

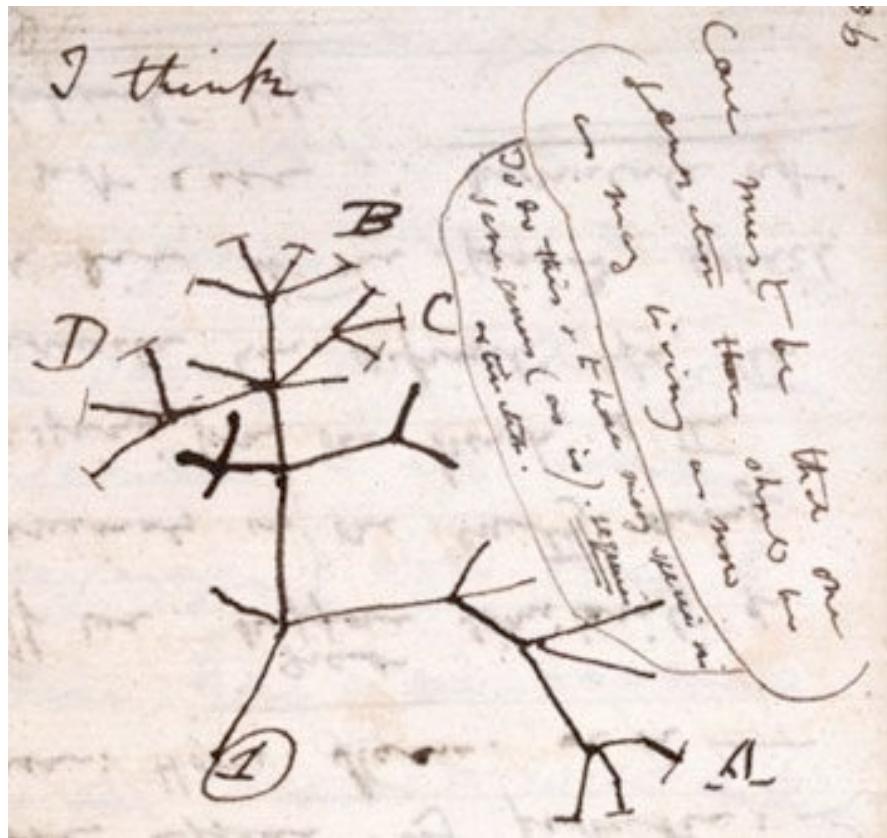
# **UBC Bioinformatics**

**Topic 9: Phylogenetic inference**

# Why do we need phylogenetics?

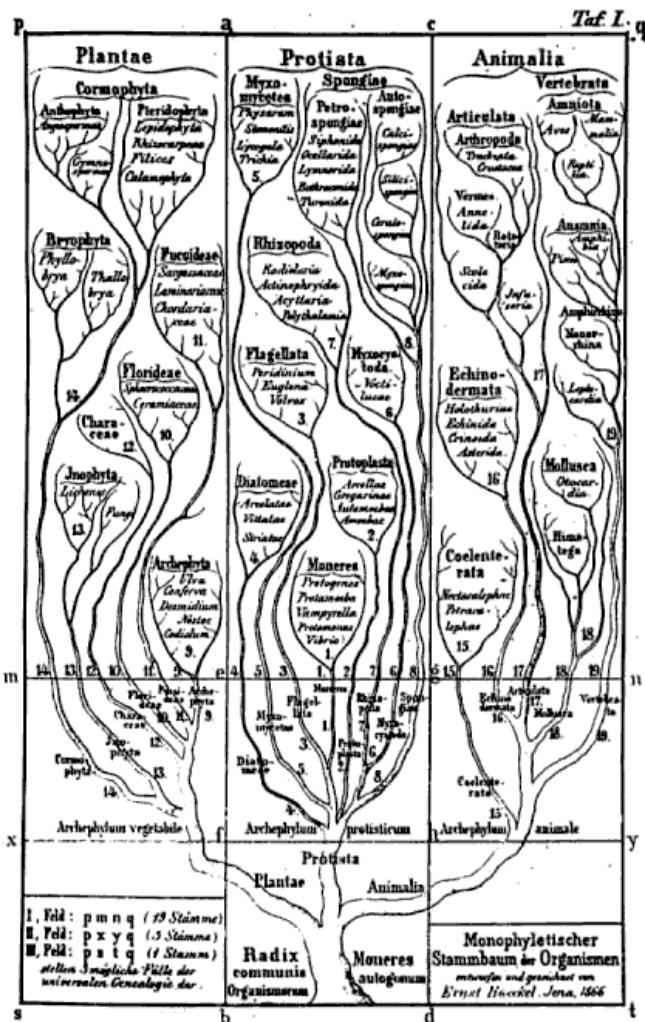
- Infer our evolutionary history
- Understand the diversity of life
- Determine how climate change may impact species distributions
- Determine how to develop drought/pest tolerant crops
- Development of new drugs and to synthesize materials
- Forensic scientists use phylogenetics in legal matters

# Thinking about life

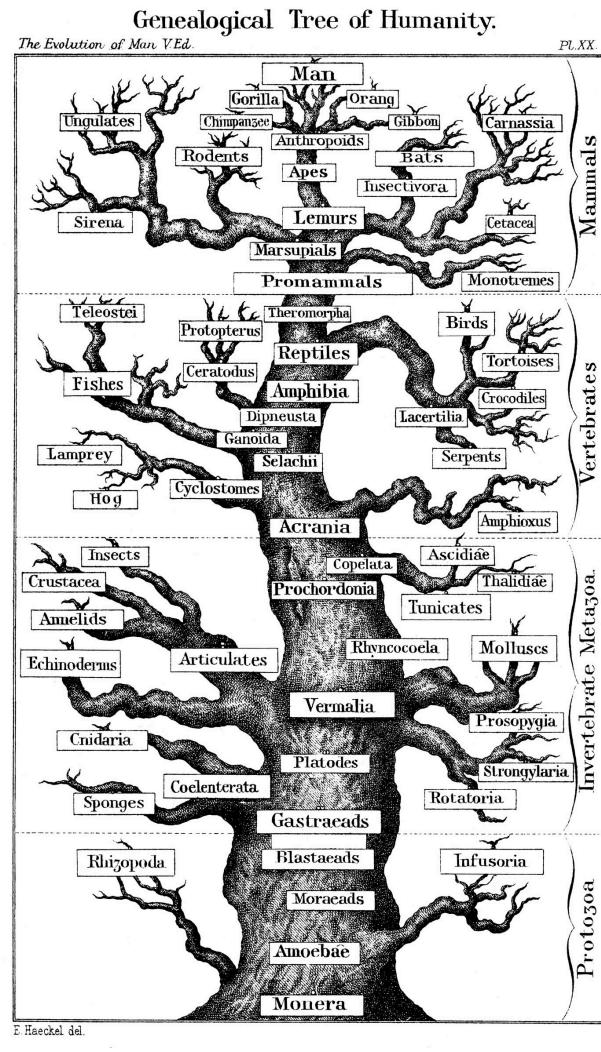


Darwin's “Transmutation of species” notebook (1837)

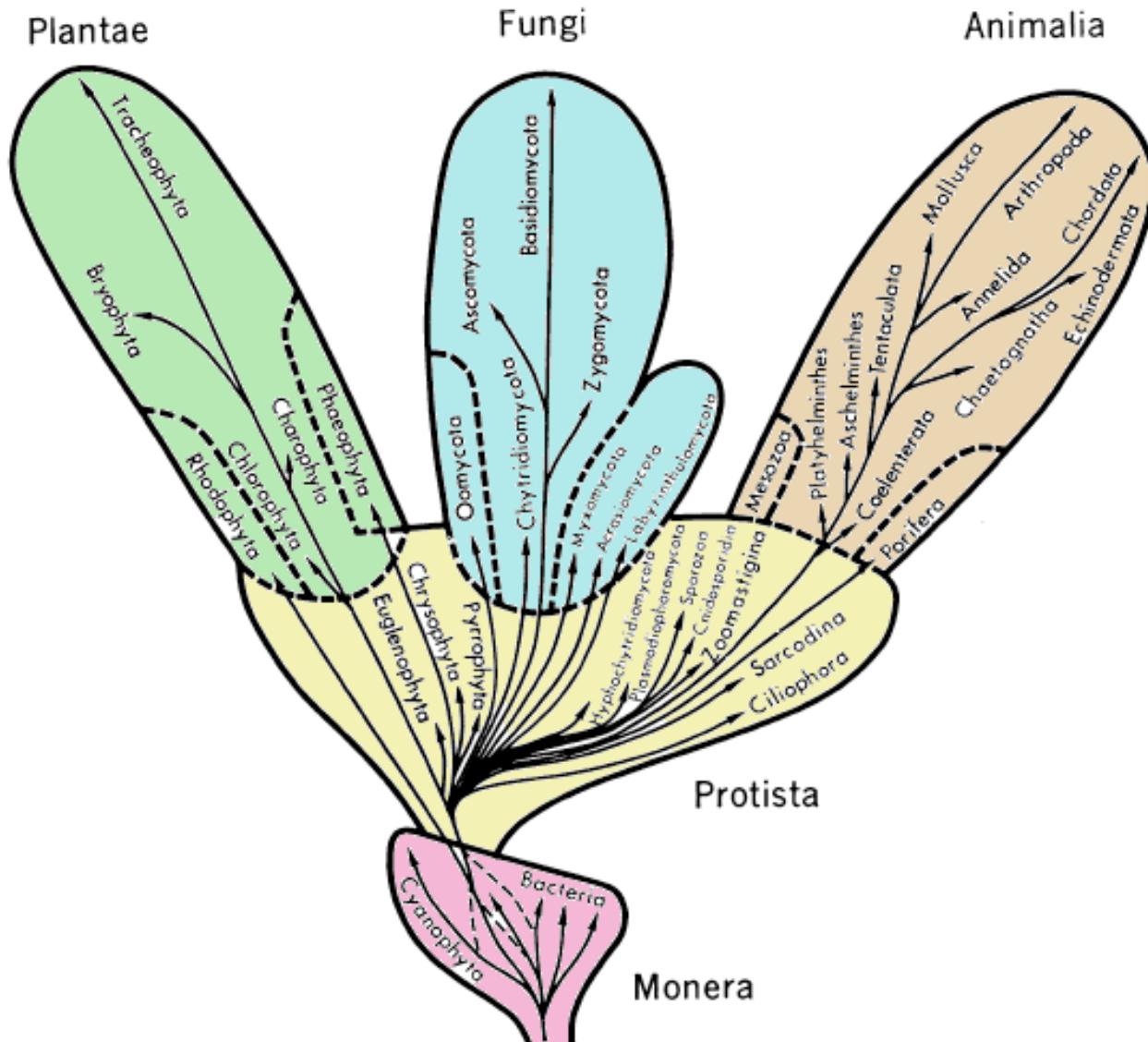
A phylogenetic view of biological diversification



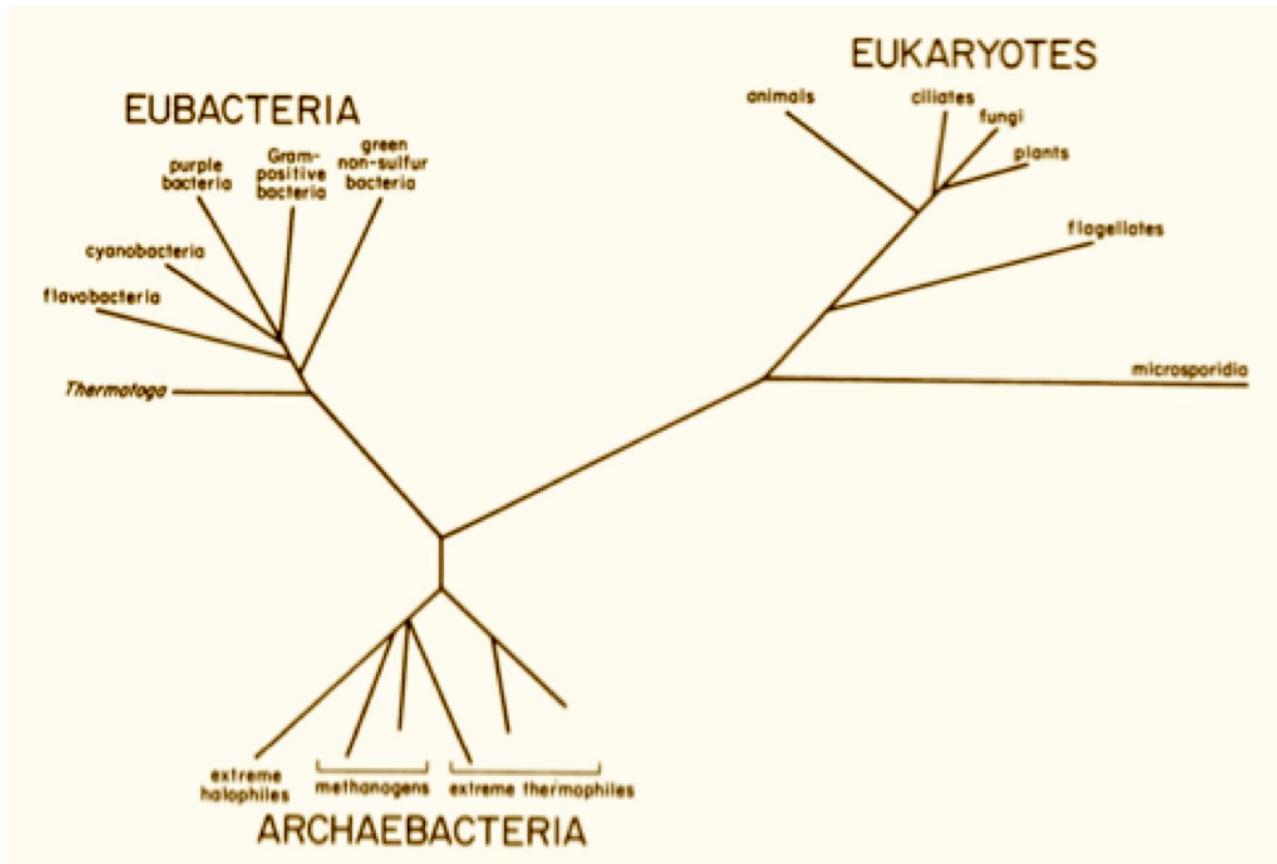
A phylogenetic view of biological diversification



# Whittaker's five kingdom system



# Carl Woese's three domain system (1977)



Analysis of ribosomal RNA sequences lead to Carl Woese to propose a three domain system (1977; PNAS 1990)

# **Basic components of phylogenetic analysis**

- All phylogenetic analyses include four basic steps:
  - 1) data acquisition

# **Basic components of phylogenetic analysis**

- All phylogenetic analyses include four basic steps:
  - 1) data acquisition
  - 2) construction of character matrix (e.g. sequence alignment)

# **Basic components of phylogenetic analysis**

- All phylogenetic analyses include four basic steps:
  - 1) data acquisition
  - 2) construction of character matrix (e.g. sequence alignment)
  - 3) phylogeny estimation

# **Basic components of phylogenetic analysis**

- All phylogenetic analyses include four basic steps:
  - 1) data acquisition
  - 2) construction of character matrix (e.g. sequence alignment)
  - 3) phylogeny estimation
  - 4) interpretation of phylogenetic inference.

