

Sequence profiles, Hidden Markov models and homology modeling

Morten Nielsen,
CBS, Department of Systems Biology, DTU
and
Instituto de Investigaciones
Biotecnológicas, Universidad de San Martín,
Argentina

Summary

- Sequence logo is a power tool to visualize (binding) motifs
 - Information content identifies essential residues for function and/or structural stability
 - Weight matrices can be derived from very limited number of data using the techniques of
 - Sequence weighting
 - Pseudo counts
-

Sequence Profiles and Weight matrices

- Alignments based on conventional scoring matrices (BLOSUM62) scores all positions in a sequence in an equal manner
 - Some positions are highly conserved, some are highly variable (more than what is described in the BLOSUM matrix)
 - Sequence profile are ideal suited to describe such position specific variations
-

Sequence alignment

- Conventional sequence alignment uses a (Blosum) scoring matrix to identify amino acids matches in the two protein sequences

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	-2	2	-3
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	-1	1
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	2	-3
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

Alignment scoring matrices

- Blosum62 score matrix. Fg=1. Ng=0?

	L	A	G	D	S	D
F						
I						
G						
D						
S						
L						

A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	-2	2	-3
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	1	1
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2
T	0	-1	0	-1	-1	-1	-2	-2	-1	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	2	-3
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

Alignment scoring matrices

- Blosum62 score matrix. Fg=1. Ng=0?

	L	A	G	D	S	D
F	0	-2	-3	-3	-2	-3
I	2 → -1	-4	-3	-2	-3	
G	-4	0	6 → -1	0	-1	
D	-4	-2	-1 → 6	0	6	
S	-2	1	0	0 → 4	0	
L	4	-1	-4	-4	-2	-4

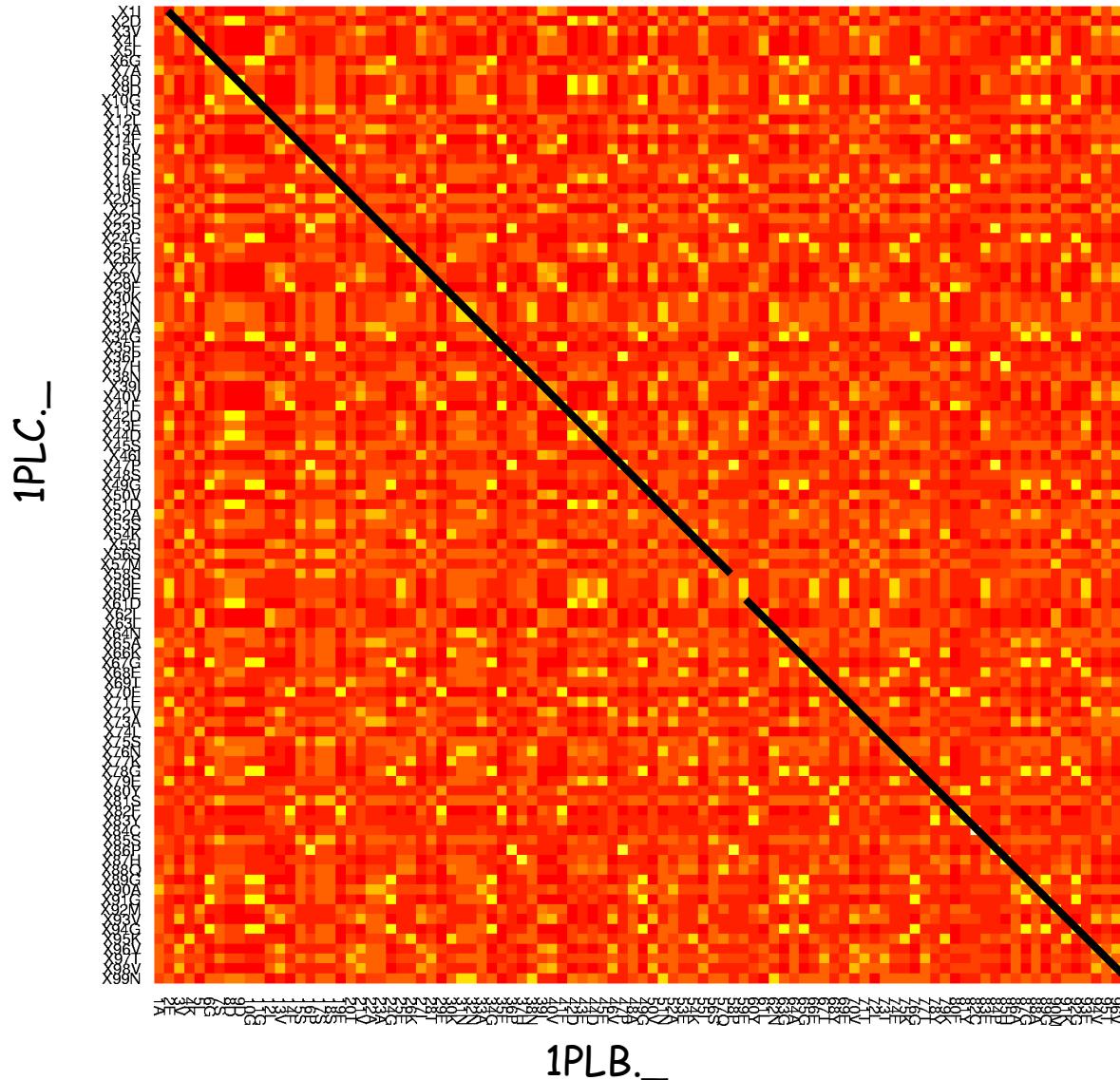
A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-2	-1	-2	-2	2	-2	-3
I	-1	-3	-3	-3	-1	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3	
L	-1	-2	-3	-4	-1	-2	-3	-4	3	2	-4	2	1	0	-3	-2	-1	-2	-1	1
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	-1	1
F	-2	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1	
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
S	1	-1	1	0	-1	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2	
T	0	-1	0	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0	
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-1	1	-4	-3	-2	11	2	-3	
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

- Score = 2-1+6+6+4=17

LAGDS

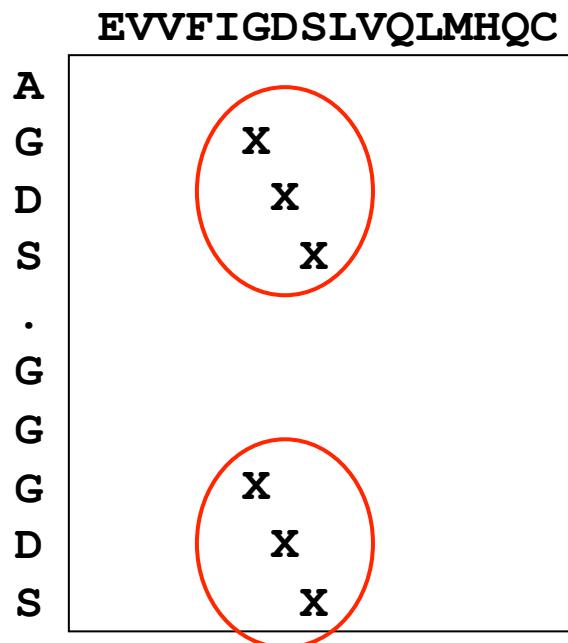
I-GDS

When Blast works!



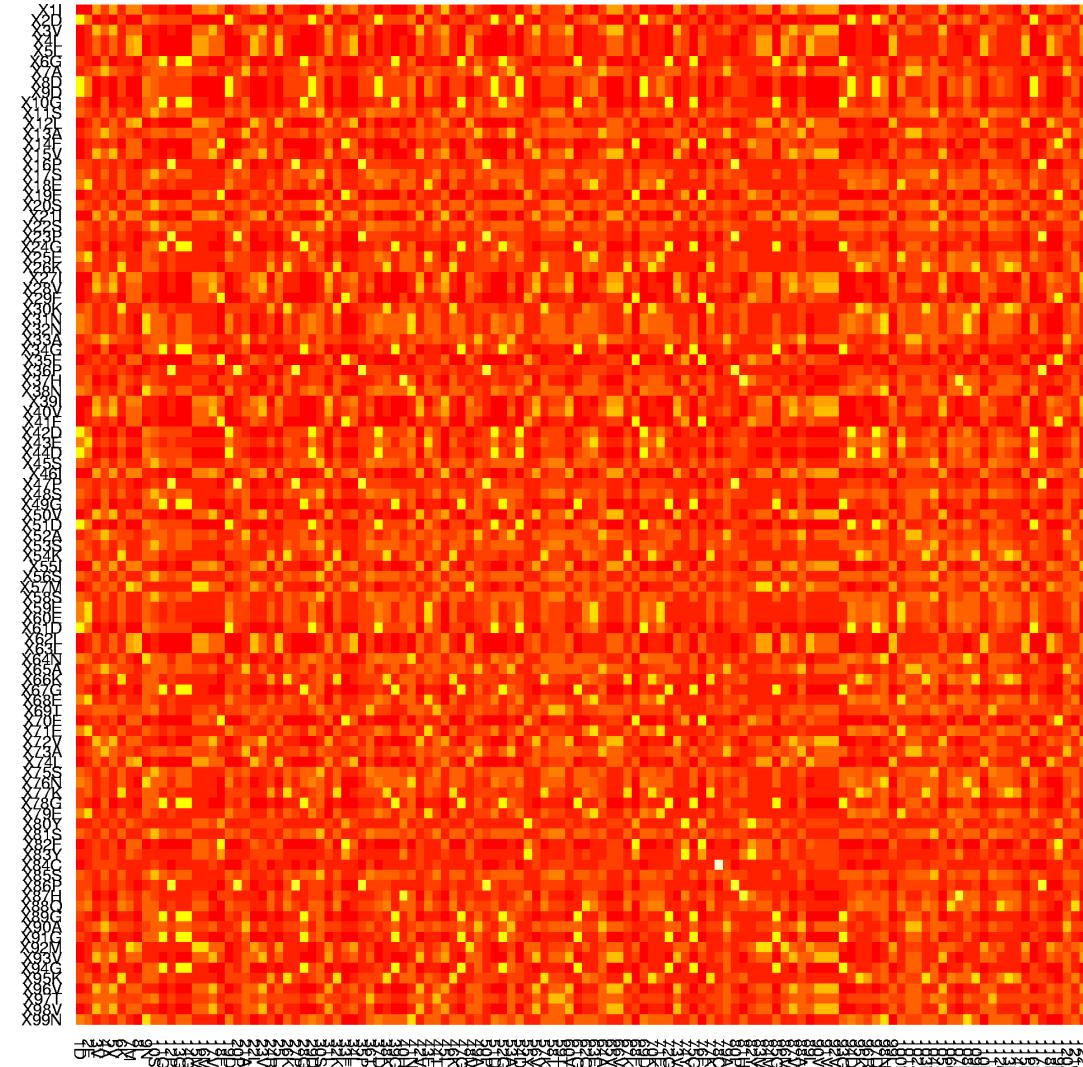
What goes wrong when Blast fails?

- Conventional sequence alignment uses a (Blosum) scoring matrix to identify amino acids matches in the two protein sequences
- This scoring matrix is identical at all positions in the protein sequence!



When Blast fails!

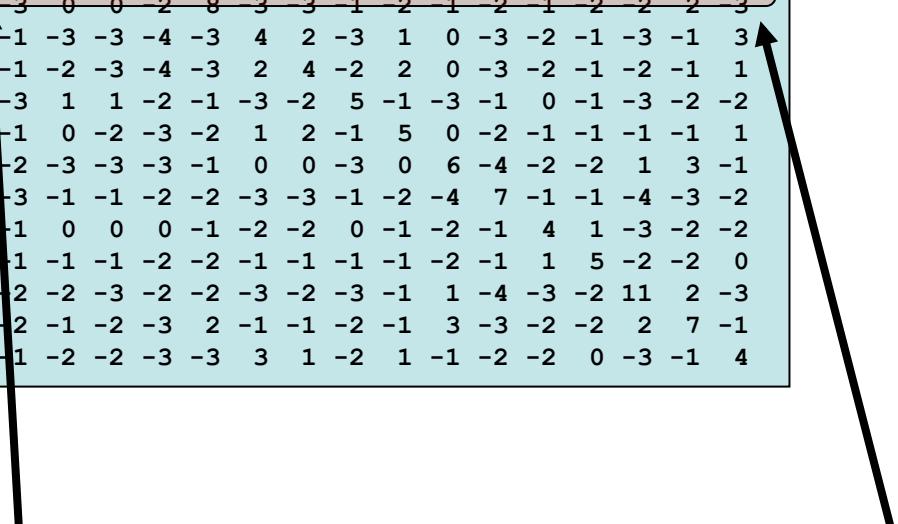
CENTER FOR BIOLOGICAL
CALCULUS AND ANALYSIS CBS



1PMY.

Alignment match scores

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3
H	2	0	1	1	3	0	0	2	8	3	3	1	2	1	2	1	2	2	2	3
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	-1	1
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
S	1	-1	1	0	-1	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2	
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	2	-3
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4



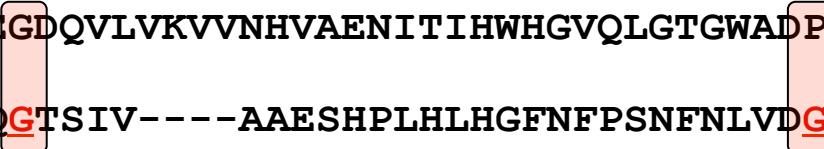
TVNGQ--FPGPRLAGVAREGDQVLVKVVNHVAENITIHWGVQLGTGWADGPAYVTQCPI

Sequence profiles

- In reality not all positions in a protein are equally likely to mutate
 - Some amino acids (active sites) are highly conserved, and the score for mismatch must be very high
 - Other amino acids can mutate almost for free, and the score for mismatch should be lower than the BLOSUM score
- Sequence profiles can capture these differences

Sequence profiles

TVNGQ--FPGPRLAGVAREGDQVLVKVVNHVAENITIHWHGVLGTGWADPPAYVTQCPITKAVVLTFTNTSVEICLVMQGTSIV----AAESHPLHLHGFNFPNSNFLVDGMERNTAGVP



Sequence profiles

Conserved Non-conserved

```
ADDGSLAFVPSEF--SISPGEKIVFKNNAGFPNIVFDEDSIPSGVDASKISMSEEDLLN
TVNGAI--PGPLIAERLKEGQNVVRTNTLDEDTSIHWHGLLVPFGMDGVPGVSFPG---I
-TSMAPAFGVQEFYRTVKQGDEVTVTIT----NIDQIED-VSHGFVVVNHGVSME---I
IE--KMKYLTPEVFYTIKAGETVYWVNGEVMPHNVAFKKGIV--GEDAFRGEMMTKD---
-TSVAPSFSQPSF-LTVKEGDEVTVIVTNLDE-----IDDLTHGFTMGNHGVAME---V
ASAETMVFEPEFLVLEIGPGDRVRFVPTHK-SHNAATIDGMVPEGVEGFKSRINDE---
TVNGQ--FPGPRLAGVAREGDQVLVKKVNHVAENITIHWGVQLGTGWADPPAYVTQCP
TKAVVLTFTNTSVEICLVMQGTSIV---AAESHPLHLHGFPNSPNFLVDGMERNTAGVP
```

Matching any thing
but $G \Rightarrow$ large
negative score

Any thing can match

$$p_a = \frac{\alpha \cdot f_a + \beta \cdot g_a}{\alpha + \beta}$$

$$g_b = \sum_a f_a \cdot q_{b|a}$$

How to make sequence profiles

1. Align (BLAST) sequence against large sequence database (Swiss-Prot)
 2. Select significant alignments and make sequence profile
 3. Use profile to align against sequence database to find new significant hits
 4. Repeat 2 and 3 (normally 3 times!)
-

