

Bayesian Phylogenetics

Analytical Paleobiology

Rachel Warnock

13.08.24



Objectives for today and tomorrow

Intro to **Bayesian phylogenetics** in paleobiology

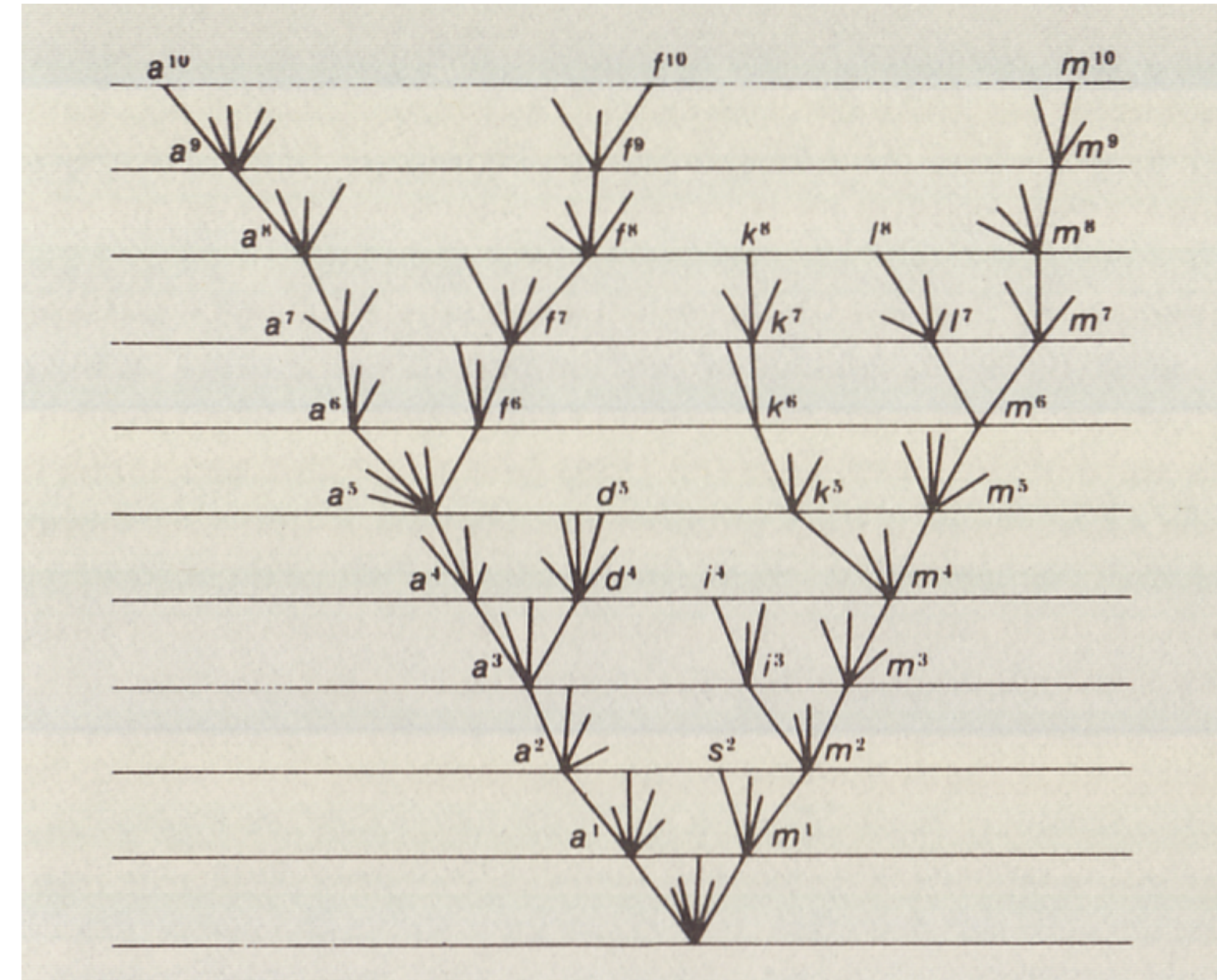
- Tree building
- Substitution models
- Dating trees
- Clock models
- Tree models
- Diversification rates
- Morphological models



Please ask questions!

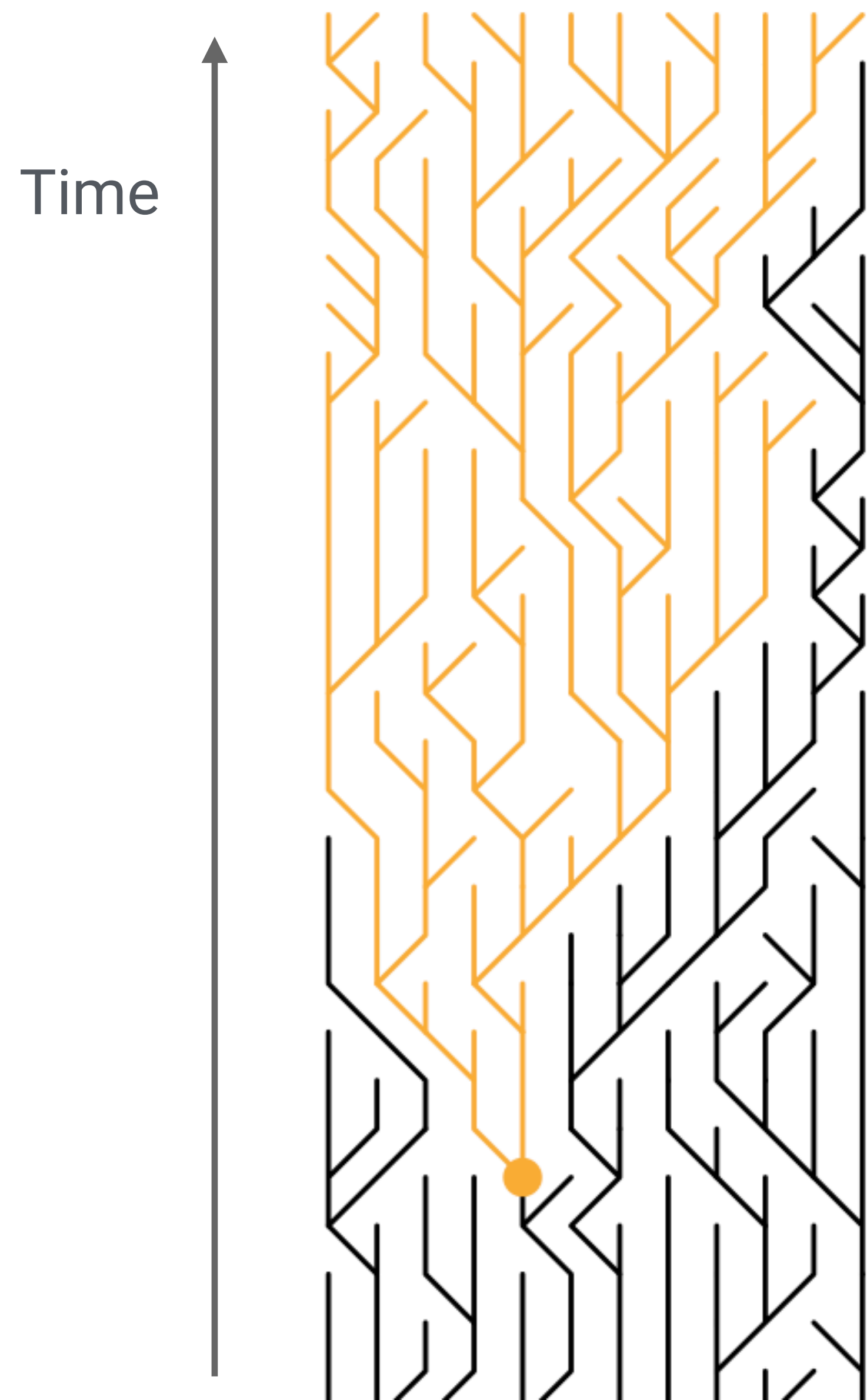
Today's objectives

- Intro to RevBayes
- Bayesian tree inference
- Morphological models



Time tree from Darwin's *Origin of Species*

What is phylogenetics?



- populations
- species
- viruses
- cells
- languages

Data

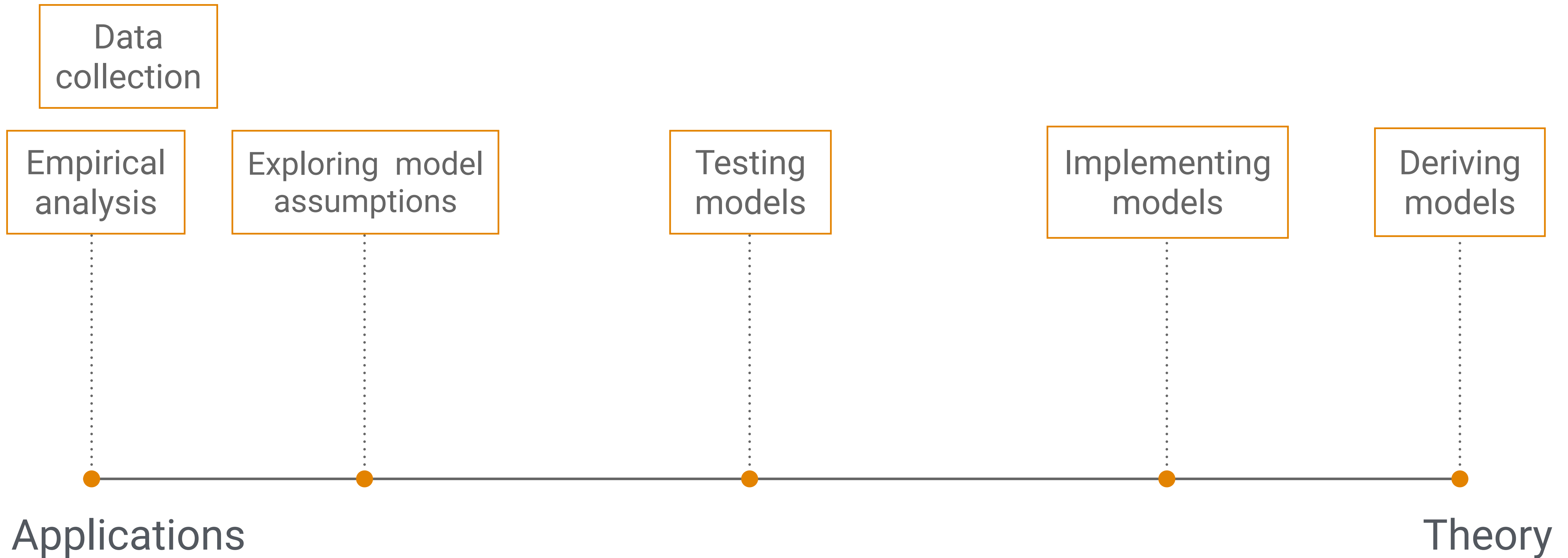
- DNA
- morphology
- words

In this course we mainly focus on trees that include **one representative per species**



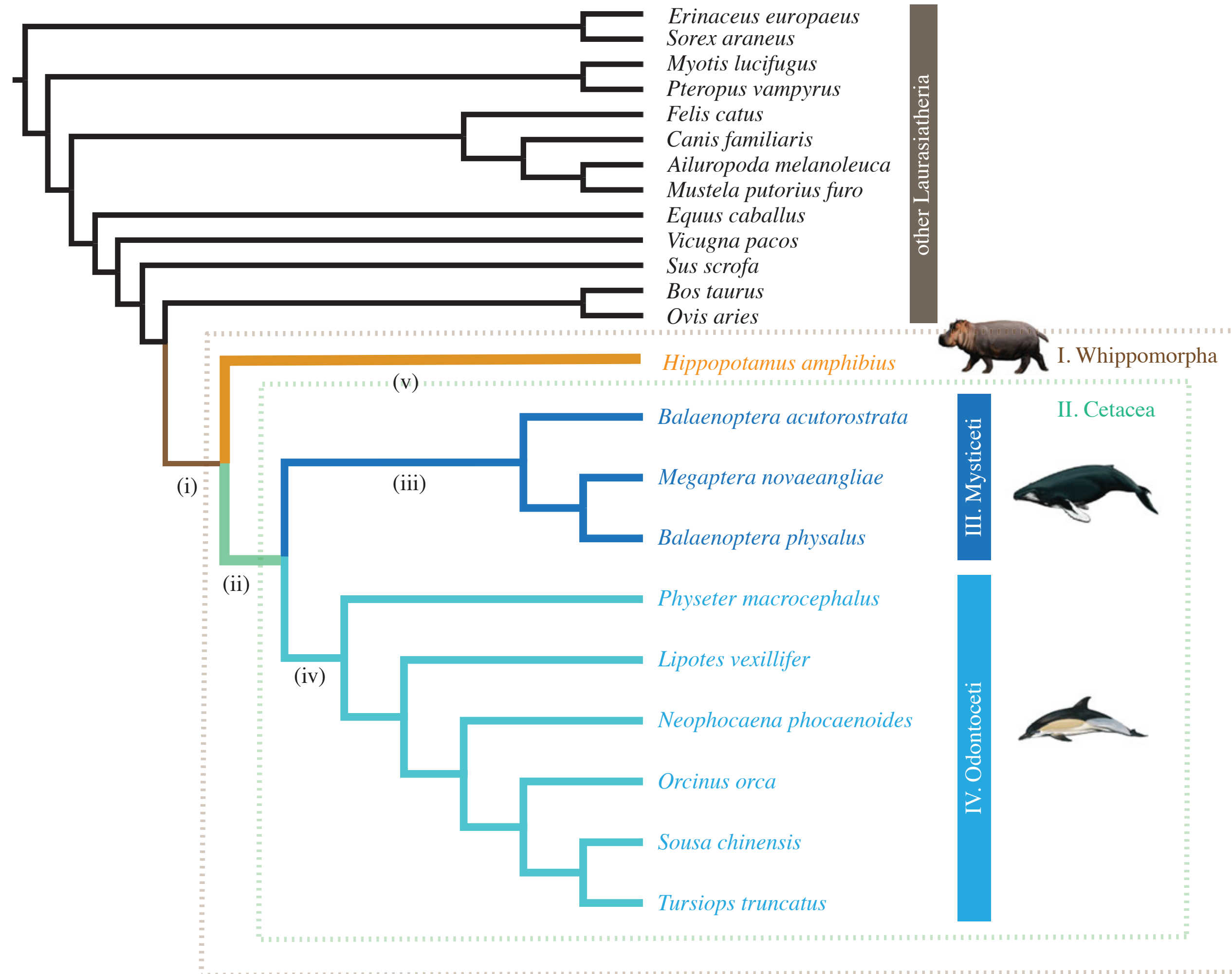
[Scots poem](#) - also the [BEAST2](#) logo!

Research topics in phylogenetics



Trees in paleobiology

What can we learn from trees?

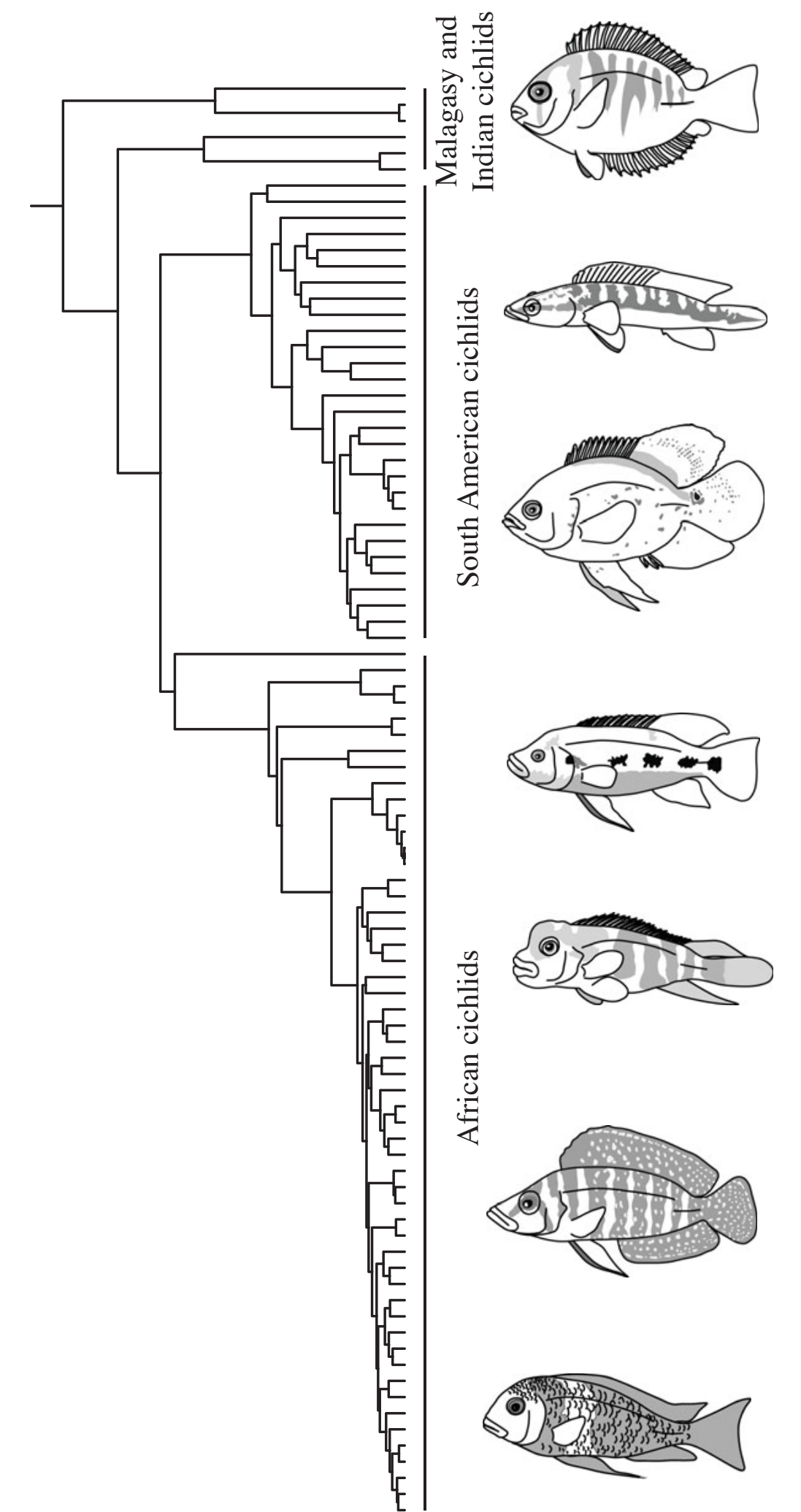


How are our favourite species related?

Does the phylogeny support the taxonomy?

What can we learn from trees?

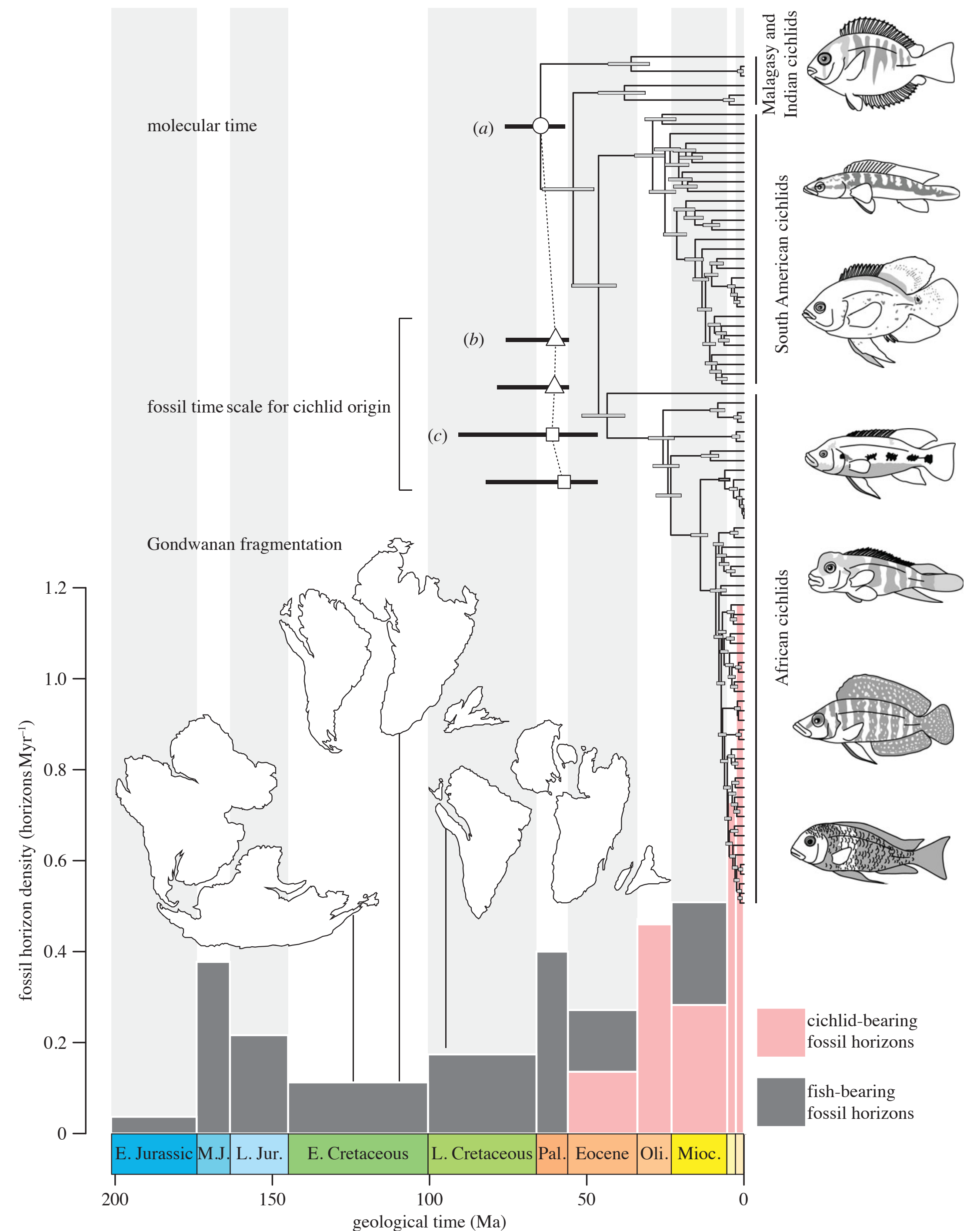
- Evolutionary relationships



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](#))



Phylogenetics

Phylogenetics aims to reconstruct the phylogeny of individual samples based on molecular or morphological character data

A phylogeny captures part of evolutionary history that is otherwise not directly observable

Phylodynamics aims to quantify the processes that gave rise to the tree, e.g., speciation, extinction

What do we mean by model?

(the following is my take on things – intended to be useful but not definitive)

What is a **statistical model**? When is an equation a model?

What is a **mechanistic model**?

What is the difference between an **algorithm** and a model?

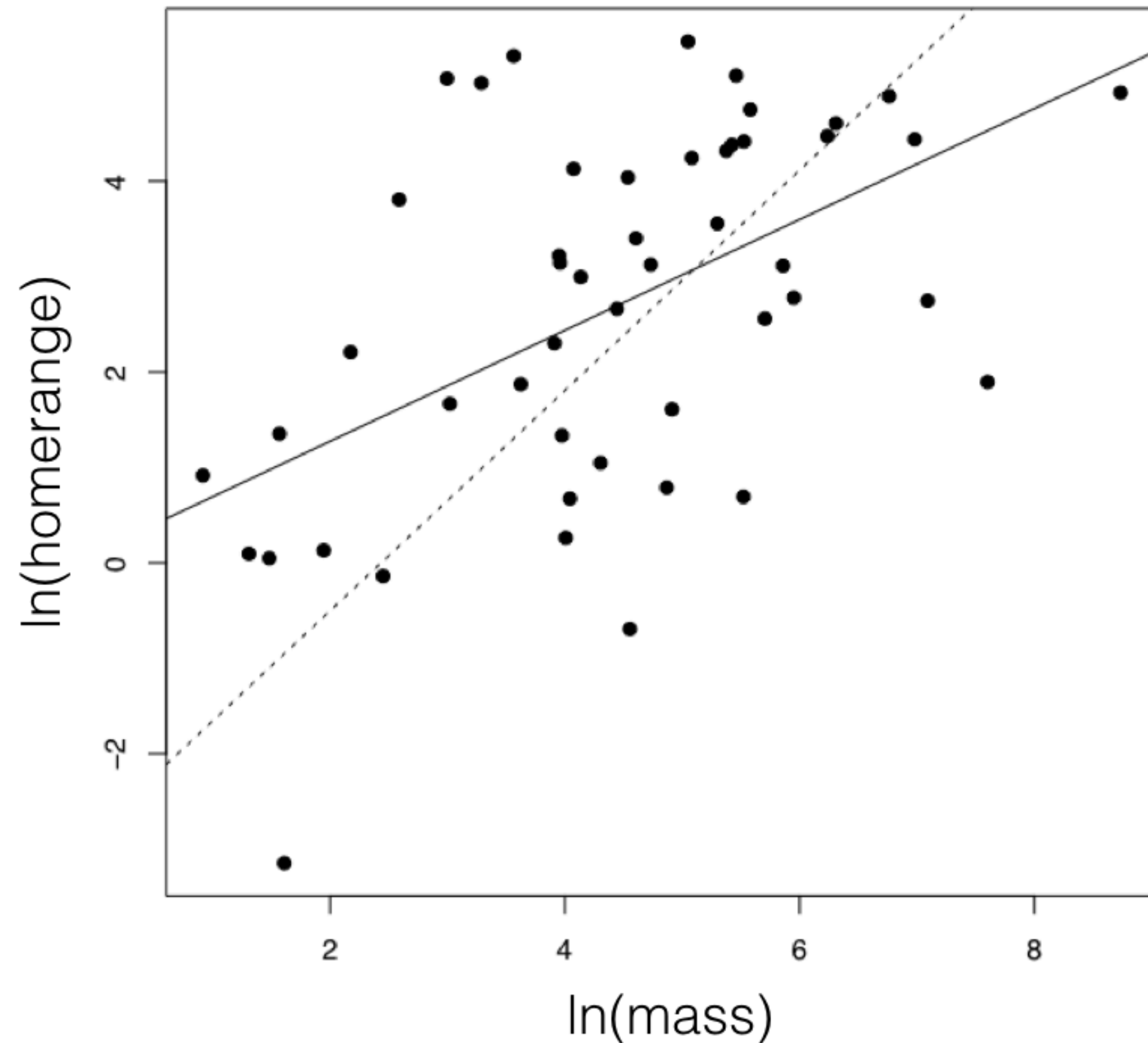
A statistical model is a type of model that includes a set of assumptions about the data-generating process

It should be possible to **simulate data** under the assumptions of the model

If we're lucky, we might *also* be able to **estimate parameters** under the model*. This isn't always possible because some models are too complex

*A fancy way of saying this is, "we can perform **inference** under the model"

An example



The solid black line is a linear regression line

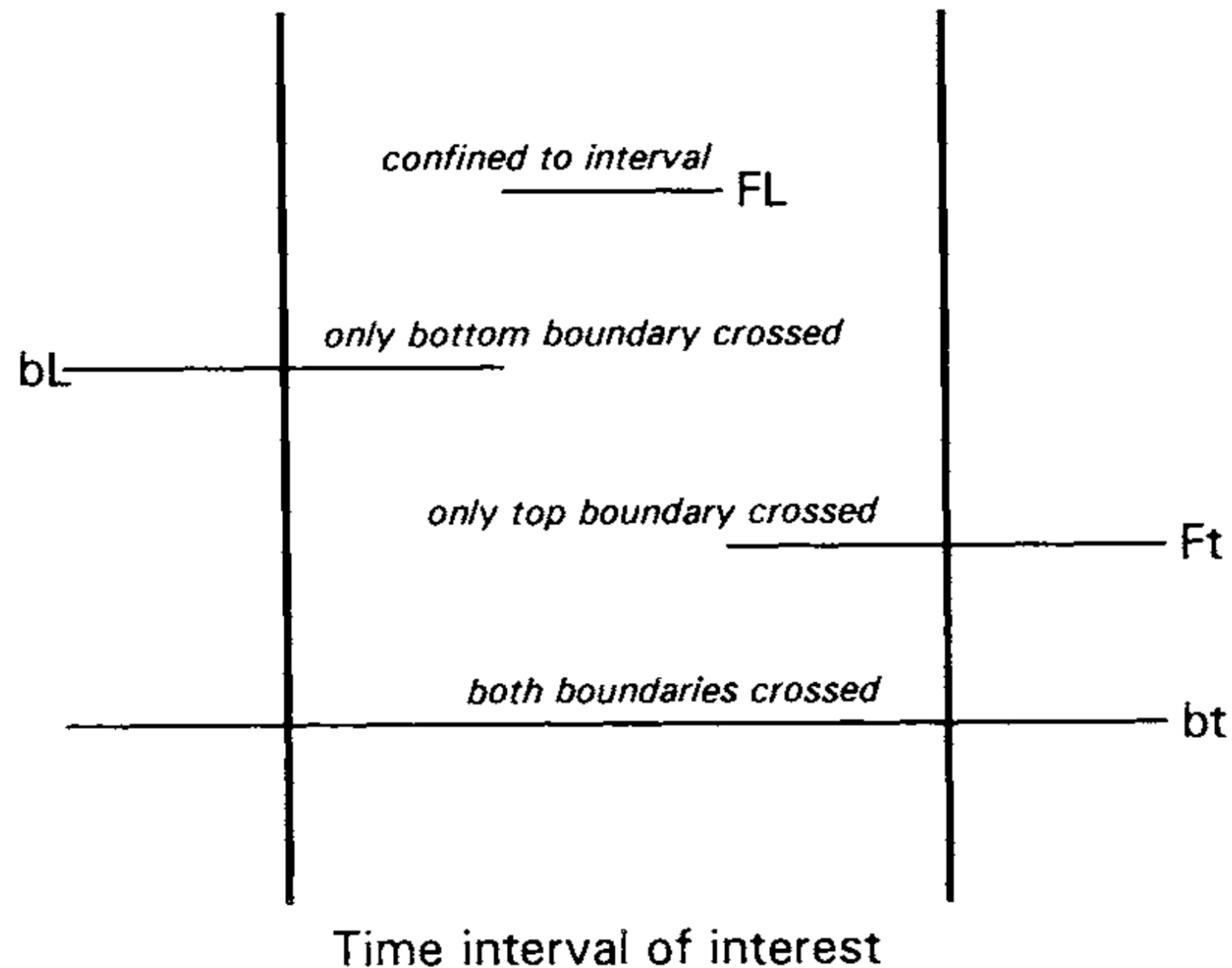
We can estimate the parameters of the regression model

$$y = X\beta + \varepsilon$$

It's also straightforward to simulate data under this model

Non model-based approaches are still useful

Four fundamental classes of taxa



The boundary-crosser and three-timer metrics are *not* models

They provide a clever way of approximating origination and extinction rates (and often perform well), but don't describe the data generating processes

