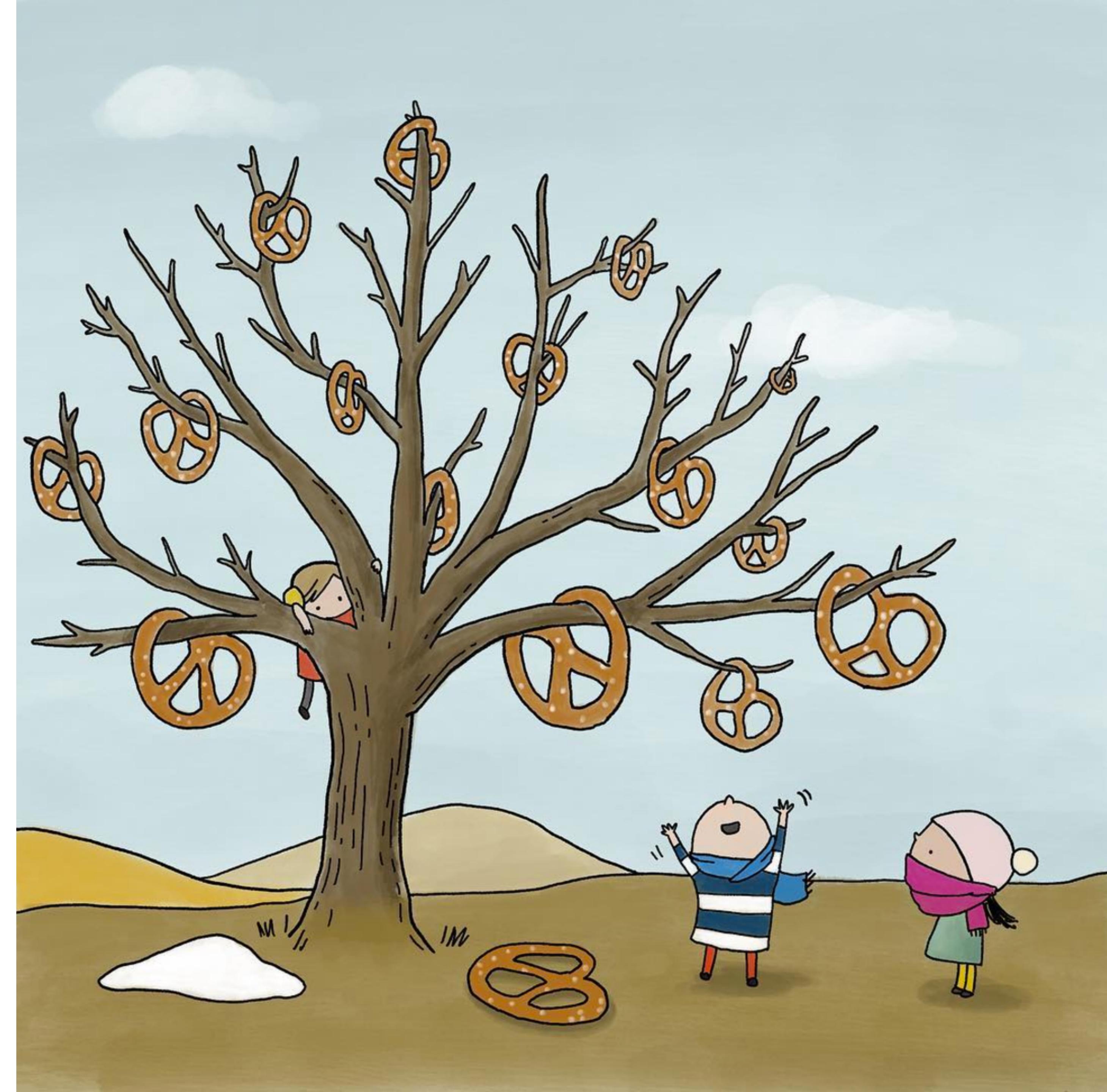


Phylogenetics workshop

Introduction to Bayesian
tree inference

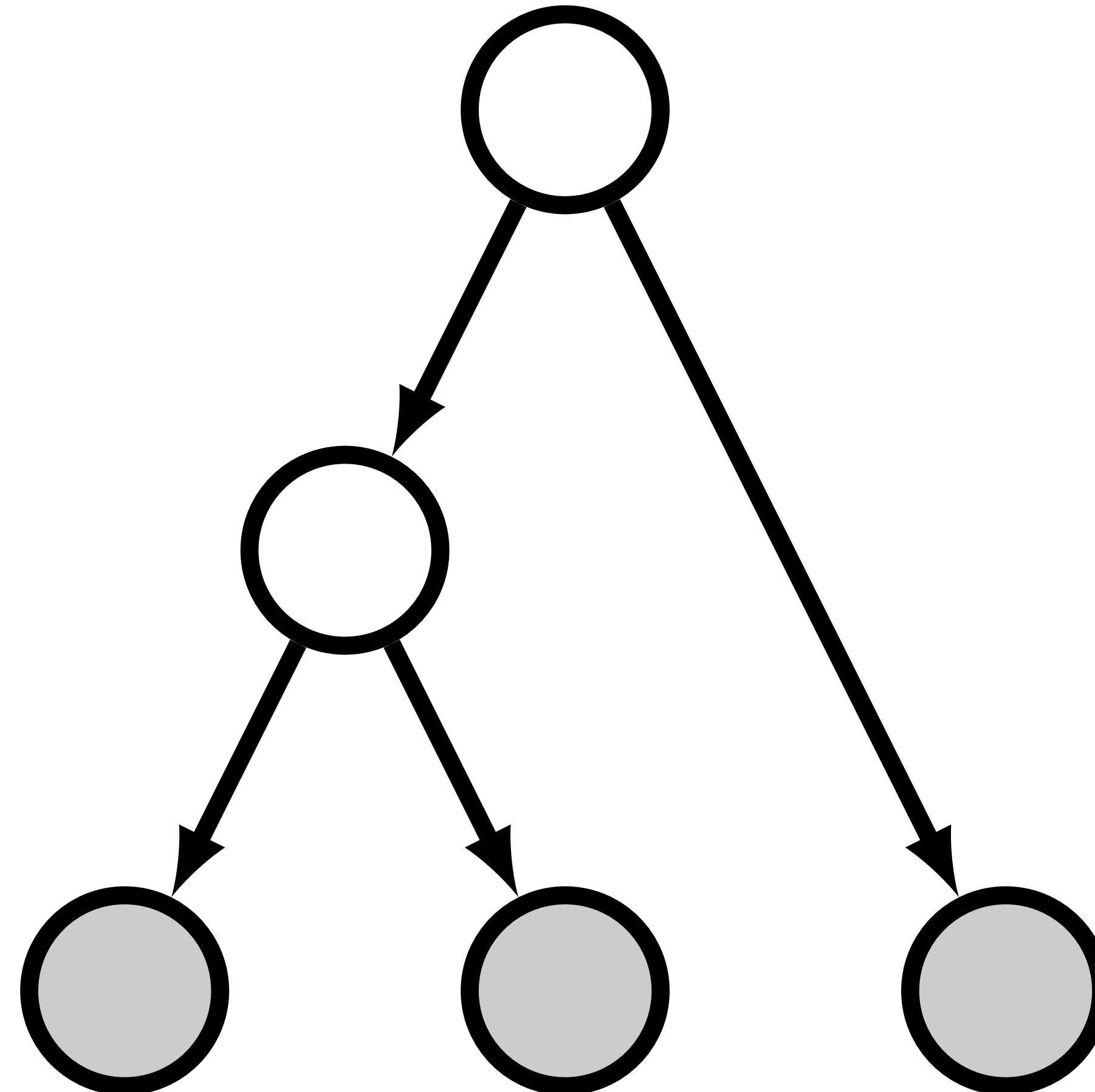
Rachel Warnock & co

10.12.24



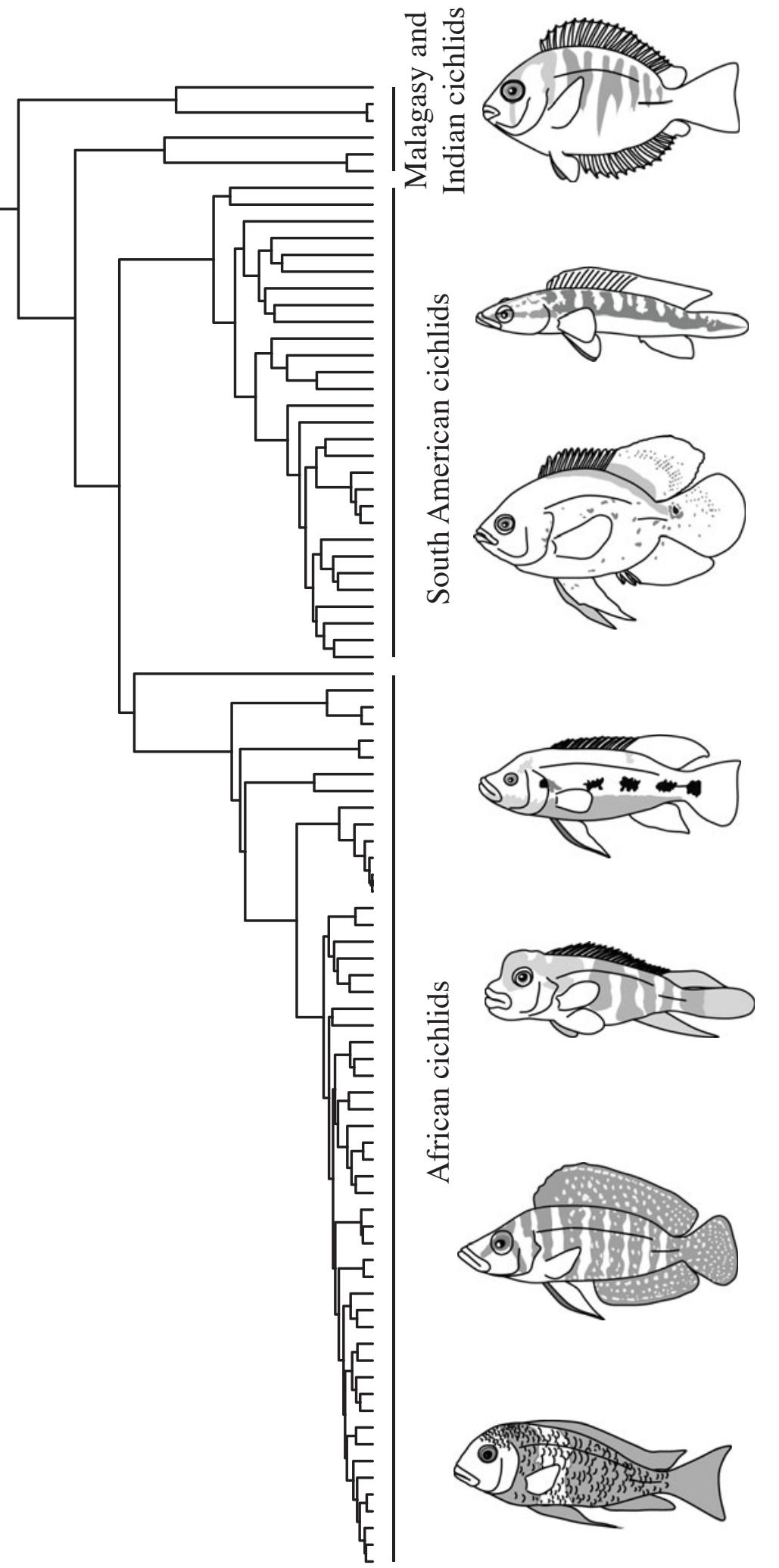
Part 1 objectives

- Bayesian inference
- MCMC
- RevBayes



What can we learn from trees?

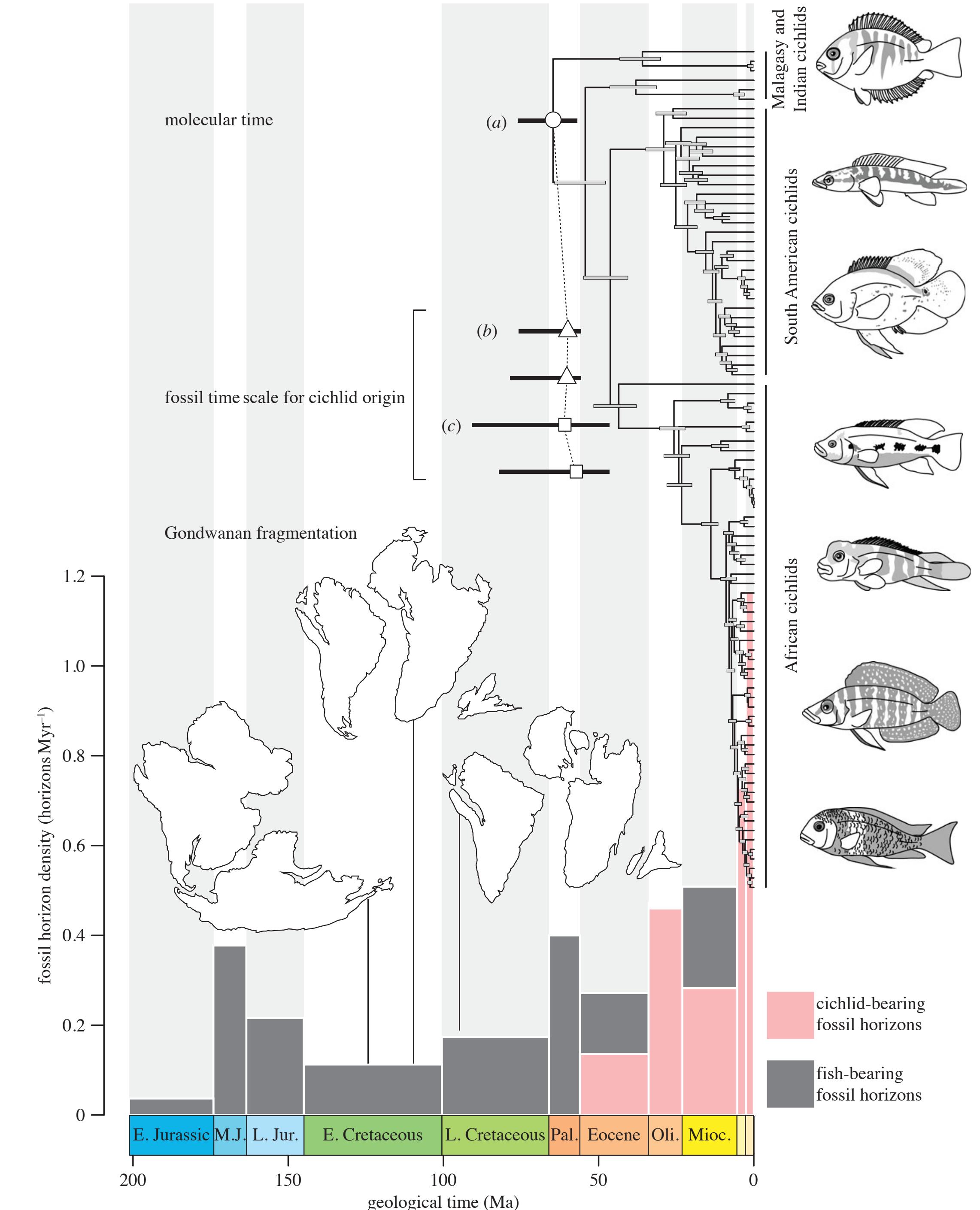
- Evolutionary relationships



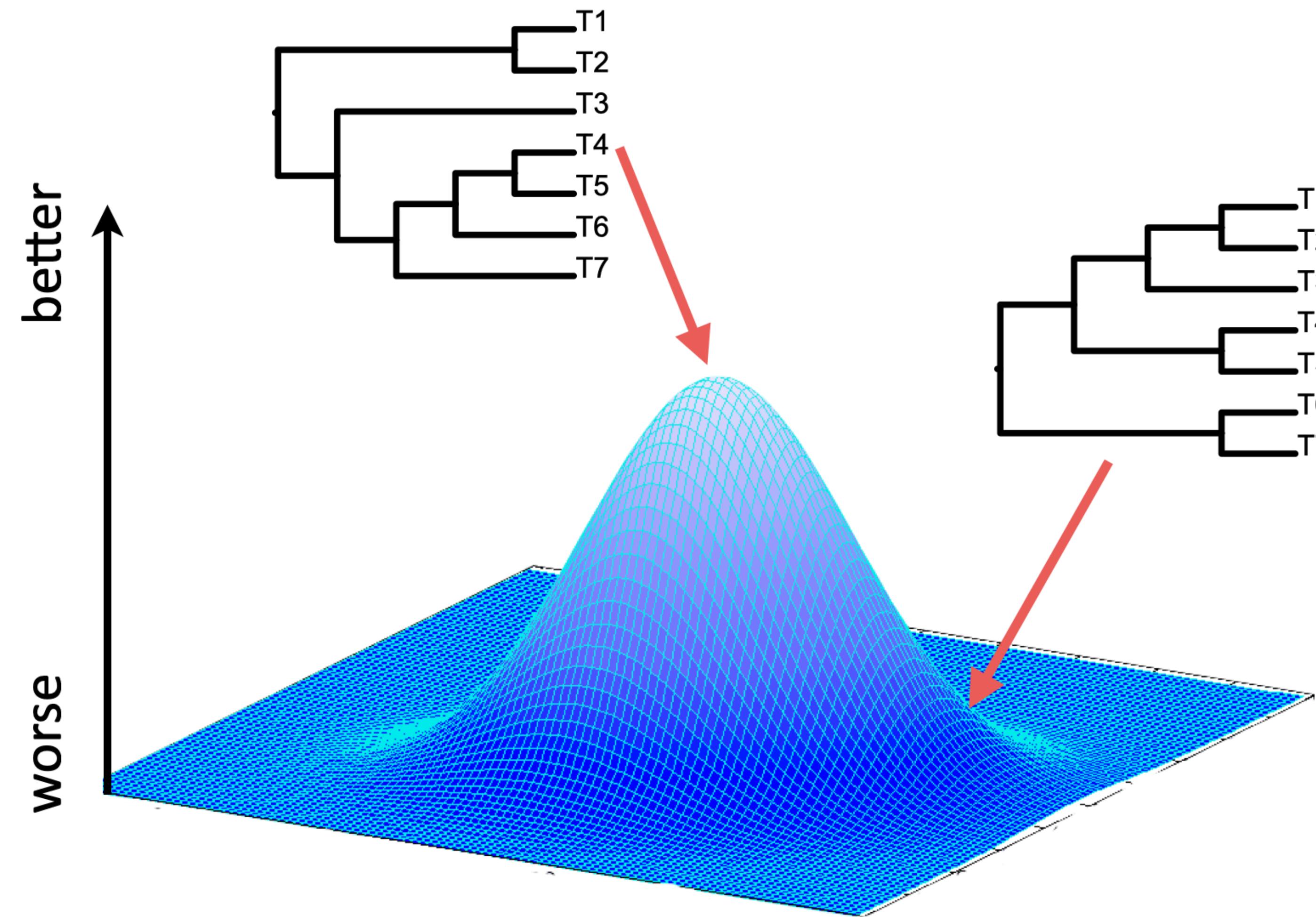
What can we learn from trees?

