

Computational Molecular Evolution

by Mátyás Pajkos

Dosztányi Lab

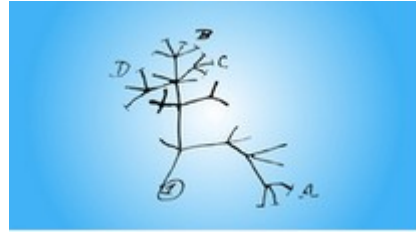
2018 May



Eötvös Loránd
University

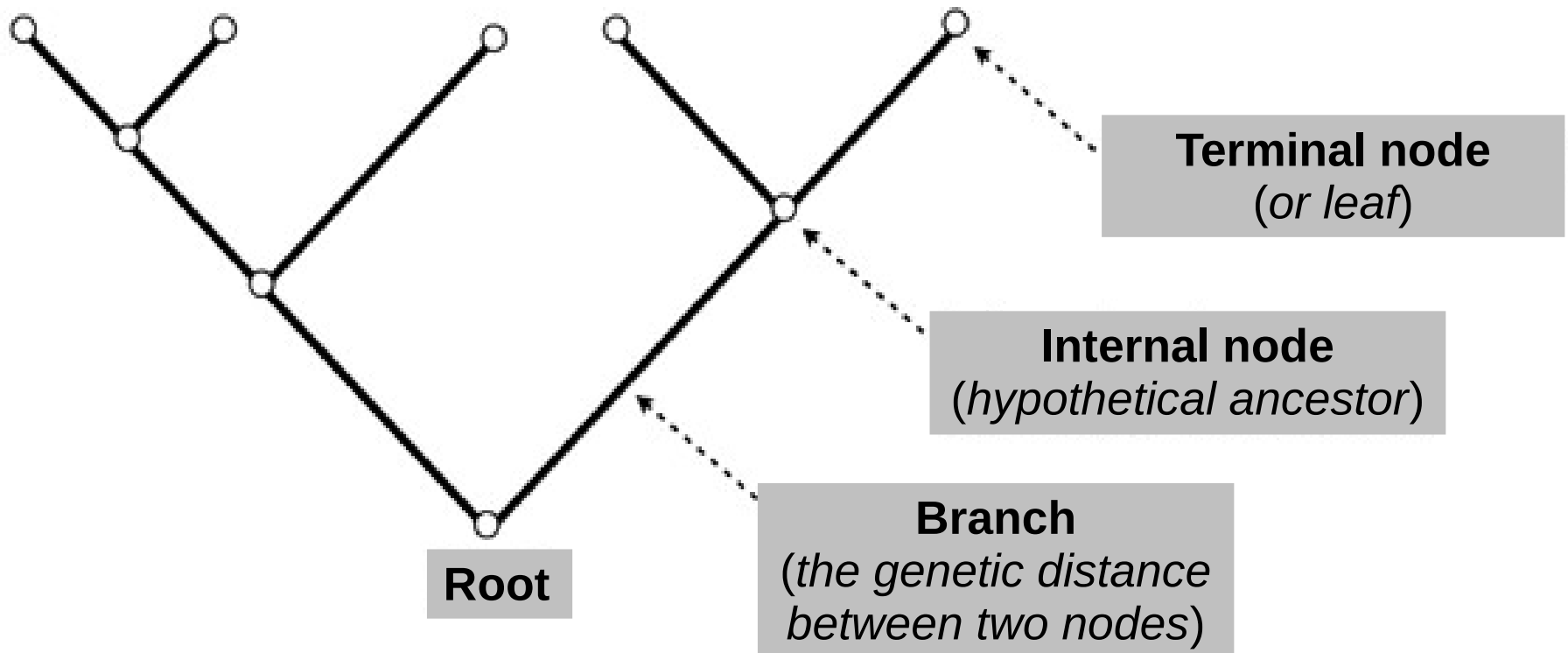
- *Phylogenetic trees: Terminology and representation*
- *Reconstructing trees using present-day data*
- *Orthologous groups*
- *Detection of molecular selection*
- *Tutorial*

Phylogenetic trees: Terminology and representation

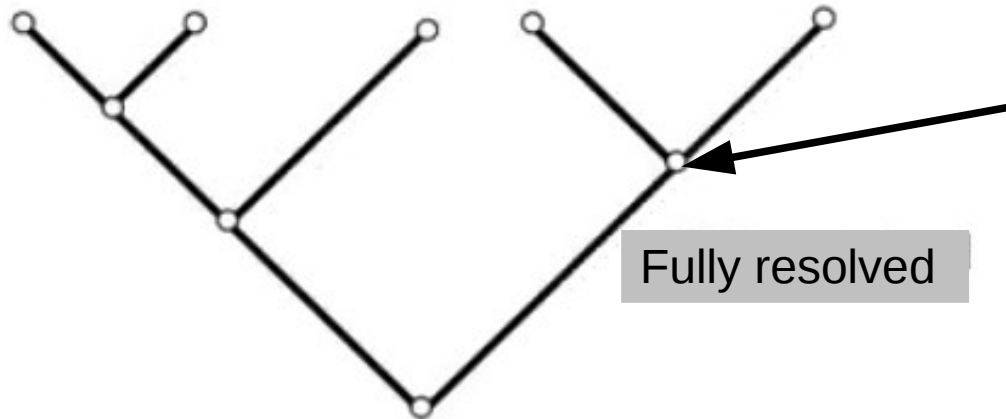


- Terminology
- Representation
- The newick format

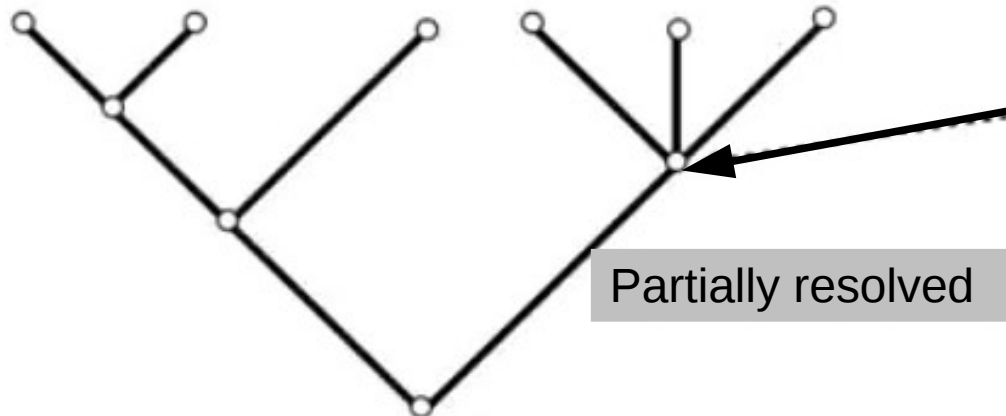
Phylogenetic trees: Terminology



Phylogenetic trees: Representation

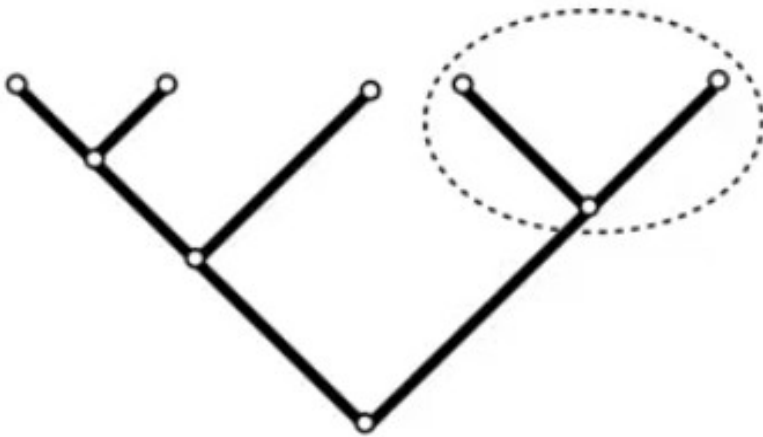


- An internal node has exactly two branches going out
- Reason: A population split into two