# Sequence archives

The mission of UniProt is to provide the sequence and functional information.

**Pfam 27.0 (March 2013, 14831 families)**

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. [More...](#)

UniProtKB	UniRef	Sequence clust	QUICK LINKS	YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...
Swiss-Prot (547,357)	Sequence search	Sequence clust	<a href="#">SEQUENCE SEARCH</a>	Analyze your protein sequence for Pfam matches
Manually annotated and reviewed.	View a Pfam family	View a Pfam family	<a href="#">VIEW A PFAM FAMILY</a>	View Pfam family annotation and alignments
TREMBL (89,451,166)	View a clan	View a clan	<a href="#">VIEW A CLAN</a>	See groups of related families
Automatically annotated and not reviewed.	View a sequence	View a sequence	<a href="#">VIEW A SEQUENCE</a>	Look at the domain organisation of a protein sequence
	View a structure	View a structure	<a href="#">VIEW A STRUCTURE</a>	Find the domains on a PDB structure
	Keyword search	Keyword search	<a href="#">KEYWORD SEARCH</a>	Query Pfam by keywords
	Jump to	Jump to	<b>JUMP TO</b>	<input type="text" value="enter any accession or ID"/> <b>Go</b> <b>Example</b>

Enter any type of accession or ID to jump to the page for a Pfam family or clan,



# Sequence archives

UniProtKB

BLAST Align Retrieve/ID Mapping

The mission of UniProt is to provide the sequence and functional information.

UniProtKB

Swiss-Prot (547,357)  
Manually annotated and reviewed.

TrEMBL (89,451,166)  
Automatically annotated and not reviewed.

UniRef

Sequence cluster

TreeFam

Home Search Browse Download Help Forum

Pfam 27.0 (March 2013, 14,400 families)

The Pfam database is a large collection of protein alignments and hidden Markov models (HMMs).

e!Ensembl

BLAST/BLAT | BioMart | Tools | Downloads | Help & Documentation | Blog | Mirrors

Search TreeFam... Examples: BRCA2,ENSP00000428982, or do a sequence search

SPECIES IN TREEFAM

Species Tree used in TreeFam 9. See full tree [here](#).

Euarchontoglires  
Mammalia  
Vertebrata  
Chordata  
Metazoa  
Eukaryota

Primates  
Glires  
Laurasiatheria  
Frogs/Lizards/E  
Tunicates  
Arthropoda  
Nematoda

Variant Effect Predictor

Ve!P

Find SNPs and other variants for my gene

GIRATAACATTC  
CTTRAAGTCTT  
CTTCTAAATTGT

Compare genes across species

Insert into tree  
Insert Into TreeFam gene tree using

Advanced

Search

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Search

Go

Examples: BRCA2,ENSP00000428982, or do a sequence search

Variant Effect Predictor

Ve!P

Find SNPs and other variants for my gene

GIRATAACATTC  
CTTRAAGTCTT  
CTTCTAAATTGT

Compare genes across species

# **Constructing a distance matrix**

**Pairwise alignment** – based on similarity or distance

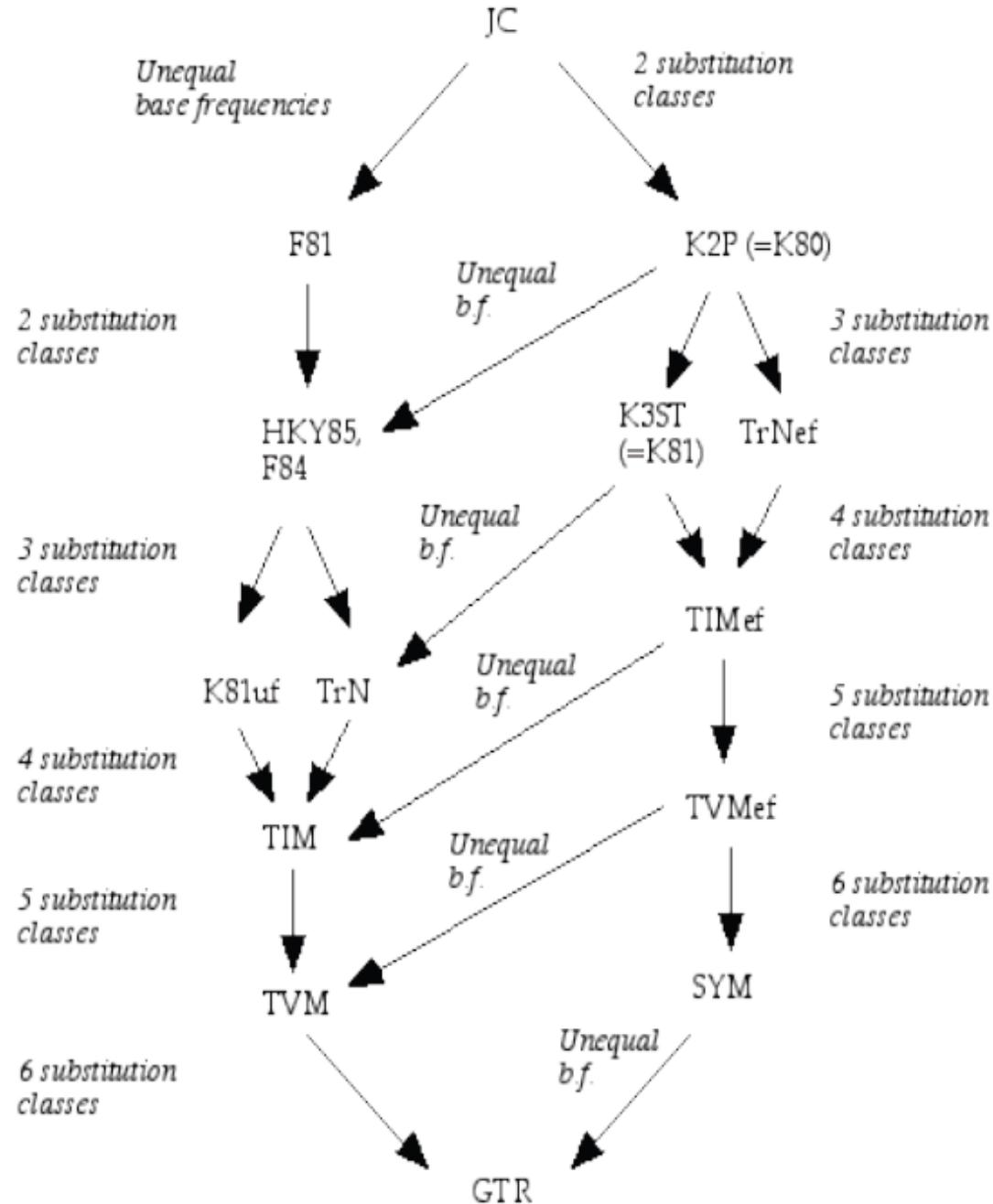
**Edit distance** – commonly, Hamming or Levenshtein distance between two strings

What about biological alphabets?  
Are all changes equally likely?

# Markov model of molecular evolution:

$$P_{ij}(t) = e^{Qt}$$

where  $Q$  is defined by one of the nucleotide substitution models shown to the right



# Jukes-Cantor (1969) described simplest Markov model for nucleotide substitutions

$\text{Pr}(\text{sequences differ at a position } - i \neq j) = 3/4 * (1 - e^{-4/3t})$

This implies equal probabilities for all substitutions ( $Q_{ij} = Q_{ji}$ )

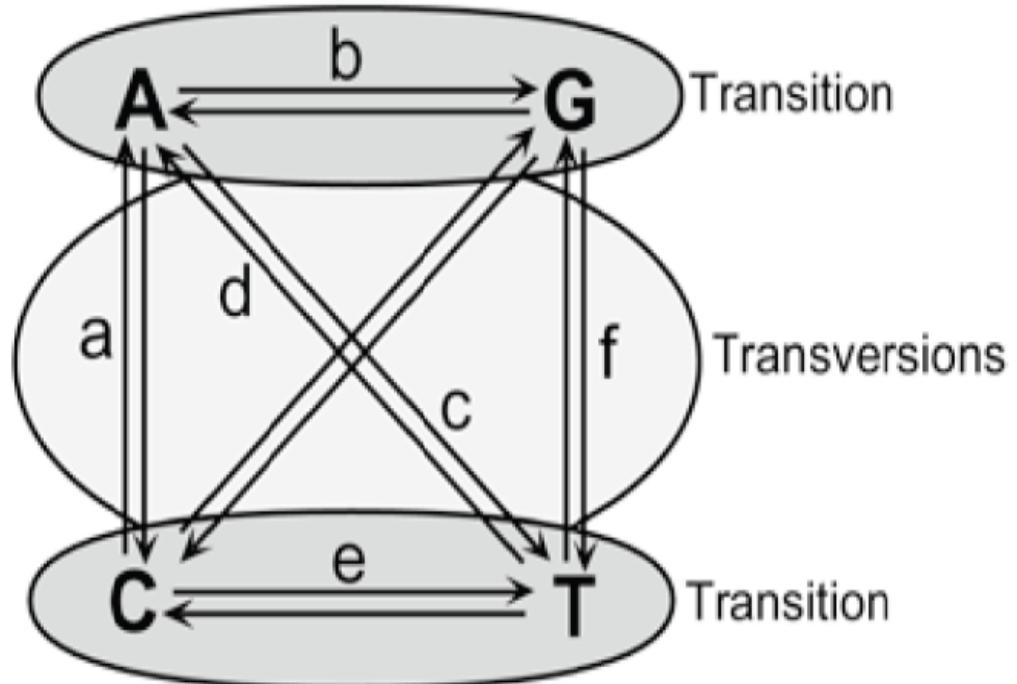
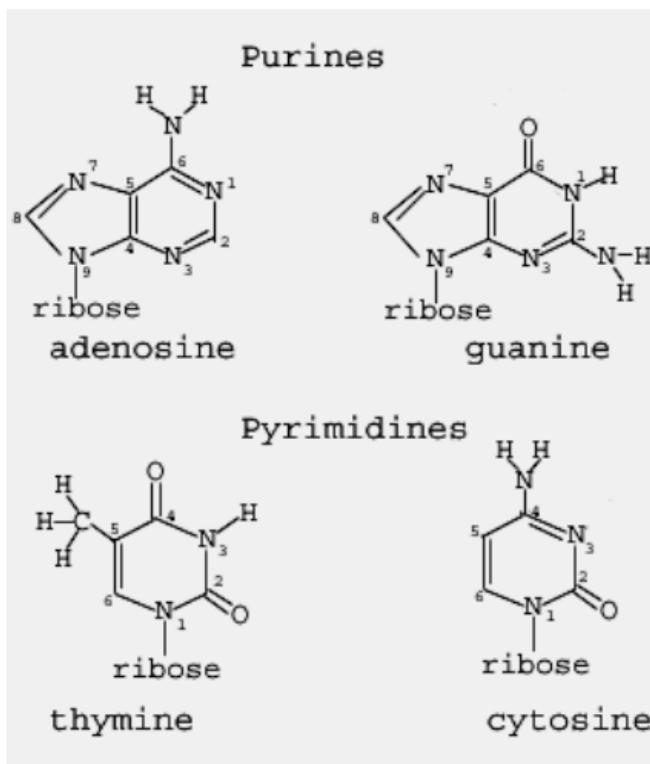
$$\mathcal{Q} = \begin{pmatrix} T & C & A & G \\ . & 1 & 1 & 1 \\ 1 & . & 1 & 1 \\ 1 & 1 & . & 1 \\ 1 & 1 & 1 & . \end{pmatrix}$$

