

Likelihood in Phylogenetics

19 July 2016

Workshop on Molecular Evolution
Woods Hole, Mass.

Paul O. Lewis

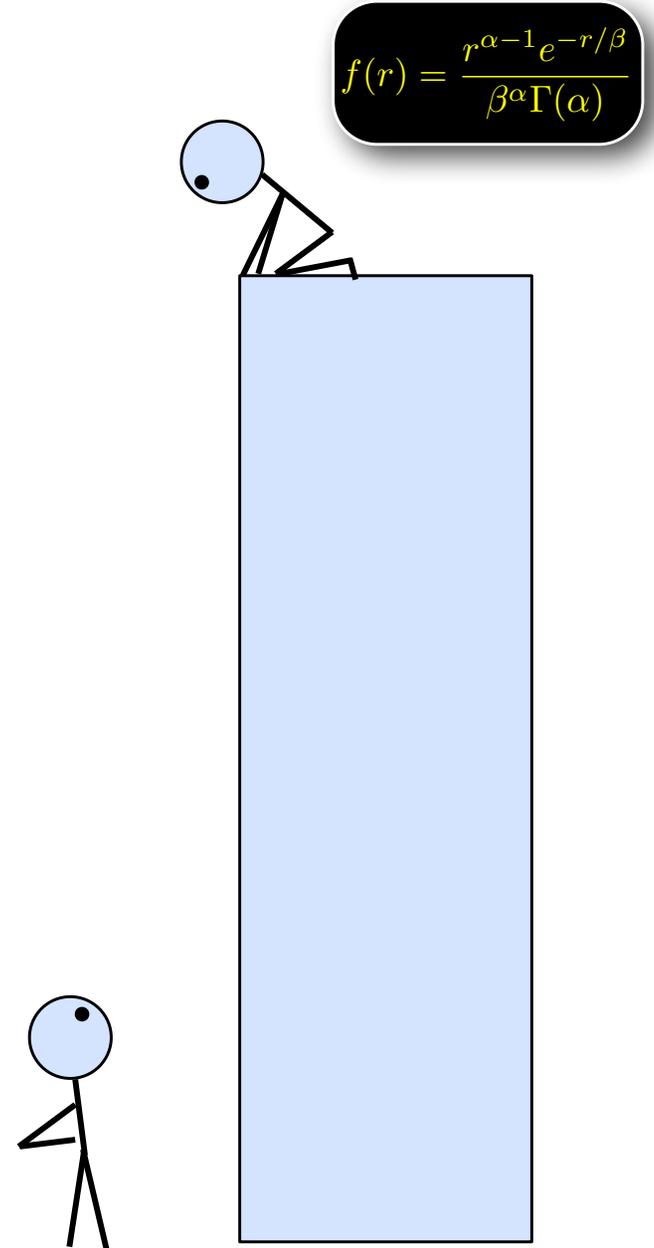
Department of Ecology & Evolutionary Biology

UConn
UNIVERSITY OF CONNECTICUT

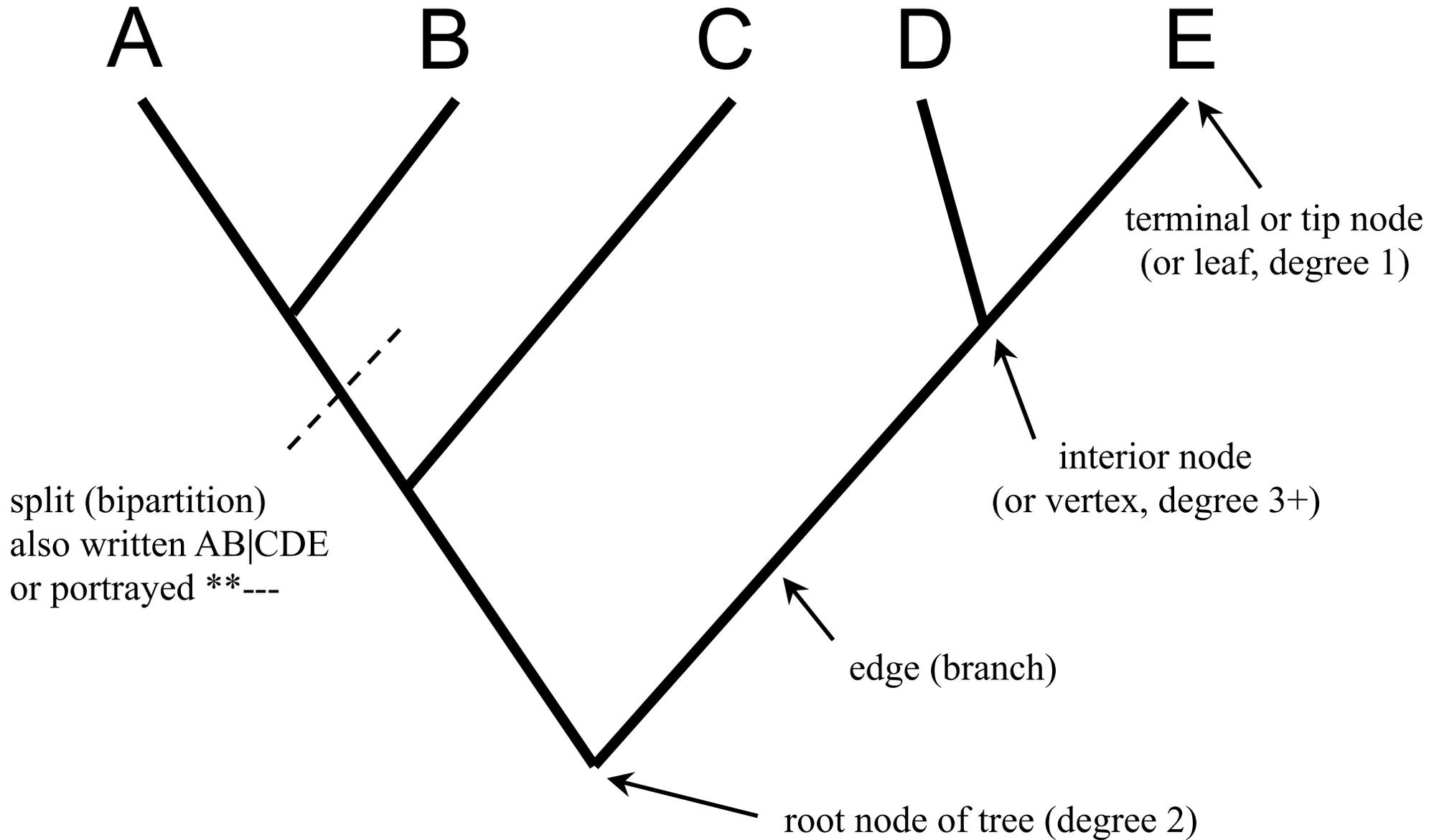


Goals

- * Explain jargon
 - * Increase comfort level
 - * Provide background
- In other words...give a hand up



Tree jargon



The Plan

- Probability review

- The AND and OR rules
- Independence of events

- Likelihood

- Substitution models

- What does it mean?
- Likelihood of a single sequence
- Maximum likelihood distances
- Likelihoods of trees

- Markov model basics
- Transition probabilities
- Survey of models
- Rate heterogeneity
- Codon models
- Amino acid models

Combining probabilities

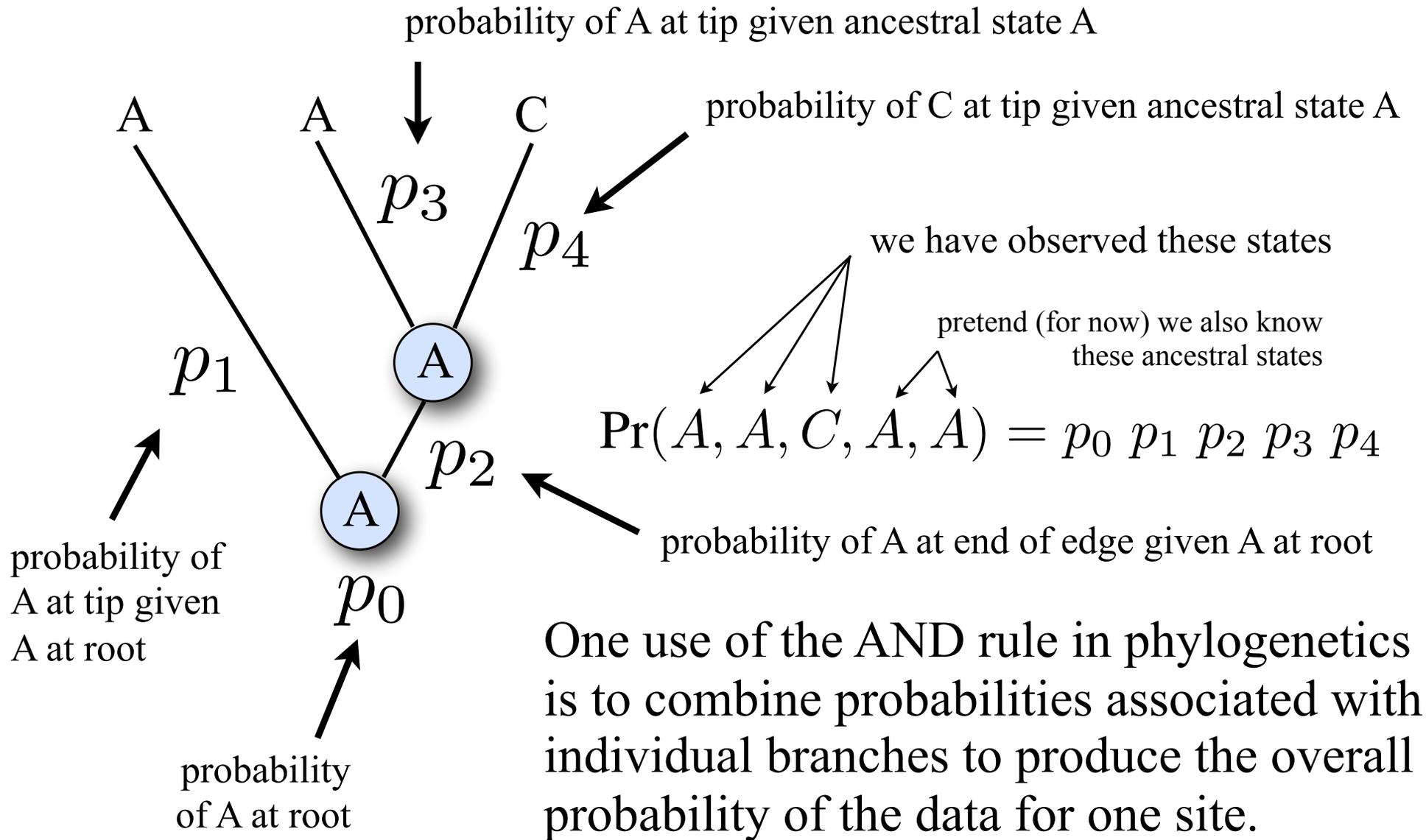
- *Multiply* probabilities if the component events must happen **simultaneously** (i.e. where you would naturally use the word AND when describing the problem)

Using 2 dice, what is the probability of



$$(1/6) \times (1/6) = 1/36$$

AND rule in phylogenetics



Combining probabilities

- *Add* probabilities if the component events are **mutually exclusive** (i.e. where you would naturally use the word OR in describing the problem)

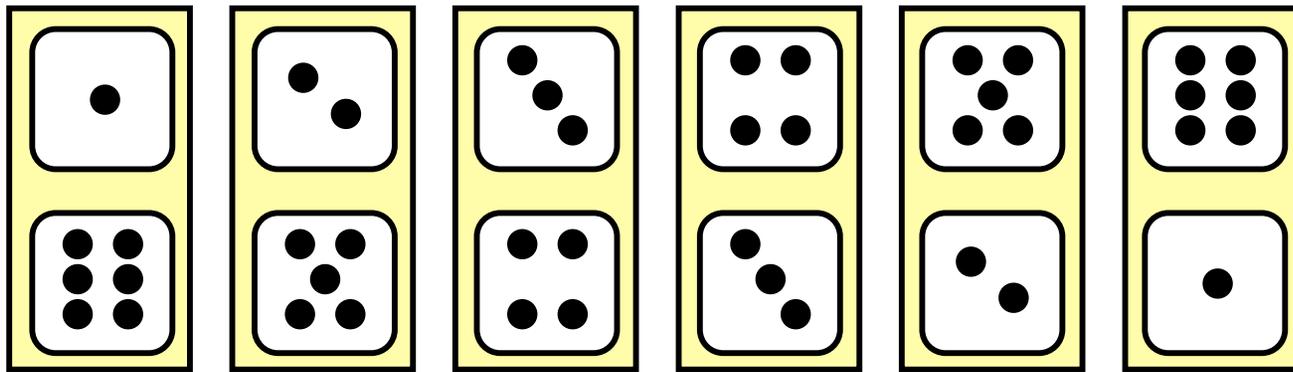
Using one die, what is the probability of



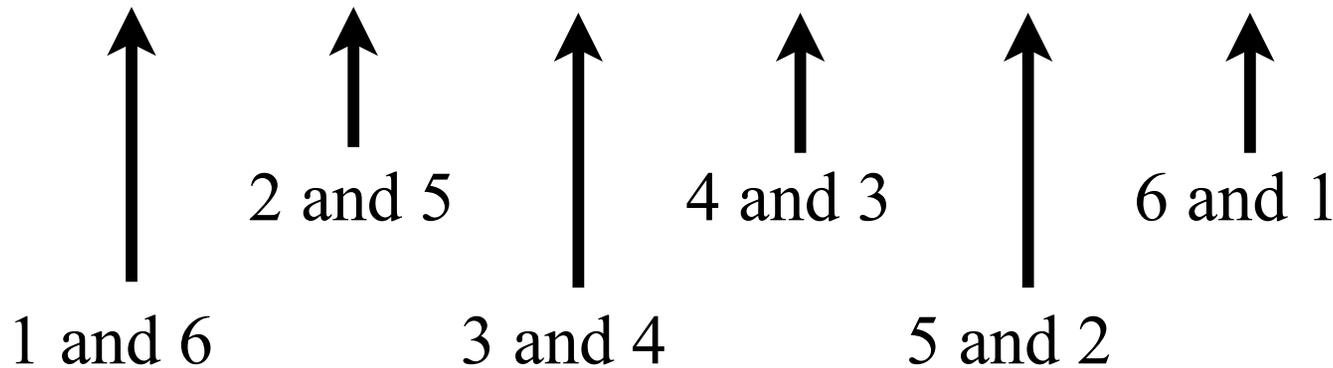
$$(1/6) + (1/6) = 1/3$$

Combining AND and OR

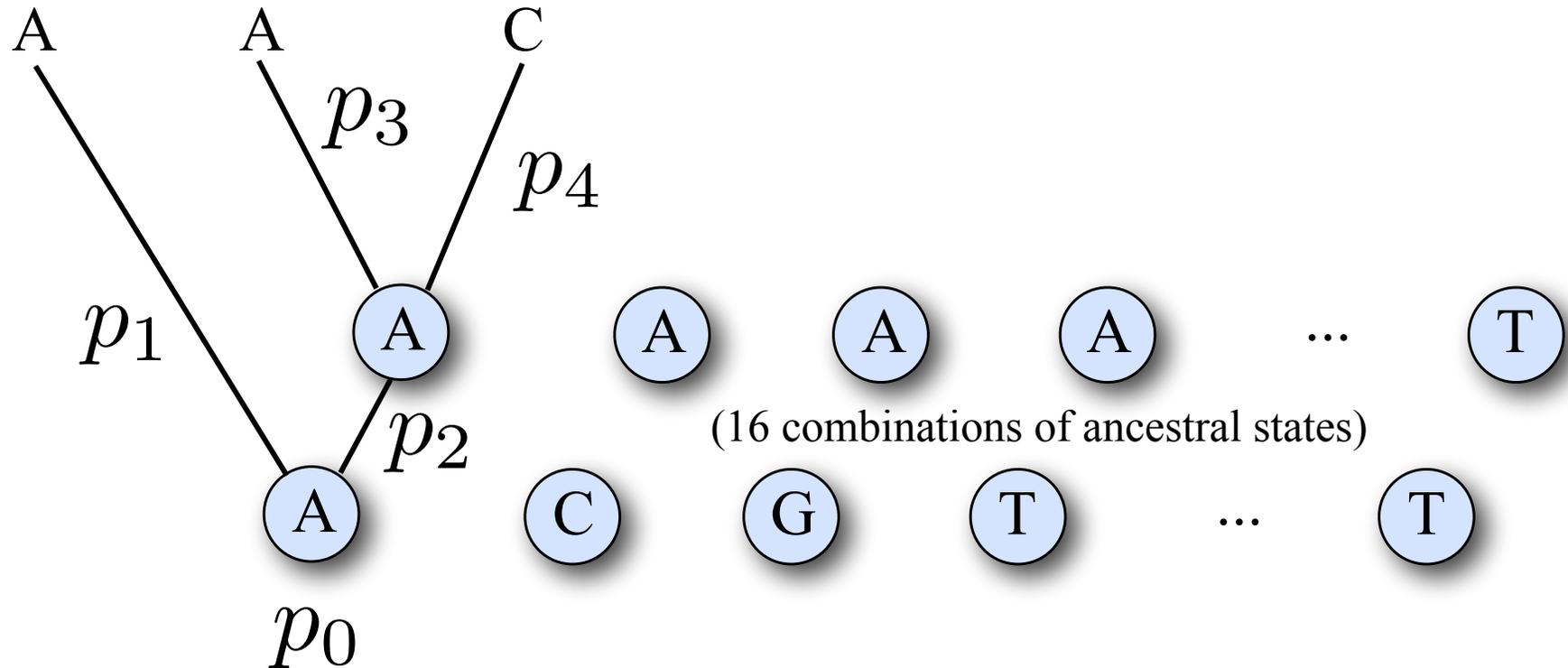
What is the probability that the sum of two dice is 7?



$$(1/36) + (1/36) + (1/36) + (1/36) + (1/36) + (1/36) = 1/6$$



Using both AND and OR in phylogenetics

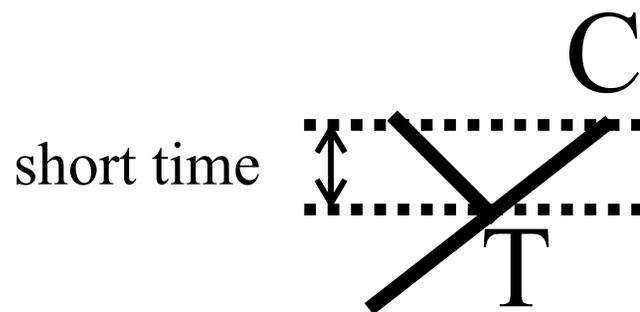


AND rule used to compute probability of the observed data for *each combination* of ancestral states.

