



Barcelona  
Biomedical  
Research  
Park



# Introduction to phylogeny, Orthology and Paralogy

Toni Gabaldón

Centre for Genomic Regulation (CRG), Barcelona  
[\(tgabaldon@crg.es\)](mailto:tgabaldon@crg.es)  
<http://gabaldonlab.crg.es>





**Barcelona**  
**Supercomputing**  
**Center**  
Centro Nacional de Supercomputación



INSTITUTE  
FOR RESEARCH  
IN BIOMEDICINE





# The ‘World’s Most Beautiful Data Center’ is a Supercomputer Housed in a Church

The MareNostrum 4 is only the world’s 25th most powerful supercomputer, but it definitely has the most style.

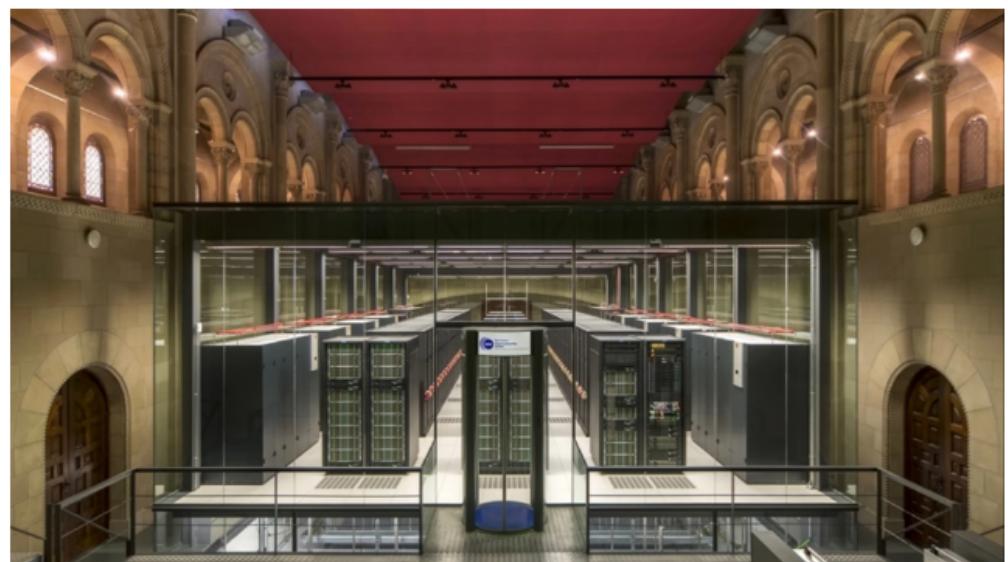
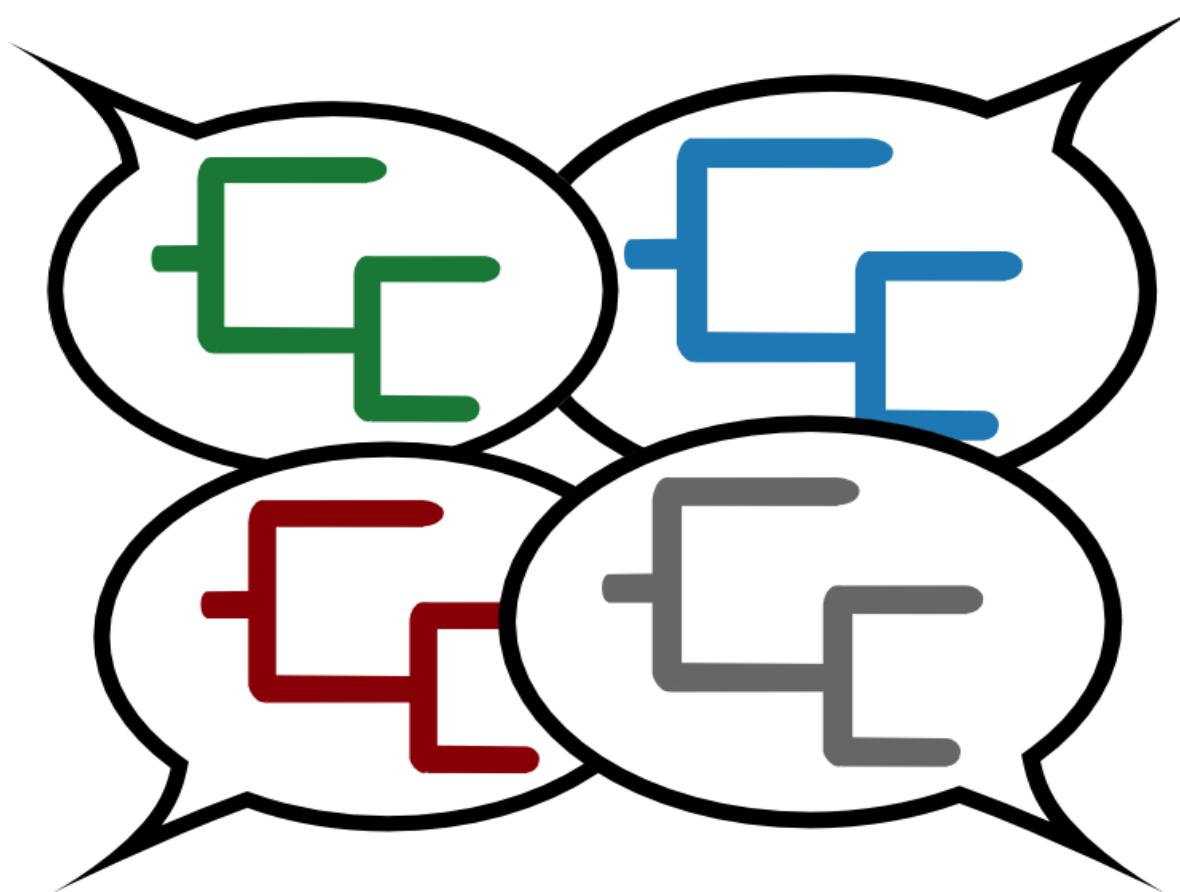


Image: Barcelona Supercomputing Center

From the outside, [Torre Girona Chapel](#) at the Polytechnic University of Catalonia in Barcelona looks like any one of the thousands of old churches that can be found throughout Spain, with a large cross mounted on the roof and a rose window perched above the entrance. Step through the chapel doors, however, and you won’t find any religious iconography or a congregation in prayer.

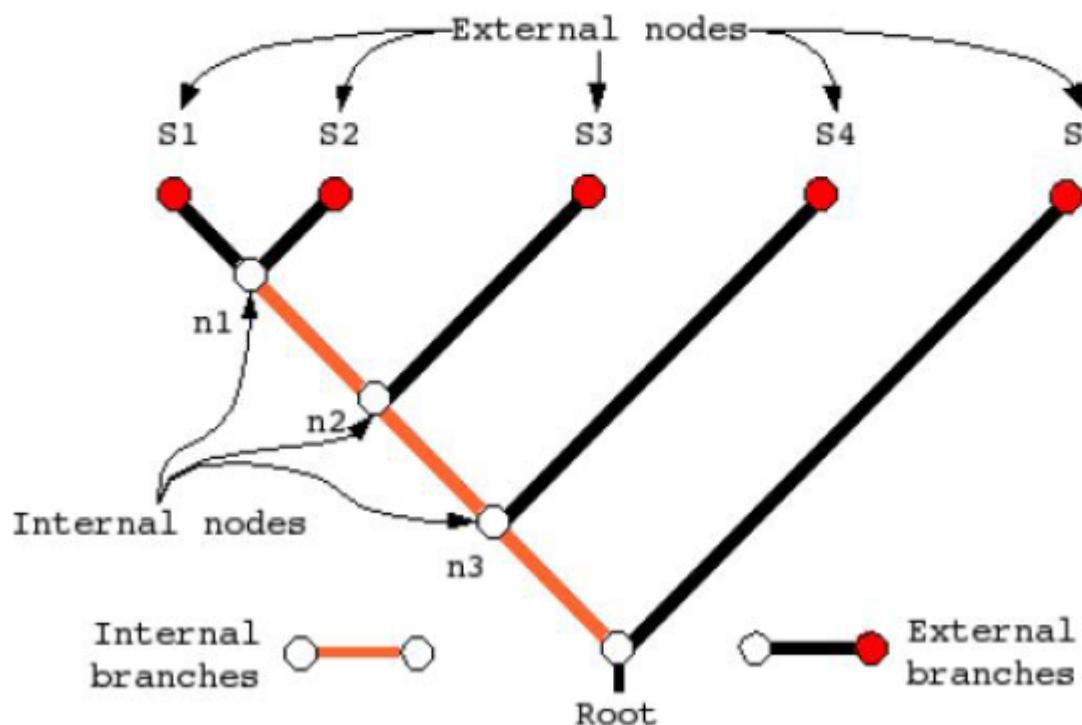
# A brief primer on phylogenetics and tree reconstruction



# A phylogenetic tree

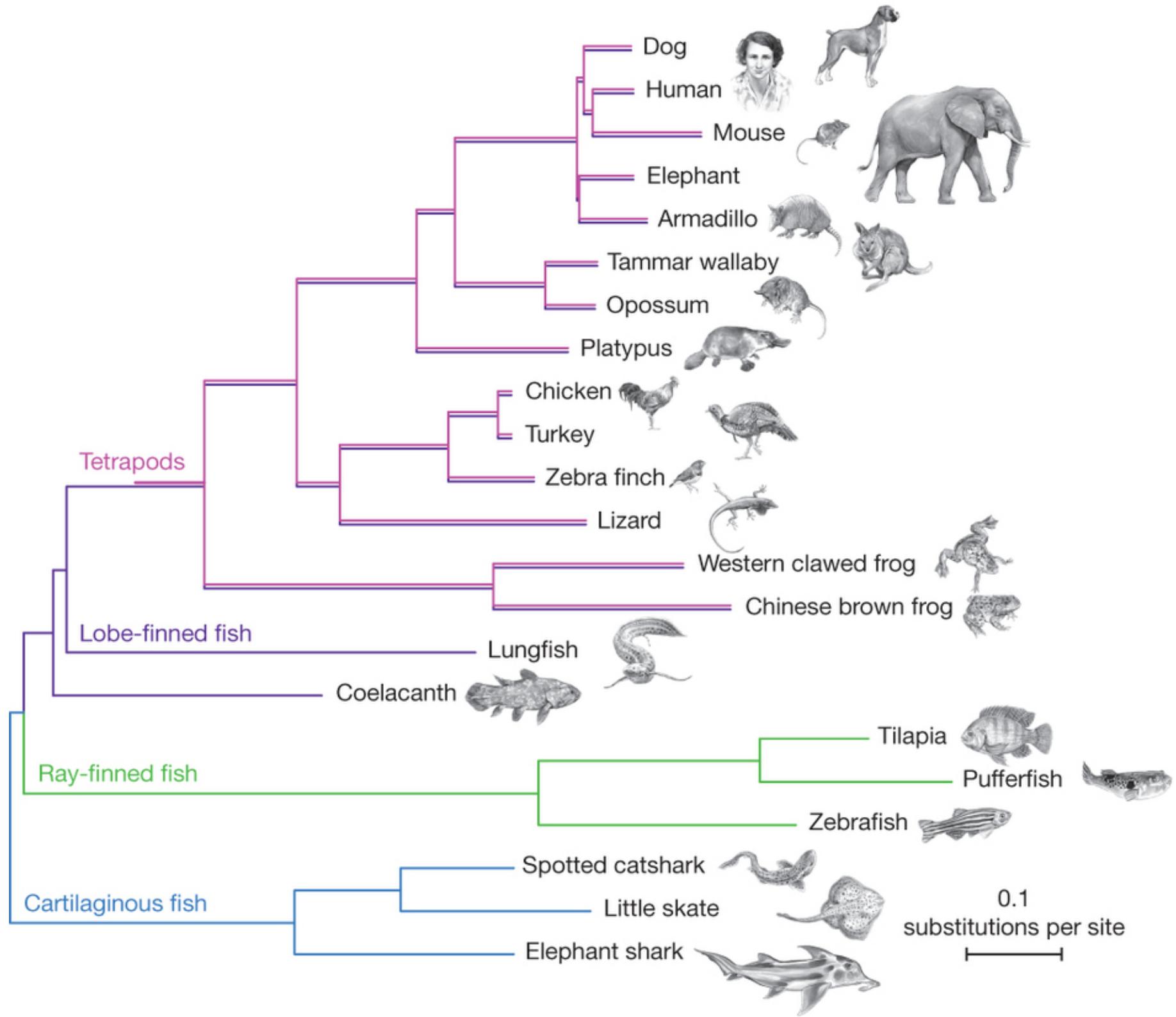
A branching diagram (**bipartite graph**) showing the inferred **evolutionary relationships** among various **biological species** or **other entities** (e.g sequences) **based on similarities and differences** in their physical and/or genetic characteristics.

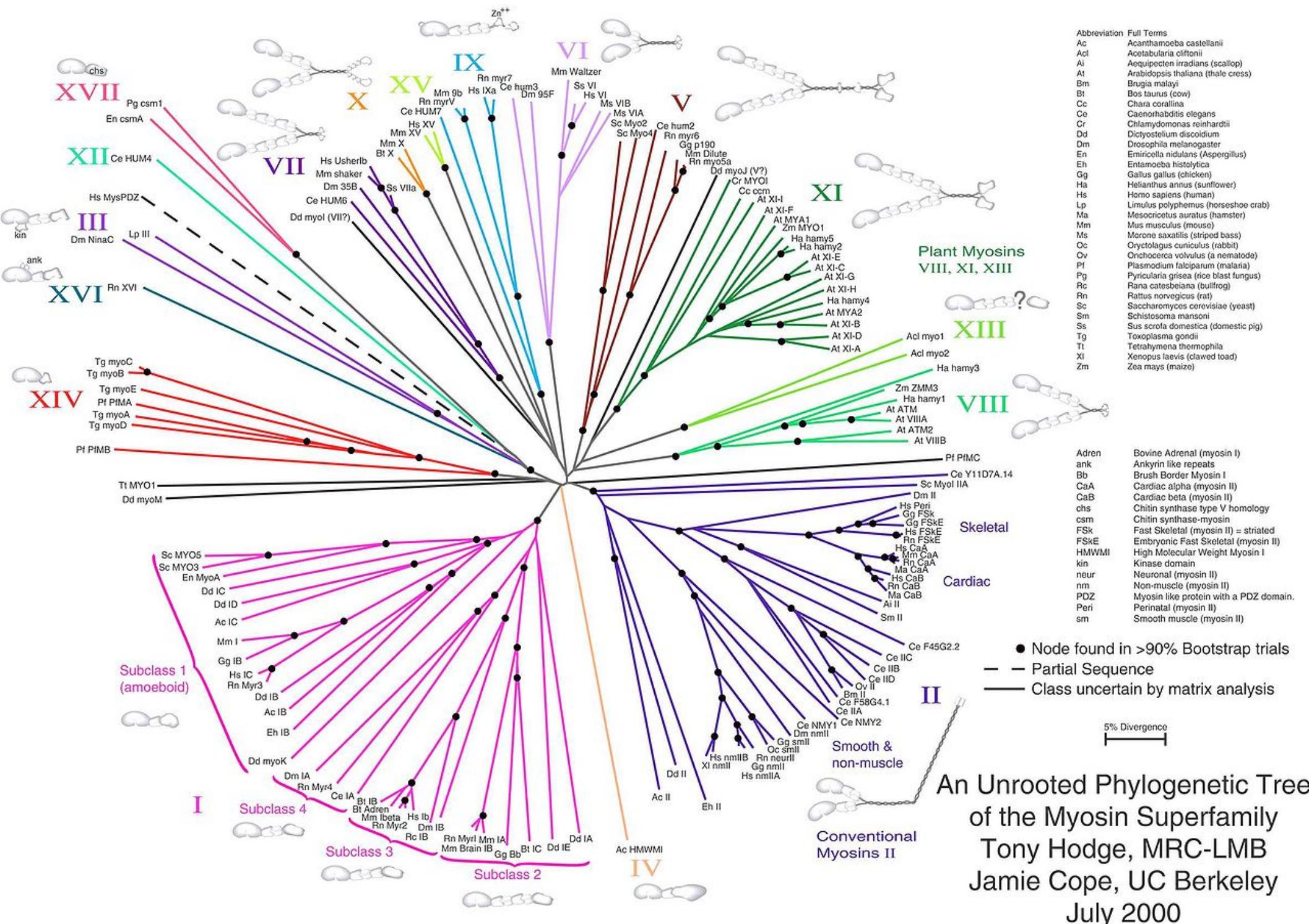
- **Nodes & branches.** Trees contain internal and external nodes and branches. In molecular phylogenetics, **external nodes** are sequences representing **genes, populations or species!**. Sometimes, **internal nodes** contain the ancestral information of the clustered species. A **branch** defines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



A phylogenetic tree is a **hypothesis** of how  
**things\*** are related through **evolution**

\* species, genes, ....





A phylogenetic tree is a **hypothesis** of how things\* are related through **evolution**

How we find our best hypothesis:

- data (i.e. sequences)
- a “model” of how this type of data evolve
- a way to assess how good our hypothesis is as compared to other possible hypotheses

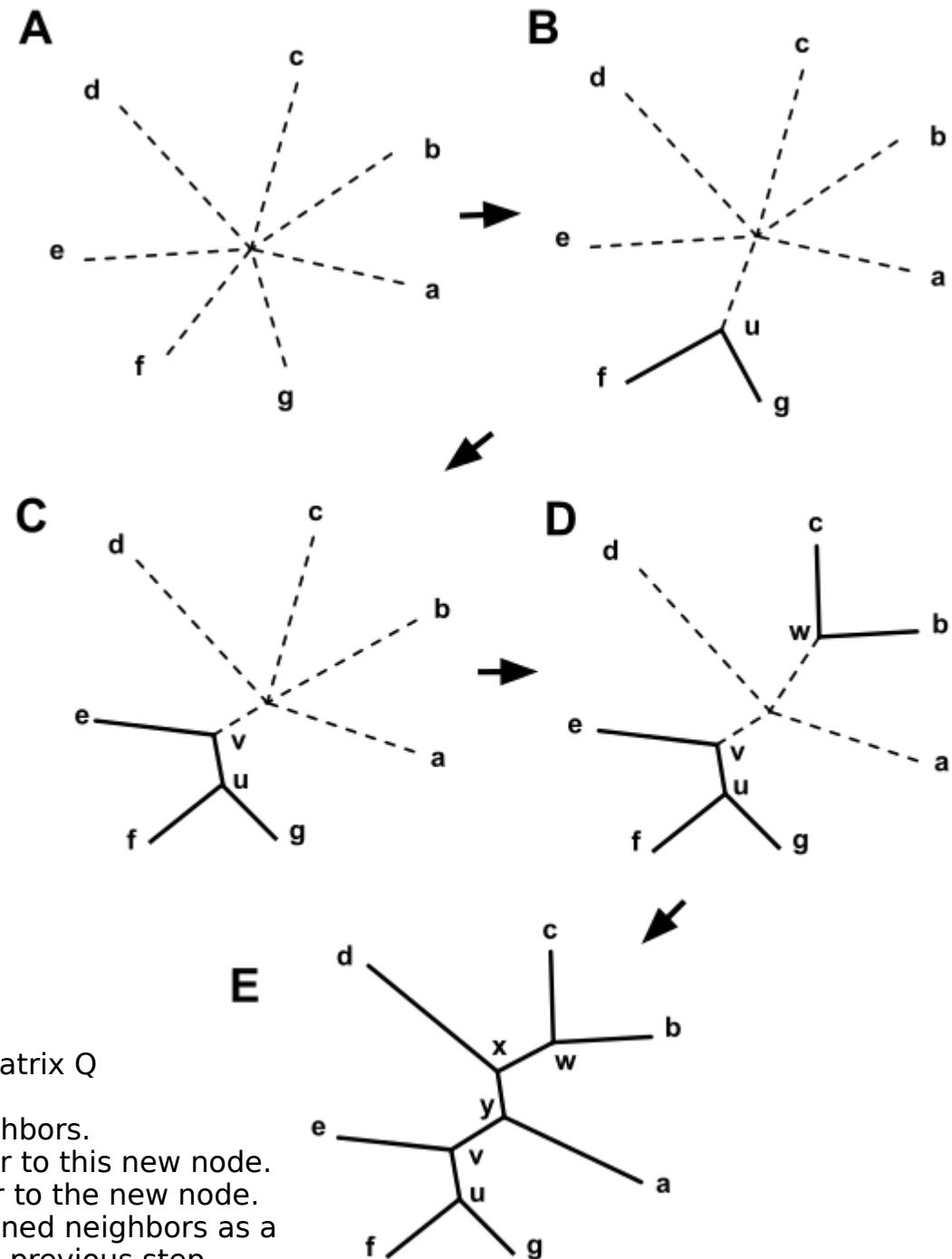
Phylogenetic approaches:

Distance methods (NJ, UPGMA)

Maximum Parsimony

Probabilistic Methods (Maximum Likelihood, Bayesian Inference, including coalescence methods)

# 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.

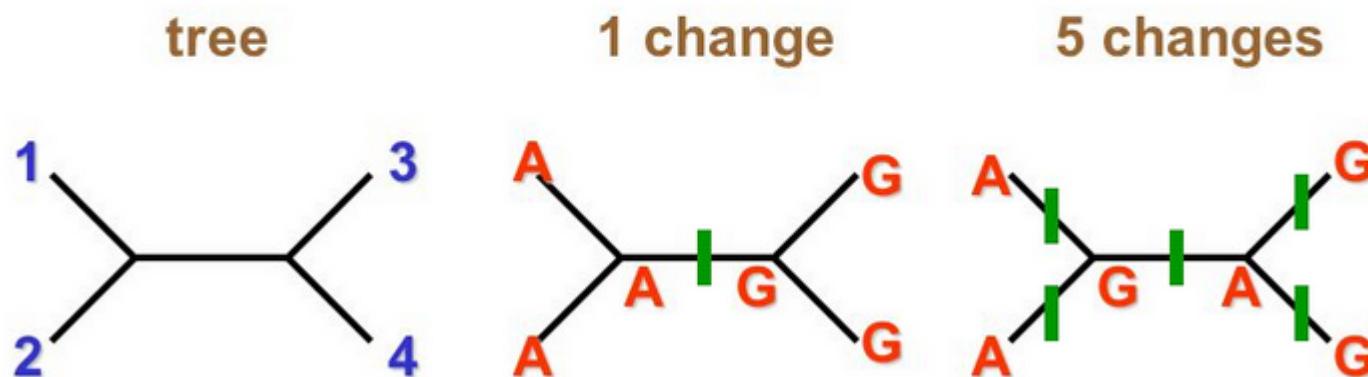
5. Start the algorithm again, considering the pair of joined neighbors as a single taxon and using the distances calculated in the previous step

# 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.



So, how to find the best tree?

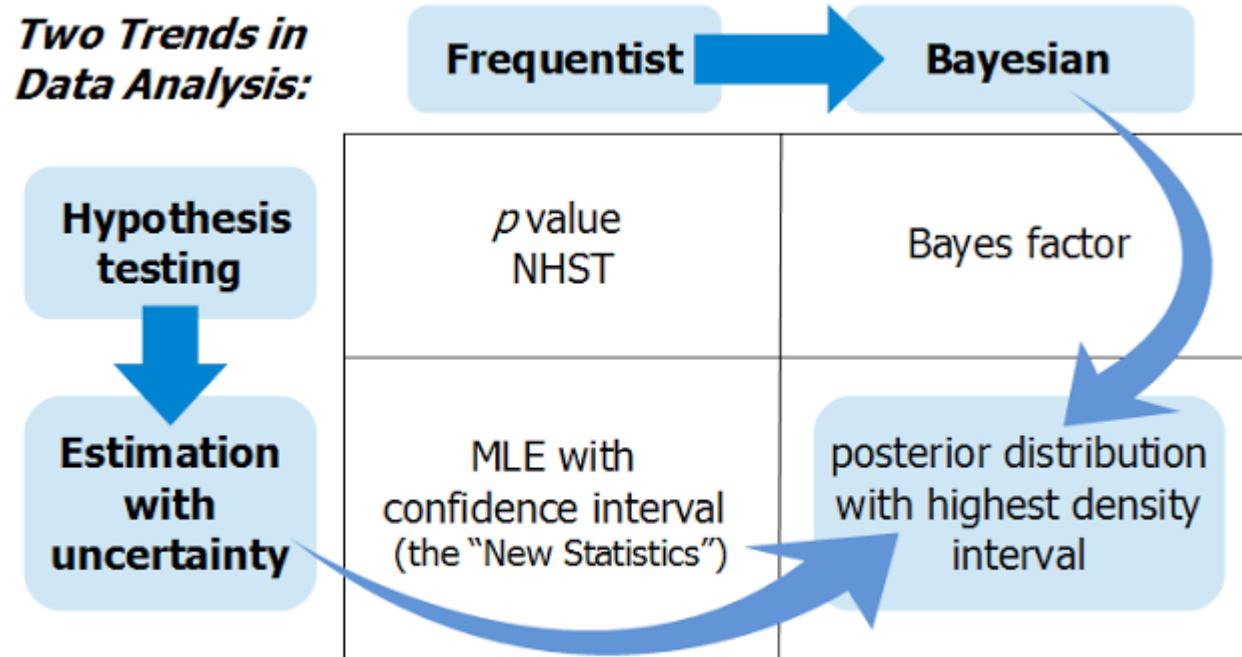
**Exhaustive** search: make ALL trees first, and then see which one best fits the data (you need an optimality criterion)

**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.

Different software differ in searching heuristics

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	$3 \times 10^{74}$

## Probabilistic methods render themselves for testing



Copyright © 2015 John K. Kruschke

## 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.



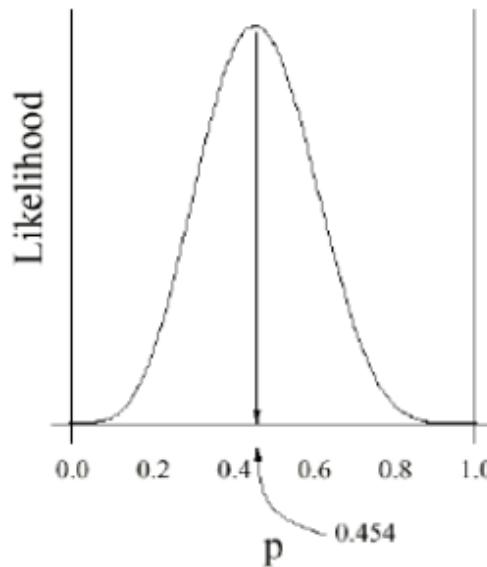
If tosses are all **independent**, and all have the same **unknown heads probability  $p$** , then the observing sequence of tosses:

**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>

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

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

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

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

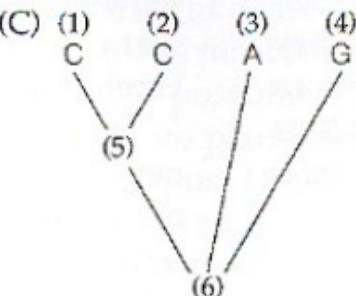
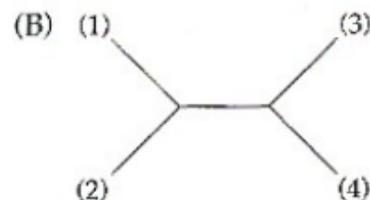
---

<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)

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

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

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

- 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.

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

Models can be very complex to capture different processes, increasing the number of parameters

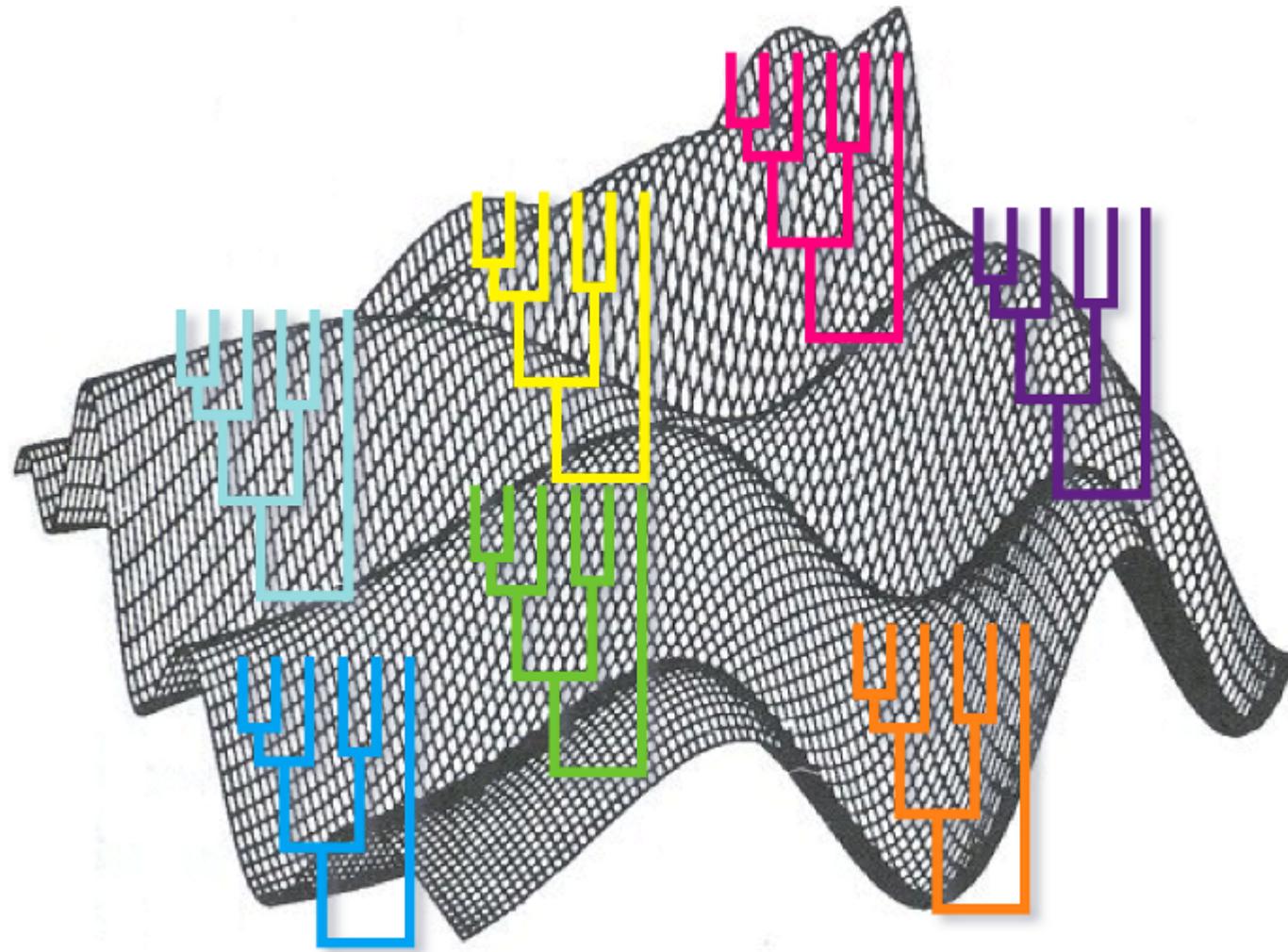
- i.e. fast and slow evolving sites

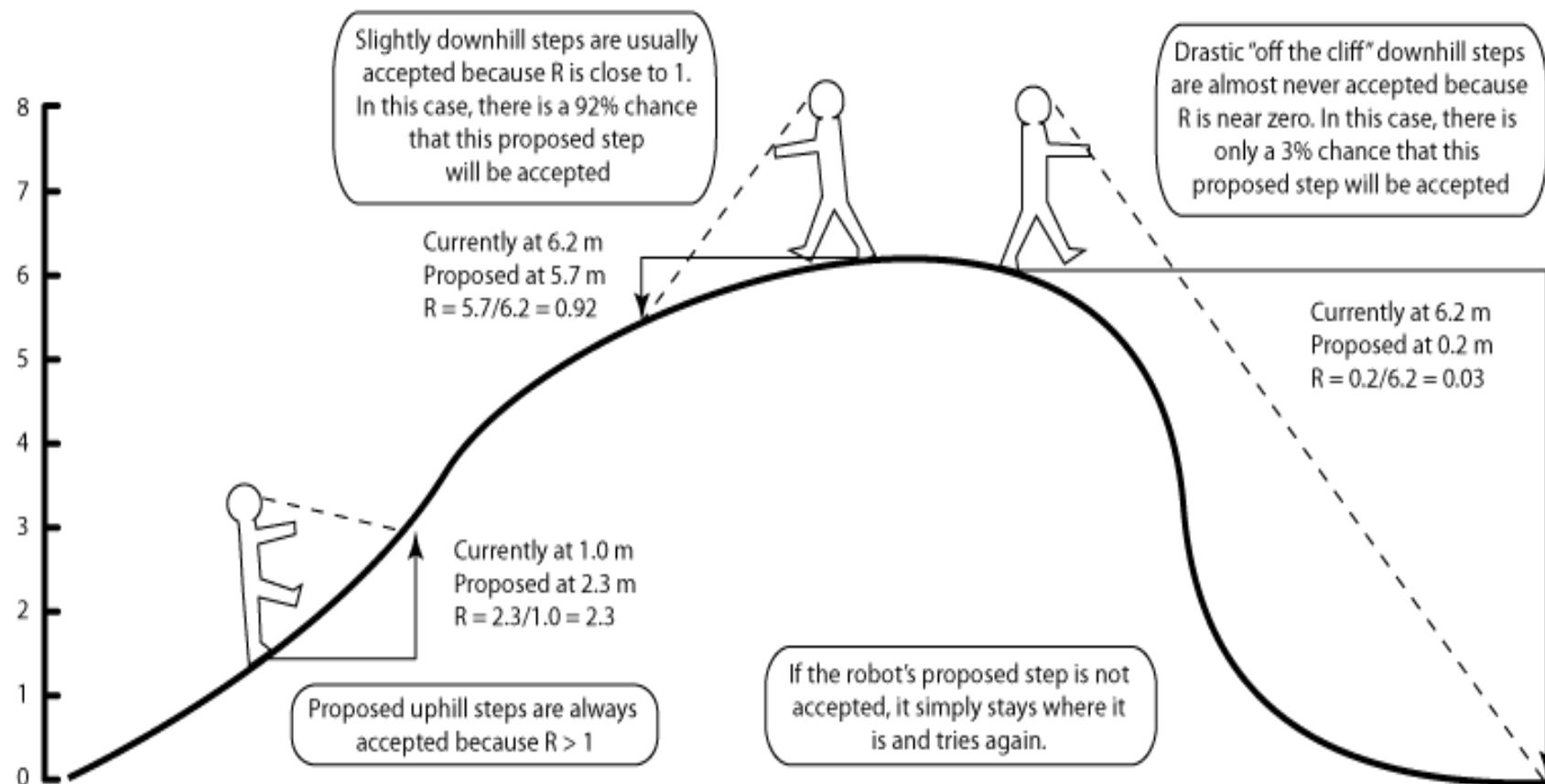
Models can be specifically made for specific groups of sequences

- i.e. for mitochondrially encoded proteins

Model choice can influence results.

