

# Codon Models

# The Genetic Code

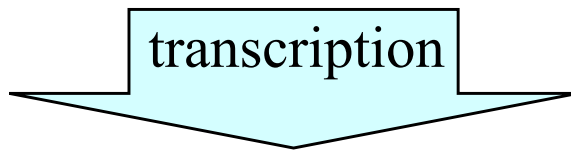
First 12 nucleotides at the 5' end of the *rbcL* gene in corn:

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

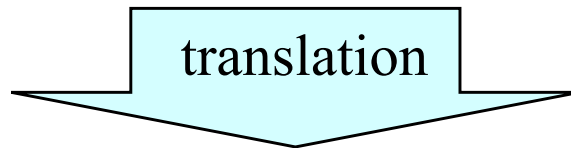
DNA double helix

## 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



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



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

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

Table 1 from: Lewis (2001b)

**Table I. Part of Muse and Gaut's 61 × 61 instantaneous rate matrix<sup>a</sup>**

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	...	---

Compare to HKY85 model:

# Interpreting codon model results

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

omega	mode of selection	example
$\omega < 1$		
$\omega = 1$		
$\omega > 1$		

# Amino Acid Models



# Amino Acid Models

AA models began with Margaret Dayhoff's 1978 PAM (Point Accepted Mutation) matrices.

*PAM: a point mutation* that results in replacement of an amino acid at a particular site in a protein sequence *and which has been accepted* by natural selection

PAM1 is a matrix recording the probability of each possible amino acid substitution for proteins that differ in 1% of their amino acid positions.

# Simulation to show why PAM1 approximates rate matrix

1 million sites simulated for two taxa under HKY:

	A	C	G	T	
A	246563	306	2366	605	kappa = 5 pi = .1 .2 .3 .4 edge length = 0.01
C	166	246110	429	3184	
G	801	317	248759	639	
T	157	1577	463	247558	

---	0.02779	0.21490	0.05495
0.01508	---	0.03896	0.28919
0.07275	0.02879	---	0.05804
0.01426	0.14323	0.04205	---

counts normalized so that sum of off-diagonal elements is 1.0

---	0.02857	0.21429	0.05714
0.01429	---	0.04286	0.28571
0.07143	0.02857	---	0.05714
0.01429	0.14286	0.04286	---

true rate matrix

PAM1 matrix comparable to the **instantaneous rate matrix** of nucleotide models:

$$\mathbf{Q}_{GTR} = \begin{pmatrix}
 \text{---} & \pi_C a & \pi_C b & \pi_T c \\
 \pi_A a & \text{---} & \pi_G d & \pi_T e \\
 \pi_A b & \pi_C d & \text{---} & \pi_T f \\
 \pi_A c & \pi_C e & \pi_G f & \text{---}
 \end{pmatrix}$$

exchangeabilities

relative frequencies

The diagram illustrates the GTR rate matrix with two types of groupings. 'exchangeabilities' are indicated by red circles around the off-diagonal elements:  $\pi_C a$ ,  $\pi_C b$ ,  $\pi_T c$ ,  $\pi_G d$ ,  $\pi_T e$ , and  $\pi_T f$ . 'relative frequencies' are indicated by blue circles around the diagonal elements:  $\pi_A c$ ,  $\pi_C e$ , and  $\pi_G f$ . Lines connect the top of the red circles to a central point labeled 'exchangeabilities', and the bottom of the blue circles to a central point labeled 'relative frequencies'.

PAM30, PAM60, etc., are comparable to **transition matrices**: they are obtained by extrapolating the rate matrix (PAM1) to larger amounts of time



# Amino acid model evolution

PAM ➤ Dayhoff et al. (1978)

- used closely-related amino-acid sequences

JTT ➤ Jones et al. (1992)

- based on much larger database

WAG ➤ Whelan and Goldman (2001)

- **Q** MLE allowed more divergent sequences

LG ➤ Le and Gascuel (2008)

- added rate heterogeneity to **Q** MLE
- based on 50,000 sequences

IQ-TREE implements 26 additional amino acid models!

