Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

- PAM matrices, Dayhoff et al, 1978
- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

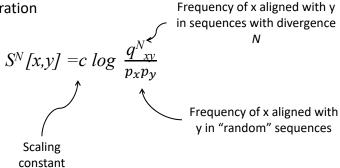
Overall strategy for both PAM and BLOSUM

- 1. Trusted amino acid alignments
- 2. Obtain amino acid pair counts (A_{xy}^N) with corrections for
 - Evolutionary divergence
 - Sample biases
- 3. Estimate substitution frequencies, q_{xy}^N , from pair counts, A_{xy}^N
- 4. Log odds substitution matrix: $S^{N}[x,y] = c \log \frac{q^{N}_{xy}}{p_{x}p_{y}}$

Log odds substitution matrices

Two sequences have N PAMs divergence, if, on average, N amino acid replacements per 100 residues occurred since their separation

Freq in se



Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

- ➤ PAM matrices, Dayhoff et al, 1978
- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

PAM Matrices

Atlas of Protein Sequence & Structure 1965 - 1978



Examined 1572 changes in 71 groups of closely related proteins



Margaret Dayhoff PhD in Chemistry, 47 Watson Computing Lab Fellow 47 - 48

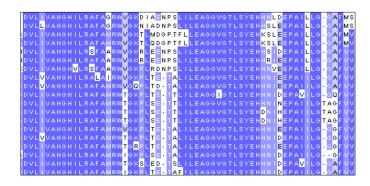
Evolutionary divergence (amino acids)

- PAM: Percent Accepted Mutation
 - Accepted Mutations are mutations that are retained and passed on to future generations
- We say the divergence between two sequences is N PAMs, if, on average, N amino acid replacements per 100 residues (including multiple substitutions) occurred since their separation.

1. Trusted multiple sequence alignments

Examined 1572 changes in 71 groups of closely related proteins

At least 85% identical



Amino Acid Substitution Matrices <u>Parameterized for evolutionary divergence (N)</u>

Overall strategy for both PAM and BLOSUM

- 1. Trusted amino acid alignments
- **2. Obtain amino acid pair counts** $\left(A_{xy}^{N}\right)$ with corrections for
 - Evolutionary divergence
 - Sample biases
- 3. Estimate substitution frequencies, q_{xy}^N , from pair counts, A_{xy}^N
- 4. Log odds substitution matrix: $S^{N}[x,y] = c \log \frac{q^{N}_{xy}}{p_{x}p_{y}}$

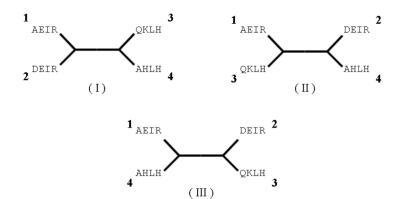
2. Obtain amino acid pair counts (A_{xy}) with corrections for evolutionary divergence and sample biases

Counting amino acid pairs on a tree:

For each unrooted tree with *k* leaves

- > Select the tree(s) that require the fewest substitutions to explain the data
- Count amino acid pairs on the branches of the tree

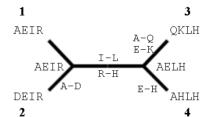
Suppose we have an alignment of four sequences. There are 3 hypotheses (i.e., 3 unrooted trees) for their evolutionary relationships



How to select the tree(s) that require the fewest substitutions to explain the data...

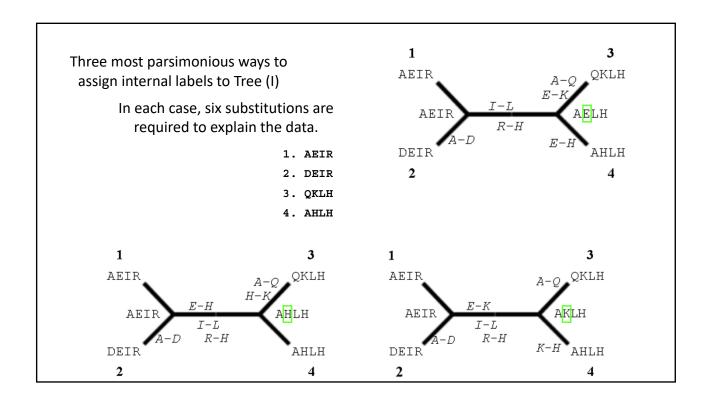
For a given a tree, assign labels to internal nodes that minimize the number of changes required to explain the data

