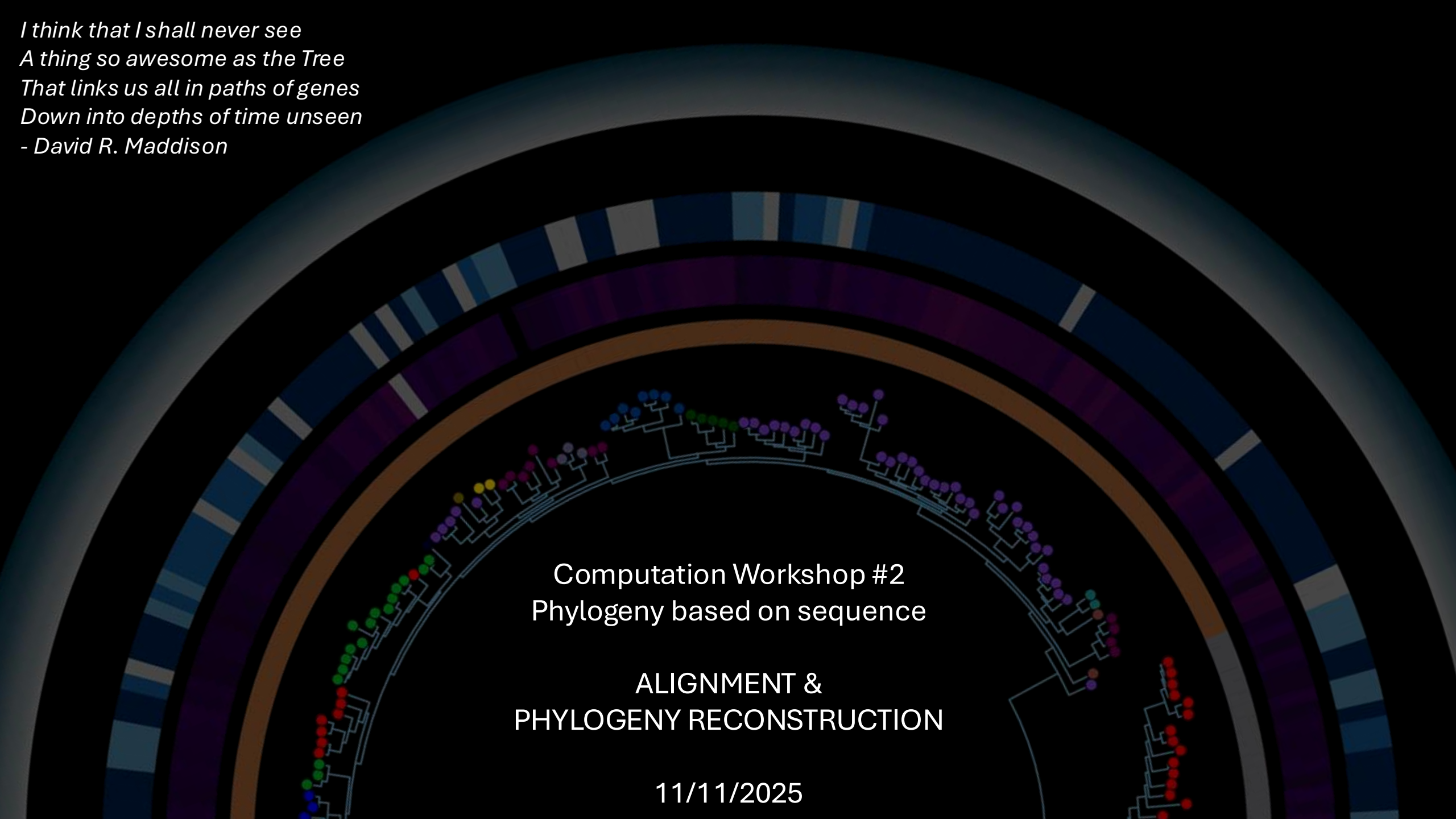*I think that I shall never see*
*A thing so awesome as the Tree*
*That links us all in paths of genes*
*Down into depths of time unseen*
*- David R. Maddison*
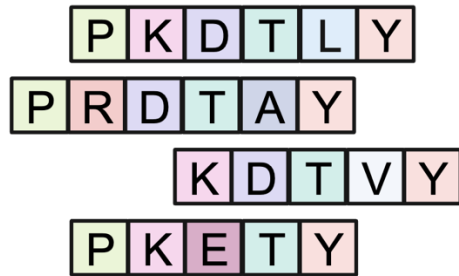
Computation Workshop #2
Phylogeny based on sequence
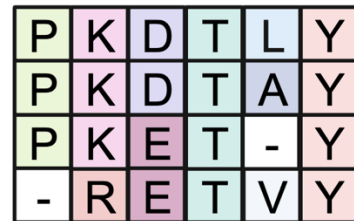
ALIGNMENT &
PHYLOGENY RECONSTRUCTION
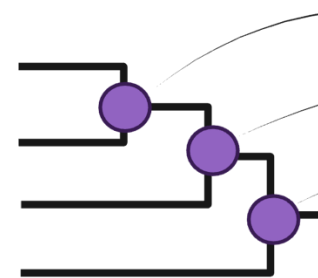
11/11/2025

