

# Phylogenetic inference and likelihood 2

Emily Jane McTavish

Life and Environmental Sciences  
University of California, Merced

`ejmctavish@ucmerced.edu`, `twitter:snacktavish`

(With thanks to Mark Holder and Paul Lewis for slides)

# Combining probabilities

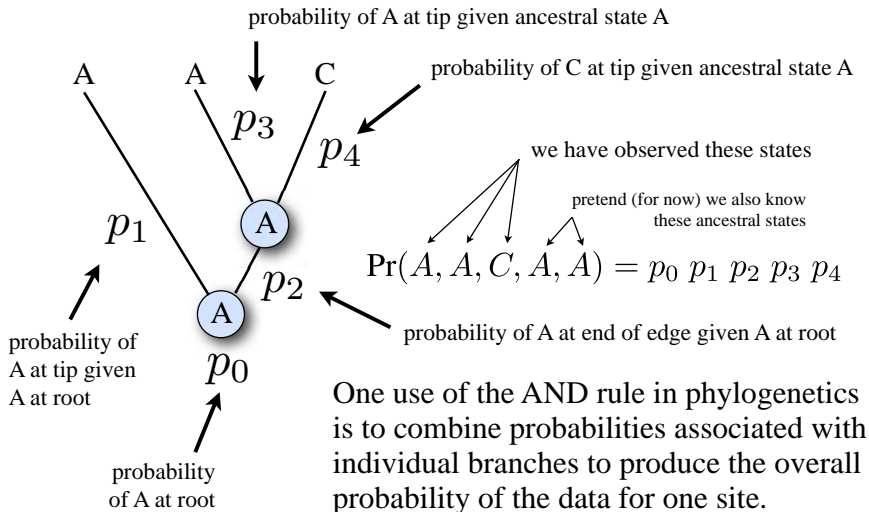
- *Multiply* probabilities if the component events must happen **simultaneously** (i.e. where you would naturally use the word AND when describing the problem)

Using 2 dice, what is the probability of



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

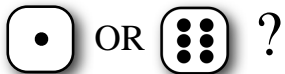
# AND rule in phylogenetics



# 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)

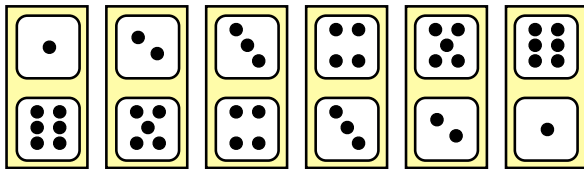
Using one die, what is the probability of



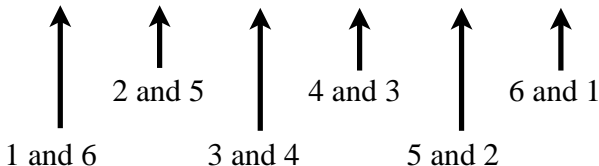
$$(1/6) + (1/6) = 1/3$$

# Combining AND and OR

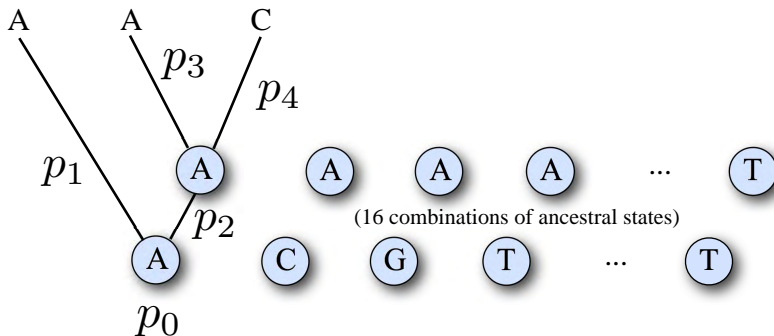
What is the probability that the sum of two dice is 7?



