

A QUICK INTRO TO:  
**PHYLOGENOMICS**

MARIA VALDEZ CABRERA - POSTDOCTORAL FELLOW

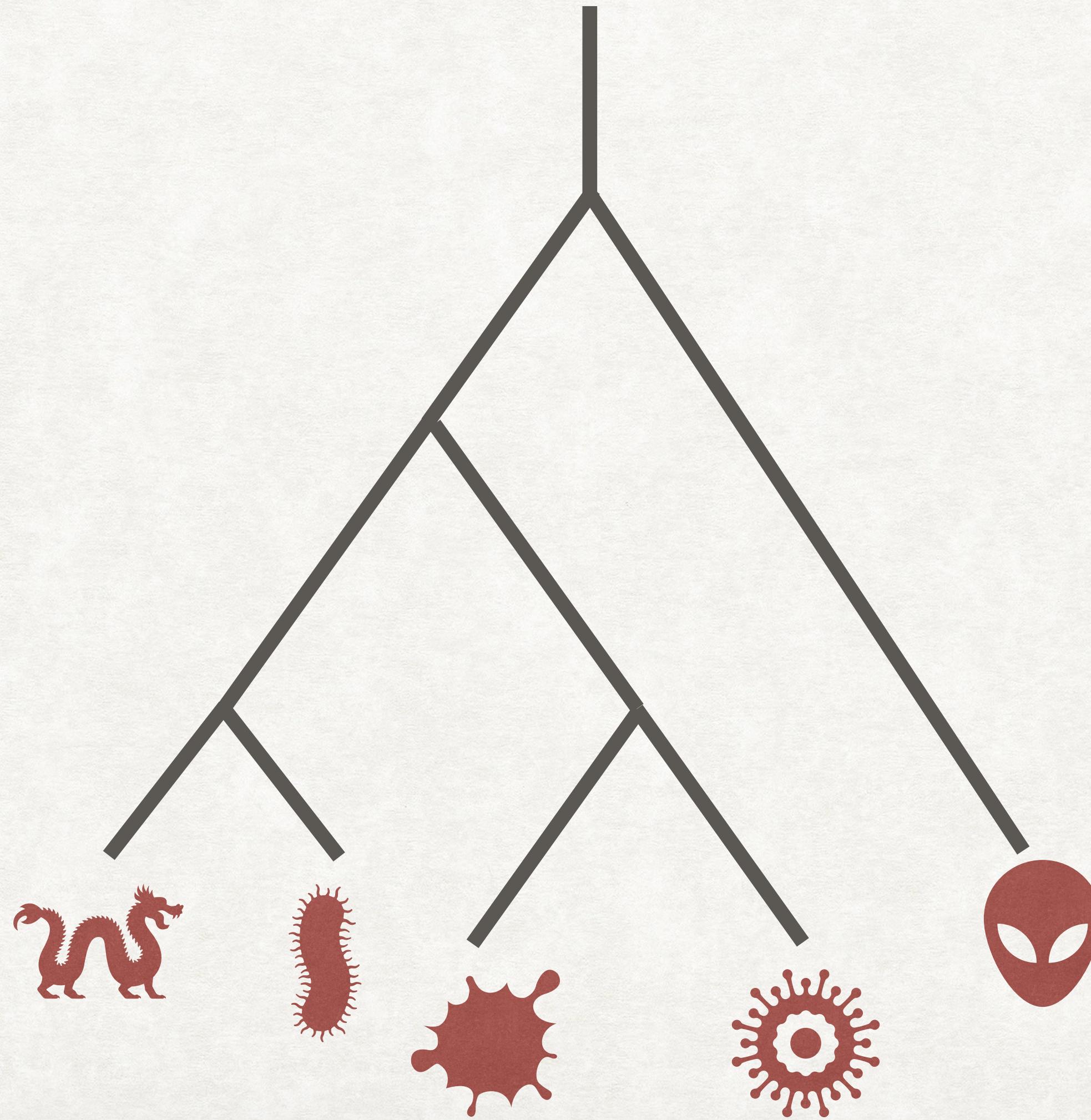
# OBJECTIVES FOR THE HOUR

- 
- Getting phylogenetic trees on your mind.
  - Giving you a statistician's point of view of these objects.
  - Discussing the potential and limitations of phylogenetic tree construction for microbiome data
  - Introducing you to the fabulous GToTree workflow by Mike Lee

# OBJECTIVES FOR THE HOUR

- Getting phylogenetic trees on your mind.
  - Giving you a statistician's point of view of these objects.
  - Discussing the potential and limitations of phylogenetic tree construction for microbiome data
- 
- GToTree** 
- Introducing you to the fabulous GToTree workflow by Mike Lee

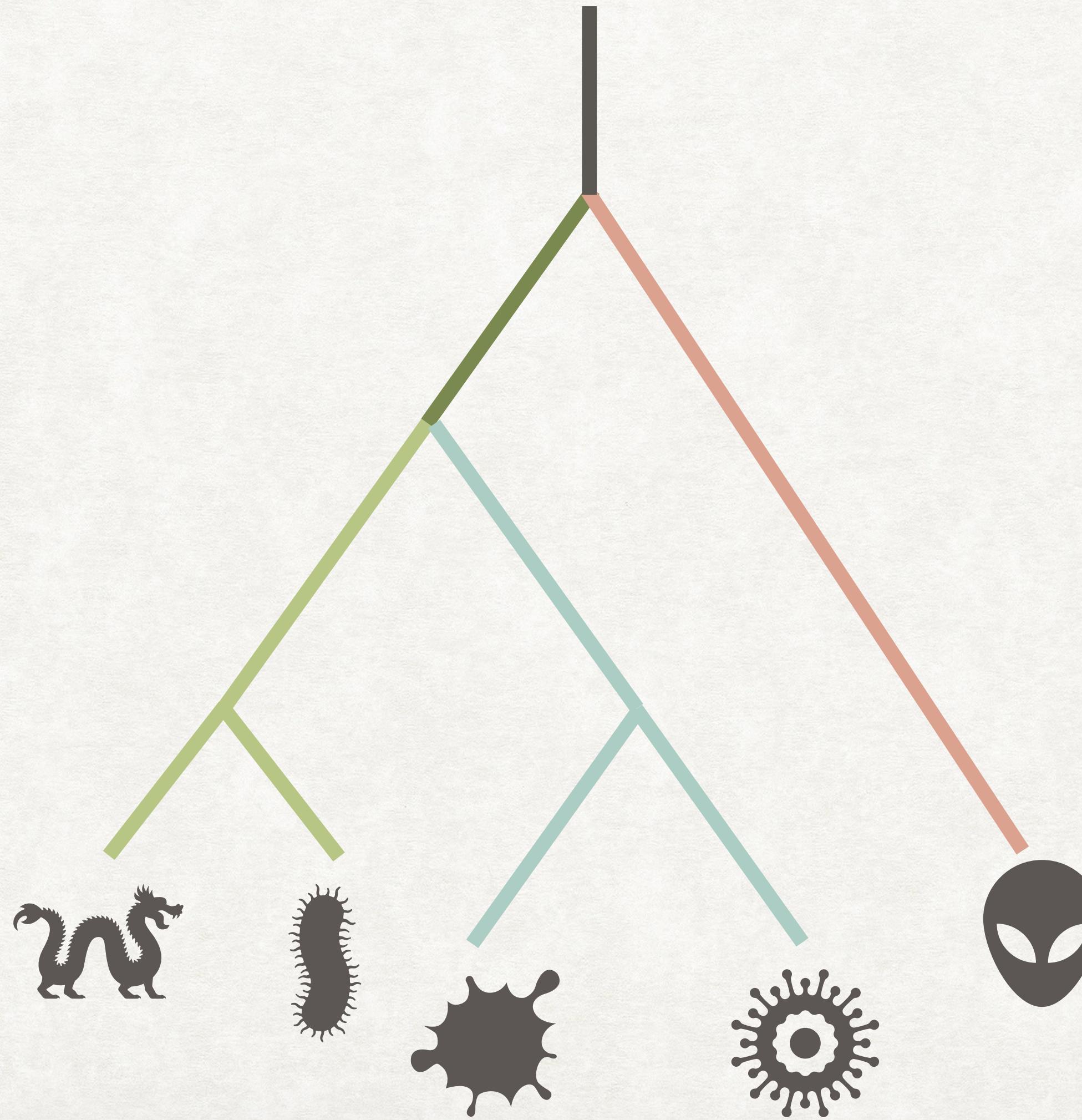
# WHAT EVEN IS A PHYLOGENETIC TREE?



