

# Phylogenetics introduction

Episode in the series on phylogenetics

Mikael Bodén

# **Part 1: Phylogenetics introduction**

Motivation

Background

Homology

1

Reference ...ACGTACGGTTACACAAAACCGTTTACGTAGTTGTGACG...

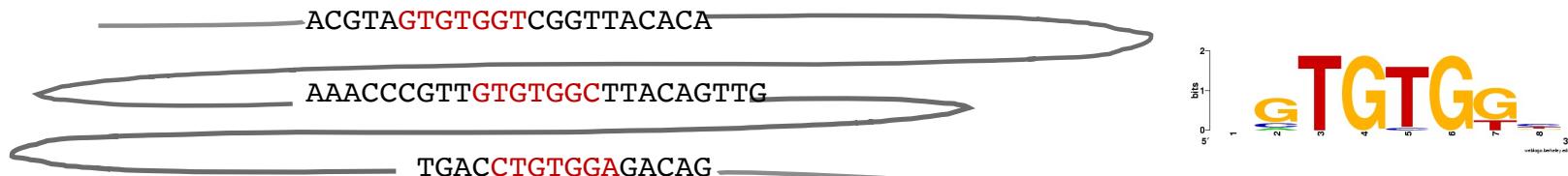
Sequence  
reads

```

TTACACAAATCCCGTTGCACGTA
TACACAAATCCCGTTGTACGTAG
ACACAAAACCCGTTGCACGTAGT
CAAATCCCGTTGCACGTAGTTGT
AAATCCCGTTGTACGTAGTTGTG

```

2

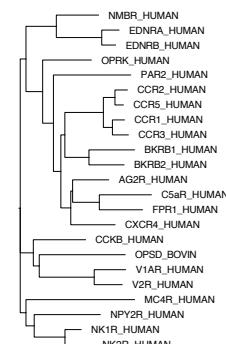


3

```

----EHLKQRREVAKTVFCLVVIPALCNPLHLSRILKKTVVNEMDKNRCELLSFLLMDYIGINLAATMSCINPIALVFVSKKFRNCQSCLCCCCC----
----DHLKQRREVAKTVFCLVLVFALCNPLHLSRILKLELYNNDNPRNCSELLFLLDYGIGNMSLSCINPIALYLVLSKFKRNCSCLCCCNCC----
----QKNKRNNDFIKLIMAIVLFFFWIPHOIFTFLDVLUIQGIIRD-CRIADIVTAMPITCICIAFFNCLNPLFYGFLGKFKRYFLOLKYP----
----QKRK----ALKNVLILLAPAFCWLPYYGISDSFILELIKOGEFENTWKWLSALAFFHCCLNPILYAFLGAKFTKSQHALTSVSR----
----PNEKK---SKAVNLFVIMIIPEFLWMPYPNNYILLSVFQDPLFHE-CEOSRHLDLAVVTEVIATHCCVNPVYIYAFVGERFRKYLROLFHRRVA-
----PSKKK-YKAIRLFVIMIIPEFLWMPYPNNYILLSVFQDPLFHE-CEOSRHLDLAVVTEVIATHCCVNPVYIYAFVGERFRKYLROLFHRRVA-
----RNEKKRRHRAVRVIVTFIMIVYFLPWMPYPNNYILLSVFQDPLFHE-CEOSRHLDLAVVTEVIATHCCVNPVYIYAFVGERFRKYLROLFHRRVA-
----RTRCGGRKDSKTALILLTLVVAFLCWCAPYHFAFLEFLFQVOAVRG-CFWDEEIDLLOLANFAFTTSSSLPVYIVFVGRLFRTKWELIKOCP-
----FKEIQTE-RATVLVVLLLFICCWLPOISEFLDTLHLGILS-CEDEILDVIQASMAXSICLNPVYVVGKRFRKYEWEVYGVCQ-
----SVKSISRAKIRTVKMTFVVIBAYIVCWCAPYFLIQOMSWDMSVITE-----SENPITITALLGSLSCNCPWYIMFSGHLLDCVORFFCON-
----EGAHVSAAVKTVRMRVVIVVVILCWCAPYFLIQOMSWDMSVITE-----LENGPFVLLMLLSASLSCTNCPWYIASFS-SSSESELRLLCAR-
----RYHEQVSAKRKVVKMIVVVCTFAICCWIPHIIFFLPYINPDLYLKK-----FIQOVVLAMMLAMSSTMNPIIYCCLNDRFLGKHAFRCCPFISA-
----ANLRHLQAMKFVKDMVLVLFTFAICCWIPHIIFFLPQEDIYCHK-----FIQOVVLAMMLAMSSTMNPIIYCCLNDRFLGKHAFRCCPFVTE-
----GSREKDRNLRITTLVLVVVVAVFVCMPIPIIFLVEALQSTSHST-----ALSSYFCIALGYTSSSLPNPIIYAFLDENRKCFRDFCPPLKMR-
----RATRSTKTLVVLVVVVASFFFIFWIPYOVGIMMSFLPEPSSPT-LLLNKLDSLCVSFAINCCINPIIYVVAQQGPFQGRLRKSLPSLR-
----LIKSSPRPLVLSFVAAAFFLCWSPYQVVVALLAATVRIRELLOG-----MYKEIGIAVDVTSALAFSICCLPMLYVFMQDFRERLHALPASLER-
----ANDHYHORRQKTLMLCVCVVVVVVBNLPLHAFQLAVDIDSQVLDLK-EWKLIFTVHILAMCSTFANPLLYGWMNSNYRKAFLSAFCEQR-
----NESEKKRRAIKLIVVLAMLICFPSNLLLVVHFLIKSQGOS-----HVYALIVACLCSSNCCLPDFVVVFSHPDFRHAKNLLCRSVR-
----ESATTQKAEKEVRWIIMVIAFLCWIPYAVFYFHQGSDFG-----IFMTPAFFATSSAnPVIYIMKOFRNCVWILLCGNKNP-
----AIRQGANMKGAIELTILIGVFVCWCAPFLHLIFYICPQNPICV-----CFMSHFNLYLILIMCSIIDPLIYALRSQERKTFEICCCYPLG-

```

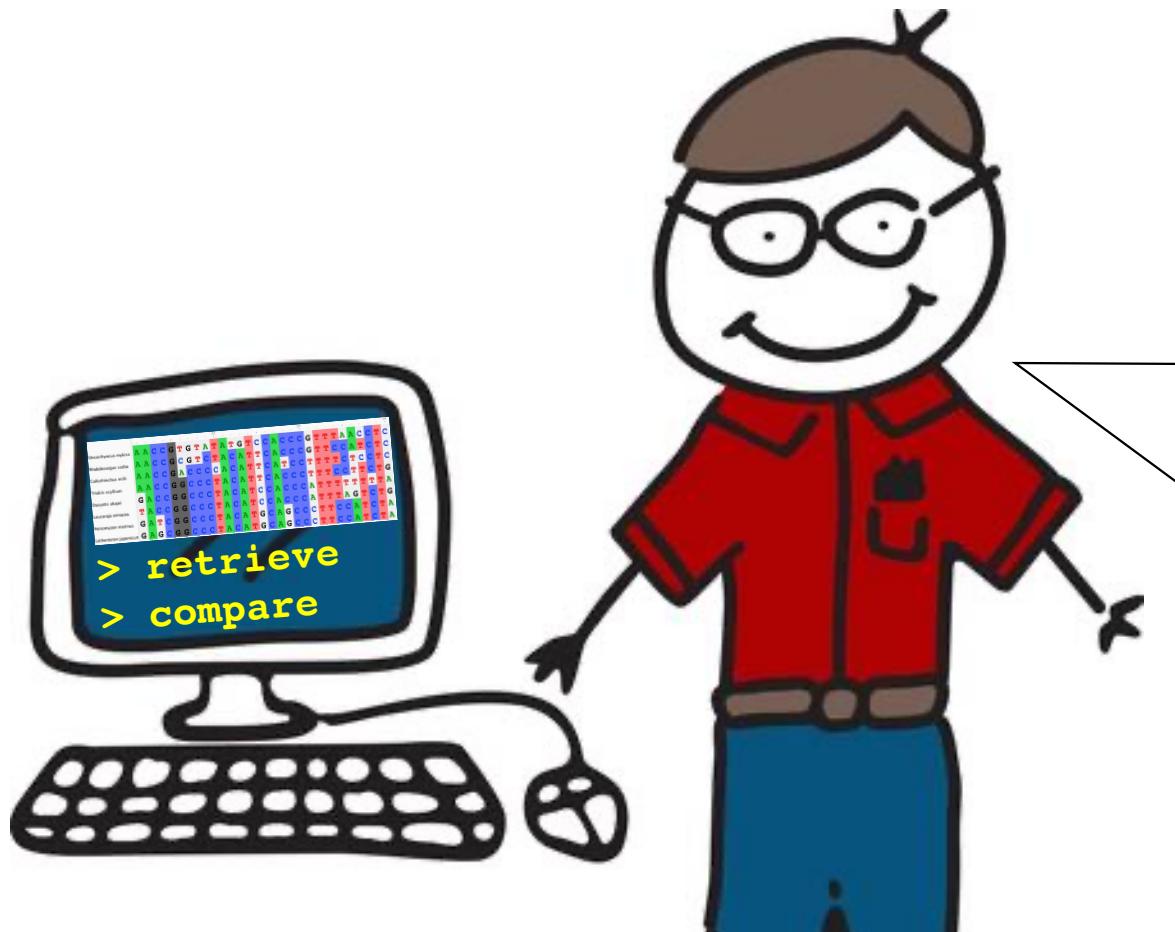


- 
- Study parts
  - Probe function and structure
  - **NEW:** Trace history

10100100100  
1010001110101010  
01001001110010000  
001000111000100101  
100110111110010  
11001001111001  
00001111100  
100001111001  
100001110101  
010100100100  
1001001000111

I am a **bioscience** student.  
Why should I pay attention?