Each model would induce a different likelihood landscape





# 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>.

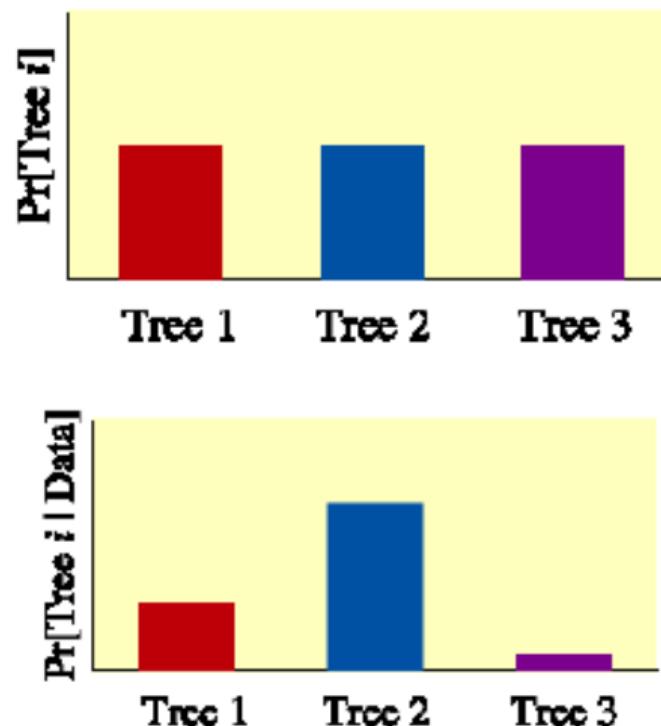
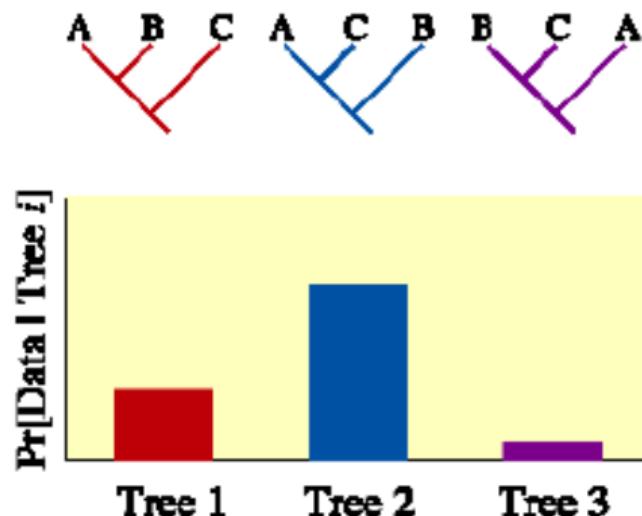
*Thomas Bayes (1702-1761)*



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

## 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.

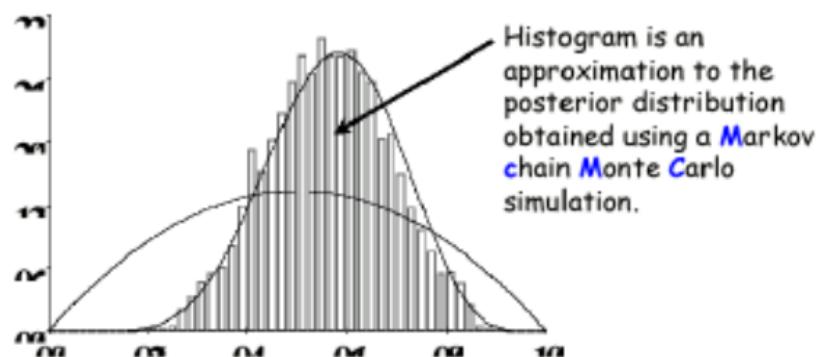


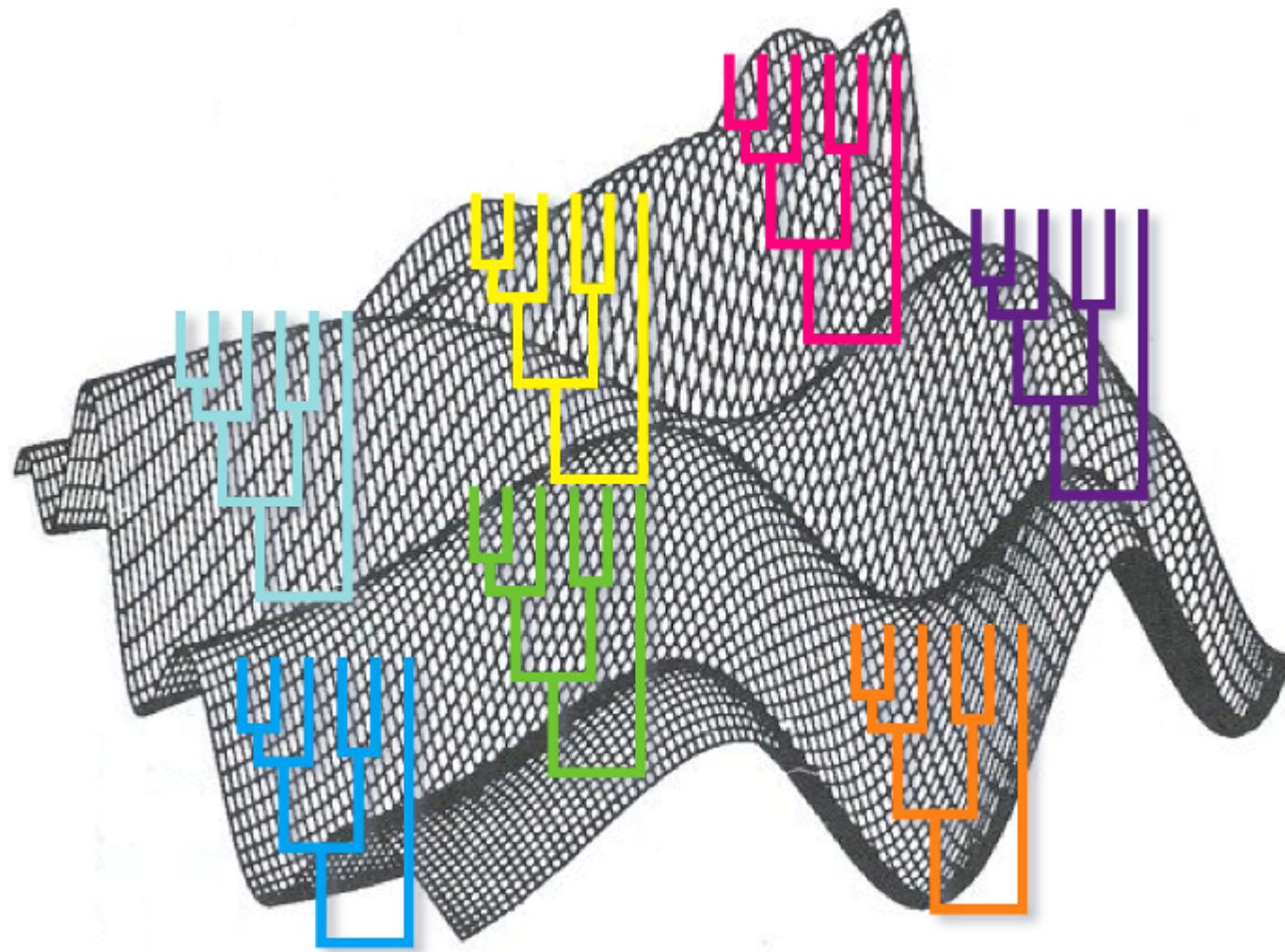
- $P(D/\theta)$  The **likelihood** is proportional to the probability of the observations (data sets) conditional on the tree.

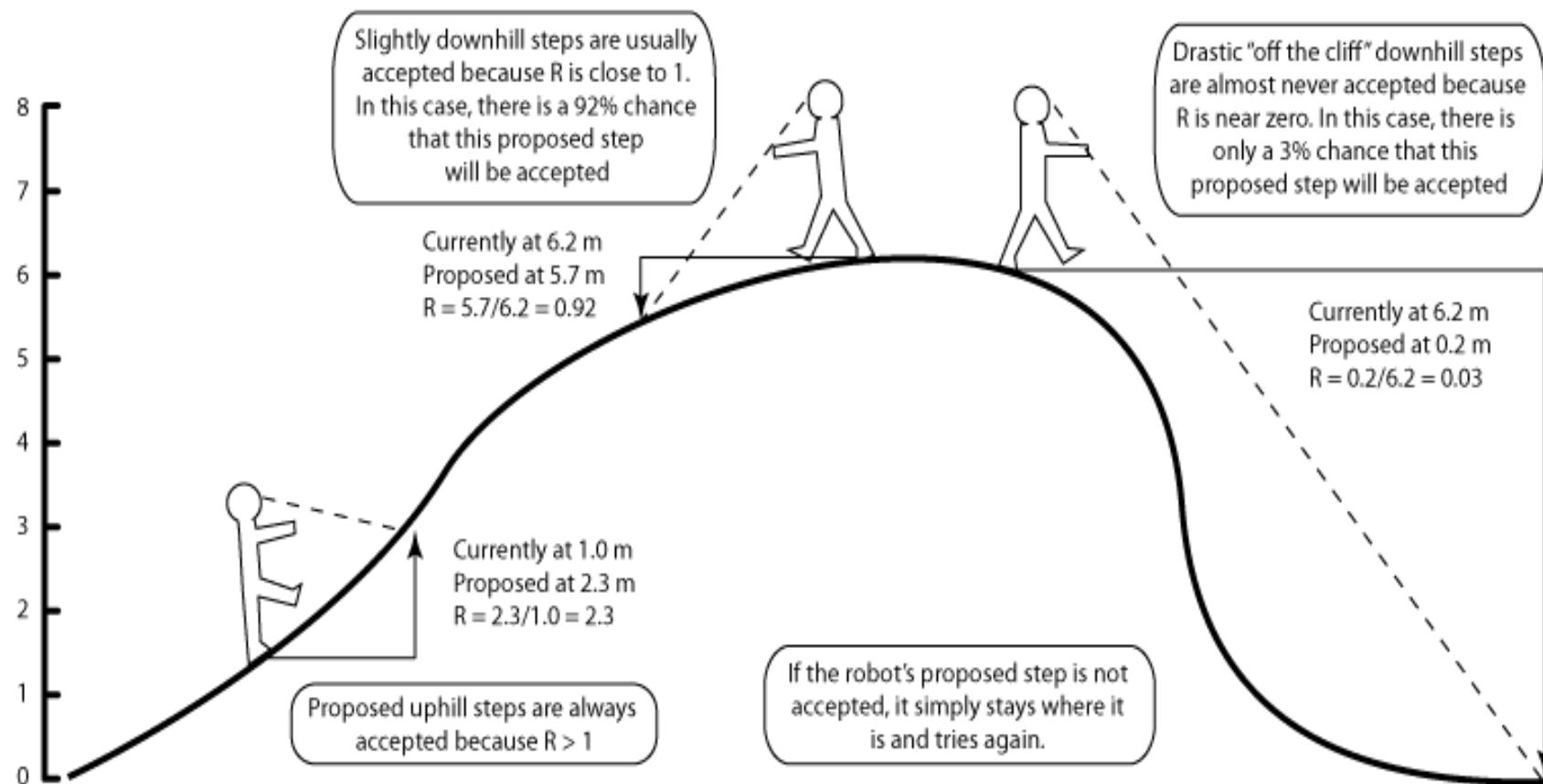
## How to find the solution

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>





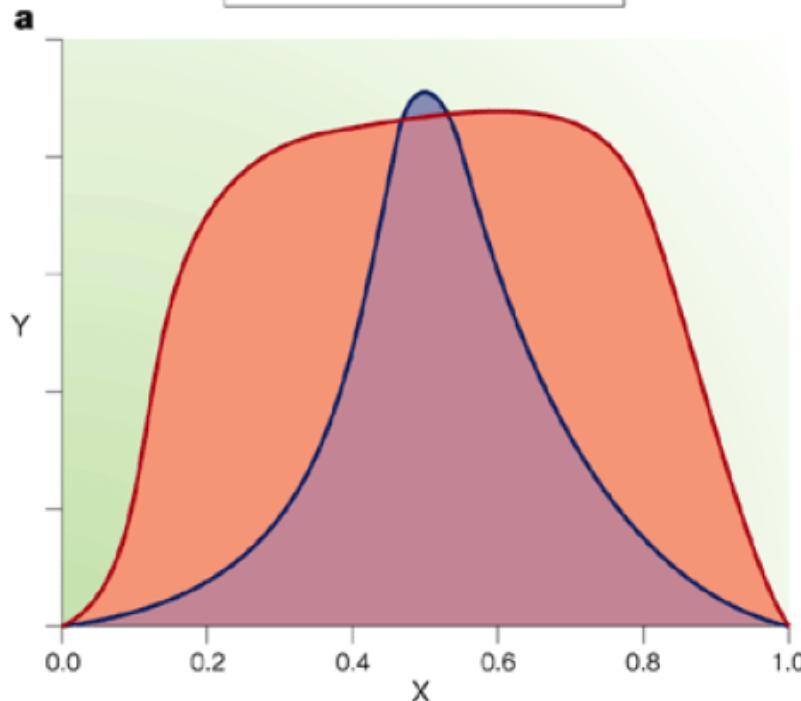


A fundamental difference:

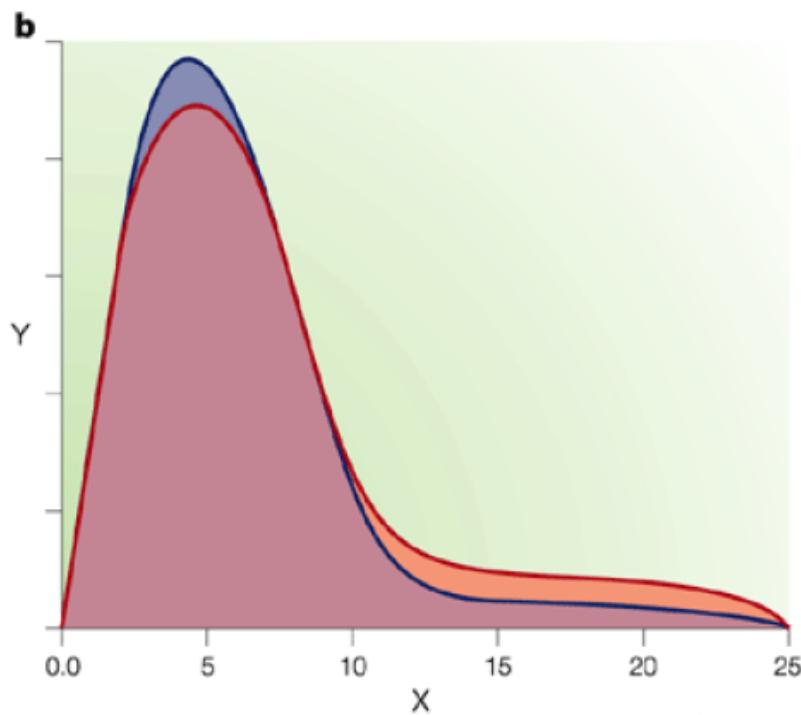
**ML** commonly uses **Joint Estimation**: finding the highest point in the parameter landscape

**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



ML favours A  
Bayesian favours B



ML favours A  
Bayesian favours B

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

# More to come in the following days

Xiaofan Zhou

Alexey Kozlov

Stephen Crotty

Laura Kubatko

Mario Dos Reis

Tracy Heath



Barcelona  
Biomedical  
Research  
Park



# Orthology Part I: concepts and implications

Toni Gabaldón

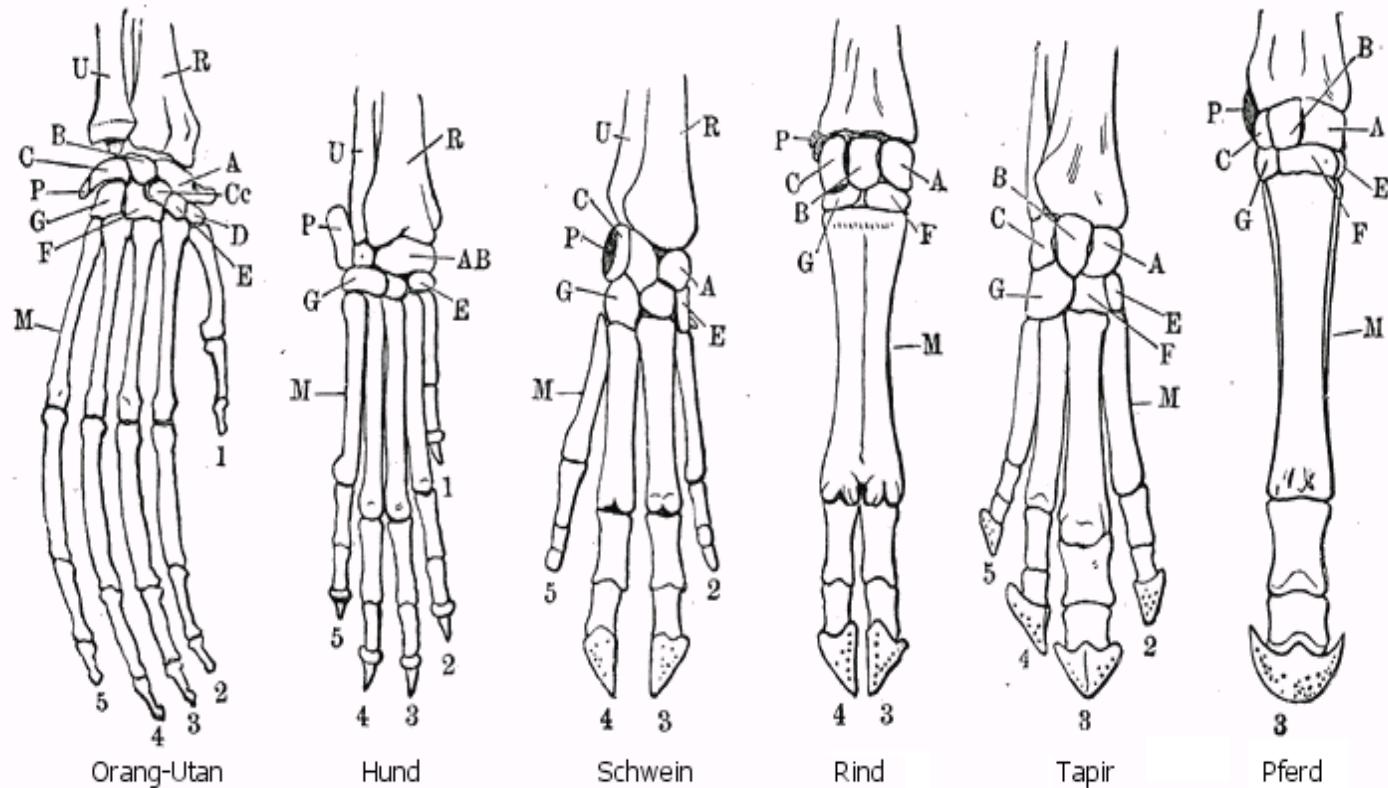
Centre for Genomic Regulation (CRG), Barcelona  
[\(tgabaldon@crg.es\)](mailto:tgabaldon@crg.es)  
<http://gabaldonlab.crg.es>





Richard Owen.

# Homology

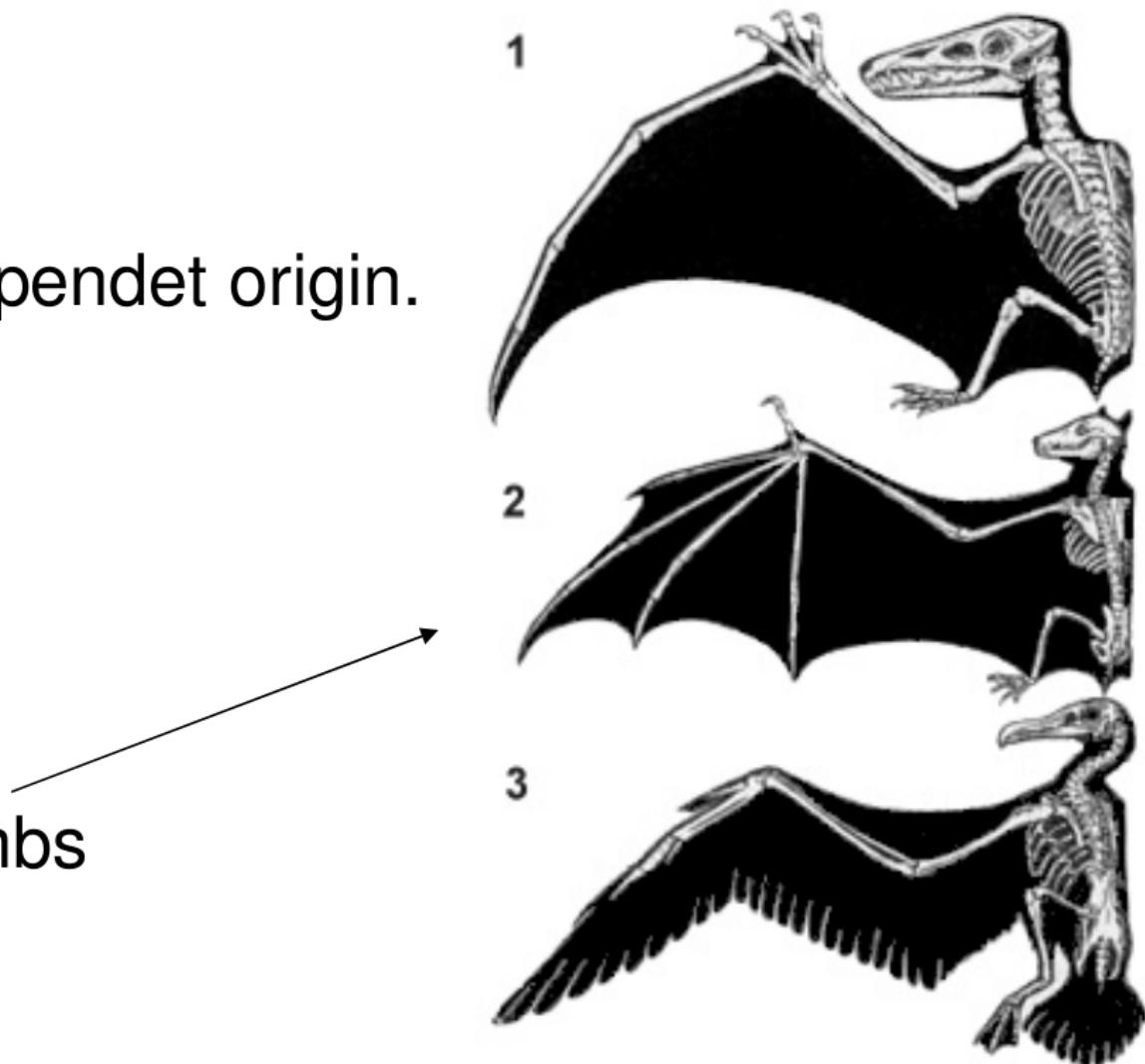


R Radius (Speiche), U Ulna (Elle), A-G, Cc, P Knochen des Carpus (Handwurzel): A Scaphoideum (Kahnbein), B Lunare (Mondbein), C Triquetrum (dreieckiges Bein), D Trapezium (großes vieleckiges Bein), E Trapezoïdes (kleines vieleckiges Bein), F Capitatum (Kopfbein), G Hamatum (Hafenbein), P Pisiforme (Erbsenbein), Cc Centrale Carpi, M Metacarpus (Mittelhand).  
Die Zahlen 1-5 bezeichnen die Finger (1 Daumen, 5 kleiner Finger).

"the same organ in different animals under every variety of form and function" R. Owen  
 → organs in two species are homologous only if the same structure was present in their last common ancestor. Homology → common ancestry

Analogous structures:  
Similar function but independent origin.

Homologous as forelimbs  
But  
Analogous as wings



Extension of the concept of homology to sequences:

*Two sequences are homologous if they share common ancestry*

AAB24882	TYHMCQFHCRYVNNHSGE <b>KLYECNERSKAFSCP</b> SHLQCHKR <b>RQIGEKTHEHNQCGKAFPT</b> 60
AAB24881	----- <b>YECNQCGKAF</b> AQHSSLKCHYRTHIGE <b>KPYECNQCGKAFSK</b> 40
	*****: . ***: * * :** * :*****, :* *****..

AAB24882	<b>PSHLQYHERHTHTGEKPYE</b> CHQCG <b>QAFKKCSLLQRHKRTHTGEKPYE-CNQCGKAF</b> AQ- 116
AAB24881	<b>HSHLQCHKR</b> THTGEKPYE <b>CNQCGKAFSQHGLLQRHKRTHTGEKPYMNVINMVKPLHNS</b> 98
	**** *:*****:*****:***.: .*****:***** : *.: :

## **Important:** Similarity and Homology

Similarity and homology are often confused. e.g.

“the sequences are 50% homologous”, “these two sequences are highly homologous”

Why is this incorrect?

Where does the confusion comes from?

# Detour

## Sequence similarity, homology detection and blast database queries

```
AAB24882  TYHMCQFHCRYVNNHSGEKLYECNERSKAFCSCPShLQCHKRRQIGEKTHEHNQCGKAFPT 60
AAB24881  -----YECNQCGKAFAQHSSLKCHYRTHIGEKPYECNQCGKAFSK 40
                         ****: .***: * *;** * ;****,:* *****.. .

AAB24882  PSHLQYHERHTHTGEKPYECHQCGQAFKKCSLLQRHKRTHTGEKPYE-CNQCGKAFAQ- 116
AAB24881  HSHLQCHKRTHTGEKPYECNQCGKAFSQHGLLQRHKRTHTGEKPYMNVINMVKPLHNS 98
                         *** *:*****:***:*. : .*****: : *.: : 
```

Are these two sequences **significantly** similar?  
(i.e. how likely is that such an alignment is the result of chance?)

>  ref|NP\_114344.1| G NADH dehydrogenase subunit 5 [Macaca sylvanus]  
Length=603

GENE ID: 803075 ND5 | NADH dehydrogenase subunit 5 [Macaca sylvanus]  
(10 or fewer PubMed links)

Score = 796 bits (2056), Expect = 0.0, Method: Compositional matrix adjust.  
Identities = 438/564 (77%), Positives = 478/564 (84%), Gaps = 0/564 (0%)

Query	24	VNPNNKKNSYPHYVKSIVASTFIISLFPTTMFMCLDQEVIISNWHWATTQTTQLSLSFKLD	83
Sbjct	24	+NPNKK+ YP+YVK+ V FI SL TT++M L+QE II +WHW TQT L+LSFKLD	
Query	84	INPNKKHLYPNVYVKTAVMYAFITSLSSTTLYMFLNQETIIWSWHWMMTQTLSTLSFKLD	83
Sbjct	84	YFSMMFIPVALFVTWSIMEFSLWYMNSDPNIQFFKYLLIFLITMLILVTANNLFQLFIG	143
Query	144	YFSMMF P+AL TWSIMEFSLWYM+SDPNI+QFFKYLLIFLITMLILVTANNLFQ FIG	
Sbjct	84	YFSMMFTPIALLTWSIMEFSLWYMSSDPNIDQFFKYLLIFLITMLILVTANNLFQFFIG	143
Query	144	YFSMMFTPIALLTWSIMEFSLWYMSSDPNIDQFFKYLLIFLITMLILVTANNLFQFFIG	
Sbjct	144	WEVGVGIMSFLLISWWYARADANTAAIQAVLYNRIGDIGFILALAWFILHSNSWDPQQMAL	203
		WEVG+GIMSFLLISWW+AR DANTAAIQA+LYNRIGDIG IL + WF+LH NSWD QQM	
		WEGMGIMSFLLISWWHARTDANTAAIQAILYNRIGDIGLILTMTWFLHYNSWDFQQMLA	203

# Alignment scores are sums of residue pairing scores according to a scoring Matrix

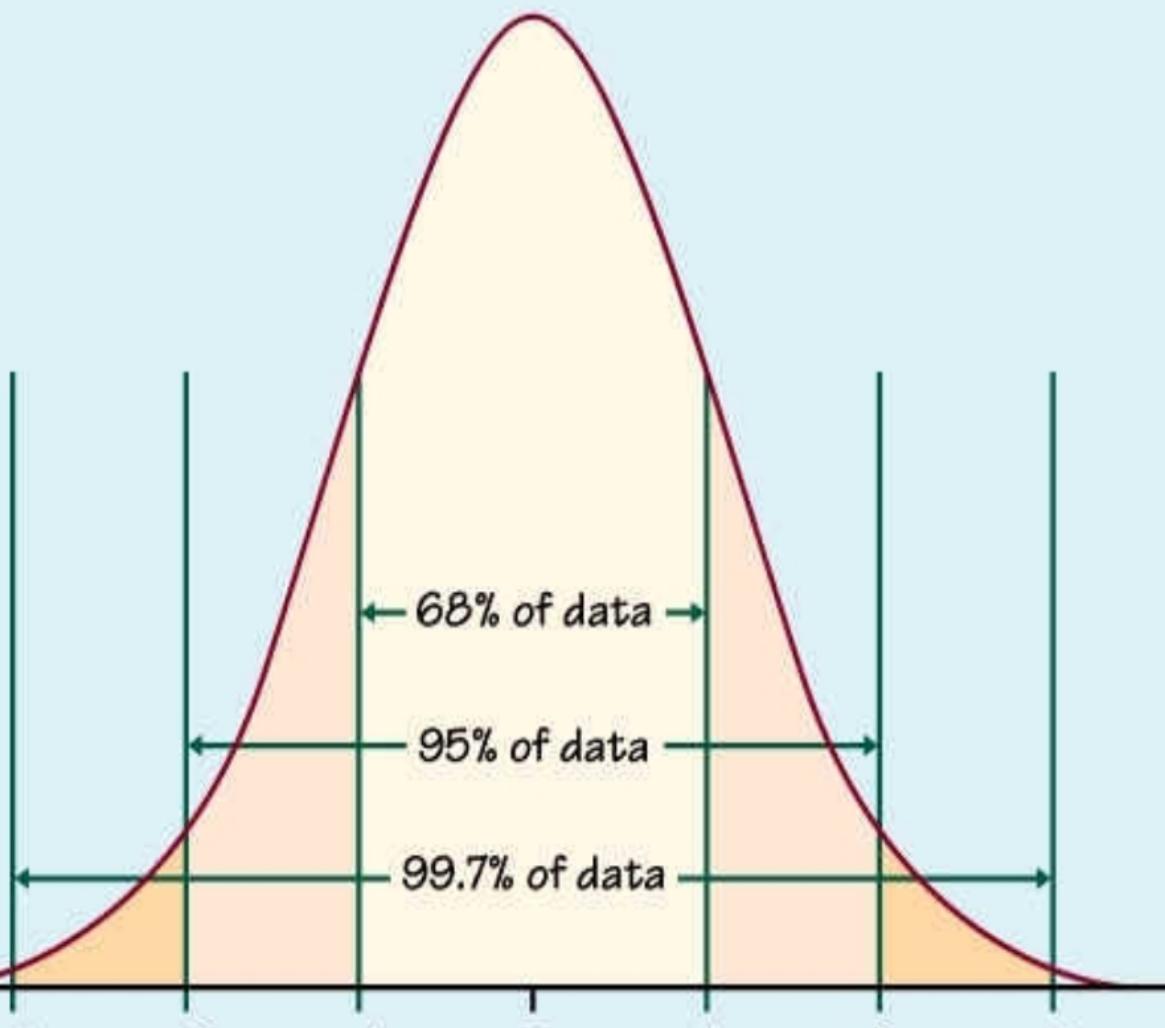
	BLOSUM62																				
A	4																				
R	-1	5																			
N	-2	0	6																		
D	-2	-2	1	6																	
C	0	-3	-3	-3	9																
Q	-1	1	0	0	-3																
E	-1	0	0	(2)	-4																
G	0	-2	0	-1	-3	-2	-2	-2	6												
H	-2	0	1	-1	-3	0	0	-2	8												
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	(4)											
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4										
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5									
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5								
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6							
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7						
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4					
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5				
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	(11)			
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7		
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4	
X	0	-1	-1	-1	-2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-2	0	0	-2	-1	-1	
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	X

Positive for chemically similar substitution

Common amino acids have low weights

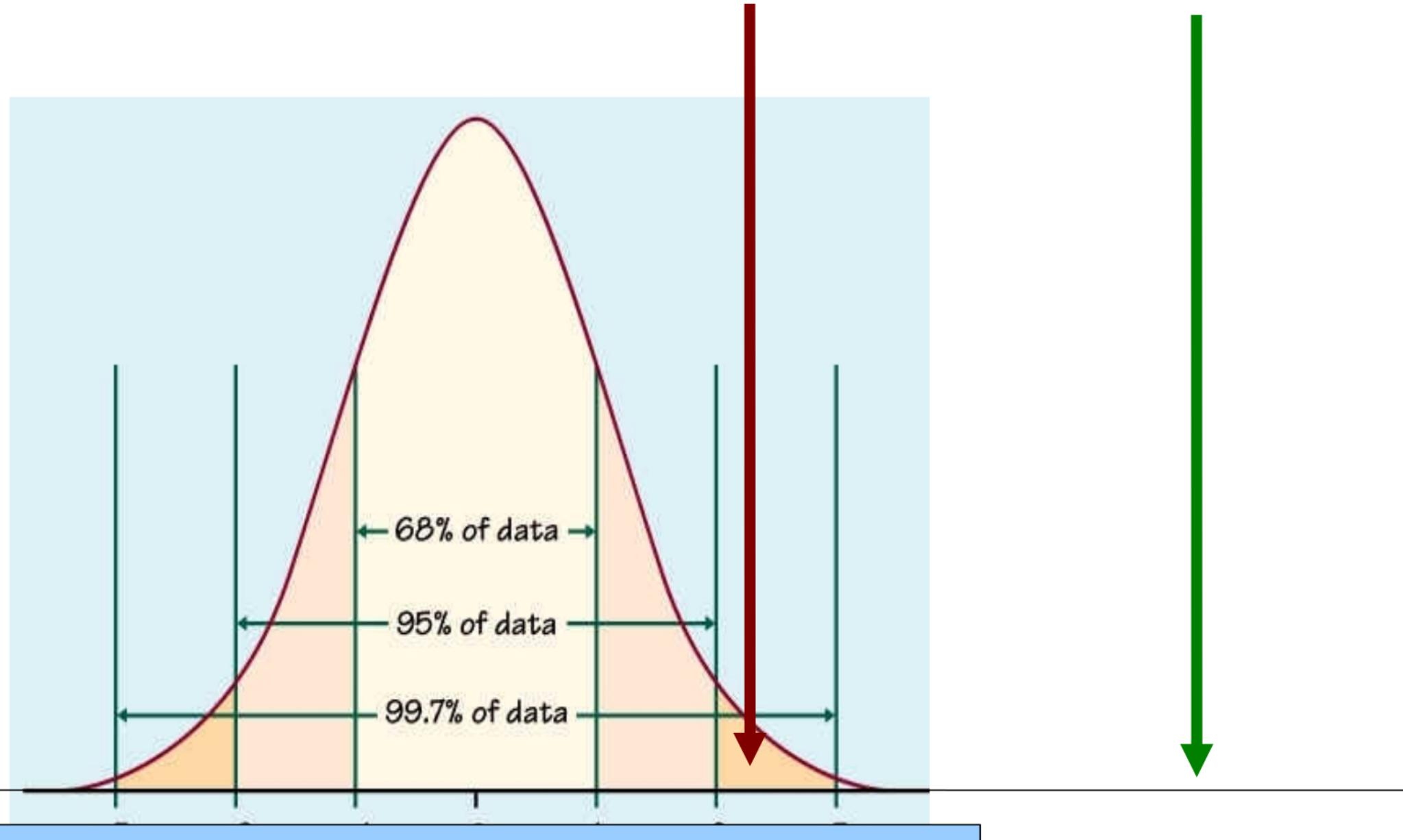
Rare amino acids have high weights

# Distribution of scores in comparisons of **random\***-sequences



\* considering the representation of the different amino acids (nucleotides) in a DataBase

Your score



$$E = m n 2^{-S'}$$

Length of hit and query

Normalized score

E-value (Expectation value)= the number of sequences that would be expected to have that **score** (or higher) if the query sequence were compared against a **database** containing unrelated sequences

>  ref|NP\_114344.1| G NADH dehydrogenase subunit 5 [Macaca sylvanus]  
Length=603

GENE ID: 803075 ND5 | NADH dehydrogenase subunit 5 [Macaca sylvanus]  
(10 or fewer PubMed links)

Score = 796 bits (2056), Expect = 0.0, Method: Compositional matrix adjust.  
Identities = 438/564 (77%), Positives = 478/564 (84%), Gaps = 0/564 (0%)

Query	24	VNPNNKKNSYPHYVKSIVASTFIISLFPTTMFMCLDQEVIISNWHWATTQTTQLSLSFKLD	83
Sbjct	24	+NPNKK+ YP+YVK+ V FI SL TT++M L+QE II +WHW TQT L+LSFKLD	83
Query	84	YFSMMFIPVALFTWSIIMEFSLWYMNSDPNIHQFFKYLLIFLITMLILVTANNLFQLFIG	143
Sbjct	84	YFSMMF P+AL TWSIIMEFSLWYM+SDPNI+QFFKYLLIFLITMLILVTANNLFQ FIG	143
Query	144	WEGVGIMSFLLISSWWYARADANIAAIQAVLYNRIGDIGFILALAWFILHSNSWDPQQMAL	203
Sbjct	144	WEG+GIMSFLLISSWW+AR DANTAAIQA+LYNRIGDIG IL + WF+LH NSWDFQQMLA	203

E-value

Coverage over the query

## From homology to orthology

- Homologs are sequences derived from a common ancestor...
- What are then orthologs?.... and paralogs?

## Are these sentences correct?

- Orthologs are homologous genes that have the same function
- Orthologs are homologous genes in different species, while paralogs are homologous genes in the same species
- The ortholog is the most similar sequence among the homologs in another species
- If gene A is orthologous to gene B, and gene B is orthologous to gene C, then A and C are orthologous to each other.
- Orthologs are genes that do not duplicate and, when they exist, they are always present in single copy
- After a duplication, the orthologous copy is the one that keeps the function of the ancestral gene



Fitch W.M.

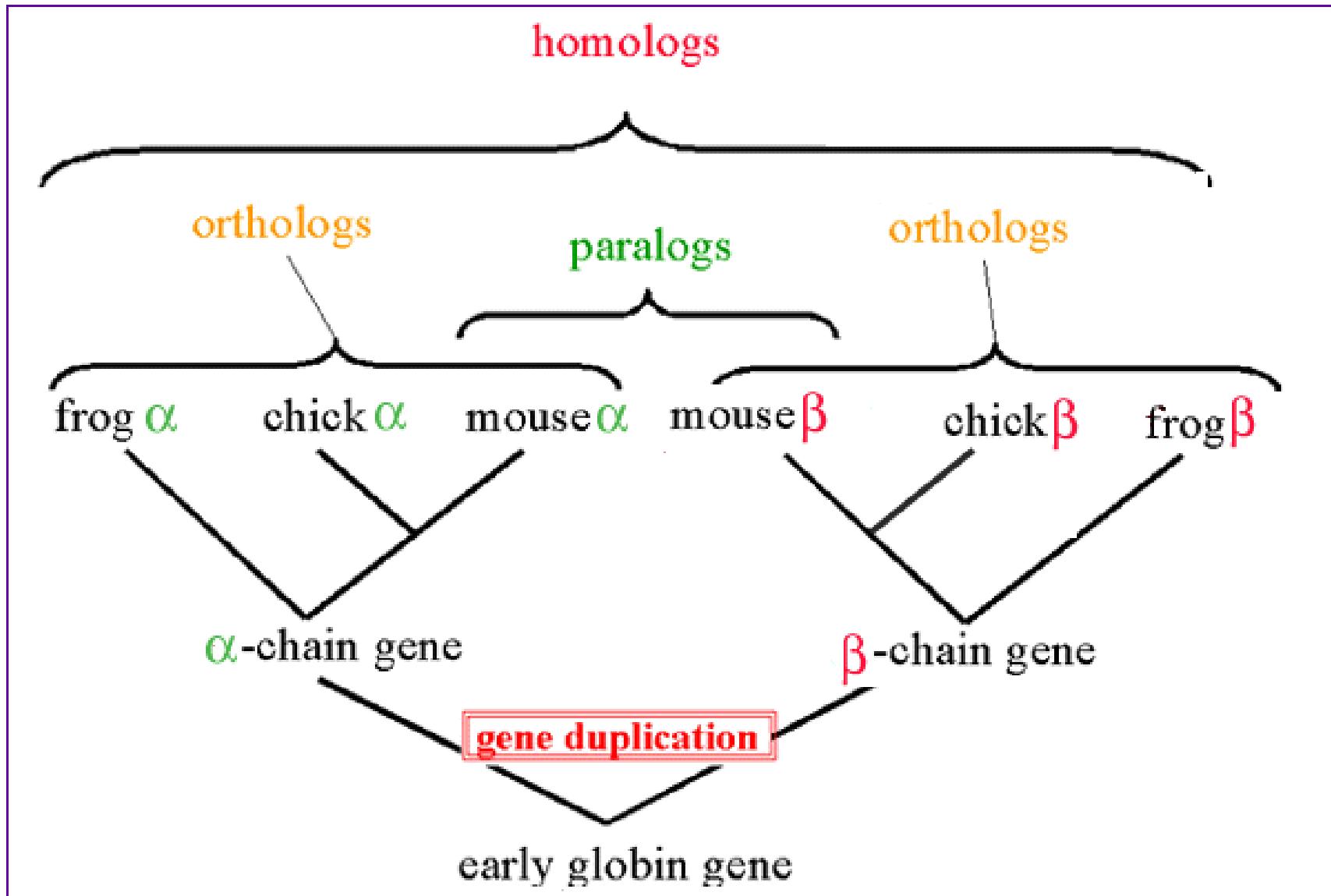
Distinguishing homologous from  
analogous proteins.

