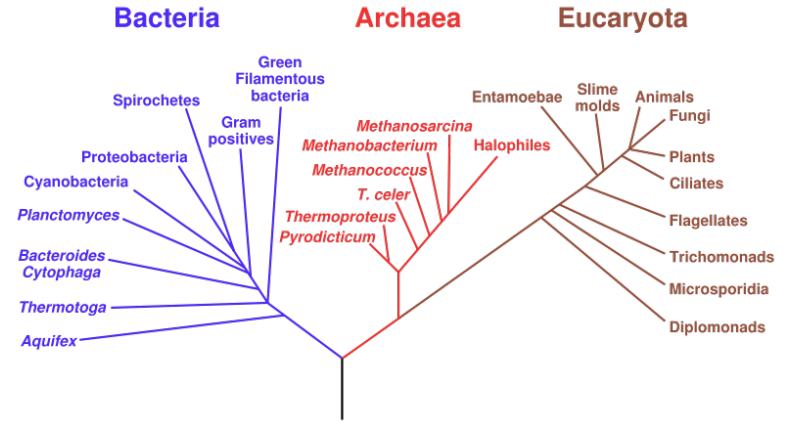




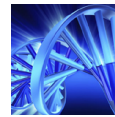
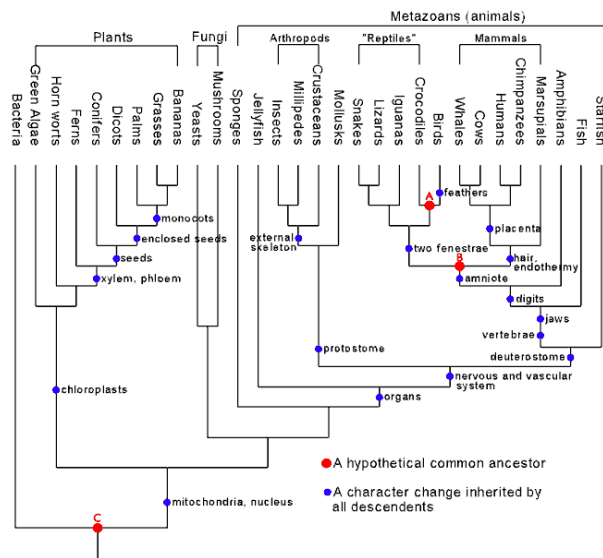
Mutations



Tree Of Life



Tree Of Life

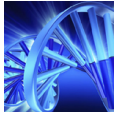


Mutations vs. Substitutions

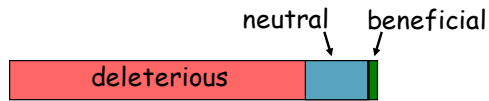
- *Mutations* are changes in DNA
- *Substitutions* are mutations that evolution has tolerated

Which rate is greater?

Replicative proofreading and DNA repair constrain the mutation rate

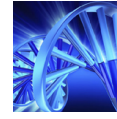


Selectionist Evolution

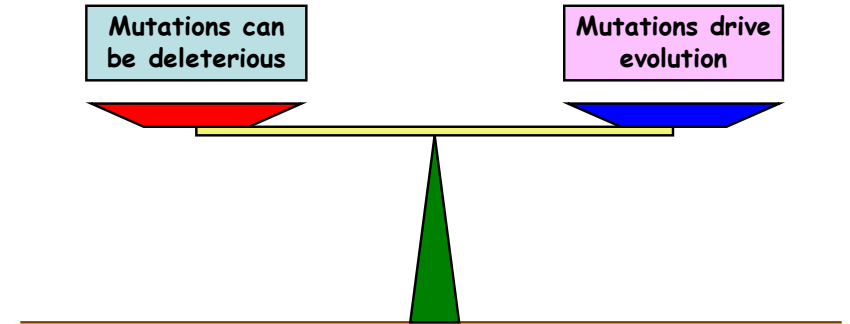


- Most mutations are deleterious; removed via negative selection
- Advantageous mutations positively selected
- Variability arises via selection