# I am a **computing** student. Why should I pay attention?

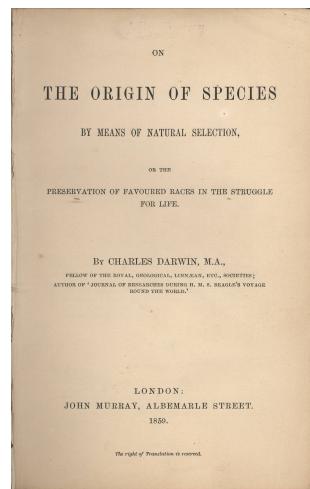
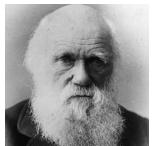


How cool is it that I will learn

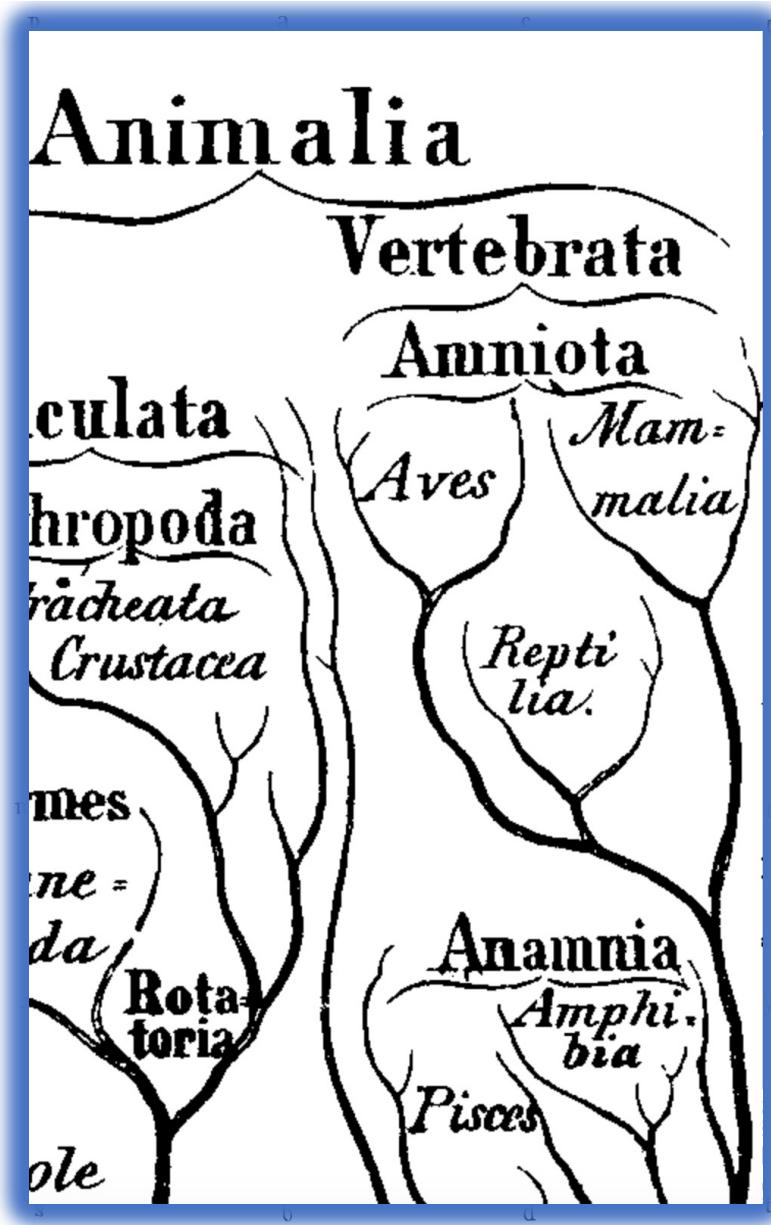
- how to determine *when* existing sequences shared an ancestor?
- how to recreate phylogenetic trees, and ancient, probably extinct sequences?
- what it takes to identify the evolutionary events that formed existing sequences?

# Phylogeny and common ancestor

- All organisms are related by genealogical descent with *modification*, as on a *branching tree*



Phylogenetic tree suggested by Haeckel (1866)



# Homology, orthology and paralogy

**Homology** means descent from a common ancestor

**Orthology** is a special case of homology, in which the descent of sequences from a common ancestral form is topologically congruent with organismal **speciation**

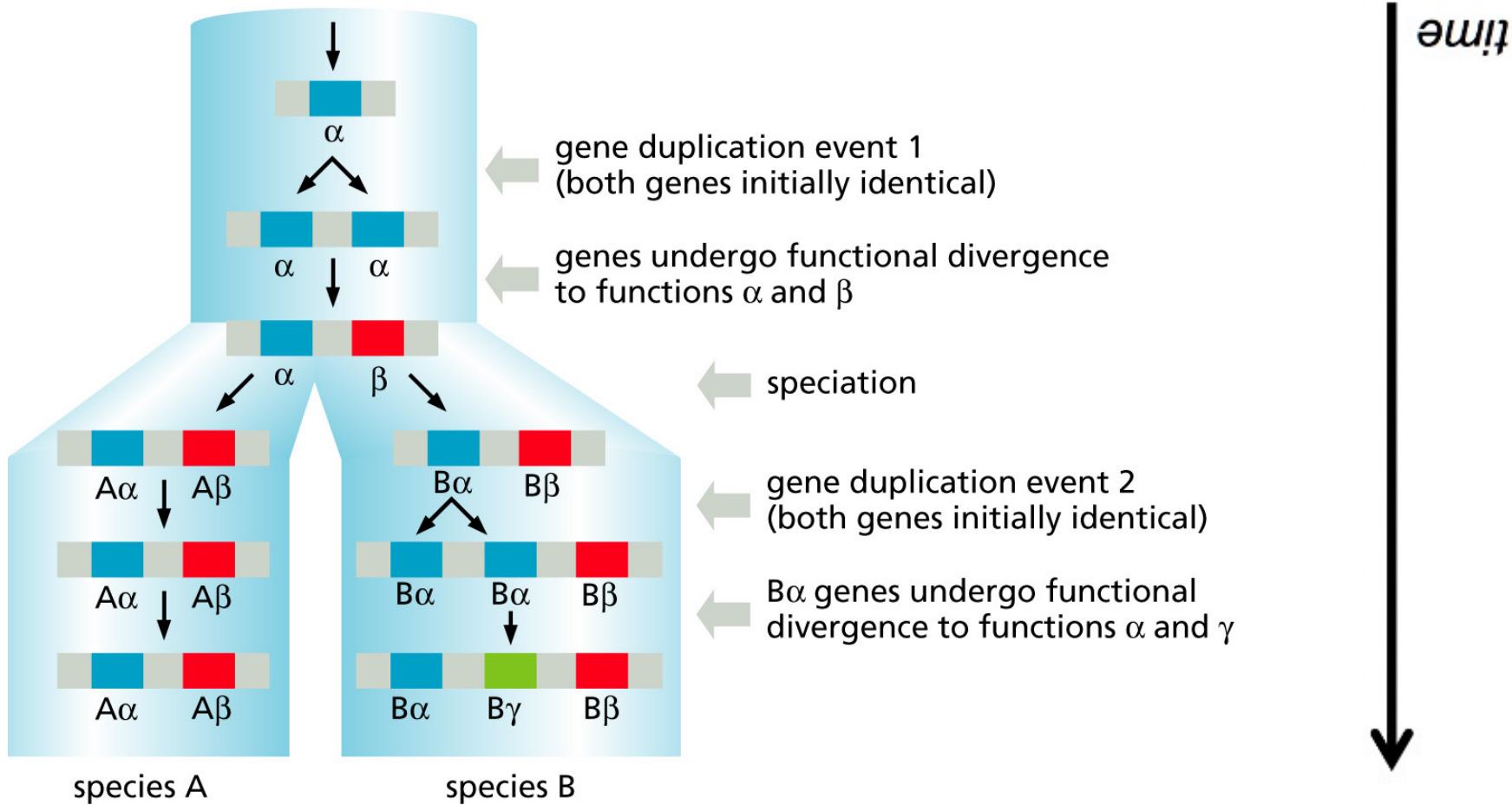
**Paralogy** is a special case of homology, in which the sequences from a common ancestral sequence arose through **duplication**

