

# **RGGS Comparative Genomics 2 – Computational Methods (Session 14)**

Jose Barba

Gerstner Scholar in Bioinformatics & Computational Biology

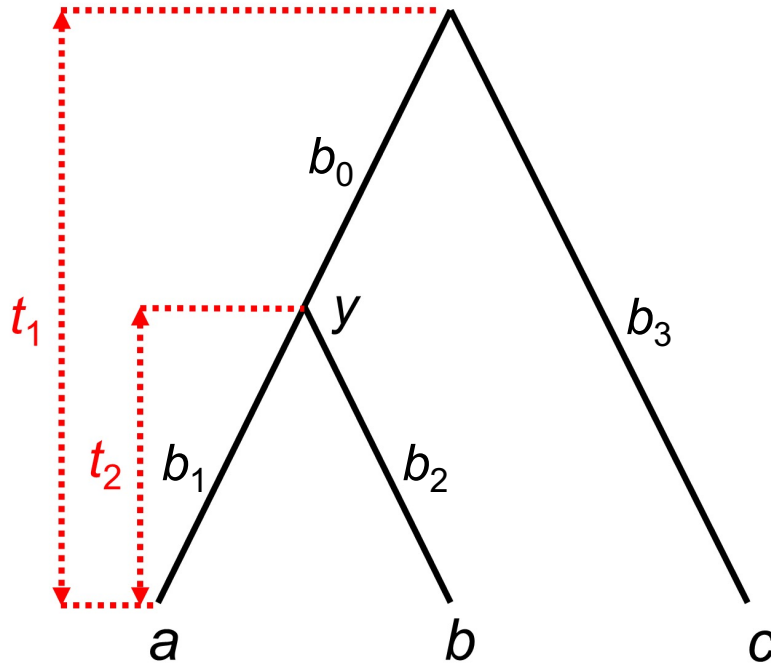
## **Session 14 outline**

- **ML and Bayesian phylogenetic inference (topics from Session 13 refer to that presentation)**
- **Molecular clock dating using phylogenomic data**
- **A tutorial for conducting a phylogenomic analysis using simulated data will be completed**
  - A comparison of phylogenetic methods will be conducted
  - A Bayesian timetree will be inferred

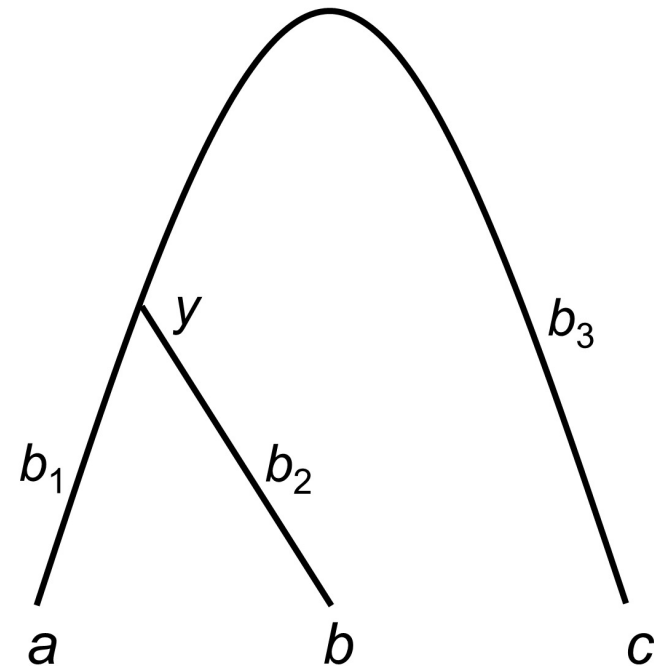
# Relative-rate test

- If substitution rate is constant over time or among lineages, molecular clock holds. For distantly related species this hypothesis should not be assumed

(a) Clock (rooted tree)