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



# 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.

# Independence

This is always true...

$$\underset{\text{joint probability}}{\Pr(A \text{ and } B)} = \Pr(A) \underset{\text{conditional probability}}{\Pr(B|A)}$$

If we can say this...

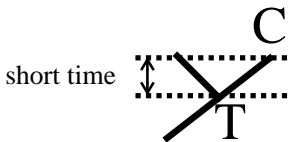
$$\Pr(B|A) = \Pr(B)$$

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

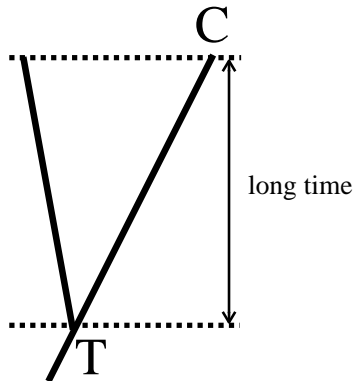
$$\Pr(A \text{ and } B) = \Pr(A) \Pr(B)$$

# 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:

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

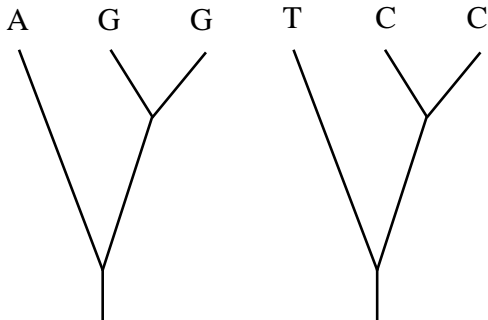
If we can say this...

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

...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

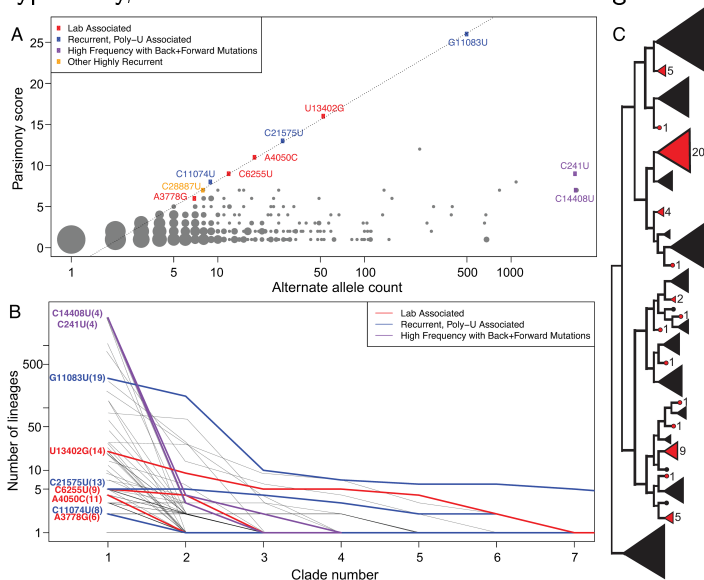


The site data patterns AGG and TCC are assumed by most models to be conditionally independent.

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

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

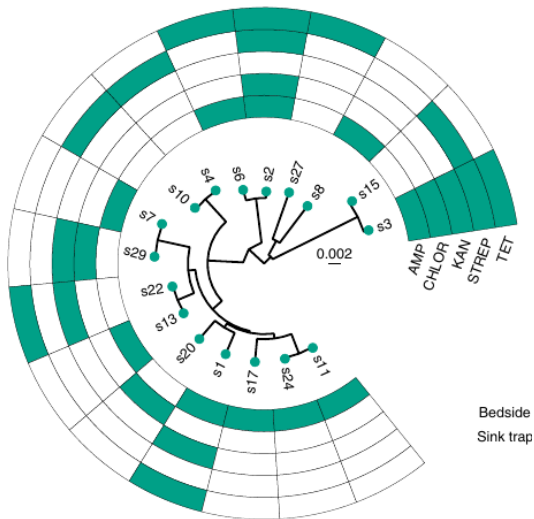
# Discussion: How does the likelihood for the data observed at different site types vary, conditioned on the tree and branch lengths?



