



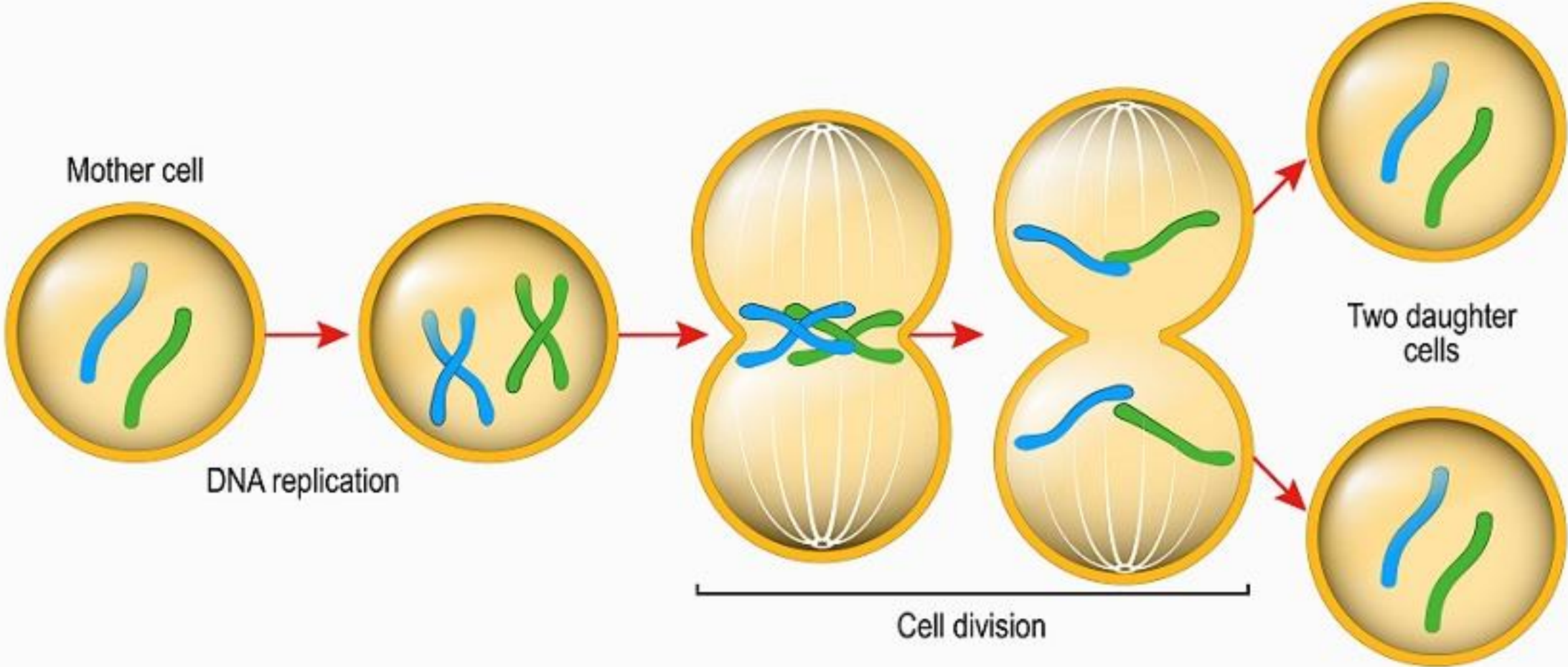
# Statistical Phylogenetics: An Introduction

