

Sequence profiles, Hidden Markov models and homology modeling

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

Identification of essential residues in protein sequences

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

TVNGQ--FPGPRLAGVAREGDQVLVKVVNHVAENITIHWGVQLGTGWADPPAYVTQCPIT
TKAVVLTFTNTSVEICLVMQGTSIV----AAESHPLHLHGFPNSNFNLVDGMERNTAGVP

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	-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

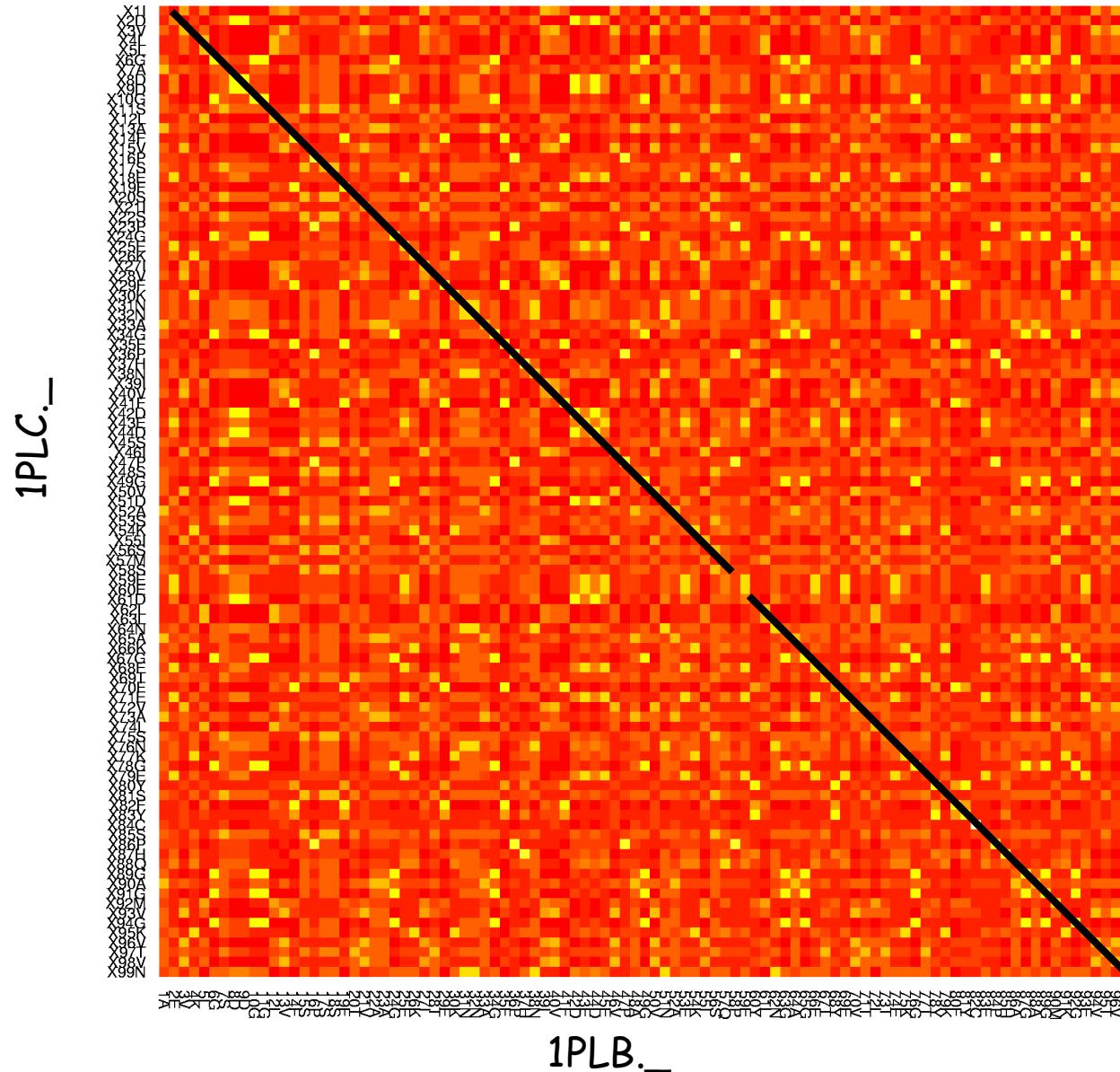
A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-2	-1	1	0	-3	-2	0	
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2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

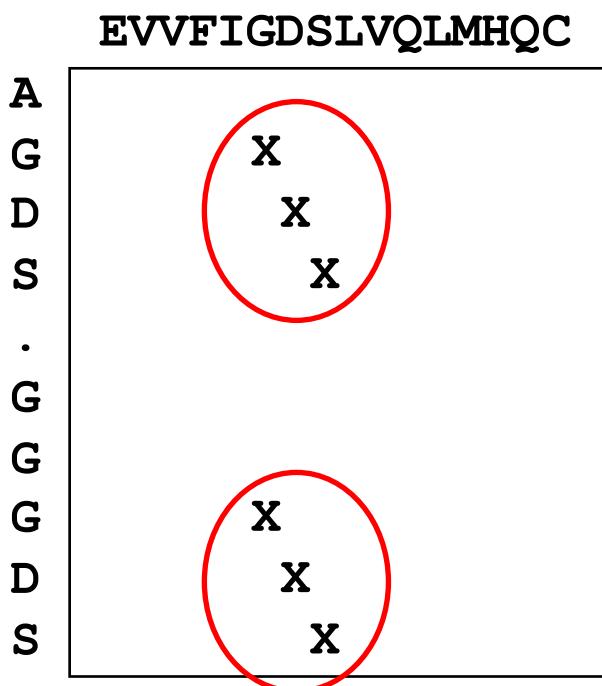
I-GDS

When Blast works!