# ASR Workflow

# 1. Sequence Curation

- Collection of homolog protein sequences from different organisms
- The most important step!

- Need to consider:
  - Which ancestor I am interested in?
  - Which group of modern organisms should I collect?
  - What can I use as outgroup?
  - What should be the size of the final dataset?
  - Single gene or multiple genes?

# 1. Sequence Curation



Query Sequence

↓

*Homology Search*

BLASTp

↓

*Multiple Sequence Alignment*

MAFFT

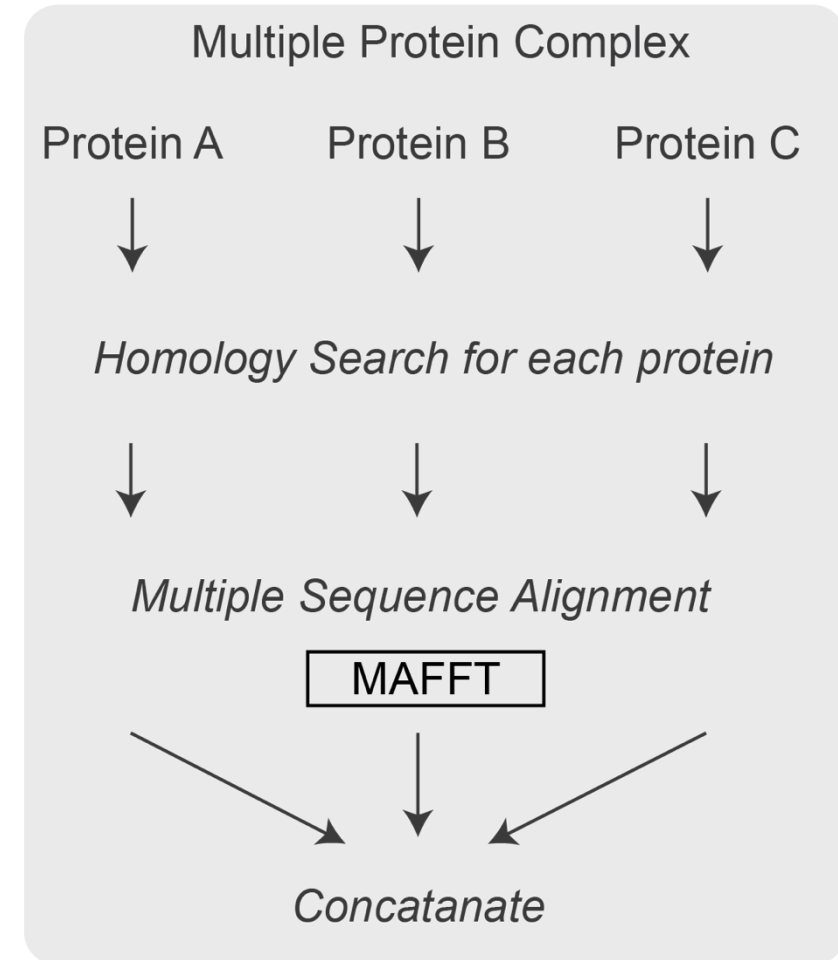*Data Curation & Filtering*

Protein BLAST
protein ▶ protein

Multiple Protein Complex

Protein A          Protein B          Protein C

↓                  ↓                  ↓

*Homology Search for each protein*

↓                  ↓                  ↓

*Multiple Sequence Alignment*

MAFFT

*Concatanate*

Garcia et al., *Methods in Mol. Bio.*, 2022

# 2. Multiple Sequence Alignment (MSA)

| P | K | D | T | L | Y |
|---|---|---|---|---|---|
| P | K | D | T | A | Y |
| P | K | E | T | - | Y |
| - | R | E | T | V | Y |

Goal: introduce gaps into sequences so that columns of alignment contain character states that are homologous.