2024-12-04, Miking Workshop

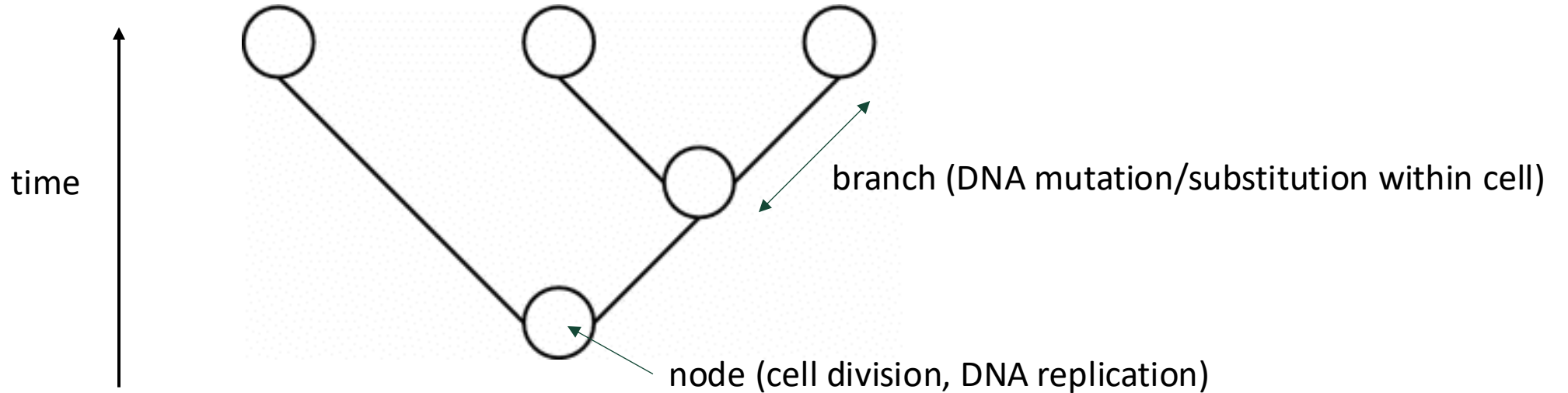
Prof. Fredrik Ronquist

Dept. Bioinformatics and Genetics

# Cell Division (Mitosis)



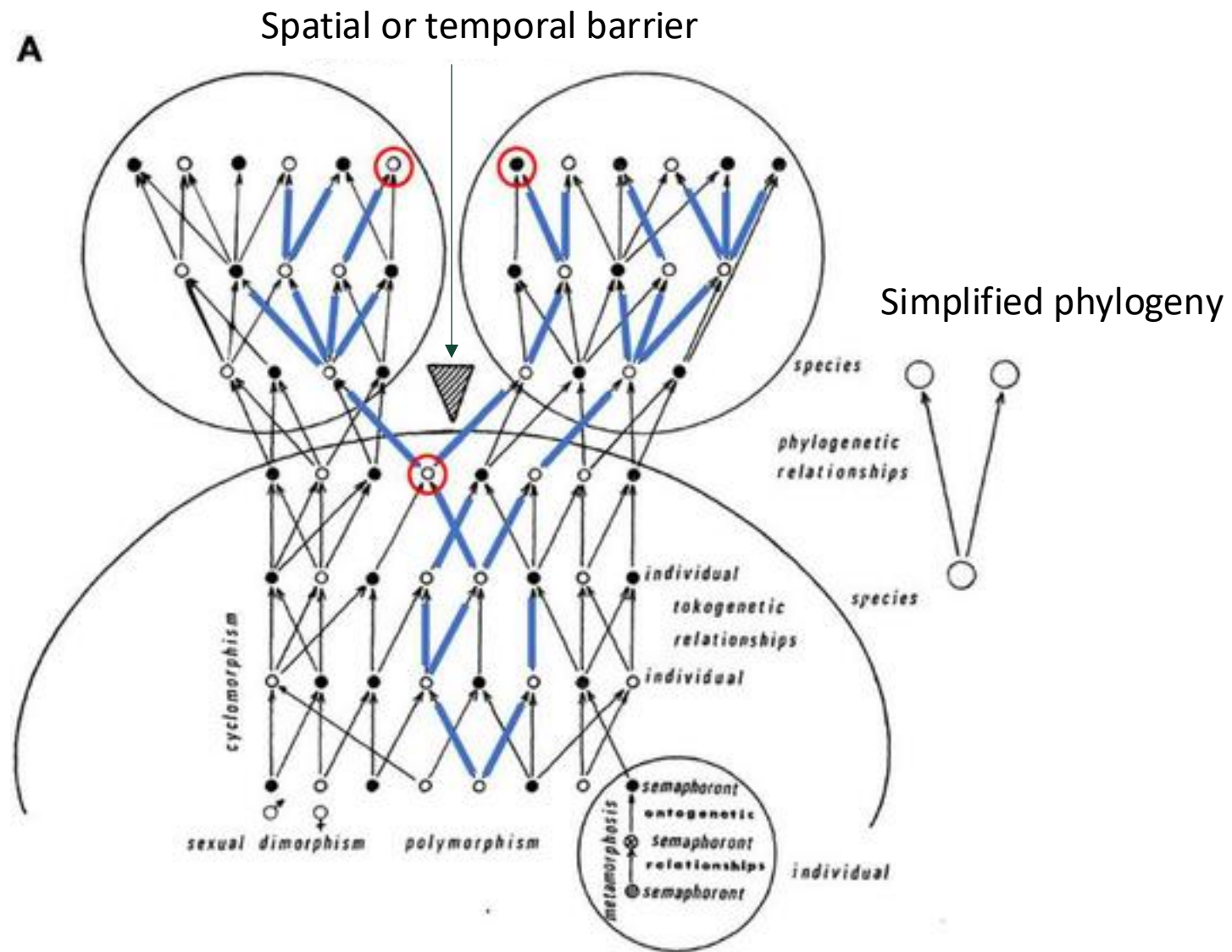
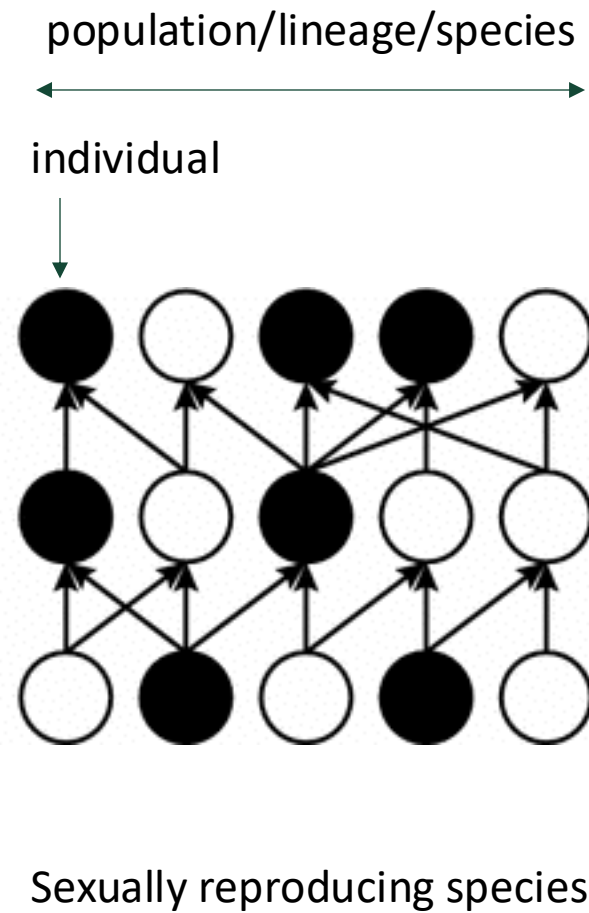
# Phylogenetic Model: The Basics



Phylogenetic tree (phylogeny)

Stochastic model:

