



How to read and make phylogenetic trees
Zuzana Starostová

How to make phylogenetic trees?

Workflow:

- ✓ obtain DNA sequence
 - quality check
 - sequence alignment
 - calculating genetic distances (optional)
 - phylogeny estimation – topology and branch length
 - reliability test (bootstrap)
 - tree visualization

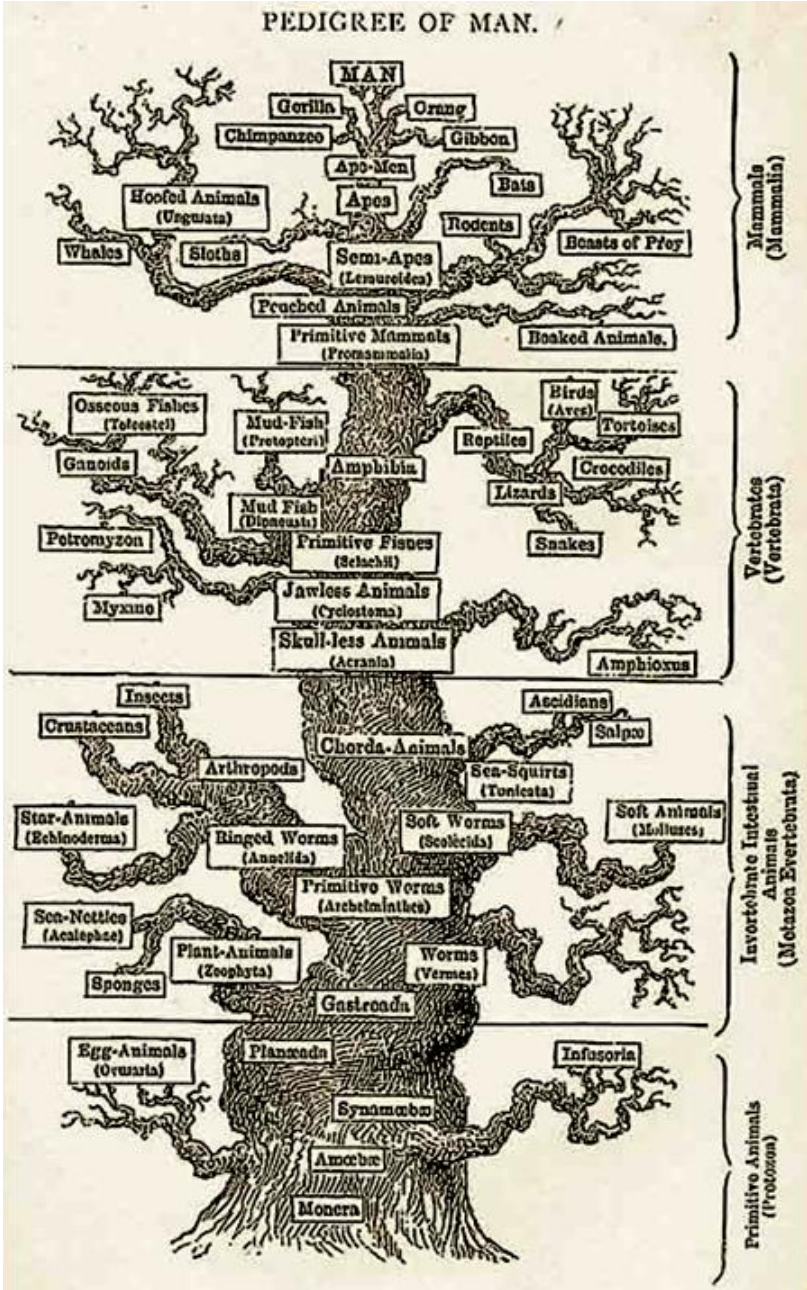
- the process of evolution produces a pattern of relationships between species - as lineages evolve and split and modifications are inherited, their evolutionary paths diverge

- this produces a branching pattern of evolutionary relationships – **phylogenetic tree**

- **phylogeny** – the evolutionary history of a species or group of related species

- phylogenies trace patterns of shared ancestry between lineages

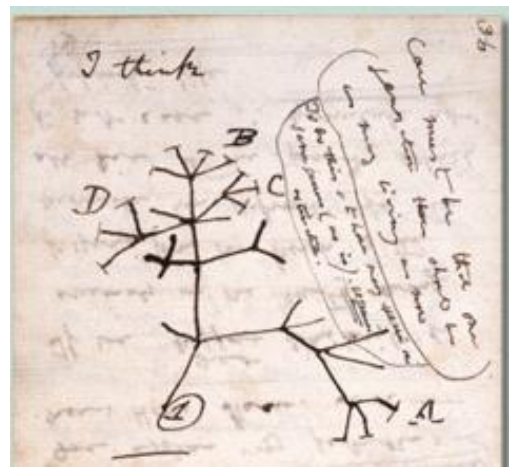
- each lineage has a part of its history that is unique to it alone and parts that are shared with other lineages



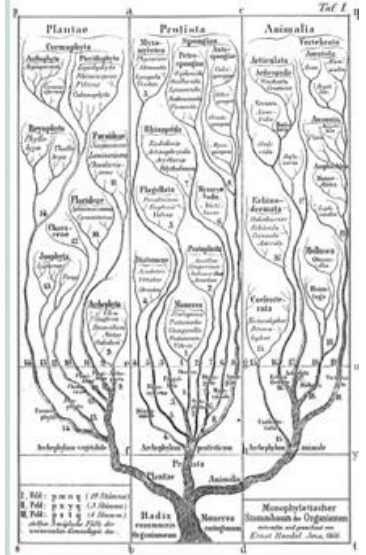
Tree of life



Charles Darwin
species share common ancestor
relationships among species compared to
“the great Tree of Life“

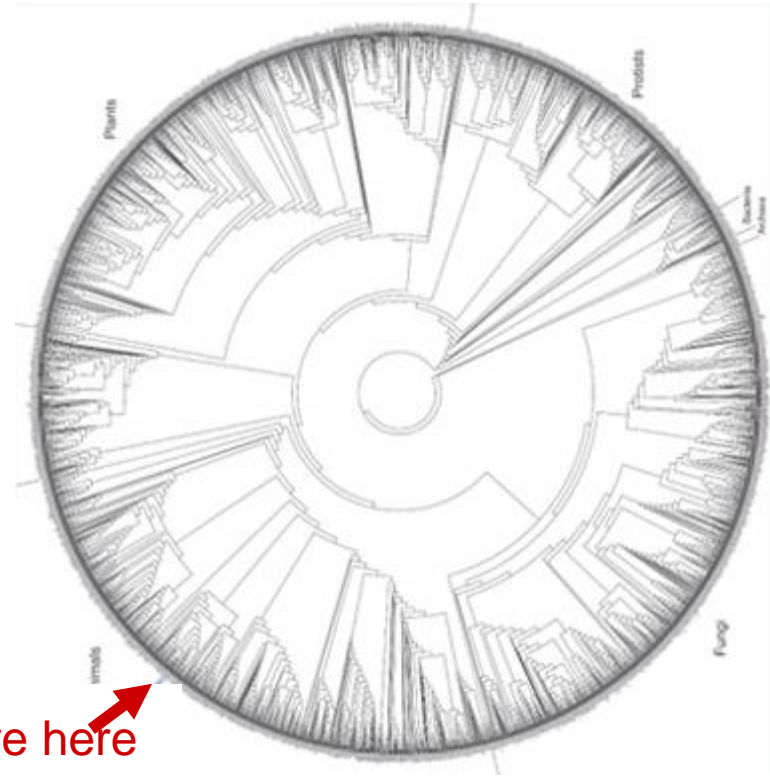
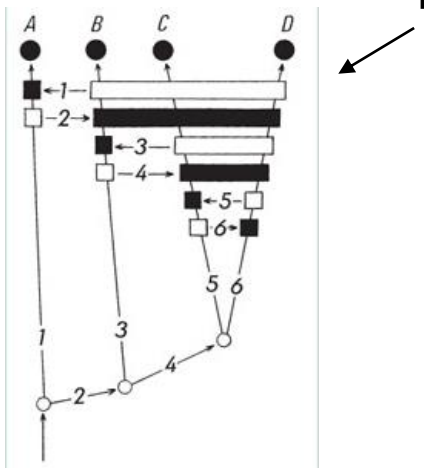


Darwin's diary, 1837



Ernst Haeckel – 1866
diagram of relationships among species
based on general similarity

Willi Hennig – 1960s
foundation of modern
phylogenetics



you are here

Why is phylogeny important?

- understanding and classifying the diversity of life on Earth

We use phylogenetic trees for:

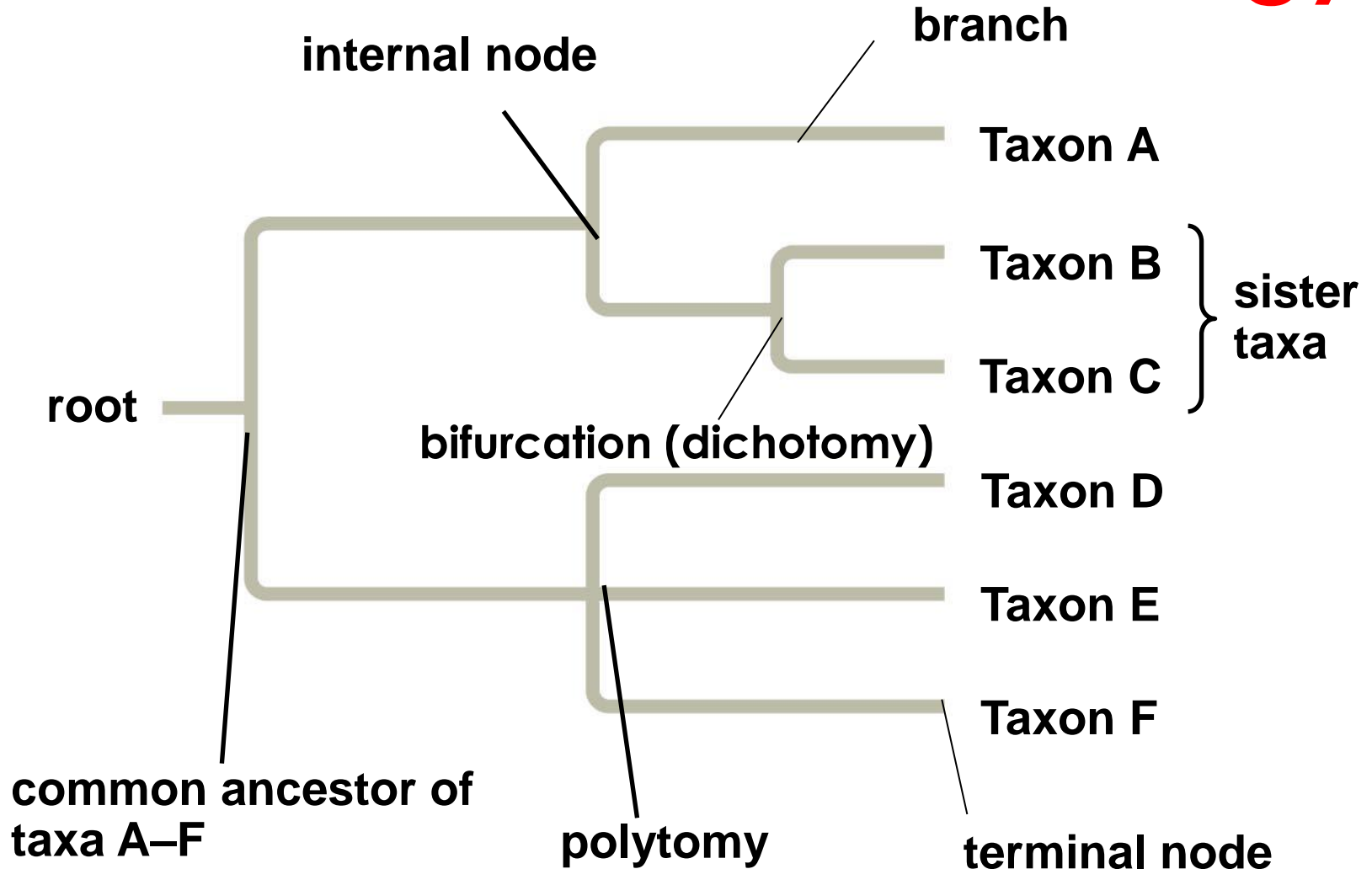
- comparative analysis and character evolution
- biogeography
- dating – age of different taxa
- genetic engineering
- disease epidemiology
- conservation
- ...



