# Introduction to Bioinformatics for Computer Scientists
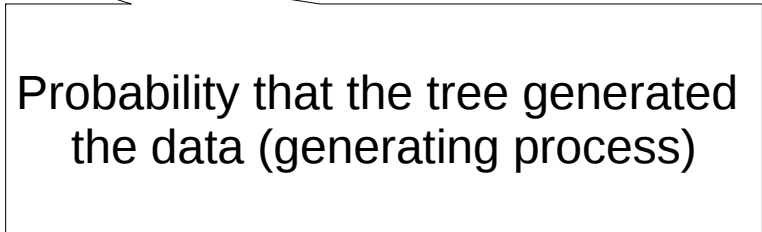
## Lecture 9b

# Likelihood

- Given:
  - MSA
  - Tree topology with branch lengths
  - Model
  - We can calculate $P_{x \to z}(b)$ for a branch length (or time) $b$
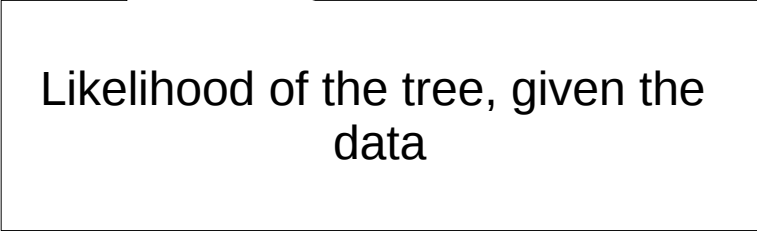
# Likelihood

- $L(T|D) = P(D|T)$

Probability that the tree generated the data (generating process)

# Likelihood

- $L(T|D) = P(D|T)$

Likelihood of the tree, given the data

# Likelihood

- L(T|D) = P(D|T)

**Likelihood:** 10 coin flips → 10 heads
What's the likelihood that the coin is fair?

**Probability:** Probability of landing heads up 10 times

# Likelihood

- $L(T|D) = P(D|T)$

- $L(T|D) = \Pi\ P(s_i|T)$

Alignment site $i$

# Likelihood

- $L(T|D) = P(D|T)$

- $L(T|D) = \Pi\, P(s_i|T)$

Alignment site *i*

What is problematic about this term?

# Likelihood

- $L(T|D) = P(D|T)$
- $L(T|D) = \prod P(s_i|T)$
- $\log(L(T|D)) = \sum \log(P(s_i|T))$

# Likelihood

- $L(T|D) = P(D|T)$

- $L(T|D) = \Pi\ P(s_i|T)$

- $\log(L(\textcolor{red}{T}|D)) = \Sigma\ \log(P(s_i|T))$

This is the model
1. Tree topology
2. Branch lengths
3. Model of nucleotide substitution
 → generally lumped into parameter vector **Θ**: $L(\textbf{Θ}|D)$

# Likelihood

- $L(T|D) = P(D|T)$

- $L(T|D) = \Pi\ P(s_i|T)$

- $\log(L(\textcolor{red}{T}|D)) = \Sigma\ \log(P(s_i|T))$

This is the model
1. Tree topology
2. Branch lengths
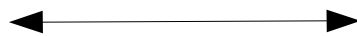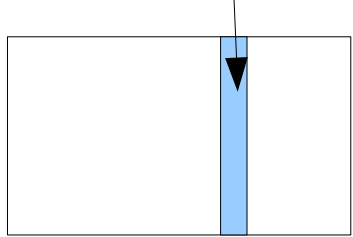3. Model of nucleotide substitution
   → generally lumped into parameter vector $\Theta$: $L(\Theta|D)$

How do we compute this?

# Likelihood of a Tree

- We assume that sites evolve independently

Likelihood of site $i$

MSA length $n$

# Likelihood of a Tree

- We assume that sites evolve independently

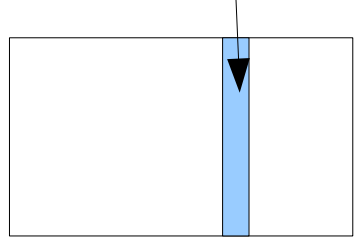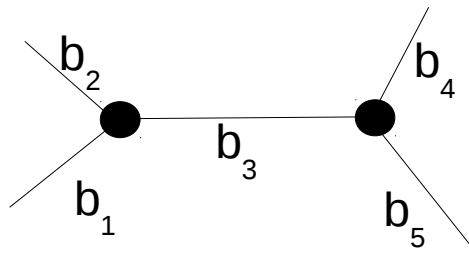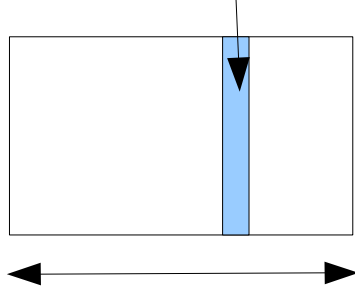Likelihood of site $i$



$b_2$

$b_4$

$b_3$

$b_1$

$b_5$

MSA length $n$

# Likelihood of a Tree

- We assume that sites evolve independently

Likelihood of site $i$

MSA length $n$

$b_2$

$b_4$

$b_3$

$b_1$

$b_5$

Model $M$

$P_{ij}(t)$

# Likelihood of a Tree

- We assume that sites evolve independently

Likelihood of site $i$

$b_2$  $b_4$

$b_3$

$b_1$  $b_5$

Model $M$
$P_{ij}(t)$

MSA length $n$

- Overall likelihood: $L := \Pi L_i$

# Likelihood of a Tree

- We assume that sites evolve independently

Likelihood of site $i$

b$_2$    b$_4$

b$_3$

b$_1$    b$_5$

Model ***M***
$P_{ij}(t)$

MSA length $n$

- Overall likelihood: $L := \Pi L_i$

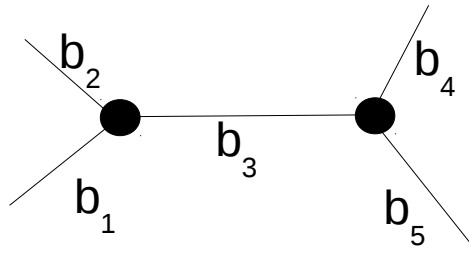- $P_{ij}(t)$ $i,j$ in $\{A, C, G, T\}$

Branch length/time
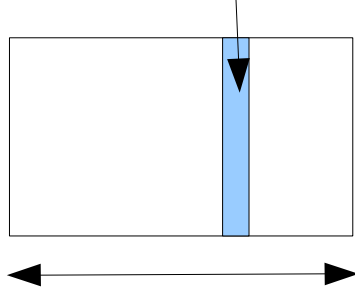
# Likelihood of a Tree