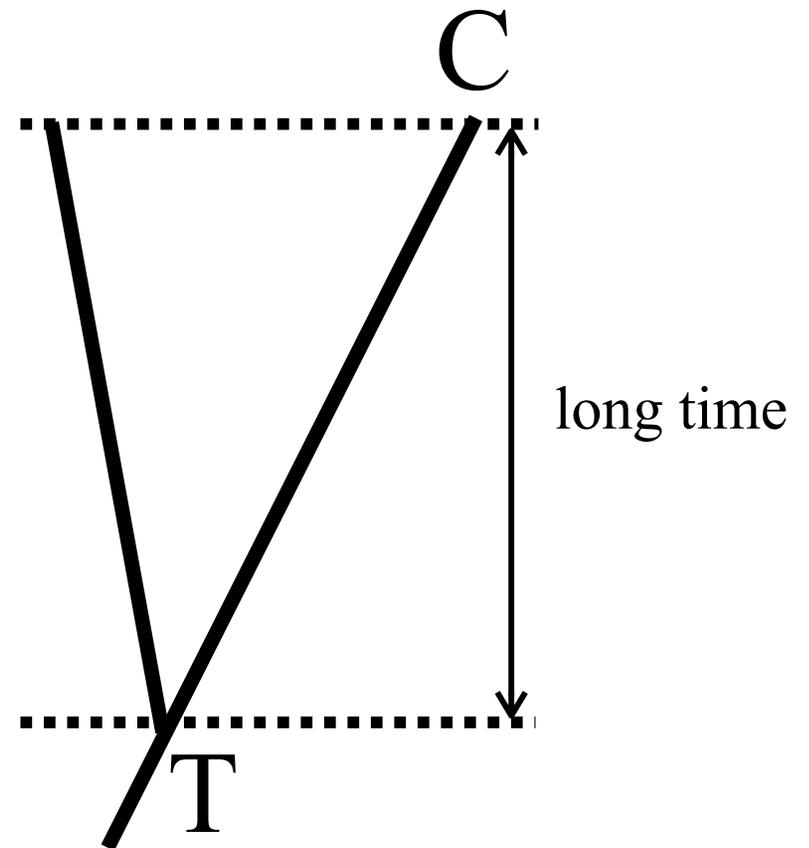
OR rule used to combine different combinations of ancestral states.

Non-independence in molecular evolution

The state present in the descendant is **not independent** of the state in the ancestor



less probable



more probable

Conditional Independence

Assume both A and B depend on C:

$$\Pr(A|C) \neq \Pr(A) \quad \Pr(B|C) \neq \Pr(B)$$

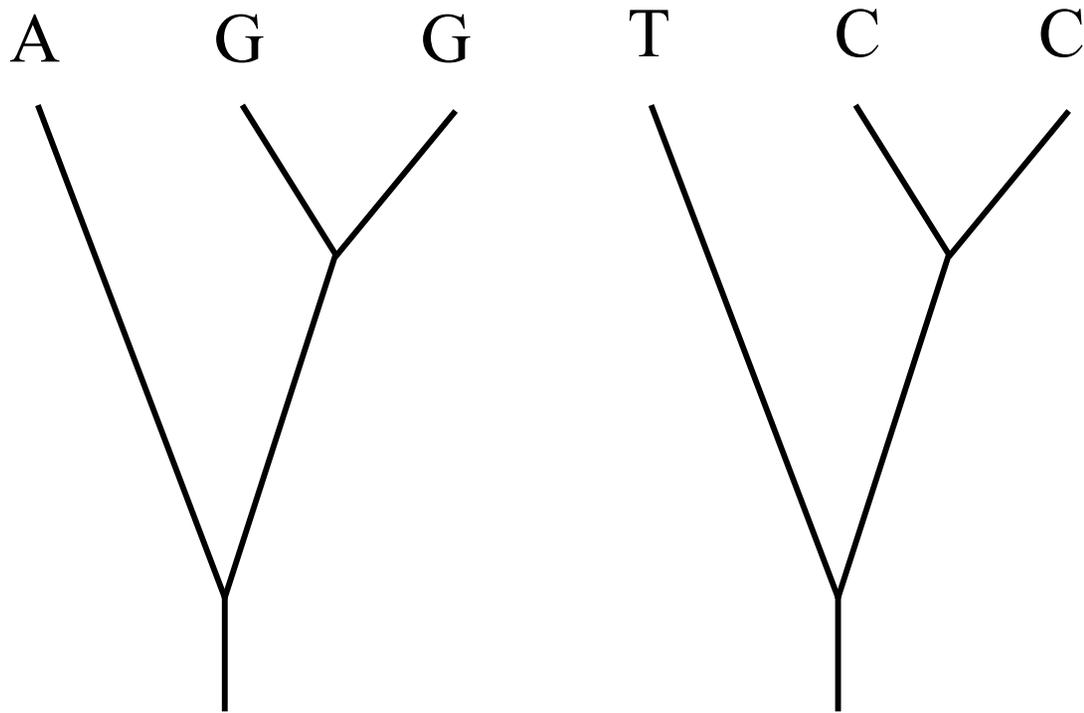
If we can say this...

$$\Pr(B|A,C) = \Pr(B|C)$$

...then events A and B are **conditionally independent** and we can express the joint (conditional) probability as the product of $\Pr(A|C)$ and $\Pr(B|C)$

$$\Pr(A \text{ and } B|C) = \Pr(A|C) \Pr(B|C)$$

Conditional independence in molecular evolution



The site data patterns AGG and TCC are assumed by most models to be conditionally independent.

The patterns both depend on the underlying tree (including edge lengths) and the substitution model.

$$\Pr(\text{AGG and TCC}|\text{tree, model}) = \Pr(\text{AGG}|\text{tree, model}) \Pr(\text{TCC}|\text{tree, model})$$

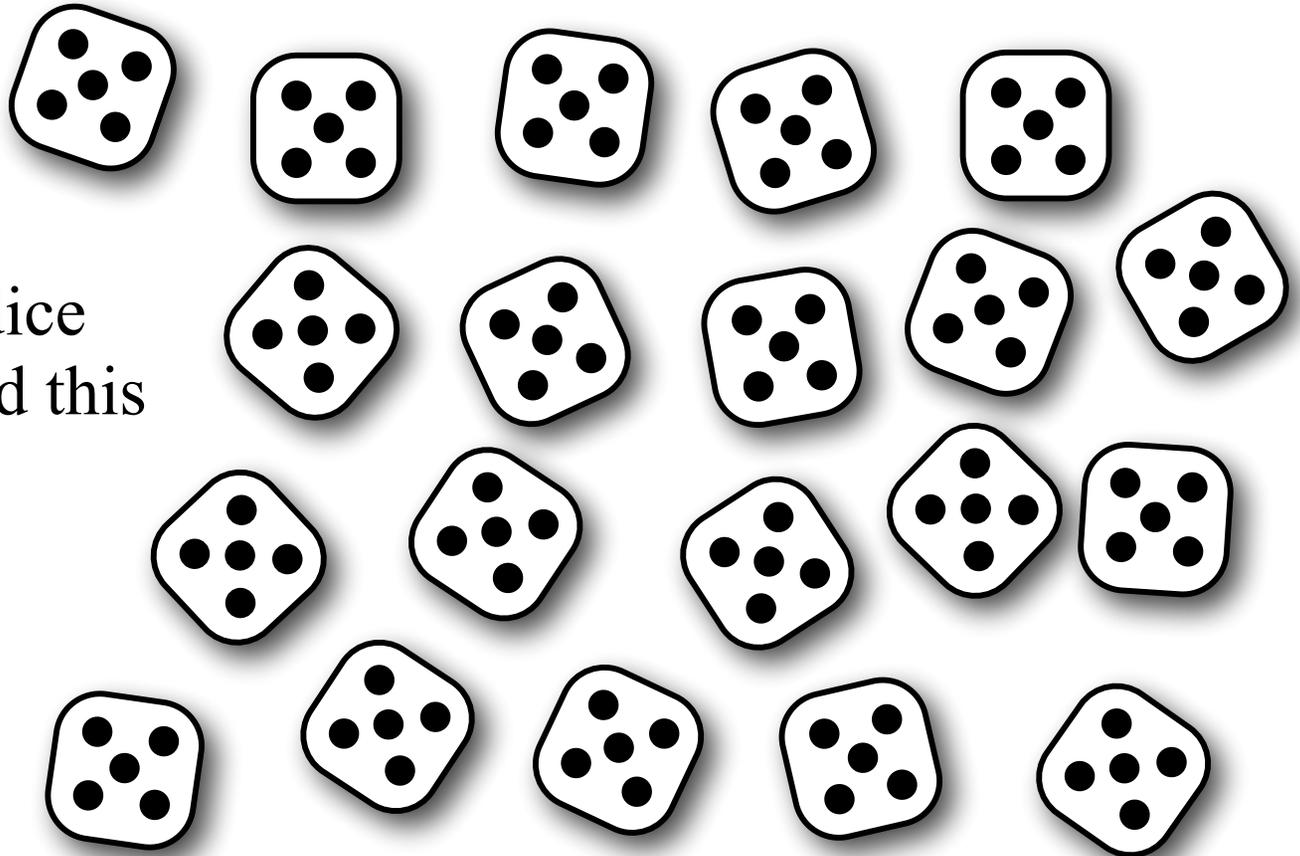
Likelihood

The Likelihood Criterion

The probability of the observations computed using a model tells us how surprised we should be.

The preferred model is the one that surprises us least.

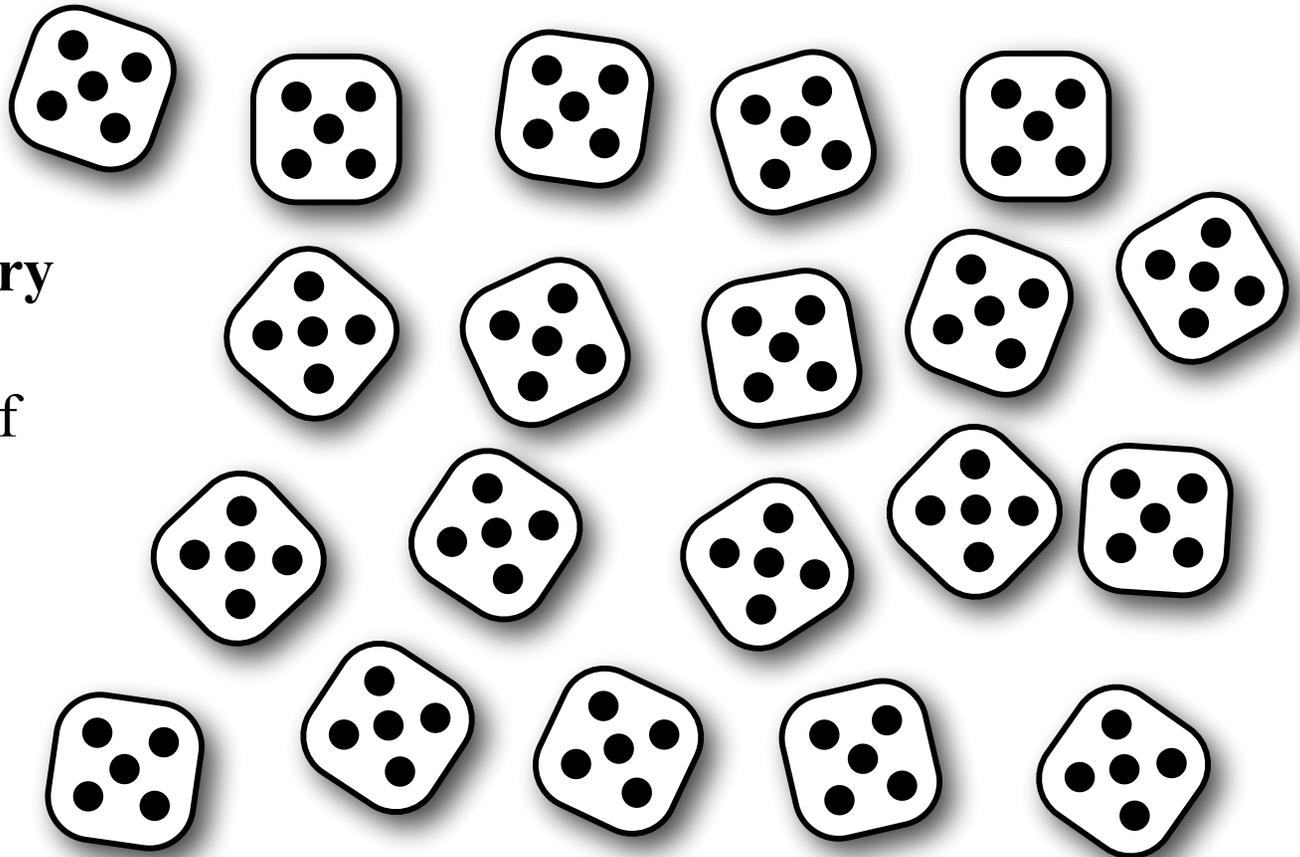
Suppose I threw 20 dice down on the table and this was the result...



The Fair Dice model

$$\Pr(\text{obs.}|\text{fair dice model}) = \left(\frac{1}{6}\right)^{20} = \frac{1}{3,656,158,440,062,976}$$

You should have been **very surprised** at this result because the probability of this event is **very small**: only 1 in 3.6 quadrillion!

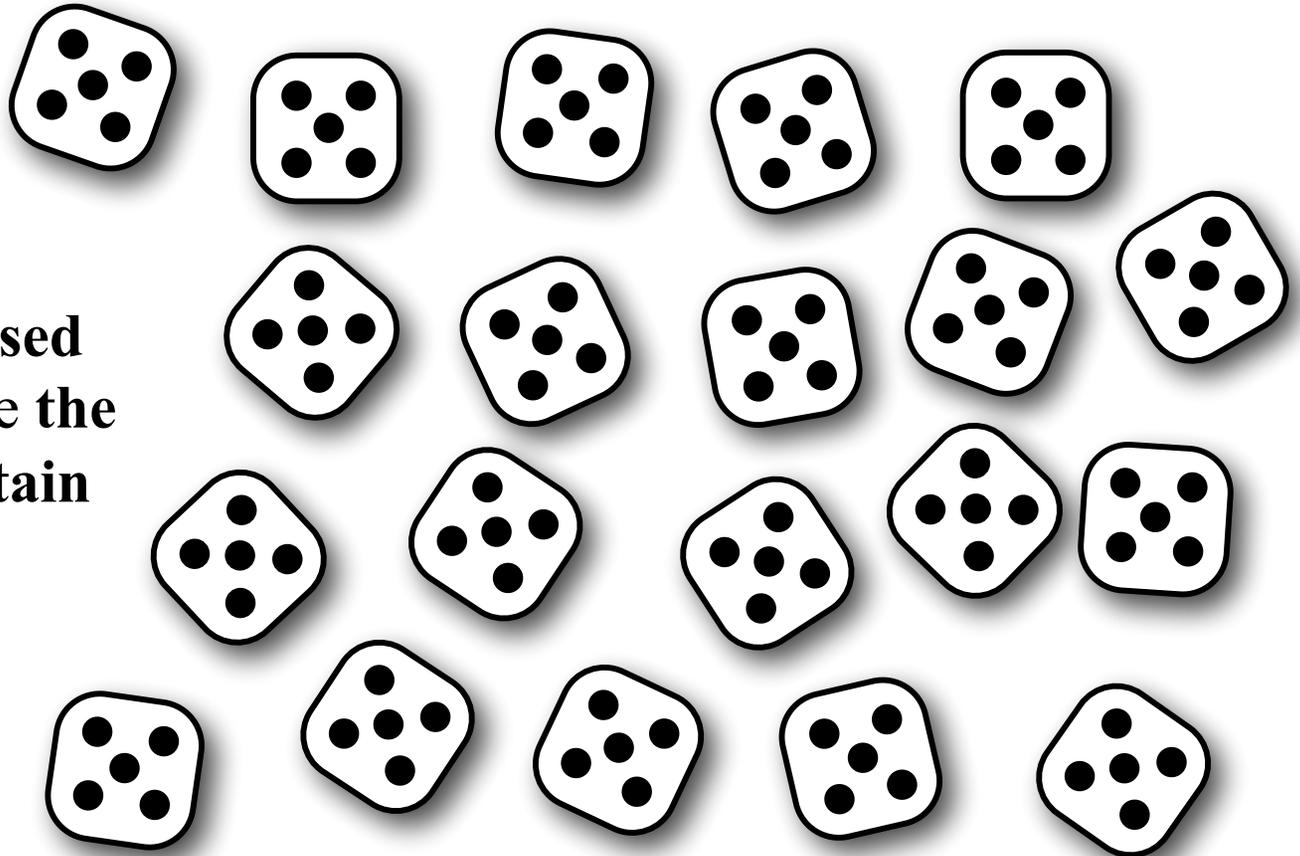


The Trick Dice model

(assumes dice each have 5 on every side)

$$\Pr(\text{obs.} | \text{trick dice model}) = 1^{20} = 1$$

You should **not be surprised at all** at this result because **the observed outcome is certain** under this model



Results

Model	Likelihood	Surprise level
Fair Dice	$\frac{1}{3,656,158,440,062,976}$	Very, <i>very</i> , <i>very</i> surprised
Trick Dice	1	Not surprised at all

winning model maximizes likelihood
(and thus minimizes surprise)

Likelihood: why a new term?