# Protein mixture models



Which amino acid equilibrium relative frequencies would provide a better fit for the site indicated?

$$\pi = \left\{ \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20}, \frac{1}{20} \right\}$$

$$\pi = \left\{ 0, \frac{1}{2}, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, \frac{1}{4}, 0, \frac{1}{4}, 0 \right\}$$

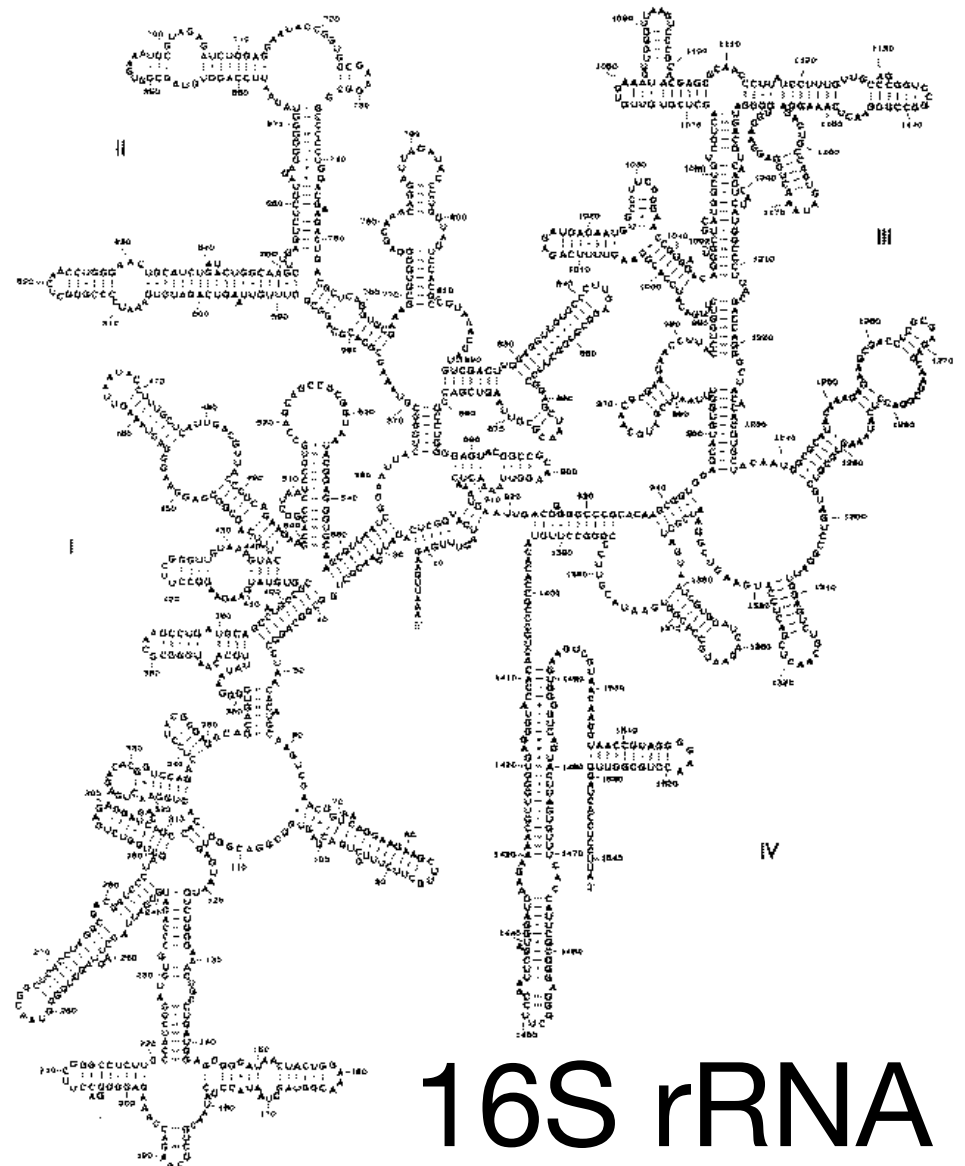
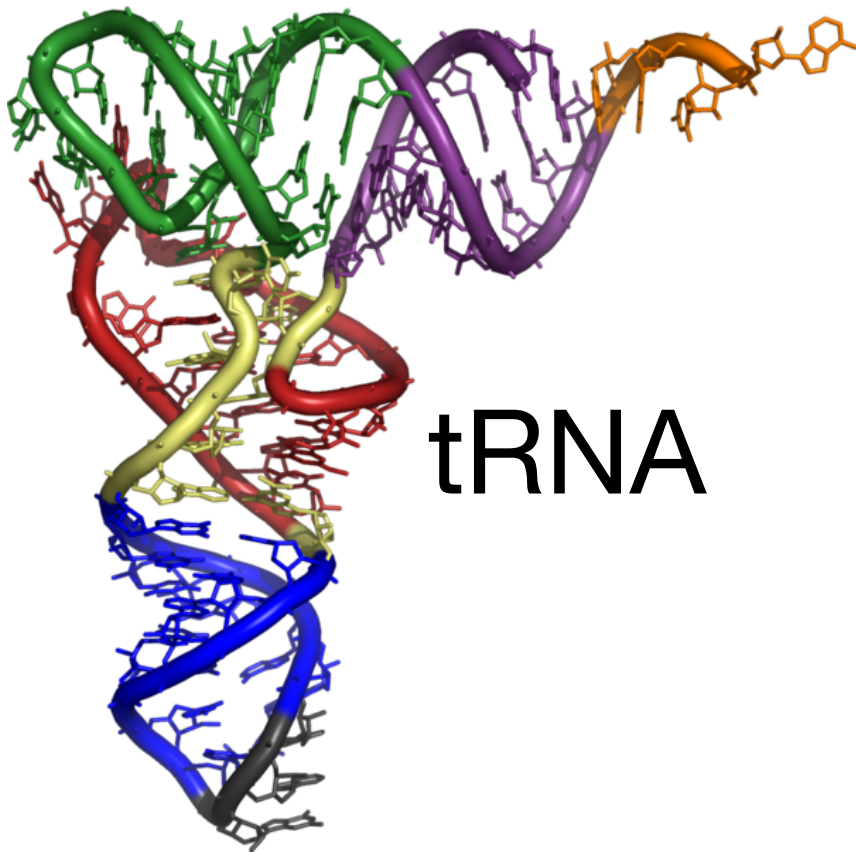
# C20 protein mixture model