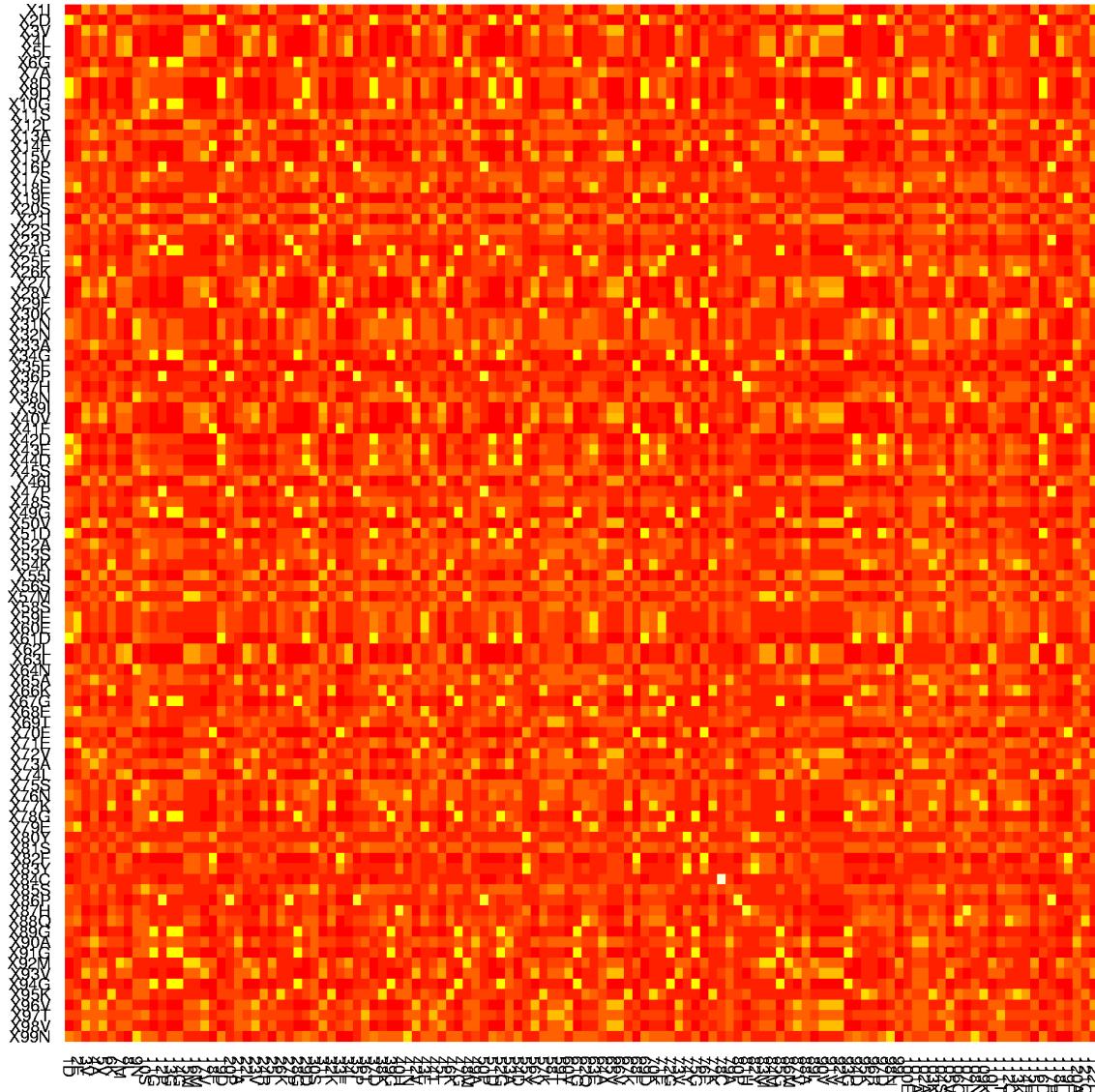
What goes wrong when Blast fails?

- 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!