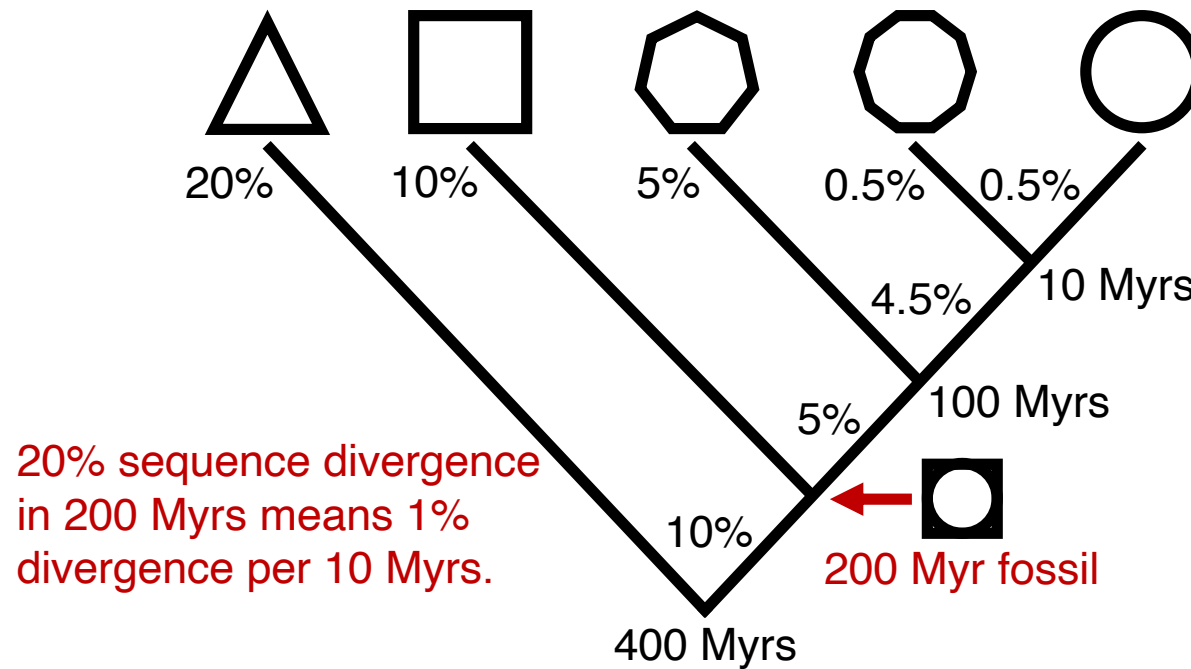


(b) No clock (unrooted tree)



# Molecular clock hypothesis

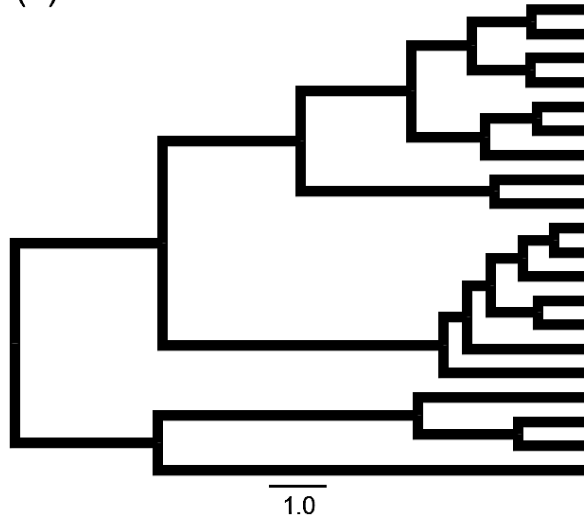
- Simple but powerful approach measuring timescale of evolutionary divergences. Expected distance between sequences grows linearly with time of divergence