- We assume that sites evolve independently



Likelihood of site $i$

MSA length $n$

Model $M$
$P_{ij}(t)$

- Overall likelihood: $L := \Pi\, L_i$

- $P_{ij}(t)$ i,j in {A, C, G, T}

  → Probability of being in state $j$ after time $t$

  → We assume that $P_{ij}(t)$ is a Markov Process
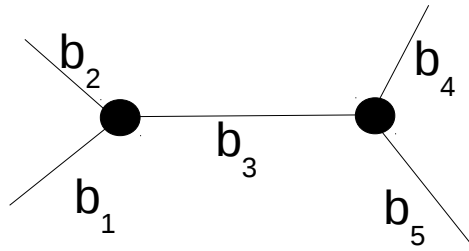
# Likelihood of a Tree

- We assume that sites evolve independently

Likelihood of site *i*



MSA length *n*

Model **M**
$P_{ij}(t)$

- Overall likelihood: $L := \Pi\, L_i$
- $P_{ij}(t)$ i,j in {A, C, G, T}
  - → Probability of being in state *j* after time *t*
  - → We assume that $P_{ij}(t)$ is a Markov Process
- Equilibrium frequency vector π = ($π_A$, $π_C$, $π_G$, $π_T$)
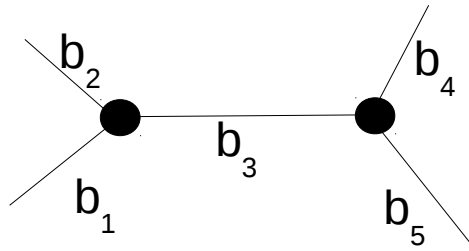
# Likelihood of a Tree

- We assume that sites evolve independently

Likelihood of site $i$



MSA length $n$

Model $M$
$P_{ij}(t)$

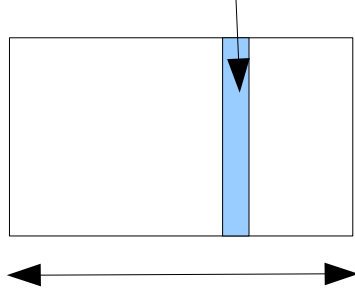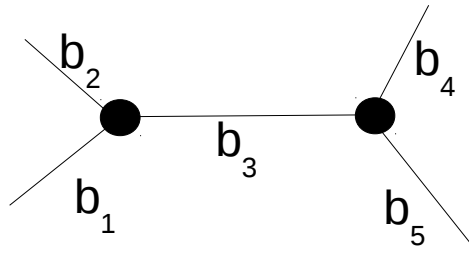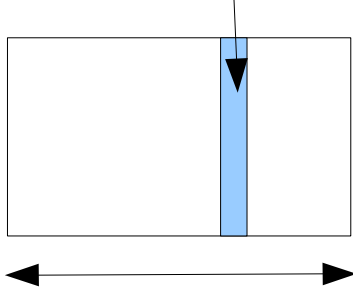- Overall likelihood: $L := \Pi\, L_i$
- $P_{ij}(t)$ $i,j$ in $\{A, C, G, T\}$
  - → Probability of being in state $j$ after time $t$
  - → We assume that $P_{ij}(t)$ is a Markov Process
- Equilibrium frequency vector $\pi = (\pi_A, \pi_C, \pi_G, \pi_T)$
- **Time reversibility:** $\pi_i P_{ij}(t) = \pi_j P_{ij}(t)$

# What's the likelihood of this tree?

# What's the likelihood of this tree?

# What's the likelihood of this tree?

# What's the likelihood of this tree?

Assume the inner states are given!
What is the likelihood of the tree if we
Interpret it as **Markov** diagram?



22

# What's the likelihood of this tree?

Assume the inner states are given!
What is the likelihood of the tree if we
Interpret it as **Markov** diagram?

$$L_i = \pi\, P_{AA}(b_1)\, P_{AA}(b_2)\, P_{AA}(b_3)\, P_{AT}(b_4) P_{TT}(b_5)\, P_{TG}(b_6)$$

# What's the likelihood of this tree?

Assume the inner states are given!
What is the likelihood of the tree if we
Interpret it as **Markov** diagram?

$L_i = \pi_A \, P_{AA}(b_1) \, P_{AA}(b_2) \, P_{AA}(b_3)$
$P_{AT}(b_4) P_{TT}(b_5) \, P_{TG}(b_6)$

*We are multiplying here, because to
observe the data at the tips, given the
tree, the initial state must be* **A** $\pi_A$

**A**

b1

b4

**A**

**T**

b2

b3

b5

b6

A

A

T

G

# What's the likelihood of this tree?

Assume the inner states are given!
What is the likelihood of the tree if we
Interpret it as **Markov** diagram?

$$L_i = \pi_A \, P_{AA}(b_1) \, P_{AA}(b_2) \, P_{AA}(b_3) \, P_{AT}(b_4) P_{TT}(b_5) \, P_{TG}(b_6)$$

*We are multiplying here, because to observe the data at the tips, given the tree, the initial state must be* **A** $\pi_A$

**AND** then this happened

# What's the likelihood of this tree?



Assume the inner states are given!
What is the likelihood of the tree if we
Interpret it as **Markov** diagram?

$$L_i = \pi_A \, P_{AA}(b_1) \, P_{AA}(b_2) \, P_{AA}(b_3) \, P_{AT}(b_4) P_{TT}(b_5) \, P_{TG}(b_6)$$

*We are multiplying here, because to observe the data at the tips, given the tree, the initial state must be* **A** $\pi_A$

