# Introduction to Phylogenetics

Heath Blackmon

19 December 2024
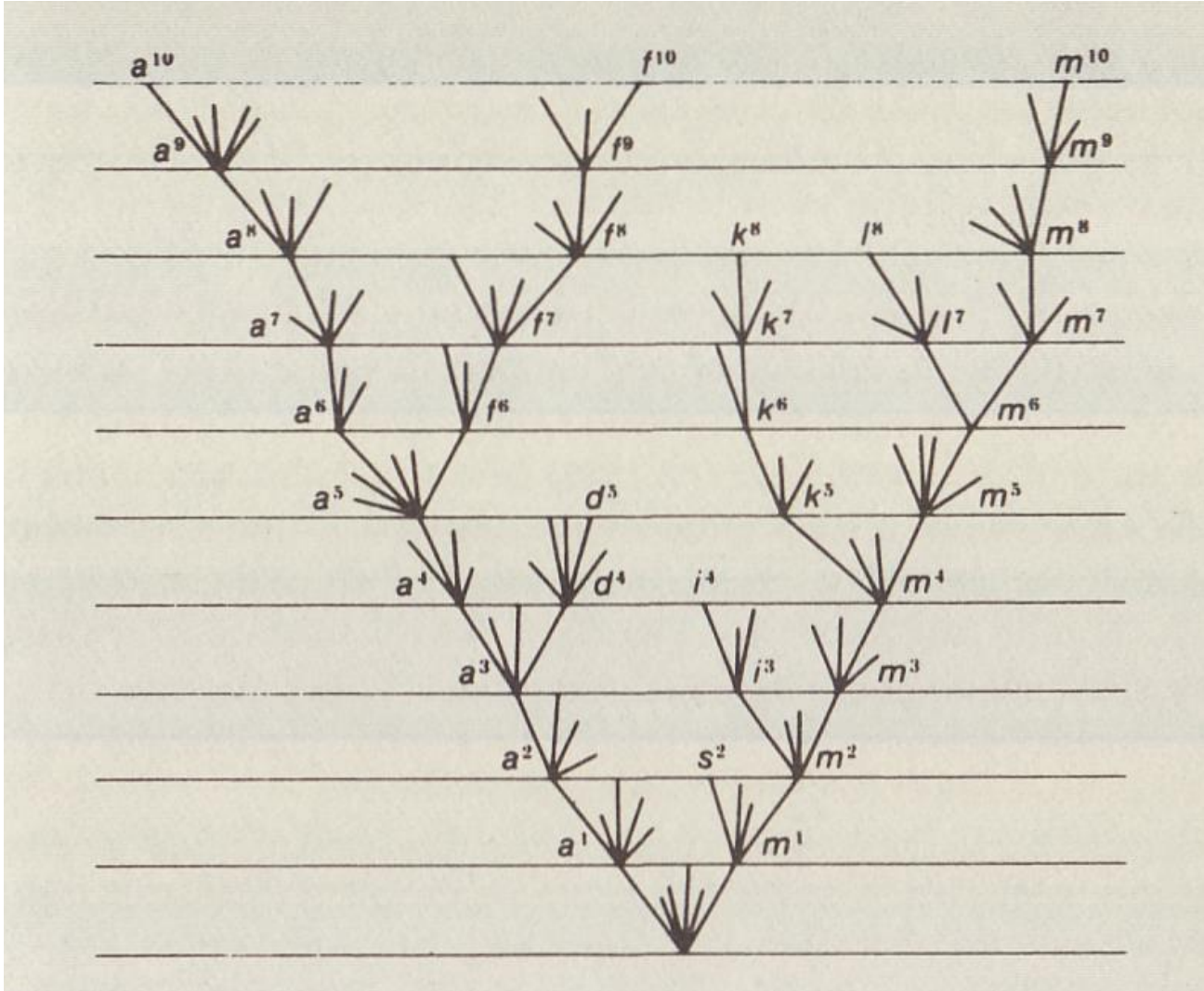
**Conceptual Foundations 1 hour**

- What is a phylogeny
- Why perform phylogenetic inference
- Data for phylogenetics
- Overview of models of evolution
- Methods of inference
- Interpreting trees