- Model of branching process
- Model of substitution process



# Applications of Statistical Phylogenetics

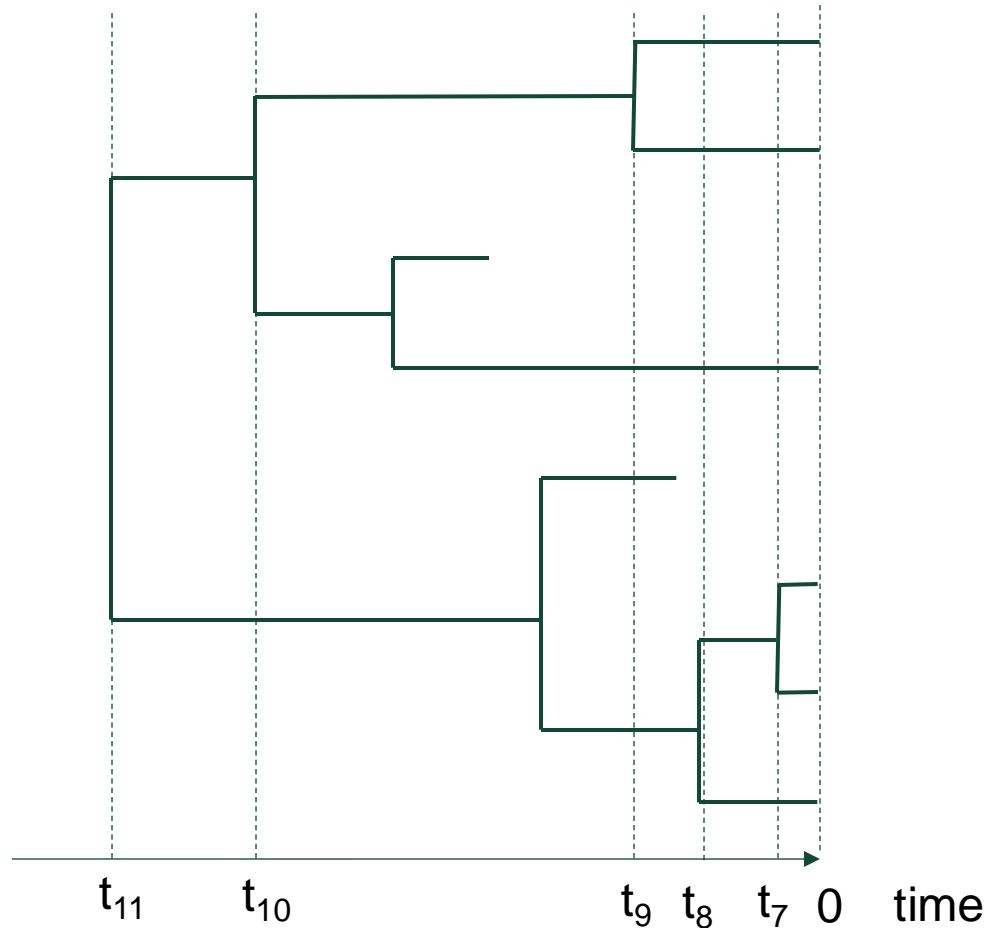
- Virus transmission pathways (HIV, SARS-CoV-2,...)
- Epidemiology
- Predicting next year's influenza outbreak
- Identification of pathogens
- Relationships among organisms
- Divergence time estimation
- Molecular evolution processes
- Selection analysis (which genes cause X)
- Patterns of diversification and extinction
- Biogeography (where does an organism come from)
- ...



**“Nothing in Biology Makes Sense Except  
in the Light of Evolution”**

Theodosius Dobzhansky (1900-1975)

# Tree model



Birth-death model:

Birth rate  $\lambda$

Death rate  $\mu$

Root time  $t_{\text{mrca}} = t_{11}$

The birth-death model induces a probability distribution on

Topology  $\mathcal{T}$

Speciation times  $t$

Given a substitution rate  $r$ , branch lengths  $b$  are given by

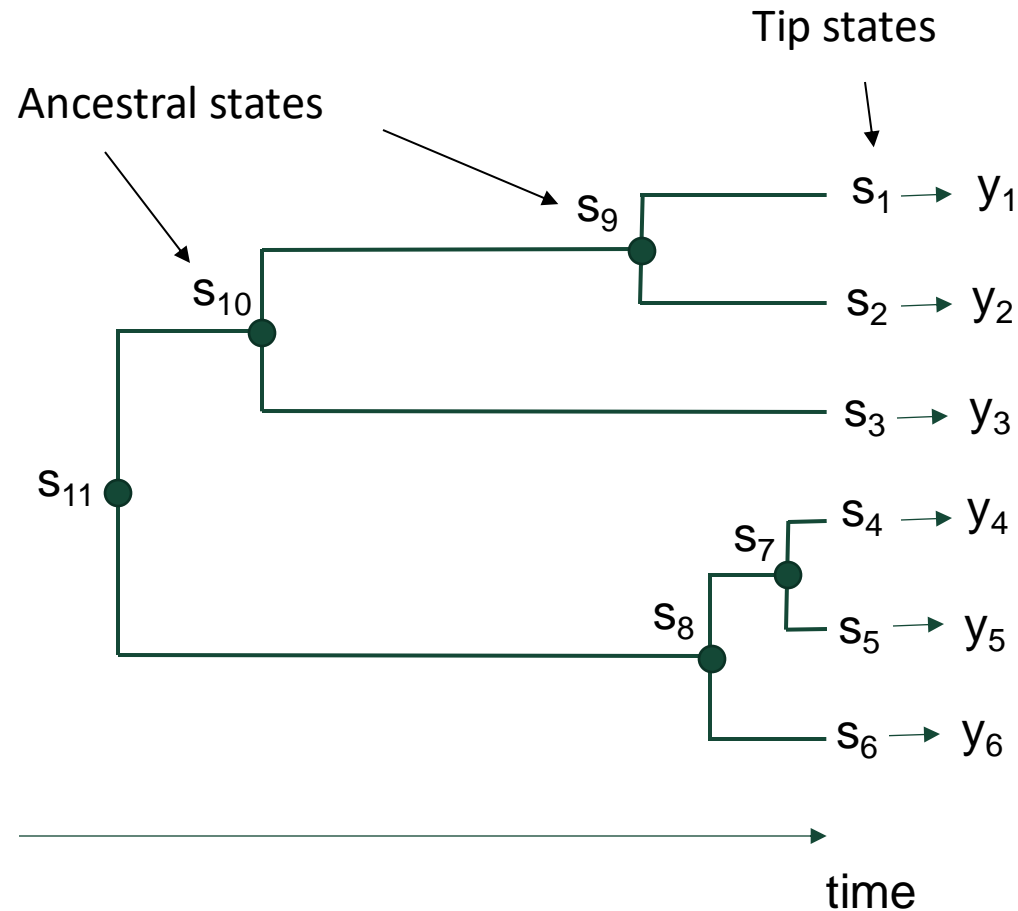
$$b_i = r(t_{a_i} - t_i)$$

# Substitution model

DNA sequence

GATAAATAATATAAGATTTTGAC...

site (s) assumed iid



Observation error usually ignored, that is, it is assumed that  $y_i = s_i$  for all leaves  $i$ .



