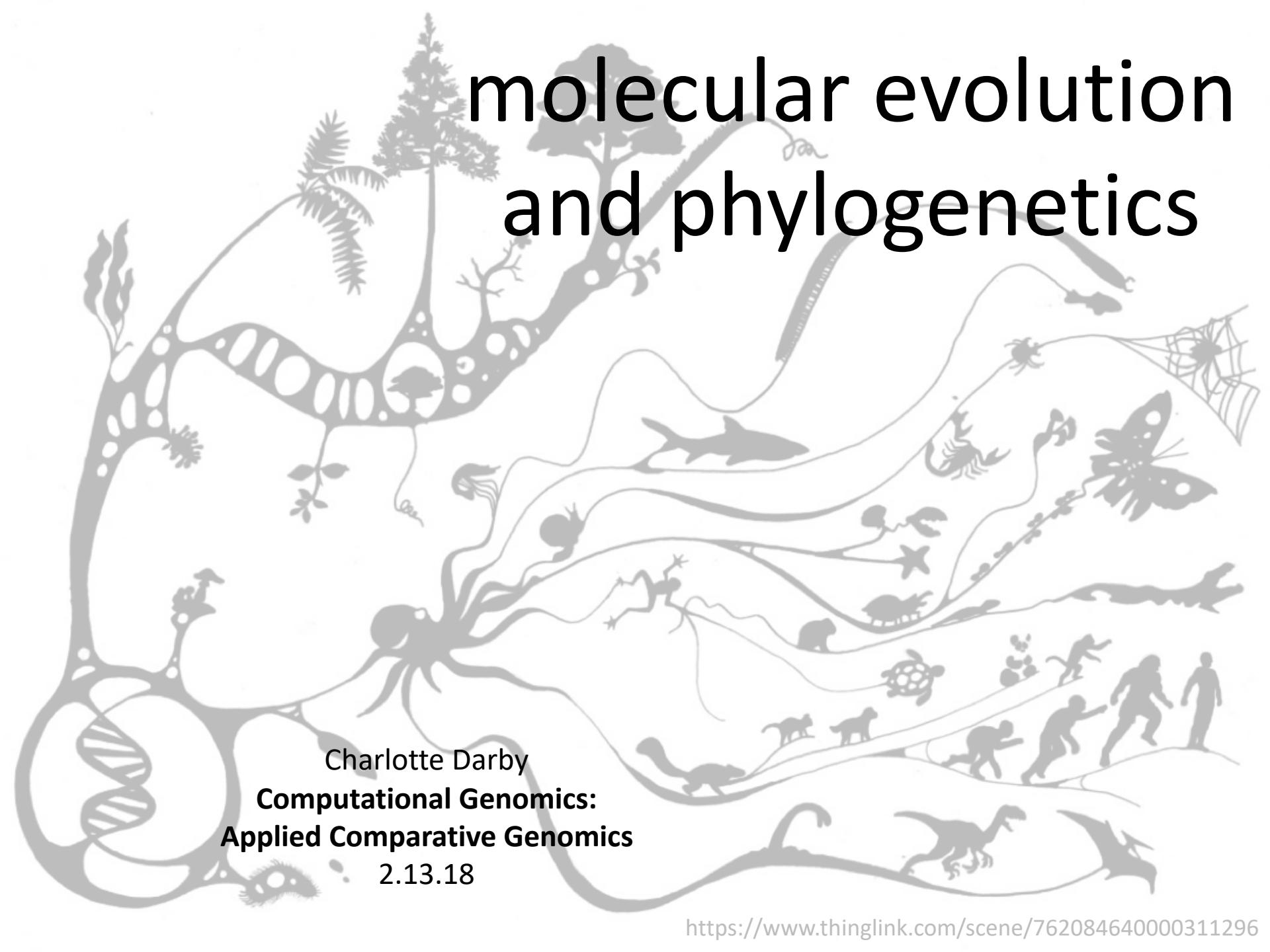
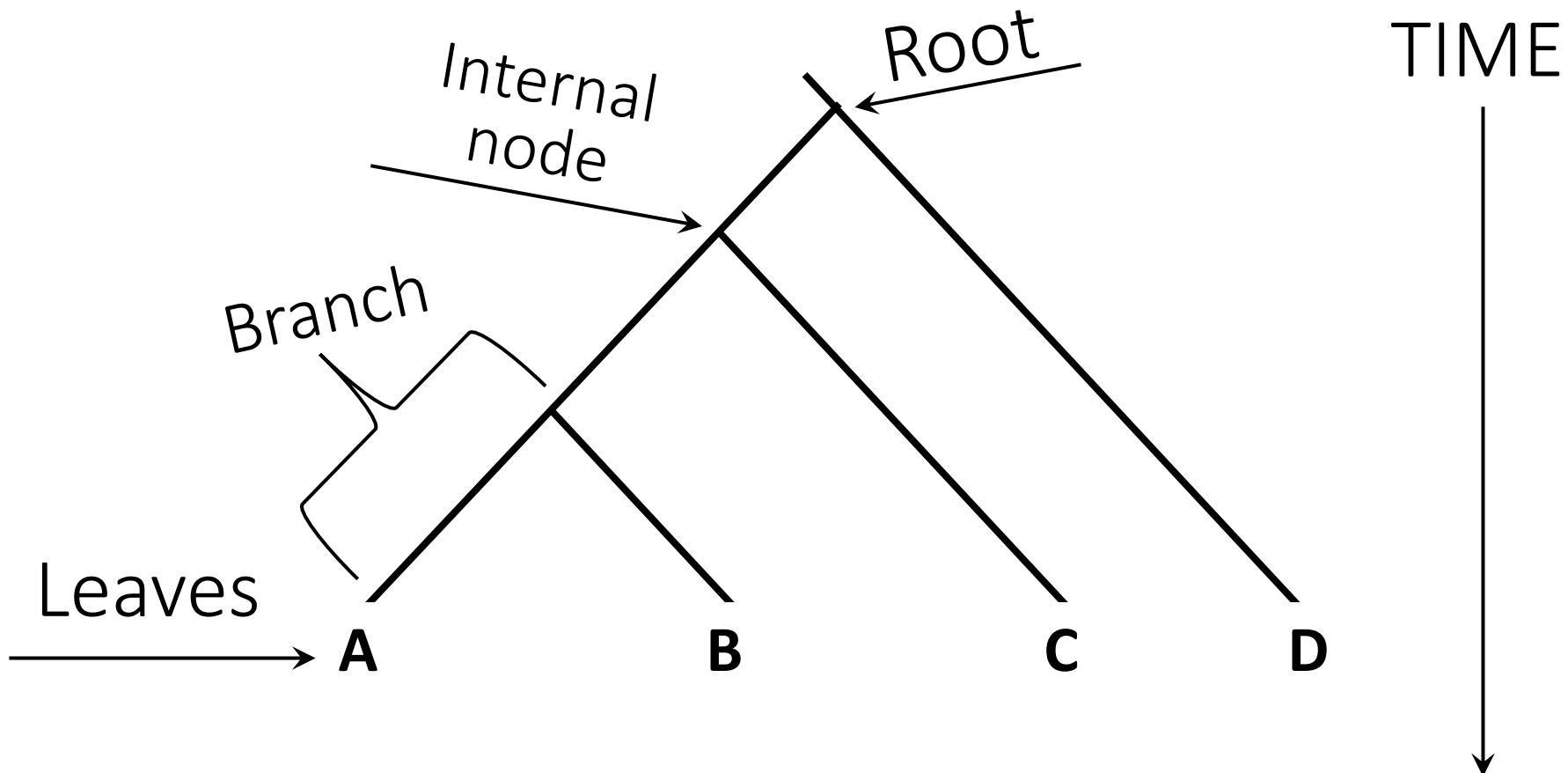


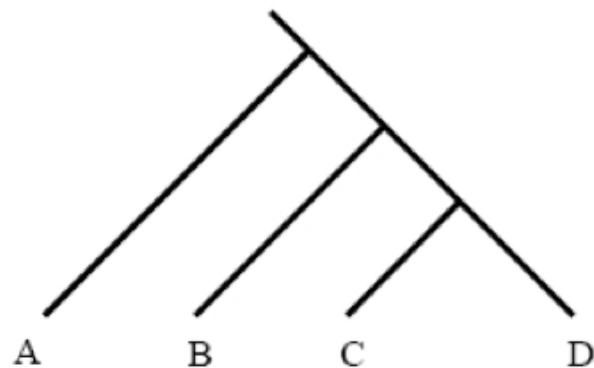
molecular evolution and phylogenetics



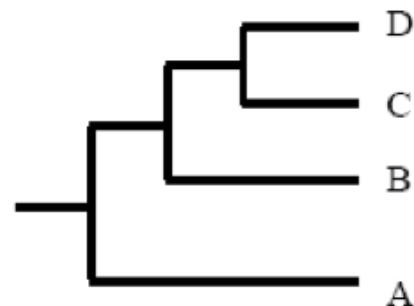
Charlotte Darby
Computational Genomics:
Applied Comparative Genomics
2.13.18



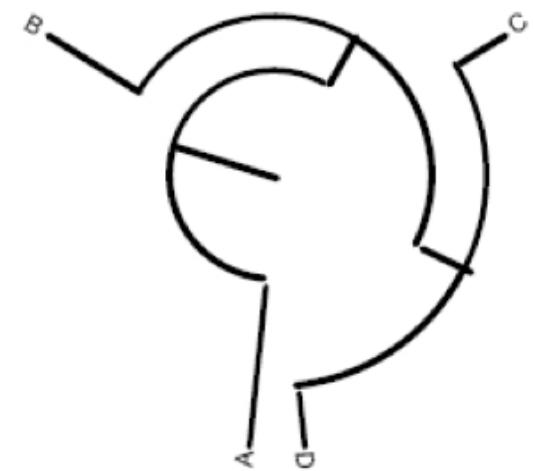
computer scientist



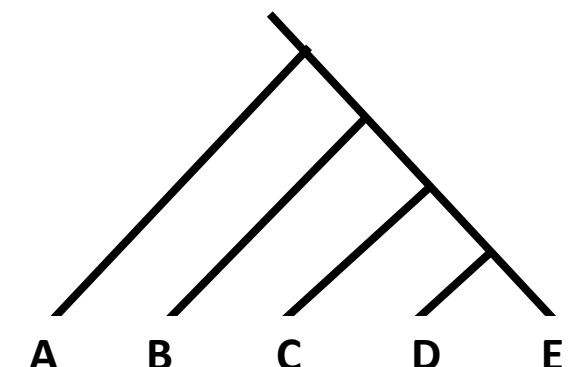
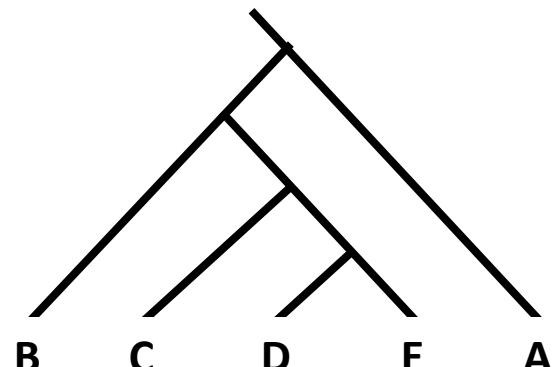
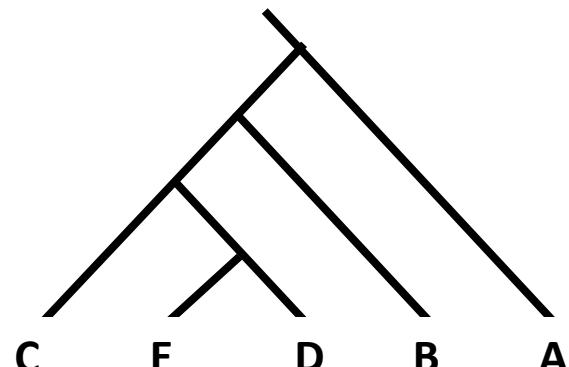
biologist

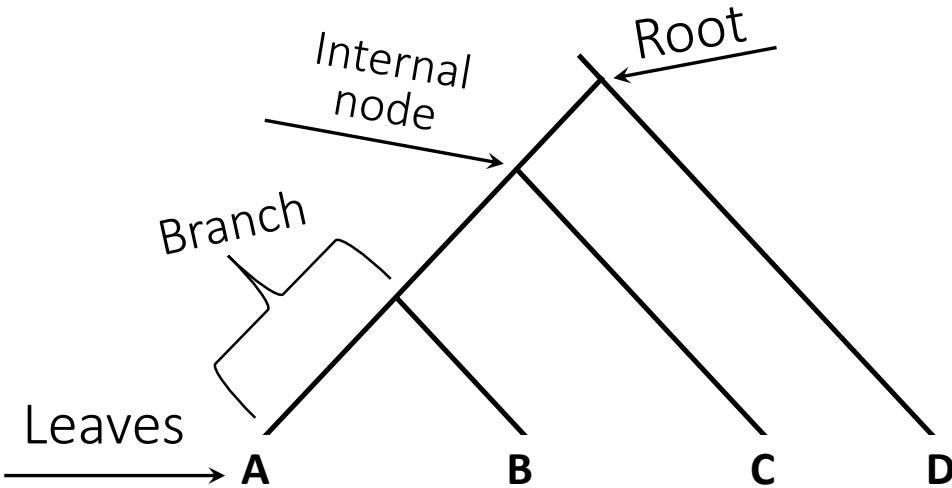


???