**Hands-On Tools and Methods 1 hour**

- Exercise 1: Running a Simple RAxML Analysis
- Exercise 2: Setting up a Basic BEAST Run
- Exercise 3: Visualizing and evaluating trees in R
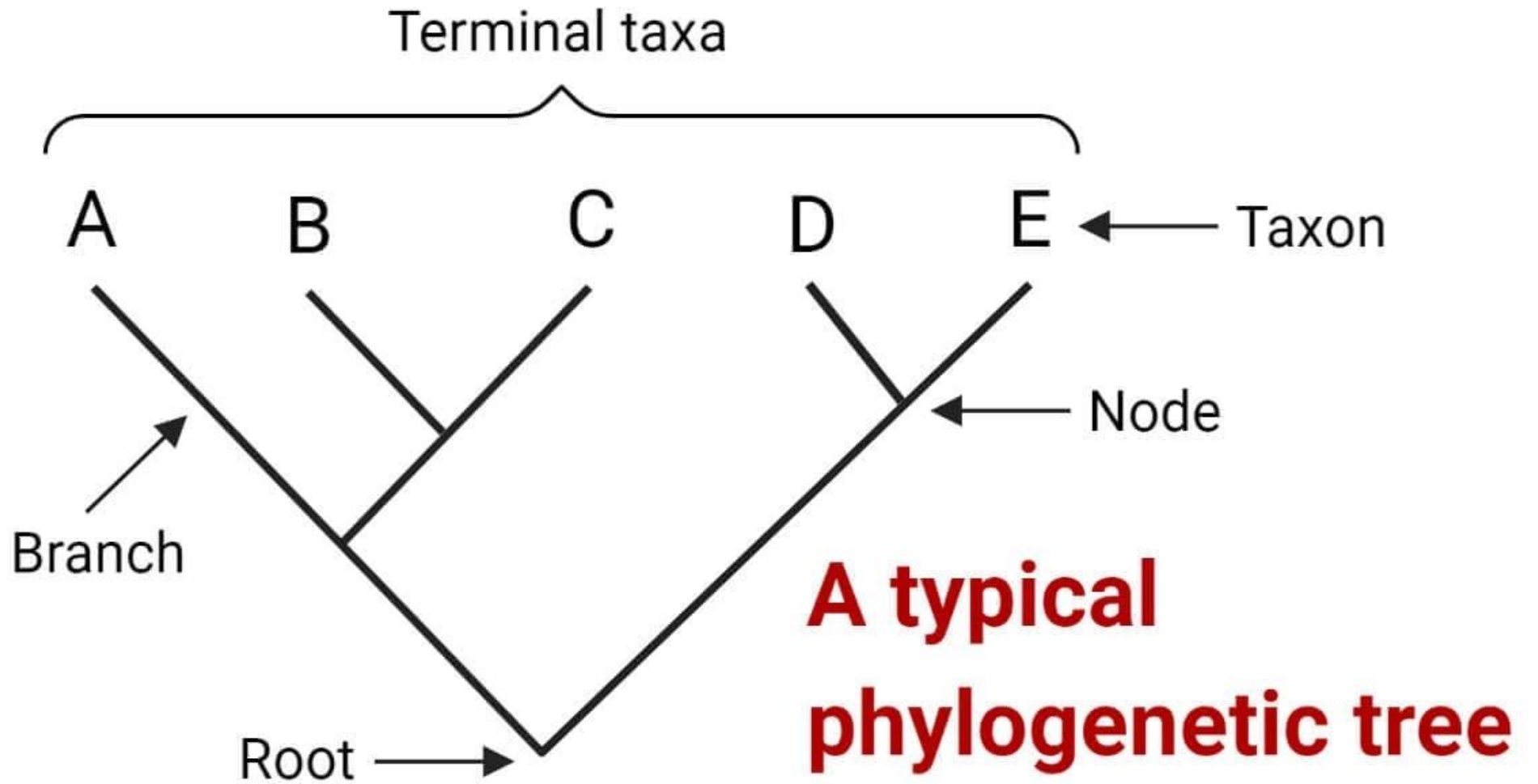
**Mini-Project and Discussion 1 hour**