DNA sequences are drawn iid from a discrete-state continuous-time Markov chain over four nucleotides: A, C, G, T

$$Q = \begin{matrix} & \begin{matrix} A & C & G & T \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \begin{matrix} -\rho_C r_{AC} & \rho_G r_{AG} & \rho_T r_{AT} & 0 \\ \rho_A r_{AC} & -\rho_G r_{CG} & \rho_T r_{CT} & 0 \\ \rho_A r_{AG} & \rho_C r_{CG} & -\rho_T r_{GT} & 0 \\ \rho_A r_{AT} & \rho_C r_{CT} & \rho_G r_{GT} & -\rho \end{matrix} \end{matrix}$$

Instantaneous rate matrix for the General Time Reversible (GTR) substitution model

$\rho$  Stationary state frequencies

$r$  Exchangeability rates

## Standard approach in statistical phylogenetics

- Specify model: script or command-line settings
- Hard-coded Bayesian MCMC (or maximum likelihood) inference machinery



# MrBayes: Bayesian Inference of Phylogeny

## Download MrBayes

MrBayes may be downloaded as a pre-compiled executable or in source form (recommended).

### Current release

The most recent release version of MrBayes is [3.2.7a](#), released March 6, 2019.

The 3.2.7a [source code](#) is available for compilation on Unix machines.

Pre-compiled (provisional) executables are available for Windows ([MrBayes-3.2.7-WIN.zip](#)). These are, however, serial versions compiled without the [Beagle](#) library. The serial version works well for smaller analyses but if you plan to run large analyses using many parallel chains, you should use the MPI version instead. Refer to the [User Manual](#), and the [INSTALL](#) document on GitHub for help with installation of the program.

MrBayes may also be installed through the Homebrew package manager on macOS, Linux, and Windows Subsystem for Linux (WSL). Please see the [INSTALL](#) document for instructions.

### Older releases

You can get access to older releases (from release 3.2.0 onwards), by browsing the [releases directory on github](#).

### Developer version

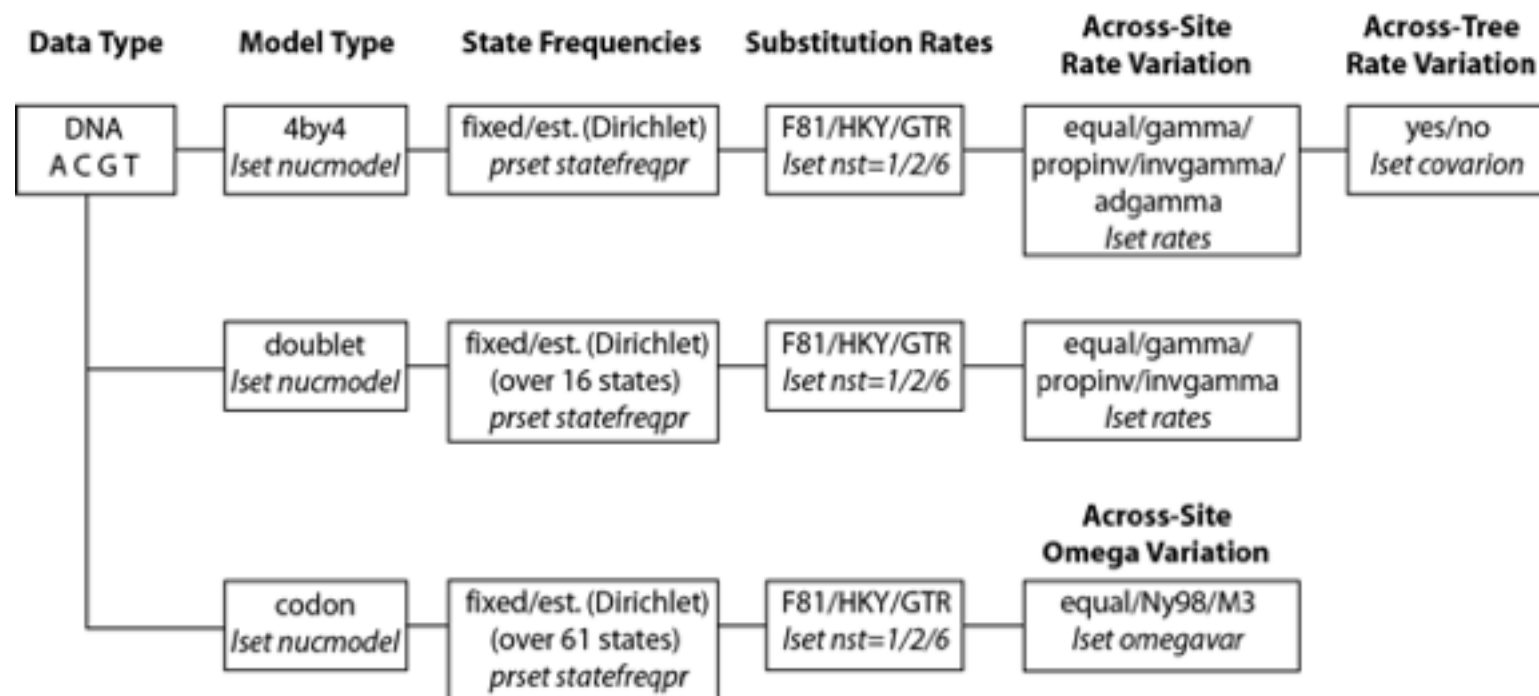
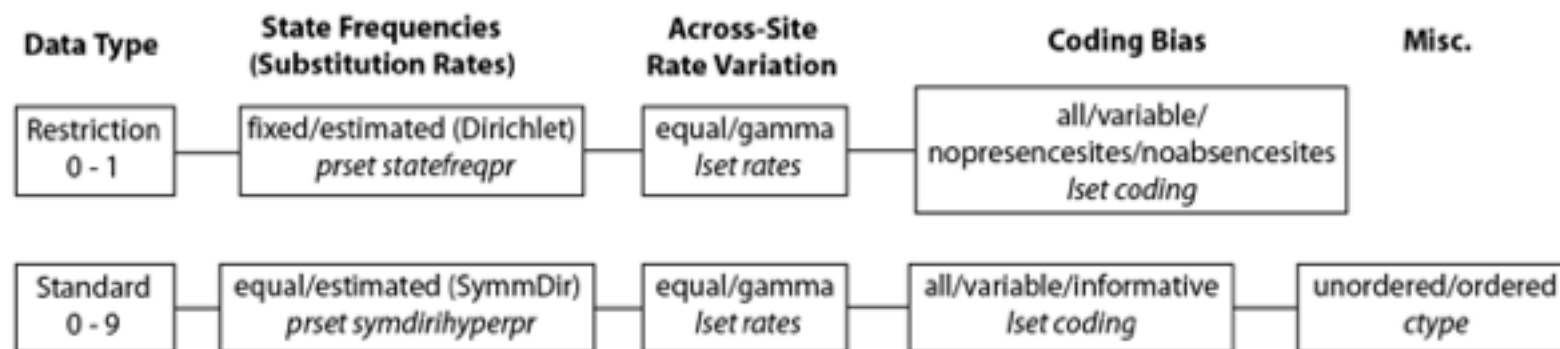
You may also download cutting edge [developer version](#) of MrBayes from the [Git repository](#) hosted at GitHub. Note that you have to compile the code yourself. Read instructions in the [INSTALL](#) file in the source code repository for further instructions.