Each set of three figures shows the same topological relationship among the leaves



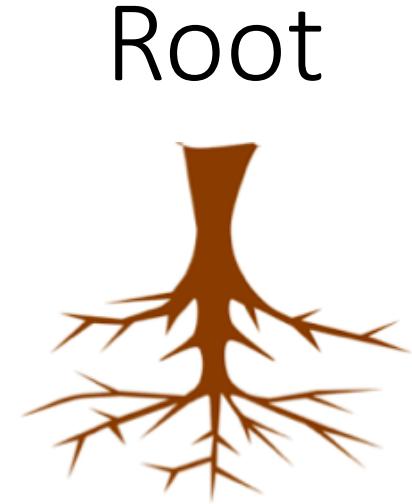
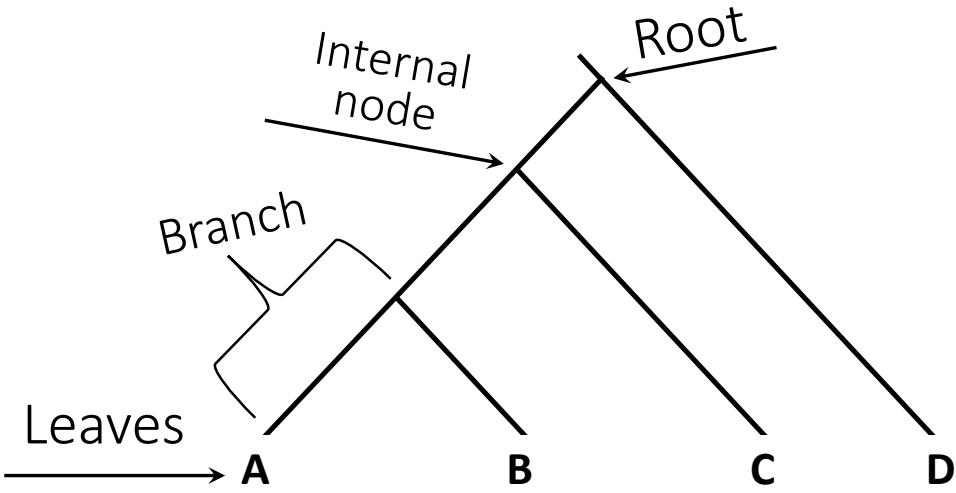


Leaves

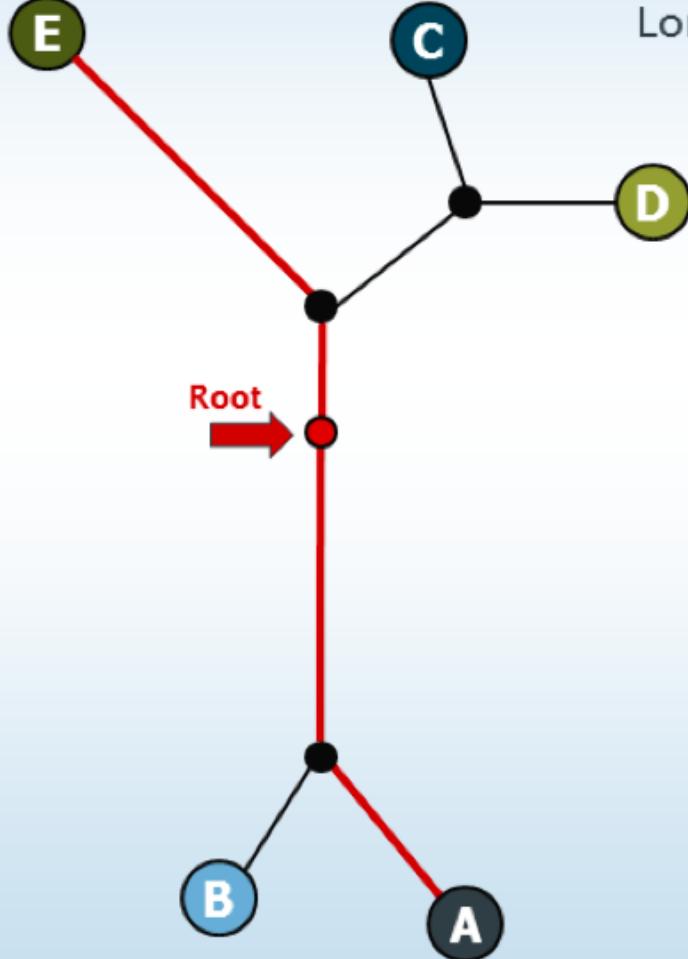


- Strains or isolates of the same species
- OR different species
- OR DNA or protein sequences
 - From the same genome or different genomes

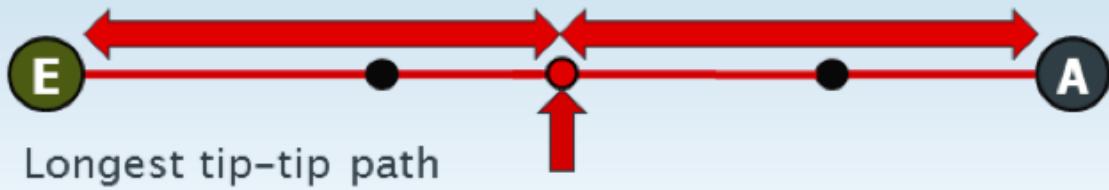
Leaves are implied to be extant
(present at the current time)



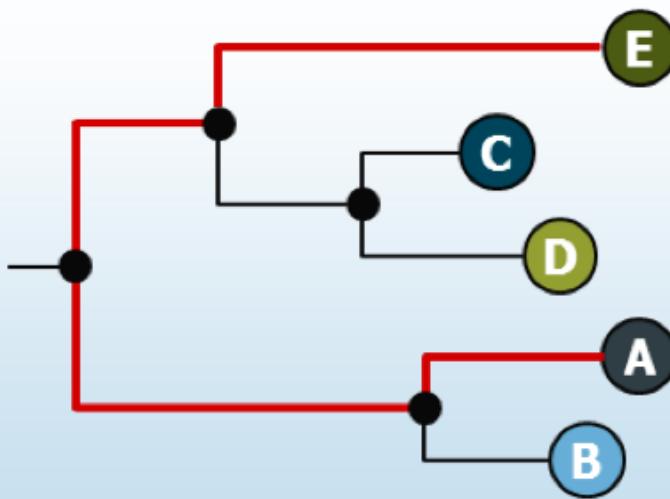
- Implies what happened first in chronological time
- Tree-building algorithms may not infer the root of a tree
 - Represented as a node with 3 outgoing branches where the root ought to be
- Could root halfway along the longest path
 - midpoint rooting
- Could root by prior knowledge of an outgroup
 - e.g. bacteria in the tree above



Unrooted

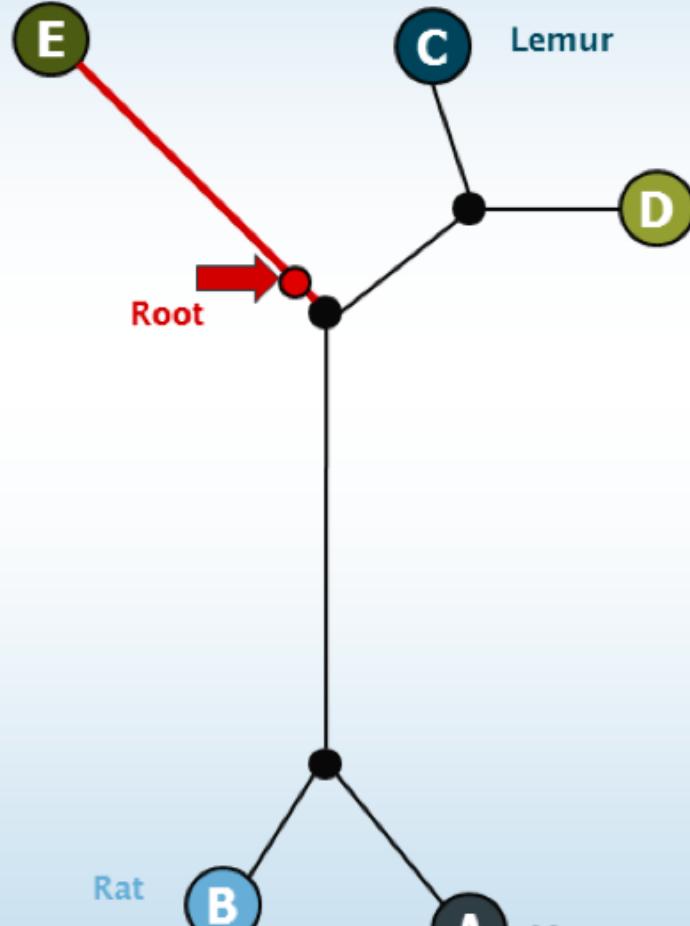


Longest tip-tip path



Mid-point
Rooted

Kangaroo



Root

Lemur

Human

Rat

B

Mouse

Unrooted

Outgroup Rooted

E

Kangaroo

C

Lemur

D

Human

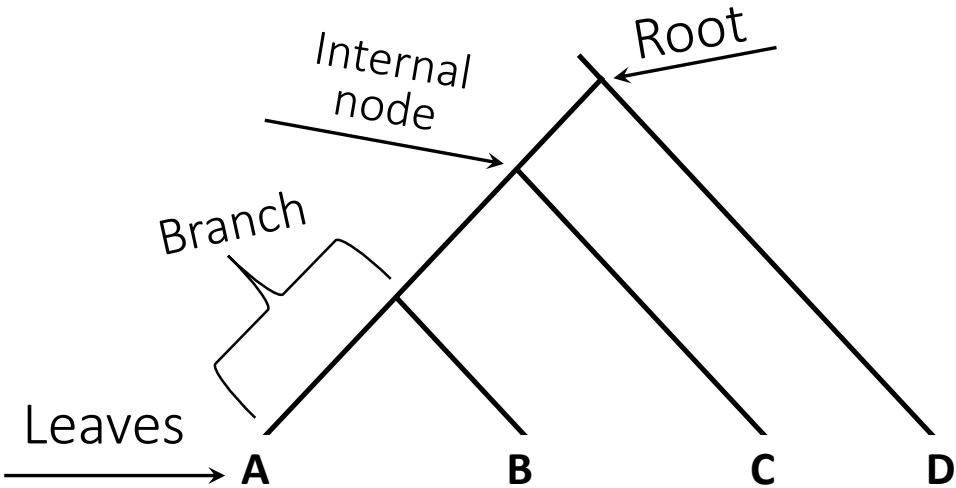
A

Mouse

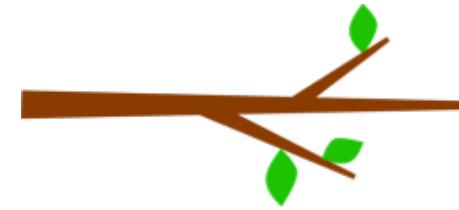
B

Rat

Mid-point
Rooted

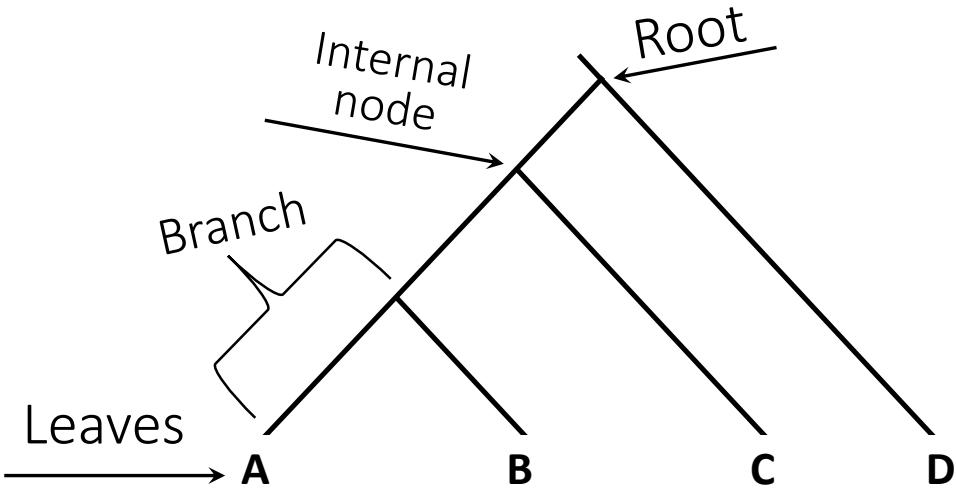


Internal node

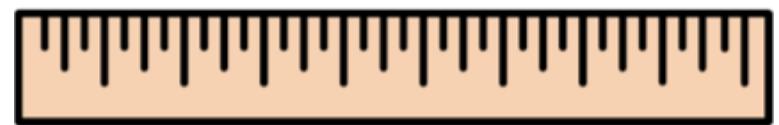


- Represents some ancestral state
- Most recent common ancestor (MRCA) of the leaves following it