Syst. Zool. 1970; 19: 99-113

Original definition of orthology and paralogy by Walter Fitch (1970, Systematic Zoology 19:99-113):

*"Where the homology is **the result of gene duplication** so that both copies have descended side by side during the history of an organism, (for example, alpha and beta hemoglobin) the genes should be called **paralogous** (para = in parallel).*

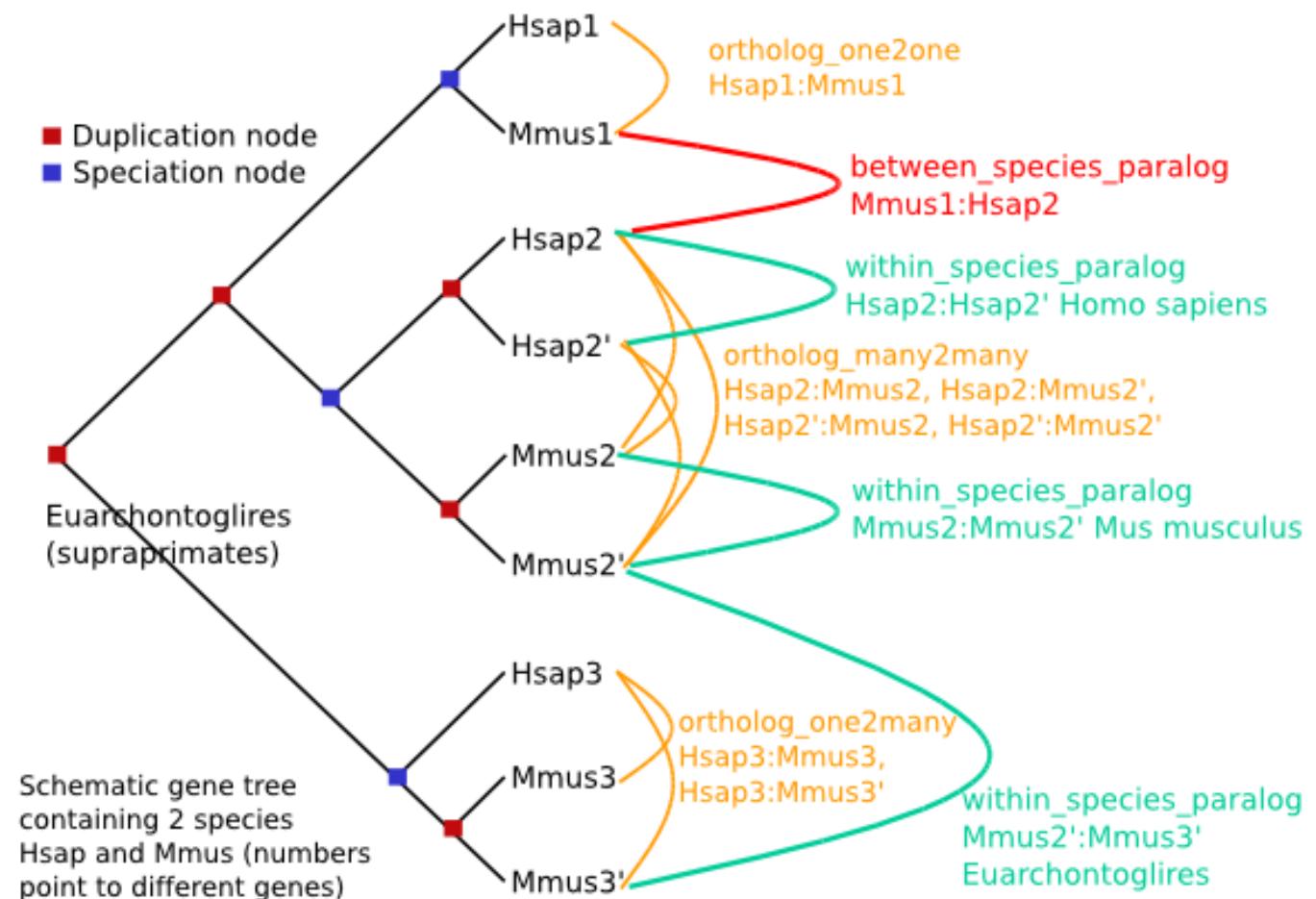
*Where the homology is **the result of speciation** so that the history of the gene reflects the history of the species (for example alpha hemoglobin in man and mouse) the genes should be called **orthologous** (ortho = exact)."*



# Corollary:

- Orthology definition is purely on evolutionary terms (not functional, not synteny...)
- There is no limit on the number of orthologs or paralogs that a given gene can have (when more than one ortholog exist, there is nothing such as “*the true ortholog*”)
- Many-to-Many orthology relationships do exist (co-orthology)
- No limit on how ancient/recent is the ancestral relationship of orthologs and paralogs
- Orthology is non-transitive (as opposed to homology)

Orthology relationships can be complex, and intricate



# Why predicting orthology is important?

- **Important implications for phylogeny:** only sets of orthologous genes are expected to reflect the underlying species evolution (although there are many exceptions)
- The most exact way of **comparing two (or more) genomes** in terms of their gene content. Necessary to uncover how genomes evolve.
- Implications for **functional inference:** orthologs, as compared to paralogs, are more likely to share the same function

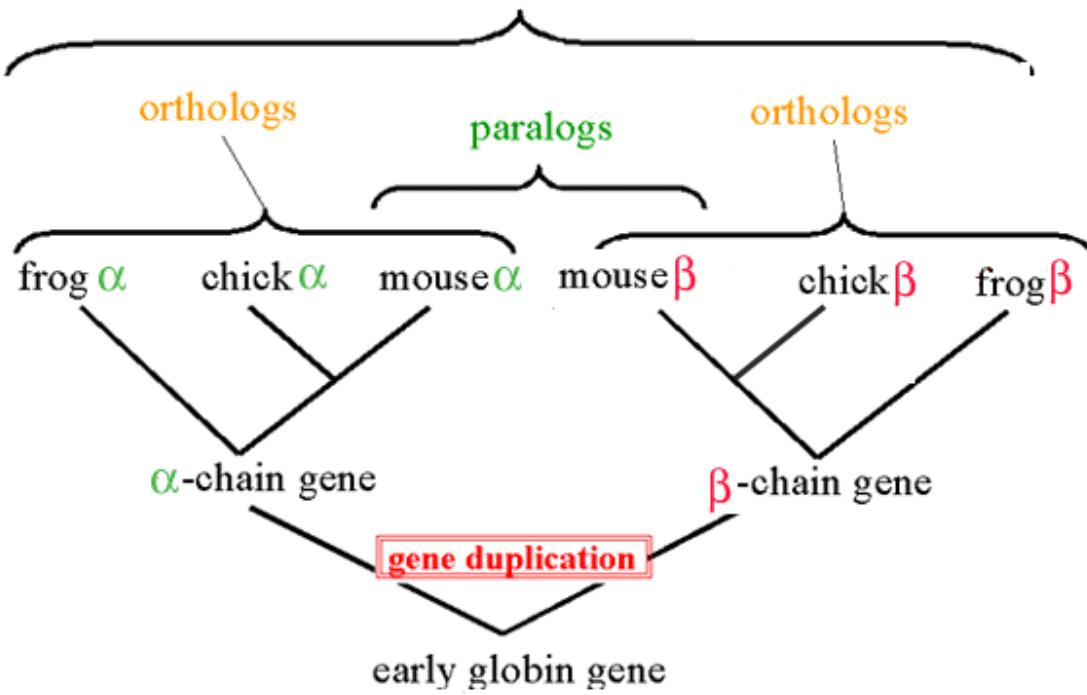
# Why predicting orthology is important?

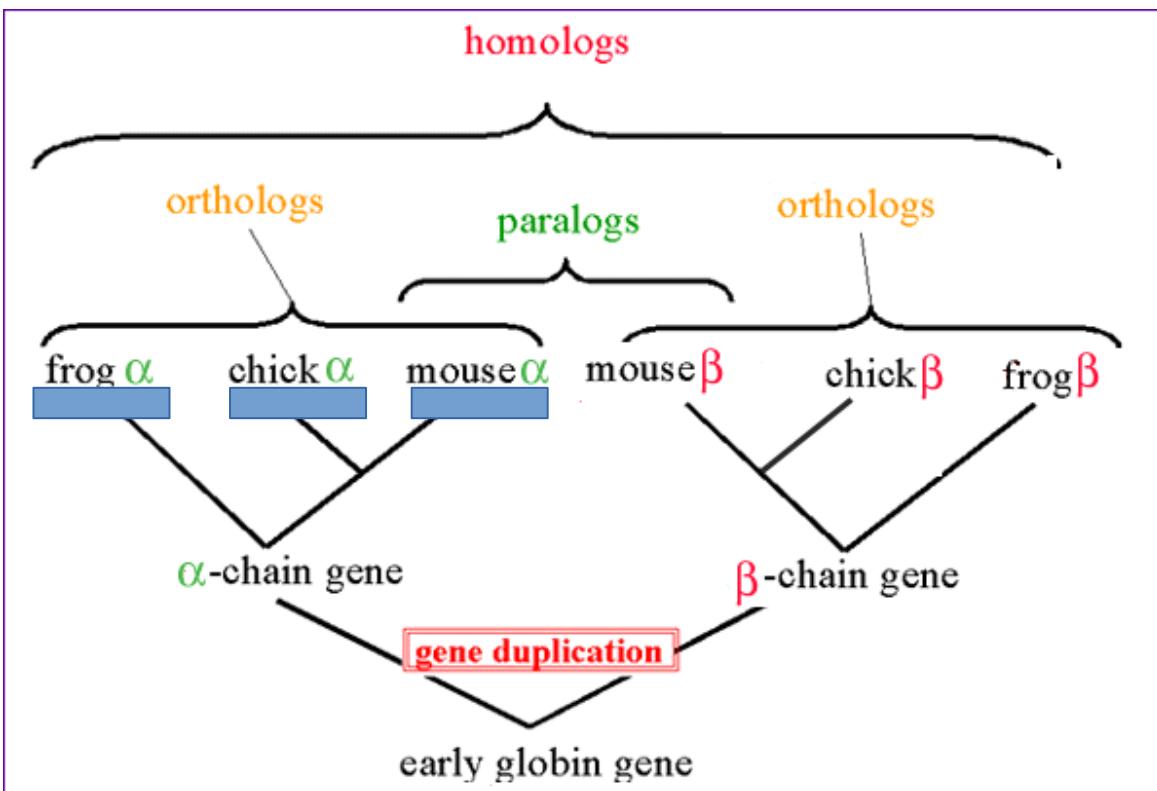
- **Important implications for phylogeny:** only sets of orthologous genes are expected to reflect the underlying species evolution (although there are many exceptions)

Where the homology is **the result of speciation** so that the history of the gene reflects the history of the species (for example alpha hemoglobin in man and mouse) the genes should be called **orthologous** (ortho = exact).

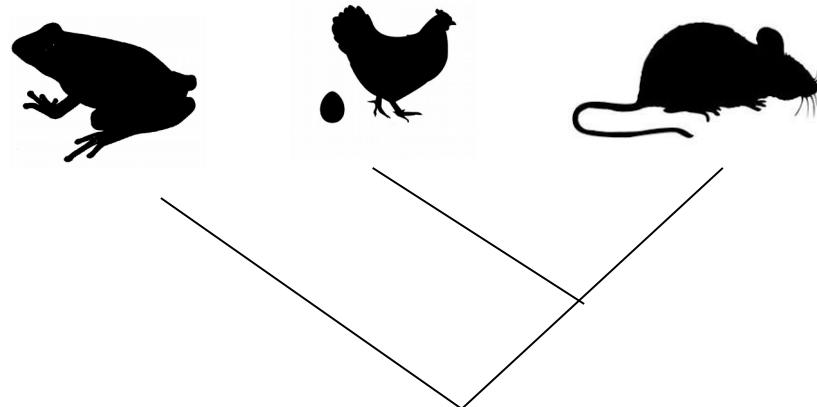


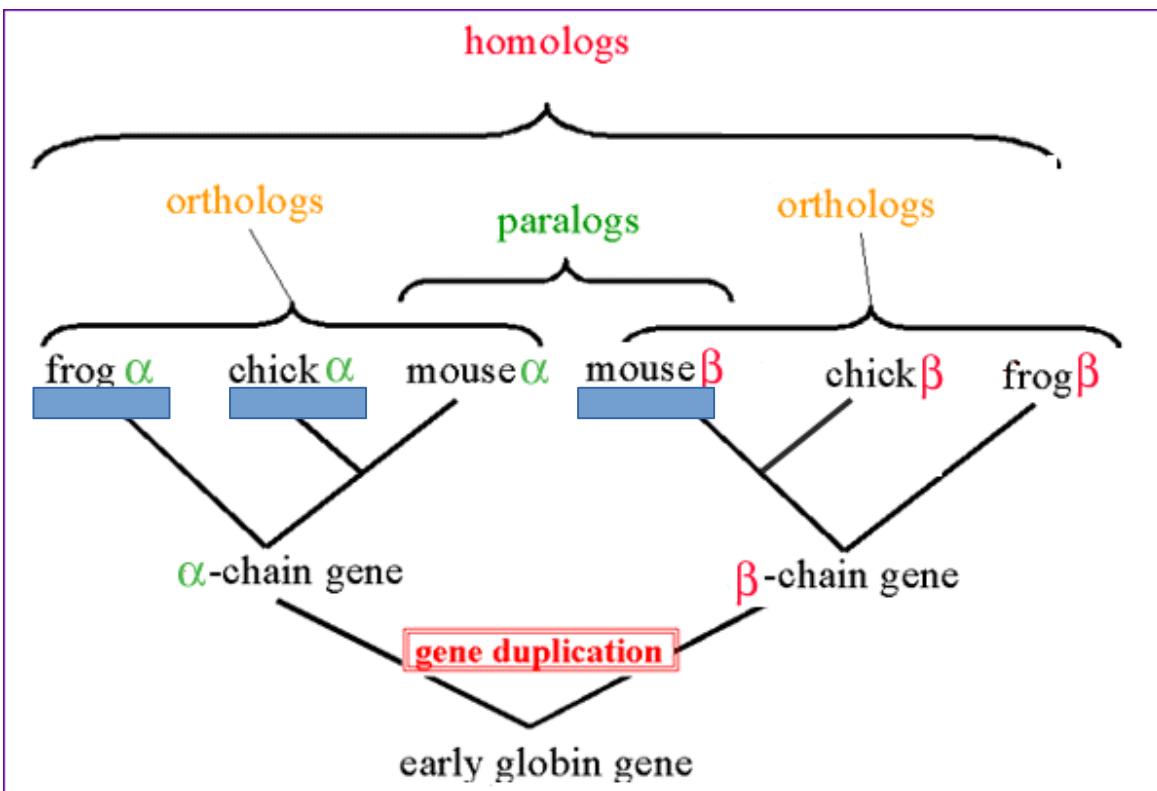
homologs



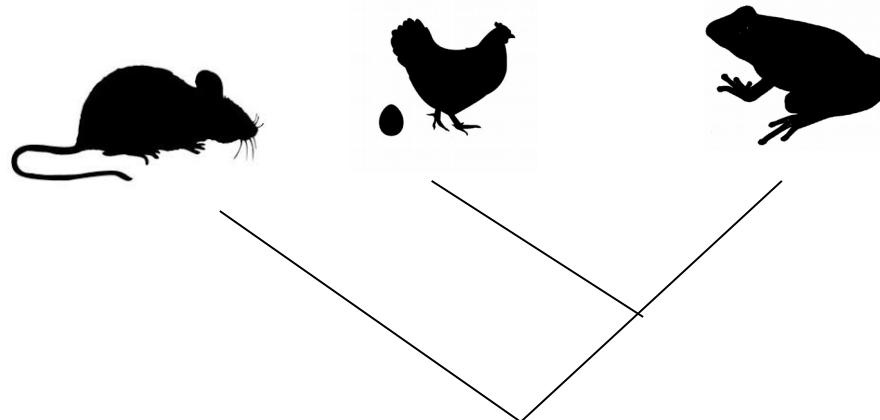


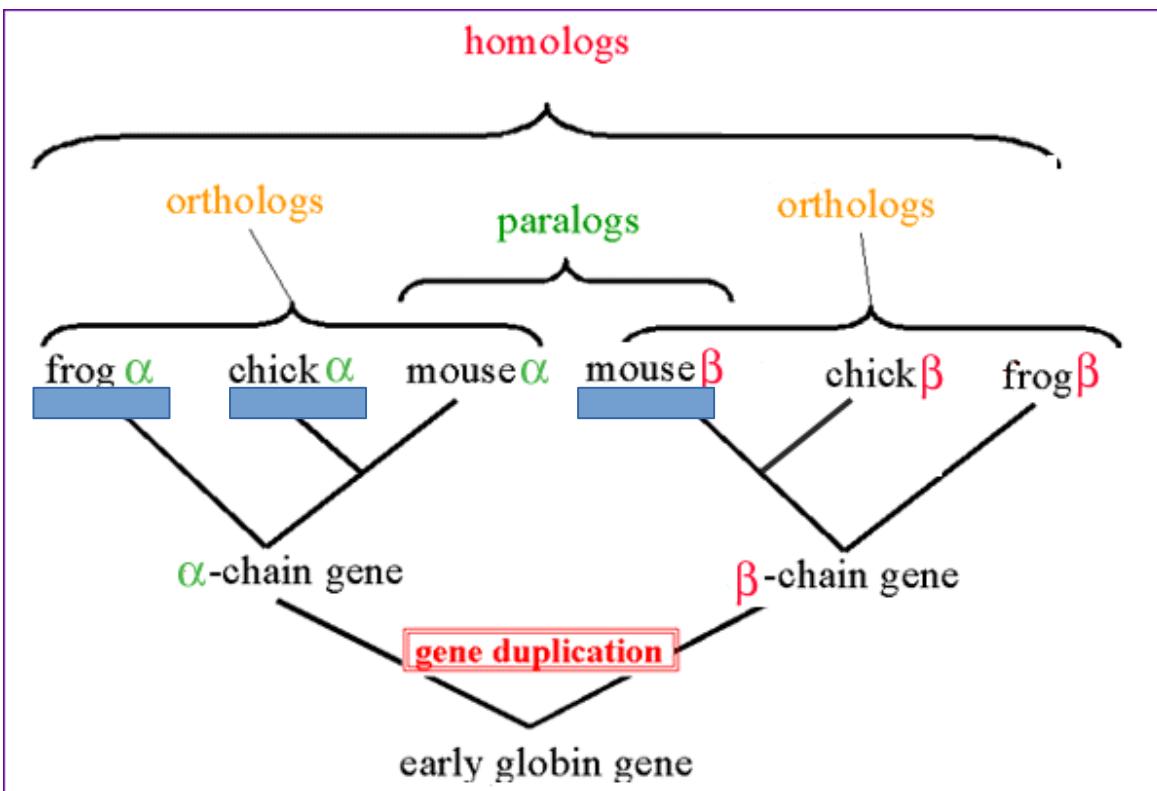
Tree from orthologous dataset:



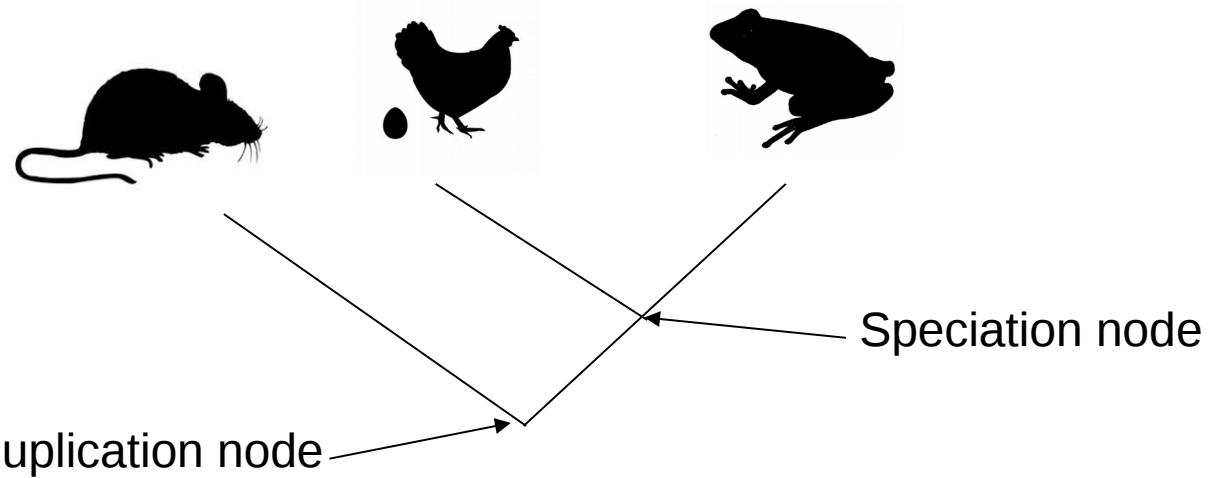


Tree from non-orthologous dataset:





Tree from non-orthologous dataset: NOT a SPECIES TREE



Seems easy, but it's not\*

\*at least not always.

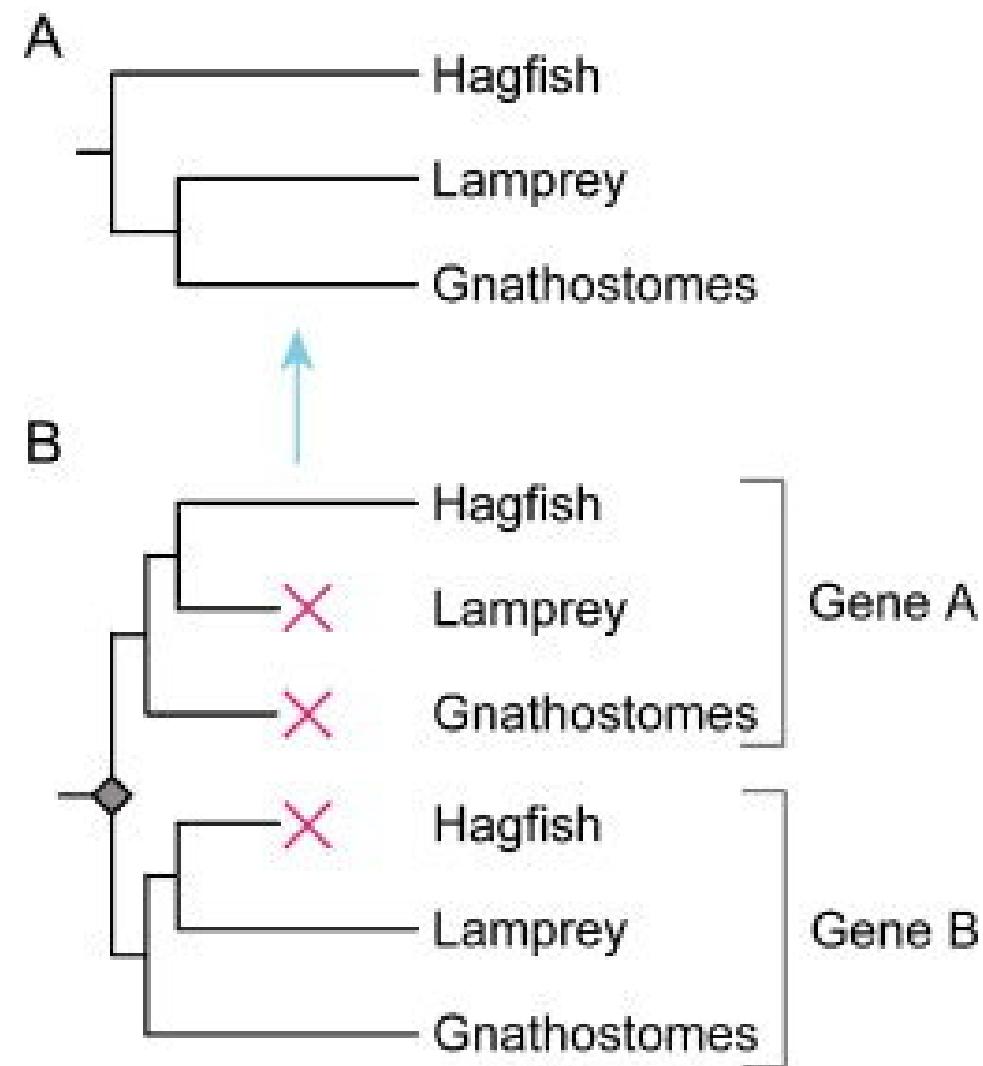
## The hidden paralogy problem

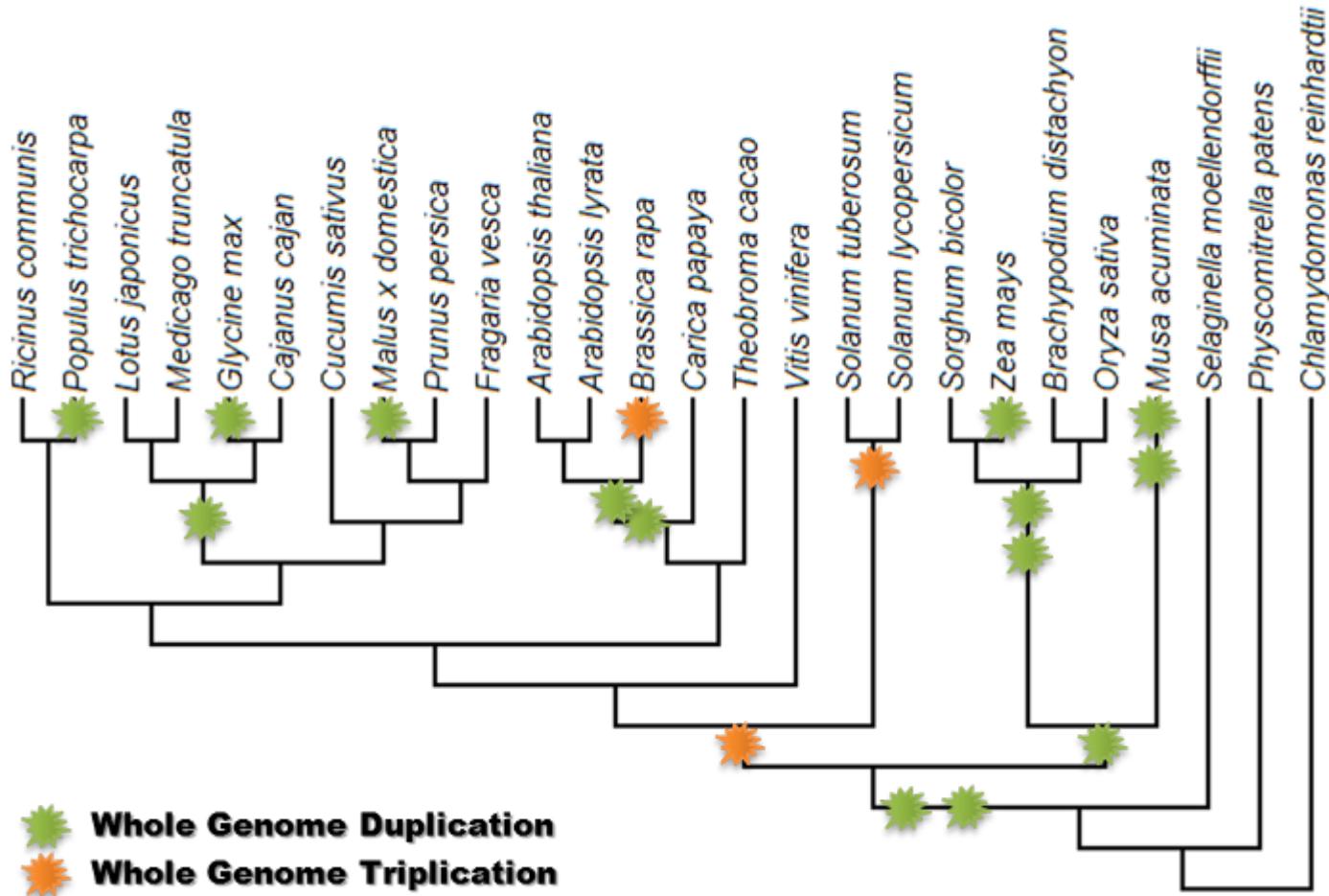
Complex duplications and loss patterns can result in paralogous genes being recovered as putative orthologs by most methods, resulting in faulty phylogenetic relationships.

This problem is exacerbated following whole genome duplication events, usually followed by massive differential gene loss.

Annotation problems and fragmented genomes could result in similar patterns.

Increasing taxonomic coverage is one approach to alleviate this problem .

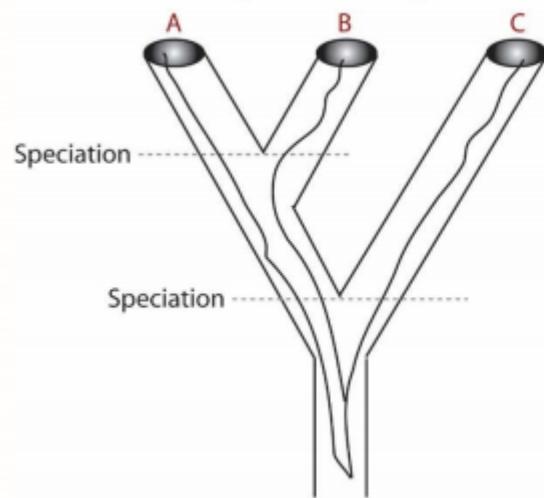




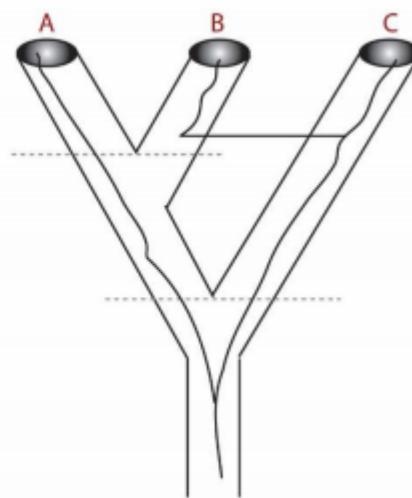
Even if you get a truly orthologous dataset.

Orthologous genes are not guaranteed to reflect the species tree!

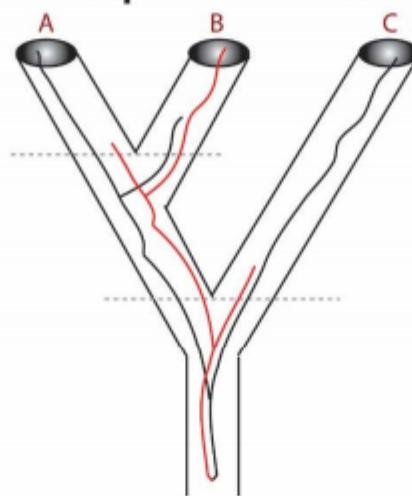
### Lineage Sorting



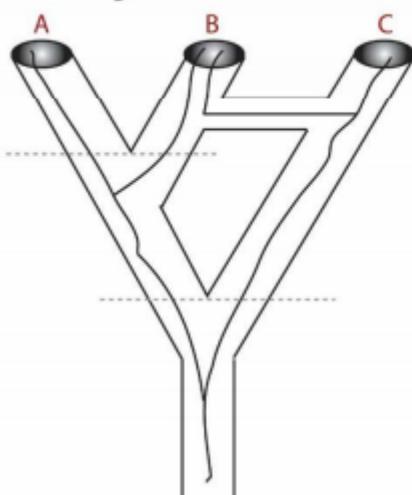
### Horizontal Gene Transfer



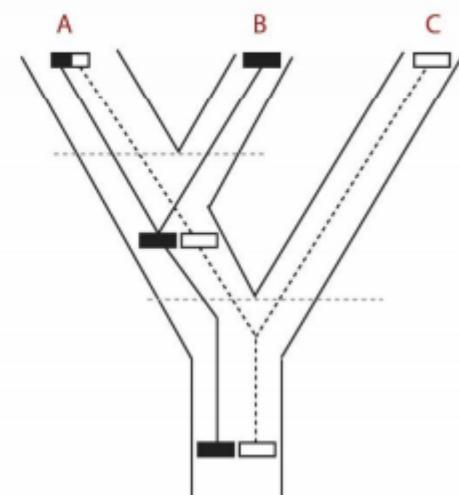
### Gene Duplication and Loss



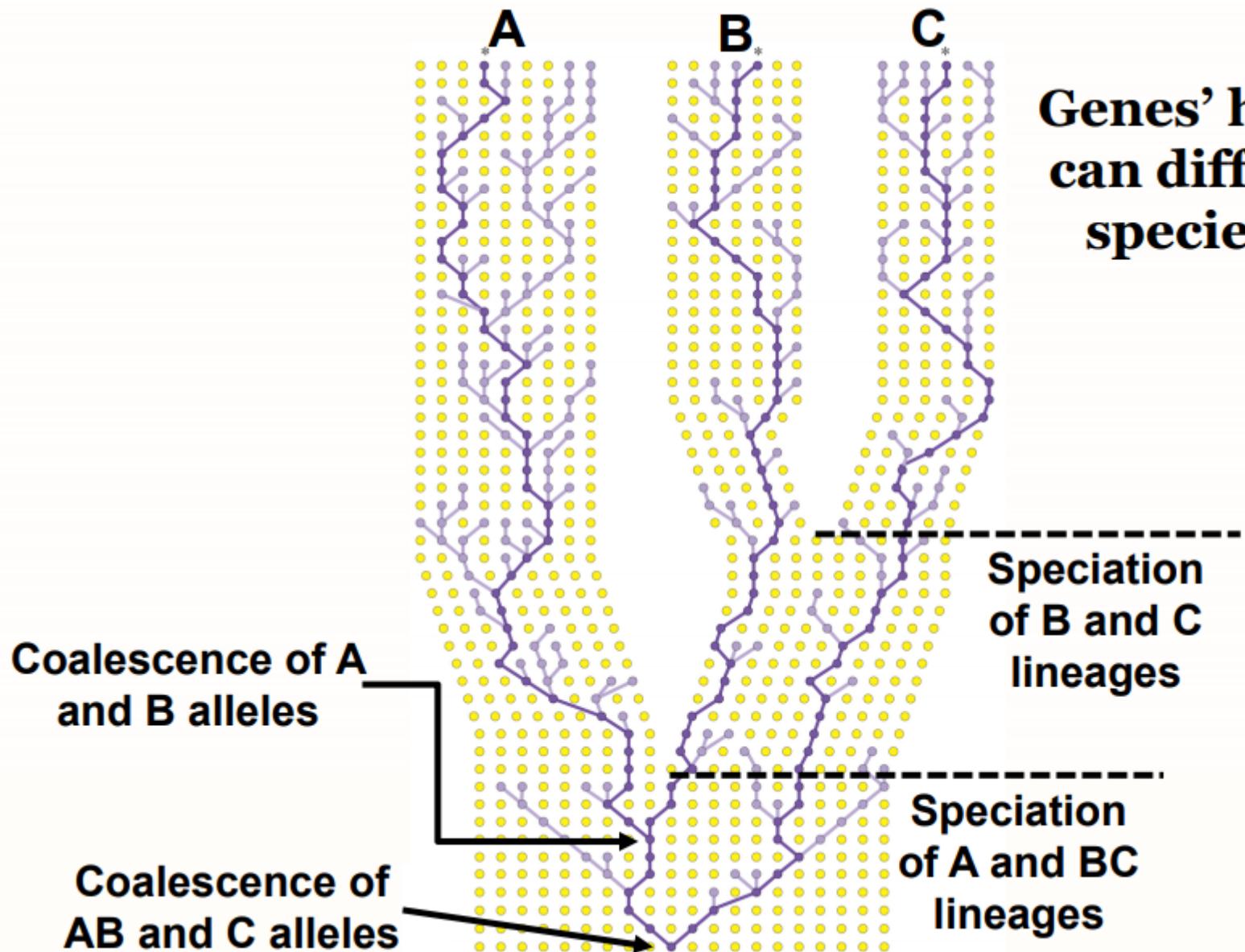
### Hybridization



### Recombination



## Incomplete lineage sorting



**Genes' histories  
can differ from  
species ones**

Incomplete lineage sorting in the primate lineage:



**Informative Sites**

**8,561 / 11,293**

(~76%)