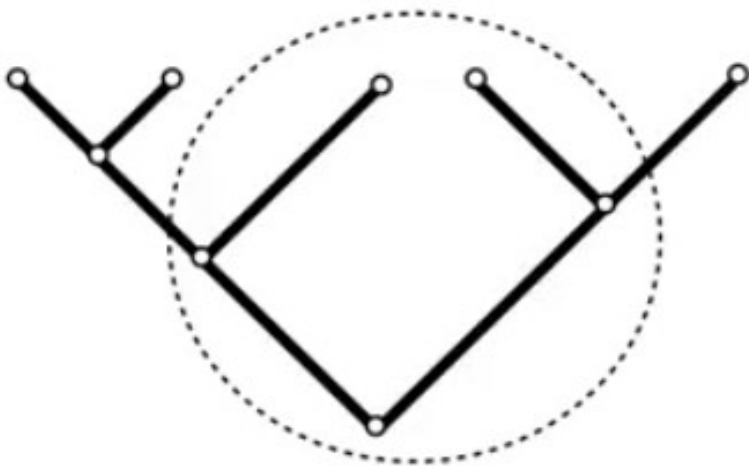
- An internal node has more than two branches going out
- Reason: We do not have enough data

Phylogenetic trees: Representation



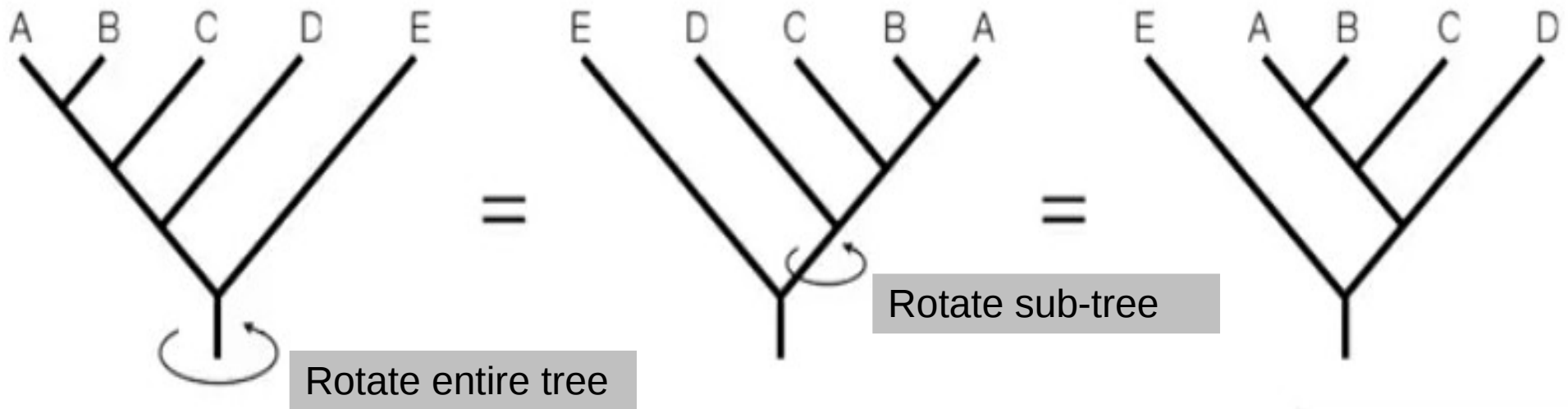
Monophyletic
(clade)

- A clade is a group of organisms that group includes all descendants of their common ancestor
- All the member of such a group, they have several shared features



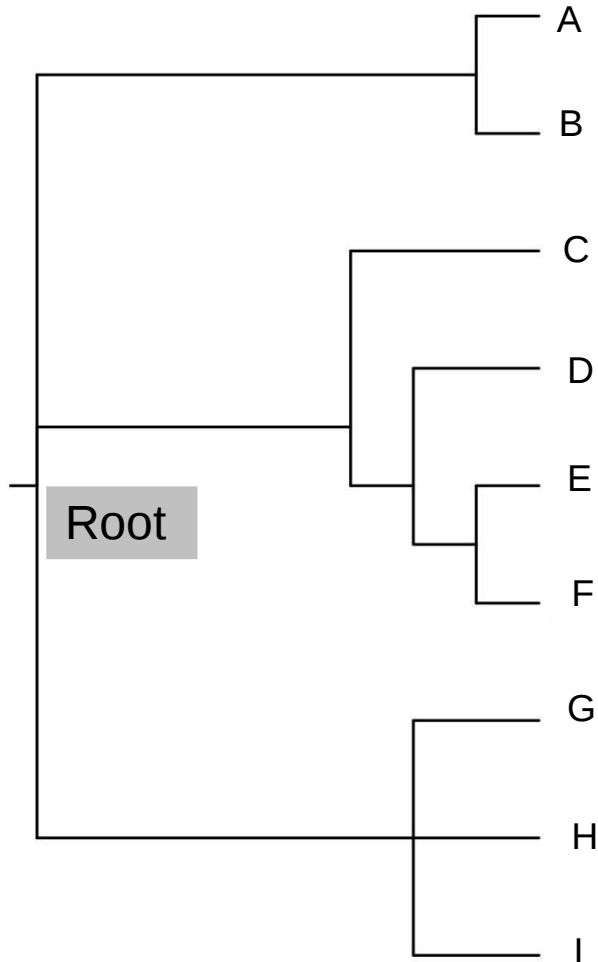
Non-monophyletic
(paraphyletic)

Phylogenetic trees: Representation of topology



- Three different representation of the same tree-topology
- It is not always true that the neighbours are closely related to each other

Phylogenetic trees: rootedness



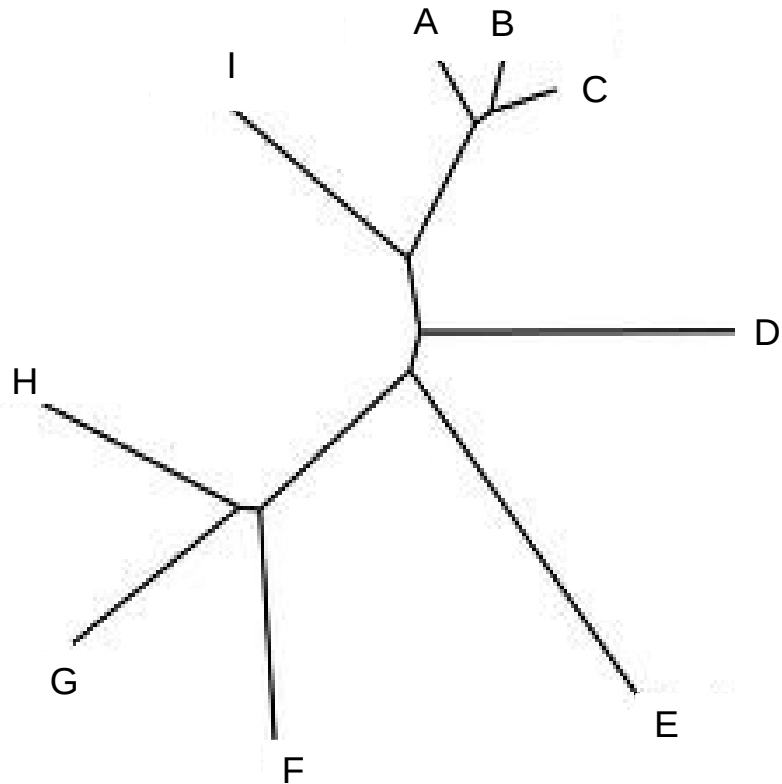
- A rooted tree has a single node (**The root**) that represents a point in time that is earlier than any other node in the tree
- A rooted tree has directionality (nodes can be ordered in terms of "earlier" or "later")
- In the rooted tree, distance between two nodes is represented along the time-axis only (the second axis just helps spread out the leaves)

Old



Young

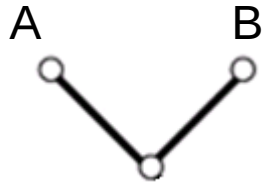
Phylogenetic trees: Unrootedness



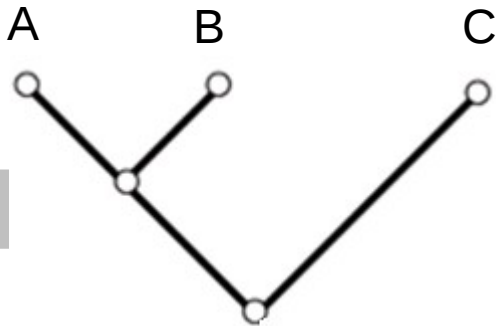
- In unrooted trees there is no directionality: we do not know if a node is younger or older than another node
- Distance along branches directly represents node distance