**MAFFT**

**PRANK**

**MUSCLE**

**Ω CLUSTAL**

# 2. Multiple Sequence Alignment (MSA)

| P | K | D | T | L | Y |
| P | K | D | T | A | Y |
| P | K | E | T | - | Y |
| - | R | E | T | V | Y |

- MSA algorithms aim to evaluate multiple possible alignments and select the one with the <u>highest score</u>.

- Aligning identical characters increases the score, while aligning different characters reduces the score.

- The extent of the score parameters depends on the substitution matrix involved.

**Alignment 1** ✅

| M | L | T | T | T | C |
| M | L | A | - | - | C |

+5   +5   -3   -4   -1   +5   = 7

**Alignment 2**

| M | L | T | T | T | C |
| M | L | - | A | - | C |

+5   +5   -4   -3   -4   +5   = 4

**Scoring Parameters**
Match= +5
Mismatch= -3
Gap open= -4
Gap extension= -1

BLOSUM Substitution Matrix

**Alignment Modulates Ancestral Sequence Reconstruction Accuracy**

Ricardo Assunção Vialle,[†,1,2,3] Asif U. Tamuri,[*,1,4] and Nick Goldman[1]
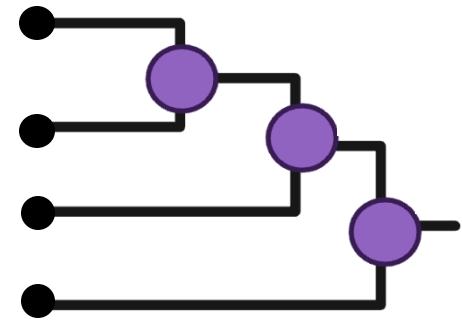
A source of ASR inaccuracy! 😰