**1,302 / 11,293**

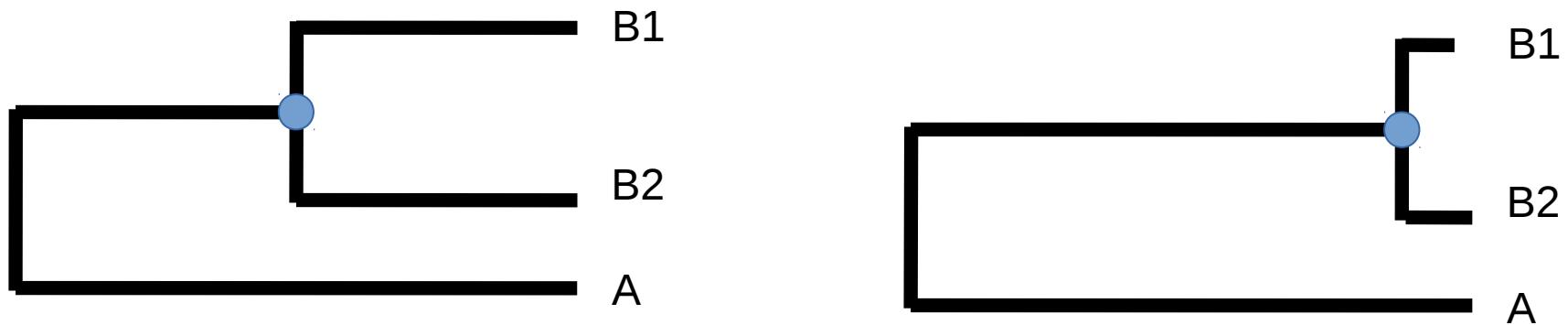
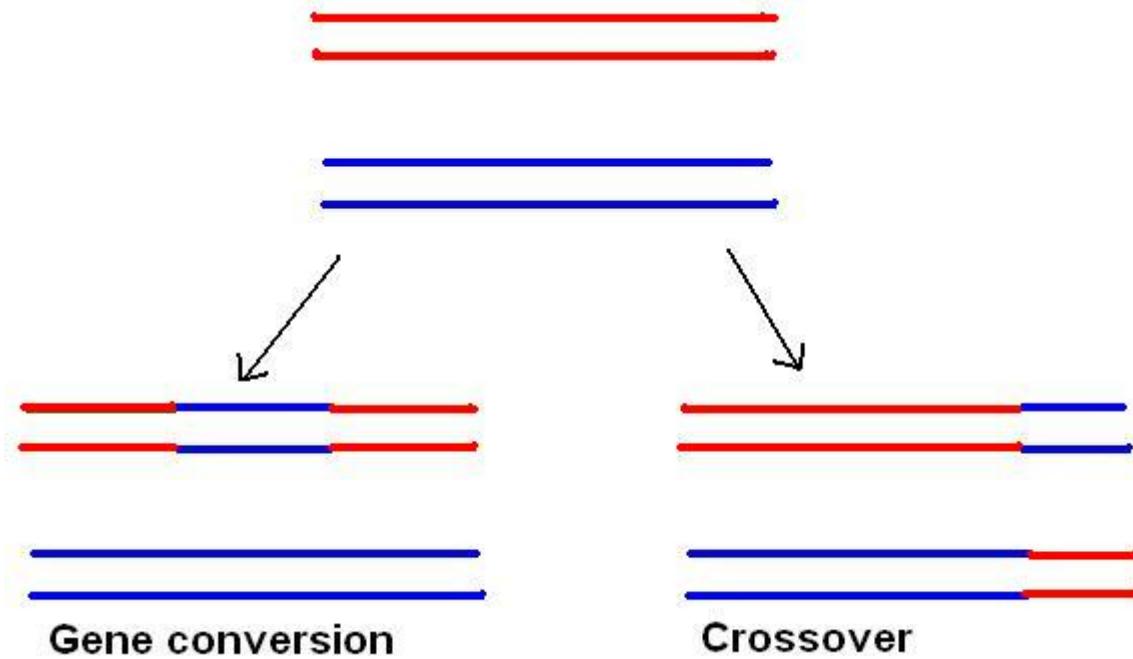
(~11.5%)



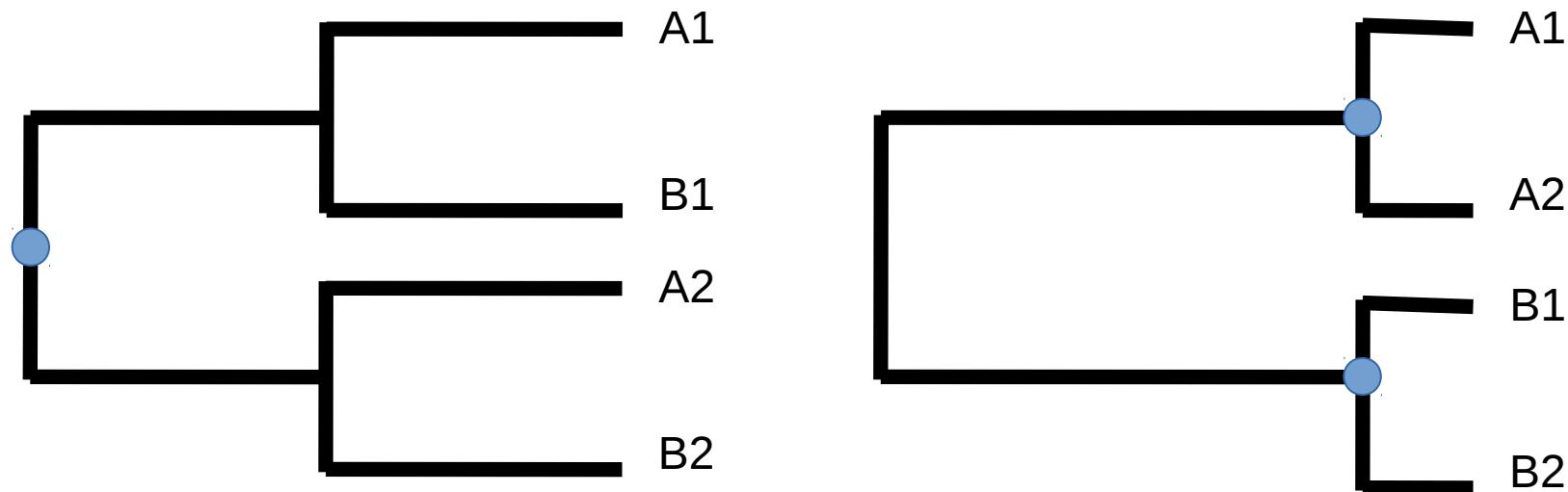
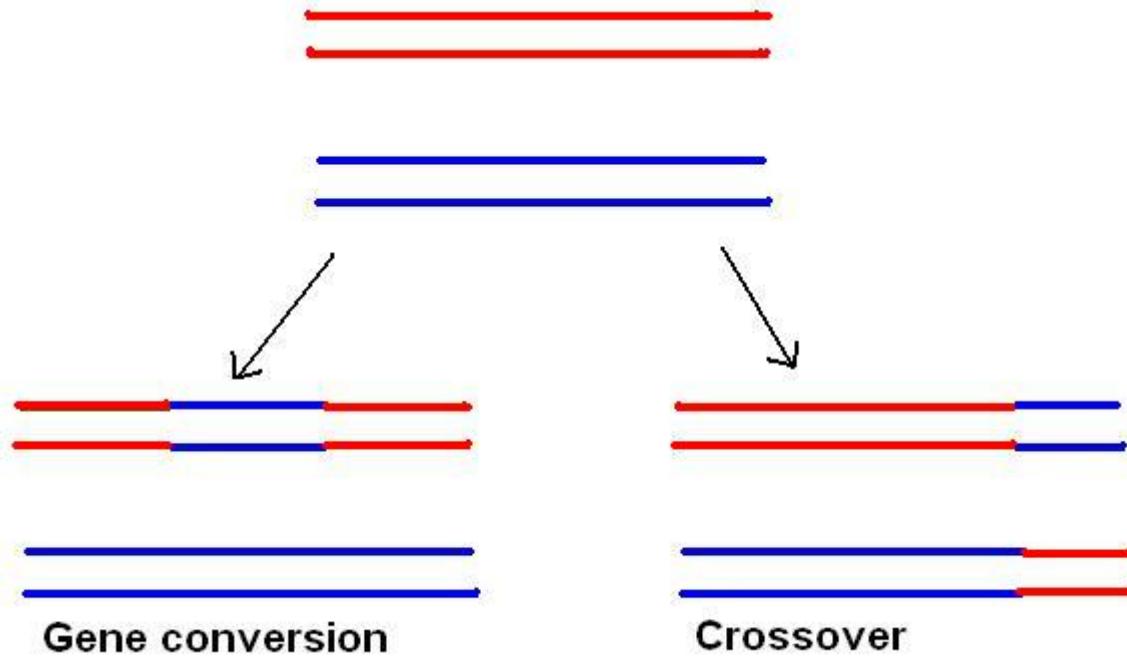
**1,430 / 11,293**

(~12.5%)

## Gene conversion



## Gene conversion

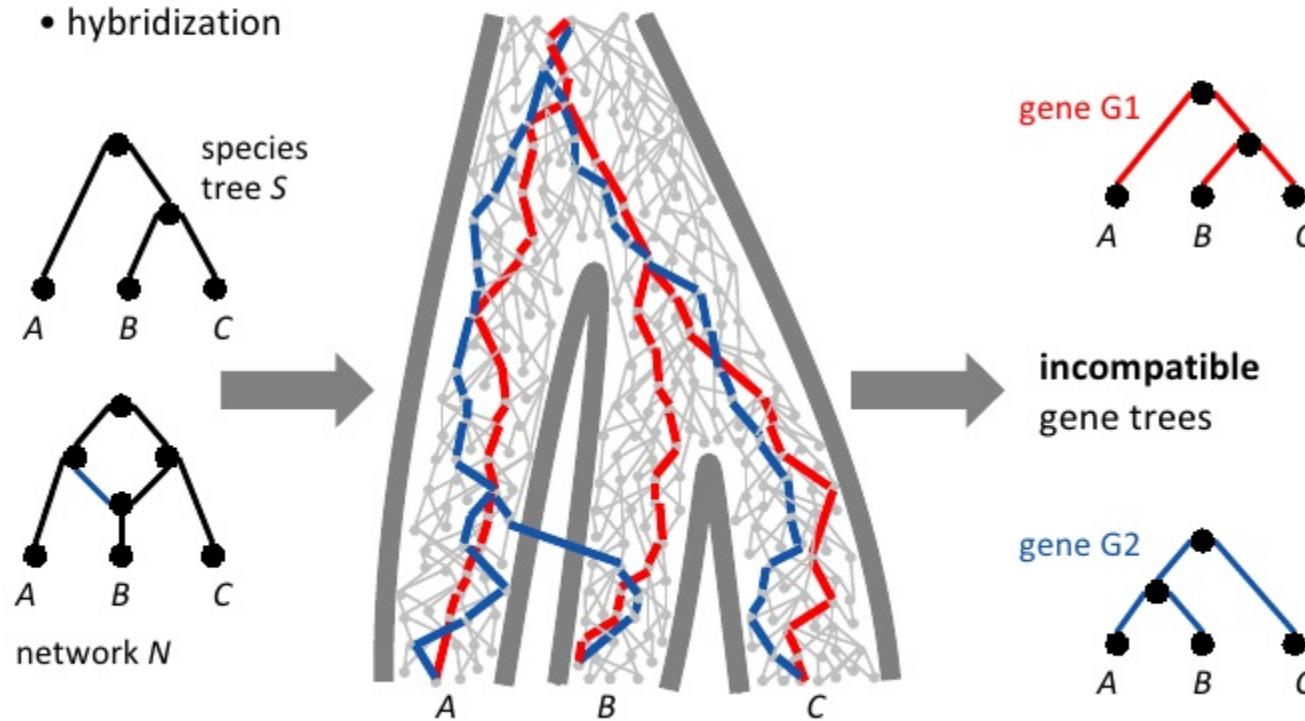


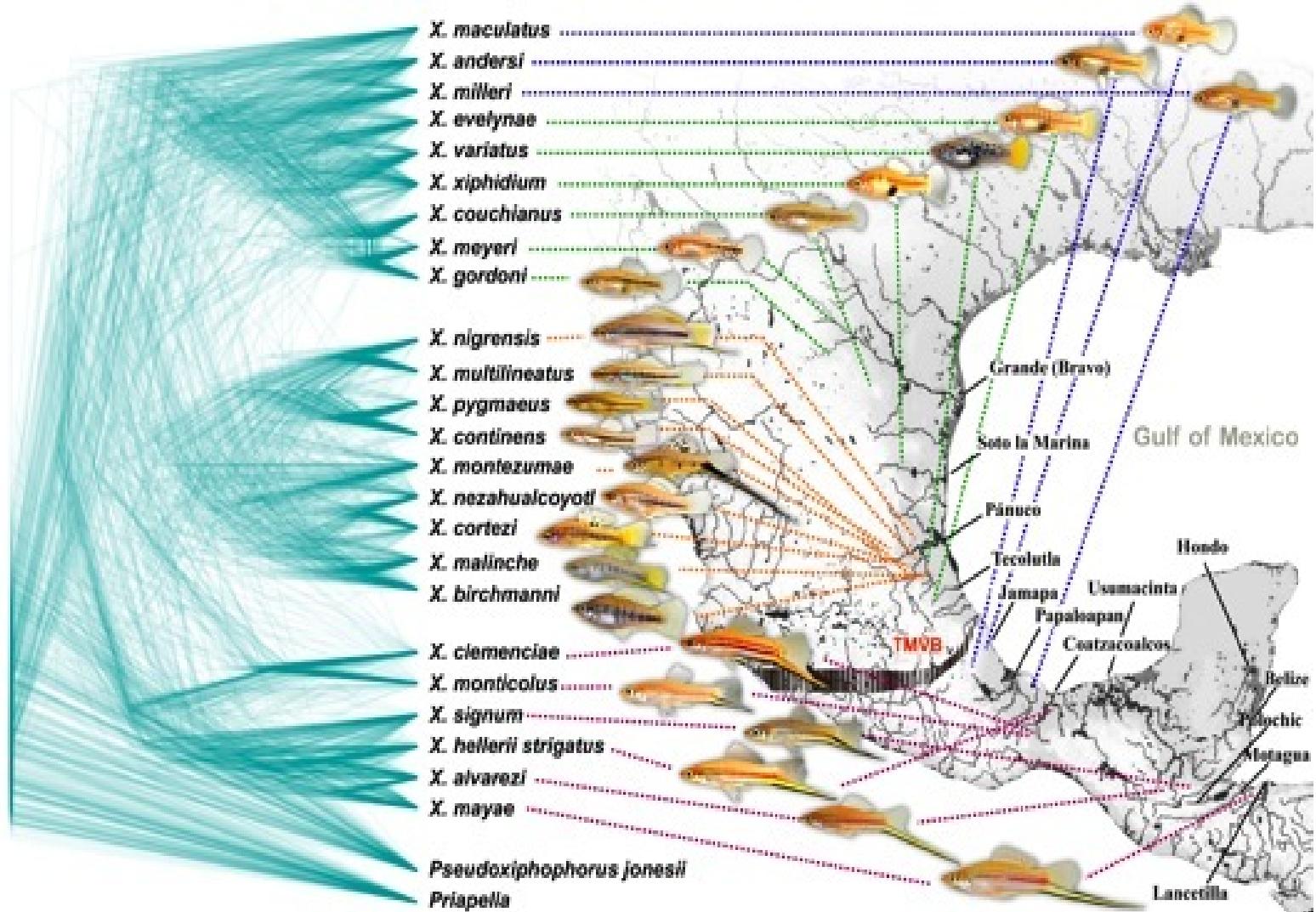
# Introgression / Hybridization

## Genetic material transfer

Genetic material transfers between coexisting species:

- horizontal gene transfer
- hybridization





# Why predicting orthology is important?

- **Important implications for phylogeny:** only sets of orthologous genes are expected to reflect the underlying species evolution (although there are many exceptions)

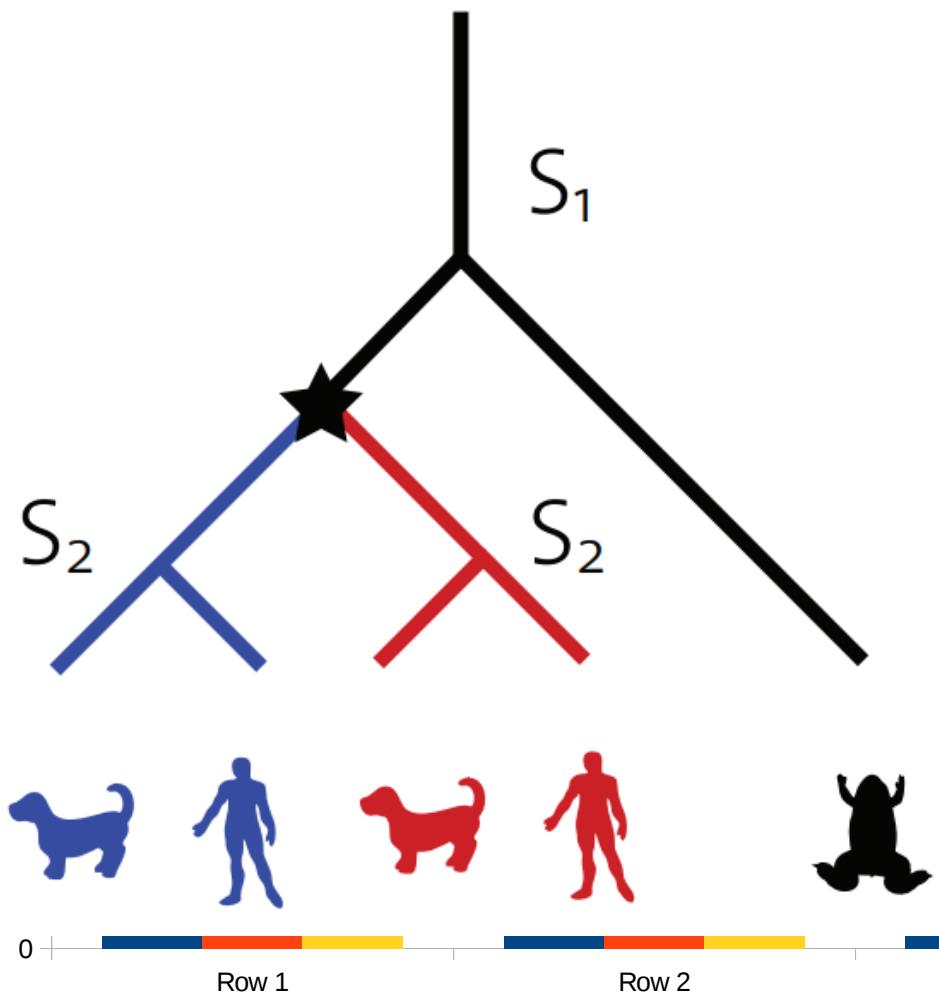
Yes, orthologs are useful to retrieve species trees.

However, even if you have the orthologs, getting the species tree is not straightforward.

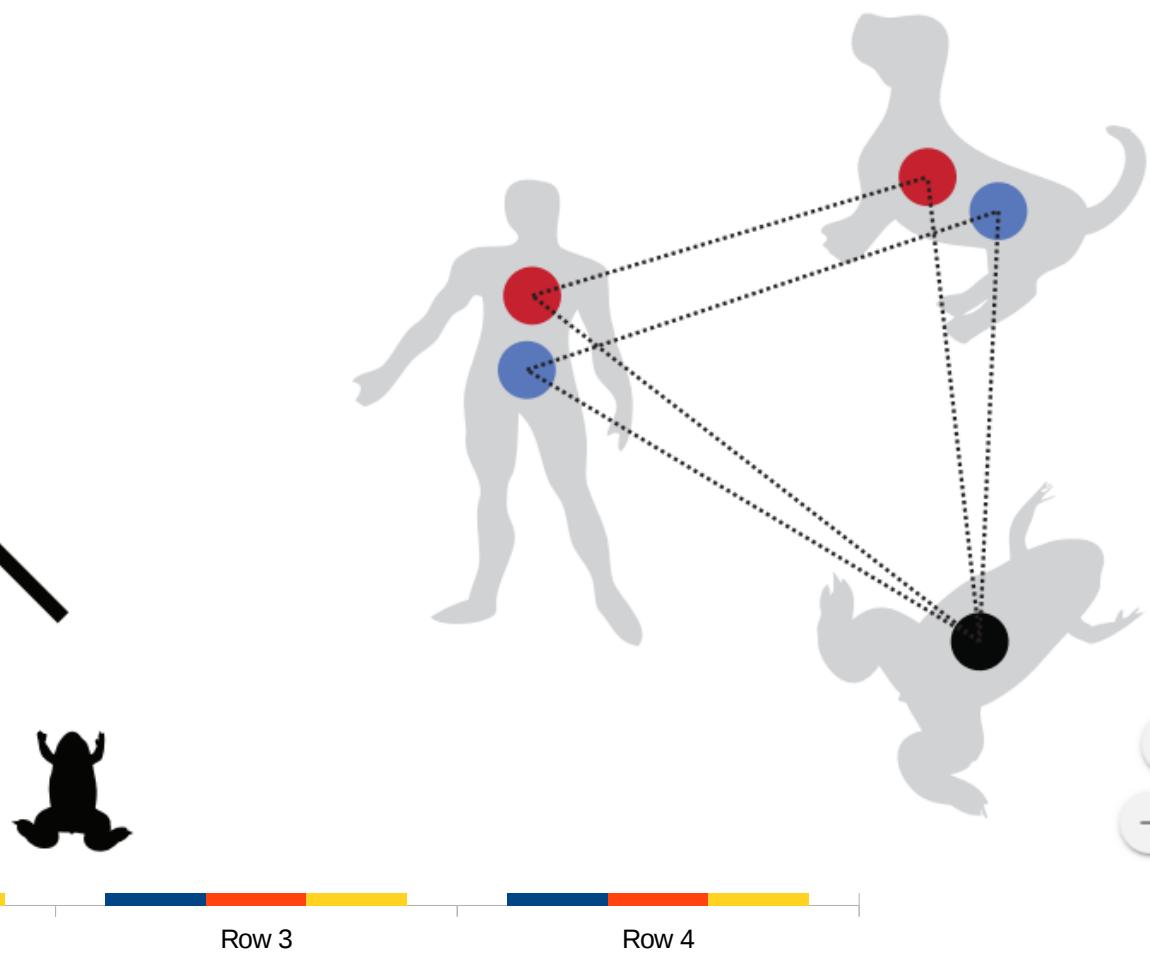
# Why predicting orthology is important?

- **Important implications for phylogeny:** only sets of orthologous genes are expected to reflect the underlying species evolution (although there are many exceptions)
- The most exact way of **comparing two (or more) genomes** in terms of their gene content. Necessary to uncover how genomes evolve.
- Implications for **functional inference:** orthologs, as compared to paralogs, are more likely to share the same function

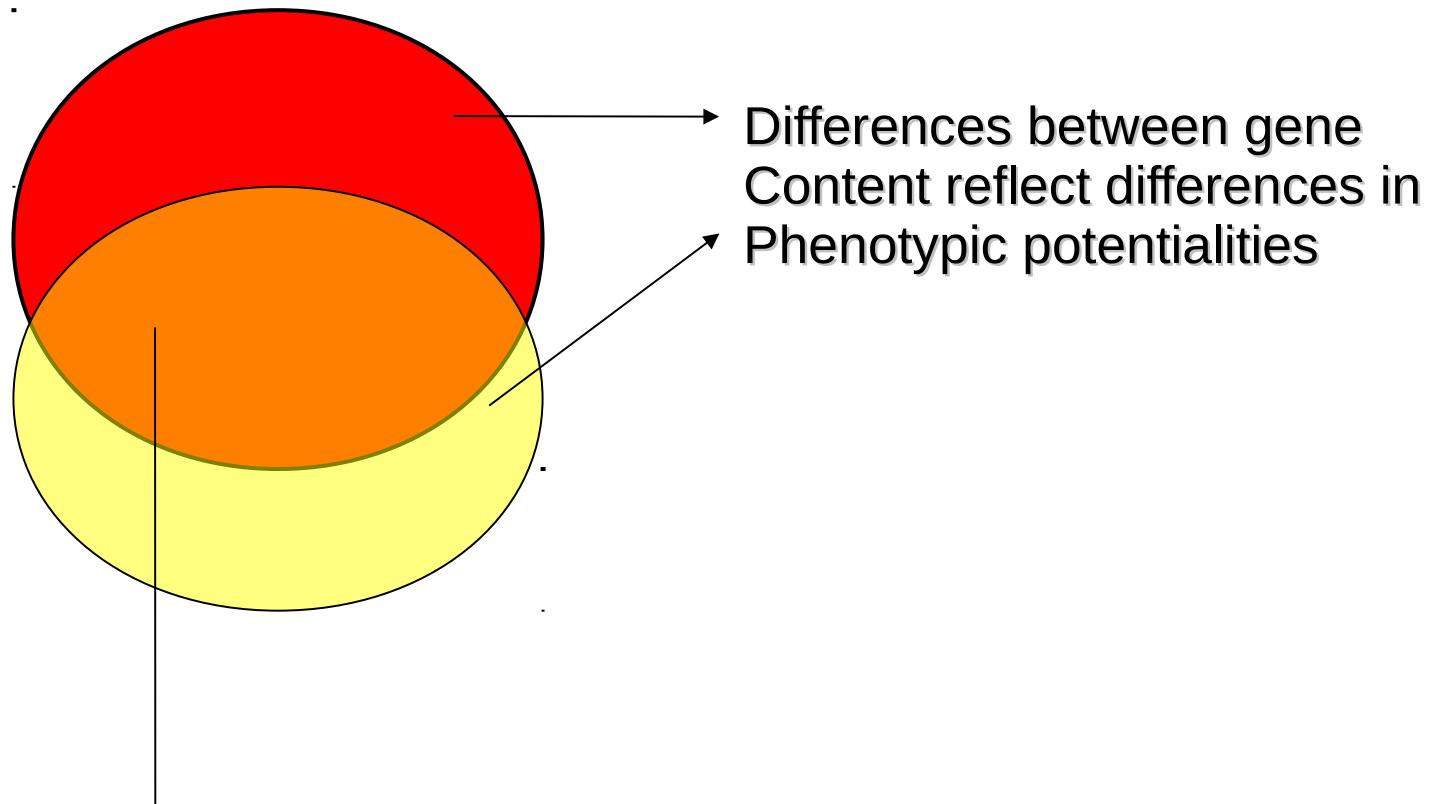
a)

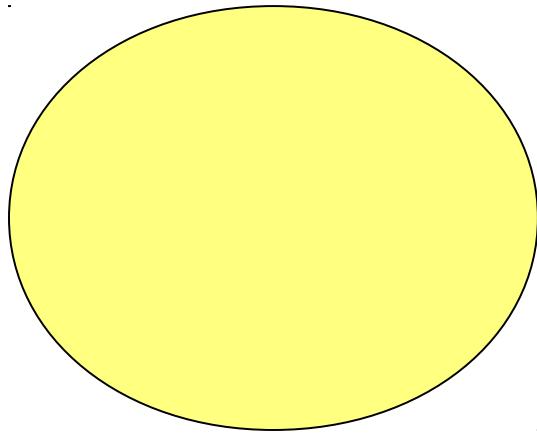


b)

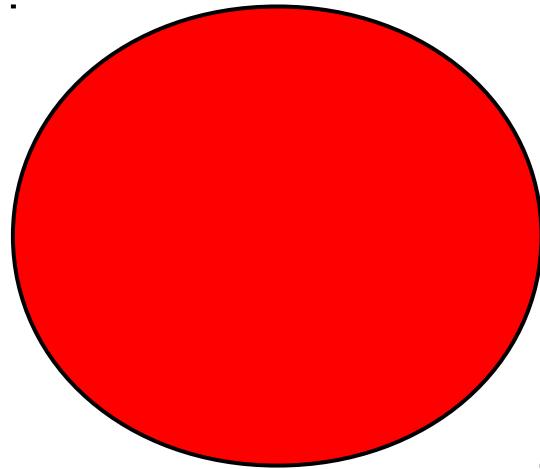


Gene content → co-evolution. (The easy case, few genomes. )



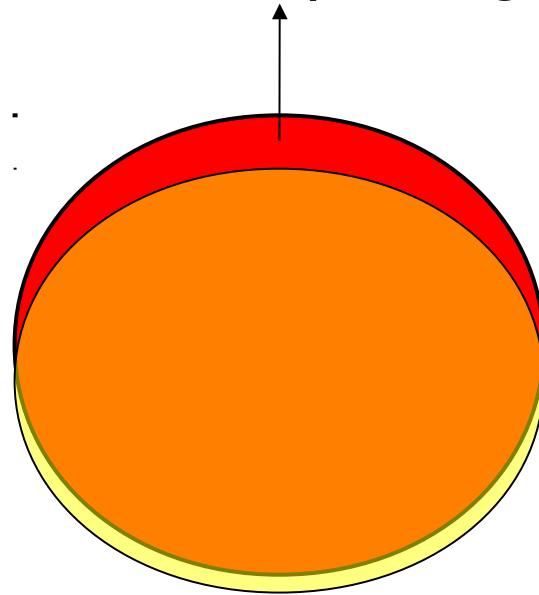


*L. innocua* (non-pathogen)



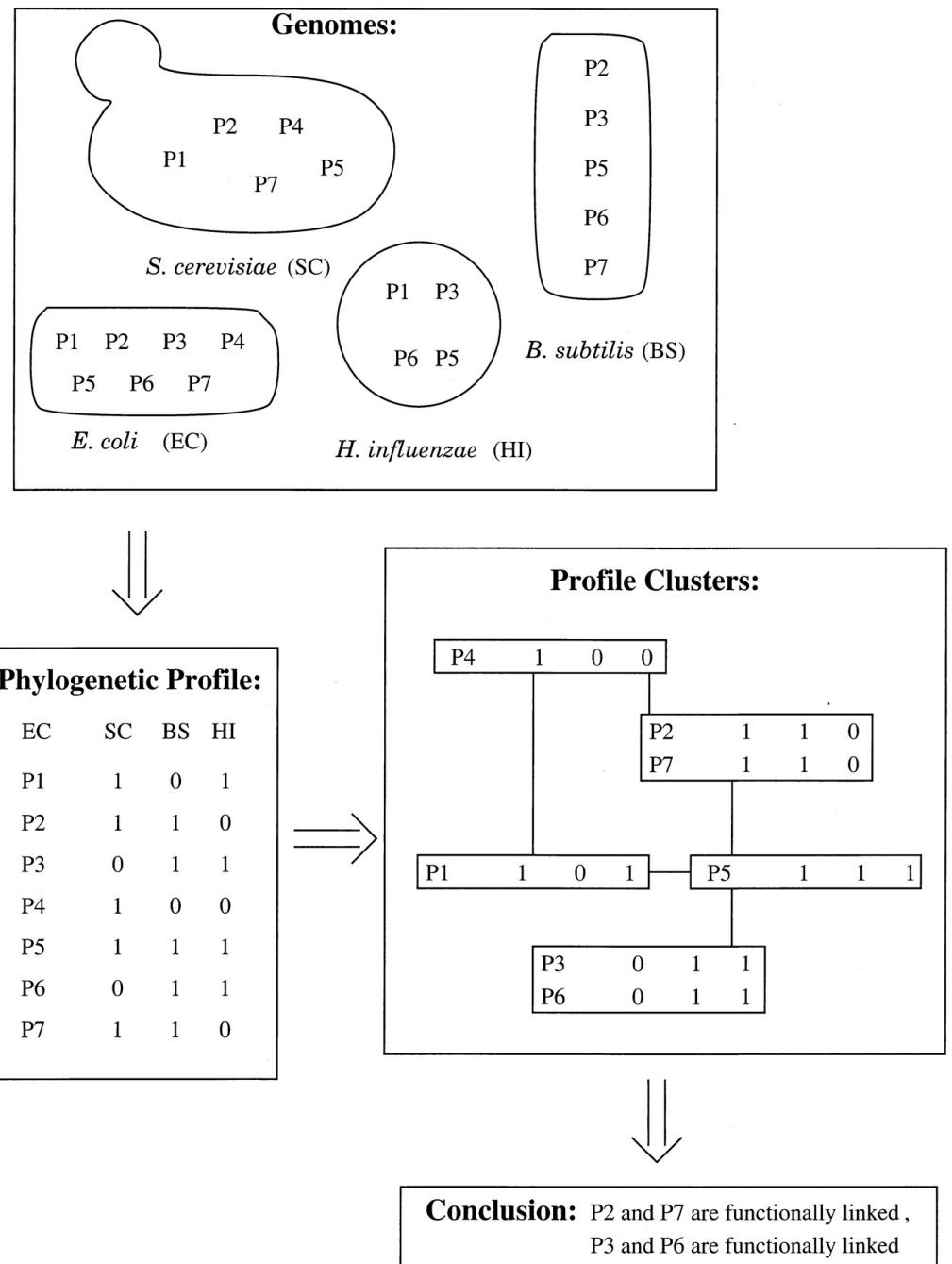
*L. monocytogenes* (pathogen)

Genes involved in pathogenicity



L. innocua (non-pathogenic) monocytogenes (pathogenic)

## More than two genomes



## Assigning protein functions by comparative genome analysis: Protein phylogenetic profiles



Matteo Pellegrini, Edward M. Marcotte, Michael J. Thompson, David Eisenberg, and Todd O. Yeates

PNAS April 13, 1999; 96 (8) 4285-4288; <https://doi.org/10.1073/pnas.96.8.4285>

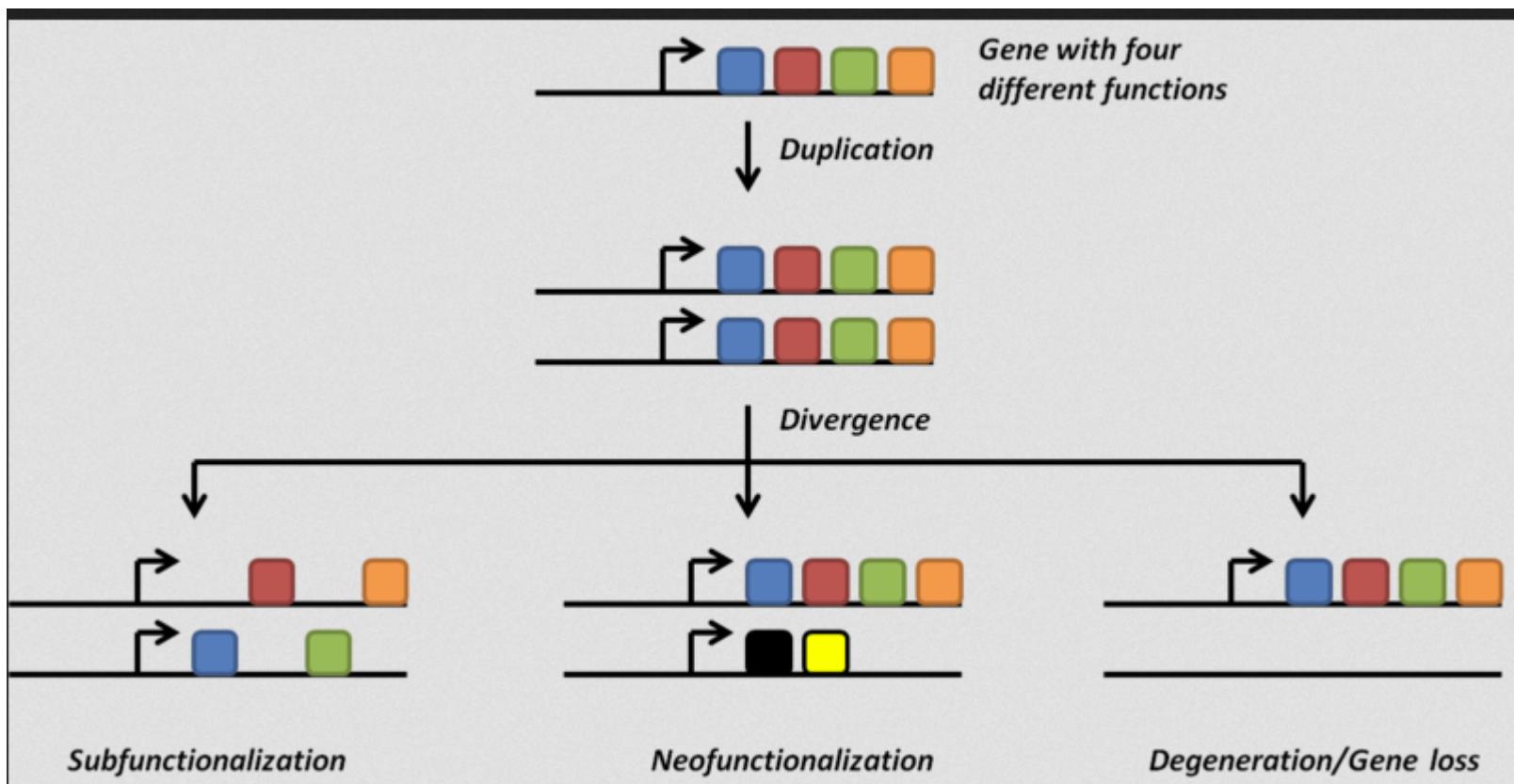
Contributed by David S. Eisenberg

# Why predicting orthology is important?

- **Important implications for phylogeny:** only sets of orthologous genes are expected to reflect the underlying species evolution (although there are many exceptions)
- The most exact way of **comparing two (or more) genomes** in terms of their gene content. Necessary to uncover how genomes evolve.
- Implications for **functional inference:** orthologs, as compared to paralogs, are more likely to share the same function

REALLY???, IS THIS TRUE IF SO, WHY IS THAT?