- Evolutionary relationships
- Timing of diversification events
- Geological context
- Rates of phenotypic evolution
- Diversification rates

Image adapted from Friedmann et al. (2013)



How do we find the ‘best’ tree?



It depends how you measure ‘best’

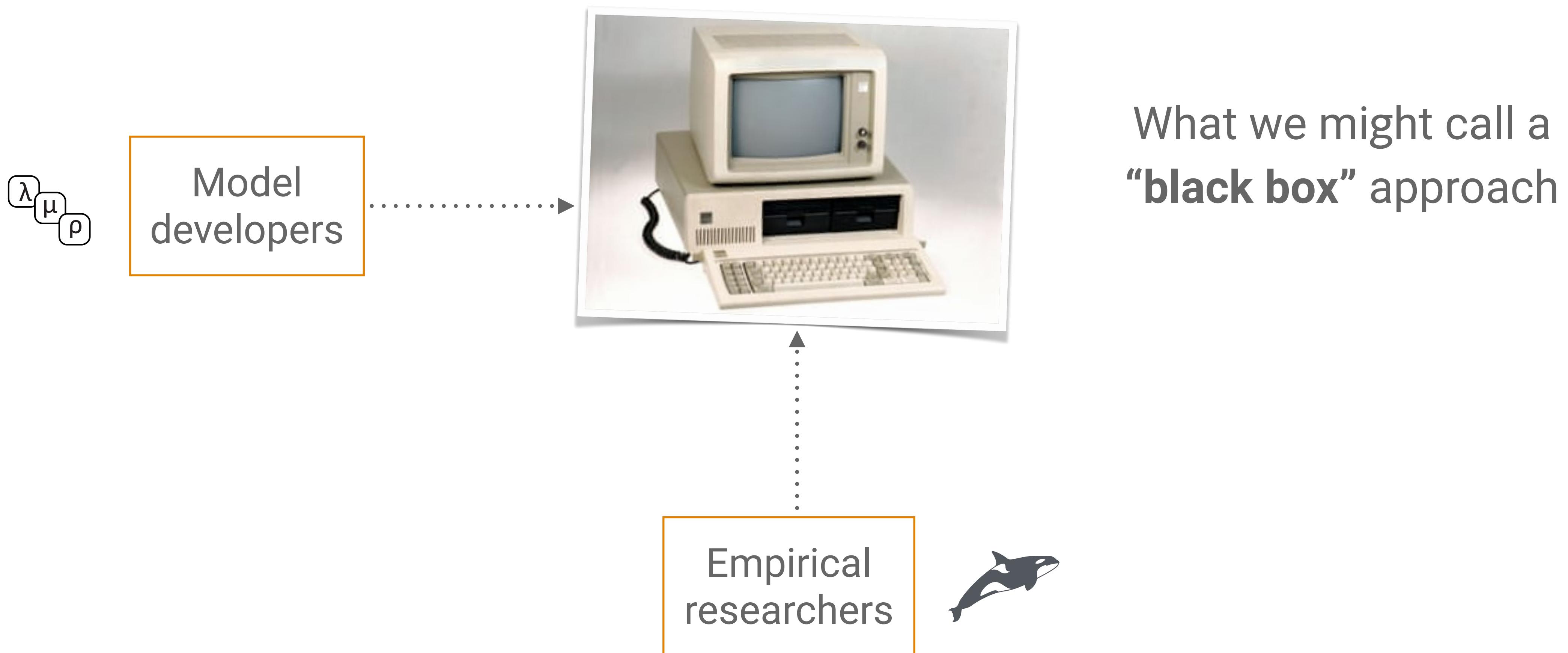
Method	Criterion (tree score)
Maximum parsimony	Minimum number of changes
.....
Maximum likelihood	Likelihood score (probability), optimised over branch lengths and model parameters
.....
Bayesian inference	Posterior probability, integrating over branch lengths and model parameters

Both maximum likelihood and Bayesian inference are model-based approaches

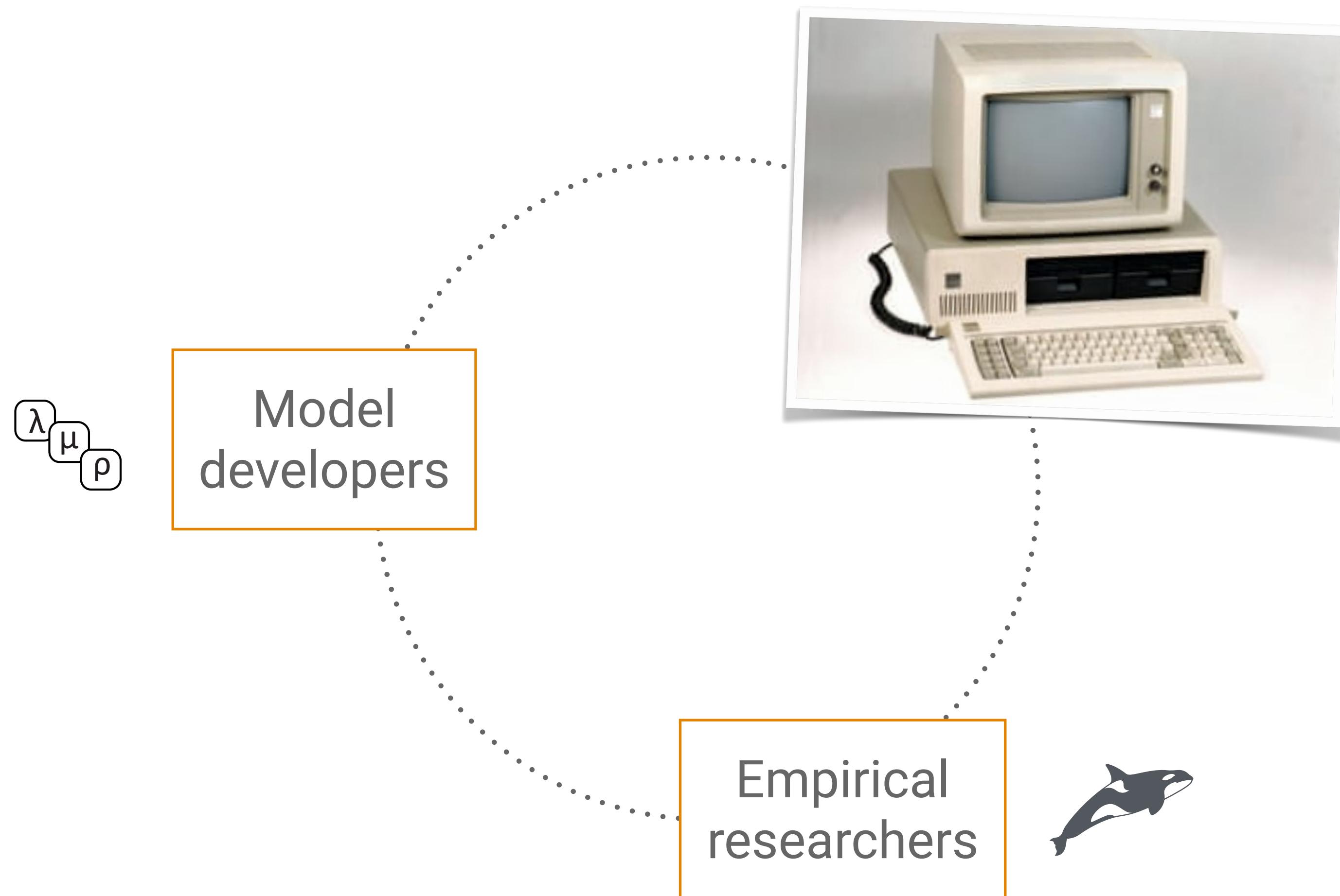
Note these are not the only approaches to tree-building but they are the most widely used

RevBayes

Phylogenetic inference – the old way



Phylogenetic inference – a better way?



The goal is to bring researchers with different expertise together, increase transparency, and do better research

The screenshot shows a web browser window displaying the RevBayes GitHub page (revbayes.github.io). The page features a navigation bar with links for Download, Tutorials, Documentation, Interfaces, Workshops, Jobs, and Developer. A small portrait of a man in historical attire is positioned next to the Download link. The main content area contains a phylogenetic tree diagram on the left and descriptive text on the right.

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](#)

Specifying graphical models using the Rev syntax

Table 1: Rev assignment operators, clamp function, and plate/loop syntax.

Operator	Variable
<code><-</code>	constant variable
<code>~</code>	stochastic variable
<code>:=</code>	deterministic variable
<code>node.clamp(data)</code>	clamped variable
<code>=</code>	inference (<i>i.e.</i> , non-model) variable
<code>for(i in 1:N){...}</code>	plate

Bayesian tree inference

Bayes' theorem

$$\Pr(\text{model} \mid \text{data}) = \frac{\Pr(\text{data} \mid \text{model}) \Pr(\text{model})}{\Pr(\text{data})}$$

Bayes' theorem

Likelihood

The probability of the data given the model assumptions and parameter values

$$\Pr(\text{model} \mid \text{data}) = \frac{\Pr(\text{data} \mid \text{model}) \Pr(\text{model})}{\Pr(\text{data})}$$

Bayes' theorem

Priors

This represents our prior knowledge of the model parameters

$$\Pr(\text{model} \mid \text{data}) = \frac{\Pr(\text{data} \mid \text{model}) \Pr(\text{model})}{\Pr(\text{data})}$$

Bayes' theorem

$$\Pr(\text{model} \mid \text{data}) = \frac{\Pr(\text{data} \mid \text{model}) \Pr(\text{model})}{\Pr(\text{data})}$$

Marginal probability

The probability of the data, given all possible parameter values. Can be thought of as a normalising constant

Bayes' theorem

posterior

Reflects our combined knowledge based on the likelihood and the priors

$\Pr(\text{model} \mid \text{data}) =$

$$\frac{\Pr(\text{data} \mid \text{model}) \Pr(\text{model})}{\Pr(\text{data})}$$

Bayes' theorem – the gist

3 points to remember

We can include **explicit models** that describe the evolutionary and sampling processes

The priors can be used to incorporate **existing knowledge**

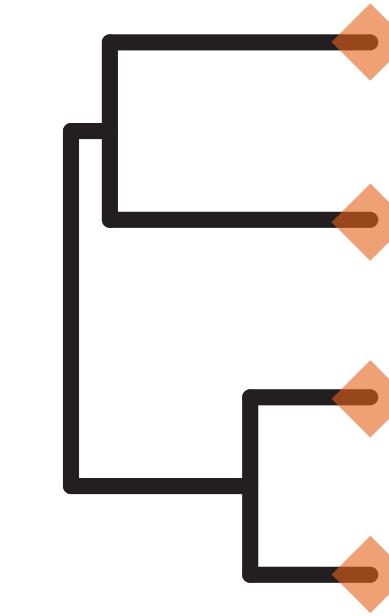
The outcome is not a point estimate, but a **distribution** of plausible parameter values and trees that reflect the uncertainty

Components used to infer trees

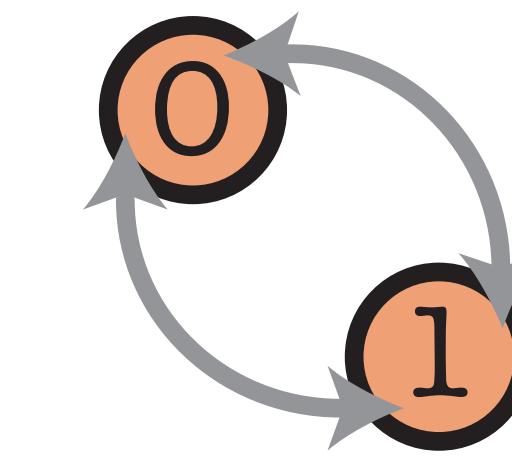
without considering time

0101...
1101...
0100...

data
sequences or
characters



tree model
topology and
branch lengths



substitution
model

Bayesian tree inference

$$\text{posterior} \quad P(E \mid \text{0101...}, \text{1101...}, \text{0100...}) = \frac{\text{likelihood} \quad P(\text{0101...} \mid E) \quad P(E)}{\text{priors} \quad P(\text{0101...}, \text{1101...}, \text{0100...})}$$

Diagram illustrating the components of Bayesian tree inference:

- posterior**: $P(E \mid \text{0101...}, \text{1101...}, \text{0100...})$
- likelihood**: $P(\text{0101...} \mid E)$
- priors**: $P(E)$
- marginal probability**: $P(\text{0101...}, \text{1101...}, \text{0100...})$

The diagram shows a phylogenetic tree with two terminal nodes. The left node is labeled with '0' and the right node with '1'. Arrows indicate the direction of evolution from root to leaves. The likelihood term $P(\text{0101...} \mid E)$ corresponds to the probability of observing the sequence '0101...' at the first node given the tree E . The prior term $P(E)$ represents the probability of the tree structure itself. The marginal probability $P(\text{0101...}, \text{1101...}, \text{0100...})$ is the joint probability of all observed sequences.

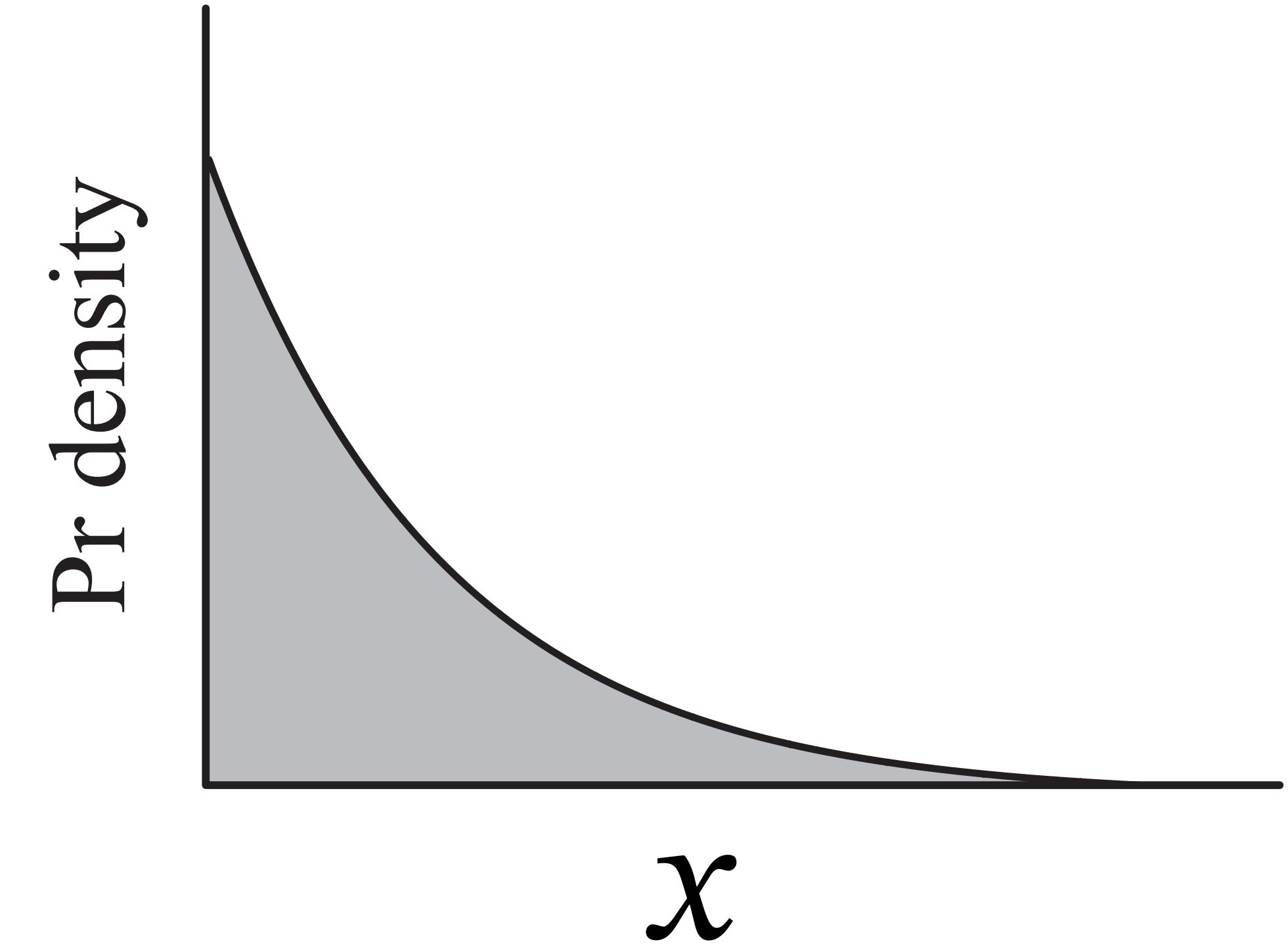
Introduction to MCMC

Probabilities vs probability densities

In phylogenetics, probabilities are not normally discrete (i.e., represented by a single value)

We're often dealing with a lot of uncertainty and typically work with **probability densities**

Probability densities introduce some complexity



The distribution height reflects the relative probability of a given range of values

Bayesian tree inference

$$\text{posterior} \quad P(E \mid \text{0101...}, \text{1101...}, \text{0100...}) = \frac{\text{likelihood} \quad P(\text{0101...} \mid E) \quad P(E)}{\text{priors} \quad P(\text{0101...}, \text{1101...}, \text{0100...})}$$

Diagram illustrating the components of Bayesian tree inference:

- posterior**: $P(E \mid \text{0101...}, \text{1101...}, \text{0100...})$
- likelihood**: $P(\text{0101...} \mid E)$
- priors**: $P(E)$
- marginal probability**: $P(\text{0101...}, \text{1101...}, \text{0100...})$

The diagram shows a phylogenetic tree with two terminal nodes. The left node is labeled '0' and the right node is labeled '1'. Arrows indicate the direction of evolution from root to leaves. The likelihood term $P(\text{0101...} \mid E)$ corresponds to the probability of observing the sequence '0101...' at the first node given the tree E . The prior term $P(E)$ corresponds to the probability of the tree E itself.

Bayesian tree inference

$$= \frac{P(\text{0101...} | \text{E} \circlearrowleft \text{O} \rightarrow \text{1}) P(\text{E} \circlearrowleft \text{O} \rightarrow \text{1})}{\int P(\text{0101...} | \text{E} \circlearrowleft \text{O} \rightarrow \text{1}) P(\text{E} \circlearrowleft \text{O} \rightarrow \text{1}) d\text{E} \circlearrowleft \text{O} \rightarrow \text{1}}$$

this part is incredibly difficult to calculate!

What is Markov chain Monte Carlo (MCMC)?