To

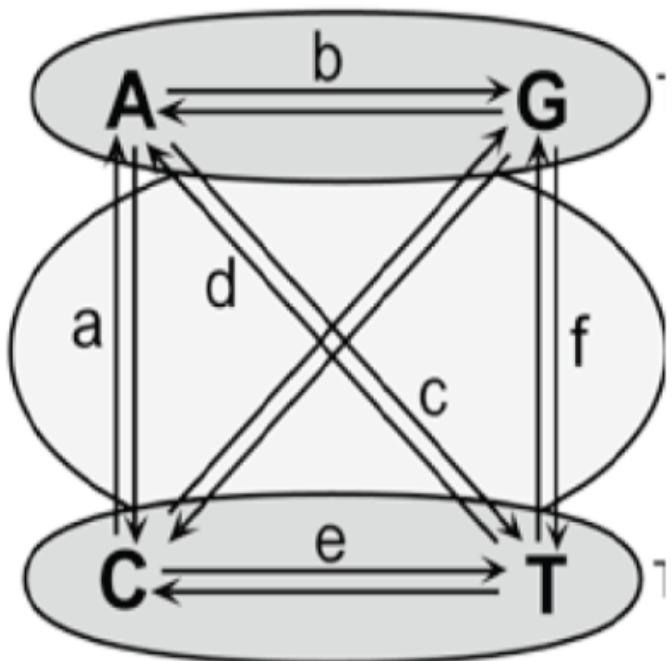
From

T  
C  
A  
G

# What about transition bias?



# Kimura (1980) modified the JC model



$$\mathcal{Q} = \begin{pmatrix} T & C & A & G \\ . & \kappa & 1 & 1 \\ \kappa & . & 1 & 1 \\ 1 & 1 & . & \kappa \\ 1 & 1 & \kappa & . \end{pmatrix} \begin{matrix} T \\ C \\ A \\ G \end{matrix}$$

# What about variation in nucleotide frequencies? – Felsenstein (1981)

$$\mathcal{Q} = \begin{pmatrix} T & C & A & G \\ \cdot & \pi_C & \pi_A & \pi_G \\ \pi_T & \cdot & \pi_A & \pi_G \\ \pi_T & \pi_C & \cdot & \pi_G \\ \pi_T & \pi_C & \pi_A & \cdot \end{pmatrix} \begin{matrix} T \\ C \\ A \\ G \end{matrix}$$

# Hasegawa, Kishino, and Yano (1985) put these together

$$\mathcal{Q} = \begin{pmatrix} T & C & A & G \\ \cdot & \kappa\pi_C & \pi_A & \pi_G \\ \kappa\pi_T & \cdot & \pi_A & \pi_G \\ \pi_T & \pi_C & \cdot & \kappa\pi_G \\ \pi_T & \pi_C & \kappa\pi_A & \cdot \end{pmatrix} \begin{matrix} T \\ C \\ A \\ G \end{matrix}$$

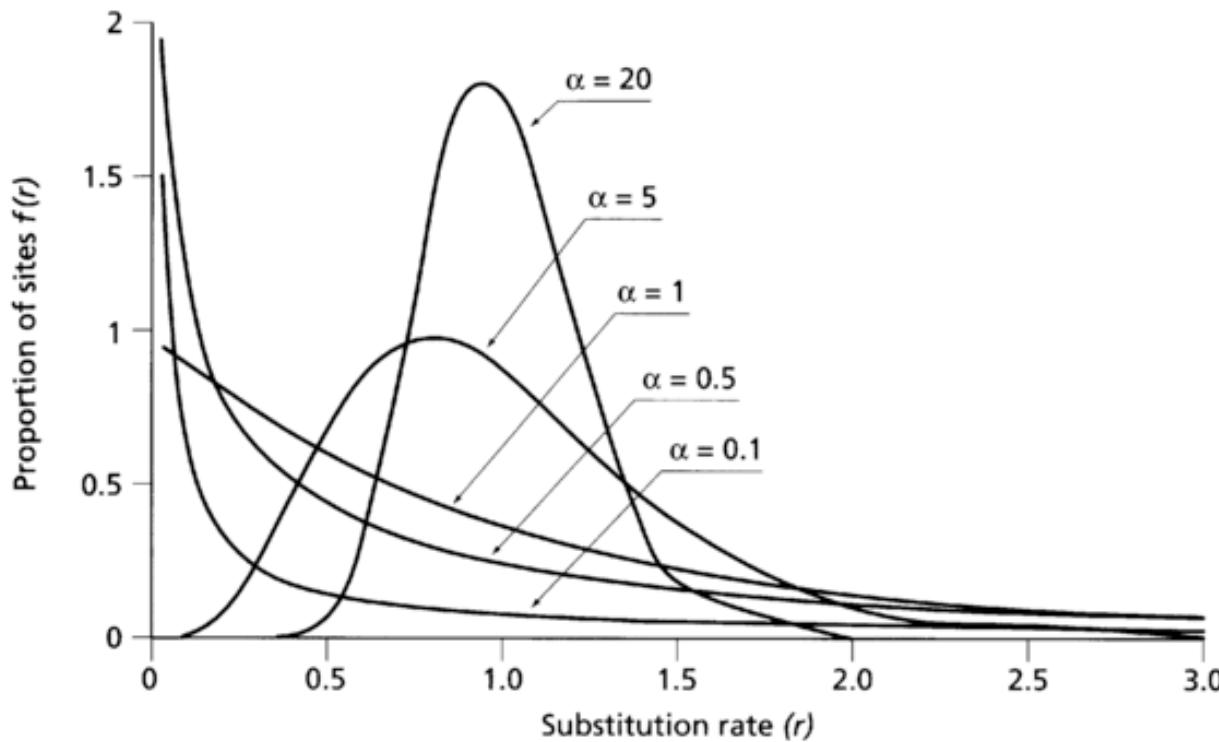
# General time reversible (GTR) model

– Tavaré 1986

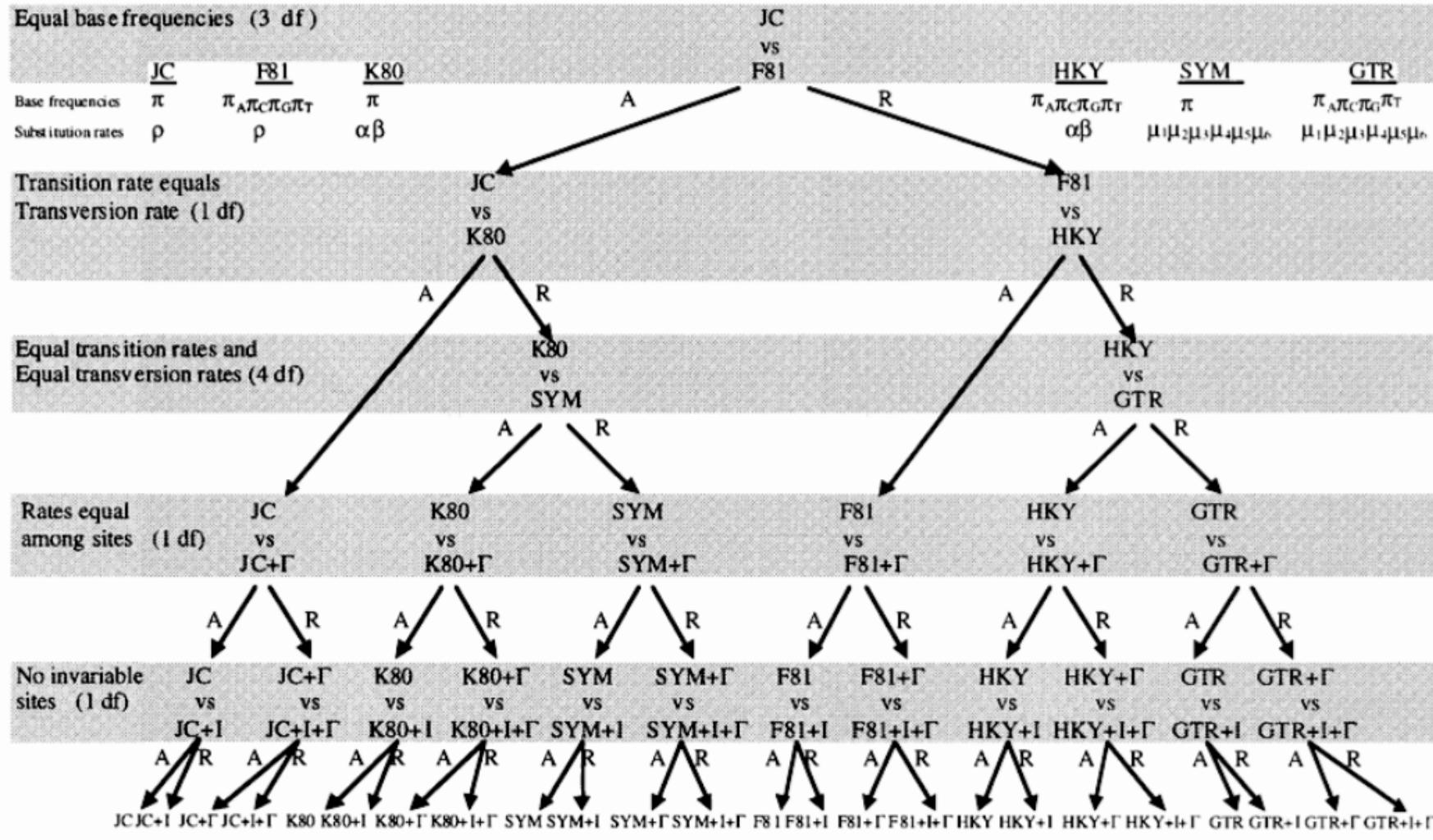
$$\mathcal{Q} = \begin{pmatrix} \cdot & a\pi_C & b\pi_A & c\pi_G \\ a\pi_T & \cdot & d\pi_A & e\pi_G \\ b\pi_T & d\pi_C & \cdot & \pi_G \\ c\pi_T & e\pi_C & \pi_A & \cdot \end{pmatrix} \begin{matrix} \text{T} \\ \text{C} \\ \text{A} \\ \text{G} \end{matrix}$$

# Can we relax the assumption of equal rates across all sites?

- Allow a certain proportion of invariable sites (I)
- Model a distribution of rates - gamma distribution



# Modeltest (Posada and Crandall 1998)



# Empirical estimation of AA transition frequencies

**PAM** (Point Accepted Mutation), **JTT** (Jones, Taylor, Thornton) and **WAG** (Whelan and Nick Goldman) matrices based on sequence comparisons in global alignments

**BLOSUM** (BLOcks of Amino Acid SUbstitution Matrix) matrices, based pairwise sequence comparisons, used for BLAST

BLOSUM 80

PAM 1

*Less divergent*

BLOSUM 62

PAM 120

← →

BLOSUM 45

PAM 250

*More divergent*

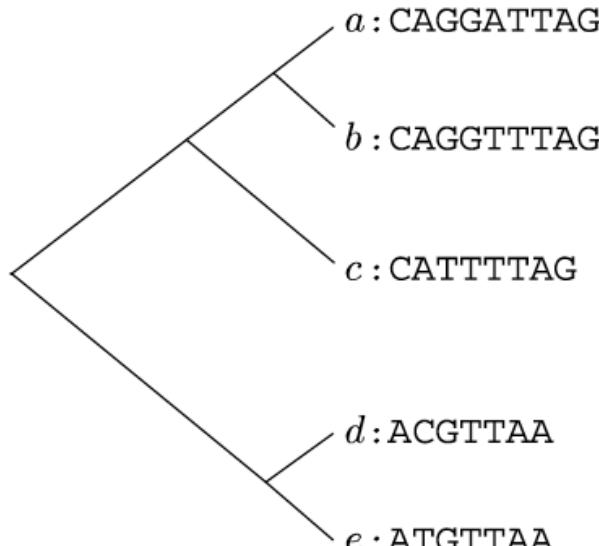
# BLOSUM62 substitution matrix

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	-2	2	-3
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	-1	1
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	2	-3
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

# Progressive alignment

*a*: CAGGATTAG  
*b*: CAGGTTTAG  
*c*: CATTTTAG  
*d*: ACGTTAA  
*e*: ATGTTAA

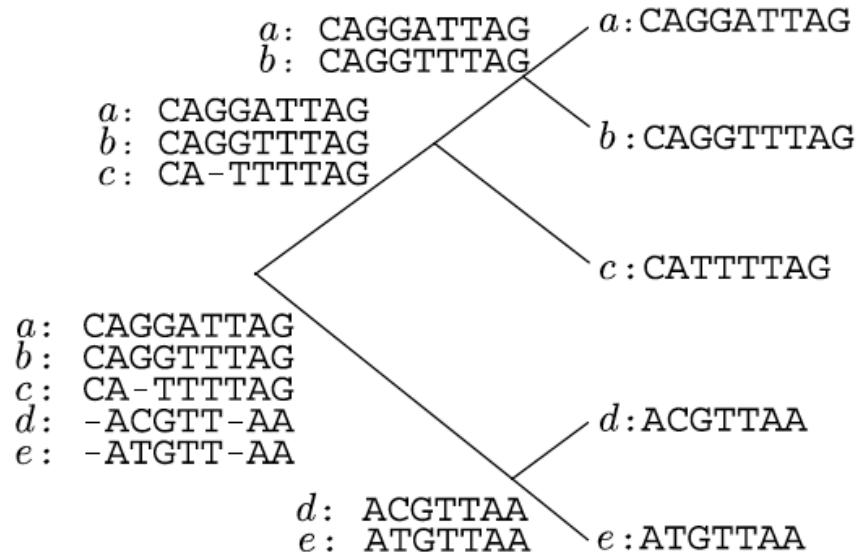
**Input**



**Guide tree**

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>
<i>a</i>	0	1	3	4	4
<i>b</i>	1	0	2	4	4
<i>c</i>	3	2	0	5	5
<i>d</i>	4	4	5	0	1
<i>e</i>	4	4	5	1	0