## After duplication: diversify or die (neofunctionalization or subfunctionalization models)



# How confident can we be that orthologs are similar, but paralogs differ?

Romain A. Studer and Marc Robinson-Rechavi

Department of Ecology and Evolution, Biophore, Lausanne University, CH-1015 Lausanne, Switzerland and Swiss Institute of Bioinformatics, CH-1015 Lausanne, Switzerland

OPEN  ACCESS Freely available online

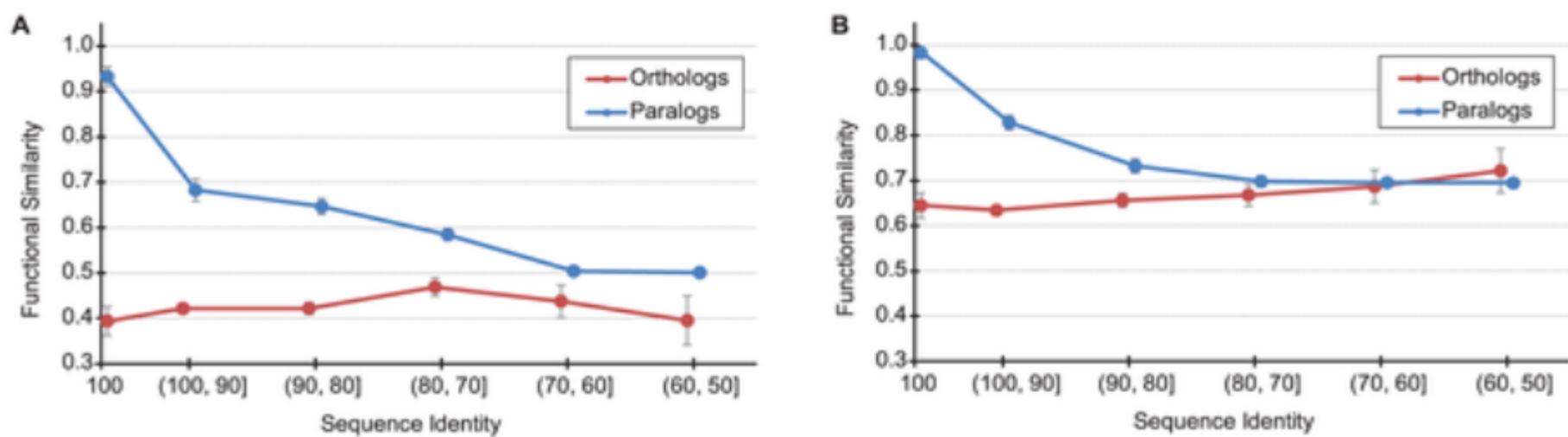
PLOS COMPUTATIONAL BIOLOGY

## Testing the Ortholog Conjecture with Comparative Functional Genomic Data from Mammals

Nathan L. Nehrt<sup>1,3</sup>, Wyatt T. Clark<sup>1,3</sup>, Predrag Radivojac<sup>1,\*</sup>, Matthew W. Hahn<sup>1,2\*</sup>

**1** School of Informatics and Computing, Indiana University, Bloomington, Indiana, United States of America, **2** Department of Biology, Indiana University, Bloomington, Indiana, United States of America

**Figure 1. The relationship between functional similarity and sequence identity for human-mouse orthologs (red) and all paralogs (blue).**



Nehrt NL, Clark WT, Radivojac P, Hahn MW (2011) Testing the Ortholog Conjecture with Comparative Functional Genomic Data from Mammals. PLoS Comput Biol 7(6): e1002073. doi:10.1371/journal.pcbi.1002073  
<http://www.ploscompbiol.org/article/info:doi/10.1371/journal.pcbi.1002073>

# On the Use of Gene Ontology Annotations to Assess Functional Similarity among Orthologs and Paralogs: A Short Report

Paul D. Thomas<sup>1\*</sup>, Valerie Wood<sup>2</sup>, Christopher J. Mungall<sup>3</sup>, Suzanna E. Lewis<sup>3</sup>, Judith A. Blake<sup>4</sup> on behalf of the Gene Ontology Consortium

**1** Division of Bioinformatics, Department of Preventive Medicine, University of Southern California, Los Angeles, California, United States of America, **2** Cambridge Systems Biology Centre and Department of Biochemistry, University of Cambridge, Cambridge, United Kingdom, **3** Genomics Division, Lawrence Berkeley National Laboratory, Berkeley, California, United States of America, **4** Bioinformatics and Computational Biology, The Jackson Laboratory, Bar Harbor, Maine, United States of America

OPEN  ACCESS Freely available online

## Resolving the Ortholog Conjecture: Orthologs Tend to Be Weakly, but Significantly, More Similar in Function than Paralogs

Adrian M. Altenhoff<sup>1,2</sup>, Romain A. Studer<sup>2,3,4</sup>, Marc Robinson-Rechavi<sup>2,3</sup>, Christophe Dessimoz<sup>1,2,5\*</sup>

**1** ETH Zurich, Department of Computer Science, Zürich, Switzerland, **2** Swiss Institute of Bioinformatics, Lausanne, Switzerland, **3** Department of Ecology and Evolution, University of Lausanne, Lausanne, Switzerland, **4** Institute of Structural and Molecular Biology, Division of Biosciences, University College London, London, United Kingdom, **5** FMRI -European Bioinformatics Institute, Hinxton, Cambridgeshire, United Kingdom

*Nature Reviews Genetics* | AOP, published online 4 April 2013; doi:10.1038/nrg3456

## PERSPECTIVES

BRIEFINGS IN BIOINFORMATICS, VOL 12, NO 5, 442–448  
Advance Access published on 22 April 2011

doi:10.1093/bib/bbr022

### OPINION

#### Functional and evolutionary implications of gene orthology

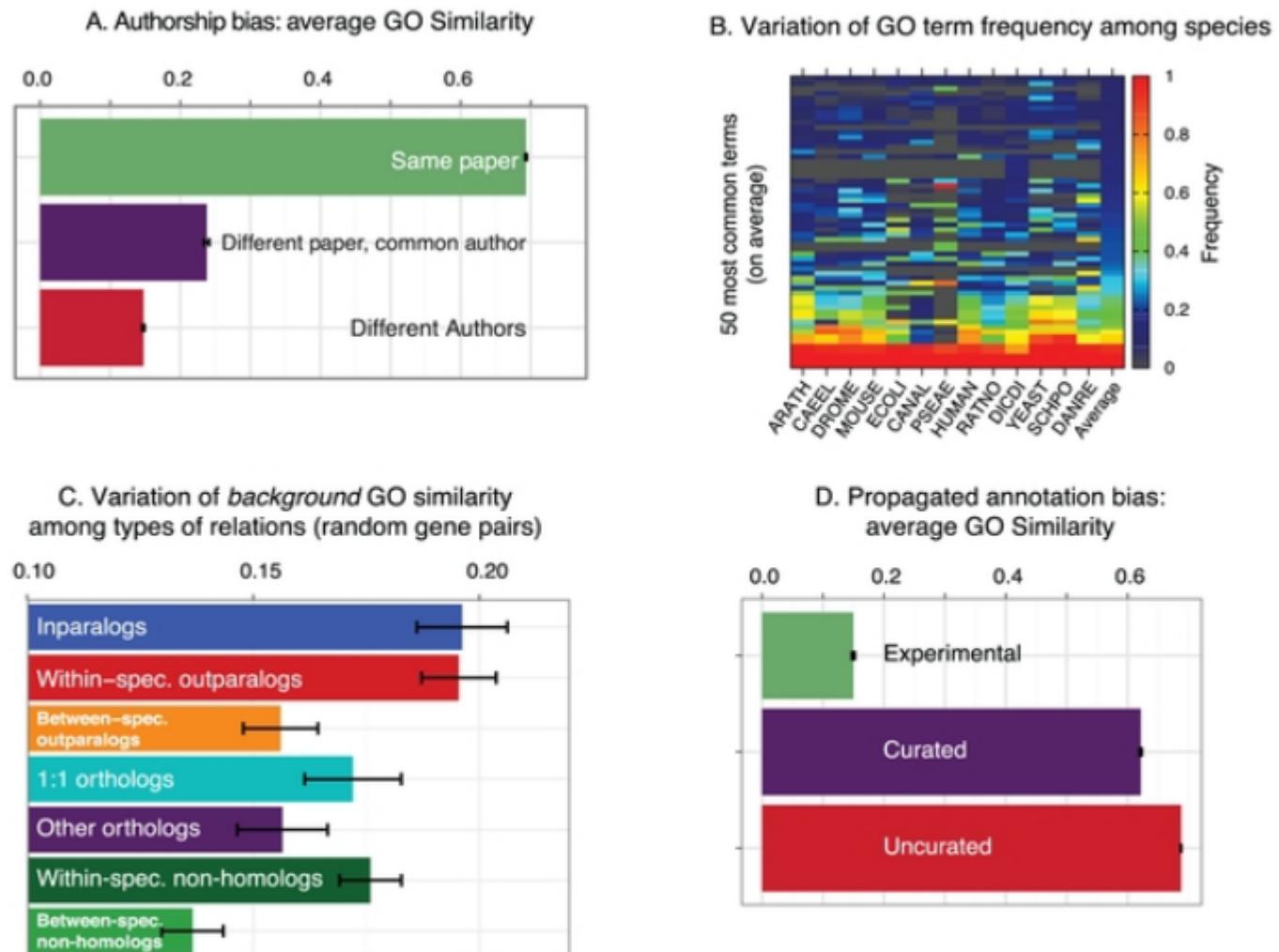
Toni Gabaldón and Eugene V. Koonin

#### Evidence for short-time divergence and long-time conservation of tissue-specific expression after gene duplication

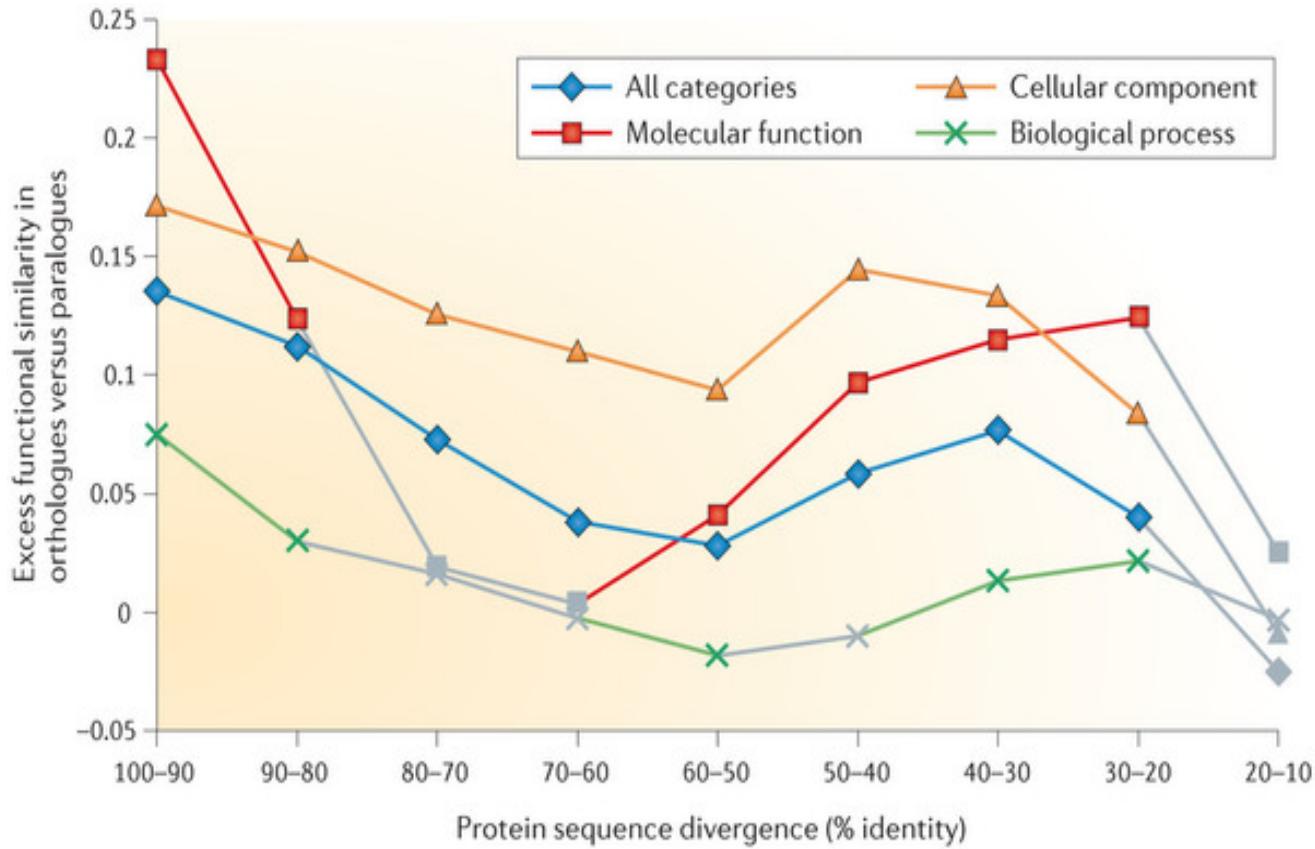
Jaime Huerta-Cepas, Joaquín Dopazo, Martijn A. Huynen and Toni Gabaldón

Submitted: 19th January 2011; Received (in revised form): 22nd March 2011

**Figure 1. Potential confounding factors in GO analyses.**



Altenhoff AM, Studer RA, Robinson-Rechavi M, Dessimoz C (2012) Resolving the Ortholog Conjecture: Orthologs Tend to Be Weakly, but Significantly, More Similar in Function than Paralogs. *PLoS Comput Biol* 8(5): e1002514.  
doi:10.1371/journal.pcbi.1002514  
<http://www.ploscompbiol.org/article/info:doi/10.1371/journal.pcbi.1002514>



Nature Reviews | Genetics

**Questions about this part?**



Barcelona  
Biomedical  
Research  
Park



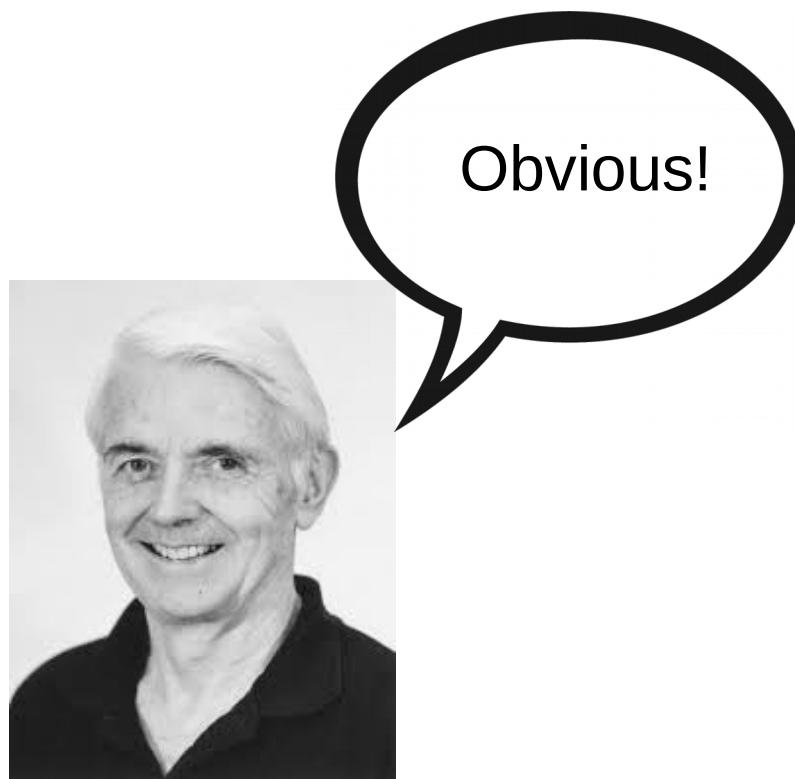
# Orthology Part II

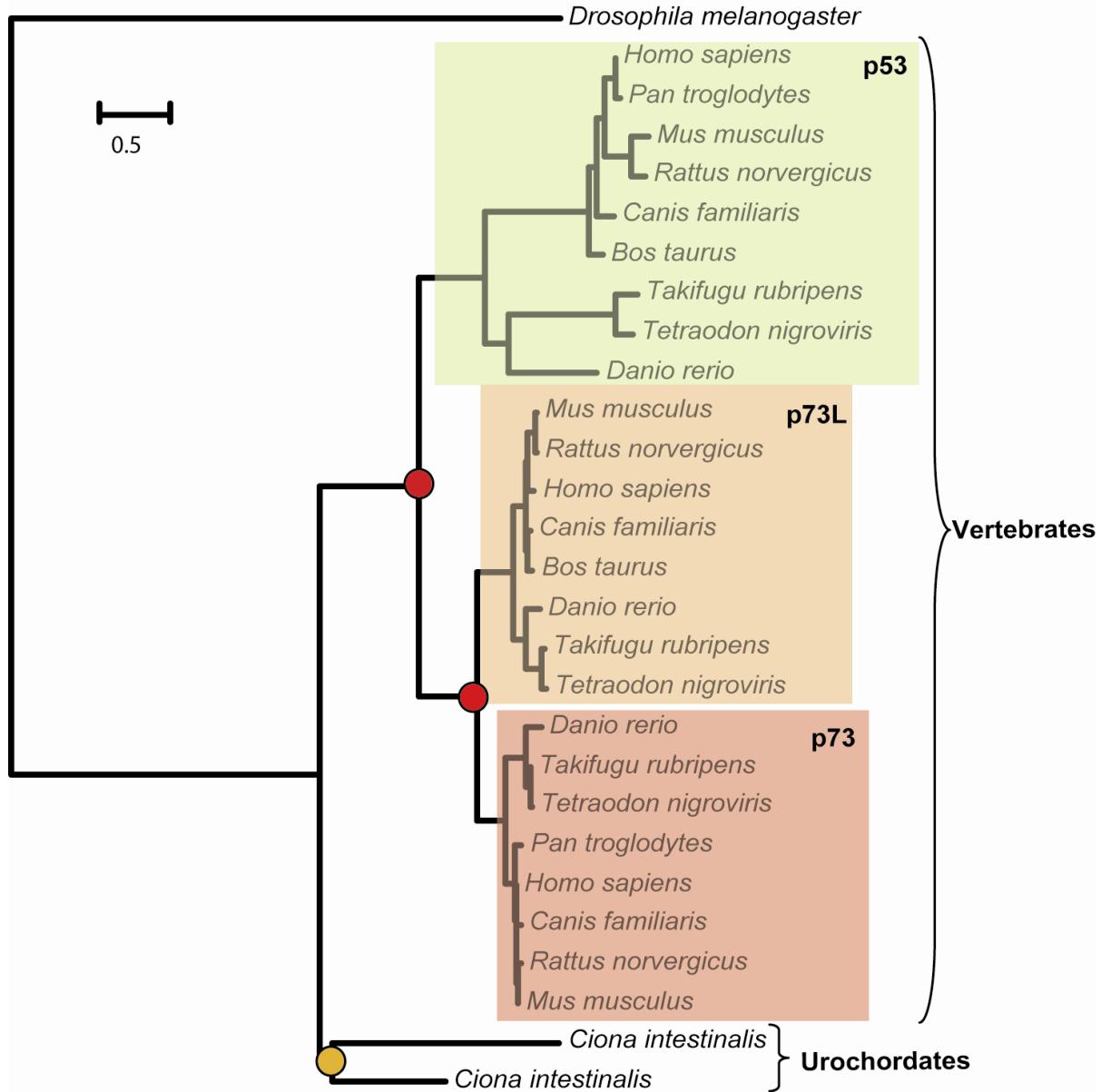
## Orthology prediction methods

Toni Gabaldón  
Centre for Genomic Regulation (CRG), Barcelona

## **Classical approach: phylogenetic inference**

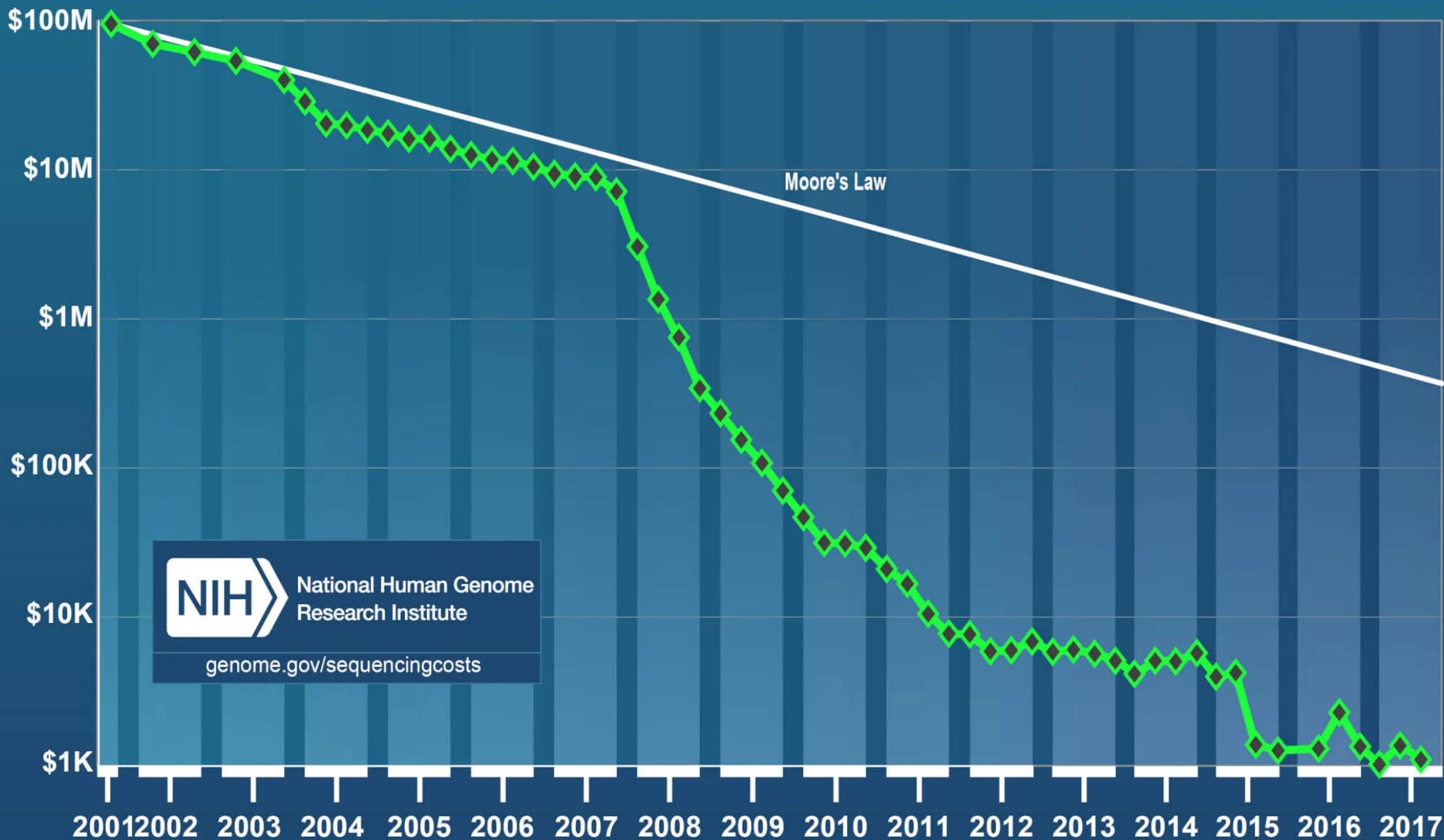
- Build a gene tree
- Compare to the species tree
- Infer duplications and speciation events
- Assign orthology and paralogy relationships accordingly

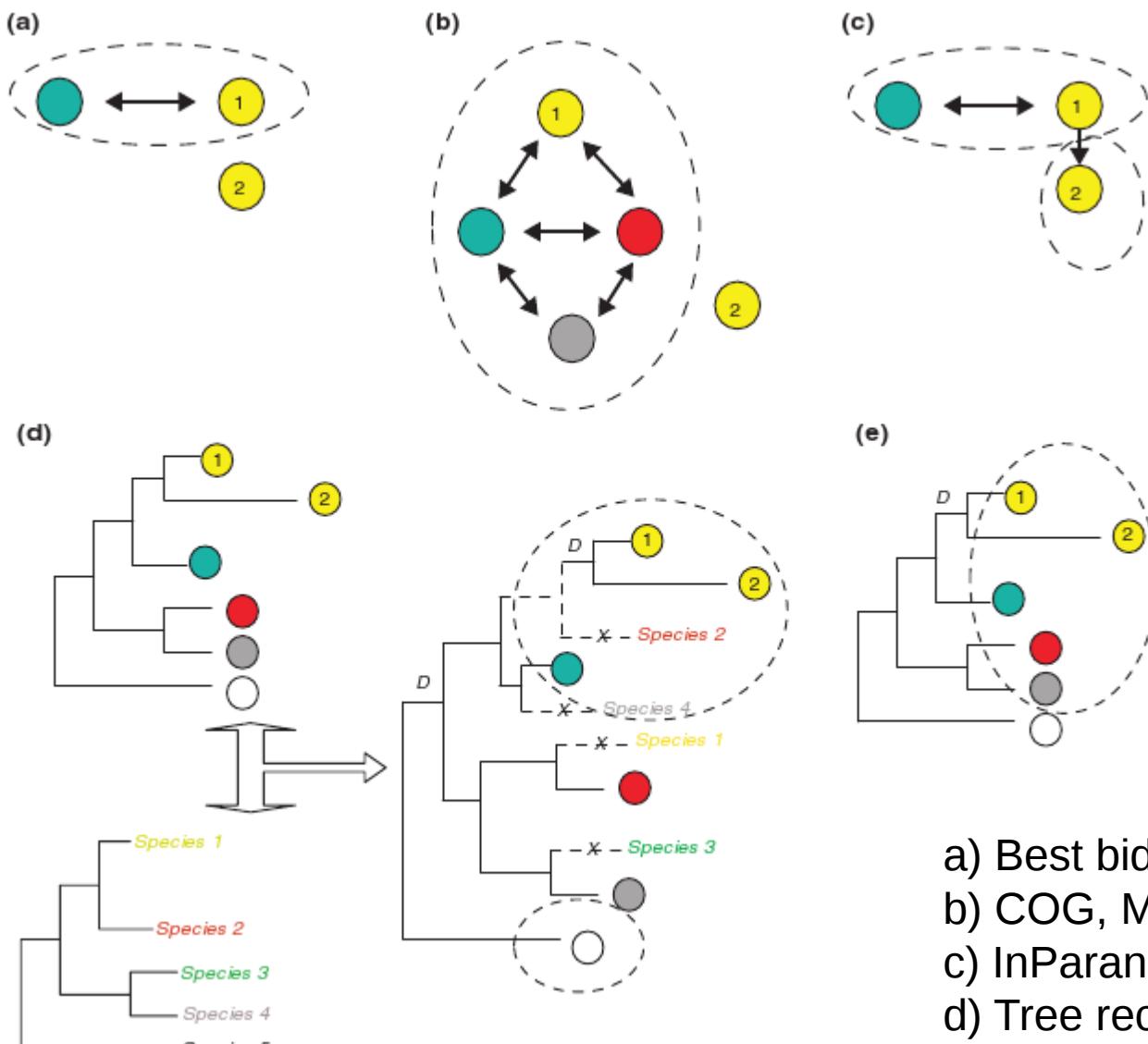




**Going genome-wide scale:**  
Everything must be done automatic and “blind”

# *Cost per Genome*





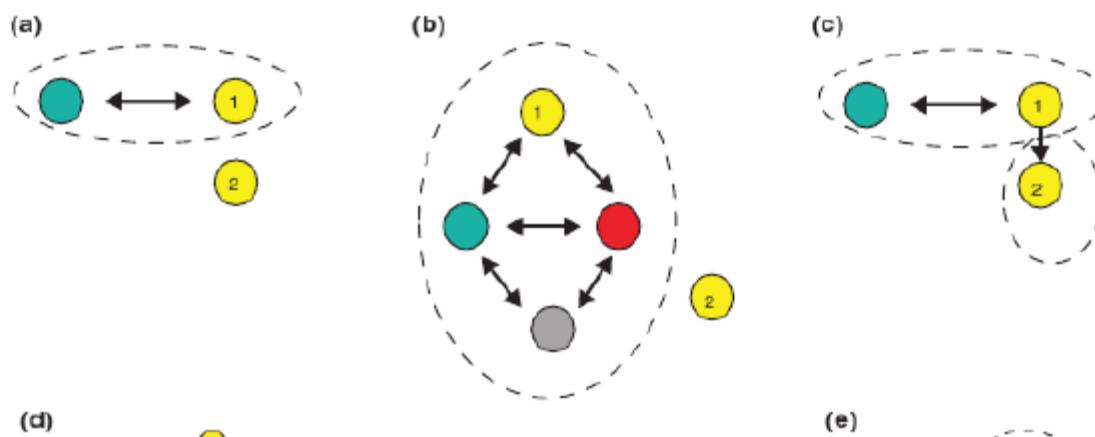
- Best bidirectional hits
- COG, MCL-clustering approach
- InParanoid
- Tree reconciliation
- Species-overlap (PhylomeDB)

## Similarity-based approaches (many more approaches):

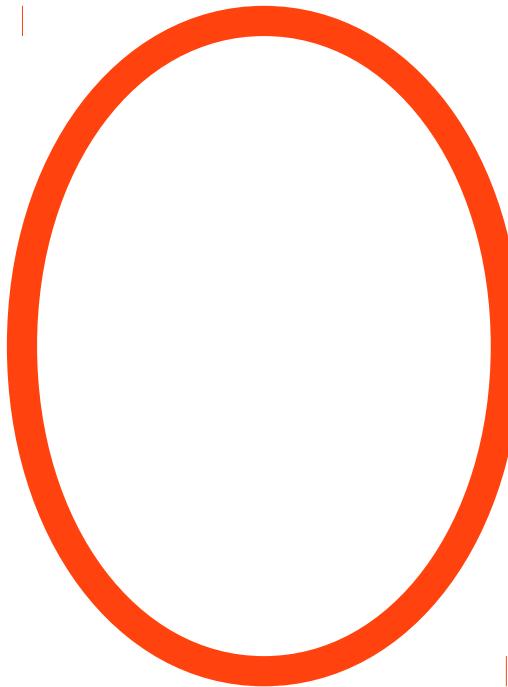
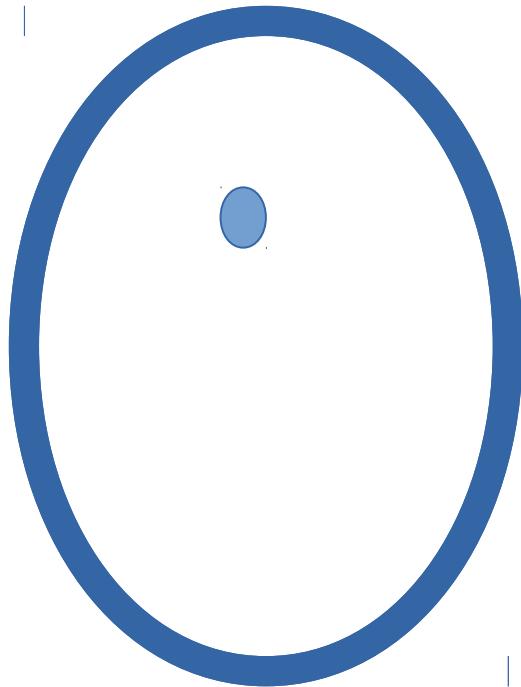
### - Best Reciprocal Hits

Detects all orthologies as one-to one. Highly affected by paralogy. Low rate of false positives but high rates of false negatives.

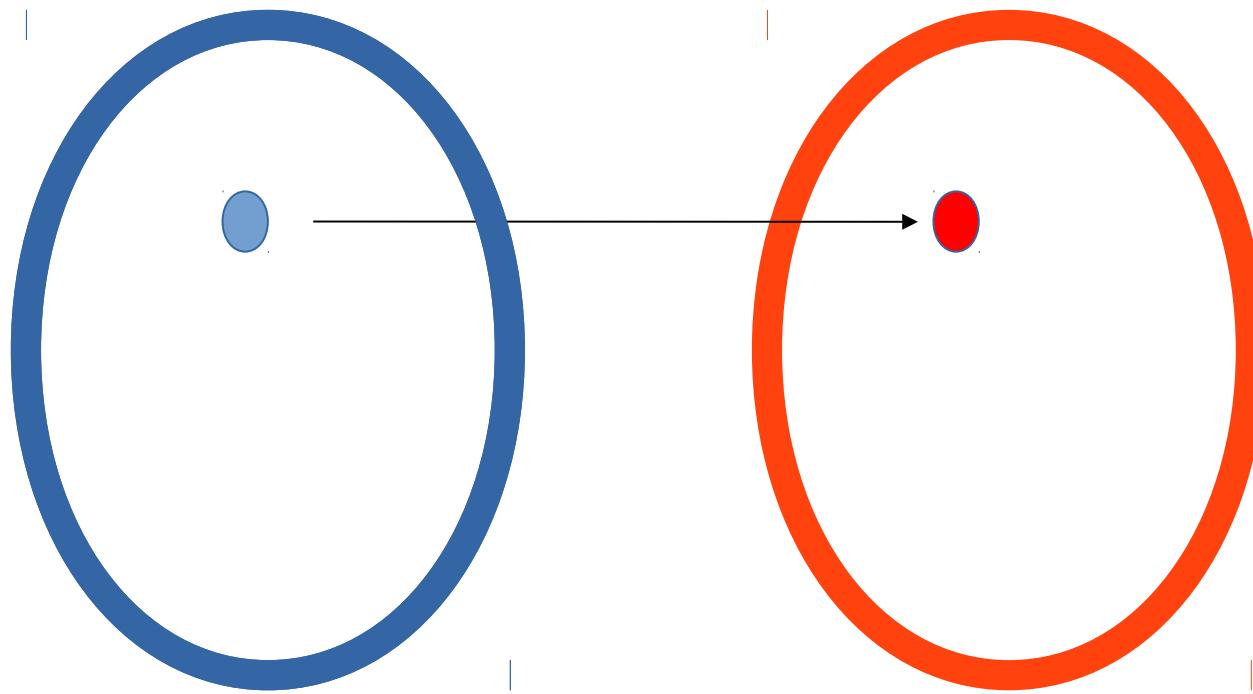
The simplest and fastest method, still widely used



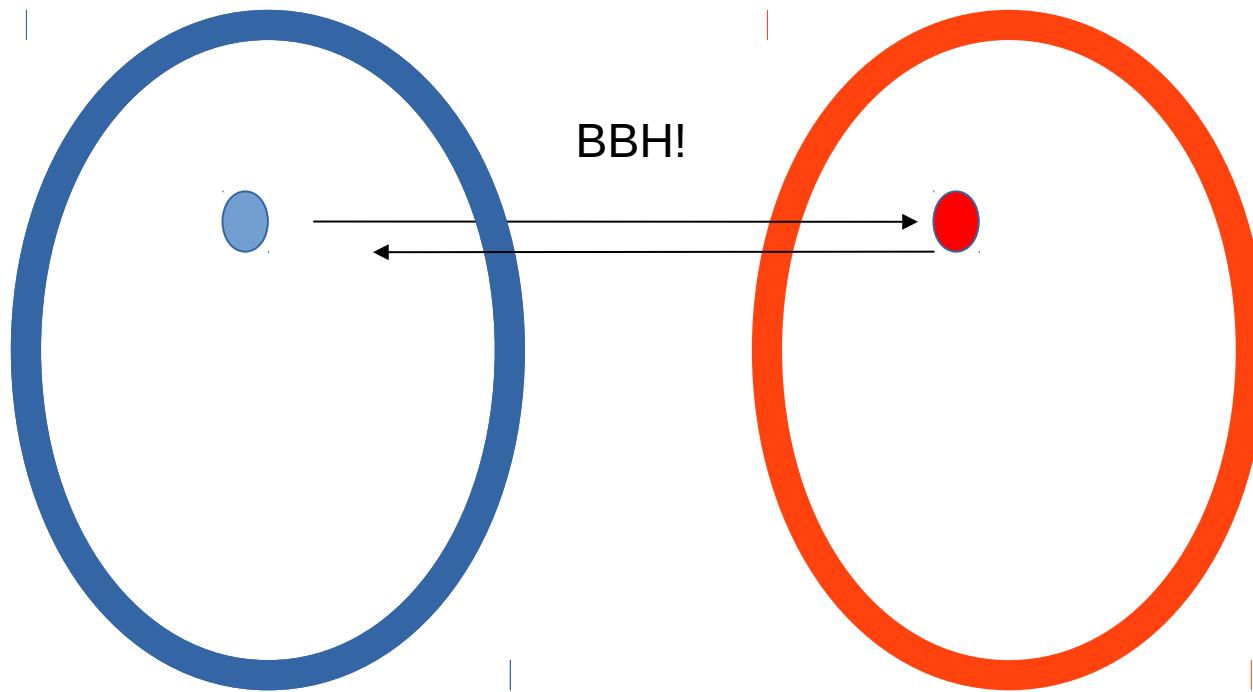
Best bidirectional hit (BBH), Best reciprocal hits (BRH)