The Tree of Life Web Project (ToL) is a collaborative effort of biologists and nature enthusiasts from around the world. On more than 10,000 World Wide Web pages, the project provides information about biodiversity, the characteristics of different groups of organisms, and their evolutionary history.

<http://tolweb.org/tree/>

how to read trees - terminology

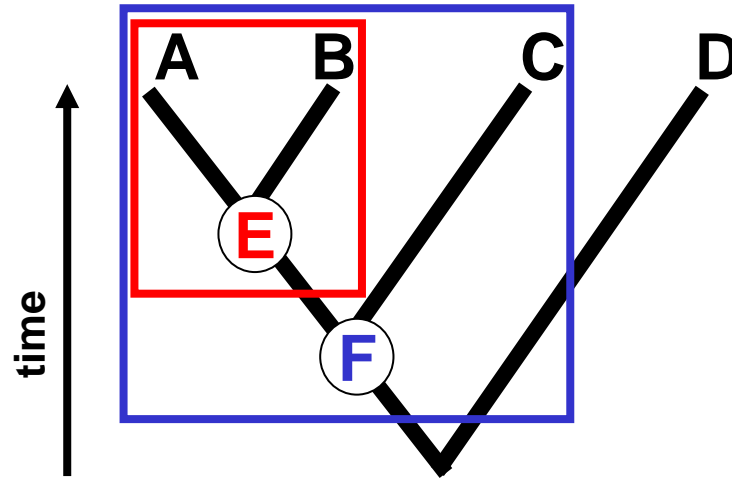


- root of the tree represents the ancestral lineage, and the tips of the branches represent the descendants of that ancestor
- as one moves from the root to the tips = moving forward in time

phylogenetic tree

branching diagram showing relationships between taxa based on their shared common ancestors

taxa (e.g.: species):



A and B are more closely related because they share a common ancestor (here “E”) that C and D do not share

A+B+C are more closely related to each other than to D because they share a common ancestor (“F”) that D does not share

Phylogeny and classification

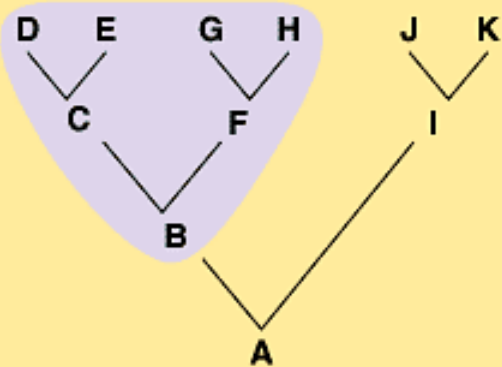
phylogenetic (cladistic) classification reflects evolutionary history

the only valid group for classification is **monophyletic**

group = clade = group that includes a common ancestor and all the descendants (living and extinct) of that ancestor.

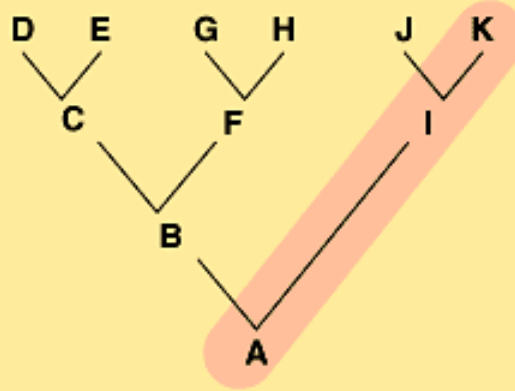


Taxon 1
(monophyletic)



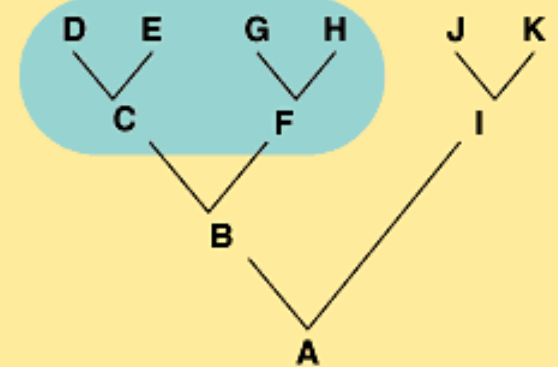
(a) Monophyletic

Taxon 2
(paraphyletic)



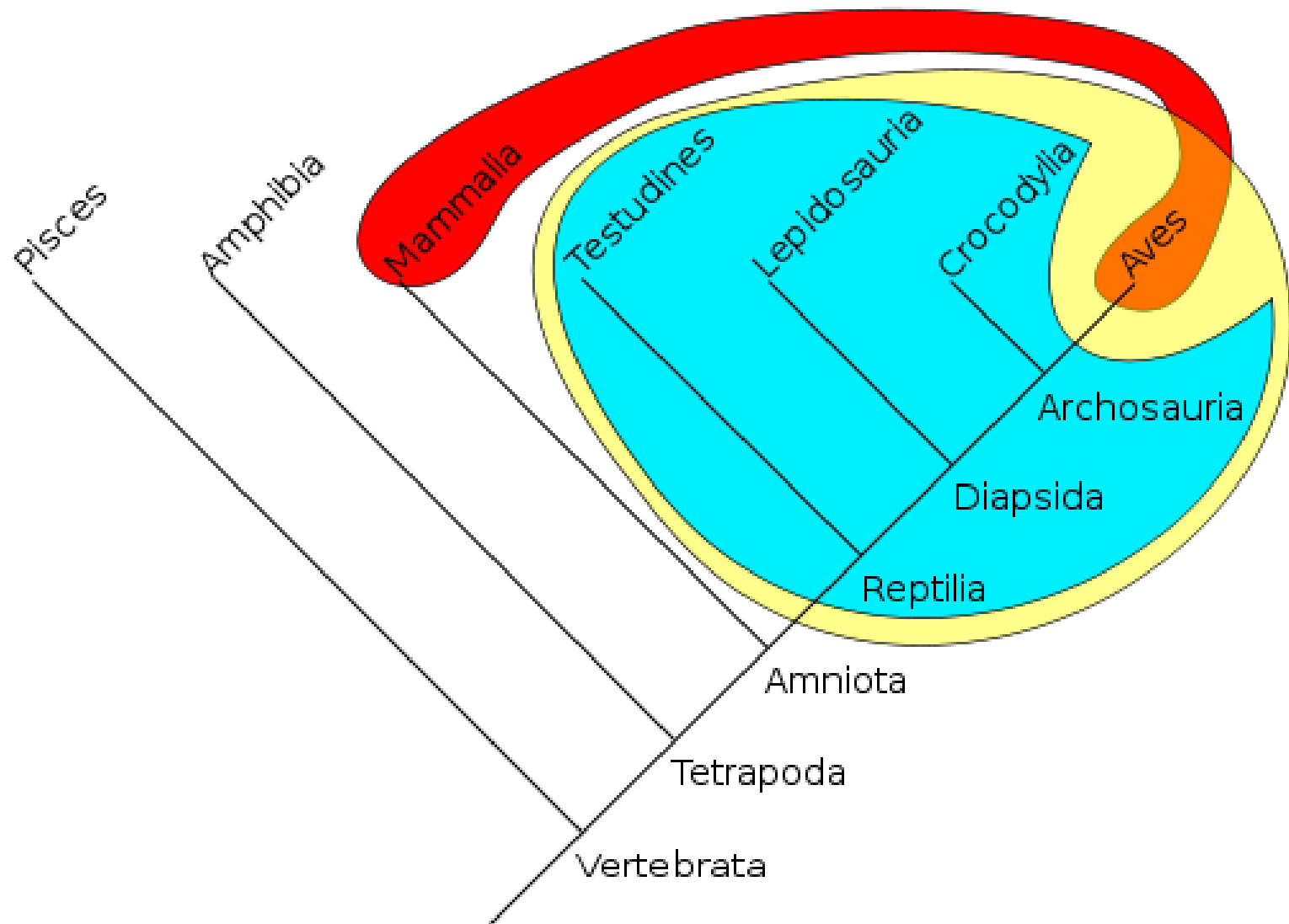
(b) Paraphyletic

Taxon 3
(polyphyletic)



(c) Polyphyletic

- Monophyly
- Paraphyly
- Polyphyly



rooted vs. unrooted phylogenetic trees

rooted tree – root leads to the common ancestor of all studied taxa (e.g.: species)

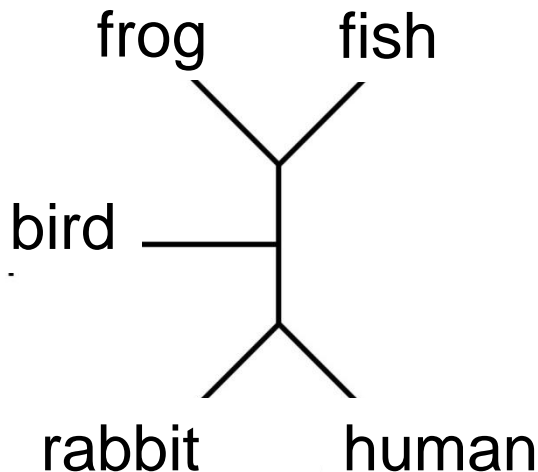
rooting the tree = indicating the direction of the evolutionary process

helps to determine what is plesiomorphic and apomorphic

How to root a tree?