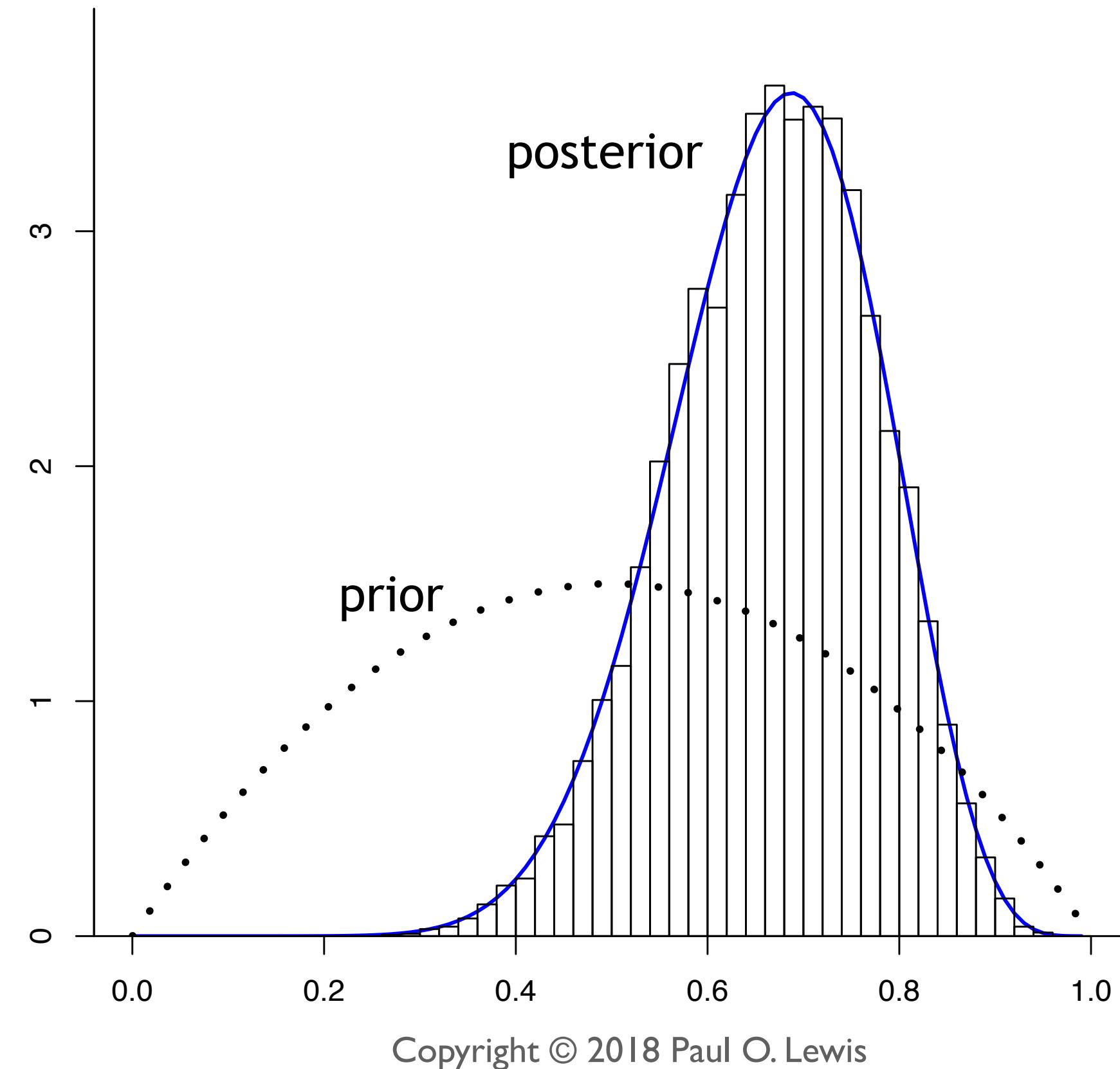
A group of algorithms for approximating the posterior distribution

Markov chain means the progress of the algorithm doesn't depend on its past

Monte Carlo (named for the casino in Monaco) methods estimate a distribution via random sampling

We use this algorithm to visit different regions the parameter space. The number of times a given region is visited will be in proportion to its posterior probability

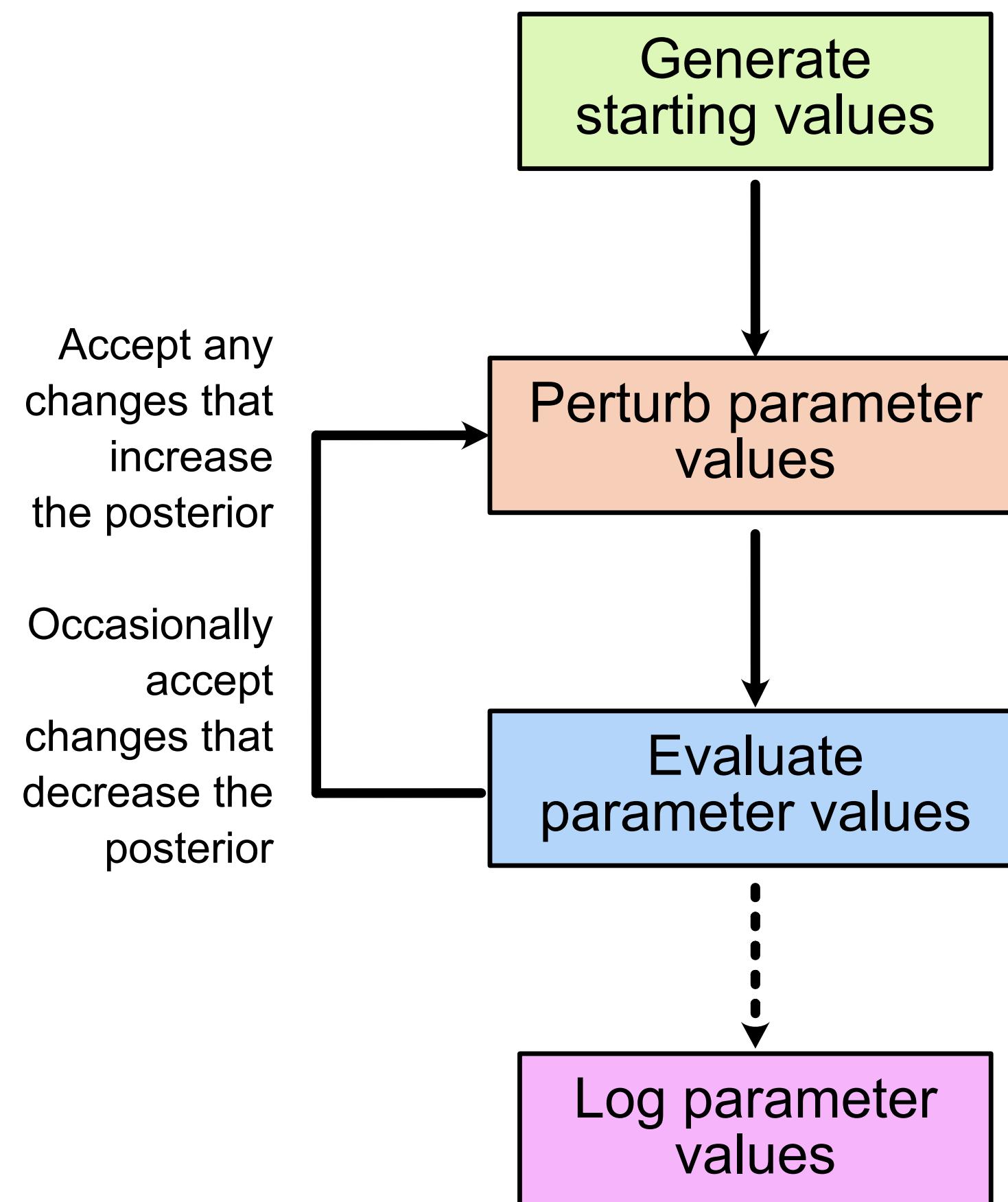
What is Markov chain Monte Carlo (MCMC)?



The aim is to produce a
histogram that provides a good
approximation of the posterior

The Metropolis-Hastings algorithm

Flowchart



Pseudocode

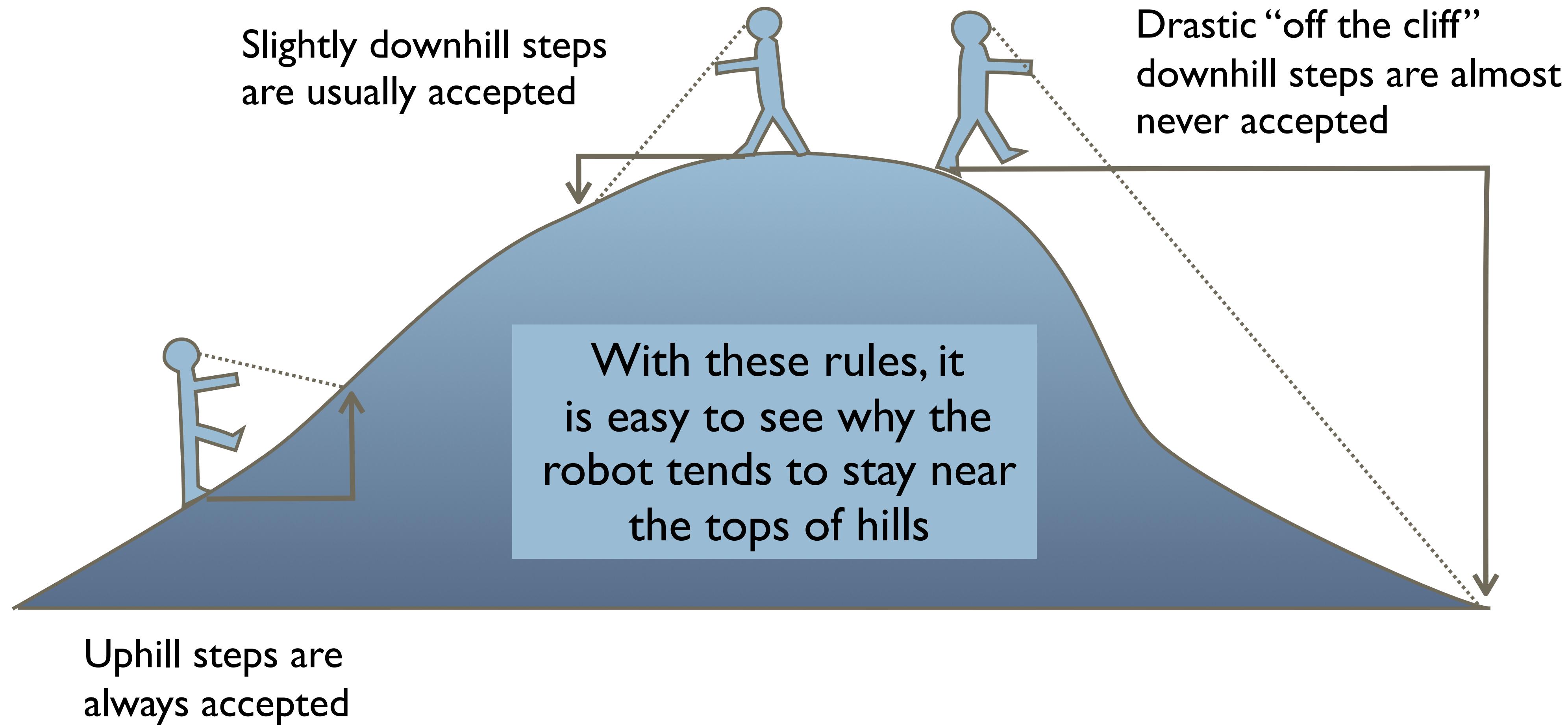
```
initialize starting values;

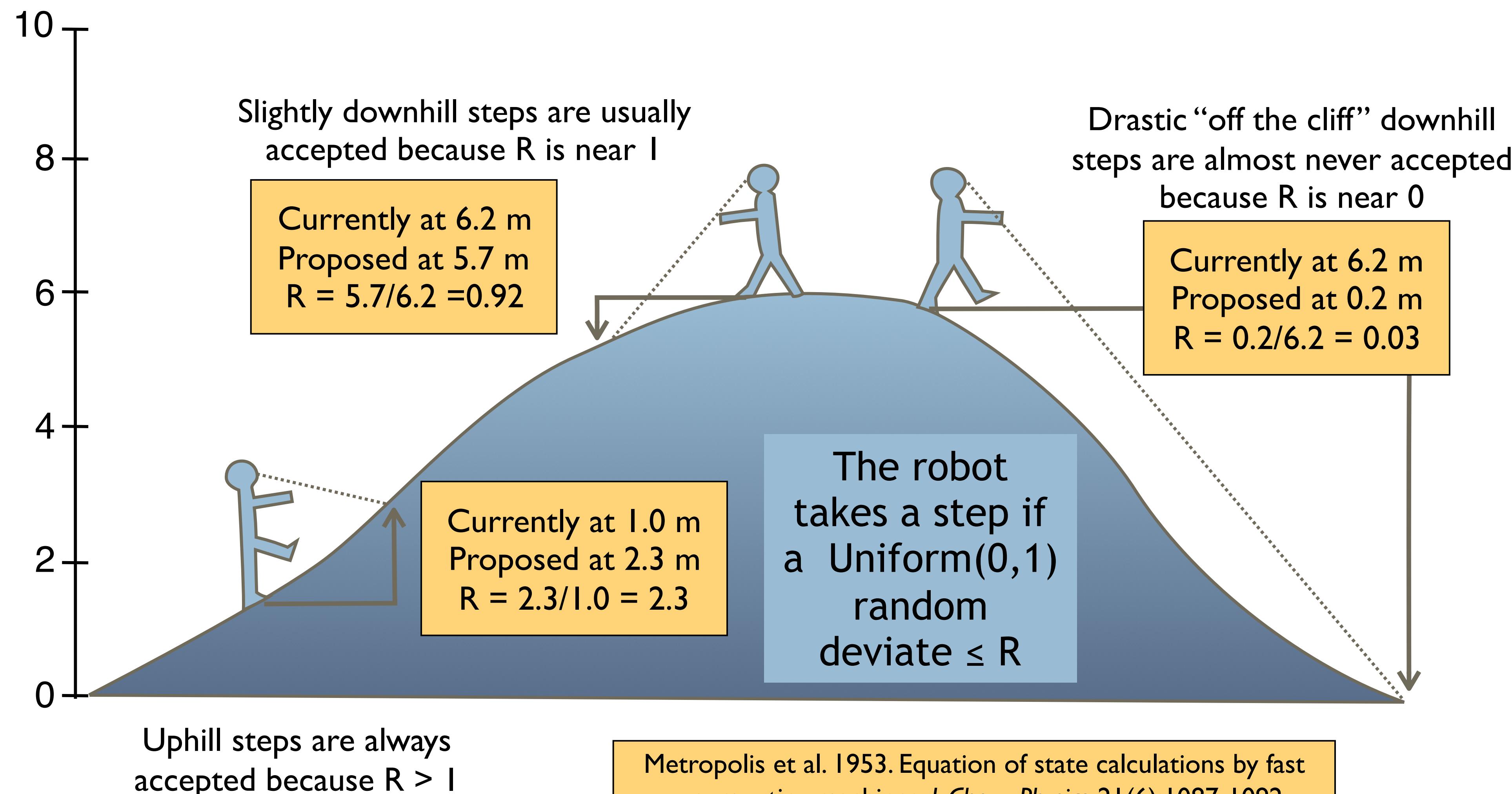
for i in mcmc steps
do
    propose new parameter values;
    calculate the Hastings ratio R;

    if( R > 1 )
        accept the new values;
    else
        accept the new values with Pr = R;

    store the values with frequency j;
done
```

MCMC robot's rules

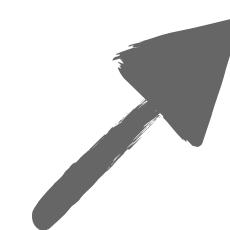




Hastings ratio

new parameter
values

$$R = \frac{P(\text{E}^* | \text{0101... 1101... 0100...})}{P(\text{E} | \text{0101... 1101... 0100...})}$$



=

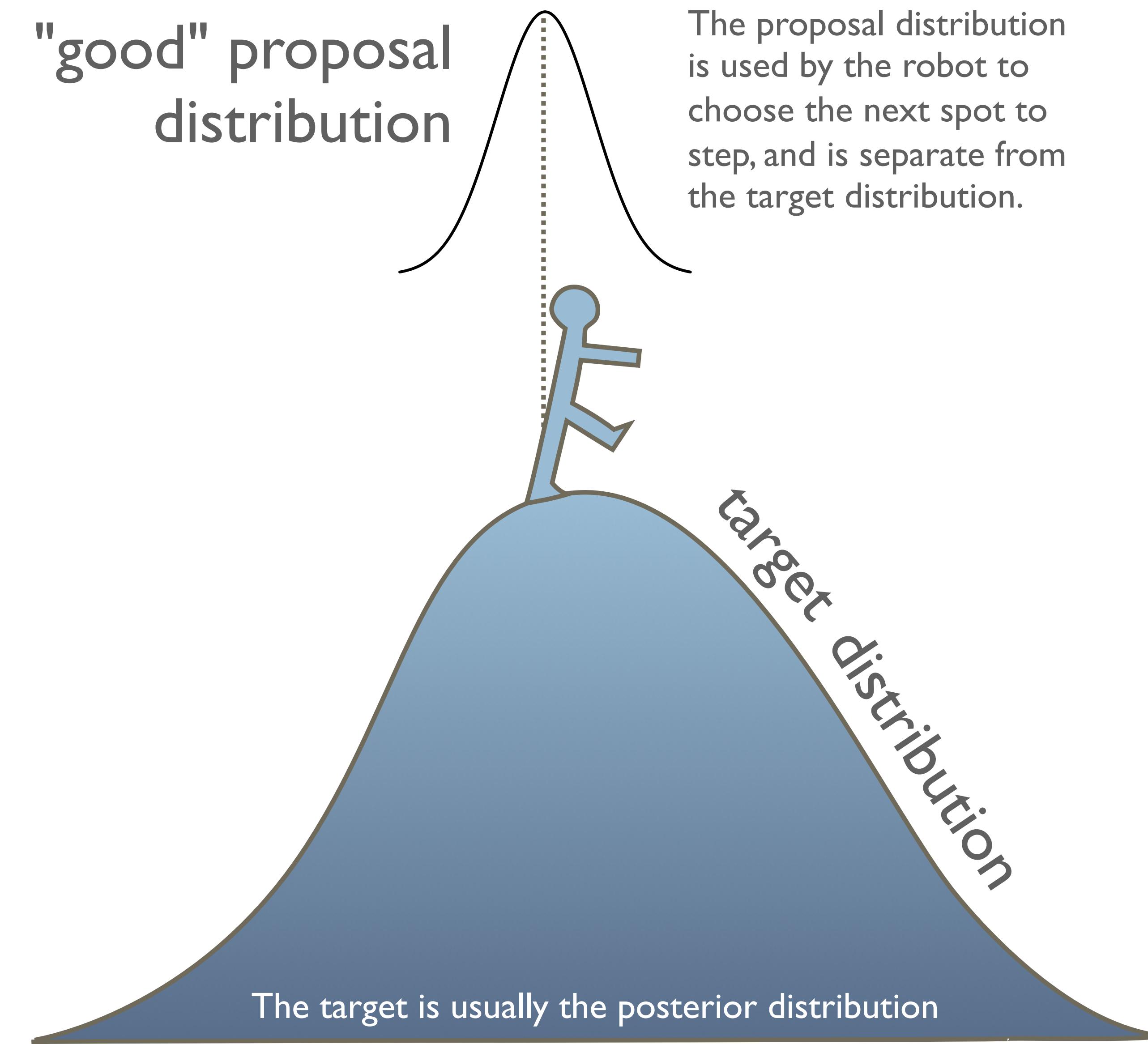
=

$$\frac{\cancel{P(\text{0101... 1101... 0100...})} P(\text{E}^* | \text{0101... 1101... 0100...}) P(\text{E}^*)}{\cancel{P(\text{0101... 1101... 0100...})} P(\text{E} | \text{0101... 1101... 0100...}) P(\text{E})}$$