## Components:

- Leaves
- Shape (topology)
- Branch lengths
- Root

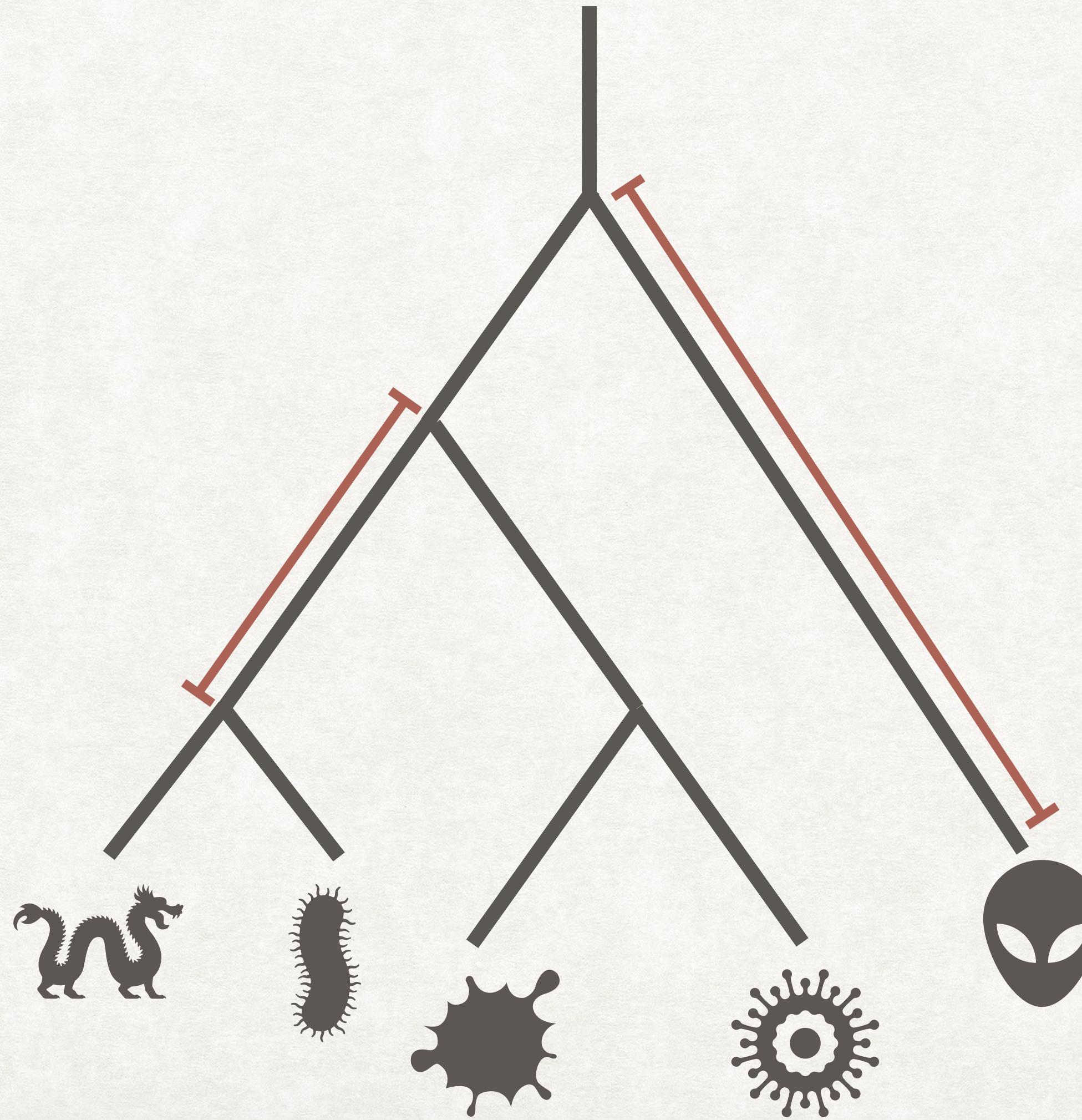
# WHAT EVEN IS A PHYLOGENETIC TREE?



## Components:

- Leaves
- Shape (topology)
- Branch lengths
- Root

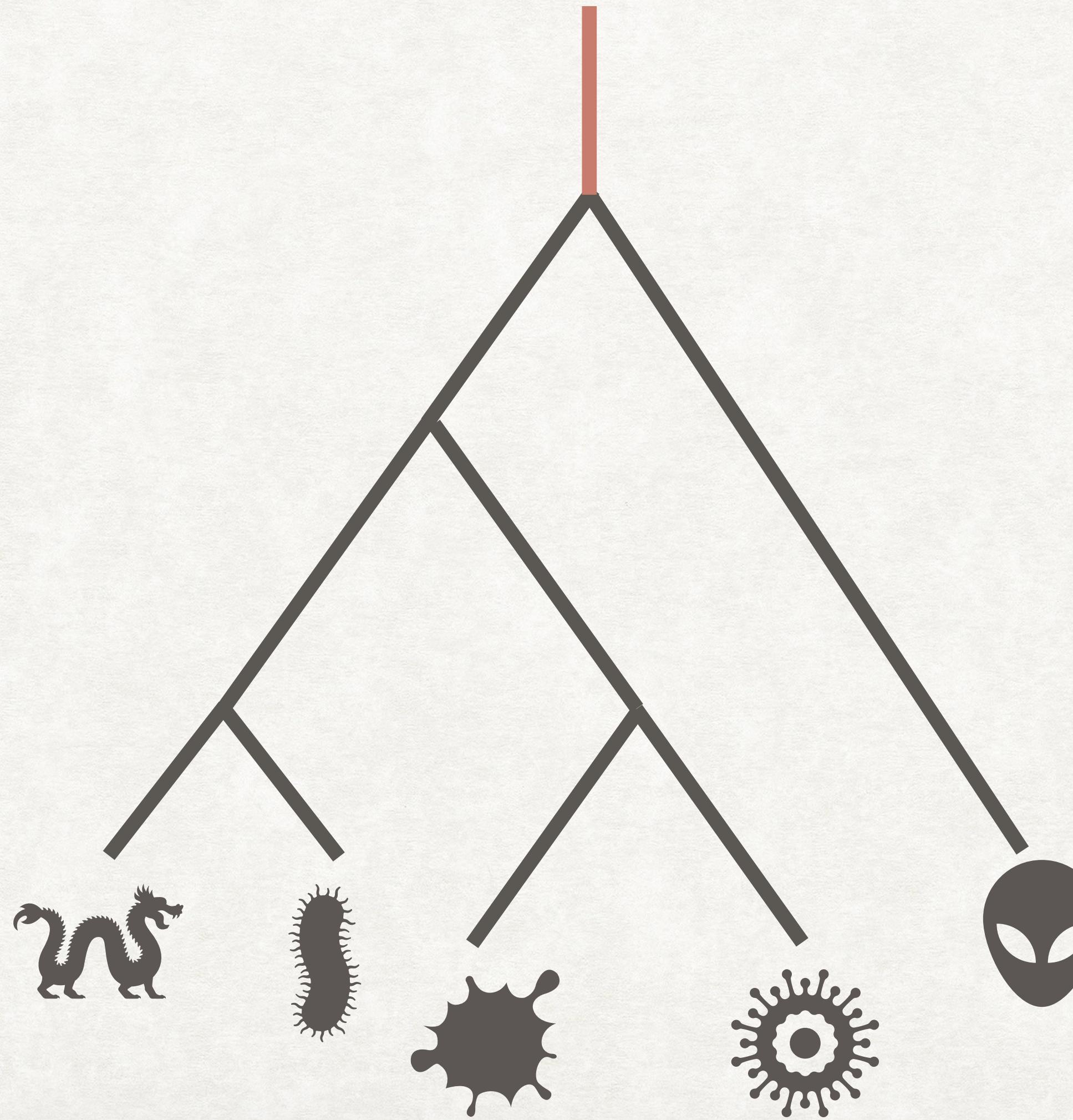
# WHAT EVEN IS A PHYLOGENETIC TREE?



## Components:

- Leaves
- Shape (topology)
- **Branch lengths**
- Root

# WHAT EVEN IS A PHYLOGENETIC TREE?



## Components:

- Leaves
- Shape (topology)
- Branch lengths
- Root

## ACTIVITY ALERT

What scientific questions can be explored with Trees?

# ★ ACTIVITY ALERT ★

## What scientific questions can be explored with Trees?

### Taxonomy

Contents lists available at ScienceDirect

**Diagnostic Microbiology and Infectious Disease**

journal homepage: [www.elsevier.com/locate/diagmicrobio](http://www.elsevier.com/locate/diagmicrobio)

ELSEVIER

Review Article

**Taxonomic update on proposed nomenclature and classification changes for bacteria of medical importance, 2015**

J. Michael Janda \*

Public Health Laboratory, Department of Public Health, Kern County, Bakersfield, CA 93306-3302

**ARTICLE INFO**

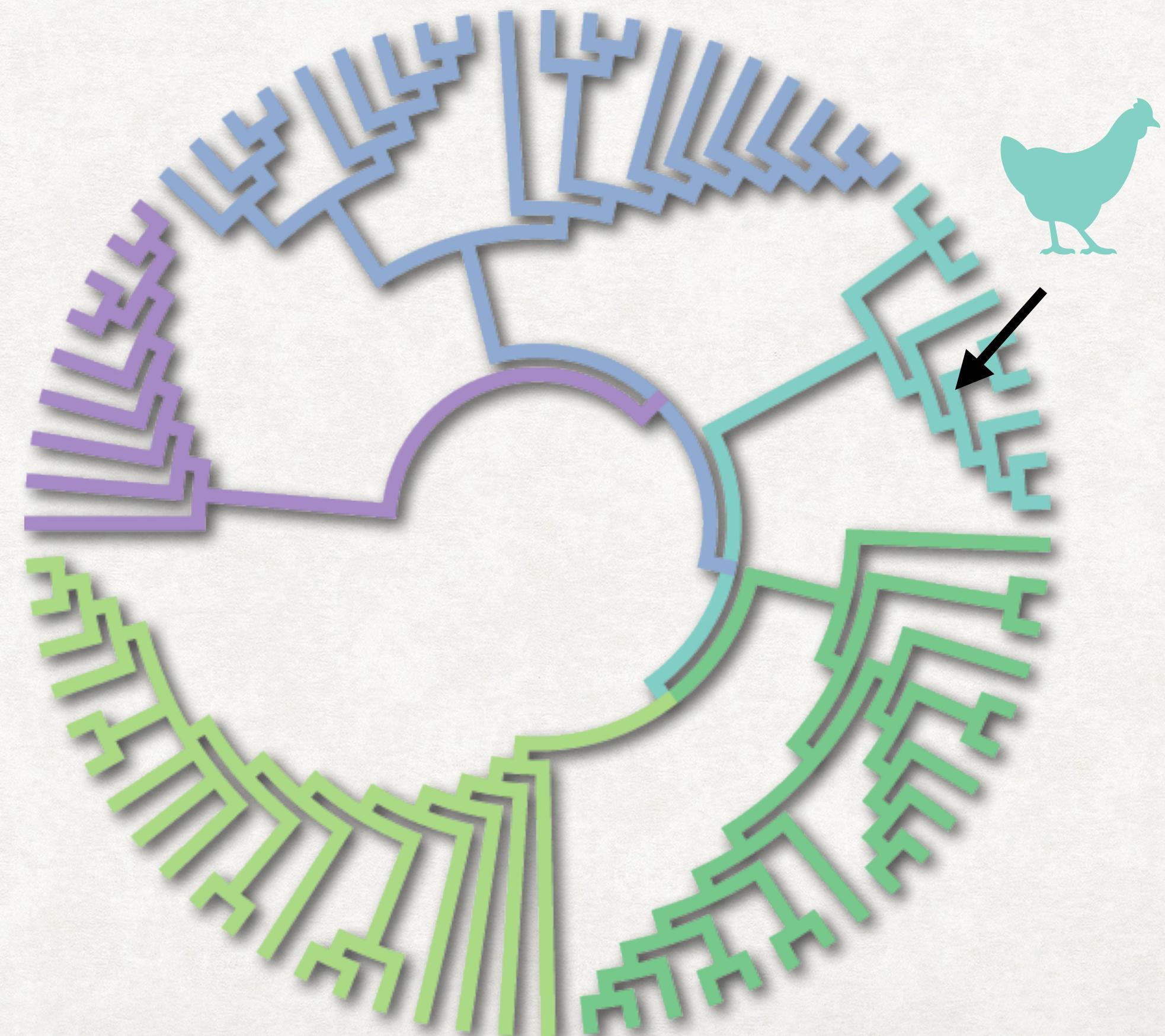
**Article history:**  
Received 10 May 2016  
Received in revised form 28 June 2016  
Accepted 30 June 2016  
Available online 4 July 2016