[Home](#)  
» [Download](#) «  
[Manual](#)  
[Bug Report](#)  
[Authors](#)  
[Links](#)



<https://nbisweden.github.io/MrBayes/download.html>

Models supported by MrBayes 3 (simplified)



# MrBayes script

```
#NEXUS
```

```
begin mrbayes;
```

```
    execute data.nex;
```

```
    outgroup Ibalia;
```

```
    charset morphology = 1-166;
```

```
    charset molecules = 167-3246;
```

```
    charset COI = 167-1244;
```

```
    charset EF1a = 1245-1611;
```

```
    charset LWRh = 1612-2092;
```

```
    charset 28S = 2093-3246;
```

```
    partition favored= 5: morphology, COI, EF1a, LWRh, 28S;
```

```
    set partition=favored;
```

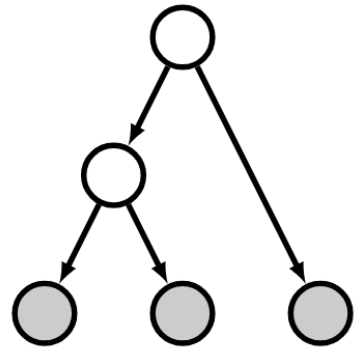
```
    lset applyto=(1) rates=gamma;
```

```
    lset applyto=(2,3,4,5) rates=invgamma nst=mixed;
```

```
    unlink revmat=(all) pinvar=(all) shape=(all) statefreq=(all);
```

```
    prset ratepr=variable;
```

```
end;
```



# RevBayes

Bayesian phylogenetic inference using probabilistic graphical models and an interpreted language

## About

RevBayes provides an interactive environment for statistical computation in phylogenetics. It is primarily intended for modeling, simulation, and Bayesian inference in evolutionary biology, particularly phylogenetics. However, the environment is quite general and can be useful for many complex modeling tasks.

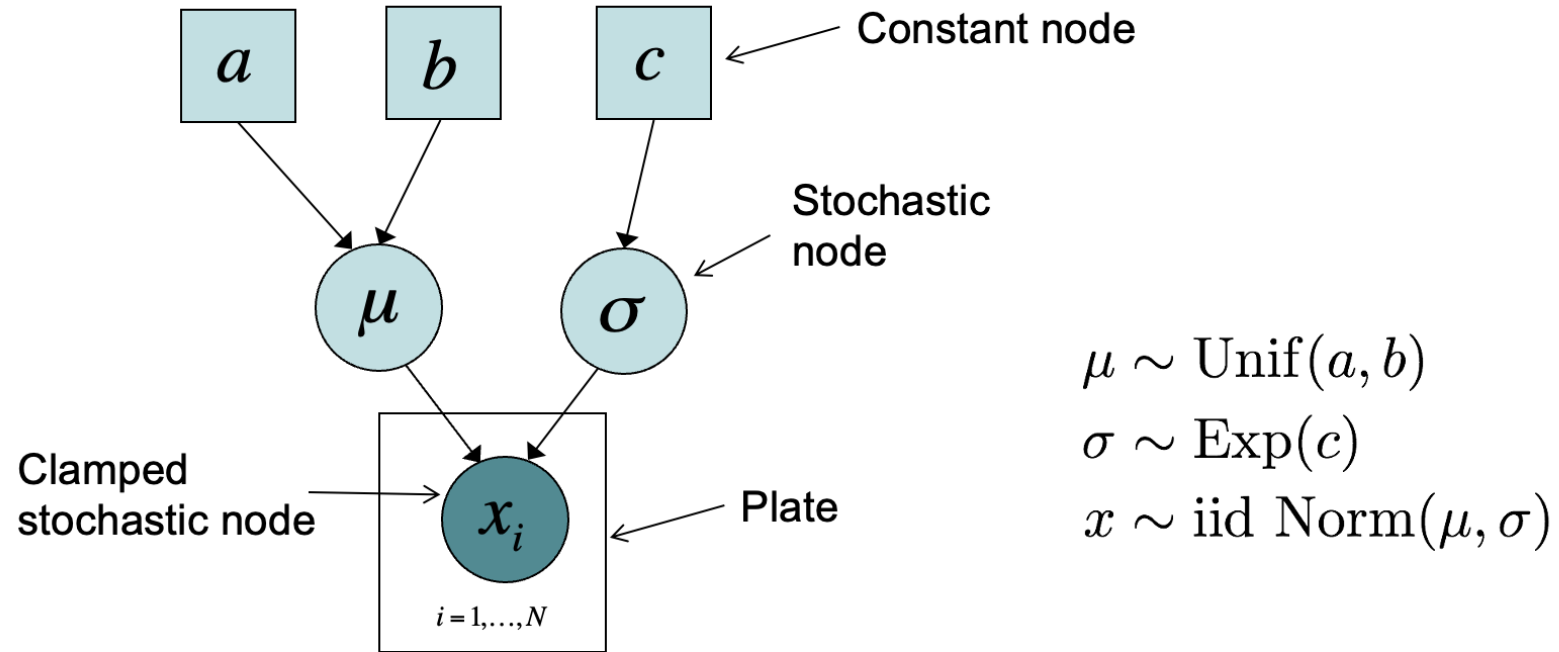
RevBayes uses its own language, Rev, which is a probabilistic programming language like [JAGS](#), [STAN](#), [Edward](#), [PyMC3](#), and related software. However, phylogenetic models require inference machinery and distributions that are unavailable in these other tools.