Phylogenetic trees: The Newick format

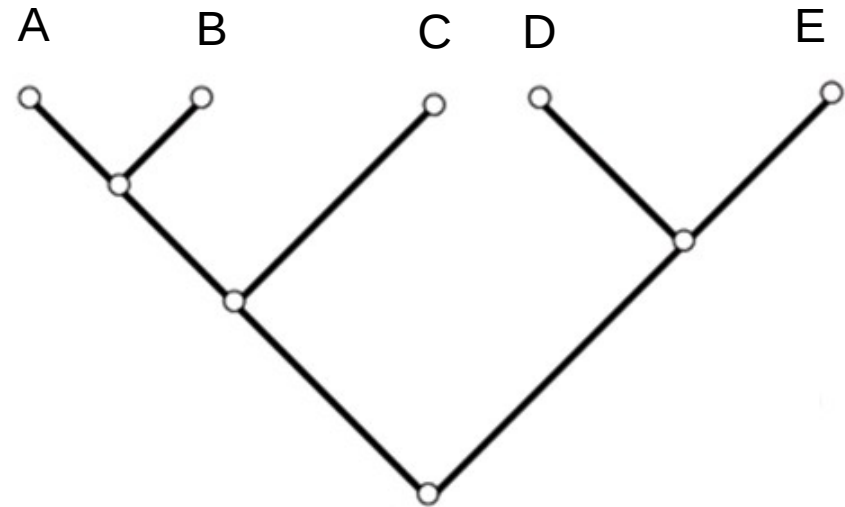
`(A , B);`



`((A , B) , C);`



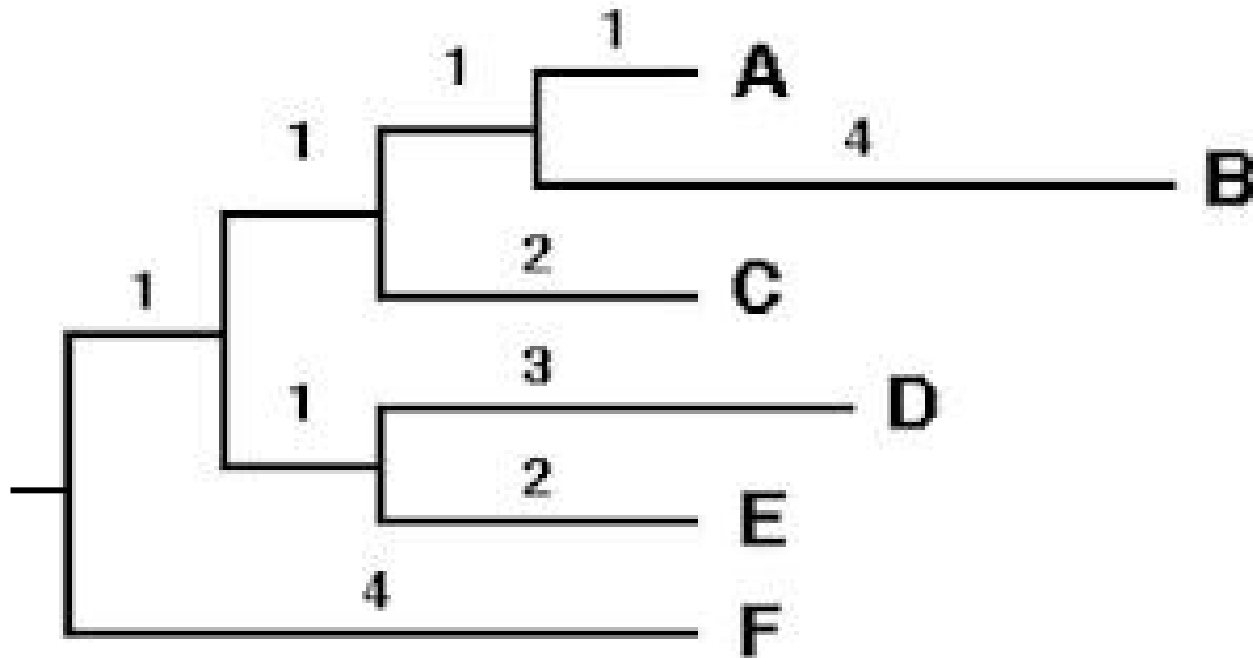
`(((A , B) , C) , (D , E));`



- Standard computer-readable format
- It is based on nested brackets, commas and a terminal semicolon

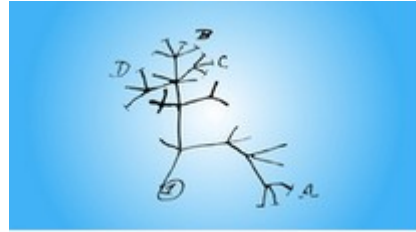
Phylogenetic trees: The Newick format

```
(((( (A:1 , B:4 ):1 , C:2 ):1 ) , (D:3 , E:2 ):1 ):1 , F:4);
```



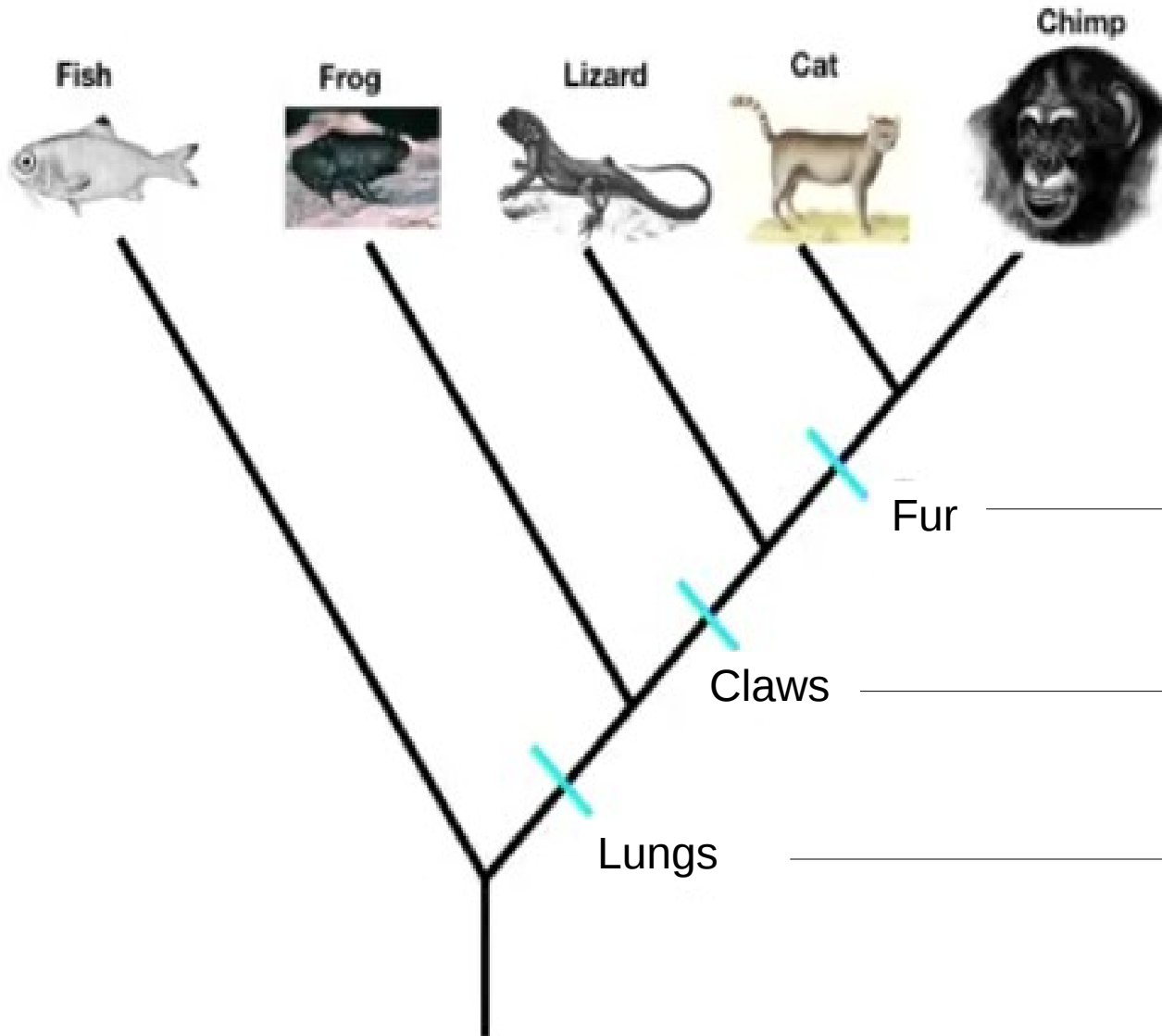
- In the Newick format the branch lengths can be indicated