Stability of SARS-CoV-2 phylogenies Turakhia Y, De Maio N, Thornlow B, Gozashti L, Lanfear R, et al. (2020)  
 Stability of SARS-CoV-2 phylogenies. PLOS Genetics 16(11): e1009175.  
<https://doi.org/10.1371/journal.pgen.1009175>

Discussion: How does the likelihood for the observed data antibiotic resistance vary, conditioned on the tree and branch lengths?

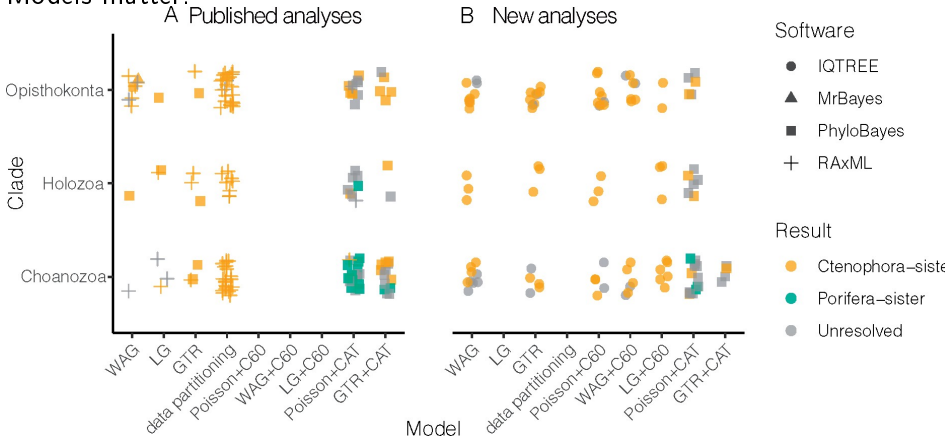
**d**



# Models

- Models help us intelligently **interpolate between our observations** for purposes of **making predictions**
- **Adding parameters** to a model generally increases its fit to the data
- **Underparameterized** models lead to poor fit to observed data points
- **Overparameterized** models lead to poor prediction of future observations
- Criteria for choosing models include likelihood ratio tests, AIC, BIC, Bayes Factors, etc.
  - all provide a way to choose a model that is neither underparameterized nor overparameterized

# Models matter!



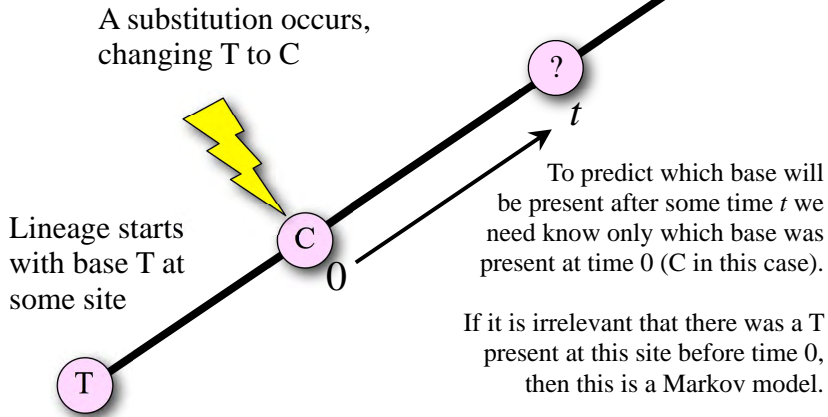
Yuanning Li, Xing-Xing Shen, Benjamin Evans, Casey W Dunn, Antonis Rokas, Rooting the Animal Tree of Life, Molecular Biology and Evolution, Volume 38, Issue 10, October 2021, Pages 4322–4333, <https://doi.org/10.1093/molbev/msab170>

# 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  
( $r_{A \leftrightarrow C} = r_{A \leftrightarrow G} = r_{A \leftrightarrow T} = r_{C \leftrightarrow G} = r_{C \leftrightarrow T} = r_{G \leftrightarrow T} = \alpha$ )
- 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...