Outcome	Fair coin model	Two-heads model
H	0.5	1
T	0.5	0
	1	1

Likelihoods of models given one particular data outcome are *not* expected to sum to 1.0

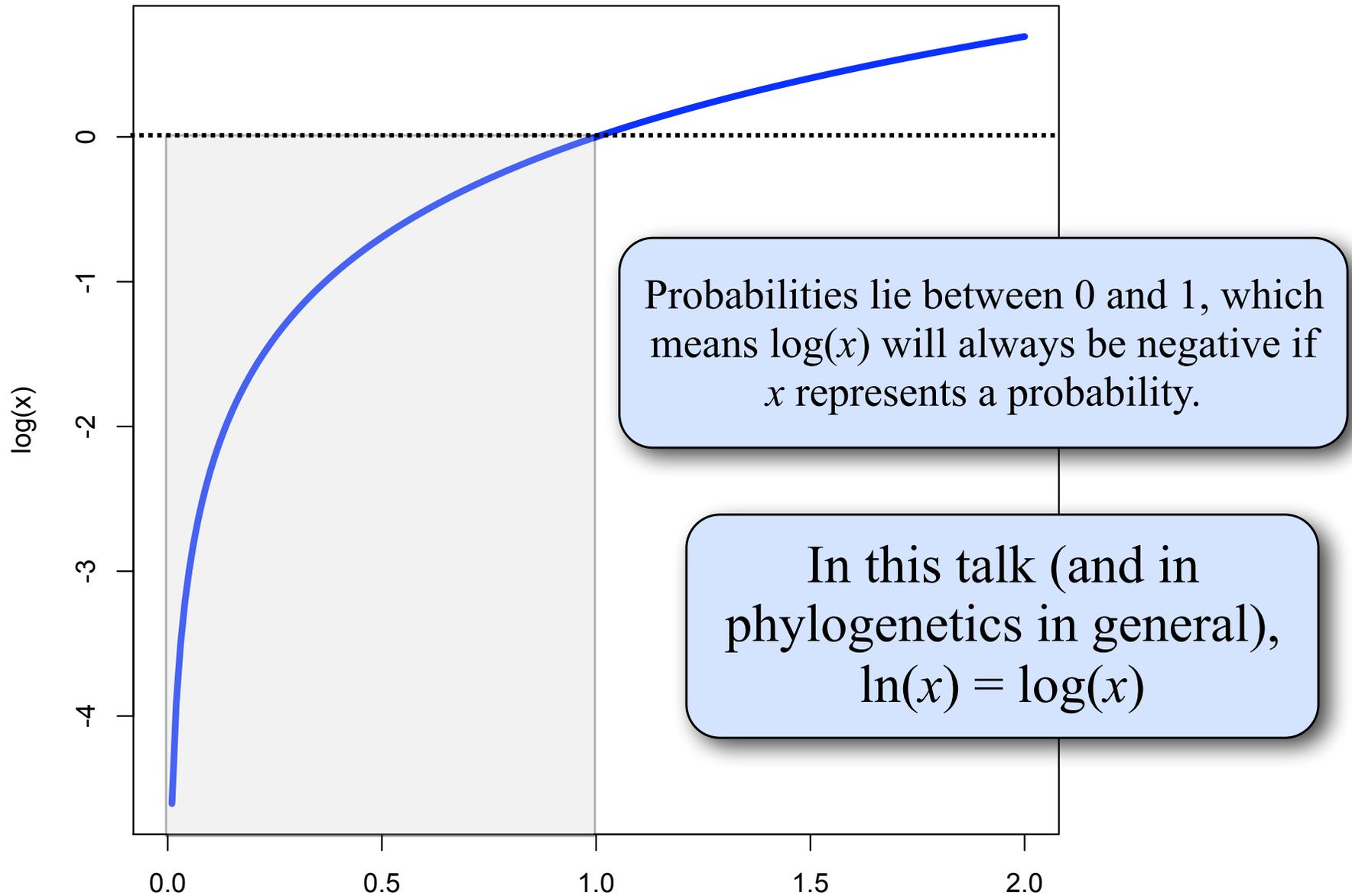
Probabilities of data outcomes given one particular model sum to 1.0

Likelihood and model comparison

- Analyses using likelihoods ultimately involve **model comparison**
- The models compared can be **discrete** (as in the fair vs. trick dice example)
- More often the models compared differ **continuously**:
 - Model 1: branch length is 0.05
 - Model 2: branch length is 0.06

Rather than having an infinity of models, we instead think of the branch length as a **parameter** within one model

Likelihoods vs. log-likelihoods



Likelihood of a single sequence

First 32 nucleotides of the $\psi\eta$ -globin gene of gorilla:

GAAGTCCTTGAGAAATAAACTGCACACTGG

$$L = \pi_G \pi_A \pi_A \pi_G \pi_T \pi_C \pi_C \pi_T \pi_T \pi_G \pi_A \pi_G \pi_A \pi_A \pi_A \pi_T \pi_A \pi_A \pi_A \pi_C \pi_T \pi_G \pi_C \pi_A \pi_C \pi_A \pi_C \pi_A \pi_C \pi_T \pi_G \pi_G$$
$$= \pi_A^{12} \pi_C^7 \pi_G^7 \pi_T^6$$

Note that we are assuming independence among sites here

$$\log L = 12 \log(\pi_A) + 7 \log(\pi_C) + 7 \log(\pi_G) + 6 \log(\pi_T)$$

We can already see by eye-balling this that a model allowing **unequal** base frequencies will **fit better** than a model that assumes **equal** base frequencies because there are about twice as many As as there are Cs, Gs and Ts.

Model ranking using LRT or AIC

Likelihood Ratio Tests (LRT) and the Akaike Information Criterion (AIC) provide two ways to evaluate whether an **unconstrained** model fits the data significantly better than a **constrained** version of the same model.

Find *maximum* logL under the *unconstrained* model:

$$\begin{aligned}\log L_{\text{unconstrained}} &= 12 \log(\pi_A) + 7 \log(\pi_C) + 7 \log(\pi_G) + 6 \log(\pi_T) \\ &= 12 \log(0.375) + 7 \log(0.219) + 7 \log(0.219) + 6 \log(0.187) \\ &= -43.1\end{aligned}$$

This model has 3 estimated parameters

Find *maximum* logL under the *constrained* model:

$$\begin{aligned}\log L_{\text{constrained}} &= 12 \log(\pi_A) + 7 \log(\pi_C) + 7 \log(\pi_G) + 6 \log(\pi_T) \\ &= 12 \log(0.25) + 7 \log(0.25) + 7 \log(0.25) + 6 \log(0.25) \\ &= -44.4\end{aligned}$$

This model has 0 estimated parameters

Likelihood Ratio Test (LRT)

Calculate the likelihood ratio test statistic:

$$\begin{aligned} R &= -2 [\log(L_{\text{constrained}}) - \log(L_{\text{unconstrained}})] \\ &= -2 [-44.4 - (-43.1)] \\ &= 2.6 \end{aligned}$$

(Note that the log-likelihoods used in the test statistic have been *maximized* under each model separately)

“unconstrained” does fit better than “constrained” ($-43.1 > -44.4$), but not significantly better ($P = 0.457$, chi-squared with 3 d.f.*)

*The number of degrees of freedom equals the difference between the two models in the number of estimated parameters. In this case, unconstrained has 3 parameters and constrained has 0, so $\text{d.f.} = 3 - 0 = 3$

Akaike Information Criterion (AIC)

Calculate AIC for each model:

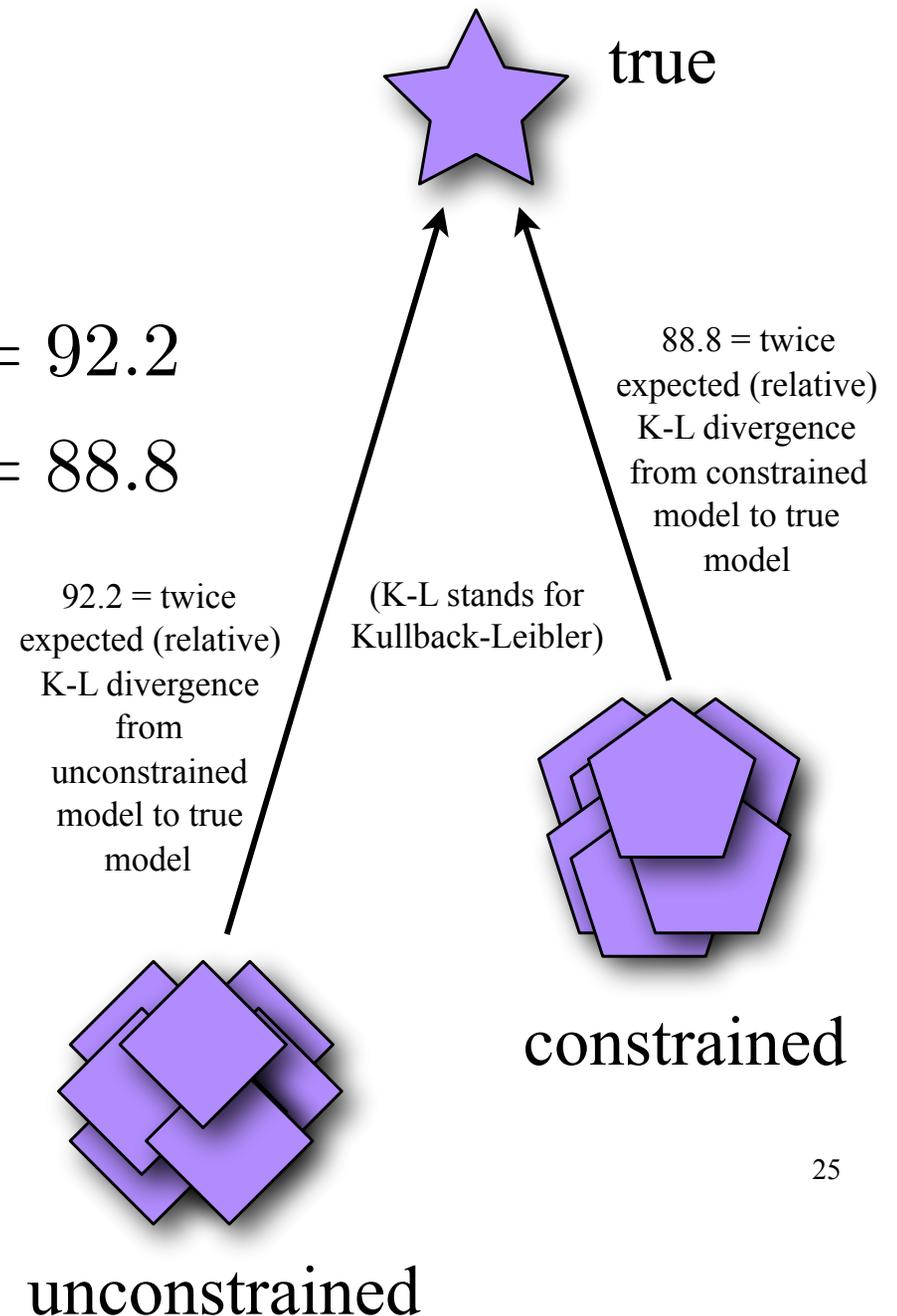
$$AIC = 2k - 2 \log(\max(L))$$

$$AIC_{\text{unconstrained}} = 2(3) - 2(-43.1) = 92.2$$

$$AIC_{\text{constrained}} = 2(0) - 2(-44.4) = 88.8$$

The constrained model is a better choice than the unconstrained model according to AIC

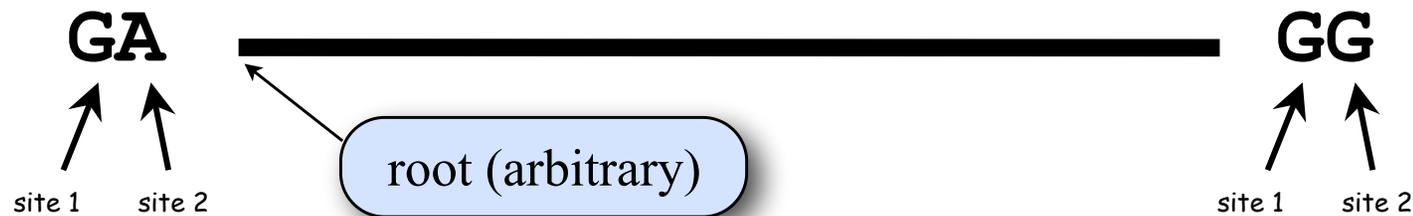
Dave Swofford will give you a much more complete explanation of LRT and AIC this afternoon



Likelihood of the simplest tree

sequence 1  sequence 2

To keep things simple, assume that the sequences are only 2 nucleotides long:



$$L = L_1 L_2$$

$$= \left[\begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \end{pmatrix} \left(\frac{1}{4} + \frac{3}{4} e^{-4\alpha t} \right) \right] \left[\begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \end{pmatrix} \left(\frac{1}{4} - \frac{1}{4} e^{-4\alpha t} \right) \right]$$

Pr(G)

Pr(G|G, αt)

Pr(A)

Pr(G|A, αt)

Note that we are NOT assuming independence here

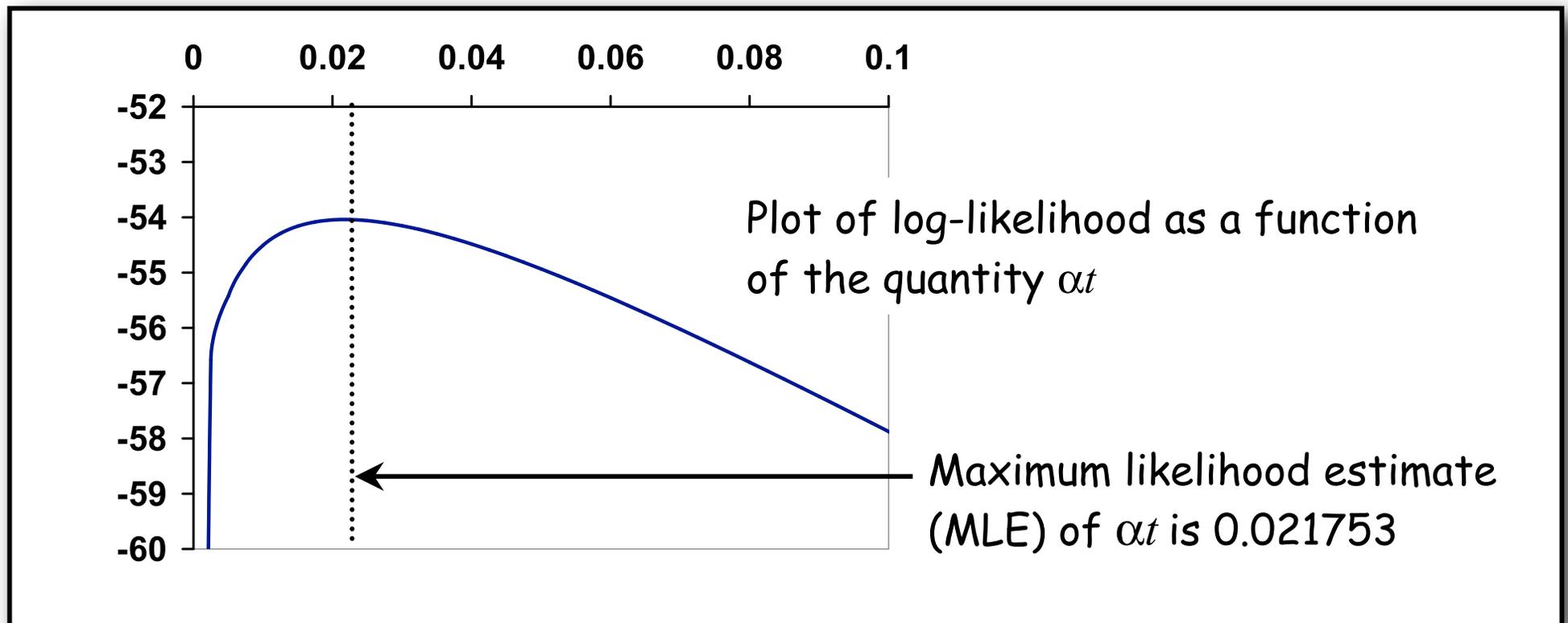
Maximum likelihood estimation

First 32 nucleotides of the $\psi\eta$ -globin gene of gorilla and orangutan:

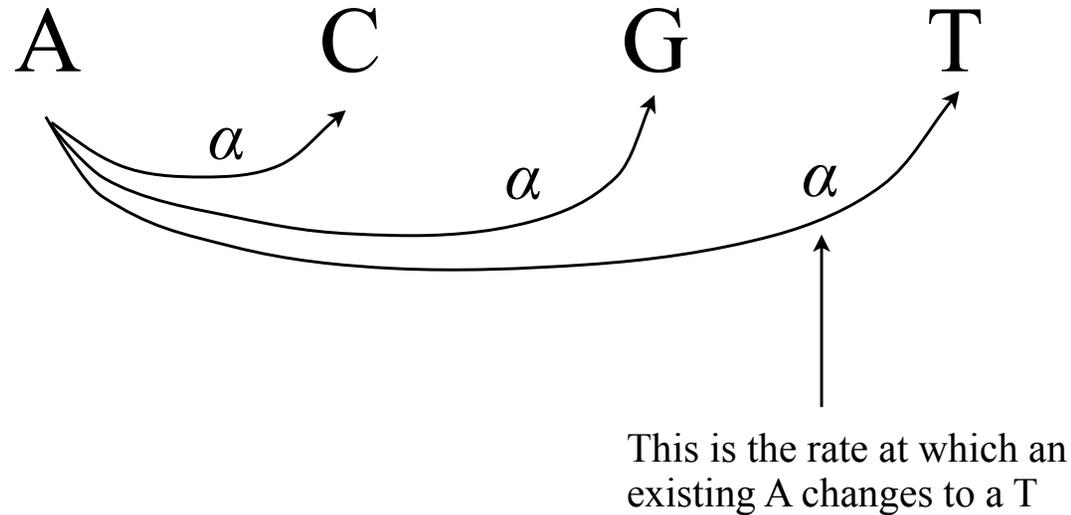
gorilla **GAAGTCCTTGAGAAATAAACTGCACACACTGG**

orangutan **GGACTCCTTGAGAAATAAACTGCACACACTGG**

$$L = \left[\binom{1}{4} \left(\frac{1}{4} + \frac{3}{4} e^{-4\alpha t} \right) \right]^{30} \left[\binom{1}{4} \left(\frac{1}{4} - \frac{1}{4} e^{-4\alpha t} \right) \right]^2$$



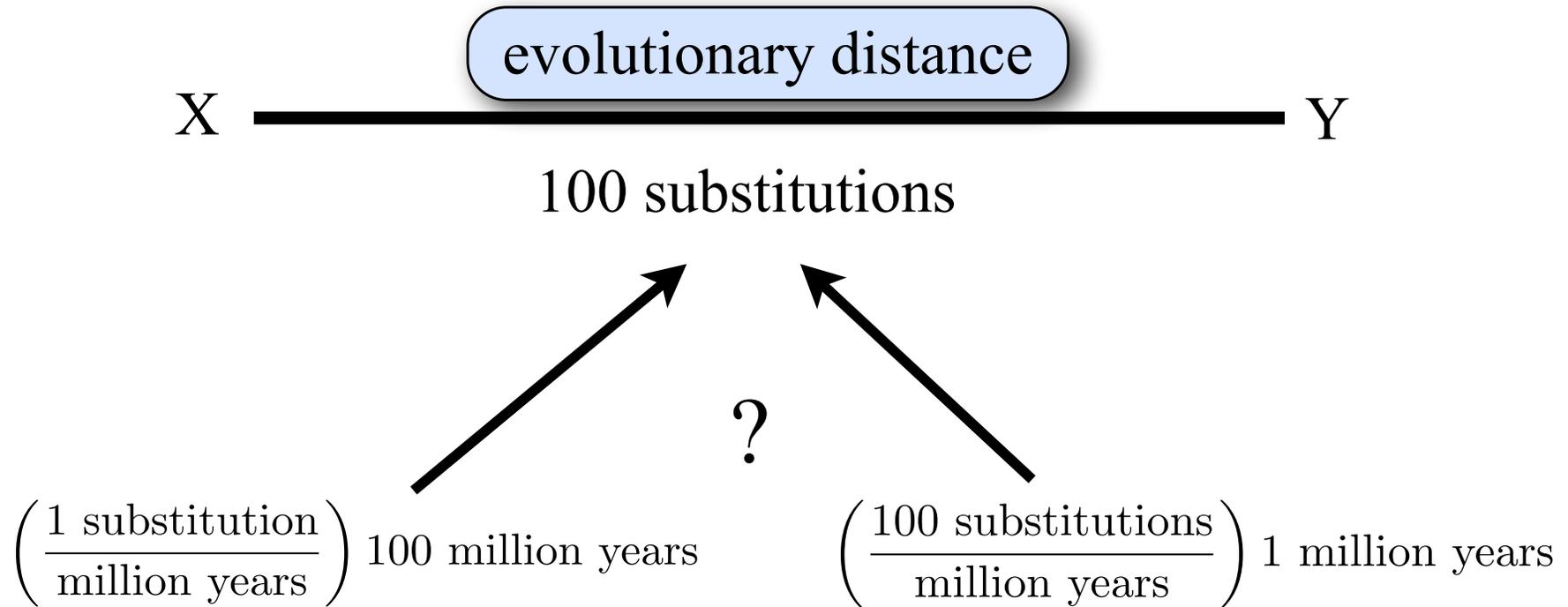
number of substitutions = rate \times time



Overall substitution rate is 3α , so the expected number of substitutions (v) is

$$v = 3\alpha t$$

Rate and time are confounded



On Friday, Tracy Heath will introduce models that allow separate estimation of rates and times, but without extra information/constraints, sequence data allow only estimation of the **number** of substitutions.