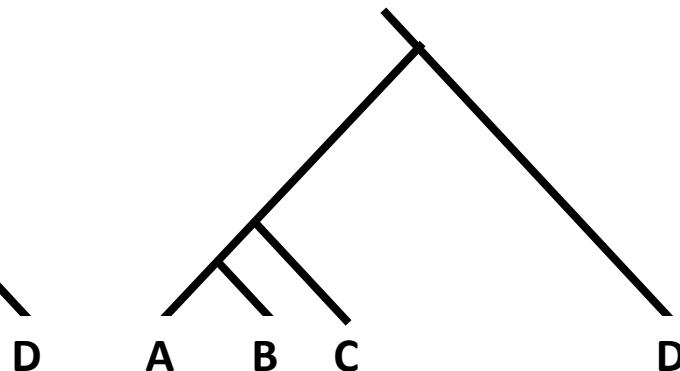
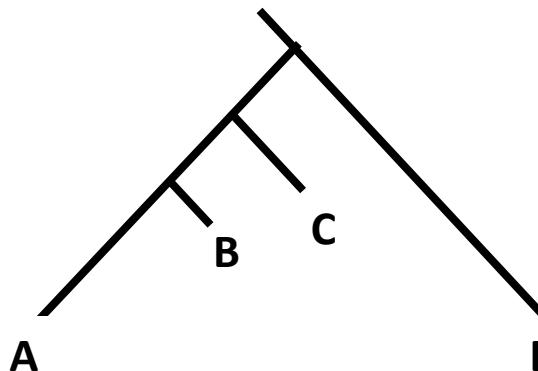
Internal nodes are implied to NOT be extant



Branch length



- Chronological time
- OR number of changes to the DNA / protein sequence
 - Depending on assumptions you make about rate of evolution when building tree, may not be ultrametric (every root-leaf path has same length)
- OR no distance implied – just defines a branching order



Why visualize biological information as trees?

(besides having a colorful Figure 1 for your Nature paper!)

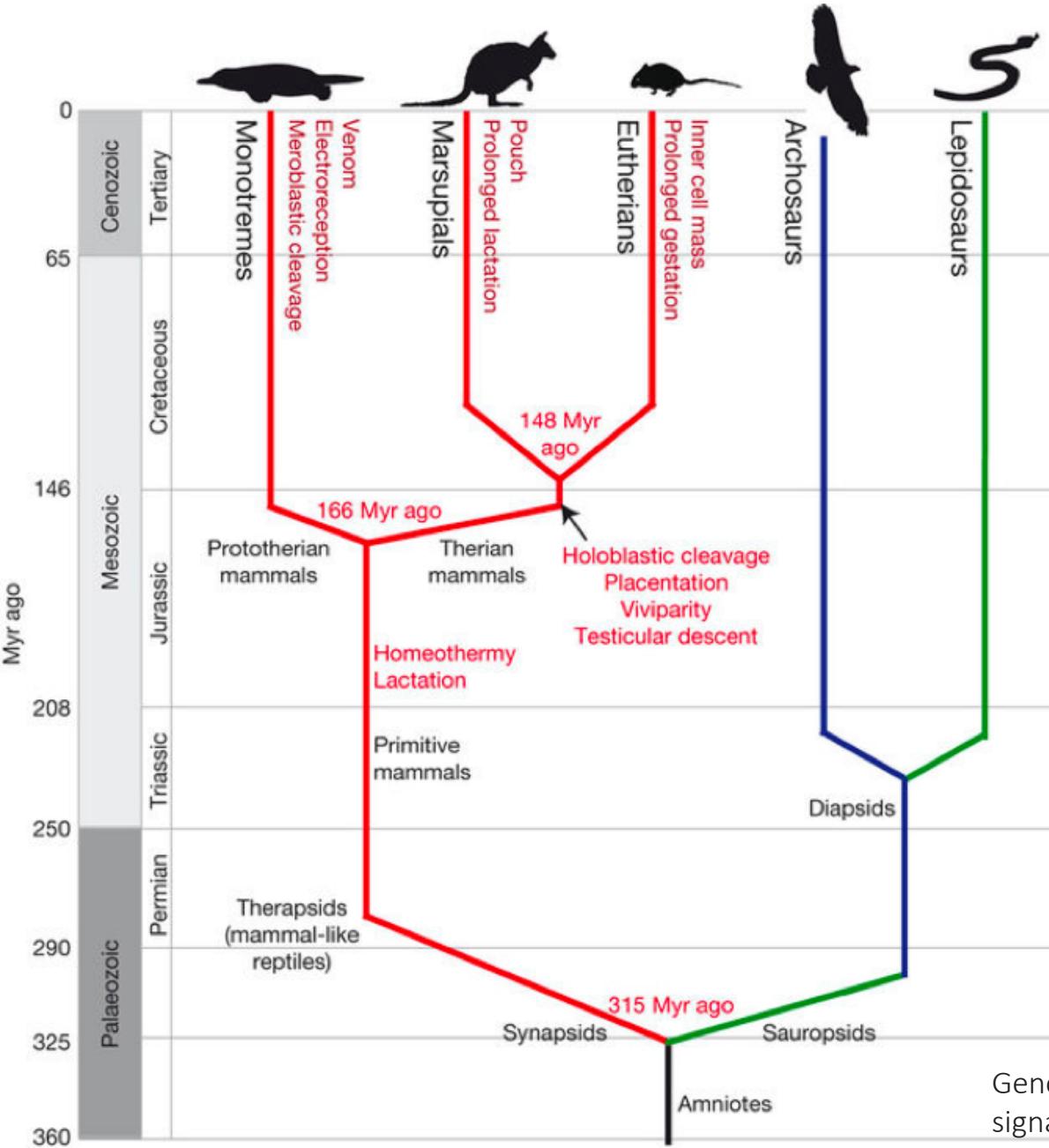
Trees show fundamental structure in your data.

A picture's worth a thousand words!

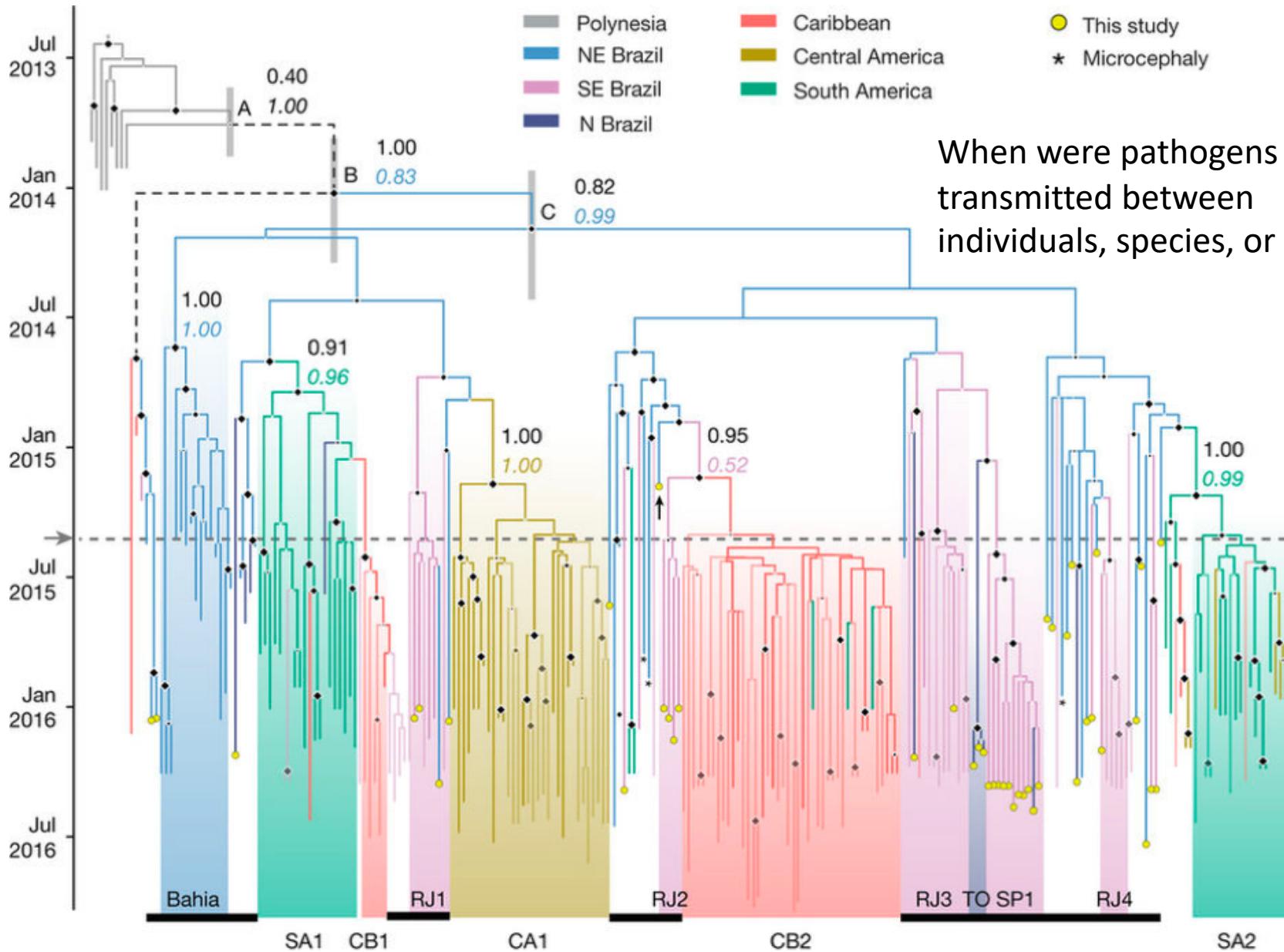
The last amniote common ancestor of synapsids and sauropsids is placed 315 MYA. The sauropsid lineage evolved into archosaurs (such as birds, crocodilians, and dinosaurs) and lepidosaurs (such as snakes and lizards), diverging in the Triassic period. In the synapsid lineage, these primitive mammals acquired homothermy and lactation before the divergence of protetherian and therian mammals 166 MYA in the Jurassic period.