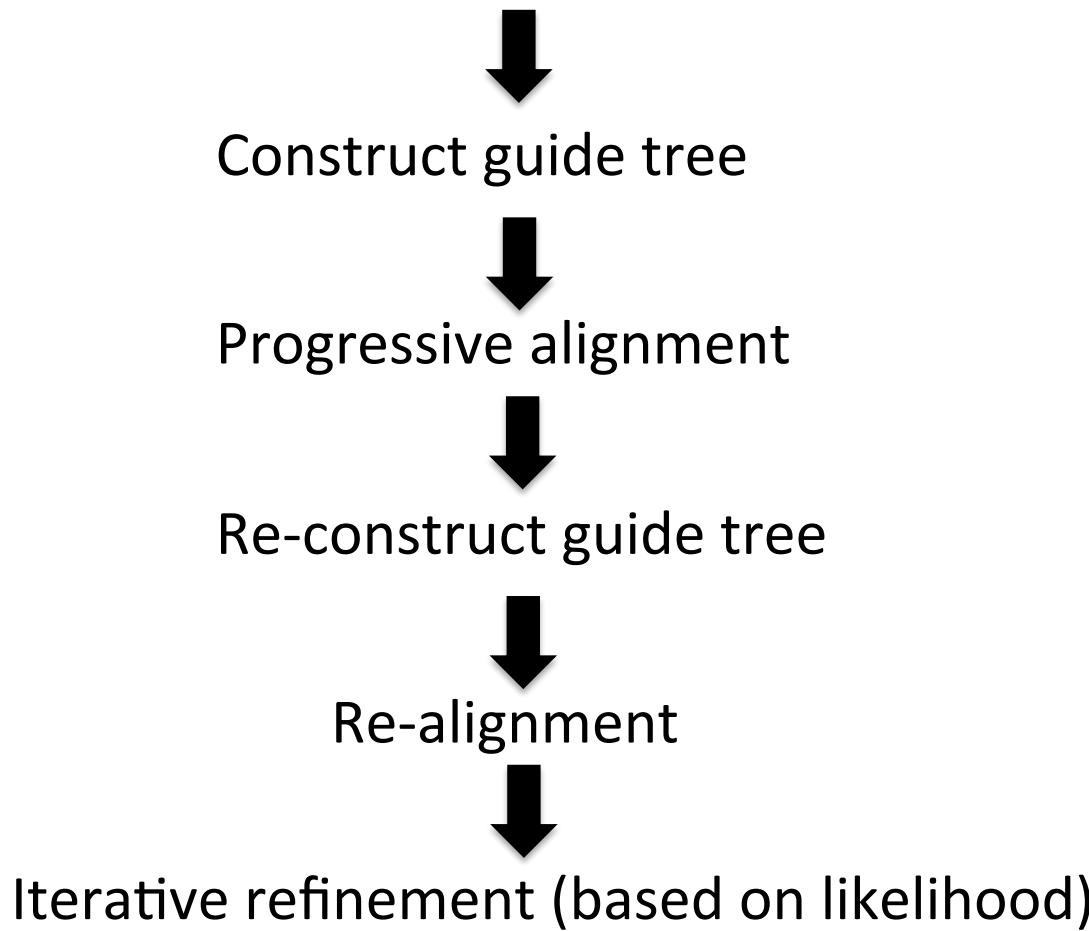
**Distance matrix**



**Progressive alignment**

# Iterative alignment

Distance matrix based on the number of shared seeds



# Alignment methods (condensed)

## Progressive

- ClustalW
- T-Coffee
- LALIGN
- PSAlign

## Iterative

- MUSCLE
- MAFFT
- DIALIGN

## Other

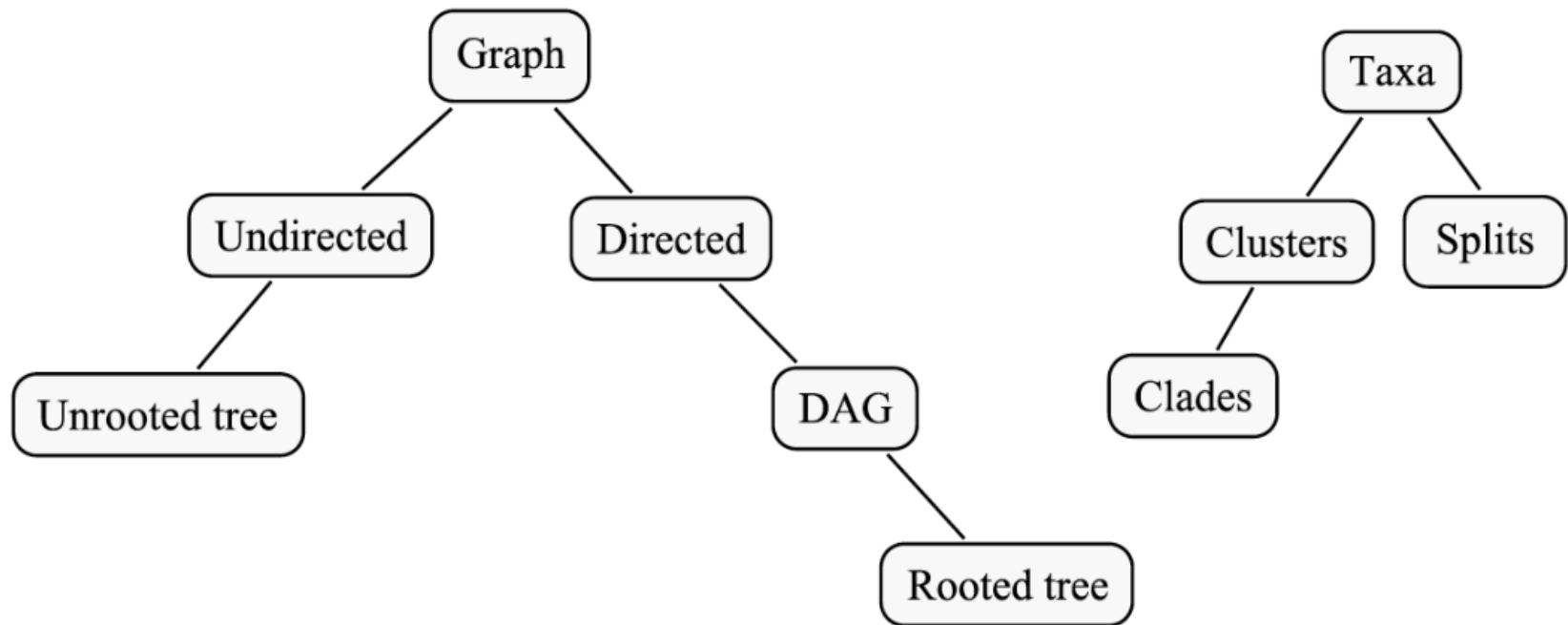
- Saté
- MAFFT
- Prank

Slow,  
inaccurate

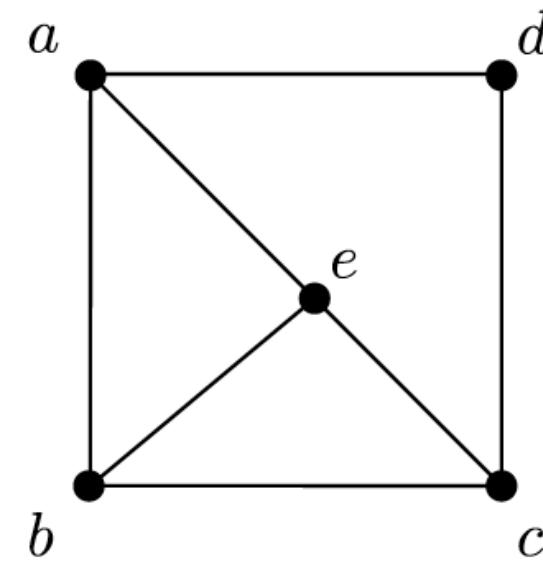
Fast,  
moderate-sized  
data

Large data,  
accurate

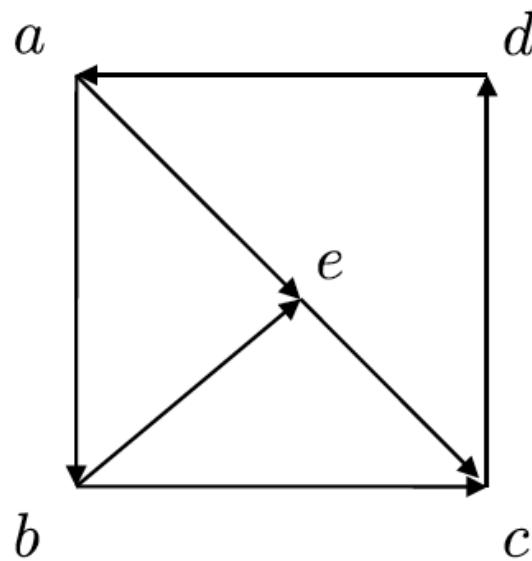
# Basic graph types



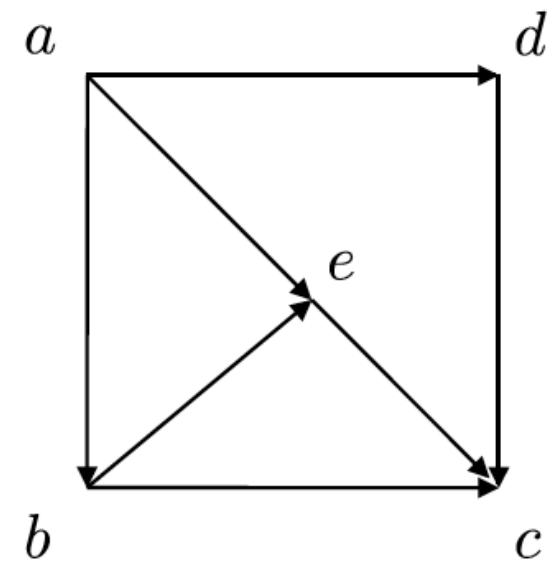
# Basic graph types



Undirected graph

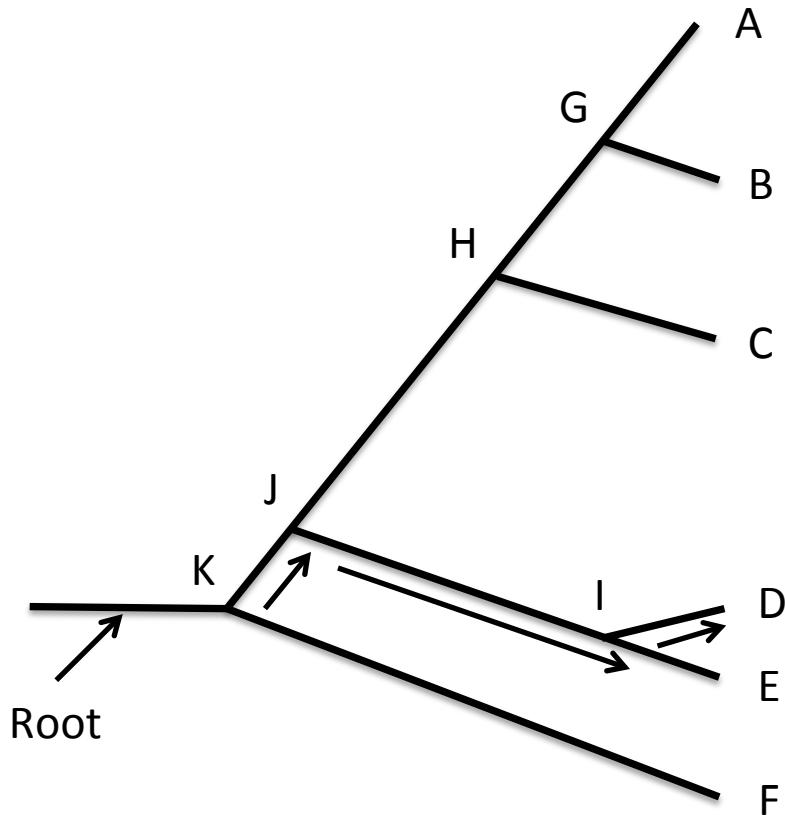


Directed graph

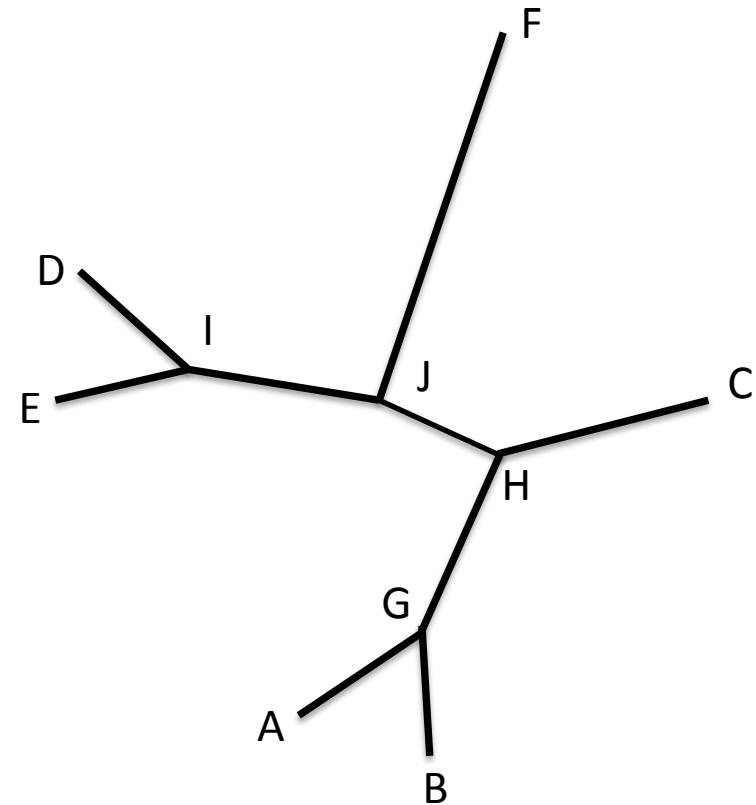


DAG

# Basic tree types



Rooted tree



Unrooted tree

# Phylogeny inference

There are four primary types of phylogenetic inference:

