

Phylogenetic trees

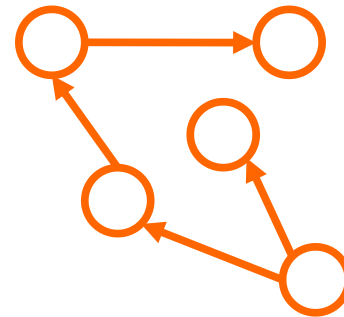
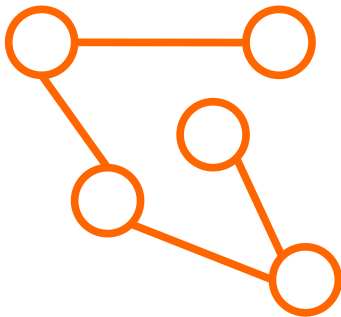
An introduction

Outline

- What are phylogenetic trees?
- Phylogenetic tree terminology
 - extant, ancestral, branch point, branch length, etc.
- Why do we construct/estimate phylogenetic trees?
- What is the difficulty of estimating phylogenetic trees?

What is a tree?

- Graph theoretically:
 - Undirected case: graph without cycles
 - Directed case: underlying undirected graph is a tree
 - Often it is required that $\text{indegree}(v) \leq 1$ for all v



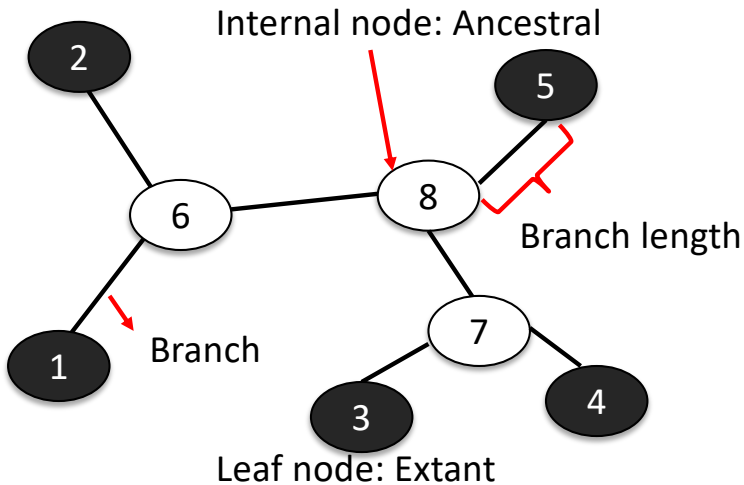
What are phylogenetic trees?

- A tree that describes evolutionary relationships between biological entities
 - Example entities: species, genes, strains
- This relationship is called *phylogeny*
- *Phylogenetics*: the task of inferring the true phylogeny from observations in existing organisms
- Meaning of nodes in a phylogenetic tree:
 - Leaves represent *extant* (current day) entities
 - Internal nodes represent ancestral entities

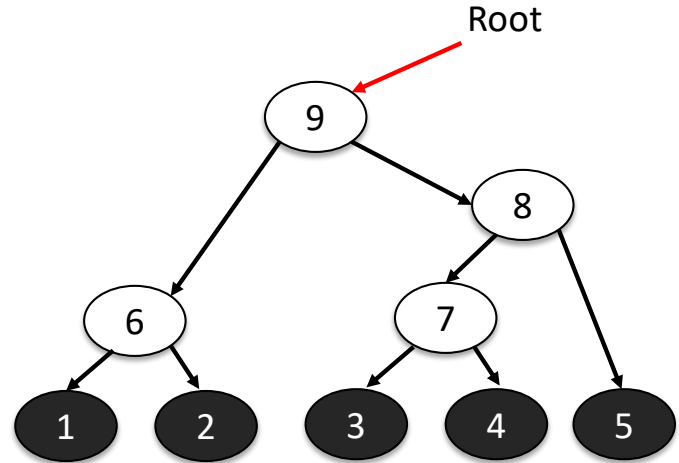
Phylogenetic tree basics

- Leaves represent entities (genes, species, individuals/strains) being compared
 - the term taxon (*taxa* plural) is used to refer to these when they represent species and broader classifications of organisms
 - For example if taxa are species, the tree is a species tree
- Phylogenetic trees can be rooted or unrooted
 - Rooted trees: the root node represents the common ancestor
- In a *rooted* tree, path from root to a node represents an evolutionary path
 - Gives directionality to evolutionary time
- An *unrooted* tree specifies relationships among entities, but lacks directionality information