**Keywords:**  
Nomenclature  
Bacterial taxonomy  
Classification  
Updates

**ABSTRACT**

A key aspect of medical, public health, and diagnostic microbiology laboratories is the accurate and rapid reporting and communication regarding infectious agents of clinical significance. Microbial taxonomy in the age of molecular diagnostics and phylogenetics creates changes in taxonomy at a rapid rate further complicating this process. This update focuses on the description of new species and classification changes proposed in 2015.

© 2016 Elsevier Inc. All rights reserved.



# ACTIVITY ALERT

## What scientific questions can be explored with Trees?

### ARTICLE

<https://doi.org/10.1038/s41467-019-13443-4>

OPEN

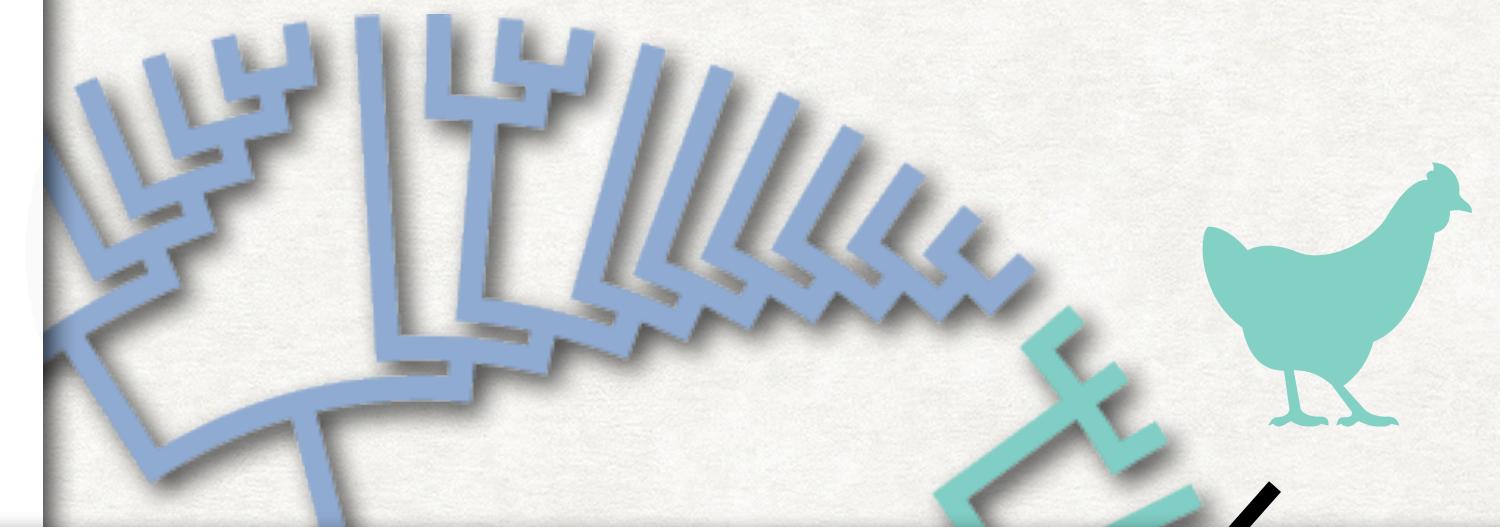
### Phylogenomics of 10,575 genomes reveals evolutionary proximity between domains Bacteria and Archaea

Qiyun Zhu<sup>1,19</sup>, Uyen Mai<sup>2,19</sup>, Wayne Pfeiffer<sup>3</sup>, Stefan Janssen<sup>1,4</sup>, Pedro Belda-Ferre<sup>1</sup>, Gabriel A. Al-Ghalith<sup>6</sup>, Evguenia Kopylova<sup>1</sup>, John B. Yin<sup>8,9</sup>, Shi Huang<sup>1,10</sup>, Nimaichand Salam<sup>11</sup>, Jian-Yu Jiao<sup>11</sup>, Zijun Yimeng Yang<sup>6</sup>, Erfan Sayyari<sup>8</sup>, Maryam Rabiee<sup>2</sup>, James T. Morton<sup>12</sup>, Wen-Jun Li<sup>11</sup>, Curtis Huttenhower<sup>14,15</sup>, Nicola Segata<sup>1,5</sup>, Larry Rob Knight<sup>1,2,16,18\*</sup>

Available online 4 July 2016

Keywords:  
Nomenclature  
Bacterial taxonomy  
Classification  
Updates

© 2016



### Insights into the phylogeny and coding potential of microbial dark matter

**Authors:** Christian Rinke, Patrick Schwientek, Alexander Sczyrba, Natalia N. Ivanova, Iain J. Anderson and Jan-Fang Cheng