- 1. AEIR
- 2. DEIR
- 3. QKLH
- 4. AHLH

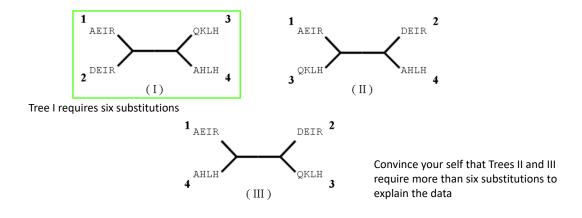


There may be more than one set of labels that satisfies this criterion

Tree I requires six substitutions



Select the most parsimonious tree; i.e., the tree that requires the fewest substitutions to explain the data.

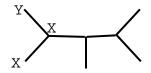


2. Obtain amino acid pair counts (A_{xy}) with corrections for evolutionary divergence and sample biases

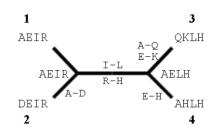
Counting amino acid pairs on a tree:

For each unrooted tree with *k* leaves

- Select the tree(s) that require the fewest substitutions to explain the data
- Count amino acid pairs on the branches of the tree
- · For each branch,
 - if labeled x _____y, $A_{xy}^N=A_{xy}^N$ +1 and $A_{yx}^N=A_{yx}^N$ +1 if labeled x _____x, $A_{xx}^N=A_{xx}^N$ +2



Impact of counting pairs on a tree: some examples



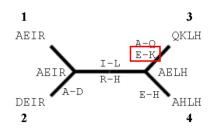
- 1. AEIR
- 2. DEIR
- 3. QKLH
- 4. AHLH

Counts along tree branches

$$A_{DQ} = 0$$

Pairwise counts in the MSA

Impact of counting pairs on a tree: some examples



- 1. AEIR
- 2. D<mark>E</mark>IR
- 3. QKLH
- 4. AHLH

Counts along tree branches

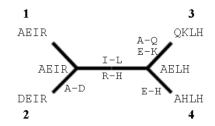
$$A_{DQ} = 0$$

$$A_{EK} = 1$$

Pairwise counts in the MSA

$$A_{DQ} = 1$$

Impact of counting pairs on a tree: some examples



Counts along tree branches

$$A_{DQ} = 0$$

$$A_{EK} = 1$$

- 1. AEIR
- 2. DEIR
- 3. QKLH
- 4. AHLE

Pairwise counts in the MSA

$$A_{DQ} = 1$$

$$A_{EK} = 2$$

$$A_{II} = 4$$

Amino Acid Substitution Matrices <u>Parameterized for evolutionary divergence (N)</u>

Overall strategy for both PAM and BLOSUM

- 1. Trusted amino acid alignments
- 2. Obtain amino acid pair counts $\left(A_{xy}^{N}\right)$ with corrections for
 - Evolutionary divergence
 - Sample biases
- 3. Estimate substitution frequencies, q_{xy}^N , from pair counts, A_{xy}^N
- 4. Log odds substitution matrix: $S^N[x,y] = c \log \frac{q^N}{p_x p_y}$

- 3. Estimate substitution frequencies, q_{xy}^N , from pair counts, A_{xy}
 - Markov model with 20 states (A, C, D, E ... Y)
 - Estimate 1 PAM transition matrix P^1 from $A_{\chi y}$
 - N-PAM transition matrix: $P^N = (P^1)^N$
 - $q_{xy}^N = p_x P_{xy}^N$
 - $S^N[x,y] = c \log \frac{q^N_{xy}}{p_x p_y}$

Is P_{xy}^N a symmetric matrix?

Is $S^N[x,y]$ a symmetric matrix?