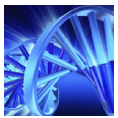
B - 5



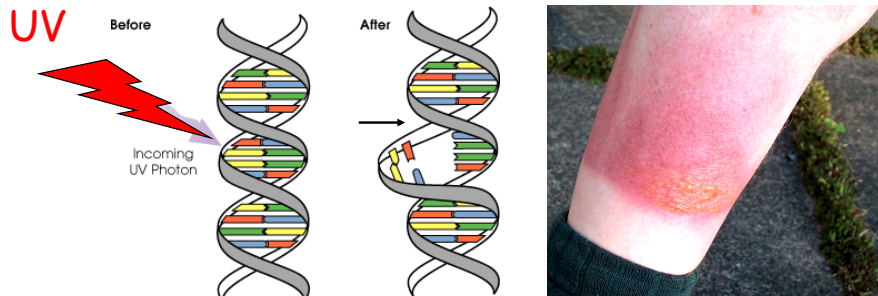
Why Are Mutations Important?



B - 6



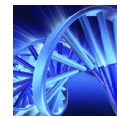
UV Damage to DNA



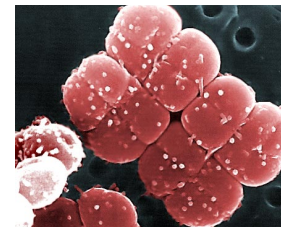
Thymine dimers

What happens if damage is not repaired?

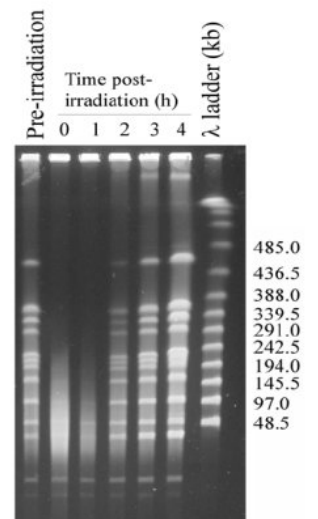
B - 7



Radiation Resistant

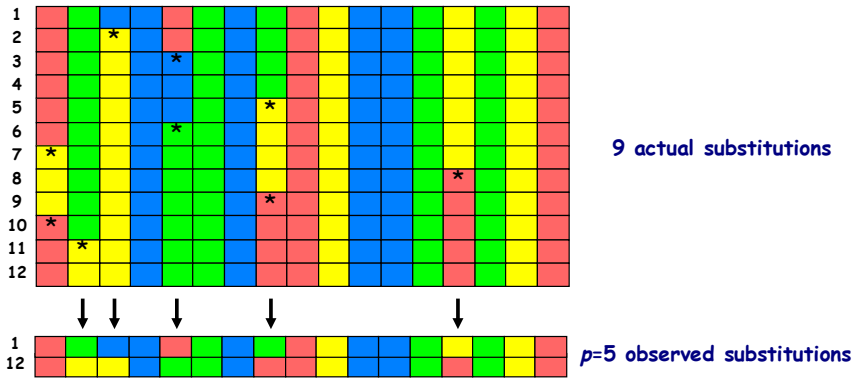


- 10 Gray will kill a human
- 60 Gray will kill an *E. coli* culture
- *Deinococcus* can survive 5000 Gray