## Conceptual Foundations 1 hour

- What is a phylogeny
- Why perform phylogenetic inference
- Data for phylogenetics
- Overview of models of evolution
- Methods of inference
- Interpreting trees

## Hands-On Tools and Methods 1 hour

- Exercise 1: Running a Simple RAxML Analysis
- Exercise 2: Setting up a Basic BEAST Run
- Exercise 3: Visualizing and evaluating trees in R

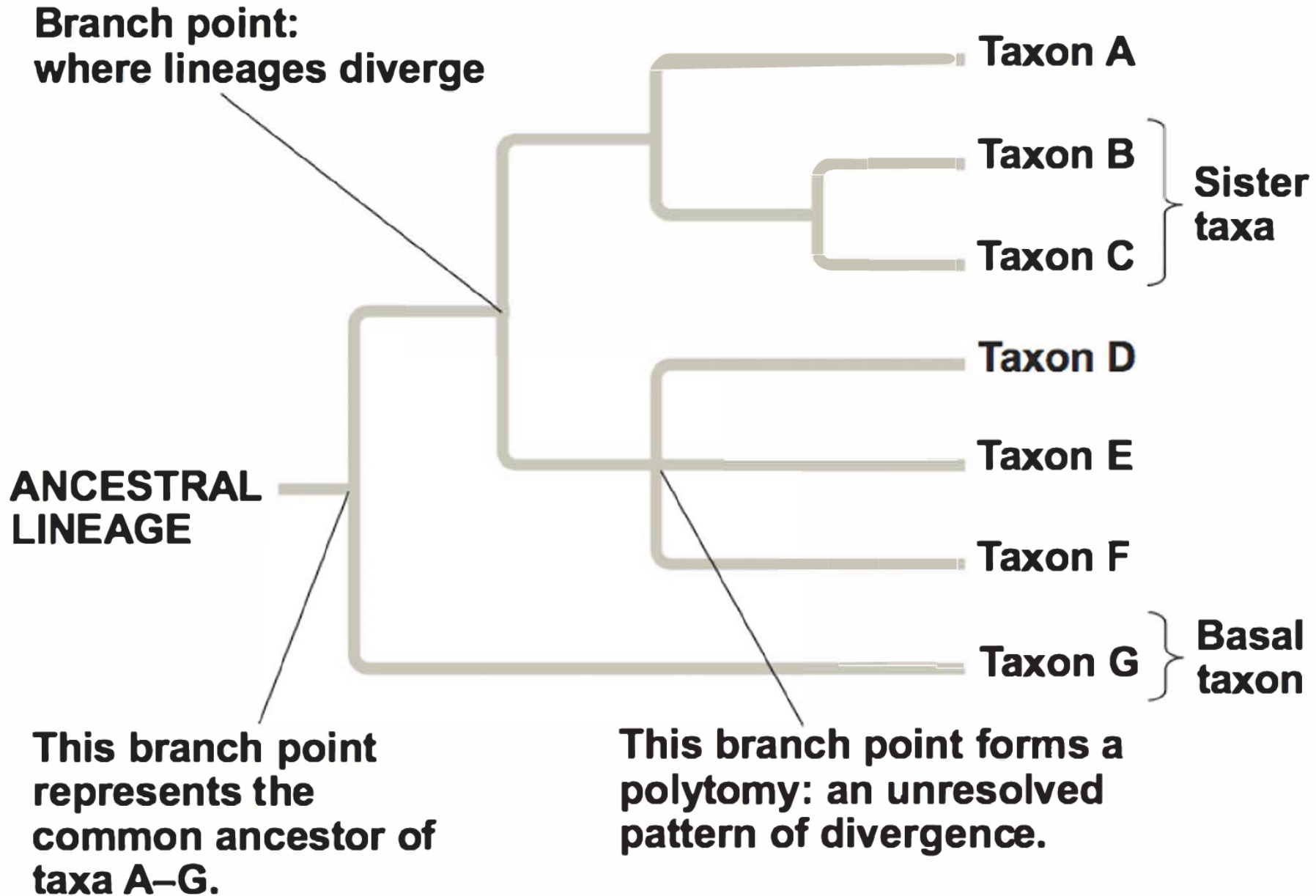## Mini-Project and Discussion 1 hour

There is grandeur in this view of life, with its several powers, having been originally breathed into a few forms or into one; and that, whilst this planet has gone cycling on