1PLC._



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	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--FPGPRLAGVAREGDQVLVKVVNHVAENITIHWHGVQLGTGWADGPAYVTQCPI

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

Conserved Non-conserved

ADDGSLAFVPSEF--SISPGEKIVFKNNAGFPNIVFDED S I P S G V D A S K I S M S E E D L L N
TVNGAI--PGPLIAERLKEGQNVRVTNTLDEDTSIHWHGLLVPFGMDGVPGVS F P G --- I
-TSMAPAFGVQE FYRTVKQGD EVTVTIT----NIDQIED-VSHGFVVVN HGV S M E --- I
IE---KMKYLTP E VFYTIKAGE TVYWVN G E V M P H N V A F K K G I V -- G E D A F R G E M M T K D ---
-TSVAPSFSQPSF-LTVKEGDEVTVI VTNLDE-----IDDLTHGFTMGNHGVAME---V
ASAETMVFE PDFLVLEIGPGDRVRFVPTHK-SHNAATIDGMVPEGVEGFKSRINDE---
TVNGQ--FPGPRLAGVAREGDQQLVKVVNHVAENITIHWGVQLGTGWADGPAYVTQCPITKAVVLTFTNTSVEICLVMQGTSIV----AAESHPLHLHGFNFPSNFNLVDPMERNTAGVP

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

Any thing can match

How to make sequence profiles

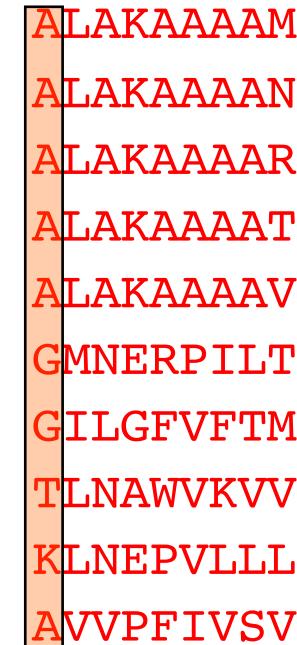
1. Align (BLAST) sequence against large sequence database (Swiss-Prot)
 2. Select significant alignments and make sequence profile
 3. Use profile to align against sequence database to find new significant hits
 4. Repeat 2 and 3 (normally 3 times!)
-

Sequence logos. Visualization of sequence profiles

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

$$\begin{aligned} P_A &= 6/10 = 0.6 \\ P_G &= 2/10 = 0.2 \\ P_T &= P_K = 1/10 = 0.1 \\ P_C &= P_D = \dots P_V = 0.0 \end{aligned}$$

$$\begin{aligned} q_A &= 0.07 \\ q_G &= 0.07 \\ q_T &= 0.05 \\ q_K &= 0.06 \end{aligned}$$