The Rev language is similar to the language used in R. Like the R language, Rev is designed to support interactive analysis. It supports both functional and procedural programming models, and makes a clear distinction between the two. Rev is also more strongly typed than R.

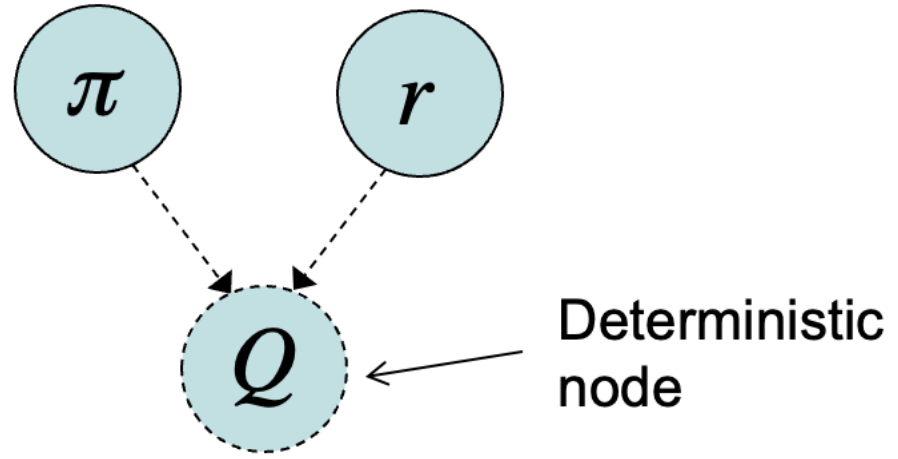
RevBayes is a collaboratively [developed](#) software project.

[GitHub](#) | [License](#) | [Citation](#) | [Users Forum](#)