Reconstructing a tree using present-day data



- Homology
- Homologous alignment characters

Reconstructing a tree using present-day data: The basic idea



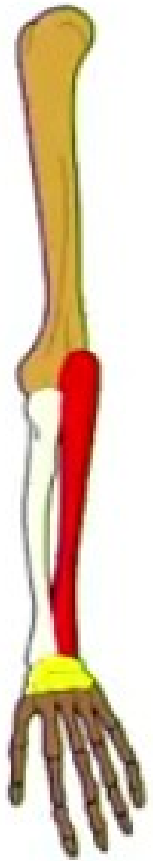
- Trying to group animals that share the largest number of derived traits

- From this point the organisms have lungs, claws and fur

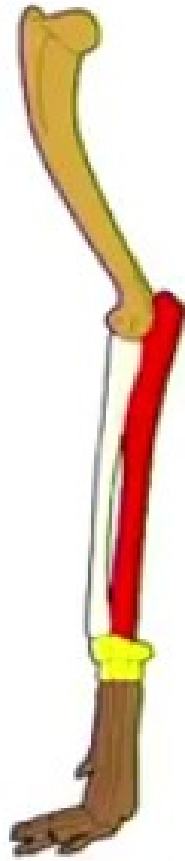
- From this point the organisms have lungs, claws

- From this point the organisms have lungs

Reconstructing a tree using present-day data: Homology



Human



Dog



Bird

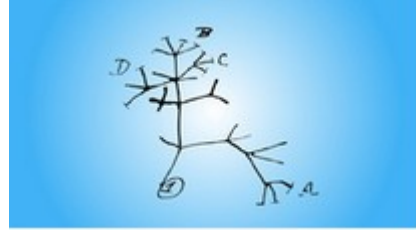
- Homologous trait means the organisms have derived from a common ancestor
- Human, dog, bird have derived from a common ancestor and they all have the same bone structure that is a homologous feature

Reconstructing a tree using present-day data: Molecular phylogenetic

A	A	G	C	G	T	T	G	G	G	C	A	A
B	A	G	C	G	T	T	T	G	G	C	A	A
C	A	G	C	T	T	T	G	T	G	C	A	A
D	A	G	C	T	T	T	T	T	G	C	A	A
				1			2	3				

- Homology means the homologous characters
- Homologous characters mean columns in alignment

Phylogenetic tree building methods



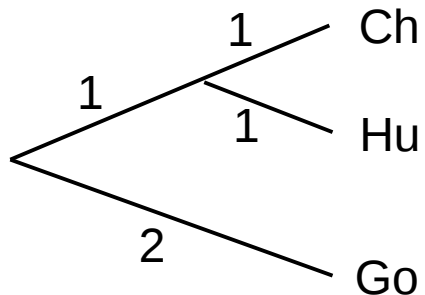
- Distance-based method
- Maximum parsimony method
- Maximum likelihood method

Phylogenetic tree building methods: Distance matrix

Gorilla: ACGT**CGTA**
 Human: ACGTTCCT
 Chimp: ACGT**TCG**

↓ ↓ ↓ ↓
 ↑ ↑

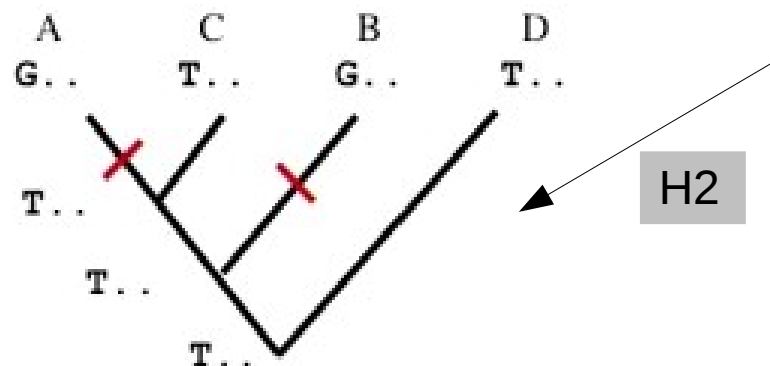
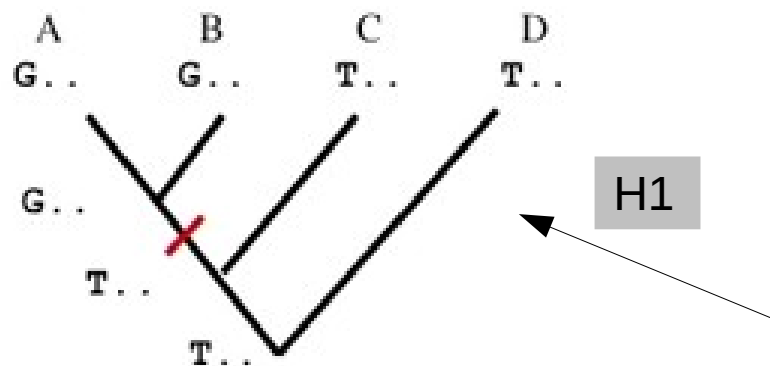
	Go	Hu	Ch
Go	-	4	4
Hu		-	2
Ch			-



- Count the number of substitutions among the sequences
- Write these number in a matrix to get the distance matrix
- According to the matrix the phylogenetic tree can be built

Phylogenetic tree building methods: Maximum Parsimony

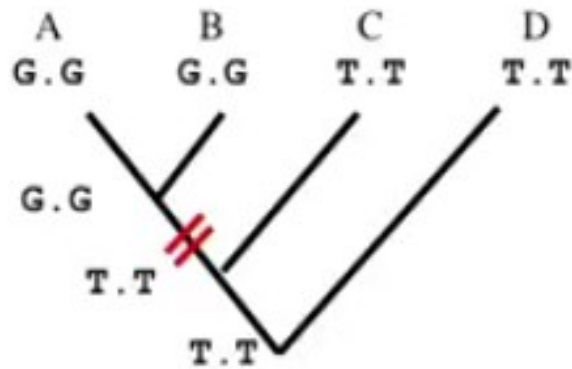
- Maximum parsimony: choose the simplest possible hypothesis



Taxon	Nucleotide position		
	1	2	3
A	G	G	G
B	G	T	G
C	T	G	T
D	T	T	T

- H1 is the simplest possible hypothesis
- The tree has 1 mutation

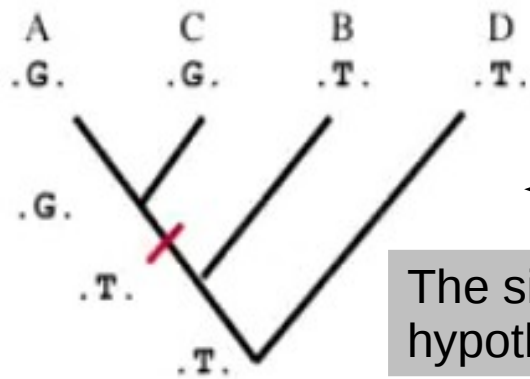
Phylogenetic tree building methods: Maximum Parsimony



- This is the same as the first column
- The tree has 2 mutations

Taxon	Nucleotide position		
	1	2	3
A	G	G	G
B	G	T	G
C	T	G	T
D	T	T	T

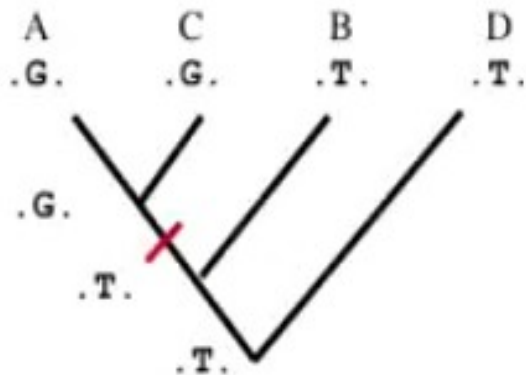
Phylogenetic tree building methods: Maximum Parsimony