according to the fixed law of gravity, from so simple a beginning endless forms most beautiful and most wonderful have been, and are being, evolved.
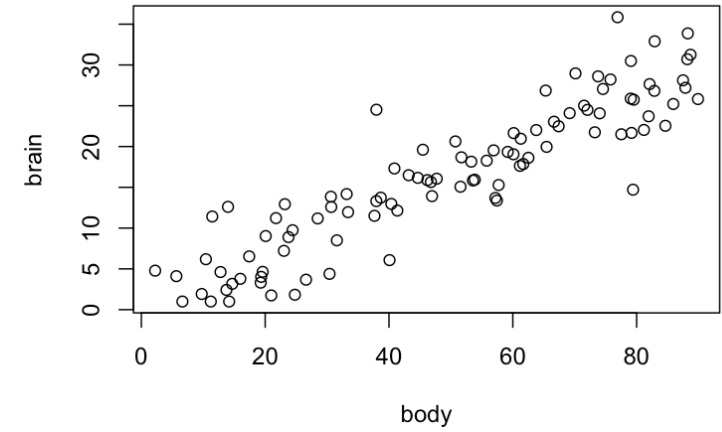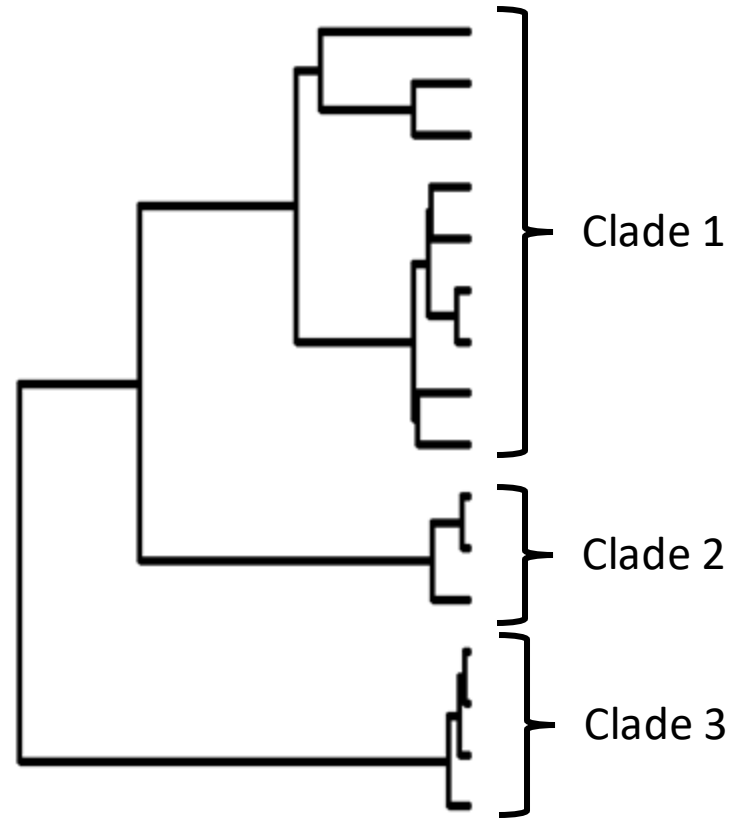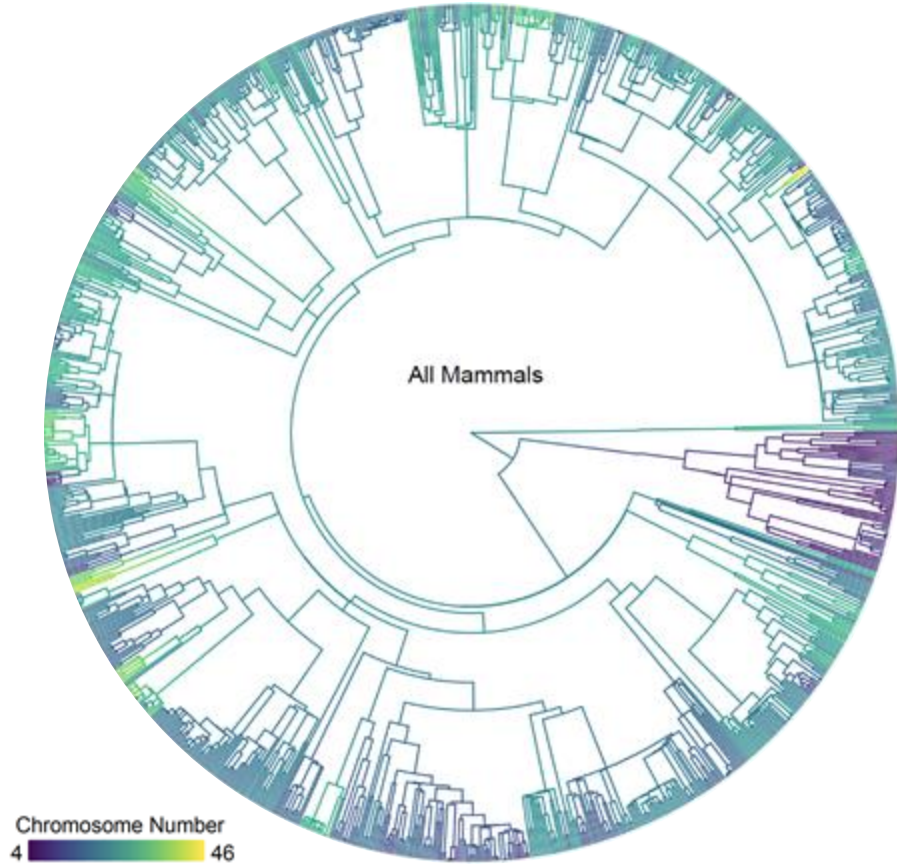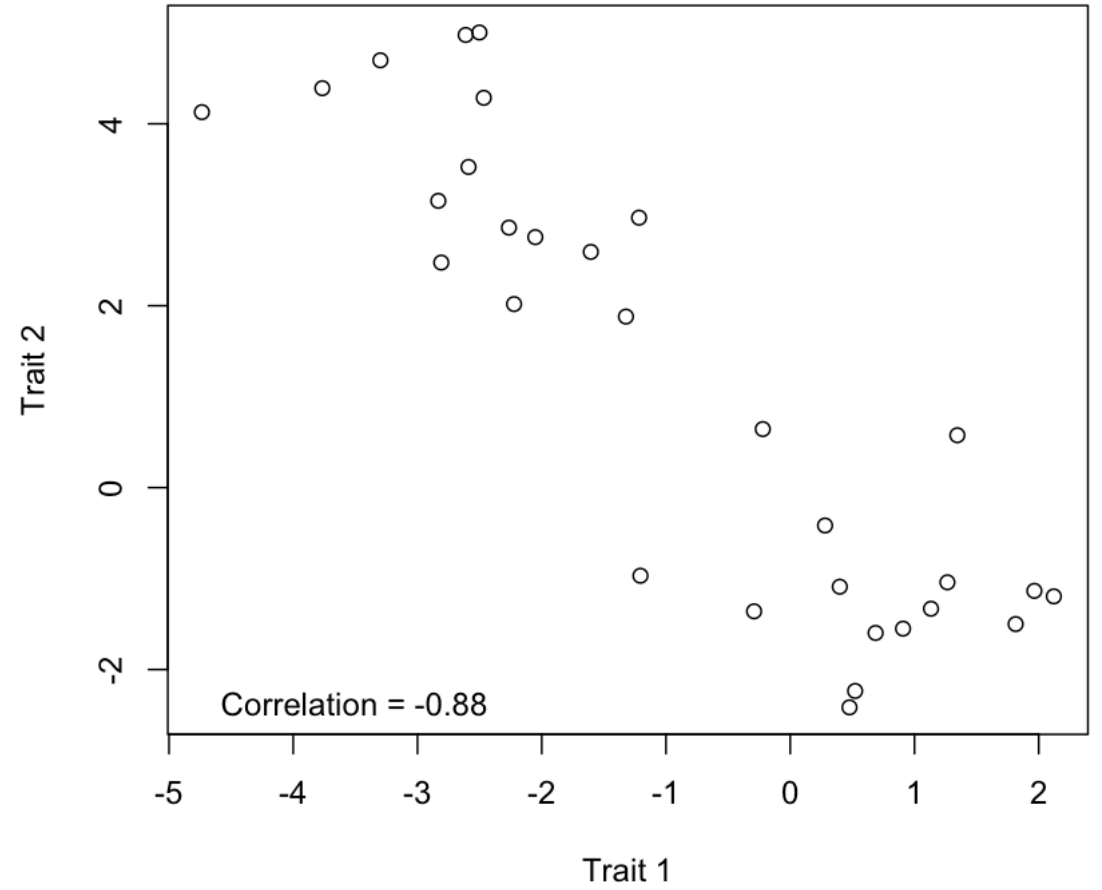
Darwin 1859