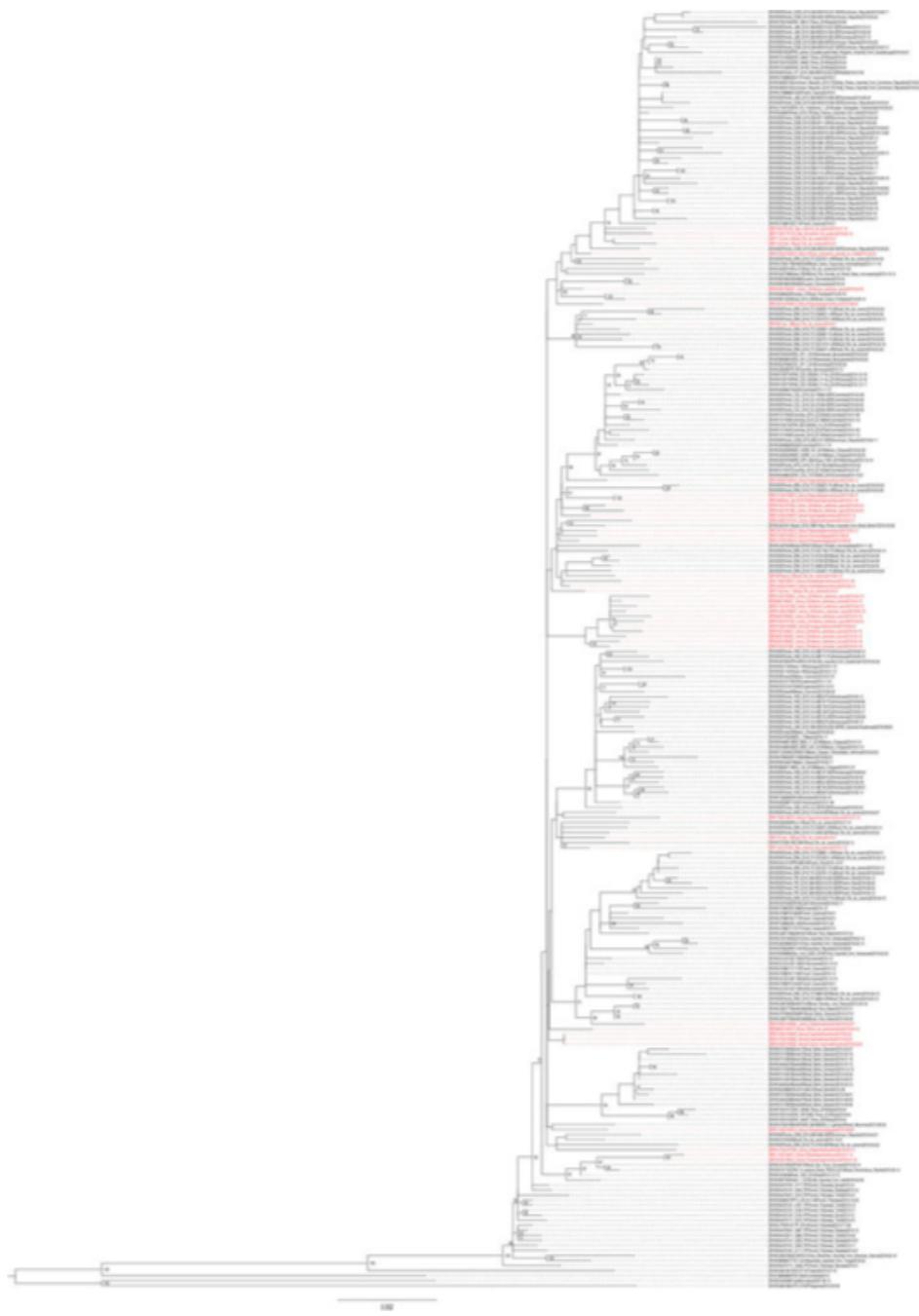
While descendant species of the protetherian mammals - extant monotremes such as the platypus *Ornithorhynchus anatinus* - have characteristics such as venom, electroreception, and meroblastic cleavage, therian mammals evolved holoblastic cleavage, placentation, viviparity, and testicular descent before their divergence 148 MYA into marsupials and eutherians. Further diversification leads us to observe a pouch and prolonged lactation in extant marsupials, and inner cell mass and prolonged gestation in eutherians.

species phylogeny



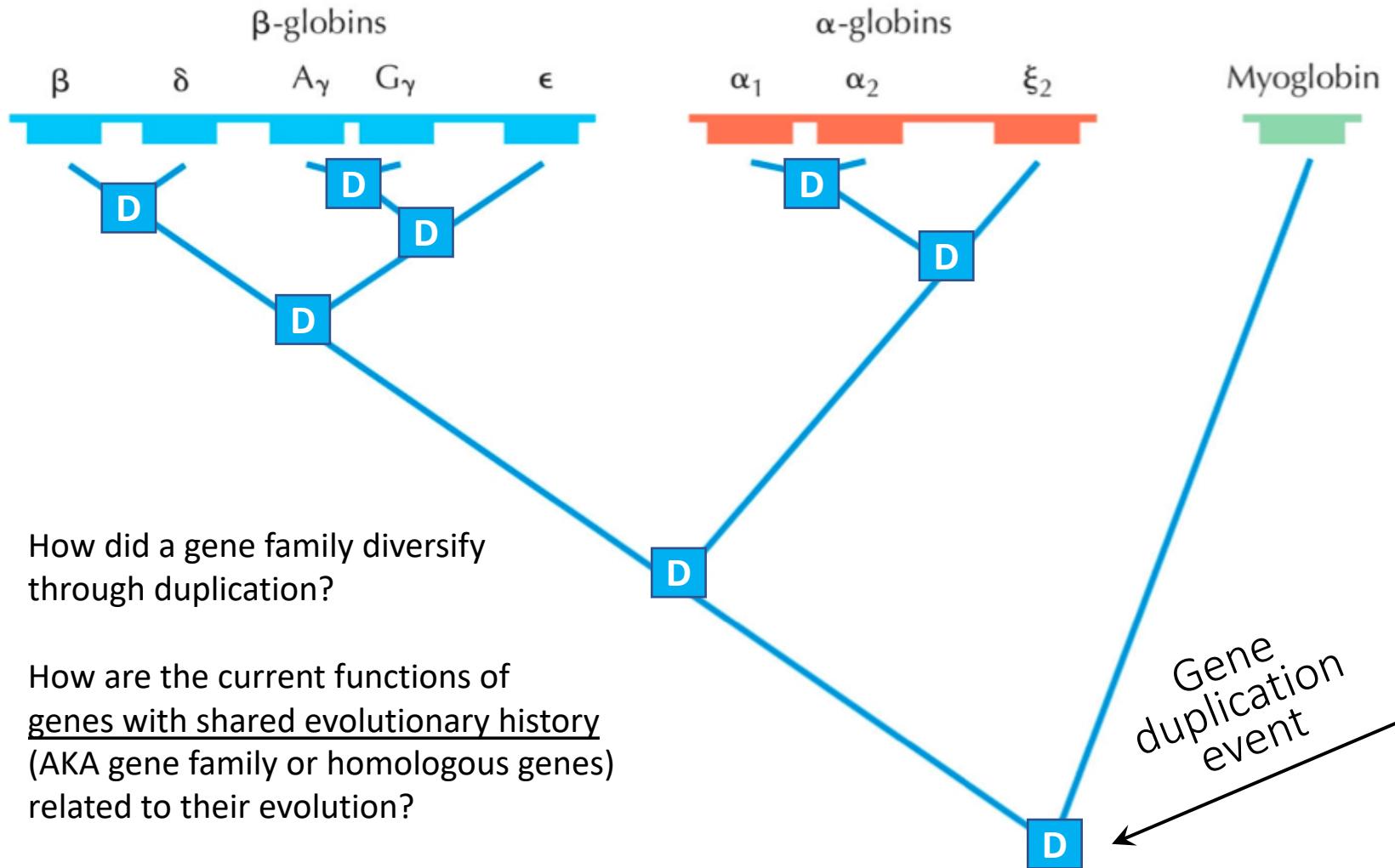
strain phylogeny



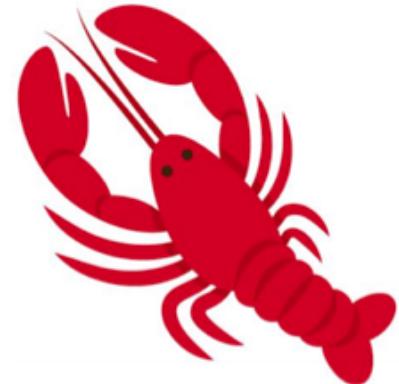


Establishment and cryptic transmission of Zika virus in Brazil and the Americas. Faria et al., Nature 2017

gene phylogeny



$((A,B),C),(D,E))$



“The Newick Standard was adopted 26 June 1986 by an informal committee meeting convened by me [Joe Felsenstein] during the Society for the Study of Evolution meetings in Durham, New Hampshire and consisting of James Archie, William H.E. Day, Wayne Maddison, Christopher Meacham, F. James Rohlf, David Swofford, and myself. (The committee was not an activity of the SSE nor endorsed by it). The reason for the name is that the second and final session of the committee met at Newick's restaurant in Dover, New Hampshire, and we enjoyed the meal of lobsters.”

```
((raccoon:19.19959,bear:6.80041):0.84600[50], ((sea_lion:11.99700,  
seal:12.00300):7.52973[100], ((monkey:100.85930,cat:47.14069):20.59  
201[80], weasel:18.87953):2.09460[75]):3.87382[50], dog:25.46154);  
  
(((A:5,B:5)f:3,C:9)g:1,(D:2,E:6)h:2);
```

Demo: <http://tree.bio.ed.ac.uk/software/figtree/>

Building Trees

You need a principled way of choosing the tree that best represents the relationship among leaves

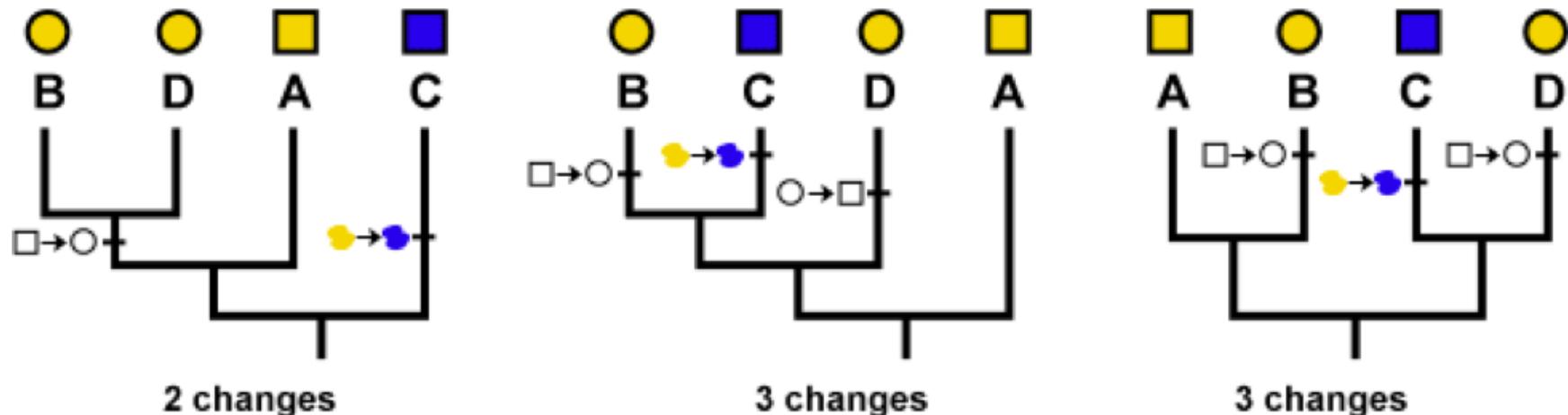
There are $(2n-3)!!$ rooted trees for n taxa