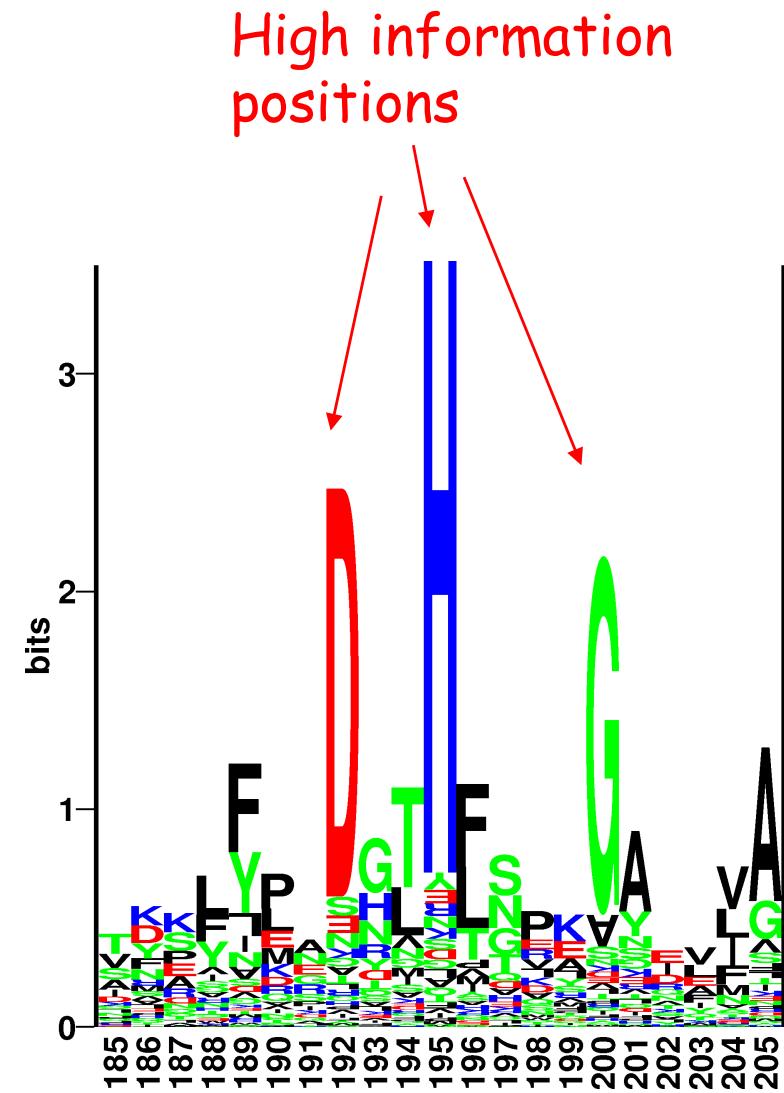


ALAKAAAAM
ALAKAAAAN
ALAKAAAAR
ALAKAAAAT
ALAKAAAAV
GMNERPILT
GILGFVFTM
TLNAWVKVV
KLNEPVLLL
AVVPFIVSV

Sequence logos

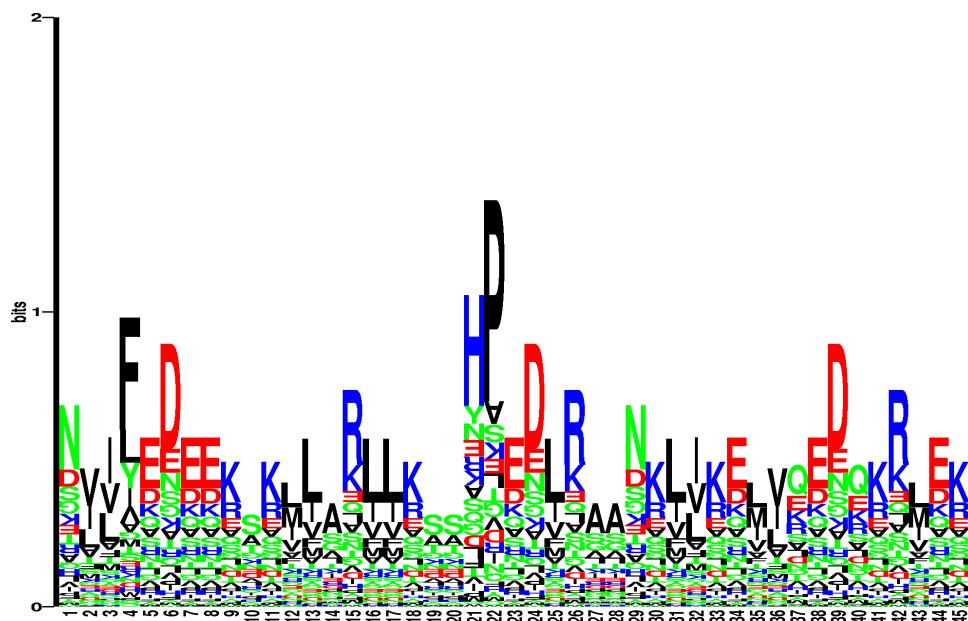
$$I = \sum_a p_a \log\left(\frac{p_a}{q_a}\right)$$

- 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)

IFEDEEKS**K**MLARLLKSSHPEDLRAA**N**KLIKELVQED**Q**KRLEK
Blosum62

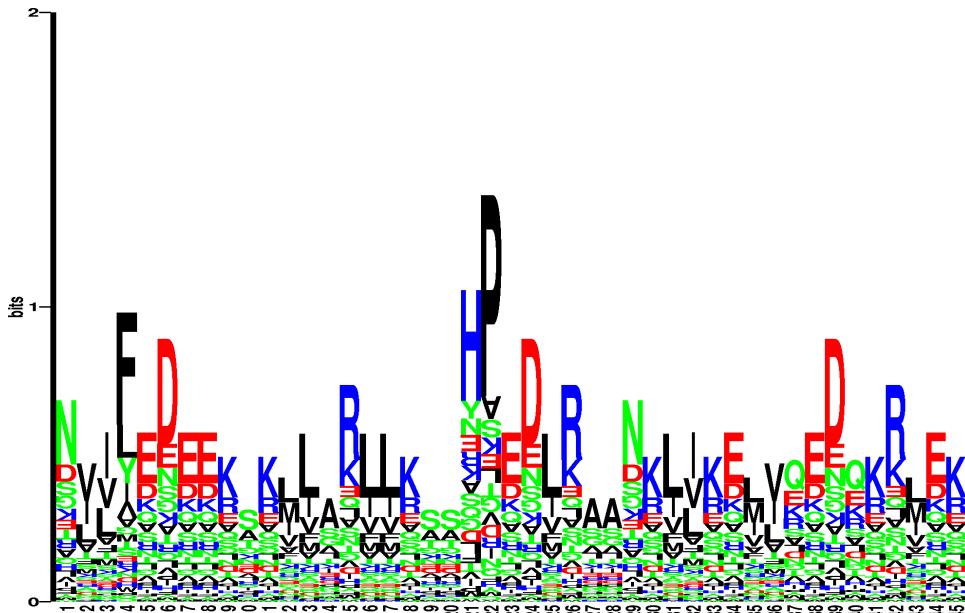


	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2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

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

Sequence profiles (1J2J.B)