**AND** then this happened
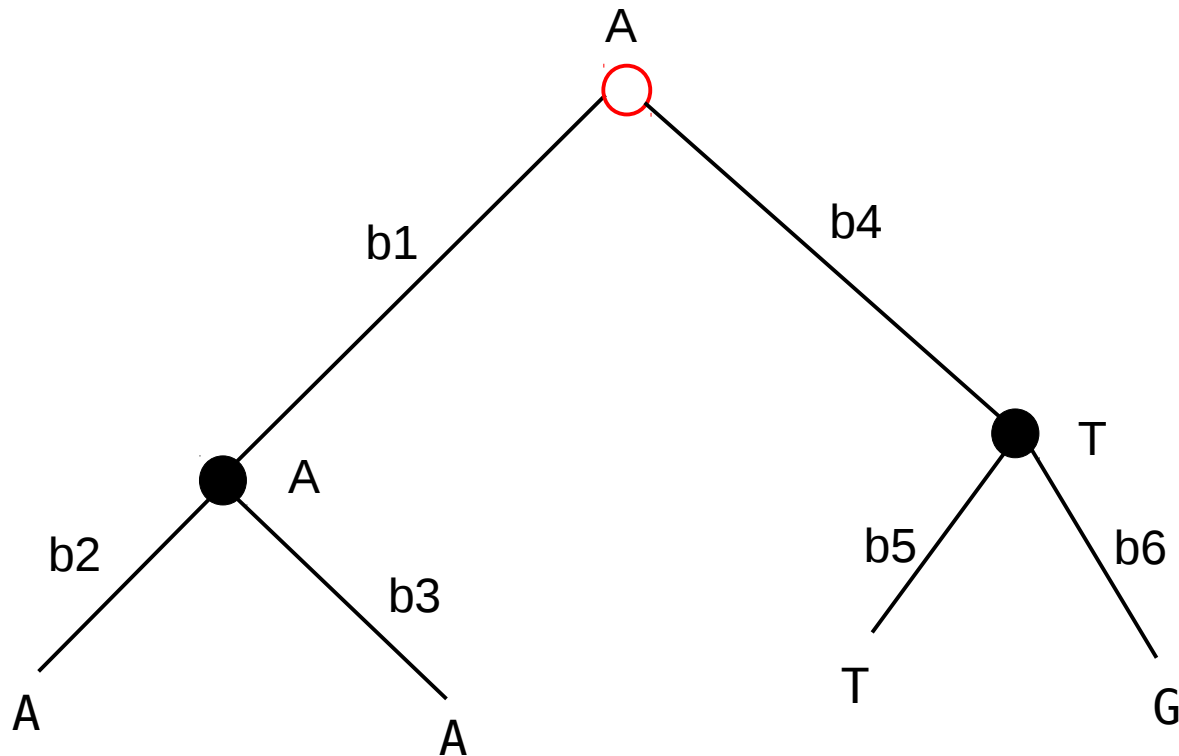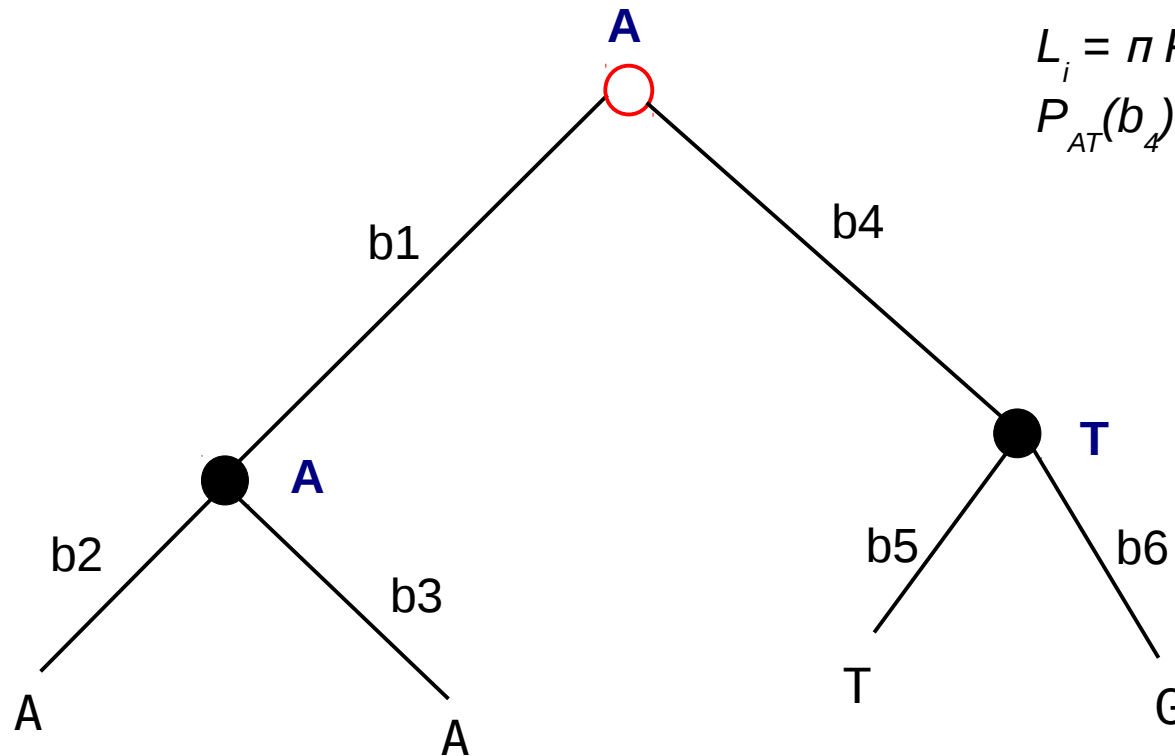**AND** this

# What's the likelihood of this tree?

Assume the inner states are given!
What is the likelihood of the tree if we
Interpret it as **Markov** diagram?

$$L_i = \pi_A \, P_{AA}(b_1) \, P_{AA}(b_2) \, P_{AA}(b_3) \, P_{AT}(b_4) P_{TT}(b_5) \, P_{TG}(b_6)$$

*We are multiplying here, because to observe the data at the tips, given the tree, the initial state must be* **A** $\pi_A$