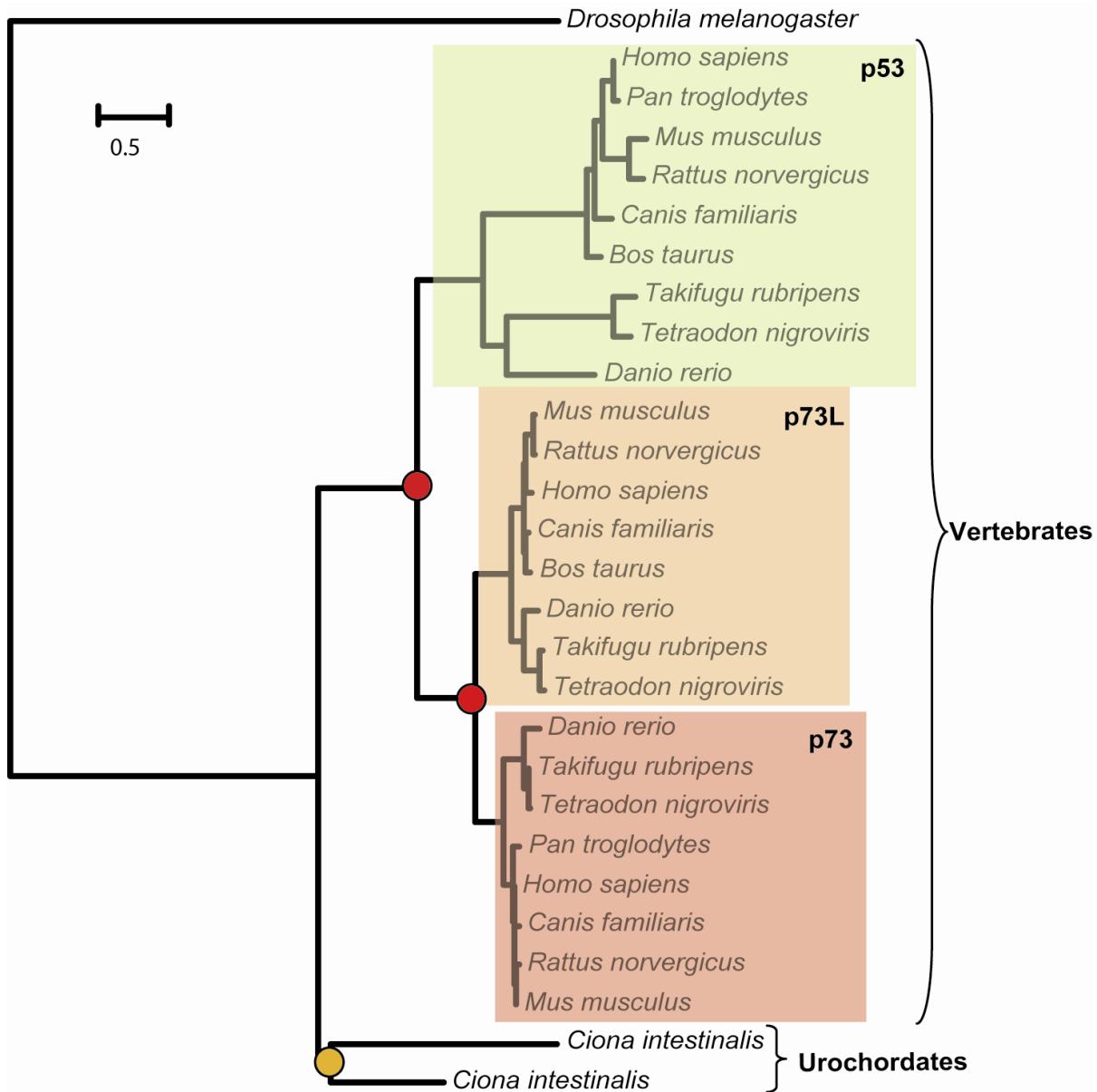


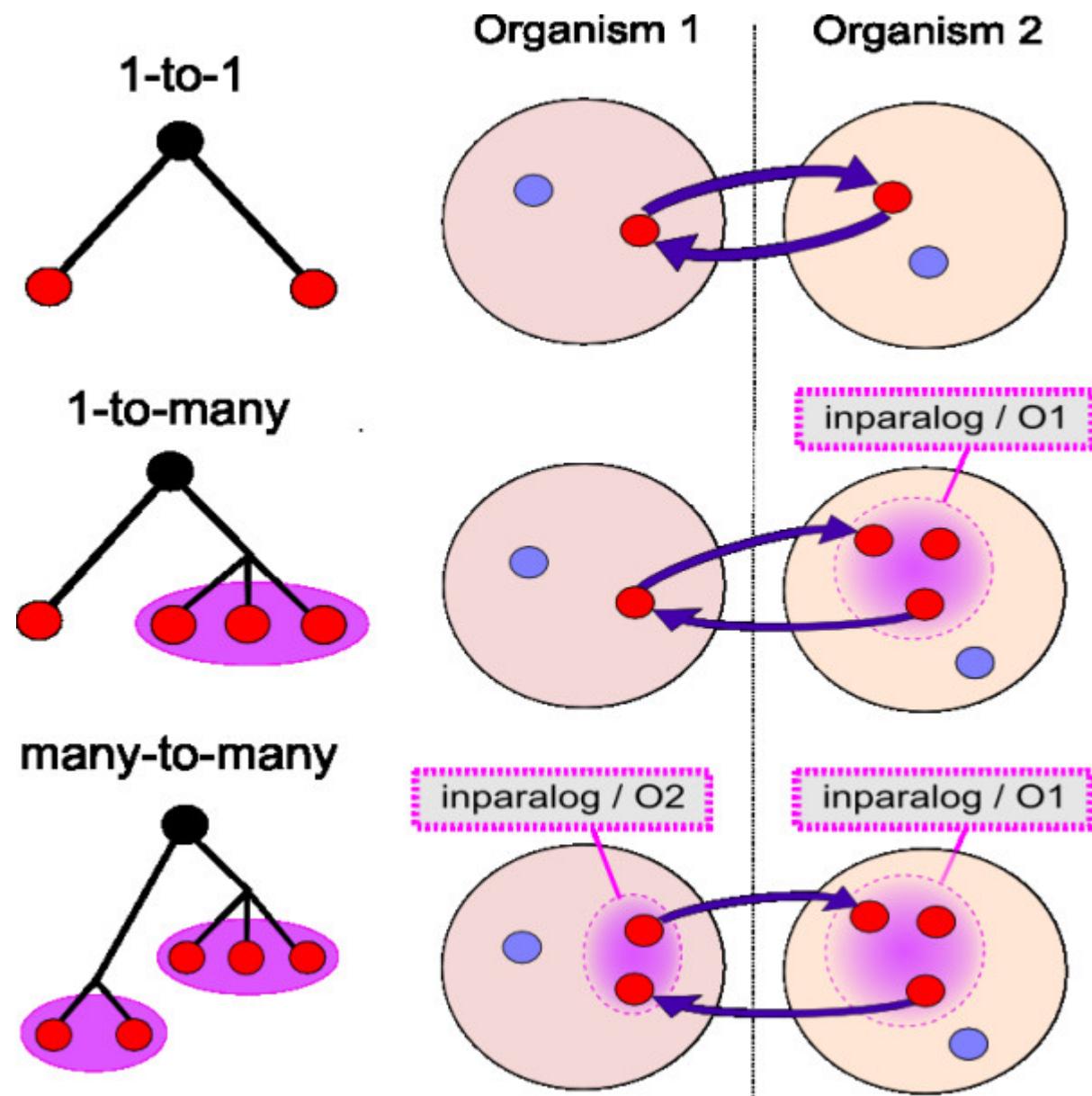
Best bidirectional hit (BBH), Best reciprocal hits (BRH)

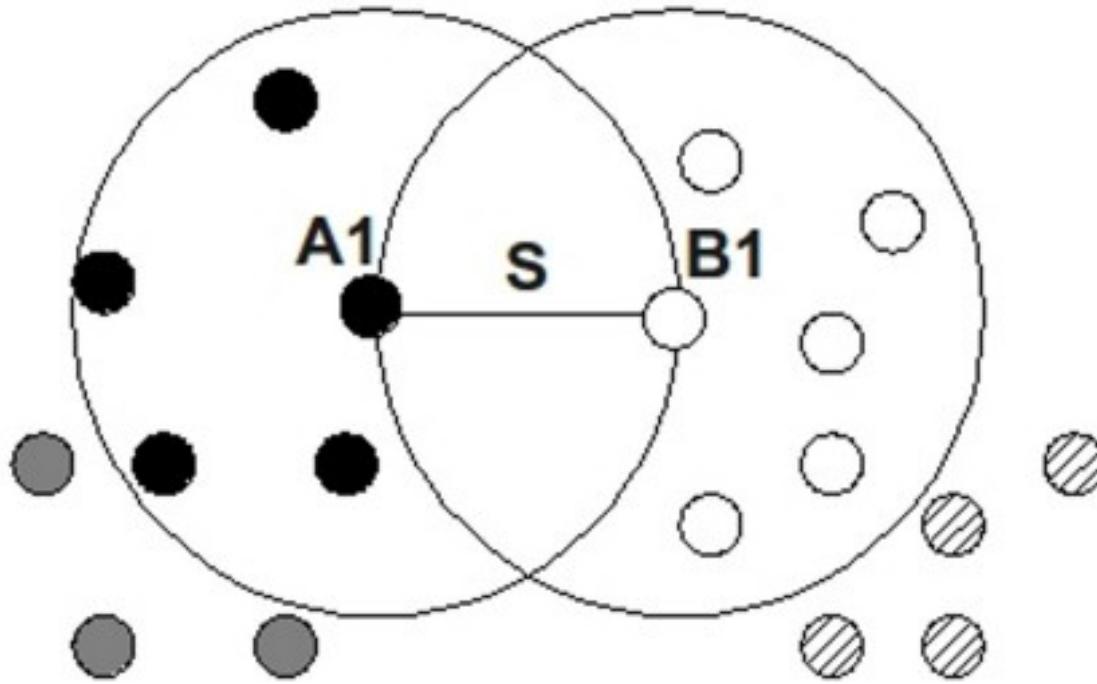


## Best bidirectional hit (BBH), Best reciprocal hits (BRH)









- A1** Seed ortholog species A
- B1** Seed ortholog species B
- S** Score seed orthologs

- Inparalogs species A
- Inparalogs species B
- Outparalogs species A
- Outparalogs species B

**In-paranoid:** improved BRH to detect in-paralogs as well. Works well at the pairwise level.

Note:

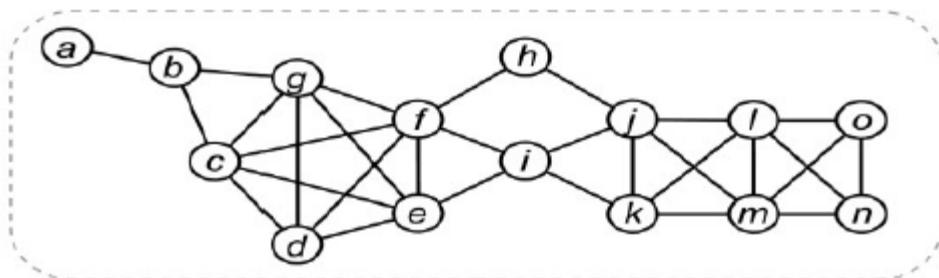
Definition of **in-** and **out-paralogues** require the specification of a given **speciation-node** of reference

## **COG-like** (used by many DBs like STRING)

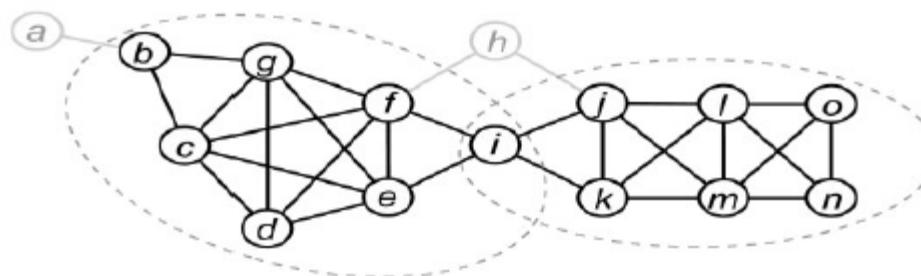
Exploits multi-species information.  
Predicts clusters of orthologous  
groups (in-paralogs) not all pairs in  
a cluster are paralogs.

Can be used at different stringent  
levels

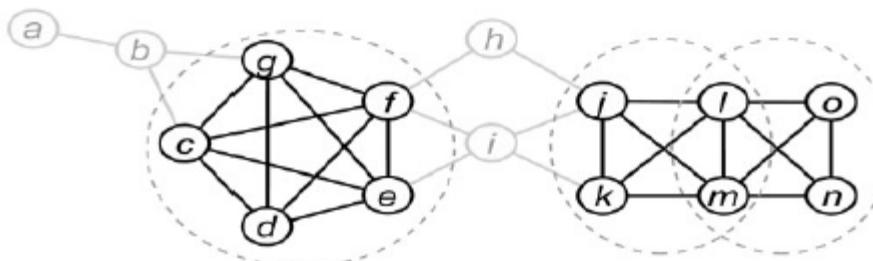
**2**



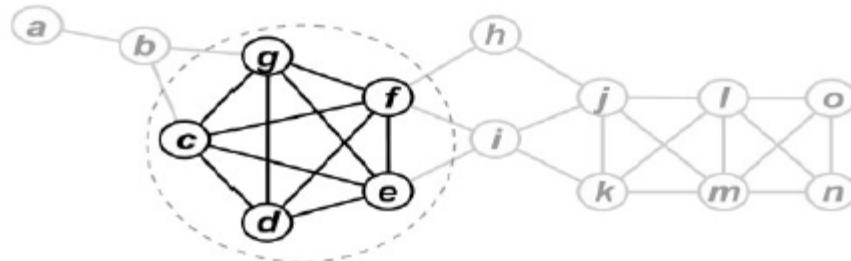
**3**



**4**



**5**



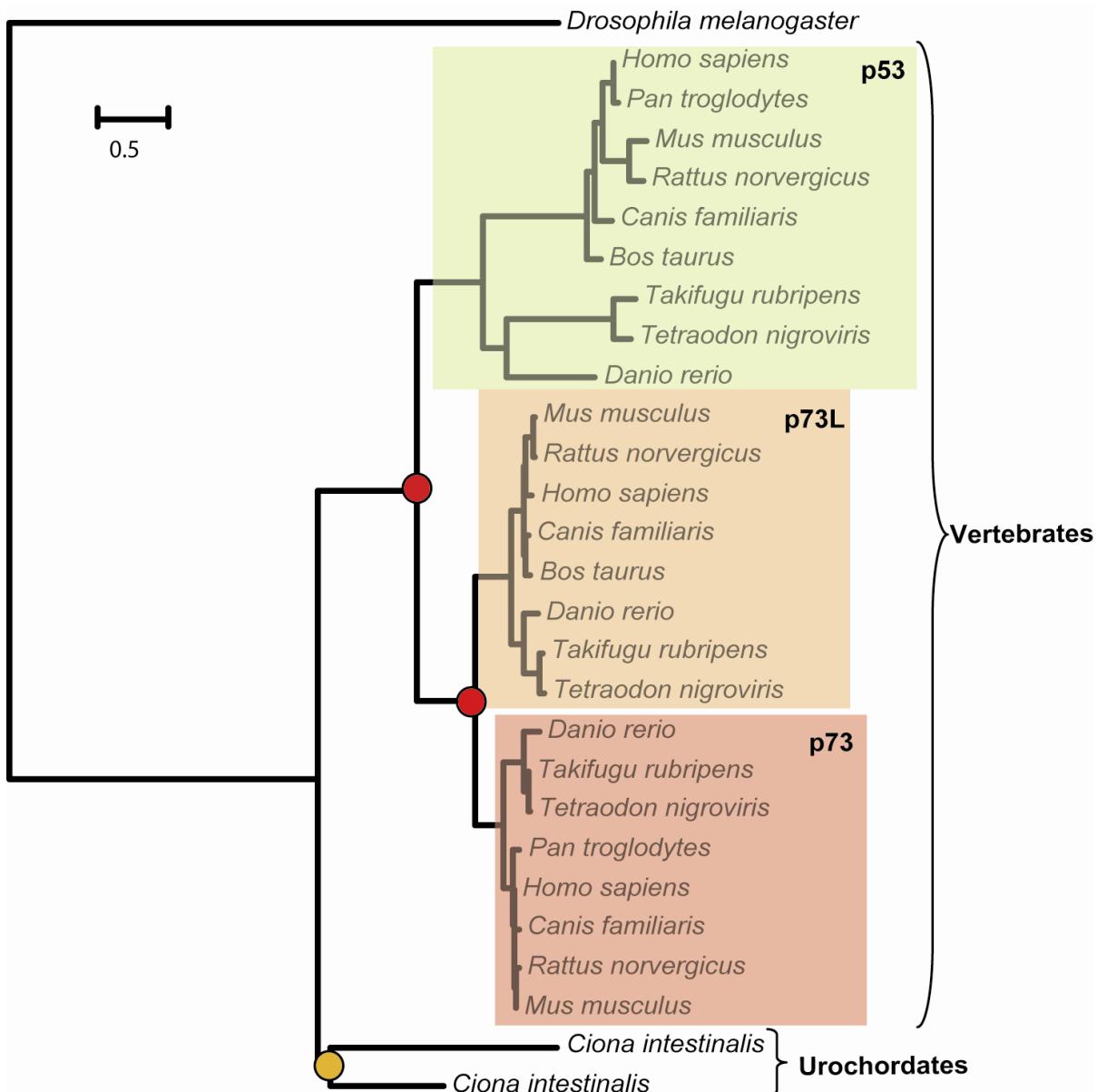
## Clustering methods produce: **orthologous groups**

Equivalent to the earlier concept of **sub-family**

Orthologous groups = Group of sequences derived from a single gene in a common ancestor. They may include orthologs and in-paralogues.

Each orthologous group has implicit the specification of an ancestral species of reference ( a speciation node).

How many orthologous groups? 3 at the level of vertebrates, 1 at the level of chordates



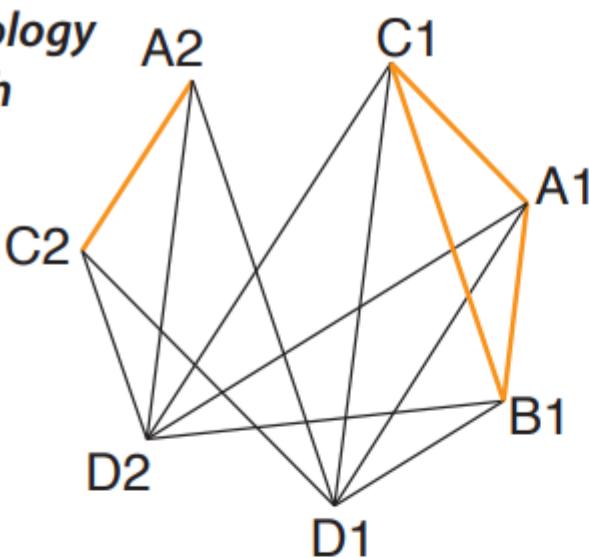
## Additional useful definitions

- **In-paralogs and out-paralogs** (Sohnhammer and koonin): It is defined relative to a given speciation event. In-paralogs are derived from duplications occurred subsequent to the speciation event and are therefore specific of one lineage. Out-paralogs are paralogs emerged from duplications occurred before the speciation. (Important: if you change the speciation events these relationships change)
- **Orthologous group (~Orthogroup):** Also defined relative to a speciation event. It is the complete set of genes in one of the lineages formed by a speciation event. (it includes orthologs and in-paralogs, so not all the genes in an orthologous group are orthologs to each other)

The definition of a reference ancestral species is just an approximation to the inherently hierarchical nature of gene family evolution: and is thus incomplete.

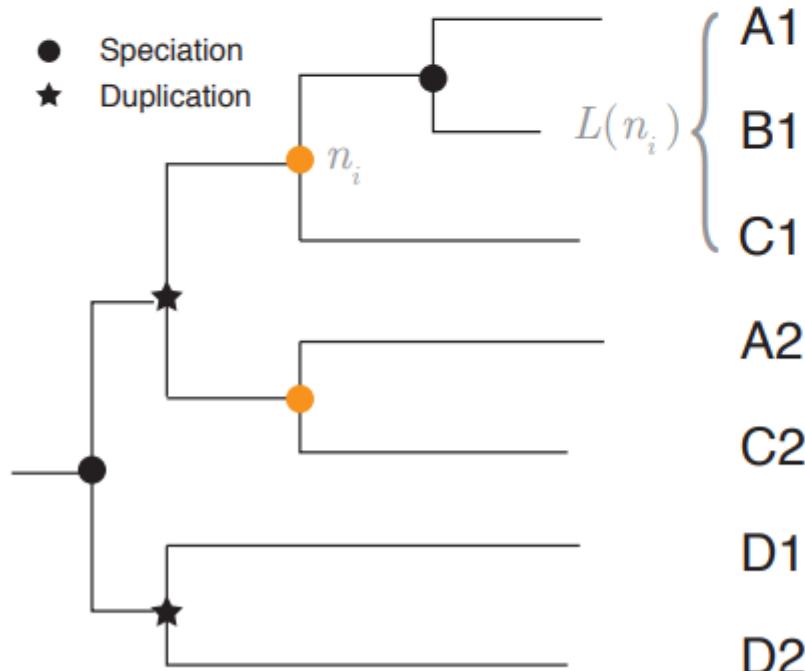
To alleviate this, many databases define orthologous groups at various hierarchical levels (e.g Metazoa, Vertebrates, Mammals, Primates)

*Orthology  
Graph*



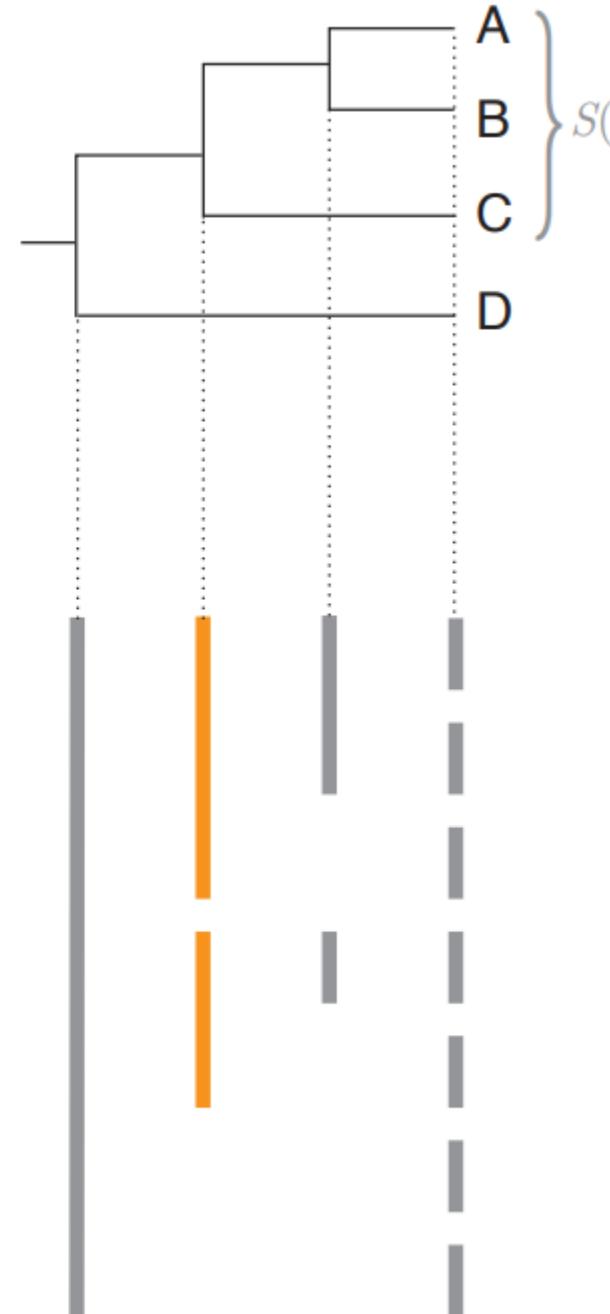
- Orthologs of induced subgraph  $G S(n_i)$
- Other orthologs

- Speciation
- ★ Duplication



*Gene Tree*

*Species Tree*



*Hierarchical Groups*

Methods based on phylogeny were not used at a large scale due to limitations in computational power (phylogenetics is costly).

However, these have changed recently, fast pipelines and algorithms are available:

Ensembl trees, PhylomeDB, TreeFam, etc..

Review

# Large-scale assignment of orthology: back to phylogenetics?

Toni Gabaldón

Bioinformatics and Genomics Program, Center for Genomic Regulation, Doctor Aiguader, 88, 08003 Barcelona, Spain.  
Email: tgabaldon@crg.es

Published: 30 October 2008

*Genome Biology* 2008, **9**:235 (doi:10.1186/gb-2008-9-10-235)

## Abstract

---

Reliable orthology prediction is central to comparative genomics. Although orthology is defined by phylogenetic criteria, most automated prediction methods are based on pairwise sequence comparisons. Recently, automated phylogeny-based orthology prediction has emerged as a feasible alternative for genome-wide studies.

---

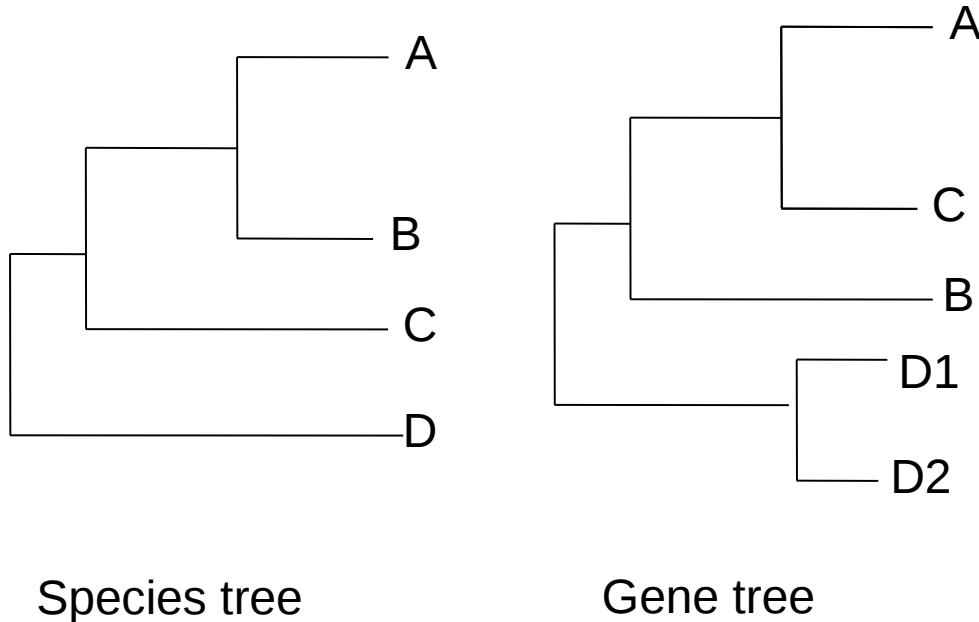
# Phylogeny-based methods

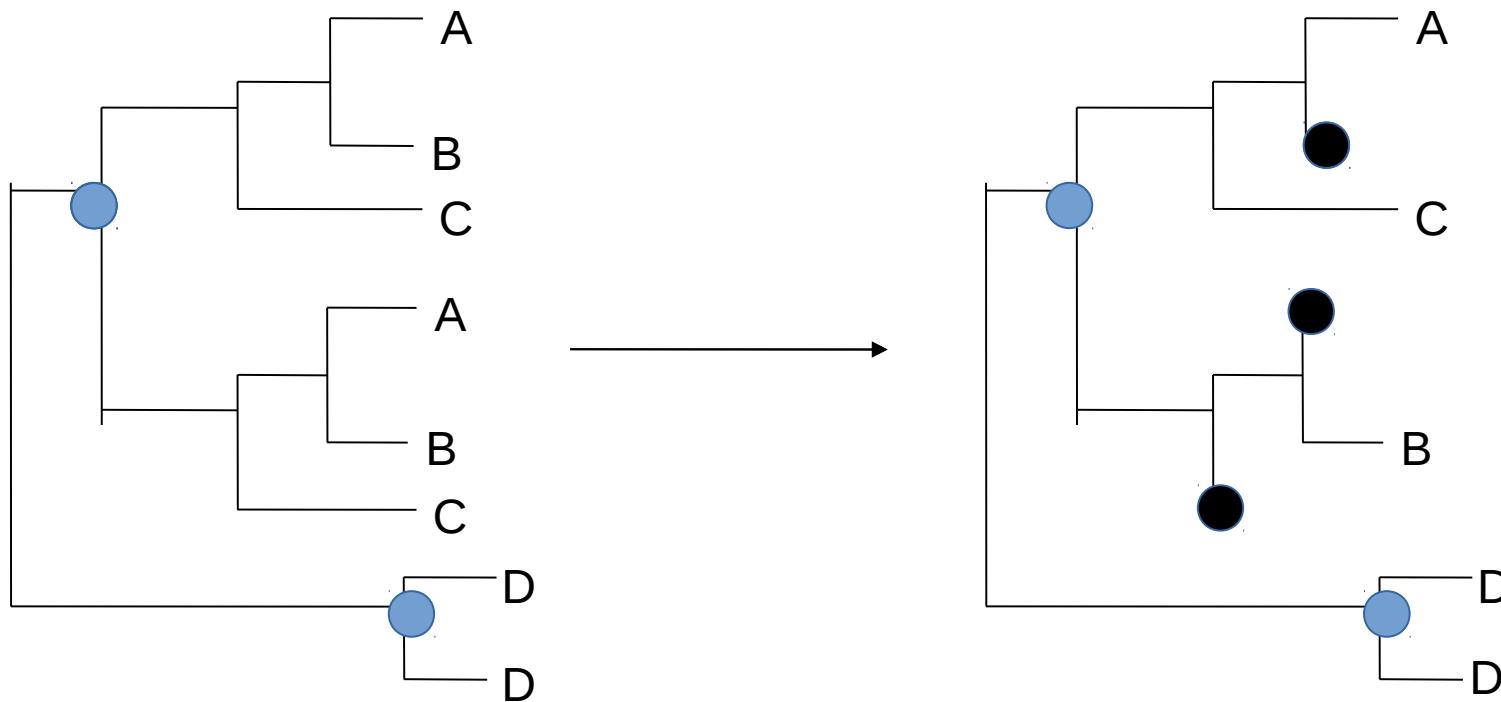
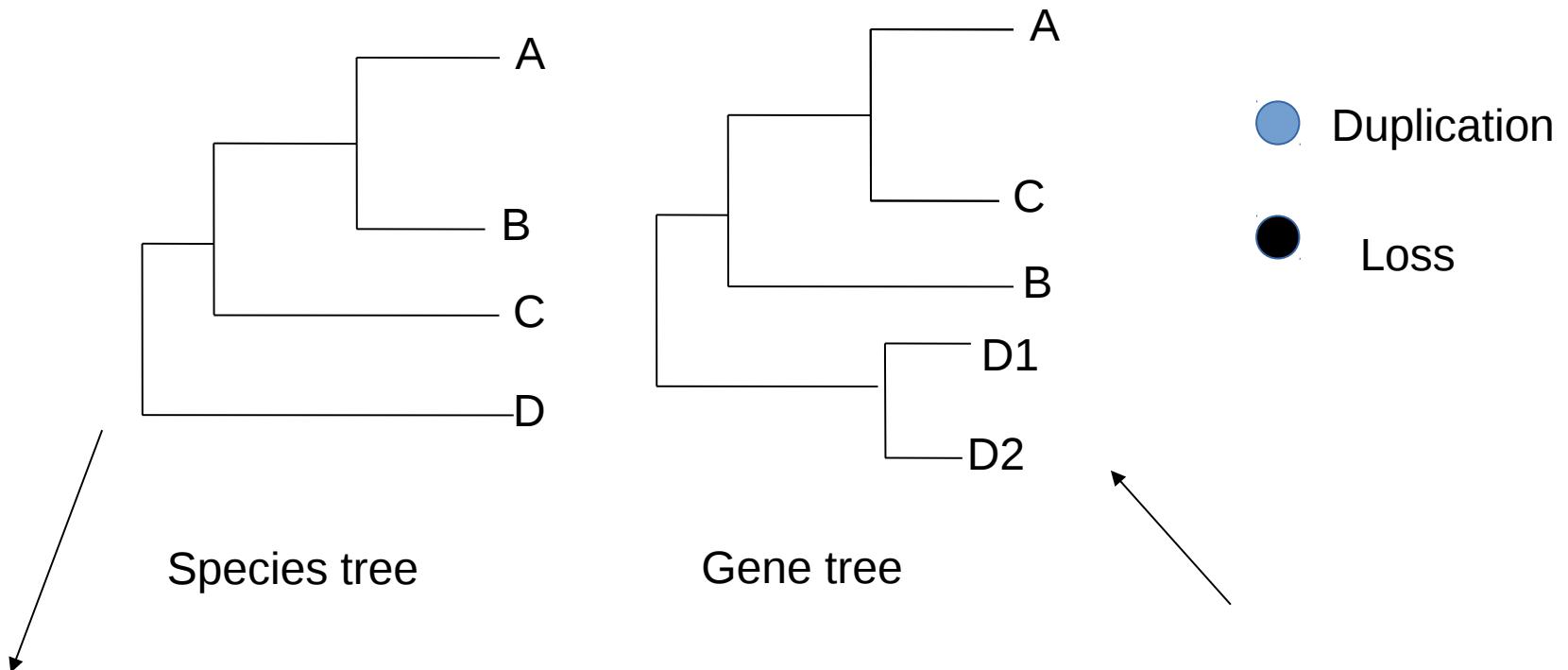
- General procedure: reconstruct the evolution of a gene family (phylogenetics), detect duplication and speciation nodes and predict orthology and paralogy accordingly.
- Two main methods for predicting duplication and speciation nodes from a tree:
  - Species tree reconciliation (RIO, Ensembl)
  - Species-overlap algorithms

## 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



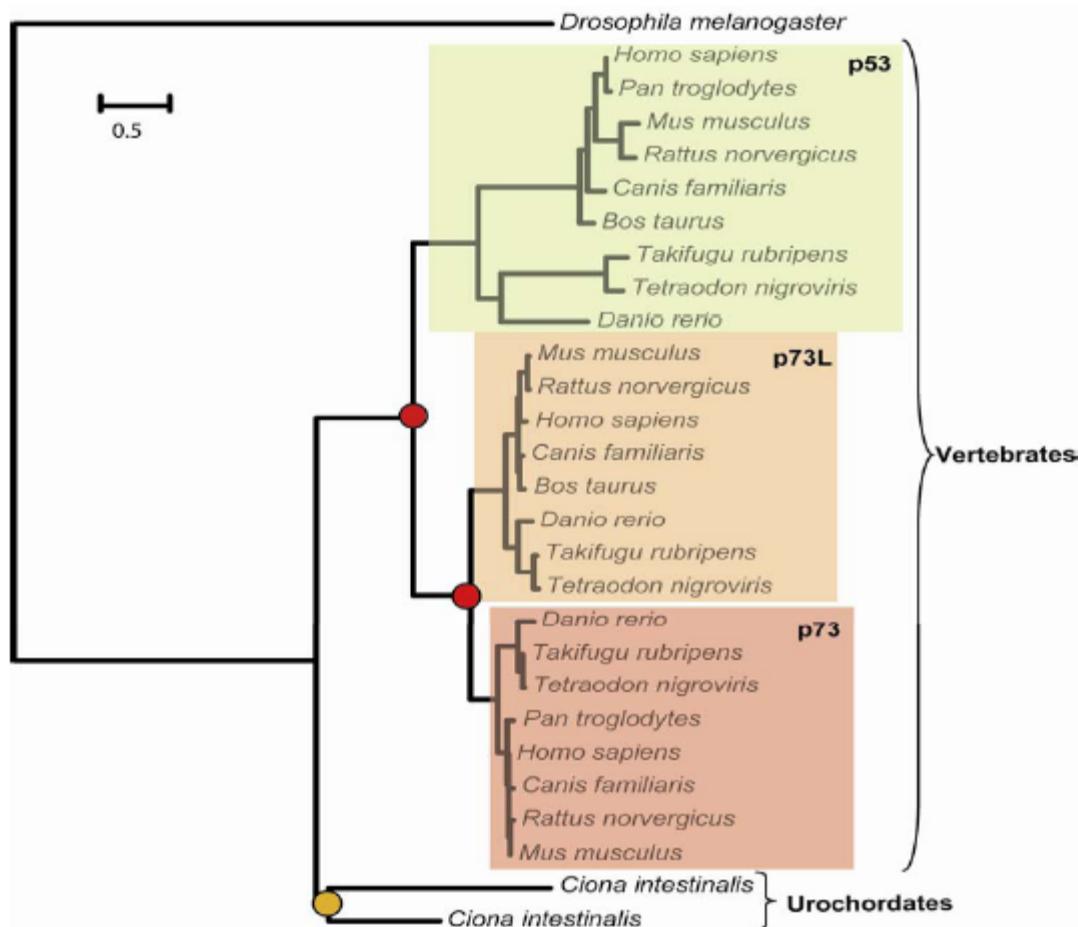


Reconciliation with the species tree readily provides you information on speciation and duplication nodes in a tree

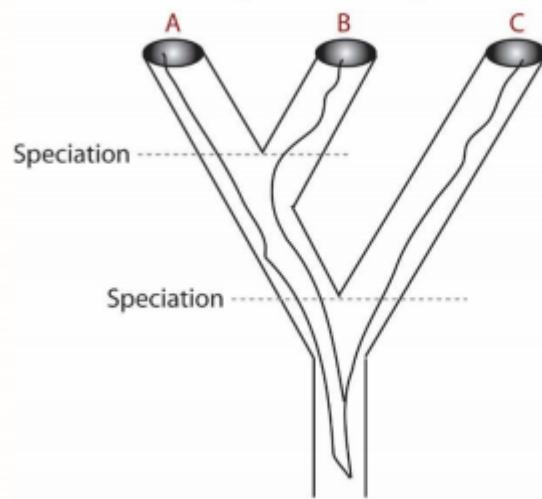
It works when these two assumptions are correct:

A) We know the true species tree

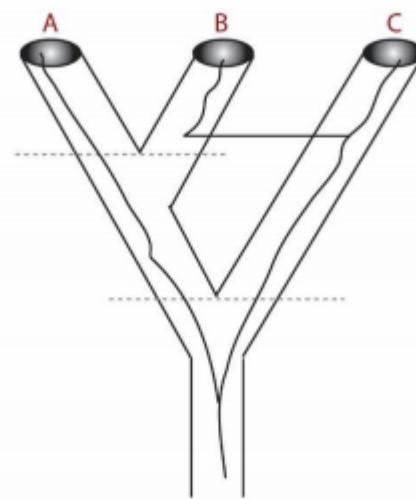
B) The gene tree is correct and reflects the species evolution