- **Distance based** – minimize distance between clusters
- **Parsimony** – minimum number of mutations
- **Maximum Likelihood** – maximizes likelihood of obtaining data, under a model of evolution
- **Bayesian methods** – compute posterior distribution of trees based on the data, model of evolution, and assumed prior dist. of trees

Support obtained from posterior distribution, or bootstrapping

# Maximum Parsimony

- **Maximum Parsimony (MP)** – aims to find the tree with the fewest changes
- Calculate the minimum number of changes at each position along the branches to explain the observed terminal nodes
- *Parsimony length* – sum of scores for all positions

# Maximum Parsimony

## Caveats:

Always underestimates the real divergence between distantly related taxa

There is no way to correct for multiple substitutions, this may lead to long branch attraction\*

\*Substitution models can account for increased prob. change on long branches

# **Distance-based methods**

**Unweighted pair group method with arithmetic means (UPGMA)** – Clustering by searching for the smallest distance in the pairwise matrix

- repeats by finding distance to new clusters
- averages distance between original clusters
- assumes the rate is the same on all branches

# **Distance-based methods**

**Neighbor-Joining (NJ)** – tree constructed by sequentially finding pairs

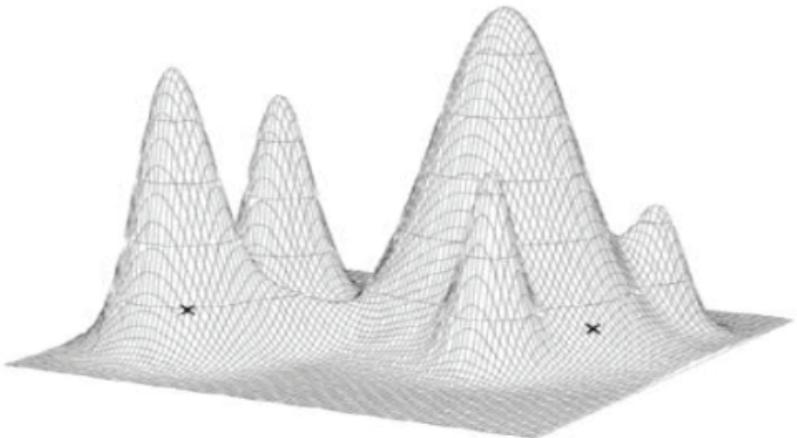
- doesn't attempt to cluster based on similarity
- minimizes length of all internal branches (and tree)
- starts with a star-like tree, adds a branch, recalculates tree

How to assess confidence in the resulting tree?

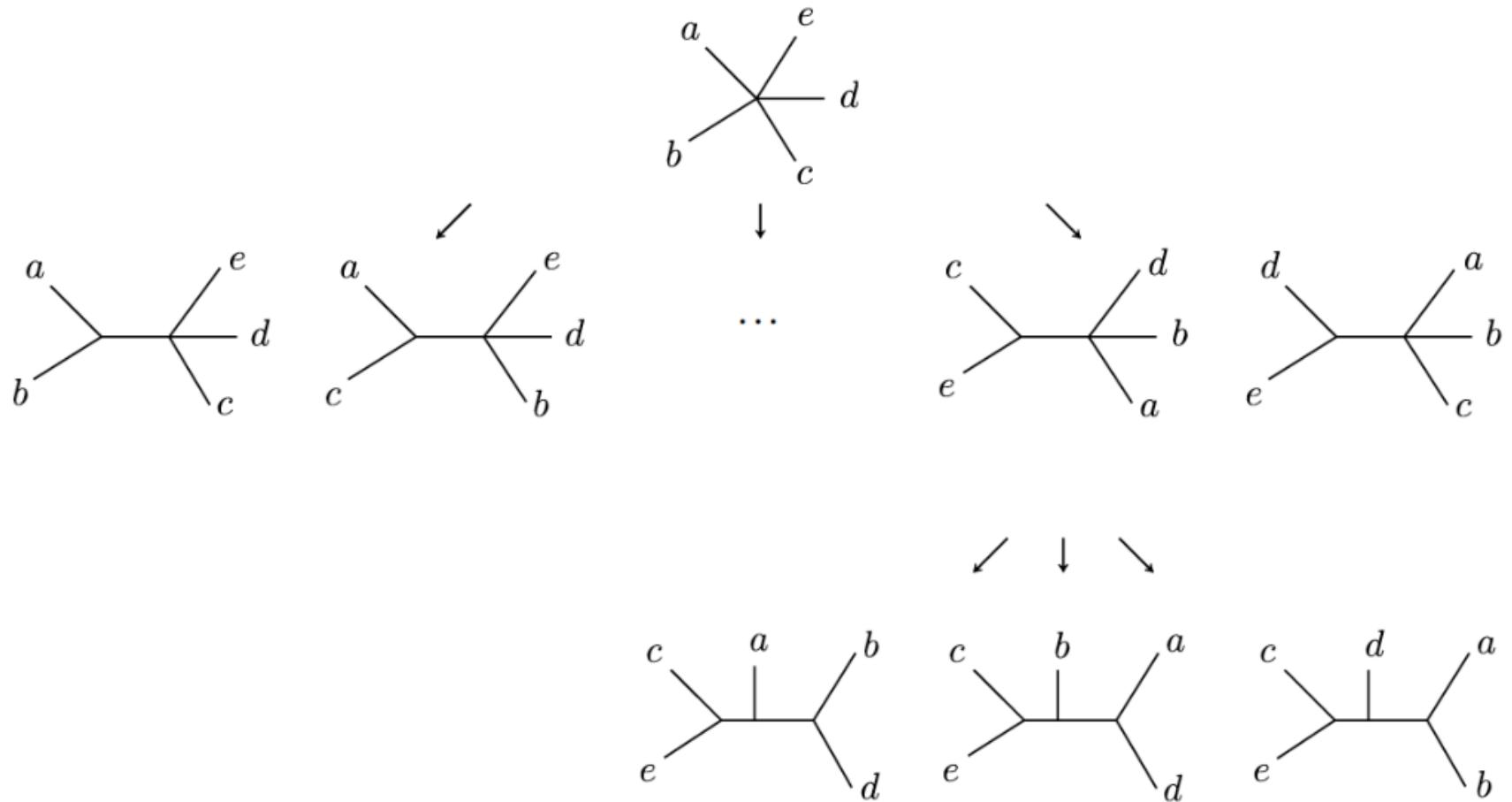
# Tree search problem

Number of taxa	Number of trees
10	$2 \times 10^6$
22	$3 \times 10^{23}$
50	$3 \times 10^{74}$
100	$2 \times 10^{182}$
1,000	$2 \times 10^{2,860}$
10,000	$8 \times 10^{38,658}$
100,000	$1 \times 10^{486,663}$
1,000,000	$1 \times 10^{5,866,723}$
10,000,000	$5 \times 10^{68,667,340}$

Likelihood ( or inv [sum  
of branch lengths])



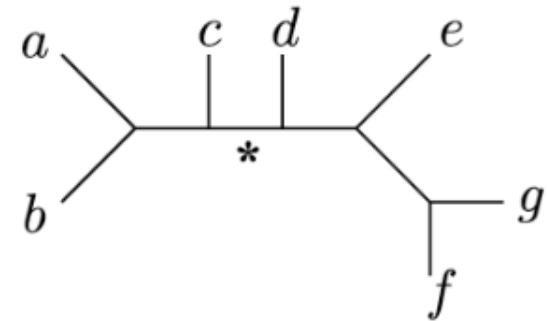
# Star-like decomposition



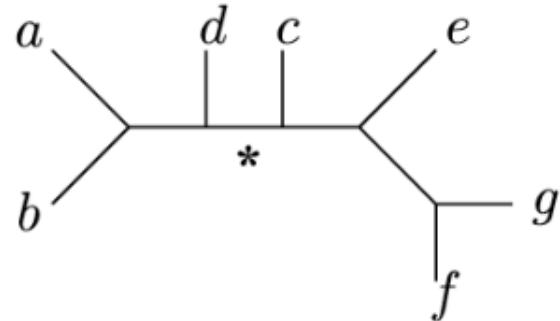
**Branch swapping can refine trees initially found in heuristic search**

- Nearest-Neighbor Interchange (NNI)
- Subtree Pruning & Regrafting (SPR)
- Tree Bisection & Reconnection (TBR)

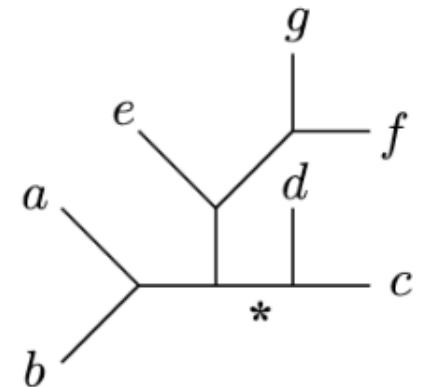
# Nearest neighbor interchange (NNI)



Phylogenetic tree

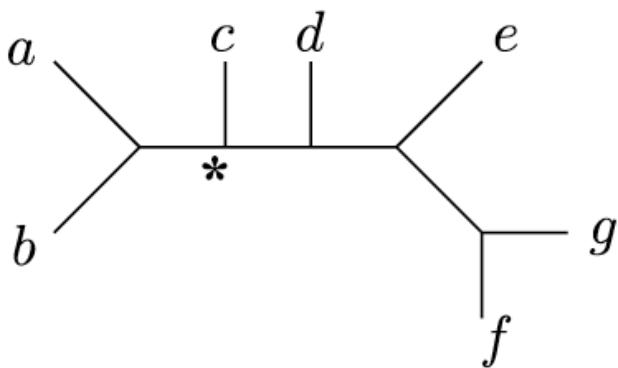


One NNI tree

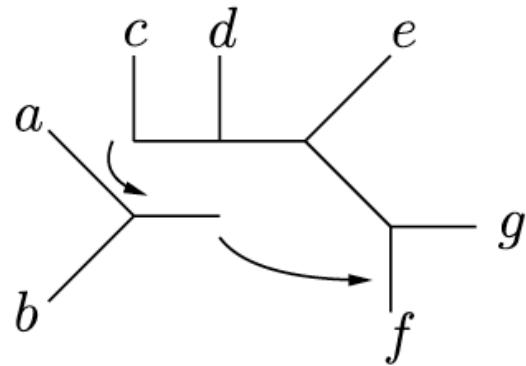


The other NNI tree

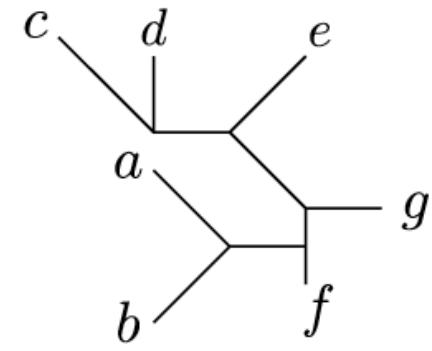
# Subtree prune and regraft (SPR)



Phylogenetic tree

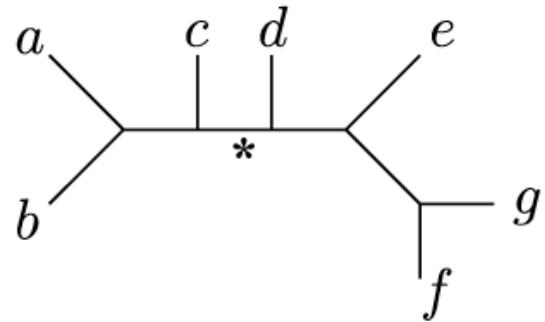


Subtree prune...

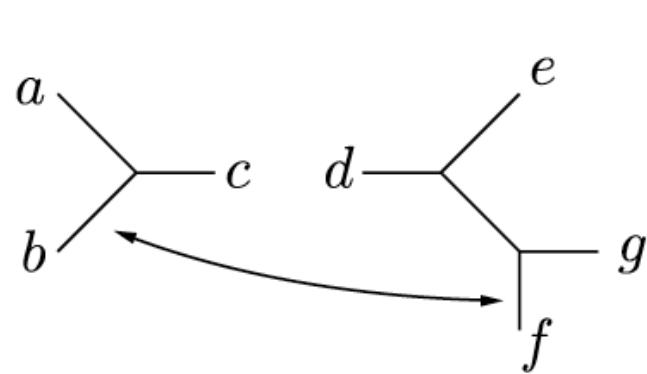


Regraft

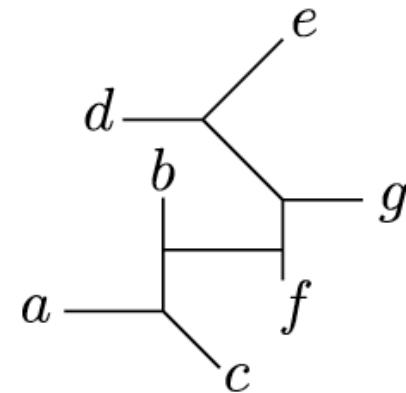
# Tree bisection and reconnection (TBR)



Tree bisection...



Reconnection choice...