A Sequence Mutating at Random

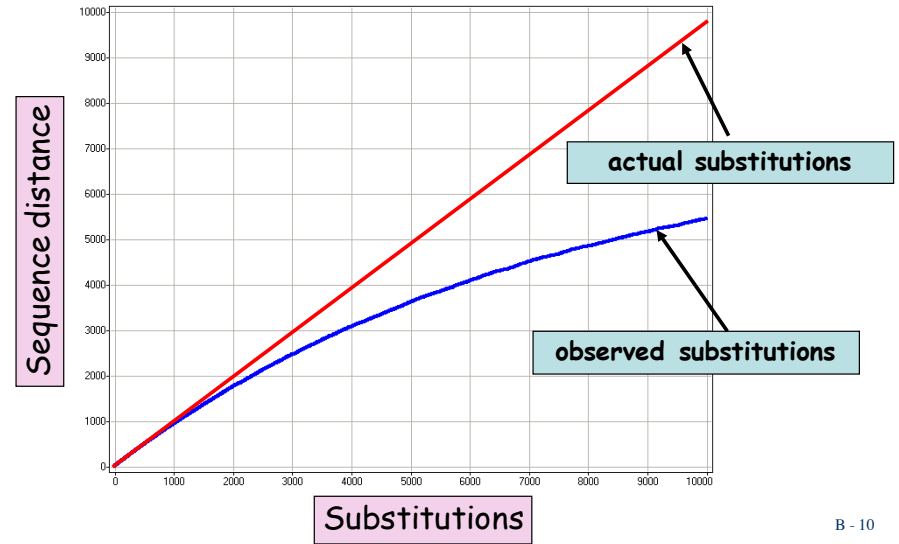


Multiple substitutions at one site can cause underestimation of number of actual substitutions

B - 9



Simulating Random Mutations



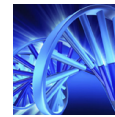
B - 10



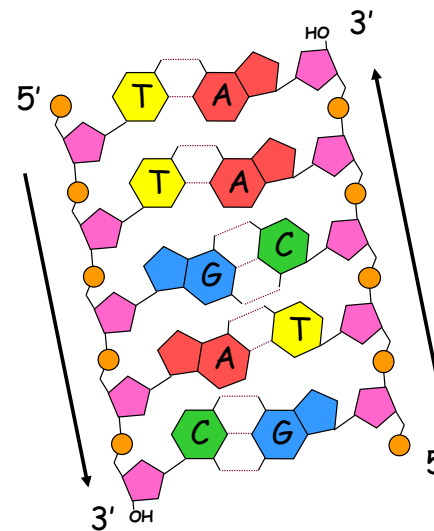
Measuring Sequence Divergence: Why Do We Care?

- Inferring phylogenetic relationships
- Dating divergence, correlating with fossil record
- Use in sequence alignments and homology searches of databases*

* Comparative genomics is an important field. Determining not only how many substitutions exist between two sequences but how similar two sequences are.^{B-11}



DNA Structure



G-C: 3 hydrogen bonds
A-T: 2 hydrogen bonds

Two base types:
- Purines (A, G)
- Pyrimidines (T, C)

B - 12



Not All Base Substitutions Are Created Equal

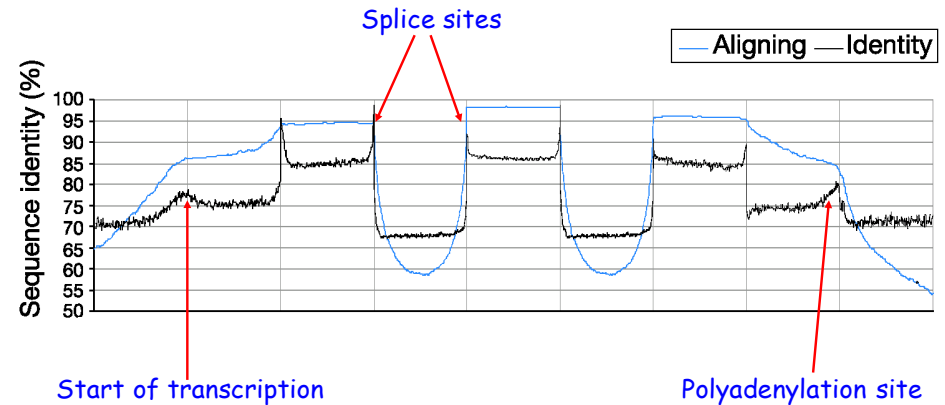
- **Transitions**
 - Purine to purine ($A \rightarrow G$ or $G \rightarrow A$)
 - Pyrimidine to pyrimidine ($C \rightarrow T$ or $T \rightarrow C$)
- **Transversions**
 - Purine to pyrimidine ($A \rightarrow C$ or T ; $G \rightarrow C$ or T)
 - Pyrimidine to purine ($C \rightarrow A$ or G ; $T \rightarrow A$ or G)