**AND** then this happened
**AND** this
**AND** this

**A**

b1

b4

**A**

**T**

b2

b3

b5

b6

A

A

T

G

# What's the likelihood of this tree?

Assume the inner states are given!
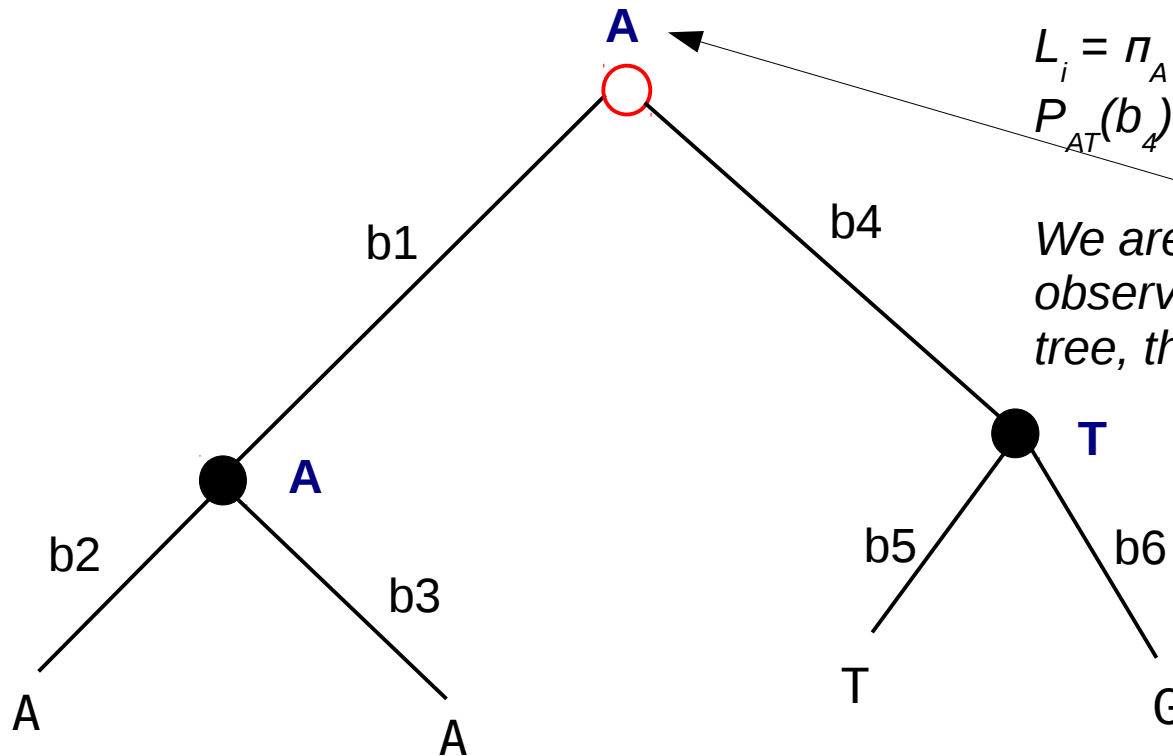What is the likelihood of the tree if we
Interpret it as **Markov** diagram?

$$L_i = \pi_A \, P_{AA}(b_1) \, P_{AA}(b_2) \, P_{AA}(b_3) \, P_{AT}(b_4) P_{TT}(b_5) \, P_{TG}(b_6)$$

*We are multiplying here, because to observe the data at the tips, given the tree, the initial state must be* **A** $\pi_A$

**AND** then this happened
**AND** this
**AND** this
**AND** this

A

b1          b4

A          T

b2     b3        b5     b6
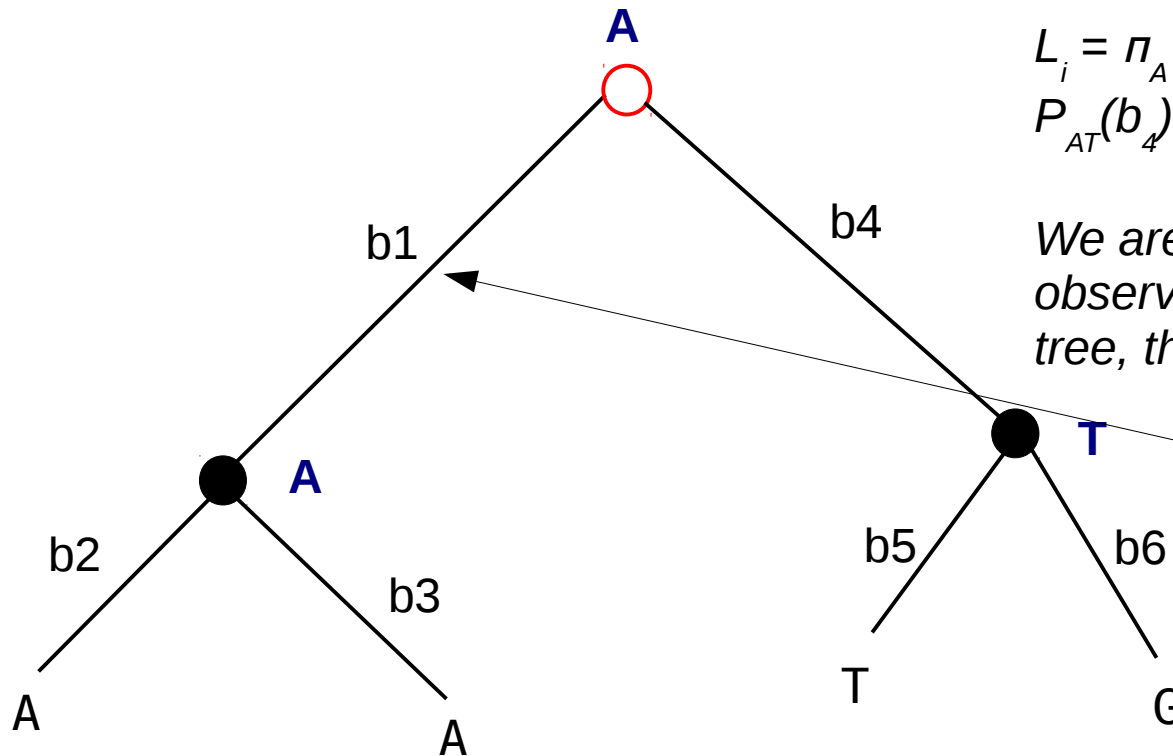
A          A      T          G

# What's the likelihood of this tree?

Assume the inner states are given!
What is the likelihood of the tree if we
Interpret it as **Markov** diagram?

$$L_i = \pi_A \, P_{AA}(b_1) \, P_{AA}(b_2) \, P_{AA}(b_3) \, P_{AT}(b_4) P_{TT}(b_5) \, P_{TG}(b_6)$$