*Good reads:*

Koonin EV (2001) *Genome Biol* 2(4); <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC138920/>

Jensen RA (2001) *Genome Biol* 2(8); <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC138949/>

(A)



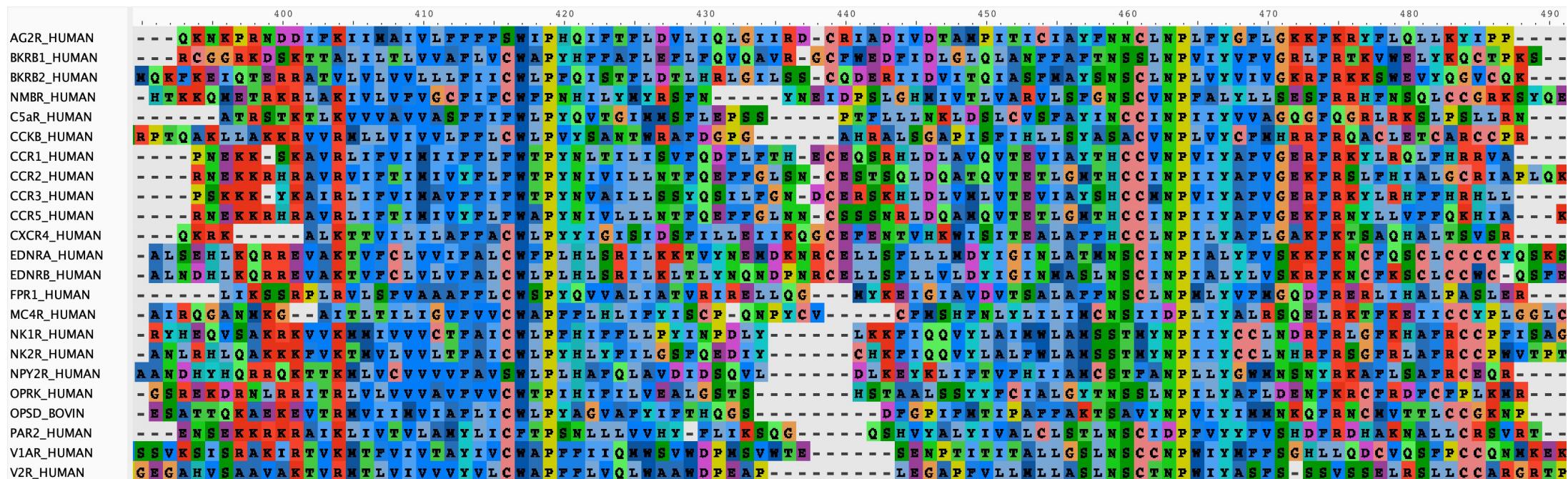
# **Part 2: Phylogenetics introduction**