Transition rate $\sim 2x$ transversion rate

B - 13



Substitution Rates Differ Across Genomes



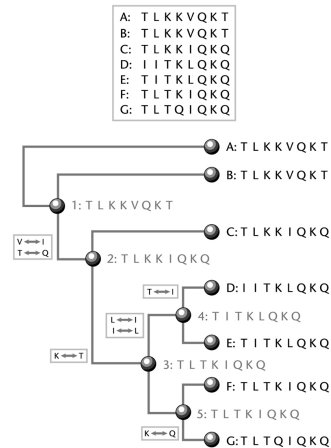
Alignment of 3,165 human-mouse pairs

B - 14

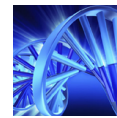


The PAM Model of Protein Sequence Evolution

- Empirical data-based substitution matrix
- Global alignments of 71 families of closely related proteins.
- Constructed hypothetical evolutionary trees
- Built matrix of 1572 amino acid point accepted mutations



B - 15



Original PAM Substitution Matrix

		Original amino acid																				
		A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
Replacement amino acid	A																					
	R	30																				
	N	109	17																			
	D	154	0	532																		
	C	33	10	0	0																	
	Q	93	120	50	76	0																
	E	266	0	94	831	0	422															
	G	579	10	156	162	10	30	112														
	H	21	103	226	43	10	243	23	10													
	I	66	30	36	13	17	8	35	0	3												
	L	95	17	37	0	0	75	15	17	40	253											
	K	57	477	322	85	0	147	104	60	23	43	39										
	M	29	17	0	0	0	20	7	7	0	57	207	90									
	F	20	7	7	0	0	0	17	20	90	167	0	17									
	P	345	67	27	10	10	93	40	49	50	7	43	43	4	7							
	S	772	137	432	98	117	47	86	450	26	20	32	168	20	40	269						
	T	590	20	169	57	10	37	31	50	14	129	52	200	28	10	73	696					
	W	0	27	3	0	0	0	0	0	3	0	13	0	0	10	0	17	0				
	Y	20	3	36	0	0	0	0	0	0	13	23	10	0	260	0	22	23	6			
	V	365	20	33	17	33	27	37	97	30	661	303	17	77	10	50	43	186	0	17		
A	Ala	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V		
R	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val			

Dayhoff, 1978

Count number of times residue i was replaced with residue j

B - 16



Deriving PAM Matrices

For each amino acid, calculate its *relative mutability*, i.e., the likelihood that the amino acid will mutate:

$$m_j = \frac{\text{\# times amino acid } j \text{ mutated}}{\text{total occurrences of amino acid } j}$$

TABLE 3-1 Relative Mutabilities of Amino Acids

The value of alanine is arbitrarily set to 100.

Asn	134	His	66
Ser	120	Arg	65
Asp	106	Lys	56
Glu	102	Pro	56
Ala	100	Gly	49
Thr	97	Tyr	41
Ile	96	Phe	41
Met	94	Leu	40
Gln	93	Cys	20
Val	74	Trp	18

Source: From Dayhoff (1978). Used with permission.

B - 17



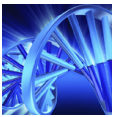
Deriving PAM Matrices

Calculate mutation probabilities for each possible substitution

$M_{i,j}$ = relative mutability \times proportion of all substitutions to j by changing to i