Foote (2000)

Mechanistic or process based models are based on 'physical principles'. They describe the data as a function of a set of parameters that have a tangible biological or geological meaning

A regression model is not mechanistic – it describes the relationship between x and y but the parameters don't have a biological meaning

Many models used in phylogenetics are mechanistic, e.g., they might include parameters for origination, extinction, or sampling

An [algorithm](#) is a precise rule (or set of rules) specifying how to solve some problem

```
i = 1
while i < 11:
    print(i)
    i = i + 1
```

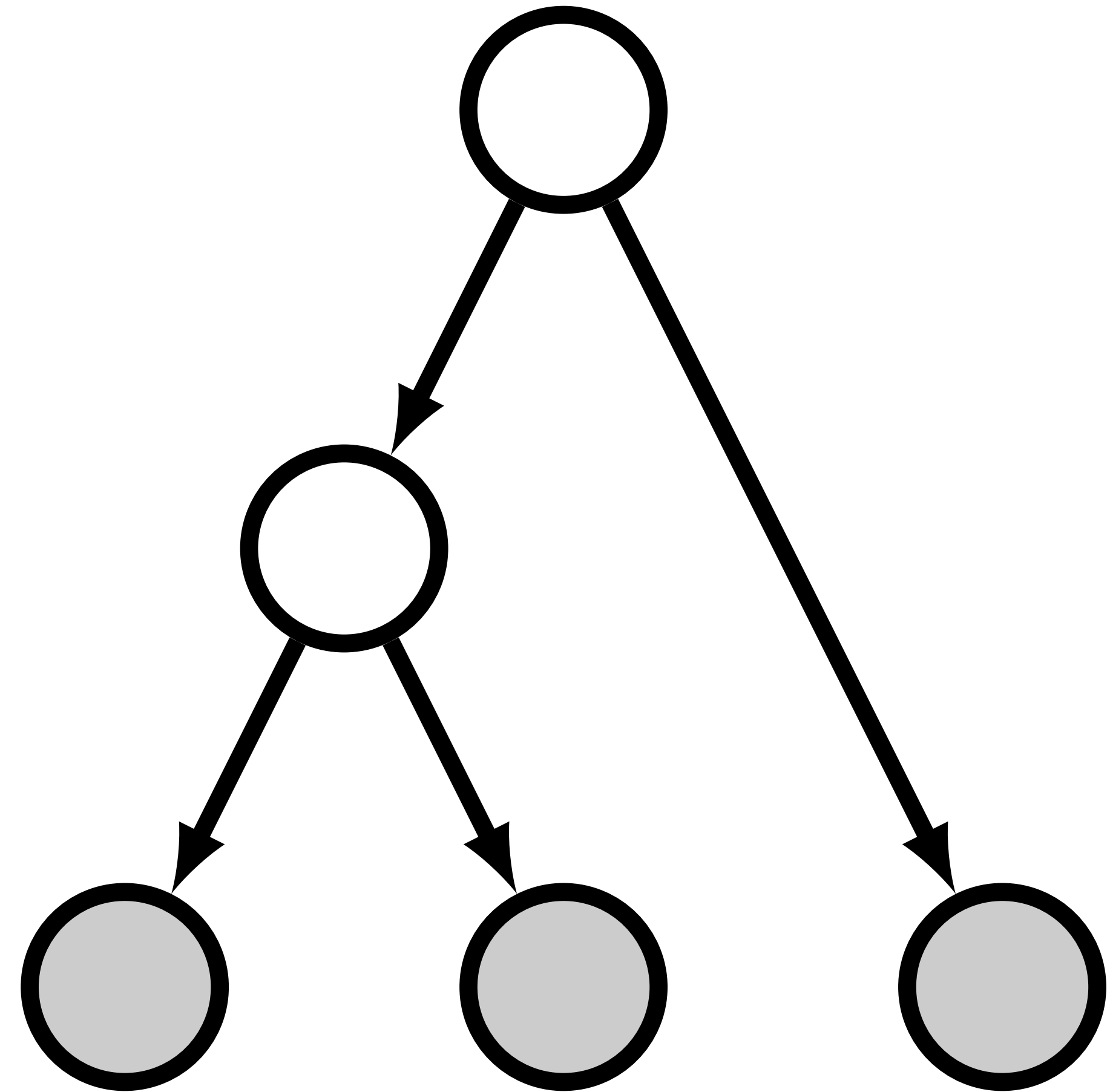
```
for i in range(1,11):
    print(i)
```

Used in phylogenetics for all sorts of tasks, e.g., traversing tree space

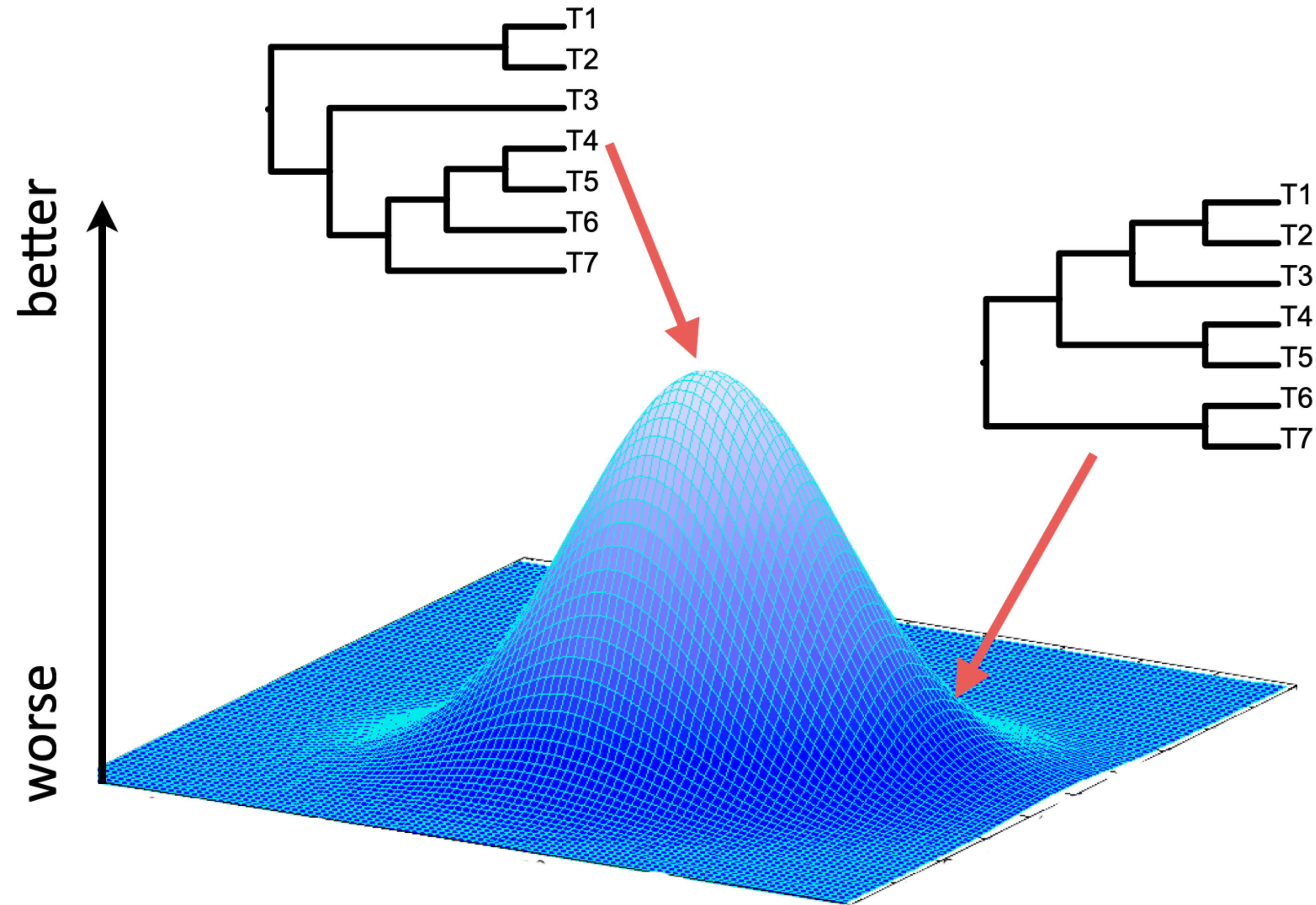
Mini reading group

Next

- graphical models
- RevBayes
- Bayesian inference
- MCMC



How do we find the 'best' tree?



It depends how you measure 'best'

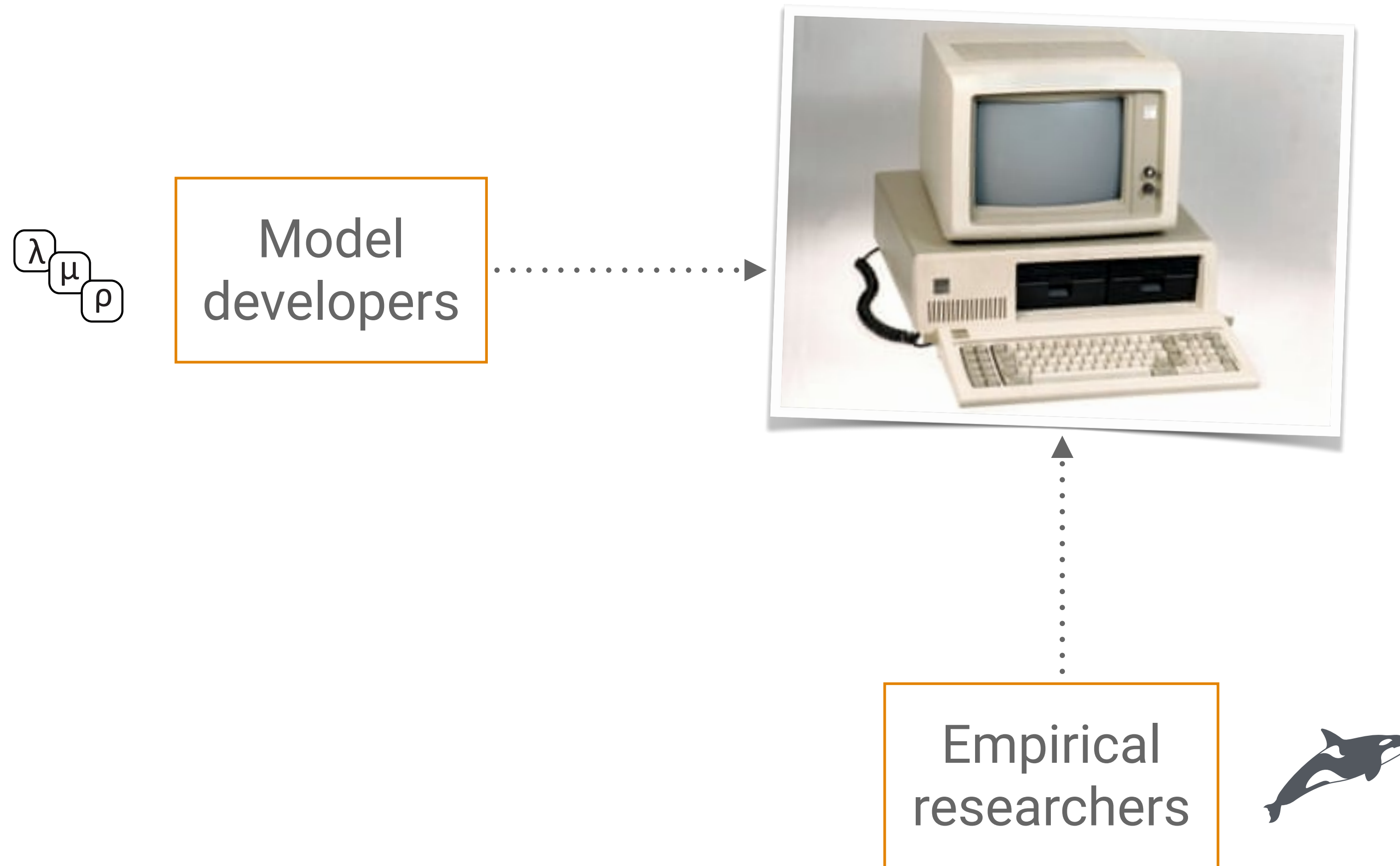
Method	Criterion (tree score)
Maximum parsimony	Minimum number of changes
Maximum likelihood	Likelihood score (probability), optimised over branch lengths and model parameters
Bayesian inference	Posterior probability, integrating over branch lengths and model parameters

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

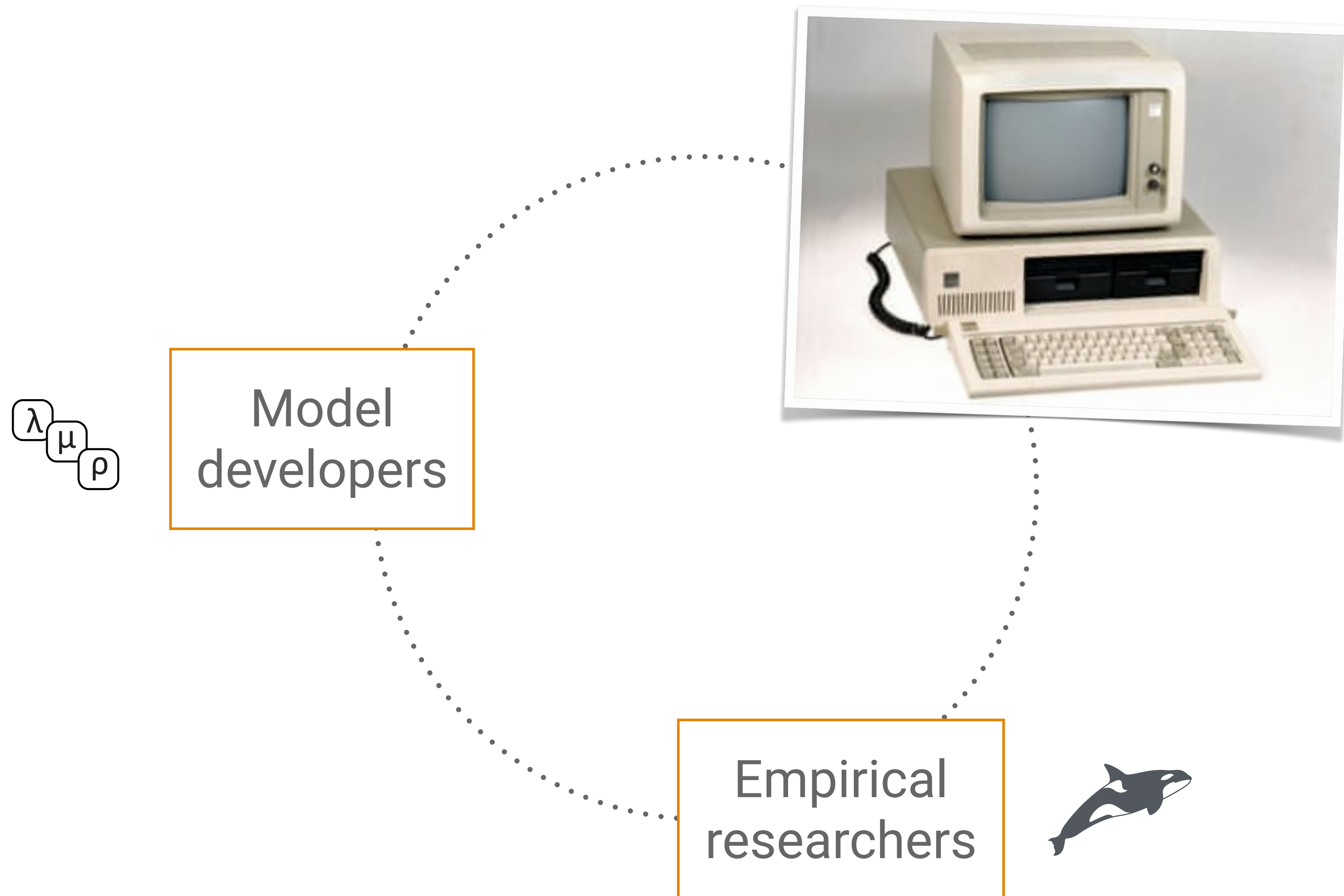
RevBayes

Phylogenetic inference – the old way

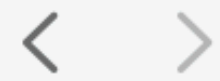


What we might call a
“**black box**” approach

Phylogenetic inference – a better way?



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



revbayes.github.io



Download

Tutorials

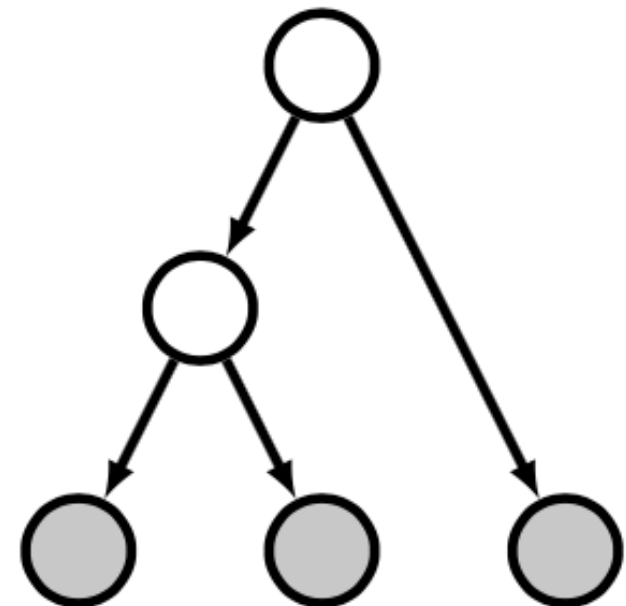
Documentation

Interfaces

Workshops

Jobs

Developer



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.

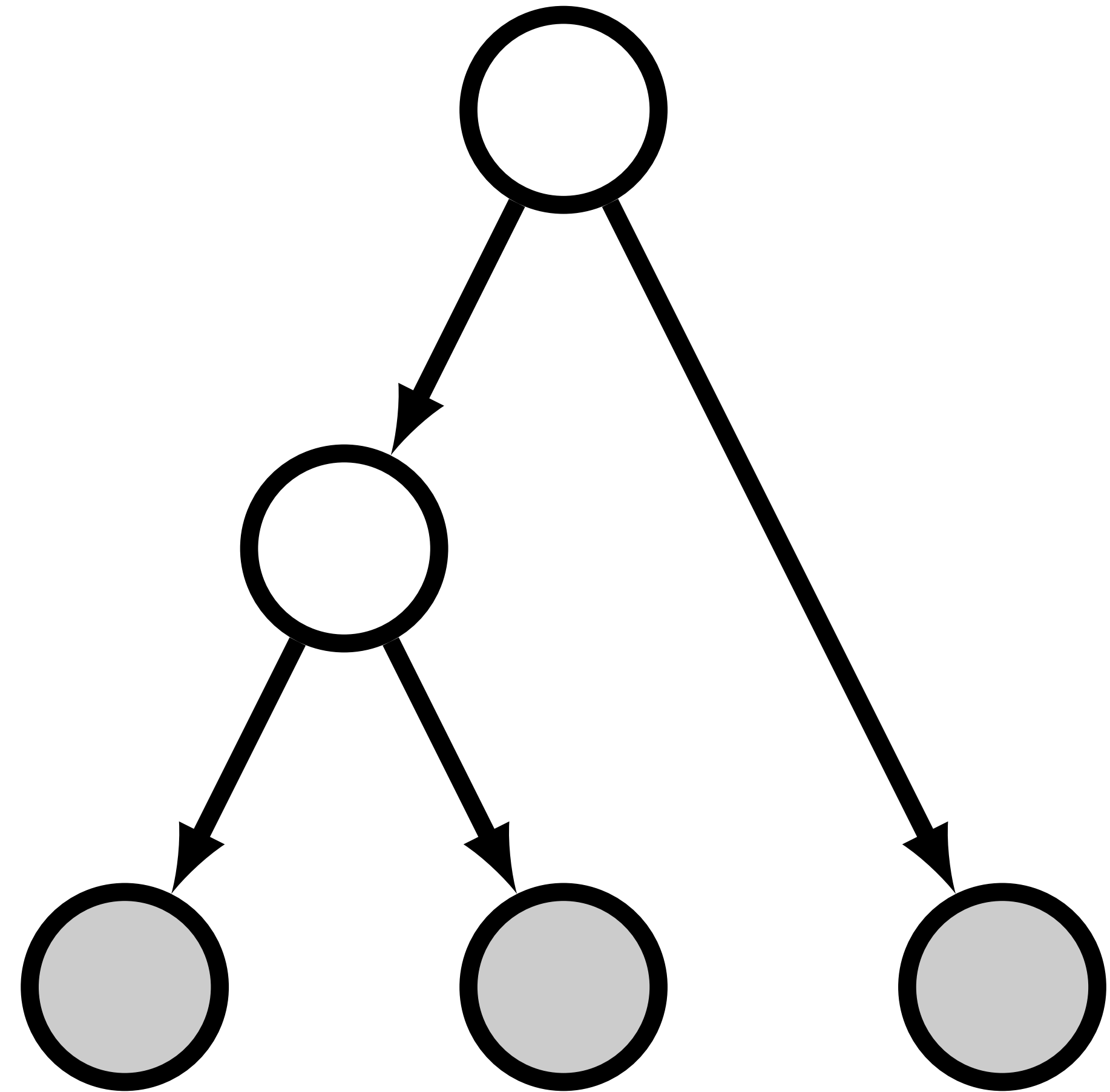
[GitHub](#) | [License](#) | [Citation](#) | [Users Forum](#)

Graphical models

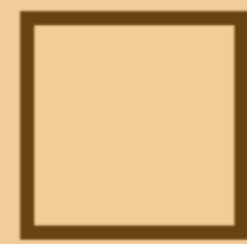
Graphical models

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

a. fixed value variables



b) Stochastic node

b. random variables that depend on other variables



c) Deterministic node

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



d) Clamped node
(observed)

d. observed stochastic variables (data)



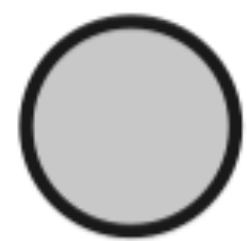
a) Constant node



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.

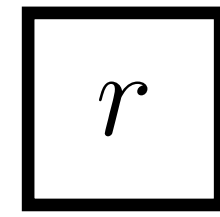
Operator	Variable
<code><-</code>	constant variable
<code>~</code>	stochastic variable
<code>:=</code>	deterministic variable
<code>node.clamp(data)</code>	clamped variable
<code>=</code>	inference (<i>i.e.</i> , non-model) variable
<code>for(i in 1:N){...}</code>	plate

a)

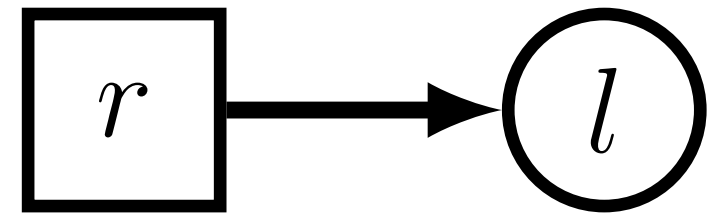
r

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

a)

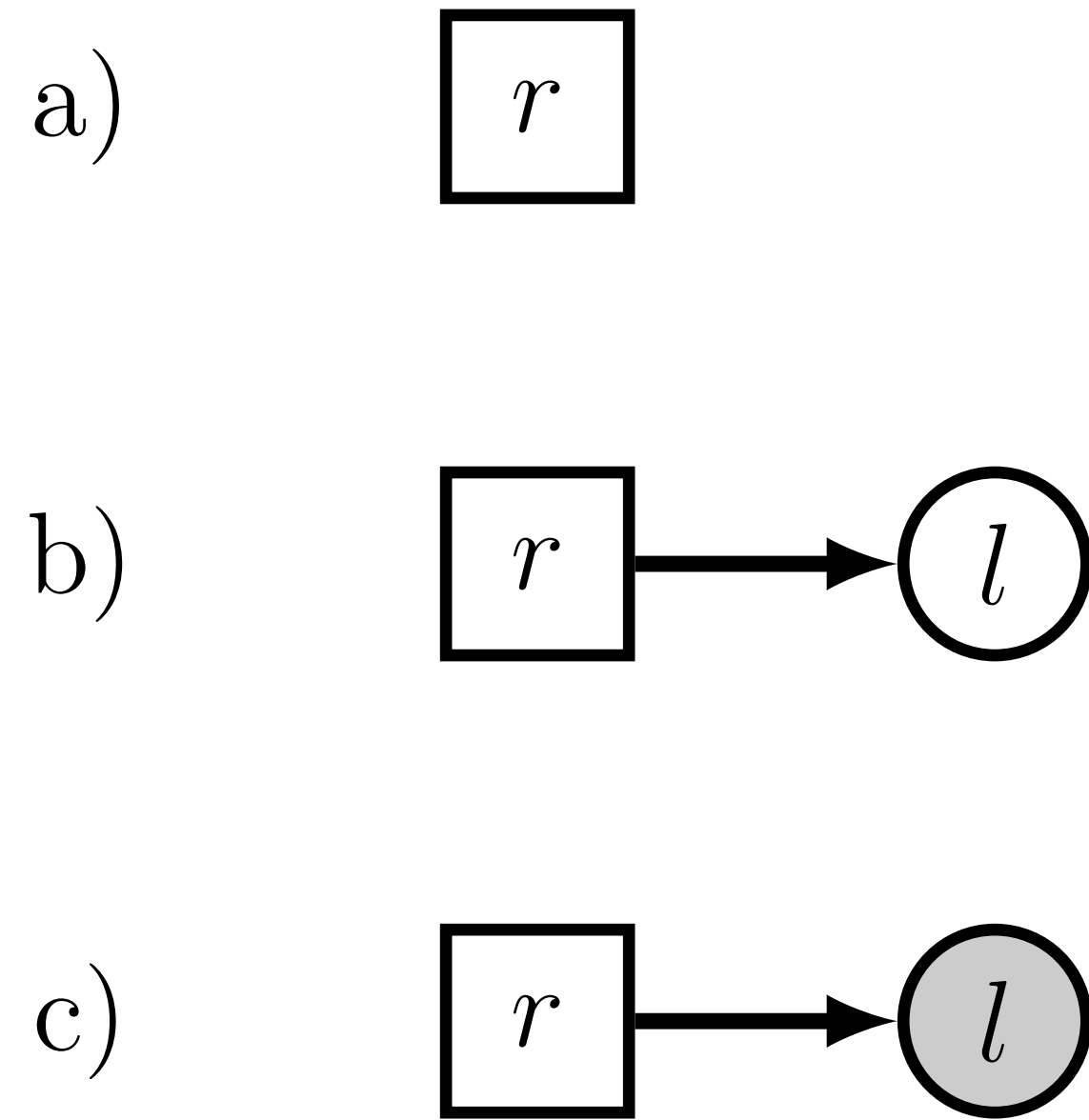


b)



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

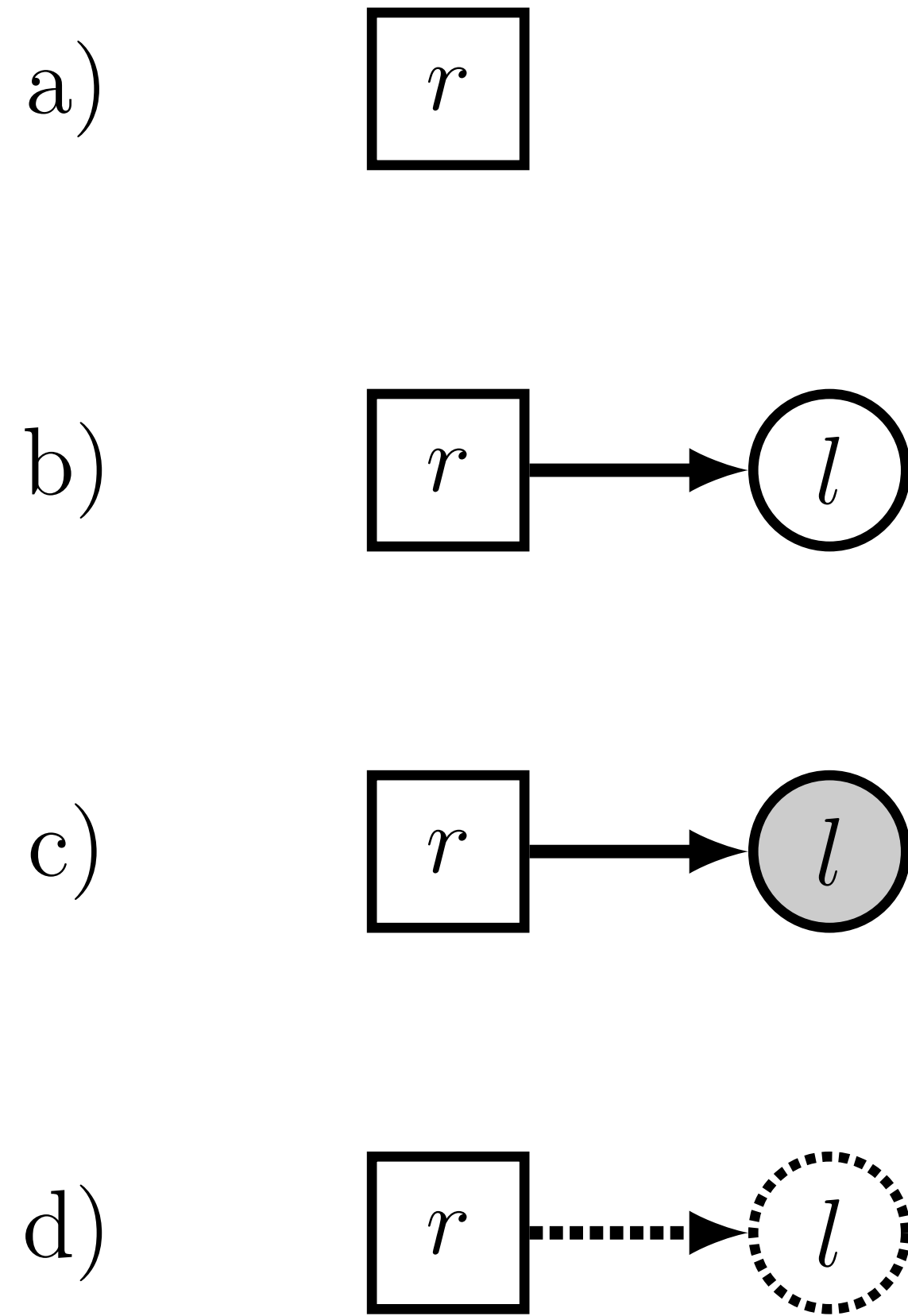
```
# stochastic node  
l ~ dnExp(r)
```