Sequence alignment

Phylogenetic trees, including terminology

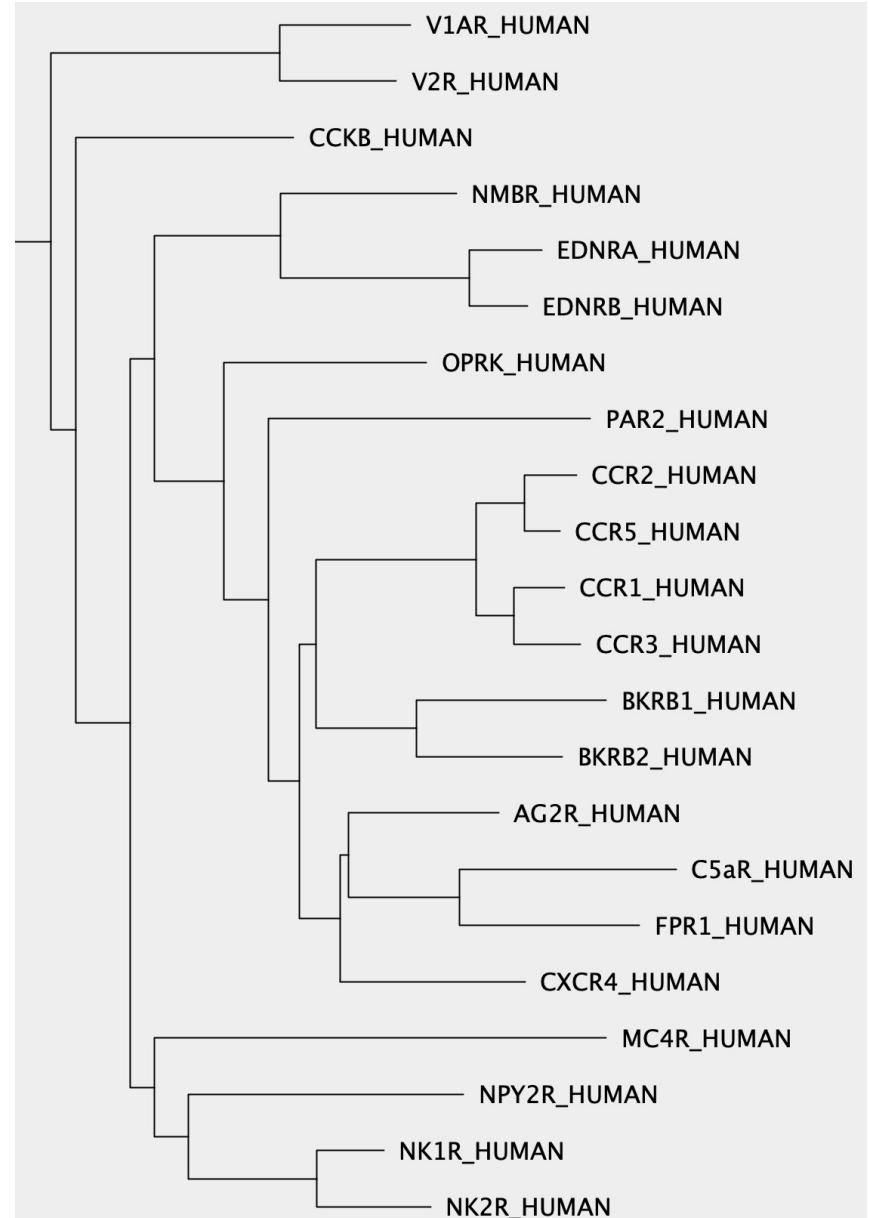
Species v. gene trees

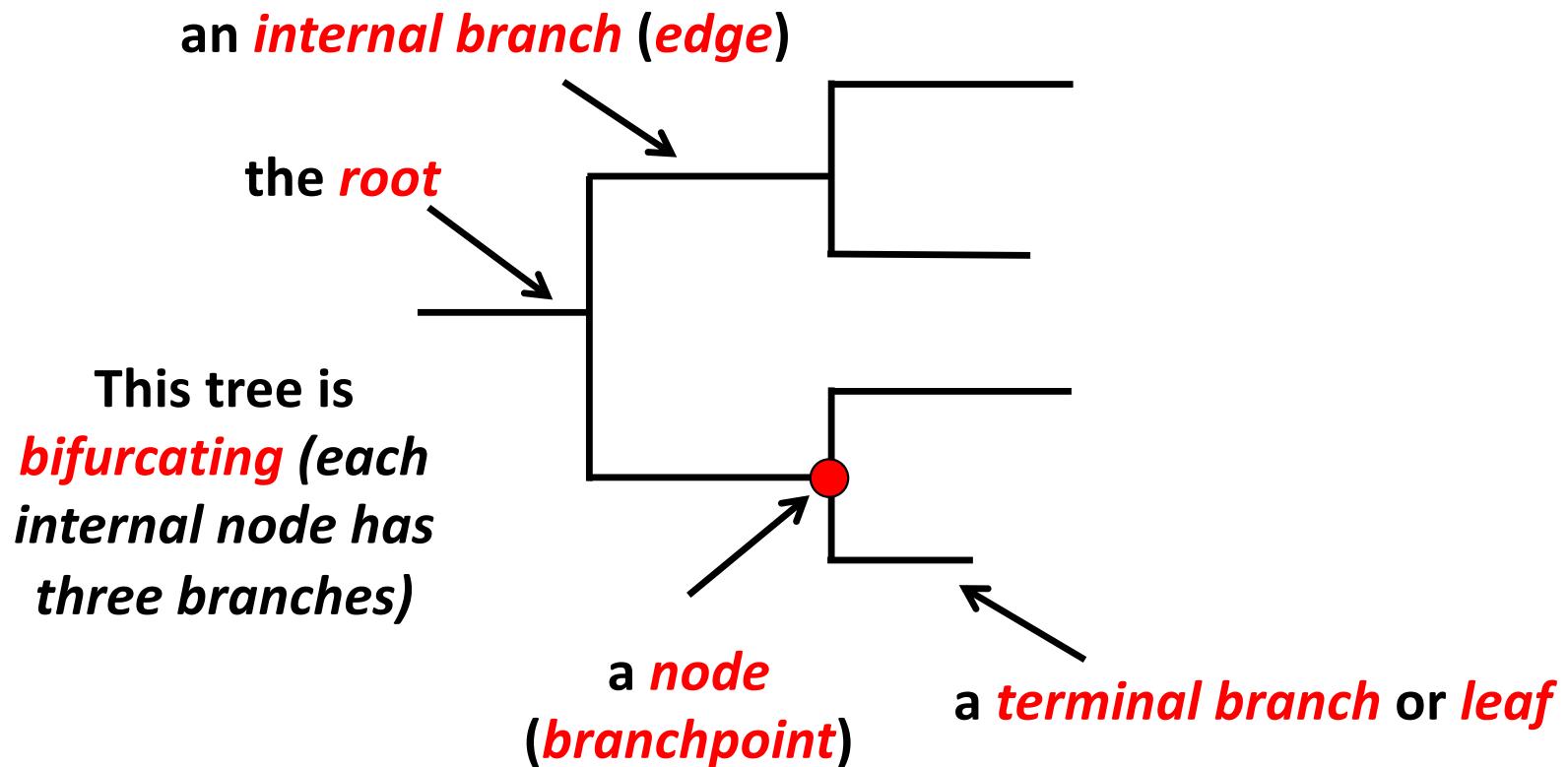
# Multiple sequence alignment is input for phylogenetic analysis

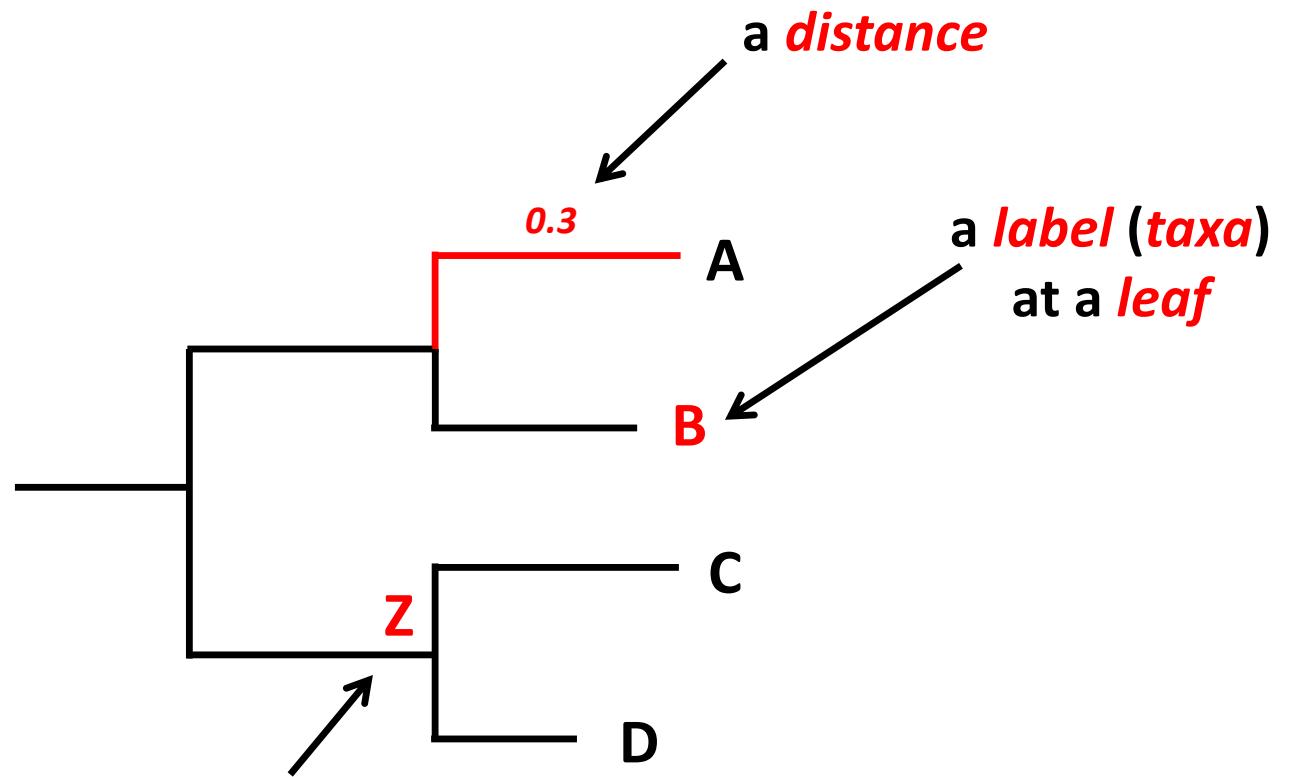


# Trees

- Trees are connected graphs that do not contain circuits
- Trees are important as data structures
- A phylogenetic tree is a hypothesis about the evolutionary relationships between a set of objects (taxa); interpretation is subject to the type of data

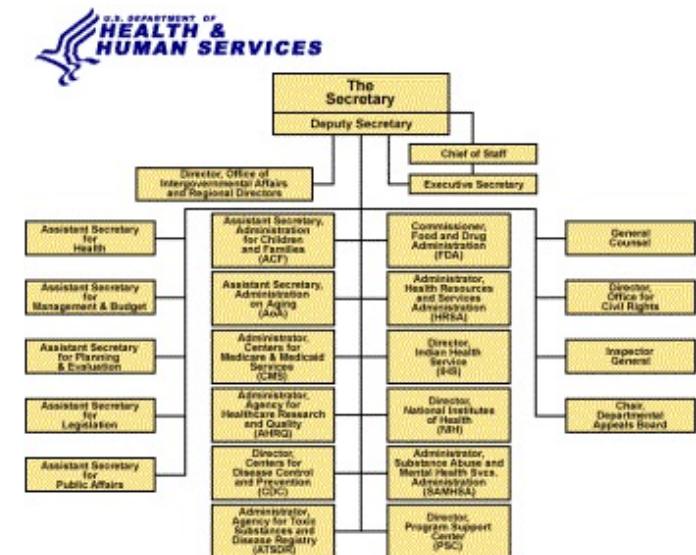
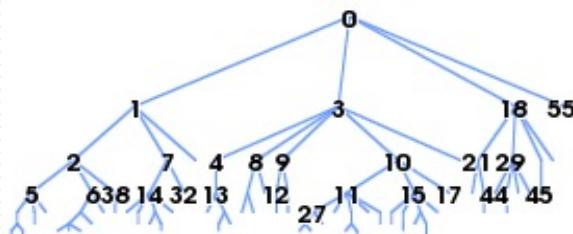
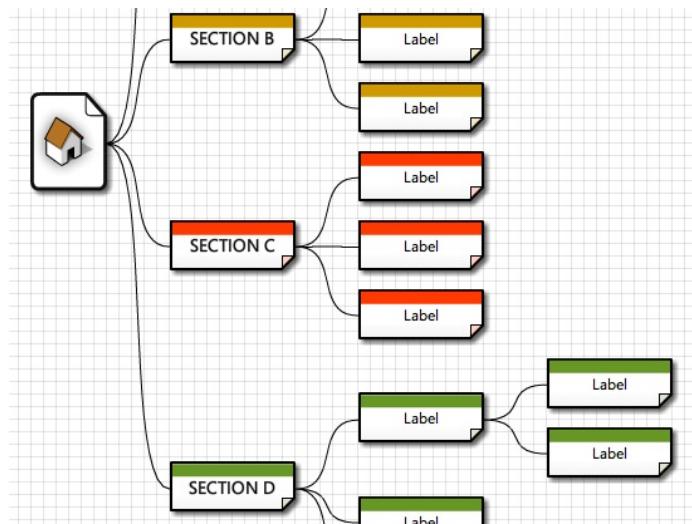
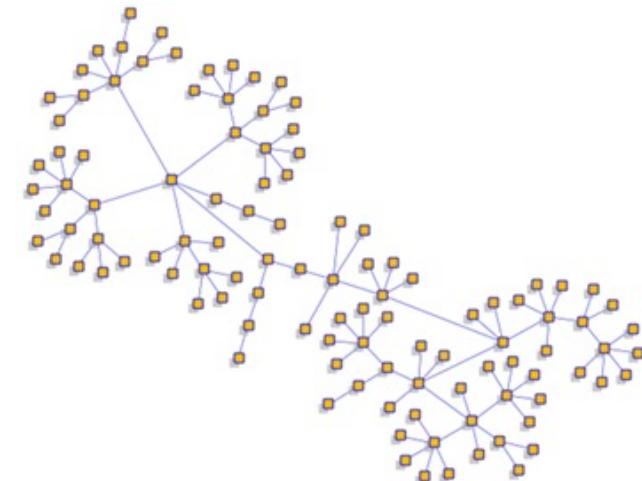
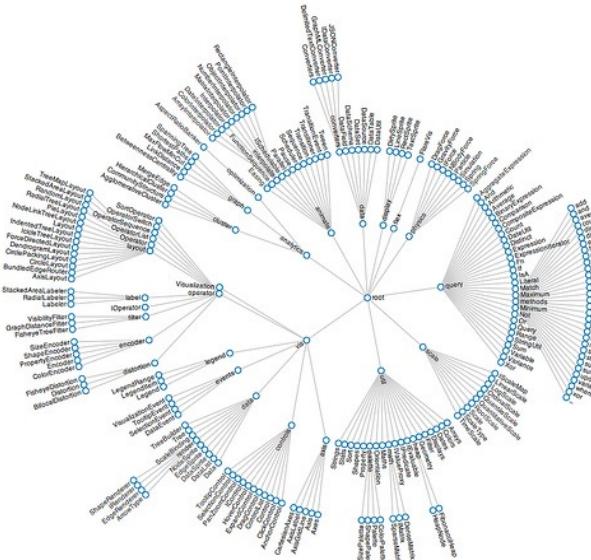
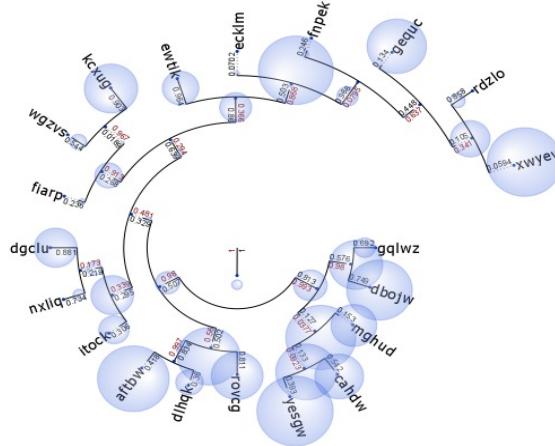