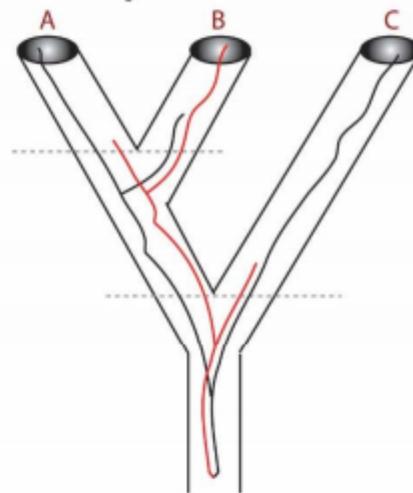
### Lineage Sorting



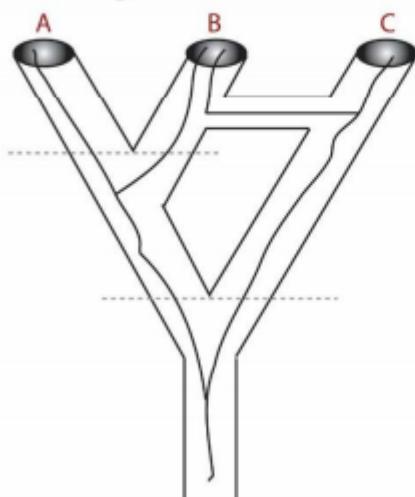
### Horizontal Gene Transfer



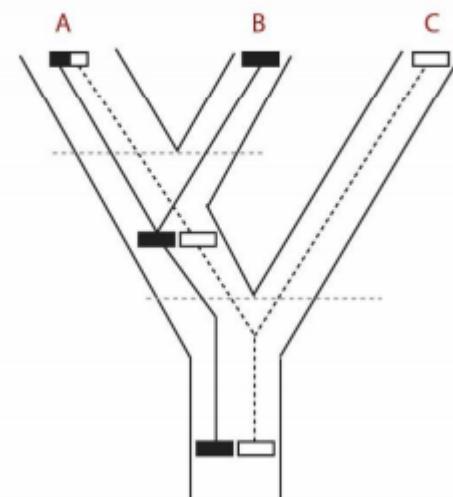
### Gene Duplication and Loss



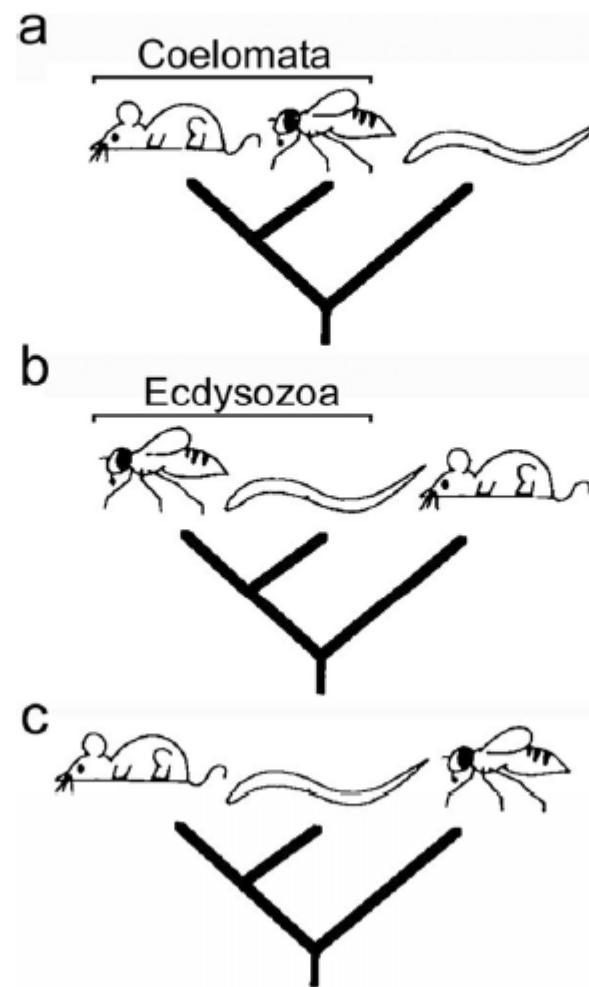
### Hybridization



### Recombination



## Uncertainty in species trees and topological variability in gene trees



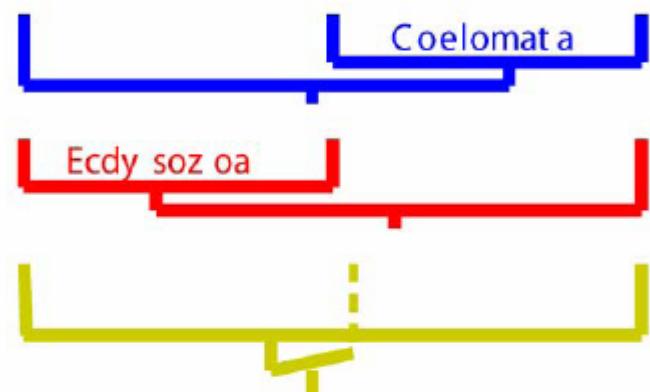
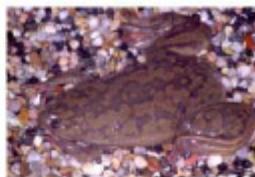
Nematodes



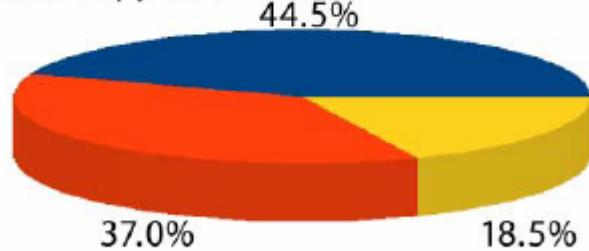
Arthropods



Chordates



Phylome supp ort:



**What percentage of gene trees from the human phylome support each topology?**

Similar results for

Primates

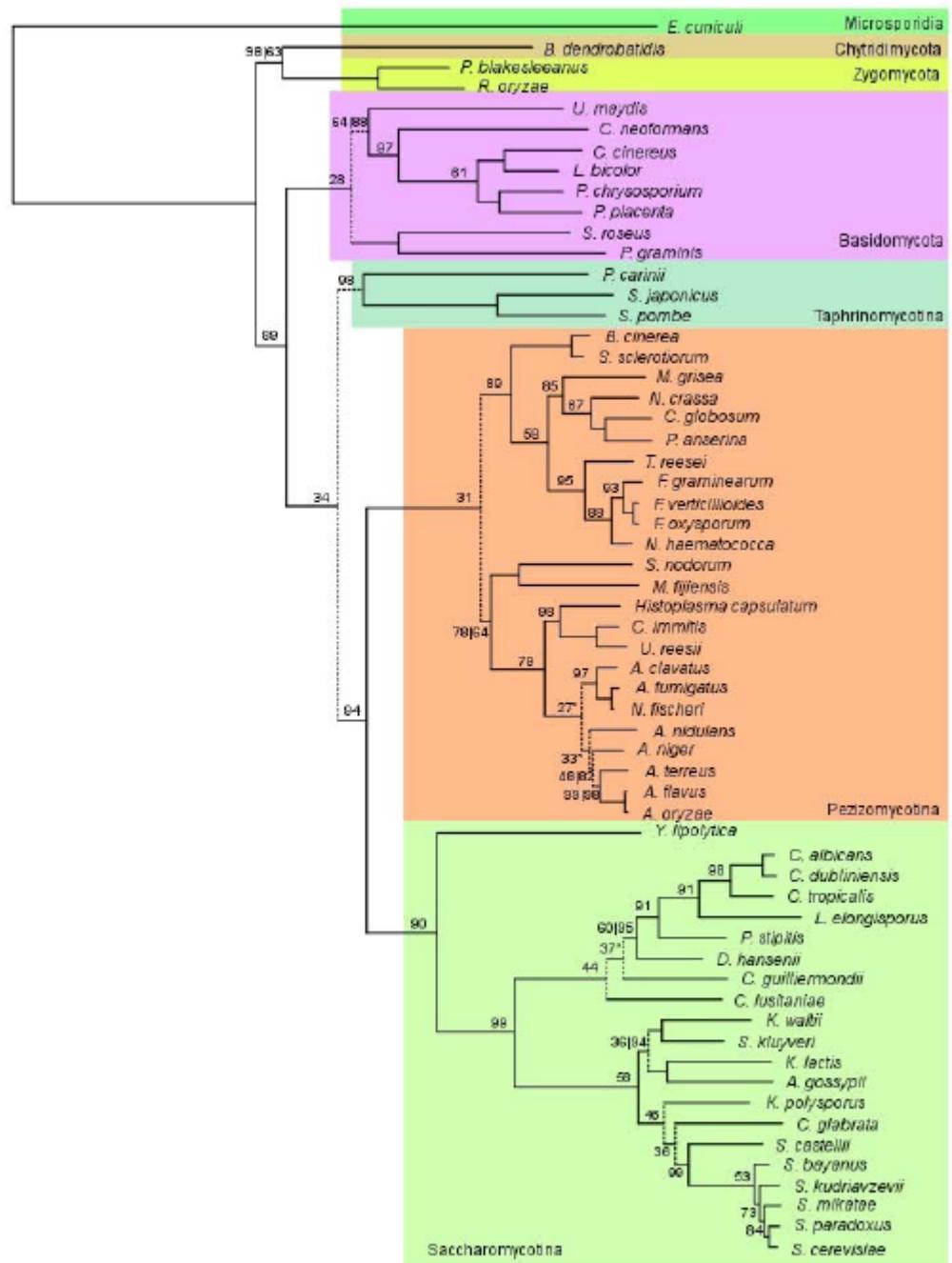
Rodents

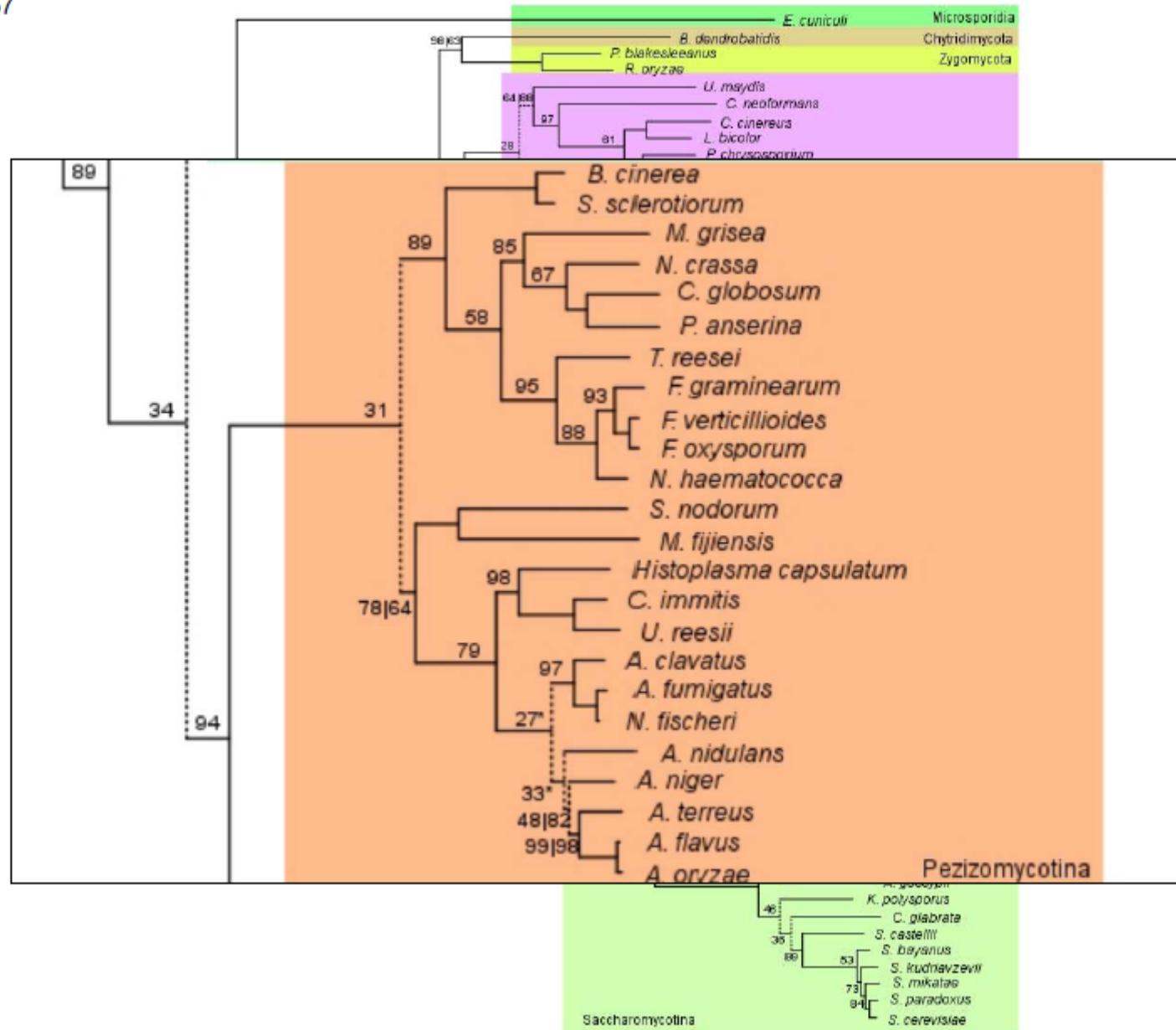
laurasatheria

The tree vs the forest:

Comparison of a fungal species tree with the topological variability of the fungal phylome

Marcet-Houben M and Gabaldón T,  
2009  
PLoS ONE 4(2): e4357





This large-degree of topological variability might be in part due to phylogenetic artifacts, insufficient phylogenetic signal, etc. But also to real evolutionary processes that render a gene tree different from a species tree: lineage sorting, gene conversion, etc

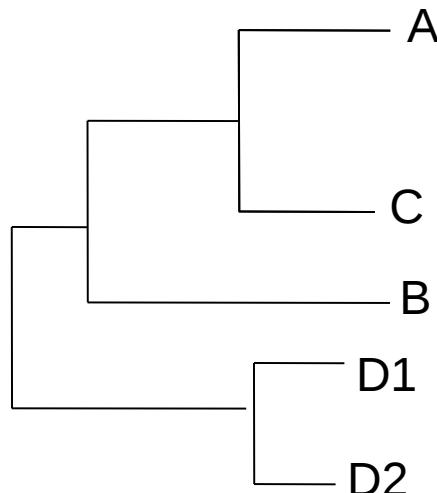
In any case: strict interpretation of gene and species trees will result in many incorrect predictions

## Species overlap algorithm.

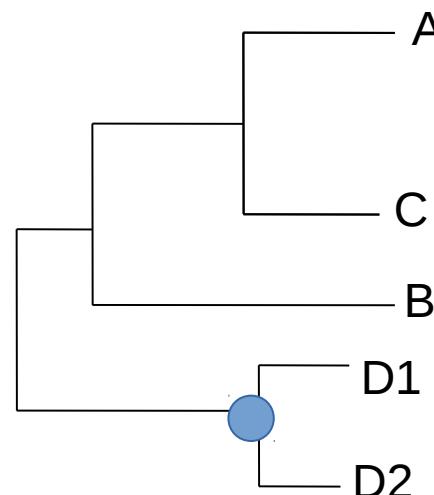
It does not require a species-tree but needs to know the species to which The genes belong

In essence can be seen as a reconciliation with an unresolved species tree

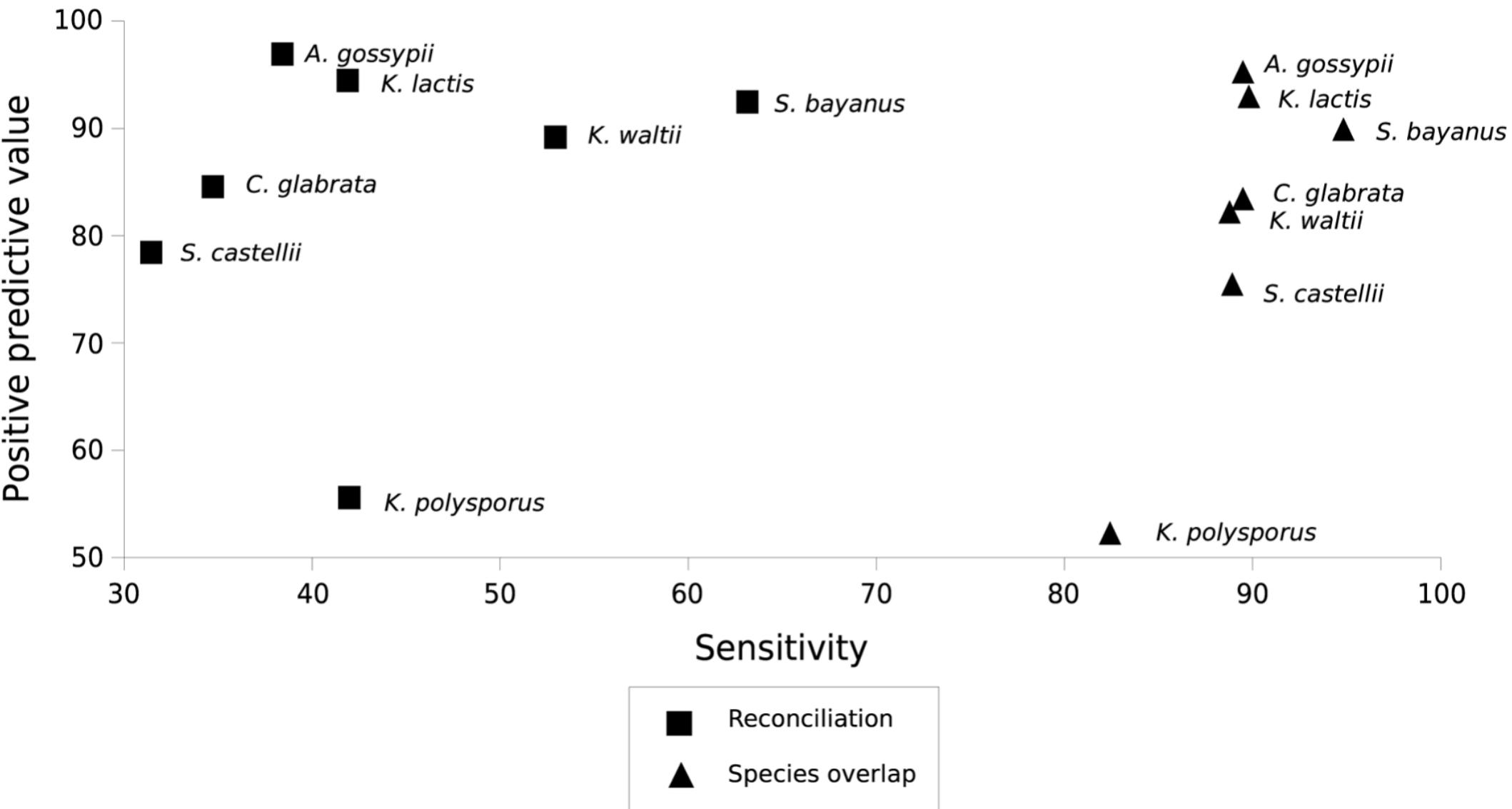
For every node in the gene tree evaluate whether the daughter partitions share any species. If the overlap (number of species shared over total number of species ) is higher than the given threshold. Inpute a duplication at that node.



Gene tree



# T60 orthology prediction benchmark



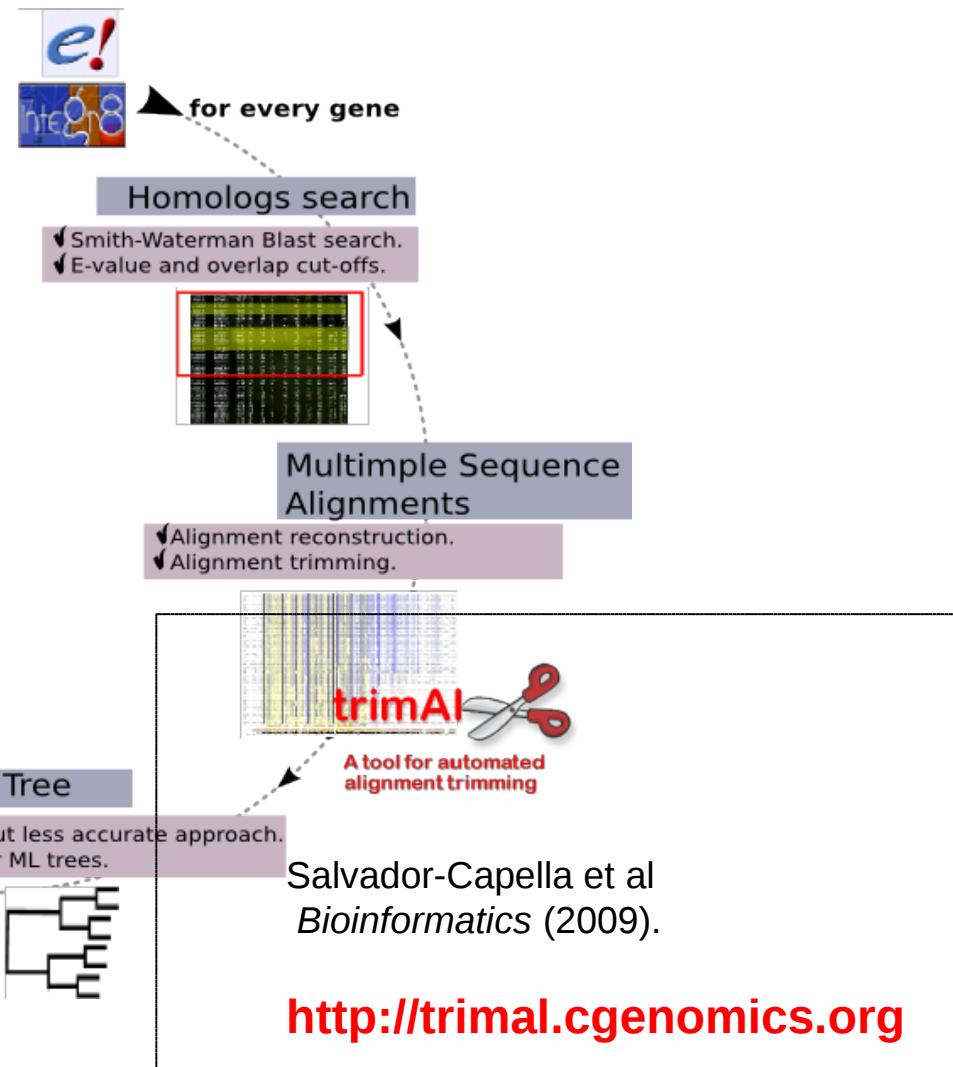
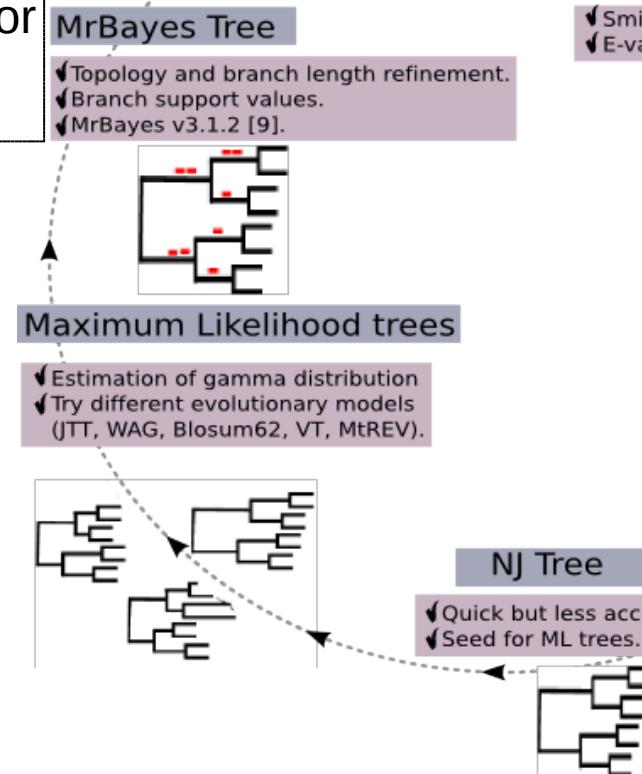
Benchmark on YGOB (Genome alignments and synteny information)

# Our pipeline:

Huerta-Cepas et al.  
*Nucleic Acids Res.* (2008)  
[www.phylomedb.org](http://www.phylomedb.org)



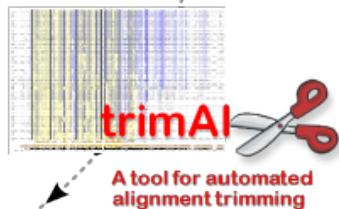
ETE: Environment for  
Tree Exploration  
[ete.cgenomics.org](http://ete.cgenomics.org)



Pipeline described in Huerta-Cepas et al NAR (2011)

## Multiple Sequence Alignments

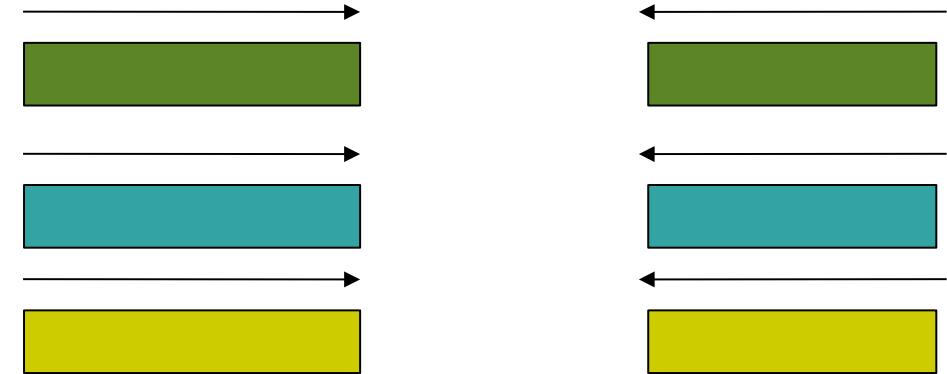
- ✓ Alignment reconstruction.
- ✓ Alignment trimming.



A tool for automated alignment trimming

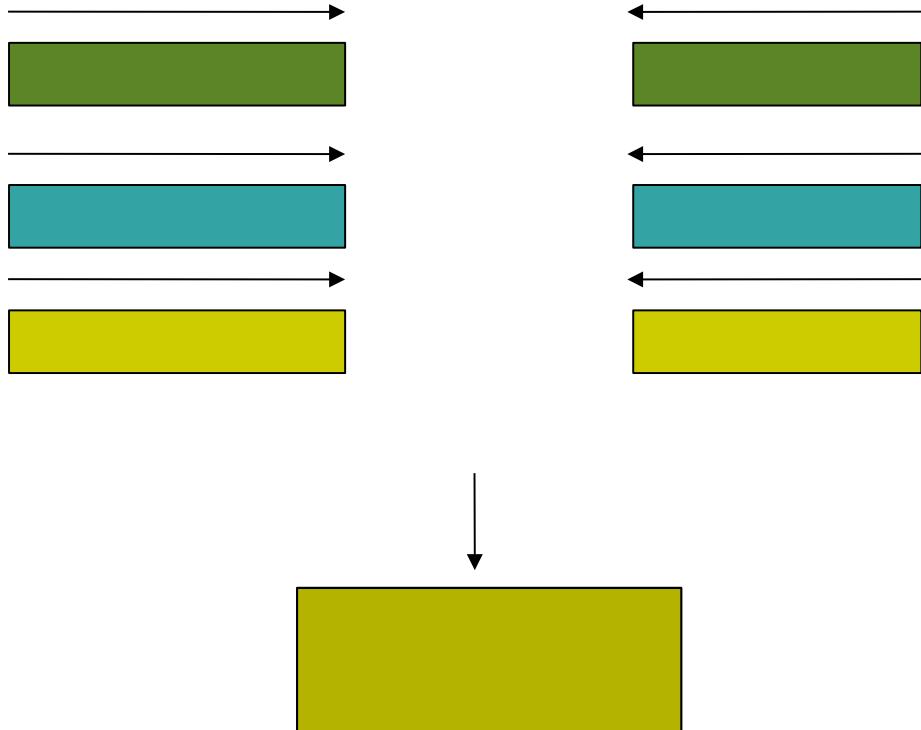
e  
s accurate approach.  
rees.

The set of homologous  
Sequences are aligned by 3 different aligners  
in forward and reverse modes (Head or Tails  
approach)



The set of homologous Sequences are aligned by 3 different aligners in forward and reverse modes (Head or Tails approach)

A consensus is built



sw\_DSBA\_PSESM/1 ---MRNLI I SAALVAASLFGMSAQAAEPIESGKQYV-ELTSAVPV  
sw\_DSBA\_SALTY/1 ---MKKIWLADLAGMVLAFSASAQISDGKQYI-TLDKP--V  
sw\_DSBA\_ENTAM/3 AKWINSIFKSVVLTAAALALPFTAS-AFTE-GTDYM-VLEKP---  
sw\_DSBA\_LEGPN/1 -----LMPMTALATQFIE-GKDYQTVASAQ-LS

cons



sw\_DSBA\_PSESM/1 AVPGK-IEVIELFWYGCPhCYAFEPTI---NPWVEKLPSDVNFVR  
sw\_DSBA\_SALTY/1 --AGE-PQVLEFFSFYCPHcyQFEEVLHVSDNVKKLPEGTKMTE  
sw\_DSBA\_ENTAM/3 -IPDADKTLIKVFSYACPFCYKDYKAVT--GPVADKVADLVTFVP  
sw\_DSBA\_LEGPN/1 TNKDKTPLITEFFSYGCPWCYKIDAPLN--D-WATRMGKGHLER

cons

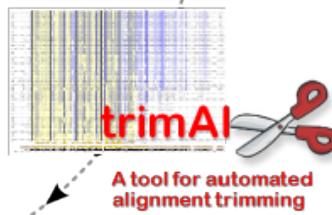


Ht

## Multimle Sequence Alignments

Alignment reconstruction.

Alignment trimming.



A tool for automated alignment trimming

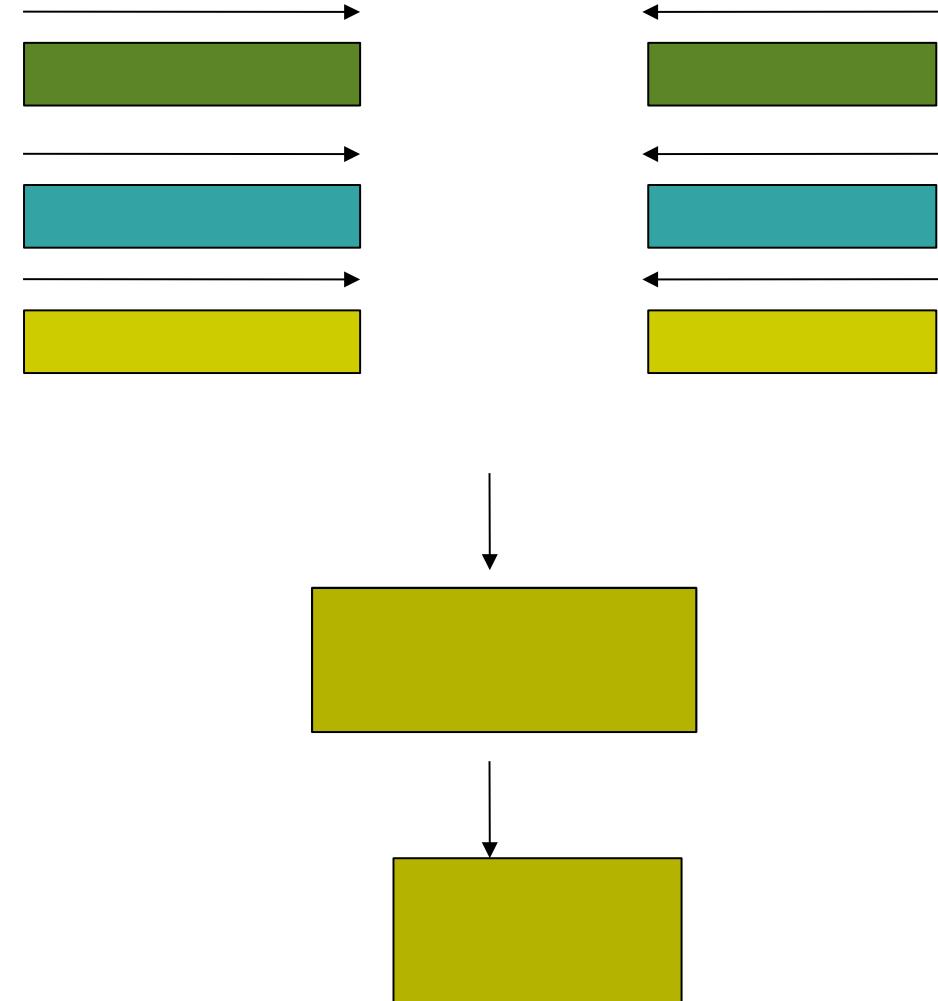
e approach.

The set of homologous  
Sequences are aligned by 3 different aligners  
in forward and reverse modes (Head or Tails  
approach)

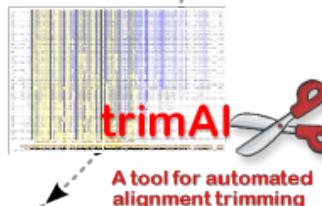
A consensus is built

The consensus is trimmed (trimAl)  
based on:

- consistency across the 6 alignments
- gap content



Alignment reconstruction.  
Alignment trimming.



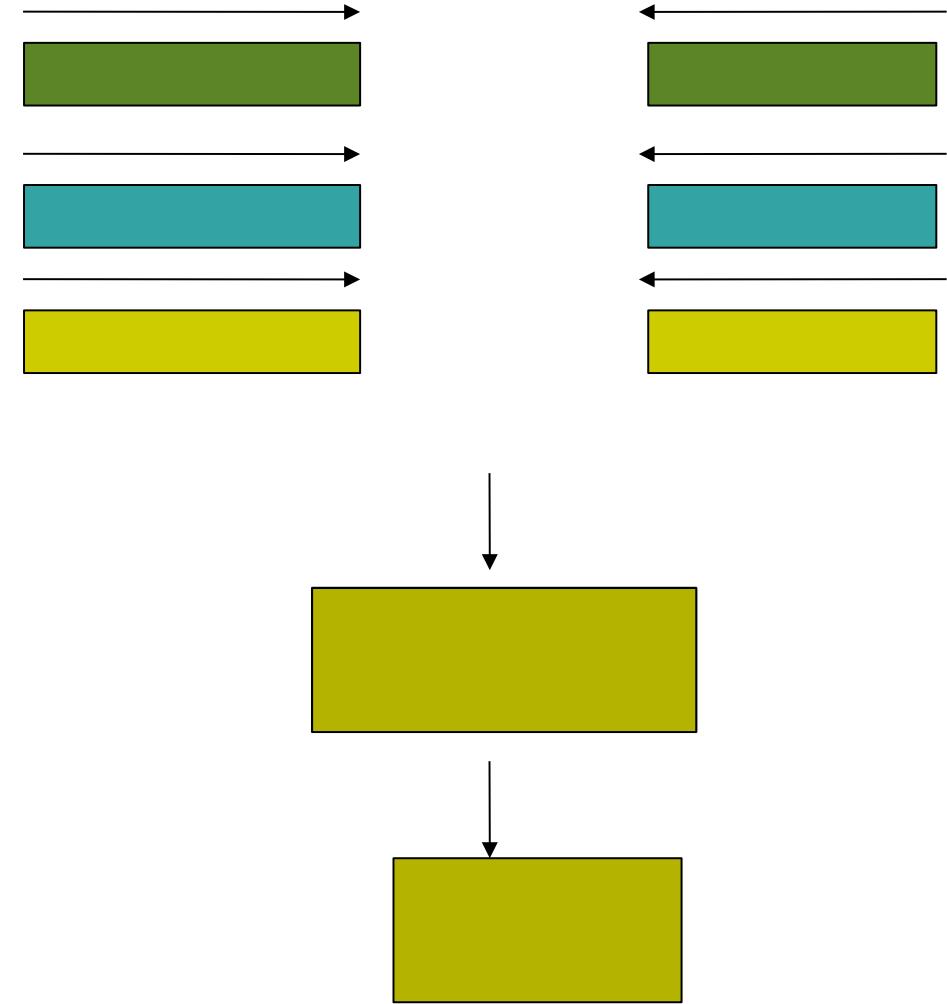
urate approach.

The set of homologous  
Sequences are aligned by 3 different aligners  
in forward and reverse modes (Head or Tails  
approach)

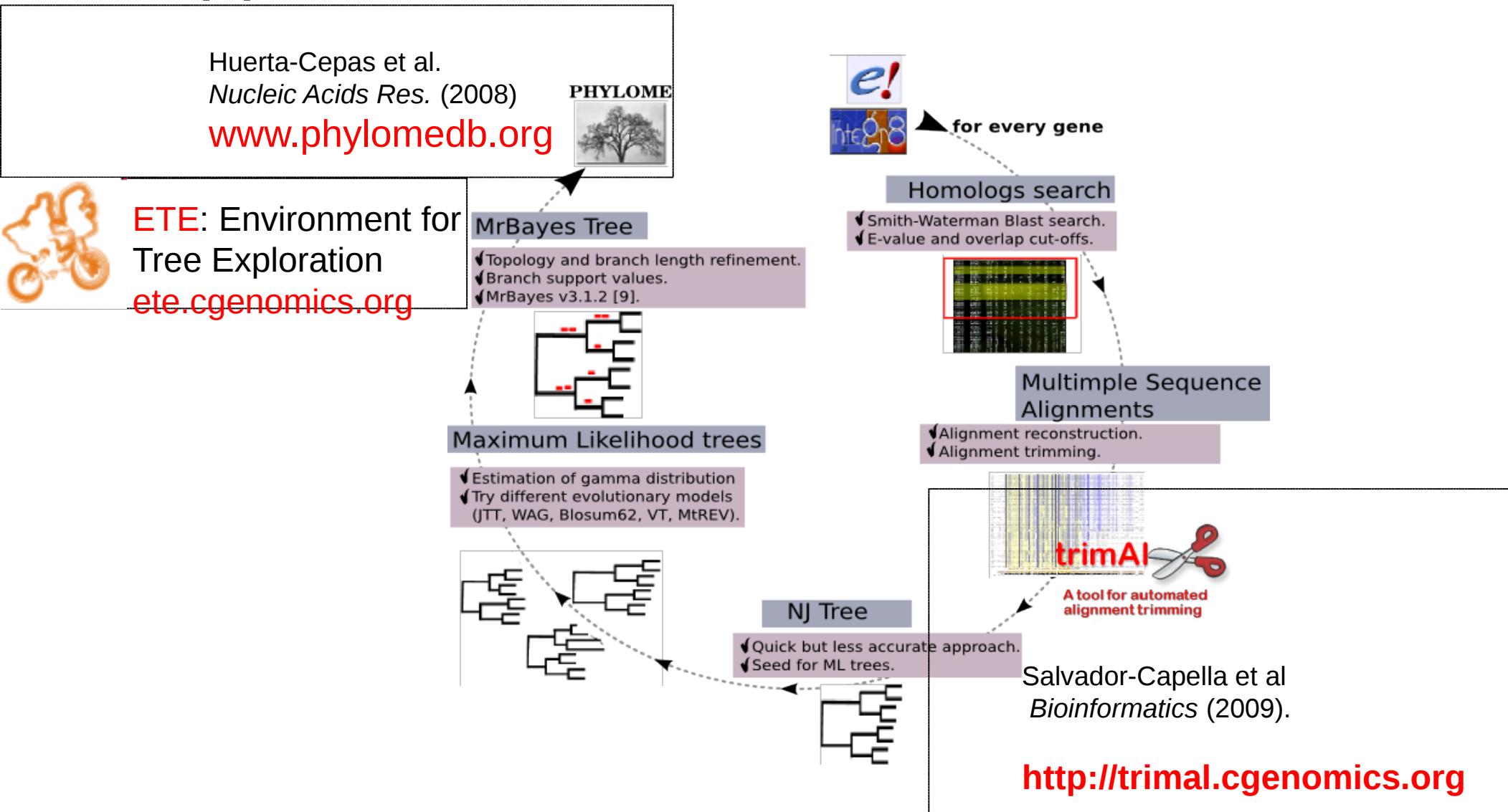
A consensus is built

The consensus is trimmed (trimAI)  
based on:

- consistency across the 6 alignments
- gap content



# Our pipeline:



Pipeline described in Huerta-Cepas et al NAR (2011)

**Search in PhylomeDB**(i.e. ENSG00000139618, YBL058W,  
TP53 )[RandomTree!](#)[BLAST search](#)**Latest Phylomes**

Arxula adeninivorans	2014
Beta vulgaris	2013
Clogmia albipunctata	2013
Penicillium digitatum	2012
Schistosoma mansoni	2012

[see all phylomes](#)**PhylomeDB uses****PhylomeDB cross linking**

# Welcome to PhylomeDB 4!

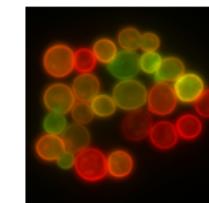
PhylomeDB is a public database for complete **catalogs of gene phylogenies** (phylomes). It allows users to interactively explore the evolutionary history of genes through the visualization of phylogenetic trees and multiple sequence alignments. Moreover, phylomeDB provides genome-wide orthology and paralogy predictions which are based on the analysis of the phylogenetic trees. The automated pipeline used to reconstruct trees aims at providing a high-quality phylogenetic analysis of different genomes, including Maximum Likelihood tree inference, **alignment trimming** and evolutionary model testing.