$$M_{i,j} = \frac{m_j \times A_{i,j}}{\sum_i A_{i,j}}$$

B - 18

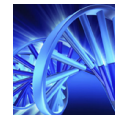


PAM1 Mutation Probability Matrix

		Original amino acid																			
		A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly	H His	I Ile	L Leu	K Lys	M Met	F Phe	P Pro	S Ser	T Thr	W Trp	Y Tyr	V Val
A	9867	2	9	10	3	8	17	21	2	6	4	2	6	2	22	35	32	0	2	18	
R	1	9913	1	0	1	10	0	0	10	3	1	19	4	1	4	6	1	8	0	1	
N	4	1	9822	36	0	4	6	6	21	3	1	13	0	1	2	20	9	1	4	1	
D	6	0	42	9859	0	6	53	6	4	1	0	3	0	0	1	5	3	0	0	1	
C	1	1	0	0	9973	0	0	0	1	1	0	0	0	0	1	5	1	0	3	2	
Q	3	9	4	5	0	9876	27	1	23	1	3	6	4	0	6	2	2	0	0	1	
E	10	0	7	56	0	35	9865	4	2	3	1	4	1	0	3	4	2	0	1	2	
G	21	1	12	11	1	3	7	9935	1	0	1	2	1	1	3	21	3	0	0	5	
H	1	8	18	3	1	20	1	0	9912	0	1	1	0	2	3	1	1	1	4	1	
I	2	2	3	1	2	1	2	0	0	9872	9	2	21	7	0	1	7	0	1	33	
L	3	1	3	0	0	6	1	1	4	22	9947	2	45	13	3	1	3	4	2	15	
K	2	37	25	6	0	12	7	2	2	4	1	9926	20	0	3	8	11	0	1	1	
M	1	1	0	0	0	2	0	0	0	5	8	4	9874	1	0	1	2	0	0	4	
F	1	1	1	0	0	0	0	1	2	8	6	0	4	9946	0	2	1	3	28	0	
P	13	5	2	1	1	8	3	2	5	1	2	2	1	1	9926	12	4	0	0	2	
S	28	11	34	7	11	4	6	16	2	2	1	7	4	3	17	9840	38	5	2	2	
T	22	2	13	4	1	3	2	2	1	11	2	8	6	1	5	32	9871	0	2	9	
W	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	9976	1	0	
Y	1	0	3	0	3	0	1	0	4	1	1	0	0	21	0	1	1	2	9945	1	
V	13	2	1	1	3	2	2	3	3	57	11	1	17	1	3	2	10	0	2	9901	

Dayhoff, 1978

B - 19



Deriving PAM Matrices

Calculate log odds ratio to convert mutation probability to substitution score

$$S_{i,j} = 10 \times \log_{10} \left(\frac{(M_{i,j})}{f_i} \right)$$

Mutation probability
(Prob. substitution from j to i is an accepted mutation)

Frequency of residue i
(Probability of amino acid i occurring by chance)

B - 20



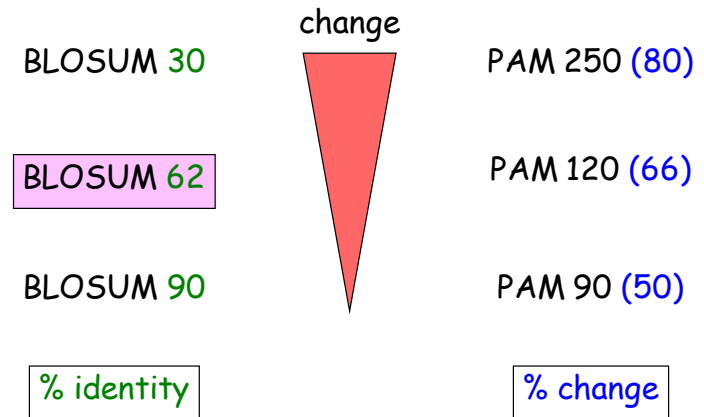
BLOSUM Uses Clustering To Reduce Sequence Bias

- Cluster the most similar sequences together
- Reduce weight of contribution of clustered sequences
- BLOSUM number refers to clustering threshold used (e.g. 62% for BLOSUM 62 matrix)

B - 25

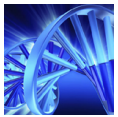


BLOSUM and PAM Substitution Matrices



BLAST algorithm uses BLOSUM 62 matrix

B - 26



PAM and BLOSUM

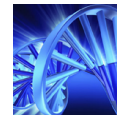
PAM

- Smaller set of closely related proteins - short evolutionary period
- Use global alignment
- More divergent matrices extrapolated
- Errors arise from extrapolation

BLOSUM

- Larger set of more divergent proteins-longer evolutionary period
- Use local alignment
- Each matrix calculated separately
- Clustering to avoid bias
- Errors arise from alignment errors

B - 27



Importance of Scoring Matrices

- Scoring matrices appear in all analyses involving sequence comparison
- The choice of matrix can strongly influence the outcome of the analysis

B - 28