- introducing outgroup

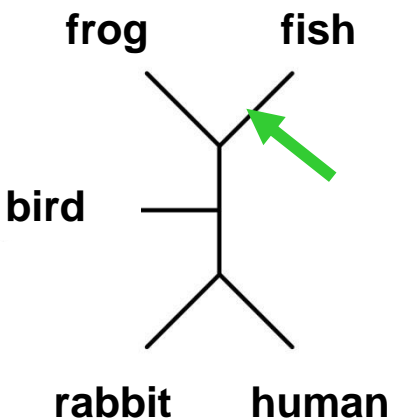
outgroup – a species or group of species that is clearly distant from all the species of interest (ingroup), but still closely related



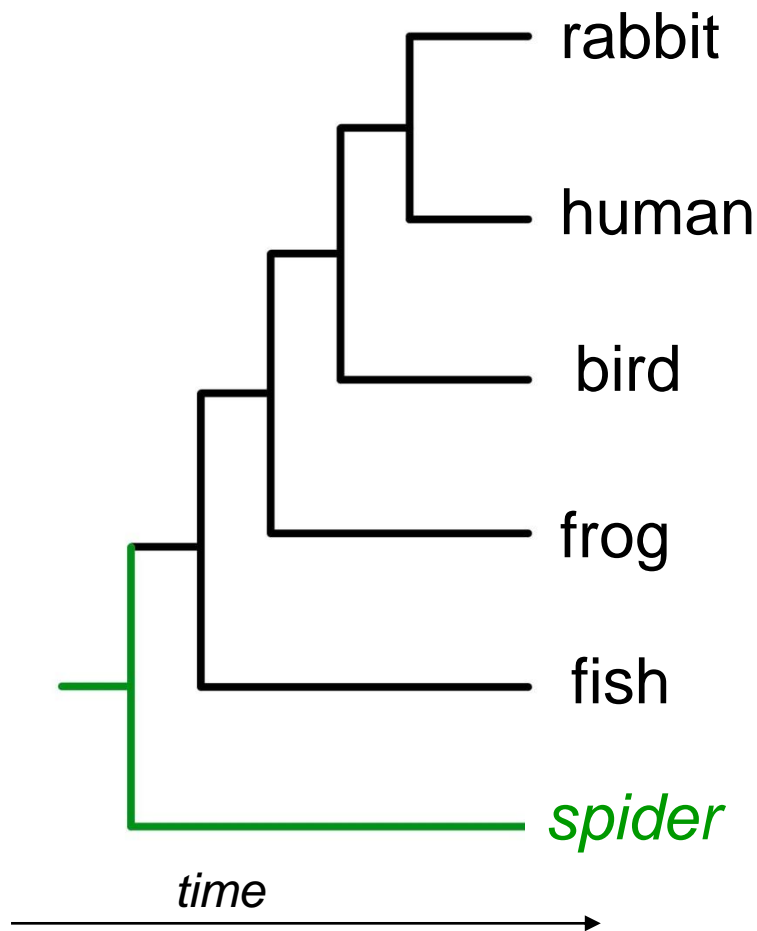
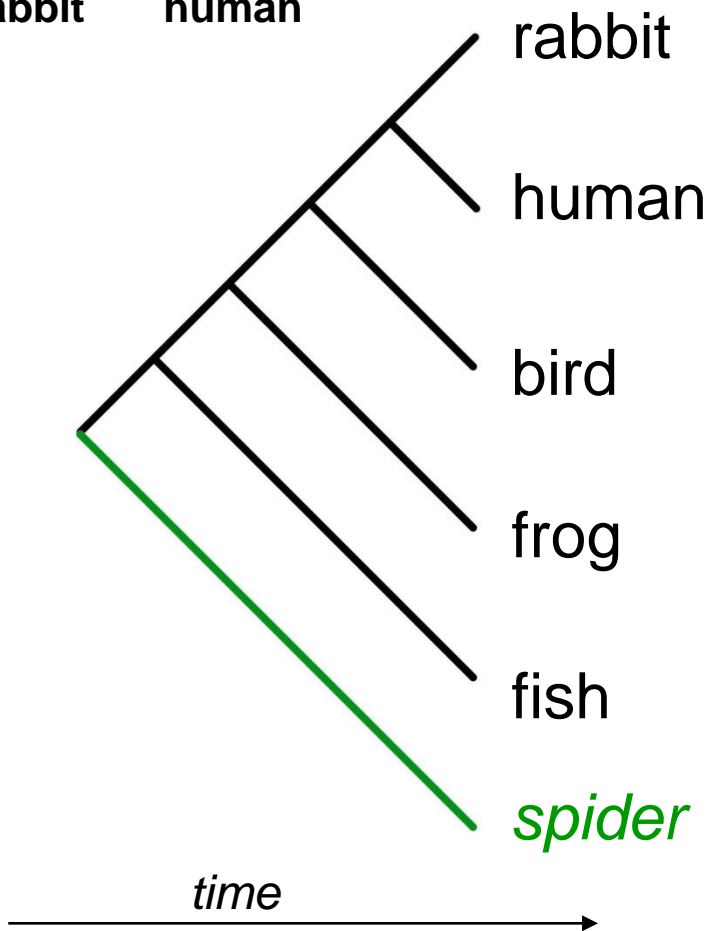
here we will use an invertebrate species to root the tree

(since these species are all Craniata it would be better to use Urochordata (tunicates) to root the tree)

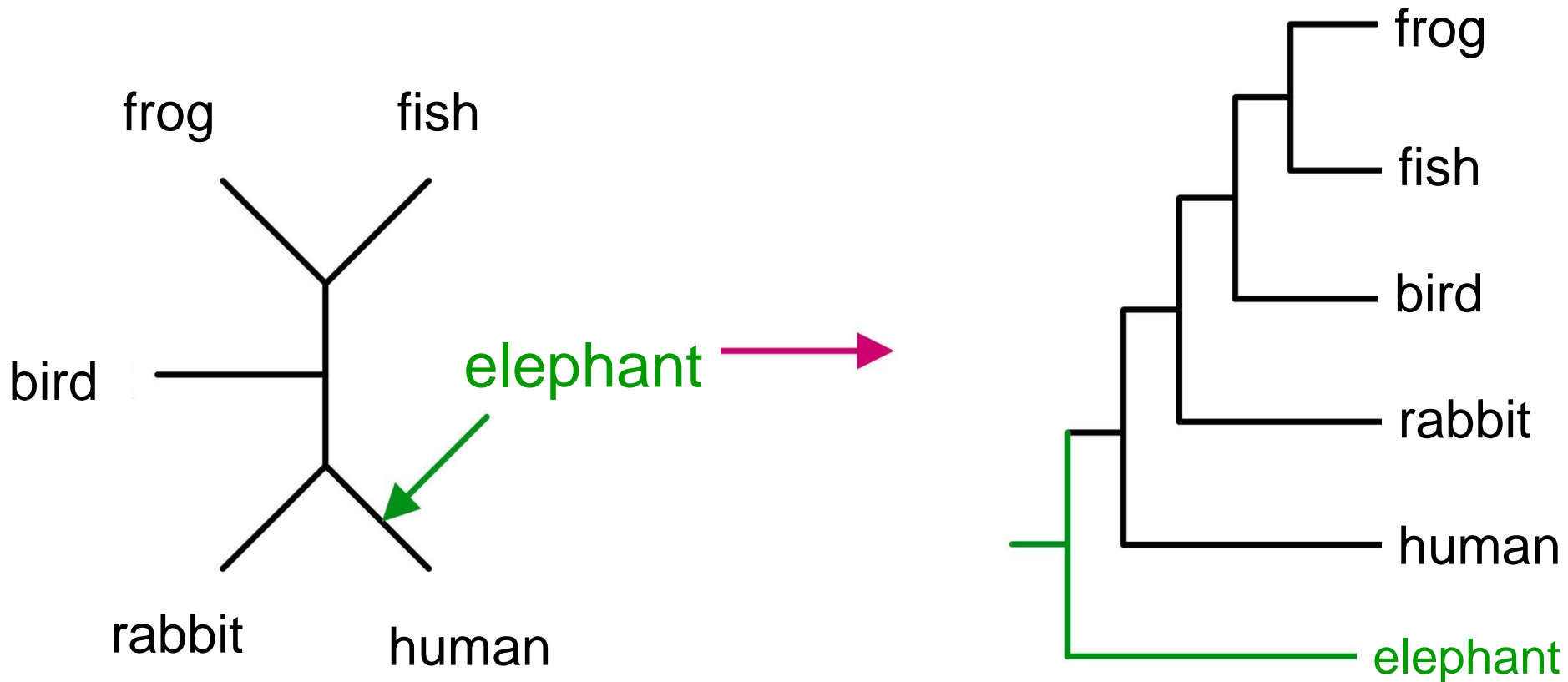
unrooted tree



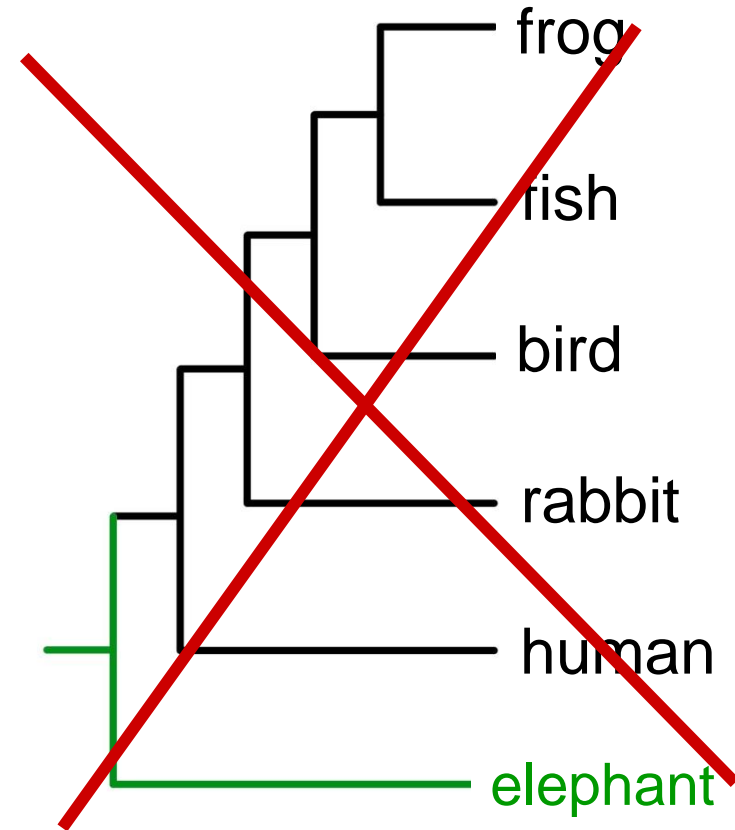
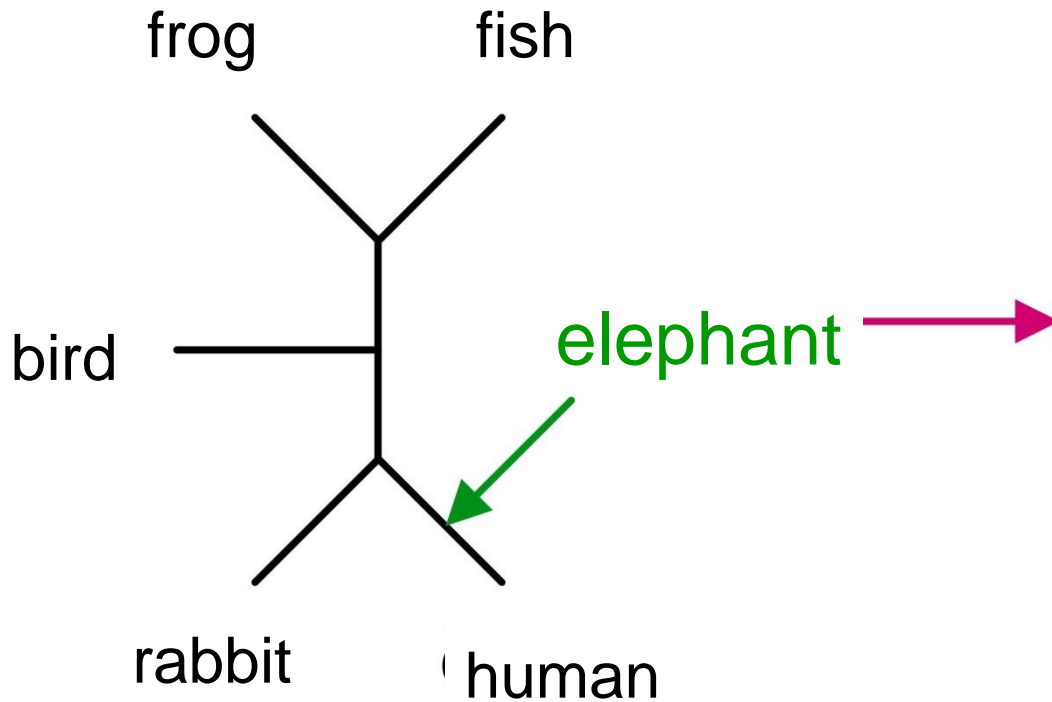
we place the root to the middle of the branch connecting **outgroup** with the rest of the tree



But what if our outgroup is wrong?