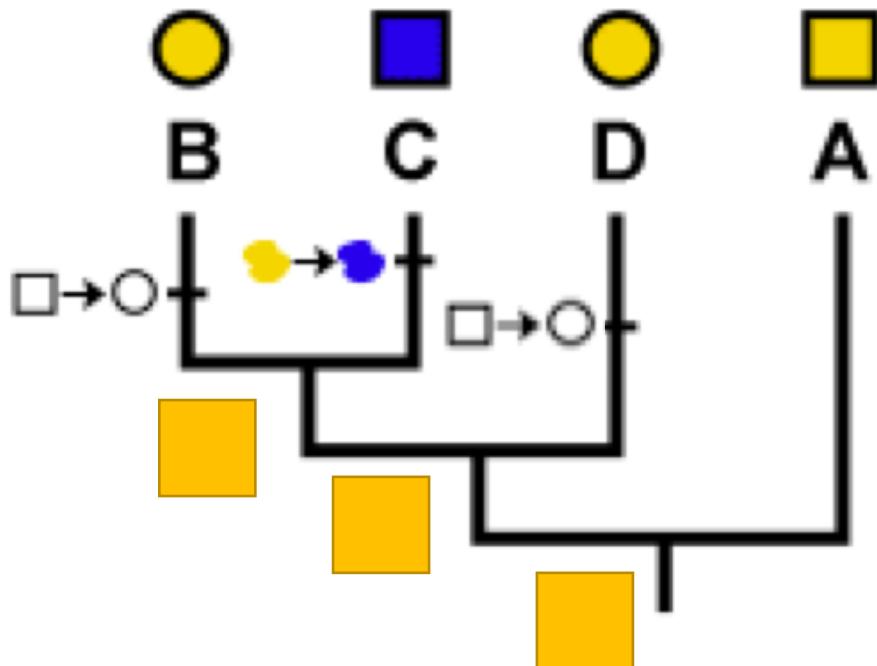
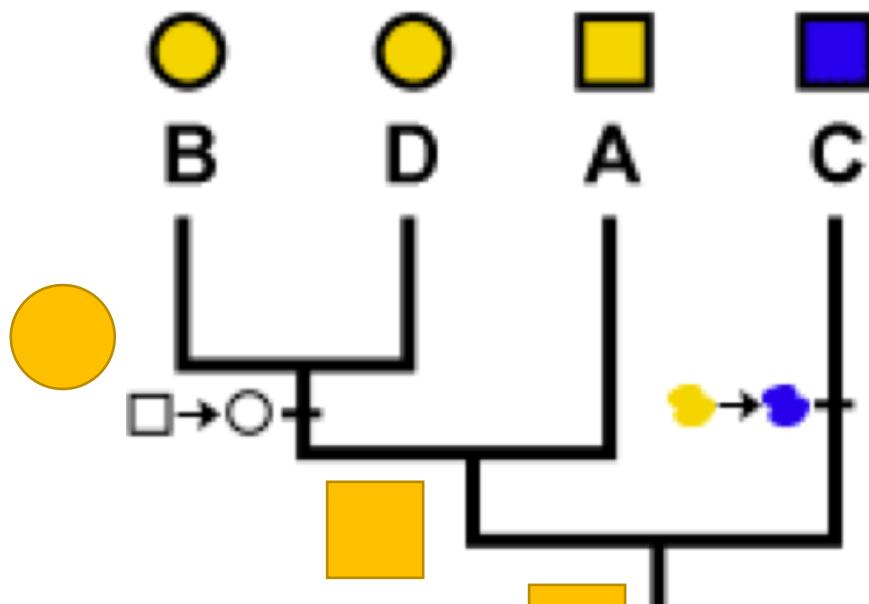
- 3 taxa: 3 trees
- 5 taxa: 105 trees
- 10 taxa: 34,459,425 trees
- 20 taxa: 8.2×10^{21} trees
- 50 taxa: 2.8×10^{76} trees
- 100 taxa: 3.3×10^{184} trees

The **double factorial**, symbolized by two exclamation marks (!!), is a quantity defined for all integers greater than or equal to -1. For an even integer n , the **double factorial** is the product of all even integers less than or equal to n but greater than or equal to 2.

Parsimony: scenario that requires the fewest changes is the best

taxon	characters	
	shape	color
A	□	yellow
B	○	yellow
C	□	blue
D	○	yellow



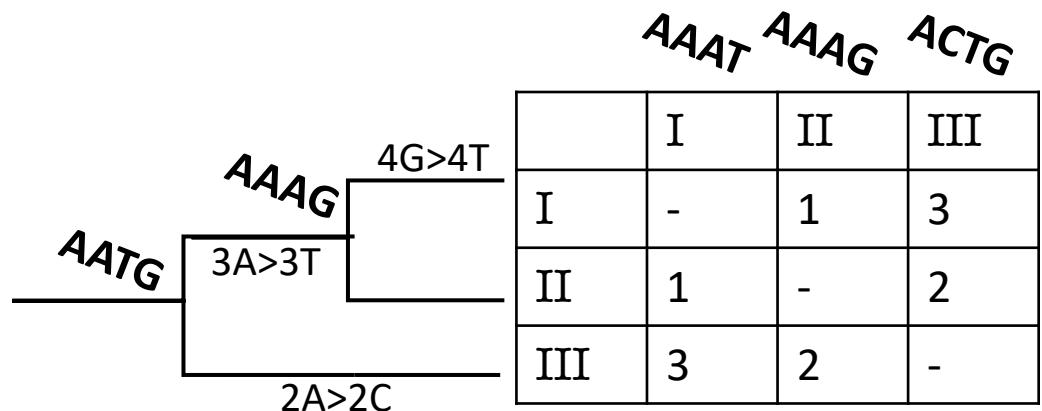


Clustering a distance matrix

- Calculate all pairwise distances
 - Use a mathematical model of protein / nucleotide evolution to calculate the distance between two sequences
- Here: Hamming distance

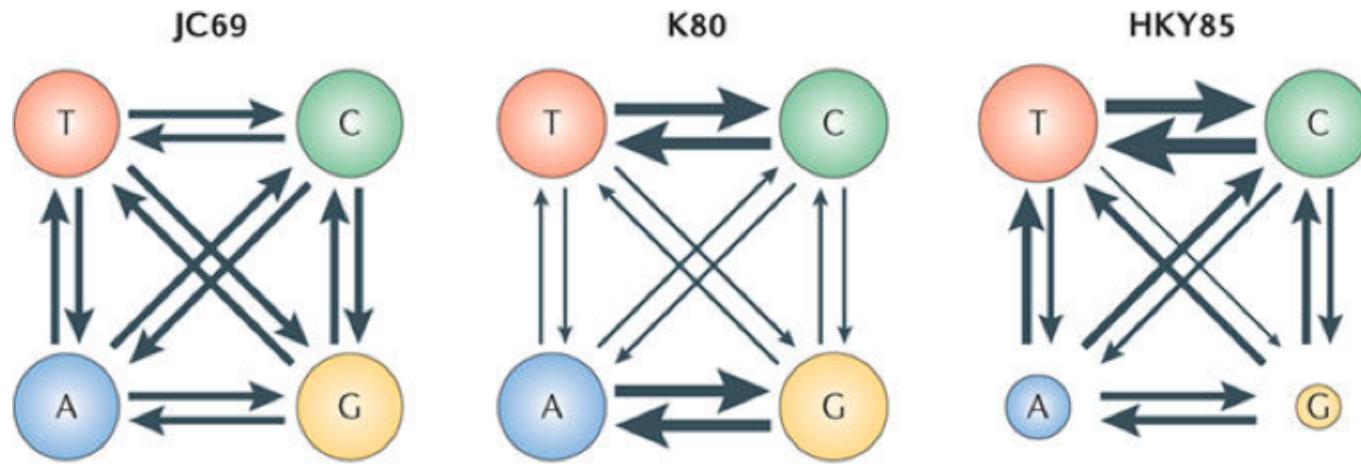
Example algorithms:

- Least squares
- Neighbor joining
- Minimum evolution
- UPGMA (Unweighted Pair Group Method with Arithmetic Mean)



Software: PHYLIP

Estimating distance between sequences: nucleotide evolution models



Nature Reviews | Genetics

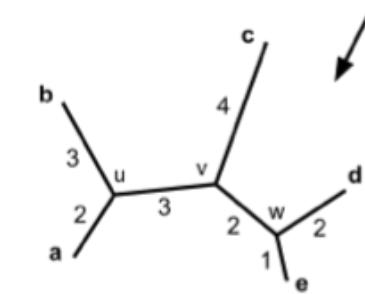
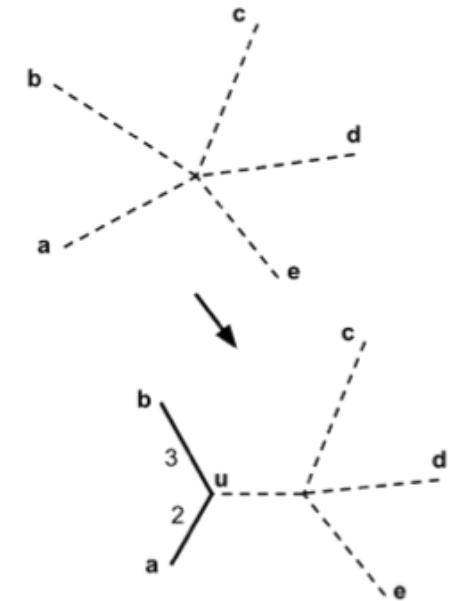
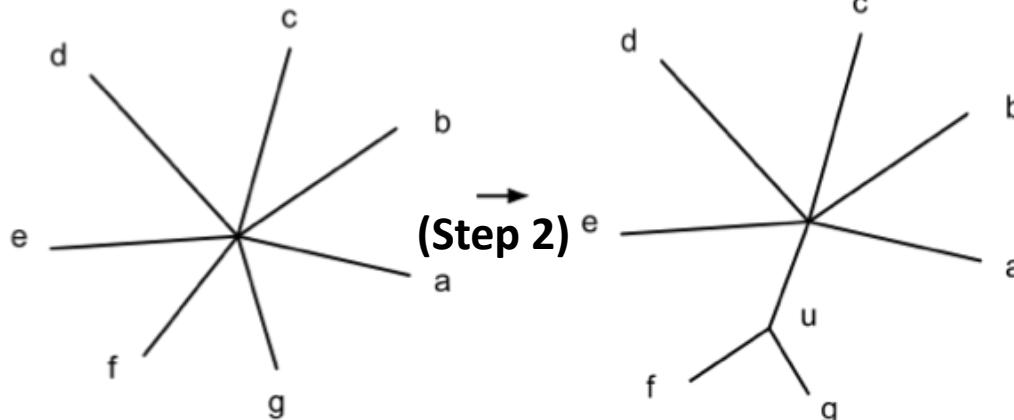
“The thickness of the arrows indicates the substitution rates of the four nucleotides (T, C, A and G), and the sizes of the circles represent the nucleotide frequencies when the substitution process is in equilibrium. Note that both JC69 and K80 predict equal proportions of the four nucleotides.”

Clustering: Neighbor joining

0. Start with a star graph
1. Calculate distance matrix
2. Find the minimum entry (closest leaves) and join those leaves with a new internal node, with branch lengths based on the values in the matrix
3. Replace the rows/columns of the leaves with a single row/column of the new internal node
4. Recalculate distance matrix