A convenient convention

Because rate and time are confounded, it is convenient to arbitrarily standardize things by setting the rate to a value such that **one substitution** is expected to occur in **one unit of time** for each site.

This results in “time” (the length of a branch) being measured in units of **evolutionary distance (expected number of substitutions per site)** rather than years (or some other calendar unit).

evolutionary distance $v = 3\alpha t$

$$v = 3 \left(\frac{1}{3} \right) t$$

Setting $\alpha=1/3$ results in v equalling t

Evolutionary distances for several common models

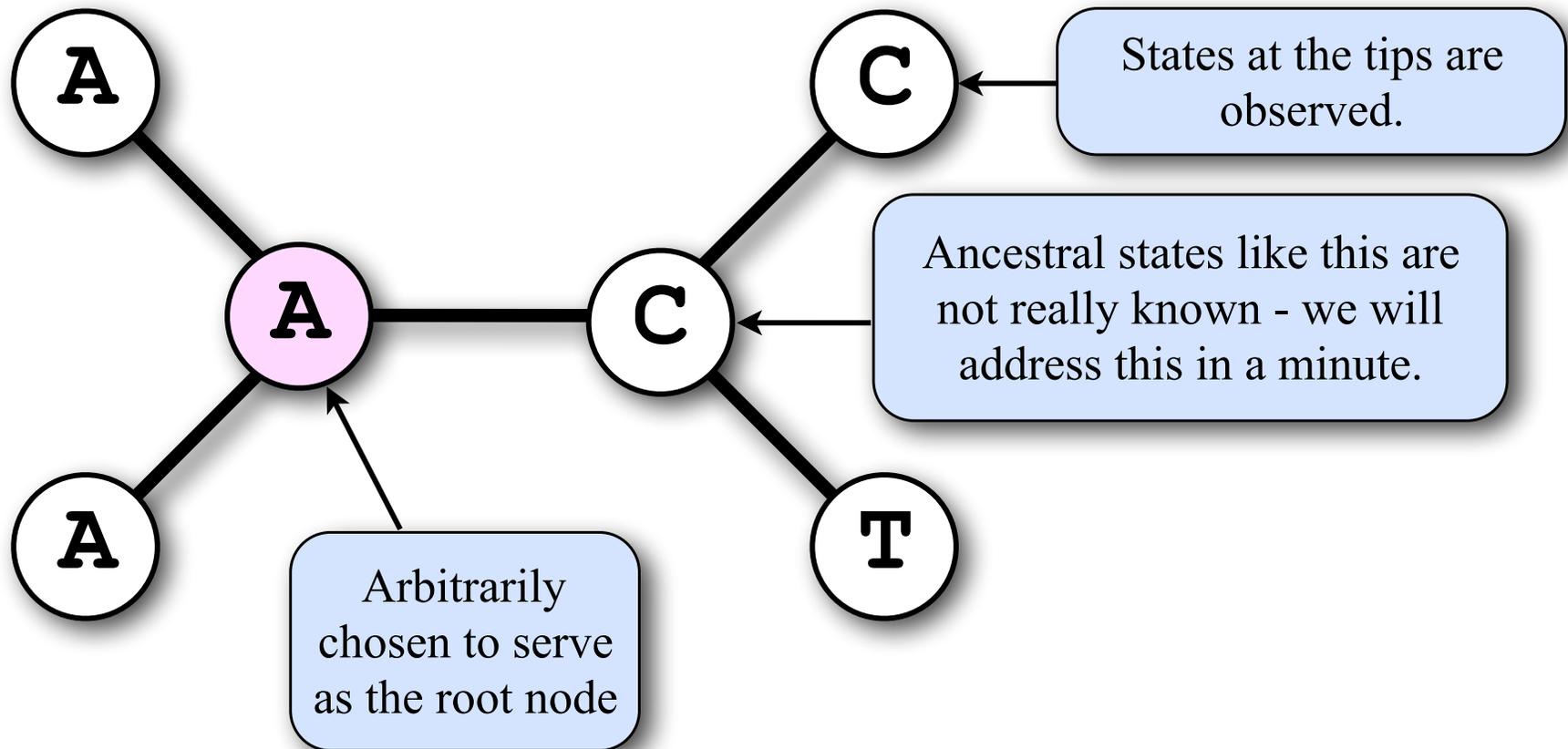
Model	Expected no. substitutions: $v = \{r\}t$
JC69	$v = \{3\alpha\}t$
F81	$v = \{2\mu(\pi_R\pi_Y + \pi_A\pi_G + \pi_C\pi_T)\}t$
K80	$v = \{\beta(\kappa + 2)\}t$
HKY85	$v = \{2\mu[\pi_R\pi_Y + \kappa(\pi_A\pi_G + \pi_C\pi_T)]\}t$

In the formulas above, the overall rate r (in curly brackets) is a function of all parameters in the substitution model.

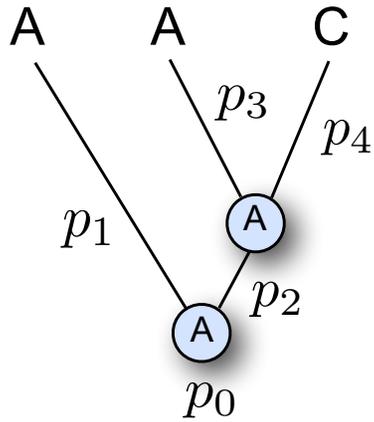
One substitution model parameter is always determined from the edge length (using our convention that $v = t$); the others are usually global (i.e. same value applies to all edges).

Likelihood of an unrooted tree

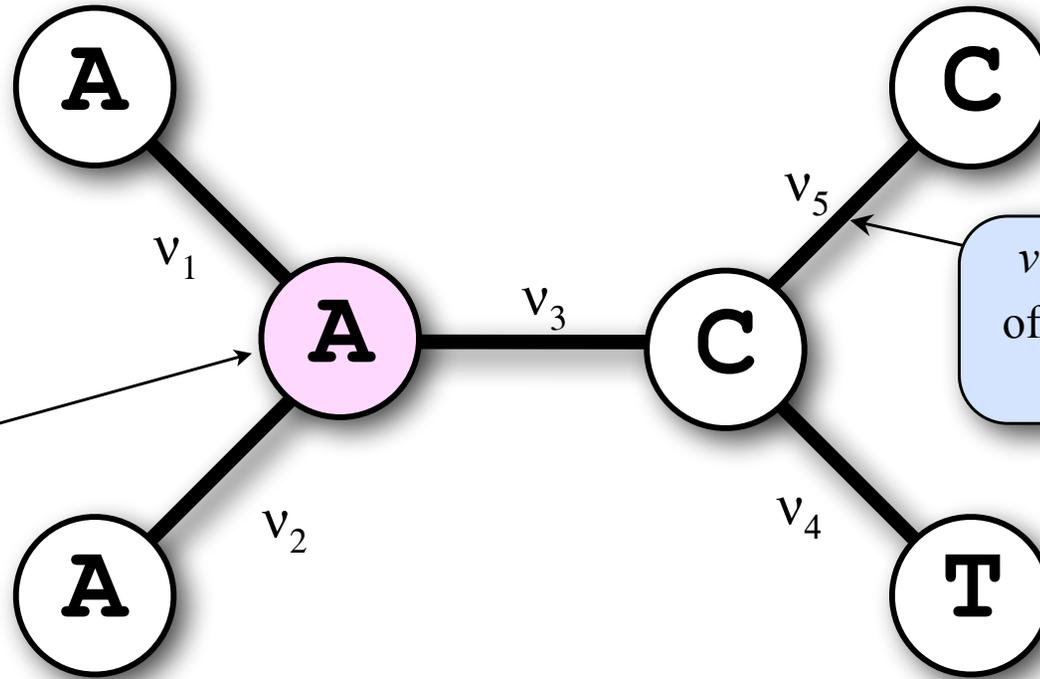
(data shown for only one site)



From slide 6



Likelihood for site k



v_5 is the expected number of substitutions for just this one branch

π_A

$$L_k = \frac{1}{4} \left[\frac{1}{4} + \frac{3}{4} e^{-4v_1/3} \right] \left[\frac{1}{4} + \frac{3}{4} e^{-4v_2/3} \right] \left[\frac{1}{4} - \frac{1}{4} e^{-4v_3/3} \right] \left[\frac{1}{4} - \frac{1}{4} e^{-4v_4/3} \right] \left[\frac{1}{4} + \frac{3}{4} e^{-4v_5/3} \right]$$

$P_{AA}(v_1)$

$P_{AA}(v_2)$

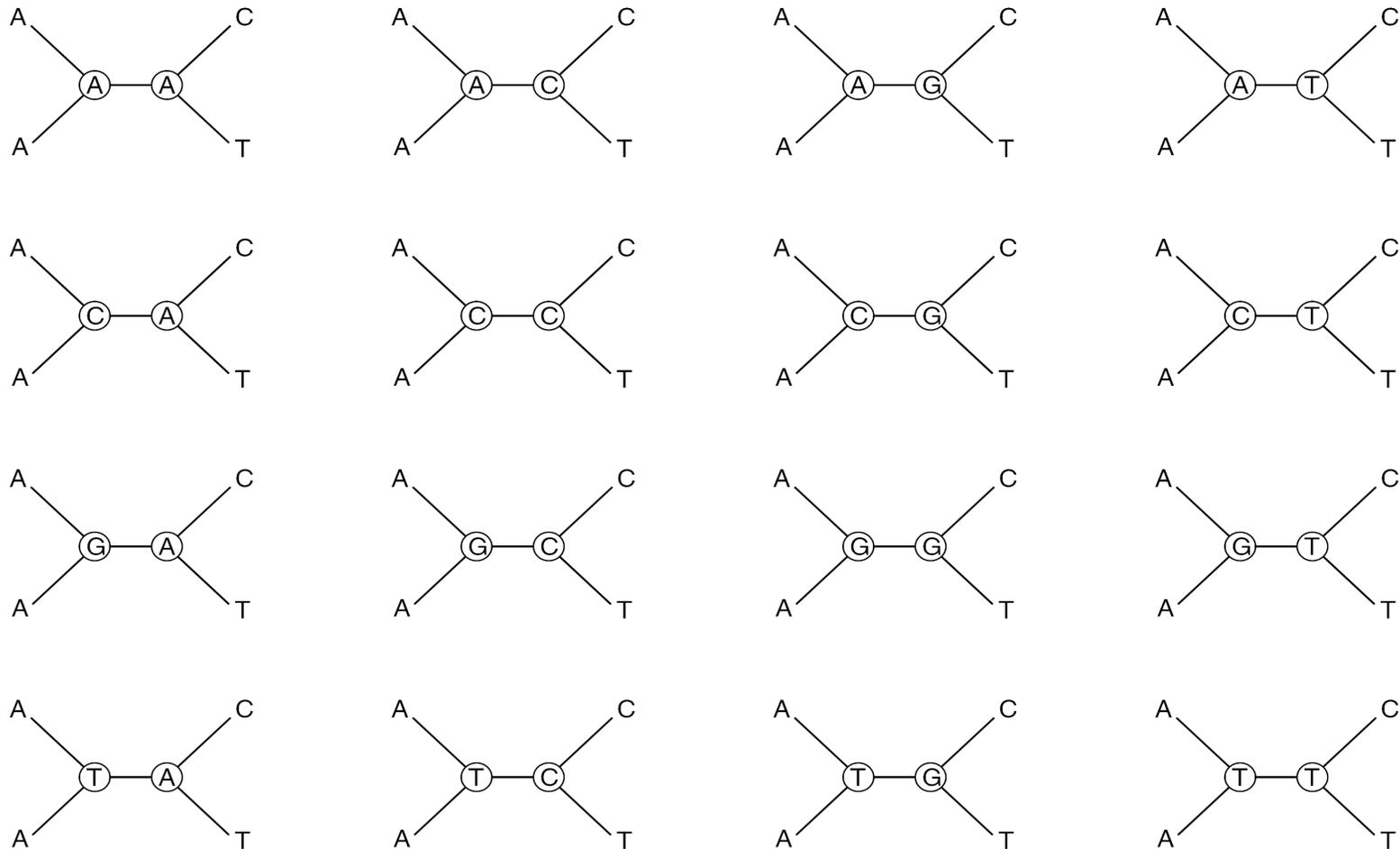
$P_{AC}(v_3)$

$P_{CT}(v_4)$

$P_{CC}(v_5)$

Note use of the AND probability rule

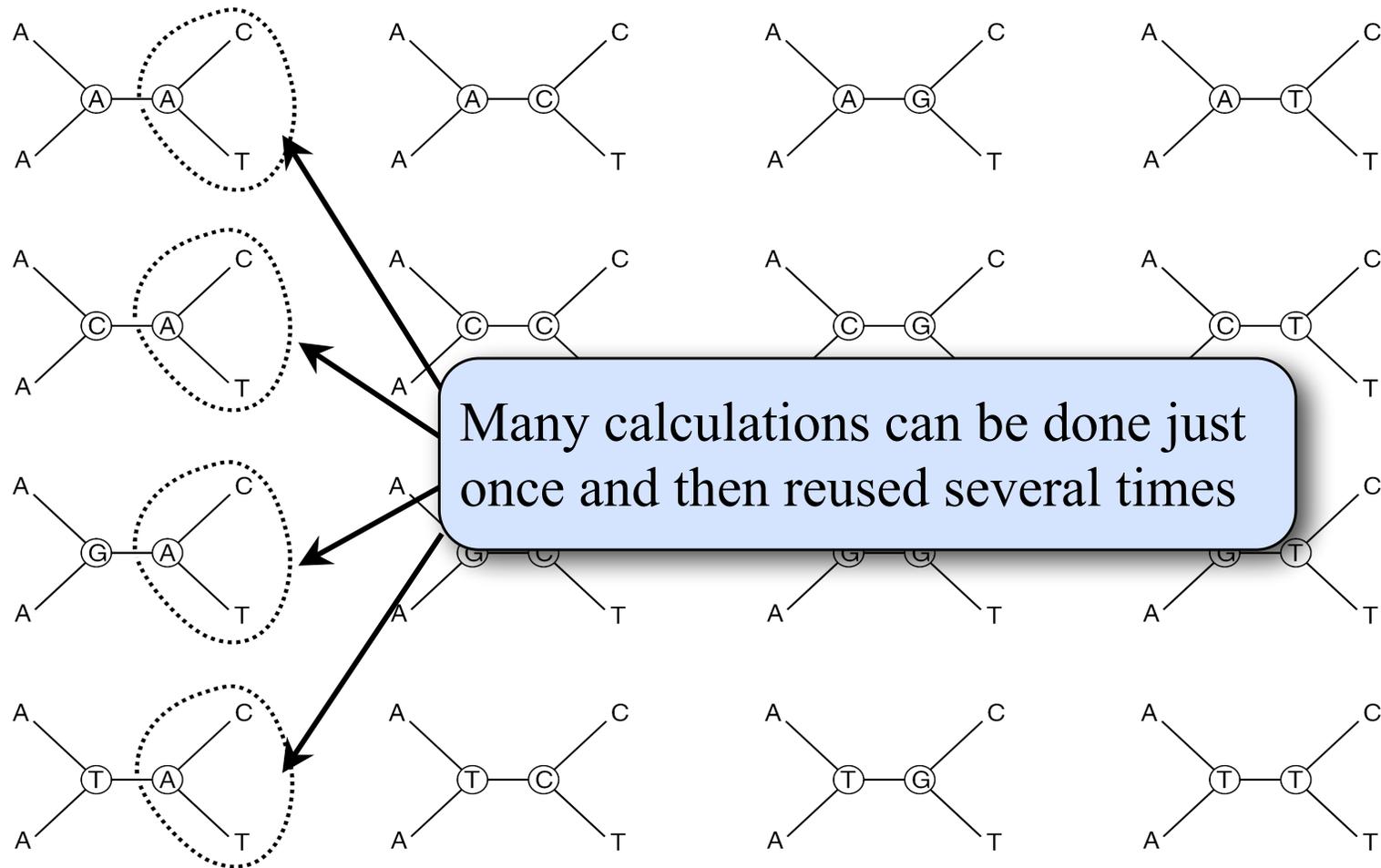
Brute force approach would be to calculate L_k for all 16 combinations of ancestral states and sum them



Note use of the OR probability rule

Pruning algorithm

(same result, less time)



Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of Molecular Evolution* 17:368-376

Substitution Models

Jukes-Cantor (JC69) model

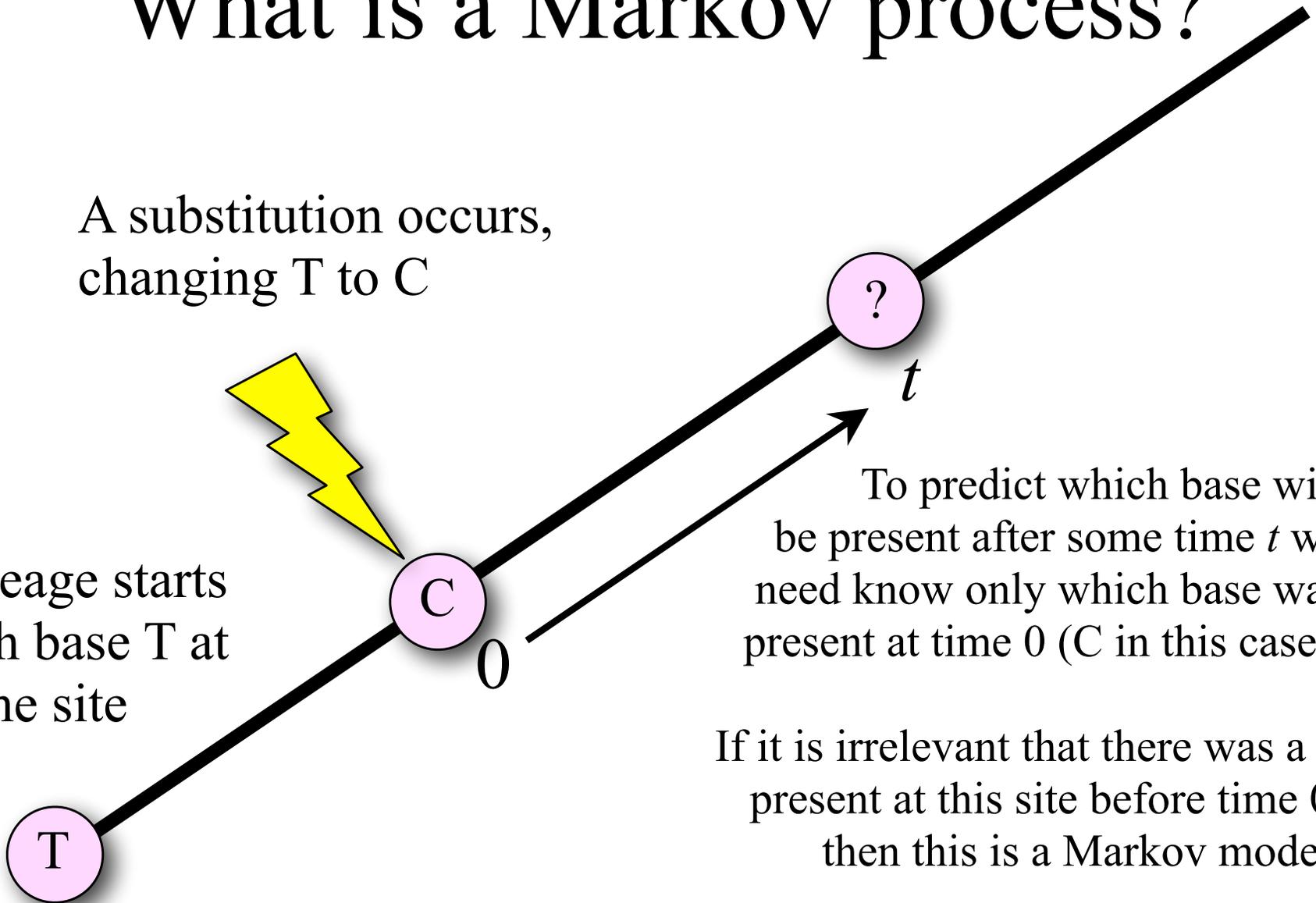
- The four bases (A, C, G, T) are expected to be **equally frequent** in sequences ($\pi_A = \pi_C = \pi_G = \pi_T = 0.25$)
- Assumes **same rate** for all types of substitution ($r_{A \rightarrow C} = r_{A \rightarrow G} = r_{A \rightarrow T} = r_{C \rightarrow G} = r_{C \rightarrow T} = r_{G \rightarrow T} = \alpha$)
- Usually described as a **1-parameter** model (the parameter being the edge length)
 - Remember, however, that each edge in a tree can have its own length, so there are really as many parameters in the model as there are edges in the tree!
- Assumes substitution is a **Markov** process...