**Date:** July 25, 2013

**From:** Nature(Vol. 499, Issue 7459)

**Publisher:** Nature Publishing Group

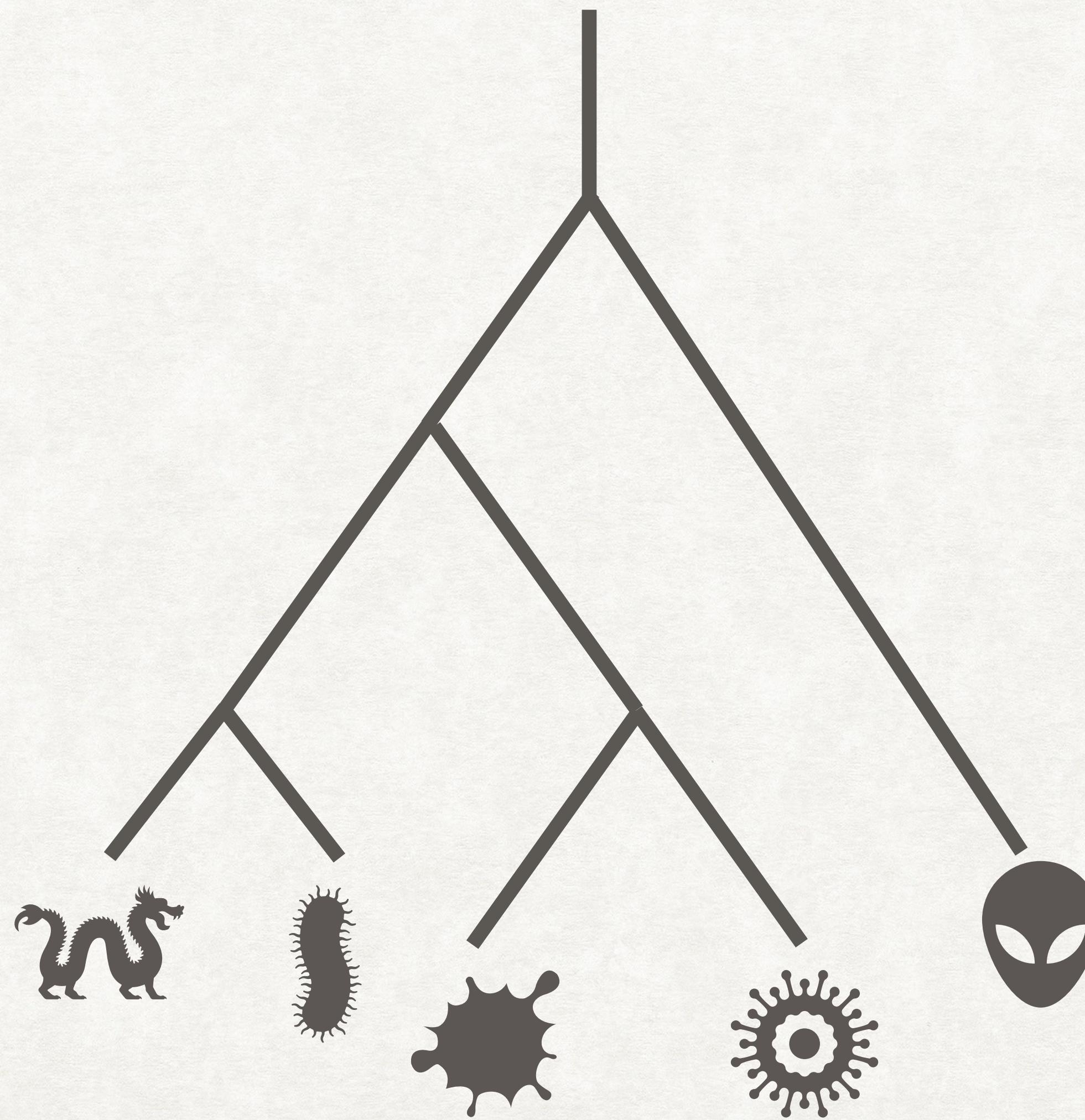
**Document Type:** Report

**Length:** 5,379 words

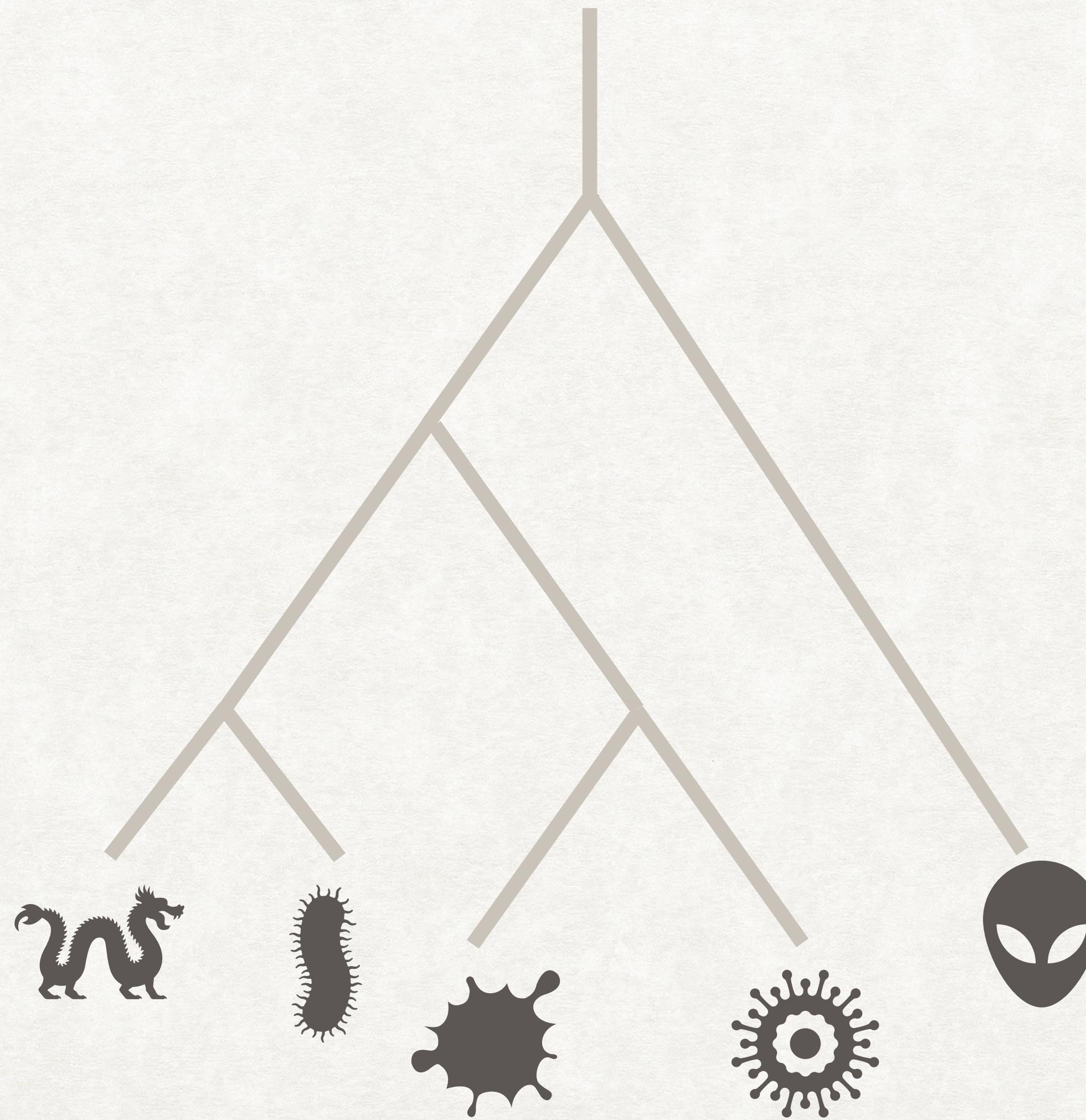
**DOI:** <http://dx.doi.org/10.1038/nature12352>

Abstract:

# WHAT EVEN IS A PHYLOGENETIC TREE?

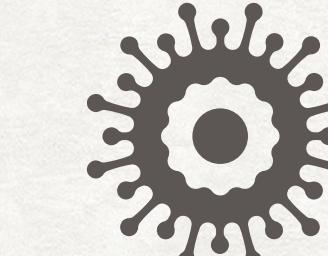
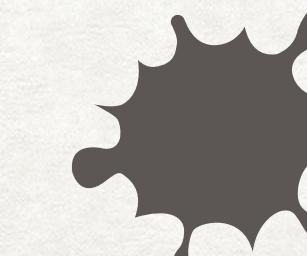
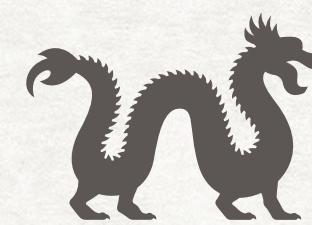


# WHAT EVEN IS A PHYLOGENETIC TREE?



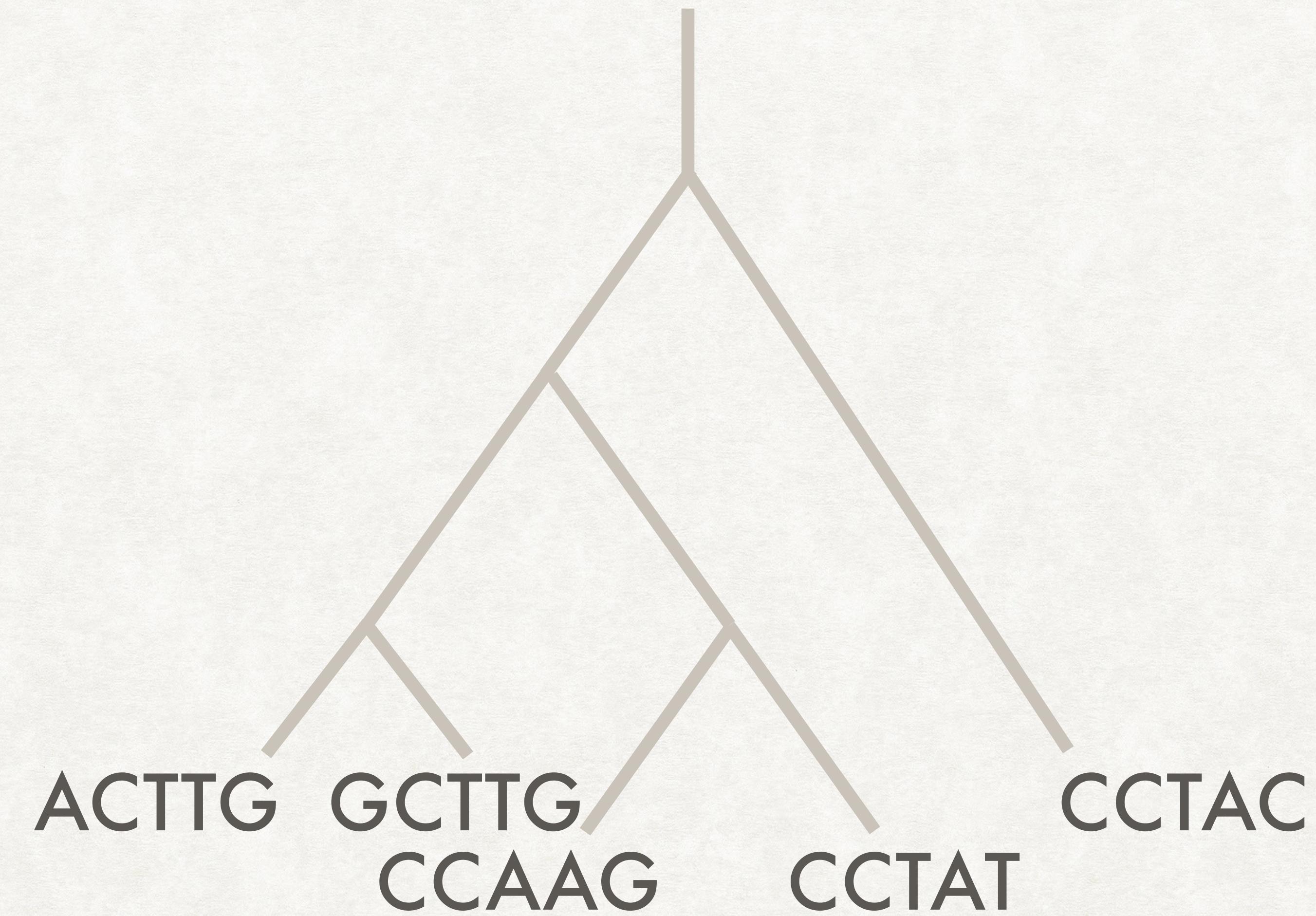
# WHAT EVEN IS A PHYLOGENETIC TREE?