PhylomeDB includes also a public download section with the complete set of trees, alignments and orthology predictions, as well as a **web API** that facilitates cross linking trees from external sources. Finally, phylomeDB provides an advanced tree visualization interface based on the **ETE toolkit**, which integrates tree topologies, taxonomic information, domain mapping and alignment visualization in a single and interactive tree image.

[What's new in phylomeDB 4?](#)**Latest story****New Zygomycete phylome: the human pathogen *Lichtheimia corymbifera***

Mon, 09/15/2014 - 21:09

PhylomeDB extends its repertoire of fungal phylomes with that of a genome of a poorly sample clade, that of the basal group zygomycetes. In this case the phylome (245) of the human pathogenic mucorales *Lichtheimia corymbifera* has served to reveal extensive past gene duplications in this group. *Lichtheimia* species are the second most important cause of mucormycosis in Europe. The sequencing of its genome and the comparison with other Zygomycete species, particularly of *Rhizopus delemar*, the main

**Popular Phylome Collections****Human****Fungi****Plants****Model Species****Latest News**

 **New Zygomycete phylome: the human pathogen *Lichtheimia corymbifera***

Mon, 09/15/2014 - 21:09

 **New hemiascomycete phylome: *Blastobotrys (Axula) adeninivorans*, a yeast of biotechnological interest.**

Mon, 05/19/2014 - 11:01

 **Help us to improve PhylomeDB: complete our survey.**

Thu, 02/20/2014 - 16:22

[show all](#)**PhylomeDB Twitter****Tweets**[Follow](#)

 **phylomedb** @phylomedb

23 Oct

New birds, crocs, and fungal phylomes to come soon at phylomeDB. stay tuned!  
[Expand](#)

 **phylomedb** @phylomedb

15 Sep

New Zygomycete phylome: the human fungal pathogen *Lichtheimia corymbifera*  
[phylomedb.org/?q=node/537](#)  
[Expand](#)

 **phylomedb** @phylomedb

26 Aug

NOTICE: PhylomeDB will be down due to MAINTENANCE



[Login] Home

Collections All phylomes Downloads Help FAQ About

## Search in PhylomeDB

(i.e. ENSG00000139618, YBL058W,

TP53 )

Search

RandomTree!

## BLAST search

## Latest Phylomes

Clogmia albipunctata	2013
Penicillium digitatum	2012
Schistosoma mansoni	2012
<a href="#">Cucumis melo</a>	2012

[see all phylomes](#)

## PhylomeDB uses



## TP53 tree in phylome 218

AS seed in Rat phylome

JTT (Ik:-18130.4)

-- in collateral trees --

Tree features

Search

Clear search

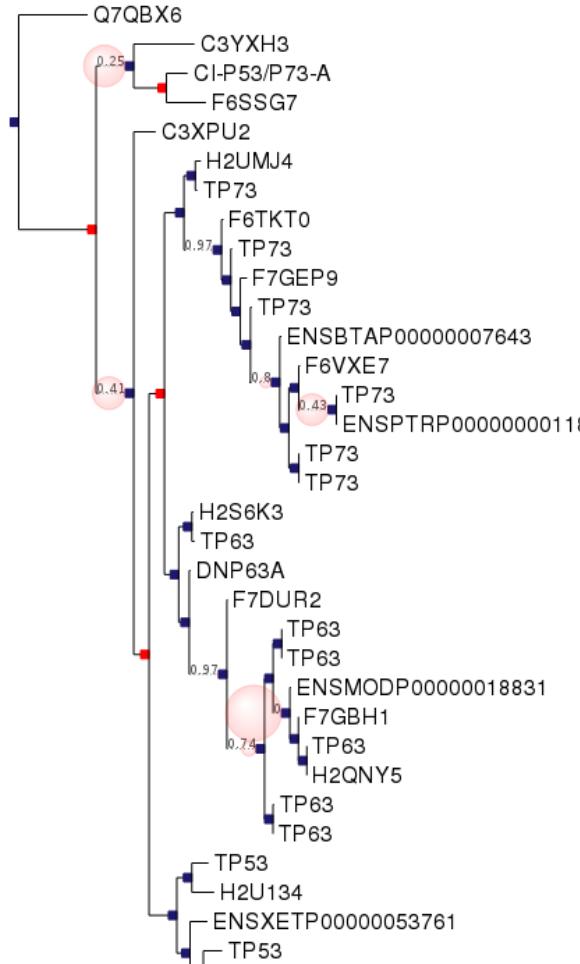
Image

Hard link

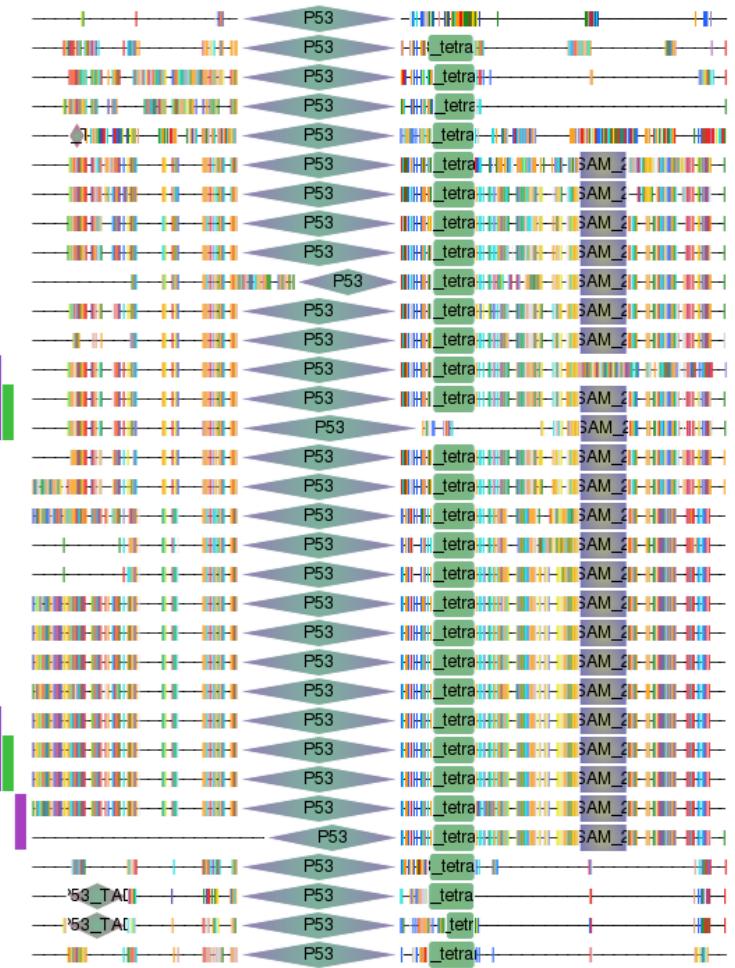
Download OrthoXML

See alignments

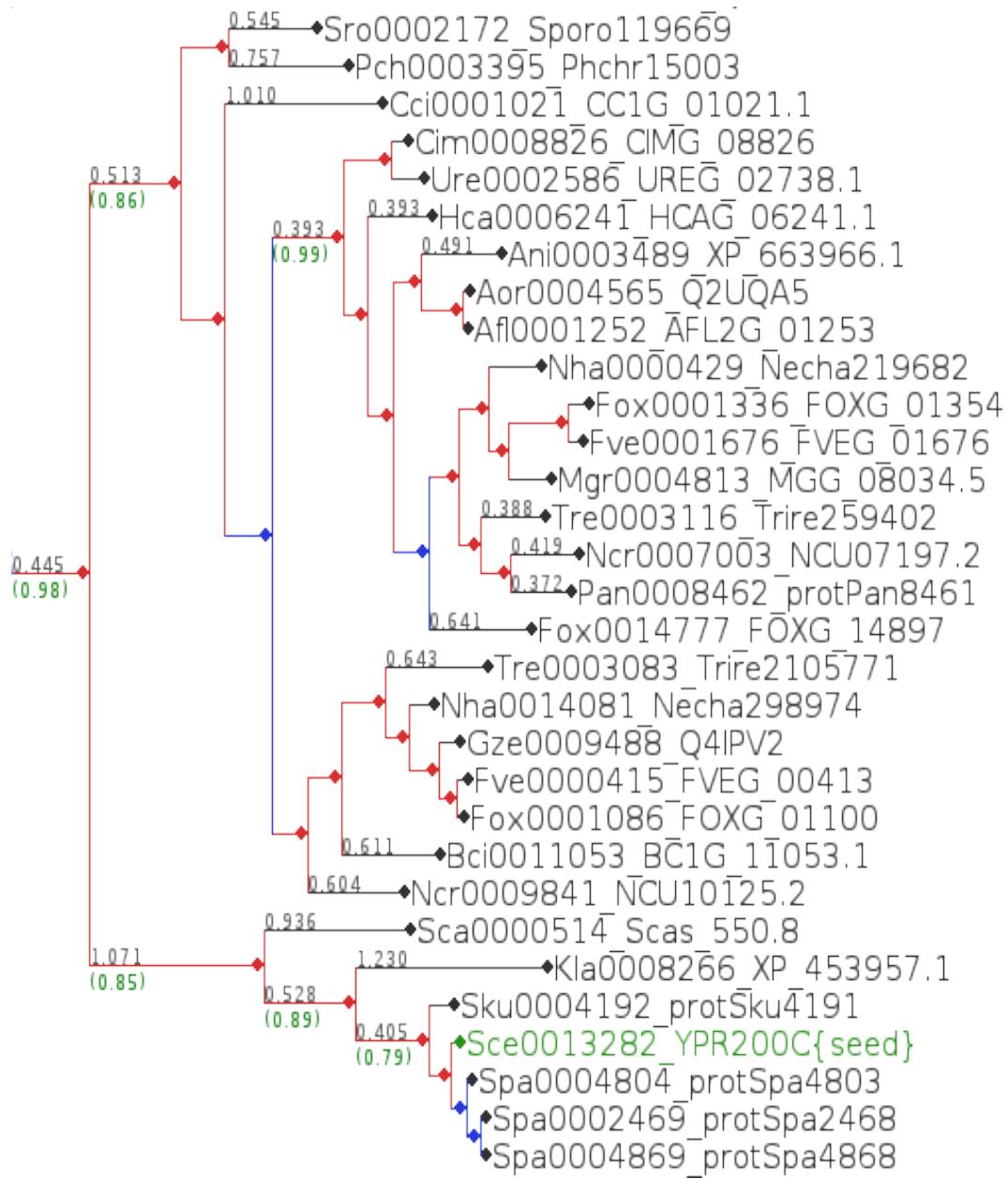
Download data.tar.gz



Anopheles gambiae  
Branchiostoma floridae  
Ciona intestinalis  
Ciona intestinalis  
Branchiostoma floridae  
Takifugu rubripes  
Danio rerio  
Xenopus tropicalis  
Gallus gallus  
Monodelphis domestica  
Canis familiaris  
Bos taurus  
Macaca mulatta  
Homo sapiens  
Pan troglodytes  
Mus musculus  
Rattus norvegicus  
Takifugu rubripes  
Danio rerio  
Gallus gallus  
Ornithorhynchus anatinus  
Rattus norvegicus  
Mus musculus  
Monodelphis domestica  
Macaca mulatta  
Homo sapiens  
Pan troglodytes  
Canis familiaris  
Bos taurus  
Danio rerio  
Takifugu rubripes  
Xenopus tropicalis  
Gallus gallus



**These phylomes can now be interrogated in many ways**

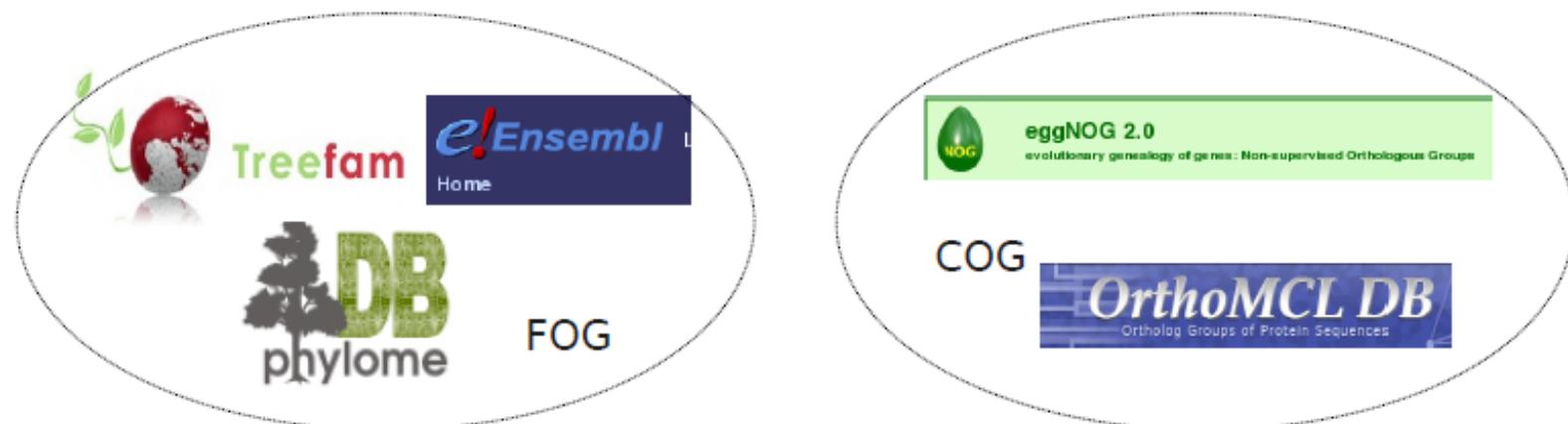


- Families that show a particular topology
  - Detect and date duplication events
  - Genes that have accelerated evolutionary rates at a particular lineage (positive/relaxed selection)
  - Families expanded at particular lineages
  - Footprints of horizontal gene transfer, lineage sorting, gene conversion and other evolutionary processes
  - Search for co-evolving genes
  - predict functional properties
  - across-species prediction of orthology and paralogy



# MetaPhOrs

(Meta-Phylogeny-Based-Orthologs)

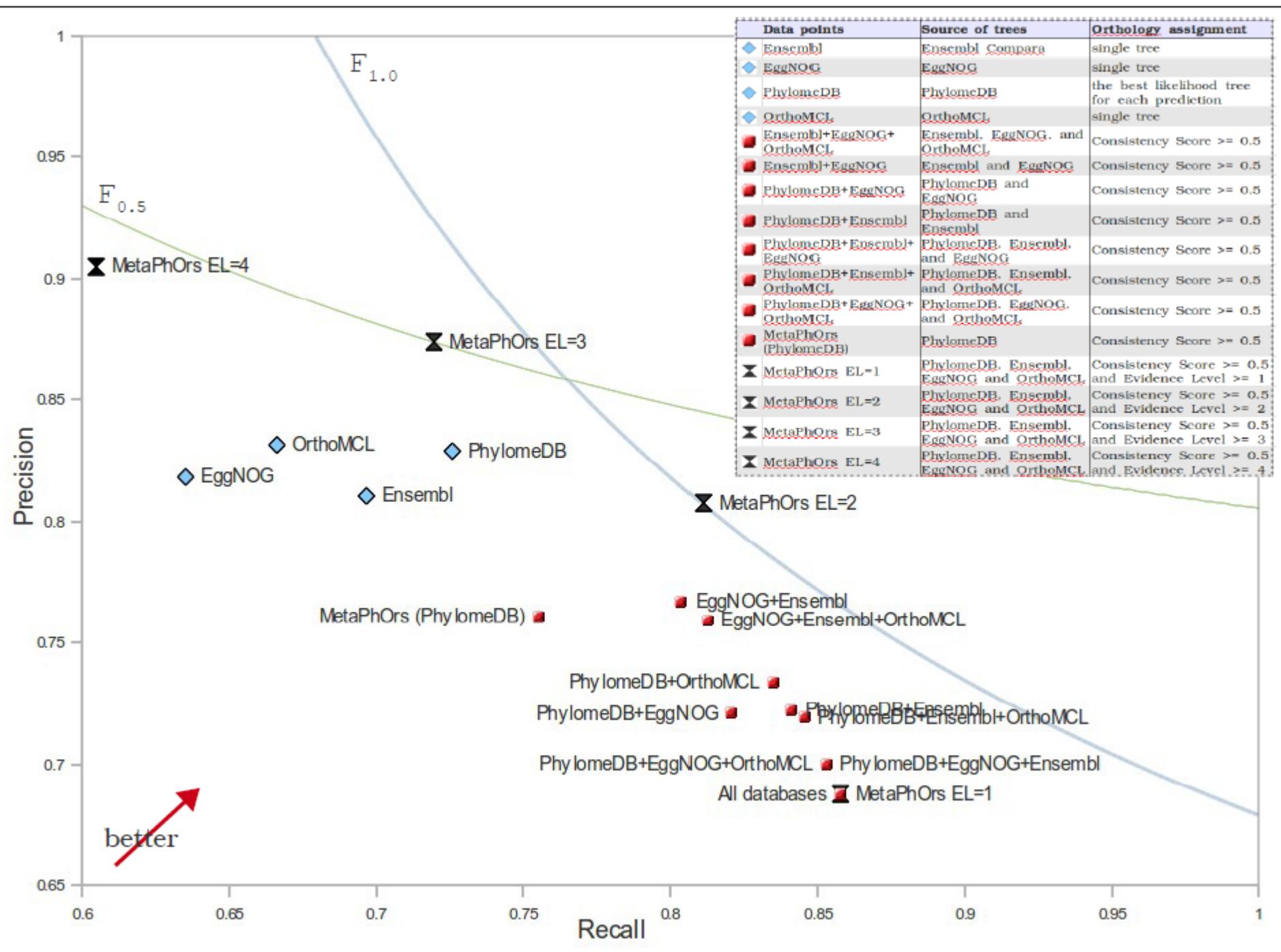


Use existing tree repositories

Reconstruct trees for orthologous groups

Integrate and use consistency across datasets as a proxy of reliability

result: phylogeny-based predictions across 800 genomes with a confidence score



# A plethora of methods for ortholog prediction

<http://questfororthologs.org>

## QUEST FOR ORTHOLOGS

ORTHOLOGY DATABASES  
DOCUMENTS (INTRANET)  
MAILING-LIST & CONTACT

\*\*\* More info on Quest for Orthologs 5 in Los Angeles, 8-10 June 2017 \*\*\*

### Welcome

This is the site of the Quest for Orthologs consortium. Proteins and functional modules are evolutionarily conserved even between distantly related species, and allow knowledge transfer between well-characterized model organisms and human. The underlying biological concept is called 'Orthology' and the identification of gene relationships is the basis for comparative studies.

More than 30 phylogenomic databases provide their analysis results to the scientific community. The content of these databases differs in many ways, such as the number of species, taxonomic range, sampling density, and applied methodology. What is more, phylogenomic databases differ in their concepts, making a comparison difficult – for the benchmarking of analysis results as well as for the user community to select the most appropriate database for a particular experiment.

The Quest for Orthologs (QfO) is a joint effort to benchmark, improve and standardize orthology predictions through collaboration, the use of shared reference datasets, and evaluation of emerging new methods.

The main sections of this site are:

- [Meetings](#)
- [Community Standards](#) (Reference proteome, standardized formats, benchmarking, etc..)
- [Working groups](#)
- [Orthology databases](#)
- [Documents \(Intranet\)](#)
- [Mailing-List and Contact](#)

To contribute to this website, please create an account (see below) and [contact us!](#)

[ Back to top | Sitemap ]

[ Log In | Old revisions ]

prsnl10 on DW under the hood | home.txt · Last modified: 2017/03/06 21:19 by Christophe Dessimoz

Database	Description / Scientific focus applications (Max. 2 sentences)	Last updated	Update frequency	QfO Prote
DIOPT	Integrative ortholog prediction tool of 10 algorithms	2016		partial
eggNOG	A database for phylogenetically refined Orthologous Groups and functional annotation.	2016	biennial	no
Ensembl Compara	<a href="#">Evolutionary relationships among Ensembl species genes</a> ; Projection of	2016	4-5x / year	no

all domains of life through 6 divisions (sets of  
eukaryotes, 352 viruses

66 chordates and 240 others

yes

yes

¿With over 30 orthology databases, based on various methods, which ones to choose?

- Different taxonomic focuses
- Different methodologies
- Different outputs (pairwise relationships, groups, etc)
- Different interfaces
- Different accuracies (**how to benchmark this?**)

## **Final warnings:**

Most methods assume the complete, fully (and correctly) annotated genome for each of the compared species is available.

## Deserve special considerations:

- Working with highly fragmented/incomplete genomes or transcriptomes
- Working with bacteria (pangenome concept, rampant HGT)

## The “boundaries” of the orthology concept

*Where the homology is the result of gene duplication so that both copies have descended side by side during the history of an organism, (for example, alpha and beta hemoglobin) the genes should be called paralogous (para = in parallel).*

*Where the homology is the result of speciation so that the history of the gene reflects the history of the species (for example alpha hemoglobin in man and mouse) the genes should be called orthologous (ortho = exact)."*

## The “boundaries” of the orthology concept

*Where the homology is the result of **gene duplication** so that both copies have descended side by side during the history of an **organism**, (for example, alpha and beta hemoglobin) the **genes** should be called **paralogous** (para = in parallel).*

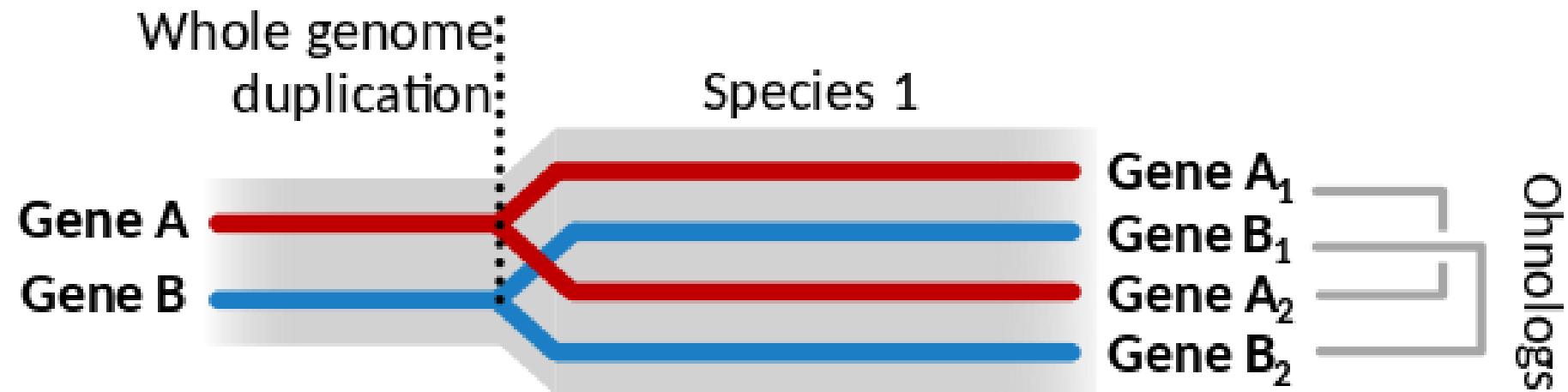
*Where the homology is the result of **speciation** so that the history of the **gene** reflects the history of the **species** (for example alpha hemoglobin in man and mouse) the genes should be called **orthologous** (ortho = exact)."*

**GENE, SPECIATION, DUPLICATION**

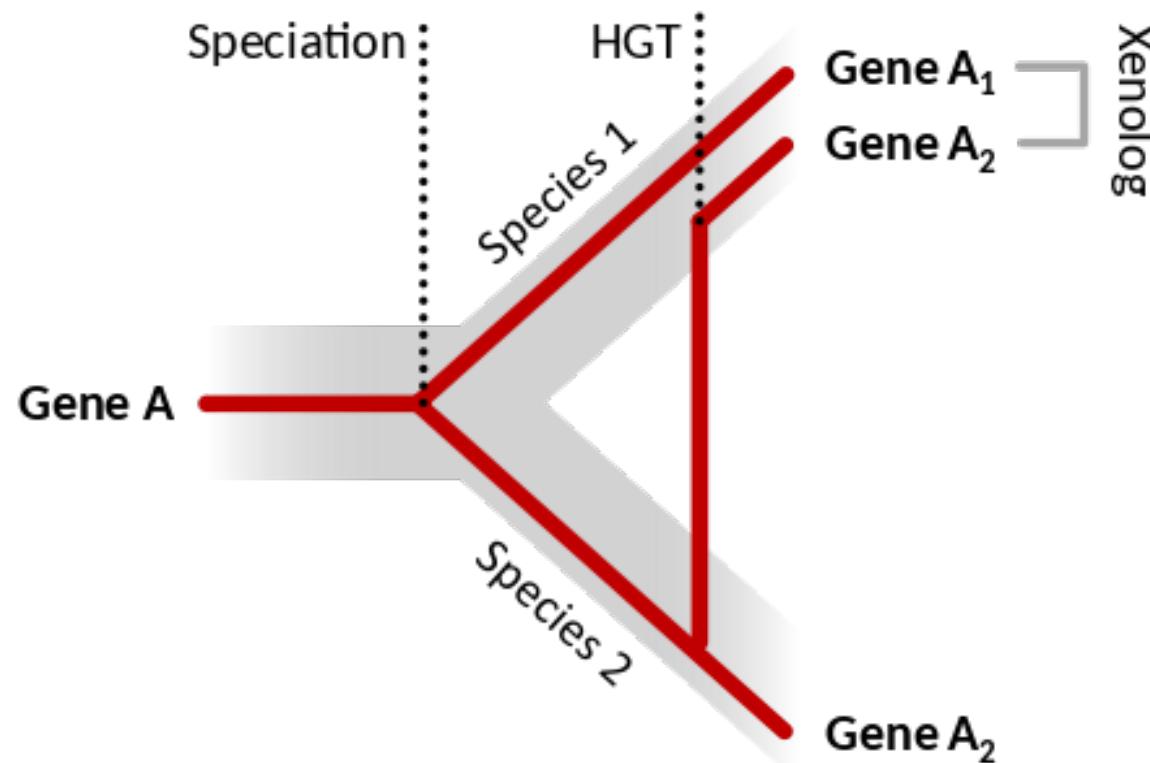
**Can we accomodate orthology to evolutionary processes other than speciation and duplication?**

**Ohnologs, Xenologs, Homeologs**

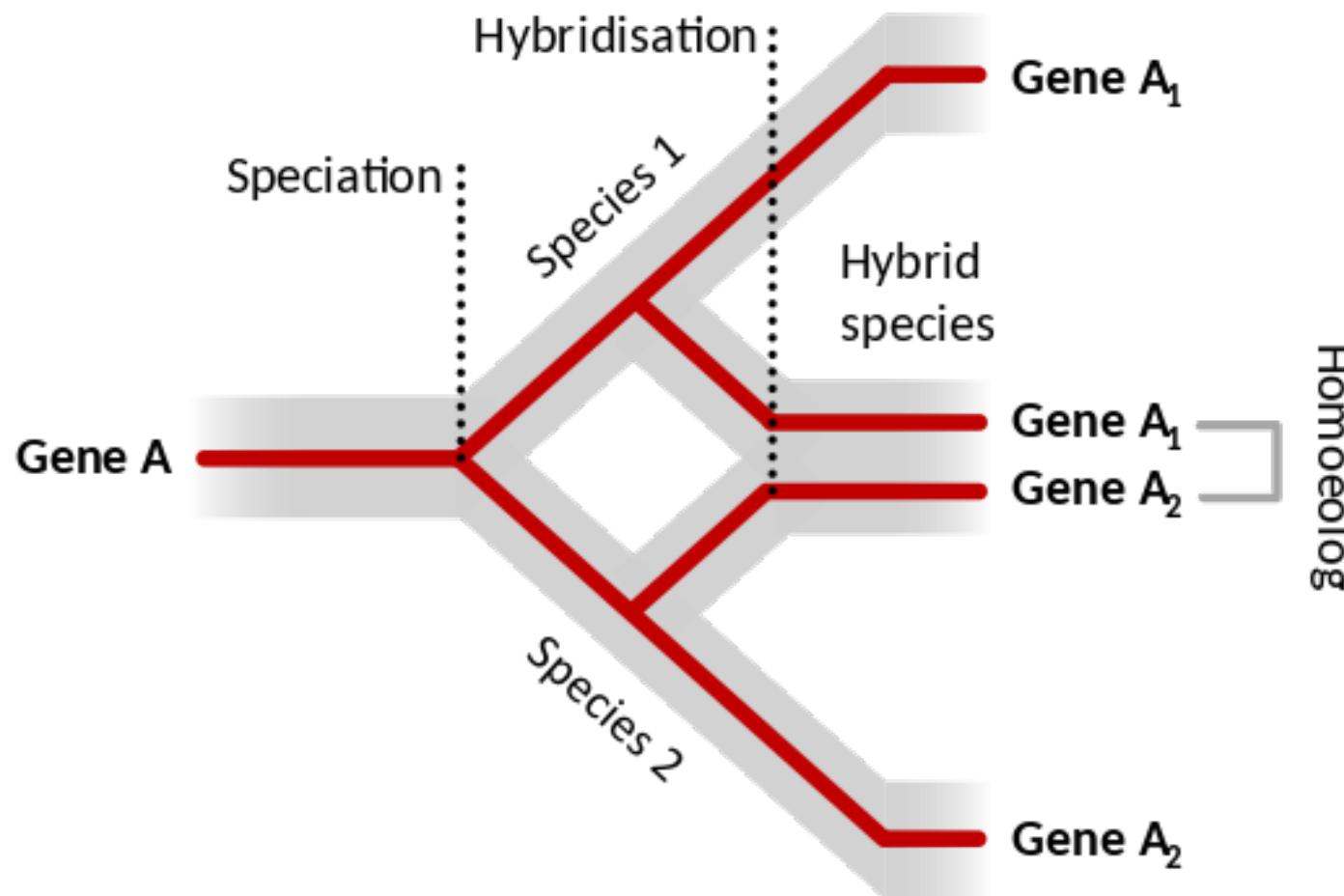
## Ohnologs, Xenologs, and Homeologs



## Ohnologs, Xenologs, and Homeologs



## Ohnologs, Xenologs, and Homeologs



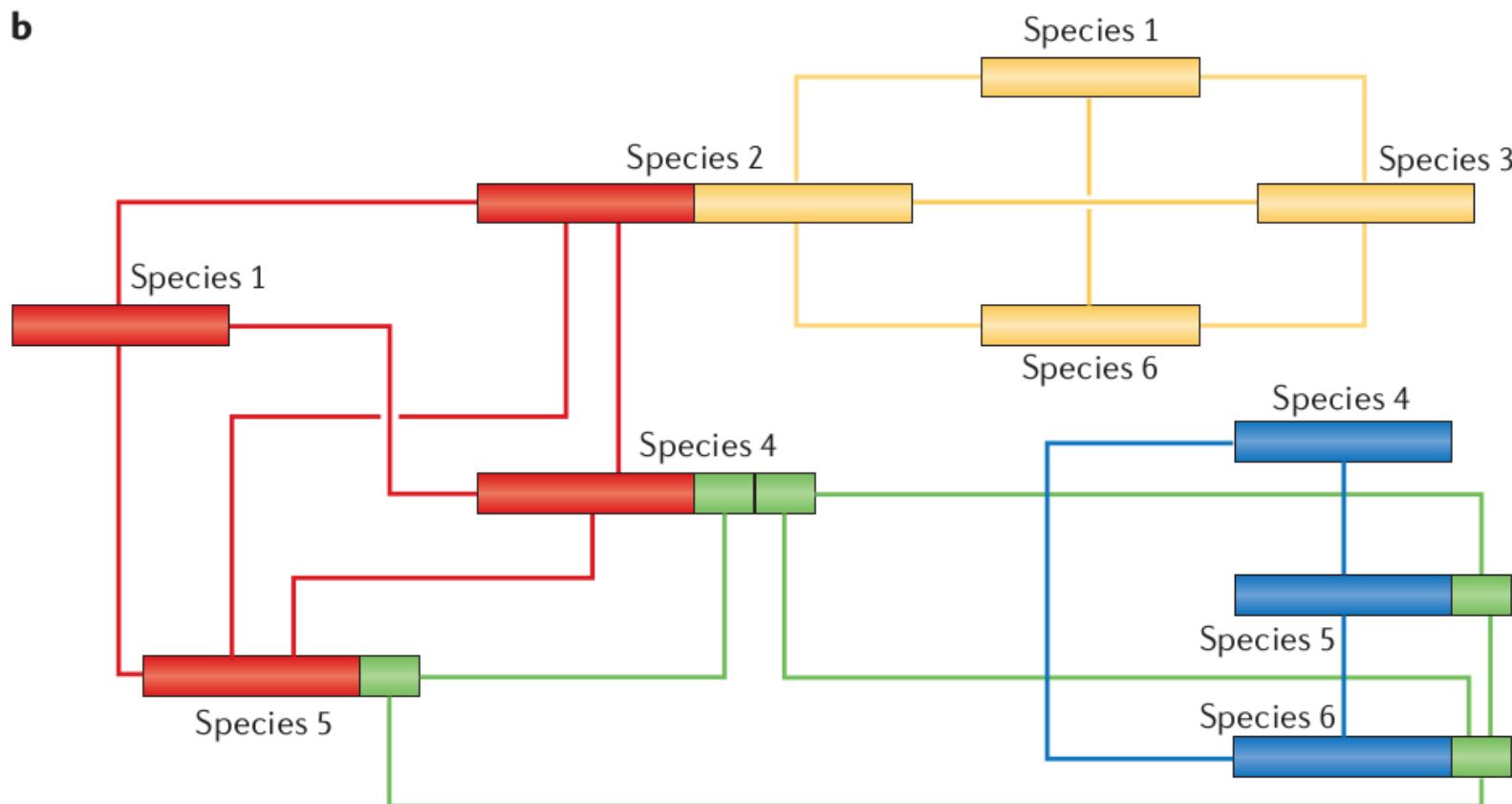
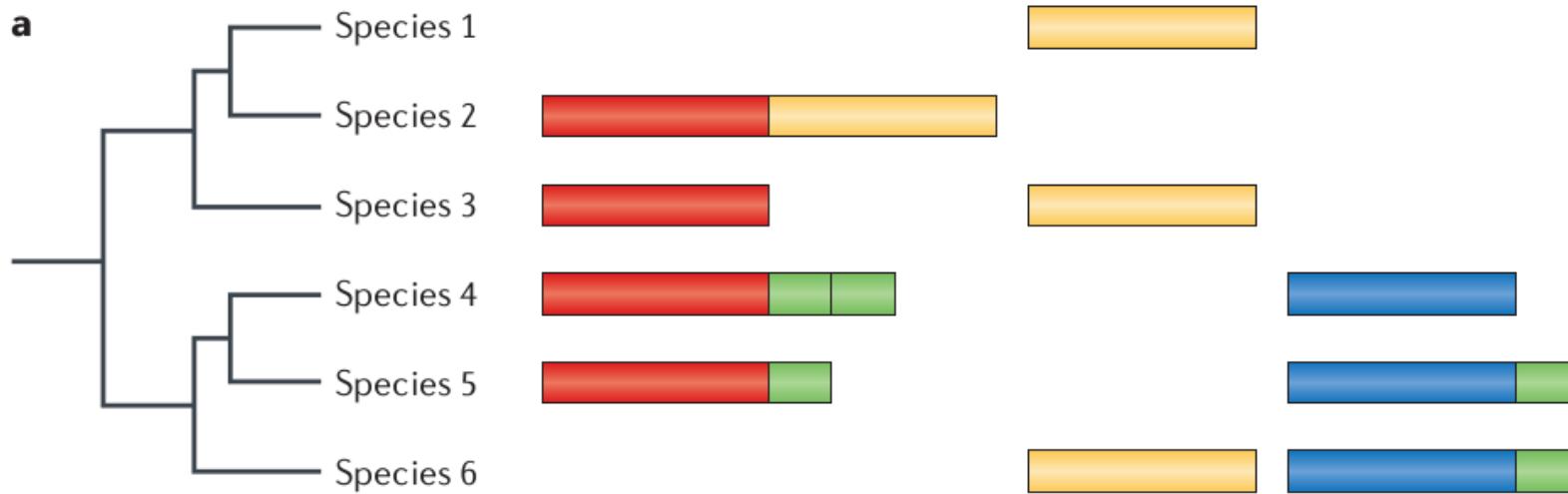
## **Can orthology be defined beyond genes?**

In principle, the concept can simply be extrapolated to any loci that involves through duplication and speciation, but where to set the level of resolution?

**Domains?**

**Single nucleotides?**

## Box 2 | Units of orthology



## **And what about the species boundary?**

Two alleles of the same gene segregating in a population are diverging from each other (they are clearly homologs) but there is no speciation event separating them, they are still “the same gene in the same species”.. however, they can diverge and even change their chromosomal location

This can be particularly problematic in microbial organisms (pan/core genome, rampant gene flow, etc).

# **Questions?**