Sequence logos. Visualization of sequence profiles

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

$$I = \sum_a p_a \log\left(\frac{p_a}{q_a}\right)$$

$$P_A = 6/10 = 0.6$$

$$P_G = 2/10 = 0.2$$

$$P_T = P_K = 1/10 = 0.1$$

$$P_C = P_D = \dots P_V = 0.0$$

$$q_A = 0.07$$

$$q_G = 0.07$$

$$q_T = 0.05$$

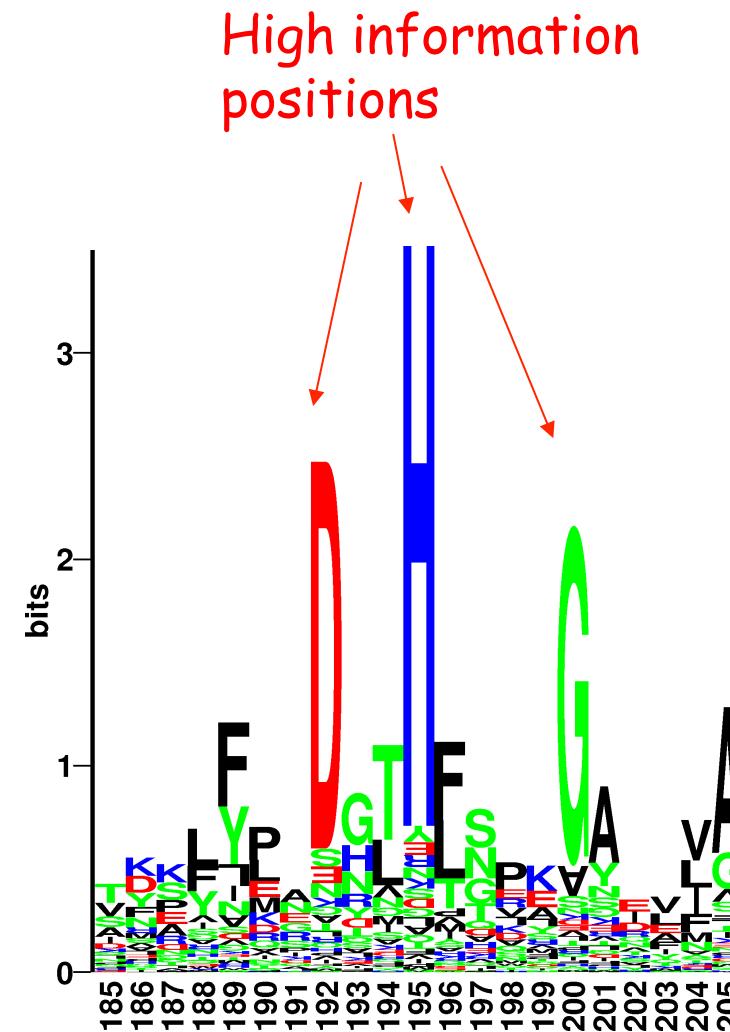
$$q_K = 0.06$$

ALAKAAAAM
ALAKAAAAN
ALAKAAAAR
ALAKAAAAT
ALAKAAAAV
GMNERPILT
GILGFVFTM
TLNAWVKVV
KLNEPVLLL
AVVPFIVSV

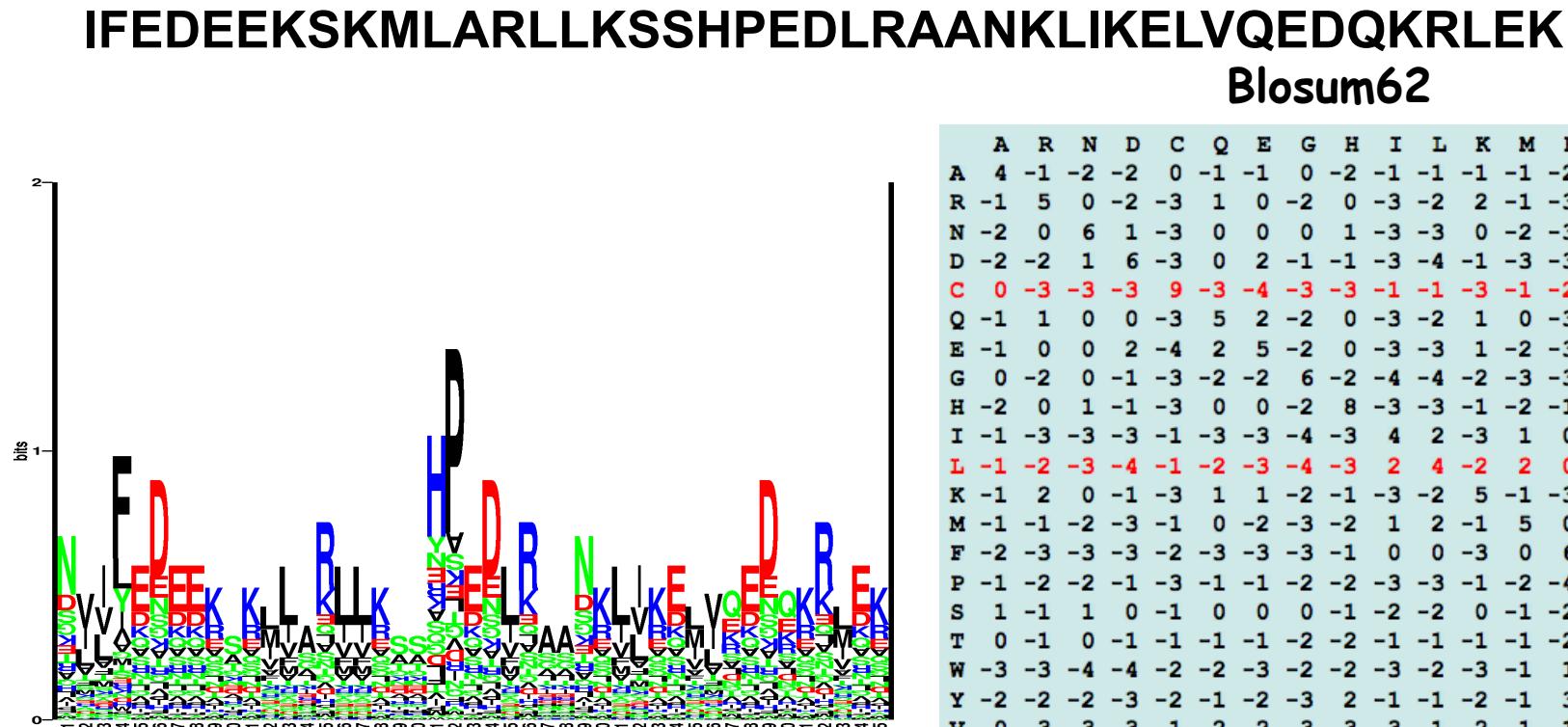
Sequence logos

$$I = \sum_a p_a \log\left(\frac{p_a}{q_a}\right)$$

- Height of a column equal to I
- Relative height of a letter is p
(letters are upside down if $q>p$)

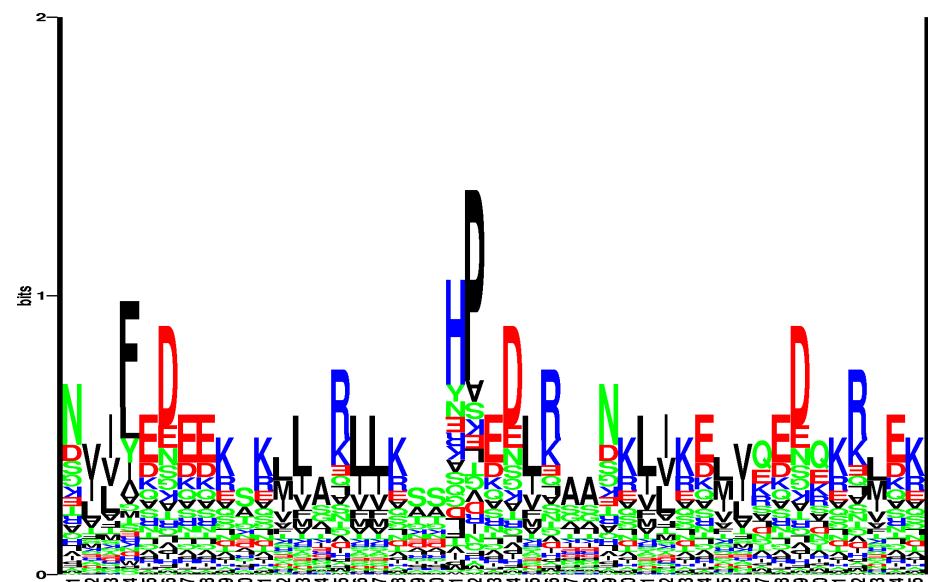


Sequence profiles (1J2J.B)



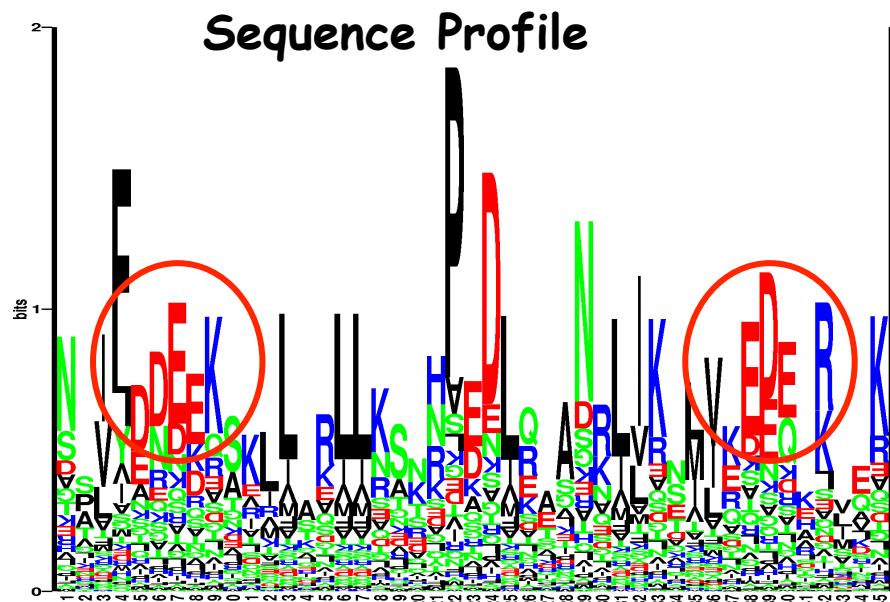
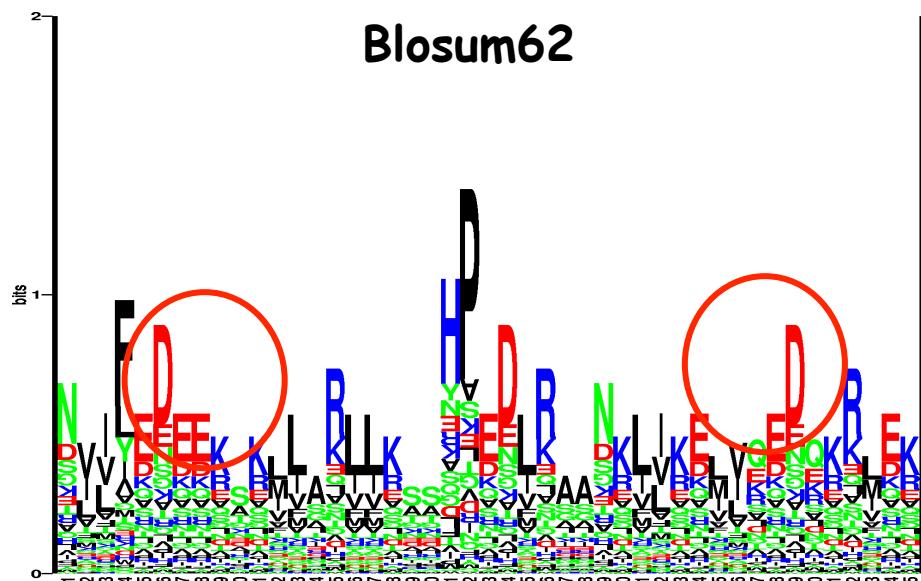
$$W_{ij} = \log(p_{ij}/q_j)$$

Sequence profiles (1J2J.B)



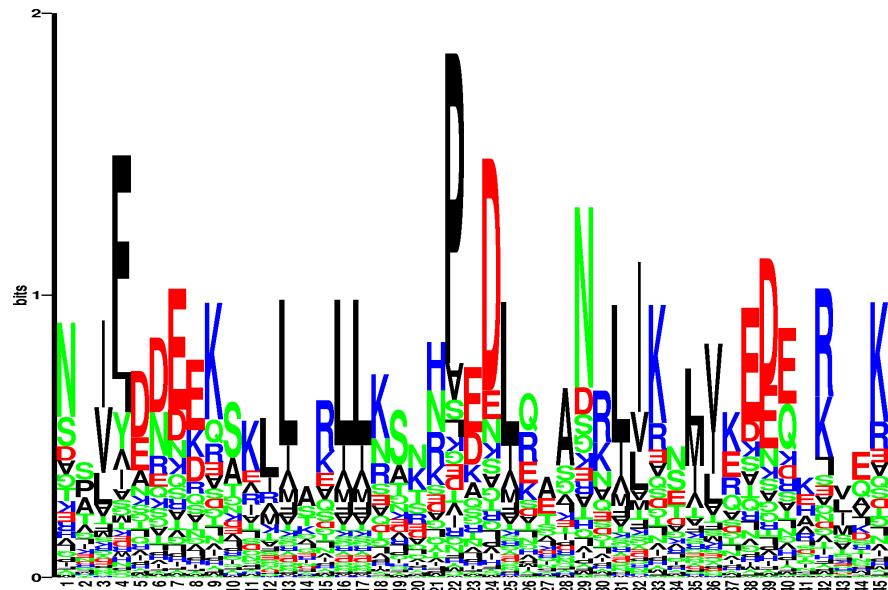
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
0 I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
1 F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1
2 E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
3 D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
4 E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
5 E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
6 K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2
7 S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2
8 K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2
9 M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	1	
10 L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
11 A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
12 R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
13 L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1

Sequence profiles (1J2J.B)



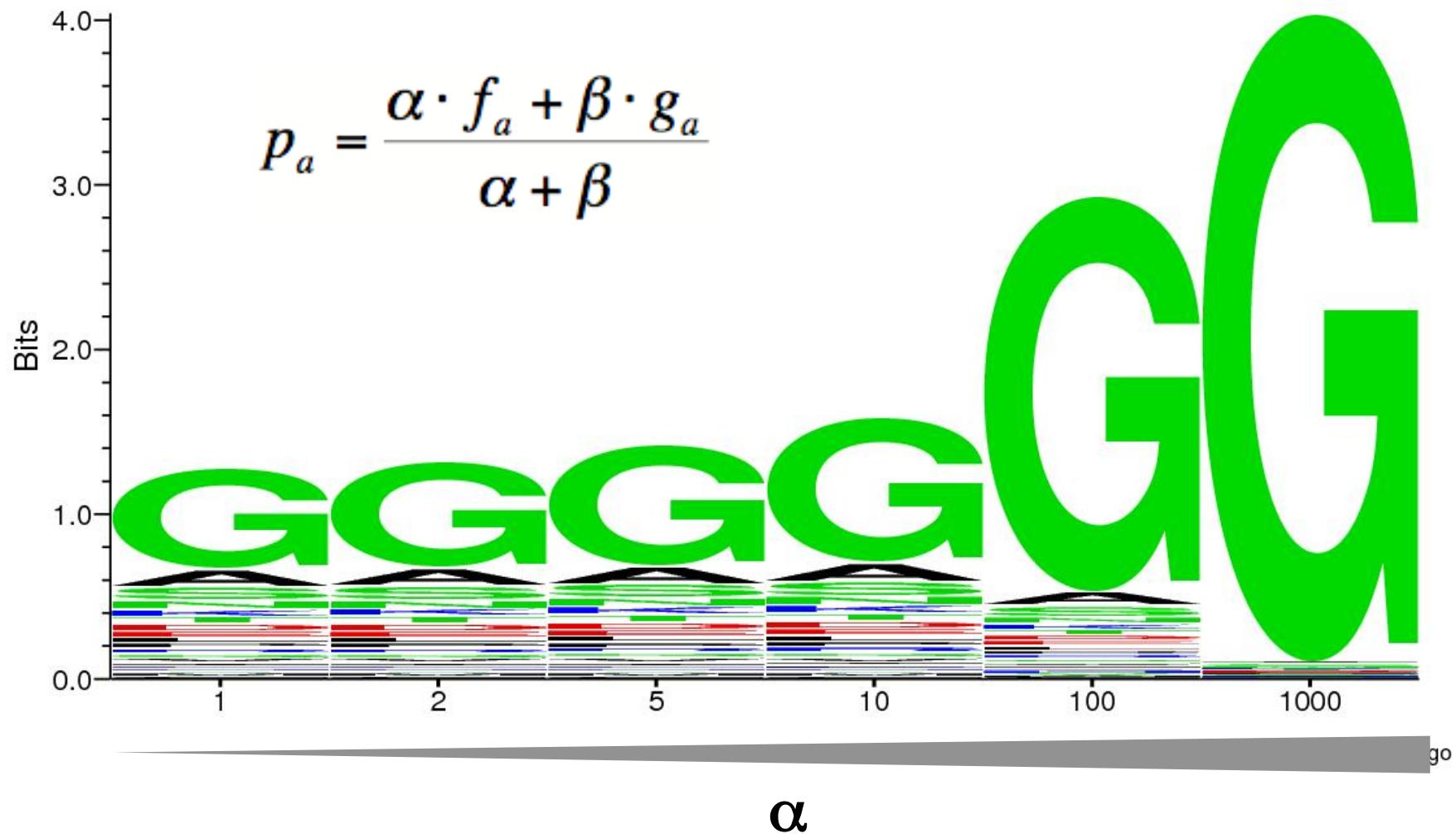
Sequence profiles (1J2J.B)