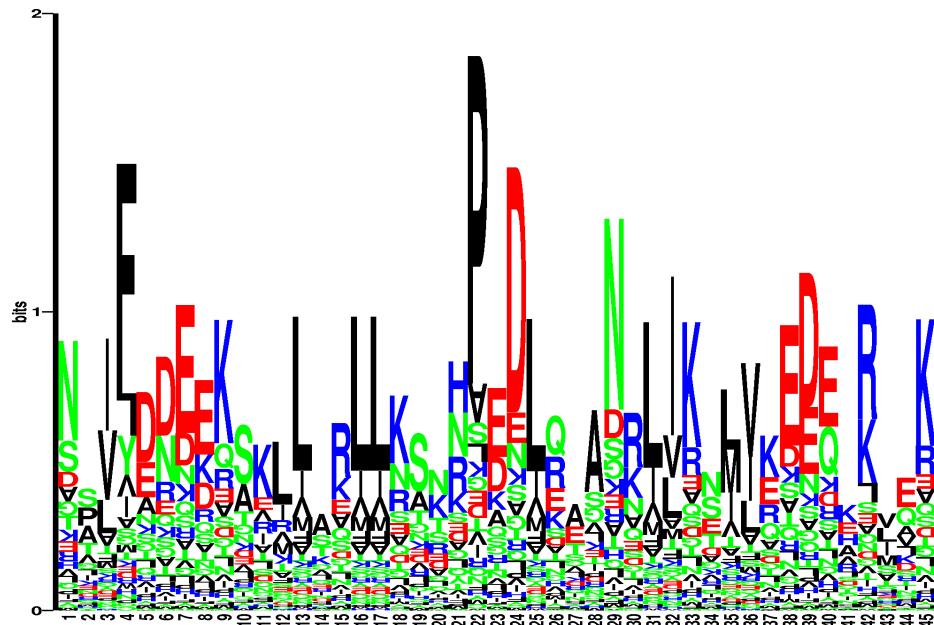
IFEDEEKS**K**MLARLLKSSHPEDLRAA**N**KL**I**KELVQ**E**D**Q**KR**L**E**K**



	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	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	-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)

IFEDEEKS**KMLARLLKSSHPEDLRAAANKLIKELVQEDQKRLEK**



	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	3
2	F	-3	-3	-3	-4	-3	-4	-4	-3	-1	0	0	-3	0	7	-4	-3	-2	1	3	-1
3	E	-2	-1	1	5	-4	1	4	-2	-1	-4	-4	0	-3	-4	-2	0	-1	-4	-3	-3
4	D	-2	-2	1	6	-4	-1	1	-2	-1	-3	-4	-1	-3	-4	-2	0	-1	-5	-3	-3
5	E	-1	0	0	1	-4	2	5	-2	0	-4	-3	1	-2	-4	-1	0	-1	-3	-2	-3
6	E	-1	0	0	1	-4	2	5	-2	0	-4	-3	1	-2	-4	-1	0	-1	-3	-2	-3
7	K	-1	2	0	-1	-3	1	1	-2	-1	-3	-3	5	-2	-3	-1	0	-1	-3	-2	-3
8	S	1	-1	0	0	-1	0	0	0	-1	-3	-3	0	-2	-3	-1	5	1	-3	-2	-2
9	K	-1	2	0	-1	-3	1	1	-2	-1	-3	-3	5	-2	-3	-1	0	-1	-3	-2	-3
10	M	-1	-2	-3	-4	-2	-1	-3	-4	-2	1	3	-2	5	0	-3	-2	-1	-2	-1	1
11	L	-2	-2	-4	-4	-2	-2	-3	-4	-3	1	4	-3	2	0	-3	-3	-1	-2	-1	1
12	A	4	-2	-1	-2	-1	-1	-1	0	-2	-1	-2	-1	-1	-2	-1	1	2	-3	-2	0

Blast2logo

Blast2logo 1.0 Server

[Instructions](#)

[Output format](#)

SUBMISSION

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

```
>Ex
VALAELYIPEVARRLGQGWHEDECTFAEVТИCTARLQAILRDIATWSADEGGMRDGPALVLLPPG
EQHTLGMAMAVAKLRLGVSVCLMSTGPAELRELFGKRRFDAIMISLAHAEMLEVGRKLVKTLKD
MTGGRIPVAMGGALFLDGTEASIPADIVTNDIEALQ
```