The marginal probability of the data cancels out

All we're left to calculate is the likelihood ratio and the prior odds ratio

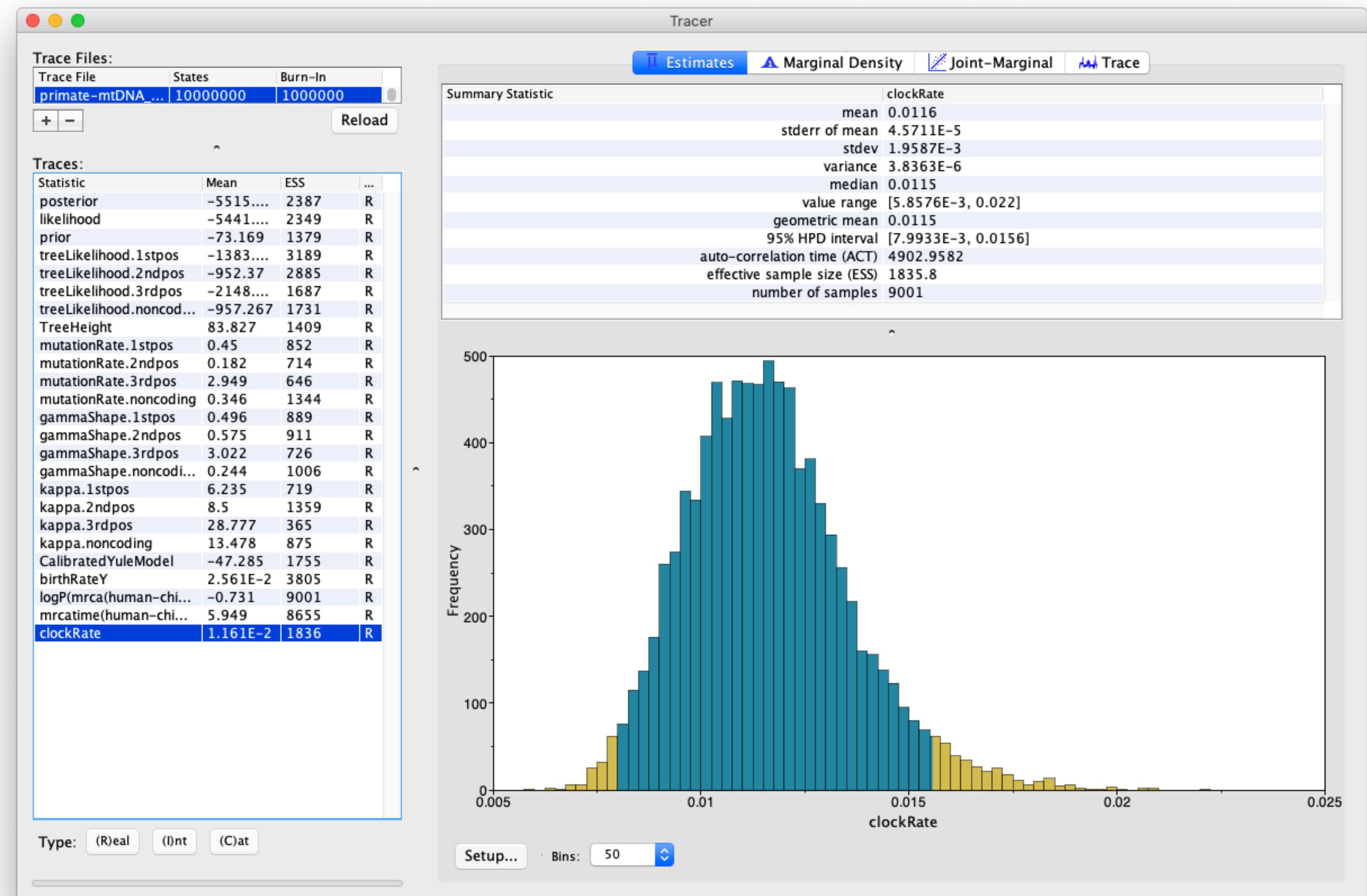
Proposals



Summarising the posterior

Tracer

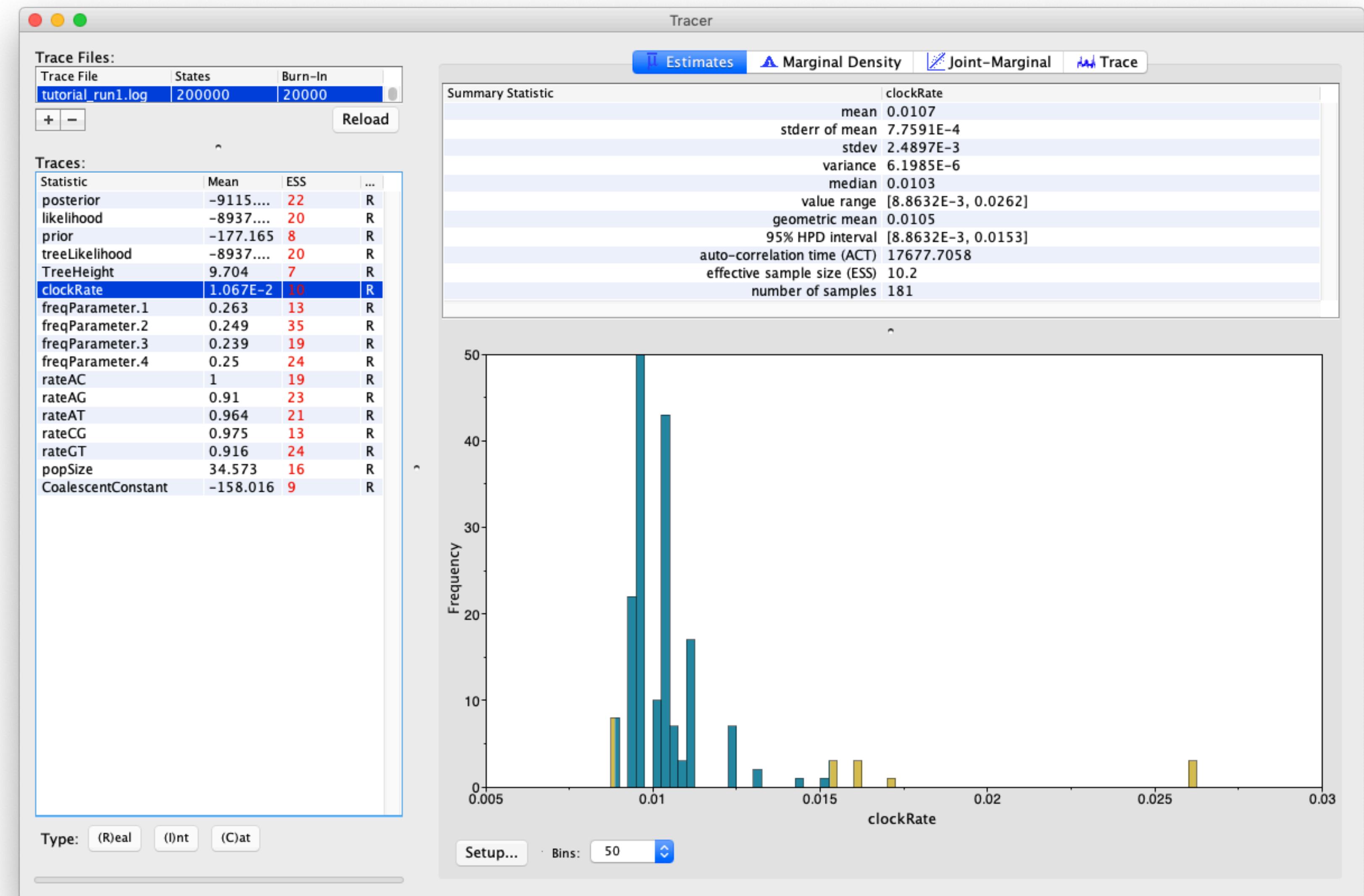
This is an example of good convergence



Summarising the posterior

Tracer

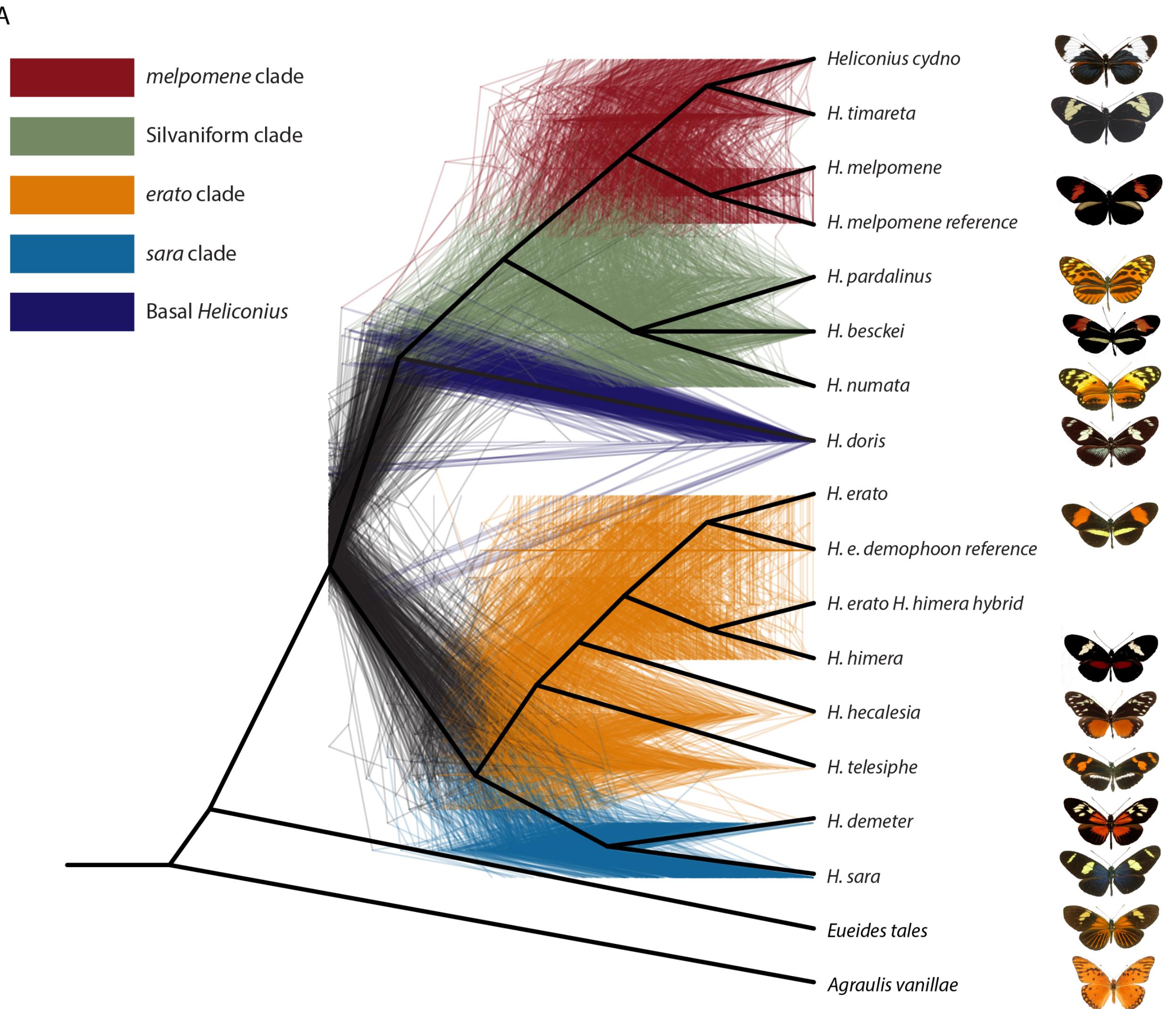
This is an example of poor convergence



Summarising the posterior

Summarising trees is
much more challenging

Presenting a single
summary tree can be
misleading



Summarising the posterior

Maximum clade credibility (MCC) tree – the tree in the posterior sample that has the highest posterior probability (i.e., clade support) across all nodes

The **95% highest posterior density (HPD)** – the shortest interval that contains 95% of the posterior probability. The Bayesian equivalent of the 95% confidence interval