IFEDEEKSKMLARLLKSSHPEDLRAAANKLIKELVQEDQKRLEK



	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
1	I	-1	-3	-3	-3	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	-1	-3	-2
2	F	-3	-3	-3	-4	-3	-4	-4	-3	-1	0	0	-3	0	7	-4	-3	-2	1	3
3	E	-2	-1	1	5	-4	1	4	-2	-1	-4	-4	0	-3	-4	-2	0	-1	-4	-3
4	D	-2	-2	1	6	-4	-1	1	-2	-1	-3	-4	-1	-3	-4	-2	0	-1	-5	-3
5	E	-1	0	0	1	-4	2	5	-2	0	-4	-3	1	-2	-4	-1	0	-1	-3	-2
6	E	-1	0	0	1	-4	2	5	-2	0	-4	-3	1	-2	-4	-1	0	-1	-3	-2
7	K	-1	2	0	-1	-3	1	1	-2	-1	-3	-3	5	-2	-3	-1	0	-1	-3	-2
8	S	1	-1	0	0	-1	0	0	0	-1	-3	-3	0	-2	-3	-1	5	1	-3	-2
9	K	-1	2	0	-1	-3	1	1	-2	-1	-3	-3	5	-2	-3	-1	0	-1	-3	-2
10	M	-1	-2	-3	-4	-2	-1	-3	-4	-2	1	3	-2	5	0	-3	-2	-1	-2	1
11	L	-2	-2	-4	-4	-2	-2	-3	-4	-3	1	4	-3	2	0	-3	-3	-1	-2	-1
12	A	4	-2	-1	-2	-1	-1	-1	0	-2	-1	-2	-1	-1	-2	-1	1	2	-3	-2

Sequence profiles or Gaining confidence



Example.

>1K7C.A

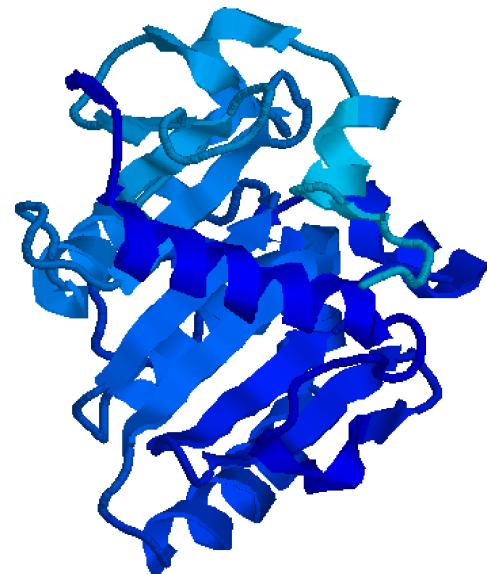
TTVYLAGDSTMAKNGGGSGTNGWGEYLASYLSATVVNDAVAGRSARSYTREGRFENIADV
VTAGDYVIVEFGHNDGGSLSTDNGRTDCSGTGAEVCYSVYDGVNETILTFPAYLENAAKL
FTAKGAKVILSSQTPNNPWETGTFVNSPTRFVEYAAEVAGVEYVDHWSYVDSIYETL
GNATVNSYFPIDHTHTSPAGAEVVAEAFLKAVVCTGTSLKSVLTTSFEGTCL

- What is the function
 - Where is the active site?
-

What would you do?

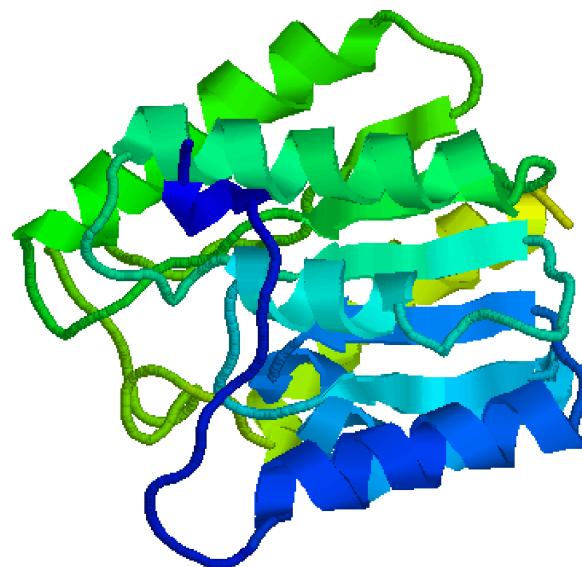
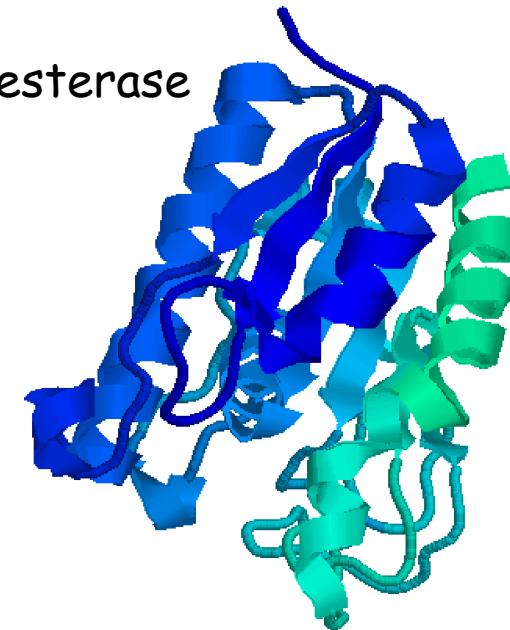
- Function
 - Run Blast against PDB
 - No significant hits
 - Run Blast against NR (Sequence database)
 - Function is Acetyl esterase?
 - Where is the active site?
-

Example. Where is the active site?



1USW Hydrolase

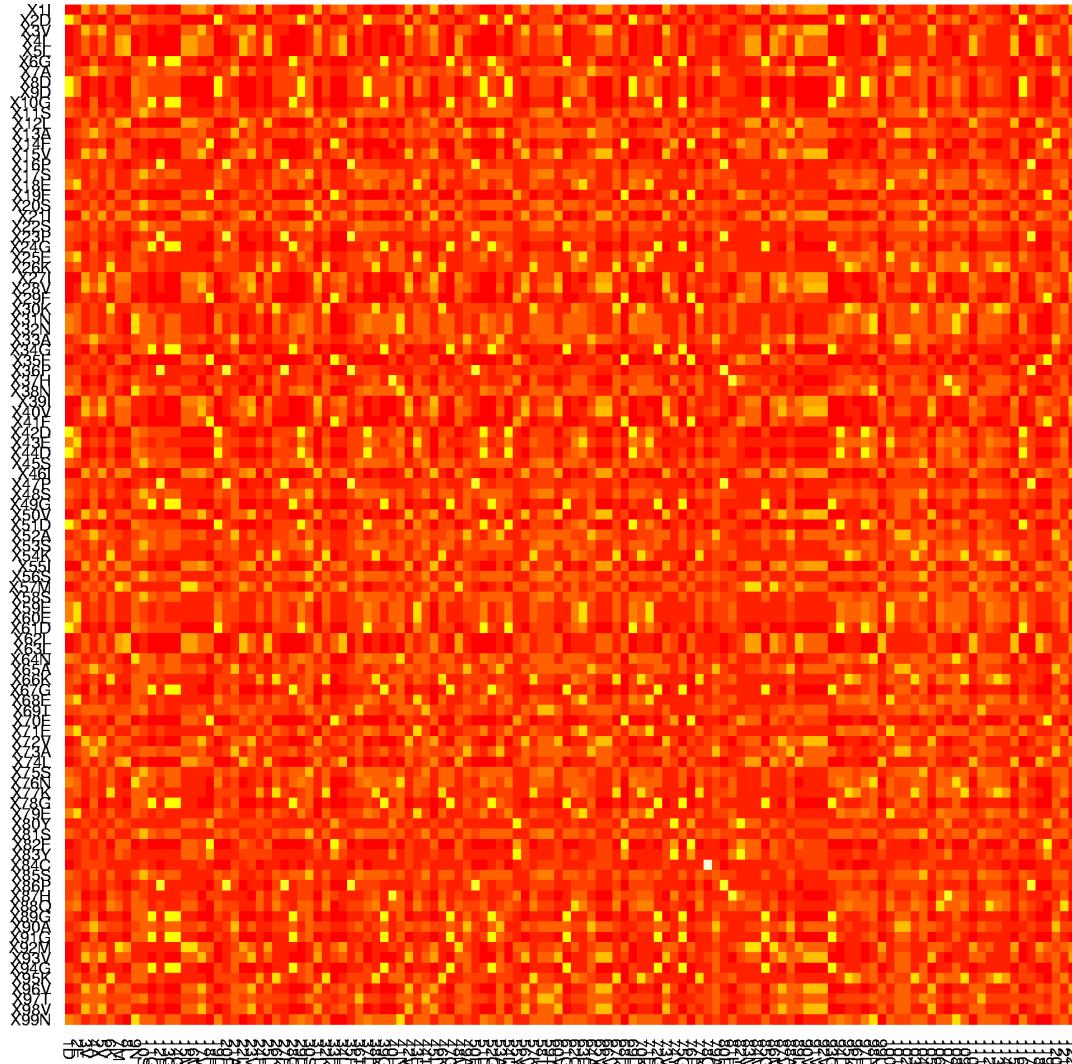
1G66 Acetylxyran esterase



1WAB Acetylhydrolase

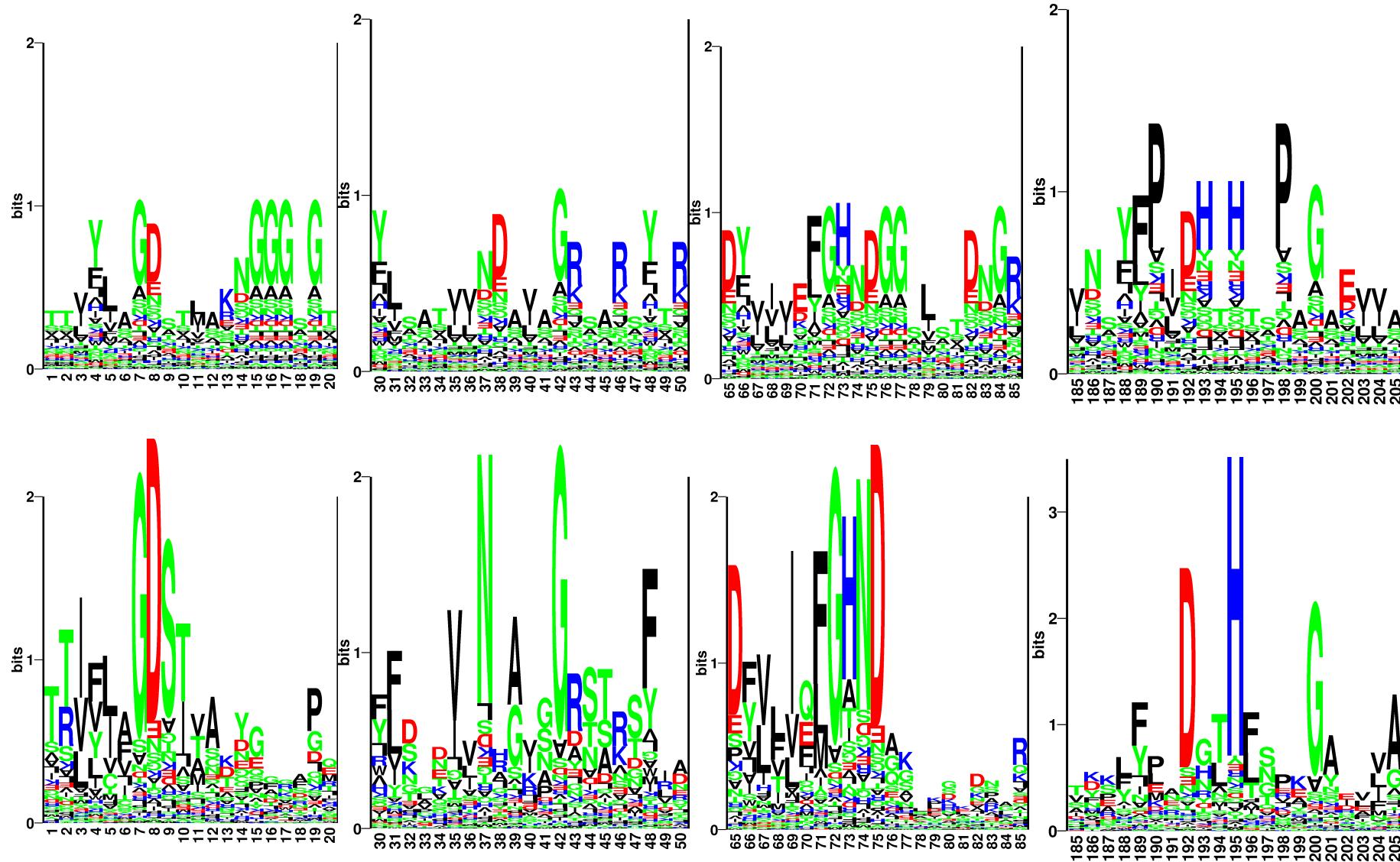
When Blast fails!