Tree basics



Unrooted tree



Rooted tree

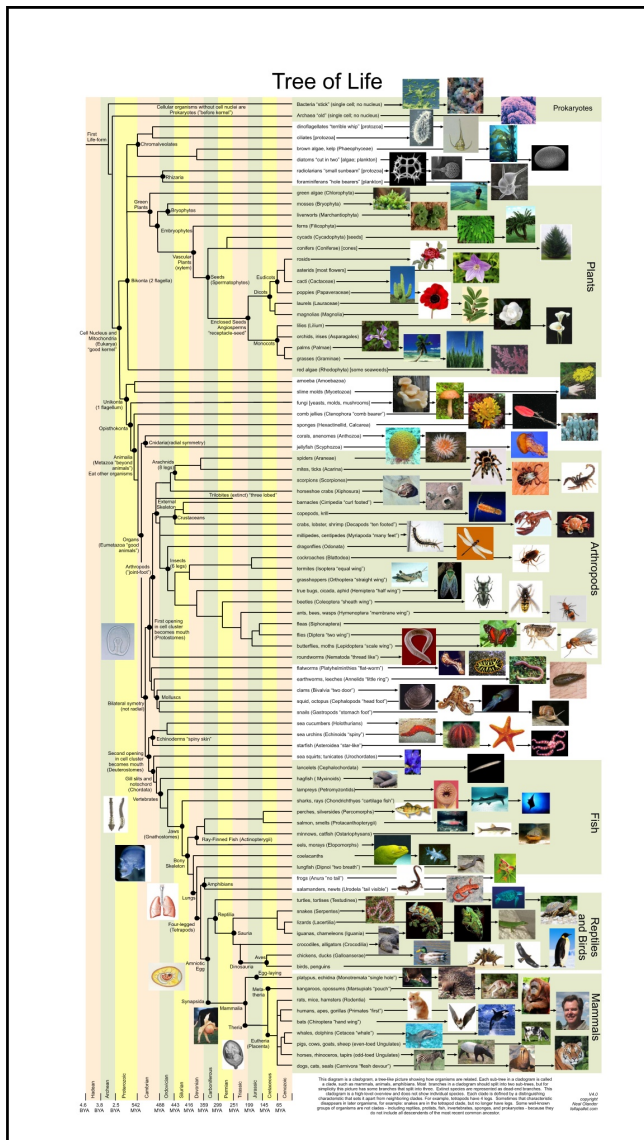
Each tree topology represents a different evolutionary history