Jukes, T. H., and C. R. Cantor. 1969. Evolution of protein molecules. Pages 21-132 in H. N. Munro (ed.), *Mammalian Protein Metabolism*. Academic Press, New York.

What is a Markov process?

A substitution occurs,
changing T to C

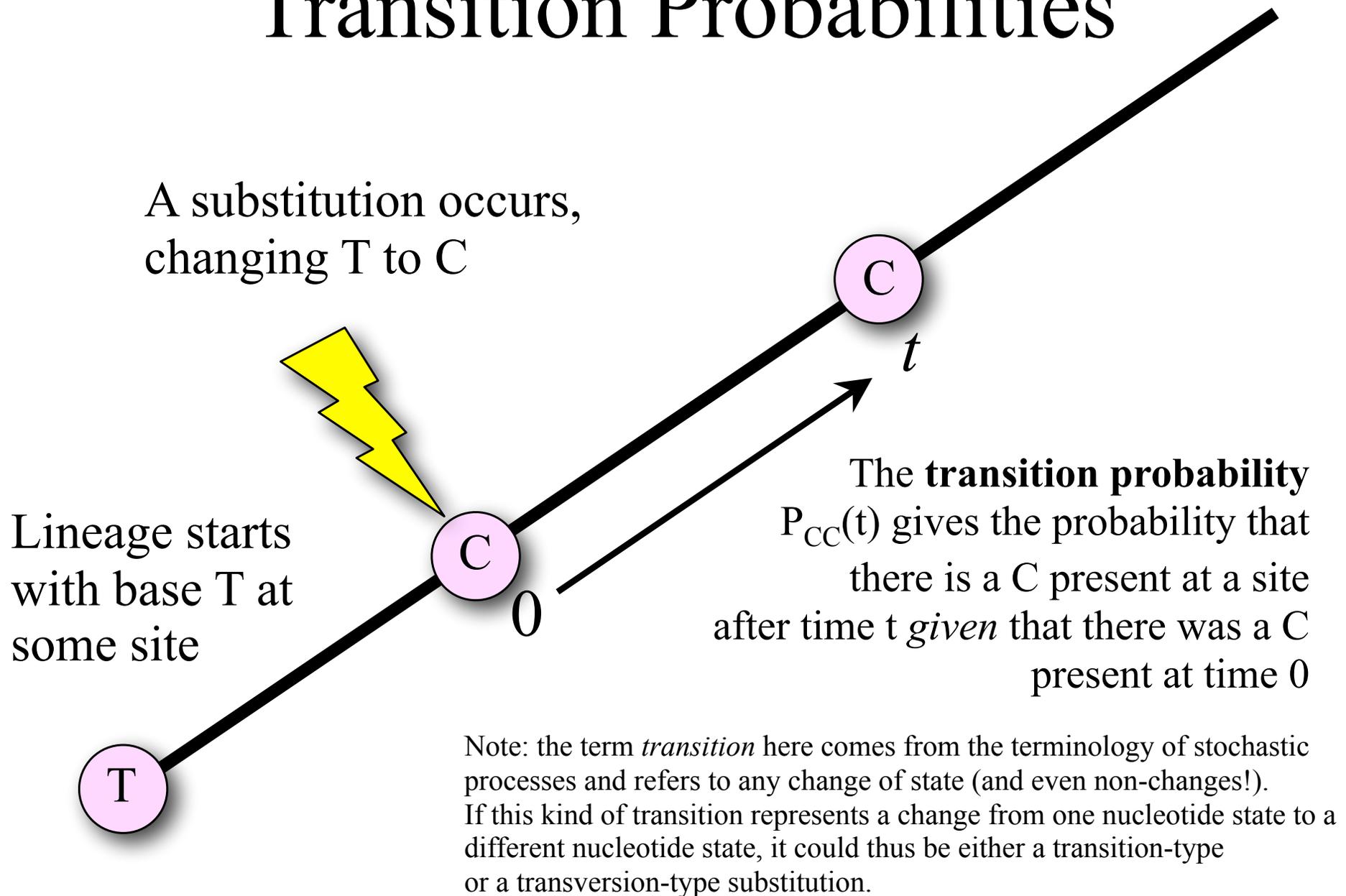
Lineage starts
with base T at
some site



To predict which base will
be present after some time t we
need know only which base was
present at time 0 (C in this case).

If it is irrelevant that there was a T
present at this site before time 0,
then this is a Markov model.

Transition Probabilities



Jukes-Cantor transition probabilities

Here is the probability that a site starting in state T will end up in state G after time t when the individual substitution rates are all α :

$$P_{TG}(t) = \frac{1}{4} (1 - e^{-4\alpha t})$$

The JC69 model has only one unknown quantity: αt

(The symbol e represents the base of the natural logarithms: its value is 2.718281828459045...)

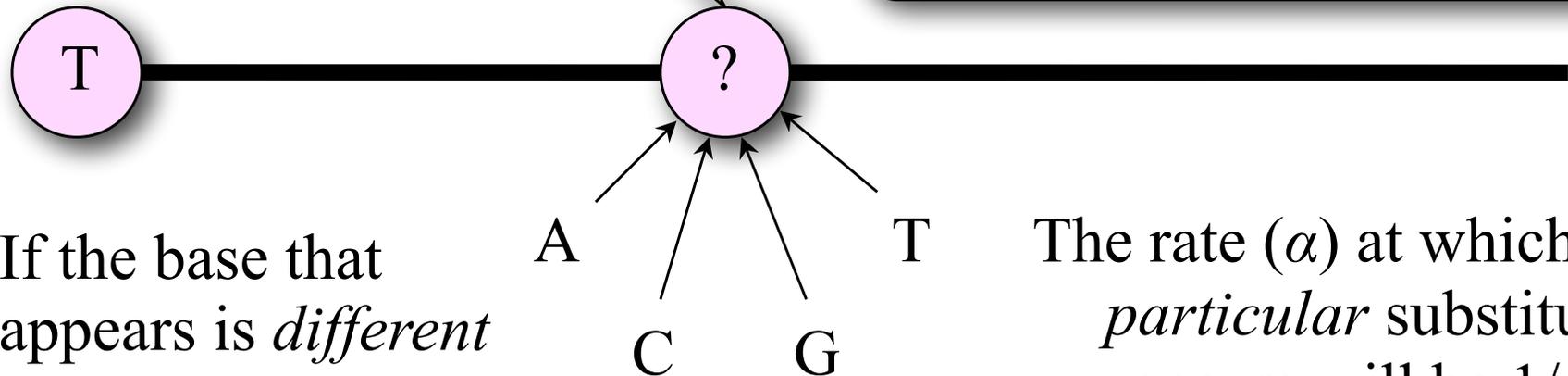
Where does a transition probability formula such as this come from?

"ACHNyons" vs. substitutions

ACHN =
"Anything
Can Happen
Now"

When an *achnyon* occurs, any base can appear in a sequence.

Note: *achnyon* is *my term* for this make-believe event. You will not see this term in the literature.



If the base that appears is *different* from the base that was already there, then a **substitution** event has occurred.

The rate (α) at which any *particular* substitution occurs will be 1/4 the *achnyon* rate (μ).
That is, $\alpha = \mu/4$
(or $\mu = 4\alpha$)

Deriving a transition probability

Calculate the probability that a site currently T will change to G over time t when the rate of this particular substitution is α :

$$\Pr(\text{zero achnyons}) = e^{-\mu t} \quad (\text{Poisson probability of zero events})$$

$$\Pr(\text{at least 1 achnyon}) = 1 - e^{-\mu t}$$

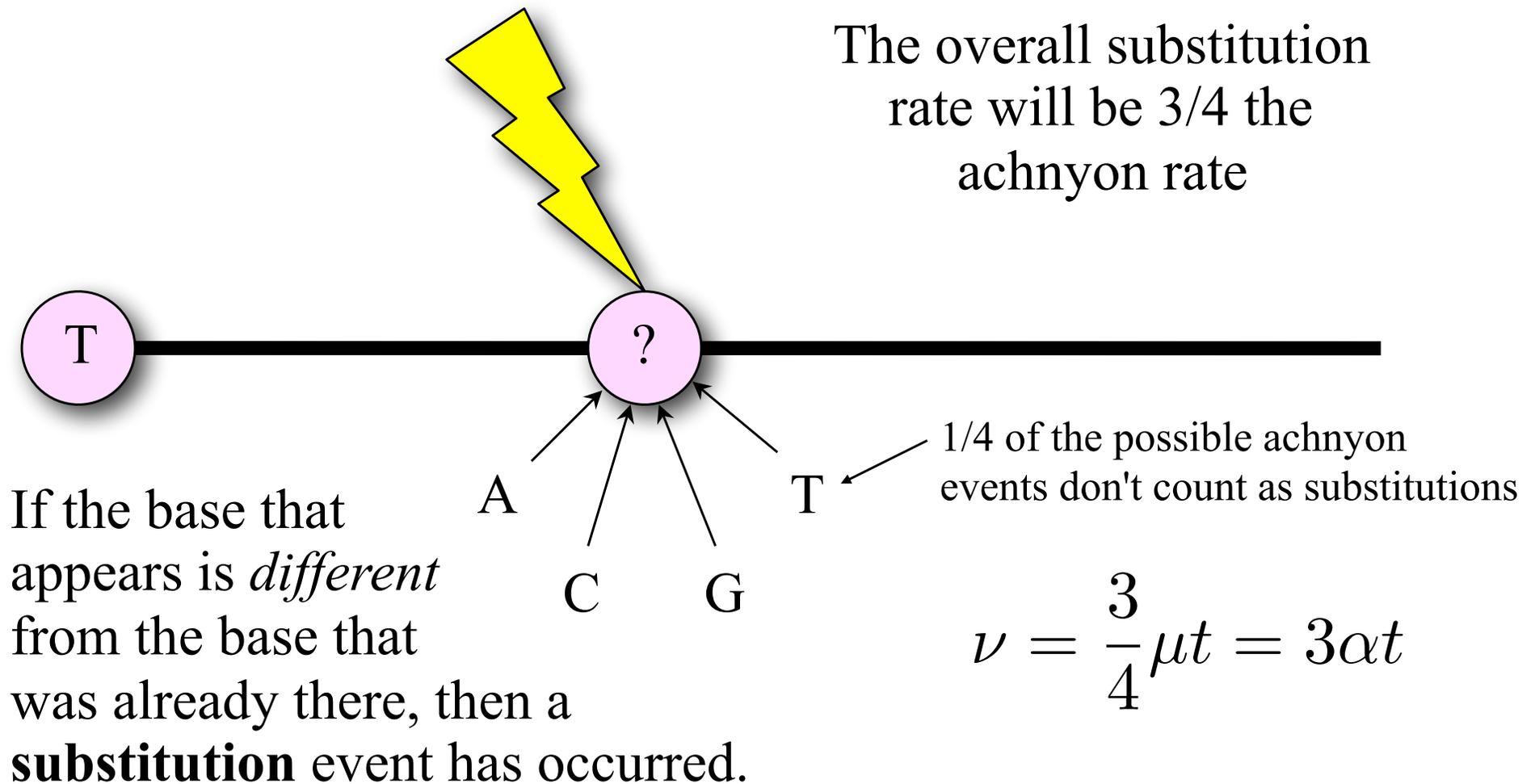
$$\Pr(\text{last achnyon results in base G}) = \frac{1}{4}$$

$$\Pr(\text{end in G} \mid \text{start in T}) = \frac{1}{4} (1 - e^{-\mu t})$$

Remember that the rate (α) of any particular substitution is one fourth the achnyon rate (μ):

$$P_{GT}(t) = \frac{1}{4} (1 - e^{-4\alpha t})$$

Expected number of substitutions



Transition Probabilities: Remarks

$$P_{TA}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TC}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TG}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TT}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$1 - e^{-4\alpha t}$$

These should add to 1.0 because T *must* change to something!

Doh! Something must be wrong here...

Transition Probabilities: Remarks

$$P_{TA}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TC}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TG}(t) = \frac{1}{4}(1 - e^{-4\alpha t})$$

$$P_{TT}(t) = \frac{1}{4}(1 - e^{-4\alpha t}) + e^{-4\alpha t}$$

Forgot to account for the possibility of *no* acnyons over time t

Coffee Break

Equilibrium frequencies

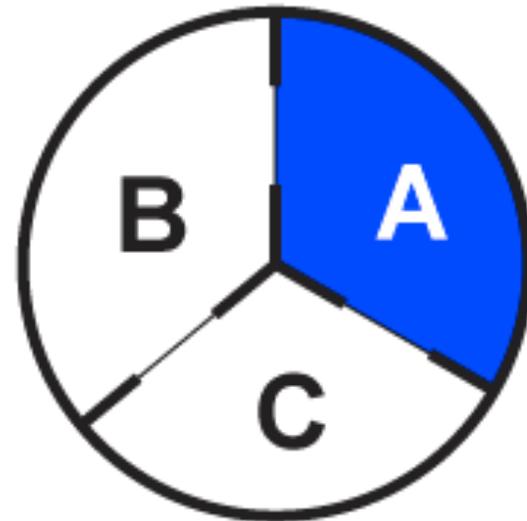
- The JC69 model assumes that the frequencies of the four bases (A, C, G, T) are equal
- The equilibrium relative frequency of each base is thus 0.25
- Why are they called *equilibrium* frequencies?

Equilibrium Frequencies

Imagine a bottle of perfume has been spilled in room A.

The doors to the other rooms are closed, so the perfume has, thus far, not been able to spread.

What would happen if we opened all the doors?

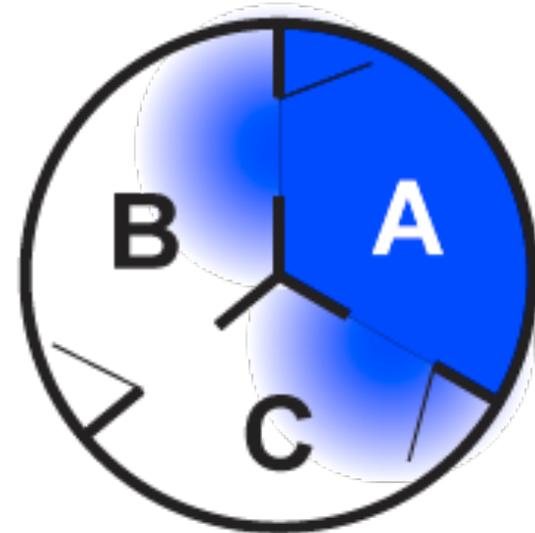


Equilibrium Frequencies

If the doors are suddenly opened, the perfume would begin diffusing from the area of highest concentration to lowest.

Molecules of perfume go both ways through open doors, but more pass one way than another, leading to a net flow from room A to rooms B and C.

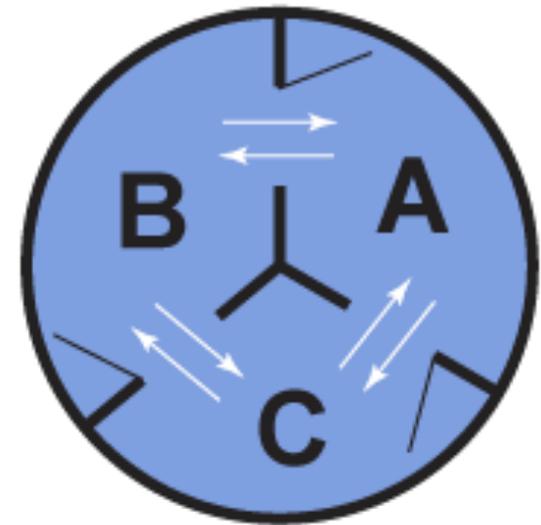
In the instant that the doors are opened, A is losing perfume molecules at *twice the rate* each of the other rooms is gaining molecules. As diffusion progresses, however, the rate of loss from A drops, approaching an equilibrium.



Equilibrium Frequencies

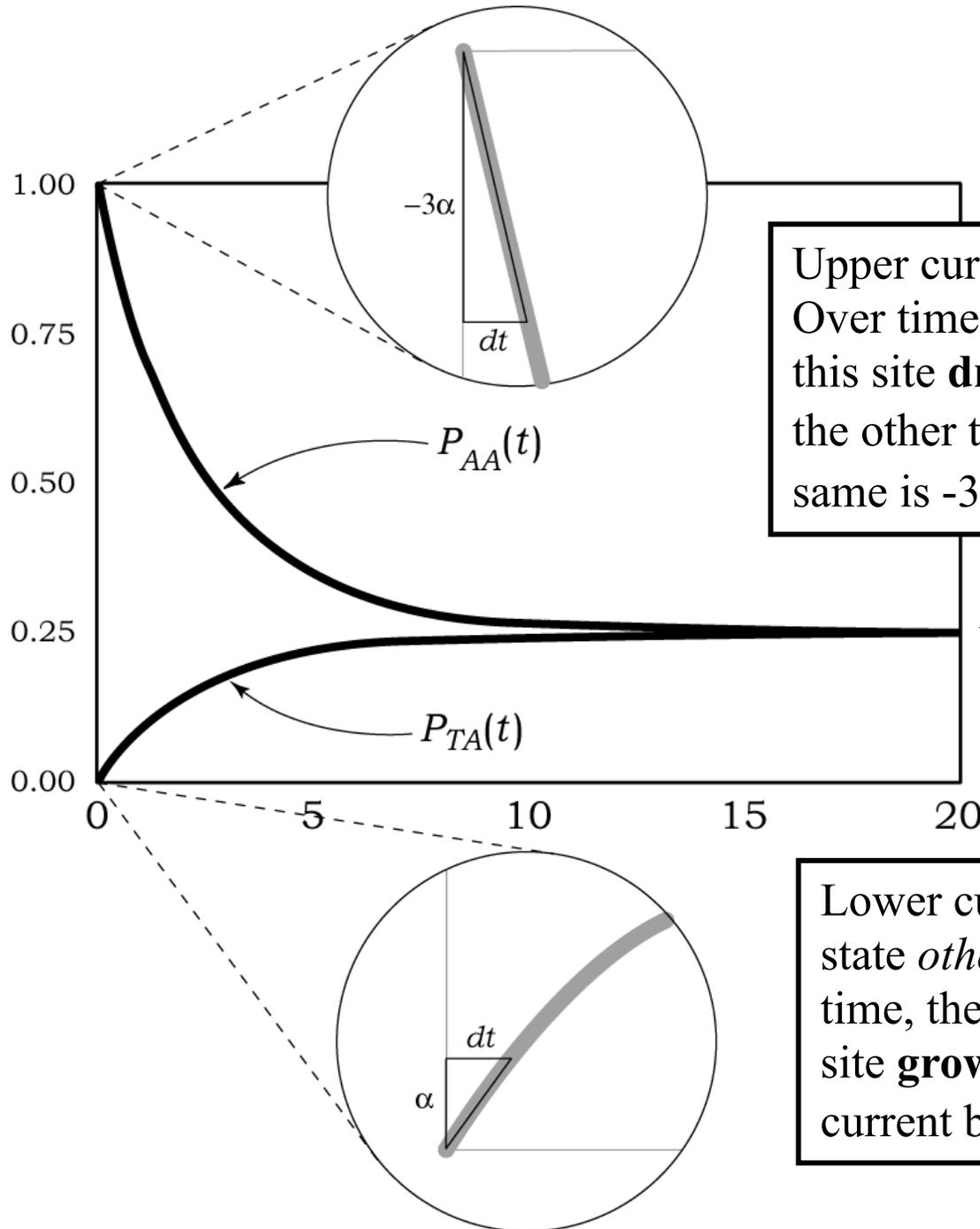
Eventually, all 3 rooms have essentially the same concentration of perfume.