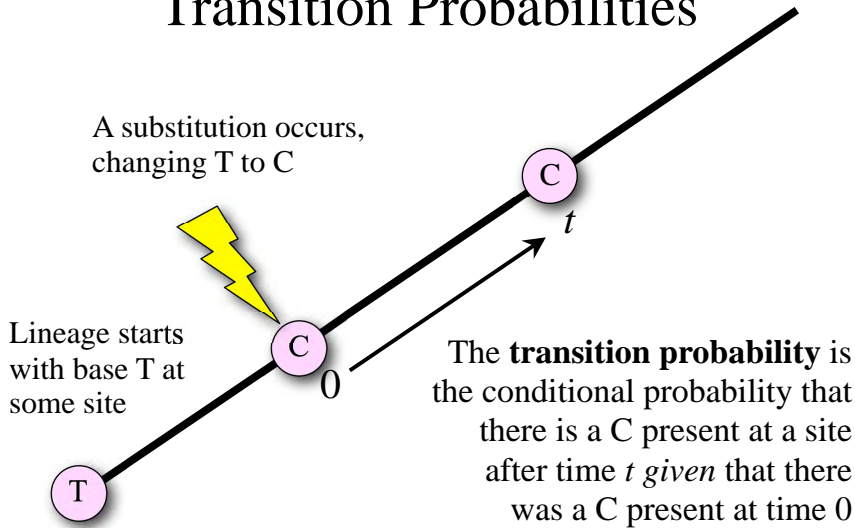
Jukes, T. H., and C. R. Cantor. 1969. Evolution of protein molecules. Pages 21-132 in H. N. Munro (ed.), *Mammalian Protein Metabolism*. Academic Press, New York.

# What is a Markov process?





# 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  $\alpha$ :

$$P_{TG}(t) = \frac{1}{4} (1 - e^{-4\alpha t}) = \Pr(G|T, \alpha t)$$

The JC69 model has only one unknown quantity:  $\alpha t$

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

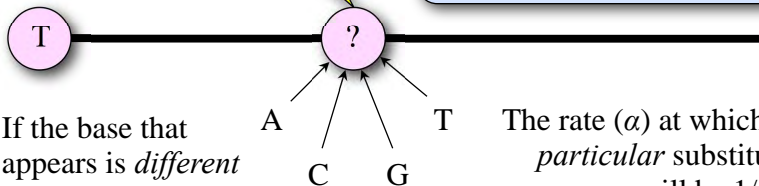
Where does a transition probability formula such as this come from?

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

# The Poisson distribution

---

Probability distribution on the number of events when:

1. events are assumed to be independent,
2. the *rate* of events some constant,  $\mu$ , and
3. the process continues for some duration of time,  $t$ .

The expectation of the number of events is  $\nu = \mu t$ .

Note that  $\nu$  can be any non-negative number, but the Poisson is a discrete distribution – it gives the probabilities of the number of events (and this number will always be a non-negative integer).

Poisson distribution can be used to explain statistical regularities of rare events

$$P(k \text{ events in interval}) = \frac{e^{-\nu} \nu^k}{k!}$$

- ▶  $\nu$  is the average number of events per interval (rate times time)
- ▶  $e$  is the number 2.71828... (Euler's number) the base of the natural logarithms
- ▶  $k$  takes values 0, 1, 2, ...
- ▶  $k! = k * (k - 1) * (k - 2) * \dots * 2 * 1$  is the factorial of  $k$ .

from wikipedia

$$P(\text{k events in interval}) = \frac{e^{-\nu} \nu^k}{k!}$$
$$P(0 \text{ events}) = \frac{e^{-\nu} \nu^0}{0!} = e^{-\nu} = e^{-\mu t}$$
$$P(\geq 1 \text{ events}) = 1 - e^{-\mu t}$$