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 SP

Number of Blast iterations

Blast E-value cutoff

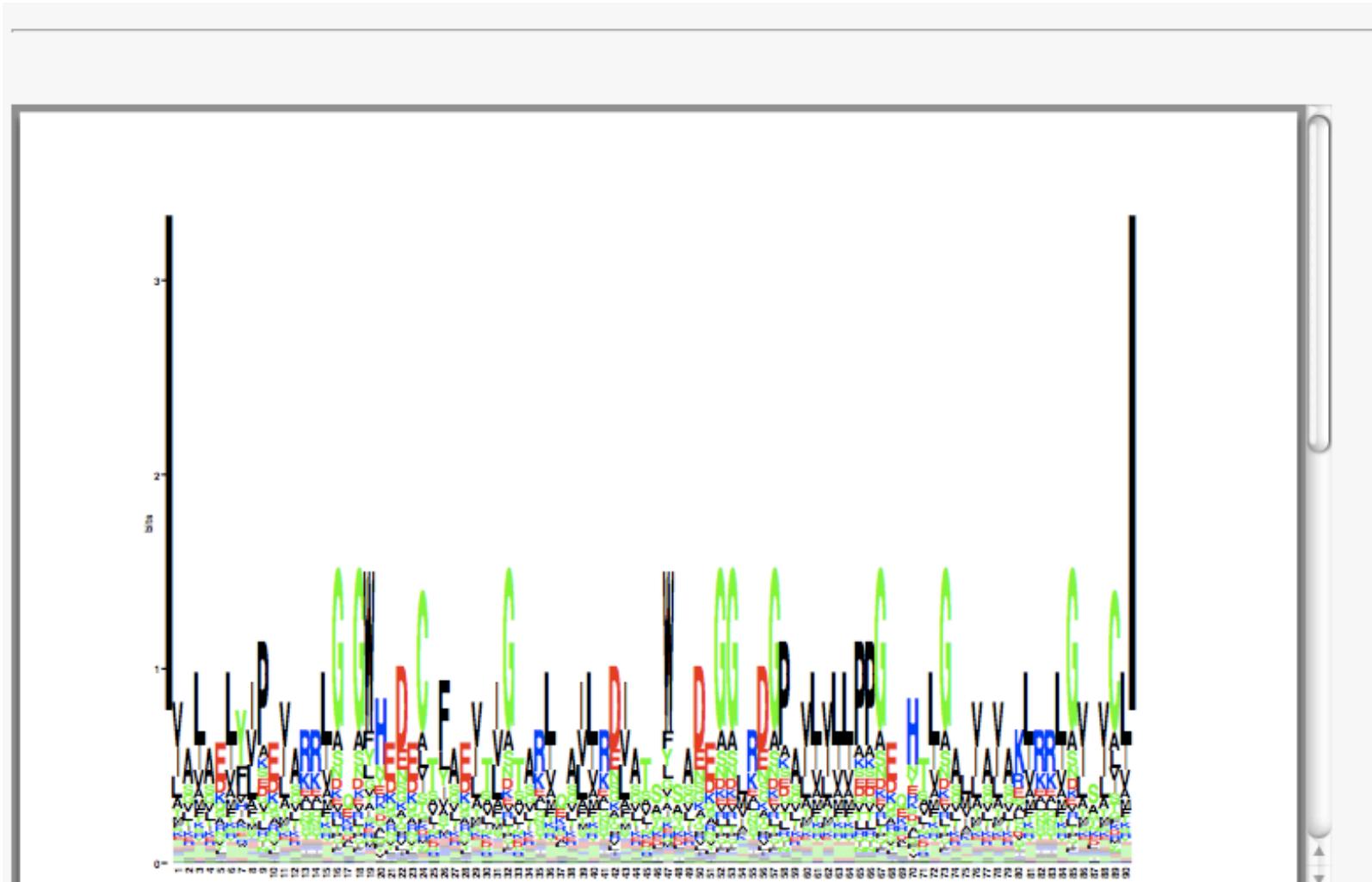
Stack Linesize

Plot Kullback-Leibler logo

File format for logo file PDF

Submit Clear fields

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	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
VALAELYIPEVARRLGQQGWHEDECTFAEVТИCTARLQAILRDIATWSADEGGMRDGPALVLLPPG
EQHTLGAMVAVAKLRRLGVSVCLRMSTGPAELRELFGRKRRFDAIMISLAHAEMLLEVGRKLVKTLKD
MTGGRIPVAMGGALFLDGTEAASIPEADIVTNIDIEALQ
```