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

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

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

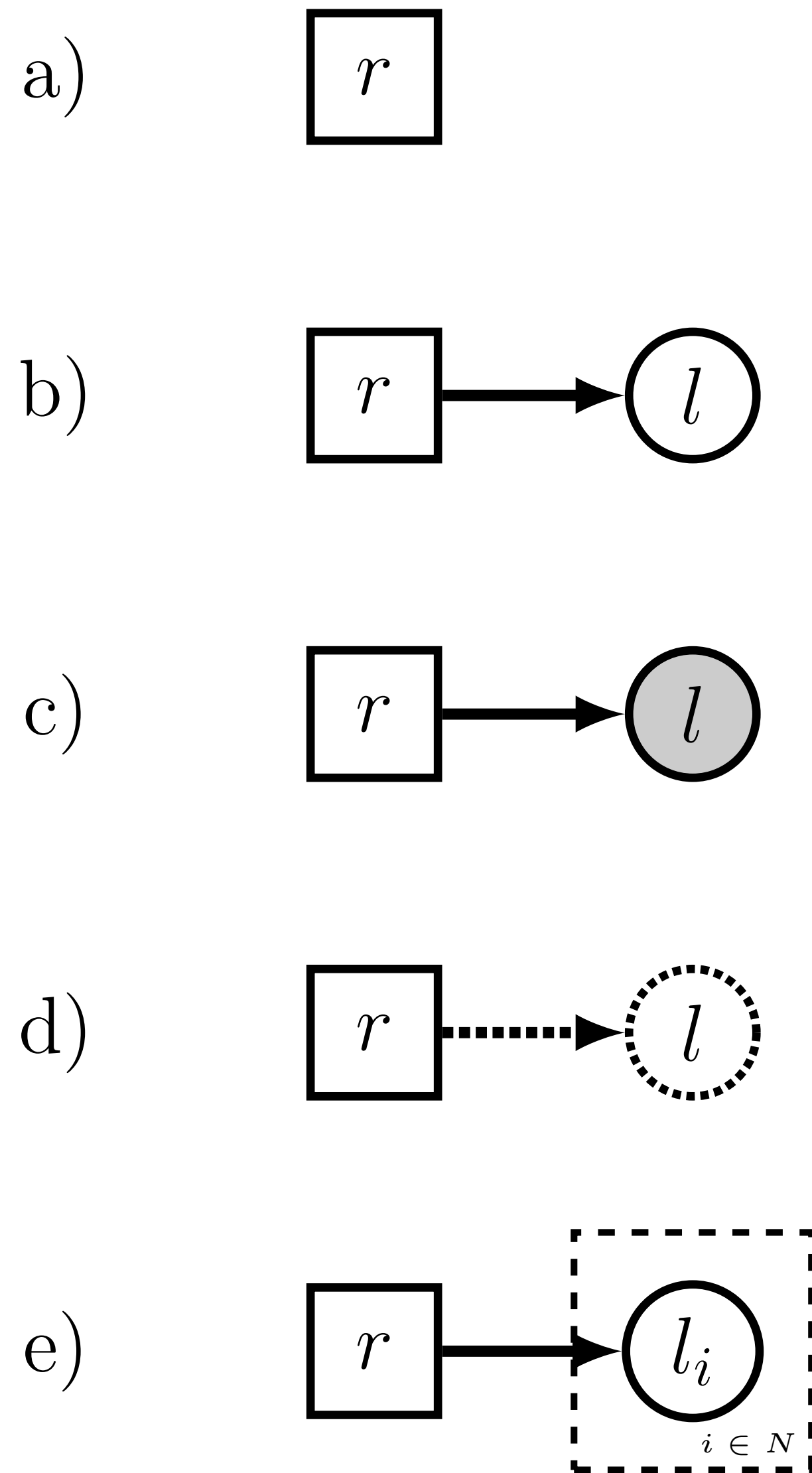


```
# constant node  
r <- 10
```

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

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

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



```
# 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)
}
```

Exercise

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

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

$\Pr(\text{data})$

Marginal probability

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

Bayes' theorem

Reflects our combined knowledge based on the likelihood and the priors

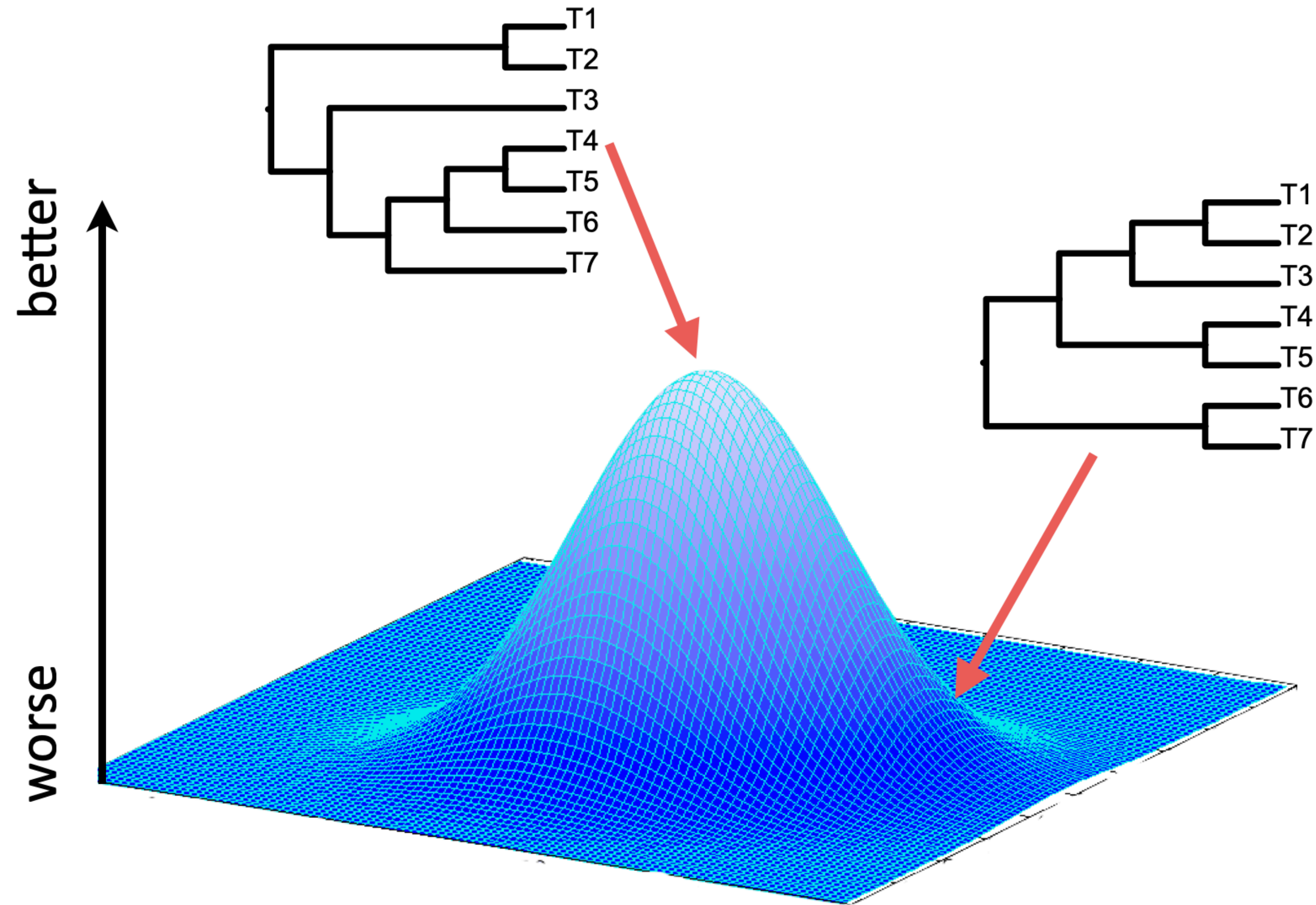
posterior

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

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

$\Pr(\text{data})$

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

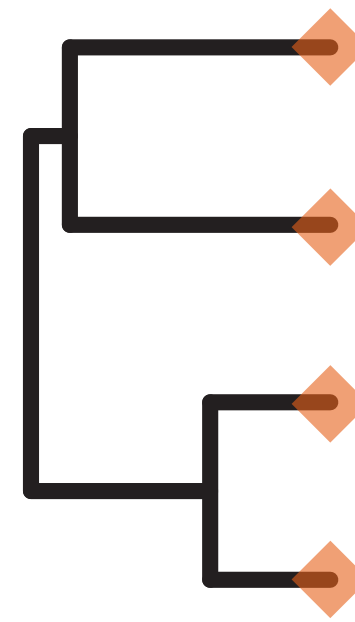
Note these are not the only approaches to tree-building but they are the most widely used

Components used to infer trees

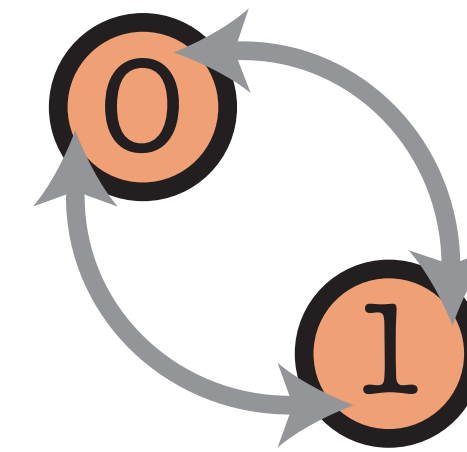
without considering time

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

data
sequences or
characters



tree
topology and
branch lengths



substitution
model

Bayesian tree inference

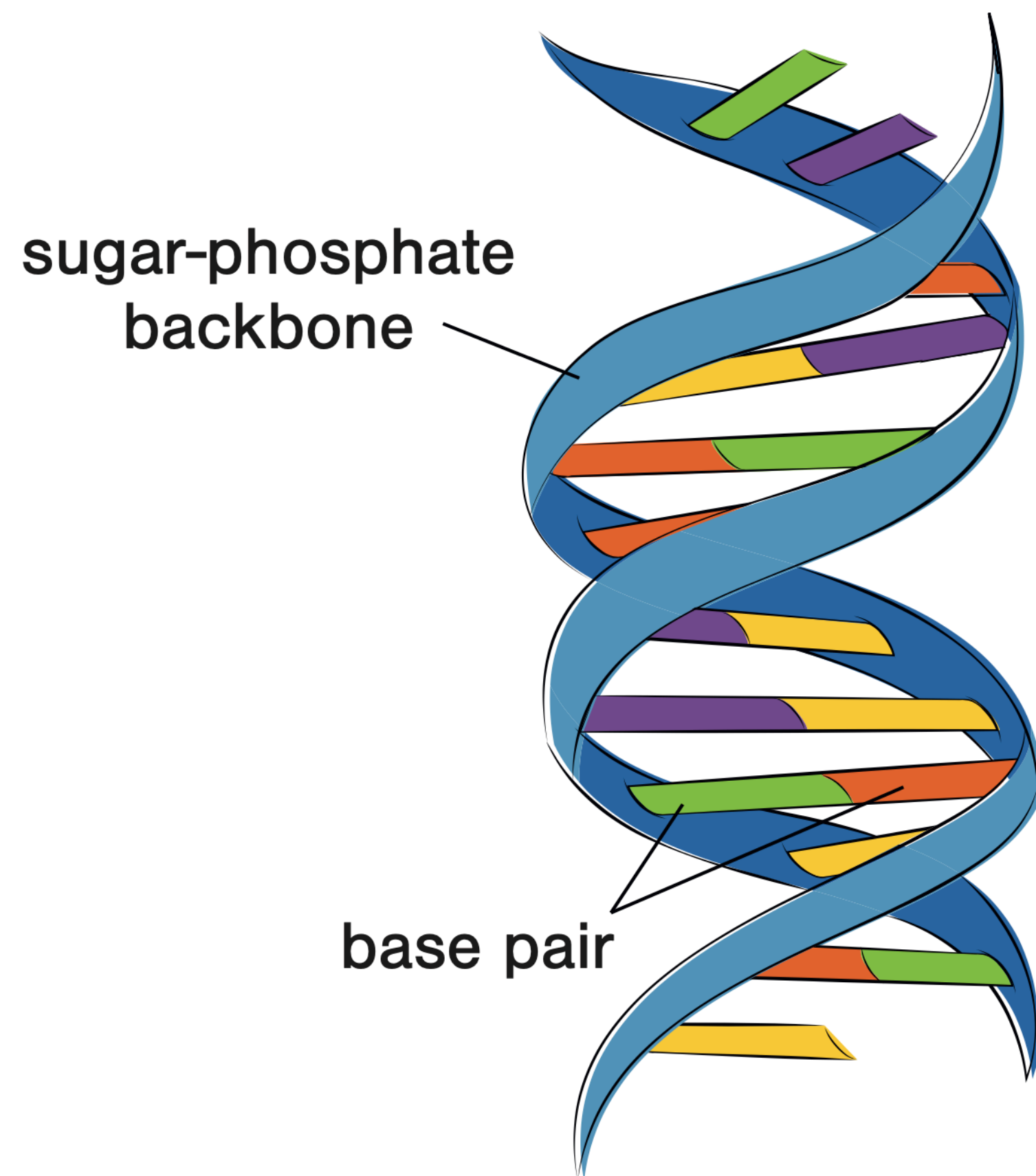
$$\begin{array}{c} \text{posterior} \\ \boxed{\phantom{\text{posterior}}} \\ P(\text{tree} \mid \text{data}) \end{array} = \frac{\begin{array}{c} \text{likelihood} \\ \boxed{\phantom{\text{likelihood}}} \\ P(\text{data} \mid \text{tree}) \end{array} \begin{array}{c} \text{priors} \\ \boxed{\phantom{\text{priors}}} \\ P(\text{tree}) \end{array}}{\begin{array}{c} \text{marginal probability} \\ \boxed{\phantom{\text{marginal probability}}} \\ P(\text{data}) \end{array}}$$

The diagram shows the Bayesian inference equation for a tree. The posterior probability $P(\text{tree} \mid \text{data})$ is equal to the product of the likelihood $P(\text{data} \mid \text{tree})$ and the prior probability $P(\text{tree})$, divided by the marginal probability $P(\text{data})$. Each term is accompanied by a small tree diagram with nodes labeled 0 and 1.

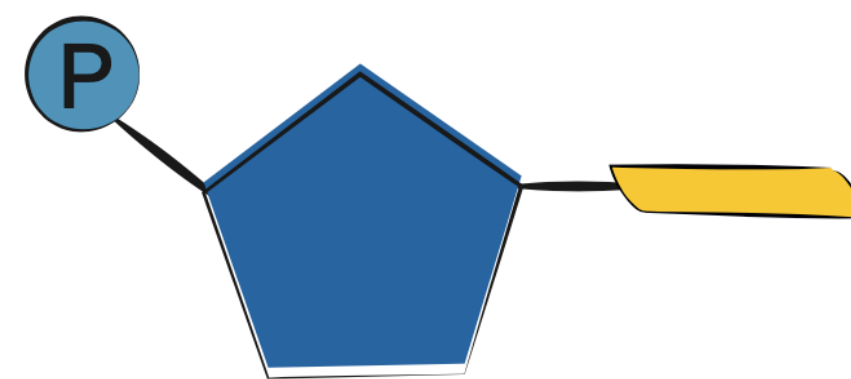
Likelihood and substitution models

Molecular evolution

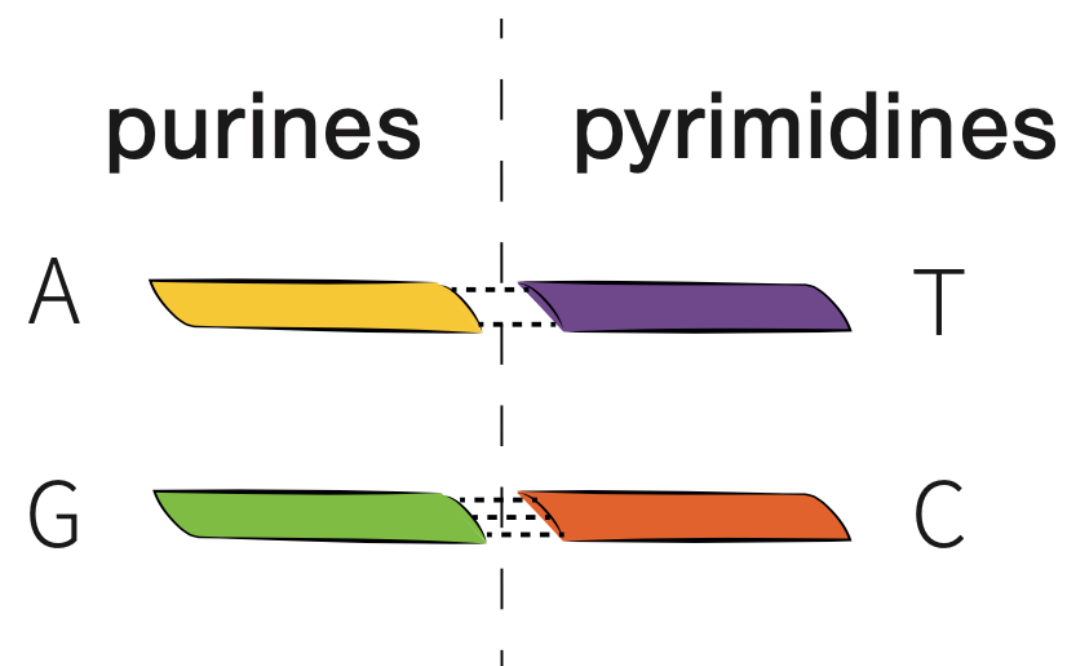
Deoxyribonucleic acid (DNA)



Nucleotide:



phosphate + sugar + nitrogenous base



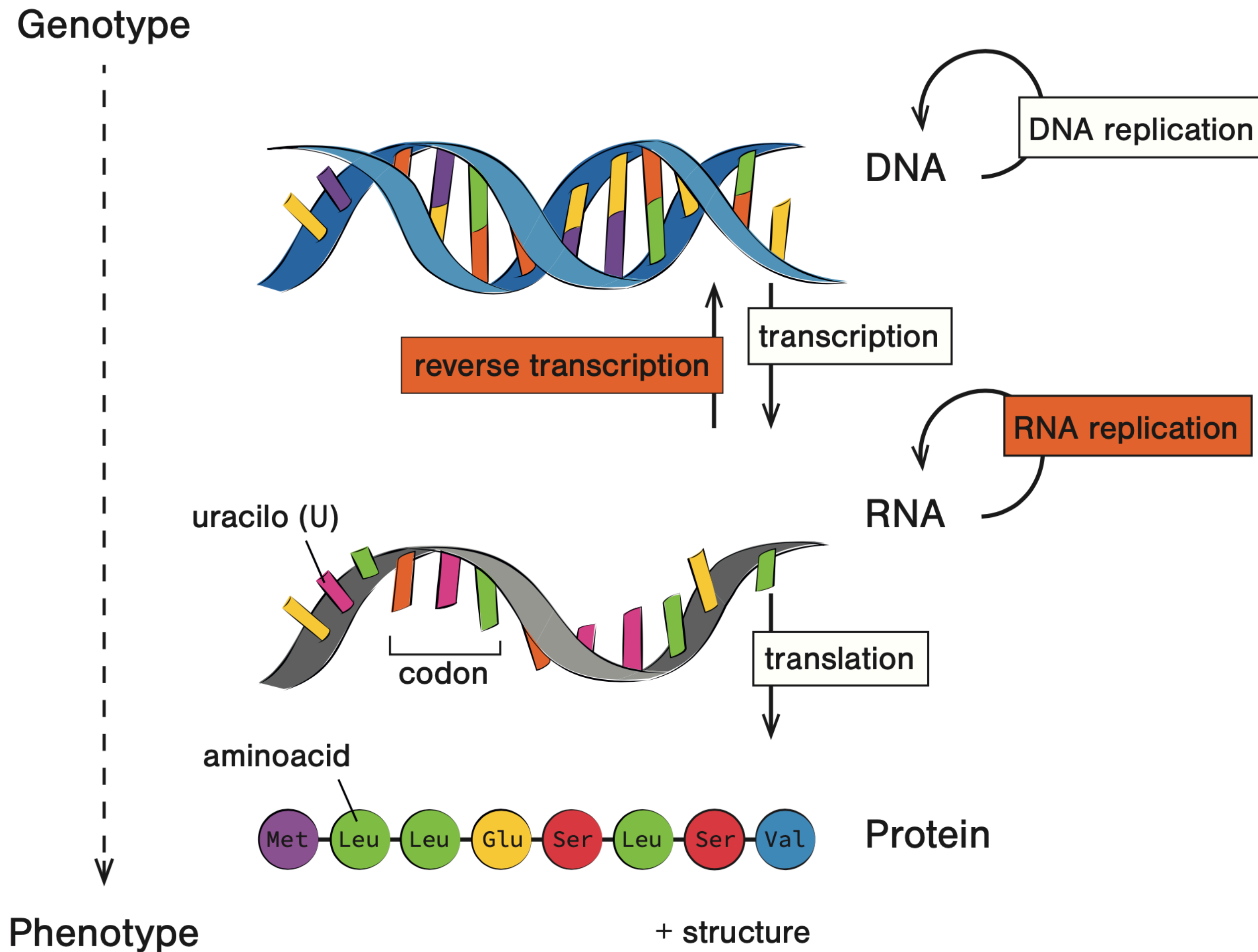
Purines

- Adenine (**A**)
- Guanine (**G**)

Pyrimidines

- Cytosine (**C**)
- Thymine (**T**)*

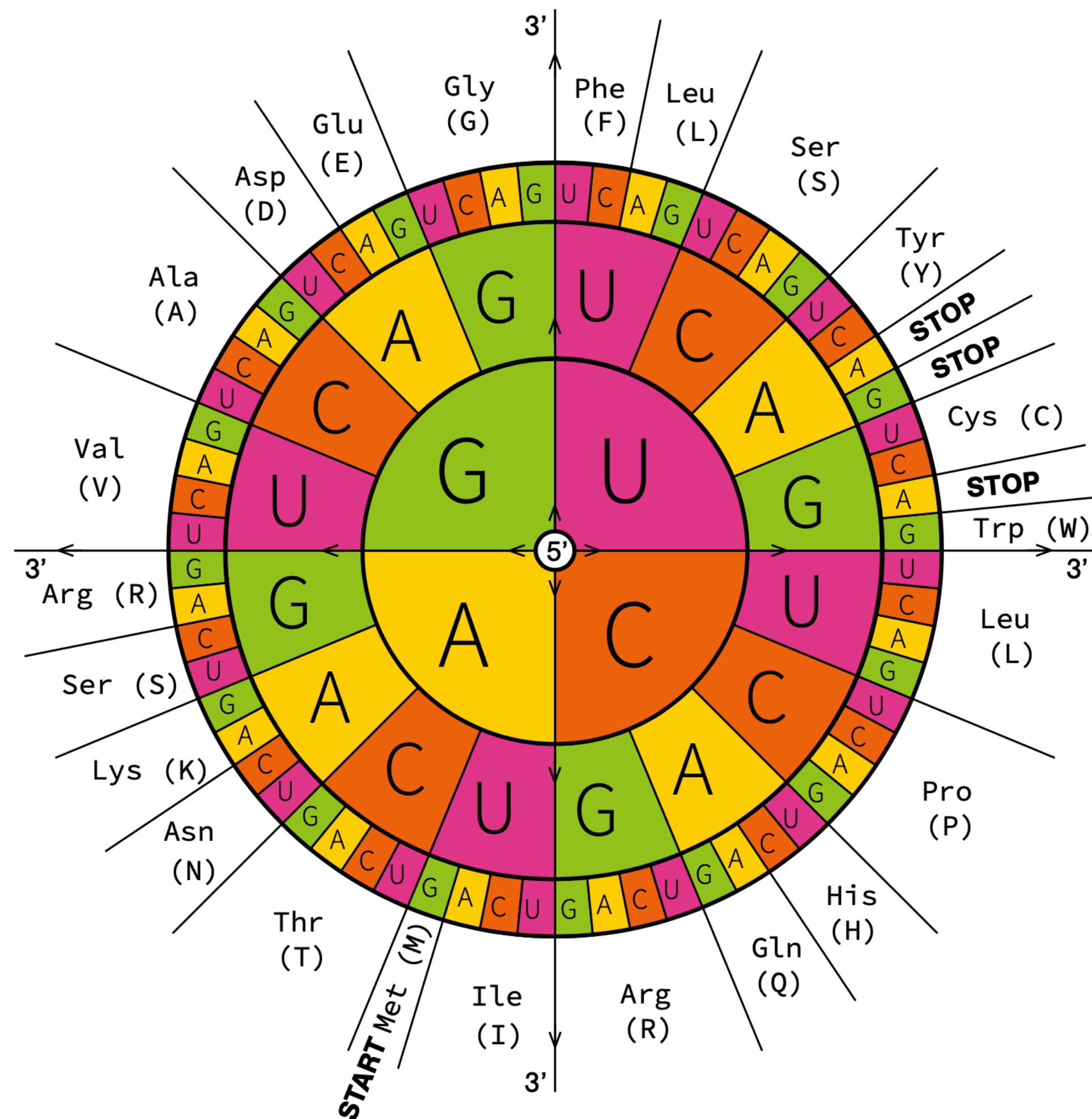
The central dogma of biology



DNA → RNA → protein

Each group of 3 successive nucleotides in a gene is a codon that encodes an amino acid (or terminate translation)

The universal genetic code



Amino acid	3-letter code	1-letter code
Alanine	Ala	A
Arginine	Arg	R
Asparagine	Asn	N
Aspartic acid	Asp	D
Cysteine	Cys	C
Glutamic acid	Glu	E
Glutamine	Gln	Q
Glycine	Gly	G
Histidine	His	H
Isoleucine	Ile	I
Leucine	Leu	L
Lysine	Lys	K
Methionine	Met	M
Phenylalanine	Phe	F
Proline	Pro	P
Serine	Ser	S
Threonine	Thr	T
Tryptophan	Trp	W
Tyrosine	Tyr	Y
Valine	Val	V
STOP		
STOP		
STOP		

$4^3 = 64$ combinations

3 terminate translation

21 amino acids

Mutation vs. substitution

Variation in genotypes (and in phenotypes) is due to errors that arise during DNA replication

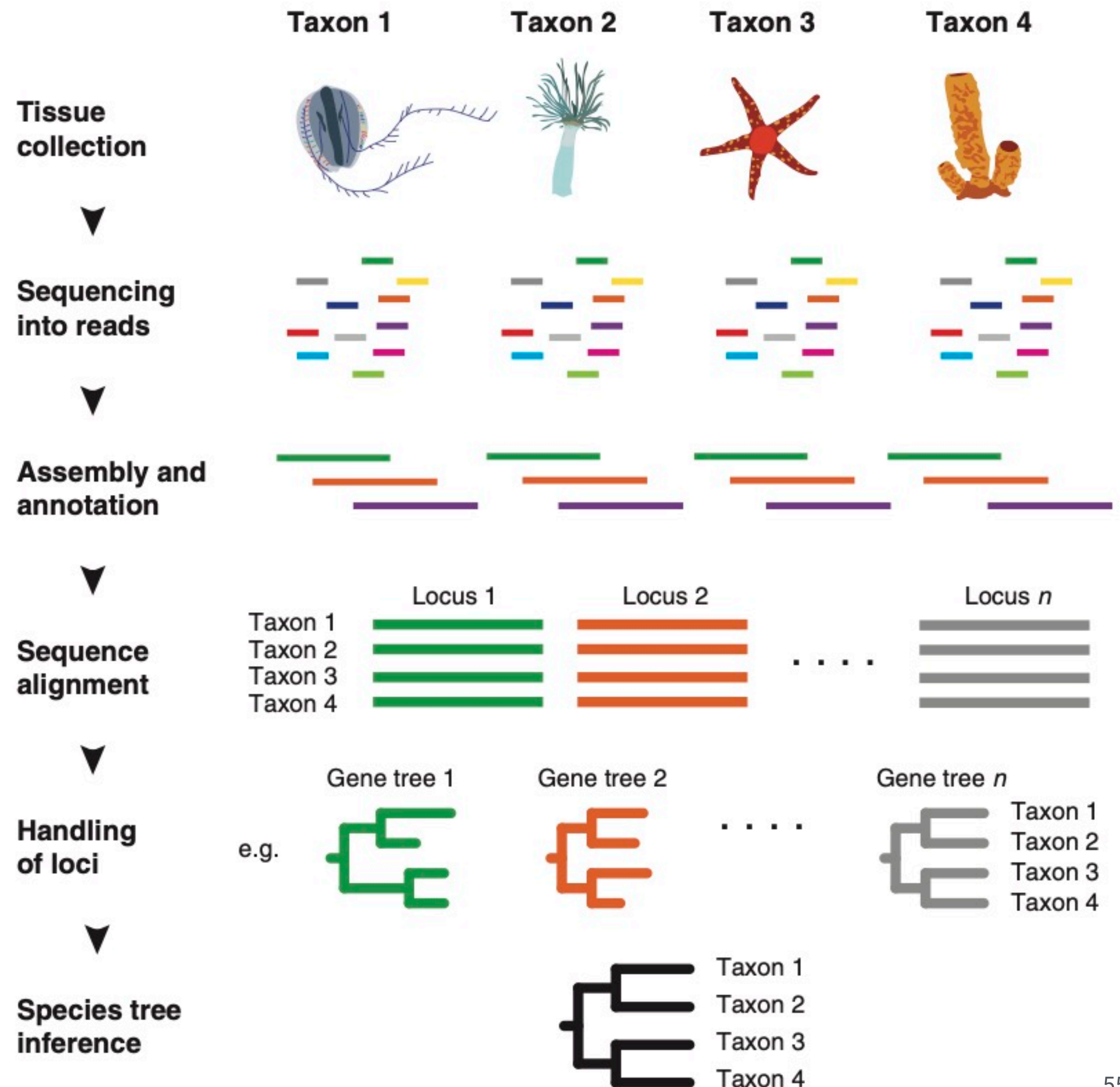
Individuals of the same species have identical characters at *most* positions in their genome (only 0.1% vary among humans)