The simplest possible hypothesis

		Nucleotide position		
Taxon		1	2	3
→	A	G	G	G
	B	G	T	G
→	C	T	G	T
	D	T	T	T

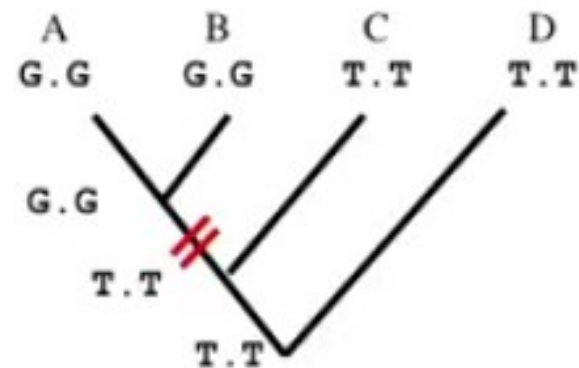
The best hypothesis differs from the others



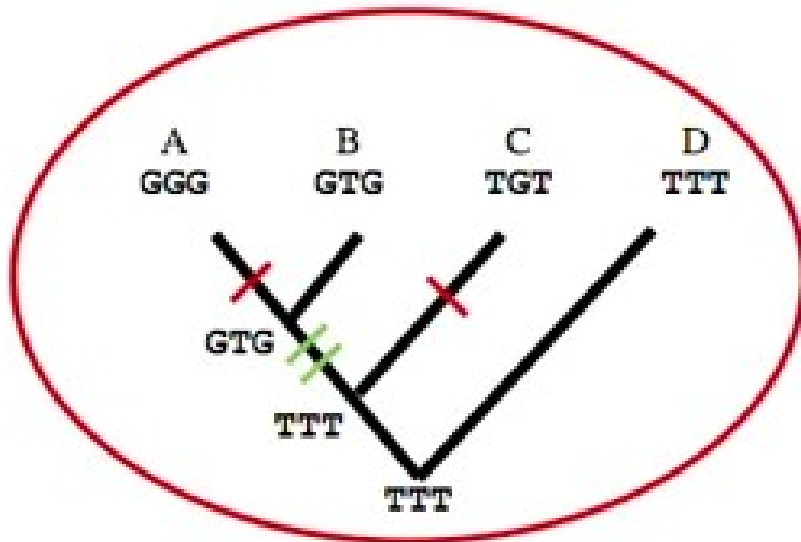
?



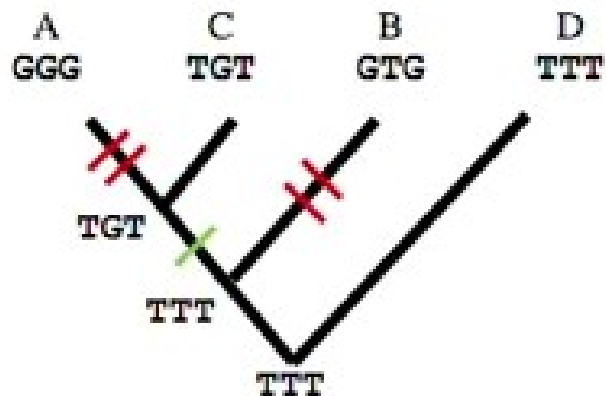
Conflict:



Phylogenetic tree building methods: Maximum Parsimony



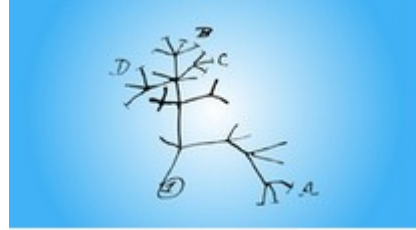
Total length of tree: 4



Total length of tree: 5

- The best tree: the smallest number of the mutations
- Count the total number of the mutations for the two versions
- Compare them and choose the smaller
- In this case we have to reject the best hypothesis at position 2 in order to get the best tree

Evolutionary implications of gene orthology

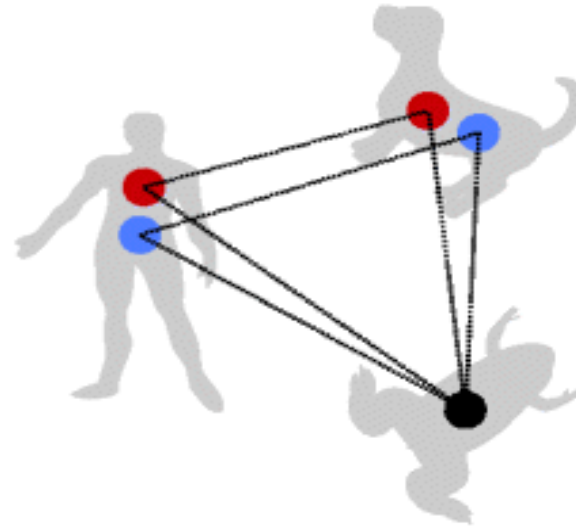


- Orthologous groups
- Prediction

Orthologous Group

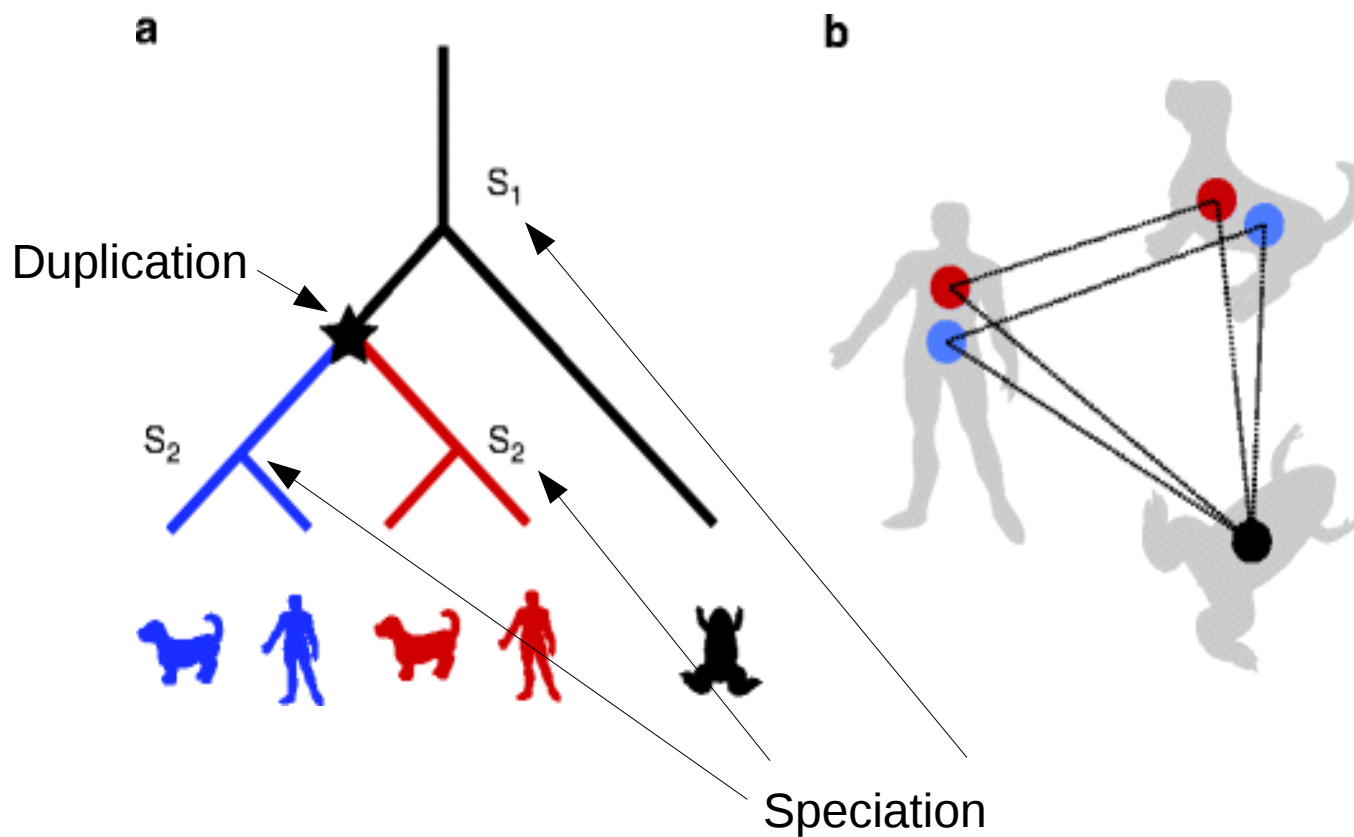
Hierarchical orthologous groups are defined as sets of genes that have descended from a single common ancestor within a taxonomic range of interest