**AN ESTIMATE!!**  
**BUT OF WHAT EXACTLY?\* ...**

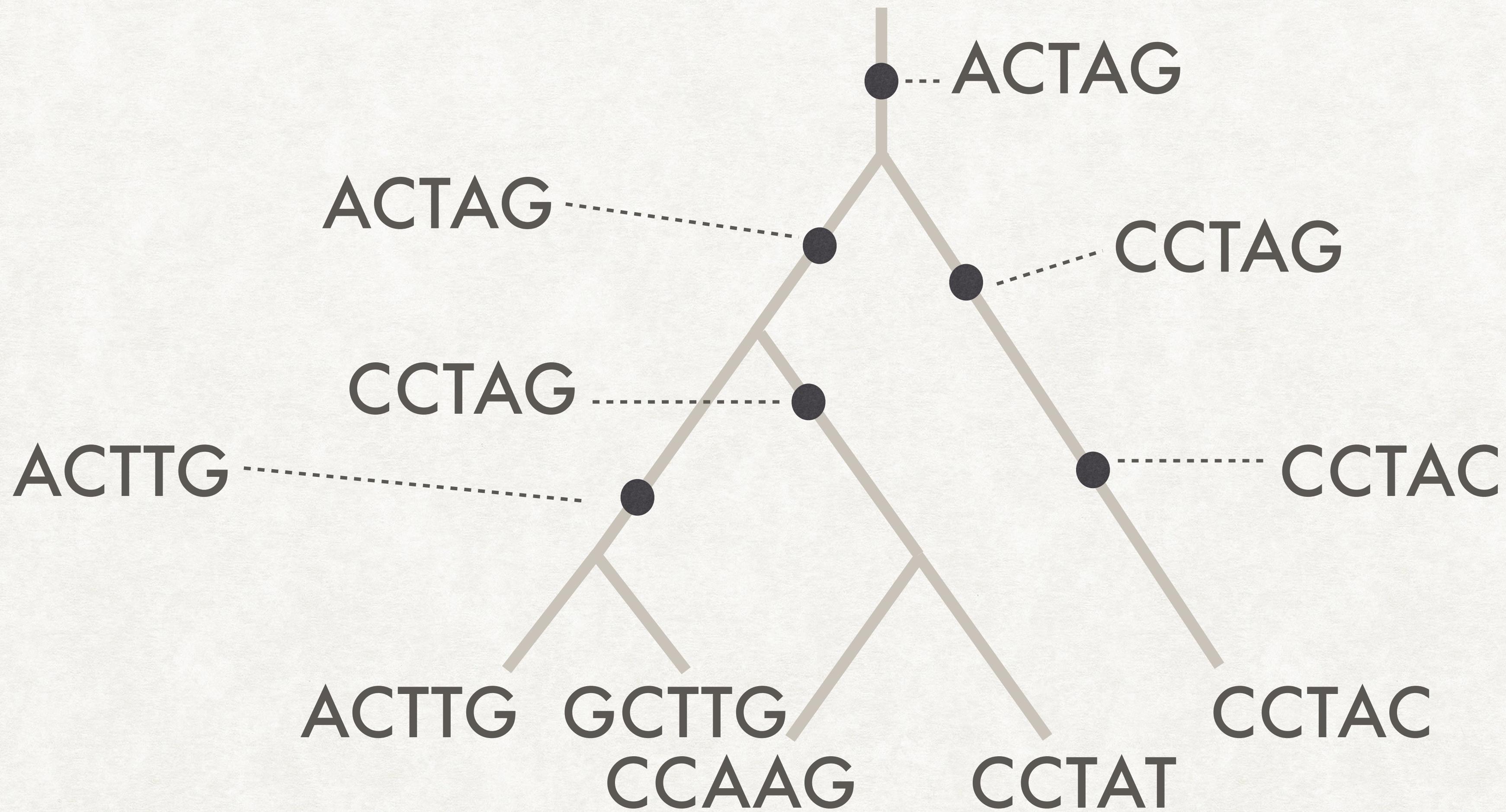


\*Not a one correct answer, and very complex when talking about microorganisms

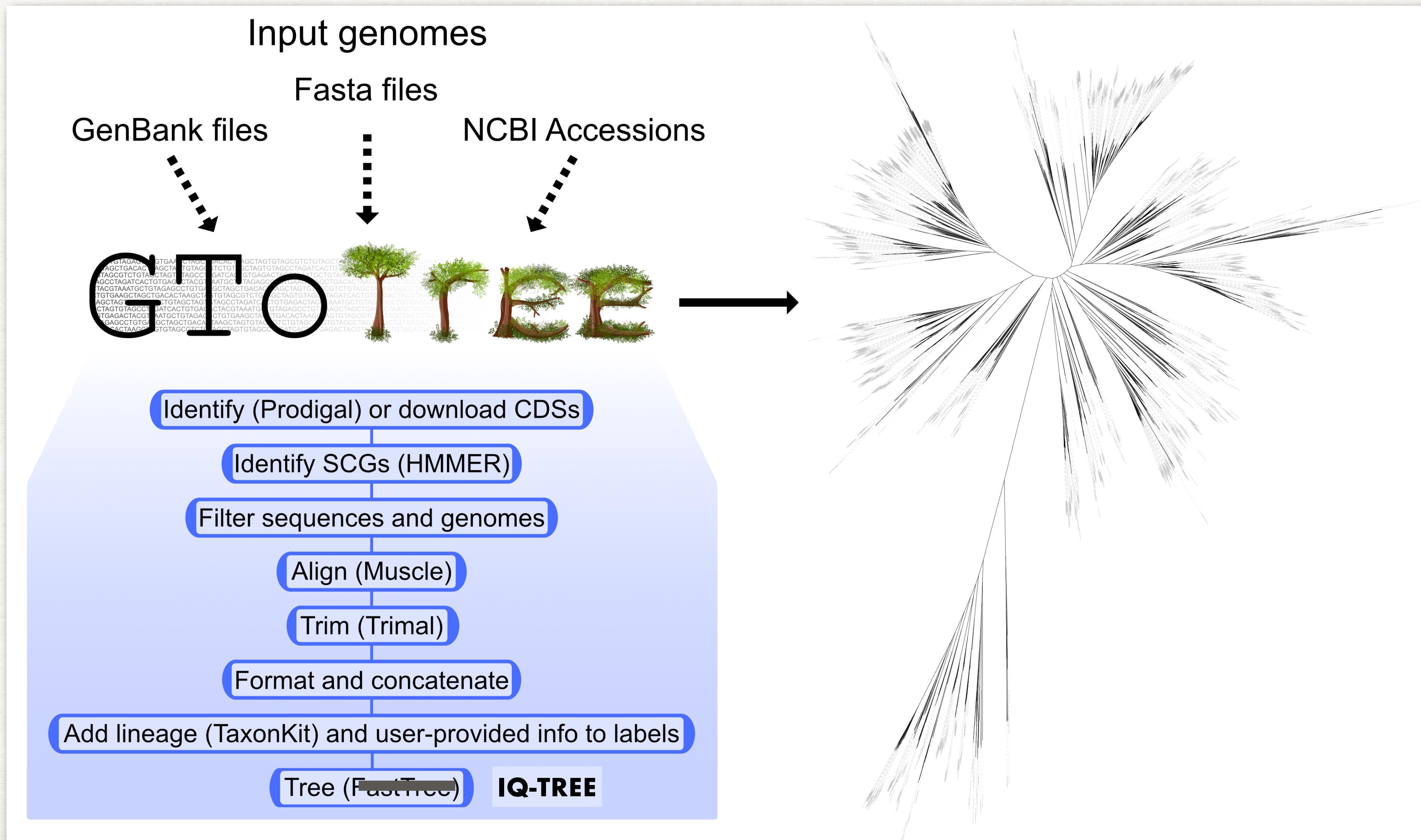