a **label** at an  
ancestor

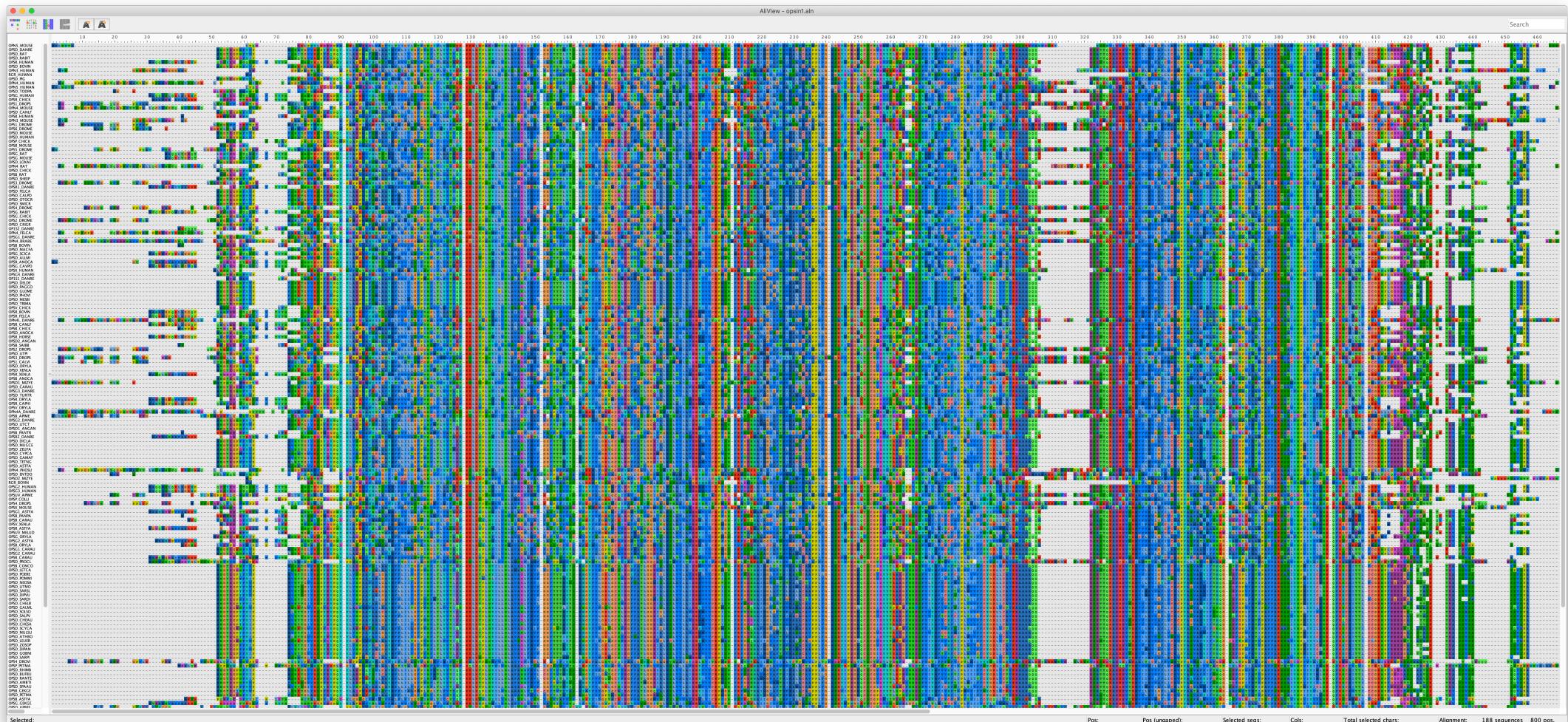
*The vertical lines are NOT edges – they're purely for the viewer's convenience*