Most mutations are repaired but can persist across generations

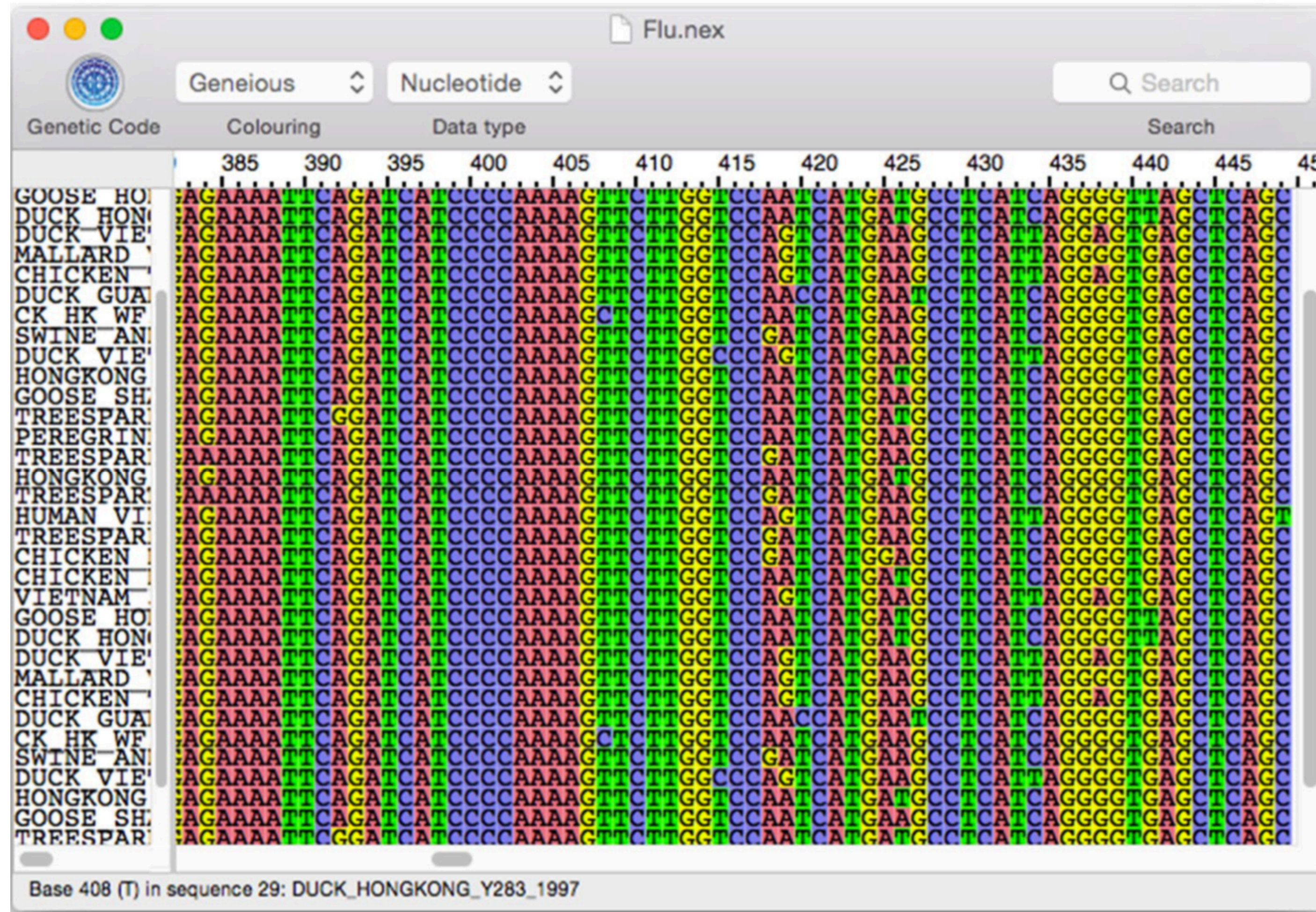
Mutations that spread throughout a population and become 'fixed' called **substitutions**

DNA sequencing

Multiple sequence alignment software establishes homology across sites from different species



Multiple sequence alignment



#NEXUS

[Cytochrome oxidase B genes - bears]

[Data source: <https://revbayes.github.io/tutorials/dating/>]

BEGIN DATA;

DIMENSIONS NTAX=10 NCHAR=1000;

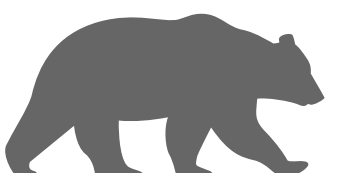
FORMAT DATATYPE = DNA MISSING=? GAP=- ;

MATRIX

Ailuropoda_melanoleuca	ATGATCAACATCCGAAAAACTCATCCATTAGTTAAAATTATCAACAACACTCATTCATTGACCT...
Arctodus_simus	ATGACCAACATCCGAAAGACTCACCCACTGGCCAAAATTATCAATAACTCATTCATCGACCT...
Helarctos_malayanus	ATGACCAACATCCGAAAAACCCACCCATTAGCTAAAATCATTAACAACACTCACTTATTGACCT...
Melursus_ursinus	ATGACCAACATCCGAAAAACCCACCCACTAGCTAAAATCATTAACAACACTCACTCATTTGACCT...
Ursus_americanus	ATGACCAACATCCGAAAAACCCACCCATTAGCTAAAATCATCAACAACACTCACTTATTGATCT...
Ursus_arctos	ATGACCAACATCCGAAAAACCCACCCATTAGCTAAAATCATCAACAACACTCATTTATTGACCT...
Ursus_maritimus	ATGACCAACATCCGAAAAACCCACCCATTAGCTAAAATCATCAACAACACTCATTTATTGATCT...
Ursus_thibetanus	ATGACCAACATCCGAAAAACCCATCCATTAGCCAAAATCATCAACAACACTCACTCATTTGATCT...
Ursus_spelaeus	ATGACCAACATCCGAAAAACCCATCCACTAGCTAAAATCATCAACAACACTCATTCATTGACCT...
Tremarctos_ornatus	ATGACCAACATCCGAAAAACTCACCCACTAGCTAAAATCATCAACAACACTCATTCATCGACCT...

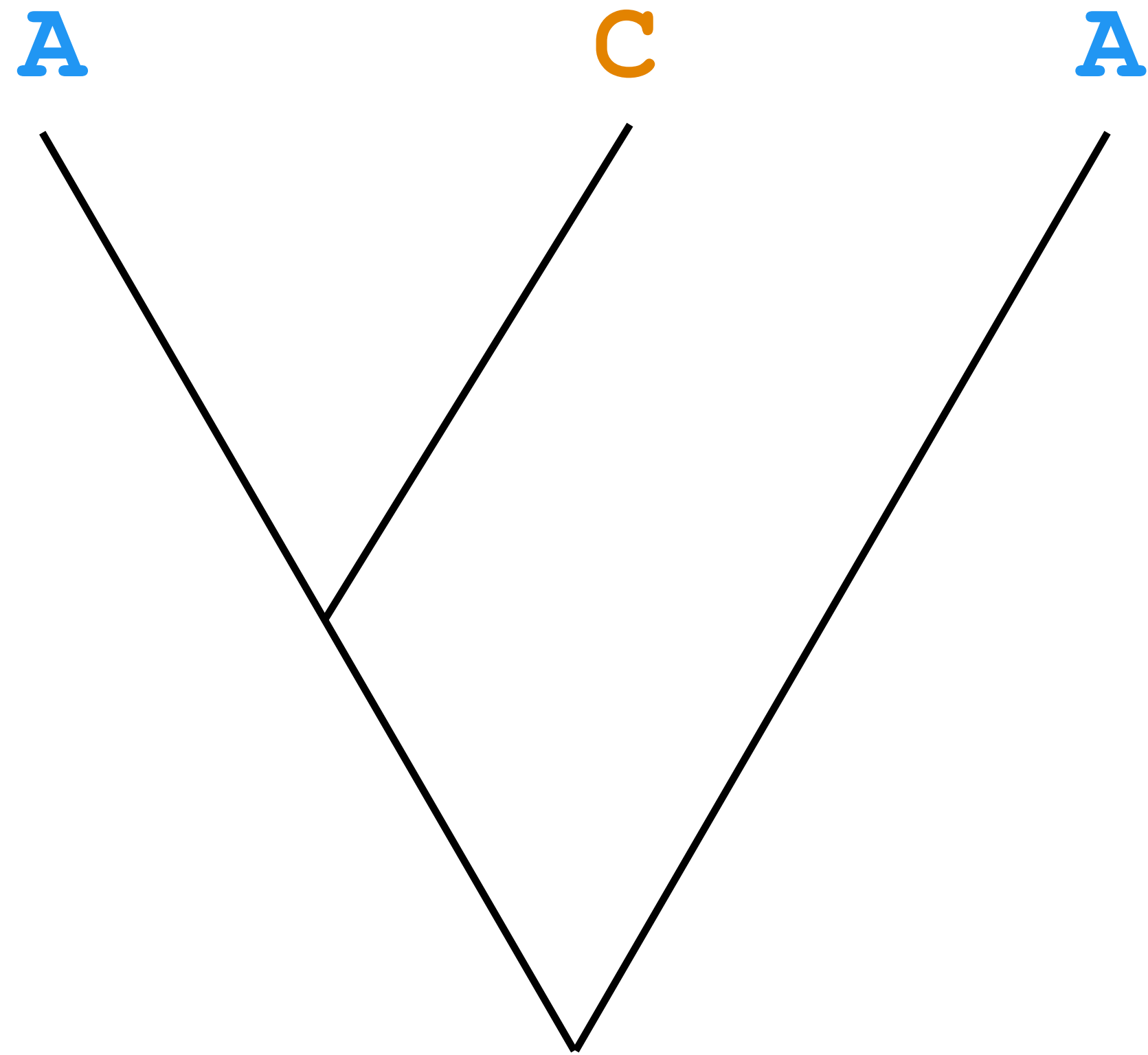
;

END;

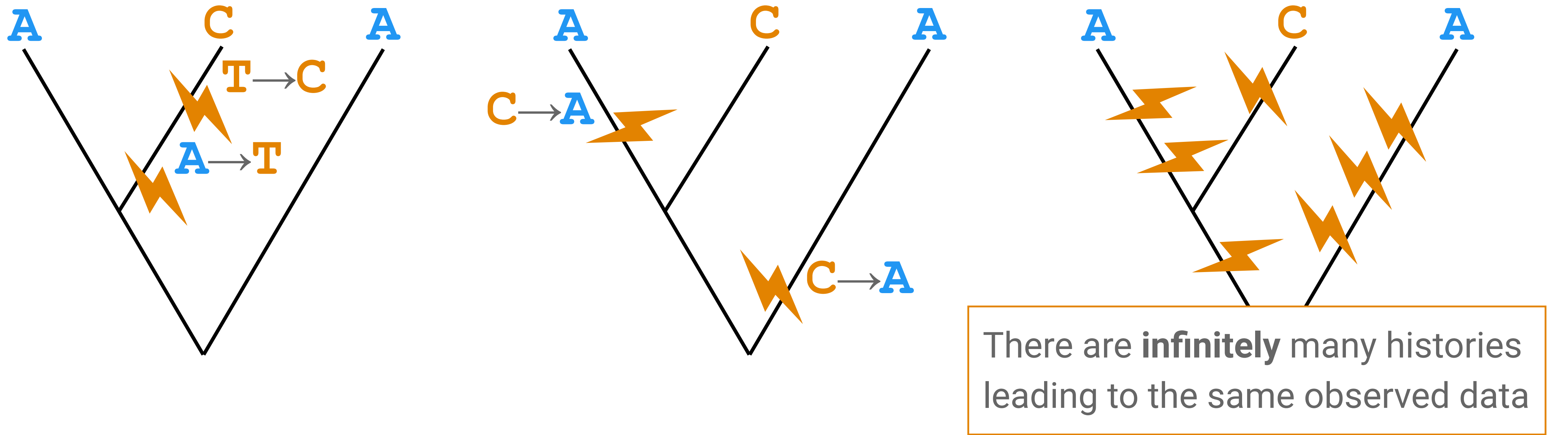
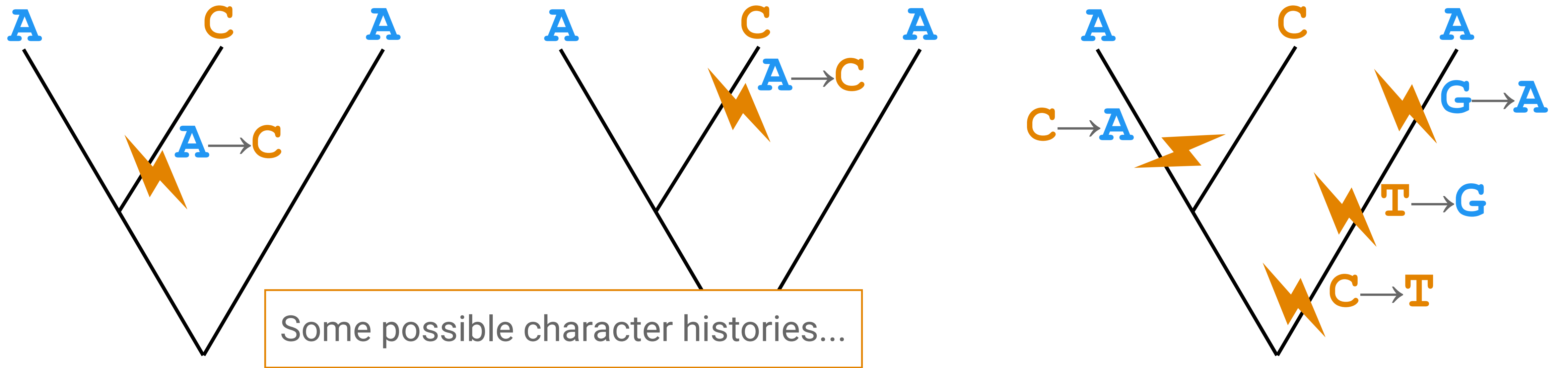


The data are the observed states at the tips

How probable is our data, given my tree?

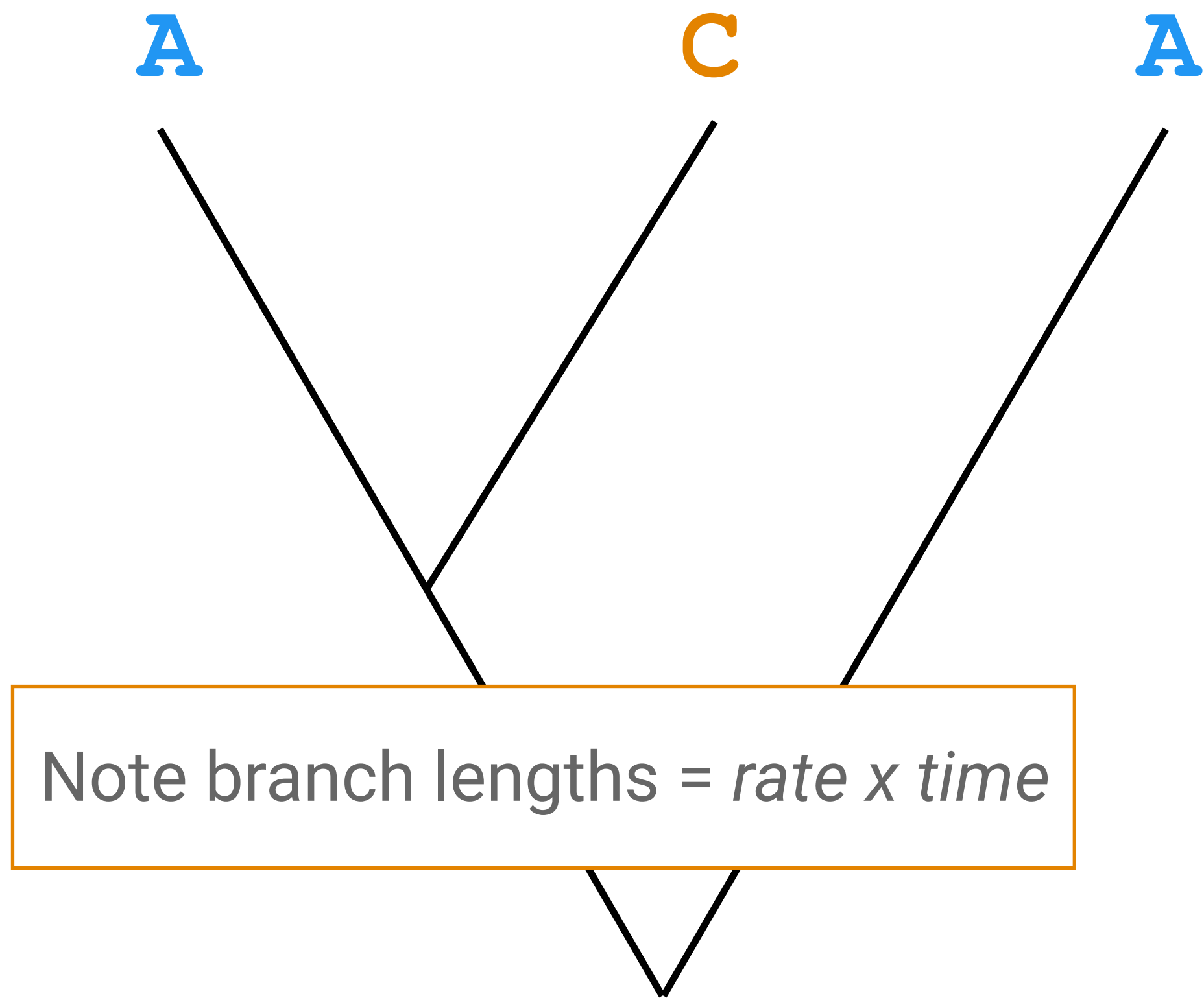


To apply a model based approach we need to be able to compute the probability of our sequence alignment (or character matrix)



The data are the observed states at the tips

How probable is our data, given my tree?



To compute P , we need:

- A model of sequence (or character) evolution
- A way of calculating the probability for given a phylogeny (tree topology + branch lengths)

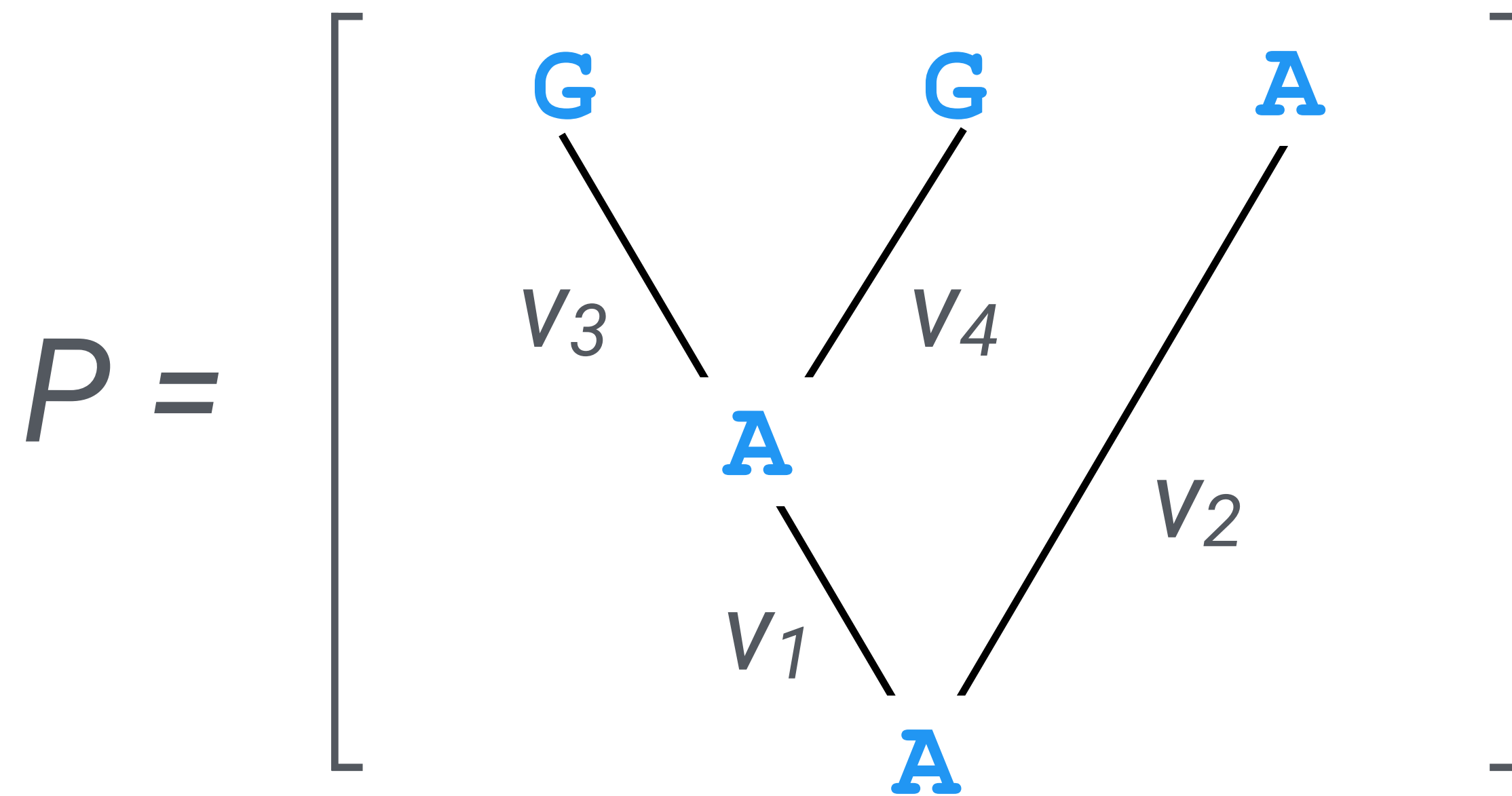
Substitution models

Models of molecular sequence evolution

Also known as substitution / site / character models

They allow us to compute the probability of changing from one state to another over branch length v

Computing the probability of the observed data



Just suppose for now
we know the ancestral
states at internal nodes

$$P_{AA}(v_1) \times P_{AA}(v_2) \times P_{AG}(v_3) \times P_{AG}(v_4)$$

$P_{ij}(v)$ — transition probabilities

Rate matrix

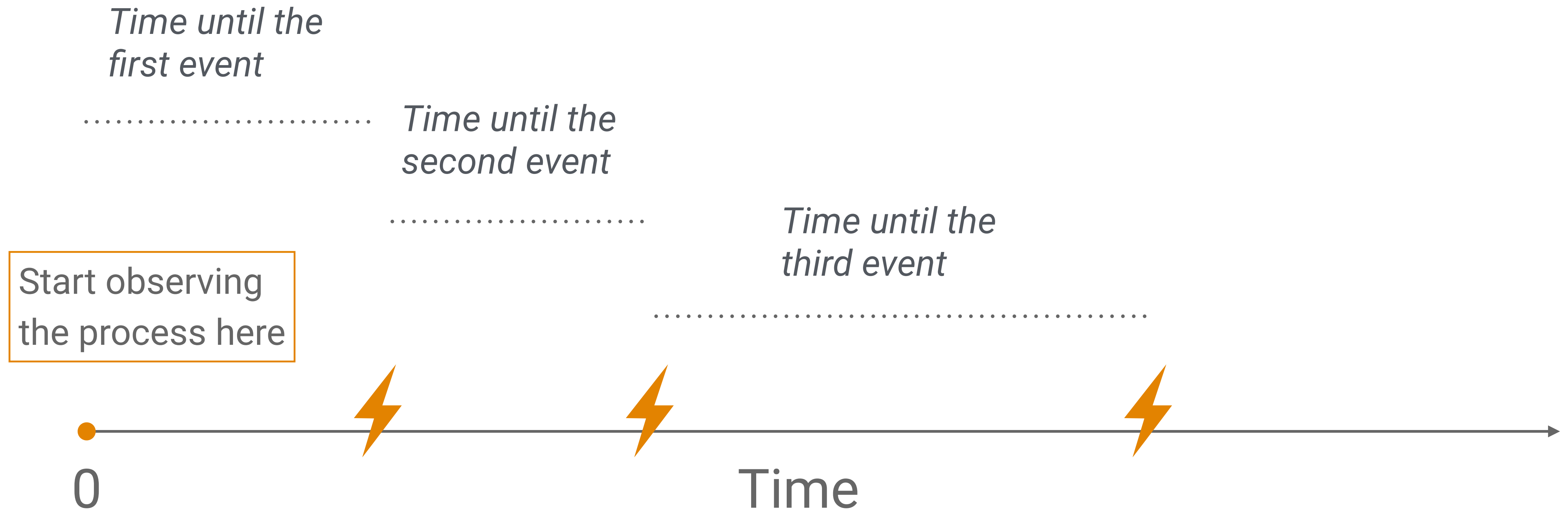
$$Q = \begin{matrix} & \mathbf{A} & \mathbf{T} & \mathbf{G} & \mathbf{C} \\ \mathbf{A} & \left[\begin{array}{cccc} -3\lambda & \lambda & \lambda & \lambda \\ \lambda & -3\lambda & \lambda & \lambda \\ \lambda & \lambda & -3\lambda & \lambda \\ \lambda & \lambda & \lambda & -3\lambda \end{array} \right] \\ \mathbf{T} & & & & \\ \mathbf{G} & & & & \\ \mathbf{C} & & & & \end{matrix}$$

In this model, we only have one parameter, substitution rate parameter λ

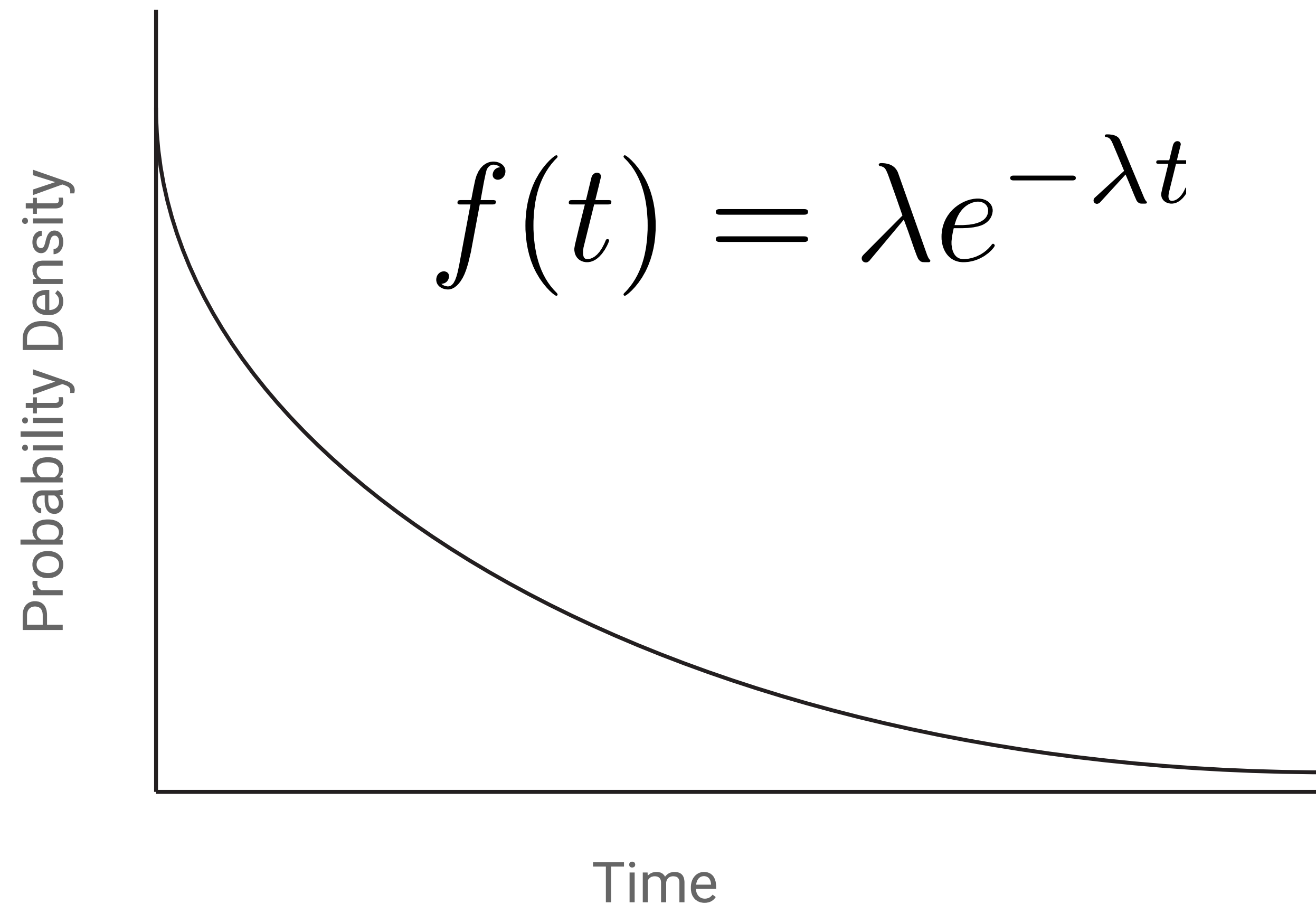
This is the Jukes-Cantor (1969) or JC69 model