1K7A.A



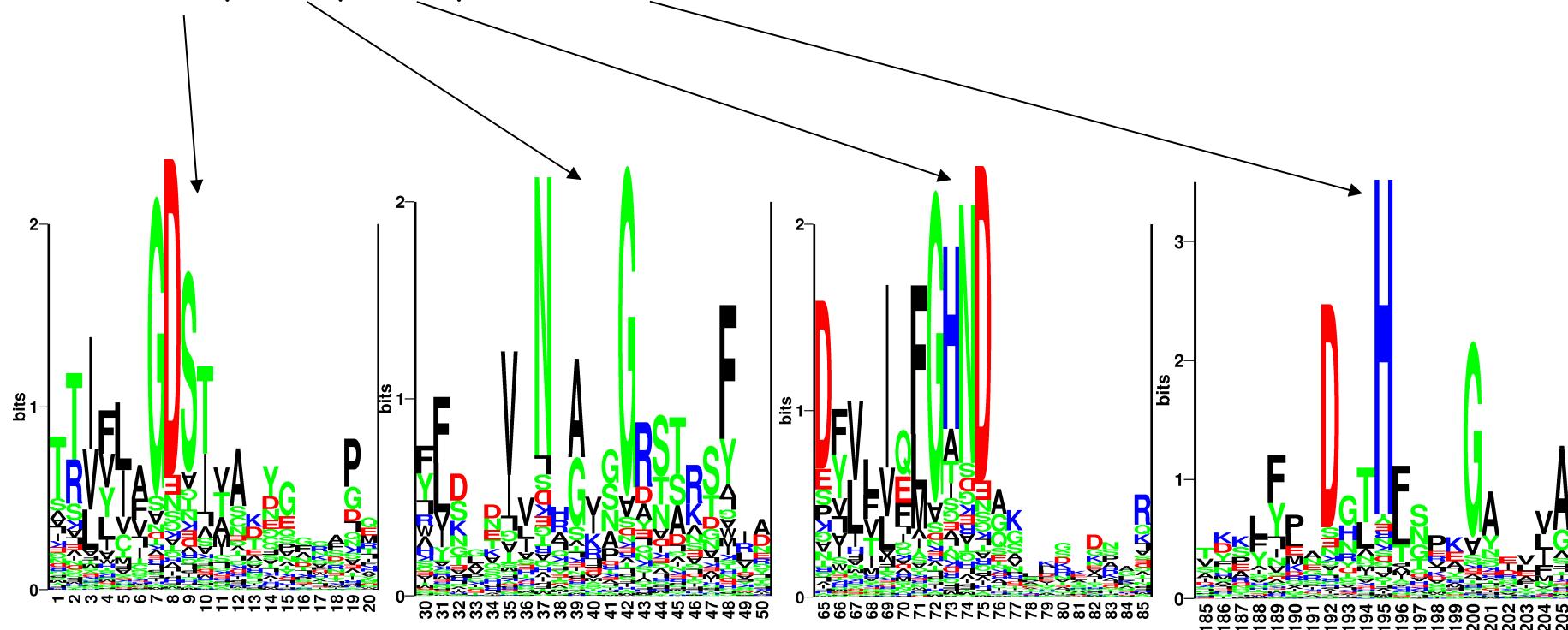
1WAB._

Example. (SGNH active site)

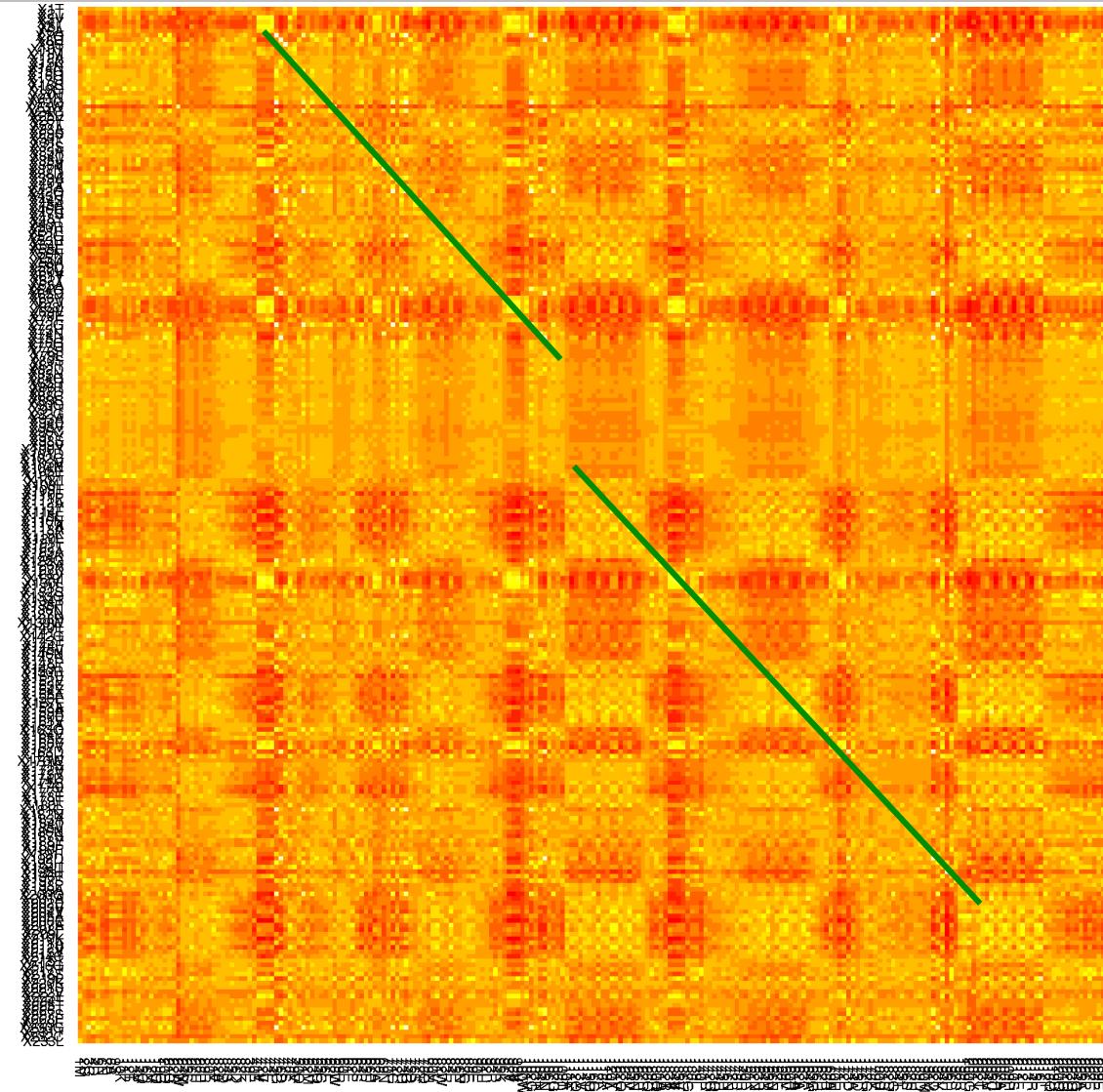


Example. Where is the active site?

- Sequence profiles might show you where to look!
- The active site could be around
 - S9, G42, N74, and H195



Profile-profile scoring matrix



Example. Where is the active site?

Align using sequence profiles

ALN 1K7C.A 1WAB._ RMSD = 5.29522. 14% ID

1K7C.A TVYLAGD**S**TMAKNGGGSGTNGGEYLASYLSATVVNDAVA**G**RSARSARYTREGRFENIADVVTAGDYVIVEFGH**N**DGGSLSTDN
1WAB._ EVVFIGD**S**LVQLMHQCE---IWRELFS---PLHALNFGIGG**G**DSTQHVLW--RLENGELEHIRPKIVVVVWGT**N**NHG-----

1K7C.A GRTDCSGTGAEVCYSVYDGNETILTFPAYLEAAKLFTAK--GAKVILSSQTPNNPWETGTFVNSPTRFVEYael-AAEVA
1WAB._ -----HTAEQVTGGIKAIQLVNERQPQARVVVLGLLPRGQ-HPNPLREKNRRVNELVRAALAGHP

1K7C.A GVEYVDHWSYVDSIYETLGNATVNSYFPIDHT**H**TSPAGAEVVAEAFLKAVVCTGTSL
1WAB._ RAHFILDADPG---FVHSDG--TISHHDMDYL**H**LSRLGYTPVCRALHSLLLRL---L

Handout exercise

Using Psi-Blast Profiles

Blast2logo

Blast2logo 1.0 Server

[Instructions](#)

[Output format](#)

SUBMISSION

Paste a single sequence in [FASTA](#) format into the field below:

```
>Ex
VALAELYIPEVARRLLGQGWHEDECTFAEVТИGTLQAILRDIATWSADEGGMRDGPALVLLPPG
EQHTLGAMVAVAKLRLGVSVCLRMSTCPAELRELFGRKRRDAIMISLAHAEMLEVGRKLVKTLKD
MTGGRIPVAMCGALFLDGTEAASIPEADIVTNDEAALQ
```

Submit a file in [FASTA](#) format directly from your local disk:

no file selected

Upload a file in [BLAST PROFILE](#) format:

no file selected

Blast Database

Number of Blast iterations

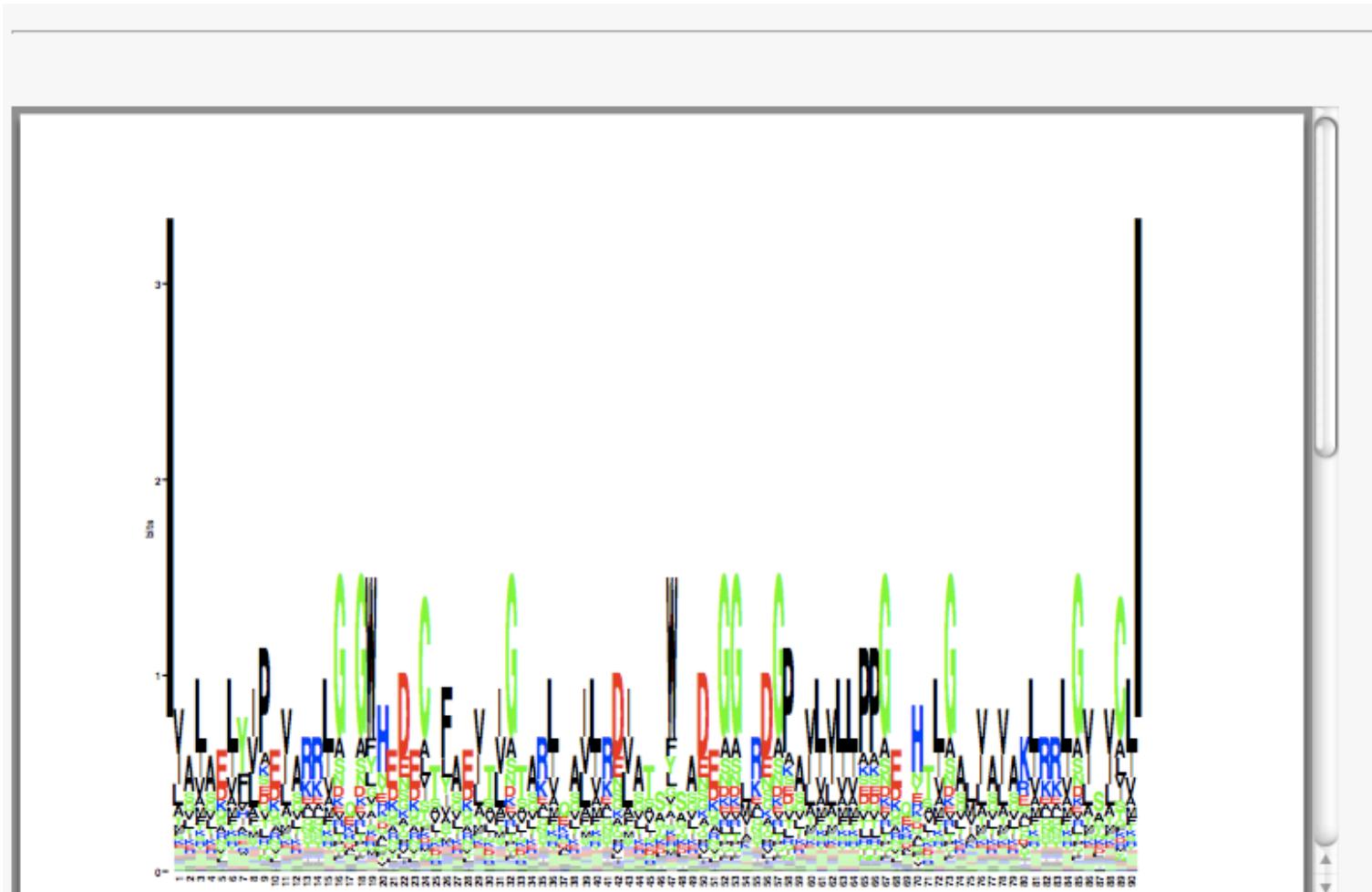
Blast E-value cutoff

Stack Linesize

Plot Kullback-Leibler logo

File format for logo file

Blast2logo



Download logo file [Logo](#)

Link to Blastprofile output file [Blast.prof](#)

Blast2logo

Last position-specific scoring matrix computed

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
1	V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
2	A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
3	L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
4	A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
5	E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
6	L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
7	Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
8	I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
9	P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
10	E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
.	

A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-2	-1	1	0	-3	-2	0	
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	-2	2	-3
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	1	
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2
T	0	-1	0	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0	
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	2	-3
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

Blast2logo

Blast2logo 1.0 Server

Instructions

Output format

SUBMISSION

Paste a single sequence in [FASTA](#) format into the field below:

```
>Ex
VALAELYIPEVARRLGQGWHEDECTFAEVТИGTRLQAILRDIATWSADEGGMRDGPALVLLPPG
EQHTLGMAMAVAKLRLGVSVCLRMSTCPAELRELFGKRRFDAIMSLAHAEMLLEVGRKLVKTLKD
MTGGRIPVAMGGALFLDGTEASIPADIVTNDEAALQ
```

Submit a file in [FASTA](#) format directly from your local disk:

Choose File no file selected

Upload a file in [BLAST PROFILE](#) format:

Choose File no file selected

Blast Database NR70

Number of Blast iterations

Blast E-value cutoff

Stack Linesize

Plot Kullback-Leibler logo

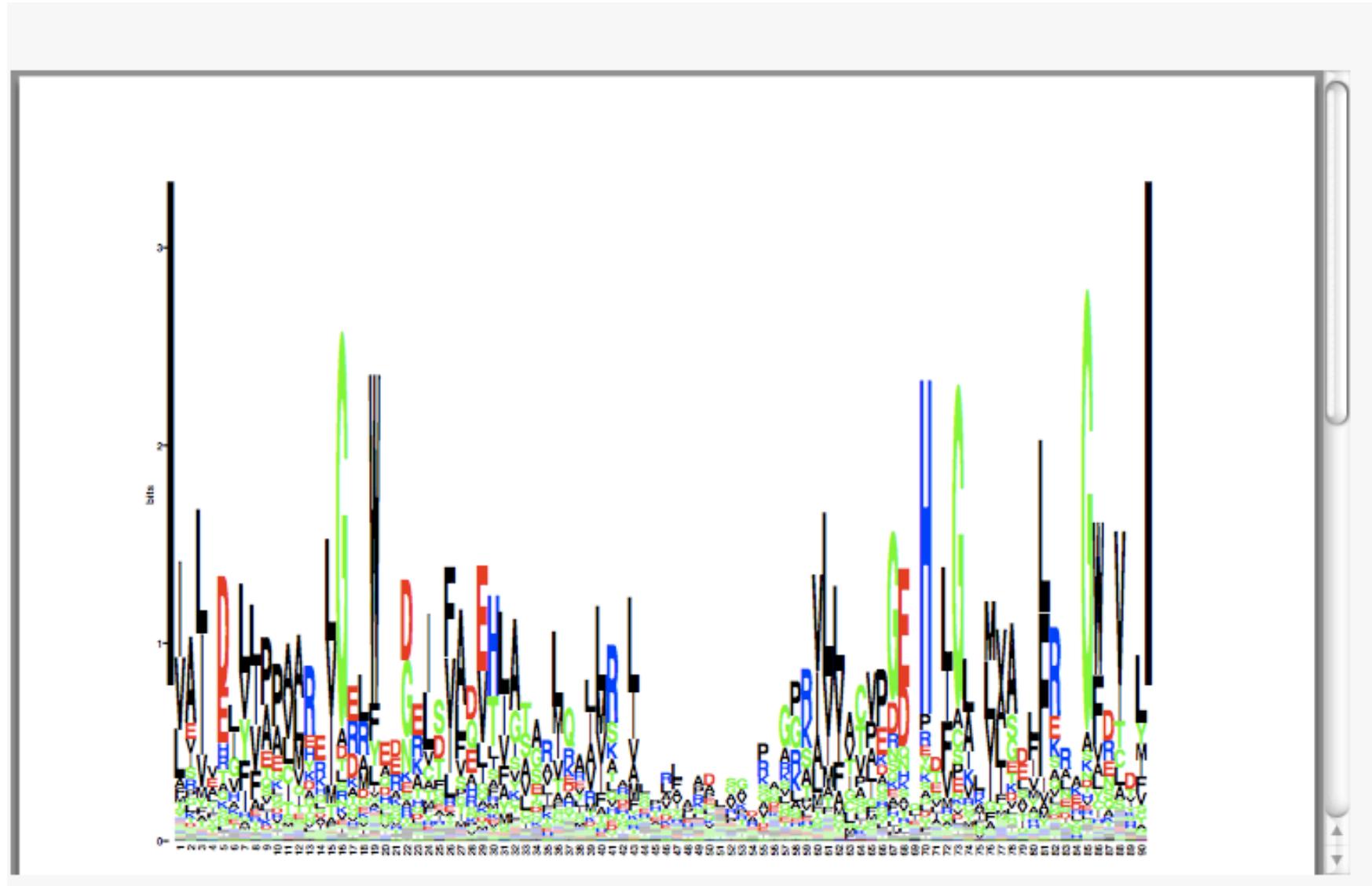
File format for logo file PDF

Submit Clear fields

Restrictions:

At most 1 sequences per submission; each sequence not more than 20,000 amino acids.

