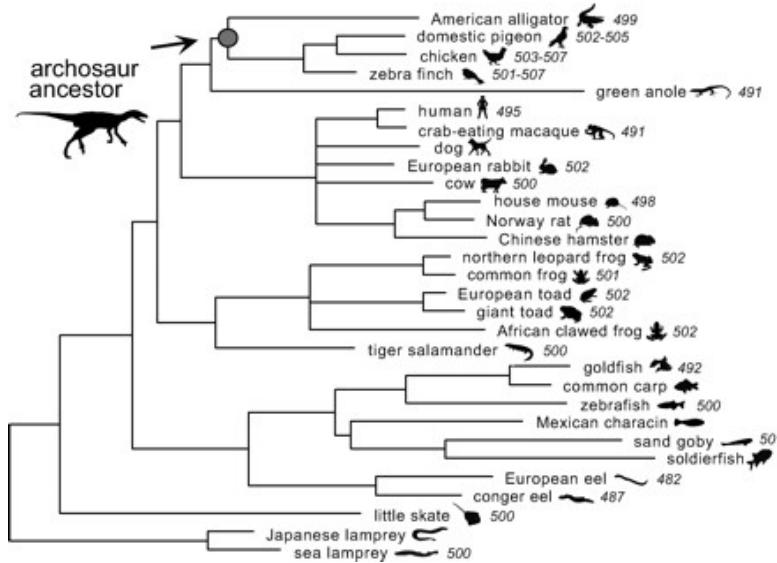
**Resurrected protein suggests that crocodiles' ancestors roamed at night.**

Helen Pearson

A reptile from the Triassic period may have done its stalking at night. So suggest scientists who have resurrected a 240-million-year-old eye protein that sees dim light<sup>1</sup>.

Such a molecule may have been found in the eyes of the earliest archosaurs, which were predecessors of the dinosaurs. Similar proteins, called rhodopsins, perceive low levels of light in humans and other animals.

Thomas Sakmar of Rockefeller University in New York and his colleagues used a computer program to extrapolate the DNA sequence of the ancient rhodopsin from known sequences in alligator, birds, frogs and fish.



Gene reconstruction gives researchers a dim view of the distant past.

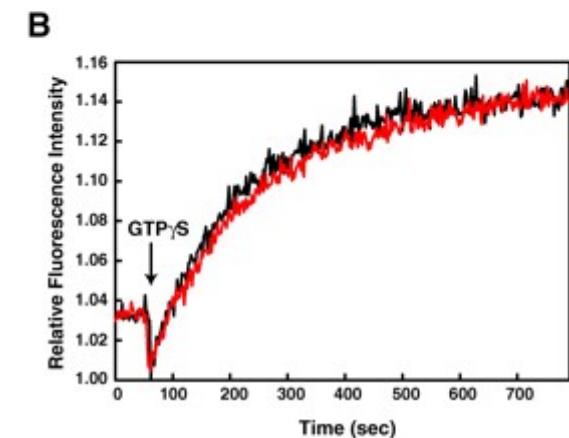
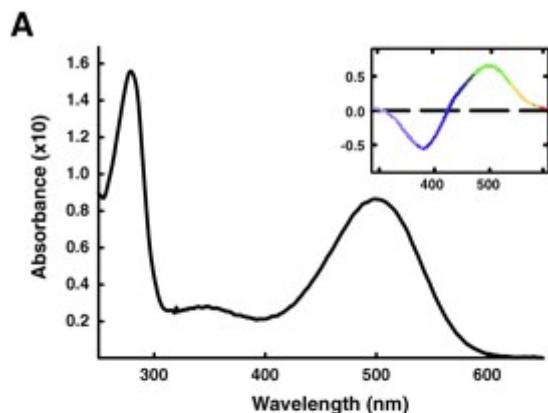
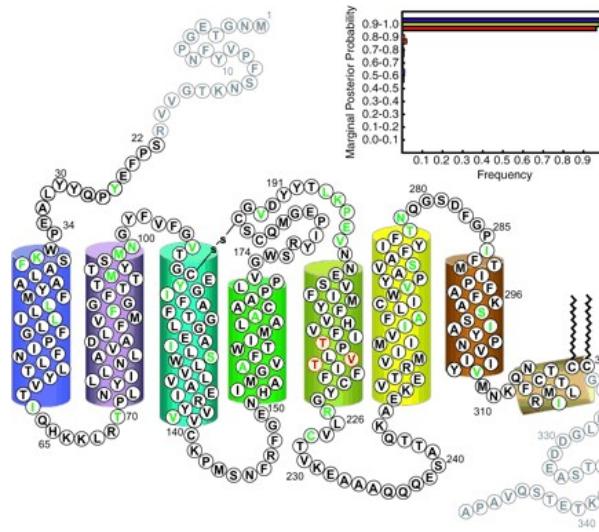
© SPL

# Recreating a Functional Ancestral Archosaur Visual Pigment FREE

Belinda S. W. Chang, Karolina Jönsson, Manija A. Kazmi, Michael J. Donoghue, Thomas P. Sakmar

*Molecular Biology and Evolution*, Volume 19, Issue 9, 1 September 2002, Pages 1483–1489,  
<https://doi.org/10.1093/oxfordjournals.molbev.a004211>

Published: 01 September 2002 Article history ▾



# **THANKS**