Continuous time Markov chain

Nucleotide substitutions (events) occur at a constant rate



The poisson process



The waiting times are **exponentially** distributed random variables

We can use this to calculate the probability of change over time (or branch length v)

The longer the interval of time, the more likely we are to observe change

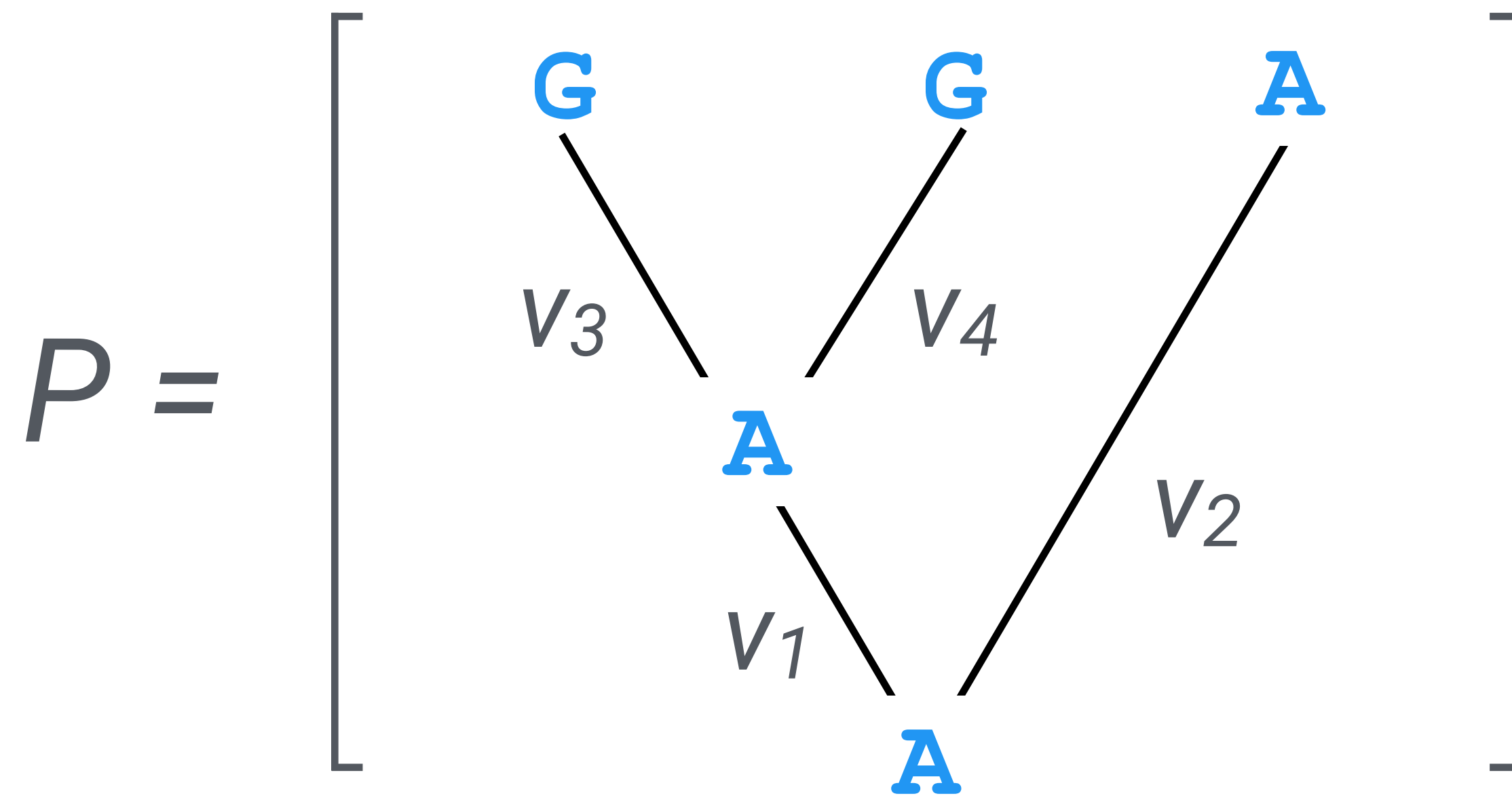
Exercise

Jukes-Cantor model transition probability applet

Felsenstein's pruning algorithm

The following slides are adapted from John Huelsenbeck (c/o Sebastian Höhna)

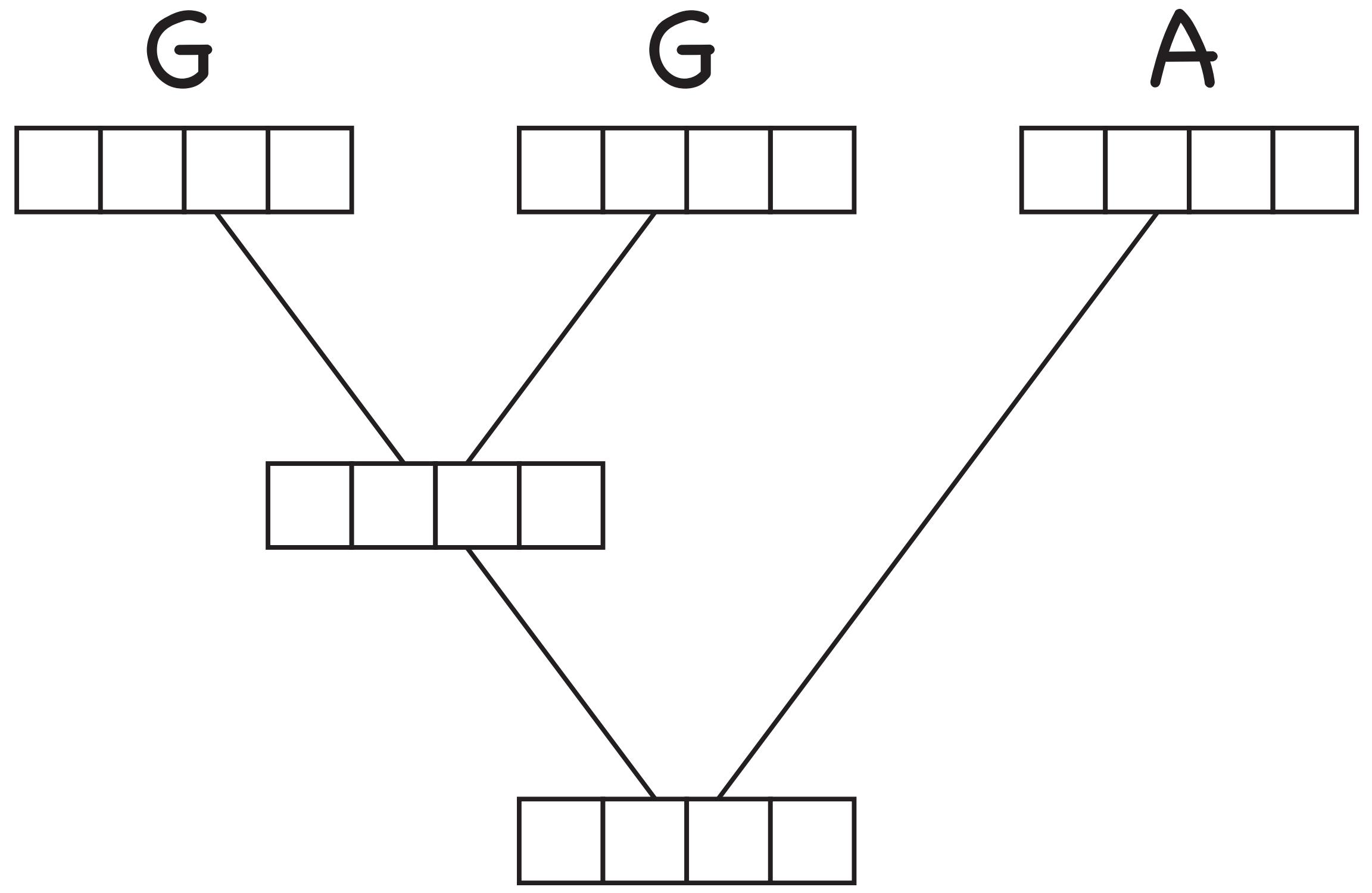
Computing the probability of the observed data



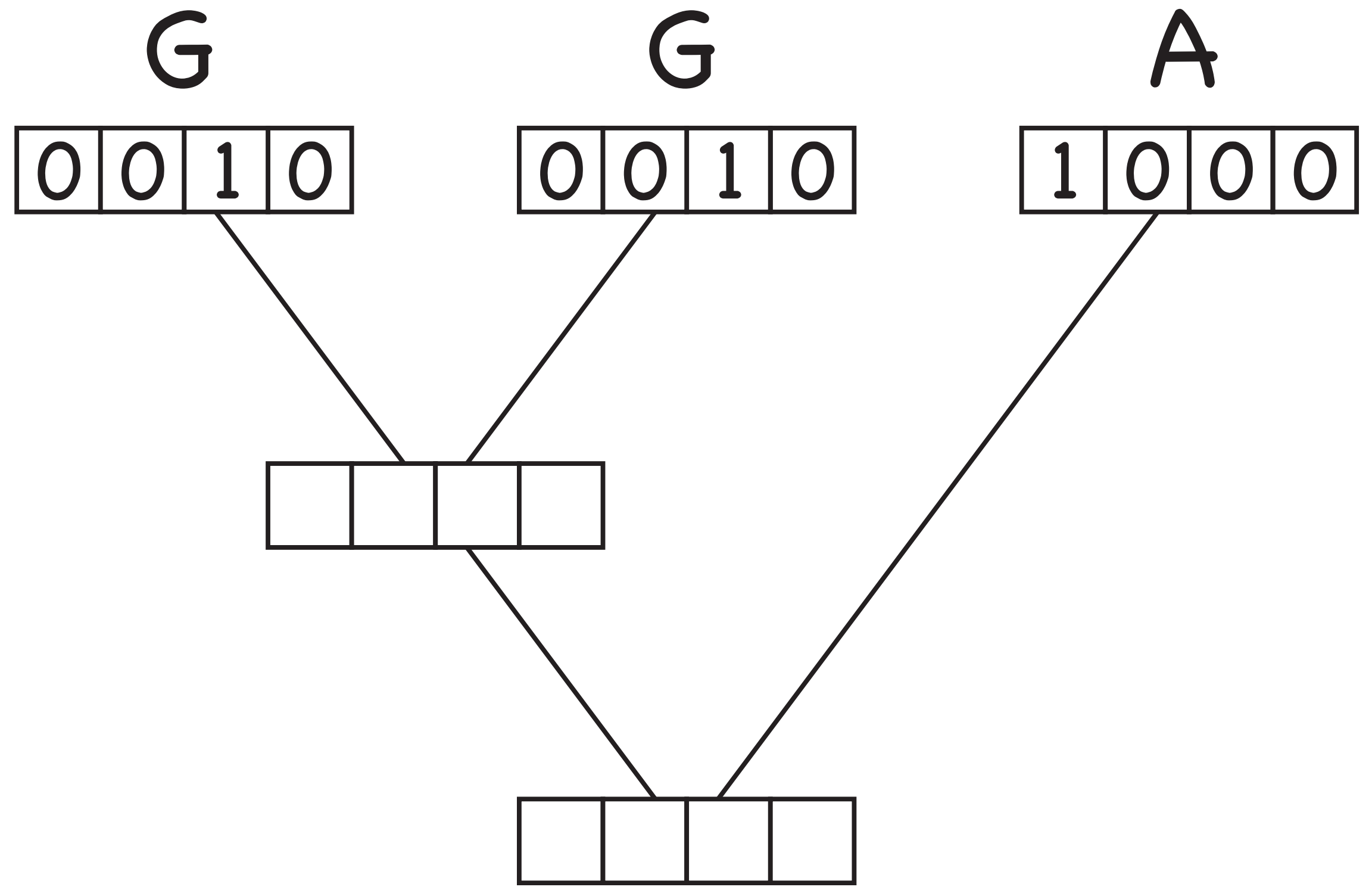
Just suppose for now we know the ancestral states at internal nodes

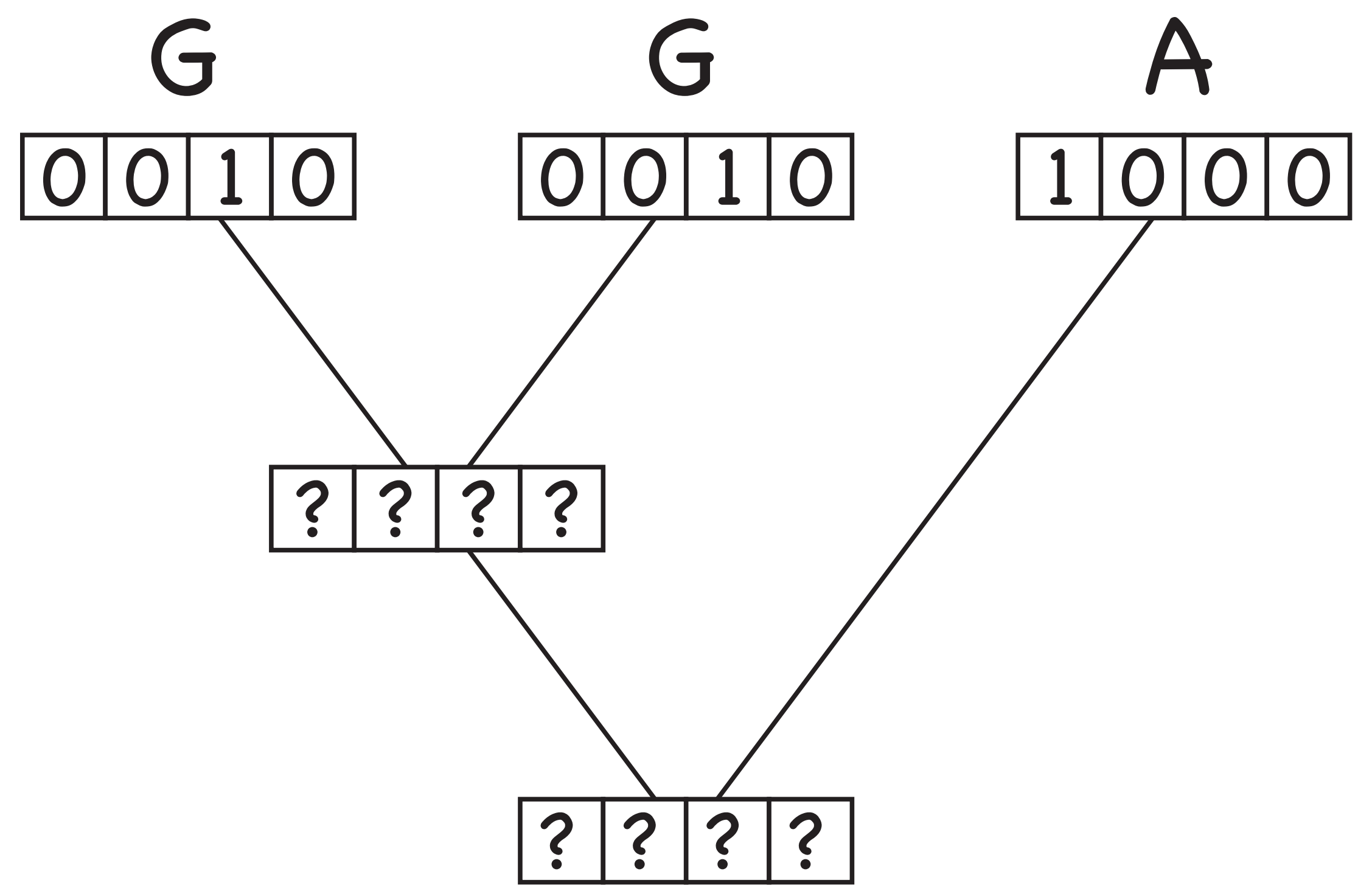
$$\pi_A \times P_{AA}(v_1) \times P_{AA}(v_2) \times P_{AG}(v_3) \times P_{AG}(v_4)$$

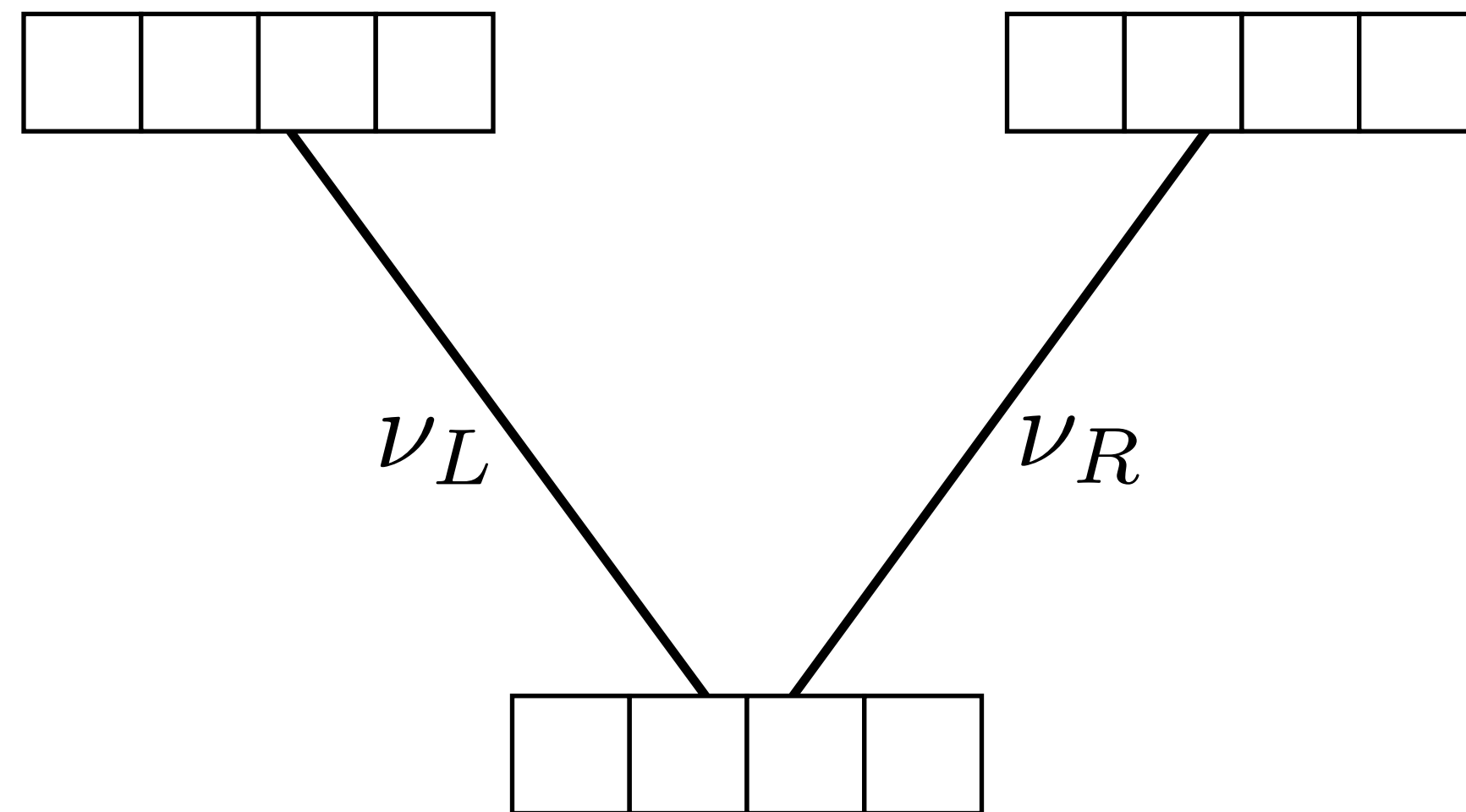
$P_{ij}(v)$ – transition probabilities
 π_i – stationary frequencies



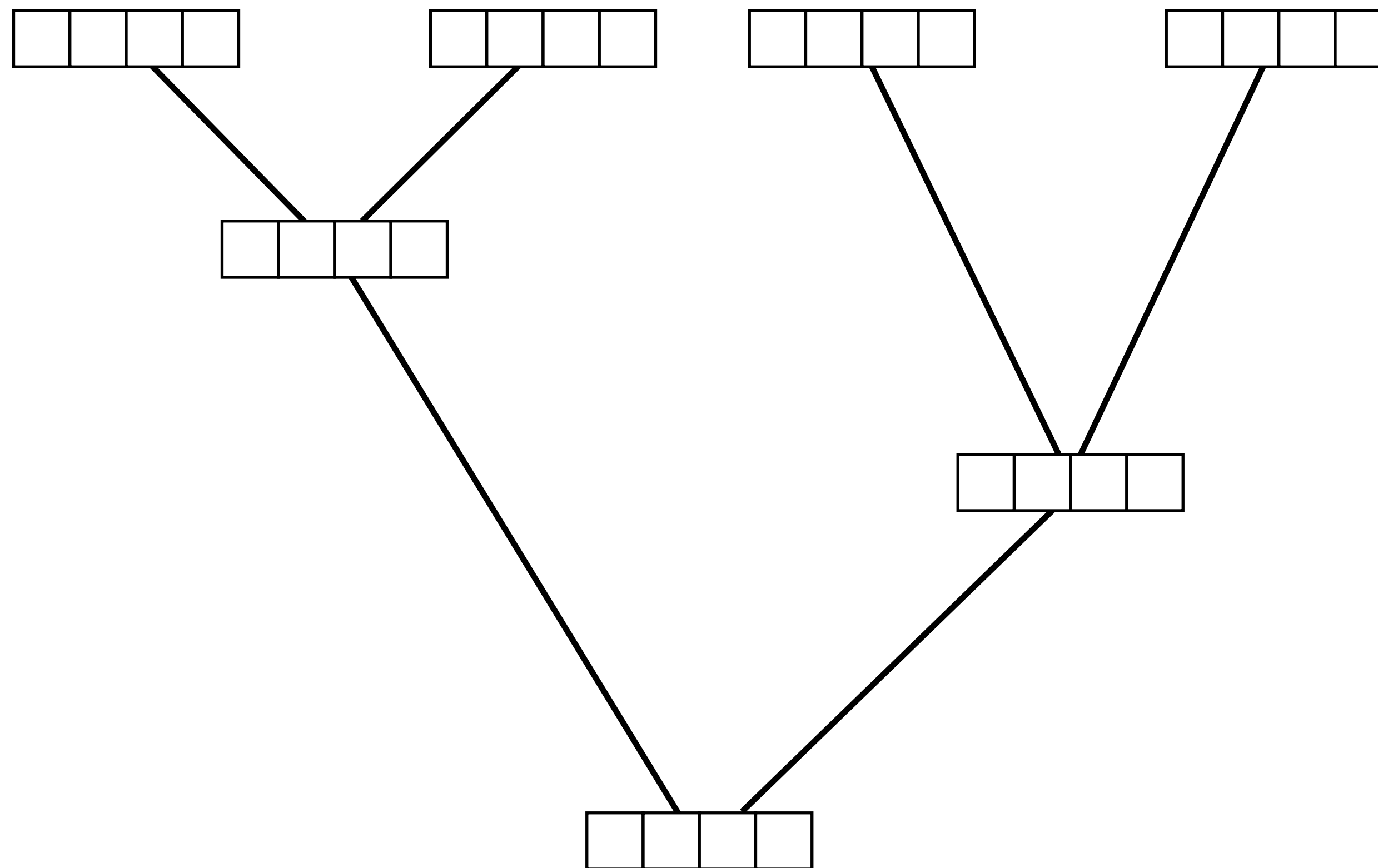
Felsenstein, J. 1981. Evolutionary trees from DNA sequences: A maximum likelihood approach.



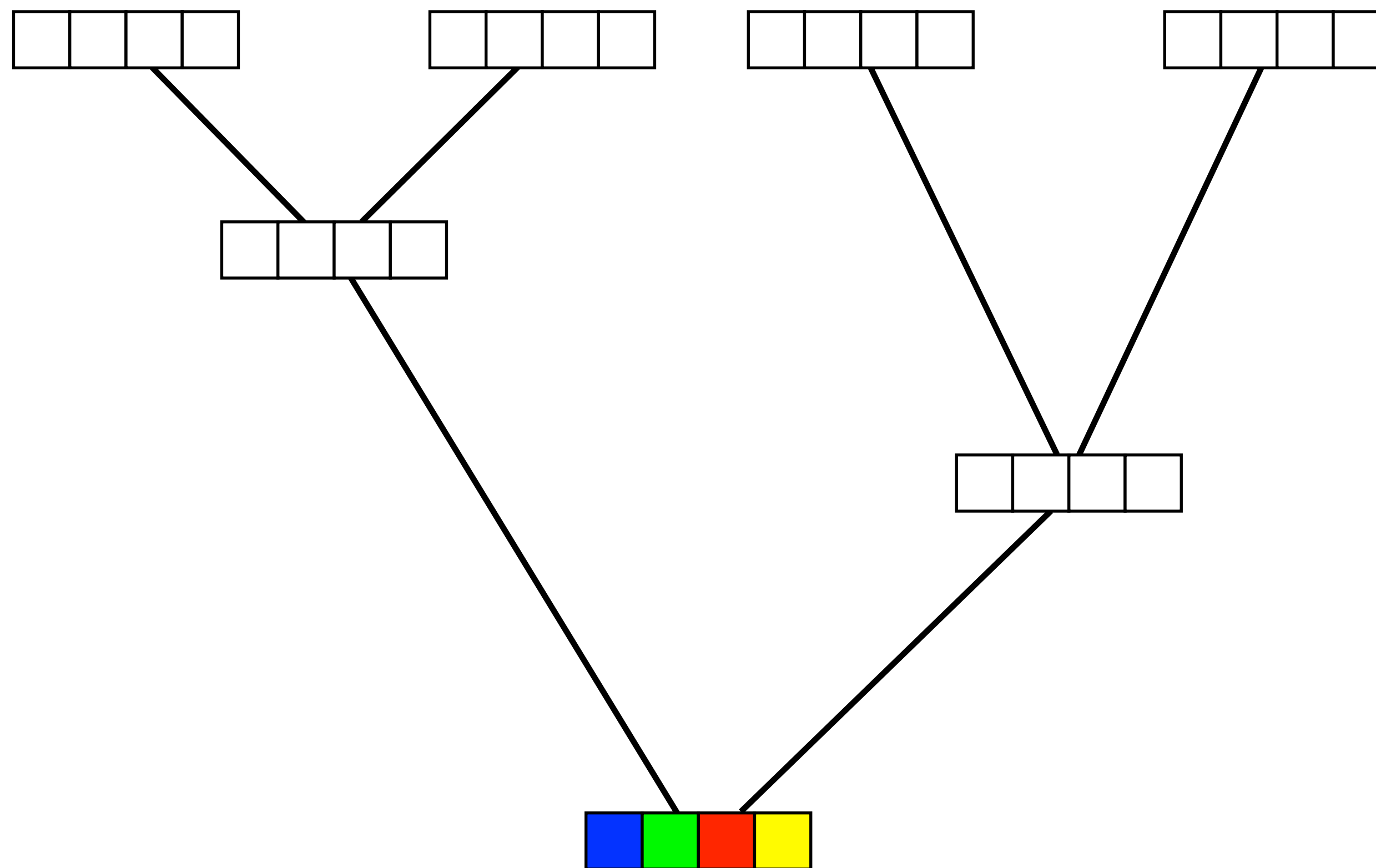




$$\ell_i = \left(\sum_j p_{ij}(\nu_L) \ell_j^L \right) \times \left(\sum_j p_{ij}(\nu_R) \ell_j^R \right)$$



$$l_{\text{Site}} = \pi_A \times l_A^{\text{Root}} + \pi_C \times l_C^{\text{Root}} + \pi_G \times l_G^{\text{Root}} + \pi_T \times l_T^{\text{Root}}$$