Blast2logo



Blast2logo

Last position-specific scoring matrix computed,

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
1	V	-2	-4	-4	-5	-2	-4	-4	-5	-4	5	2	-4	0	-1	-4	-3	-2	-4	-2	4
2	A	5	0	-3	-3	-3	-2	1	-2	-3	0	-3	-2	-2	-4	0	0	-2	-4	-3	0
3	L	-4	-5	-6	-6	-4	-5	-5	-6	-5	5	4	-5	1	-2	-5	-5	-3	-4	0	1
4	A	1	-4	-1	-1	3	-1	2	-4	-3	0	-1	-2	-3	1	-4	0	0	-4	2	2
5	E	-2	0	-2	6	-6	0	<u>4</u>	-4	2	-5	-5	-2	-5	-6	<u>-4</u>	-2	0	-6	-4	-5
6	L	-1	-2	-4	-4	-4	-2	-1	2	3	3	2	-1	0	-2	-5	-1	-1	-5	-3	1
7	Y	-4	-5	-5	-6	-4	-5	-5	-4	0	1	4	-5	-1	3	-5	-5	-4	-3	5	3
8	I	-1	-2	-5	-5	-4	-5	-2	-6	-5	4	3	-5	-1	3	-5	-4	-2	-4	-1	3
9	P	3	-4	-4	-3	-4	1	1	-4	-2	-2	-3	-2	-4	-5	6	-1	0	-5	-5	-2
10	E	2	-2	-3	-2	-3	0	<u>1</u>	-1	-3	-4	-3	-1	-1	-4	<u>6</u>	-2	-2	-4	-4	-3

Sequence profiles take home message

- Blast will often fail to recognize sequence relationships for low homology sequence pairs
 - Sequence profiles contain information on conserved/variable residues in a protein sequence
 - Sequence profiles are calculated from (multiple) sequence alignments
 - Iterative Blast enables homology recognition also for low sequence similarity
 - Sequence profiles give information on residues essential for protein function and protein structure
 - Can be used to predict impact of SNP's on protein function
 - This is often done using the Blosum matrix, but profiles are much more precise
-

Summary

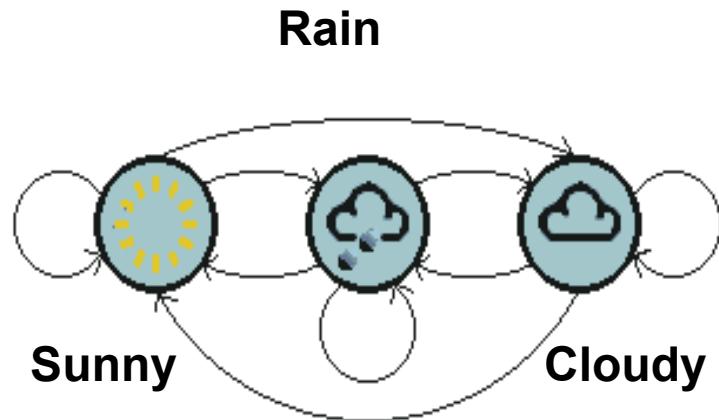
- Sequence logo is a power tool to visualize (binding) motifs
 - Information content identifies essential residues for function and/or structural stability
 - Weight matrices and sequence profiles can be derived from very limited number of data using the techniques of
 - Sequence weighting
 - Pseudo counts
 - Weight matrices and sequences profiles can accurately describe binding motifs, **sequence conservation, active sites, functional impact of SNP's, ...**
-

Hidden Markov Models, HMM's

Objectives

- Introduce Hidden Markov models and understand that they are just weight matrices with gaps
 - How to construct an HMM
 - How to “align/score” sequences to HMM’s
 - Viterbi decoding
 - Forward decoding
 - Backward decoding
 - Posterior Decoding
 - Use and construct Profile HMM
 - HMMer
-

Markov Chains



States : Three states - sunny, cloudy, rainy.

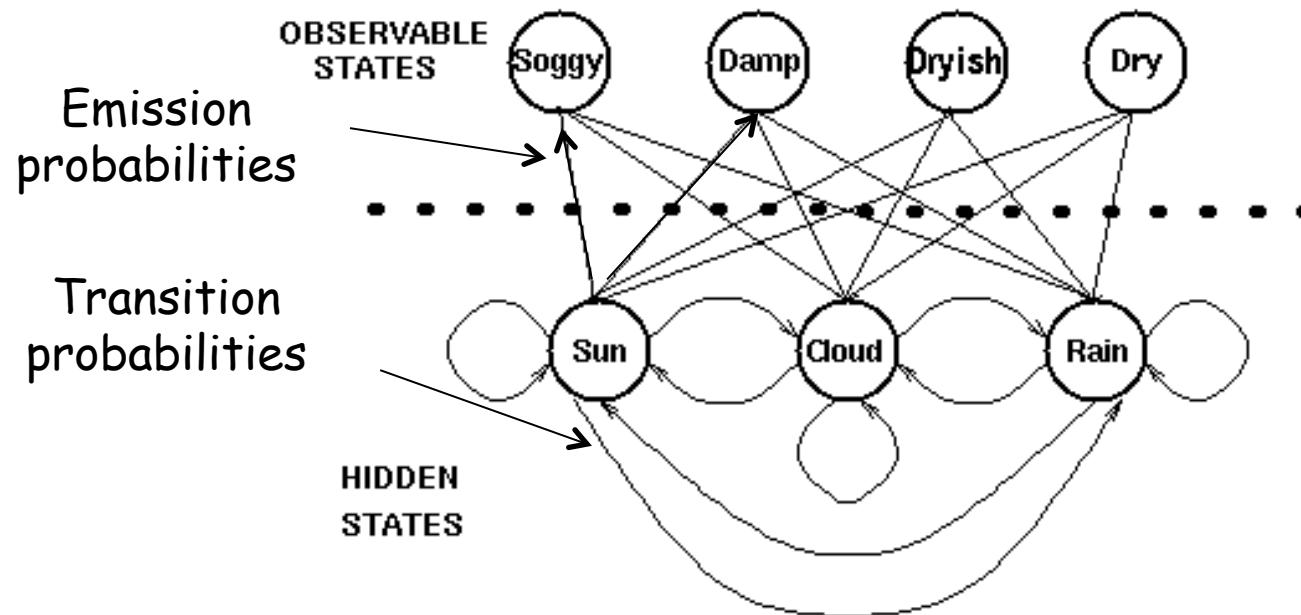
			weather today		
		Sun	Cloud	Rain	
weather yesterday	Sun	0.5	0.25	0.25	
	Cloud	0.375	0.125	0.375	
		Rain	0.125	0.625	0.375

State transition matrix : The probability of the weather given the previous day's weather.

$$\begin{pmatrix} \text{Sun} & \text{Cloud} & \text{Rain} \\ 1.0 & 0.0 & 0.0 \end{pmatrix}$$

Initial Distribution : Defining the probability of the system being in each of the states at time 0.

Hidden Markov Models



Hidden states : the (TRUE) states of a system that may be described by a Markov process (e.g., the weather).

Observable states : the states of the process that are 'visible' (e.g., seaweed dampness).

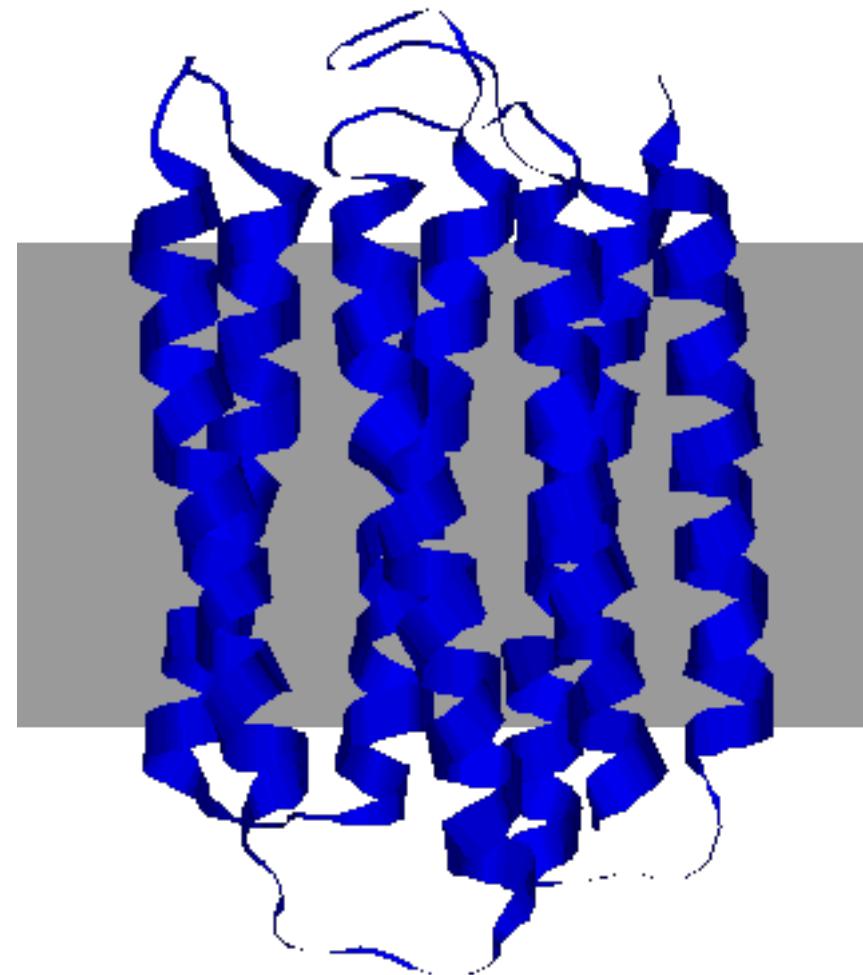
TMHMM (trans-membrane HMM) (Sonnhammer, von Heijne, and Krogh)

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

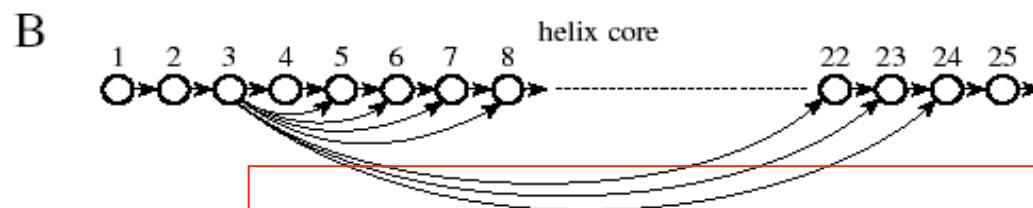
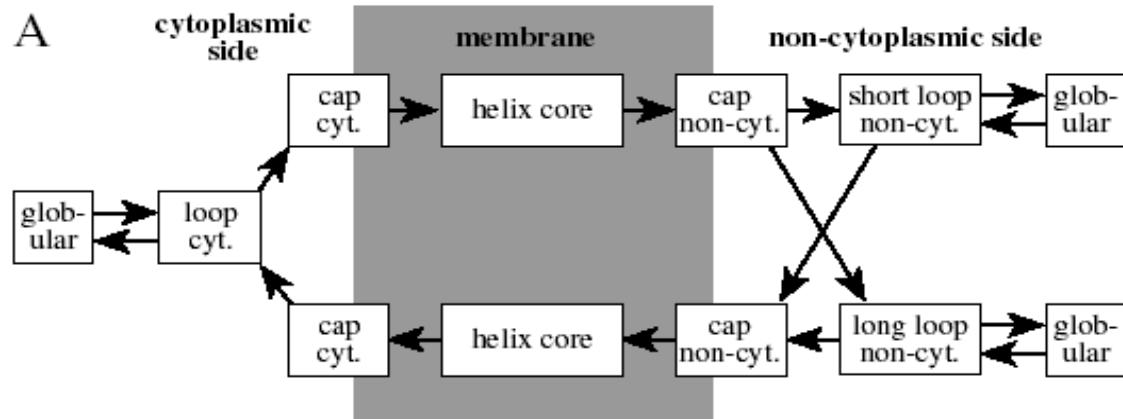
Extra cellular

Trans membrane

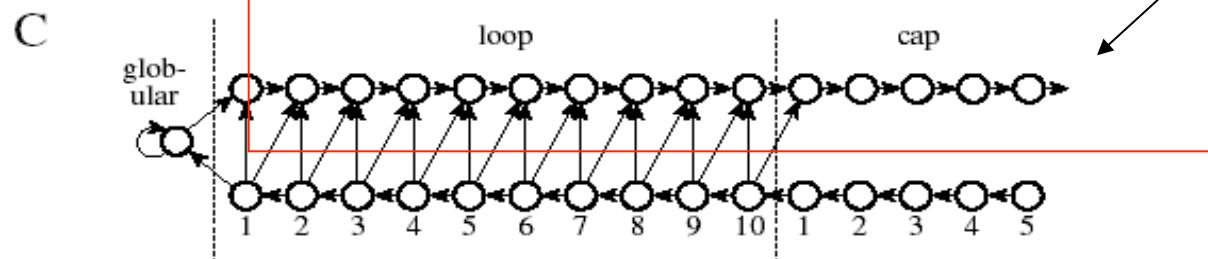
Intra cellular



TMHMM (trans-membrane HMM) (Sonnhammer, von Heijne, and Krogh)



**Model TM length distribution.
Power of HMM.
Difficult in alignment.**



ALLYVDWQILPVIL

Weight matrix construction

SLLPAIVEL YLLPAIVHI TLWVDPYEV GLVPFLVSV KLLEPVILL LLDVPTAAV LLDVPTAAV LLDVPTAAV
LLDVPTAAV VLFRGGPRG MVDGTLLLL YMNGTMSQV MLLSVPLLL SLLGLLVEV ALLPPINIL TLIKIQHTL
HLIDYLVTS ILAPPVVKL ALFPQLVIL GILGFVFTL STNRQSGRQ GLDVLTAKV RILGAVAKV QVCERIPTI
ILFGHENRV ILMEHIHKL ILDQKINEV SLAGGIIGV LLIENVASL FLLWATAEA SLPDFGISY KKREEAPSL
LERPGGNEI ALSNLEVKL ALNELLQHV DLERKVESL FLGENISNF ALSDHHIYL GLSEFTEYL STAPPAHGV
PLDGEYFTL GVLVGVALI RTLDKVLEV HLSTAFARV RLDSYVRSL YMNGTMSQV GILGFVFTL ILKEPVHGV
ILGFVFVTLT LLFGYPVYV GLSPTVWLS WLSLLVPFV FLPSDFFPS CLGGLLTMV FIAGNSAYE KLGEFYNQM
KLVALGINA DLMGYIPLV RLVTLKDIV MLLAVLYCL AAGIGILTV YLEPGPVT A LLDGTATLR ITDQVPFSV
KTWGQYWQV TITDQVPFS AFHHVAREL YLNKIQNSL MMRKLAILS AIMDKNIIL IMDKNIIK SMVGNWAKV
SLLAPGAKQ KIFGSLAFL ELVSEFSRM KLTPLCVTL VLYRYGSFS YIGEVLVSV CINGVCWTV VMNILLQYV
ILTVELGVL KVLEYVIKV FLWGPRALV GLSRYVARL FLLTRILTI HLGNVKYLV GIAGGLALL GLQDCTMLV
TGAPVTYST VIYQYMDL VLPDVFIRC VLPDVFIRC AVGIGI^A VV LVVLLGLLAV ALGLGLLPV GIGIGVLA
GAGIGVAVL IAGIGILAI LIVIGILIL LAGIGLIAA VDGIGILTI GAGIGVITA AAGIGIIQI QAGIGILLA
KARDPHSGH KACDPHSGH ACDPHSGHF SLYNTVATL RGPGRAFVT NLVPMVATV GLHCYEQLV PLKQHFAQIV
AVFDRKSDA LLDVFVRFMG VLVKSPNHV GLAPPQHLLI LLGRNSFEV PLTFGWCYK VLEWRFDST TLNAWKVV
GLCTLVAML FIDSYICQV IISAVVGIL VMAGVGSPY LLWTLVVLL SVRDRLARI LLMDCSGSI CLTSTVQLV
VLHDDLLEA LMWITQCFL SLLMWITQC QLSLLMWIT LLGATCMFV RLTRFLSRV YMGTMSQV FLTPKKLQC
ISNDVCAQV VKTDGNPPE SYVDFFWL FLYGALLLA VLFSSDFRI LMWAKIGPV SLLLELEEV SLSRFWSGA
YTAFTIPS^I RLMKQDFSV RLPRIFCSC FLWGPRAYA RLLQETELV SLFEGIDFY SLDQSVVEL RLNMFTPYI
NMFTPYIGV LMIPIPLINV TLFIGSHVV SLVIVTTFV VLQWASLAV ILAKFLHWL STAPPHVN^V LLLLTVLTV
VVLGVVFGI ILHNGAYSL MIMVKCWMI MLGTHTMEV MLGTHTMEV SLADTNSLA LLWAARPRL GVALQTMKQ
GLYDGMEHL KMVELVHFL YLQLVFGIE MLMAQEALA LMAQEALAF VYDGREHTV YLSGANLNL RMFPNAPYL
EAAGIGILT TLDSQVMSL STPPPGRV KVAELVHFL IMIGVLGV^G ALCRWGLLL LLFAGVQCQ VLLCESTAV
YLSTAFARV YLLEMILWRL SLDDYNHLV RTLDKVLEV GLPVEYLQV KLIANNTRV FIYAGSLA KLVANNTRL
FLDEFMEGV ALQPGBTALL VLDGLDVLL SLYSFPEPE ALYVDSLFF SLLQHLIGL ELTLGEFLK MINAYLDKL
AAGIGILTV FLPSDFFPS SVRDRRLARL SLREWLLRI LLSAWILTA AAGIGILTV AVPDEIPPL FAYDGKD^{YI}
AAGIGILTV FLPSDFFPS AAGIGILTV FLPSDFFPS AAGIGILTV FLWGPRALV ETVSEQSNV ITLWQRPLV