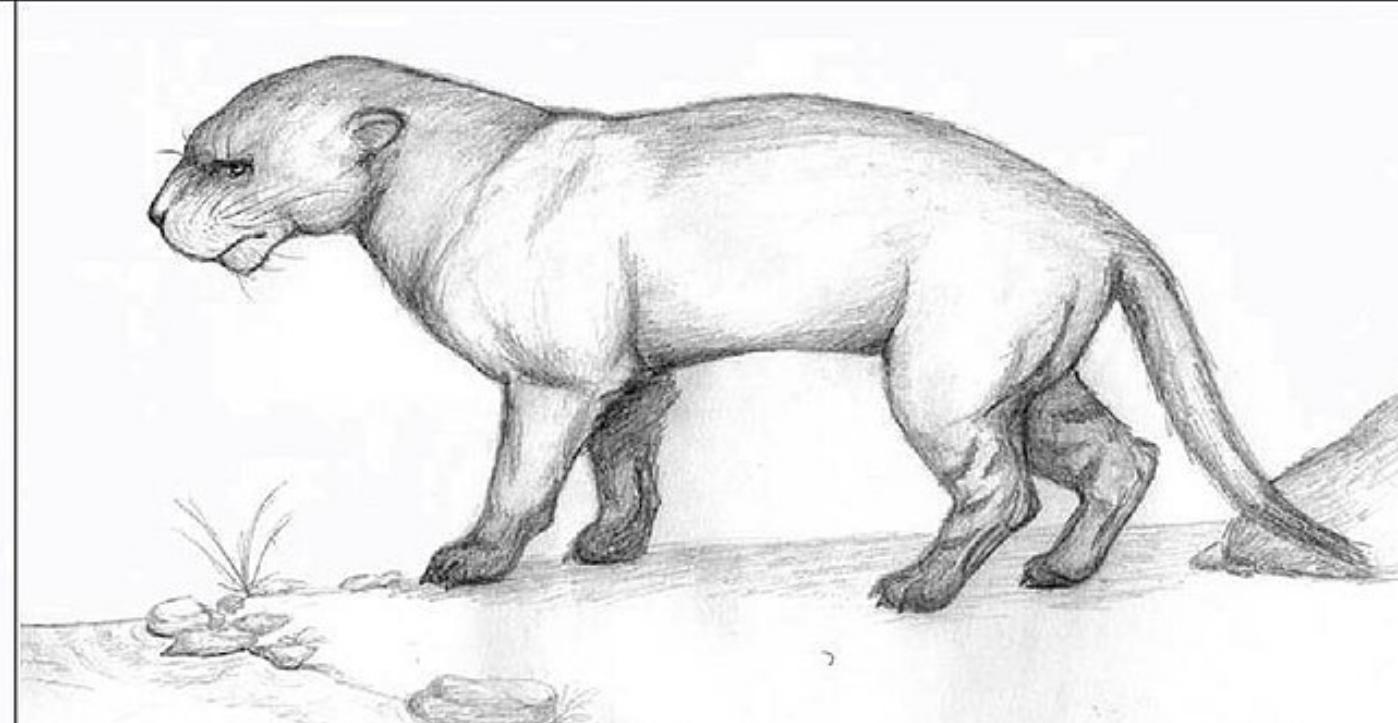
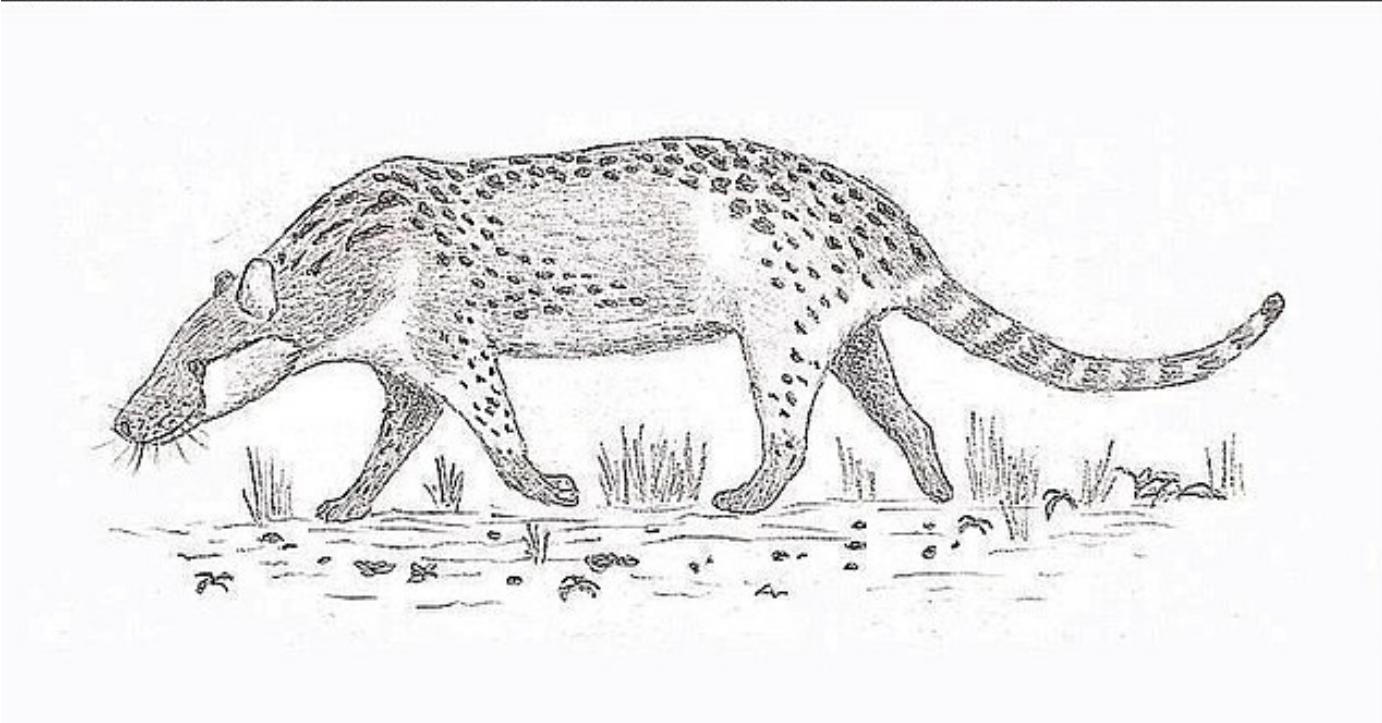
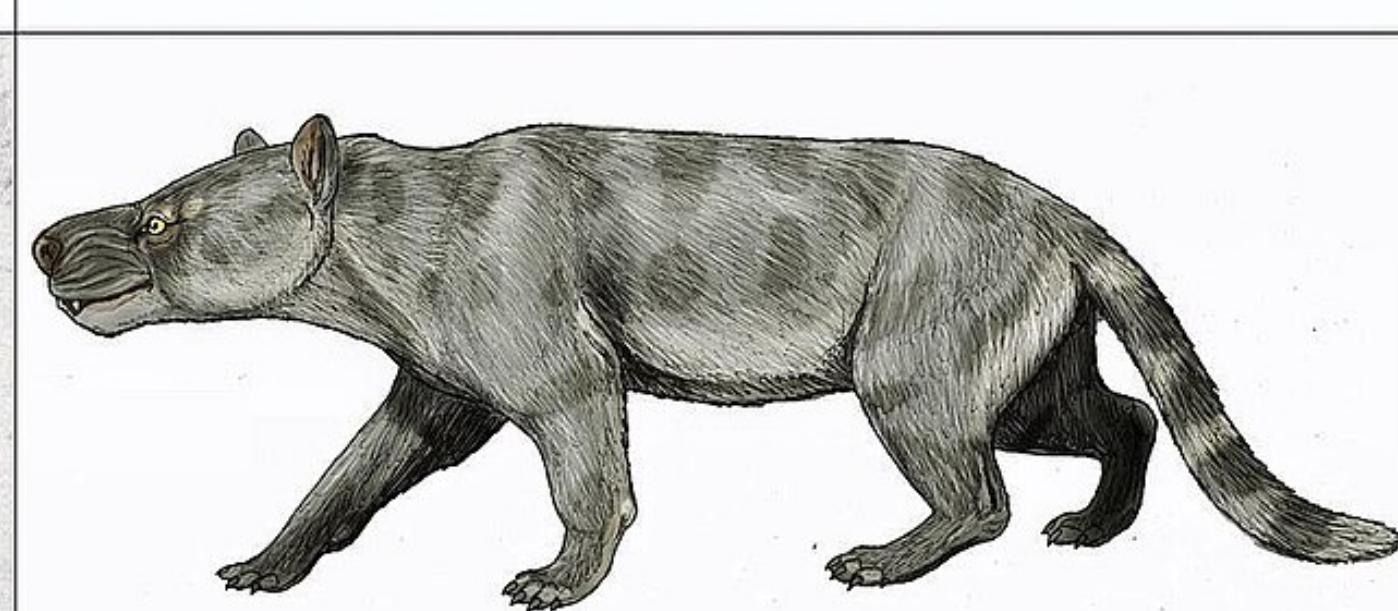
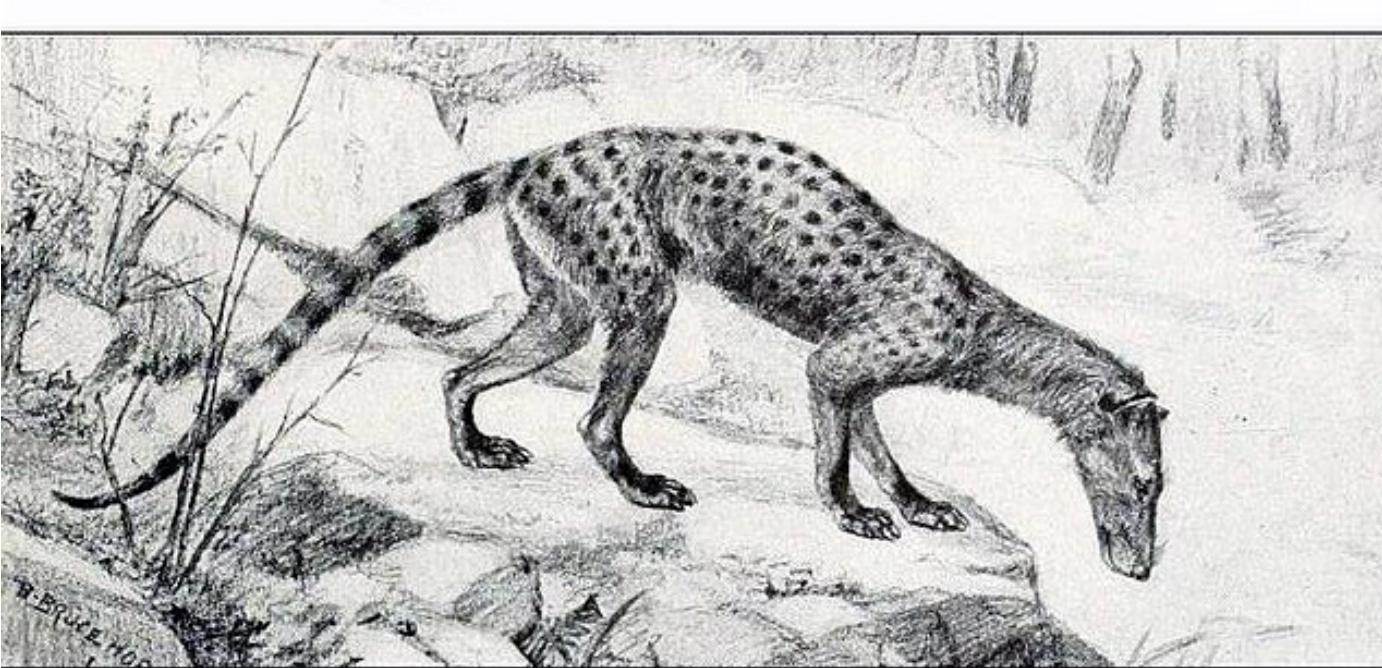
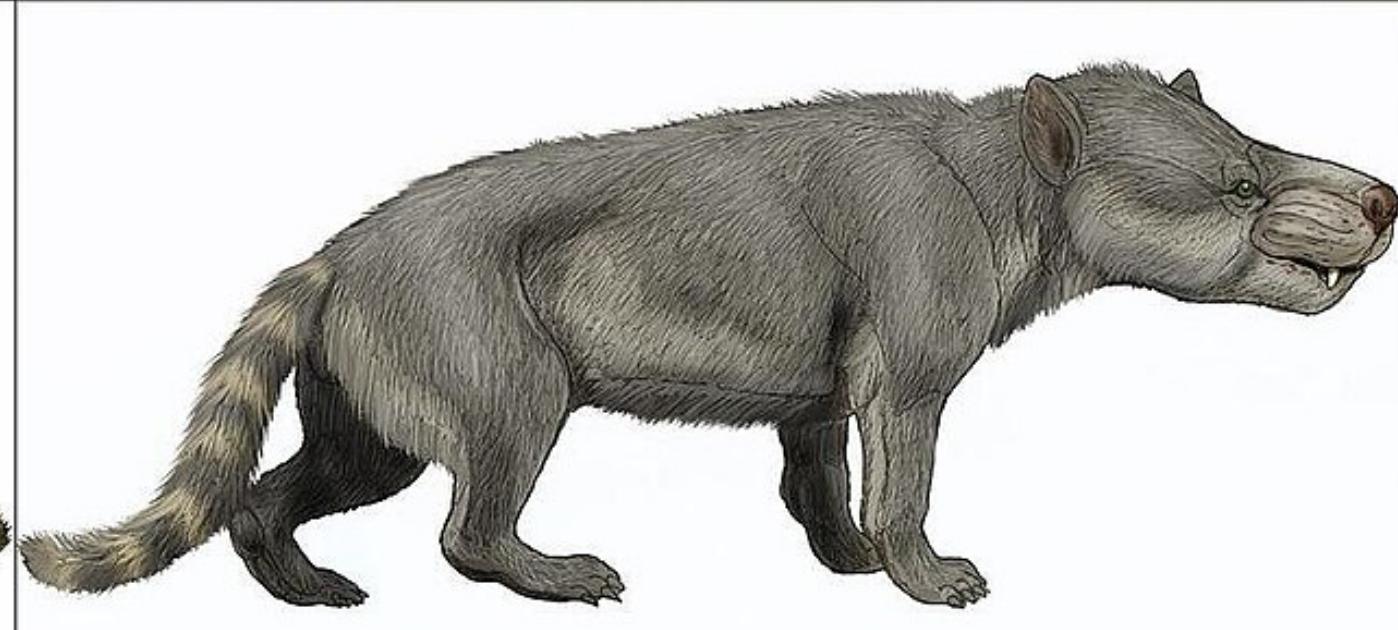
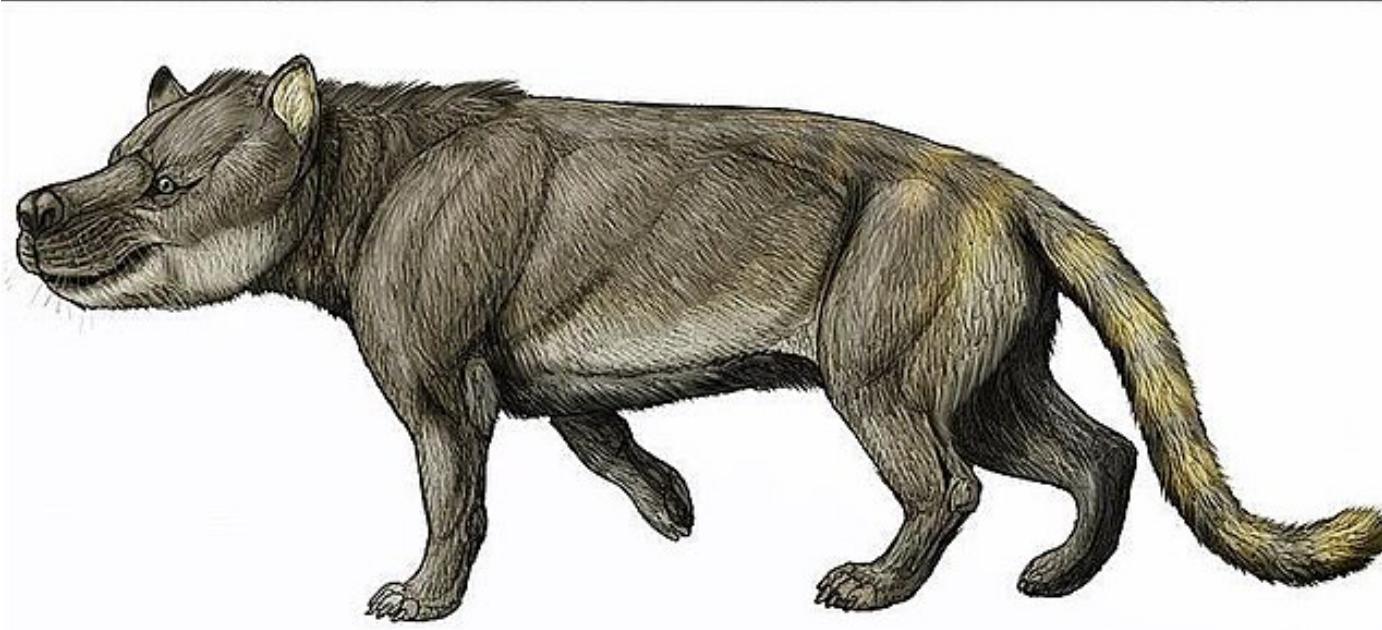
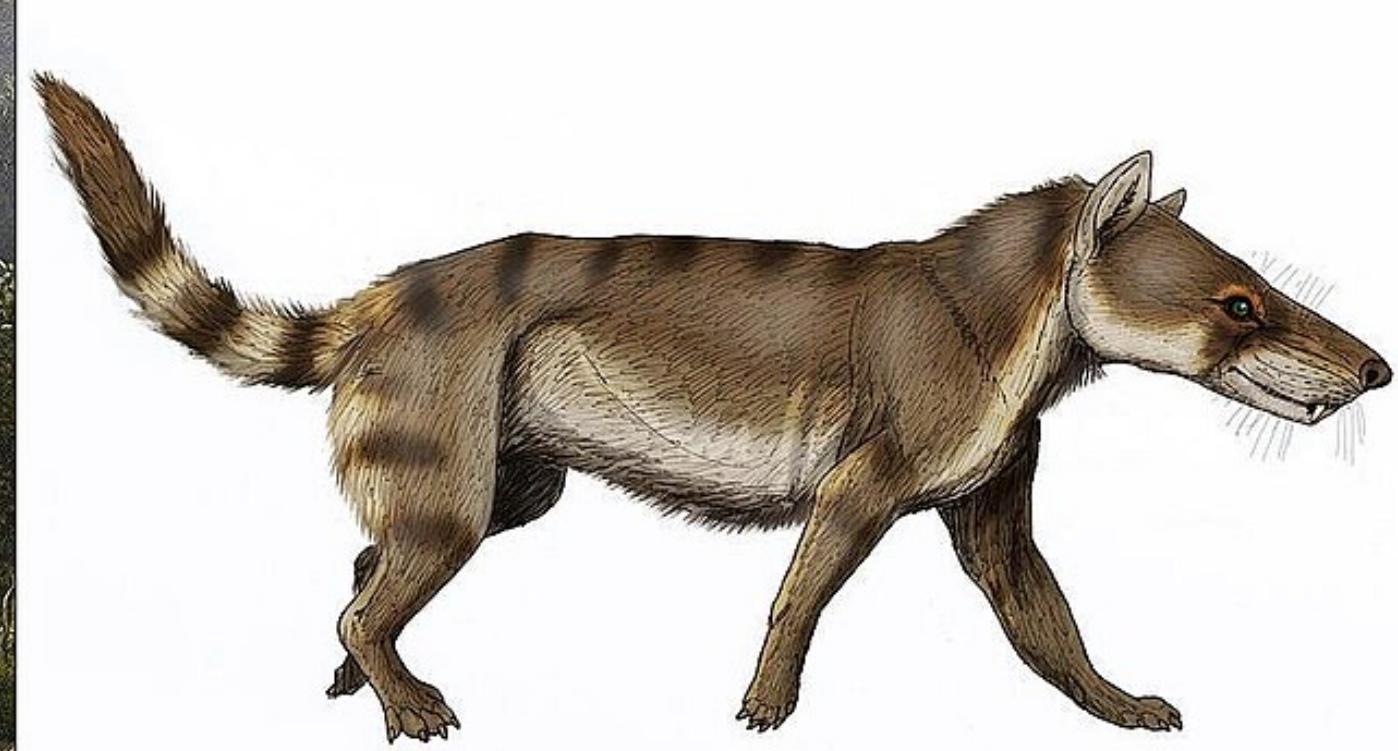
Marginal posterior density – the probability of a parameter regardless of the value of the others, represented by the histogram

Exercise overview

Hyaenodonts

12 taxa and 65 characters

Extinct order of hyper-carnivorous placentals



Data from Egi et al. (2007)
Image source Wikipedia

#NEXUS

[File downloaded from graemetlloyd.com]

BEGIN DATA;

DIMENSIONS NTAX=15 NCHAR=65;

FORMAT DATATYPE = STANDARD SYMBOLS= " 0 1 2 3 4" MISSING=? GAP=- ;

MATRIX

missing_ORIG_Procerberus_and_Cimolestes	0000000100000000001000?0010000000000011000000100000?000000000
missing_ORIG_Proviverra_and_Lesmesodon	100????0?00110001001010??0001013111021??20?1100001021010000110000
Allopteronodon_ORIG_Allopteronodon	1012111011011011111110?110011210001111010110000212100002010000
Leonhardtina_ORIG_Leonhardtina	1112121021001110001000?11000120000110212022120001120000120200001
Eurotherium_ORIG_Eurotherium	20110010101100101111102110011000011100120111122010211001022000?0
Prodiissopsalis_ORIG_Prodiissopsalis	30100010111000111?1112??0?00100011001011110201001100110200000
Cynohyaenodon_ORIG_Cynohyaenodon	11110010111102111111100011000001201020110200110101110200000
no_age_data_ORIG_Paracynohyaenodon	111100101111211111101110?1110001201020111020112110110?200000
Brychotherium_ORIG_Brychotherium	10110????11112?111111021001200100012020201210221201120102?300111
Metasinopa_ORIG_Metasinopa	21??????00112101110110??11110020001202020111121211111100200010
Dissopsalis_ORIG_Dissopsalis	42010012112201111202222101?1?02000221001111121201110101300110
Anasinopa_ORIG_Anasinopa	311??011112210101120110??0000010002020012111121211011110201010
Paratritemnodon_ORIG_Paratritemnodon	100101121101121011101122010?0?11??0021101102011212110101001000010
Masrasector_ORIG_Masrasector	1010002?00112?0111001??101101000020200111011112110001101101010
Kyawdawia_ORIG_Kyawdawia	42010212111012101111110010???0?0?1201001120111121?0001101100010

;

END;

The screenshot shows a Mac OS X Finder window with the following details:

Window Title: Hyaenodonts_analysis

Toolbar: Includes standard Mac OS X icons for file operations (red, yellow, green circles), navigation (back, forward), and search.

View Options: Shows icons for grid view, list view, and other display settings.

Table Headers: Name, Date Modified, Size, Kind.