Repeat 2,3,4 until you have an unrooted binary tree

*Resulting tree is unrooted and not guaranteed to be ultrametric

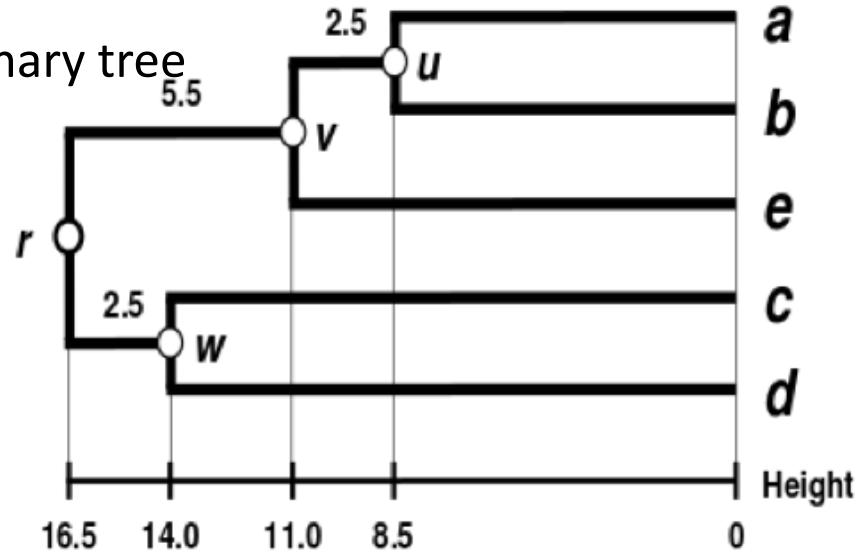


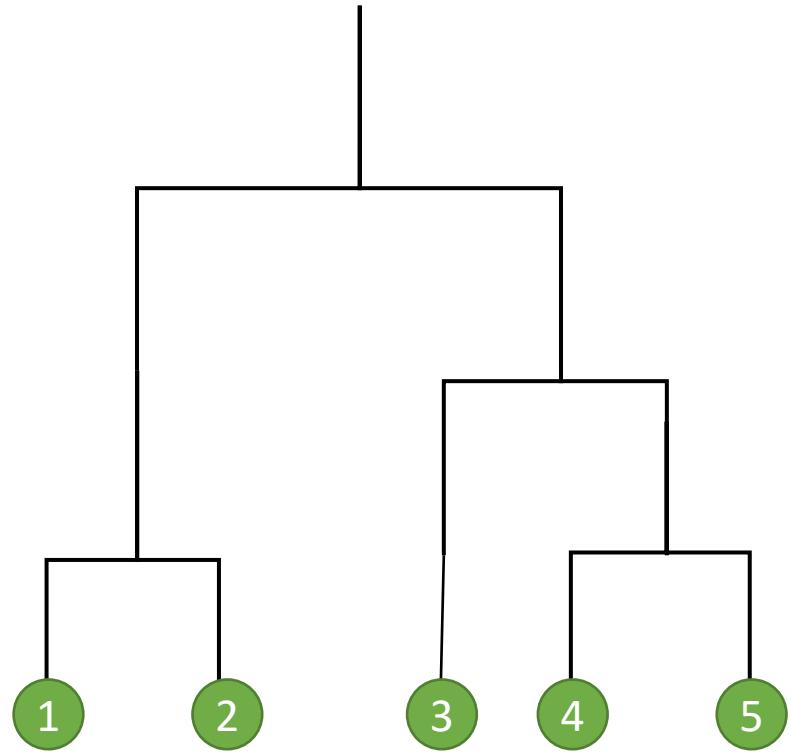
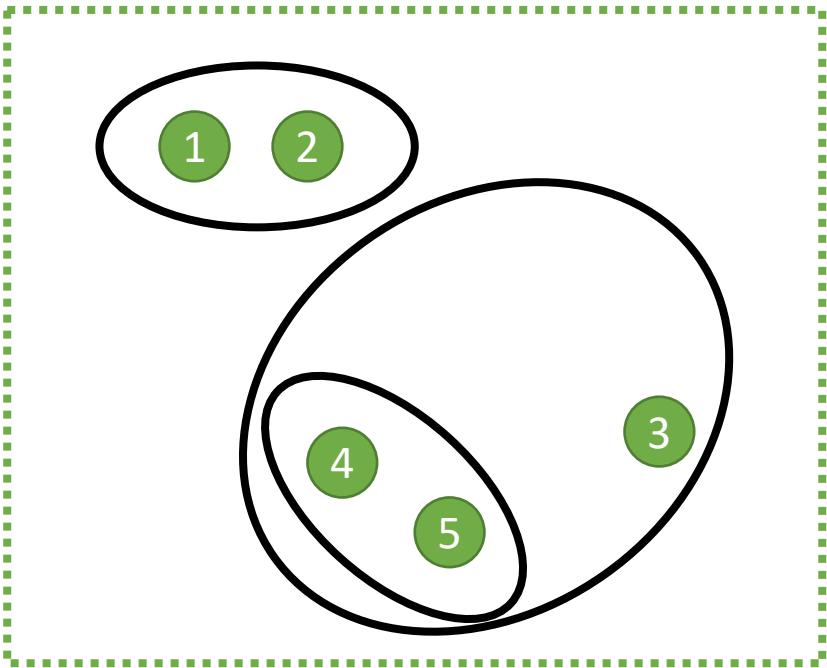
Clustering: UPGMA

0. Start with a star graph
1. Calculate distance matrix
2. Find the minimum entry (closest leaves) and join those leaves at a distance of the mean distance between their elements
3. Replace the rows/columns of the leaves with a single row/column of their union (cluster)
4. Recalculate distance matrix

Repeat 2,3,4 until you have an rooted binary tree

*Resulting tree *is* rooted and ultrametric:
assumes rate of evolution is constant
on all branches (molecular clock)





Multiple sequence alignment

Scarites	C T T A G A T C G T A C C A A - - - A A T A T T A C
Carenum	C T T A G A T C G T A C C A C A - T A C - T T T A C
Pasimachus	A T T A G A T C G T A C C A C T A T A A G T T T A C
Pheropsophus	C T T A G A T C G T T C C A C - - - A C A T A T T A C
Brachinus armiger	A T T A G A T C G T A C C A C - - - A T A T A T T T C
Brachinus hirsutus	A T T A G A T C G T A C C A C - - - A T A T A T T A C
Aptinus	C T T A G A T C G T A C C A C - - - A C A A T T T A C
Pseudomorpha	C T T A G A T C G T A C C - - - A C A A A A T T A C

Maximum likelihood estimation from an MSA

0. Start with a multiple sequence alignment
1. Build a starting tree
2. Calculate the likelihood of the tree based on an evolutionary model
3. Change the tree slightly to increase the likelihood

Repeat 2,3 until some convergence criterion met

ML is a heuristic search of “tree space”

Bootstrapping: resampling from “tree space” to gauge the quality of your solution

- Proportion of resamples that support a certain branching pattern

Software:

PhyML/RAXML

Genetic algorithm - GARLI

Bayesian estimation from an MSA

0. Start with a multiple sequence alignment
1. Build a starting tree and assume some prior probability on the distribution of trees
2. Calculate the posterior probability of the tree based on an evolutionary model
3. Sample a “nearby” tree, usually using MCMC (Markov chain Monte Carlo) algorithms
4. Accept the new tree if it is “better”, otherwise keep with the old tree

Repeat 2,3,4 until some convergence criterion met

Bayesian methods are also a heuristic search of “tree space”

The posterior probability (probability correct *given model*) gauges the quality of your solution

Software: BEAST, MrBayes

Table 2 | A summary of strengths and weaknesses of different tree reconstruction methods

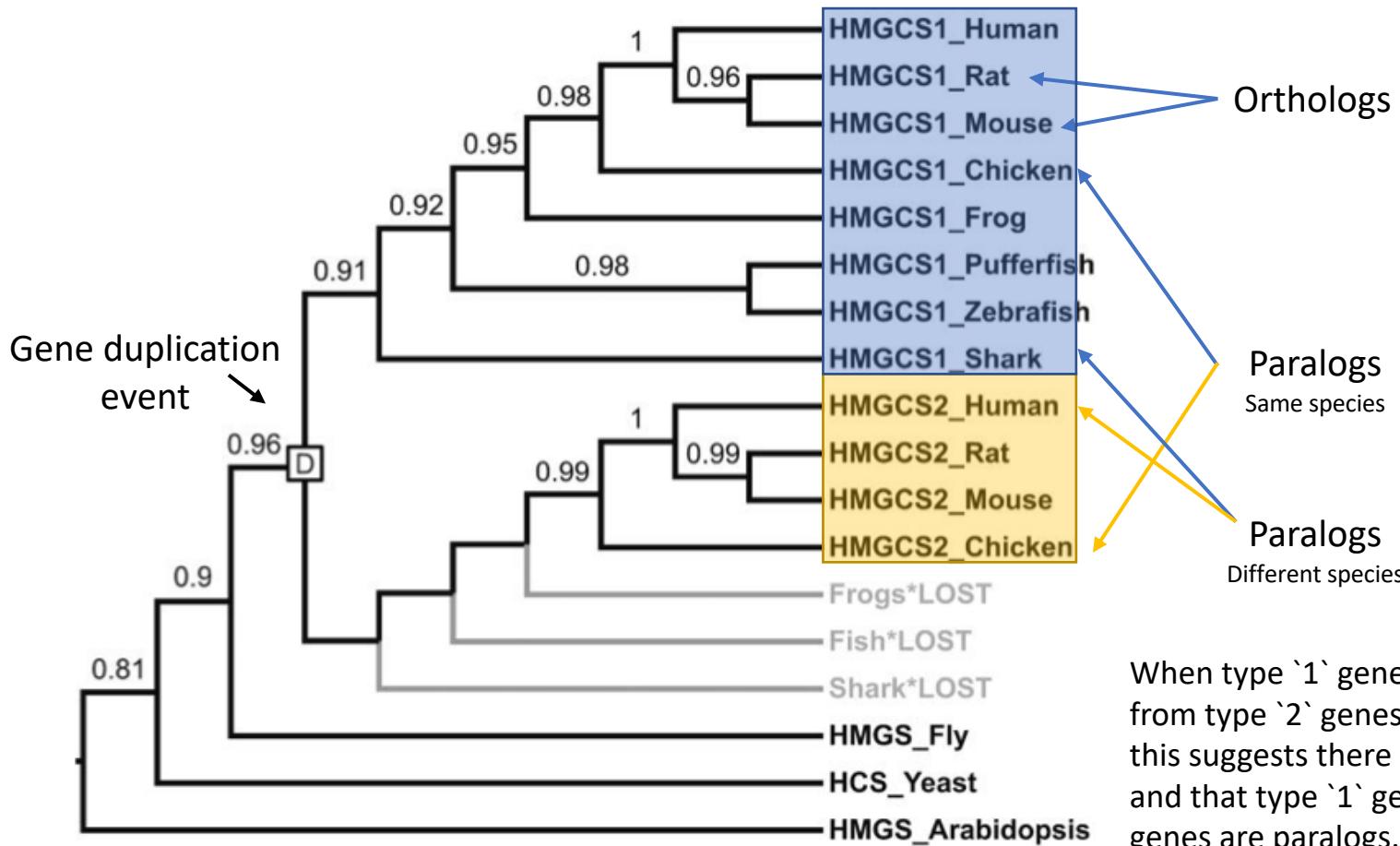
Strengths	Weaknesses
Parsimony methods	
<ul style="list-style-type: none"> • Simplicity and intuitive appeal • The only framework appropriate for some data (such as SINES and LINES) 	<ul style="list-style-type: none"> • Assumptions are implicit and poorly understood • Lack of a model makes it nearly impossible to incorporate our knowledge of sequence evolution • Branch lengths are substantially underestimated when substitution rates are high • Maximum parsimony may suffer from long-branch attraction
Distance methods	
<ul style="list-style-type: none"> • Fast computational speed • Can be applied to any type of data as long as a genetic distance can be defined • Models for distance calculation can be chosen to fit data 	<ul style="list-style-type: none"> • Most distance methods, such as neighbour joining, do not consider variances of distance estimates • Distance calculation is problematic when sequences are divergent and involve many alignment gaps • Negative branch lengths are not meaningful
Likelihood methods	
<ul style="list-style-type: none"> • Can use complex substitution models to approach biological reality • Powerful framework for estimating parameters and testing hypotheses 	<ul style="list-style-type: none"> • Maximum likelihood iteration involves heavy computation • The topology is not a parameter so that it is difficult to apply maximum likelihood theory for its estimation. Bootstrap proportions are hard to interpret
Bayesian methods	
<ul style="list-style-type: none"> • Can use realistic substitution models, as in maximum likelihood • Prior probability allows the incorporation of information or expert knowledge • Posterior probabilities for trees and clades have easy interpretations 	<ul style="list-style-type: none"> • Markov chain Monte Carlo (MCMC) involves heavy computation • In large data sets, MCMC convergence and mixing problems can be hard to identify or rectify • Uninformative prior probabilities may be difficult to specify. Multidimensional priors may have undue influence on the posterior without the investigator's knowledge • Posterior probabilities often appear too high • Model selection involves challenging computation^{138,139}

