

# The Genetic Code

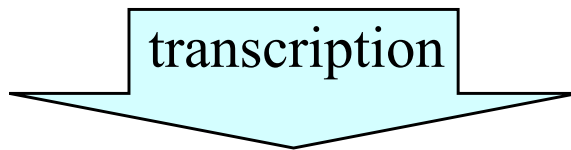
First 12 nucleotides at the 5' end of the *rbcL* gene in corn:

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

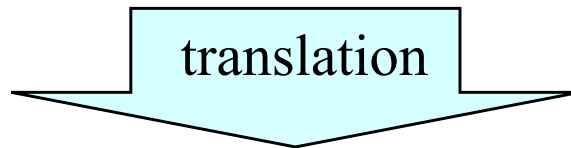
DNA double helix

## 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

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

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 shows the GTR rate matrix with exchangeabilities (red circles) and relative frequencies (blue circles) highlighted. Exchangeabilities are 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 the diagonal elements:  $\pi_A a$ ,  $\pi_A b$ ,  $\pi_A c$ ,  $\pi_C d$ ,  $\pi_C e$ ,  $\pi_G f$ , and  $\pi_G g$ .

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

# Amino acid models

WAG (Whelan And Goldman) exchangeabilities:

0.55157																			
0.50985	0.63535																		
0.73900	0.14730	5.42942																	
1.02704	0.52819	0.26526	0.03029																
0.90860	3.03550	1.54364	0.61678	0.09882															
1.58285	0.43916	0.94720	6.17416	0.02135	5.46947														
1.41672	0.58467	1.12556	0.86558	0.30667	0.33005	0.56772													
0.31695	2.13715	3.95629	0.93068	0.24897	4.29411	0.57003	0.24941												
0.19334	0.18698	0.55424	0.03944	0.17014	0.11392	0.12740	0.03045	0.13819											
0.39792	0.49767	0.13153	0.08480	0.38429	0.86949	0.15426	0.06130	0.49946	3.17097										
0.90627	5.35142	3.01201	0.47986	0.07403	3.89490	2.58443	0.37356	0.89043	0.32383	0.25756									
0.89350	0.68316	0.19822	0.10375	0.39048	1.54526	0.31512	0.17410	0.40414	4.25746	4.85402	0.93428								
0.21049	0.10271	0.09616	0.04673	0.39802	0.09992	0.08113	0.04993	0.67937	1.05947	2.11517	0.08884	1.19063							
1.43855	0.67949	0.19508	0.42398	0.10940	0.93337	0.68236	0.24357	0.69620	0.09993	0.41584	0.55690	0.17133	0.16144						
3.37079	1.22419	3.97423	1.07176	1.40766	1.02887	0.70494	1.34182	0.74017	0.31944	0.34474	0.96713	0.49391	0.54593	1.61328					
2.12111	0.55441	2.03006	0.37487	0.51298	0.85793	0.82277	0.22583	0.47331	1.45816	0.32662	1.38698	1.51612	0.17190	0.79538	4.37802				
0.11313	1.16392	0.07192	0.12977	0.71707	0.21574	0.15656	0.33698	0.26257	0.21248	0.66531	0.13751	0.51571	1.52964	0.13941	0.52374	0.11086			
0.24074	0.38153	1.08600	0.32571	0.54383	0.22771	0.19630	0.10360	3.87344	0.42017	0.39862	0.13326	0.42844	6.45428	0.21605	0.78699	0.29115	2.48539		
2.00601	0.25185	0.19625	0.15234	1.00214	0.30128	0.58873	0.18725	0.11836	7.82130	1.80034	0.30543	2.05845	0.64989	0.31489	0.23274	1.38823	0.36537	0.31473	

Determined from comparing 3905 closely-related protein sequences.

<http://www.ebi.ac.uk/goldman/WAG/wag.dat>

## Relative amino acid frequencies

0.08663	0.04397	0.03909	0.05705	0.01931	0.03673	0.05806	0.08325	0.02443	0.04847	0.08621	0.06203	0.01950	0.03843	0.04576	0.06952	0.06101	0.01439	0.03527	0.07090
Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val
A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V

# 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