Content:

Name	Date Modified	Size	Kind
data	Today at 18:28	--	Folder
Egi.nex	5. December 2024 at 20:35	2 KB	Document
main.Rev	5. December 2024 at 21:33	1 KB	Visual S...cument
output	Today at 18:27	--	Folder
scripts	Today at 18:27	--	Folder
Mk.Rev	5. December 2024 at 21:32	93 bytes	Visual S...cument

Path Bar: Macintosh HD > Users > warnock > Hyaenodonts_analysis

Files

Data

- **Egi.nex** (the nexus file, i.e., the character data)

Scripts

- **main.Rev** (used to read in the data, set up the tree model & MCMC settings)
- **Mk.Rev** (used the model of character evolution)

main.Rev

part 1

```
# read in character data
morpho <- readDiscreteCharacterData("data/Egi.nex")

# set up some useful variables
num_taxa <- morpho.ntaxa() # number of taxa
num_branches <- 2 * num_taxa - 3 # number of branches in an unrooted tree
taxa <- morpho.taxa() # list of taxon names

moves      = VectorMoves()
monitors = VectorMonitors()
```

main.Rev part 2

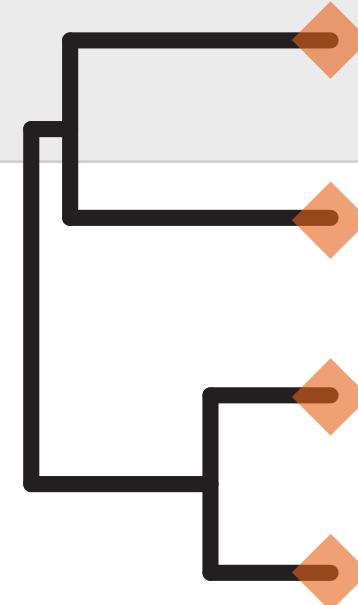
```
# prior on the tree topology
topology ~ dnUniformTopology(taxa)

moves.append( mvNNI(topology, weight = num_taxa) ) # nearest neighbour interchange
moves.append( mvSPR(topology, weight = num_taxa/10.0) ) # subtree pruning and regrafting

# prior on the branch length
for (i in 1:num_branches) {
  br_lens[i] ~ dnExponential(10.0)
  moves.append( mvScale(br_lens[i]) )
}

phylogeny := treeAssembly(topology, br_lens)
TL := sum(br_lens)
```

prior on the tree
topology &
branch lengths

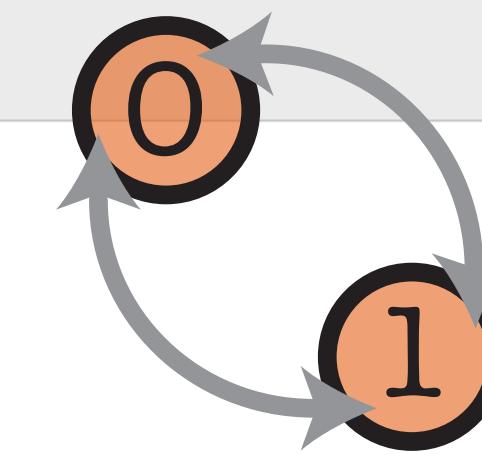


```
# define the Q matrix
Q <- fnJC(5)

# bring everything together
seq ~ dnPhyloCTMC(tree = phylogeny, Q = Q, type = "Standard")

seq.clamp(morpho)
```

Substitution
model



main.Rev

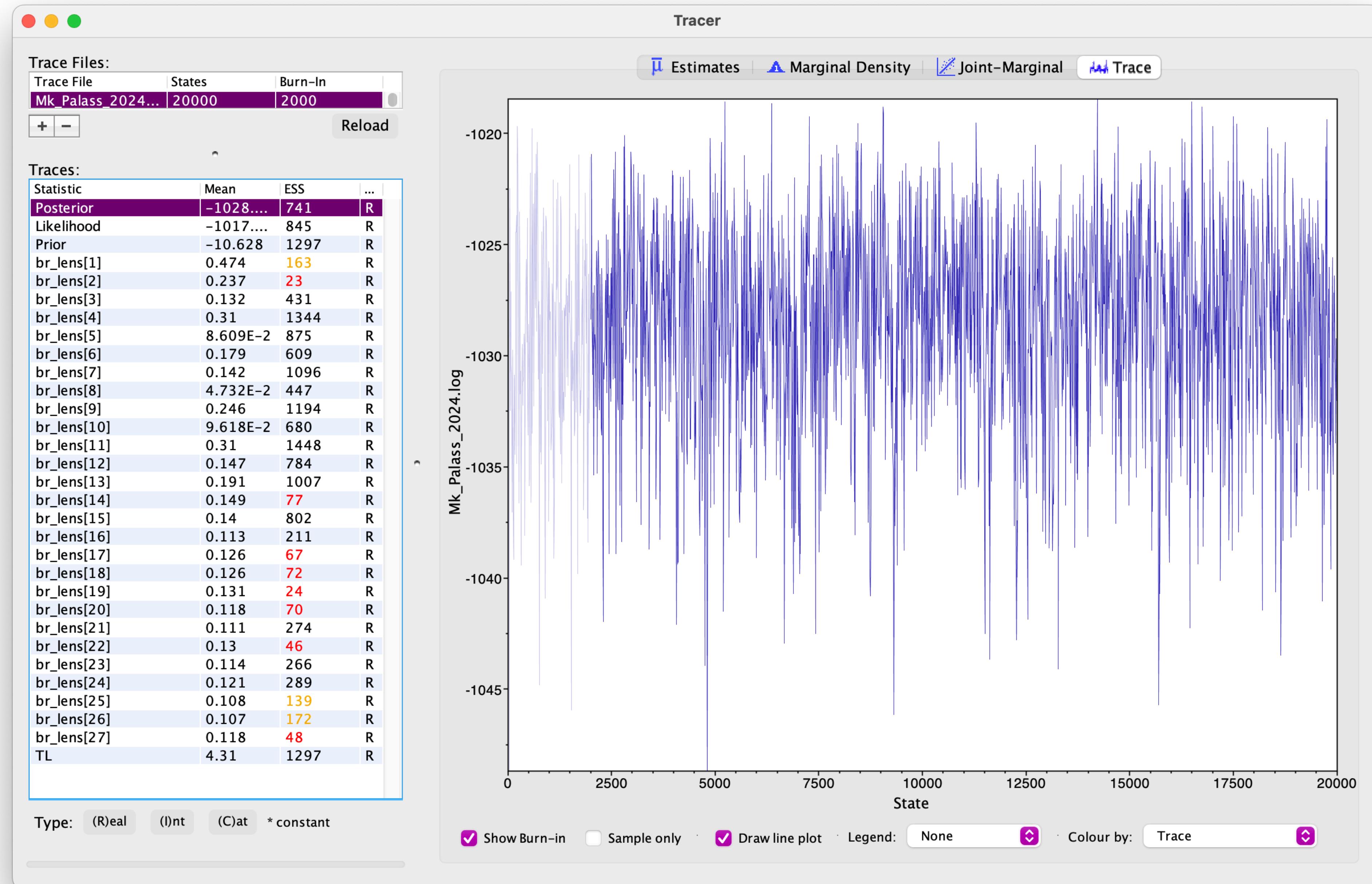
part 3

```
source("scripts/Mk.Rev")

mymodel = model(phylogeny)

# parameters printed to file
monitors.append( mnModel(filename = "output/Mk_Palass_2024.log", printgen = 10) )
# trees printed to file
monitors.append( mnFile(filename = "output/Mk_Palass_2024.trees", printgen = 10,
phylogeny) )
# parameter values printed to screen during the MCMC
monitors.append( mnScreen(printgen = 100, TL) )

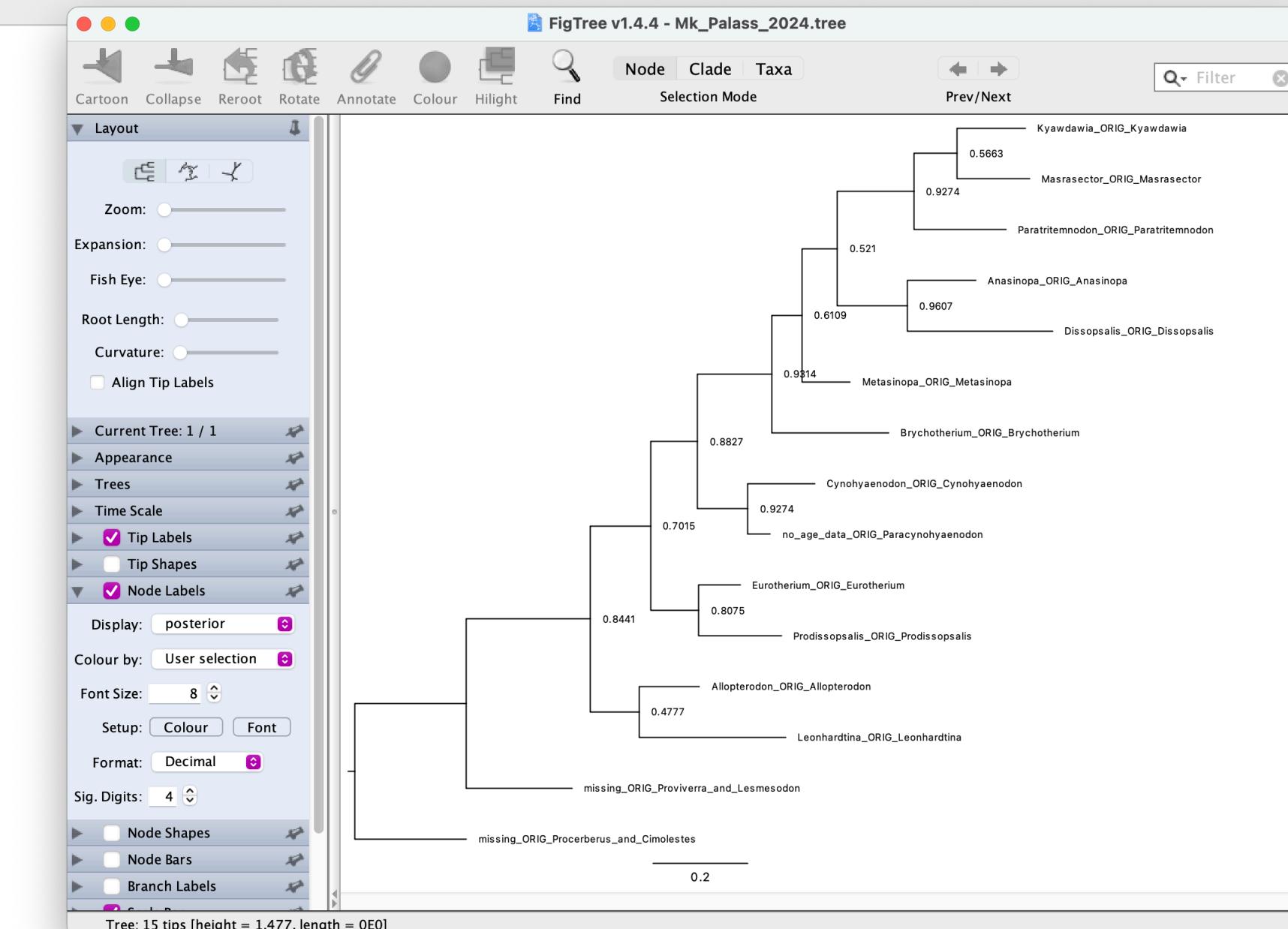
mymcmc = mcmc(mymodel, monitors, moves)
mymcmc.run(generations = 20000)
```



main.Rev part 4

```
# read the tree file back in
treetrace = readTreeTrace("output/Mk_Palass_2024.trees", treetype = "non-clock")

# generate an MAP summary tree
map_tree = mapTree(treetrace, "output/Mk_Palass_2024.tree")
```



Exercise

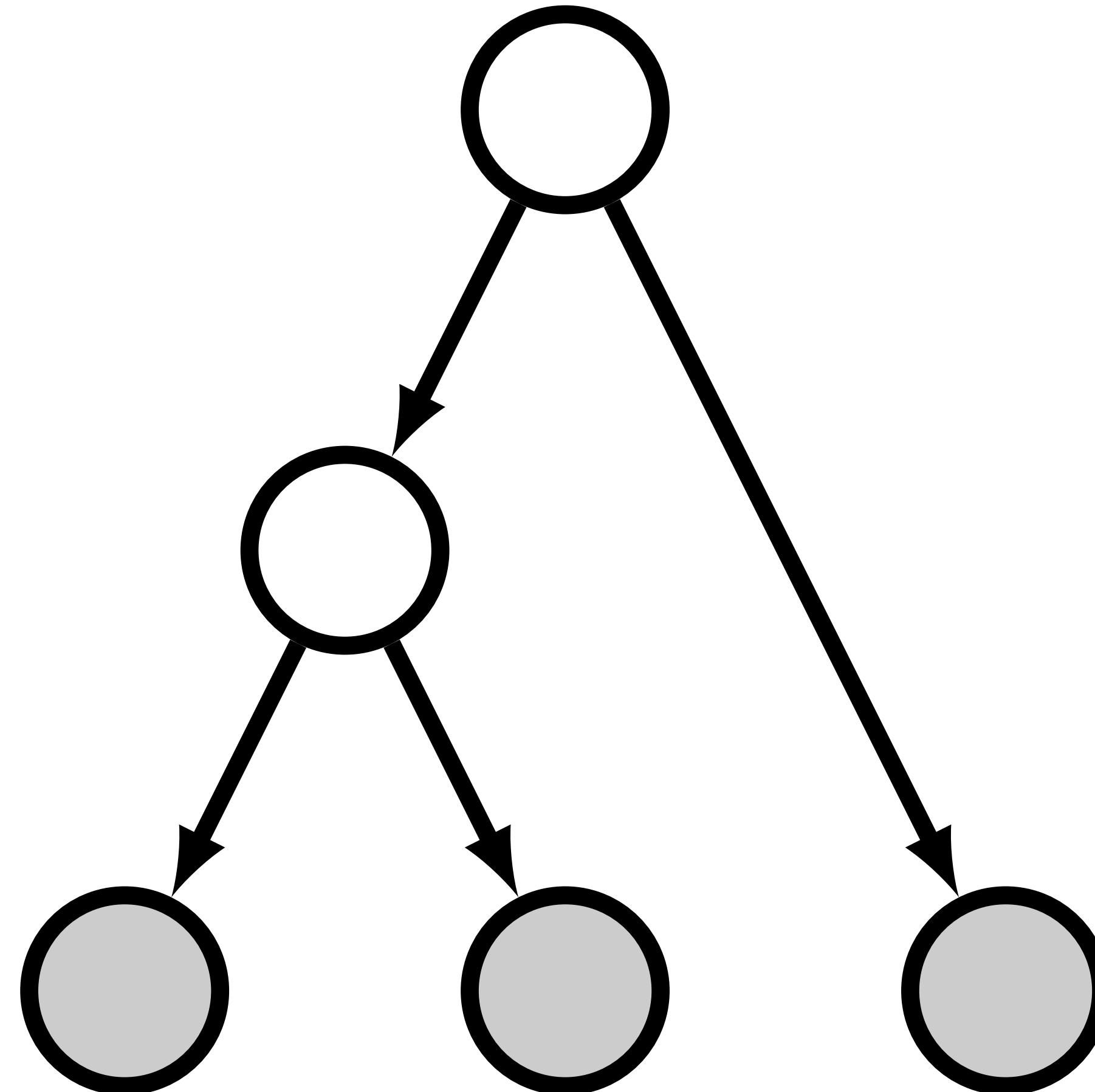
Extra slides

Graphical models

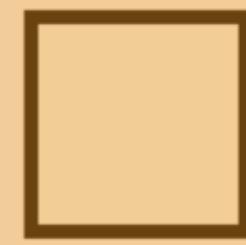
Graphical models

Provide tools for visually and computationally representing complex, parameter-rich models

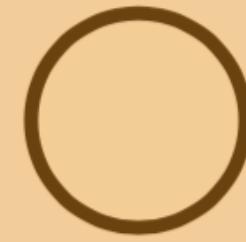
Depict the conditional dependence structure of parameters and other random variables



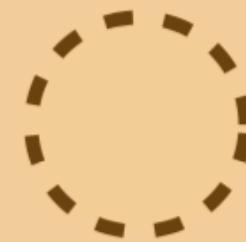
Types of variables (nodes)



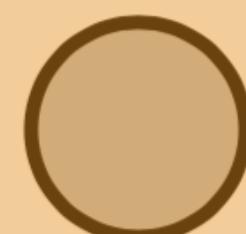
a) Constant node



b) Stochastic node



c) Deterministic node



d) Clamped node
(observed)

a. fixed value variables

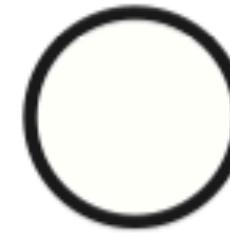
b. random variables that depend on other variables

c. variables determined by a function applied other variables (transformations)

d. observed stochastic variables (data)



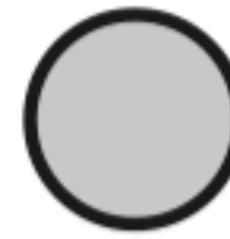
a) Constant node



b) Stochastic node



c) Deterministic node



d) Clamped node
(observed)



e) Plate

a. fixed value variables

b. random variables that depend on other variables

c. variables determined by a function applied other variables (transformations)

d. observed stochastic variables (data)

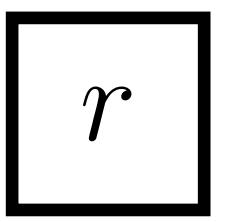
e. repetition over multiple variables (equivalent to a loop)

Specifying graphical models using the Rev syntax

Table 1: Rev assignment operators, clamp function, and plate/loop syntax.

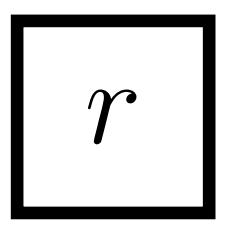
Operator	Variable
<code><-</code>	constant variable
<code>~</code>	stochastic variable
<code>:=</code>	deterministic variable
<code>node.clamp(data)</code>	clamped variable
<code>=</code>	inference (<i>i.e.</i> , non-model) variable
<code>for(i in 1:N){...}</code>	plate

a)

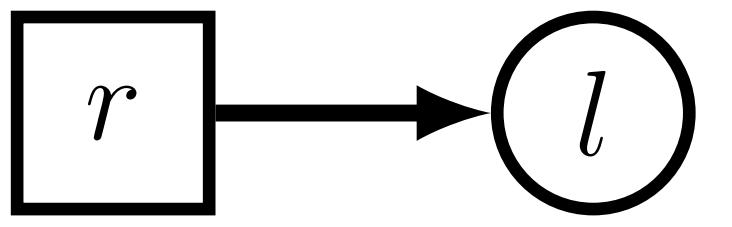


```
# constant node  
r <- 10
```

a)



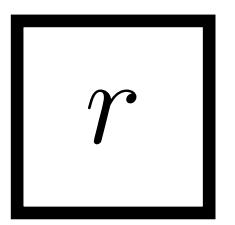
b)



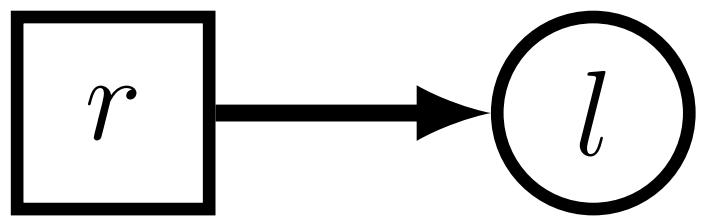
```
# constant node  
r <- 10
```

```
# stochastic node  
l ~ dnExp(r)
```

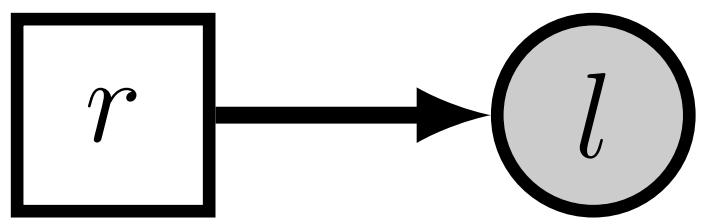
a)



b)



c)

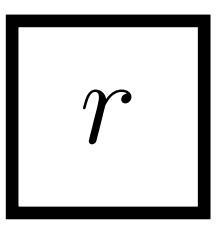


```
# constant node  
r <- 10
```

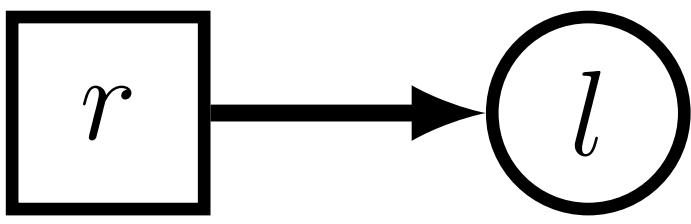
```
# stochastic node  
l ~ dnExp(r)
```

```
# stochastic node (observed)  
l.clamp(0.1)
```