3-hydroxy-3-methylglutaryl coenzyme A
synthase (HMGCS) enzyme family

Homologs: Genes with shared evolutionary history

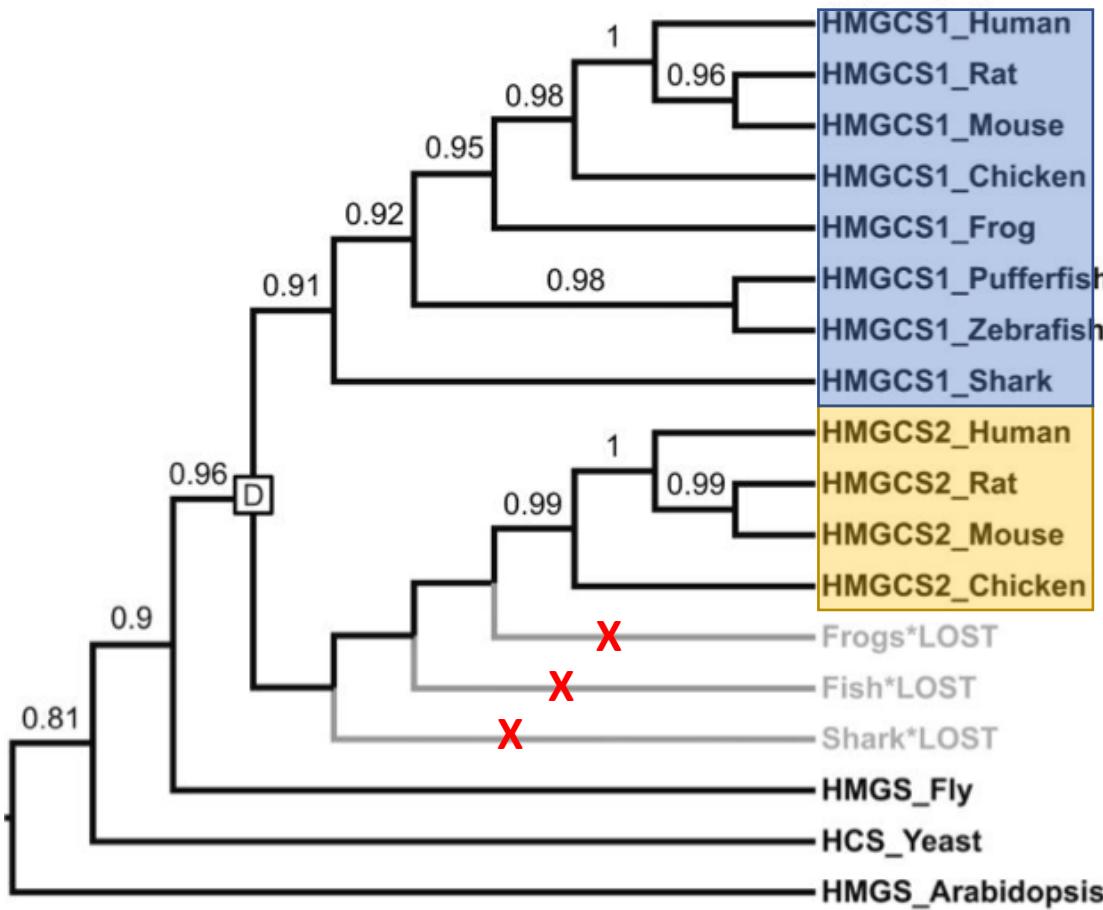
Orthologs: Genes that diverged through speciation (i.e. not duplication)

Paralogs: Genes that diverged through duplication

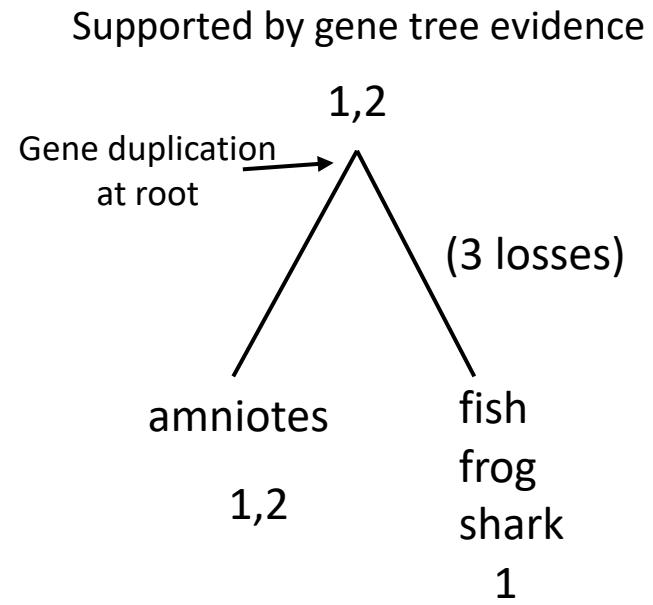
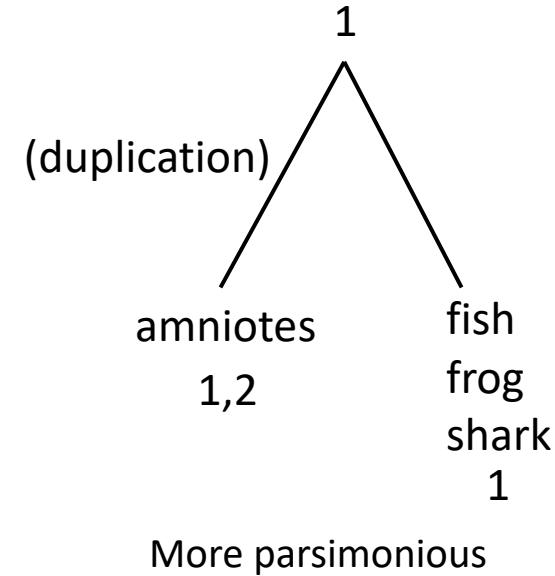


When type '1' genes cluster separately from type '2' genes in the gene tree, this suggests there was a duplication and that type '1' genes and type '2' genes are paralogs.

3-hydroxy-3-methylglutaryl coenzyme A
synthase (HMGCS) enzyme family

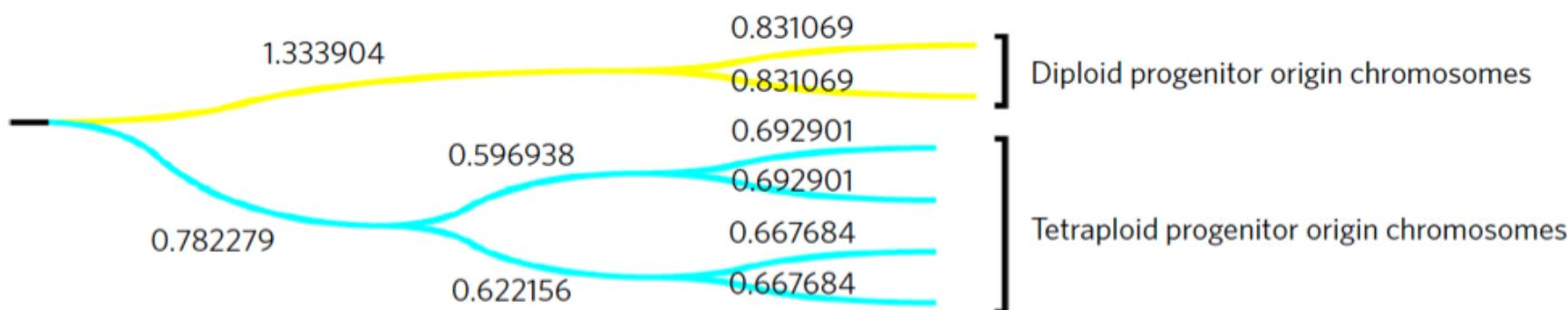


"Gene tree inferred using PhyML from sequences aligned with MAFFT and rooted using three invertebrate outgroup sequences. Branch support was assessed using aLRT scores."

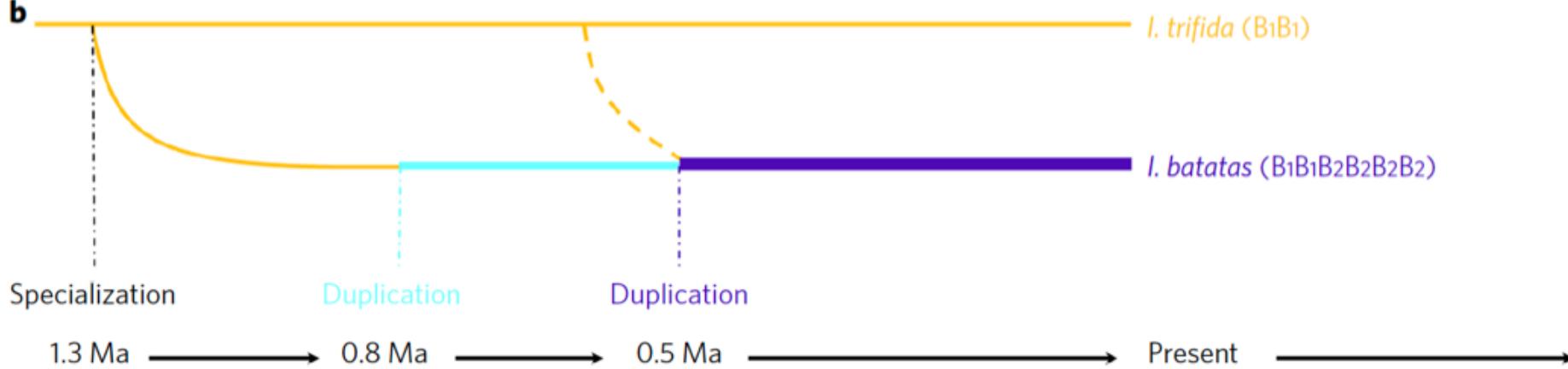


Evolution of chromosomes in hexaploid sweet potato

a



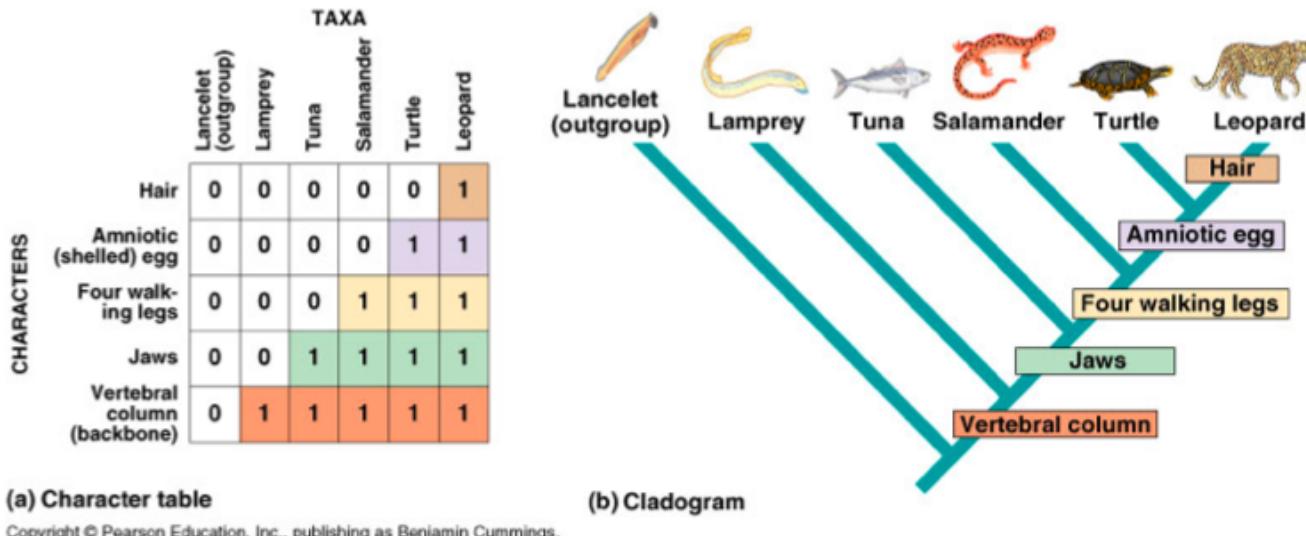
b



- a) Evolutionary history of cultivated *I. batatas* revealed by phylogenetic analysis of homologous chromosome regions. a, The dominant topology structure of all phylogenetic trees. Numbers indicate the average branch length of trees in this structure.
- b) The time points of B2 subgenome specialization and two whole-genome duplication events were estimated as 1.3, 0.8 and 0.5 million years ago (Ma). The estimation is based on 0.7% mutation rate per million years. The dashed curve indicates the crossing between diploid and tetraploid progenitors.