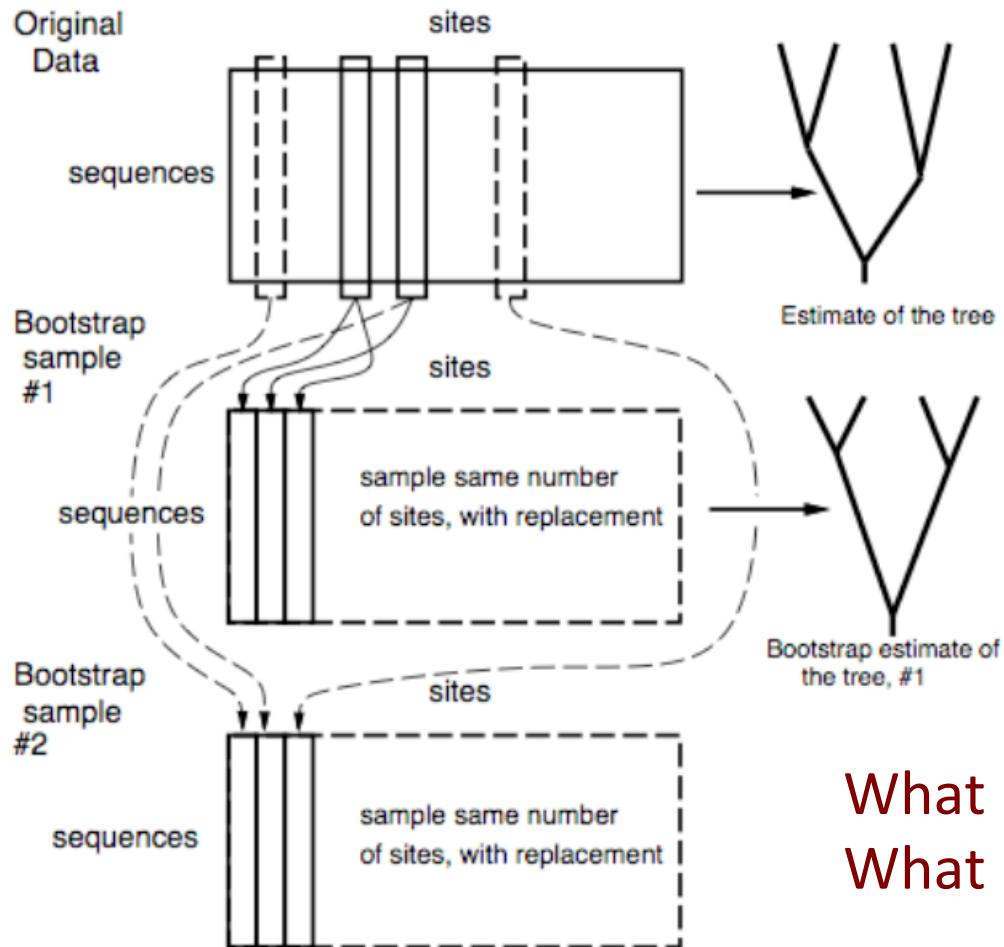
And reconnection

# Which topology is correct?

$$\text{NNI}(T) \subseteq \text{SPR}(T) \subseteq \text{TBR}(T)$$

- The size of  $\text{NNI}(T)$  grows linearly with  $n$
- The size of  $\text{SPR}(T)$  grows quadratically
- The size of  $\text{TBR}(T)$  depends not only on  $n$ , but also on the actual topology of the phylogenetic trees considered

# Nonparametric bootstrapping



What if the tree is wrong?  
What if the model is wrong?

# Likelihood methods

- In statistics we usually estimate the probability of a hypothesis given data:

$$\Pr(\text{Hypothesis} \mid \text{Data})$$

- But Likelihoods measure the probability of the data given a hypothesis:

$$L(\text{Tree} \mid \text{Data}) = \Pr(\text{Data} \mid \text{Tree})$$

- Any estimate of  $\Pr(\text{Data} \mid \text{Hypothesis})$  requires a model: e.g. the binomial for coin flips, and a substitution model for molec. phylogenetics

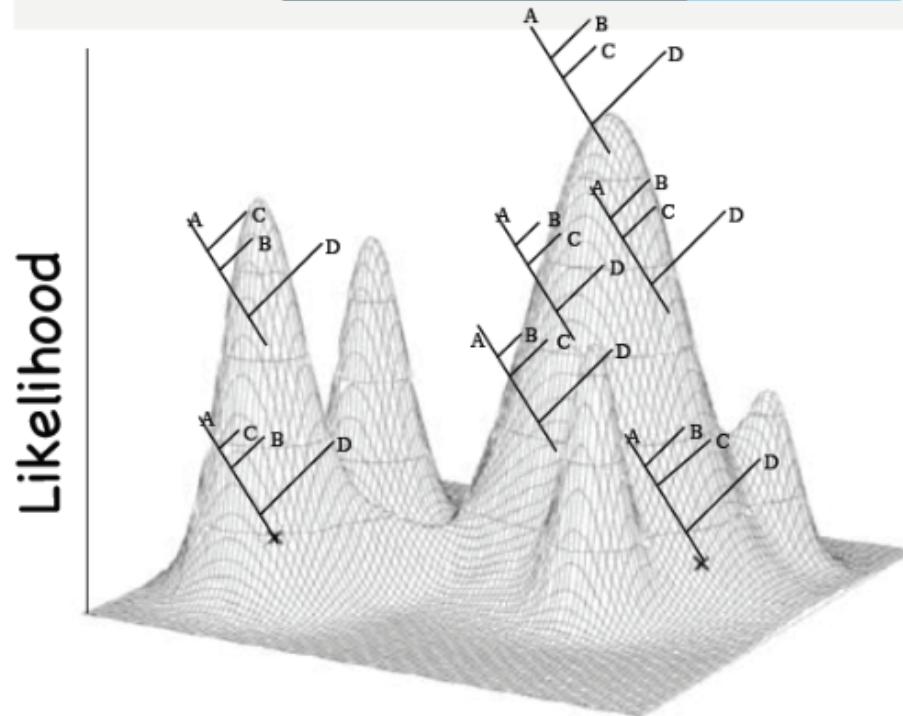
# Likelihoods in phylogenetic inference

MAX  $L(\text{Tree topology} \rightarrow \text{B} \text{ M} | \text{Data})$

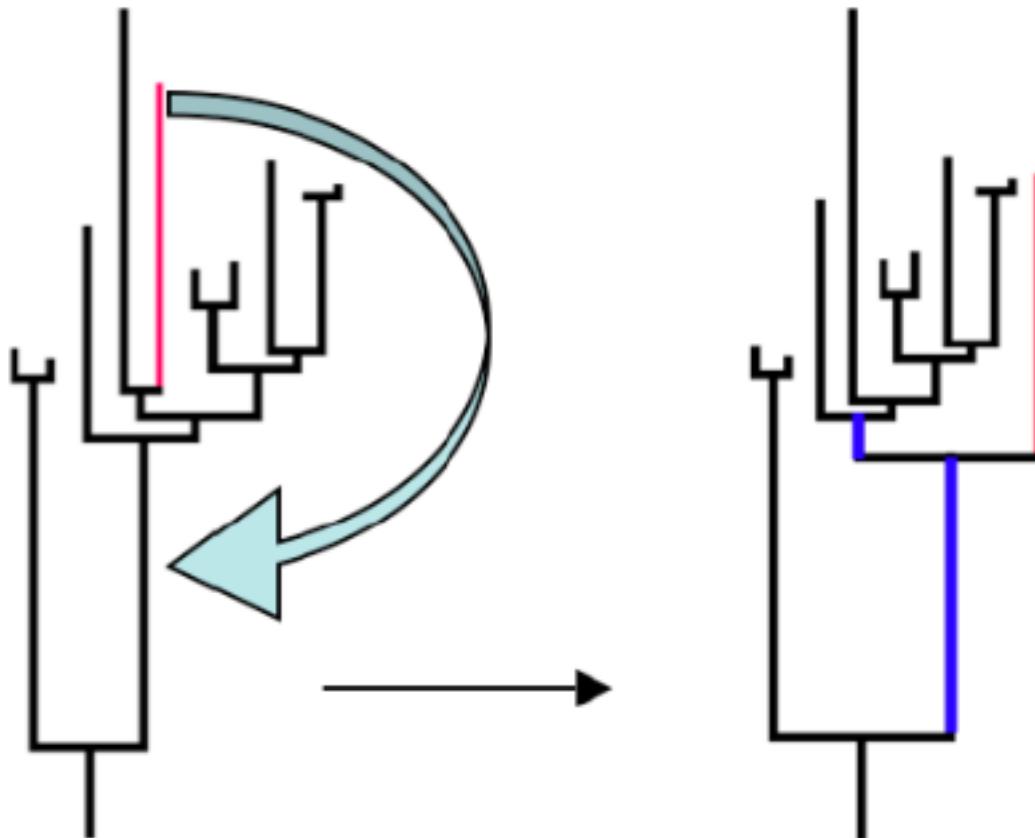
Branch-length parameters

Substitution model parameters

Tree topology

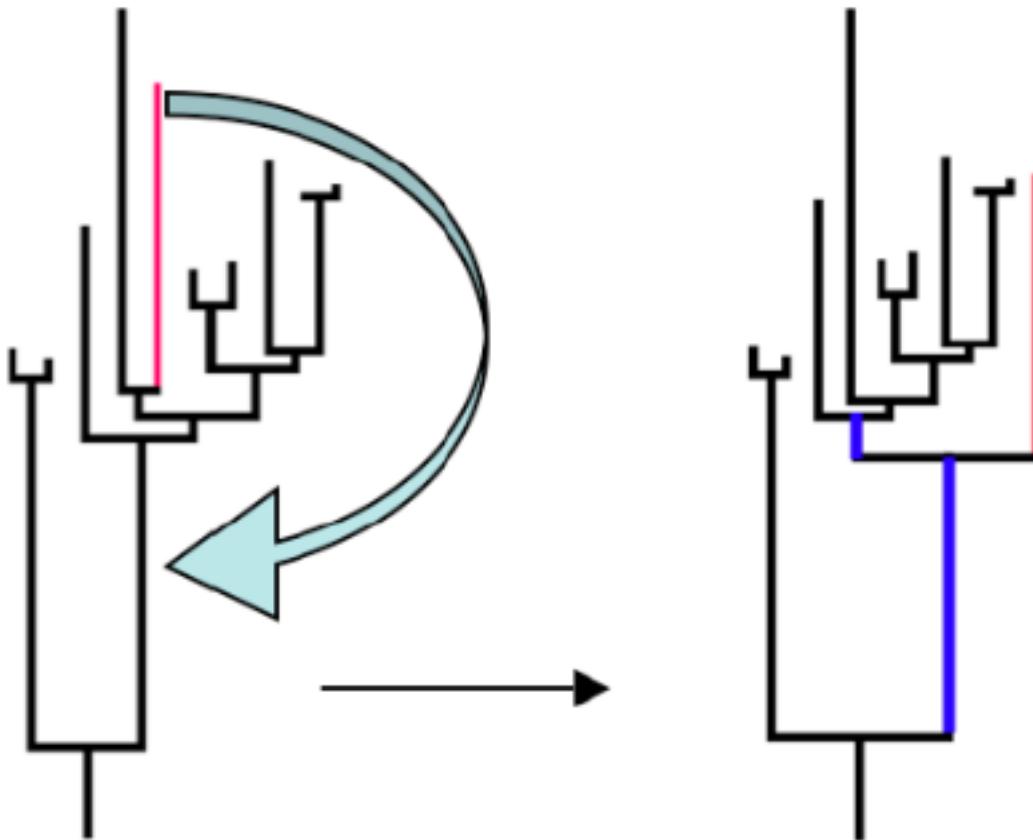


# Approximate likelihood



Garli only reoptimizes  
Branch lengths in the  
old and new  
neighborhood of  
pruned and regrafted  
branch

# RAxML employs a similar approach with SPR swapping on a MP tree



Bootstrap analysis  
with random addition  
MP tree for  
each round is highly  
recommended...??

# Bayesian inference

$X$  = Data

$\Theta$  = Model parameters

Posterior  
distribution

Prior distribution

Likelihood

$$f(\theta | X) = \frac{f(\theta)f(X | \theta)}{\int f(\theta)f(X | \theta)d\theta}$$

Normalizing constant

# Bayesian inference

“Model Parameters” may include tree topology, branch lengths and substitution model parameters

$$\text{Posterior distribution} = \frac{\text{Prior distribution} \times \text{"Likelihood"} }{\text{Normalizing constant}}$$
$$\text{PP}(T, B, M | \text{Data}) = \frac{\Pr(T) \Pr(B) \Pr(M) P(\text{Data} | T, B, M)}{\sum \int \int \Pr(T) \Pr(B) \Pr(M) P(\text{Data} | T, B, M) d(B)d(M)}$$

Diagram illustrating the components of Bayesian inference:

- Posterior distribution
- Prior distribution
- "Likelihood"
- Normalizing constant

The "Likelihood" term is highlighted in yellow.