But what if our outgroup is wrong?



rooted vs. unrooted phylogenetic trees

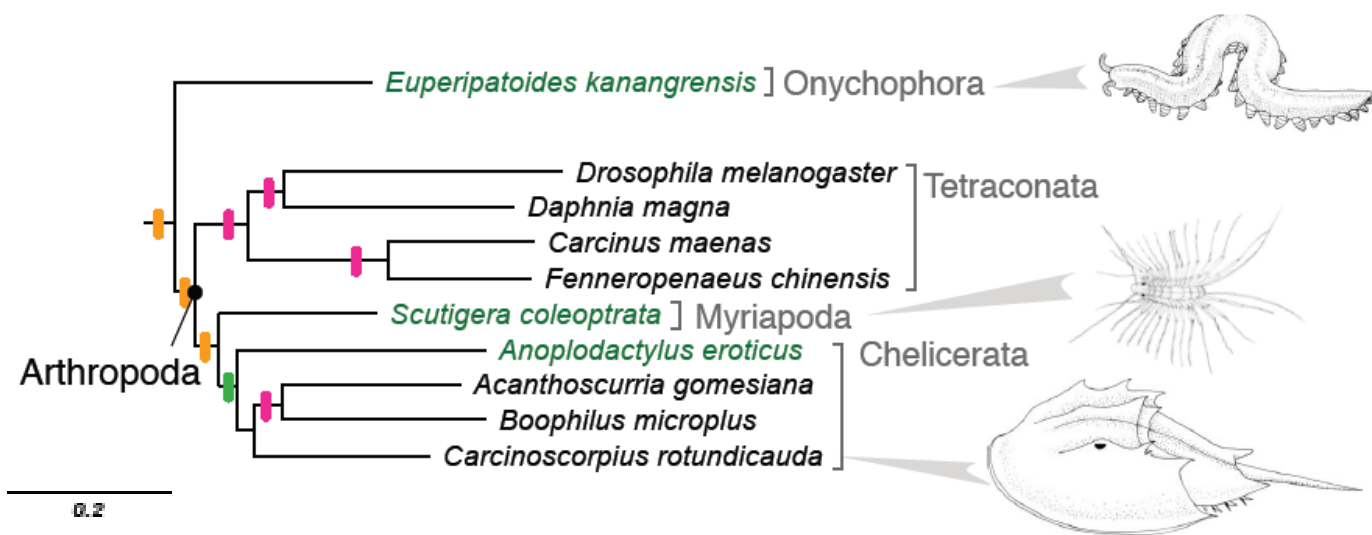
$$(2n - 3)!$$

$2^{n-2}(n-2)!$ Number of bifurcating rooted trees for n taxa (OTUs)

OTU = operational taxonomical unit

Number of bifurcating trees is increasing with number of taxa. The number of unrooted trees for n taxa is equal for the number of rooted trees for $(n-1)$ taxa.

No. of taxa	No. of unrooted trees	No. of rooted trees
3	1	3
4	3	15
5	15	105
10	2 027 025	34 459 425
11	34 459 425	654 729 075
12	654 729 075	13 749 310 575
50	1.00986×10^{57}	2.75292×10^{76}



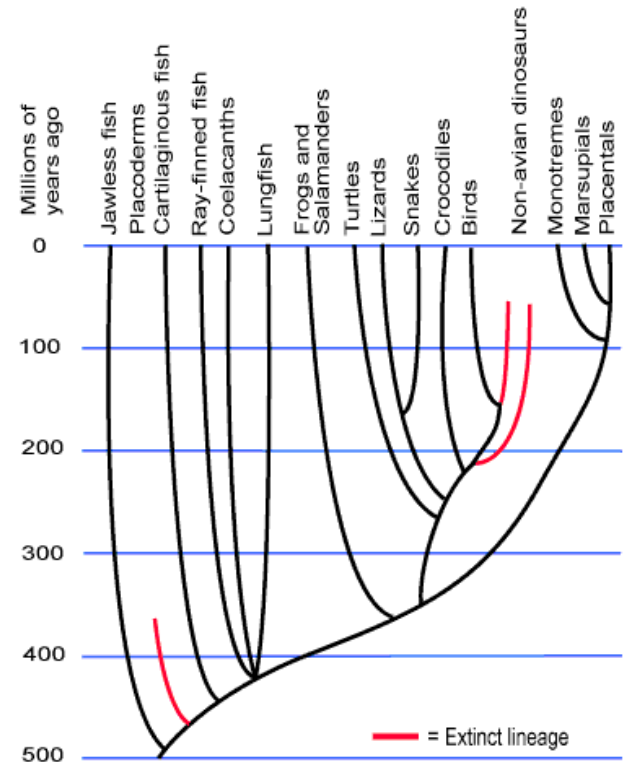
branch length = time of divergence from the common ancestor

molecular clock: sequence divergence increases over time linearly

when molecular clock holds (accumulation is linear over time) – all lineages in the tree have accumulated substitutions at the same rate

evolutionary rate dependent on metabolic rate, generation time, bottleneck events,...

we need calibration points (fossils, geological events)



How to make trees?

like family trees, phylogenetic trees represent patterns of ancestry

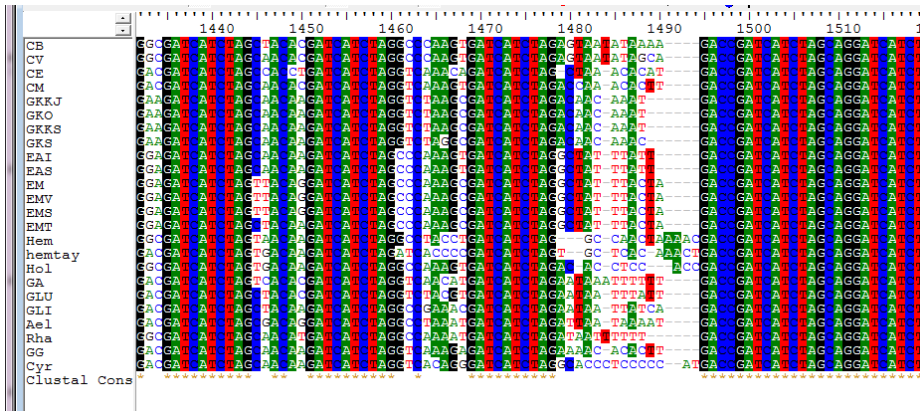
to reveal phylogenetic relationship, we have to compare characters which are **inherited from a common ancestor**

→ **homologous characters** (x analogous characters are result of convergent evolution)

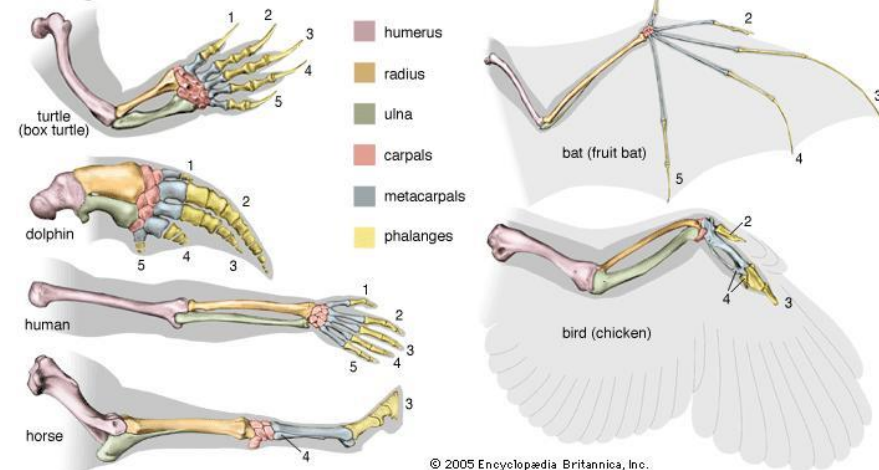
character = heritable traits that can be compared across organisms

two or more forms = **character state**

types: physical characteristics (morphology), genetic sequences, and behavioral traits



Homologies of the forelimb in six vertebrates





Pros and cons of molecular characters (mainly sequence data)

- molecular data - much more abundant (human genome 3,1 Gb, *E. coli* 4.6 Mb)
- independent (one position in the sequence of nucleotides x an eye – the eye is missing, but it means that also cornea, retina, etc. are missing)
- easy to describe (A, C, T, G at the position 175 in the cytochrome b gene) x some structure on the bone more pointed
- can resolve relationships at all different levels of organization, from species and populations to phyla and kingdoms
- less subjective
- neutral – number of shared characters mirrors phylogenetic relationship not just the same environmental selection pressures



Pros and cons of molecular characters (mainly sequence data)

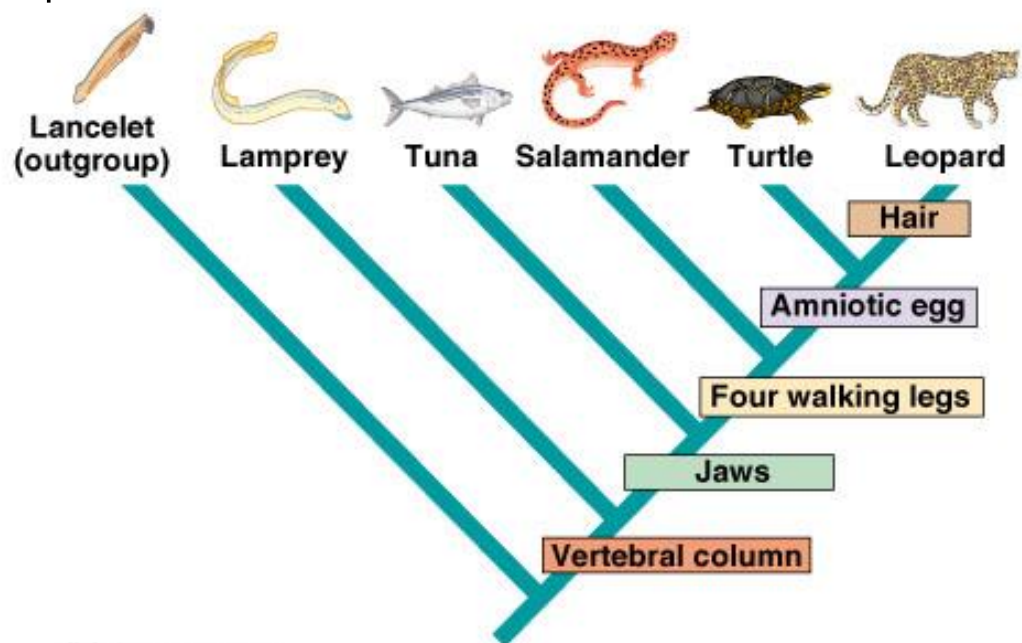
- expensive
- technically demanding
- also non-experts on a given group can use them (lack of insight)
- no information about phenotype