# 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  $\alpha$ :

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

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

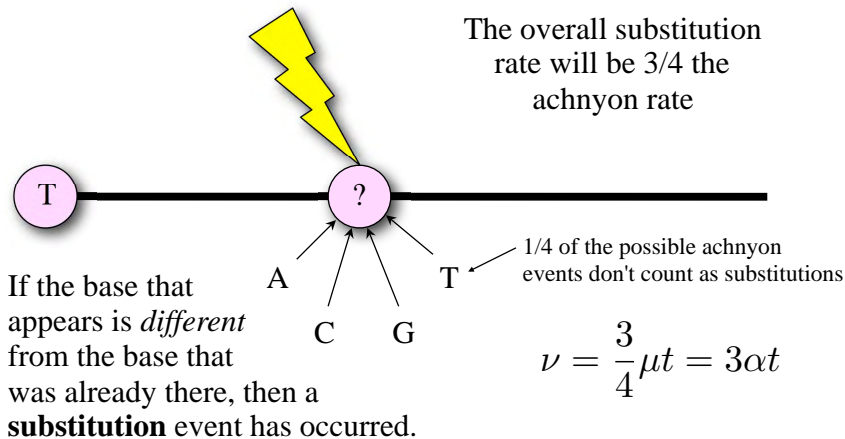
$$\Pr(\text{last achnyon results in base G}) = \frac{1}{4}$$

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

Remember that the rate ( $\alpha$ ) of any particular substitution is one fourth the achnyon rate ( $\mu$ ):

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

# Expected number of substitutions





# Transition Probabilities: Remarks

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

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

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

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

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

---

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

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  $t$

# 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?

# Models for nucleotide substitutions

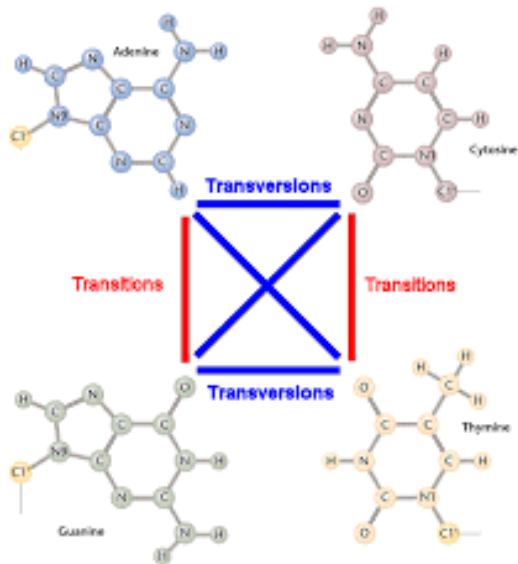
# JC69 rate matrix

1 parameter:  
 $\alpha$

		To			
		A	C	G	T
From	A	$-3\alpha$	$\alpha$	$\alpha$	$\alpha$
	C	$\alpha$	$-3\alpha$	$\alpha$	$\alpha$
	G	$\alpha$	$\alpha$	$-3\alpha$	$\alpha$
	T	$\alpha$	$\alpha$	$\alpha$	$-3\alpha$

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.

Bring in some biology!



Transitions are between two ring purines (A G) or between one ring pyrimidines (C T). they therefore involve bases of similar shape.

Transversions are interchanges of purine for pyrimidine bases.

# K80 (or K2P) rate matrix

2 parameters:

$\alpha$   
 $\beta$

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

↑ 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.