# Bayesian inference

## Pros:

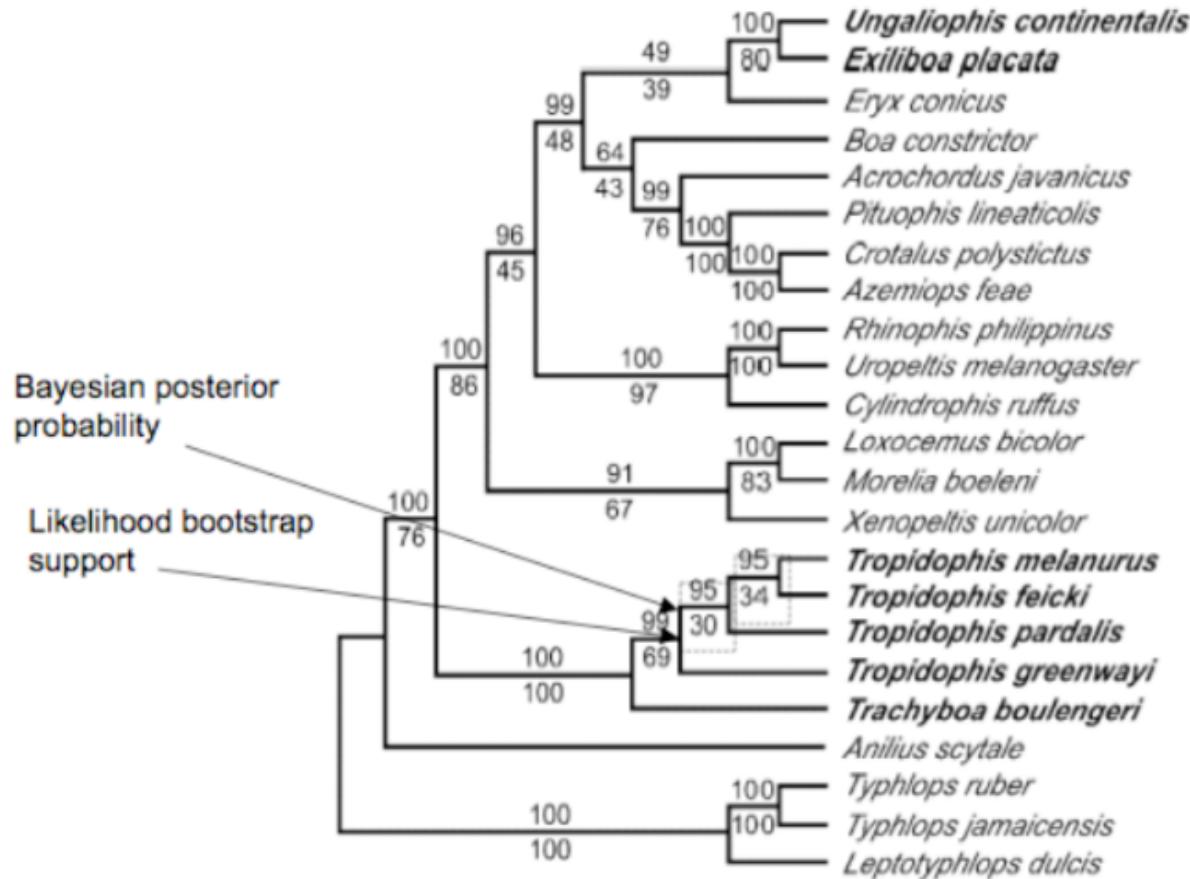
- Can incorporate prior information (divergence times, subs. rates across genes)
- Can impose parameter constraints with a prior
- Allows inference across complex tree and parameter space (rather than maximization)
  - this simplifies searches for large trees

# Bayesian inference

## Minuses:

- Must incorporate information about parameters and that is not always available (e.g., influence of branch length prior and subs. model on posterior probability)
- Numerical methods for inferring posterior probabilities are difficult to diagnose (convergence in MrBayes)

# How to assess posterior probability and ML bootstrap support?



# **KH (Kishino-Hasegawa), SH (Shamodaira-Hasegawa), AU (Approximately Unbiased) Tests - Implementations**

- Propose a set of hypothesized trees and substitution model
- Estimate likelihoods for each tree:

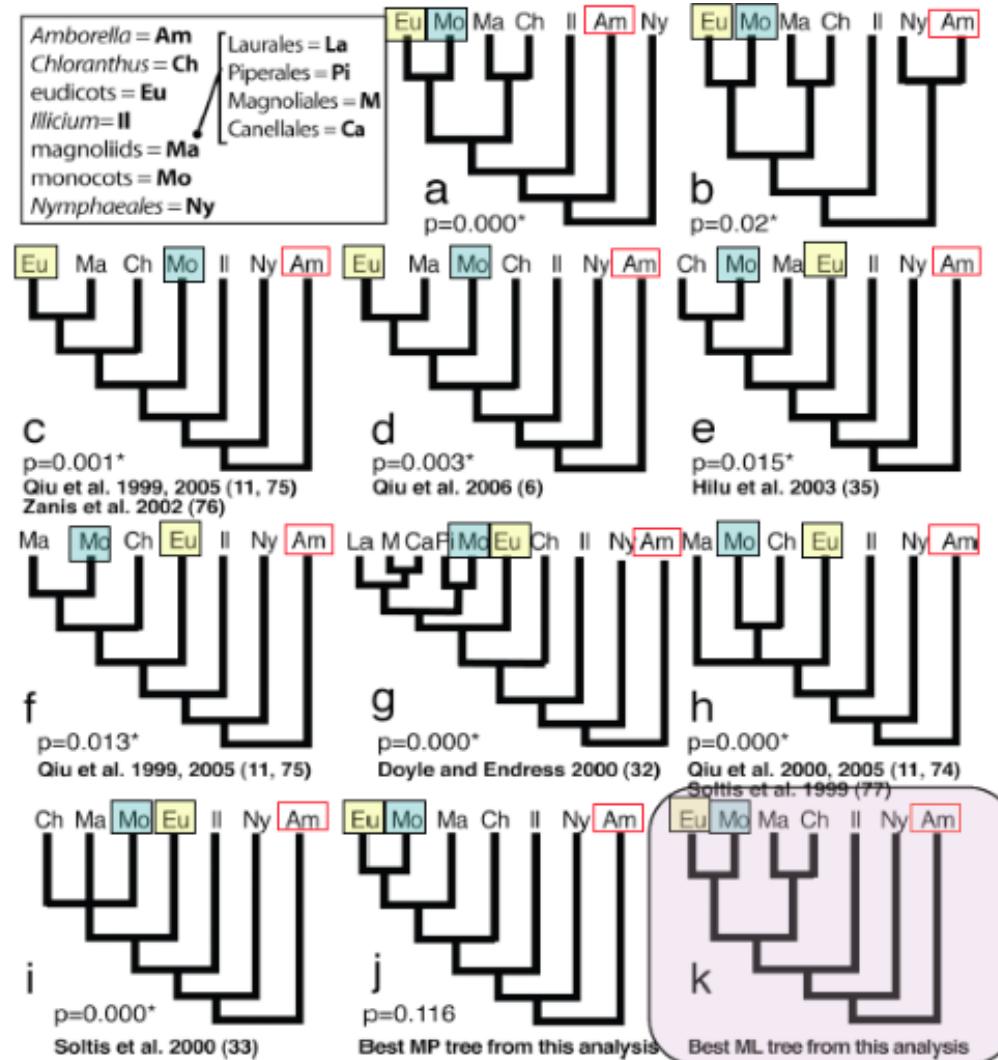
$$\text{Max } L(B, M \mid \text{Tree, Data})$$

KH test compares pairs of trees by assessing difference in likelihood

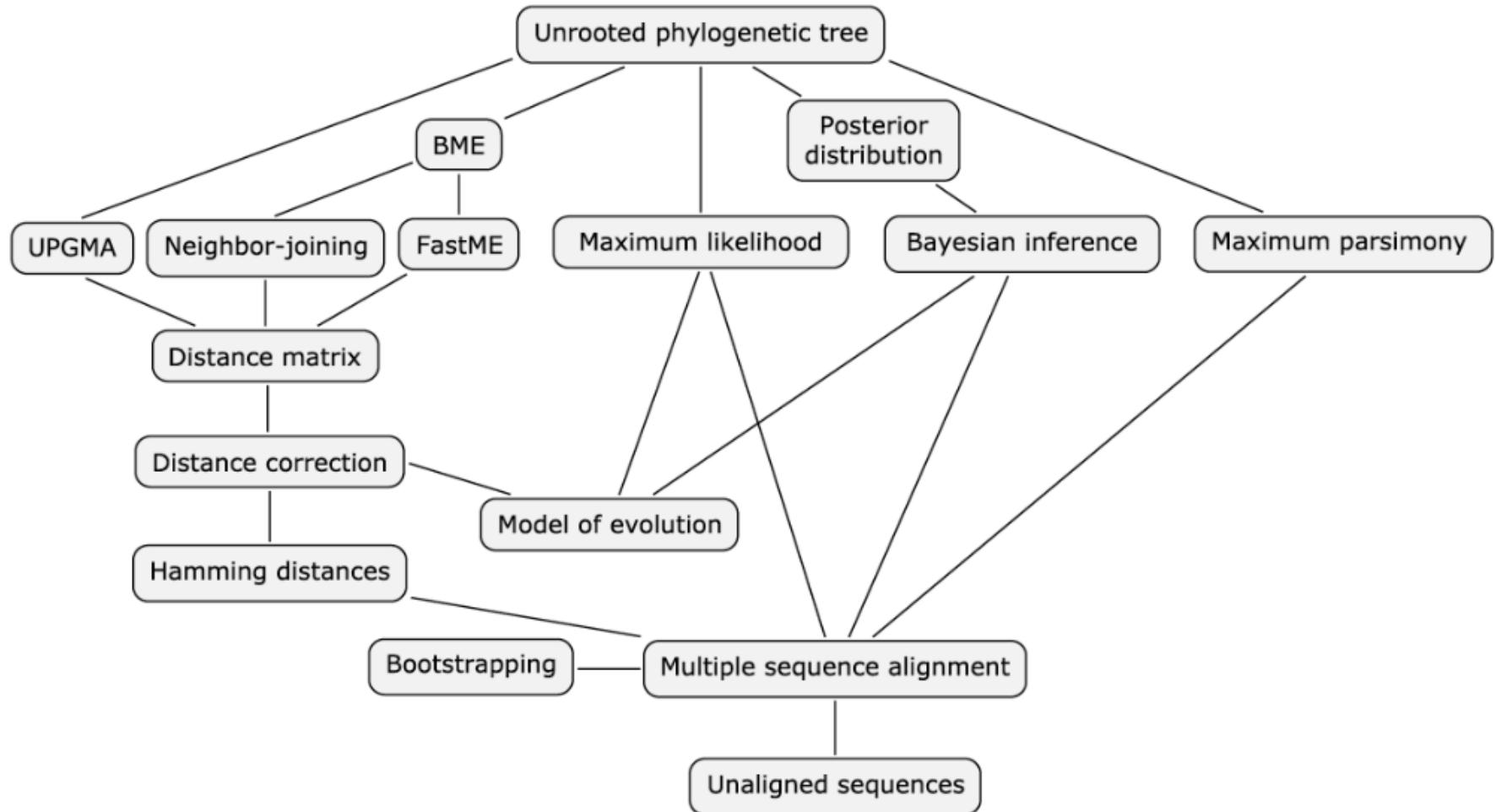
SH and AU tests bootstrap resample the original data matrix thousands of times

**Likelihoods compared for best tree vs. alternative trees**

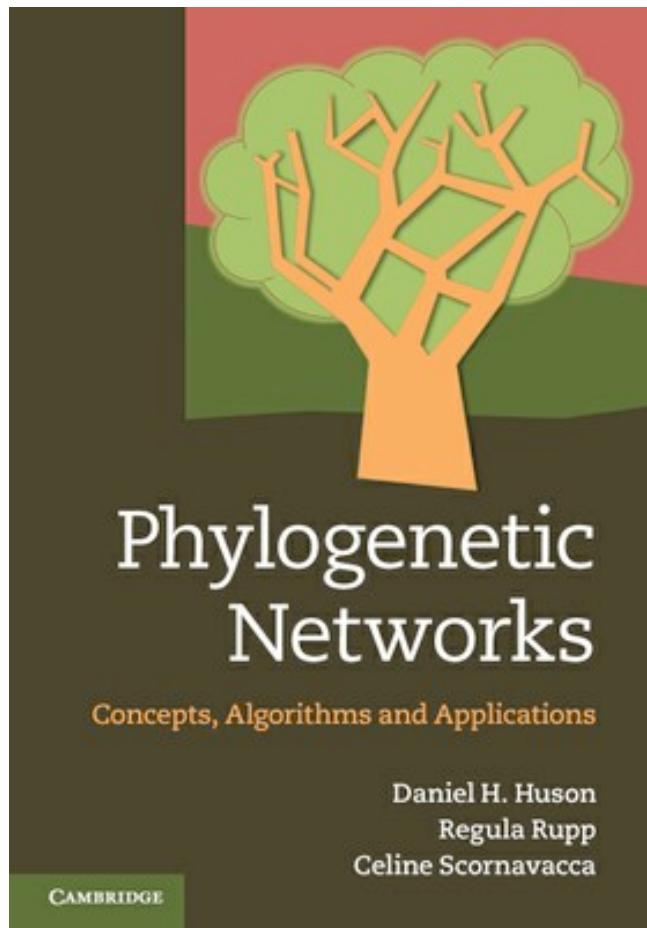
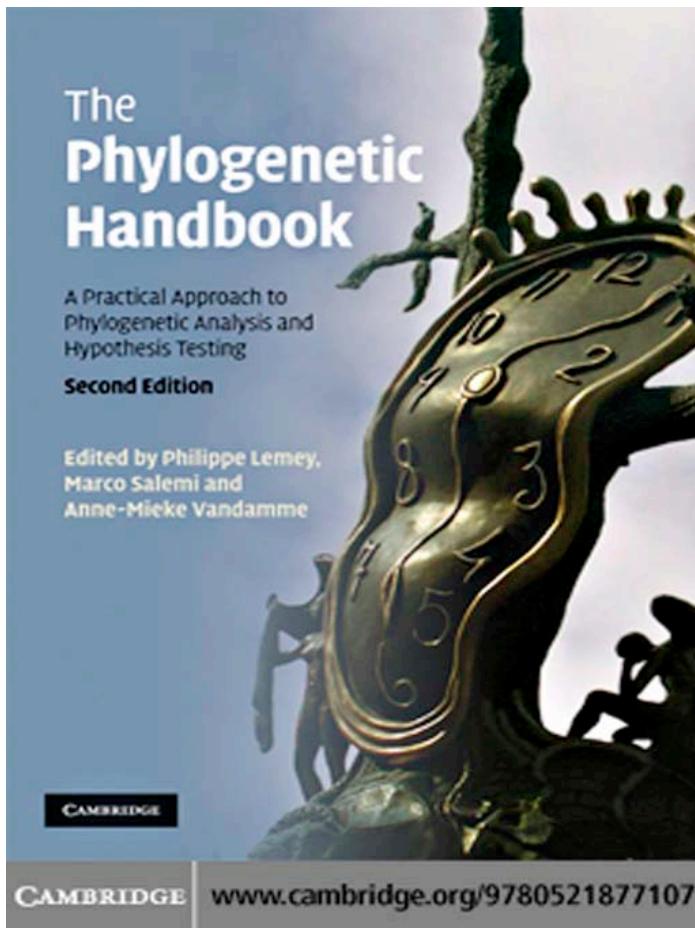
# SH, AU tests of alternative topologies