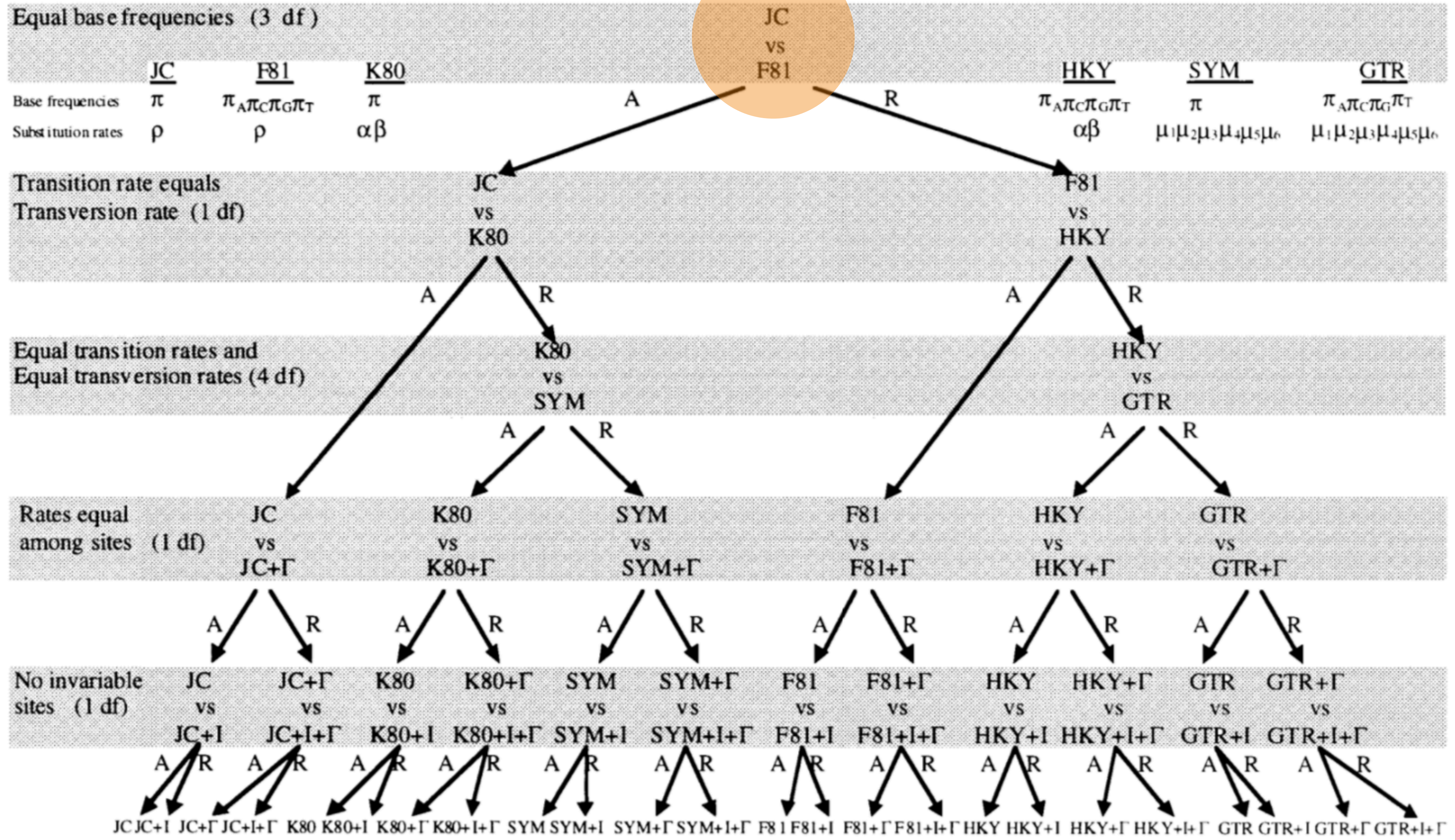


$$l_{\text{Site}} = \pi_A \times l_A^{\text{Root}} + \pi_C \times l_C^{\text{Root}} + \pi_G \times l_G^{\text{Root}} + \pi_T \times l_T^{\text{Root}}$$

$$l_i = \left(\sum_j p_{ij}(\nu_L) l_j^L \right) \times \left(\sum_j p_{ij}(\nu_R) l_j^R \right)$$

$$l_{\text{Site}} = \pi_A \times l_A^{\text{Root}} + \pi_C \times l_C^{\text{Root}} + \pi_G \times l_G^{\text{Root}} + \pi_T \times l_T^{\text{Root}}$$

Other substitution models



Base frequencies

The JC69 model assumes equal transition rates and equal base frequencies

Base frequencies are the proportion of each nucleotide in the dataset

If a given nucleotide appears in our dataset at a low frequency, we are less likely to observe a transition to that state

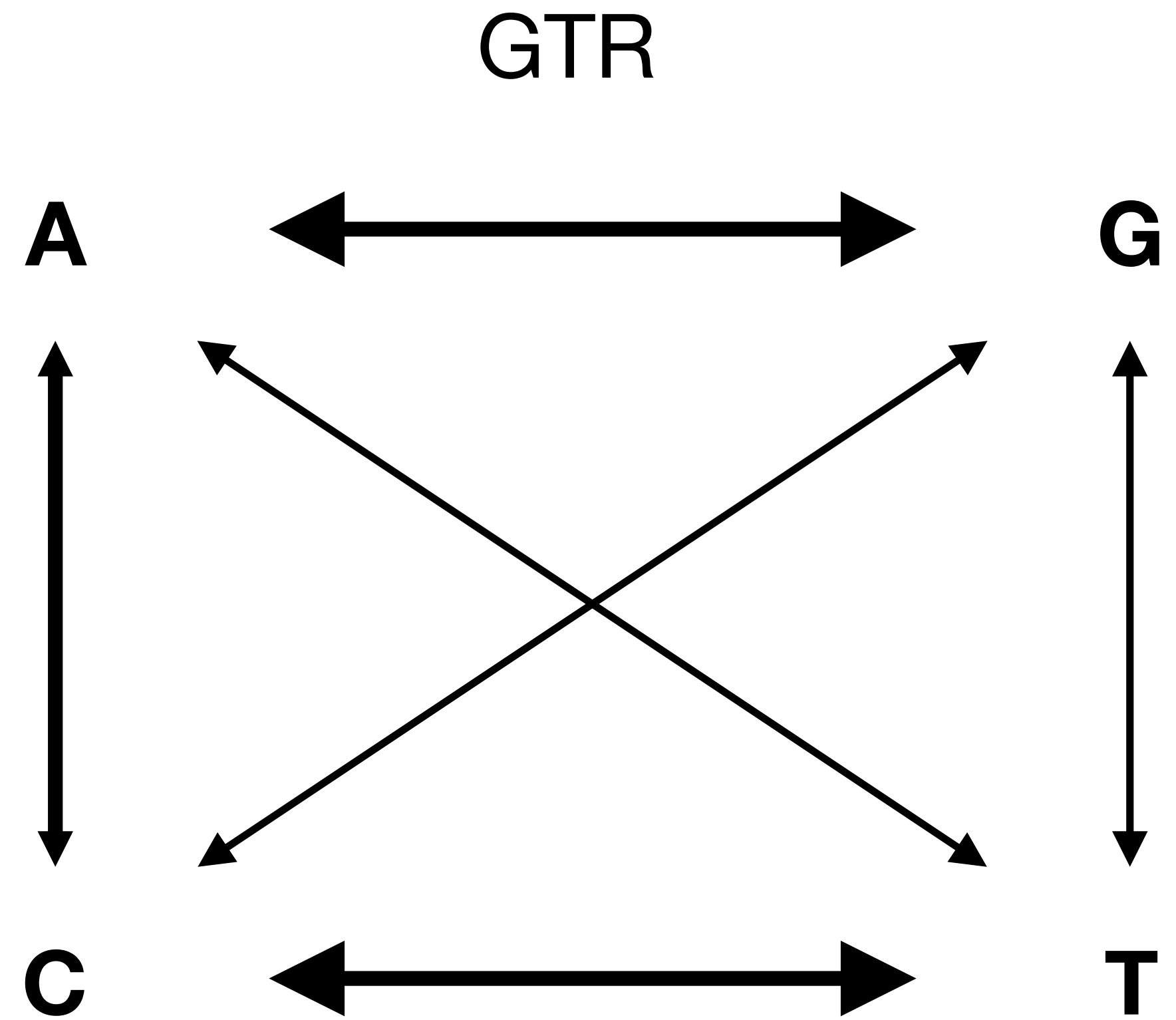
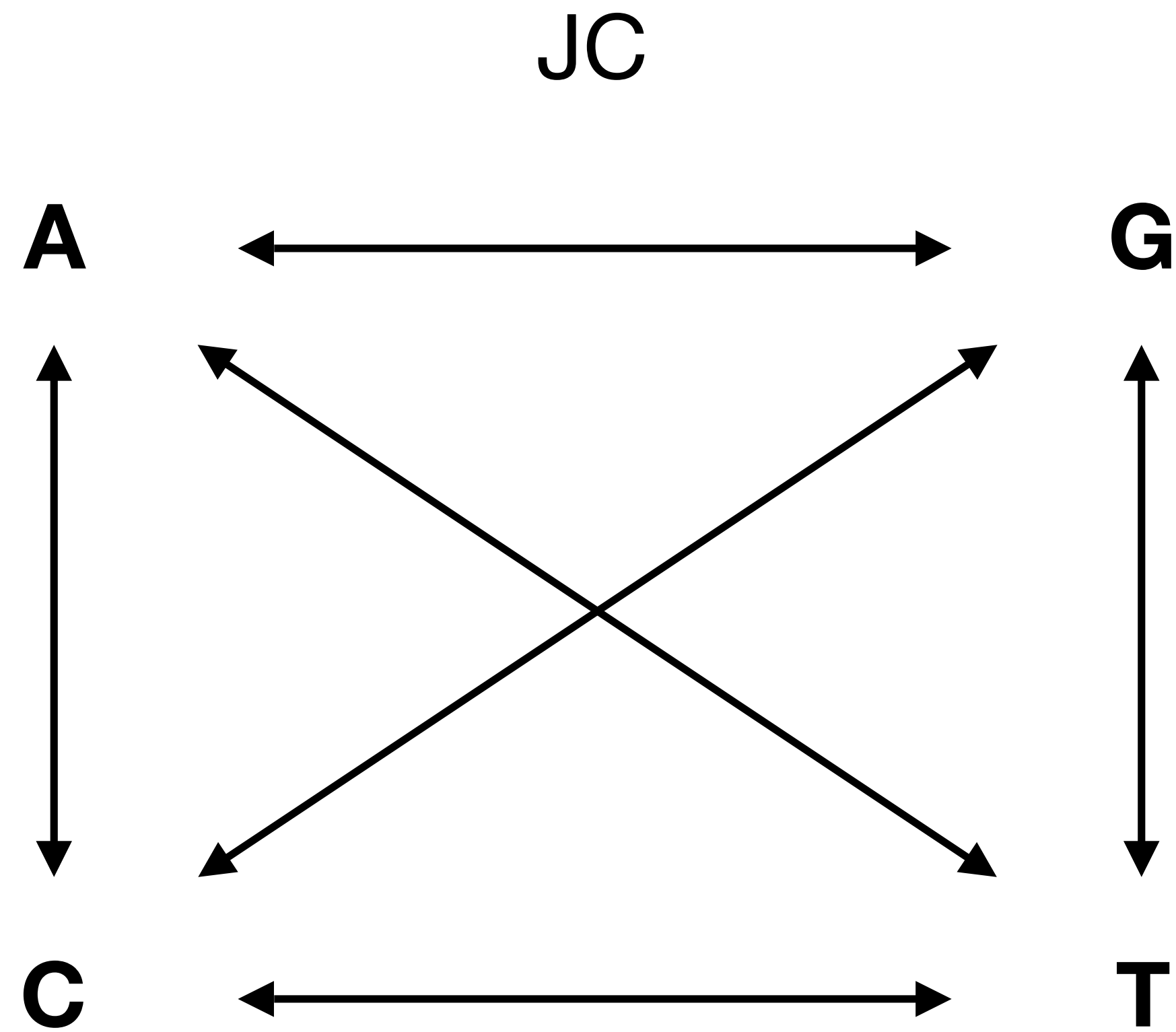
The general time reversible model

Allows for unequal transition rates (μ) and unequal base frequencies (π)

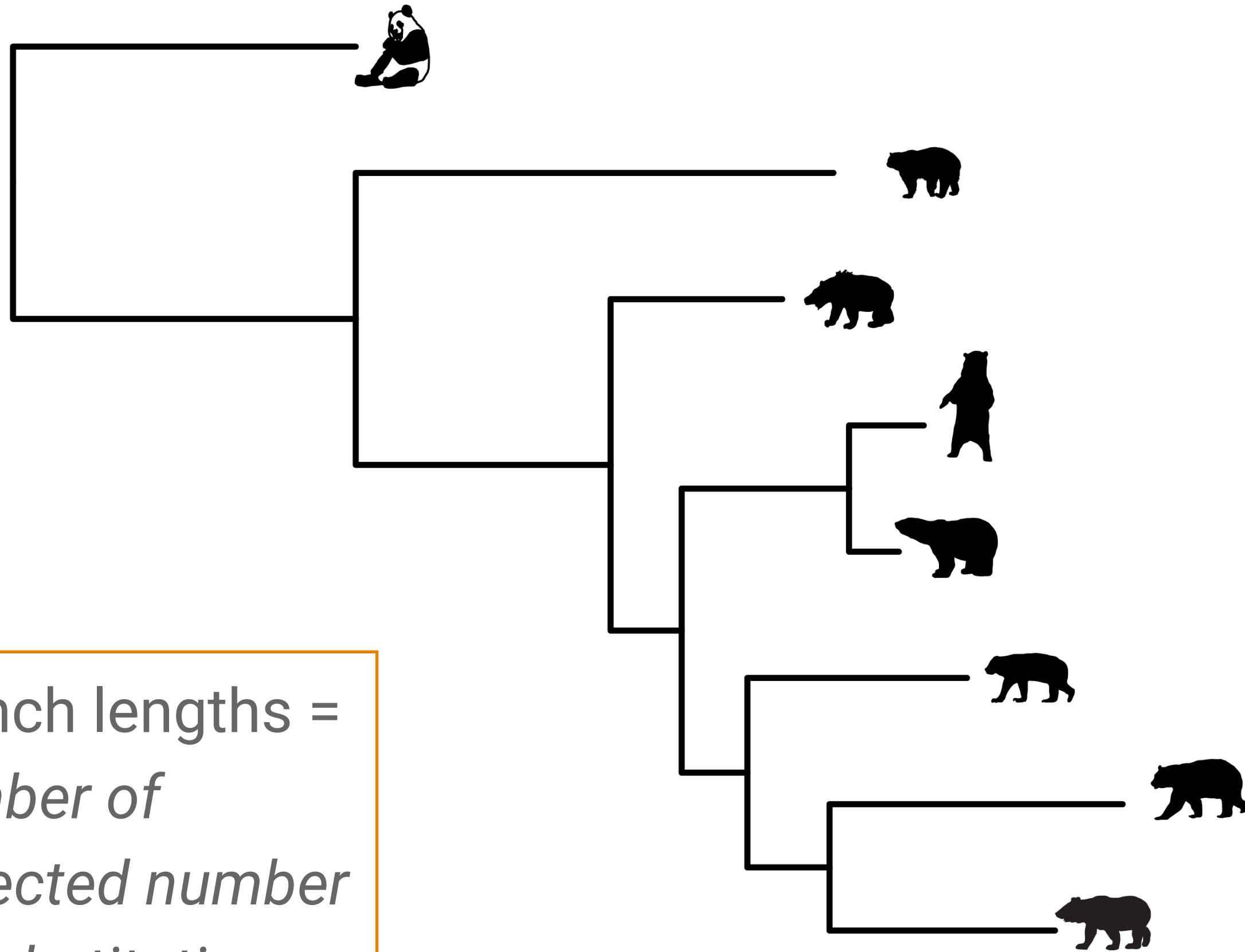
$$Q = \begin{pmatrix} * & \mu_{AG}\pi_G & \mu_{AC}\pi_C & \mu_{AT}\pi_T \\ \mu_{GA}\pi_A & * & \mu_{GC}\pi_C & \mu_{GT}\pi_T \\ \mu_{CA}\pi_A & \mu_{CG}\pi_G & * & \mu_{CT}\pi_T \\ \mu_{TA}\pi_A & \mu_{TG}\pi_G & \mu_{TC}\pi_C & * \end{pmatrix}$$

Note the rates are symmetric – e.g., the rate of change between A and T, is the same in both directions – but the frequency of each character state also affects the probability of change

The JC versus GTR models



Branch lengths



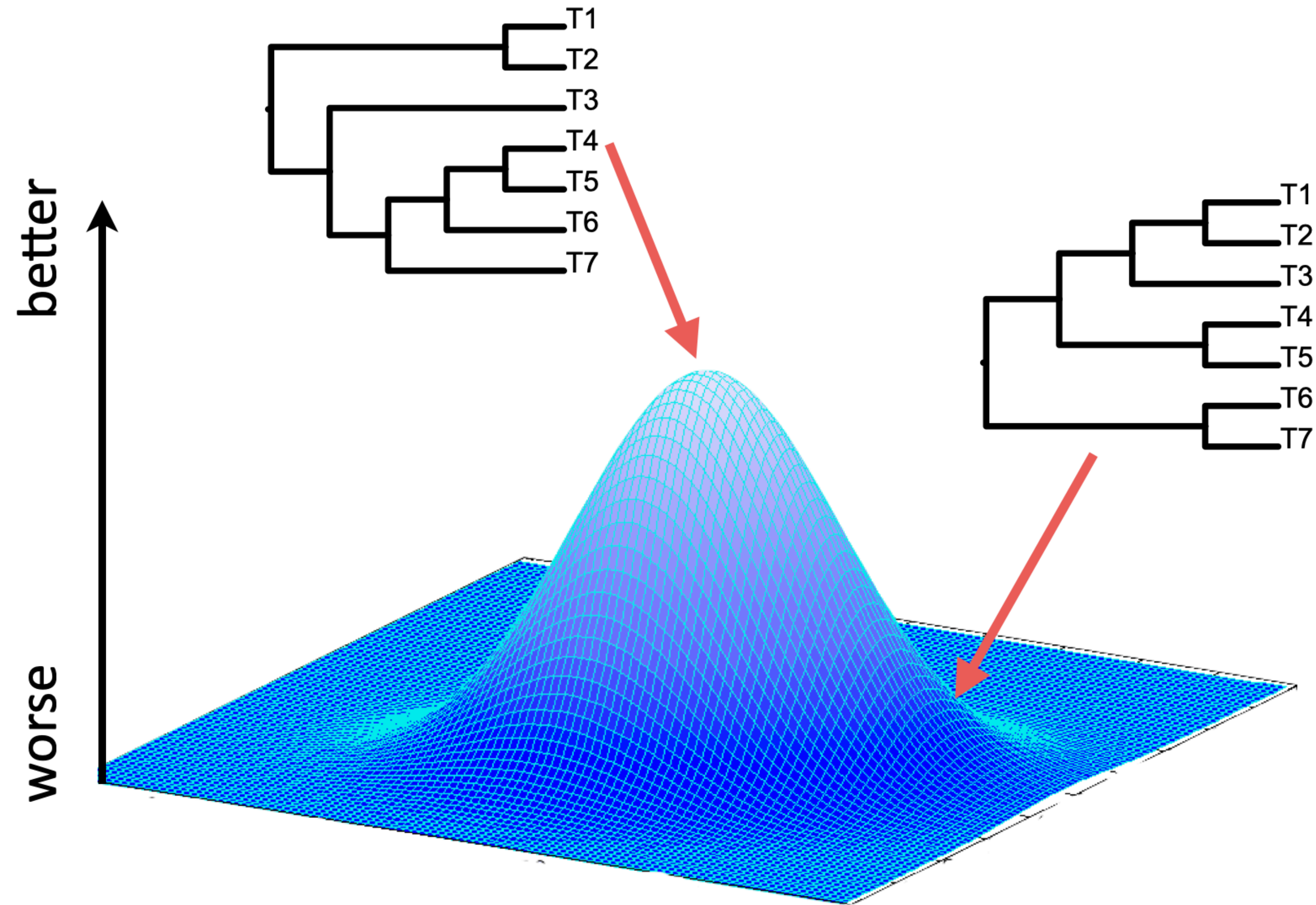
Branch lengths =
*number of
expected number
of substitutions
per site*

Branch lengths are a product
of rate and time

Without temporal information
we can only measure relative
genetic distance

Maximum likelihood

How do we find the 'best' tree?



It depends how you measure 'best'

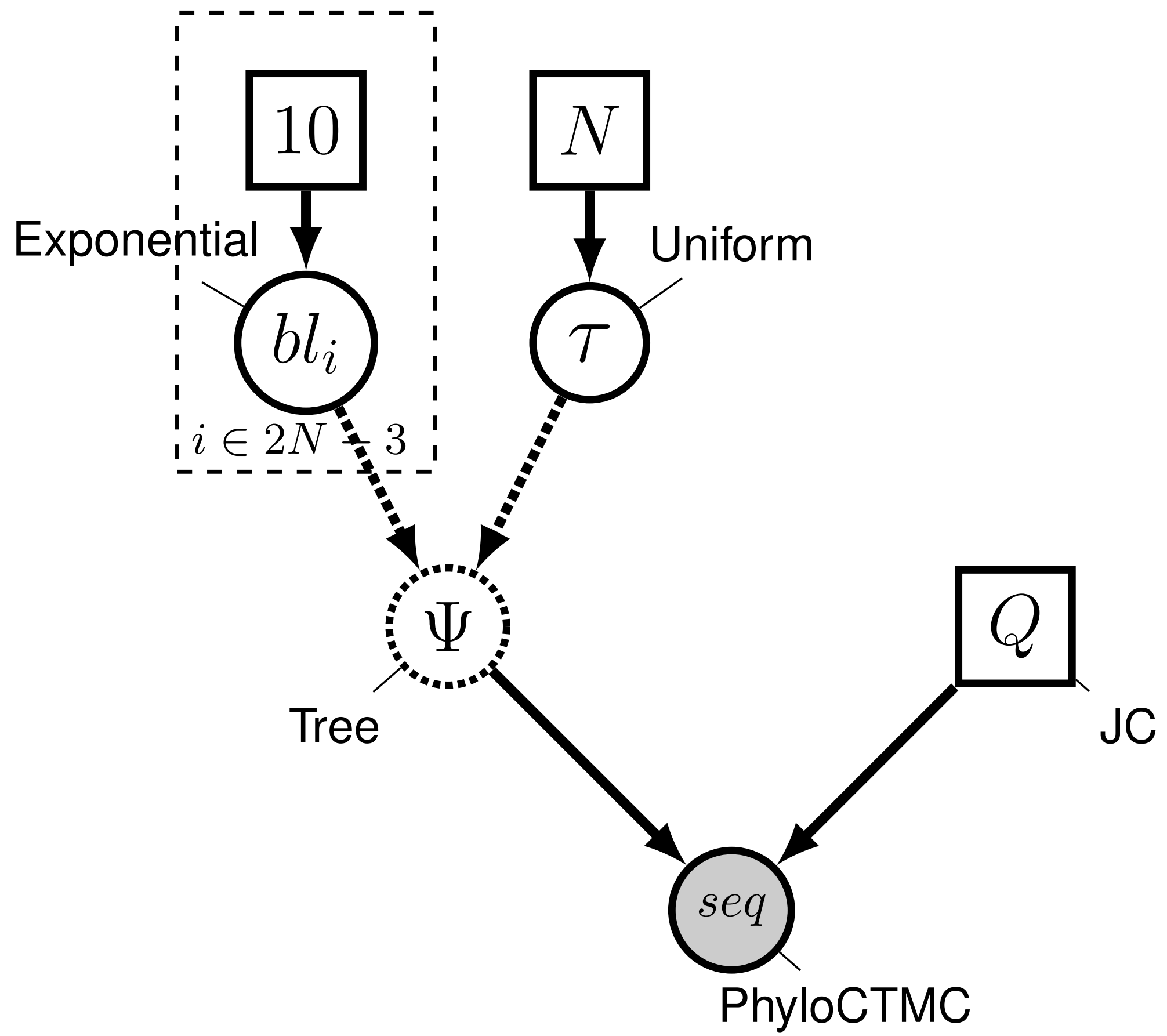
Method	Criterion (tree score)
Maximum parsimony	Minimum number of changes
Maximum likelihood	Likelihood score (probability), optimised over branch lengths and model parameters
Bayesian inference	Posterior probability, integrating over branch lengths and model parameters

Both maximum likelihood and Bayesian inference are model-based approaches

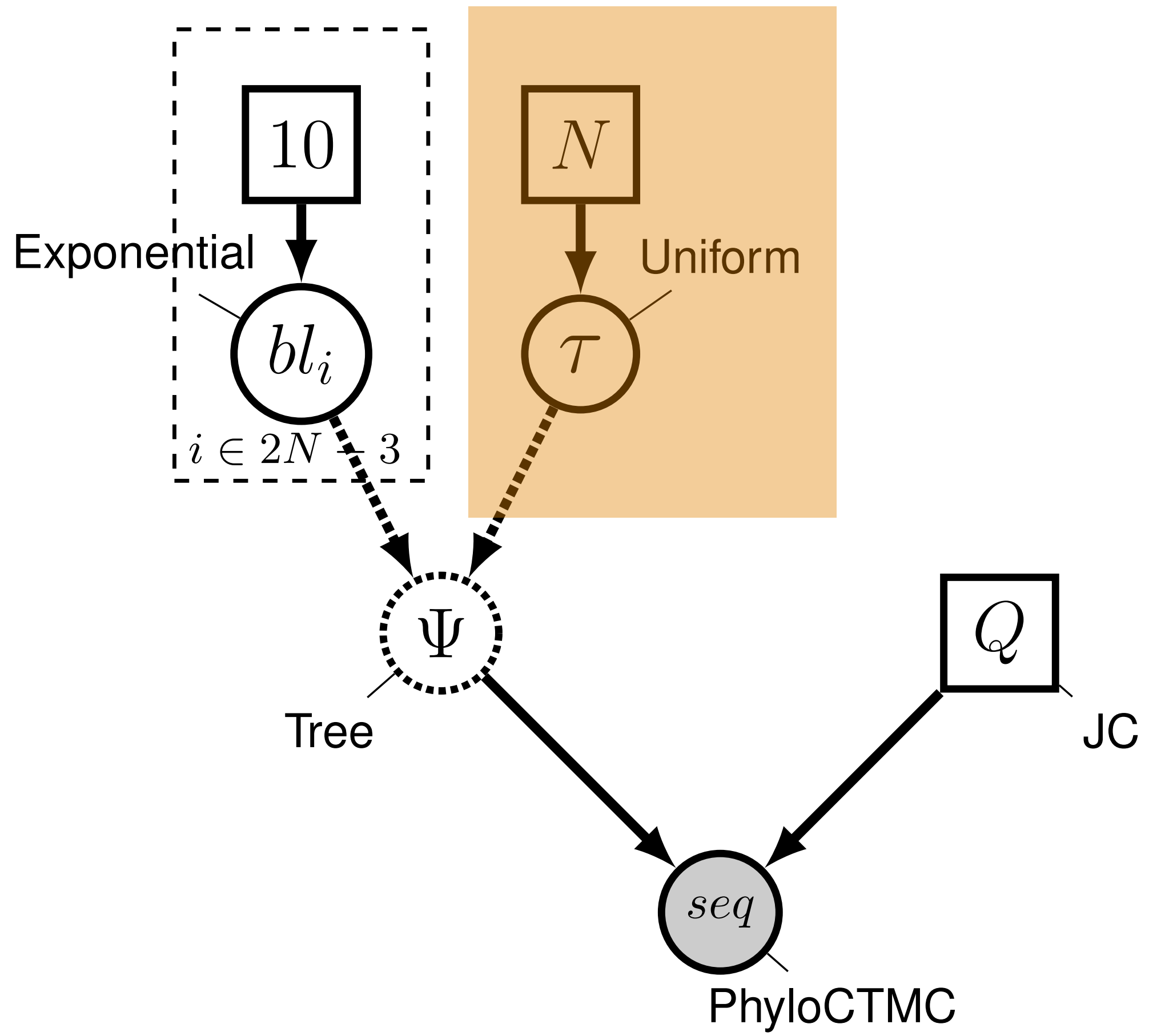
Note these are not the only approaches to tree-building but they are the most widely used

Bayesian tree inference

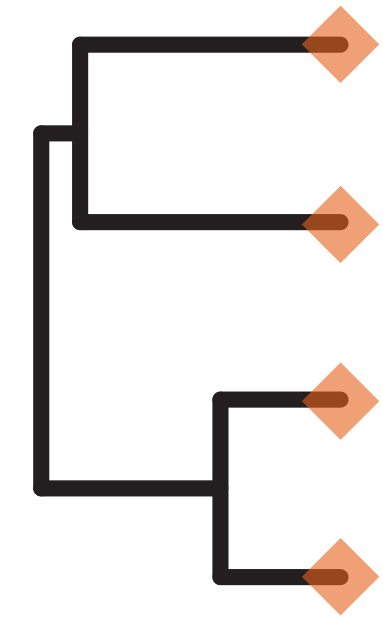
$$\begin{array}{c} \text{posterior} \\ \boxed{\phantom{\text{posterior}}} \\ P(\text{tree} \mid \text{data}) \end{array} = \frac{\begin{array}{c} \text{likelihood} \\ \boxed{\phantom{\text{likelihood}}} \\ P(\text{data} \mid \text{tree}) \end{array} \begin{array}{c} \text{priors} \\ \boxed{\phantom{\text{priors}}} \\ P(\text{tree}) \end{array}}{\begin{array}{c} \text{marginal probability} \\ \boxed{\phantom{\text{marginal probability}}} \\ P(\text{data}) \end{array}}$$