# Transition/transversion ratio (ratio) versus Transition/transversion *rate* ratio (kappa)



## 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:

$\mu$

$\pi_A$

$\pi_C$

$\pi_G$

	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

# 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:**

$\pi_A$   
 $\pi_C$   
 $\pi_G$   
 $a$   
 $b$   
 $c$   
 $d$   
 $e$   
 $f$   
 $\mu$

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)	AAAGATTACAGATTAACTTACTATACTCTGAGTATAAACTAAAGATACTGACATTTTAGCTGCAITTCGTGTAACCTCCA
Chlorella (green alga)	....C...C.T.....T...CC.C.A.....C.....C.T...A..G.C...A.G....T
Volvox (green alga)	.....TC.T.....A...C...A...C...GT.GTA.....C.....C.....A..A.G....
Conocephalum (liverwort)	.....TC.....T.....G.T...G.....G...T.....A..A.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)	.....C.....G.....T.....T..CG.C.....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

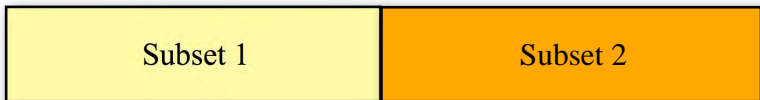
Q--L--G--V--P--P--E--E--A--G--A--A--V--A--A--E--S--S--T--G--T--W--T--T--V--W--T--D--G--L--T--S--L--D--R--Y--K--G--R--C--Y--H--I--E--	
CAACCTGSCGTTCCACCTGAAGAAGCAGGGGCTGCAGTAGCTGCAGAATCTTCTACTGGTACATGGACTACTGTTTGSACGTGACGGGATTAAGTAGTTTGGACCGATACAAAGGAAGATGCTACGATATTGAA	
....A..T.....A.....G..T...G.....A.....A.....T.....G.....T.....TC.T...T...C.C..G...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



$r_1$  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:

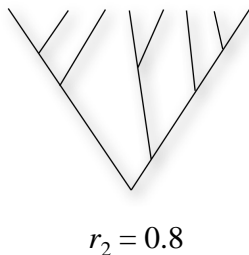
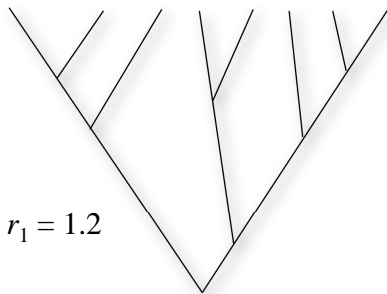
$$r_1 p(r_1) + r_2 p(r_2) = 1$$

# 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**:

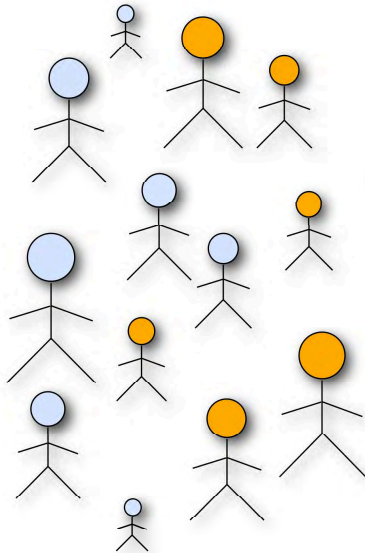
$$\begin{aligned}P_{ii}(t) &= \frac{1}{4} + \frac{3}{4}e^{-4r_1\alpha t} \\P_{ij}(t) &= \frac{1}{4} - \frac{1}{4}e^{-4r_1\alpha t}\end{aligned}$$

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

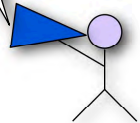
$$\begin{aligned}P_{ii}(t) &= \frac{1}{4} + \frac{3}{4}e^{-4r_2\alpha t} \\P_{ij}(t) &= \frac{1}{4} - \frac{1}{4}e^{-4r_2\alpha t}\end{aligned}$$