Molecules still move through doors, but now the rates are the same in all directions.



Back to sequence evolution: assume a sequence began with only A nucleotides (a poly-A sequence). Over time, substitution would begin converting some of these As to Cs, Gs, and Ts, just as the perfume diffused into adjacent rooms.

$\Pr(A|A)$ and $\Pr(A|T)$ as a function of time



Upper curve assumes we started with A at time 0. Over time, the probability of still seeing an A at this site **drops** because rate of changing to one of the other three bases is 3α (so rate of staying the same is -3α).

The equilibrium relative frequency of A is 0.25

Lower curve assumes we started with some state *other* than A (T is used here). Over time, the probability of seeing an A at this site **grows** because the rate at which the current base will change into an A is α .

JC69 rate matrix

1 parameter:
 α

		To			
		A	C	G	T
From	A	-3α	α	α	α
	C	α	-3α	α	α
	G	α	α	-3α	α
	T	α	α	α	-3α

Jukes, T. H., and C. R. Cantor. 1969. Evolution of protein molecules. Pages 21-132 in H. N. Munro (ed.), *Mammalian Protein Metabolism*. Academic Press, New York.

K80 (or K2P) rate matrix

2 parameters:

α
 β

		To			
		A	C	G	T
From	A	$-\alpha - 2\beta$	β	α	β
	C	β	$-\alpha - 2\beta$	β	α
	G	α	β	$-\alpha - 2\beta$	β
	T	β	α	β	$-\alpha - 2\beta$

↑ transition rate ↑ transversion rate

Kimura, M. 1980. A simple method for estimating evolutionary rate of base substitutions through comparative studies of nucleotide sequences. *Journal of Molecular Evolution* 16:111-120.

K80 rate matrix

(looks different, but actually the same)

2 parameters:

κ

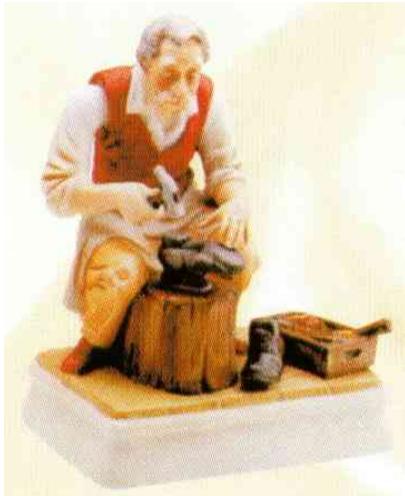
β

	A	C	G	T
A	$-\beta(\kappa + 2)$	β	$\kappa\beta$	β
C	β	$-\beta(\kappa + 2)$	β	$\kappa\beta$
G	$\kappa\beta$	β	$-\beta(\kappa + 2)$	β
T	β	$\kappa\beta$	β	$-\beta(\kappa + 2)$

All I've done is re-parameterize the rate matrix,
 letting κ equal the *transition/transversion rate ratio* $\longrightarrow \kappa = \frac{\alpha}{\beta}$

Note: the K80 model is identical to the JC69 model if $\kappa = 1$ ($\alpha = \beta$)

Transition/transversion ratio (τ) versus Transition/transversion *rate* ratio (κ)



Cobbler analogy:

- 4 cobblers in a factory make loafers
- 8 cobblers in the factory make work boots
- all cobblers produce the same number of shoes per unit time, regardless of shoe type
- what is the loafer/boot *rate ratio* and how does that compare to the loafer/boot *ratio*?

The loafer/boot *rate ratio* is 1.0 because each cobbler cranks out shoes at the same rate.

The loafer/boot *ratio*, however, is 0.5 because there are twice as many cobblers making boots as there are cobblers making loafers.

There are 8 possible transversion-type substitutions and only 4 possible transition-type substitutions: the transition/transversion ratio is thus 0.5 when the transition/transversion rate ratio is 1.

F81 rate matrix

4 parameters:

μ

π_A

π_C

π_G

	A	C	G	T
A	$-\mu(1 - \pi_A)$	$\pi_C\mu$	$\pi_G\mu$	$\pi_T\mu$
C	$\pi_A\mu$	$-\mu(1 - \pi_C)$	$\pi_G\mu$	$\pi_T\mu$
G	$\pi_A\mu$	$\pi_C\mu$	$-\mu(1 - \pi_G)$	$\pi_T\mu$
T	$\pi_A\mu$	$\pi_C\mu$	$\pi_G\mu$	$-\mu(1 - \pi_T)$

Note: the F81 model is identical to the JC69 model if all base frequencies are equal

HKY85 rate matrix

5 parameters:

κ
 β
 π_A
 π_C
 π_G

	A	C	G	T
A	—	$\pi_C \beta$	$\pi_G \beta \kappa$	$\pi_T \beta$
C	$\pi_A \beta$	—	$\pi_G \beta$	$\pi_T \beta \kappa$
G	$\pi_A \beta \kappa$	$\pi_C \beta$	—	$\pi_T \beta$
T	$\pi_A \beta$	$\pi_C \beta \kappa$	$\pi_G \beta$	—

A dash means equal to negative sum of other elements on the same row

Note: the HKY85 model is identical to the F81 model if $\kappa = 1$. If, in addition, all base frequencies are equal, it is identical to JC69.

F84 vs. HKY85

F84 model:

μ rate of process generating *all types of substitutions*

$k\mu$ rate of process generating *only transitions*

Becomes F81 model if $k = 0$

HKY85 model:

β rate of process generating *only transversions*

$\kappa\beta$ rate of process generating *only transitions*

Becomes F81 model if $\kappa = 1$

F84 first used in Felsenstein's PHYLIP package in 1984, first published by: Kishino, H., and M. Hasegawa. 1989. Evaluation of the maximum likelihood estimate of the evolutionary tree topologies from DNA sequence data, and the branching order in hominoidea. *Journal of Molecular Evolution* 29: 170-179.

GTR rate matrix

	A	C	G	T
A	—	$\pi_C a \mu$	$\pi_G b \mu$	$\pi_T c \mu$
C	$\pi_A a \mu$	—	$\pi_G d \mu$	$\pi_T e \mu$
G	$\pi_A b \mu$	$\pi_C d \mu$	—	$\pi_T f \mu$
T	$\pi_A c \mu$	$\pi_C e \mu$	$\pi_G f \mu$	—

9 parameters:

π_A
 π_C
 π_G
 a
 b
 c
 d
 e
 μ

Identical to the F81 model if $a = b = c = d = e = f = 1$. If, in addition, all the base frequencies are equal, GTR is identical to JC69. If $a = c = d = f = \beta$ and $b = e = \kappa\beta$, GTR becomes the HKY85 model.

Lanave, C., G. Preparata, C. Saccone, and G. Serio. 1984.
 A new method for calculating evolutionary substitution
 rates. *Journal of Molecular Evolution* 20:86-93.

Rate Heterogeneity

Green Plant *rbcL*

First 88 amino acids (translation is for *Zea mays*)

M--S--P--Q--T--E--T--K--A--S--V--G--F--K--A--G--V--K--D--Y--K--L--T--Y--Y--T--P--E--Y--E--T--K--D--T--D--I--L--A--A--F--R--V--T--P--
 Chara (green alga; land plant lineage) AAAGATTACAGATTAACCTACTATACTCCTGAGTATAAACTAAAGATACTGACATTTTAGCTGCATTTTCGTGTAAGTCCA
 Chlorella (green alga)C...C.T.....T..CC..C.A.....C.....T...C.T..A..G..C...A.G.....T
 Volvox (green alga)TC.T....A....C..A....GT.GTA.....C.....C.....A.....A.G.....T
 Conocephalum (liverwort)TC.....T.....G..T...G.....G..T.....A.....A.AA.G.....T
 Bazzania (moss)T.....C..T....G....A...G.G..C....G..A..T....G..A.....A.G.....C
 Anthoceros (hornwort)T.....CC.T....C....T..CG.G..C..G.....T....G..A..G.C.T.AA.G.....T
 Osmunda (fern)TC...G...C.....C..T...G.G..C..G.....T....G..A....C..AA.G.....C
 Lycopodium (club "moss") .GG.....C.T..C.....T....G..C....A..C..T..C.G..A.....AA.G.....T
 Ginkgo (gymnosperm; Ginkgo biloba)G....T.....A..C...C.....T..C..G..A....C..A.....T
 Picea (gymnosperm; spruce)T.....T.....A...C.G..C.....G..T....G..A....C..A.....T
 Iris (flowering plant)G....T.....T..CG...C.....T..C..G..A....C..A.....T
 Asplenium (fern; spleenwort)TC..C.G...T..C..C..C..A..C..G..C.....C..T..C..G..A..T..C..GA.G..C...
 Nicotiana (flowering plant; tobacco)G...A...G.....T.....CC...C..G.....T..A..G..A....C..A.....T

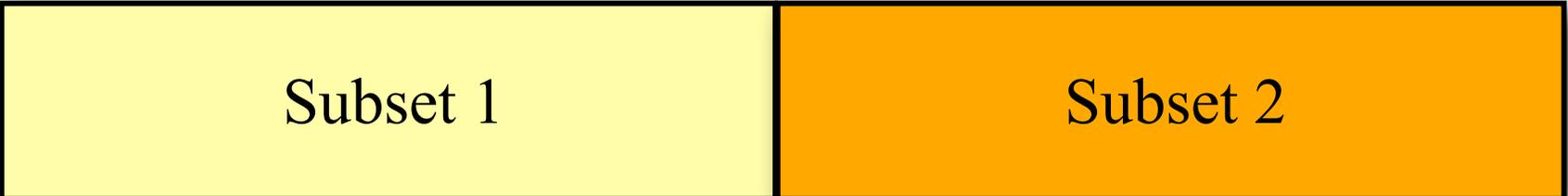
Q--L--G--V--P--P--E--E--A--G--A--A--V--A--A--E--S--S--T--G--T--W--T--T--V--W--T--D--G--L--T--S--L--D--R--Y--K--G--R--C--Y--H--I--E--
 CAACCTGGCGTTCCACCTGAAGAAGCAGGGGCTGCAGTAGCTGCAGAATCTTCTACTGGTACATGGACTACTGTTTGGACTGACGGATTAAGTACTAGTTTGGACCGATACAAAGGAAGATGCTACGATATTGAA
A..T.....A.....G..T..G.....A....A..A.....T....G....A.....T..T.....A.....T.....TC.T..T..T..C..C..G
A..T.....TGT..T....T..T....T....A..A..A.....T....A....A.....T..T.....A..C.T.....TC.T..T..T..C..C..G
 ..G...G..A...G..A.....A..A...T....T.....A.....T..TC.T...ACC.T..T..T..T.....TC.....T.G.....C
G..A..A.....A..G.....T....A..C...G...C..G.....C..T..GC.T..A..C.C..T..T.....TC.....T..C..C..
 T...A..G..G.....A..C.....T.....A.....C..T...C.T..C..CC.T...T.....TC.....C.....
C..A..A..GG...G....T..A.....G.....A....G...C....A...G..T..C.T..C..C.T..T..T..T..G..TC.....
T..A..A...C..G...G..A..C.....T....C.....C.....C.....C..T..C.T..C..C.C.T..C.....TC.G...T..A..
A..G...G...G...A...C.....C.....C.....C.....C..T..C.T..C..C.T..T..T.....G.....T..C..C..G
A..G..G..G..C..G...G..A..A.....T.....C..C.....C.....C..T..C.T.....C.T..T..T.....G..GC.....T..C..C..G
C..A...TG.....G...C..G...C.....A..A..G.....T..C.T..C..C.T..T..T.....C.....C.C..C..G
C..A..A..G.....C..A.....G..C...A.....C...G...A...G..G..C..CC.T...T...G..CC.....C..G
A.....C..G...C.....C.....A...A.....C..T..C.T..C..CC.T..T..T.....GC.....CGC..C..G

All four bases are observed at some sites...

...while at other sites, only one base is observed

Site-specific rates

Each defined subset (e.g. gene, codon position) has its own relative rate



Subset 1

Subset 2

r_1 applies to subset 1
(e.g. sites 1 - 1000)

r_2 applies to subset 2
(e.g. sites 1001-2000)

Relative rates have mean 1:

$$\frac{r_1 + r_2}{2} = 1$$

More generally:

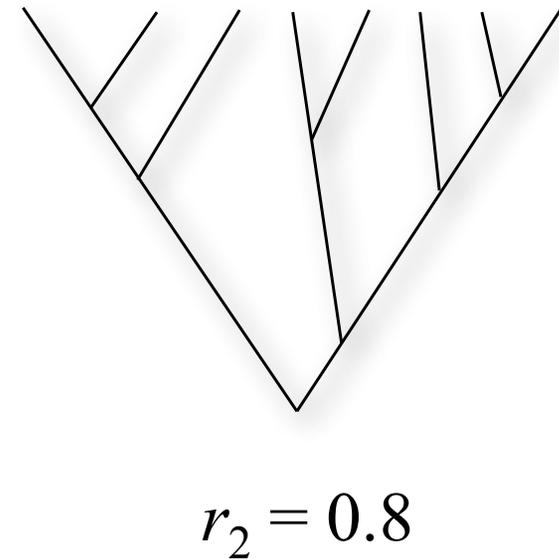
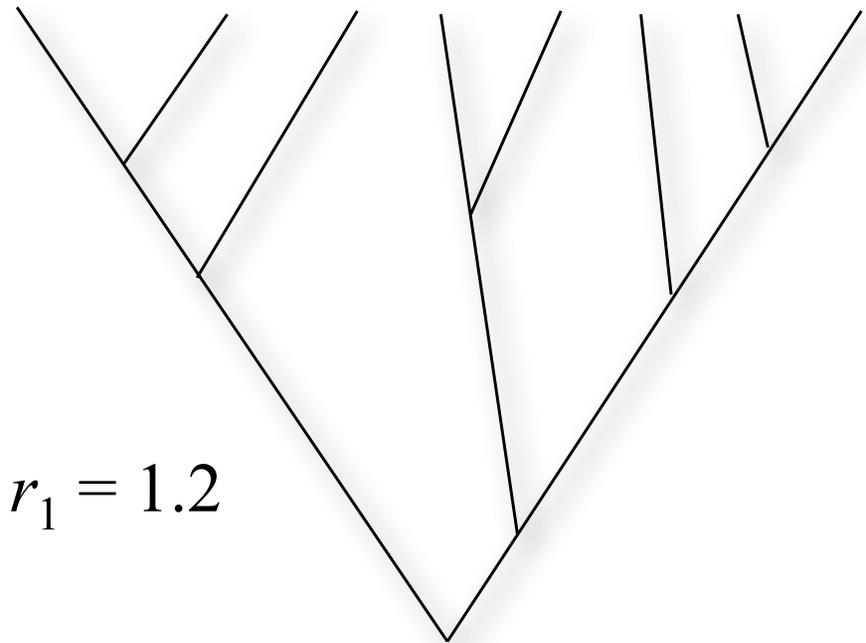
$$r_1 p(r_1) + r_2 p(r_2) = 1$$

Site-specific rates

$$L = \underbrace{\Pr(D_1|r_1) \cdots \Pr(D_{1000}|r_1)}_{\text{Gene 1}} \underbrace{\Pr(D_{1001}|r_2) \cdots \Pr(D_{2000}|r_2)}_{\text{Gene 2}}$$

Gene 1

Gene 2



Site-specific rates

JC69 transition probabilities that would be used for every site if rate *homogeneity* were assumed:

$$P_{ii}(t) = \frac{1}{4} + \frac{3}{4}e^{-4\alpha t}$$

$$P_{ij}(t) = \frac{1}{4} - \frac{1}{4}e^{-4\alpha t}$$

Site specific rates

JC69 transition probabilities that would be used for sites in **gene 1**:

$$P_{ii}(t) = \frac{1}{4} + \frac{3}{4}e^{-4r_1\alpha t}$$

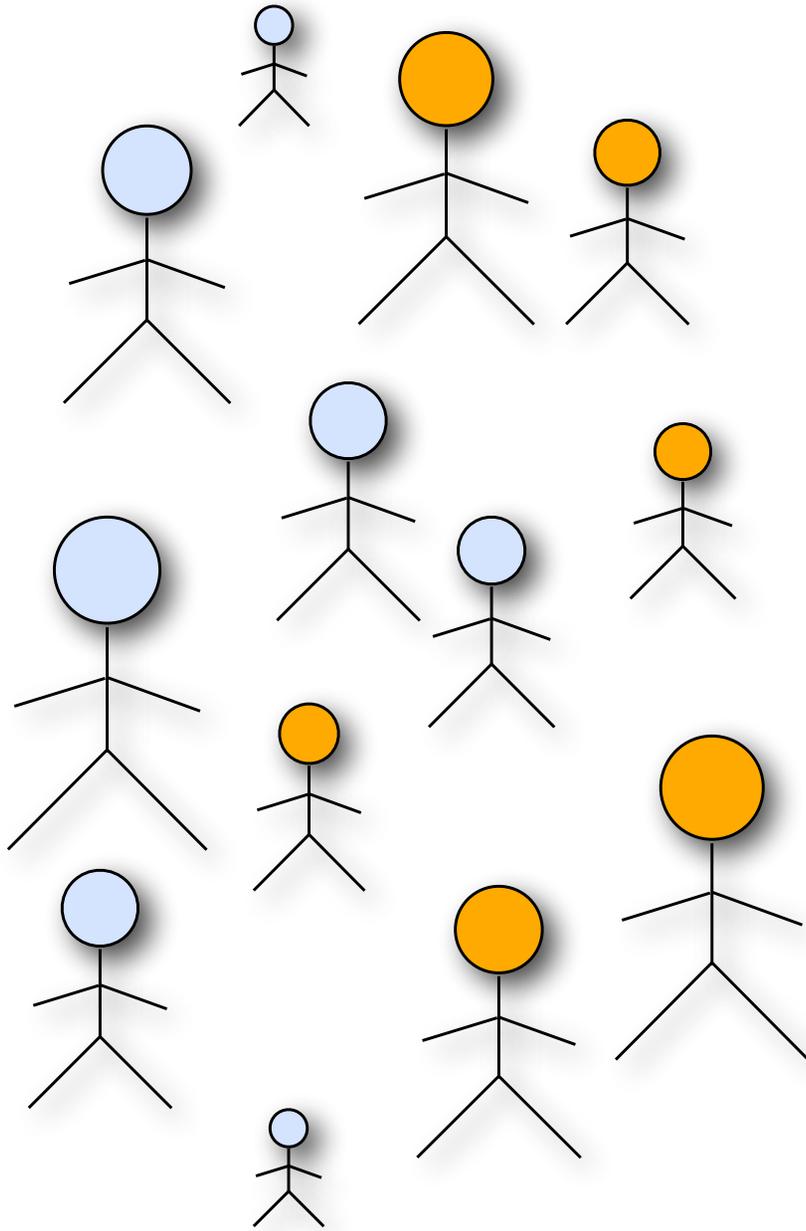
$$P_{ij}(t) = \frac{1}{4} - \frac{1}{4}e^{-4r_1\alpha t}$$