PSSM construction

- Calculate amino acid frequencies at each position using
 - Sequence weighting
 - Pseudo counts
- Define background model
 - Use background amino acids frequencies
- PSSM is

$$S(a_i) = \log \frac{p(a_i)}{q(a)}$$

More on scoring

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
1	0.6	0.4	-3.5	-2.4	-0.4	-1.9	-2.7	0.3	-1.1	1.0	0.3	0.0	1.4	1.2	-2.7	1.4	-1.2	-2.0	1.1	0.7
2	-1.6	-6.6	-6.5	-5.4	-2.5	-4.0	-4.7	-3.7	-6.3	1.0	5.1	-3.7	3.1	-4.2	-4.3	-4.2	-0.2	-5.9	-3.8	0.4
3	0.2	-1.3	0.1	1.5	0.0	-1.8	-3.3	0.4	0.5	-1.0	0.3	-2.5	1.2	1.0	-0.1	-0.3	-0.5	3.4	1.6	0.0
4	-0.1	-0.1	-2.0	2.0	-1.6	0.5	0.8	2.0	-3.3	0.1	-1.7	-1.0	-2.2	-1.6	1.7	-0.6	-0.2	1.3	-6.8	-0.7
5	-1.6	-0.1	0.1	-2.2	-1.2	0.4	-0.5	1.9	1.2	-2.2	-0.5	-1.3	-2.2	1.7	1.2	-2.5	-0.1	1.7	1.5	1.0
6	-0.7	-1.4	-1.0	-2.3	1.1	-1.3	-1.4	-0.2	-1.0	1.8	0.8	-1.9	0.2	1.0	-0.4	-0.6	0.4	-0.5	-0.0	2.1
7	1.1	-3.8	-0.2	-1.3	1.3	-0.3	-1.3	-1.4	2.1	0.6	0.7	-5.0	1.1	0.9	1.3	-0.5	-0.9	2.9	-0.4	0.5
8	-2.2	1.0	-0.8	-2.9	-1.4	0.4	0.1	-0.4	0.2	-0.0	1.1	-0.5	-0.5	0.7	-0.3	0.8	0.8	-0.7	1.3	-1.1
9	-0.2	-3.5	-6.1	-4.5	0.7	-0.8	-2.5	-4.0	-2.6	0.9	2.8	-3.0	-1.8	-1.4	-6.2	-1.9	-1.6	-4.9	-1.6	4.5

$$S = \sum_i S(a_i)$$

$$S = \sum_i \log \frac{p(a_i)}{q(a_i)}$$

$$S = \log \left(\frac{\prod_i p(a_i)}{\prod_i q(a_i)} \right)$$

Probability of observation given Model

Probability of observation given Prior
(background)

$$S = \log \left(\frac{P(a | M)}{P(a | B)} \right)$$

Hidden Markov Models

- Weight matrices do not deal with insertions and deletions
 - In alignments, this is done in an ad-hoc manner by optimization of the two gap penalties for first gap and gap extension
 - HMM is a natural framework where insertions/ deletions are dealt with explicitly
-

Multiple sequence alignment

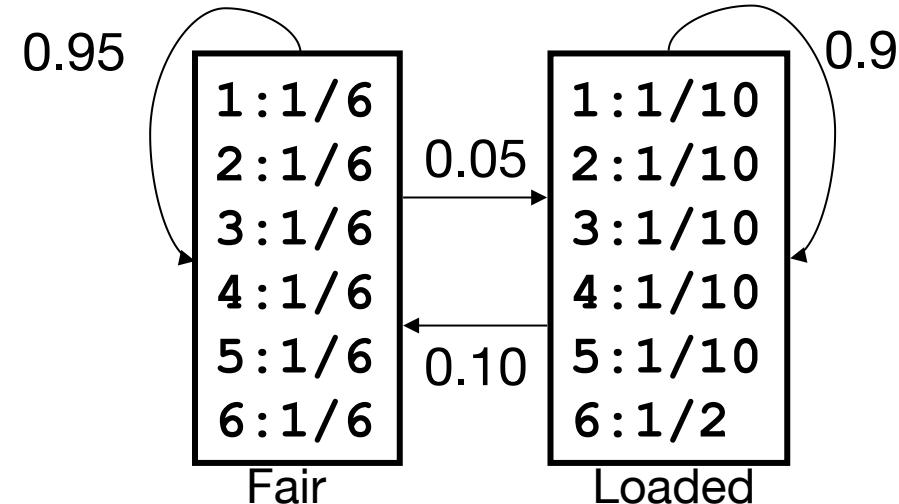
Learning from evolution



Why hidden?

The unfair casino: Loaded die $p(6) = 0.5$; switch fair to load:0.05; switch load to fair: 0.1

- Model generates numbers
 - 312453666641
- Does not tell which die was used
- Alignment (decoding) can give the most probable solution/path (Viterbi)
 - FFFFFFFLLLLL
- Or most probable set of states
 - FFFFFFFLLLLL



HMM (a simple example)

ACA---ATG

TCAACTATC

ACACC--AGC

AGA---ATC

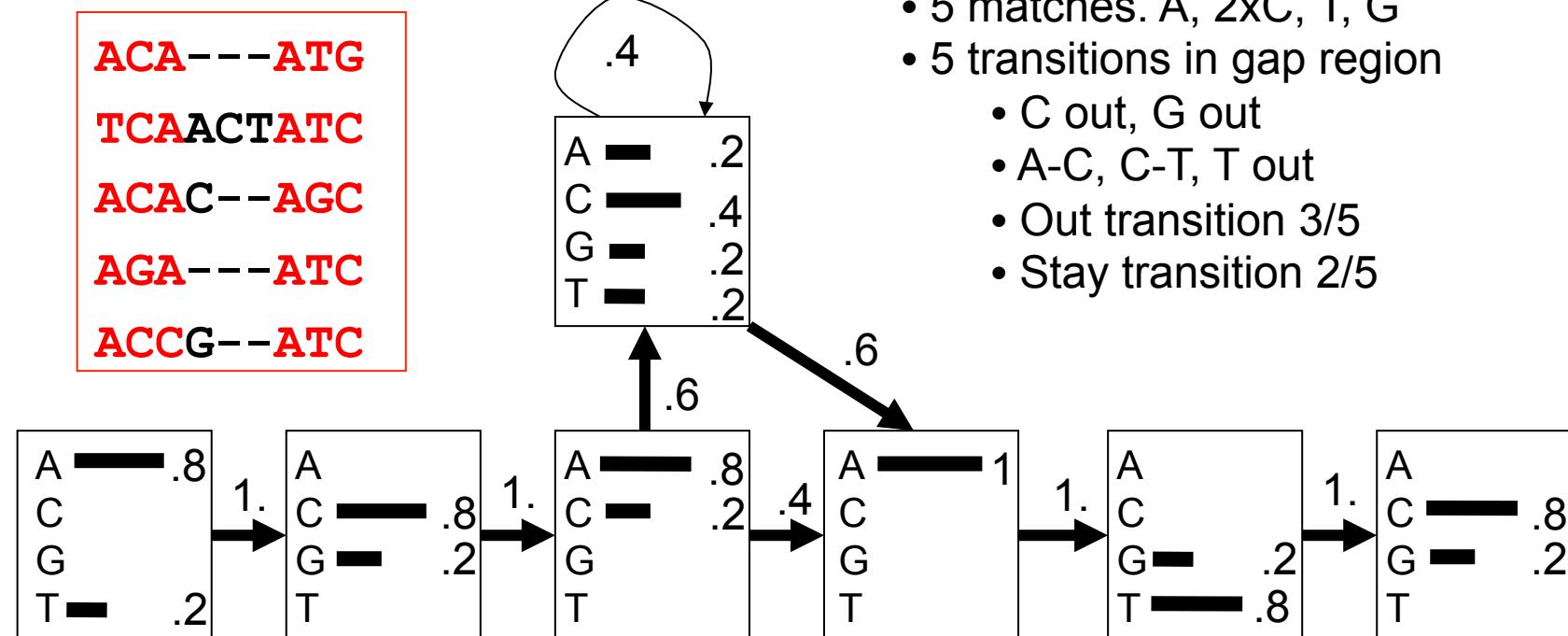
ACCG--ATC



Core of alignment

- Example from A. Krogh
- Core region defines the number of states in the HMM (**red**)
- Insertion and deletion statistics are derived from the non-core part of the alignment (black)

HMM construction (supervised learning)



$$\text{ACA---ATG} \quad 0.8 \times 1 \times 0.8 \times 1 \times 0.8 \times 0.4 \times 1 \times 1 \times 0.8 \times 1 \times 0.2 = 3.3 \times 10^{-2}$$

Scoring a sequence to an HMM

ACA---ATG $0.8 \times 1 \times 0.8 \times 1 \times 0.8 \times 0.4 \times 1 \times 0.8 \times 1 \times 0.2 = 3.3 \times 10^{-2}$

TCAACTATC $0.2 \times 1 \times 0.8 \times 1 \times 0.8 \times 0.6 \times 0.2 \times 0.4 \times 0.4 \times 0.2 \times 0.6 \times 1 \times 1 \times 0.8 \times 1 \times 0.8 = 0.0075 \times 10^{-2}$

ACAC--AGC = 1.2×10^{-2}

Consensus:

ACAC--ATC = 4.7×10^{-2} , **ACA---ATC** = 13.1×10^{-2}

Exceptional:

TGCT--AGG = 0.0023×10^{-2}

Align sequence to HMM - Null model

- Score depends **strongly** on length
 - Null model is a random model. For length L the score is 0.25^L
 - Log-odds score for sequence S
 - $\log(P(S)/0.25^L)$
 - Positive score means more likely than Null model
- This is just like we did for PSSM $\log(p/q)$!

ACA---ATG = 4.9

TCAACTATC = 3.0

ACAC--AGC = 5.3

AGA---ATC = 4.9

ACCG--ATC = 4.6

Consensus:

ACAC--ATC = 6.7

ACA---ATC = 6.3

Exceptional:

TGCT--AGG = -0.97

Note!

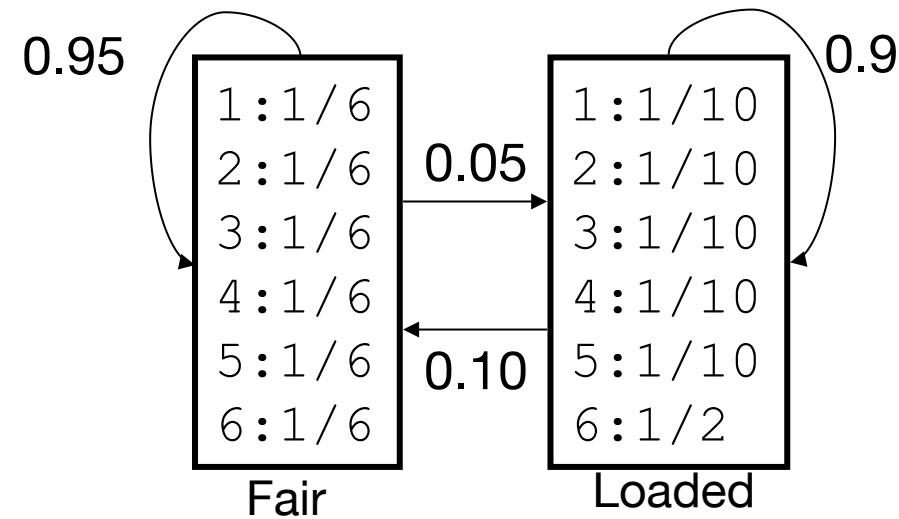
Aligning a sequence to an HMM

- Find the path through the HMM states that has the highest probability
 - For alignment, we found the path through the scoring matrix that had the highest sum of scores
 - The number of possible paths rapidly gets very large making brute force search infeasible
 - Just like checking all path for alignment did not work
 - Use dynamic programming
 - The Viterbi algorithm does the job
-

The Viterbi algorithm

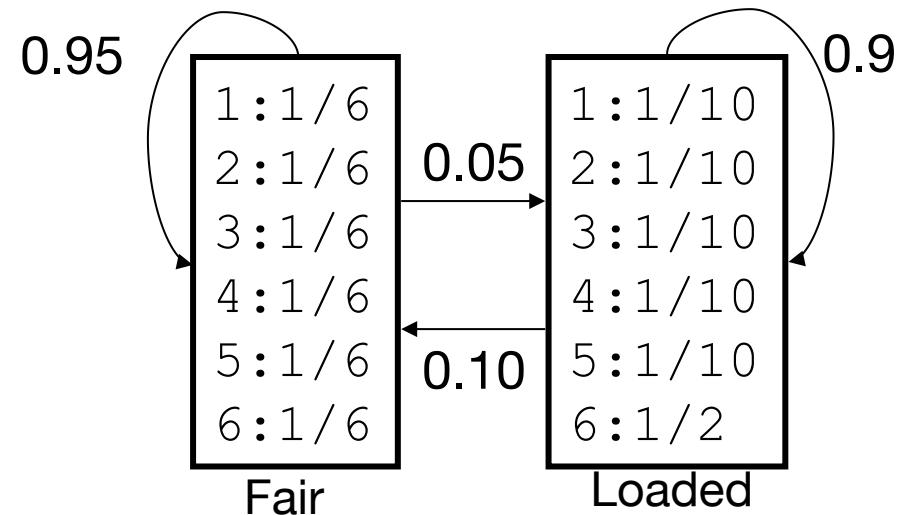
The unfair casino: Loaded dice $p(6) = 0.5$; switch fair to load:0.05; switch load to fair: 0.1

- Model generates numbers
 - 312453666641



Model decoding (Viterbi)

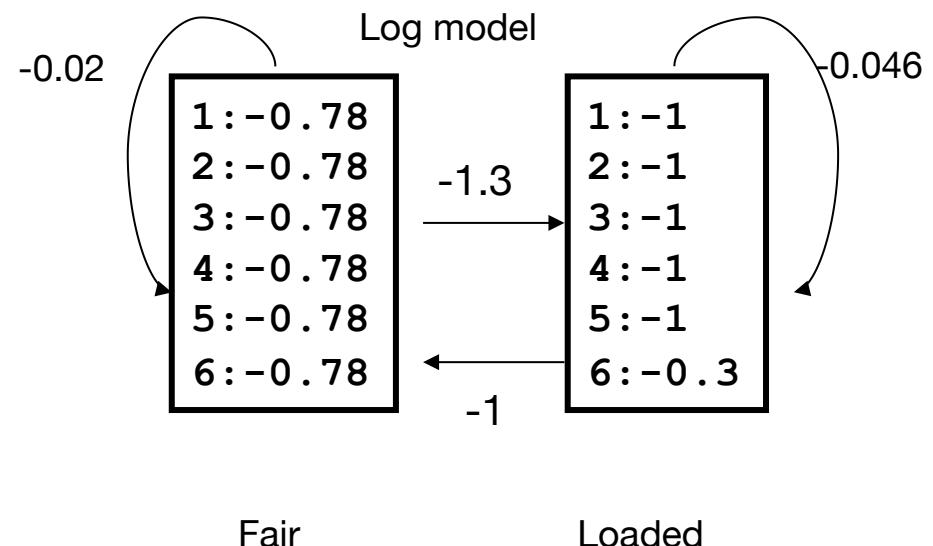
- Example: 566. What was the most likely series of dice used to generate this output?
- Use Brute force



```
FFF = 0.5*0.167*0.95*0.167*0.95*0.167 = 0.0021
FFL = 0.5*0.167*0.95*0.167*0.05*0.5 = 0.00333
FLF = 0.5*0.167*0.05*0.5*0.1*0.167 = 0.000035
FLL = 0.5*0.167*0.05*0.5*0.9*0.5 = 0.00094
LFF = 0.5*0.1*0.1*0.167*0.95*0.167 = 0.00013
LFL = 0.5*0.1*0.1*0.167*0.05*0.5 = 0.000021
LLF = 0.5*0.1*0.9*0.5*0.1*0.167 = 0.00038
LLL = 0.5*0.1*0.9*0.5*0.9*0.5 = 0.0101
```

Or in log space

- Example: 566. What was the most likely series of dice used to generate this output?