*We are multiplying here, because to observe the data at the tips, given the tree, the initial state must be* **A** $\pi_A$

**A**

b1

b4

**A**

b2

b3

**T**

b5

b6

A

A

T

G

**AND** then this happened
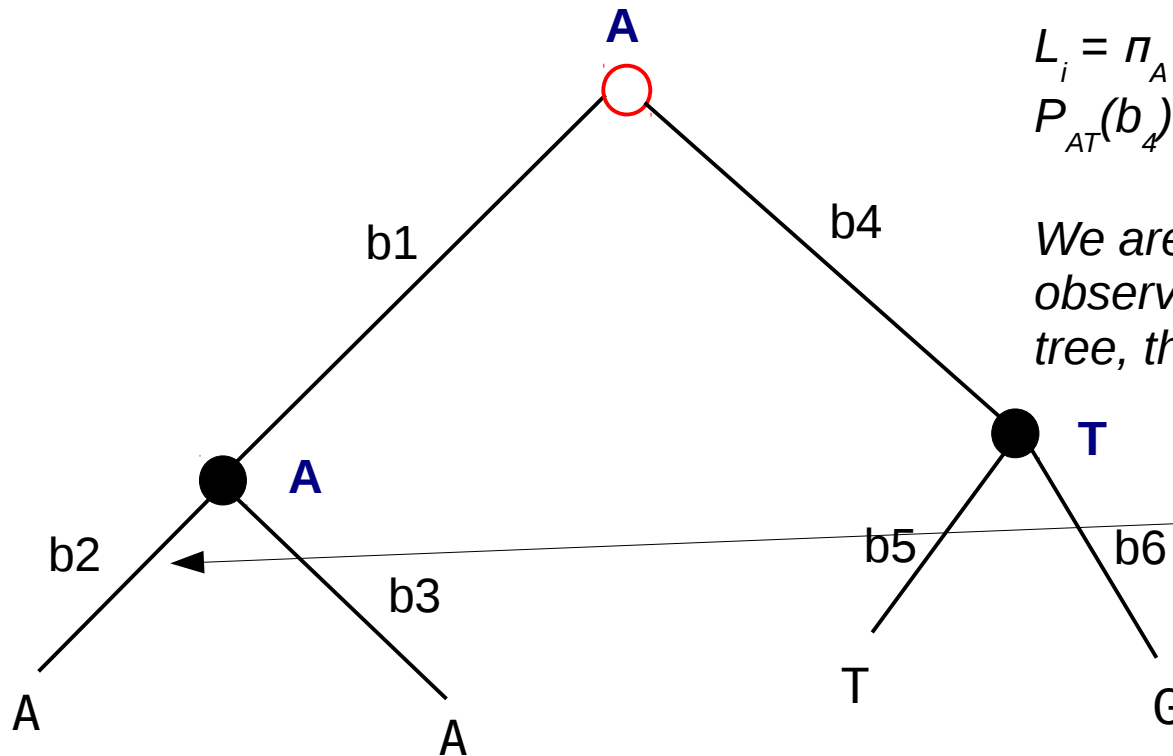**AND** this
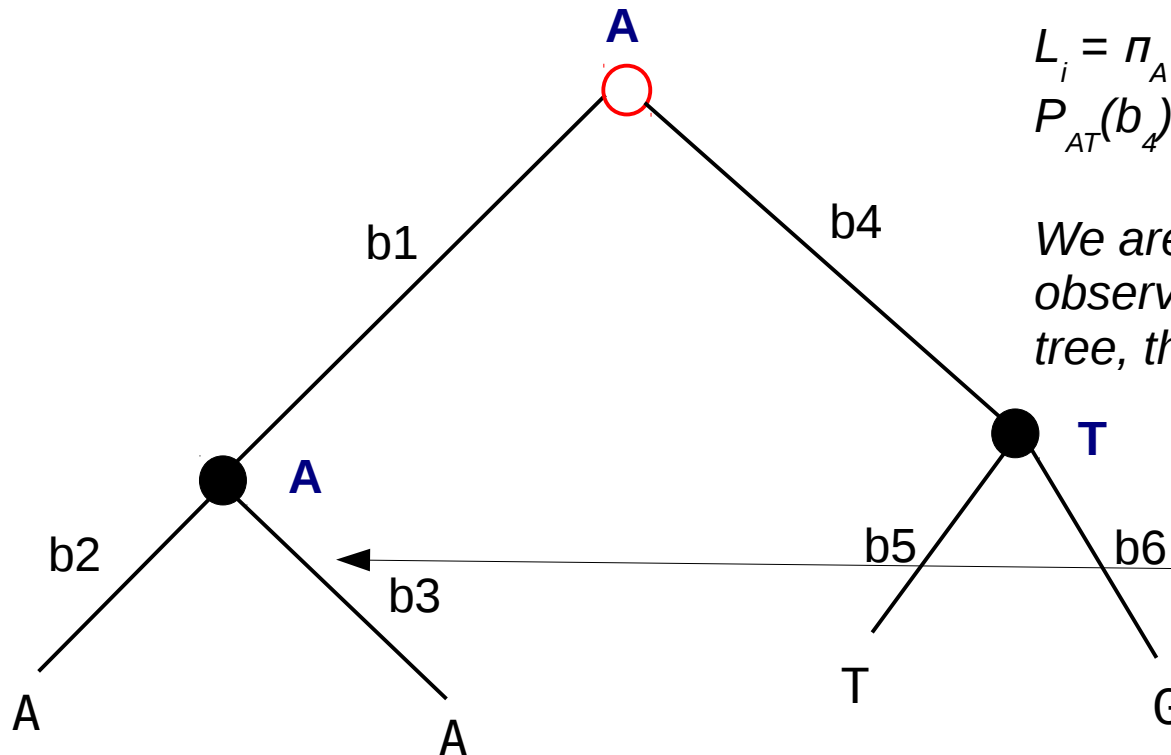**AND** this
**AND** this
**AND** this

# What's the likelihood of this tree?

Assume the inner states are given!
What is the likelihood of the tree if we
Interpret it as **Markov** diagram?

$$L_i = \pi_A \, P_{AA}(b_1) \, P_{AA}(b_2) \, P_{AA}(b_3) \, P_{AT}(b_4) P_{TT}(b_5) \, P_{TG}(b_6)$$