# TREE ESTIMATION



# TREE ESTIMATION

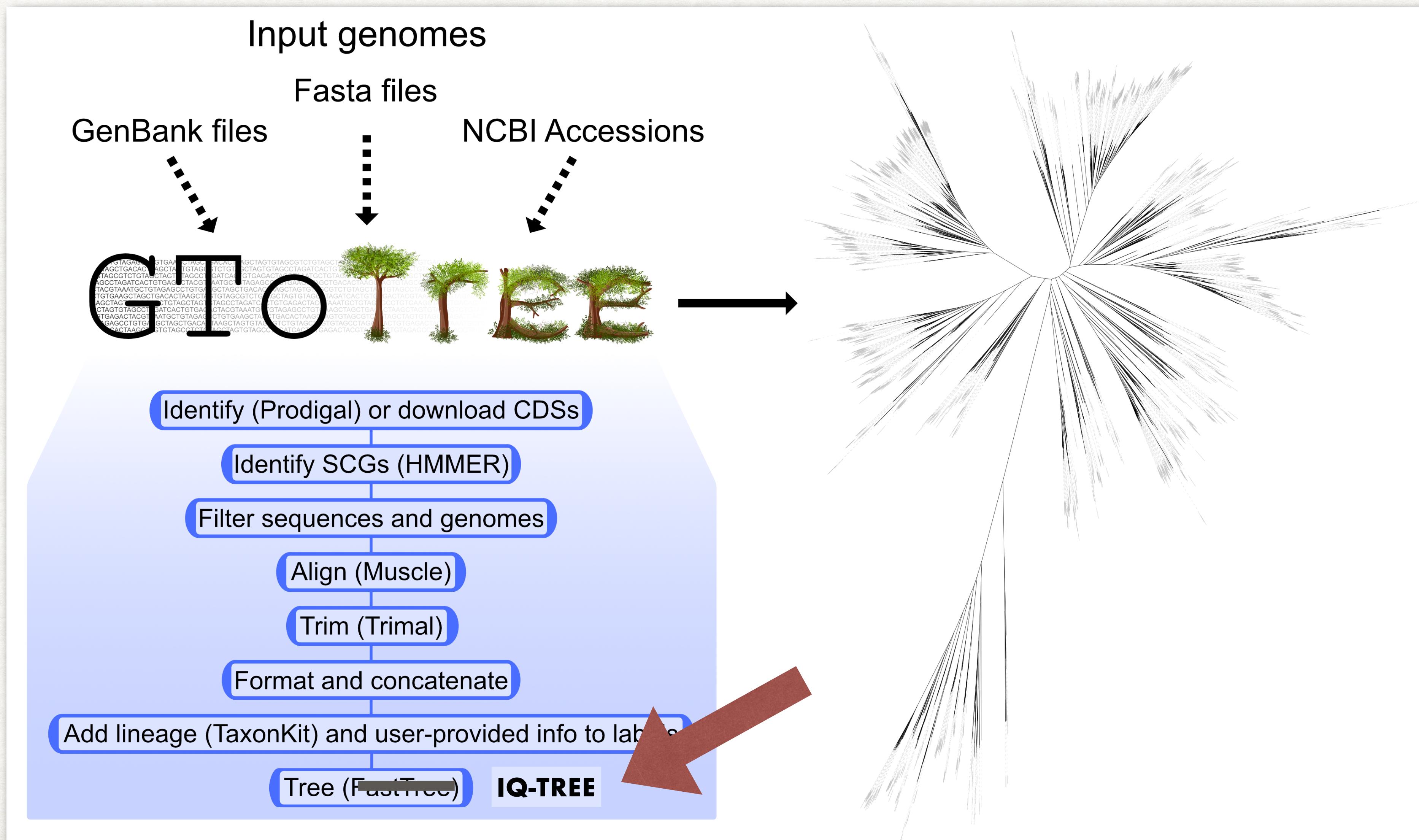


# TREE ESTIMATION

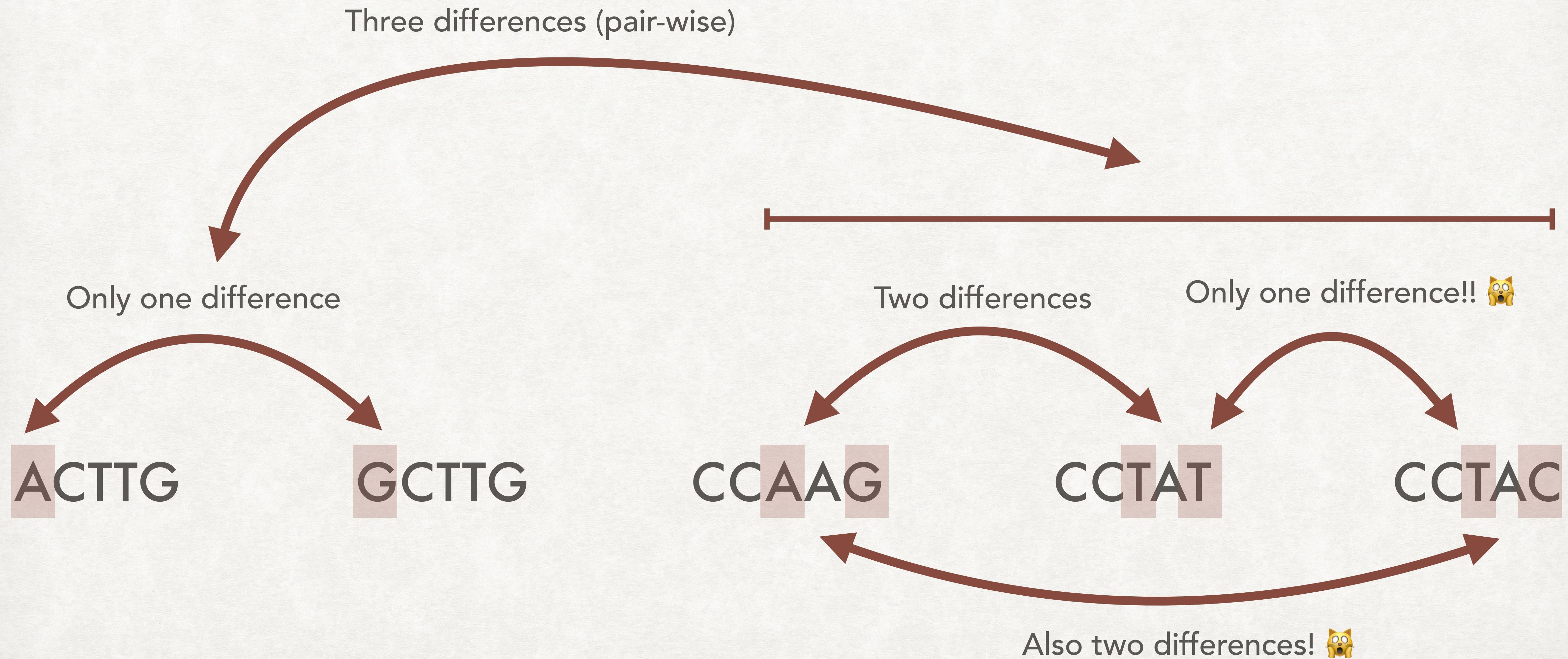


Access Lab instructions [here](#)