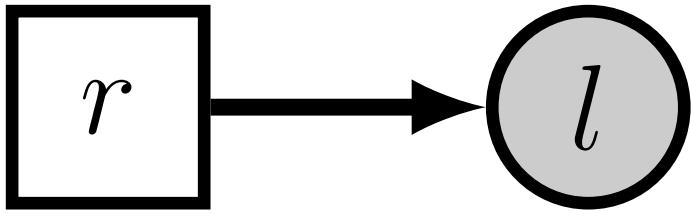
a)



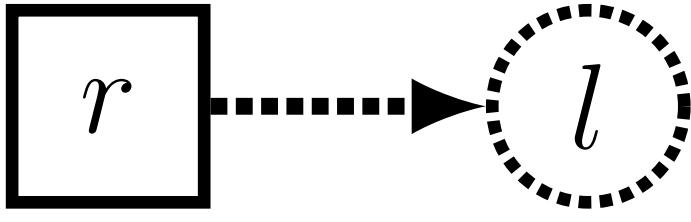
b)



c)



d)



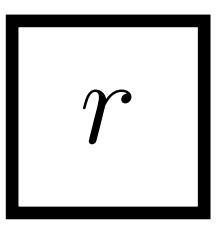
```
# constant node  
r <- 10
```

```
# stochastic node  
l ~ dnExp(r)
```

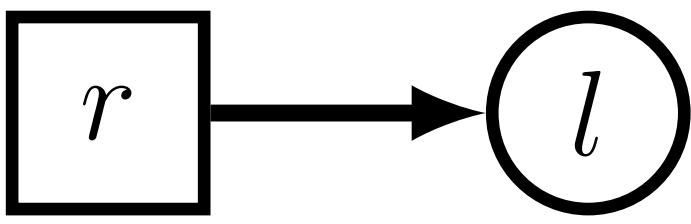
```
# stochastic node (observed)  
l.clamp(0.1)
```

```
# deterministic node  
l := exp(r)
```

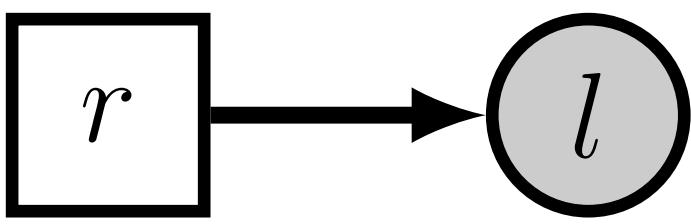
a)



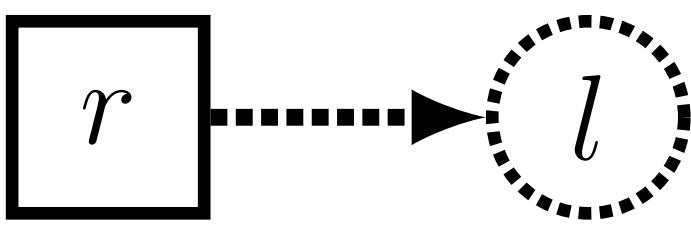
b)



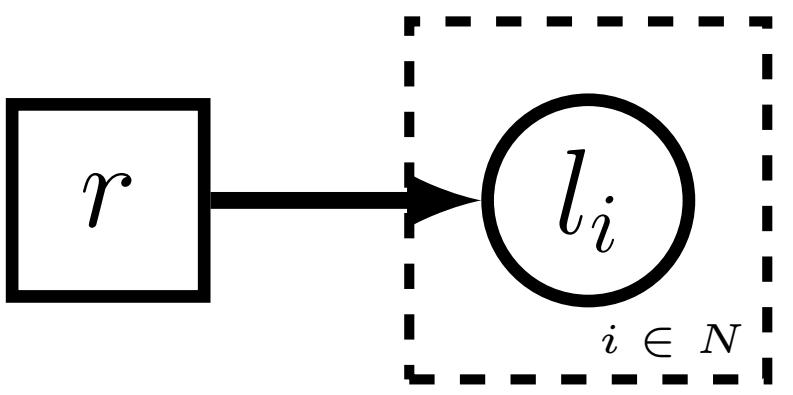
c)



d)



e)



```
# constant node  
r <- 10
```

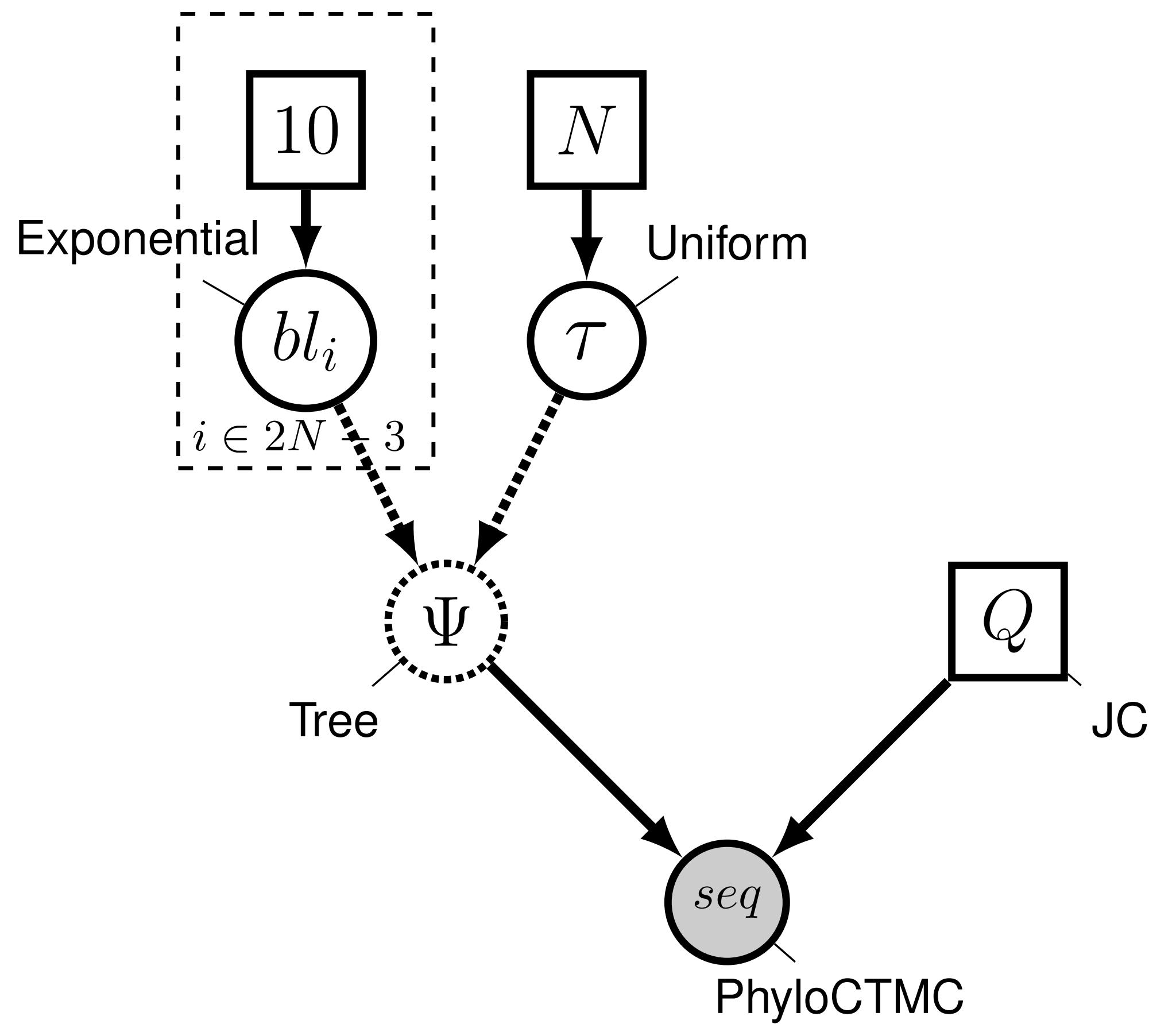
```
# stochastic node  
l ~ dnExp(r)
```

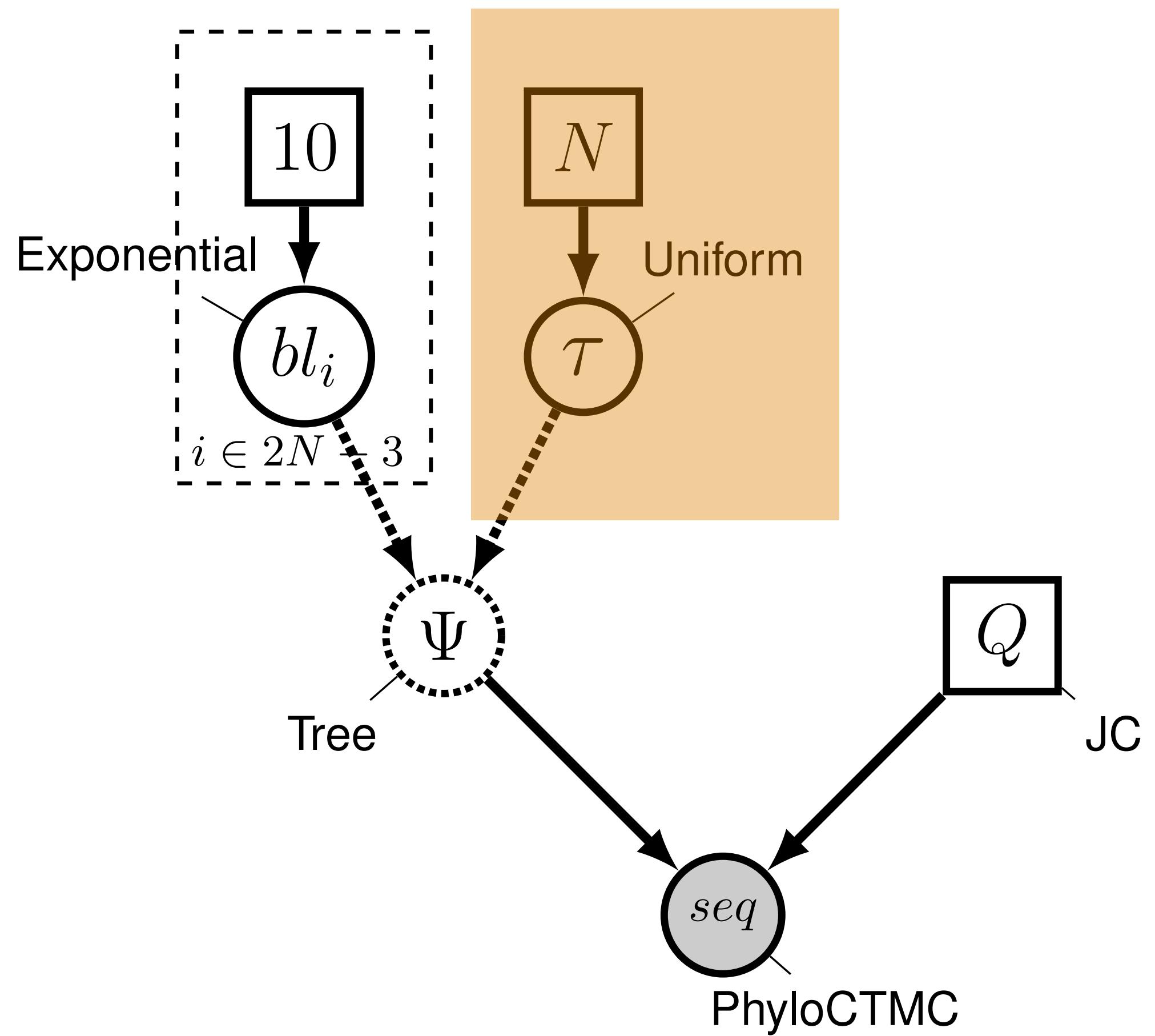
```
# stochastic node (observed)  
l.clamp(0.1)
```

```
# deterministic node  
l := exp(r)
```

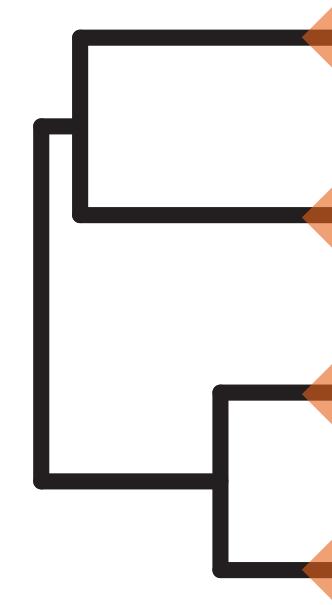
```
# stochastic nodes (iid)  
for (i in 1:N) {  
  l[i] ~ dnExp(r)  
}
```

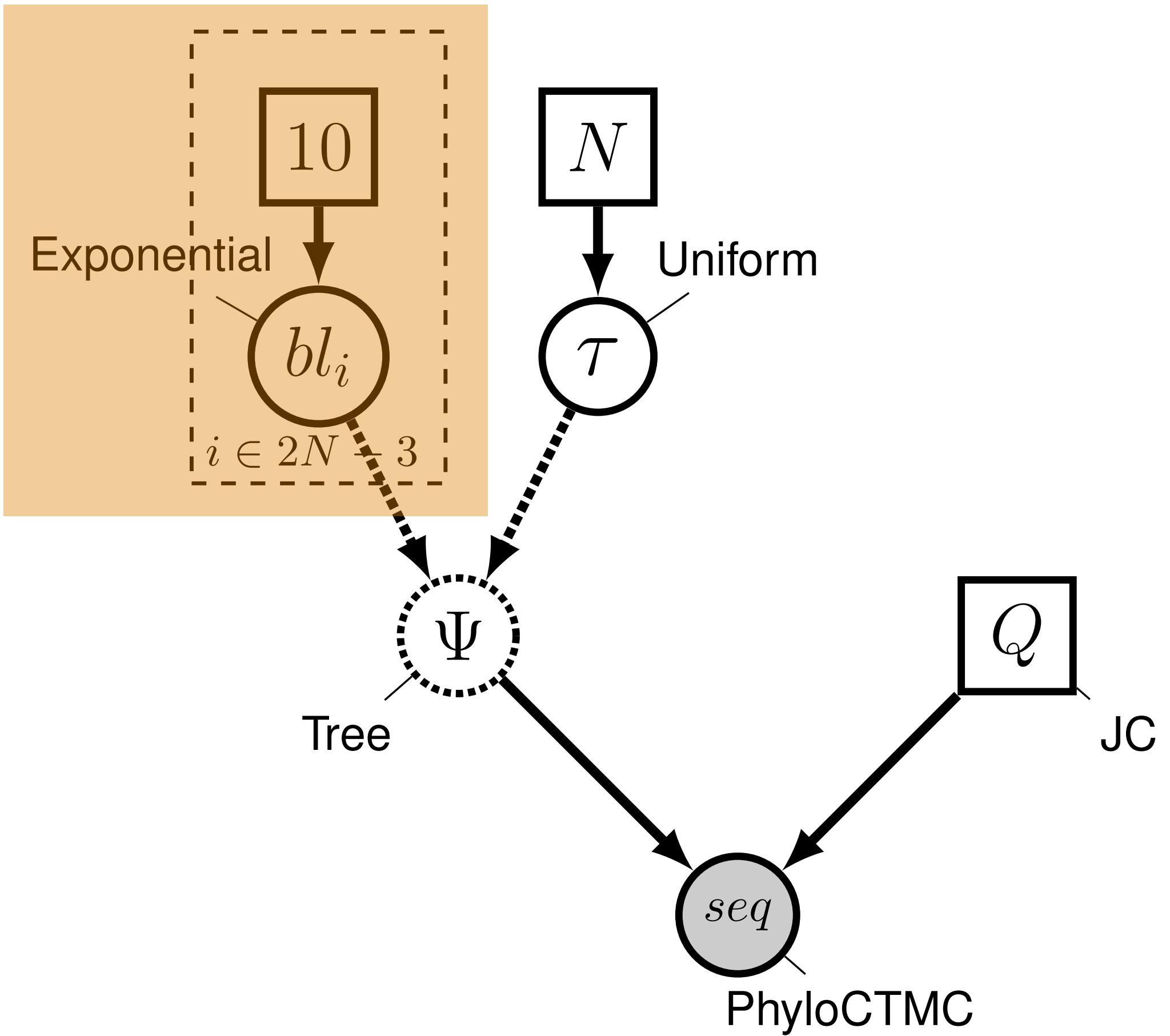
Graphical model example (DNA)