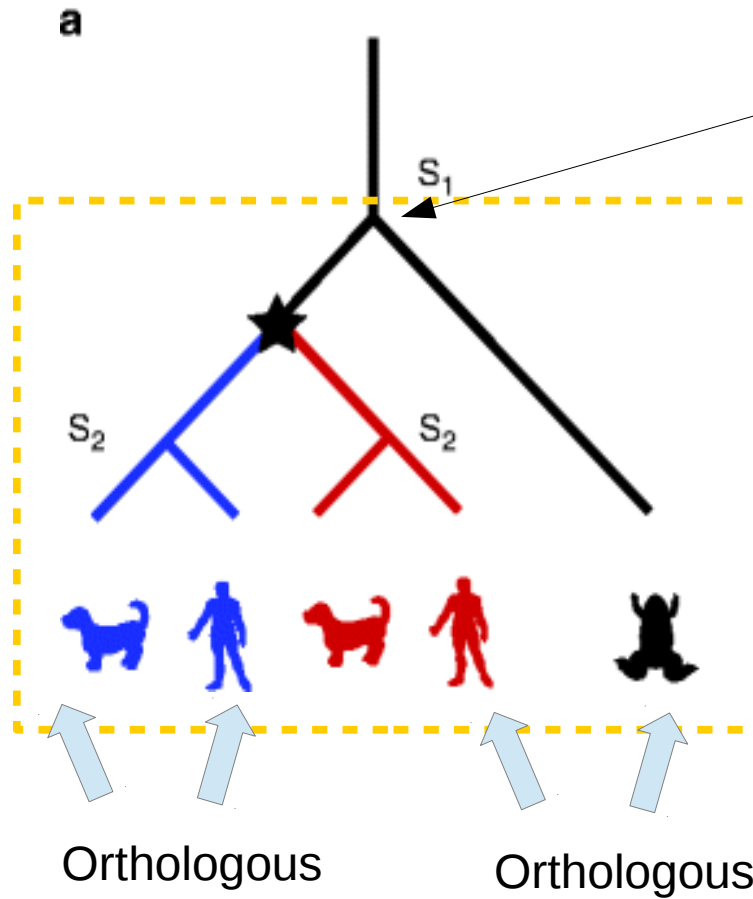
- A gene in human, frog and dog
- In human there is one copy
- But in dog and frog there are two copies



What's going on here?

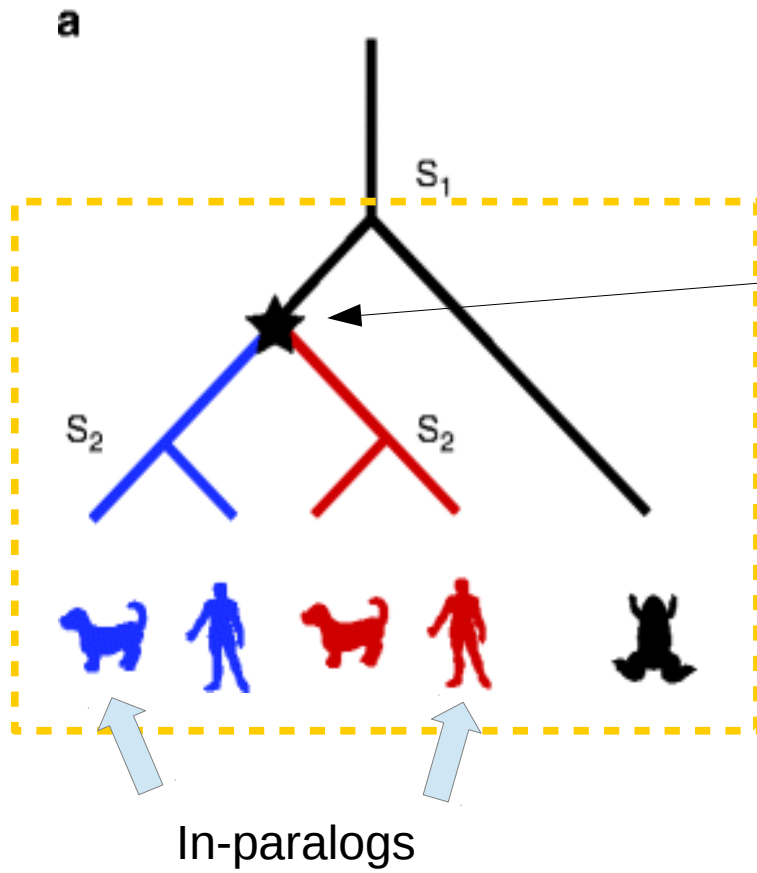
Look at the history of these five genes which is depicted in a phylogenetic tree





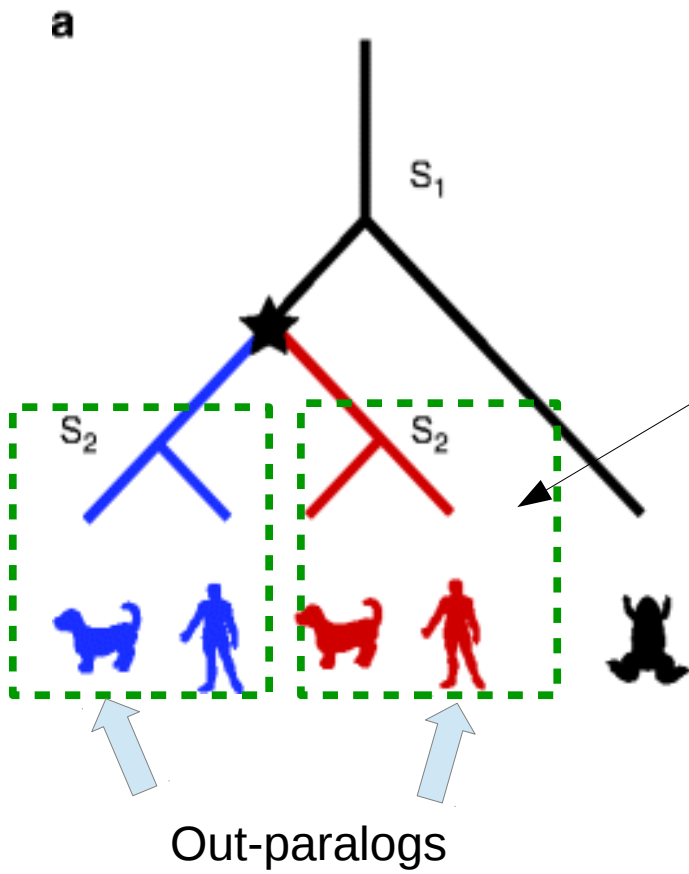
- In ancestral vertebrata: one gene
- All genes are derived from it
- This clade includes the five genes

Orthologous: Genes related by speciation



The duplication must have happened within the clade in question

In-paralogs: Genes related by duplications



- Really important to define the clade that is taxonomic level
- They started diverging at a duplications that happened before the dog-human speciation
- Different clades relative to this level