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

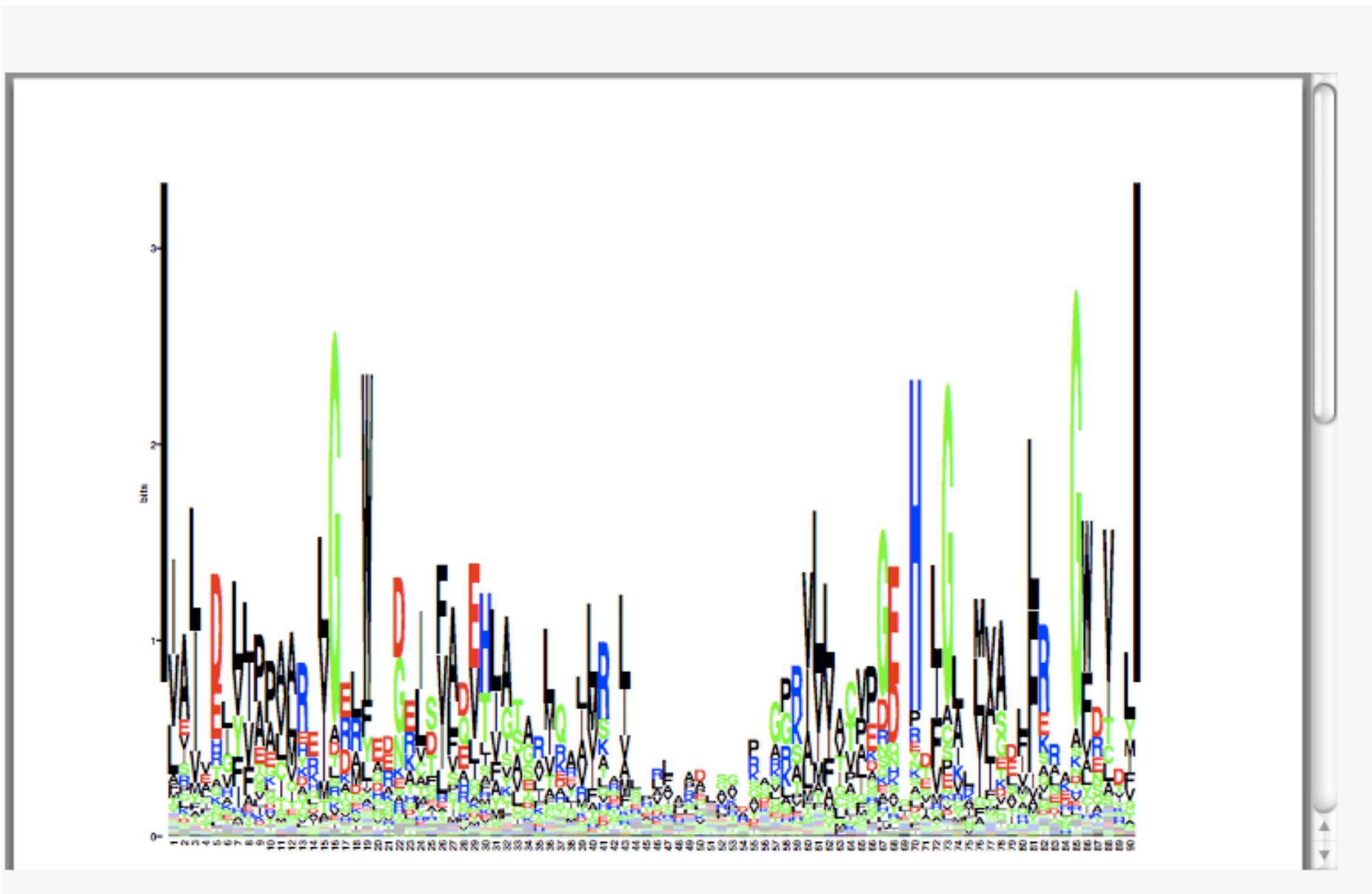
Plot Kullback-Leibler logo

File format for logo file

Restrictions:

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

Blast2logo



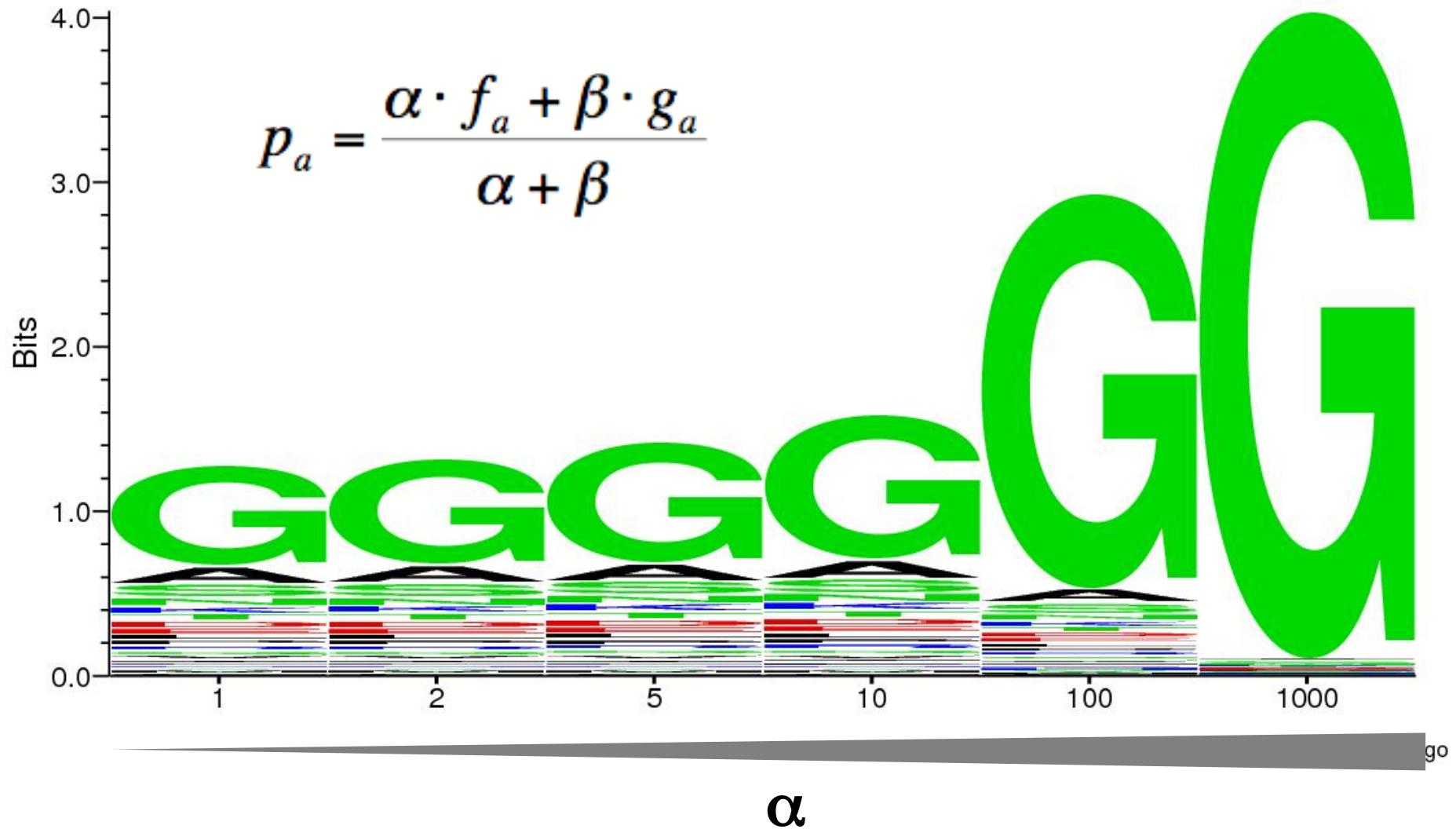
Blast2logo

CENTERFO
R BIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

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 or Gaining confidence



Example.

>1K7C.A

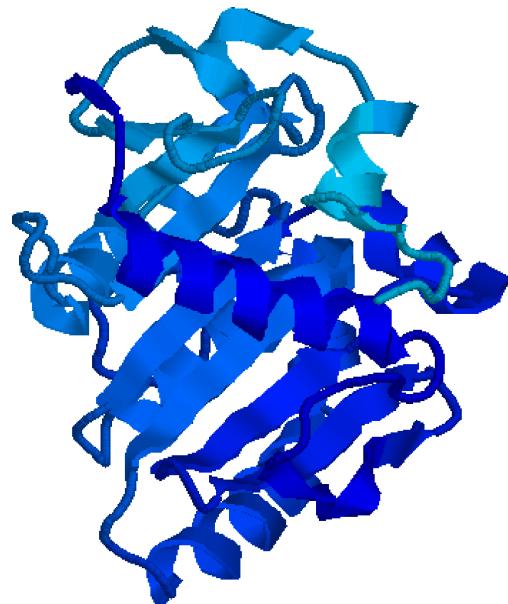
TTVYLAGDSTMAGNGGGSGTNGWGEYLASYLSATVVNDAVAGRSARSYTREGRFENIADV