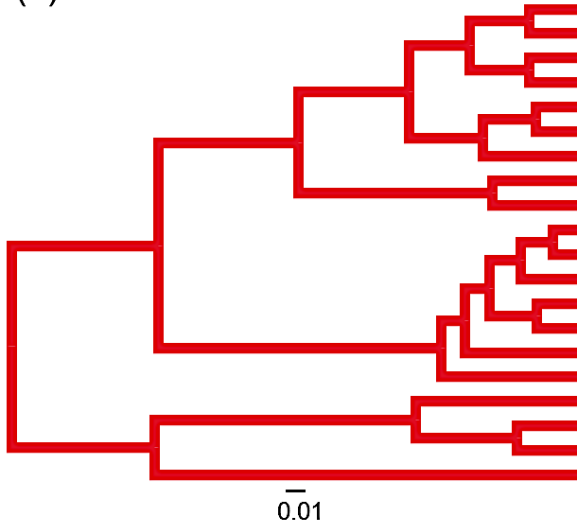
- Ages from fossil record or geological events, can be used to translate distances between sequences or tree branch lengths into absolute geological times

# Molecular clock model

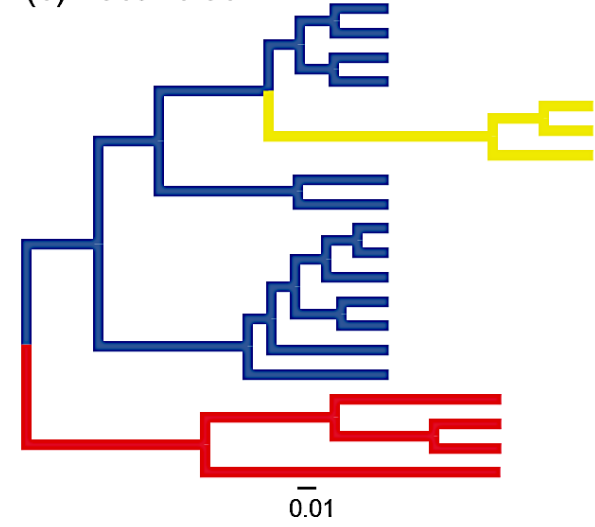
(a) Timetree



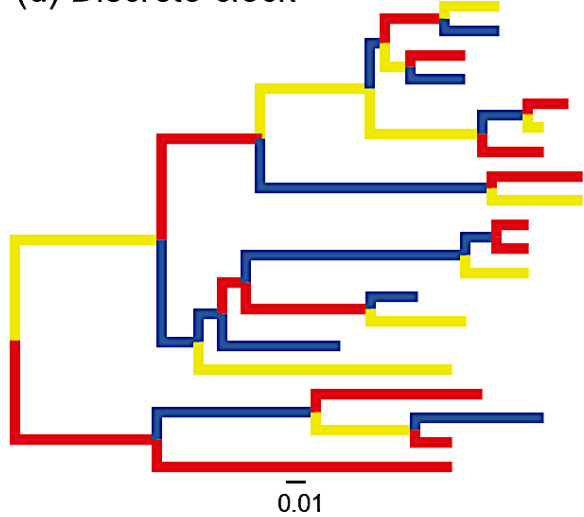
(b) Global clock



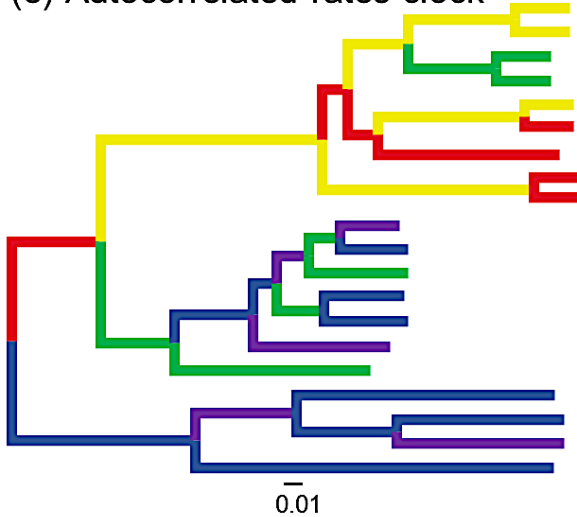
(c) Local clock



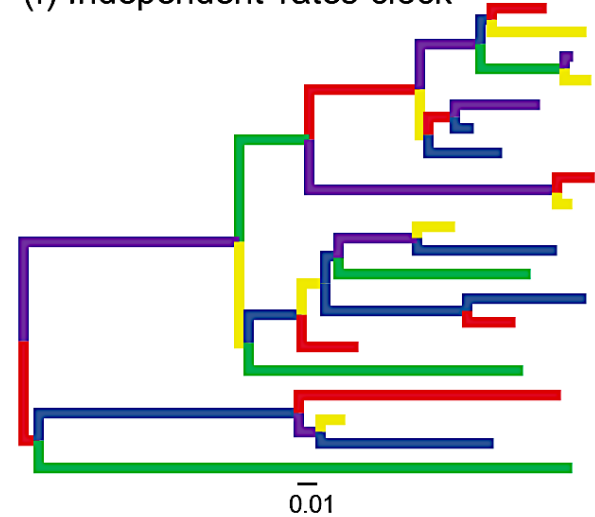
(d) Discrete clock



(e) Autocorrelated rates clock



(f) Independent rates clock



# Phylogenetic uncertainty in molecular dating

- **Types of molecular clock methods**
  - Generate reliable phylogeny before estimating divergence times
  - Jointly infer phylogeny and divergence times