For a species tree, internal nodes represent speciation events

Branch length describes the evolutionary divergence between two nodes

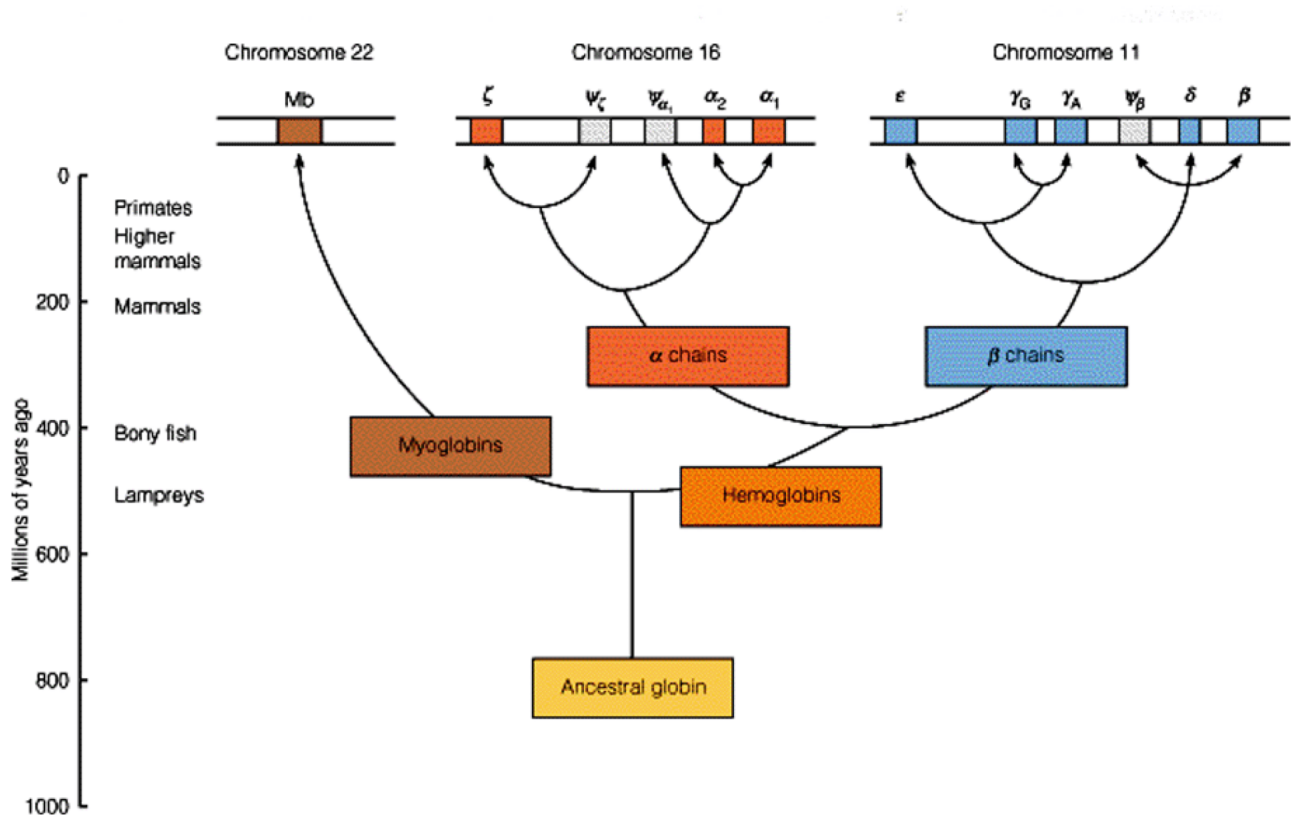
Why phylogenetic trees?

- Understand how organisms are related
 - Do humans and chimpanzees share a more recent common ancestor than do humans and gorillas? (Answer: they do)
- Ask how closely species are related
 - When did the common ancestor of humans and chimpanzees live? (Answer: around 5 million years ago)
- How specific functions/traits have evolved
 - What made us human?
- Identify signatures of conservation of sequence
- Inform multiple sequence alignments