$$L(D_k) = p(\boldsymbol{\pi}_1)p(D_k|\boldsymbol{\pi}_1) + p(\boldsymbol{\pi}_2)p(D_k|\boldsymbol{\pi}_2) + \dots + p(\boldsymbol{\pi}_{20})p(D_k|\boldsymbol{\pi}_{20})$$

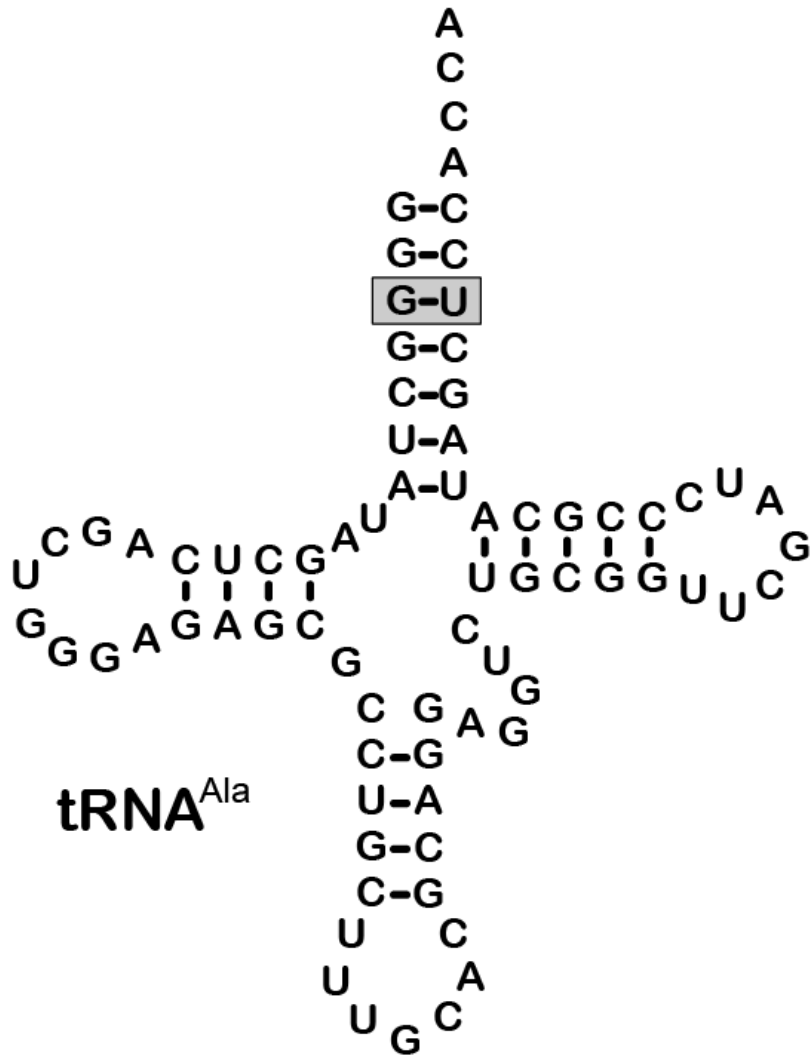
	<b>A</b>	<b>C</b>	<b>D</b>	<b>E</b>	<b>F</b>	<b>G</b>	<b>H</b>	<b>I</b>	<b>K</b>	<b>L</b>	<b>M</b>	<b>N</b>	<b>P</b>	<b>Q</b>	<b>R</b>	<b>S</b>	<b>T</b>	<b>V</b>	<b>W</b>	<b>Y</b>
1	3	1	0	0	1	0	0	38	0	9	2	0	0	0	0	0	2	42	0	0
2	17	3	0	1	1	2	0	12	1	9	3	1	2	1	1	5	12	27	0	1
3	4	1	1	4	0	2	3	1	28	3	1	3	1	9	28	4	4	2	0	1
4	1	0	0	0	3	0	0	31	0	36	5	0	0	0	0	0	1	20	0	0
5	7	0	27	34	0	3	1	0	4	1	0	4	1	6	2	4	3	1	0	0
6	9	1	3	2	1	4	1	1	3	1	1	8	2	2	2	30	26	2	0	1
7	3	0	29	8	0	8	3	0	4	0	0	25	1	3	2	8	4	0	0	1
8	9	0	8	20	0	2	2	1	14	2	1	4	1	15	9	4	5	2	0	0
9	44	4	0	1	1	11	0	1	0	1	1	1	2	1	1	21	7	5	0	0
10	1	1	0	0	39	0	3	1	0	5	1	0	0	0	0	1	0	1	5	41
11	2	1	0	0	8	0	0	12	0	49	16	0	0	1	0	1	1	6	1	1
12	5	1	1	4	2	1	3	9	8	18	6	2	1	7	9	3	6	9	1	2
13	16	0	2	4	1	3	1	1	3	3	1	2	38	3	3	9	5	3	0	1
14	5	2	2	3	10	2	16	2	3	6	2	5	1	5	5	5	3	3	3	18
15	10	1	5	11	1	2	2	6	7	5	2	3	2	7	6	7	12	11	0	1
16	4	2	0	0	20	1	2	12	1	22	4	1	1	1	1	1	2	12	3	10
17	11	1	6	3	0	47	1	0	3	1	0	8	1	2	3	8	2	0	0	0
18	18	0	8	13	0	6	2	1	8	1	1	5	2	8	6	13	7	2	0	0
19	5	1	6	5	1	4	12	1	9	2	1	16	1	10	9	8	5	1	0	2
20	12	0	13	17	0	3	1	0	7	1	0	4	20	5	3	7	4	1	0	0