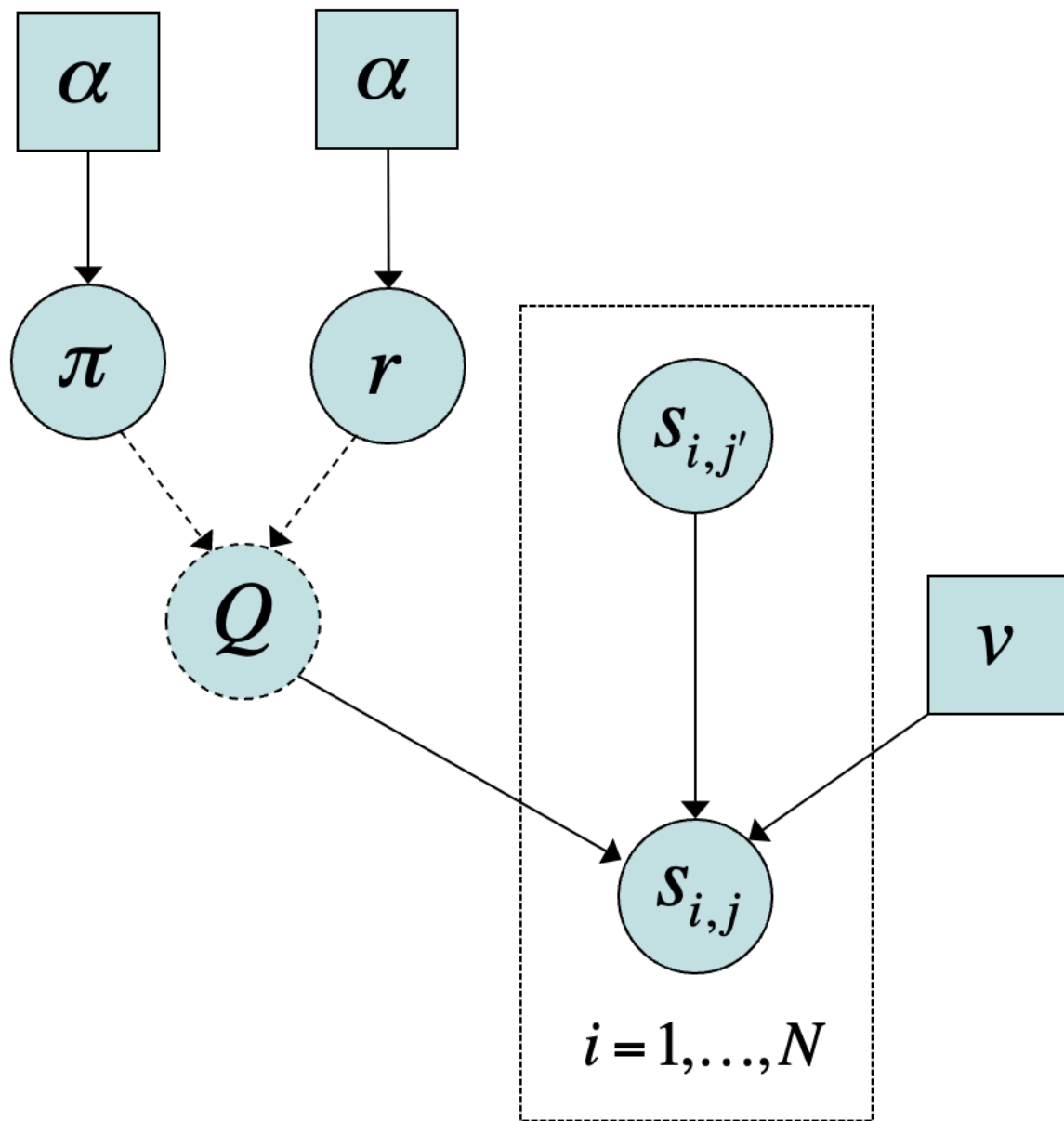
<http://www.revbayes.com>

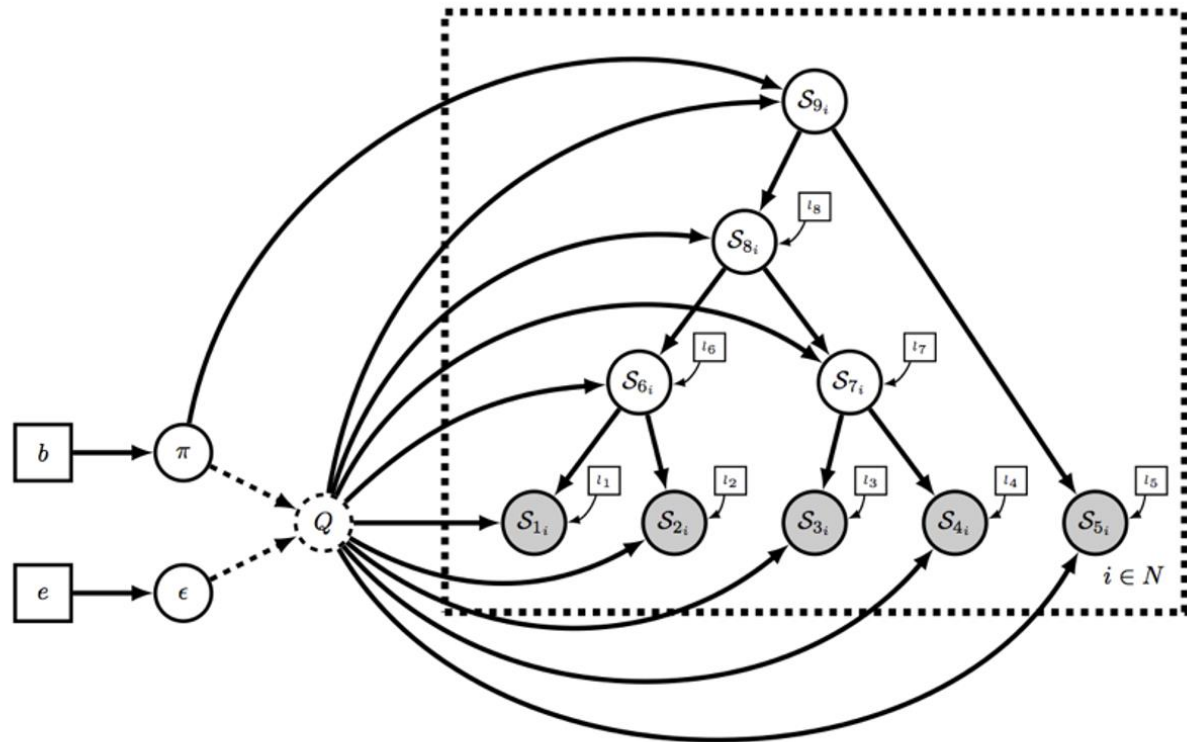


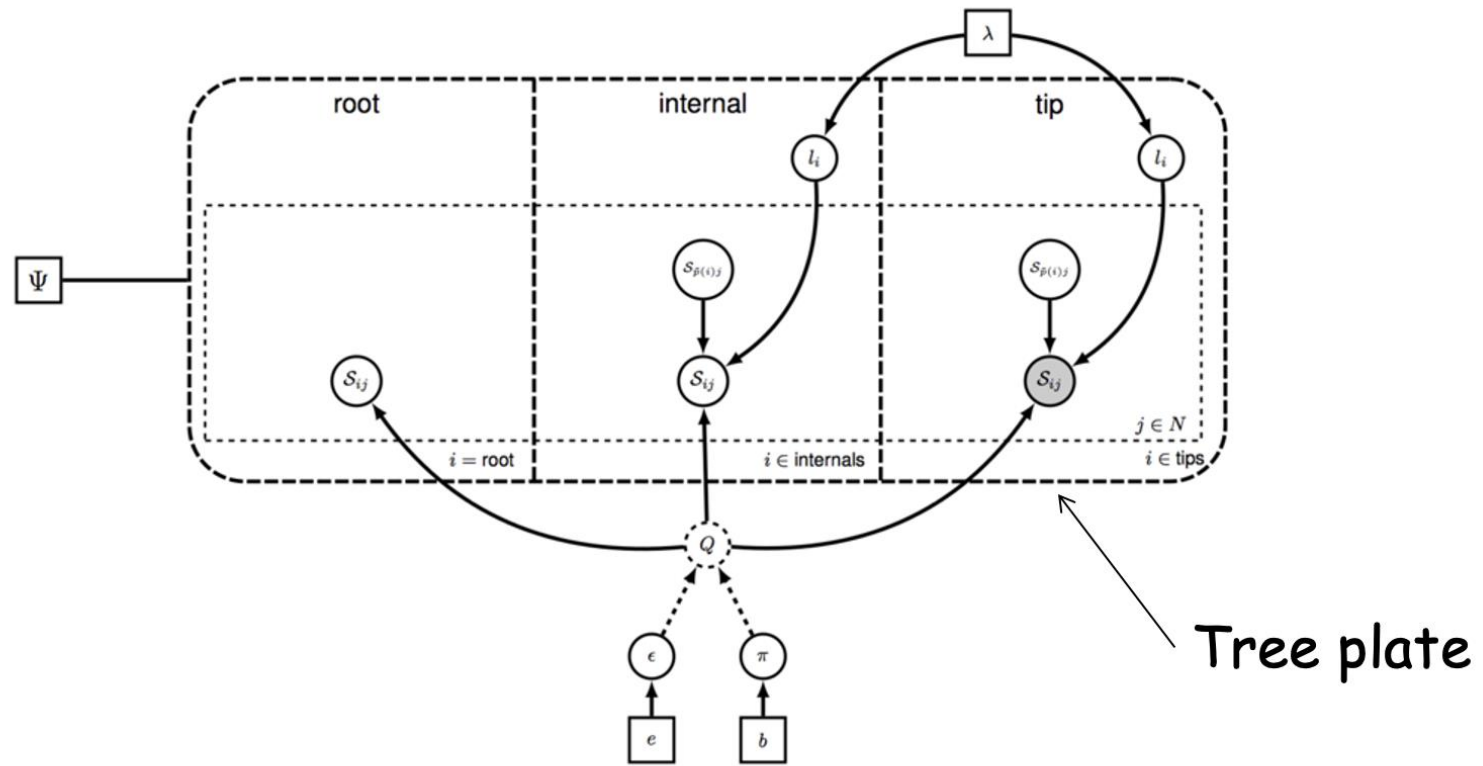
Hierarchical Normal Model

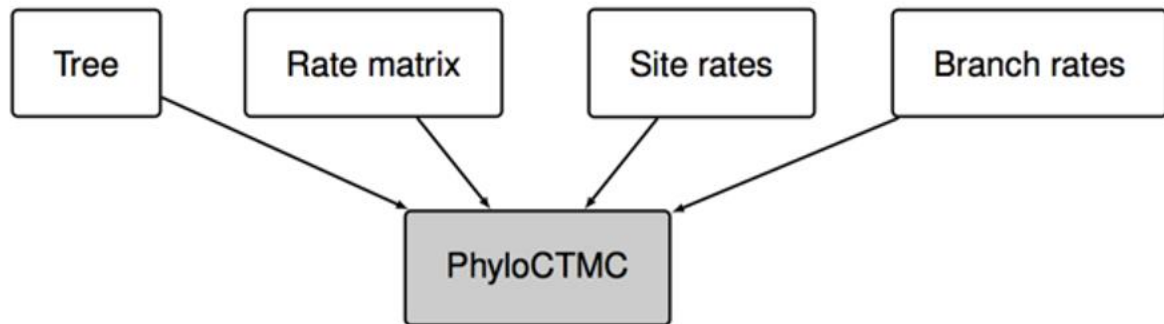




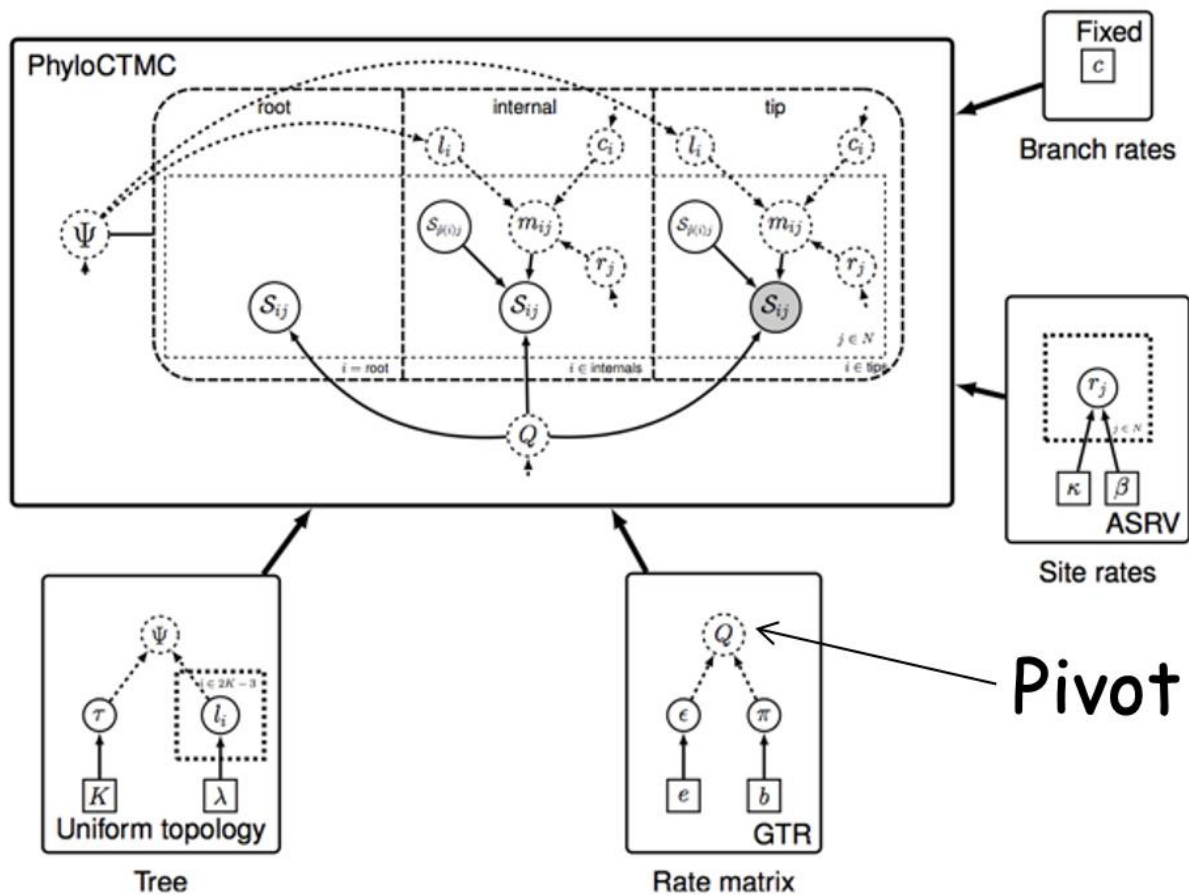








# Modular representation



# Challenges for a universal PPL for phylogenetics

- Ease of use in describing models and inference algorithms
- Computational efficiency [problem sizes are large, convergence is difficult]
  - A range of inference strategies
  - Fast computation of likelihoods/weights
  - Informative diagnostics (are we there yet?)
- Belief propagation
  - Makes a huge difference when it can be applied
  - Standard model (GTR+Gamma) requires belief propagation in two dimensions (rates across sites and ancestral nucleotide states)
  - Automated?
- Alignment
  - Required when number of random variables varies (birth-death models, sampling of change histories)
  - Automated?
- Proposal/guide distributions
  - Using quick-and-dirty approximations essential for efficient inference
  - Generic propose – update mechanism?