JC69 transition probabilities that would be used for sites in **gene 2**:

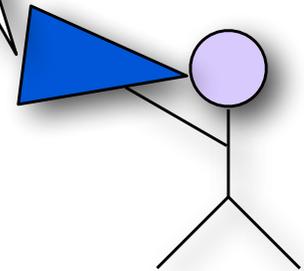
$$P_{ii}(t) = \frac{1}{4} + \frac{3}{4}e^{-4r_2\alpha t}$$

$$P_{ij}(t) = \frac{1}{4} - \frac{1}{4}e^{-4r_2\alpha t}$$

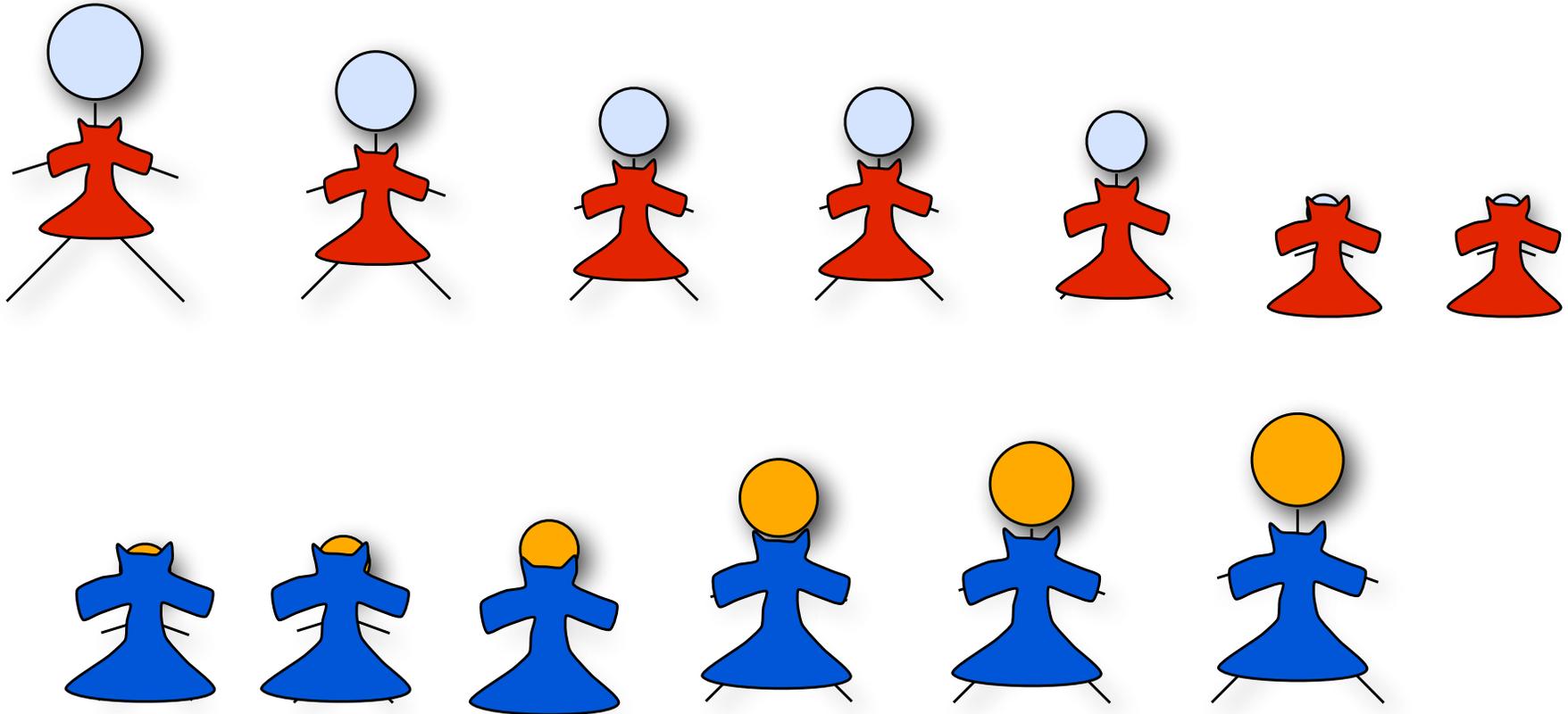
Site-specific Approach



Ok, I am going to divide you into 2 groups based on the color of your head, and everyone in each group will get a coat of the average size for their group. Very sorry if this does not work well for some people who are unusually large or small compared to their group.



Site-specific Approach

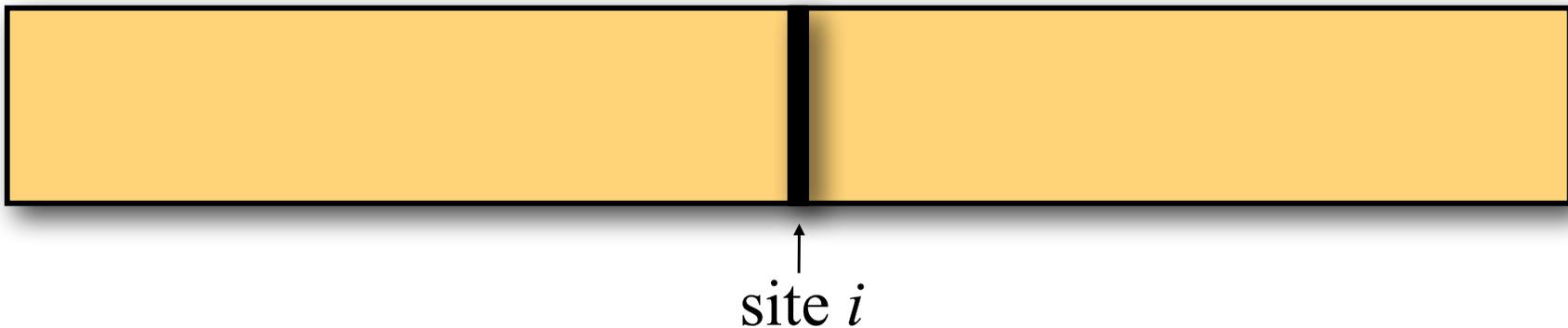


Good: costs less: need to buy just one coat for every person

Bad: every person in a group has to wear the same size coat, so the fit will be poor for some people if they are much bigger or smaller than the average size for the group in which they have been placed

Mixture Models

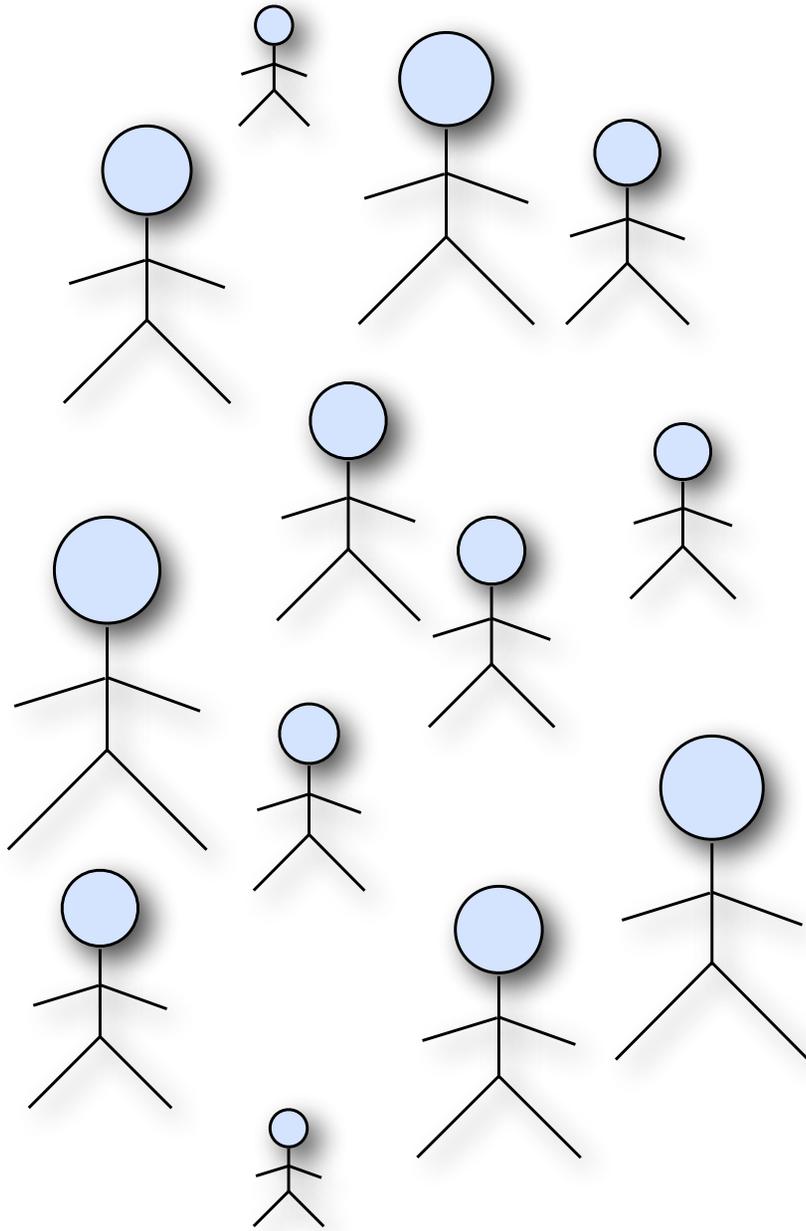
All relative rates applied to every site



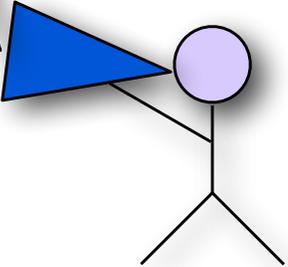
$$L_i = \Pr(D_i|r_1) \Pr(r_1) + \Pr(D_i|r_2) \Pr(r_2)$$

Common examples {
Invariable sites (I) model
Discrete Gamma (G) model

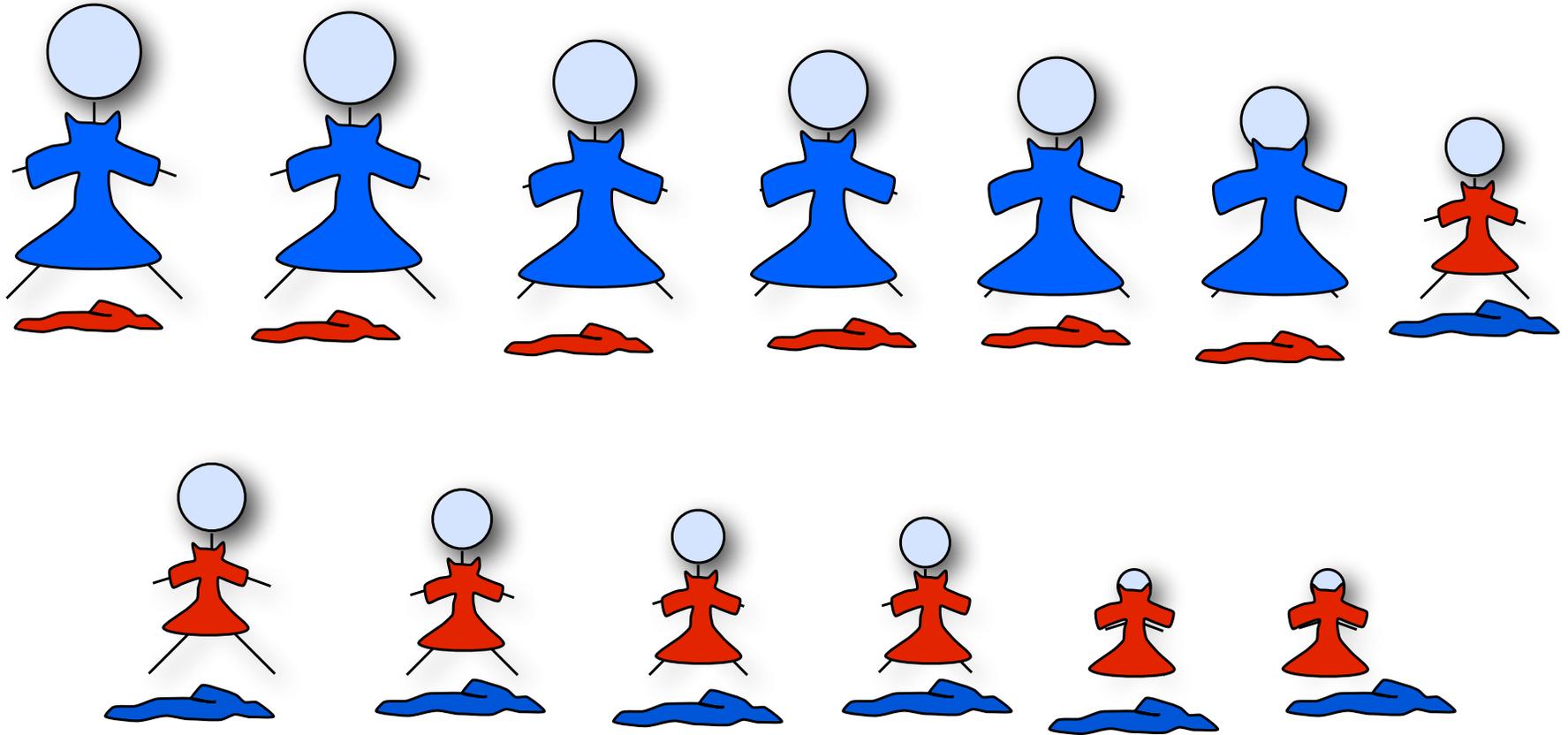
Mixture Model Approach



Ok, I am going to give each of you 2 coats: use the one that fits you best and throw away the other one. This costs twice as much for me, but on average leads to better fit for you. I have determined the two sizes of coats based on the distribution of your sizes.



Mixture Model Approach

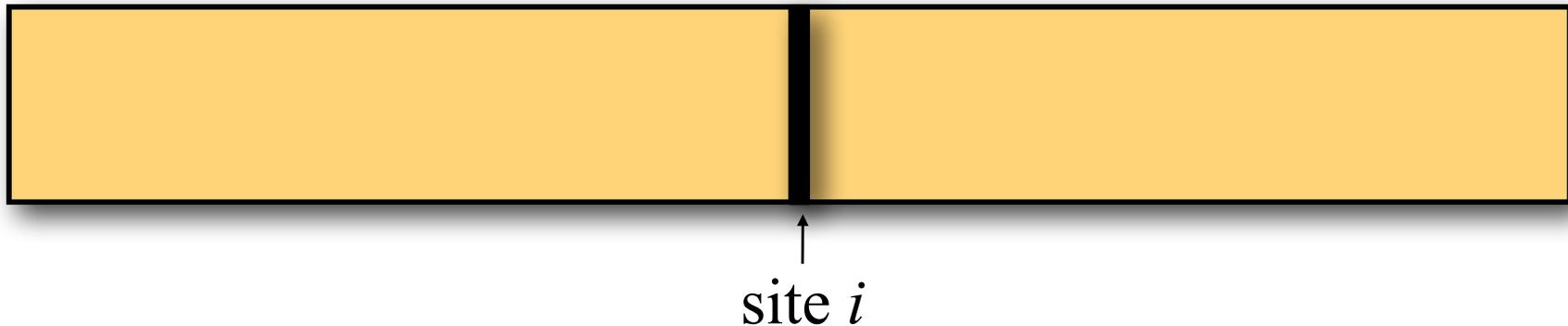


Good: every person experiences better fit because they can choose the size coat that fits best

Bad: costs more because two coats must be provided for each person

Invariable Sites Model

A fraction p_{invar} of sites are assumed to be invariable (i.e. rate = 0.0)



$$L_i = \Pr(D_i | r_1) p_{invar} + \Pr(D_i | r_2) (1 - p_{invar})$$

$$r_1 = 0.0$$

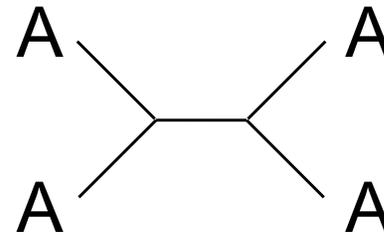
$$r_2 = \frac{1}{1 - p_{invar}}$$

Allows for the possibility that any given site could be variable or invariable

Reeves, J. H. 1992. Heterogeneity in the substitution process of amino acid sites of proteins coded for by mitochondrial DNA. *Journal of Molecular Evolution* 35:17-31.