How to make a tree - morphological example

- selected species have shared primitive and derived characters
- for reconstructing phylogeny derived shared characters = **synapomorphies** are important
- create a character table (matrix) with variable characters
- group the taxa based on synapomorphies - the more shared characters, the more closely related are the species
- for morphological data we use usually maximum parsimony method which prefers the simplest explanation of observed data

CHARACTERS	TAXA					
	Lancelet (outgroup)	Lamprey	Tuna	Salamander	Turtle	Leopard
Hair	0	0	0	0	0	1
Amniotic (shelled) egg	0	0	0	0	1	1
Four walking legs	0	0	0	1	1	1
Jaws	0	0	1	1	1	1
Vertebral column (backbone)	0	1	1	1	1	1

(a) Character table

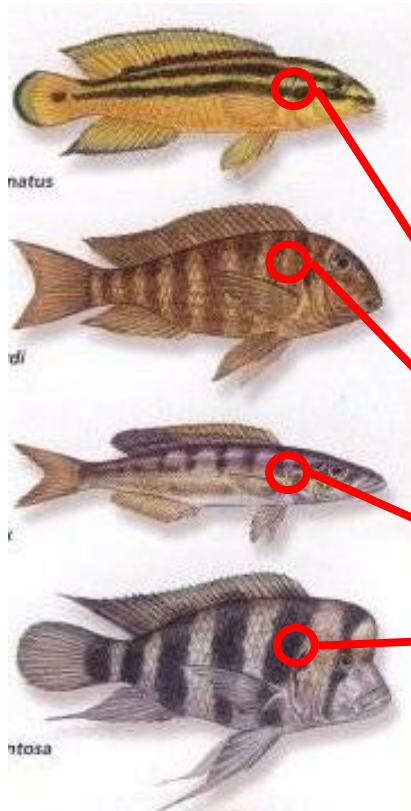


molecular characters

Outgroup
 Taxon A
 Taxon B
 Taxon C

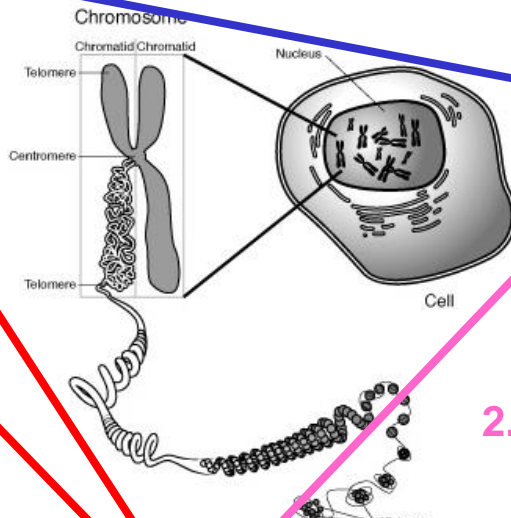
AAGCTTCATAGGAGCAACCAATTCCTAATAATAAGCCTCATAAAGCC
 AAGCTTCACCGGCGCAGTTATCCTCATAATATGCCTCATAATGCC
 GTGCTTCACCGACGCAGTTGTCCTCATAATGTGCCTCACTATGCC
 GTGCTTCACCGACGCAGTTGCCCTCATGATGAGCCTCACTATGCA

3. alignment

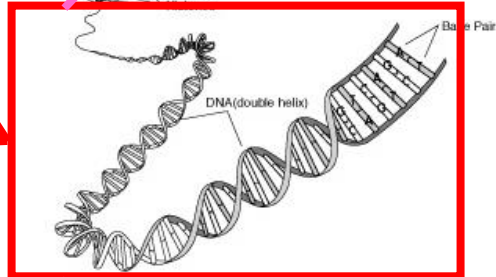


1. DNA isolation

<http://www.accessexcellence.org/AR/GG/chromosome.html>



2. sequencing



molecular characters

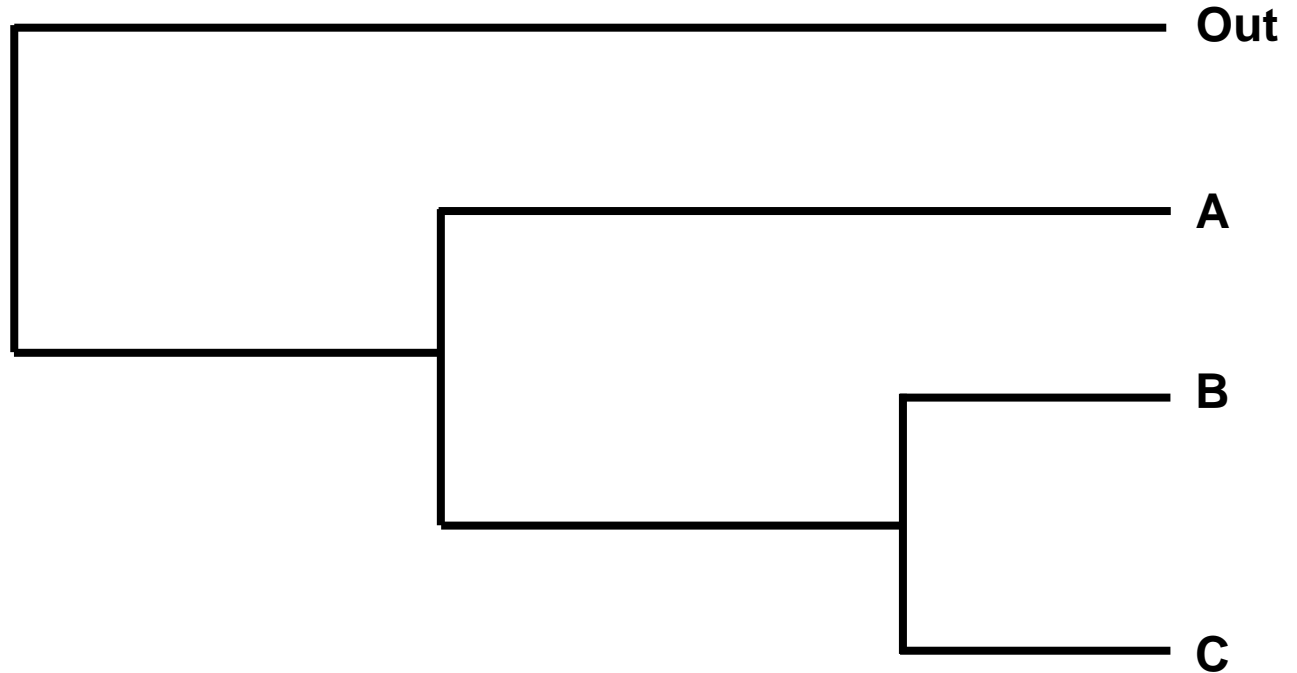
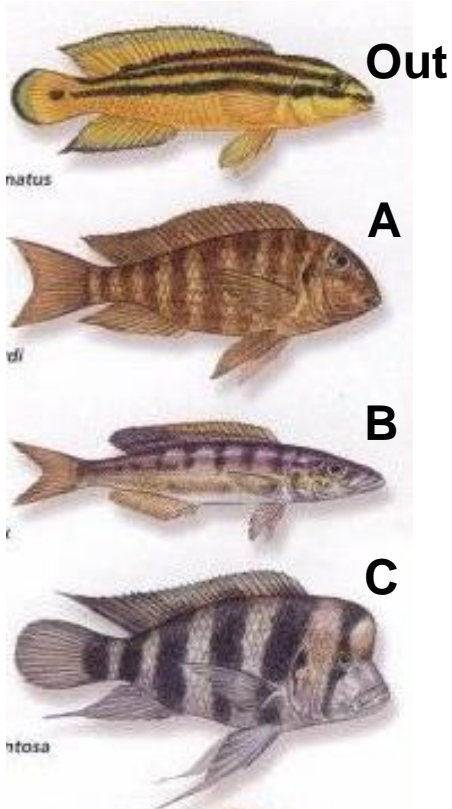
Outgroup

Taxon A

Taxon B

Taxon C

A	A	G	C	T	T	C	A	T	A	invariant sites
G	A	G	C	T	T	C	A	C	A	
G	T	G	C	T	T	C	A	C	G	
G	T	G	C	T	T	C	A	C	G	



molecular characters

Outgroup

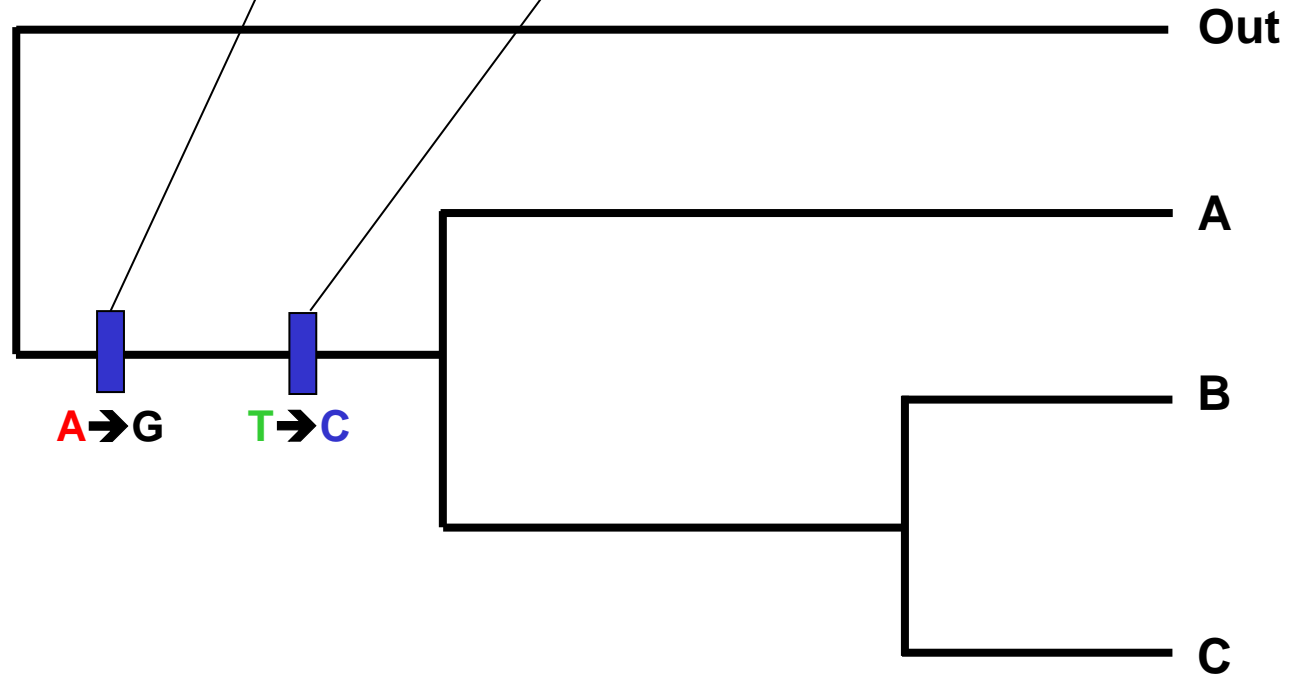
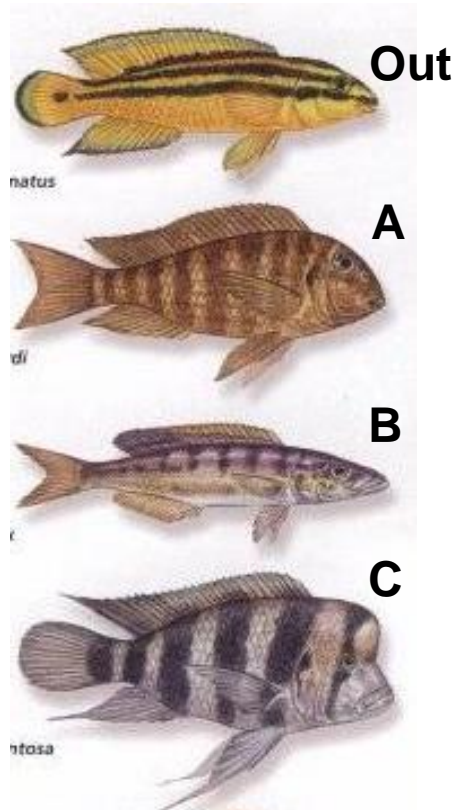
Taxon A