Out-paralogs: Genes related by duplications

Transferrin gén

Drosophila melanogaster
Drosophila melanogaster

Drosophila melanogaster

Ciona intestinalis

Danio rerio

Takifugu rubripes

Xenopus tropicalis

Gallus gallus

Gerincesek
A3

Takifugu rubripes

Danio rerio

Xenopus tropicalis

Gallus gallus

Monodelphis domestica

Homo sapiens

Mus musculus

Canis familiaris

A2.2

Canis familiaris

Mus musculus

Homo sapiens

A2.1

Gerincesek
A2

Canis familiaris

Homo sapiens

Mus musculus

A1.1

Monodelphis domestica

Gallus gallus

Xenopus tropicalis

Takifugu rubripes

Danio rerio

Gerincesek
A1

Gerinchúrosok
A

- Gén „duplication” (●) és „speciation” (□) sorozatok különböző evolúciós időben több, elkülöníthető csoportot, **referencia pontot** alakít ki

- Ortológok:** „speciation node”

Referencia pont: A3 - Gerincesek

A3 – *Danio r.* & *Gallus g.*

Referencia pont: A - Gerinchúrosok

A – *Ciona m.* & A3 – *Danio r.*

- Co – Ortológok:**

Referencia pont: A2 - Gerincesek

„one-to-many”

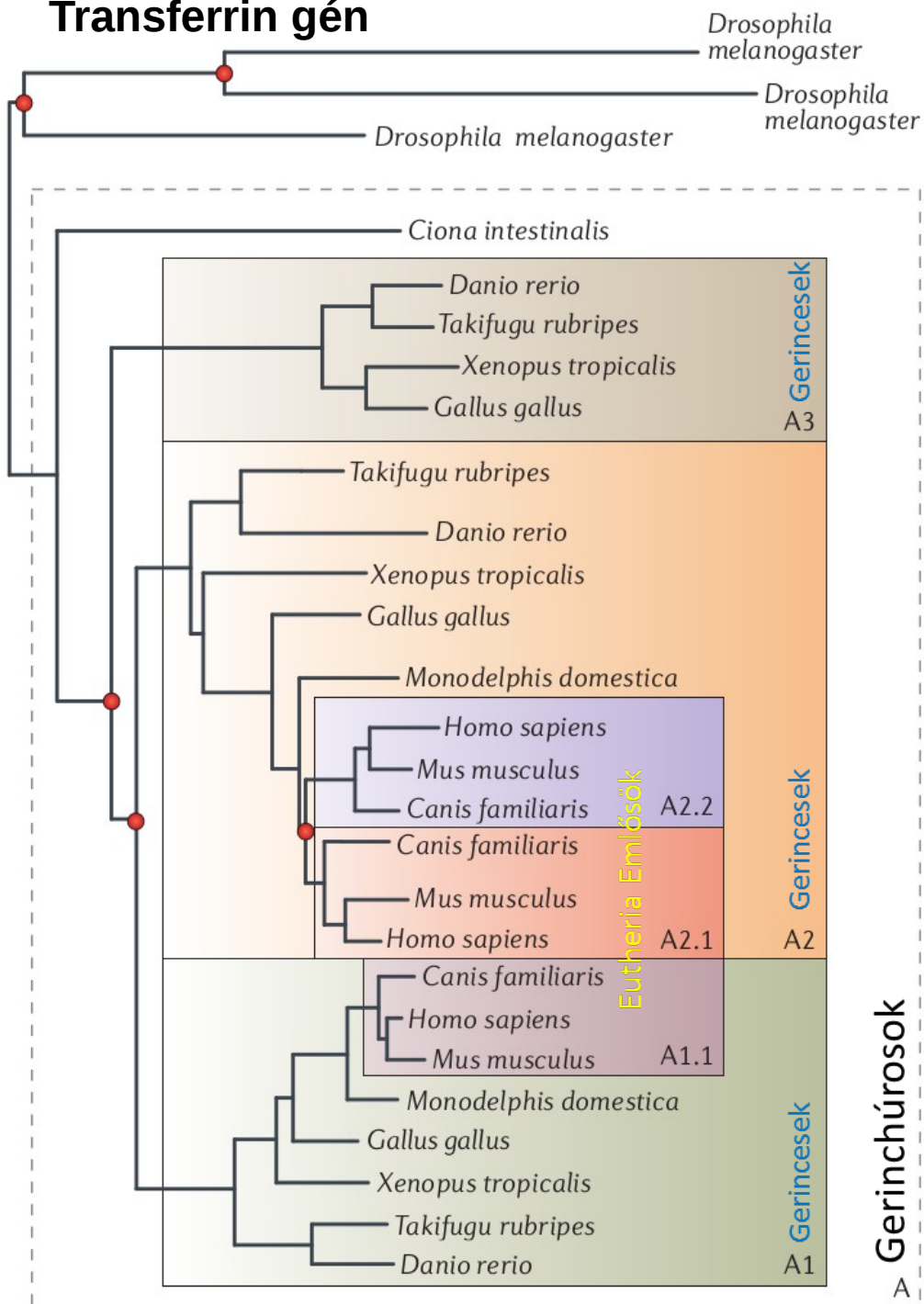
A2 – *Danio r.* & A2.1, A2.2 *Homo s.*

Referencia pont: vizsgált tartomány

„many-to-many”

Drosophila m. & *Homo s.*

Transferrin gén



- Gén „duplication” (●) és „speciation” (□) sorozatok különböző evolúciós időben több, elkülöníthető csoportot, **referencia pontot** alakít ki

- Paralógok:** „duplication node”
Referencia pont: A - Gerinchúrosok
A3 – Dani r. & A1 – Danio r.
A3 – Gallus g. & A2 – Gallus g.

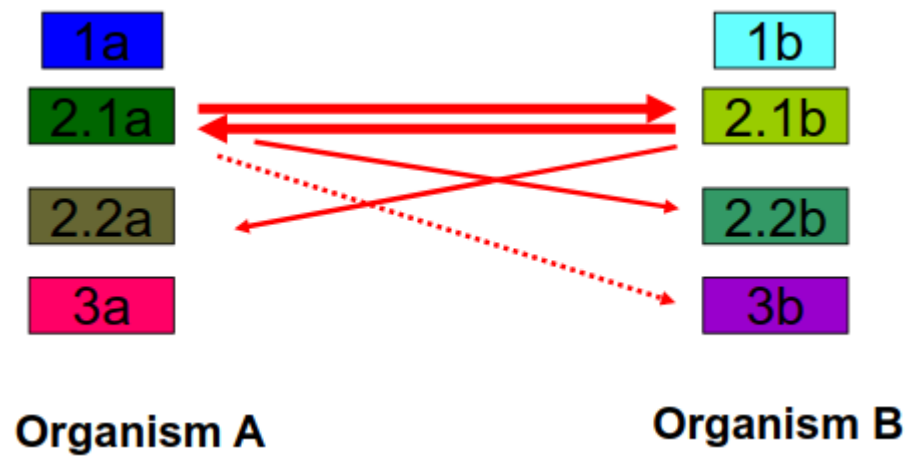
Referencia pont: A2 - Gerincesek
„In – paralogues”
A2.1 – Homo s. & A2.2 – Homo s.

Referencia pont: A1.1-A2.1 - Gerincesek
„Out - paralogues”
A2.1 – Homo s. & A1.1 – Homo s.

Referencia pont: A - Gerinchúrosok
„In – paralogues”
A2.1, A2.2, A1.1 Homo s.

Gerinchúrosok
A

- How to detect orthologous genes?
 - Easy way: **Best Reciprocal Hit** (RBH)



Detection of positive/negative molecular selection



- Substitutions
- Detection of molecular selection
- The levels of detection

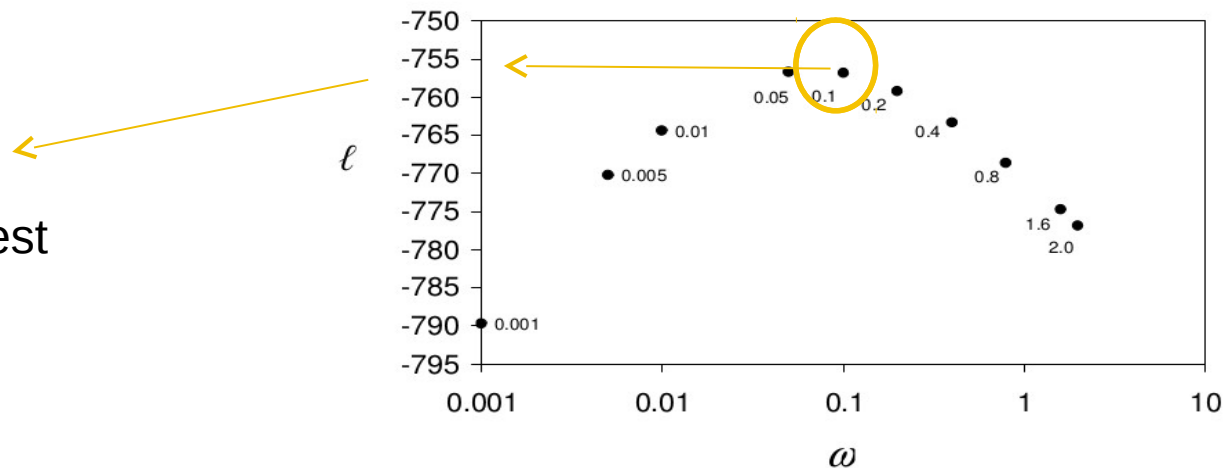
Detection of positive/negative molecular selection: substitutions