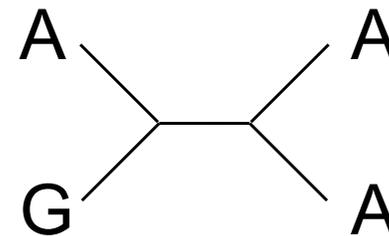
Invariable sites model

If site i is a *constant* site, both terms will contribute to the site likelihood:



$$L_i = \Pr(D_i|0.0)p_{\text{invar}} + \Pr(D_i|r_2)(1 - p_{\text{invar}})$$

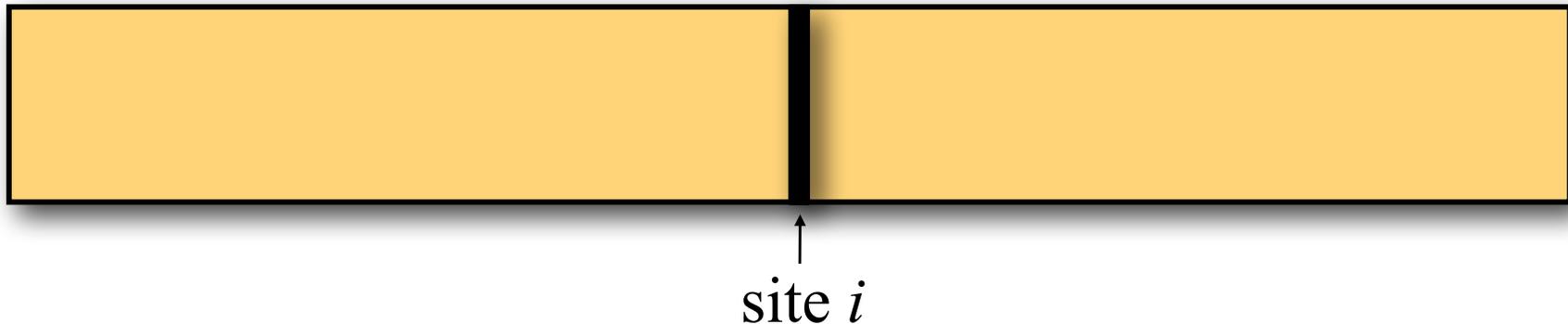
If site i is a *variable* site, there is no way to explain the data with a zero rate, so the first term is zero:



$$L_i = \cancel{\Pr(D_i|0.0)p_{\text{invar}}} + \Pr(D_i|r_2)(1 - p_{\text{invar}})$$

Discrete Gamma Model

No relative rate is exactly 0.0, and all are equally probable



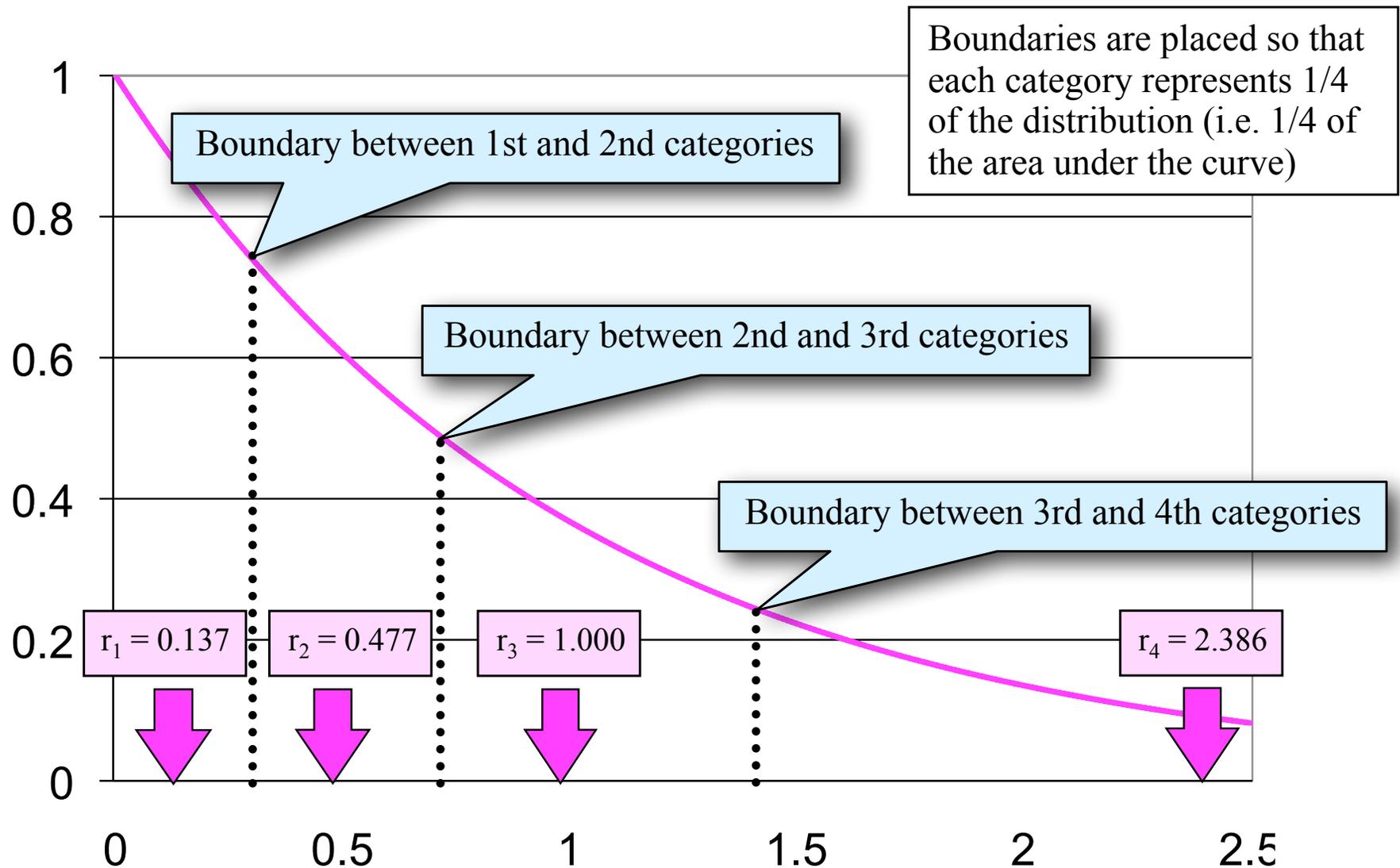
$$L = \left(\frac{1}{4}\right) \Pr(D_i|r_1) + \left(\frac{1}{4}\right) \Pr(D_i|r_2) + \left(\frac{1}{4}\right) \Pr(D_i|r_3) + \left(\frac{1}{4}\right) \Pr(D_i|r_4)$$

Relative rates are constrained to a discrete gamma distribution
Number of rate categories can vary (4 used here)

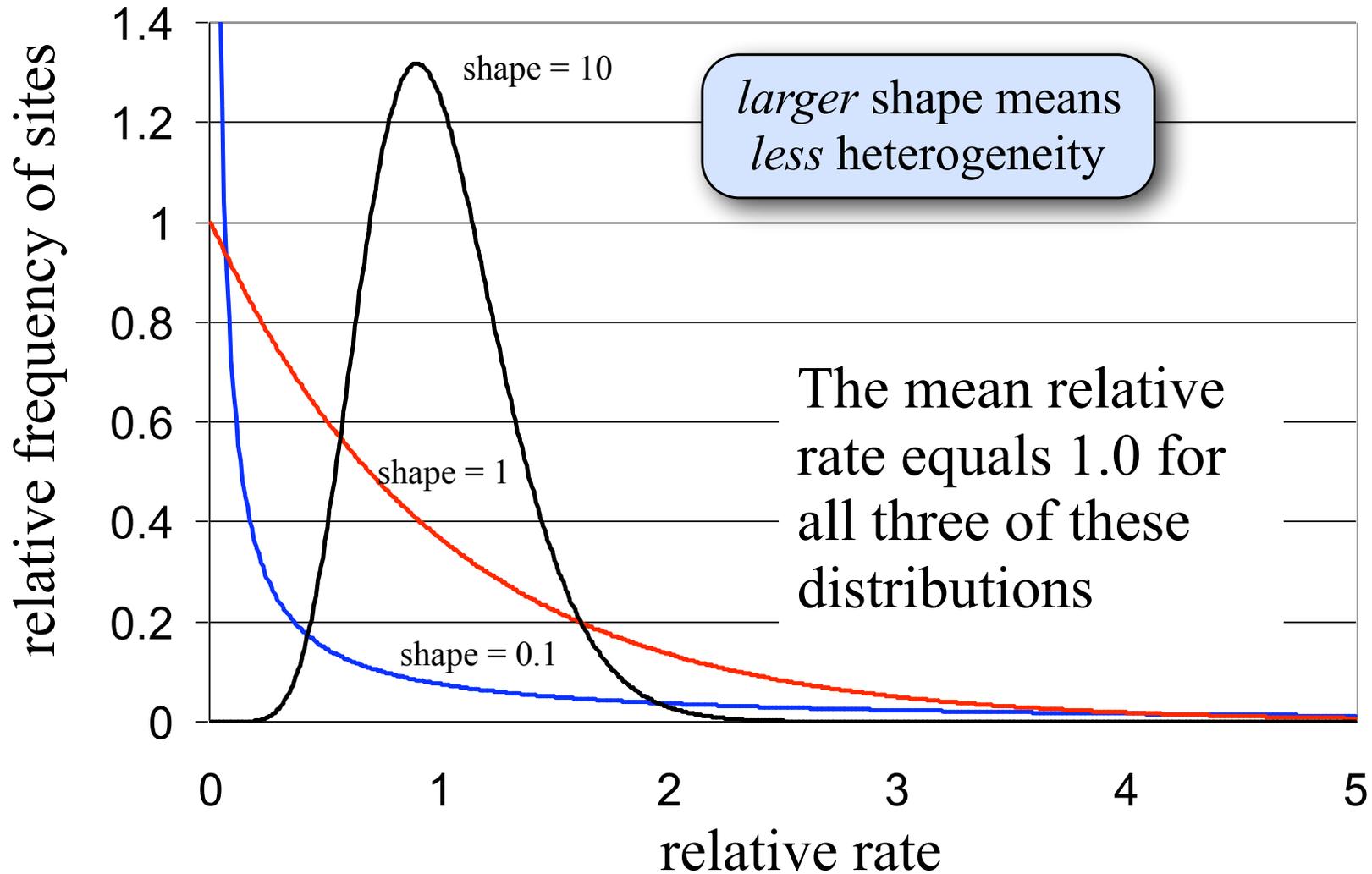
Yang, Z. 1993. Maximum-likelihood estimation of phylogeny from DNA sequences when substitution rates differ over sites. *Molecular Biology and Evolution* 10:1396-1401.

Yang, Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *Journal of Molecular Evolution* 39:306-314.

Relative rates in 4-category case



Gamma distributions



Codon models

Joe Bielawski will discuss codon models in greater detail tomorrow.

The Genetic Code

First 12 nucleotides at the 5' end of the *rbcL* gene in corn:

5' -ATG | TCA | CCA | CAA-3' coding strand
 3' -TAC | AGT | GGT | GTT-5' template strand

} DNA double helix

transcription

5' -AUG | UCA | CCA | CAA-3' mRNA

translation

N-Met | Ser | Pro | Gln-C polypeptide

Codon models treat codons as the independent units, not individual nucleotide sites.

Genetic Code

	U	C	A	G	
U	UUU Phe UUC Phe UUA Leu UUG Leu	UCU Ser UCC Ser UCA Ser UCG Ser	UAU Tyr UAC Tyr UAA Stop UAG Stop	UGU Cys UGC Cys UGA Stop UGG Trp	U C A G
C	CUU Leu CUC Leu CUA Leu CUG Leu	CCU Pro CCC Pro CCA Pro CCG Pro	CAU His CAC His CAA Gln CAG Gln	CGU Arg CGC Arg CGA Arg CGG Arg	U C A G
A	AUU Ile AUC Ile AUA Ile AUG Met	ACU Thr ACC Thr ACA Thr ACG Thr	AAU Asn AAC Asn AAA Lys AAG Lys	AGU Ser AGC Ser AGA Arg AGG Arg	U C A G
G	GUU Val GUC Val GUA Val GUG Val	GCU Ala GCC Ala GCA Ala GCG Ala	GAU Asp GAC Asp GAA Glu GAG Glu	GGU Gly GGC Gly GGA Gly GGG Gly	U C A G

<http://www.langara.bc.ca/biology/mario/Assets/Geneticcode.jpg>

First codon models

- Muse and Gaut model (MG94) is simplest
 - α = synonymous substitution rate
 - β = nonsynonymous substitution rate
 - $\pi_A, \pi_C, \pi_G, \pi_T$ = base frequencies
- Goldman and Yang model (GY94) similar
 - accounts for synon./nonsynon. *and* trs/trv bias *and* amino acid properties (later simplified, see Yang et al. 1998)

Muse, S. V., and B. S. Gaut. 1994. A likelihood approach for comparing synonymous and nonsynonymous substitution rates, with application to the chloroplast genome. *Molecular Biology and Evolution* 11:715-724.

Goldman, N., and Z. Yang. 1994. A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Molecular Biology and Evolution* 11:725-736.

Yang, Z., Nielsen, R., and Hasegawa, M. 1998. Models of amino acid substitution and applications to mitochondrial protein evolution. *Molecular Biology and Evolution* 15:1600-1611.

Table I. Part of Muse and Gaut's 61×61 instantaneous rate matrix^a

Codon before substitution (the 'from' state)	Codon after substitution (the 'to' state)							
	TTT (Phe)	TTC (Phe)	TTA (Leu)	TTG (Leu)	CTT (Leu)	CTC (Leu)	...	GGG (Gly)
TTT (Phe)	---	$\alpha\pi_C$	$\beta\pi_A$	$\beta\pi_G$	$\beta\pi_C$	0	...	0
TTC (Phe)	$\alpha\pi_T$	---	$\beta\pi_A$	$\beta\pi_G$	0	$\beta\pi_C$...	0
TTA (Leu)	$\beta\pi_T$	$\beta\pi_C$	---	$\alpha\pi_G$	0	0	...	0
TTG (Leu)	$\beta\pi_T$	$\beta\pi_C$	$\alpha\pi_A$	---	0	0	...	0
CTT (Leu)	$\beta\pi_T$	0	0	0	---	$\alpha\pi_C$...	0
CTC (Leu)	0	$\beta\pi_T$	0	0	$\alpha\pi_T$	---	...	0
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
GGG (Gly)	0	0	0	0	0	0	...	---

Note that it is still easy for the change CTT → TTA to occur, it just requires more than one instant of time

Instantaneous rate is 0.0 if two or more nucleotides must change during the codon transition

Table 1 from: Lewis, P. O. 2001. Phylogenetic systematics turns over a new leaf. Trends in Ecology and Evolution 16:30-37.

Interpreting codon model results

$\omega = \beta/\alpha$ is the nonsynonymous/synonymous rate ratio

omega	mode of selection	example(s)
$\omega < 1$	stabilizing selection (nucleotide substitutions rarely change the amino acid)	functional protein coding genes
$\omega = 1$	neutral evolution (synonymous and nonsynonymous substitutions occur at the same rate)	pseudogenes
$\omega > 1$	positive selection (nucleotide substitutions often change the amino acid)	envelope proteins in viruses under active positive selection

Amino acid models

JC69 Flashback

	A	C	G	T
A	-3α	α	α	α
C	α	-3α	α	α
G	α	α	-3α	α
T	α	α	α	-3α

← Q matrix
(instantaneous rates)

P matrix
(transition probabilities)



$\frac{1}{4} + \frac{3}{4}e^{-4\alpha t}$	$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} (1 - e^{-4\alpha t})$
$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} + \frac{3}{4}e^{-4\alpha t}$	$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} (1 - e^{-4\alpha t})$
$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} + \frac{3}{4}e^{-4\alpha t}$	$\frac{1}{4} (1 - e^{-4\alpha t})$
$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} (1 - e^{-4\alpha t})$	$\frac{1}{4} + \frac{3}{4}e^{-4\alpha t}$

A different path from Q to P

For many interesting models, it is not possible to obtain transition probabilities *analytically* (i.e. using a formula)

We can, however, obtain transition probabilities *numerically* (i.e. obtain the value of the transition probability without plugging values into a formula)

$$\mathbf{P}(t) = e^{\mathbf{Q}t}$$

$\lambda_1, \lambda_2, \lambda_3,$ and λ_4 are the eigenvalues of \mathbf{Q}

$$\mathbf{P}(t) = \mathbf{U} \begin{pmatrix} e^{\lambda_1 t} & 0 & 0 & 0 \\ 0 & e^{\lambda_2 t} & 0 & 0 \\ 0 & 0 & e^{\lambda_3 t} & 0 \\ 0 & 0 & 0 & e^{\lambda_4 t} \end{pmatrix} \mathbf{U}^{-1}$$

Inverse
of \mathbf{U}

Matrix of eigenvectors of \mathbf{Q}

Factoring a square matrix into eigenvectors and a diagonal matrix of eigenvalues is known as *diagonalization*

The elements of Q

The Q matrix is often presented in the following form, factored into a symmetric matrix R of exchangeabilities and a set of state frequencies.

Ala								
Arg	0.267828							
Asn	0.984474	0.327059						
Asp	1.199805	0.000000	8.931515					
Cys	0.360016	0.232374	0.000000	0.000000				
Gln	0.887753	2.439939	1.028509	1.348551	0.000000			
Glu	1.961167	0.000000	1.493409	11.388659	0.000000	7.086022		
Gly	2.386111	0.087791	1.385352	1.240981	0.107278	0.281581	0.811907	
His	0.228116	2.383148	5.290024	0.868241	0.282729	6.011613	0.439469	...
Ile	0.653416	0.632629	0.768024	0.239248	0.438074	0.180393	0.609526	...
Leu	0.406431	0.154924	0.341113	0.000000	0.000000	0.730772	0.112880	...
Lys	0.258635	4.610124	3.148371	0.716913	0.000000	1.519078	0.830078	...
Met	0.717840	0.896321	0.000000	0.000000	0.000000	1.127499	0.304803	...
Phe	0.183641	0.136906	0.138503	0.000000	0.000000	0.000000	0.000000	...
Pro	2.485920	1.028313	0.419244	0.133940	0.187550	1.526188	0.507003	...
Ser	4.051870	1.531590	4.885892	0.956097	1.598356	0.561828	0.793999	...
Thr	3.680365	0.265745	2.271697	0.660930	0.162366	0.525651	0.340156	...
Trp	0.000000	2.001375	0.224968	0.000000	0.000000	0.000000	0.000000	...
Tyr	0.244139	0.078012	0.946940	0.000000	0.953164	0.000000	0.214717	...
Val	2.059564	0.240368	0.158067	0.178316	0.484678	0.346983	0.367250	...
	Ala	Arg	Asn	Asp	Cys	Gln	Glu	...

R matrix
(only values below
diagonal shown)

Freq	0.087127	0.040904	0.040432	0.046872	0.033474	0.038255	0.049530	...
------	----------	----------	----------	----------	----------	----------	----------	-----

Frequencies

GTR Flashback

	A	C	G	T
A	—	$\pi_C a\mu$	$\pi_G b\mu$	$\pi_T c\mu$
C	$\pi_A a\mu$	—	$\pi_G d\mu$	$\pi_T e\mu$
G	$\pi_A b\mu$	$\pi_C d\mu$	—	$\pi_T f\mu$
T	$\pi_A c\mu$	$\pi_C e\mu$	$\pi_G f\mu$	—

The off-diagonal elements of the GTR matrix can similarly be separated into a symmetric R matrix and a diagonal matrix of frequencies.

$$\begin{pmatrix}
 - & a\mu & b\mu & c\mu \\
 a\mu & - & d\mu & e\mu \\
 b\mu & d\mu & - & f\mu \\
 c\mu & e\mu & f\mu & -
 \end{pmatrix}
 \begin{pmatrix}
 \pi_A & 0 & 0 & 0 \\
 0 & \pi_C & 0 & 0 \\
 0 & 0 & \pi_G & 0 \\
 0 & 0 & 0 & \pi_T
 \end{pmatrix}$$

R matrix
Frequencies

What does all this accomplish?

- An empirical Q matrix can be constructed from many closely-related pairwise comparisons
- A Q matrix can be extrapolated using diagonalization to generate a P matrix for any desired value of t
- This model has 0 parameters!
- Models generic features of protein evolution; Q matrix does not necessarily reflect your particular sequences
- Frequencies can be swapped with more appropriate set (locally estimated)

Ways to improve

- Base everything on a much larger protein database (JTT model)
- Avoid need to use closely-related sequence pairs by obtaining ML estimate of Q matrix (WAG model)
- Add rate heterogeneity to ML estimation of Q matrix (LG model)

JTT: Jones, D.T., Taylor, W.R., and Thornton, J.M. 1992. The rapid generation of mutation data matrices from protein sequences. *Comput. Appl. Biosci.* 8:275–282.

WAG: Whelan, S., and Goldman, N. 2001. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. *Molecular Biology and Evolution.* 18:691–699.

LG: Le, S.Q., and Gascuel, O. 2008. An improved general amino acid replacement matrix. *Molecular Biology and Evolution.* 25:1307–1320.

The End