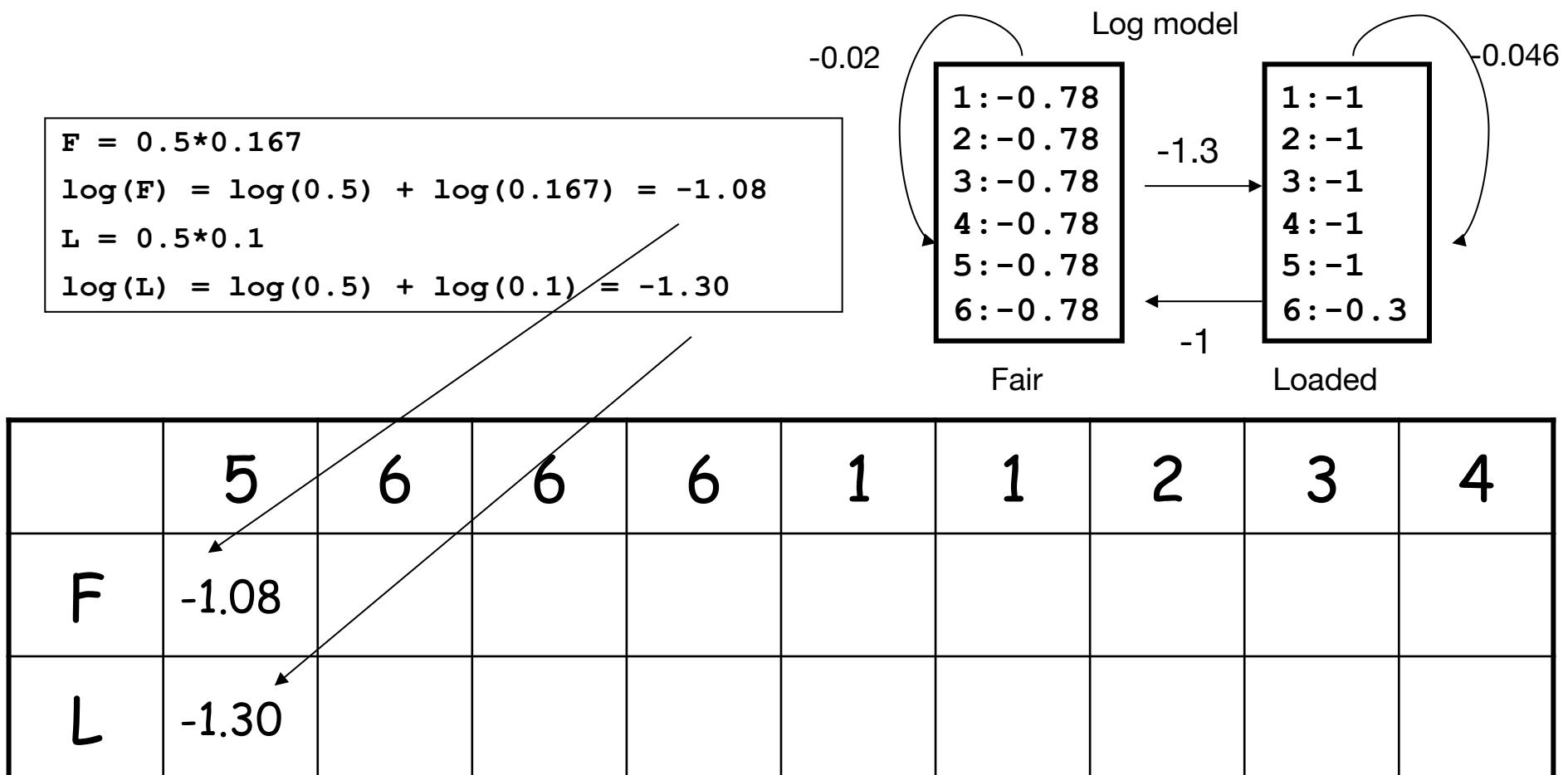
$$\text{Log}(P(\text{LLL}|M)) = \log(0.5*0.1*0.9*0.5*0.9*0.5) = \log(0.0101)$$

or

$$\begin{aligned}\text{Log}(P(\text{LLL}|M)) &= \log(0.5)+\log(0.1)+\log(0.9)+\log(0.5)+\log(0.9)+\log(0.5) \\ &= -0.3 -1 -0.046 -0.3 -0.046 -0.3 = -1.99\end{aligned}$$

Model decoding (Viterbi)

- Example: 566611234. What was the most likely series of dice used to generate this output?



Model decoding (Viterbi)

- Example: 566611234. What was the most likely series of dice used to generate this output?

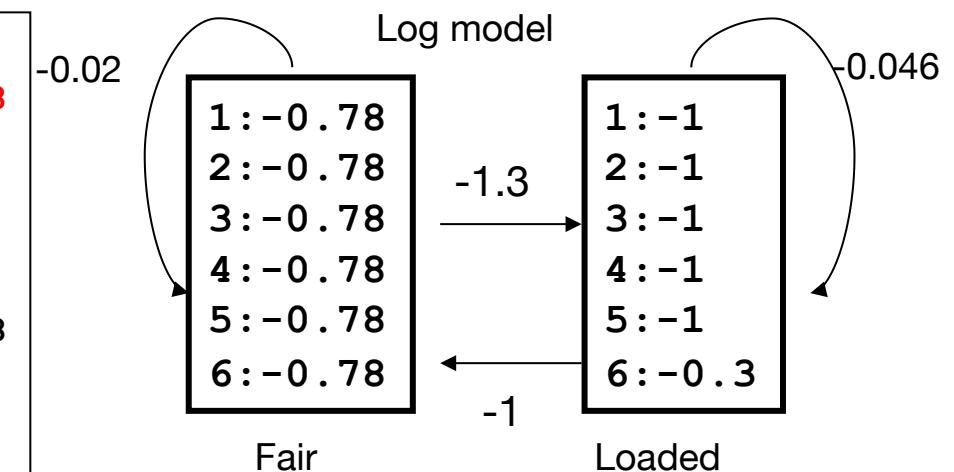
```

FF = 0.5*0.167*0.95*0.167
Log(FF) = -0.30 -0.78 - 0.02 -0.78 = -1.88

LF = 0.5*0.1*0.1*0.167
Log(LF) = -0.30 -1 -1 -0.78 = -3.08

FL = 0.5*0.167*0.05*0.5
Log(FL) = -0.30 -0.78 - 1.30 -0.30 = -2.68

LL = 0.5*0.1*0.9*0.5
Log(LL) = -0.30 -1 -0.046 -0.3 = -1.65
  
```



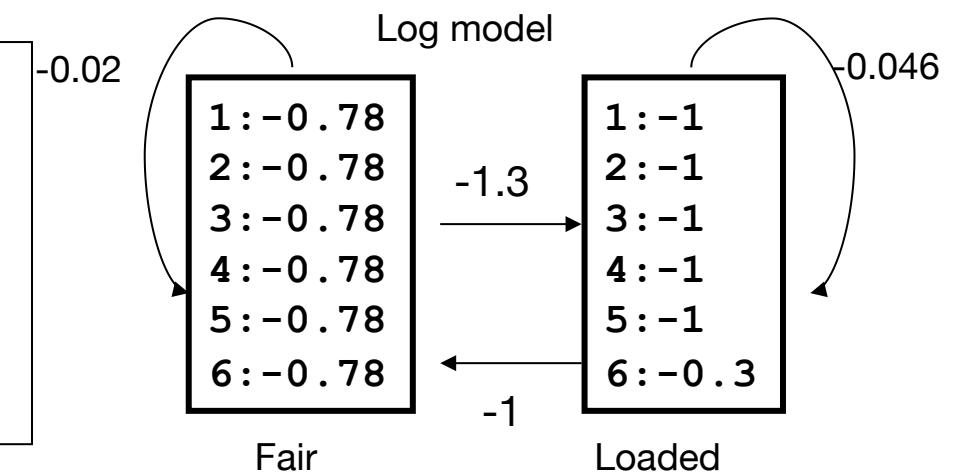
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88							
L	-1.30	-1.65							

Model decoding (Viterbi)

- Example: 566611234. What was the most likely series of dice used to generate this output?

```

FFF = 0.5*0.167*0.95*0.167*0.95*0.167 = 0.0021
FLF = 0.5*0.167*0.05*0.5*0.1*0.167 = 0.000035
LFF = 0.5*0.1*0.1*0.167*0.95*0.167 = 0.00013
LLF = 0.5*0.1*0.9*0.5*0.1*0.167 = 0.00038
FLL = 0.5*0.167*0.05*0.5*0.9*0.5 = 0.00094
FFL = 0.5*0.167*0.95*0.167*0.05*0.5 = 0.00333
LFL = 0.5*0.1*0.1*0.167*0.05*0.5 = 0.000021
LLL = 0.5*0.1*0.9*0.5*0.9*0.5 = 0.0101
  
```

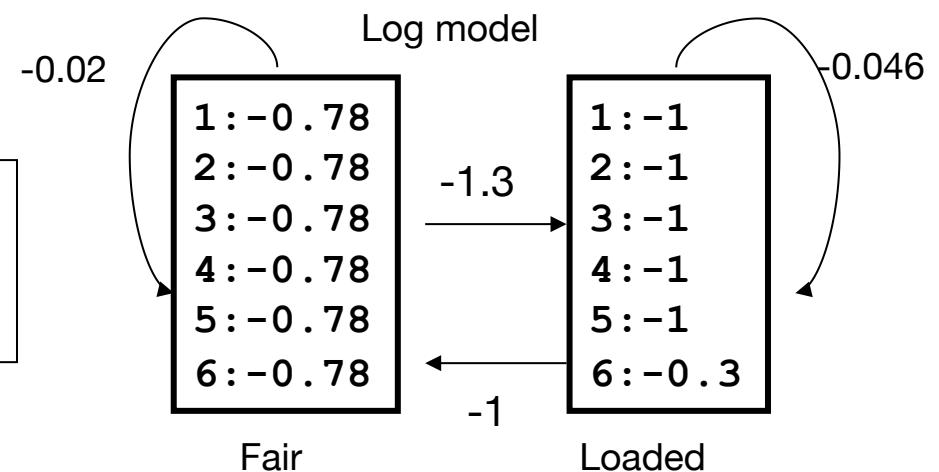


	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88							
L	-1.30	-1.65							

Model decoding (Viterbi)

- Example: 566611234. What was the most likely series of dice used to generate this output?

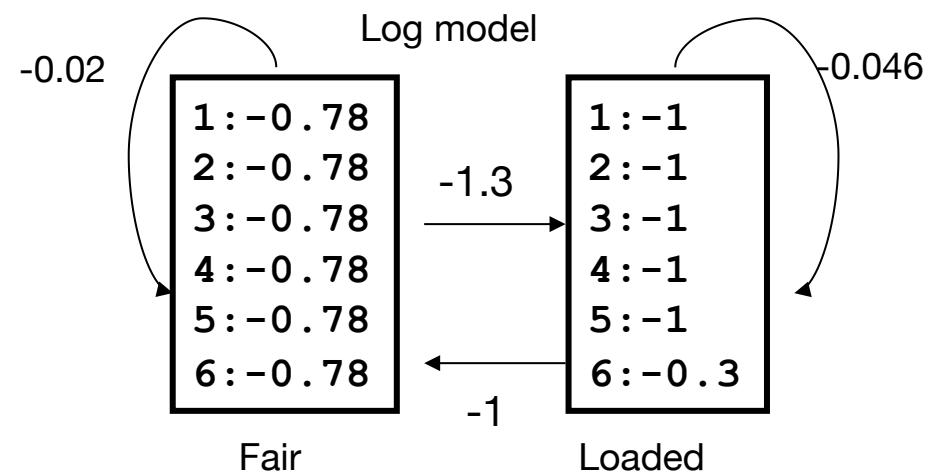
$\text{FFF} = 0.5 * 0.167 * 0.95 * 0.167 * 0.95 * 0.167 = 0.0021$
 $\text{Log}(P(\text{FFF})) = -2.68$
 $\text{LLL} = 0.5 * 0.1 * 0.9 * 0.5 * 0.9 * 0.5 = 0.0101$
 $\text{Log}(P(\text{LLL})) = -1.99$



	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68						
L	-1.30	-1.65	-1.99						

Model decoding (Viterbi)

- Example: 566611234. What was the most likely series of dice used to generate this output?