# Stems and loops

[http://en.wikipedia.org/wiki/Image:3d\\_tRNA.png](http://en.wikipedia.org/wiki/Image:3d_tRNA.png)



# Compensatory substitutions



# Muse (1995) stem model

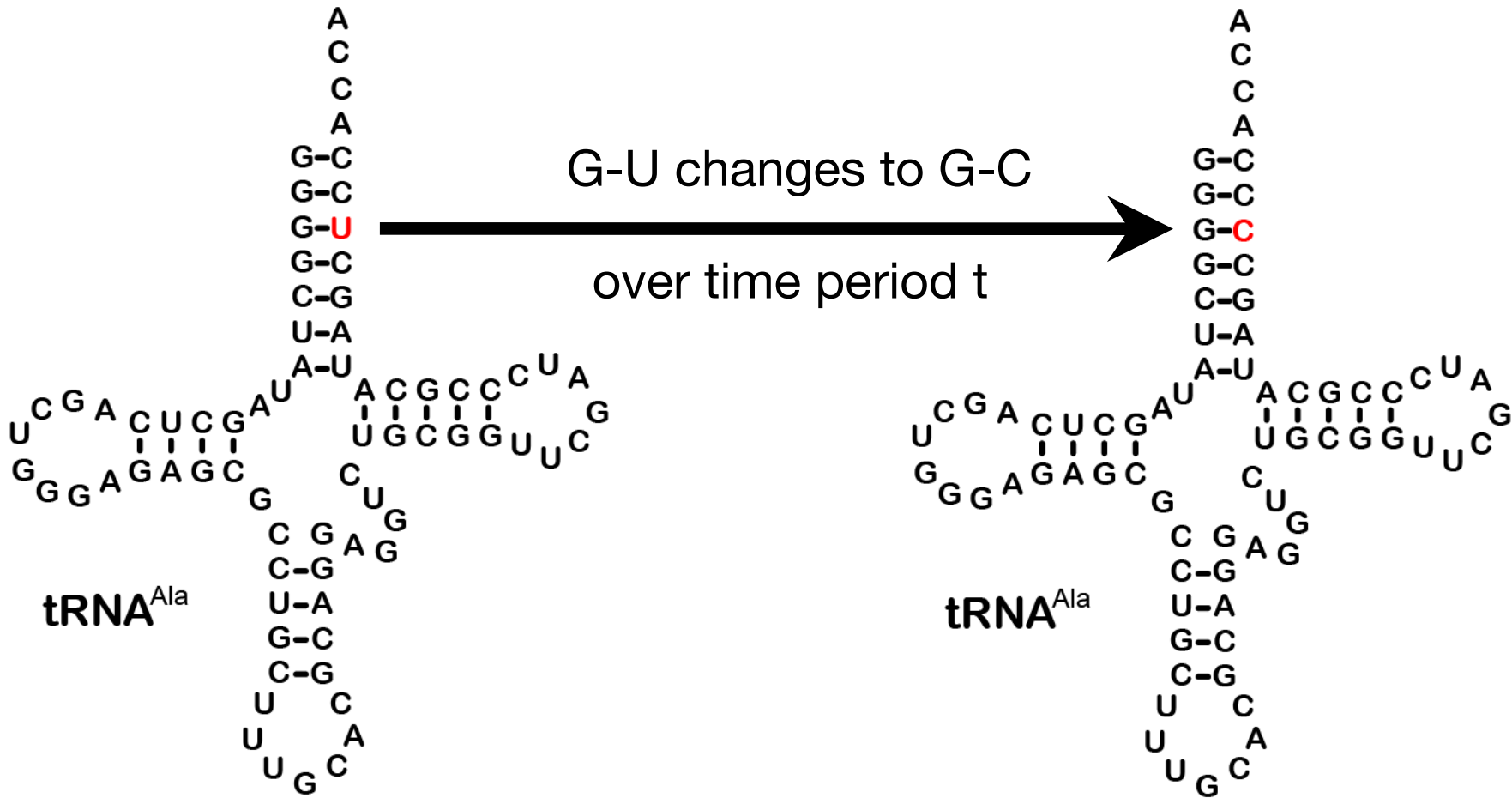
- Apply normal substitution model (e.g. JC69, K80, HKY85, etc.) to loop regions
- Apply a special "stem model" to the stem regions
- The stem model has 16 states:  
AA, AC, AG, AT, CA, CC, ..., TT  
where each state consists of one possible pairing of two nucleotides across a stem