Taxon B

Taxon C

AAGCTTCATA
GAGCTTCACA
GTGCTTCACG
GTGCTTCACG

Synapomorphies
supporting A+B+C



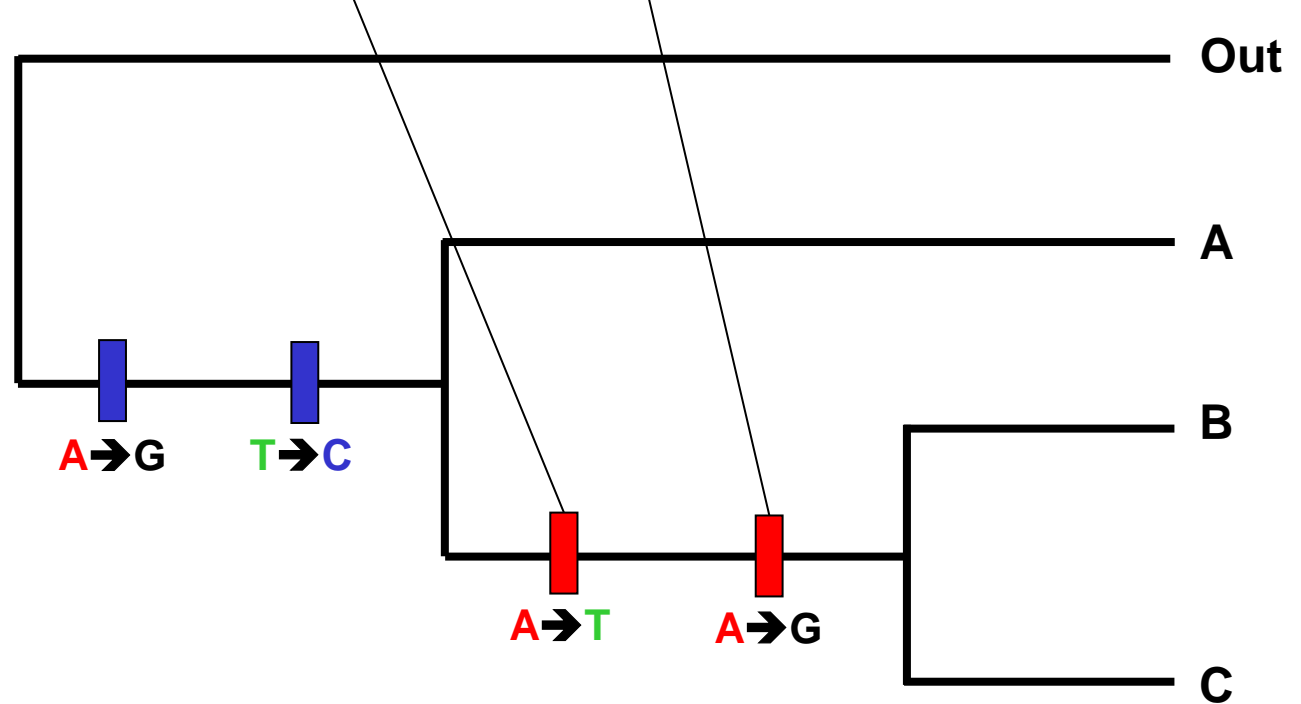
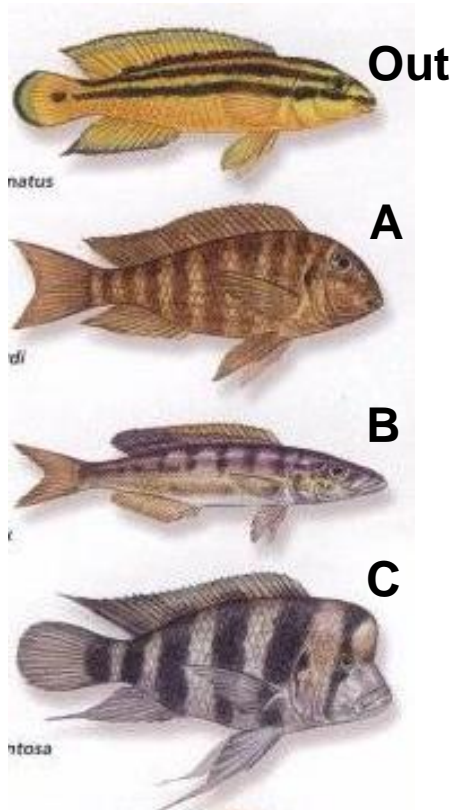
molecular characters

Outgroup
 Taxon A
 Taxon B
 Taxon C

AAGCTTCATA
 GAGCTTCACA
 GTGCTTCACG
 GTGCCTCACG

Synapomorphies
 supporting A+B+C

Synapomorphies
 supporting B+C



molecular characters

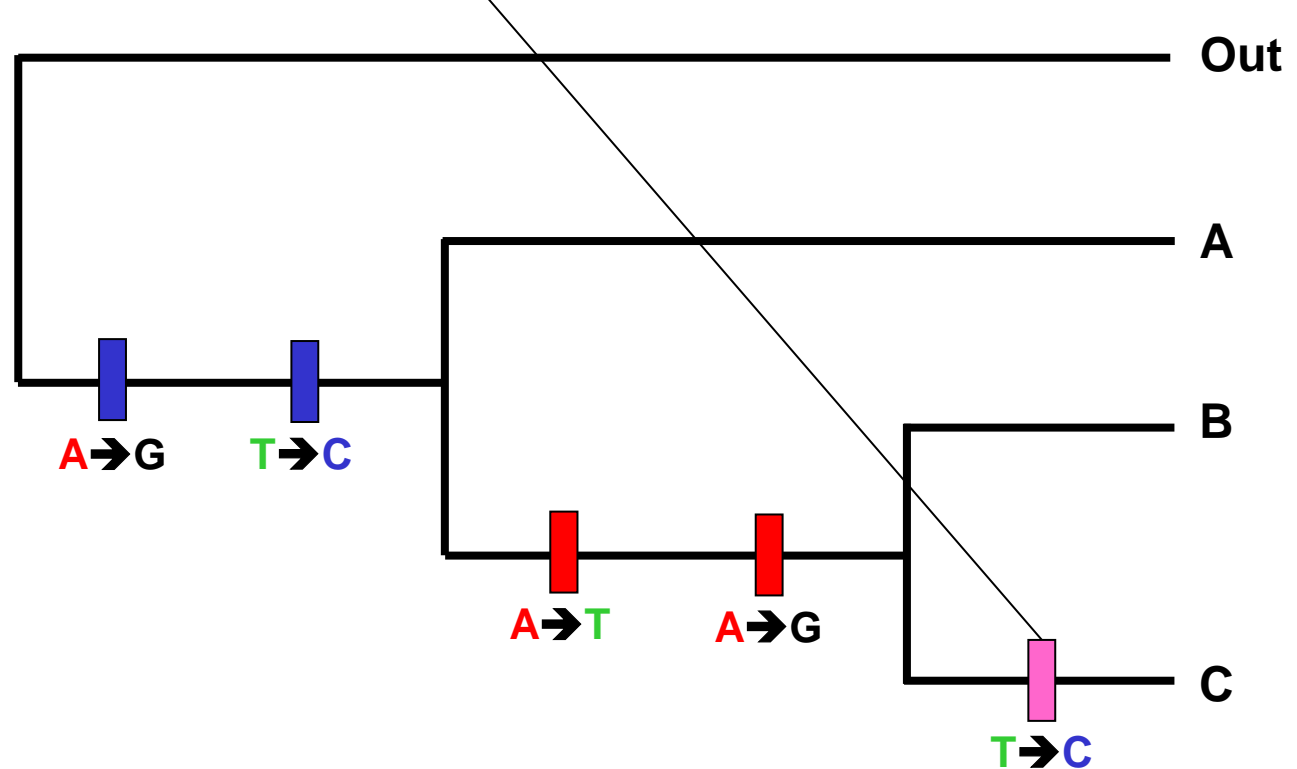
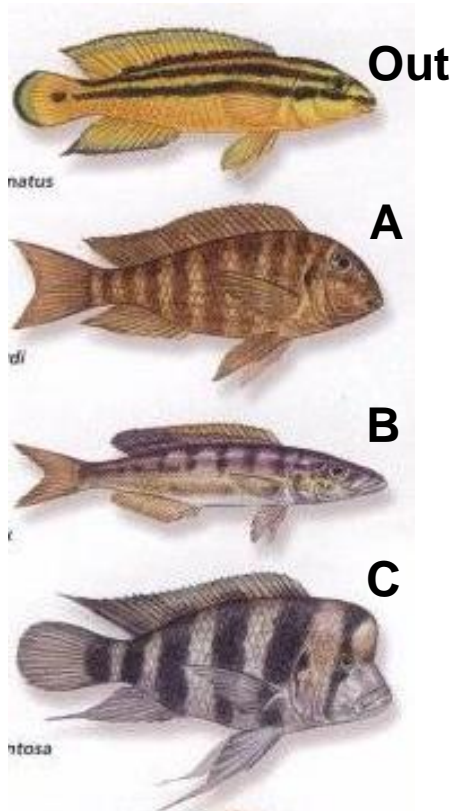
Outgroup
 Taxon A
 Taxon B
 Taxon C

AAGCTTCATA
 GAGCTTCACA
 GTGCTTCACG
 GTGCTCACG

Synapomorphies
 supporting A+B+C

Synapomorphies
 supporting B+C

Apomorphy for C



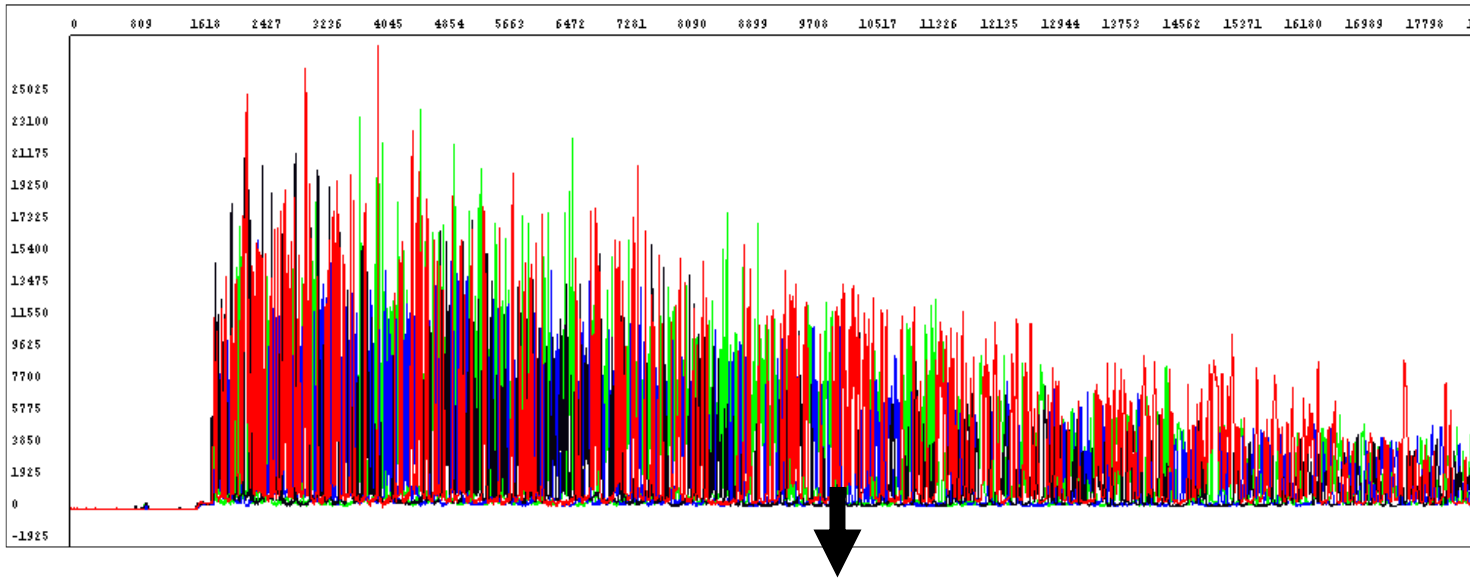
How to make phylogenetic trees?