## *Sequential Analysis*

- MCMCTree
- MultiDivTime
- PhyloBayes
- RelTime
- TreePL

## *Joint analysis*

- BEAST
- BEAST2
- MrBayes
- RevBayes
- TREETIME

# Bayesian divergence time estimation

- Using Bayesian method, can integrate fossil, molecule and clock uncertainty

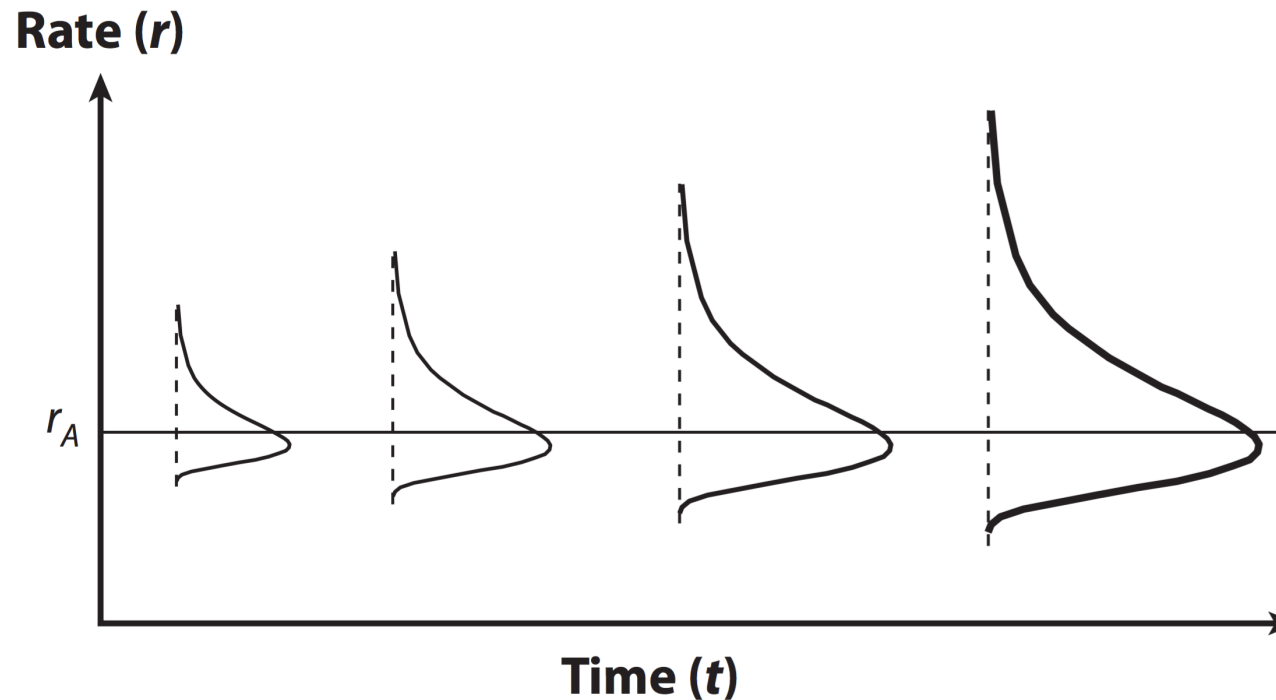
$$\underbrace{f(\mathbf{t}, \mathbf{r}, \theta | D)}_{\text{time and rate estimates}} = \underbrace{f(\theta)}_{\text{fossil uncertainty}} \underbrace{f(\mathbf{t})}_{\text{rate uncertainty}} \underbrace{f(\mathbf{r} | \mathbf{t}, \theta)}_{\text{branch length uncertainty}} \underbrace{L(D | \mathbf{t}, \mathbf{r}, \theta)}_{\text{branch length uncertainty}} / k$$