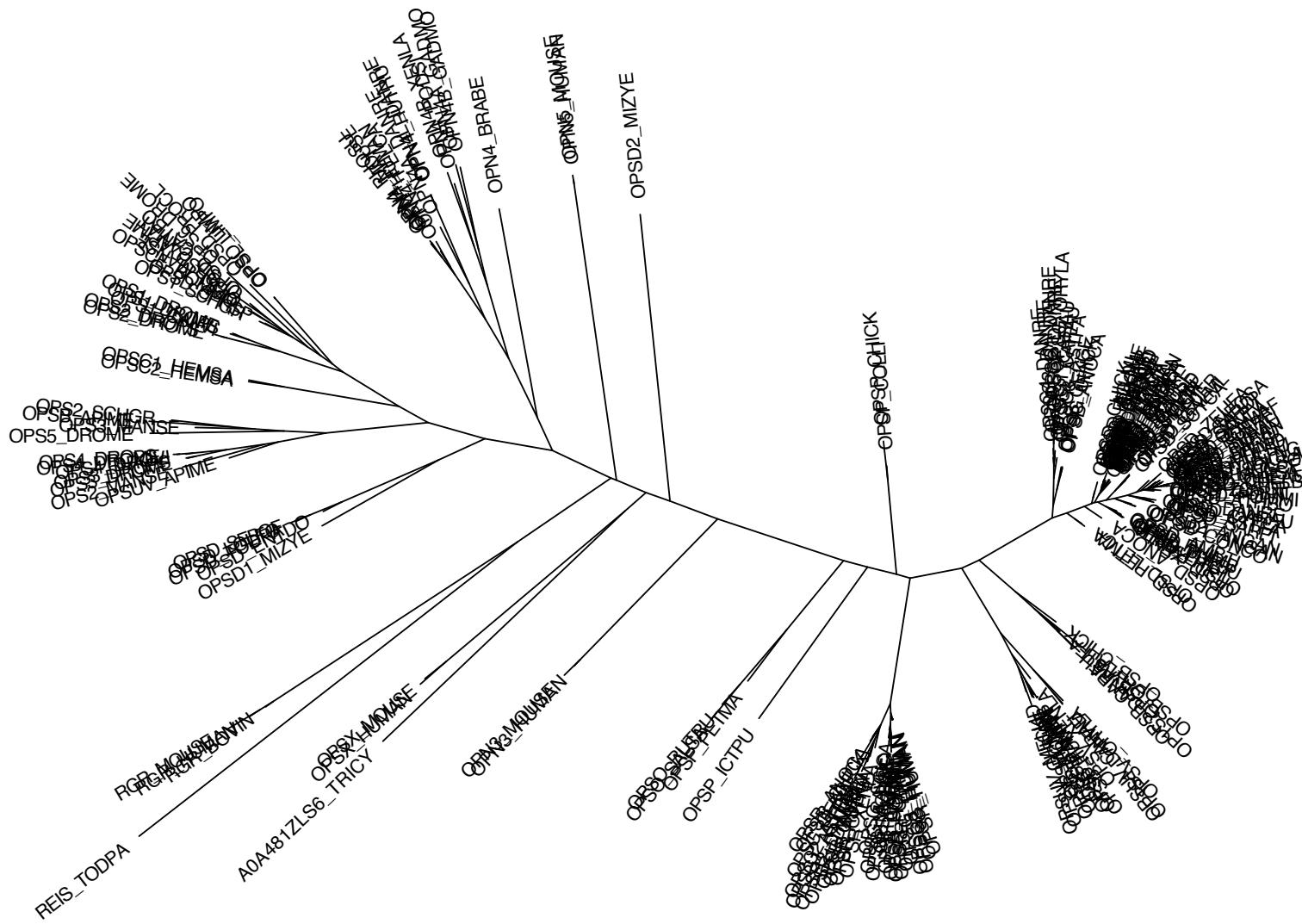


Wikimedia Commons; Mike Bostock; InviVIEW; LovelyCharts; Jaime Huerta-Cepas & Toni Gabaldón; US Dept Health & Human Services

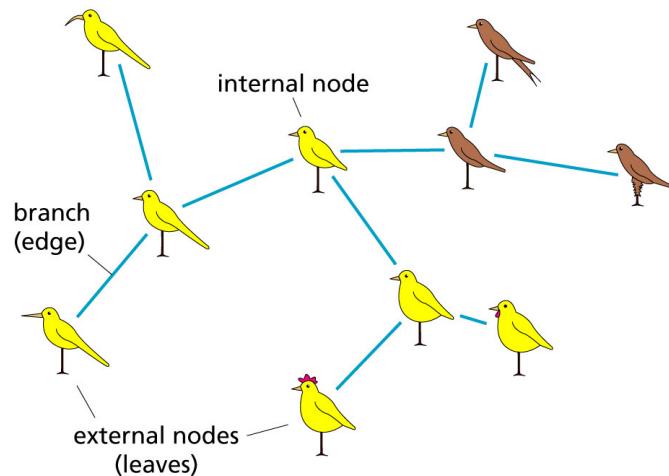
# Opsin is a GPCR-1 sub-family

- Bilateria (16464 results) 
- Chordata (14310 results) 
- Branchiostoma (10 results) 
  - Branchiostoma belcheri (Amphioxus) (5 results)
  - Branchiostoma floridae (Florida lancelet) (Amphioxus) (4 results)
  - Branchiostoma lanceolatum (Common lancelet) (Amphioxus lanceolatum) (1 results)
- Vertebrata (14300 results) 
  - Cyclostomata (jawless vertebrates) (30 results) 
  - Gnathostomata (jawed vertebrates) (14270 results) 
- Protostomia (2154 results) 
- Cnidaria (4 results) 
  - Pocilloporidae (3 results) 
  - Tripedalia cystophora (Jellyfish) (1 results)