# Tutorial time

- Running MAFFT
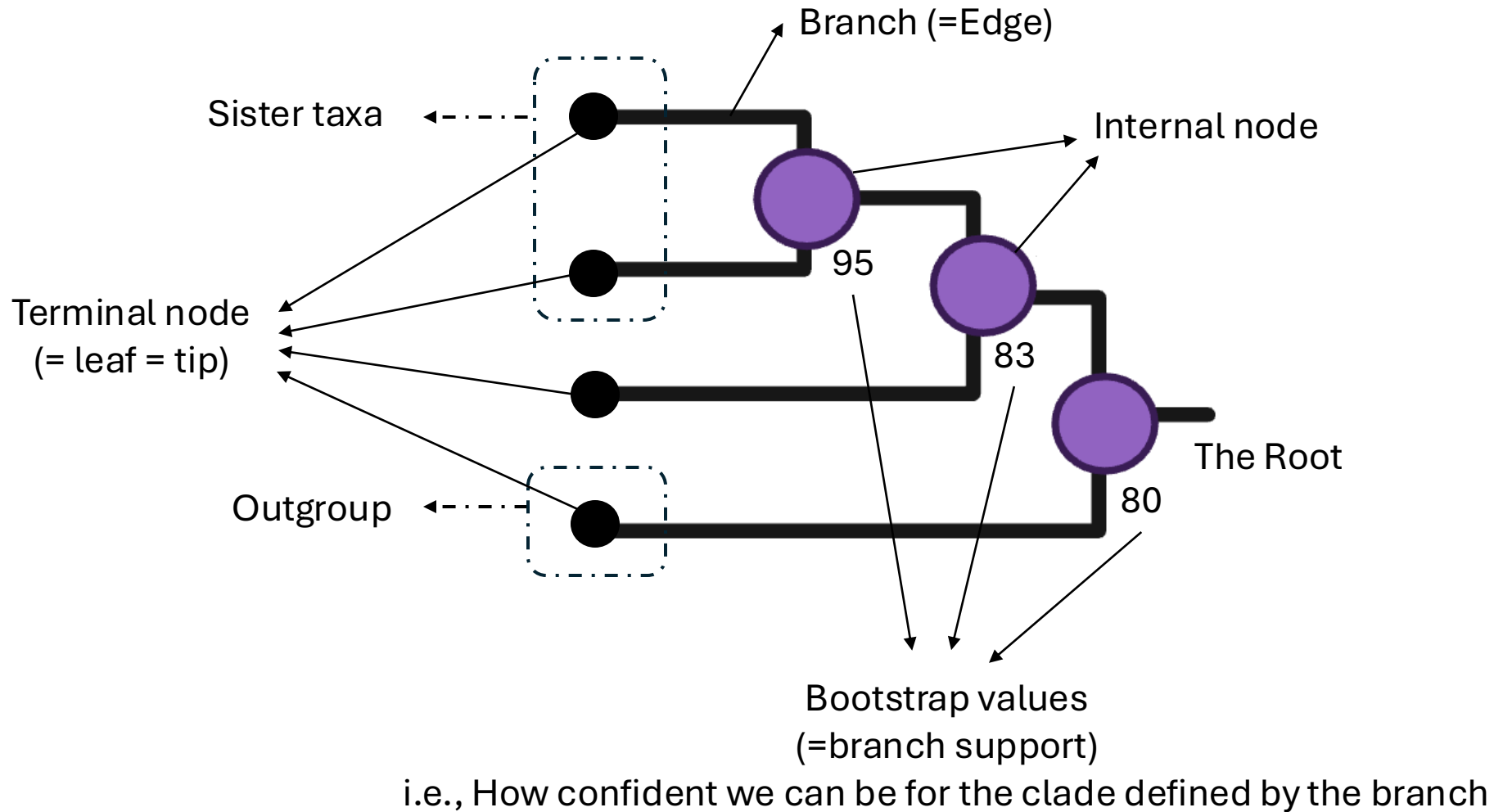
# 3.Phylogeny Reconstruction

A phylogenetic tree:

- Describes evolutionary relationships among a group of organisms based on their molecular sequences.

- Represents a model of evolutionary history depicted by ancestor-descendant relationships between organisms at different level of relatedness.

# Topology of A Phylogenetic Tree



Branch (=Edge)

Sister taxa

Internal node

Terminal node
(= leaf = tip)

95

83

The Root

80

Outgroup

Bootstrap values
(=branch support)
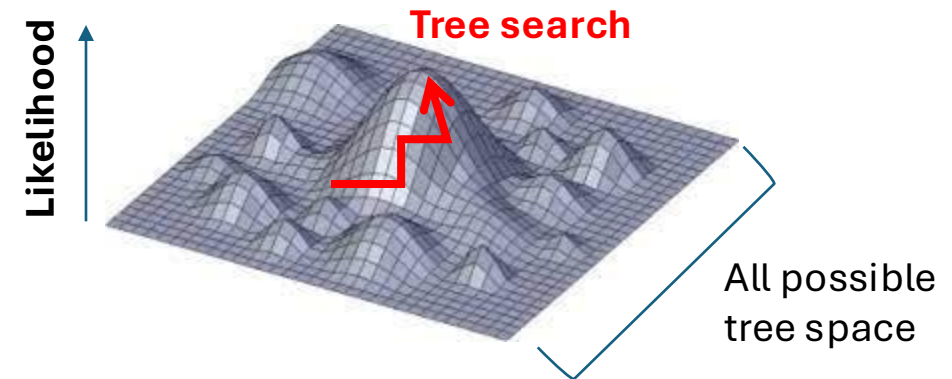i.e., How confident we can be for the clade defined by the branch

# Construction of Phylogenetic Trees

Alignment of extant molecular sequences

A model with rates of different substitutions

- Distance based = Neighbor-joining
- Character based = Maximum parsimony
- Statistical = Maximum Likelihood, Bayesian