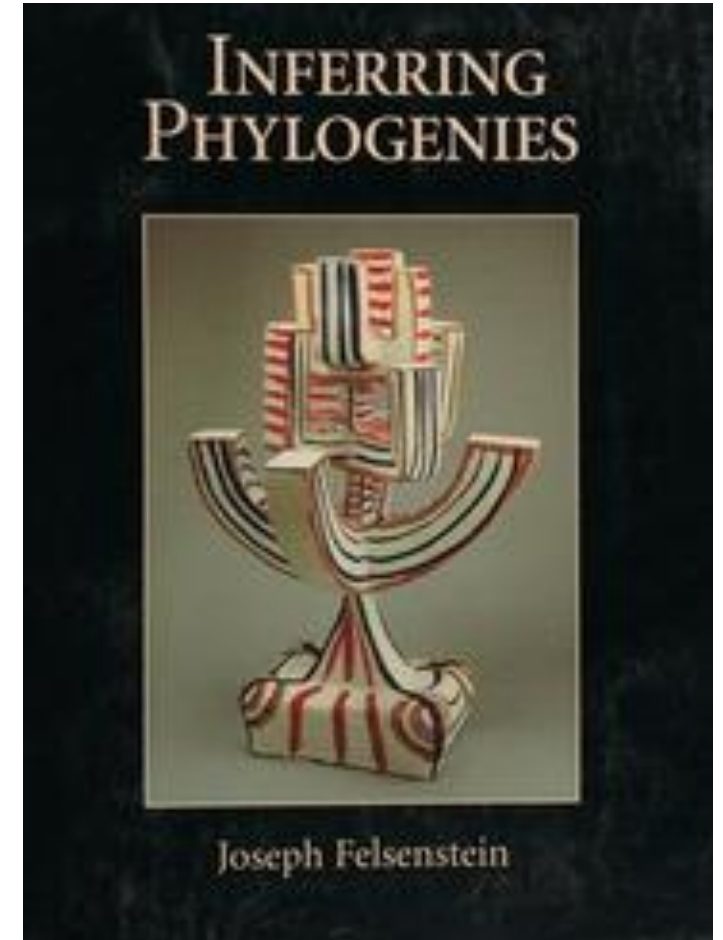A typical phylogenetic tree

# What is a phylogeny

# Why infer a phylogeny

## PHYLOGENIES AND THE COMPARATIVE METHOD

JOSEPH FELSENSTEIN

Department of Genetics SK-50, University of Washington, Seattle, Washington 98195

INFERRING
PHYLOGENIES

Joseph Felsenstein

# Data to infer a phylogeny

Data to infer a phylogeny

# Real alignments can be messy

# Alignments tell us a story that can span eons!

Our tree will only be as informative as our alignment and thus our alignments must be assessed/adjusted/QCd before we use them!