- Calculating the posterior involves a complex integral ( $k$ ), so MCMC algorithms are used to generate samples from joint posterior dist.

# Relaxed clocks and prior model of rate drift

- In AR model, rate at each node specified by conditioning on rate at ancestral node
  - Given the rate  $r_A$  at ancestral node, rate  $r$  at current node has a lognormal dist.

## Geometric Brownian motion model of rate drift





# Relaxed clocks and prior model of rate drift

- In IR model, rate for branch is a random variable drawn from a common probability distribution
- Rates effectively evolve independently on each lineage, but extent of rate variation has evolutionary constraint (imposed by prior distribution on rates)

$$f(r \mid \mu, \sigma^2) = \frac{1}{r\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (\log(r/\mu) + \frac{1}{2}\sigma^2)^2 \right\}, \quad 0 < r < \infty$$

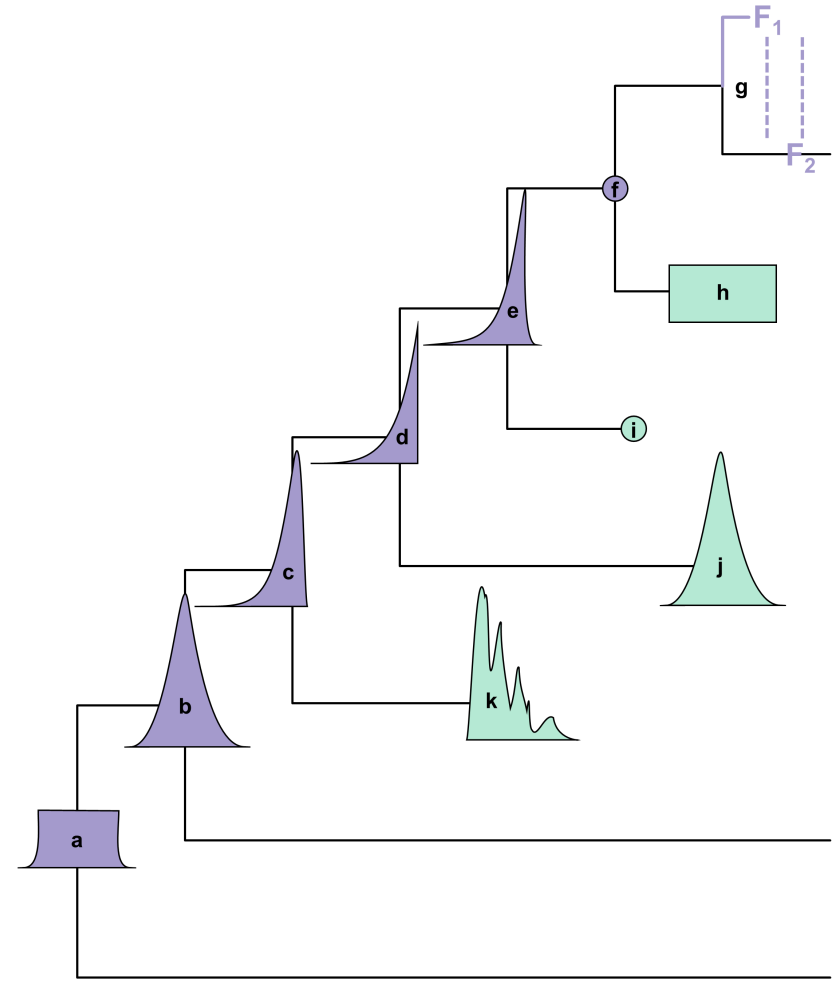
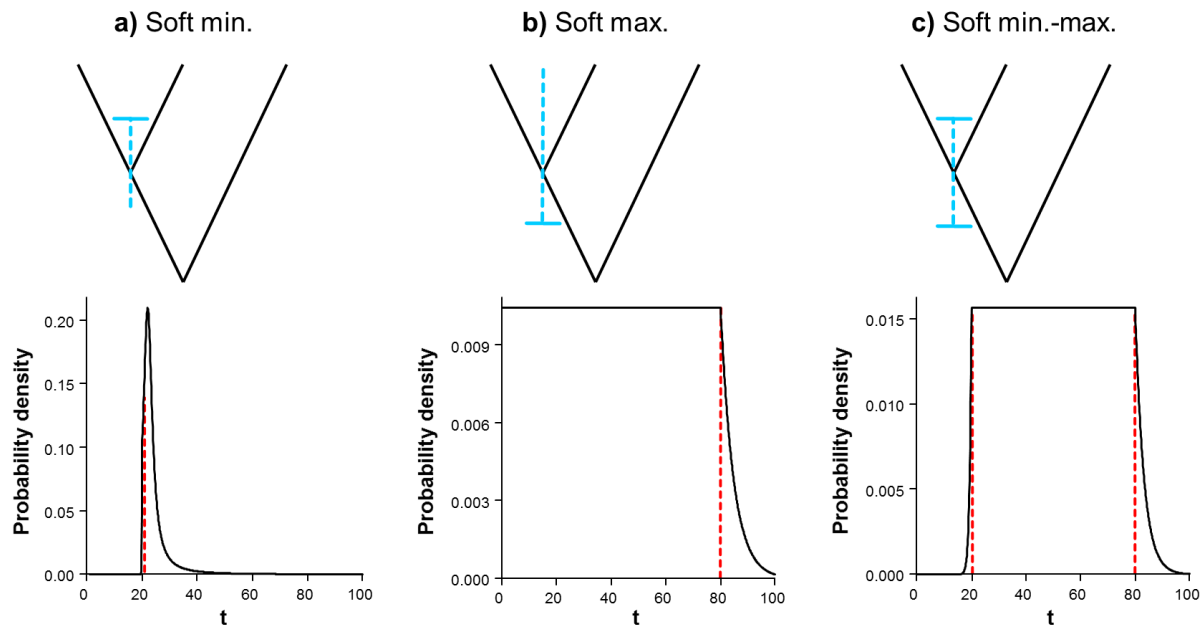
mean rate across loci