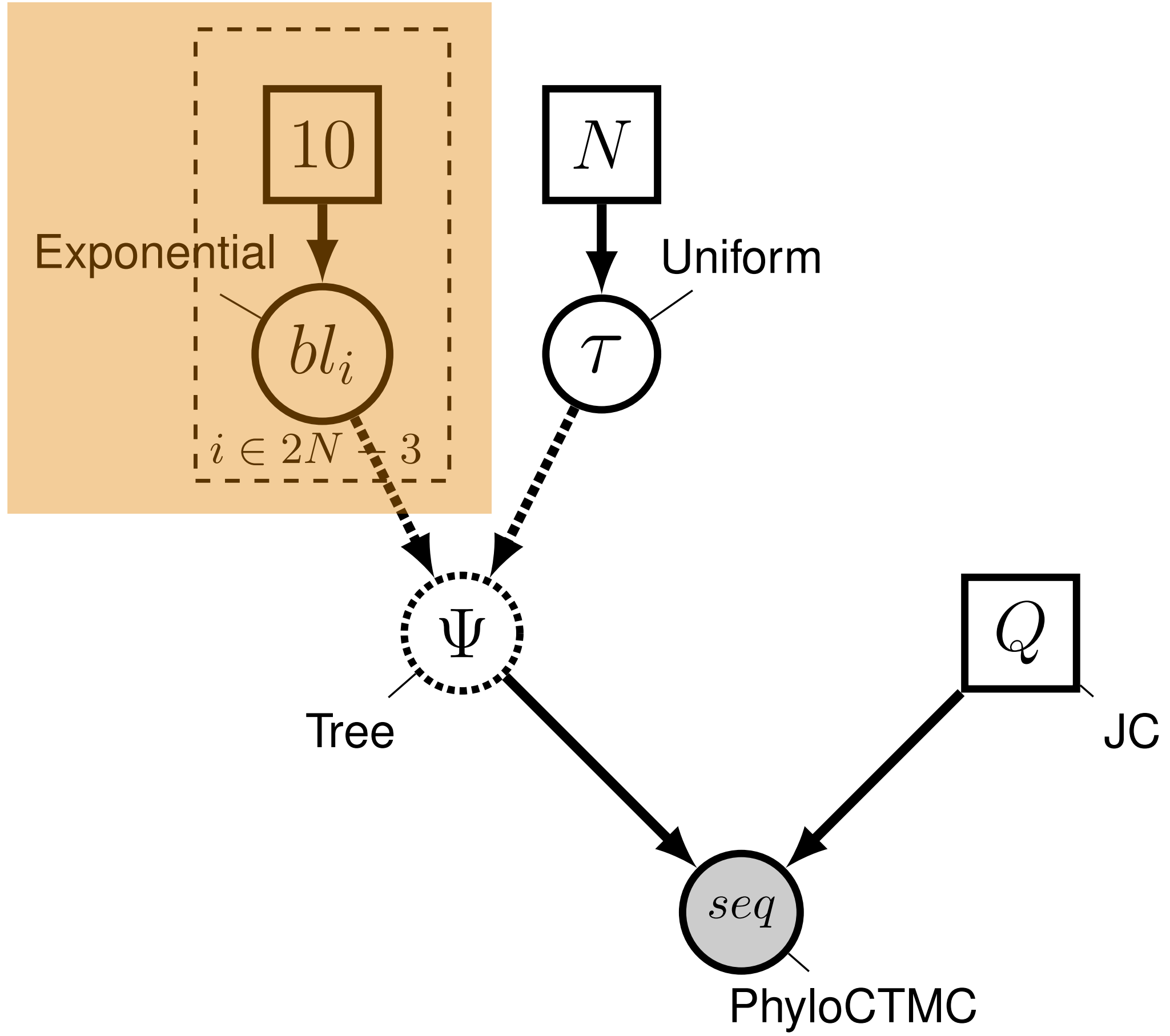


s

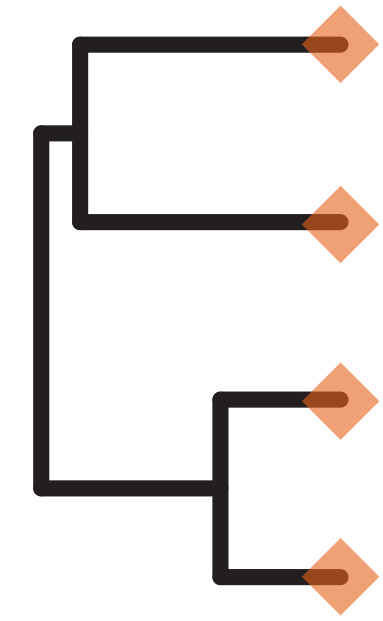


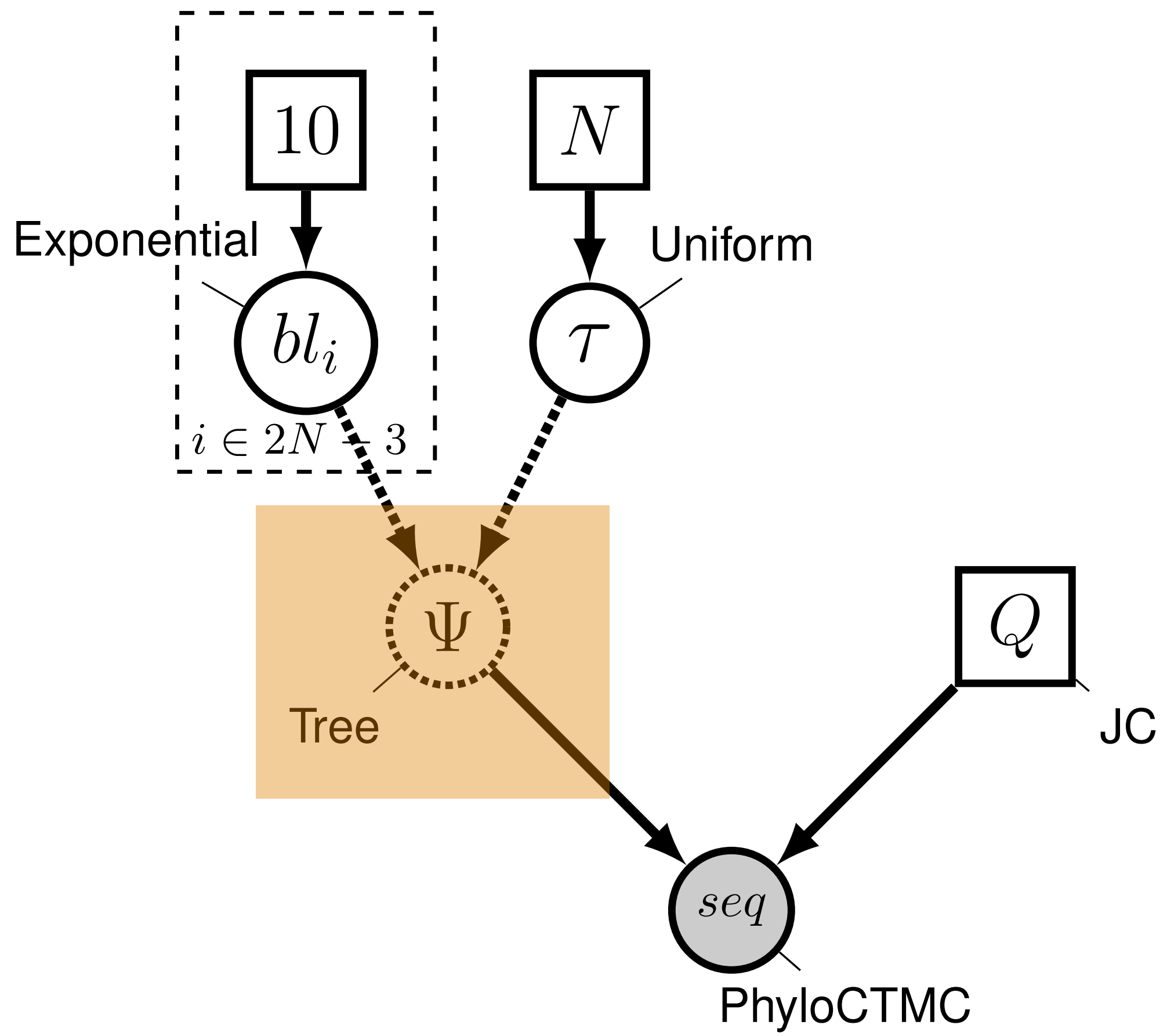
prior on the tree topology



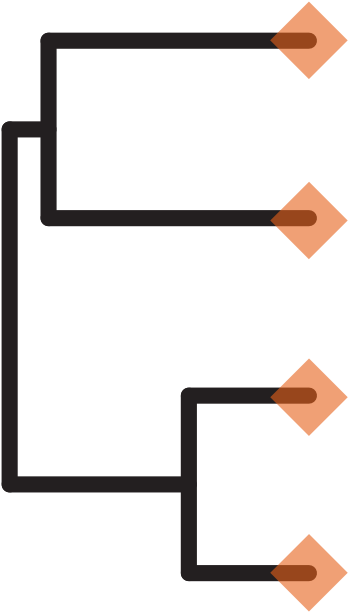


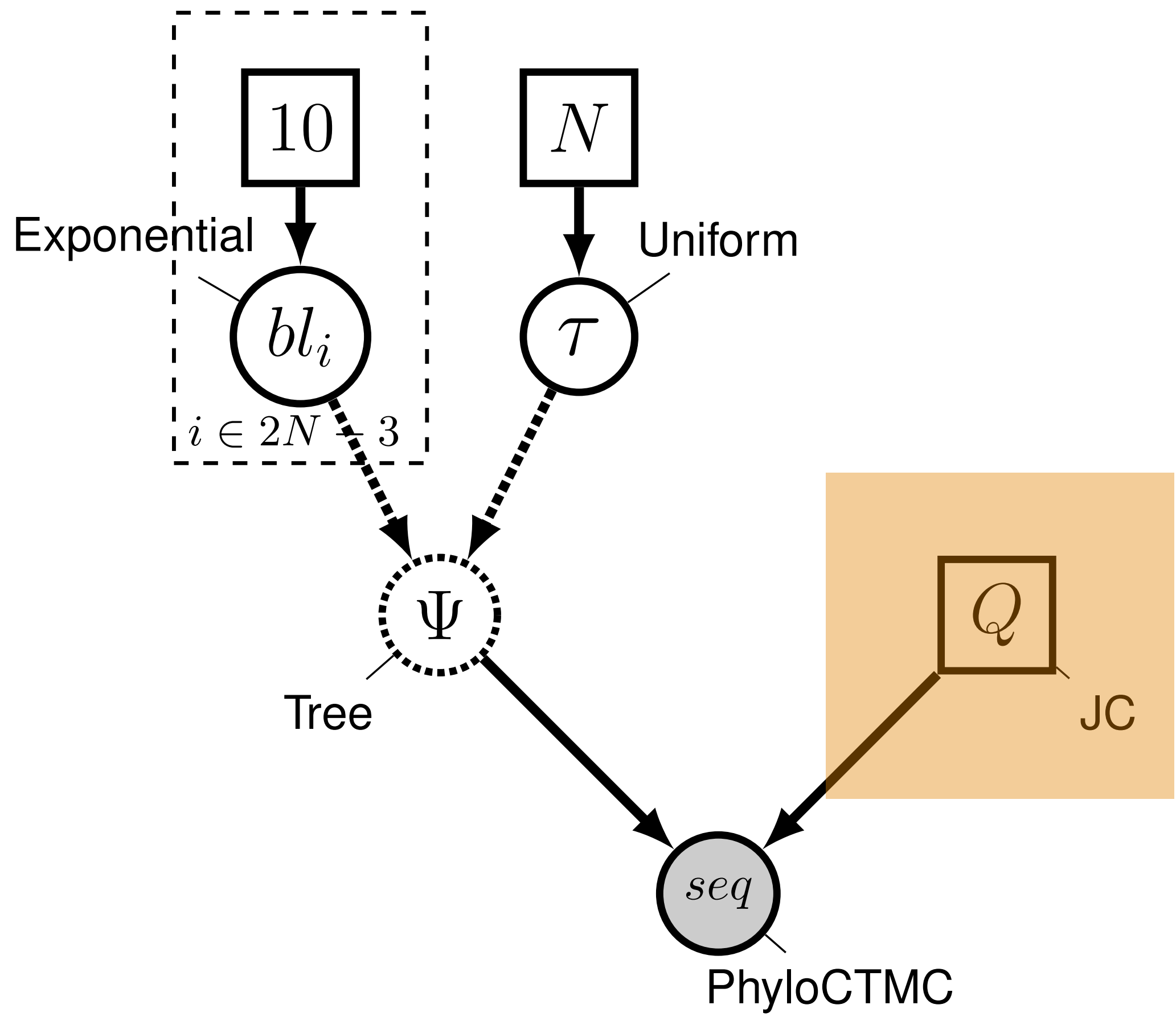
prior on the branch lengths



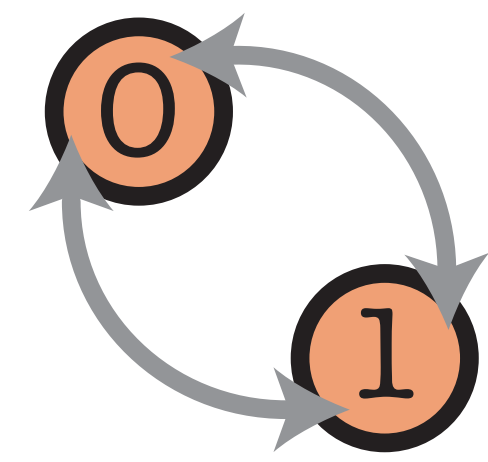


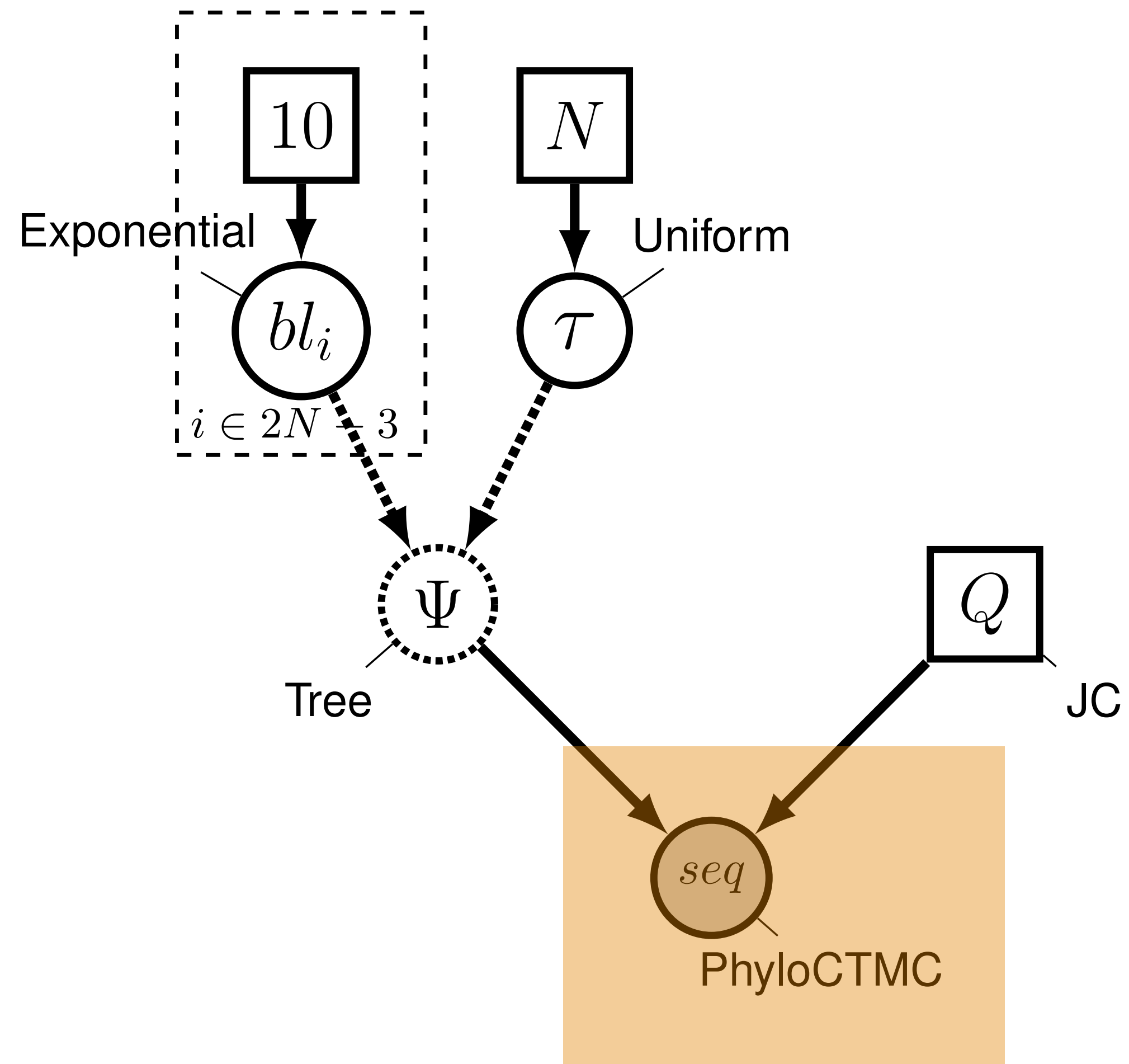
we can combine the topology and branch lengths



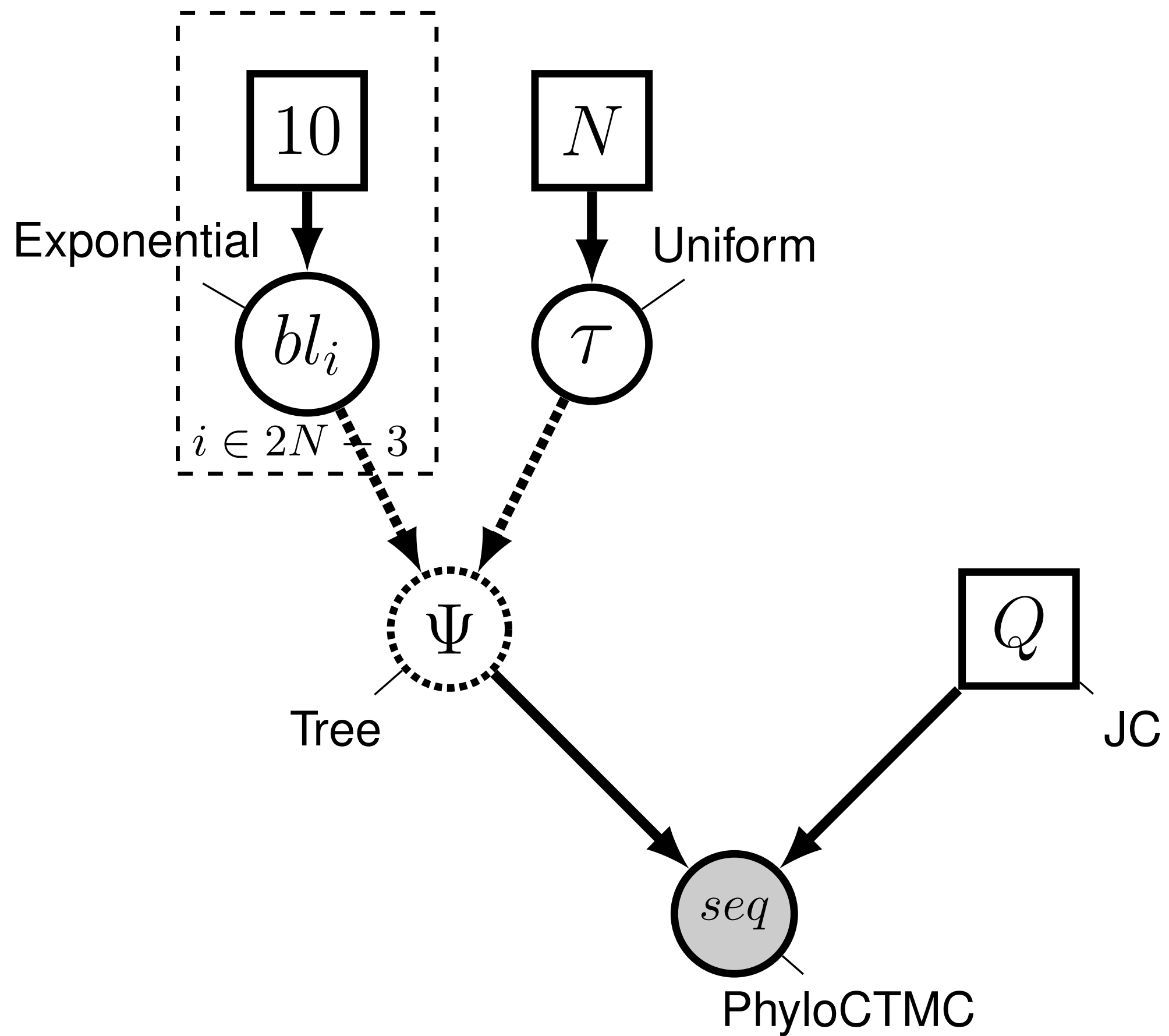


Substitution model





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



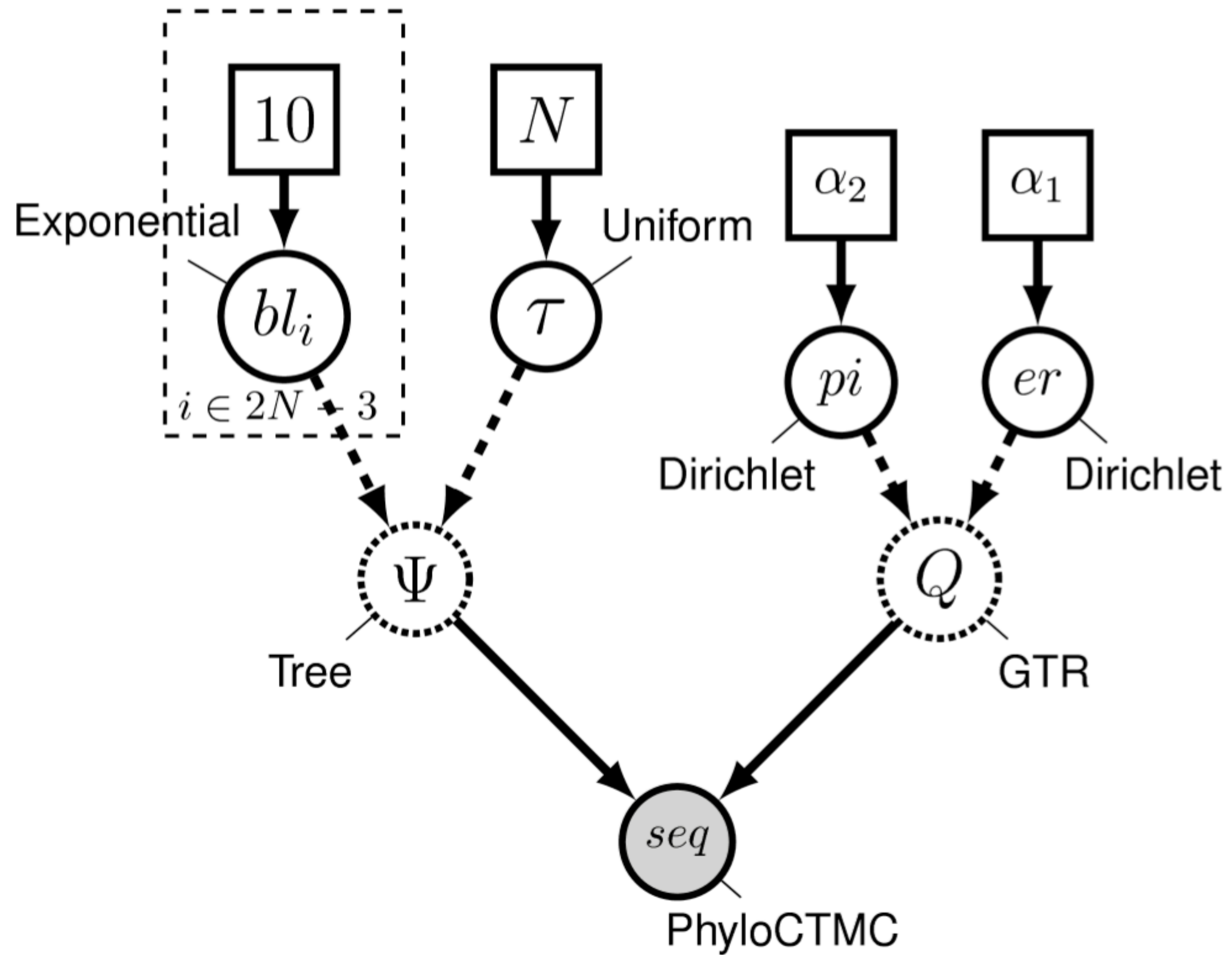
```

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)

```

Introduction to MCMC

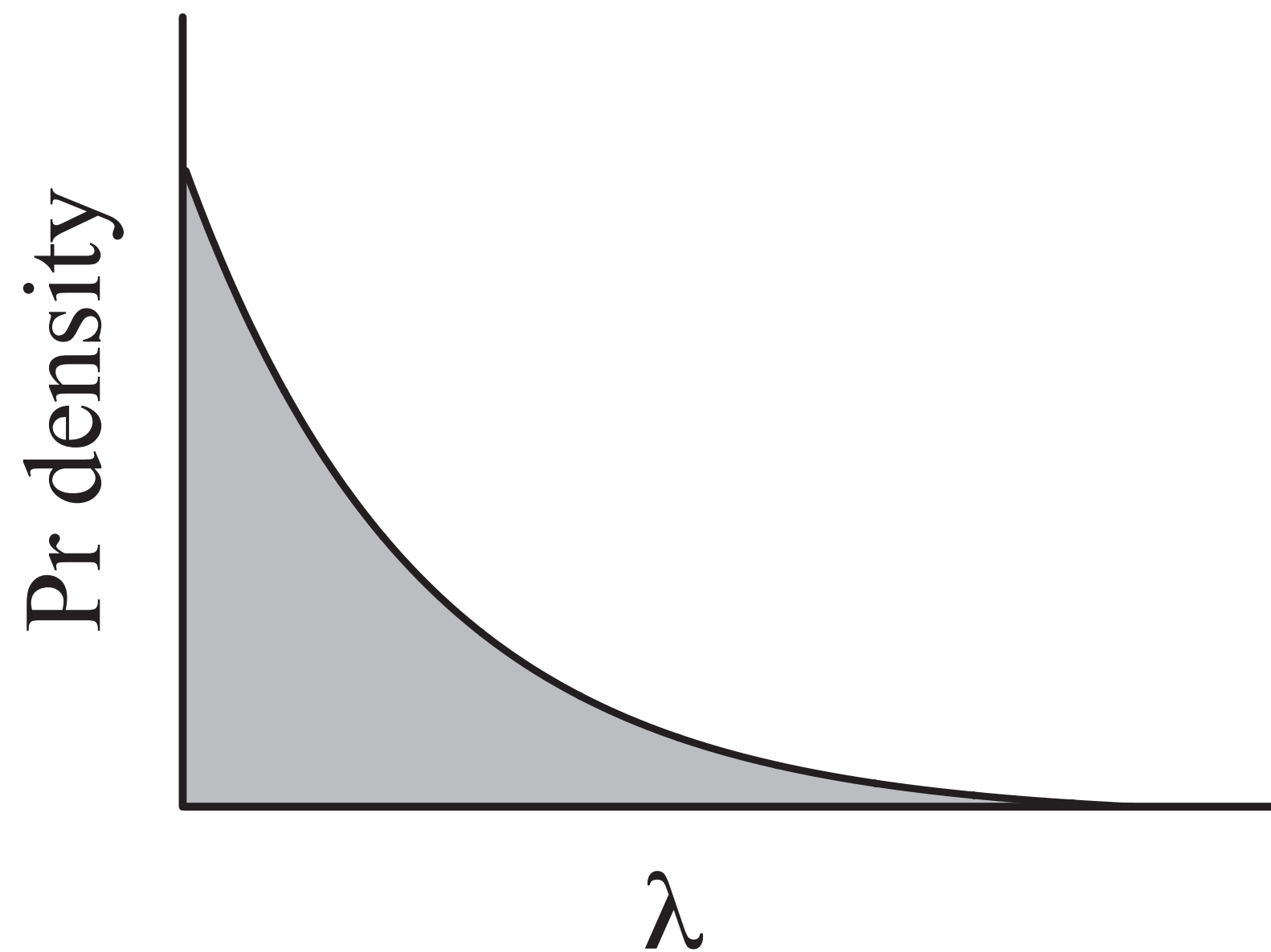
Probabilities vs probability densities

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

Probabilities vs probability densities



λ is drawn from an exponential distribution with mean δ

The x-axis represents the value of our parameter λ

The y-axis does have a value but it is not so easily interpretable

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

Bayesian tree inference

$$\begin{array}{c} \text{posterior} \\ \boxed{\phantom{\text{posterior}}} \\ P(\text{tree} \mid \text{data}) \end{array} = \frac{\begin{array}{c} \text{likelihood} \\ \boxed{\phantom{\text{likelihood}}} \\ P(\text{data} \mid \text{tree}) \end{array} \begin{array}{c} \text{priors} \\ \boxed{\phantom{\text{priors}}} \\ P(\text{tree}) \end{array}}{\begin{array}{c} \text{marginal probability} \\ \boxed{\phantom{\text{marginal probability}}} \\ P(\text{data}) \end{array}}$$

Bayesian tree inference

$$= \frac{P\left(\begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix} \mid \begin{matrix} \text{tree} \\ \text{0} \\ \text{1} \end{matrix}\right) P\left(\begin{matrix} \text{tree} \\ \text{0} \\ \text{1} \end{matrix}\right)}{\int P\left(\begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix} \mid \begin{matrix} \text{tree} \\ \text{0} \\ \text{1} \end{matrix}\right) P\left(\begin{matrix} \text{tree} \\ \text{0} \\ \text{1} \end{matrix}\right) d\begin{matrix} \text{tree} \\ \text{0} \\ \text{1} \end{matrix}}$$

this part is incredibly difficult to calculate!

What is Markov chain Monte Carlo (MCMC)?