Tree of life aims to represent the phylogeny of all species on earth

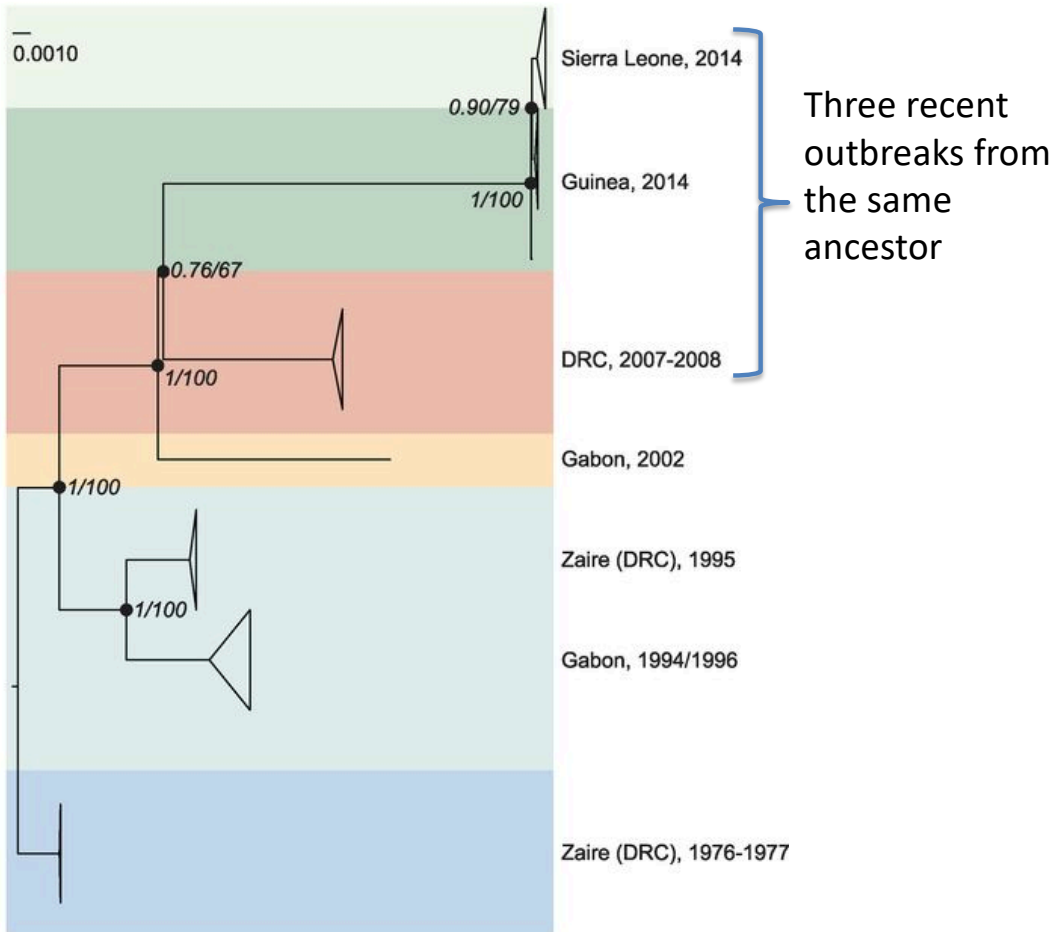
Example Gene Tree: Globins



Tracing the evolution of the Ebola virus

- Ebola virus: a lethal human pathogen, fatality rate 78%
- 2014 Ebola epidemic in Africa
 - Until recently the largest known case happened in 1976 (318 cases)
 - Outbreak reported in Feb 2014
 - 11,310 reported deaths from 2014 outbreak
 - World Health Organization ended declaration of Public Health Emergency in March 2016
- Key questions
 - Where did the pathogen come from?
 - How is it evolving?
- In a 2014 Science paper, researchers reported whole genome sequence alignment of 78 Ebola virus samples

Phylogenetic tree of the Ebola virus



Insights gained from sequence comparison

- “Genetic similarity across the sequenced 2014 samples suggests a single transmission from the natural reservoir, followed by human-to-human transmission during the outbreak”
- “..data suggest that the Sierra Leone outbreak stemmed from the introduction of two genetically distinct viruses from Guinea around the same time...”
- “..the catalog of 395 mutations, including 50 fixed nonsynonymous changes with 8 at positions with high levels of conservation across ebola viruses, provides a starting point for such studies”

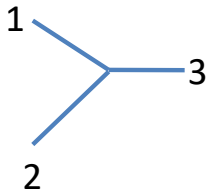
The phylogenetic tree estimation task and its difficulty