Tree search

Likelihood

All possible tree space

IQ-TREE
Efficient software for phylogenomic inference

amkozlov/raxml-ng
RAxML Next Generation: faster, easier-to-use and more flexible

stephaneguindon/ phyml
PhyML — Phylogenetic estimation using (Maximum) Likelihood

PAML: Phylogenetic Analysis by Maximum Likelihood

# Tutorial time

- Running IQTREE
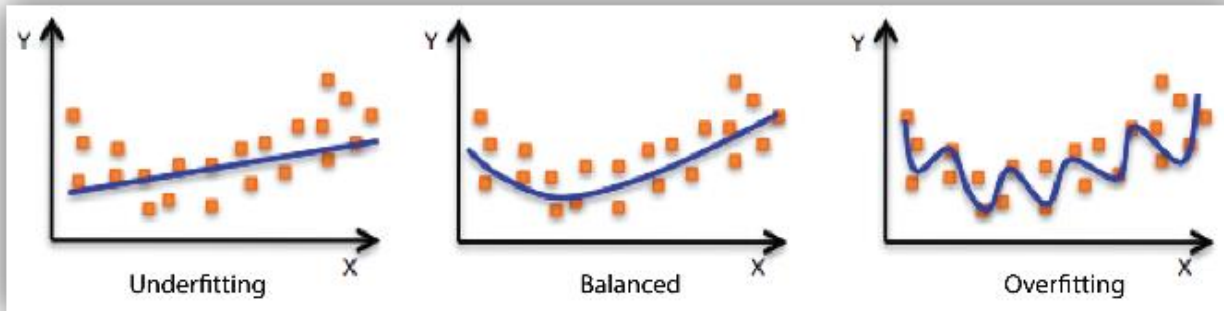
# Why do we need an evolutionary model?



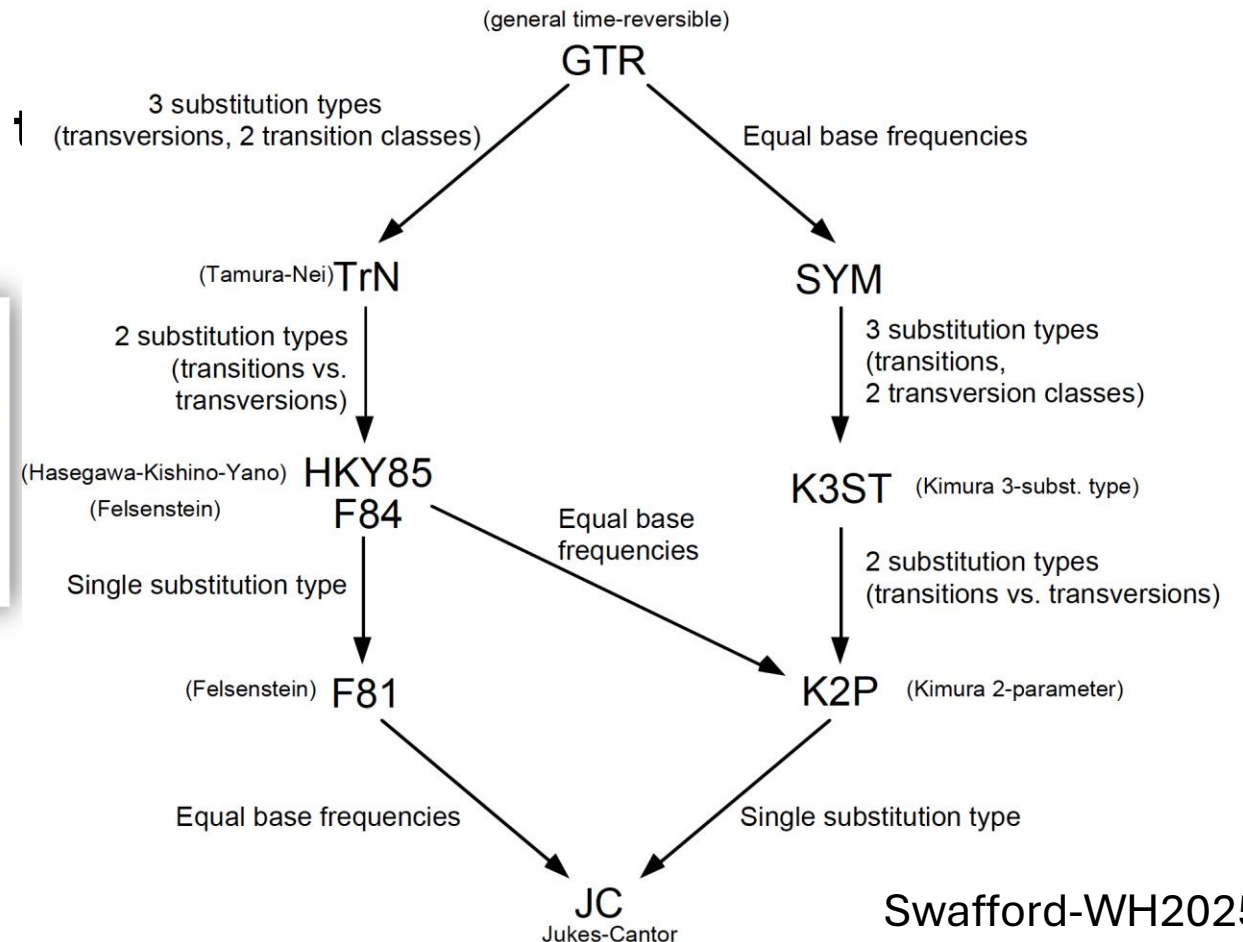**ModelFinder: fast model selection for accurate phylogenetic estimates**