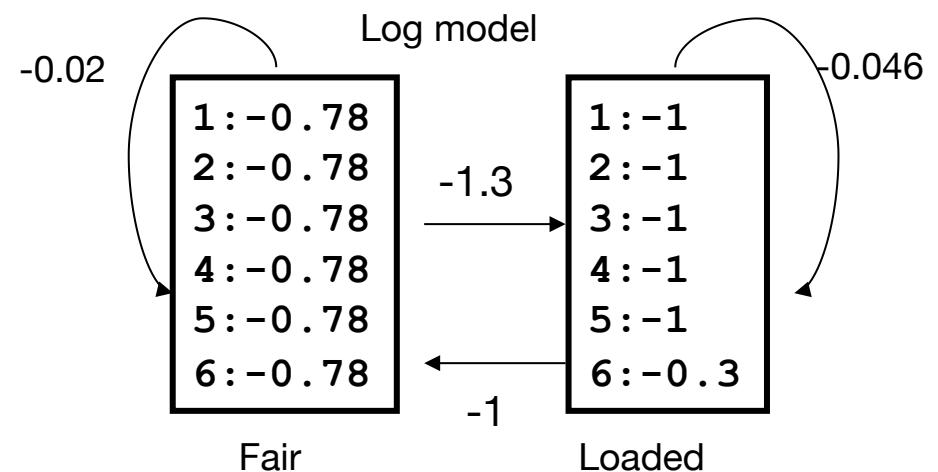
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68						
L	-1.30	-1.65	-1.99						

Model decoding (Viterbi)

- Example: 566611234. What was the most likely series of dice used to generate this output?

$$-0.78 - 0.02 - 2.68 = -3.48$$

$$-0.78 - 1 - 1.99 = -3.77$$



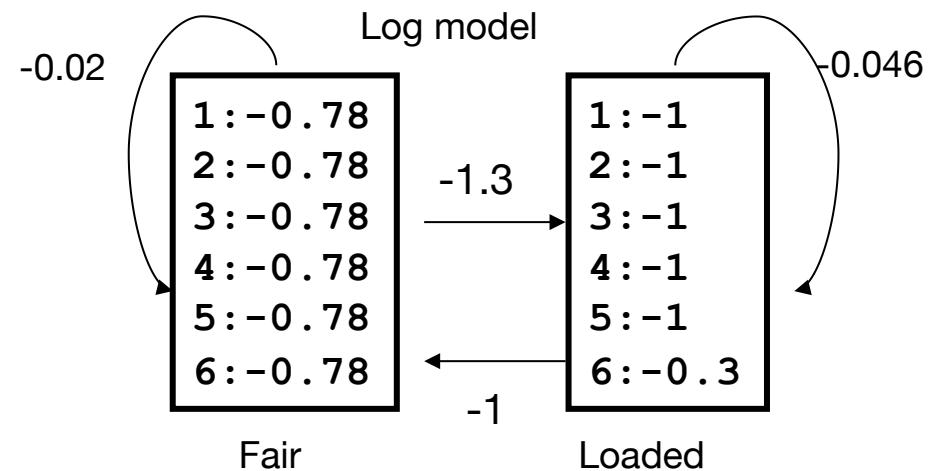
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-					
L	-1.30	-1.65	-1.99						

Model decoding (Viterbi)

- Example: 566611234. What was the most likely series of dice used to generate this output?

$$\boxed{-0.78 - 0.02 - 2.68 = -3.48}$$

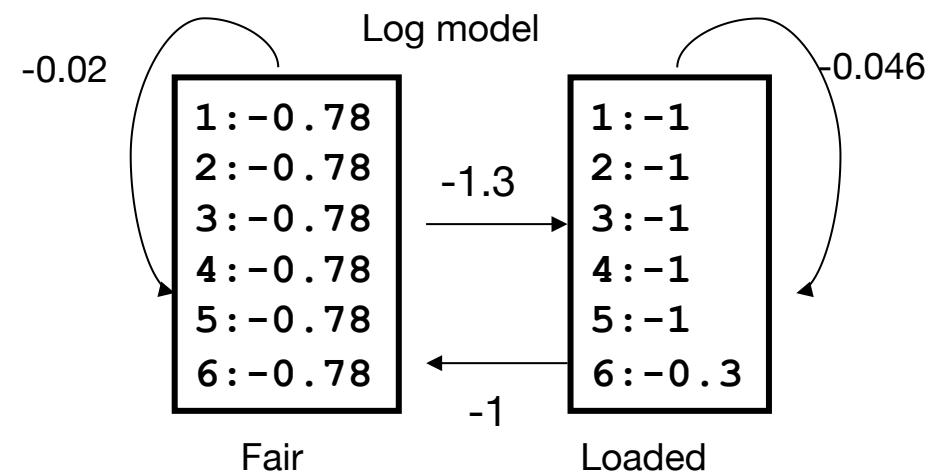
$$\boxed{-0.78 - 1 - 1.99 = -3.77}$$



	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-3.48					
L	-1.30	-1.65	-1.99						

Model decoding (Viterbi)

- Now we can formalize the algorithm!



$$P_l(i+1) = p_l(i+1) \cdot \max_k (P_k(i) \cdot a_{kl}) \quad or$$

$$\log(P_l(i+1)) = \log(p_l(i+1)) + \max_k (\log(P_k(i)) + \log(a_{kl}))$$

New match

Old max score

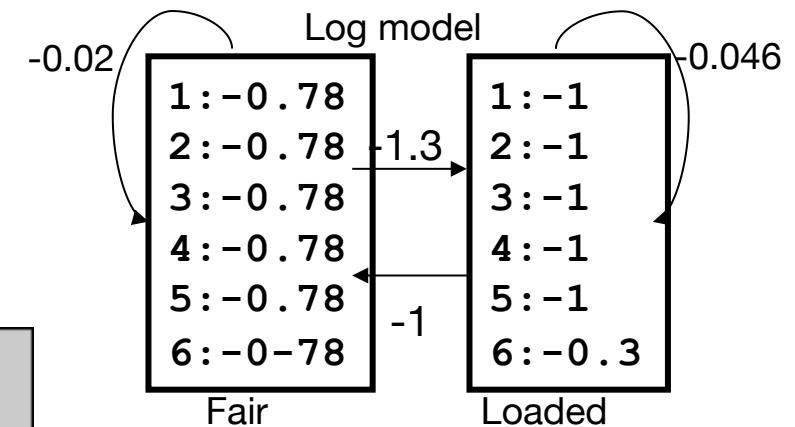
Transition

Model decoding (Viterbi). Can you do it?

- Example: 566611234. What was the most likely series of dice used to generate this output?
- Fill out the table using the Viterbi recursive algorithm
 - Add the arrows for backtracking
- Find the optimal path

$$P_l(i+1) = p_l(i+1) \cdot \max_k (P_k(i) \cdot a_{kl}) \quad \text{or}$$

$$\log(P_l(i+1)) = \log(p_l(i+1)) + \max_k (\log(P_k(i)) + \log(a_{kl}))$$



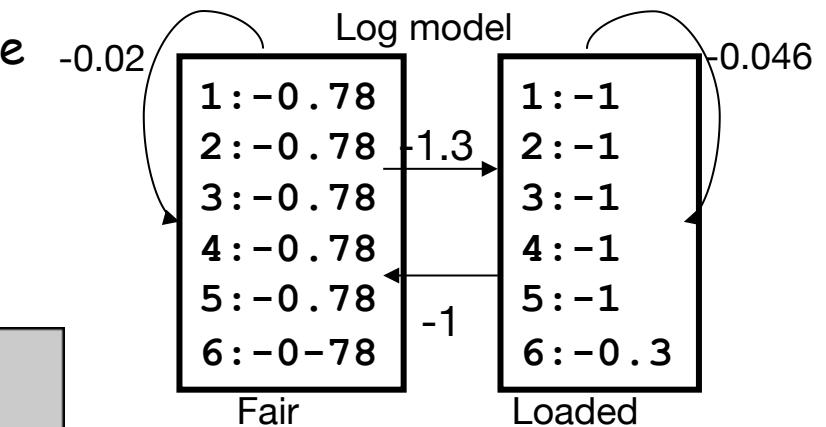
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-3.48					
L	-1.30	-1.65	-1.99						

Model decoding (Viterbi). Can you do it?

- Example: 566611234. What was the most likely series of dice used to generate this output?
- Fill out the table using the Viterbi recursive algorithm
 - Add the arrows for backtracking
- Find the optimal path

$$P_l(i+1) = p_l(i+1) \cdot \max_k (P_k(i) \cdot a_{kl}) \quad \text{or}$$

$$\log(P_l(i+1)) = \log(p_l(i+1)) + \max_k (\log(P_k(i)) + \log(a_{kl}))$$



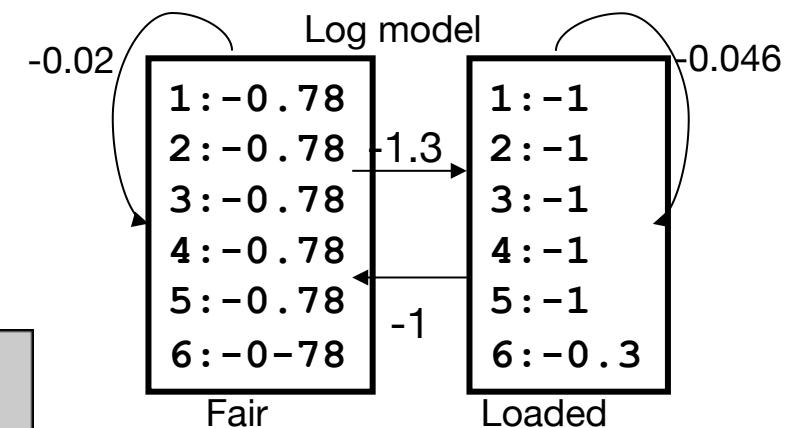
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-3.48		-4.92		-6.53	
L	-1.30	-1.65	-1.99		-3.39			-6.52	

Model decoding (Viterbi). Can you do it?

- Example: 566611234. What was the most likely series of dice used to generate this output?
- Fill out the table using the Viterbi recursive algorithm
 - Add the arrows for backtracking
- Find the optimal path

$$P_l(i+1) = p_l(i+1) \cdot \max_k (P_k(i) \cdot a_{kl}) \quad \text{or}$$

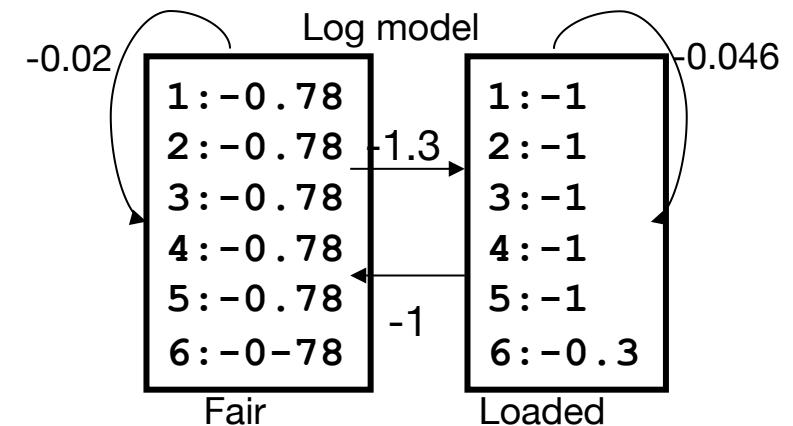
$$\log(P_l(i+1)) = \log(p_l(i+1)) + \max_k (\log(P_k(i)) + \log(a_{kl}))$$



	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-3.48	-4.12	-4.92	5.73	-6.53	-7.33
L	-1.30	-1.65	-1.99	-2.34	-3.39	-4.44	-5.48	-6.52	-7.57

Model decoding (Viterbi). Can you do it?

- Example: 566611234. What was the most likely series of dice used to generate this output?
- The most likely path is
 - LLLLFFFFFF



	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-3.48	-4.12	-4.92	5.73	-6.53	-7.33
L	-1.30	-1.65	-1.99	-2.34	-3.39	-4.32	-5.48	-6.52	-7.57

Model decoding (Viterbi).

- What happens if you have three dice?

	5	6	6	6	1	1	2	3	4
F	-1.0								
L1	-1.2								
L2	-1.3								

$$P_l(i+1) = p_l(i+1) \cdot \max_k (P_k(i) \cdot a_{kl}) \quad or$$
$$\log(P_l(i+1)) = \log(p_l(i+1)) + \max_k (\log(P_k(i)) + \log(a_{kl}))$$

And if you have a trans-membrane model

- What is the most likely path (alignment) of a protein sequence to the model

	D	G	V	L	I	M	A	D	Q
iC	-1.0								
M	-1.2								
xC	-1.3								

$$P_l(i+1) = p_l(i+1) \cdot \max_k (P_k(i) \cdot a_{kl}) \quad or$$
$$\log(P_l(i+1)) = \log(p_l(i+1)) + \max_k (\log(P_k(i)) + \log(a_{kl}))$$

The Forward algorithm

- The Viterbi algorithm finds the most probable path giving rise to a given sequence
- One other interesting question would be
 - What is the probability that a given sequence can be generated by the hidden Markov model
 - Calculated by summing over all path giving rise to a given sequence

The Forward algorithm

- Calculate summed probability over all path giving rise to a given sequence

$$P(x) = \sum_{\pi} P(x, \pi)$$

- The number of possible paths is very large making (once more) brute force calculations infeasible
 - Use dynamic (recursive) programming

The Posterior decoding (Backward algorithm)

- One other interesting question would be
 - What is the probability that an observation x_i came from a state k given the observed sequence x

or

- What is the probability that a given amino acids is part of the trans-membrane helix given the protein sequence is X ?

$$P(\pi_i = k \mid x)$$

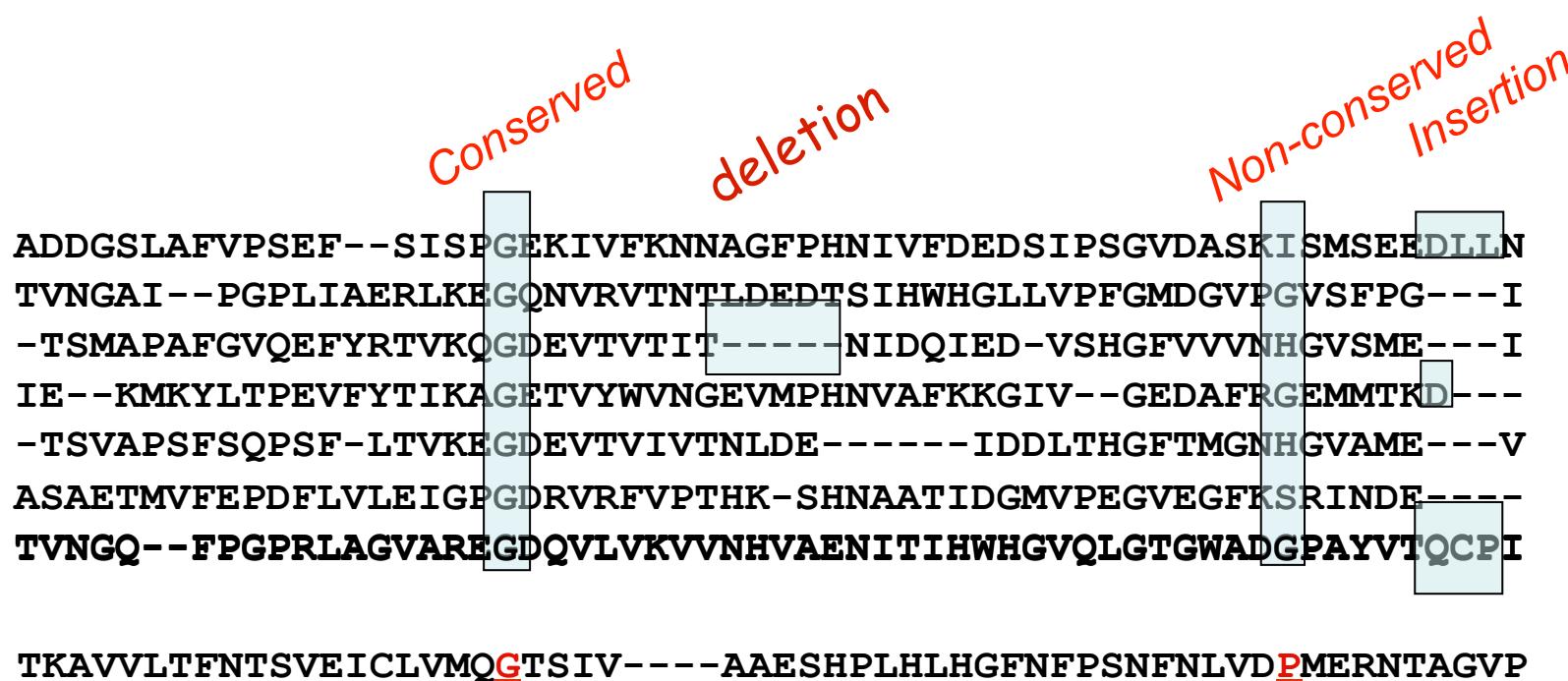
HMM's and weight matrices

- In the case of un-gapped alignments HMM's become simple weight matrices
- To achieve high performance, the emission frequencies are estimated using the techniques of
 - Sequence weighting
 - Pseudo counts

Profiles and profile HMM's

- Alignments based on conventional scoring matrices (BLOSUM62) scores all positions in a sequence in an equal manner
 - Some positions are highly conserved, some are highly variable (more than what is described in the BLOSUM matrix)
 - Profile HMM's are ideal suited to describe such position specific variations
-

Sequence profiles



Matching any thing
but $G \Rightarrow$ large
negative score

Any thing can match

HMM vs. alignment

- Detailed description of core
 - Conserved/variable positions
 - Price for insertions/deletions varies at different locations in sequence
 - These features cannot be captured in conventional alignments
-