- *Phylogenetic tree estimation task:*
- **Given:** Observed features (e.g., DNA sequences) of each extant entity
- **Do:** Construct a phylogenetic tree that “best explains” the observed features
- How hard is this?
 - Depends on the data/features and definition of “best explains”
 - Some combinatorics will tell us how many possible trees such a task will need to consider (explicitly or implicitly)

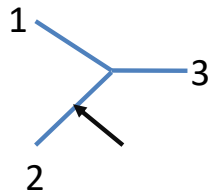
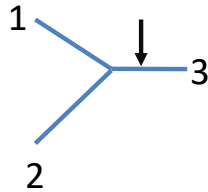
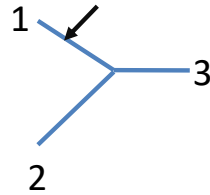
Tree counting

- A rooted binary tree with n leaf nodes has
 - $n-1$ internal nodes
 - $2n-2$ edges/branches
- An unrooted binary tree with n leaf nodes has
 - $n-2$ internal nodes
 - $2n-3$ edges/branches
 - A root can be added to any of these branches to give $2n-3$ rooted trees for any unrooted tree
- E.g. for $n=3$ there is *one* unrooted tree and *three* rooted trees

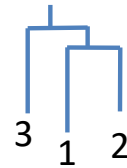
Tree counting



An unrooted tree



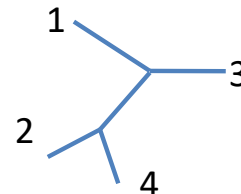
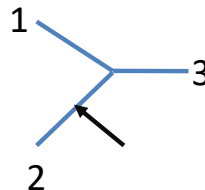
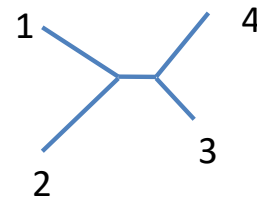
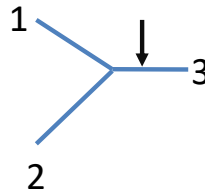
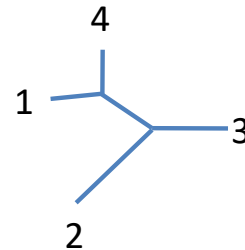
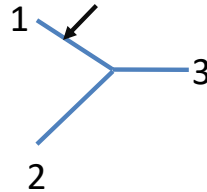
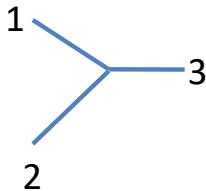
Possible positions for root



Rooted trees

Tree counting

- Instead of adding a root we could add a branch for the $n+1^{\text{th}}$ taxon



Tree counting

- A tree with 3 leaves can be grown in $(2*3)-3=3$ ways to make a tree of 4 leaves
 - 3 possible unrooted trees with 4 leaves
- Each tree with 4 leaves can be grown in $(2*4)-3=5$ ways to make a tree of 5 leaves
 - $3*5$ possible unrooted trees with 5 leaves
- Each tree of 5 leaves can be grown in $(2*5)-3=7$ ways
 - $3*5*7$ possible unrooted trees with 6 leaves
- In general for n leaves we can have
 - $(1)*(3)*(5)*...(2n-5)$ unrooted trees

Number of Possible Trees

- given n leaves, there are $\prod_{i=3}^n (2i - 5)$ possible unrooted trees
- and $(2n - 3) \prod_{i=3}^n (2i - 5)$ possible rooted trees
- This grows very fast
 - For $n=10$, we have 2 million unrooted trees
 - For $n=20$, we have $2.2 \cdot 10^{20}$

Summary

- Phylogenetic trees are representations of evolutionary relationships between biological entities
- Phylogenetic trees are used in a variety of applications including understanding the tree of life and determining the origins of viral outbreaks
- The estimation of phylogenetic trees heavily depends on the criteria being optimized
- Tree space is large - explicit consideration of all trees is not feasible