variance of the logarithm of the rate

# Fossil calibration

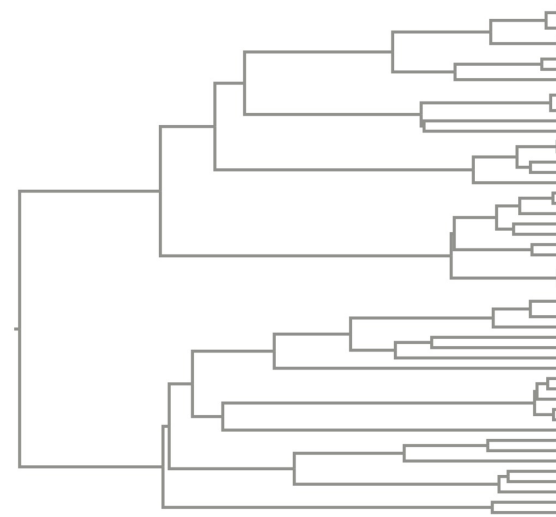
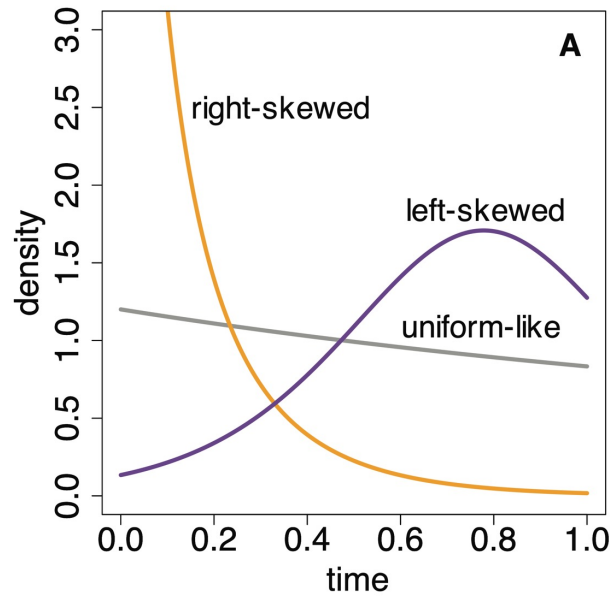
## Approaches to represent uncertainty of calibrations in phylogenetic tree

### Probability densities for describing uncertainty in fossil calibrations

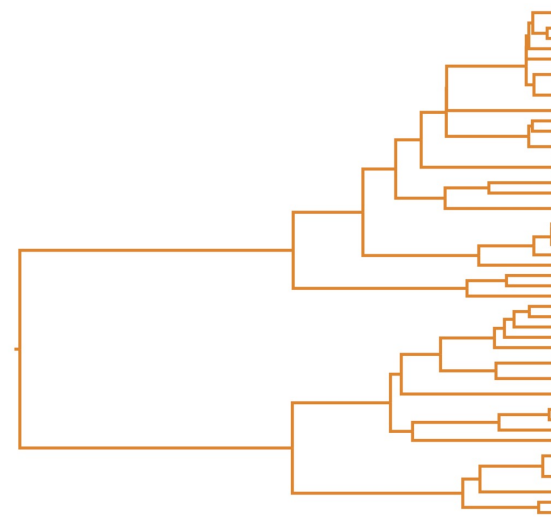


# Cladogenesis model

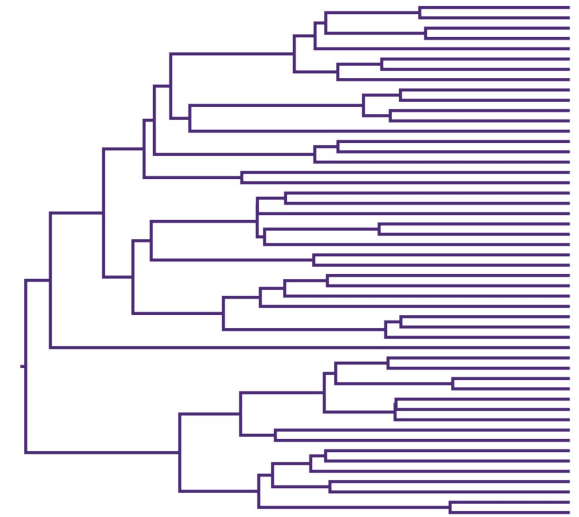
- **BD process with species sampling:**
  - Specified by a per-lineage birth rate  $\lambda$ , a per-lineage death rate  $\mu$  and a sampling fraction  $\rho$