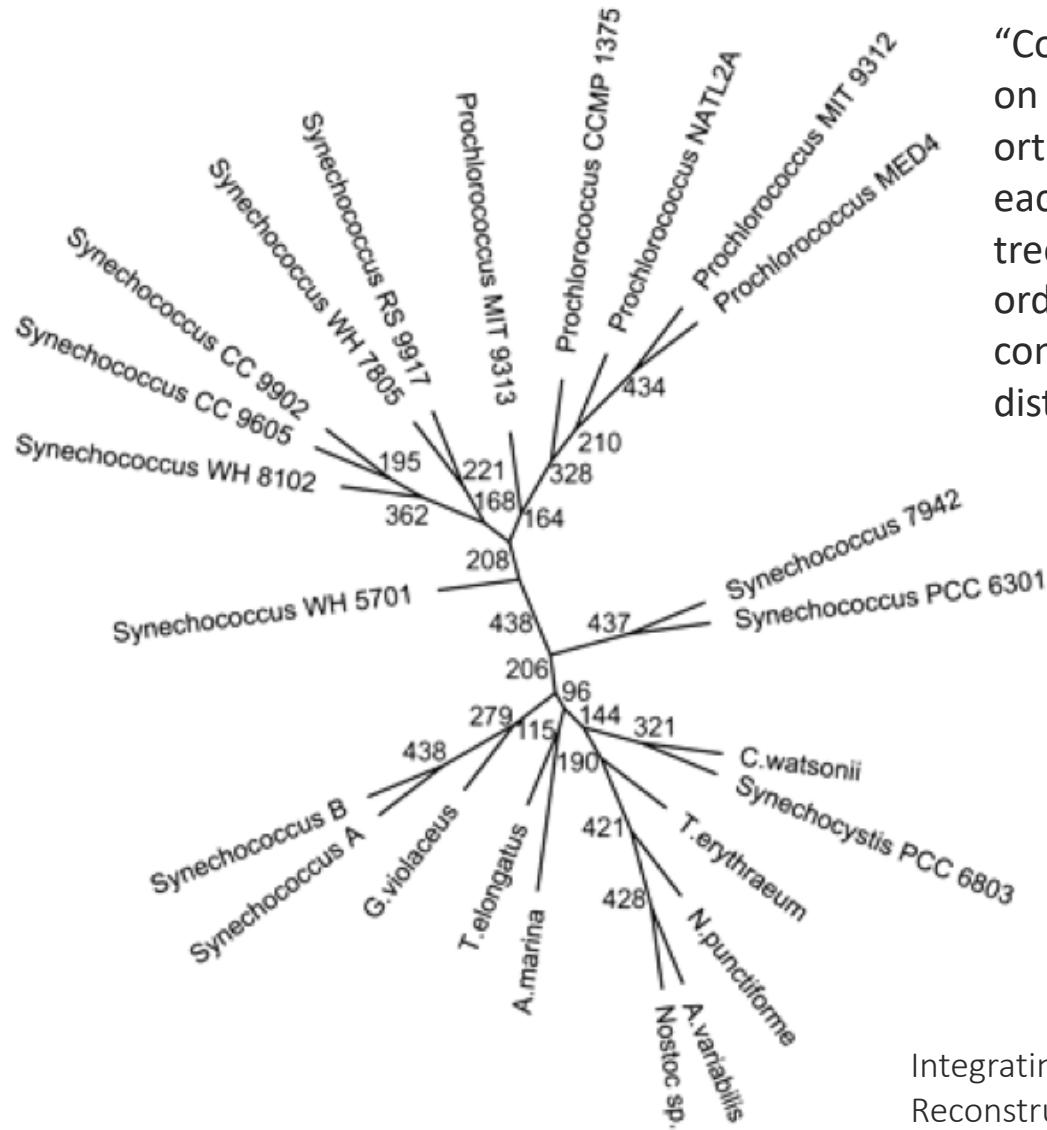
Building a species phylogeny

- Macroscopic characteristics



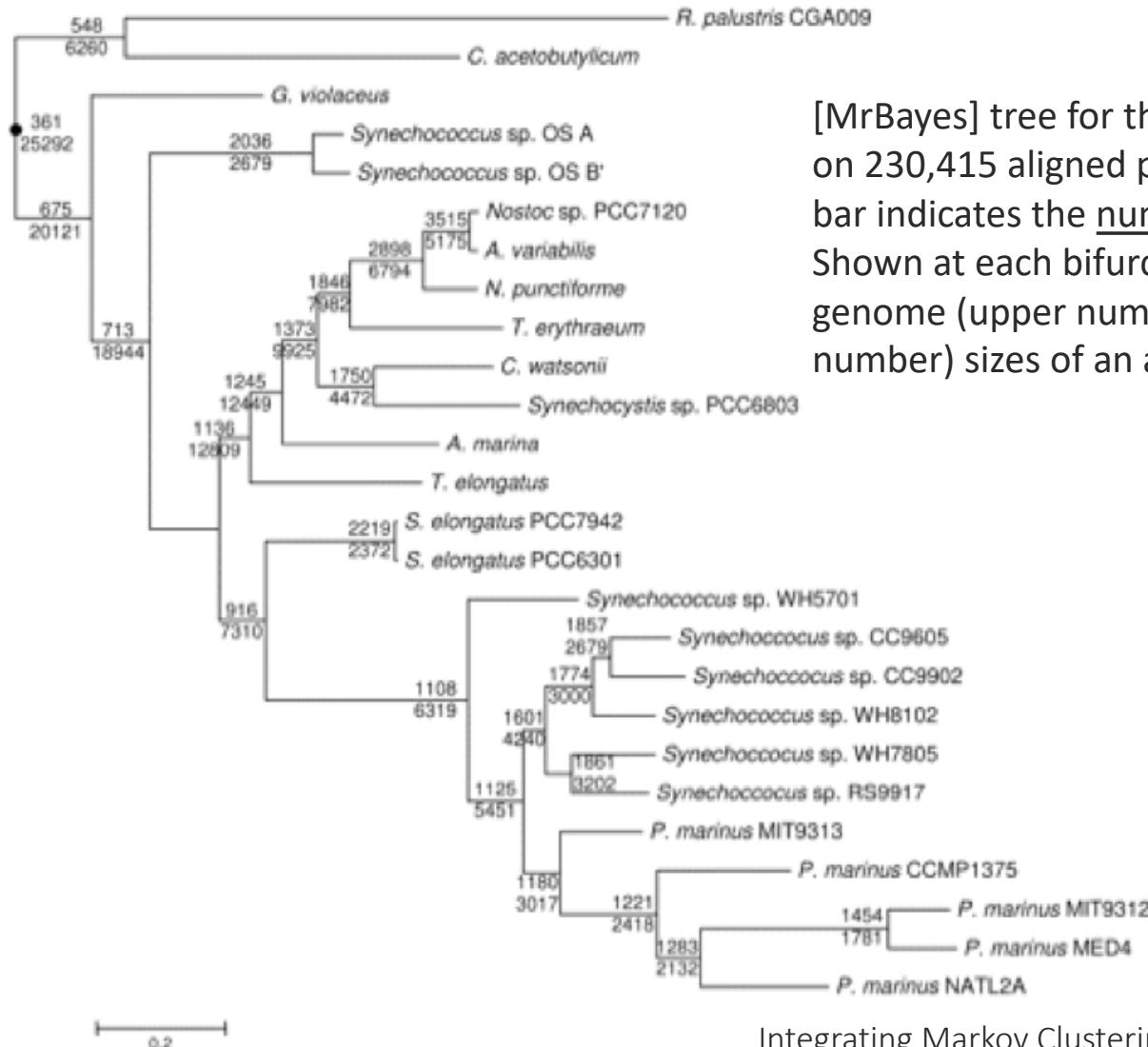
- Molecular characteristics
 - Conserved sequence, e.g. ribosome
 - Concatenation of conserved genes
 - Consensus of many individual-gene trees (supertree)

Species phylogeny of Cyanobacteria from consensus of many protein families

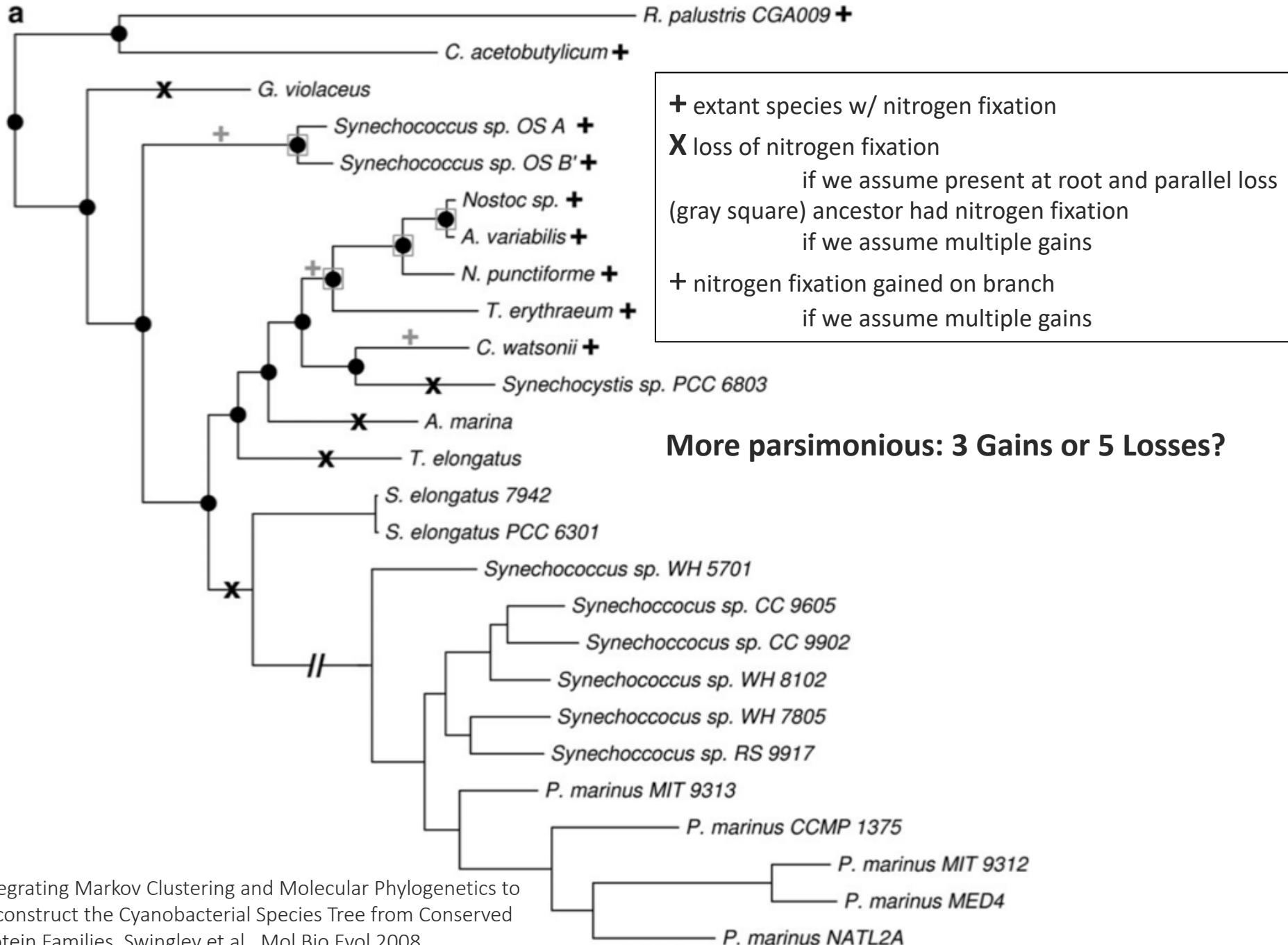


“Consensus cyanobacterial phylogeny based on maximum likelihood trees for each of 438 orthologous protein families. The numbers at each bifurcation indicate the total number of trees where that exact bifurcation/branching order is observed ... Note that as this is a consensus tree, topology is meaningful but distances are not.”

Species phylogeny of Cyanobacteria from concatenation of many multiple sequence alignments



[MrBayes] tree for the full-concatenated data set, based on 230,415 aligned positions in 26 genomes. The scale bar indicates the number of substitutions per site. Shown at each bifurcation are the predicted core-genome (upper number) and pan-genome (lower number) sizes of an ancestor at that point.



Resources

- Review article
 - Molecular phylogenetics: principles and practice. Yang and Rannala, *Nature Reviews Genetics* 2012
- Opinion piece
 - Homology: a personal view on some of the problems. Fitch, *Trends In Genetics* 2000
- Textbooks
 - *Molecular Evolution: A Phylogenetic Approach*. Page and Holmes
 - *Inferring Phylogenies*. Felsenstein

My work in evolutionary biology (horizontal gene transfer)
Xenolog classification. Darby, Stolzer et. al., *Bioinformatics* 2017