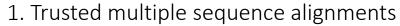
Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

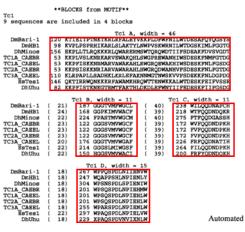
- PAM matrices, Dayhoff et al, 1978
- >BLOSUM (Block Sum) matrices, Hennikoff &Hennikoff, 1991

BLOSUM Matrices

- Trusted data
 - 2000 blocks of conserved regions in ~500 groups of proteins
- Count amino acid pairs: A_{xy}^N
 - Parameterize by evolutionary distance, N
 - Correct for sample bias
- Calculate amino acid frequencies:
 - Related pairs: q_{xy}^N
 - Background pair frequencies calculated from blocks: E_{xy}
- Log likelihood scoring matrix

•
$$S^N = 2 \log_2 \frac{q_{xy}^N}{E_{xy}}$$





~2000 blocks representing 500+ groups of proteins

Automated construction and graphical presentation of protein blocks from unaligned sequences

Steven Henikoff **h.**, Jorja G. Henikoff **1, William J. Alford **2, Shmuel Pietrokovski **3

**Basis Sainces Division, Fred Machinosa Cancer Research Crower, Sentile, IRA 80196, USA

**Howard Maghes Medical Institute, Fred Hutchinosa Cancer Research Contre, Santile, IRA 80196, USA

2. Count amino acid pairs: A_{xy}^N

Parameterize by evolutionary distance, *N* Correct for sample bias

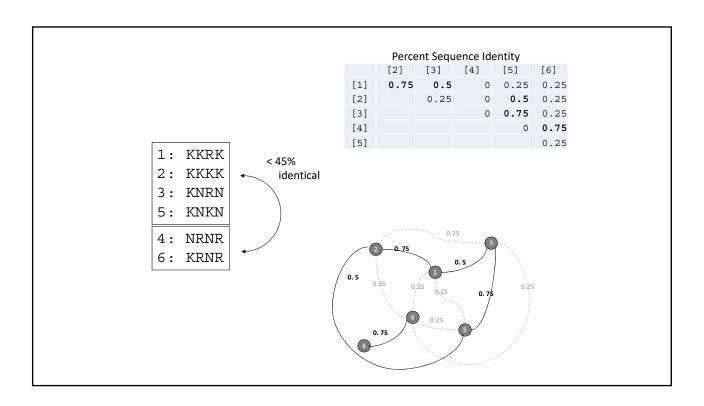
- Cluster sequences such that if s1 and s2 are in different clusters, then identity(s1, s2) < N%
- Count amino acid pairs in s1 aligned with s2 <u>only</u> if s1 and s2 are in different clusters
- Normalize for cluster size

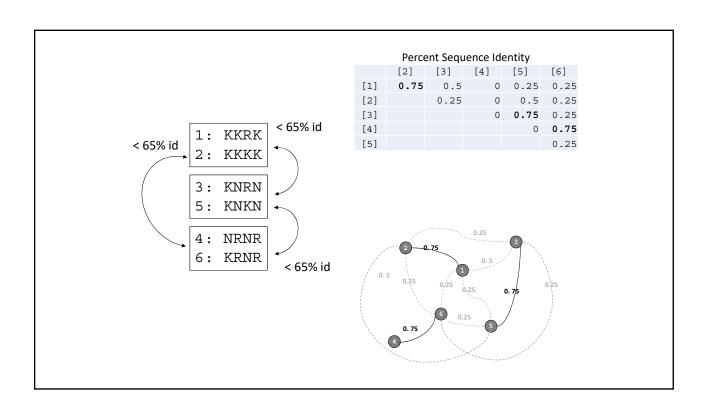
An example...

BLOSUM clustering example

			Perce	ent Sequ	ience Id	ence Identity			
1:	KKRK		[2]	[3]	[4]	[5]	[6]		
2:	KKKK	[1]	0.75	0.5	0	0.25	0.25		
2.	KNRN	[2]		0.25	0	0.5	0.25		
3:	KINKIN	[3]			0	0.75	0.25		
4:	NRNR	[4]				0	0.75		
5:	KNKN	[5]					0.25		
6:	KRNR								

Unclustered sequences: Every sequence is at least 25% identical





Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

- PAM matrices, Dayhoff et al, 1978
- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

Similarities and differences between PAM and BLOSUM

	PAM	BLOSUM
Evolutionary model	Explicit evolutionary model	None
Data	Full length MSAs of closely related sequences.	Conserved blocks. i.e., ungapped local MSAs
Bias correction	Trees	Clustering
Multiple substitutions	Markov model: $P^n = (P^1)^n$	Implicitly represented in data (clustering)
Evolutionary distance	Markov model: $P^n = (P^1)^n$	Clustering
Matrices	Transition and log odds scoring matrices	Log odds scoring matrix only.
Parameter n	Distance increases with n	Distance decreases with <i>n</i>
Biophysical properties	Derived indirectly from data	Derived indirectly from data