Subha Kalyaanamoorthy, Bui Quang Minh, Thomas K F Wong, Arndt von Haeseler & Lars S Jermiin

**Evolutionary models account for**

- Underfitting models oversimplify molecular evolution
- Overfitting models has too much parameter, takes too much time to calculate likelihood

**Jukes & Cantor (1969)**

$$Q = \begin{pmatrix} -1 & 1/3 & 1/3 & 1/3 \\ 1/3 & -1 & 1/3 & 1/3 \\ 1/3 & 1/3 & -1 & 1/3 \\ 1/3 & 1/3 & 1/3 & -1 \end{pmatrix}$$

**Kimura (1980)**

$$Q = \begin{pmatrix} -1 & 1/(\kappa+2) & \kappa/(\kappa+2) & 1/(\kappa+2) \\ 1/(\kappa+2) & -1 & 1/(\kappa+2) & \kappa/(\kappa+2) \\ \kappa/(\kappa+2) & 1/(\kappa+2) & -1 & 1/(\kappa+2) \\ 1/(\kappa+2) & \kappa/(\kappa+2) & 1/(\kappa+2) & -1 \end{pmatrix}$$

**Hasegawa, Kishino, and Yano (1985)**

$$Q = \begin{pmatrix} - & \pi_C & \kappa\pi_G & \pi_T \\ \pi_A & - & \pi_G & \kappa\pi_T \\ \kappa\pi_A & \pi_C & - & \pi_T \\ \pi_A & \kappa\pi_C & \pi_G & - \end{pmatrix}\mu$$

**GTR (Tavare, 1986)**

$$Q = \begin{pmatrix} - & r_{AC}\pi_C & r_{AG}\pi_G & r_{AT}\pi_T \\ r_{AC}\pi_A & - & r_{CG}\pi_G & r_{CT}\pi_T \\ r_{AG}\pi_A & r_{CG}\pi_C & - & \pi_T \\ r_{AT}\pi_A & r_{CT}\pi_C & \pi_G & - \end{pmatrix}\mu$$

There are two criteria for choosing the best model:
- Akaike Information Criterion (AIC)
- Bayesian Information Criterion (BIC)

Both try to find the balance between underfitting and overfitting models

- Akaike information criterion (AIC)

$$AIC_i = -2 \ln L_i + 2k$$

where *k* is the number of free parameters estimated

- AICc (corrected AIC)