*We are multiplying here, because to observe the data at the tips, given the tree, the initial state must be* **A** $\pi_A$
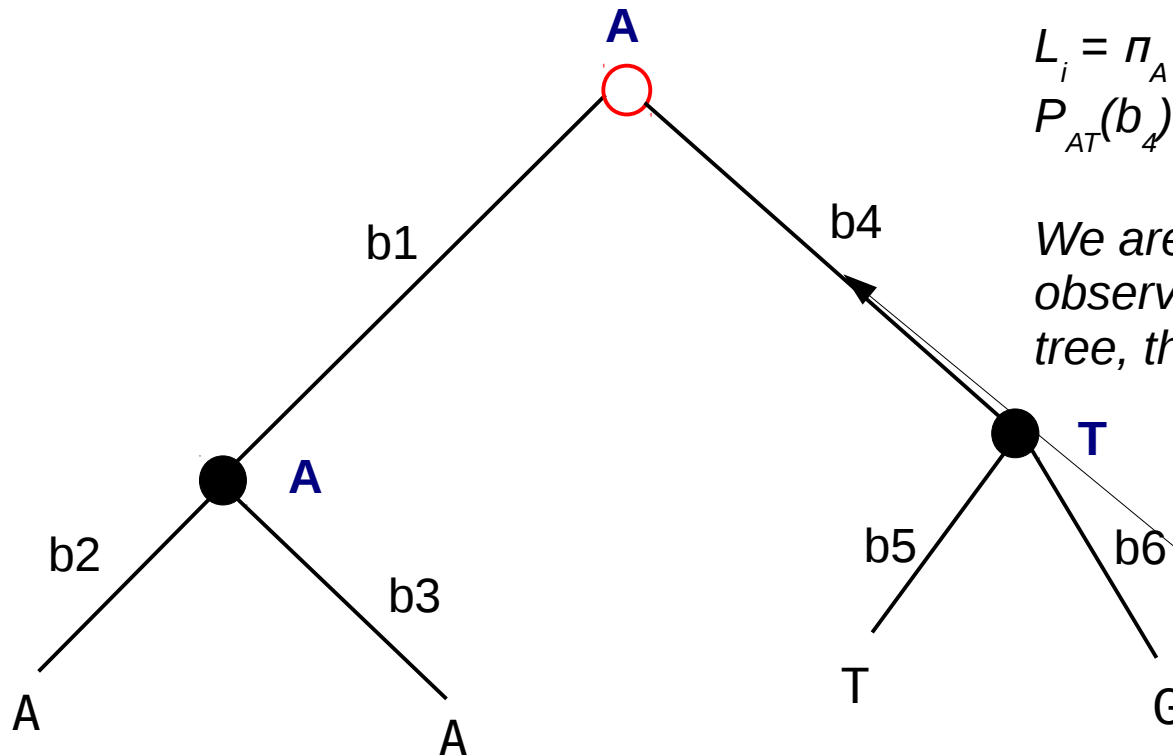
**AND** then this happened
**AND** this
**AND** this
**AND** this
**AND** this
**AND** this



**A**

b1

b4

**A**

**T**

b2

b3

b5

b6

A

A

T

G

# What's the likelihood of this tree?

However, we don't know the inner states :-(
So the question is: What are the possible
evolutionary histories that could have given
rise (generated) to the data we observe at
the tips?

**I1**

**I2**

**I3**

b1

b4

b2

b3

b5

b6

A

A

T

G

# What's the likelihood of this tree?

It could be this

# What's the likelihood of this tree?

It could be this
**OR** this

**A**

b1

b4

**A**

**C**

b2

b3

b5

b6

A

A

T

G

33

# What's the likelihood of this tree?

It could be this
**OR** this
**OR** this

# What's the likelihood of this tree?

It could be this
**OR** this
**OR** this
**OR** this



35

# What's the likelihood of this tree?

It could be this
**OR** this
**OR** this
**OR** this
**OR** this



36

# What's the likelihood of this tree?

It could be this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this



37

# What's the likelihood of this tree?

It could be this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this

**A**

b1        b4

**C**        **A**

b2        b3        b5        b6

A        A        T        G

38

# What's the likelihood of this tree?

It could be this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this

**A**

b1

b4

**G**

b2

b3

**A**

b5

b6

A

A

T

G

39

# What's the likelihood of this tree?

It could be this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this
**OR** this
…
**OR** this

**T**

b1

b4

**T**

**T**

b2

b3

b5

b6

A

A

T

G

40

# What's the likelihood of this tree?

So the likelihood of the tree is the sum (**OR!**) over the likelihoods of all possible assignments of A, C, G, and T (all possible evolutionary histories) to the inner nodes *I1, I2, I3* of the tree.

# What's the likelihood of this tree?

So the likelihood of the tree is the sum (**OR!**) over the likelihoods of all possible assignments of A, C, G, and T (all possible evolutionary histories)
to the inner nodes *I1, I2, I3* of the tree.

There are 4 x 4 x 4 possible assignments in our example
→ this sounds very compute-intensive :-(

# The Felsenstein Pruning Algorithm

I1

Post order traversal

b1

b4

I2

I3

b2

b3

b5

b6

A

A

T

G

43

# Felsenstein Pruning

conditional likelihood vectors

b1

b4

b2

b3

b5

b6

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

A

A

T

G

# Felsenstein Pruning

$P_{AA}(b1)$ P(A)

# Felsenstein Pruning



$P_{AA}(b1)\ P(A)$ **OR**
$P_{AC}(b1)\ P(C)$

P(A)
P(C)
P(G)
P(T)

b4

b1

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

b2

b3

b5

b6

A

A

T

G

# Felsenstein Pruning



$P_{AA}(b1)\ P(A)$ **OR**

$P_{AC}(b1)\ P(C)$ **OR**

$P_{AG}(b1)\ P(G)$

P(A)
P(C)
P(G)
P(T)

b4

b1

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

b2

b3

b5

b6

A

A

T

G

# Felsenstein Pruning

$P_{AA}(b1)\ P(A)$ **OR**

$P_{AC}(b1)\ P(C)$ **OR**

$P_{AG}(b1)\ P(G)$ **OR**

$P_{AT}(b1)\ P(T)$

b1

b4

b2

b3

b5

b6

A

A

T

G

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

48

# Felsenstein Pruning



**AND!**

P(A)
P(C)
P(G)
P(T)