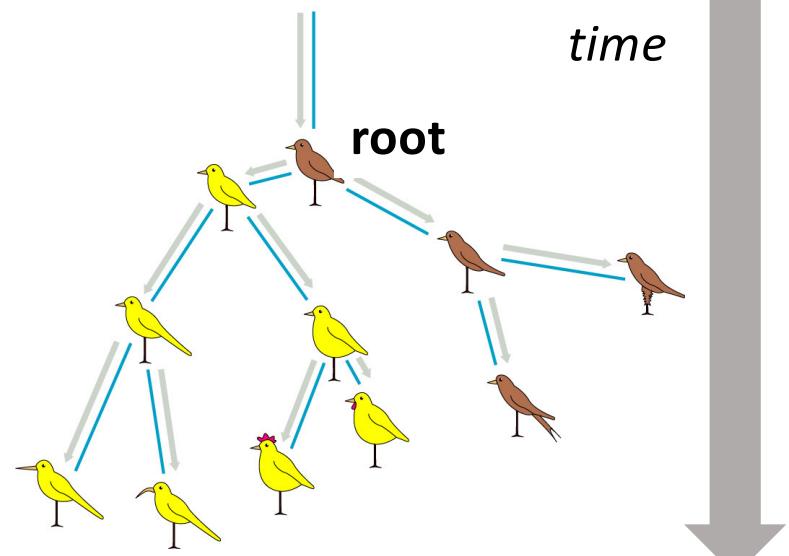


## Unrooted

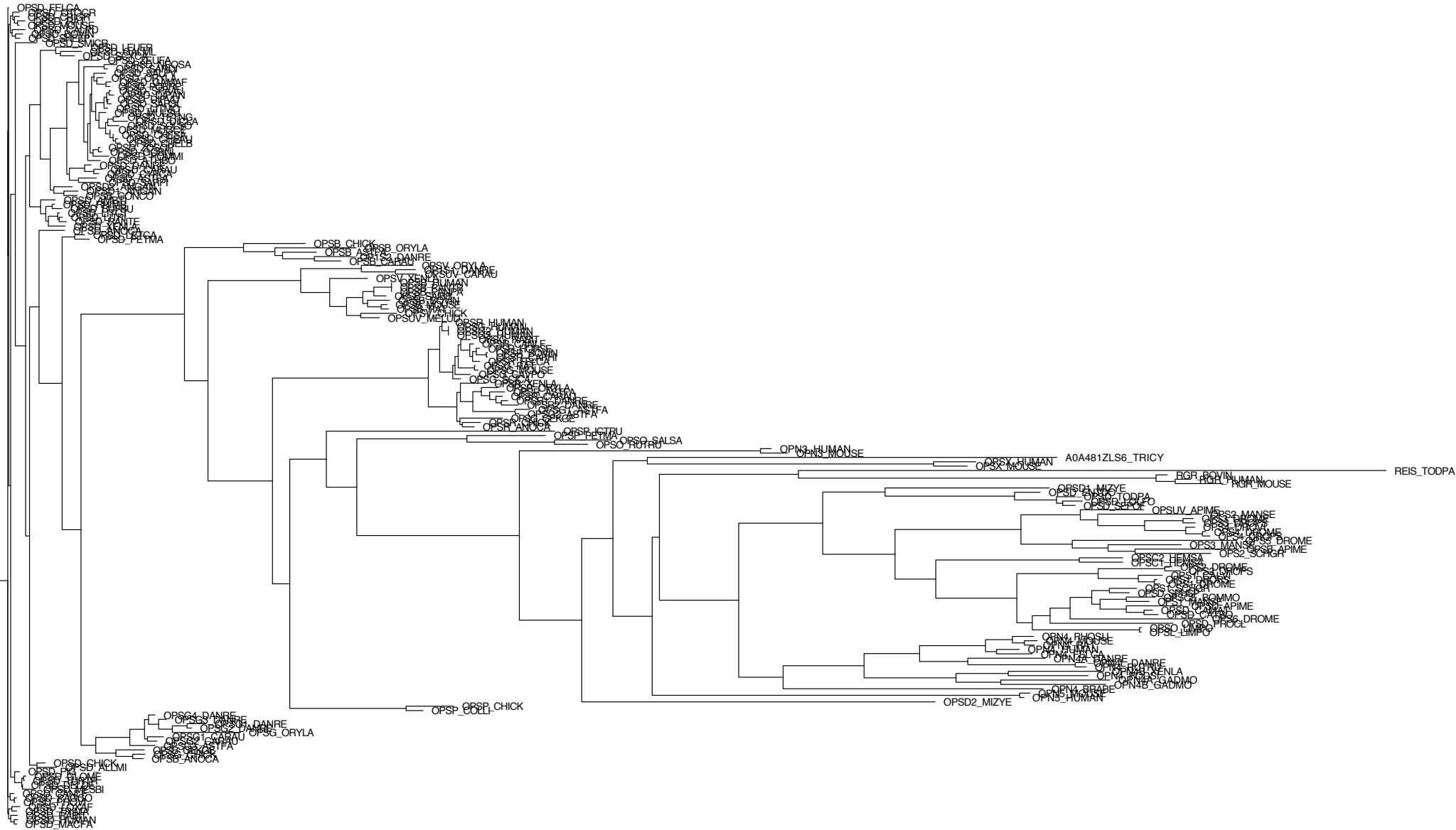


Distances are  
undirected

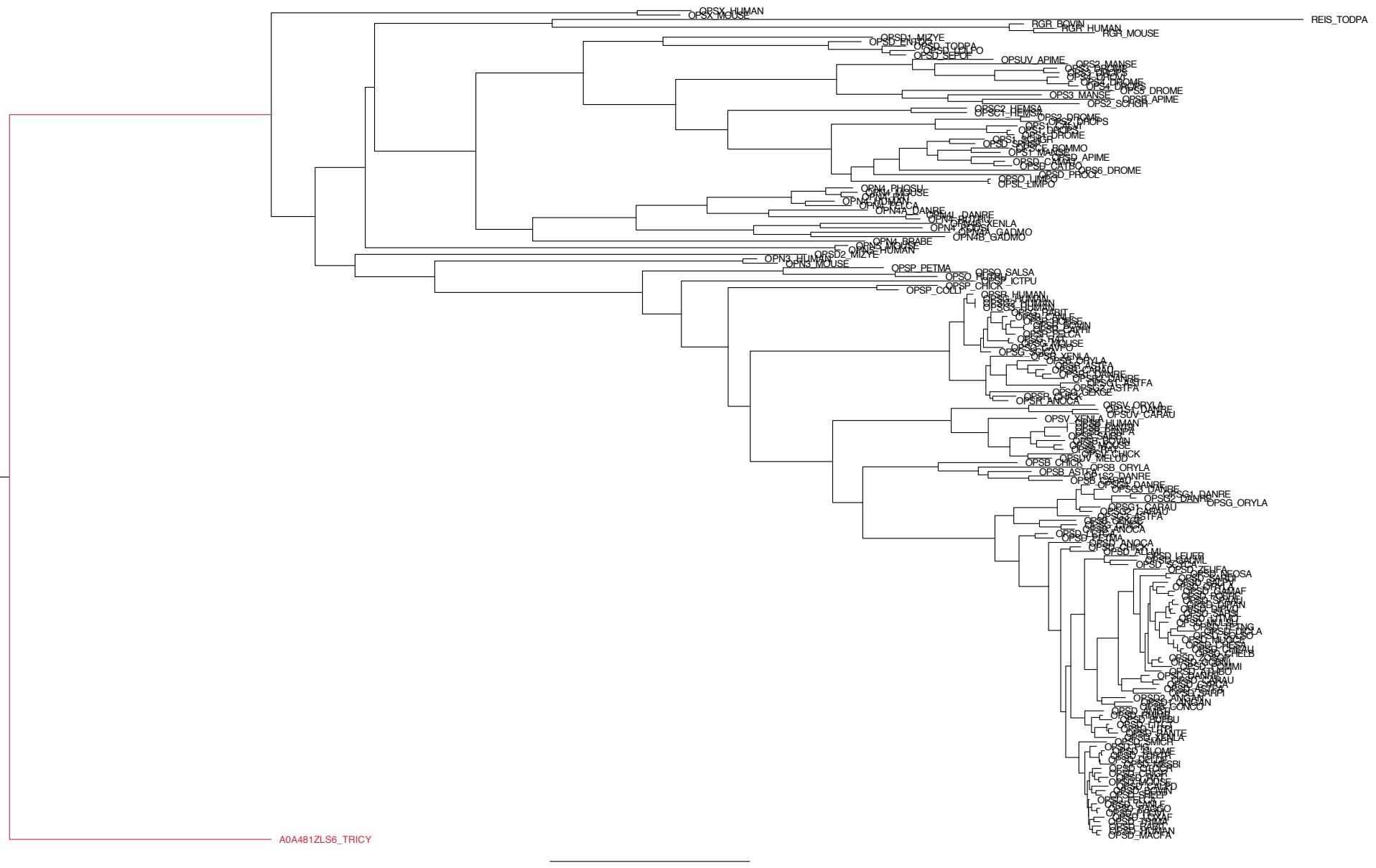
## Rooted

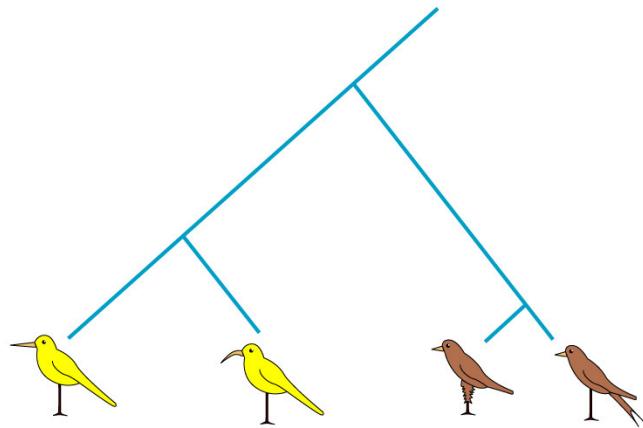


Distances are  
directed



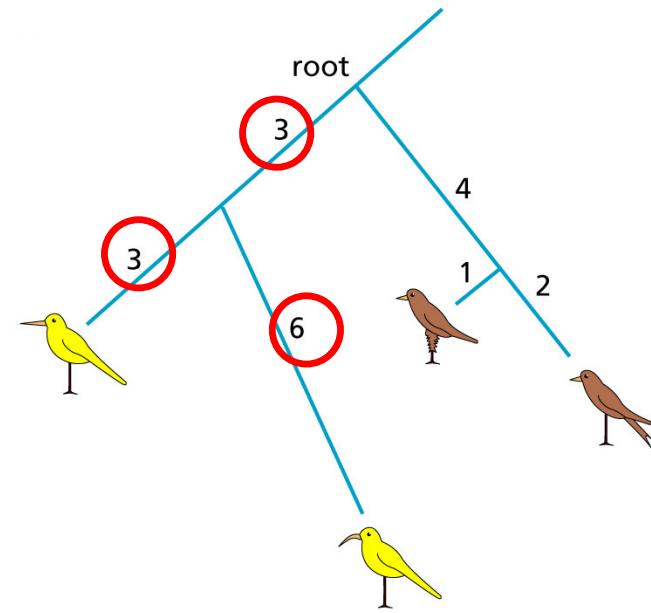
0.