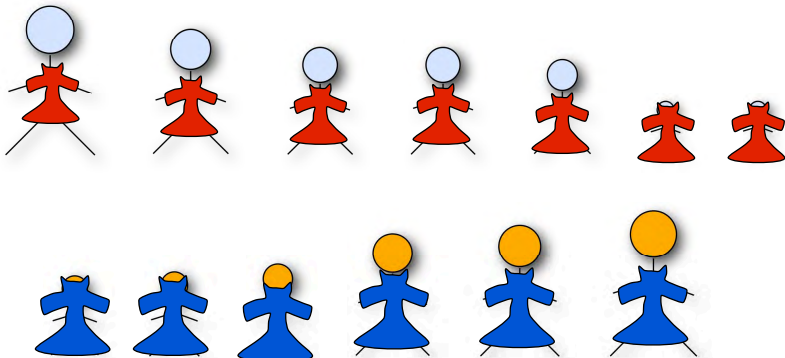
## Site-specific Approach



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

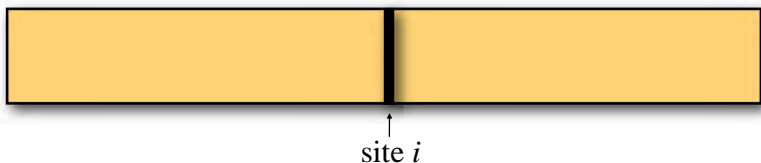


Good: costs less: need to buy just one coat for every person

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

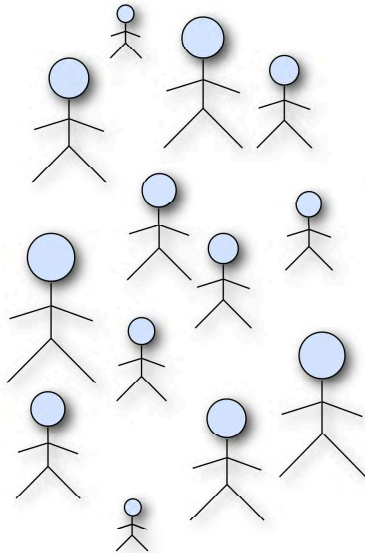
All relative rates applied to every site



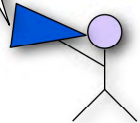
$$L_i = \Pr(D_i|r_1) \Pr(r_1) + \Pr(D_i|r_2) \Pr(r_2)$$

Common examples  $\left\{ \begin{array}{l} \text{Invariable sites (I) model} \\ \text{Discrete Gamma (G) model} \end{array} \right.$

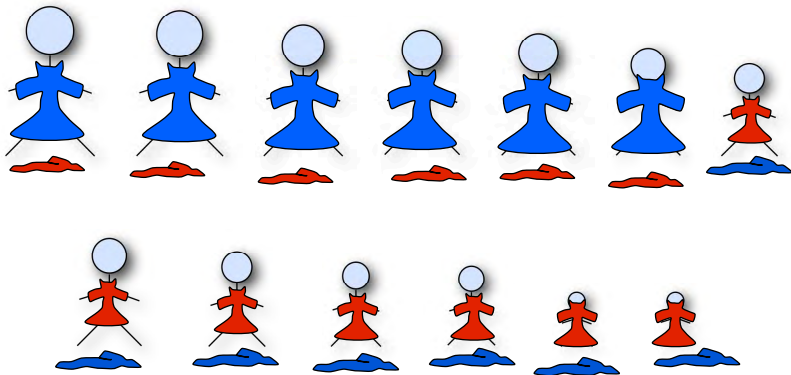
## Mixture Model Approach



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

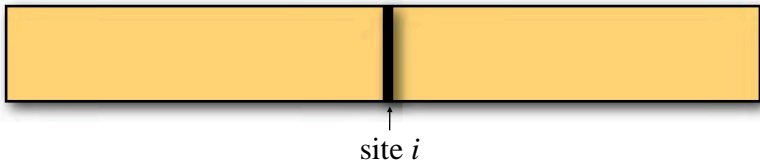


Good: every person experiences better fit because they can choose the size coat that fits best