$$AIC_c = AIC + \frac{(2k(k+1))}{(n-k-1)}$$

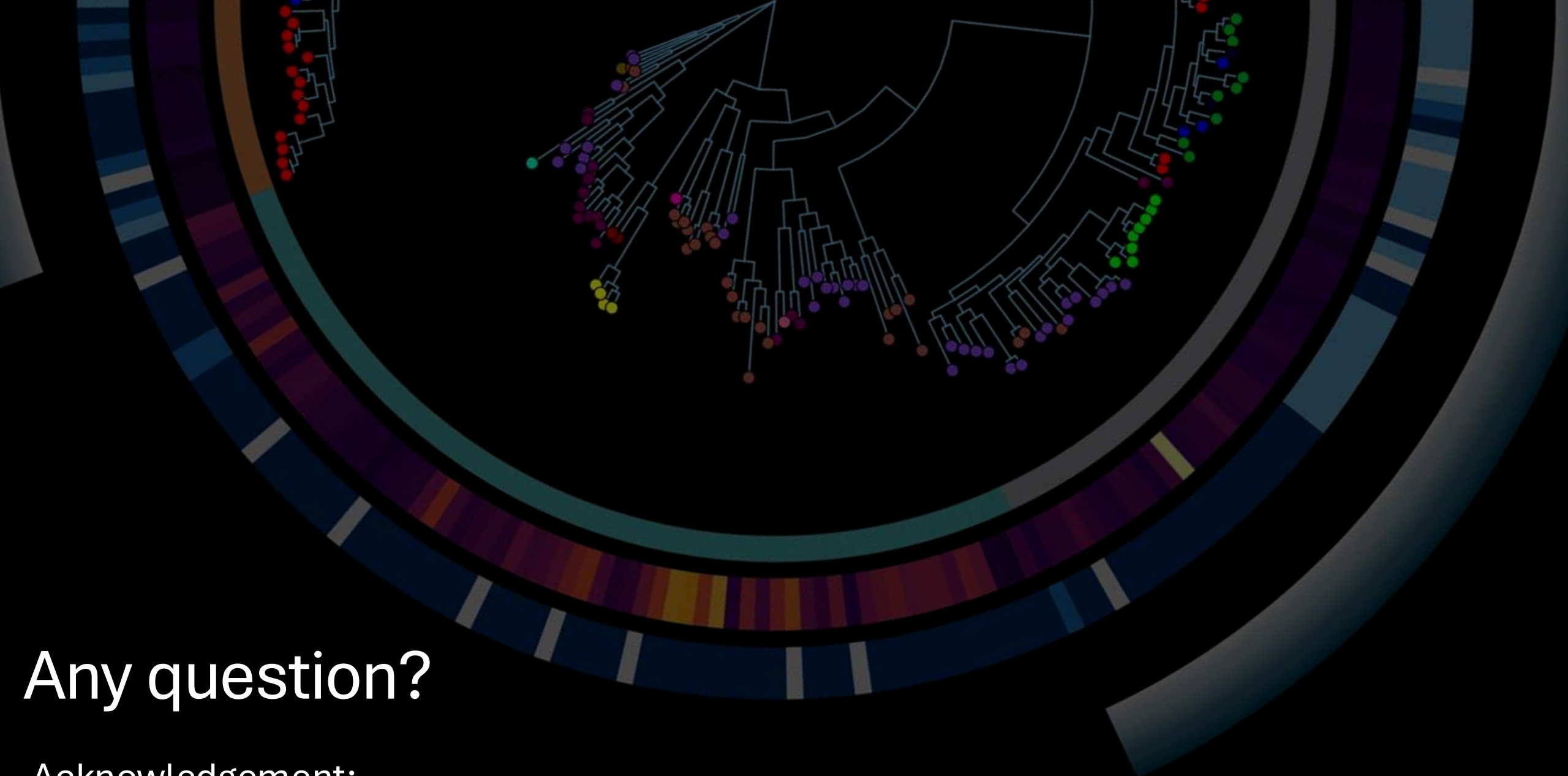- Bayesian information criterion (BIC)

$$BIC_i = -2 \ln L_i + k \ln n$$

where *k* is the number of free parameters estimated and *n* is the "sample size" (typically number of sites)

# AIC(c) vs. BIC

- BIC performs well when true model is contained in model set, and among a set of simple-ish models, AIC often selects a more complex model than the truth (indeed, AIC is formally statistically inconsistent)

- But in phylogenetics, no model is as complex as the truth, and the true model will never be contained in the model set.

- BIC often chooses models that seem *too* simple!.

# Tutorial time

- Navigating output

# Any question?

Acknowledgement:

- Some slides are modified from Dr. Amanda Garcia's ASR Workshop presentation