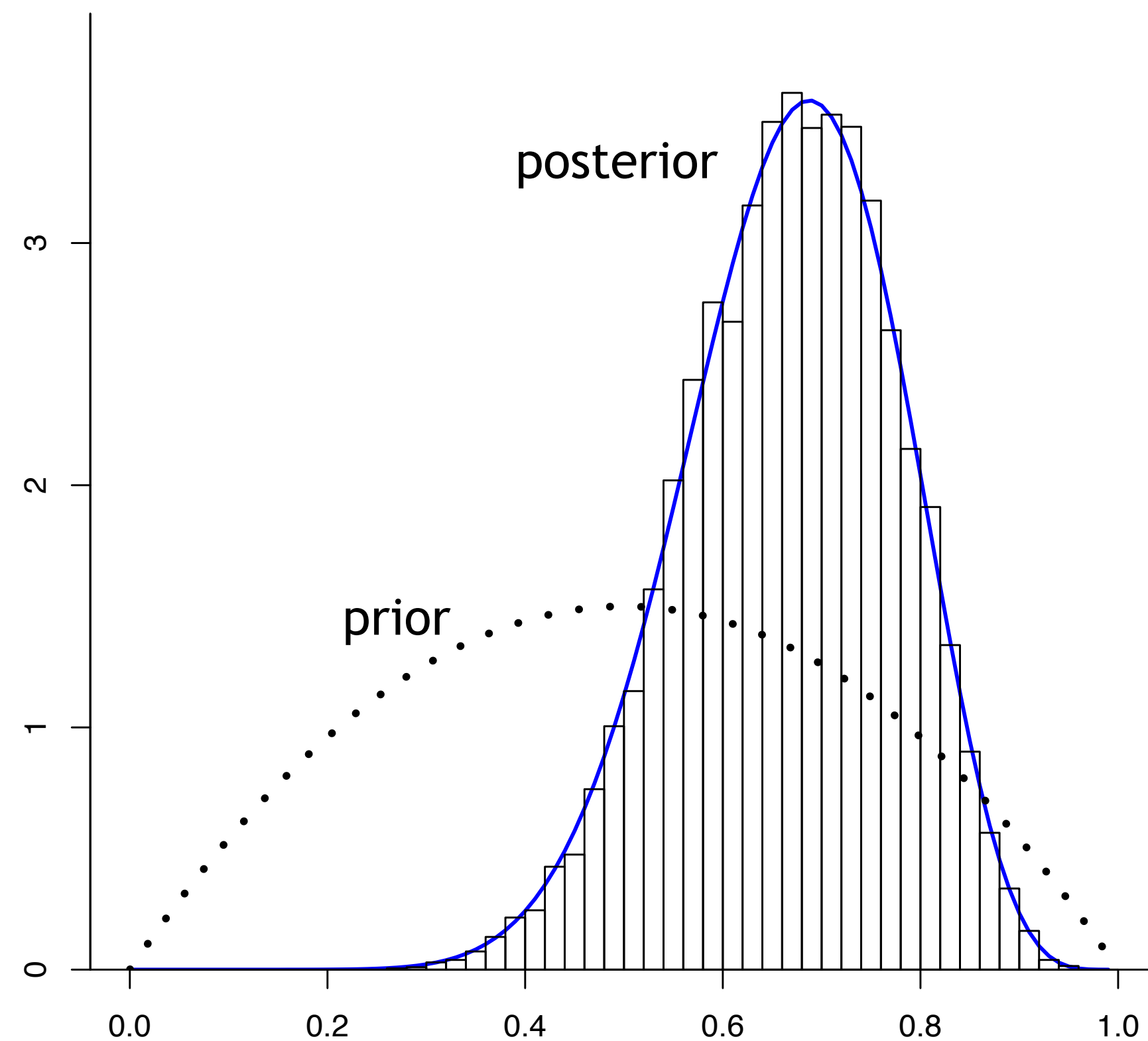
A group of algorithms for approximating the posterior distribution (also known as samplers)

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)?

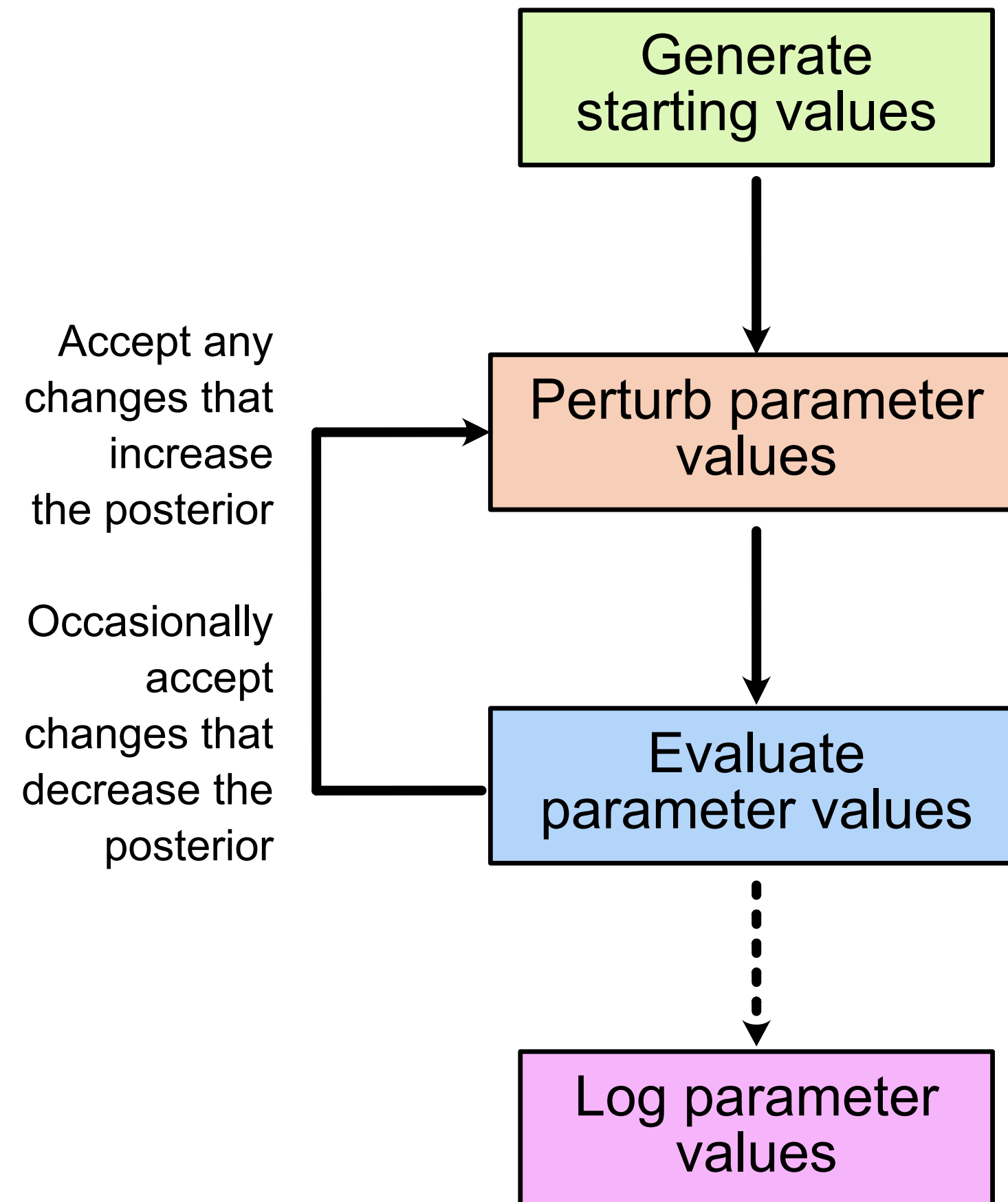


Copyright © 2018 Paul O. Lewis

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

The Metropolis-Hastings algorithm

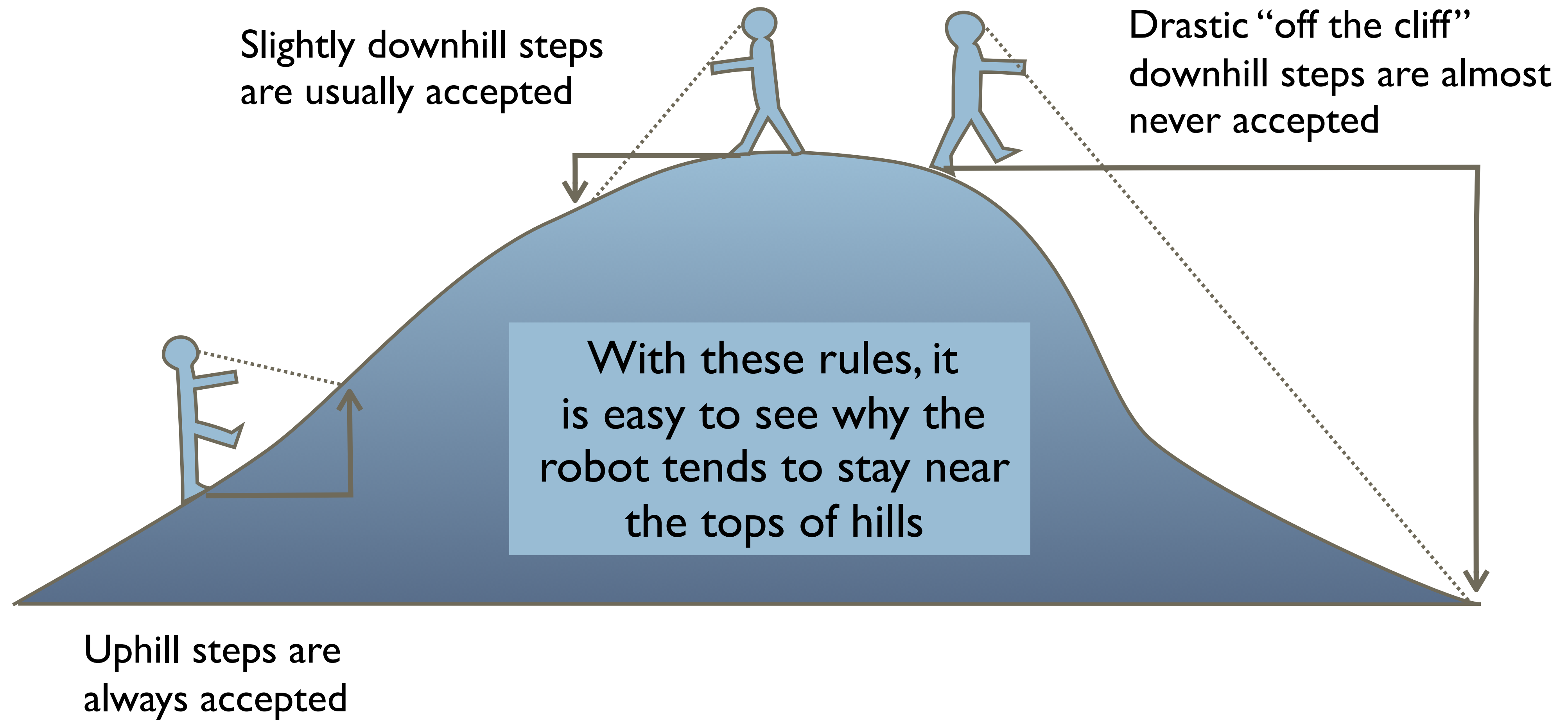
Flowchart

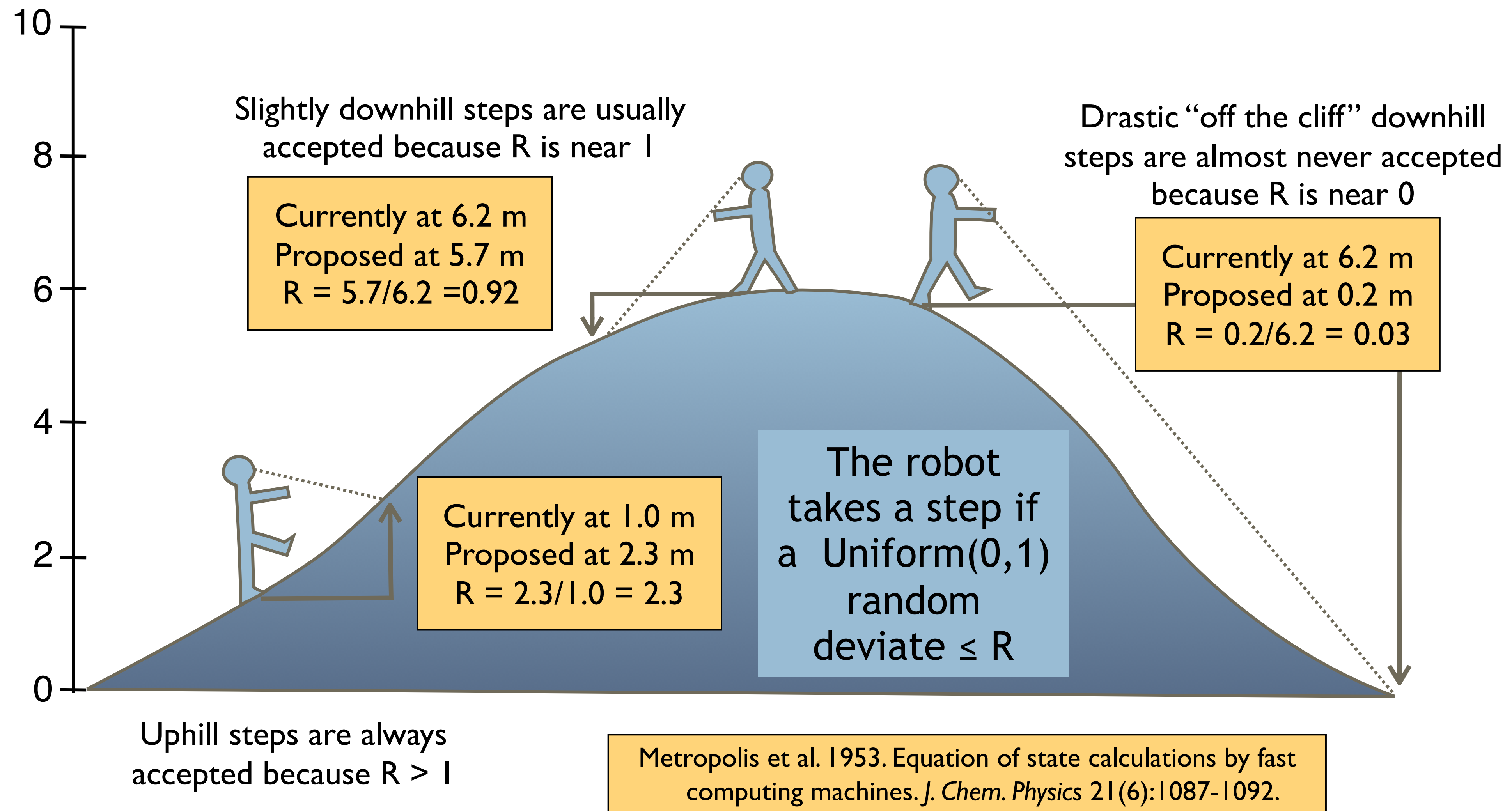


Pseudocode

```
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





When calculating the ratio (R) of posterior densities, the marginal probability of the data cancels.

$$\frac{p(\theta^* | D)}{p(\theta | D)} = \frac{\frac{p(D | \theta^*) p(\theta^*)}{\cancel{p(D)}}}{\frac{p(D | \theta) p(\theta)}{\cancel{p(D)}}} = \frac{p(D | \theta^*) p(\theta^*)}{p(D | \theta) p(\theta)}$$

Posterior
odds

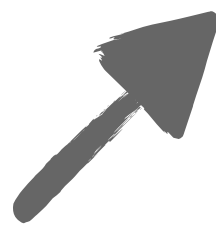
Apply Bayes' rule to
both top and bottom

Likelihood
ratio

Prior
odds

Hastings ratio

new parameter values



$$R = \frac{P(\text{Diagram with } * \mid \begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix})}{P(\text{Diagram without } * \mid \begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix})}$$

=

$$\frac{P(\begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix} \mid \text{Diagram with } *) P(\text{Diagram with } *)}{P(\begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix})}$$

$$\frac{P(\begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix} \mid \text{Diagram without } *) P(\text{Diagram without } *)}{P(\begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix})}$$

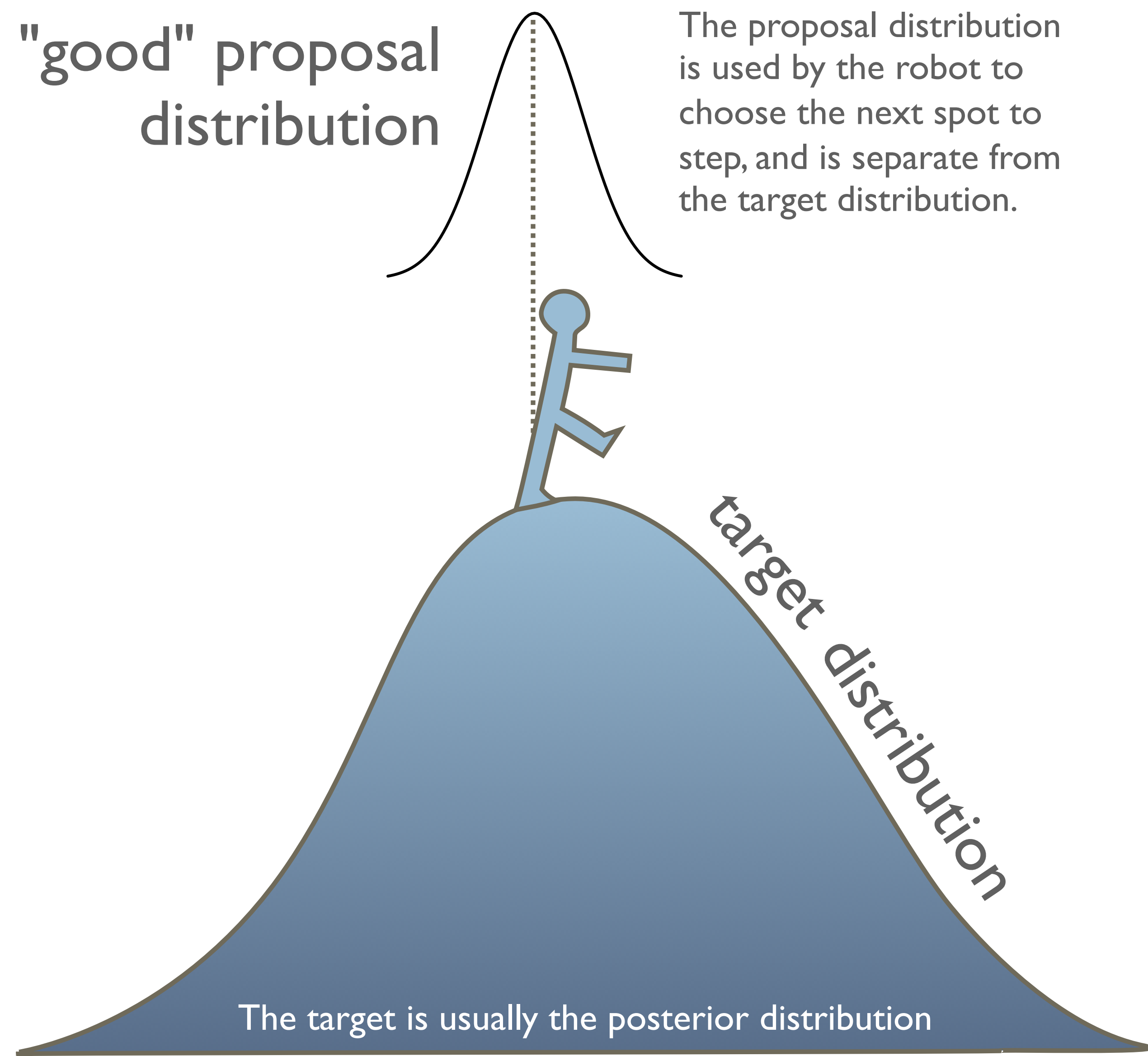
=

$$\frac{P(\begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix} \mid \text{Diagram with } *) P(\text{Diagram with } *)}{P(\begin{matrix} 0101\dots \\ 1101\dots \\ 0100\dots \end{matrix} \mid \text{Diagram without } *) P(\text{Diagram without } *)}$$

The marginal probability of the data cancels out

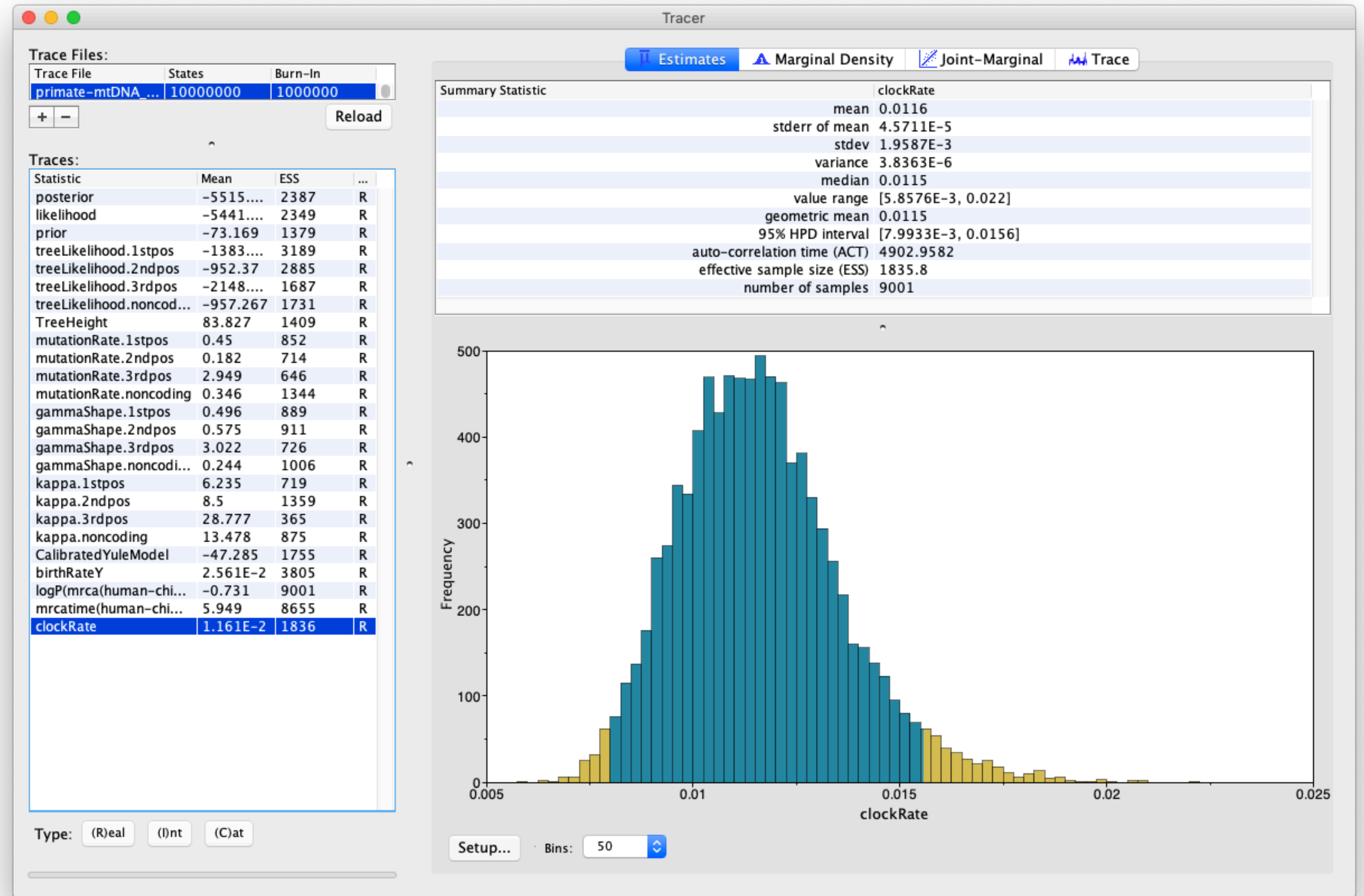
All we're left to calculate is the likelihood ratio and the prior odds ratio

Proposals



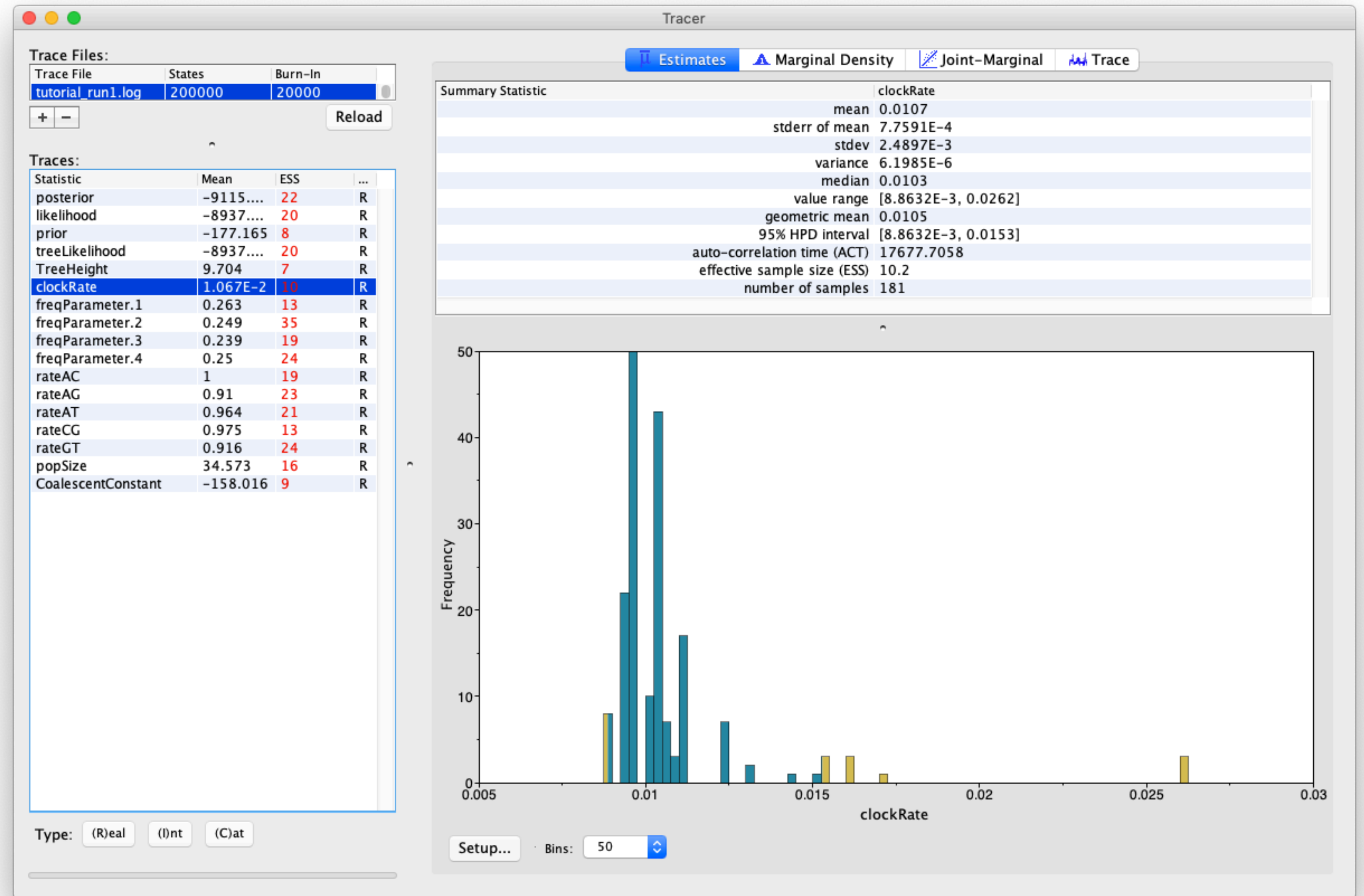
Summarising the posterior

Tracer is an amazing program for exploring MCMC output



Summarising the posterior

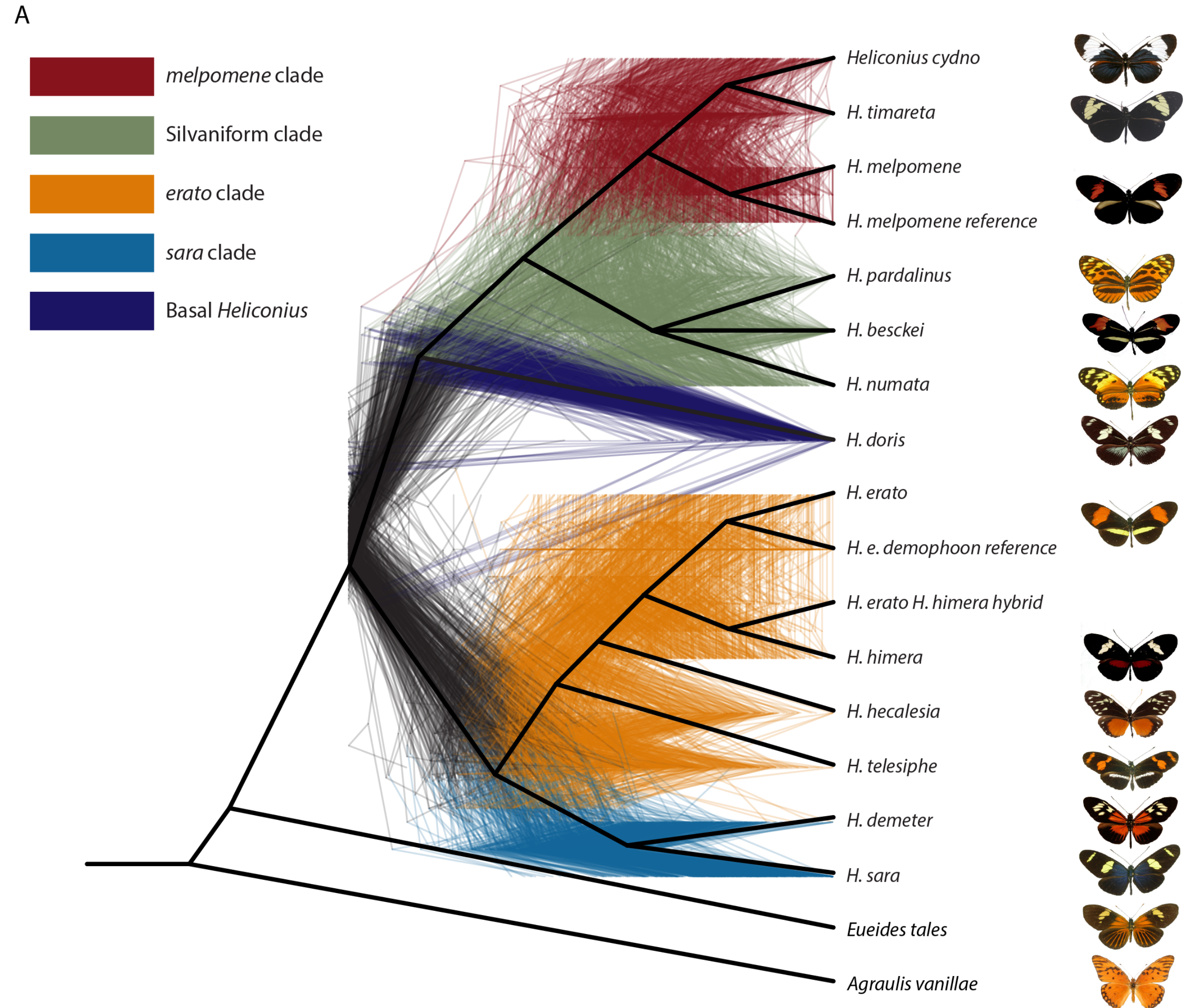
Tracer is an amazing program for exploring MCMC output



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