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

# Invariable Sites Model

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



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

$$r_1 = 0.0$$

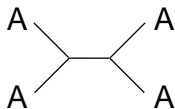
$$r_2 = \frac{1}{1 - p_{\text{invar}}}$$

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

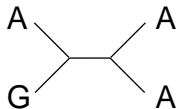
If site  $i$  is a *constant* site, both terms will contribute to the site likelihood:



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

---

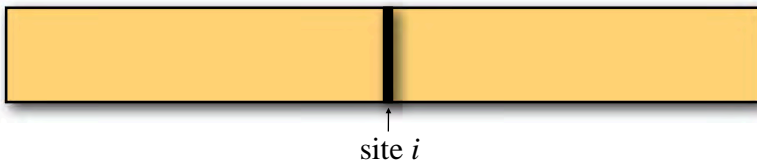
If site  $i$  is a *variable* site, there is no way to explain the data with a zero rate, so the first term is zero:



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

# Discrete Gamma Model

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



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

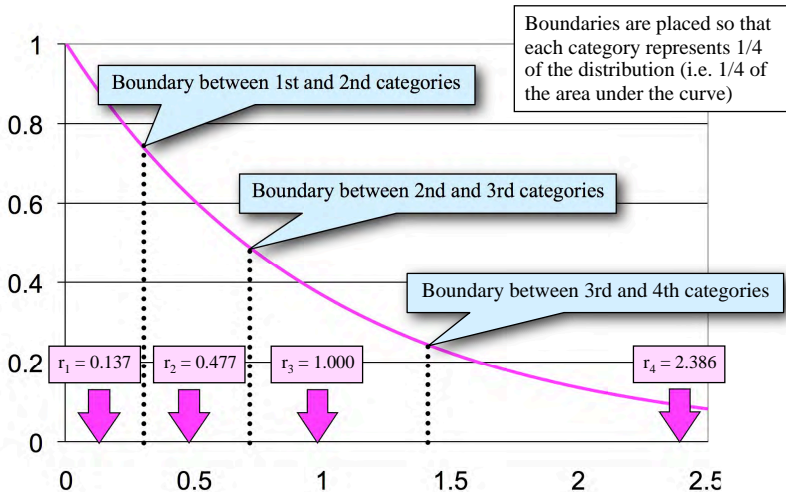
Relative rates are constrained to a discrete gamma distribution  
Number of rate categories can vary (4 used here)

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

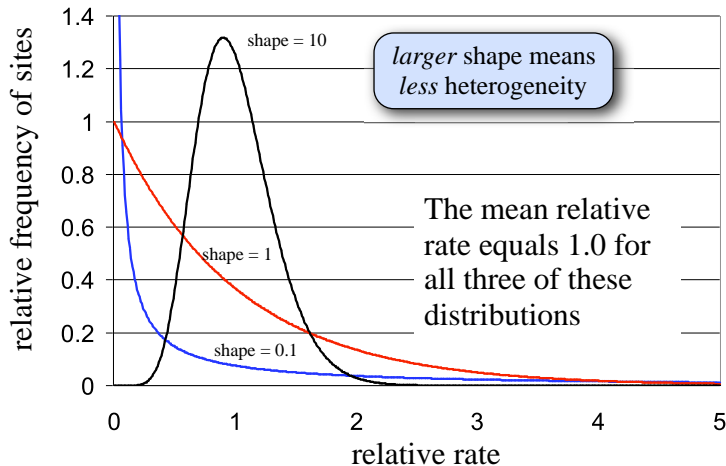
Yang, Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *Journal of Molecular Evolution* 39:306-314.



# Relative rates in 4-category case



# Gamma distributions



## Models for other data types

- ▶ Amino acids (as discussed in the Li and Dunn papers)
- ▶ Morphological data
- ▶ Expression data

How would a model for Copy Number Variant data differ from nucleotide models?

What are some components that would make biological sense in a model for CVNs?