# High level view of phylogeny inference



# Resources



# Practical session

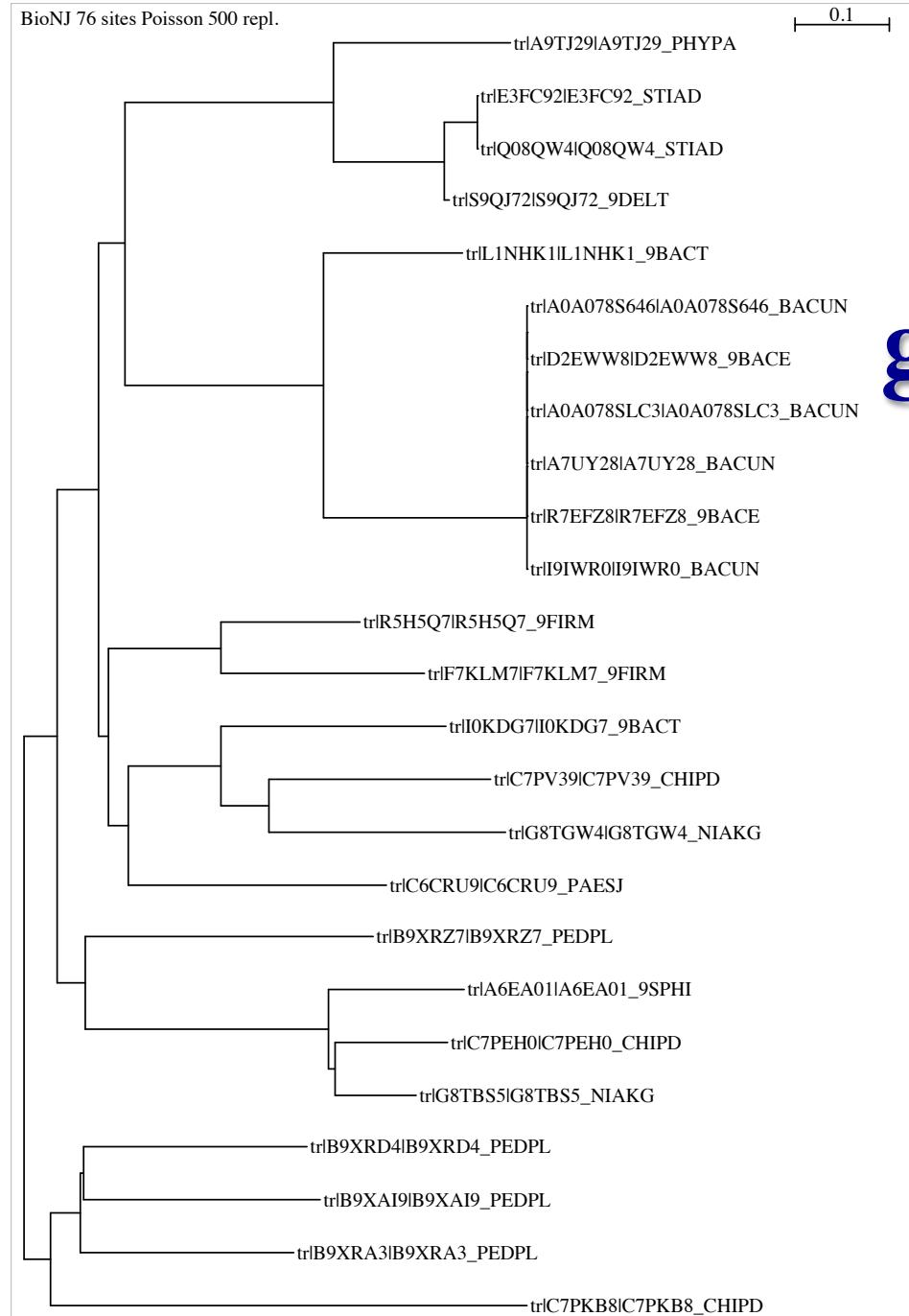
- Exercise list is online:

[github.com/UBCBio525/Bio525D/Day8](https://github.com/UBCBio525/Bio525D/Day8)

# Review of yesterday's findings

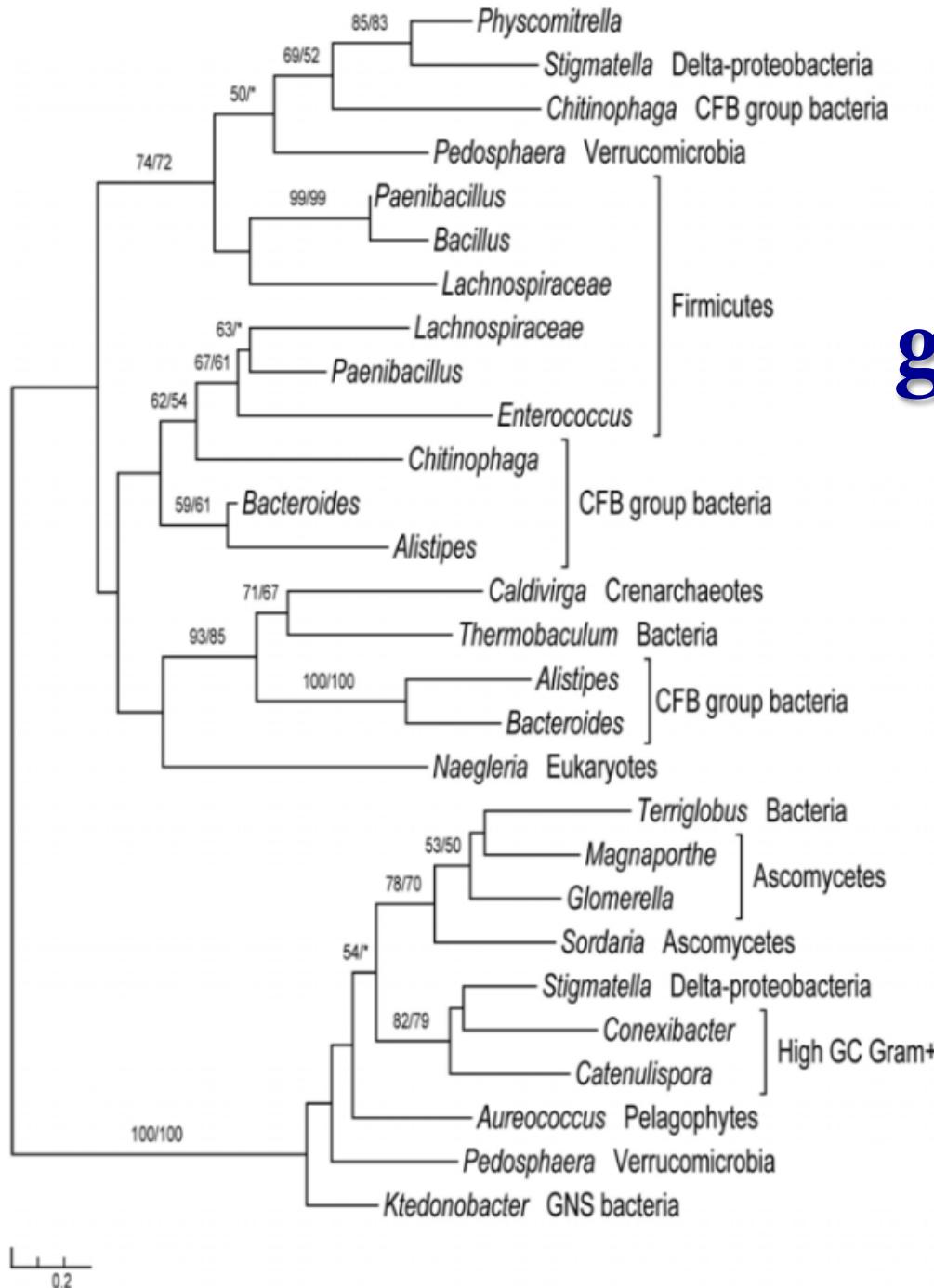
- Species identification – *Physcomitrella patens*
- BLAST hit summary
- Enrichment findings
- Evolutionary hypothesis?
- Results of phylogenetic inference....

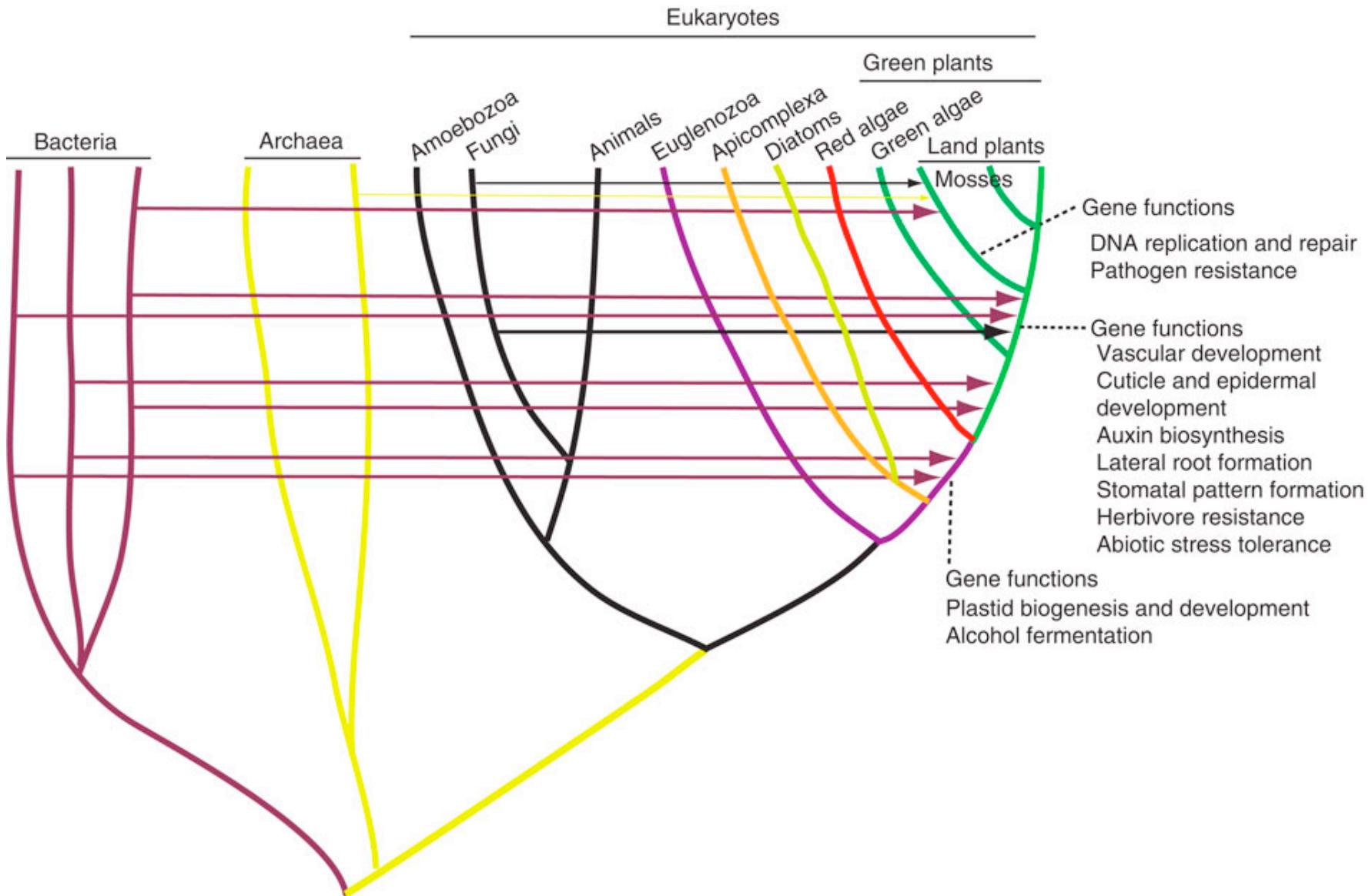
BioNJ 76 sites Poisson 500 repl.



# Phylogeny of glycoside hydrolase

# Phylogeny of glycoside hydrolase





# Reference (and data source)



## ARTICLE

Received 10 Jun 2012 | Accepted 20 Sep 2012 | Published 23 Oct 2012

DOI: 10.1038/ncomms2148

## Widespread impact of horizontal gene transfer on plant colonization of land

Jipei Yue<sup>1,2</sup>, Xiangyang Hu<sup>1,3</sup>, Hang Sun<sup>1</sup>, Yongping Yang<sup>1,3</sup> & Jinling Huang<sup>2</sup>