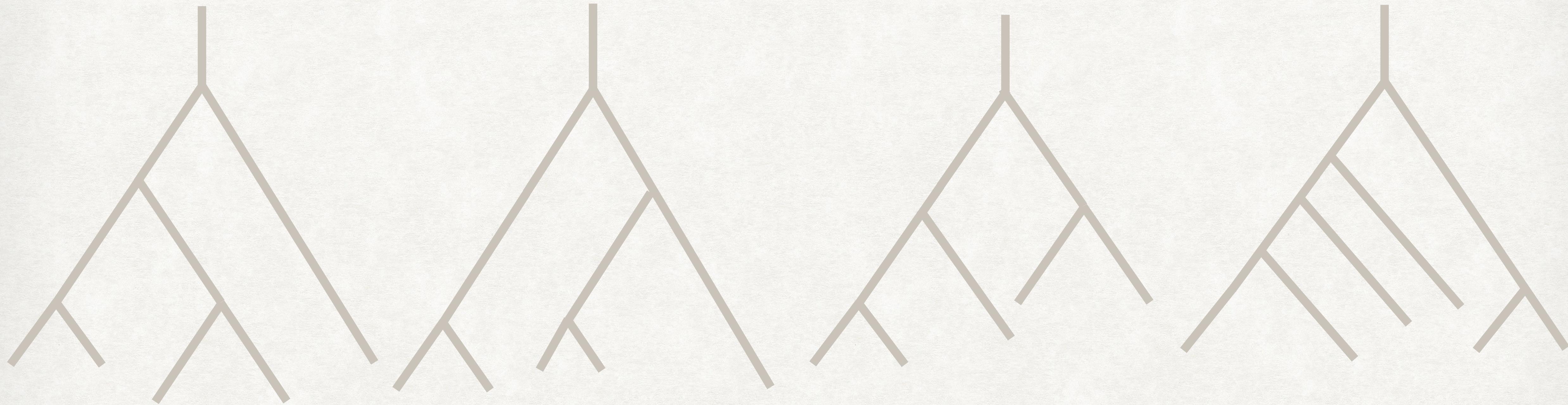
# TREE ESTIMATION



# TREE ESTIMATION



# TREE ESTIMATION



ACTTG

GCTTG

CCAAG

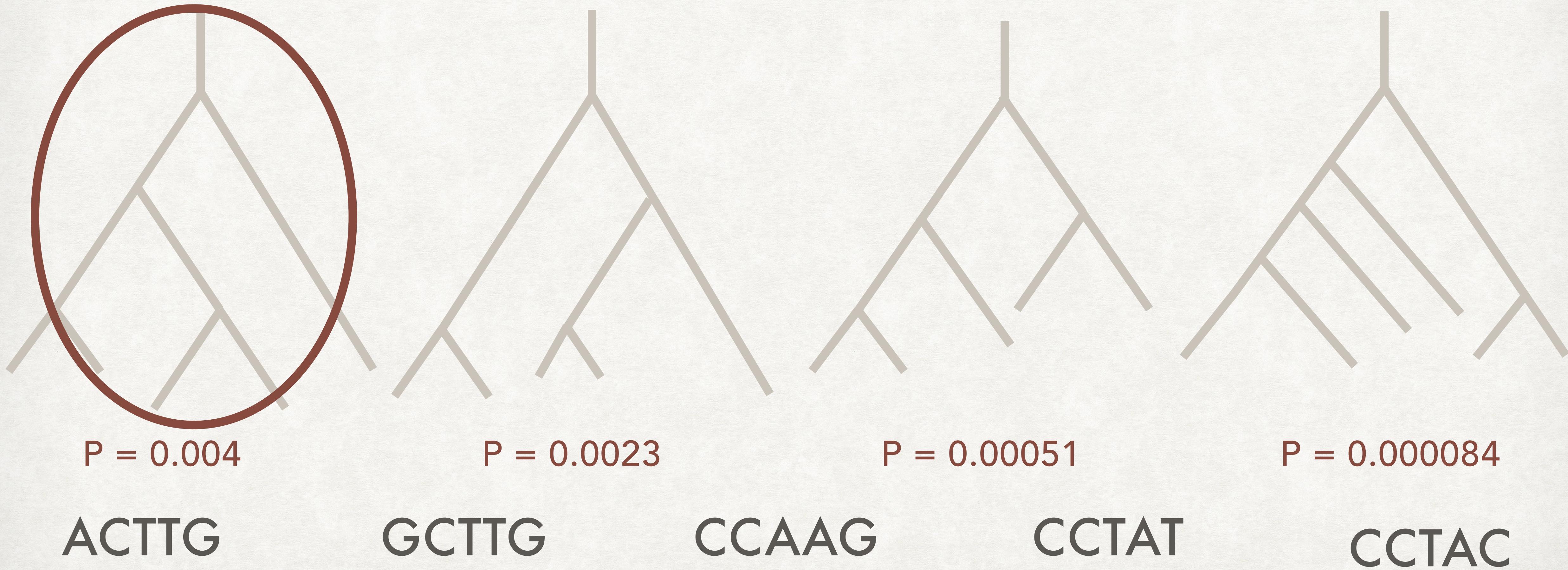
CCTAT

CCTAC

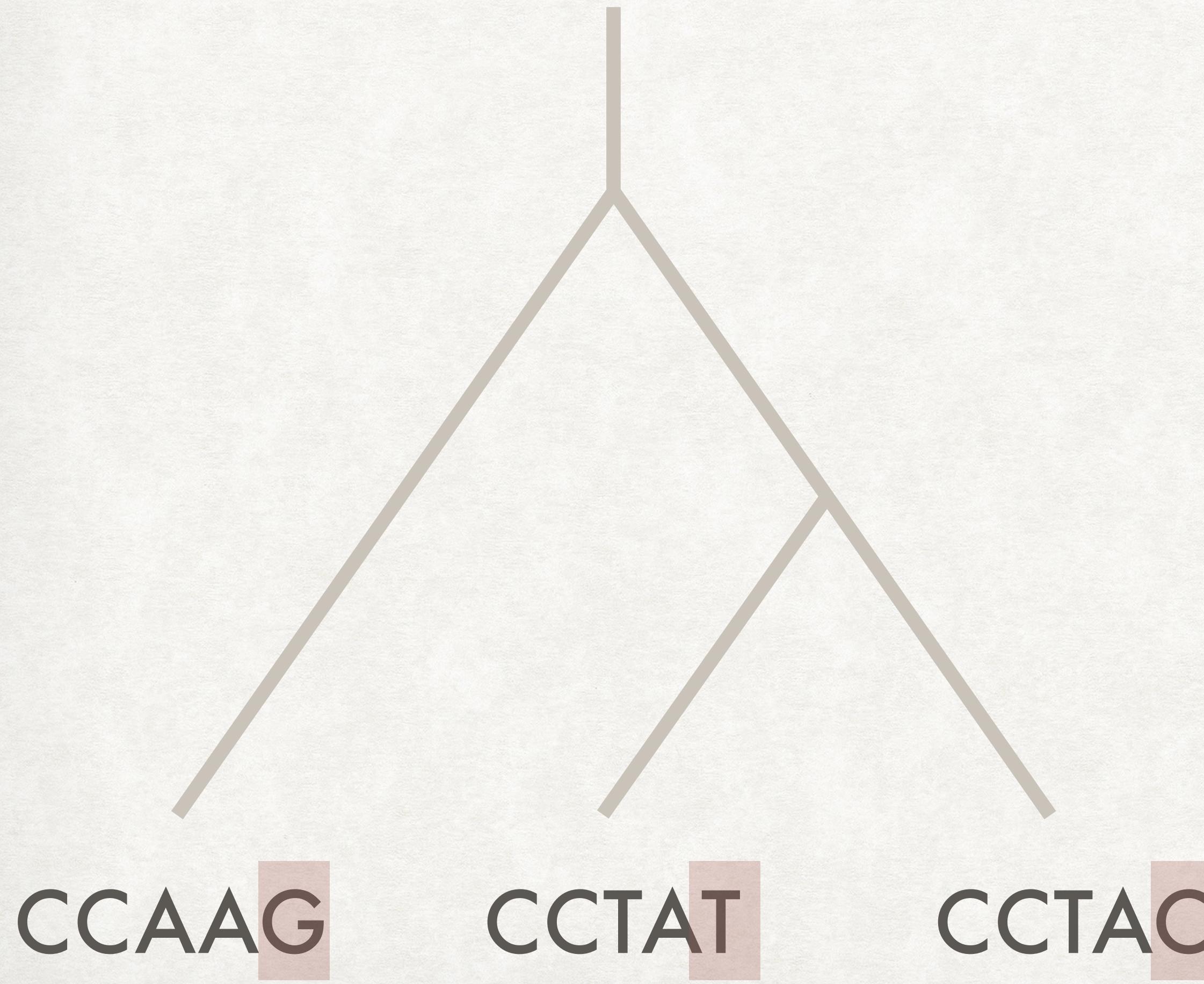
# TREE ESTIMATION

Which of these trees is more "likely" to have produced the genomes?

$$P(\text{DNA}, \dots | \text{Tree}, \Theta)$$



# TREE ESTIMATION

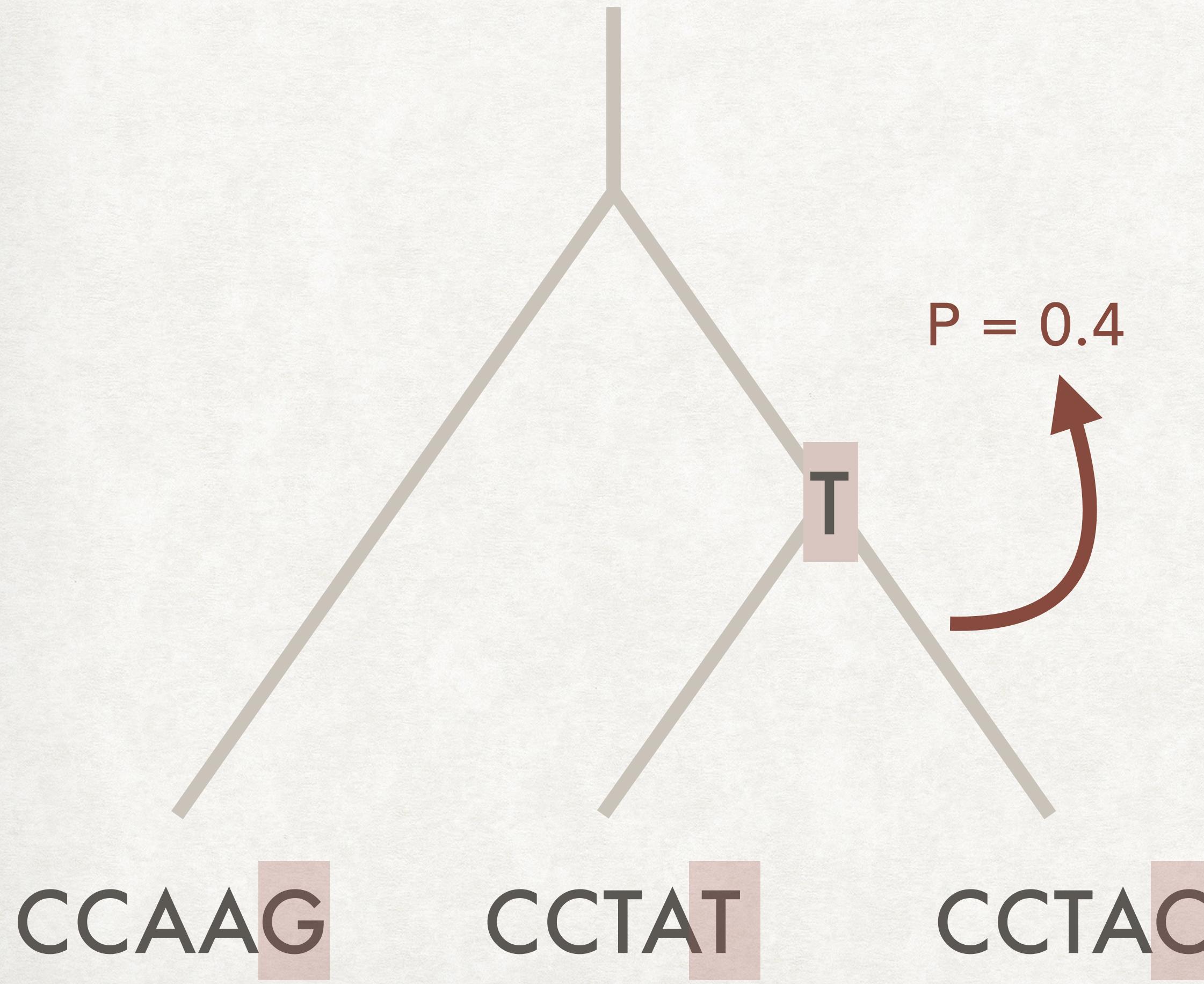


$\mu$  : Rate of change ( $\mu \times$  branch length = expected changes in one site)

$r(A \leftrightarrow C), r(A \leftrightarrow G), \dots, r(G \leftrightarrow T)$ :  
Proportion of change between different nucleotides

$\pi_A, \pi_C, \pi_G, \pi_T$ : "Equilibrium" frequencies of each nucleotide.