- The main parameter: ω (omega)

$$\omega = dN/dS$$

- **dN**: rate of non-synonymous substitutions
- **dS**: rate of synonymous substitutions

Select the smallest likelihood value



- $\omega = 1$ → neutral selection
- $\omega < 1$ → negative selection
- $\omega > 1$ → positive selection

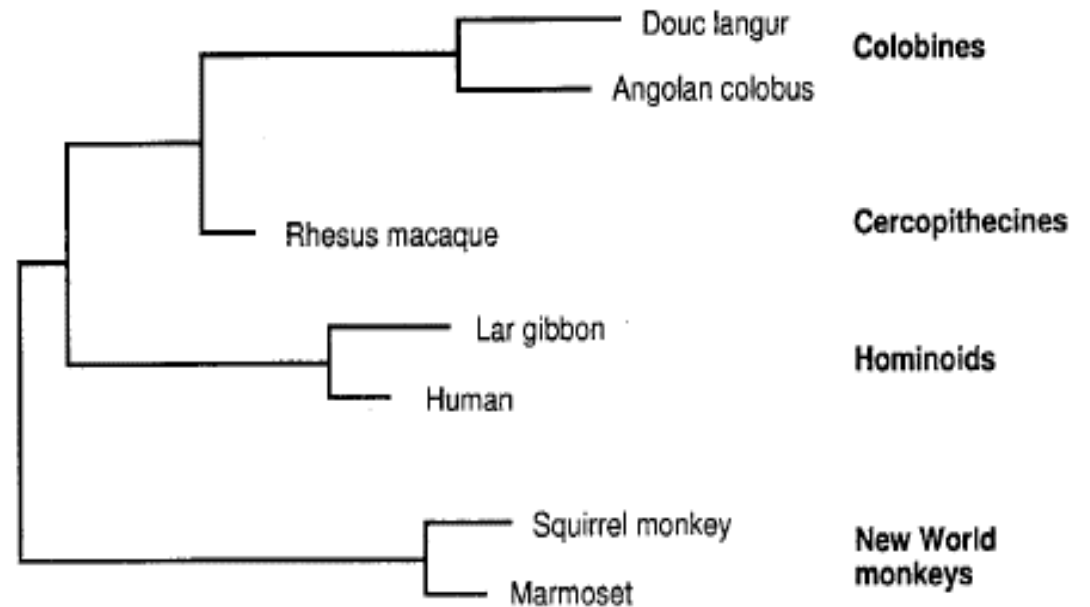
Detection of positive molecular selection: Models

- **Null-model:**

- Global average omega value

- It describes the evolution of the entire tree

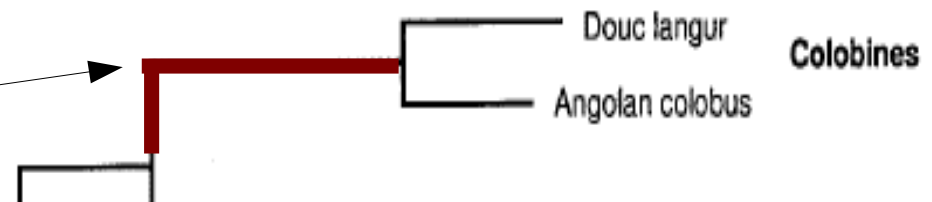
$\omega = 0.34$ → The given gene spent the overwhelmed majority of time under negative selection



- **Branch-model:**

- Partial omega value
- Global omega value

- It describes the evolution of a given branch

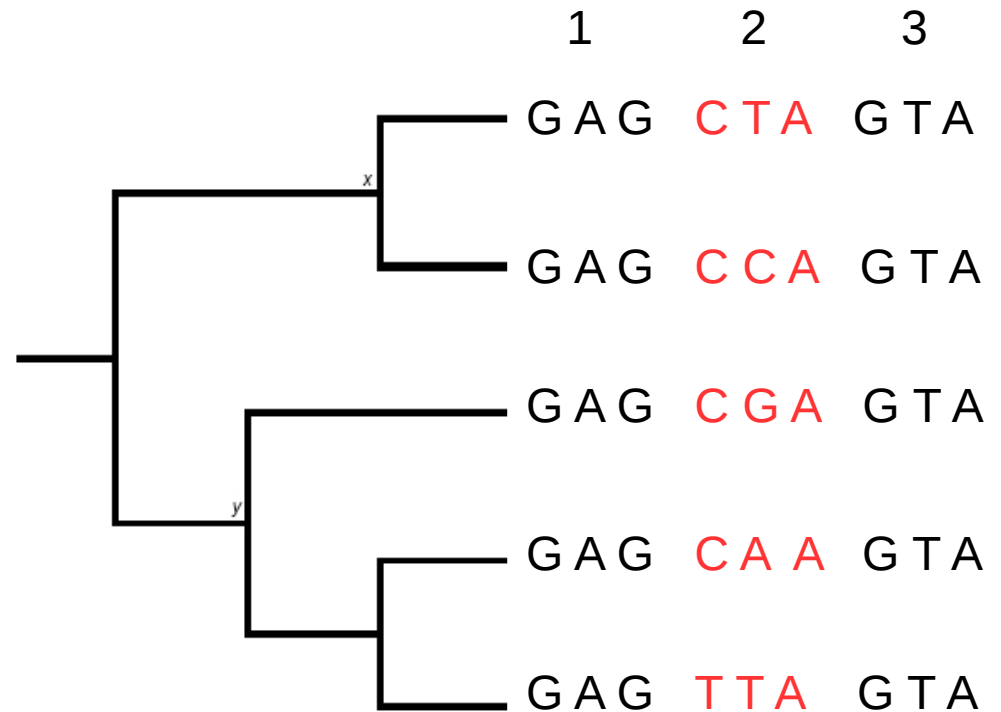


$\omega = 2.21$ → Positive selection can be observed on this branch

Detection of positive molecular selection: Models

- **Site-model:**

- Global average omega value for each codon independently
- It describes the evolution of each codon



$\omega_1 = 1$ → There is not any selection

$\omega_2 = 2.81$ → There is positive selection

$\omega_3 = 1$ → There is not any selection

Detection of positive molecular selection: Likelihood Ratio Test

- We have to declare hypotheses to calculate some kind of statistics

Hypothesis 1:

- $\omega < 1$ or $\omega > 1$
- Likelihood value 1

Hypothesis 0:

- $\omega = 1$
- Likelihood value 2

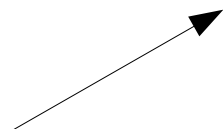
- Using the two likelihood values we can decide whether the selection is statistically significant or not

Computation 1:

$\omega = 2.3$
 $L = -745$

Computation 2:

$\omega = 1$ (fixed)
 $L = -973$

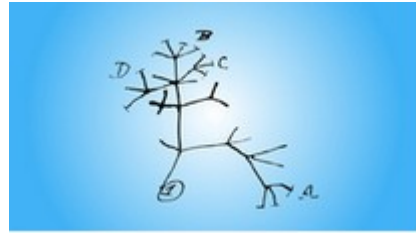


Likelihood Ratio
Test



P-value: 0.00034

Tutorial



- Prepare and view trees in FigTree viewer
- Prepare distance matrix
- Computational molecular evolution

Please download the files below:

- <http://dlab.elte.hu/index.php/education/>
 - MolEvo_Tutorial_2018_byPajkos.pdf
 - Codeml.zip