(B) uniform-like



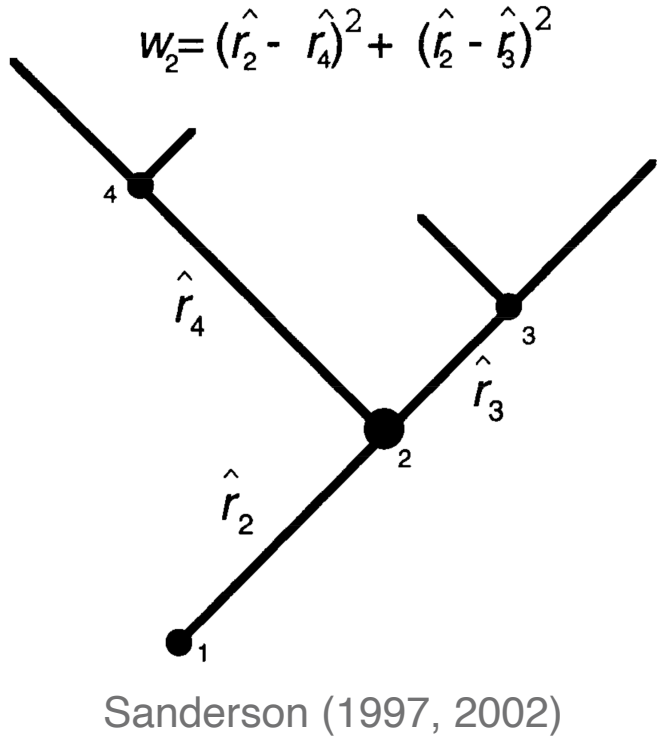
(C) right-skewed



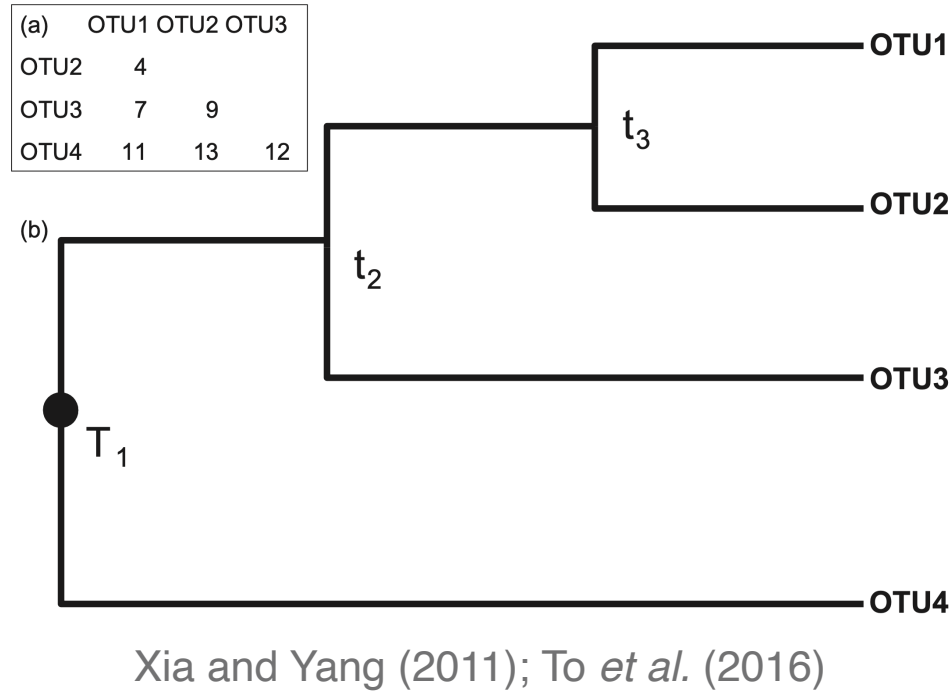
(D) left-skewed

## Rapid relaxed clock methods

## Penalized likelihood



# Least-squares



## Relative rate framework