# TREE ESTIMATION

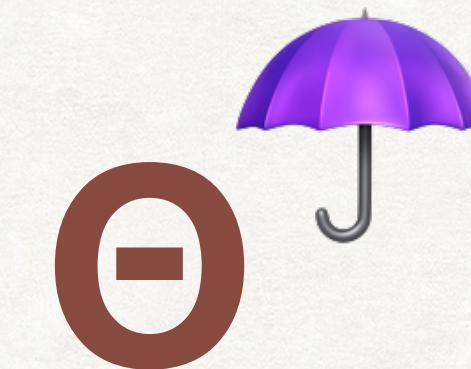
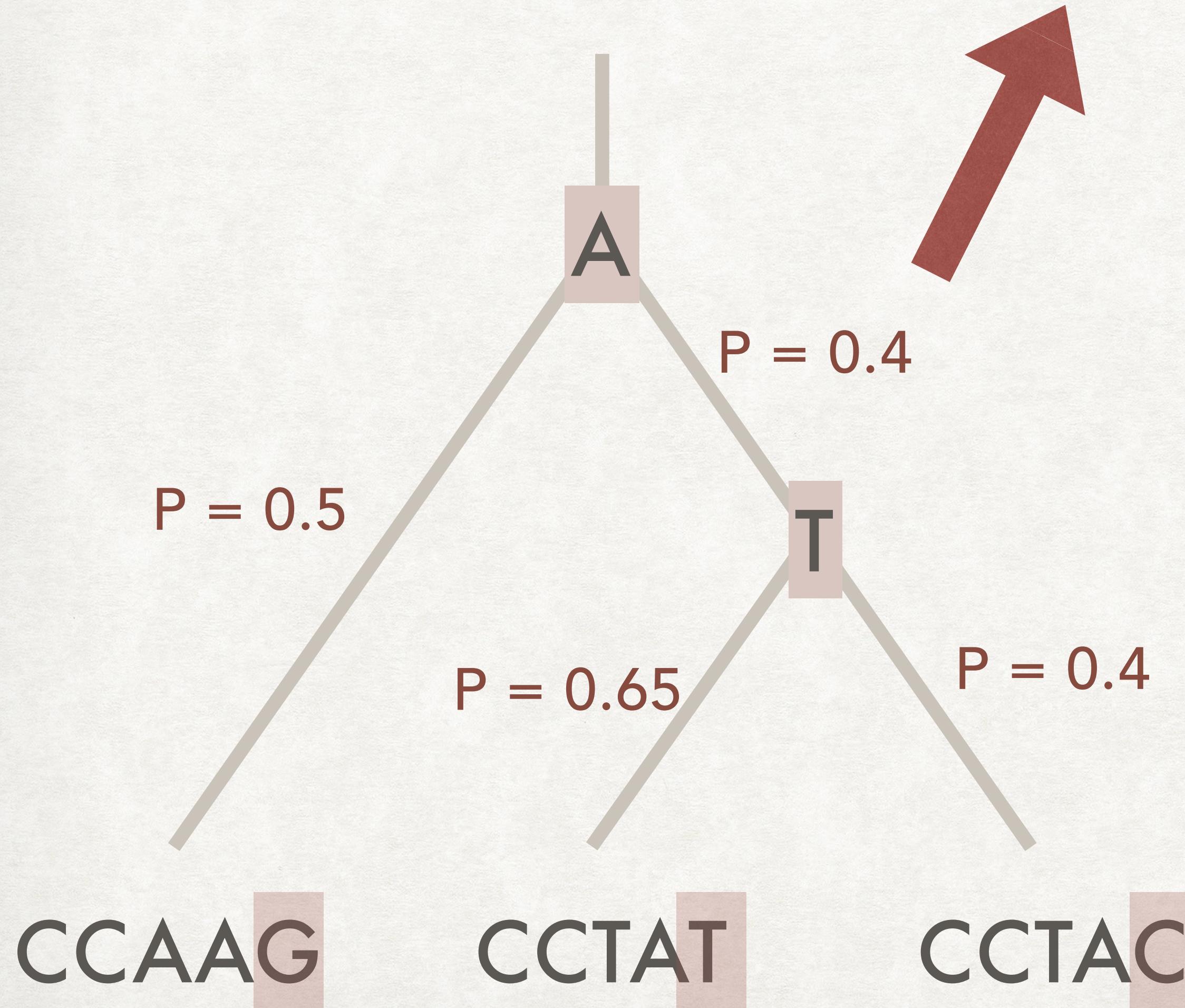


$\mu$  : Rate of change ( $\mu \times$  branch length = expected changes in one site)

$r(A \leftrightarrow C), r(A \leftrightarrow G), \dots, r(G \leftrightarrow T)$ :  
Proportion of change between different nucleotides

$\pi_A, \pi_C, \pi_G, \pi_T$ : "Equilibrium" frequencies of each nucleotide.

# TREE ESTIMATION

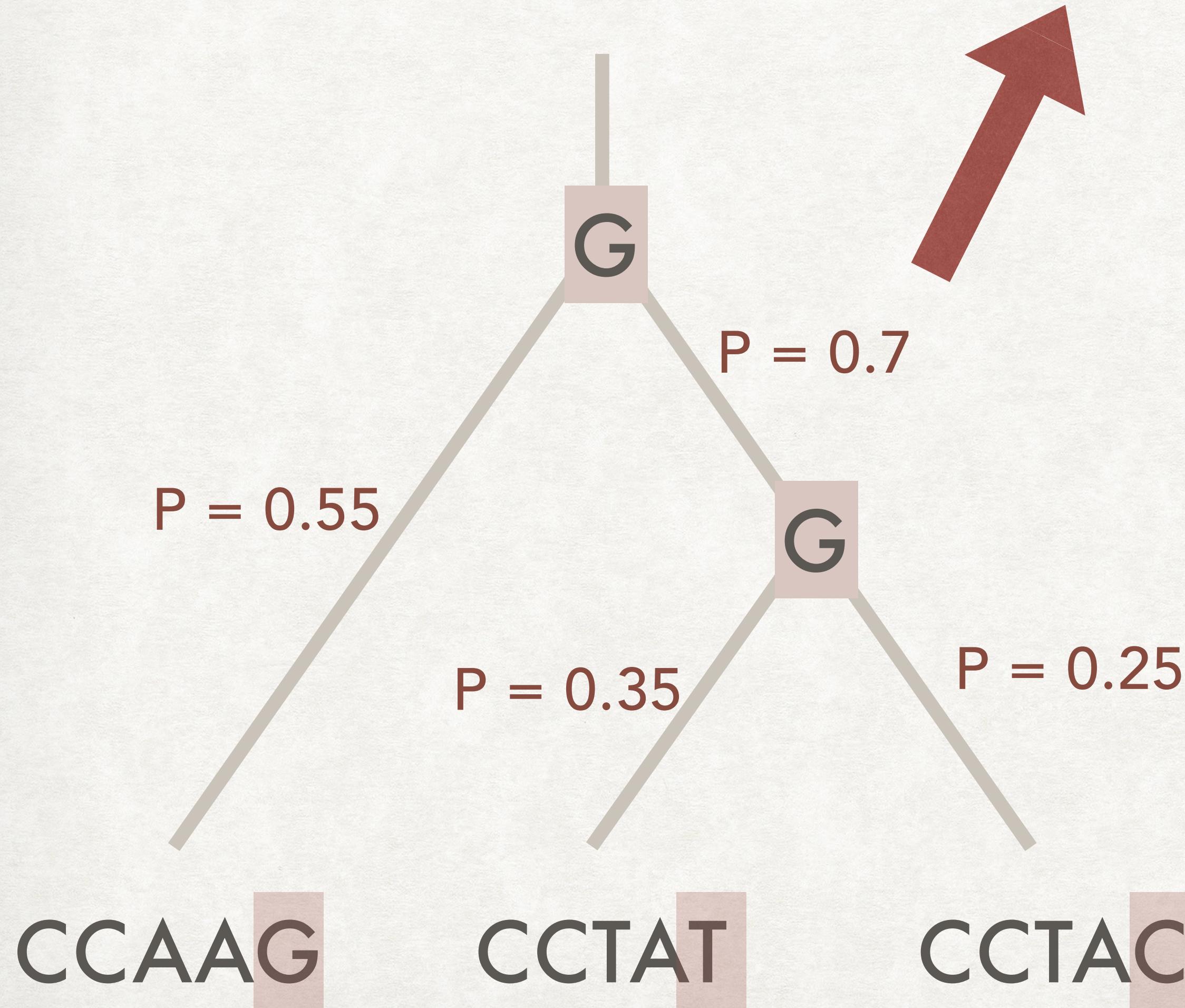


$\mu$  : Rate of change ( $\mu \times$  branch length = expected changes in one site)

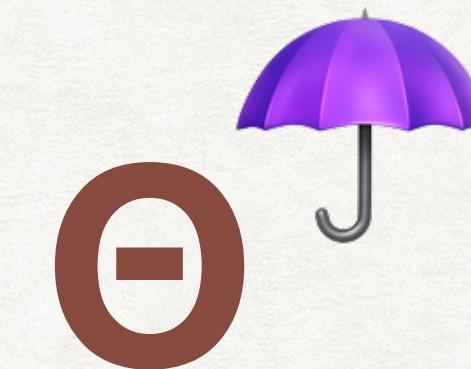
$r(A \leftrightarrow C), r(A \leftrightarrow G), \dots, r(G \leftrightarrow T)$ :  
Proportion of change between different nucleotides

$\pi_A, \pi_C, \pi_G, \pi_T$ : "Equilibrium" frequencies of each nucleotide.

# TREE ESTIMATION



$$\text{Total Prob} = 0.7 \times 0.25 \times 0.35 \times 0.55 = 0.0337\dots$$



$\mu$  : Rate of change ( $\mu \times$  branch length = expected changes in one site)

$r(A \leftrightarrow C), r(A \leftrightarrow G), \dots, r(G \leftrightarrow T)$ :  
Proportion of change between different nucleotides

$\pi_A, \pi_C, \pi_G, \pi_T$ : "Equilibrium" frequencies of each nucleotide.

# ASSUMPTIONS MADE BY THE MODEL

SO... SO... MANY  
BUT LET'S LIST SOME

- Each site (nucleotide) changes independently (questionable?)
- The proportion of transitions from A to T is the same as from T to A (fair as far as I understand)
- The rate of change/frequencies/proportions of transitions are the same across sites (questionable again) and branches of the tree (super questionable).
- When applied to aligned genes across the genome, **ALL SHARE THE SAME EVOLUTIONARY HISTORY**

# CONCLUSIONS

- Estimation of phylogenetic trees: truly complex problem to attack.
- The current methods used are strict parametric models, which we know oversimplify reality a lot.
- When talking about phylogenomics for microbes, maybe we should talk about phylogenetics instead.

P.D. There are methods not discussed here that take gene trees and return an estimate for the species trees (ASTRAL, BEAST,...). But it is still not entirely clear to me what we mean by “species tree” when talking on this scale of life.