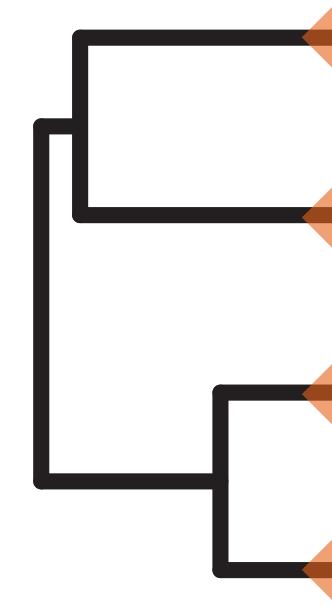


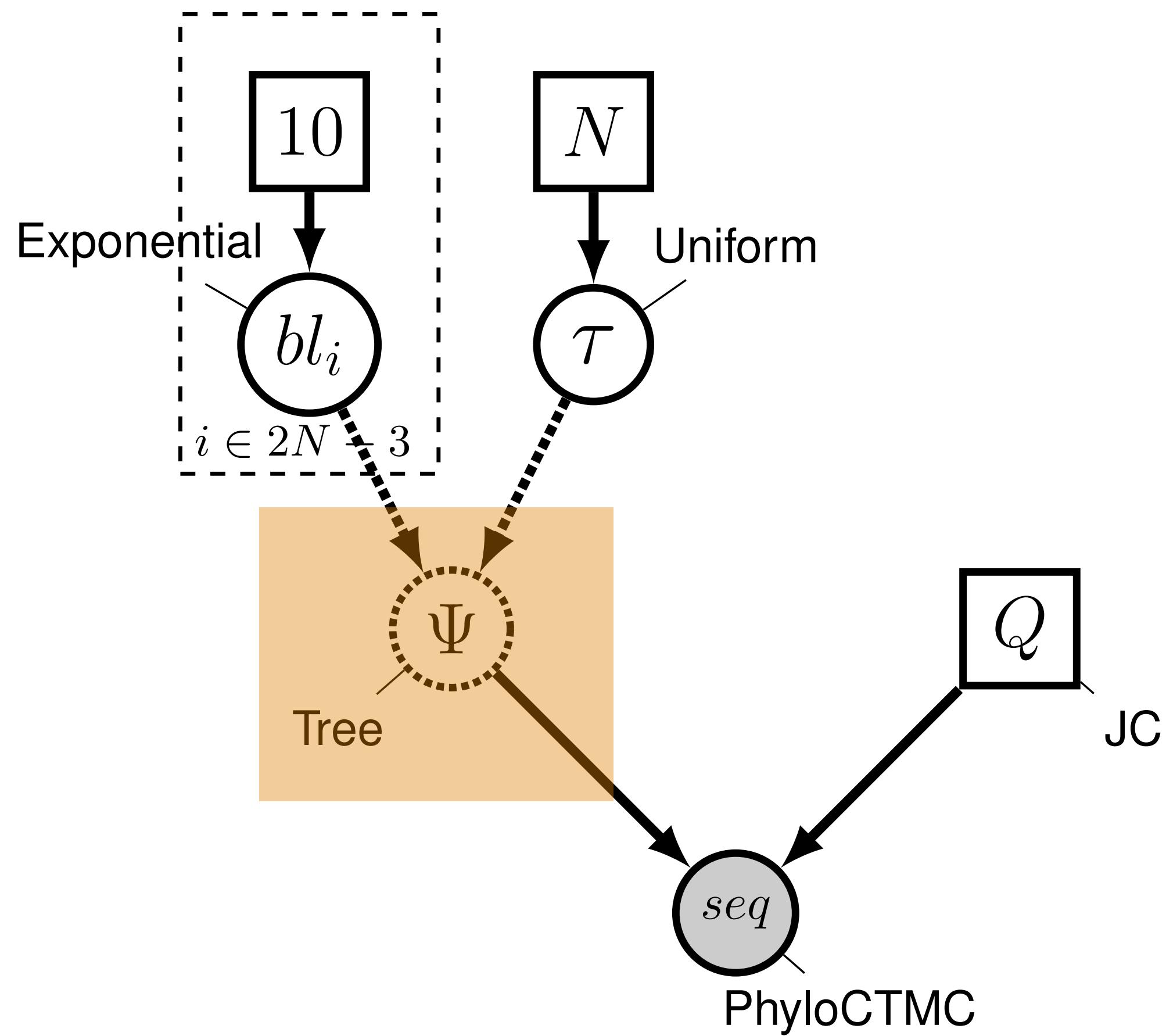
prior on the tree
topology



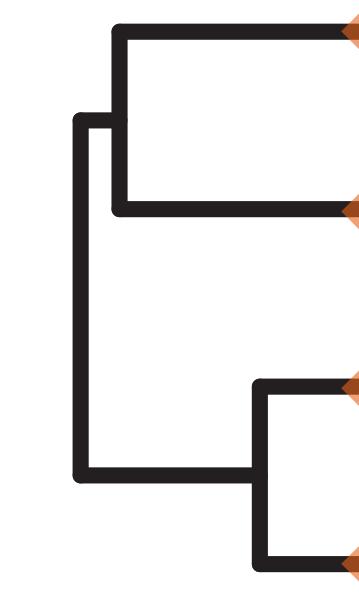


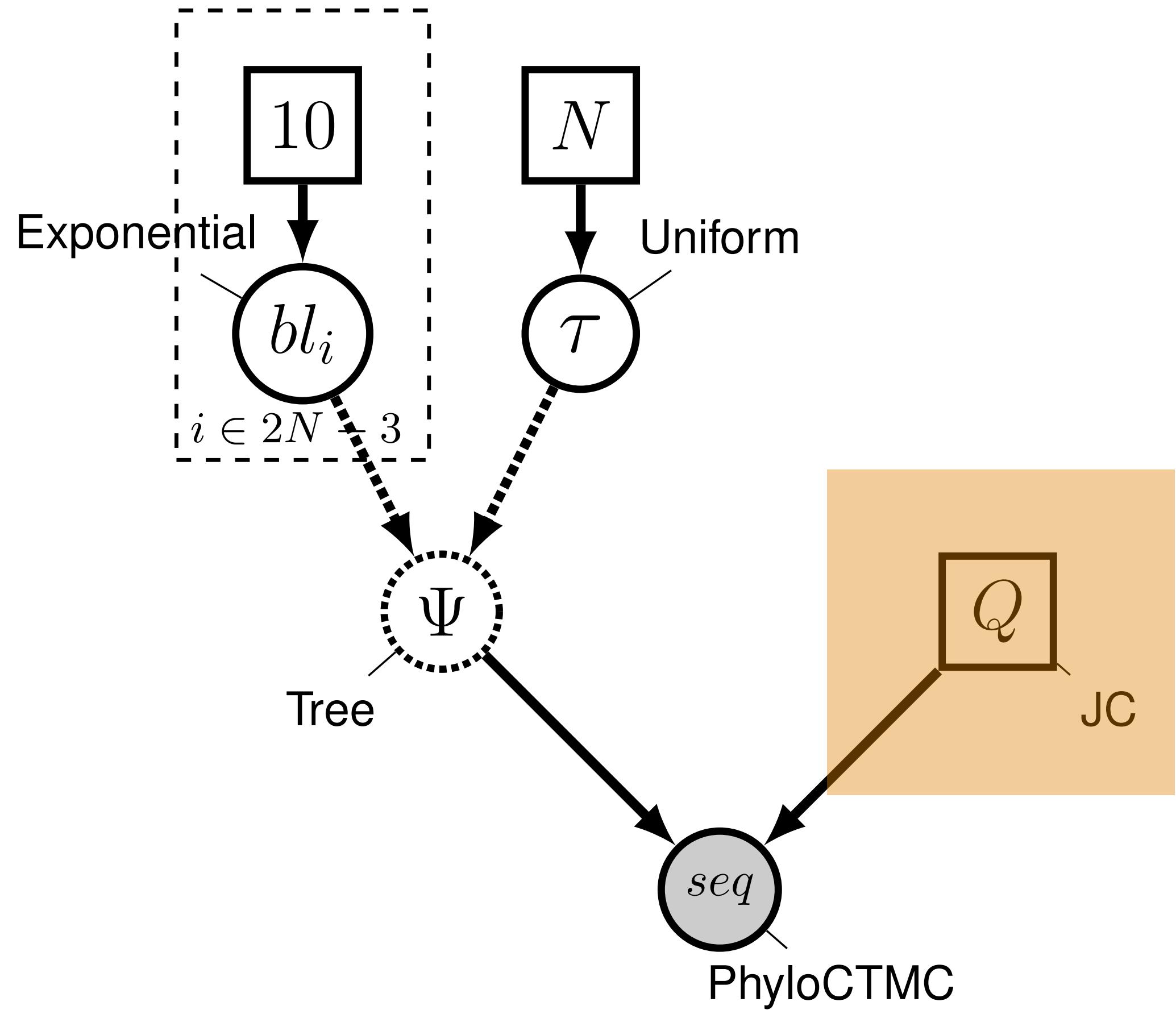
prior on the
branch lengths



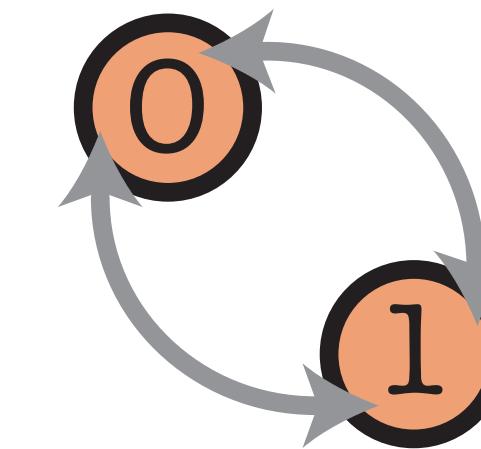


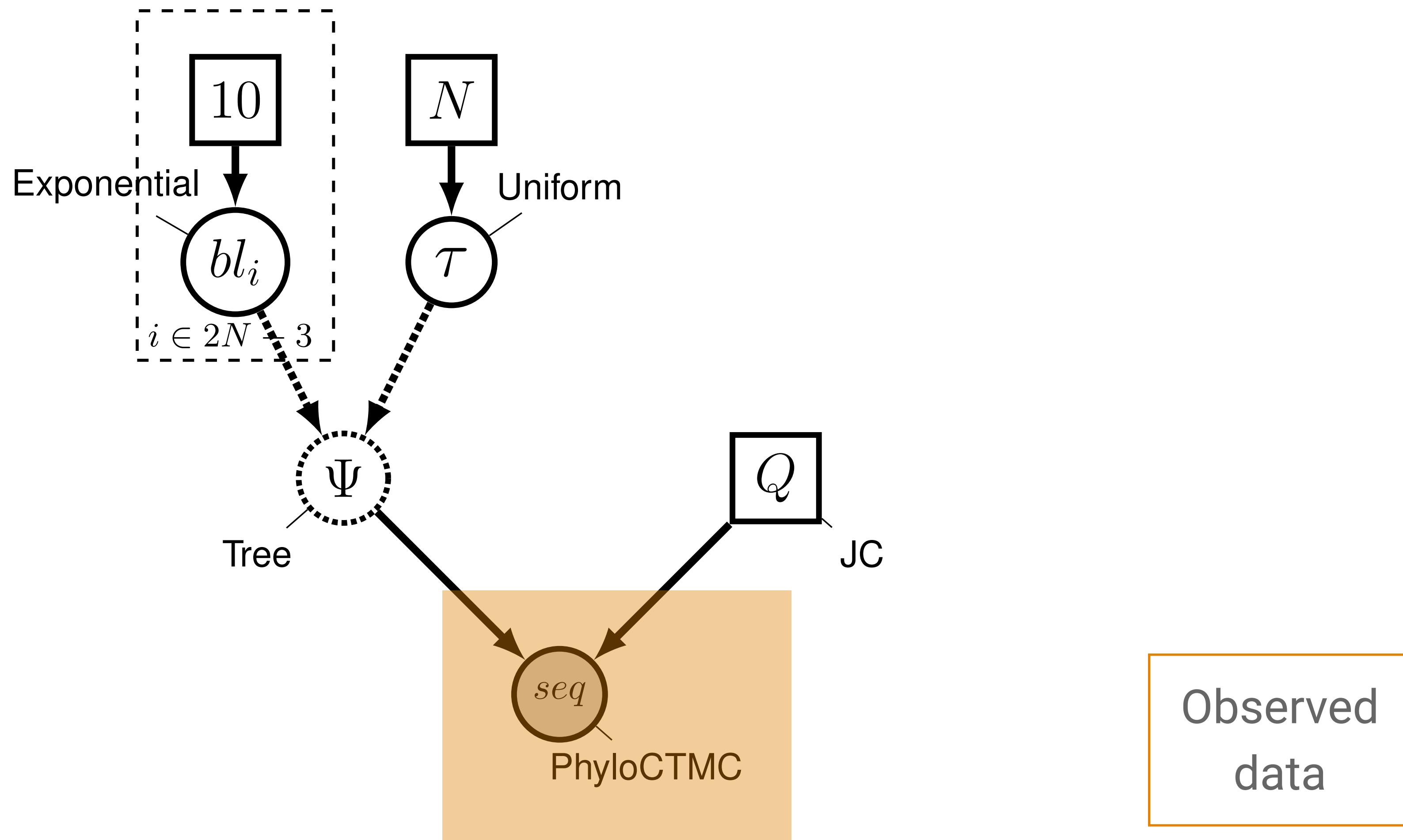
we can combine
the topology and
branch lengths



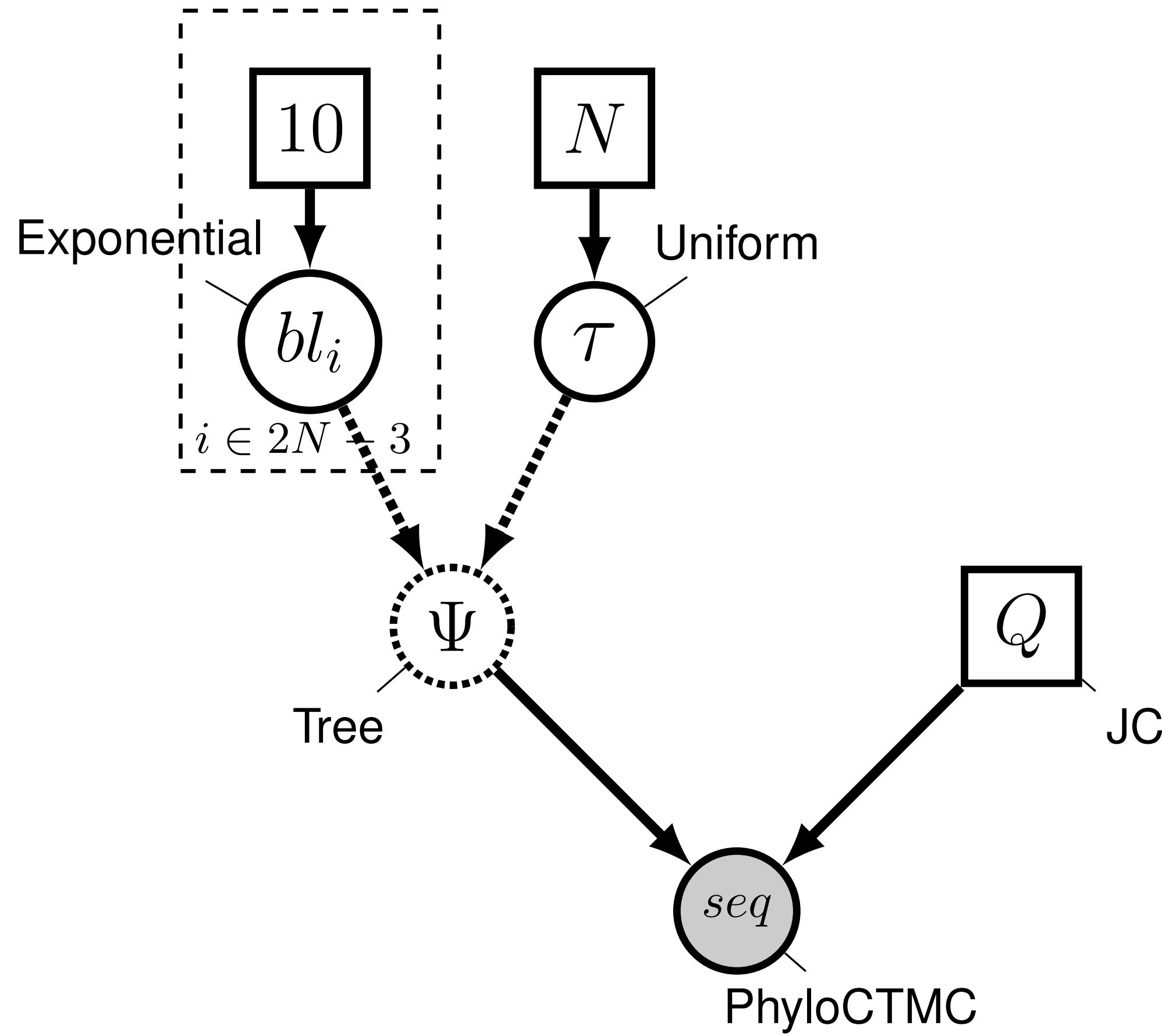


Substitution
model





0101...
1101...
0100...



```

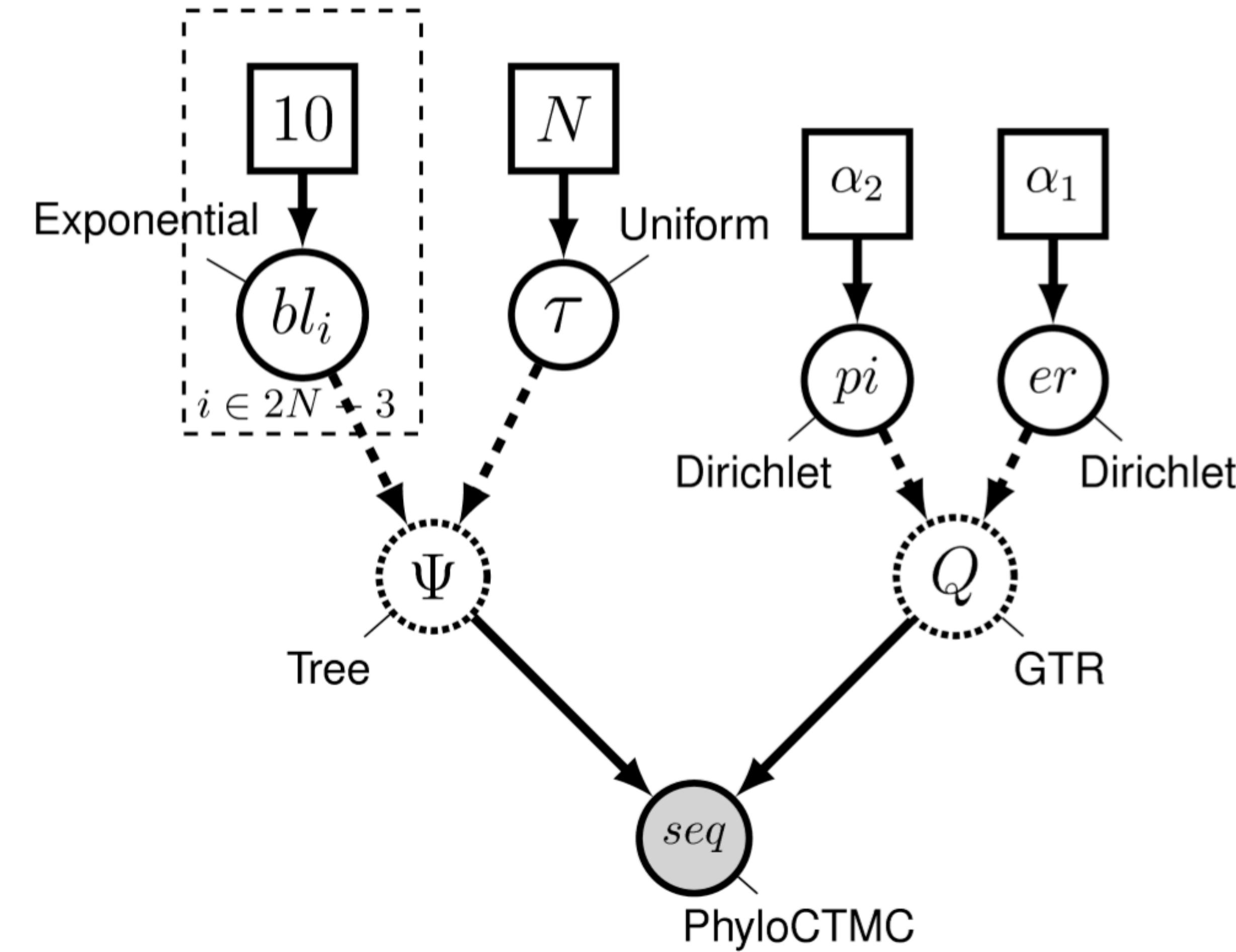
for (i in 1:n_branches) {
  bl[i] ~ dnExponential(10.0)
}
topology ~ dnUniformTopology(taxa)
psi := treeAssembly(topology, bl)

Q <- fnJC(4)

seq ~ dnPhyloCTMC( tree=psi, Q=Q, type="DNA" )

seq.clamp( data )

```



```

for (i in 1:n_branches) {
    bl[i] ~ dnExponential(10.0)
}
topology ~ dnUniformTopology(taxa)
psi := treeAssembly(topology, bl)

alpha1 <- v(1,1,1,1,1,1)
alpha2 <- v(1,1,1,1)
er ~ dnDirichlet(alpha1)
pi ~ dnDirichlet(alpha2)
Q := fnGTR(er, pi)

```