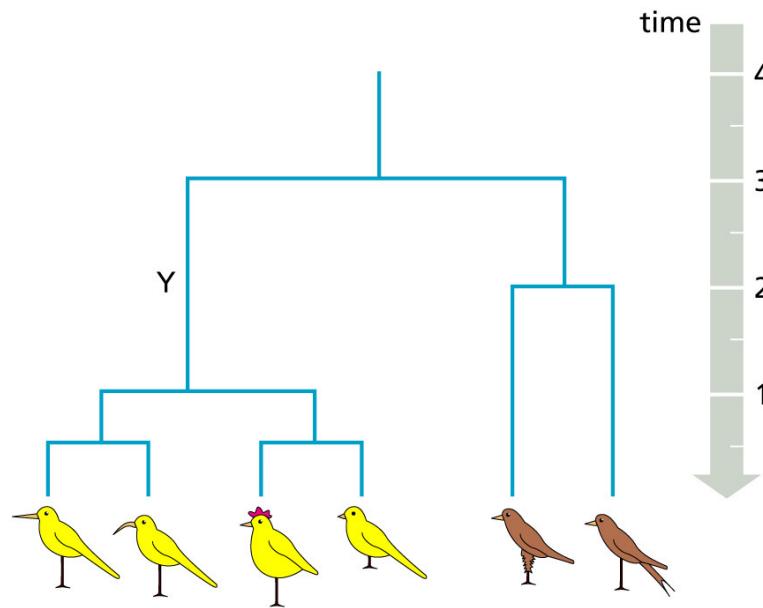
## Cladogram

branch lengths have *qualitative*  
but not *quantitative* meaning

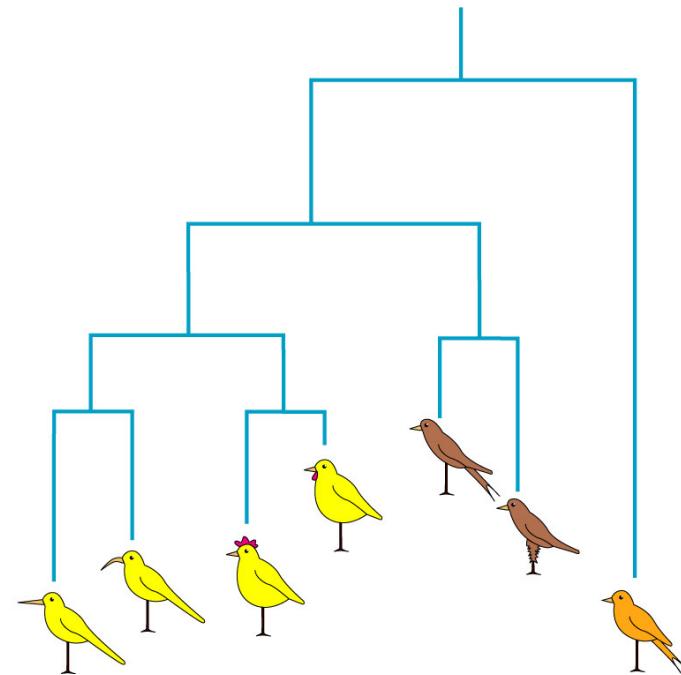


## Additive tree

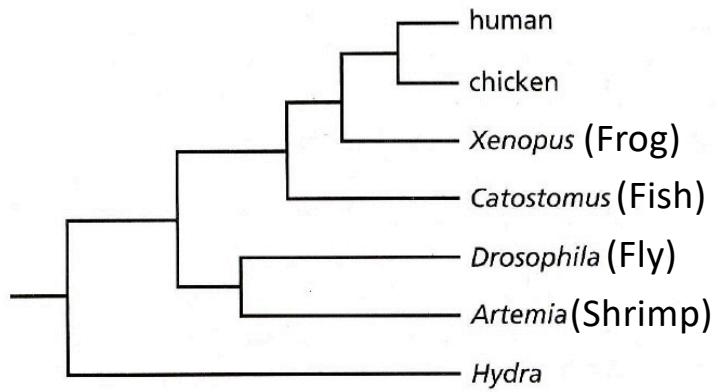
branches quantify the evolutionary  
distance, the progress of divergence



**Ultrametric tree**  
(constant rate of mutation;  
*molecular clock*)



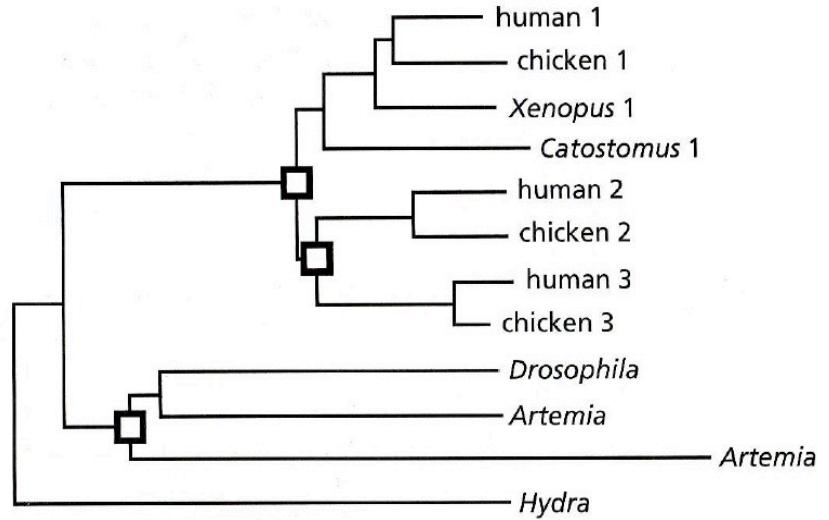
**Additive tree with  
outgroup**



### Species tree

(evolution of species)

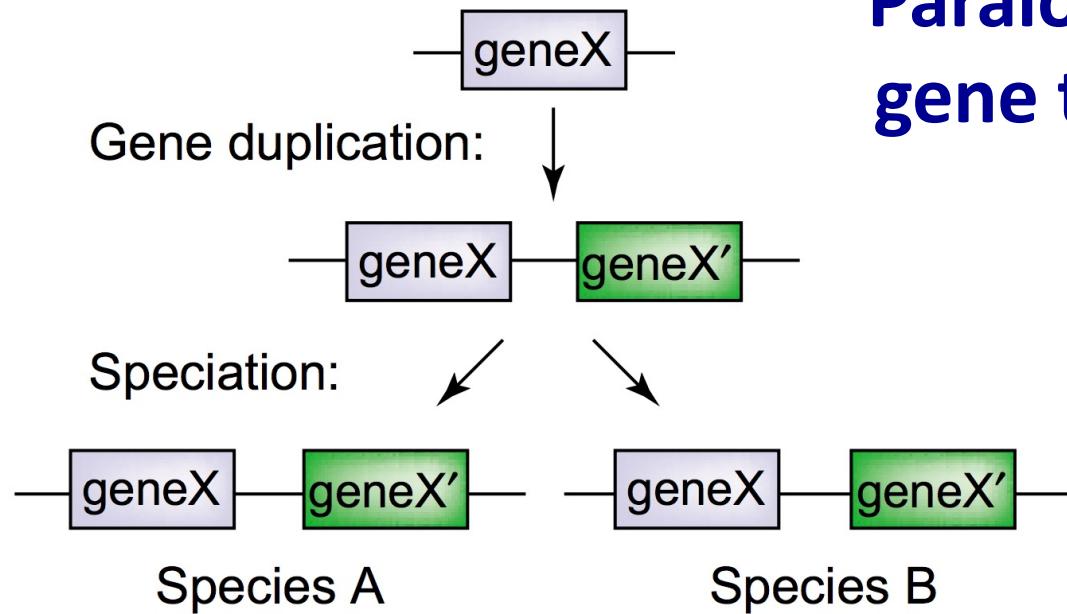
Ideally derived from **strictly orthologous** sequences



### Gene tree

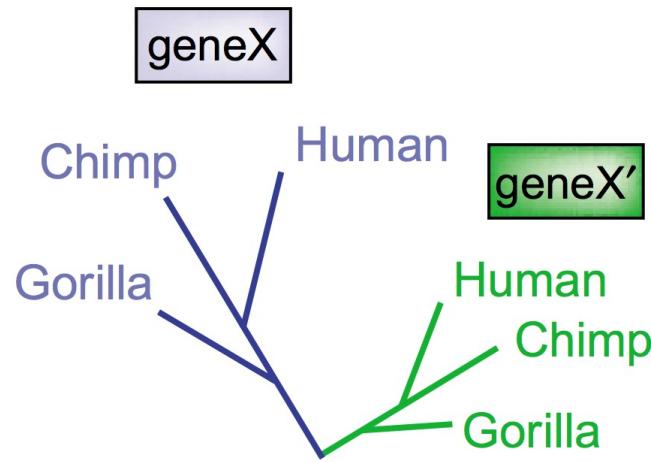
(evolution of homologous genes,  
in this case a particular family of  
membrane proteins)

## Paralogs in gene trees

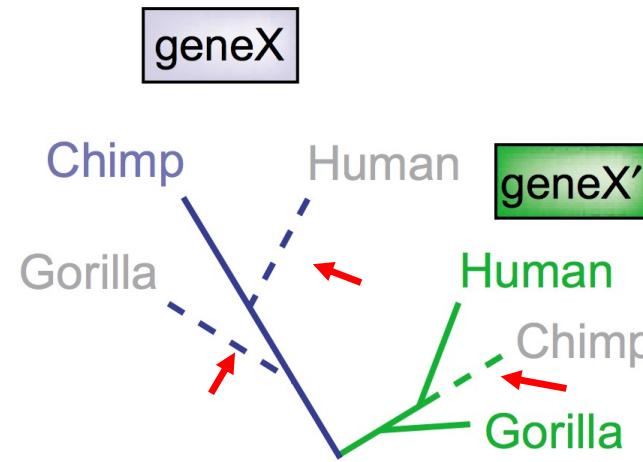


Gene X is duplicated prior to speciation. All subsequent species inherit both copies of the gene (unless one or the other is lost somewhere along the way).

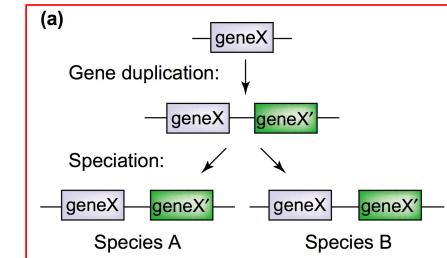
Baldauf, SL (2003) *Trends in Genetics* 19: 345-351.



All sequences of gene X are **orthologs** of each other, and all the sequences of gene X' are **orthologs** of each other. However, X and X' are **paralogs**. Both the X and X' subtrees show the true relationships among the three species.



A tree of the X/X' gene family can be misleading if not all the sequences are included (because of **incomplete sampling** or **gene loss**). If the broken branches are missing, then the true species relationships are misrepresented.



# **Read more about this?**

- Some of the material is based on Zvelebil and Baum's textbook “Understanding Bioinformatics” (chapters 7 and 8)