- Alignment tool: MAFFT, T-Coffee, Clustal, Muscle, PRANK,

- Assessment: Visually inspect alignments (MEGA, Geneious, Jalview, Clustal X)

- Algorithmically prune sites: Gblocks, TrimAI, T-Coffee + TCS

# Models of sequence evolution describe how DNA changes?

|   | A | G | C | T |
|---|---|---|---|---|
| A | - | α | β | β |
| G | α | - | β | β |
| C | β | β | - | α |
| T | β | β | α | - |

## Kimura's two-parameter model

α represents transitions A<->G or C<->T (purine to purine or pyrimidine to pyrimidine)
β represents transversions transitions between purines and pyrimidines

|   | A | G | C | T |
|---|---|---|---|---|
| A | - | $r_1$ | $r_2$ | $r_3$ |
| G | $r_1$ | - | $r_4$ | $r_5$ |
| C | $r_2$ | $r_4$ | - | $r_6$ |
| T | $r_3$ | $r_5$ | $r_6$ | - |

## General Time Reversible Model

Each type of transition has its own rate but reverse rates equal forward rates

# How do we handle models?

- Just use GTR and forget about it.[*]

- Pick the best model (jModel Test, Model Test NG, IQ-Tree)

- Average across models based on probability (MrBayes, RevBayes)

* Abadi, S., Azouri, D., Pupko, T. and Mayrose, I., 2019. Model selection may not be a mandatory step for phylogeny reconstruction. *Nature communications*, *10*(1), p.934.

# How do we actually choose among trees and branch lengths

| Parsimony | Maximum Likelihood | Bayesian |
|---|---|---|
| | | |

$$P\left( \; \vcenter{\hbox{}} \; \middle| \; \vcenter{\hbox{}} \; \vcenter{\hbox{}} \; \right)$$

1) Pick a random starting value for all parameters
2) Calculate the likelihood above
3) Make a small change to one of the parameters
4) Keep that change if the likelihood of the alignment improved
5) Repeat until you get no further improvement in likelihood

Simple Likelihood Surface

Simple Likelihood Surface

$$P\left(\m\right) = \frac{P\left(\mimage \mid \mimage\right) P\left(\mimage\right)}{P\left(\mimage\right)}$$



Impractical to calculate for non-trivial problems but can be avoided via MCMC

These are the priors

Reverend Thomas Bayes (1701–1761) was an English statistician, philosopher, and Presbyterian minister best known for Bayes' Theorem, which provides a mathematical framework for updating probabilities based on new evidence. His posthumously published work laid the foundation for Bayesian inference, a cornerstone of modern statistics, machine learning, and phylogenetics.

# MCMC algorithm

1) Pick starting values
2) Calculate the probability
3) Make a small change to one parameter
4) Calculate the new probability
5) Accept the changed parameter with this probability
6) Return to step 3

Repeat steps 3-6 1,000,0000s of times

Simple Likelihood Surface

# MCMC Run

Typically, branch lengths will either be expected substitutions per site or time (MY)

Name two species that are sister.

Which species is most closely related to t18?

What does the red dot indicate?

What are the tradeoffs between ML and Bayesian

Why do we need phylogenies?

# If this is your future what next?
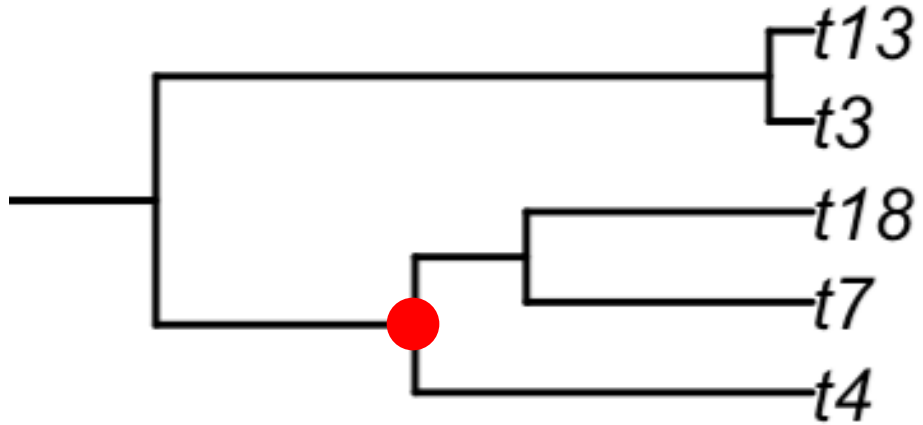
- Read inferring phylogenies!
- Talk to your mentors about good labs that might interest you.
- As you start grad school look for workshops (MBL Mol. Evol. Course)
- Amazing youtube vidoes

# If this is your future what next?

- Break
- Guided walkthrough of running raxml
- Visualizing trees