# Muse (1995) stem model

- Idea is to compare null (independence) model to the alternative (dependence) model using a likelihood ratio test
- **Independence** model assumes each site evolves independently
- **Dependence** model allows rate of evolution to be *higher* for changes that *improve stem stability* and lower for changes that destroy stability

# Independence model

(applied to one stem site)



$$L_{\text{stem "site"}} = \left[ \begin{pmatrix} 1 \\ 4 \end{pmatrix} \left( \frac{1}{4} + \frac{3}{4} e^{-4\beta t} \right) \right] \left[ \begin{pmatrix} 1 \\ 4 \end{pmatrix} \left( \frac{1}{4} - \frac{1}{4} e^{-4\beta t} \right) \right]$$

$G \rightarrow G$  (no change)
 $U \rightarrow C$  (change)

# Independence model rate matrix\*

	AA	AC	AG	AU	CA	CC	CG	CU	GA	...
AA	-	$\beta$	$\beta$	$\beta$	$\beta$	0	0	0	$\beta$	...
AC	$\beta$	-	$\beta$	$\beta$	0	$\beta$	0	0	0	...
AG	$\beta$	$\beta$	-	$\beta$	0	0	$\beta$	0	0	...
AU	$\beta$	$\beta$	$\beta$	-	0	0	0	$\beta$	0	...
CA	$\beta$	0	0	0	-	$\beta$	$\beta$	$\beta$	0	...
CC	0	$\beta$	0	0	$\beta$	-	$\beta$	$\beta$	0	...
CG	0	0	$\beta$	0	$\beta$	$\beta$	-	$\beta$	0	...
CU	0	0	0	$\beta$	$\beta$	$\beta$	$\beta$	-	0	...
GA	$\beta$	0	0	0	0	0	0	0	-	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

# Dependence model rate matrix

	AA	AC	AG	AU	CA	CC	CG	CU	GA	...
AA	-	$\beta$	$\beta$	$\lambda\beta$	$\beta$	0	0	0	$\beta$	...
AC	$\beta$	-	$\beta$	$\lambda\beta$	0	$\beta$	0	0	0	...
AG	$\beta$	$\beta$	-	$\lambda\beta$	0	0	$\lambda\beta$	0	0	...
AU	$\beta/\lambda$	$\beta/\lambda$	$\beta/\lambda$	-	0	0	0	$\beta/\lambda$	0	...
CA	$\beta$	0	0	0	-	$\beta$	$\lambda\beta$	$\beta$	$\beta$	...
CC	0	$\beta$	0	0	$\beta$	-	$\lambda\beta$	$\beta$	0	...
CG	0	0	$\beta/\lambda$	0	0	$\beta/\lambda$	-	$\beta/\lambda$	$\beta/\lambda$	...
CU	0	0	0	$\lambda\beta$	0	$\beta$	$\lambda\beta$	-	$\beta$	...
GA	$\beta$	0	0	0	0	0	0	0	-	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

# Muse (1995) dependence model

- The instantaneous rate of those changes that **increase pairing** are multiplied by parameter  $\lambda$
- The instantaneous rate of those changes that **decrease pairing** are divided by parameter  $\lambda$
- The parameter  $\lambda$  thus measures the **degree to which pairing is favored in evolution**
- Setting  **$\lambda = 1$  equals independence**
- LRT statistic =  $-2 (\log L_{indep} - \log L_{dep})$   
1 d.f. ( $\lambda$  is the only extra parameter)
- Muse (1995) found that  $\lambda = 5.29$  for ribonuclease P  
(LRT = 224.30, P=0.0)