Workflow:

- ✓ obtain DNA sequence
 - quality check
 - sequence alignment
 - calculating genetic distances
 - phylogeny estimation – topology and branch length
 - reliability test (bootstrap)
 - tree visualization

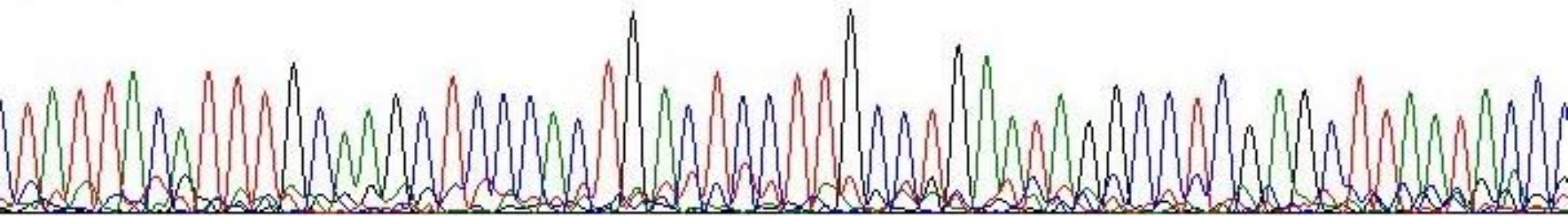
Chromatogram files:

Can be opened in different software- Chromas, BioEdit, DNASTAR , ...



PT_ND2-ND2-f File: D:\Plocha\macrogen_zuzal\MAcrogen_2004_05_05\APT_ND2-ND2-f.ab1

150 160 170 180 190 200
C T A T T A C A T T T G C A A G C T C C C A C T G A C T C C T T G C C T G A A T A G G C C T C G A G C T T A A T A C C C





Checking sequences and chromatograms
programs: Chromas, BioEdit, Geneious

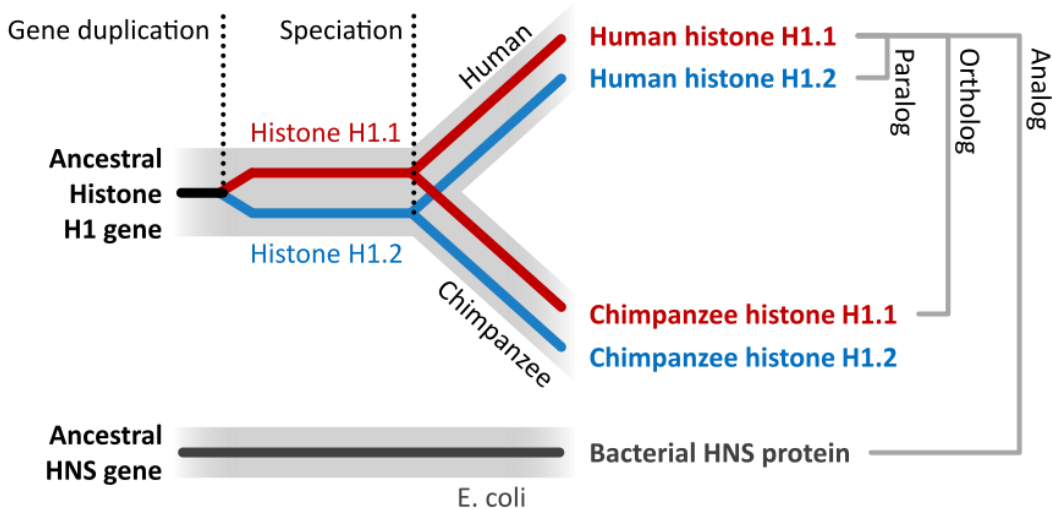
homology of molecular characters

it is necessary to compare sequences (of e.g. genes) which are **orthologous** = inferred to be descended from the same ancestral sequence separated by a speciation event

X

xenolog – sequence of gene incorporated from other species by horizontal transfer

pseudogene – sequence copied from the mitochondrial genome to the nuclear DNA



Gene phylogeny as red and blue branches within grey species phylogeny. Top: An ancestral gene duplication produces two **paralogs** (histone H1.1 and 1.2). A speciation event produces **orthologs** in the two daughter species (human and chimpanzee). Bottom: in a separate species (E. coli), a gene has a similar function (histone-like nucleoid-structuring protein) but has a separate evolutionary origin and so is an **analog**.

Additional sources of DNA sequences:

- public databases



A screenshot of the National Center for Biotechnology Information (NCBI) homepage. The page features a navigation menu on the left with categories like 'NCBI Home', 'Site Map (A-Z)', 'All Resources', 'Chemicals & Bioassays', 'Data & Software', 'DNA & RNA', 'Domains & Structures', 'Genes & Expression', 'Genetics & Medicine', 'Genomes & Maps', 'Homology', 'Literature', 'Proteins', 'Sequence Analysis', 'Taxonomy', 'Training & Tutorials', and 'Variation'. The main content area includes a 'Welcome to NCBI' message, a 'Get Started' section with links to 'Tools', 'Downloads', 'How-To's', and 'Submissions', and an 'Education Resources' section. A search bar at the top right contains the word 'Nucleotide'. The right sidebar lists 'Popular Resources' such as BLAST, Bookshelf, Gene, Genome, Nucleotide, OMM, Protein, PubChem, PubMed, PubMed Central, and SNP, as well as 'NCBI News' with a 'New NCBI Newsletter' link.

<http://www.ncbi.nlm.nih.gov/>

BLAST - **Basic Local Alignment Search Tool**

- algorithm for searching in databases for similar sequences

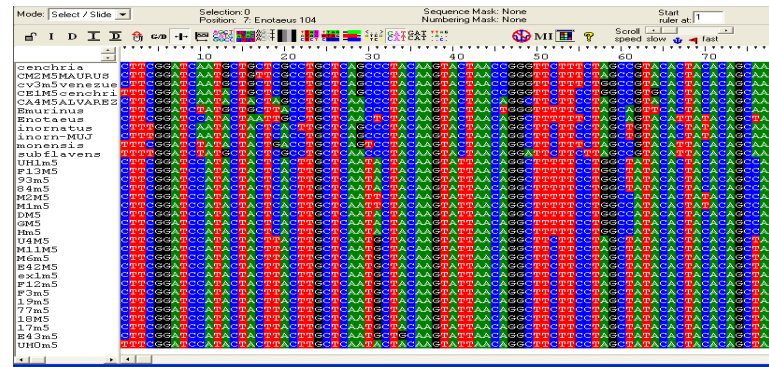
A screenshot of the BLAST (Basic Local Alignment Search Tool) website interface. The top navigation bar includes 'Home', 'Recent Results', 'Saved Strategies', and 'Help'. The main content area features a search box with the text 'BLAST finds regions of similarity between biological sequences. [more...](#)'. Below the search box is a 'New' button and a text input field containing 'DELTA-BLAST, a more sensitive protein-protein search' with a 'Go' button. The right sidebar contains sections for 'Your Recent Results' with a 'New!' indicator and a link to 'All Recent results...', and a 'News' section with a link to 'BLAST 2.3.20 released'.

How to make phylogenetic trees?

Workflow:

- ✓ obtain DNA sequence
- ✓ quality check
 - sequence alignment
 - calculating genetic distances
 - phylogeny estimation – topology and branch length
 - reliability test (bootstrap)
 - tree visualization

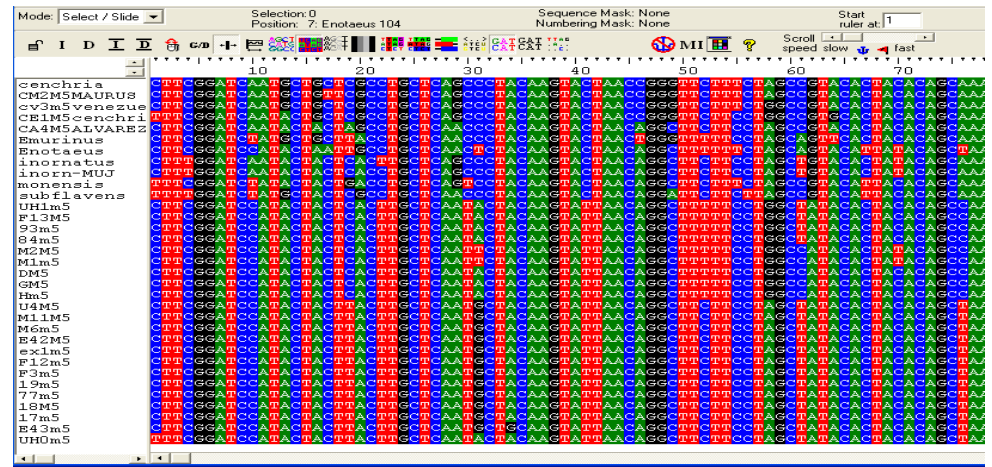
Where sequences differ and where are the same?