$P_{AA}(b4)\ P(A)$ **OR**

$P_{AC}(b4)\ P(C)$ **OR**

$P_{AG}(b4)\ P(G)$ **OR**

$P_{AT}(b4)\ P(T)$

b4

b1

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

b2

b3

b5

b6

A

A

T

G

# Felsenstein Pruning

**AND** (left branch/right branch)

$$\vec{L}_A^{(k)}(c) = \Big( \sum_{S=A}^{T} P_{AS}(b_i) \vec{L}_S^{(i)}(c) \Big) \Big( \sum_{S=A}^{T} P_{AS}(b_j) \vec{L}_S^{(j)}(c) \Big)$$

L^(k)

P(b_i)

A C G T

| | | | | |
|---|---|---|---|---|
| A | | | | |
| C | | | | |
| G | | | | |
| T | | | | |

P(A)
P(C)
P(G)
P(T)

P(b_j)

A C G T

| | | | | |
|---|---|---|---|---|
| A | | | | |
| C | | | | |
| G | | | | |
| T | | | | |

b_i

b_j

L^(i)

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

L^(j)

Position *c*

50

# Felsenstein Pruning

**OR** (along left branch)

$$\vec{L}_A^{(k)}(c) = \left( \sum_{S=A}^{T} P_{AS}(b_i) \vec{L}_S^{(i)}(c) \right) \left( \sum_{S=A}^{T} P_{AS}(b_j) \vec{L}_S^{(j)}(c) \right)$$

$L\hat{\ }(k)$

P(b_i)

P(b_j)

```
  A C G T
A
C
G
T
```

```
  A C G T
A
C
G
T
```

P(A)
P(C)
P(G)
P(T)

b_i

b_j

$L\hat{\ }(i)$

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

$L\hat{\ }(j)$

Position *c*

# Felsenstein Pruning

OR (along right branch)

$$\vec{L}_A^{(k)}(c) = \Big( \sum_{S=A}^{T} P_{AS}(b_i) \vec{L}_S^{(i)}(c) \Big) \Big( \sum_{S=A}^{T} P_{AS}(b_j) \vec{L}_S^{(j)}(c) \Big)$$

L^(k)

P(b_i)

P(b_j)

| | A | C | G | T |
|---|---|---|---|---|
| A | | | | |
| C | | | | |
| G | | | | |
| T | | | | |

P(A)
P(C)
P(G)
P(T)

b_i

b_j

| | A | C | G | T |
|---|---|---|---|---|
| A | | | | |
| C | | | | |
| G | | | | |
| T | | | | |

L^(i)

P(A)
P(C)
P(G)
P(T)

P(A)
P(C)
P(G)
P(T)

L^(j)

Position *c*

52

# Felsenstein Pruning



53

# Felsenstein Pruning

Likelihood at the root: $L_i = \pi_A\ P(A) + \pi_C\ P(C) + \pi_G\ P(G) + \pi_T\ P(T)$



54

# Why is time-reversibility important?

$$L = \sum_{S_4=A}^{T} \pi_{S_4} \sum_{S_3=A}^{T} P_{S_4 S_3}(b_1) L_{S_3}^{(3)} \sum_{S_5=A}^{T} P_{S_4 S_5}(b_4) L_{S_5}^{(5)}$$

$L_{S5}$

$b_1$   $b_4$

$L_{S3}$

# Why is time-reversibility important?

$$L = L' = \sum_{S_4=A}^{T} \pi_{S_4} \sum_{S_3=A}^{T} P_{S_4 S_3}(b_1 + x) L_{S_3}^{(3)} \sum_{S_5=A}^{T} P_{S_4 S_5}(b_4 - x) L_{S_5}^{(5)}$$

$L_{S3}$

$b_1'$

$b_4'$

$x$

$L_{S5}$

# Why is time-reversibility important?

$$L = L' = \sum_{S_4=A}^{T} \pi_{S_4} \sum_{S_3=A}^{T} P_{S_4 S_3}(b_1 + x) L_{S_3}^{(3)} \sum_{S_5=A}^{T} P_{S_4 S_5}(b_4 - x) L_{S_5}^{(5)}$$

$b_4':=b_1+b_4$

$b_1' := 0$

$L_{S5}$

$L_{S3}$

$x$

# Why is time-reversibility important?

This observation can be applied recursively to the tree

$\rightarrow$

It does not matter at all where we place the root!

$$L = L' = \sum_{S_4=A}^{T} \quad \cdots \quad \sum_{S_5=A}^{T} P_{S_4 S_5}(b_4 - x) L_{S_5}^{(5)}$$

$b_1' := 0$

$b_4' := b_1 + b_4$

$L_{S3}$

$L_{S5}$

$x$

# What's in the black box $P_{ij}(t)$?

Instantaneous rate matrix *R*!

# What's in the black box $P_{ij}(t)$?

What about the probabilities of staying in the current state?
→ they are given by the properties of continuous Markov chains!
e.g., $\lambda_{AA} = -(\lambda_{AC} + \lambda_{AG} + \lambda_{AT})$ rows in the *R* matrix need to sum to **0**

# What's in the black box $P_{ij}(t)$?

$$
\begin{array}{c c c c c}
 & A & C & G & T \\
A & * & \lambda_{AC} & \lambda_{AG} & \lambda_{AT} \\
C & & * & \lambda_{CG} & \lambda_{CT} \\
G & & & * & \lambda_{GT} \\
T & & \text{Symmetric} & & * \\
\end{array}
$$

# What's in the black box $P_{ij}(t)$?

Diagonal values are
given by the off-diagonal
values (R matrix property)
$\lambda_{AA} = -(\lambda_{AC} + \lambda_{AG} + \lambda_{AT})$

$$
\begin{array}{c c c c c}
 & A & C & G & T \\
A & * & \lambda_{AC} & \lambda_{AG} & \lambda_{AT} \\
C & & * & \lambda_{CG} & \lambda_{CT} \\
G & & & * & \lambda_{GT} \\
T & \text{Symmetric} & & & * \\
\end{array}
$$