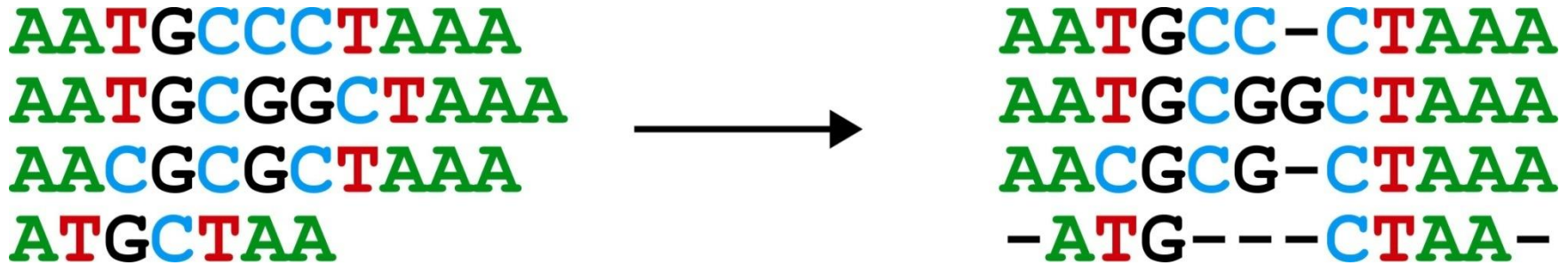
Alignment

- a way of arranging the sequences of DNA (also e.g. amino acid in the protein sequence) to identify regions of similarity
- start of every phylogenetic analysis
- assessing of position homology of each base in the sequence
- each position (column in the alignment) in the sequence represents character potentially useful for the phylogenetic analysis
- different programs for calculating and editing alignments
 - manual: BioEdit, Macaw
 - automatic – different algorithms

Clustal X, PileUp, Multalin, Mafft – often online



Alignment - pairwise alignment (two sequences)
 - multiple alignment (more sequences)



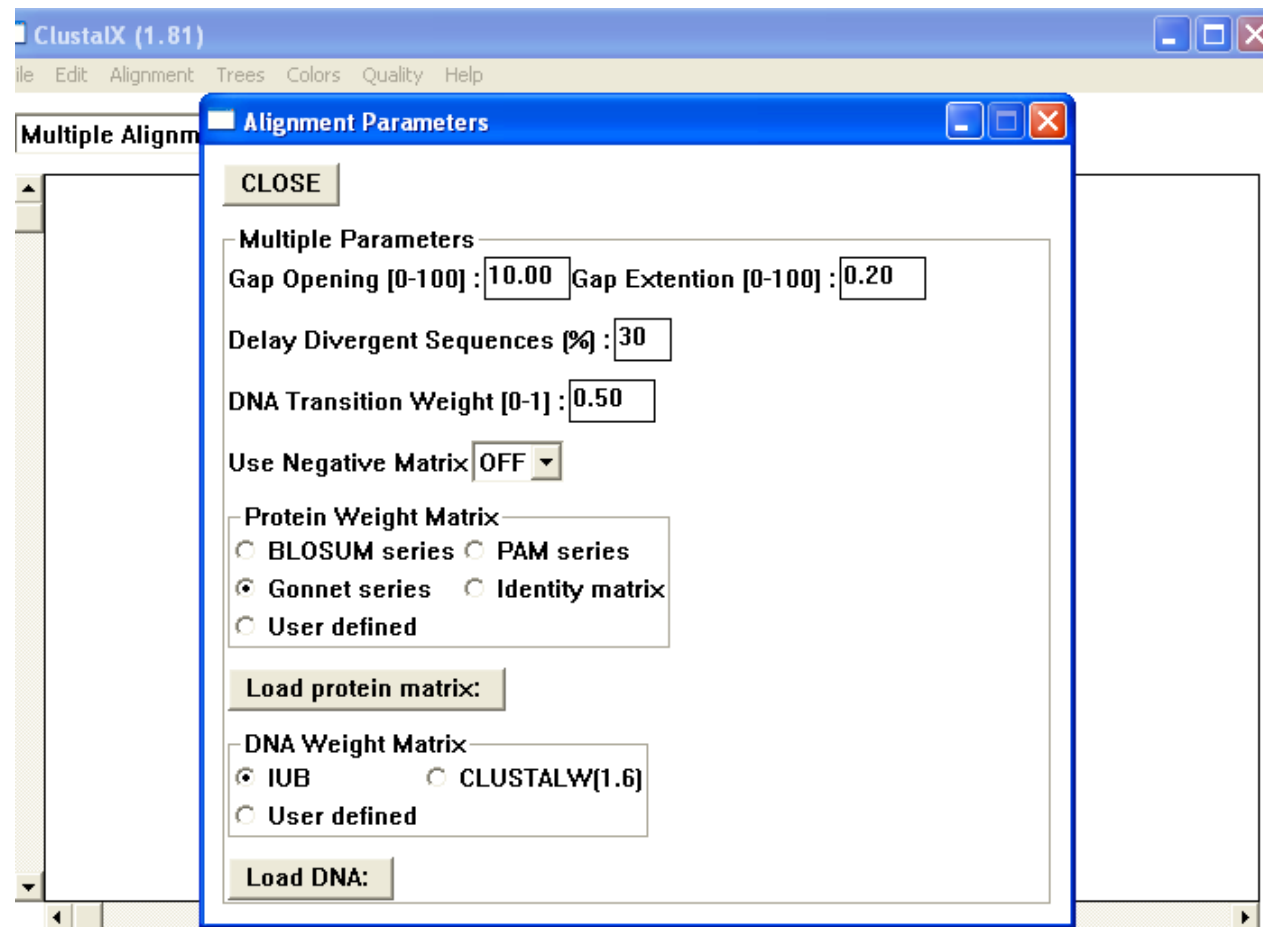
gaps are inserted between the bases so that identical or similar characters are aligned in successive columns



Alignment:

Clustal W (Clustal X)

www.clustal.org – frequently used software (Thompson et al. 1994 a 1997)



What can help:

- first use default parameters, later increase penalty for gap opening and decrease penalty for gap extension
- coding genes align as amino acids
- use information about sequence secondary structure (genes for 12S and 16S RNA) – database with alignments with secondary structure in mind (<http://www.arb-silva.de/>)
- delete ambiguously aligned positions

MOLECULAR PHYLOGENETICS AND EVOLUTION
Vol. 4, No. 1, March, pp. 1–9, 1995

**Elision: A Method for Accommodating Multiple Molecular
Sequence Alignments with Alignment-Ambiguous Sites**

WARD C. WHEELER,* JOHN GATESY,† AND ROB DeSALLE‡

MOLECULAR PHYLOGENETICS AND EVOLUTION
Vol. 2, No. 2, June, pp. 152–157, 1993

**Alignment-Ambiguous Nucleotide Sites and the Exclusion of
Systematic Data**

JOHN GATESY,* ROB DeSALLE,† AND WARD WHEELER*



C:\Users\Zuz\syncho-stary notas\synchro1\workshop\ali-podle-sek-struk\alisdek10-t=0,5.aln

Courier New 11 B 25 total sequences shade threshold 60%

Mode: Select / Slide Selection: 0 Position: Sequence Mask: None Numbering Mask: None Start ruler at: 1

Scroll speed slow fast

	1440	1450	1460	1470	1480	1490	1500	1510	1
CB	GGCGATCATCTAGCTACACGATCATCTAGGCCCAAGTGA	TCACTCTAGAGTAAATATAAAA	---	GA	CCGATCATCTAGCAGGATCATCT				
CV	GGCGATCATCTAGCAACCGATCATCTAGGCCCAAGTGA	TCACTCTAGAGTAAATATAAGCA	---	GA	CCGATCATCTAGCAGGATCATCT				
CE	GACGATCATCTAGCACTGATCATCTAGGTCAAA	CAGATCATCTAGCTAAACACAT	---	GA	CCGATCATCTAGCAGGATCATCT				
CM	GACGATCATCTAGCAACCGATCATCTAGGTCAAA	GTGATCATCTAGAGCAAACACCT	---	GA	CCGATCATCTAGCAGGATCATCT				
GKKJ	GAAAGATCATCTAGCAACAGATCATCTAGGTC	TAAAGCGATCATCTAGAGAACAAAT	---	GA	CCGATCATCTAGCAGGATCATCT				
GKO	GAAAGATCATCTAGCAACAGATCATCTAGGTC	TAAAGCGATCATCTAGAGAACAAAT	---	GA	CCGATCATCTAGCAGGATCATCT				
GKKS	GAAAGATCATCTAGCAACAGATCATCTAGGTC	TAAAGCGATCATCTAGAGAACAAAT	---	GA	CCGATCATCTAGCAGGATCATCT				
GKS	GAAAGATCATCTAGCAACAGATCATCTAGGTC	TAAAGCGATCATCTAGAGAACAAAT	---	GA	CCGATCATCTAGCAGGATCATCT				
EAI	GGAGATCATCTAGCAACAGATCATCTAGCCCAAAG	TGATCATCTAGGCTATTTAT	---	GA	CCGATCATCTAGCAGGATCATCT				
EAS	GGAGATCATCTAGCAACAGATCATCTAGCCCAAAG	TGATCATCTAGGCTATTTAT	---	GA	CCGATCATCTAGCAGGATCATCT				
EM	GGAGATCATCTAGTTACAGGATCATCTAGCCCAAAG	CGATCATCTAGGCTATTTACTA	---	GA	CCGATCATCTAGCAGGATCATCT				
EMV	GGAGATCATCTAGTTACAGGATCATCTAGCCCAAAG	CGATCATCTAGGCTATTTACTA	---	GA	CCGATCATCTAGCAGGATCATCT				
EMS	GGAGATCATCTAGTTACAGGATCATCTAGCCCAAAG	CGATCATCTAGGCTATTTACTA	---	GA	CCGATCATCTAGCAGGATCATCT				
EMT	GGAGATCATCTAGTTACAGGATCATCTAGCCCAAAG	CGATCATCTAGGCTATTTACTA	---	GA	CCGATCATCTAGCAGGATCATCT				
Hem	GGCGATCATCTAGTTACAGGATCATCTAGGCC	TACCTGATCATCTAGGC	CAACCTAAAAC	GA	CCGATCATCTAGCAGGATCATCT				
hemtay	GACGATCATCTAGTGCAACAGATCATCTAGAT	ACCCCGATCATCTAGTGC	TCACAAACT	GA	CCGATCATCTAGCAGGATCATCT				
Hol	GGCGATCATCTAGTGCAACAGATCATCTAGGCC	CAAAGTGA	CTAC	GA	CCGATCATCTAGCAGGATCATCT				
GA	GACGATCATCTAGTCAACCGATCATCTAGGTC	AAACATGATCATCTAGAGTAAATTTT	TT	GA	CCGATCATCTAGCAGGATCATCT				
GLU	GACGATCATCTAGTTACAGGATCATCTAGGTC	TACGTTGATCATCTAGAGTAAATTTT	TT	GA	CCGATCATCTAGCAGGATCATCT				
GLI	GACGATCATCTAGTTACAGGATCATCTAGGTC	GAAACGATCATCTAGAGTAAATTTT	TCA	GA	CCGATCATCTAGCAGGATCATCT				
Ael	GACGATCATCTAGTCAACAGGATCATCTAGGCC	TAAATGATCATCTAGAGTAAATTTT	TT	GA	CCGATCATCTAGCAGGATCATCT				
Rha	GGCGATCATCTAGCAACATGATCATCTAGGCC	AAAATGATCATCTAGAGTAAATTTT	TT	GA	CCGATCATCTAGCAGGATCATCT				
GG	GACGATCATCTAGCAACAGGATCATCTAGGTC	AAAAGATCATCTAGAGAAC	ACACCT	GA	CCGATCATCTAGCAGGATCATCT				
Cyr	GACGATCATCTAGCAACAGGATCATCTAGGTC	ACAGGATCATCTAGGACCC	TCCCCC	AT	GA	CCGATCATCTAGCAGGATCATCT			
Clustal Cons	*****								



Phylogeny estimation

character based (maximum parsimony, maximum likelihood, Bayesian analysis)

two types of methods

distance based (Neighbour-joining, UPGMA)

Two different approaches:

algorithm – number of specific steps resulting in one best tree
methods: UPGMA, Neighbour-joining

optimality criterion – consider and compare all theoretically possible trees based on selected criteria- number of evolutionary steps, likelihood value