# What's in the black box $P_{ij}(t)$?

|   | A | C | G | T |
|---|---|---|---|---|
| A | $*$ | $\lambda_{AC}$ | $\lambda_{AG}$ | $\lambda_{AT}$ |
| C |   | $*$ | $\lambda_{CG}$ | $\lambda_{CT}$ |
| G |   |   | $*$ | $\lambda_{GT}$ |
| T |   | Symmetric |   | $*$ |

Equilibrium frequency vector $\pi = (\pi_A, \pi_C, \pi_G, \pi_T)$ where $\pi_A + \pi_C + \pi_G + \pi_T = 1$

# The Jukes-Cantor model

$$
\begin{array}{c c c c c c}
 & A & C & G & T \\
A & * & \lambda & \lambda & \lambda \\
C & & * & \lambda & \lambda \\
G & & & * & \lambda \\
T & & & & *
\end{array}
$$

*Π = ( 1/4, 1/4, 1/4, 1/4)*

# Felsenstein 81

|   | A | C | G | T |
|---|---|---|---|---|
| A | * | λ | λ | λ |
| C |   | * | λ | λ |
| G |   |   | * | λ |
| T |   |   |   | * |

$$\Pi_i \neq \Pi_j$$

# Kimura 2-parameter model 1980

$$
\begin{array}{c c c c c}
 & A & C & G & T \\
A & * & \lambda & \zeta & \lambda \\
C & & * & \zeta & \lambda \\
G & & & * & \zeta \\
T & & & & * \\
\end{array}
$$

*Π = ( 1/4, 1/4, 1/4, 1/4)*

# HKY85

|   | A | C | G | T |
|---|---|---|---|---|
| A | $*$ | $\lambda$ | $\zeta$ | $\lambda$ |
| C |   | $*$ | $\zeta$ | $\lambda$ |
| G |   |   | $*$ | $\zeta$ |
| T |   |   |   | $*$ |

$$\Pi_i \neq \Pi_j$$

# GTR 1986

$$
\begin{array}{c|cccc}
 & A & C & G & T \\
\hline
A & * & \alpha & \beta & \gamma \\
C & & * & \delta & \varepsilon \\
G & & & * & \zeta \\
T & & & & * \\
\end{array}
$$

$$\Pi_i \neq \Pi_j$$

# GTR 1986

|   | A | C | G | T |
|---|---|---|---|---|
| A | * | α | β | γ |
| C |   | * | δ | ε |
| G |   |   | * | ζ |
| T |   |   |   | * |

$\Pi_i \neq \Pi_j$

Note that these are **relative** rates, their Values only matter relative to each other, so we can set $\zeta := 1.0$ by default

# GTR 1986

$$
\begin{array}{c c c c c}
 & A & C & G & T \\
A & * & \alpha & \beta & \gamma \\
C & & * & \delta & \epsilon \\
G & & & * & 1.0 \\
T & & & & * \\
\end{array}
$$

$\Pi_i \neq \Pi_j$

Note that these are **relative** rates, their values only matter relative to each other, so we can set $\zeta := 1.0$ by default. Although the GTR model has 6 rates, it only has 5 free parameters!

# Model Hierarchy

# GTR 1986

This is a rate matrix, time reversibility would Require $\pi_i r_{ij} = \pi_j r_{ji}$

$$
\begin{array}{c c c c c}
 & A & C & G & T \\
A & * & \alpha & \beta & \gamma \\
C & & * & \delta & \varepsilon \\
G & & & * & 1.0 \\
T & & & & * \\
\end{array}
$$

$\Pi_i \neq \Pi_j$

# GTR 1986

This is a rate matrix, time reversibility would Require $\pi_i r_{ij} = \pi_j r_{ji}$ Solution: introduce a $Q$ matrix $Q := diag(\pi)\ R$

$$
\begin{array}{c c c c c}
 & A & C & G & T \\
A & * & \alpha & \beta & \gamma \\
C & & * & \delta & \varepsilon \\
G & & & * & 1.0 \\
T & & & & * \\
\end{array}
$$

$$
\begin{pmatrix}
\pi_A & & & \\
& \pi_C & & \\
& & \pi_G & \\
& & & \pi_T \\
\end{pmatrix}
$$

$\Pi_i \neq \Pi_j$

# GTR 1986

This is a rate matrix, time reversibility would Require $\pi_i r_{ij} = \pi_j r_{ji}$
Solution: introduce a $Q$ matrix $Q := diag(\pi)\, R$

$$
\begin{array}{c|cccc}
 & A & C & G & T \\
\hline
A & * & \alpha & \beta & \gamma \\
C & & * & \delta & \varepsilon \\
G & & & * & 1.0 \\
T & & & & *
\end{array}
$$

$\Pi_i \neq \Pi_j$

$$
\begin{pmatrix}
\Pi_A & & & \\
 & \Pi_C & & \\
 & & \Pi_G & \\
 & & & \Pi_T
\end{pmatrix}
$$

Then, $\pi_i r_{ij} = \pi_j r_{ji}$ holds

# So how do we compute P(t) from Q?

- As we have seen in the lecture on Markov chains:

  $P(t) = e^{Qt} = I + Qt + 1/2! (Qt)^2 + 1/3! (Qt)^3 + \ldots$

- but this is unfortunately a matrix eponential :-(

- I will spare you the details, but in general, e.g., for GTR we need to apply an egienvector/eigenvalue decomposition of Q to calculate:

  $P(t) = U \, exp(diag(\lambda_i)t) \, U^{-1}$

Matrix and inverse matrix of eigenvectors of $Q$

# So how do we compute P(t) from Q?

- As we have seen in the lecture on Markov chains:

  $P(t) = e^{Qt} = I + Qt + 1/2! (Qt)^2 + 1/3! (Qt)^3 + \ldots$

- but this is unfortunately a matrix exponential :-(

- I will spare you the details, but in general, e.g., for GTR we need to apply an egienvector/eigenvalue decomposition of Q to calculate:

  $P(t) = U \exp(diag(\lambda_i)t) U^{-1}$

Diagonal matrix of eigenvalues of *Q*, here the exponential function *exp()* is invoked on scalar values!

# Likelihood Calculations

- So far, we have only seen how to calculate **a** likelihood on a

    - given, fixed tree topology

    - with given fixed branch lengths

    - and given, fixed remaining model parameters

- Computing the **maximum** likelihood score, is much more complicated as it requires functions for optimizing continuous parameters and functions for searching the discrete space of trees !