VTAGDYVIVEFGHNDGGSLSTDNGRTDCSGTGAEVCYSYDGVNETILTFPAYLENAAKL
FTAKGAKVILSSQTPNNPWETGTFVNSPTRFVEYAAEVAGVEYVDHWSYVDSIYETL
GNATVNSYFPIDHTHTSPAGAEVVAEFLKAVVCTGTSKSVLTTSFEGTCL

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

What would you do?

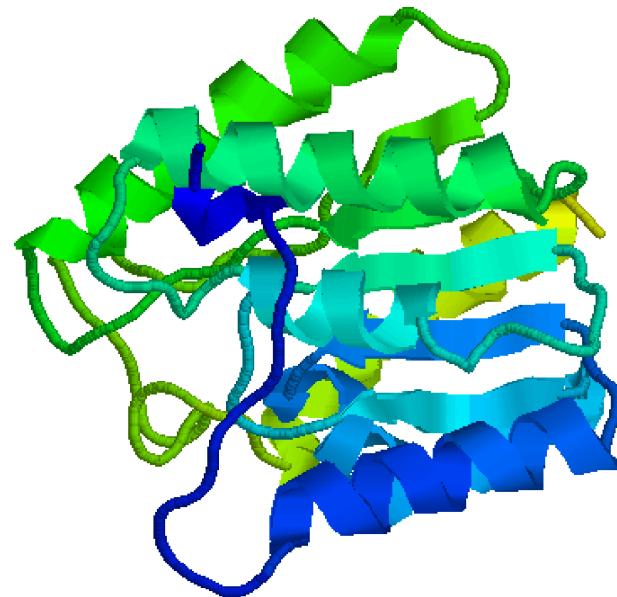
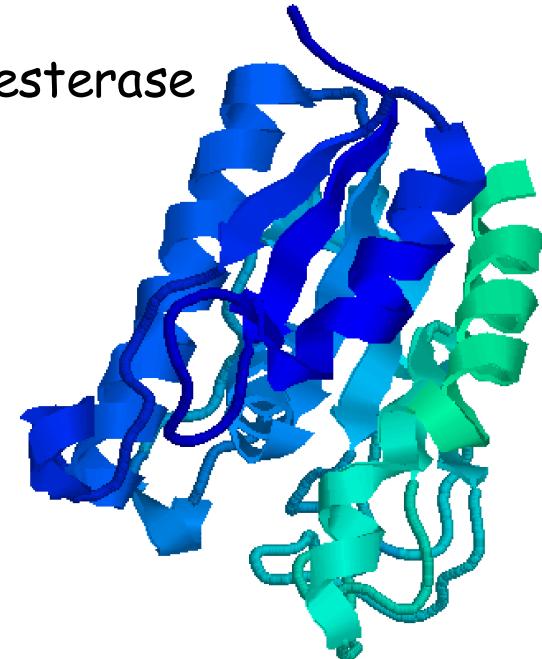
- Function
 - Run Blast against PDB
 - No significant hits
 - Run Blast against NR (Sequence database)
 - Function is Acetyl esterase?
- Where is the active site?

Example. Where is the active site?



1USW Hydrolase

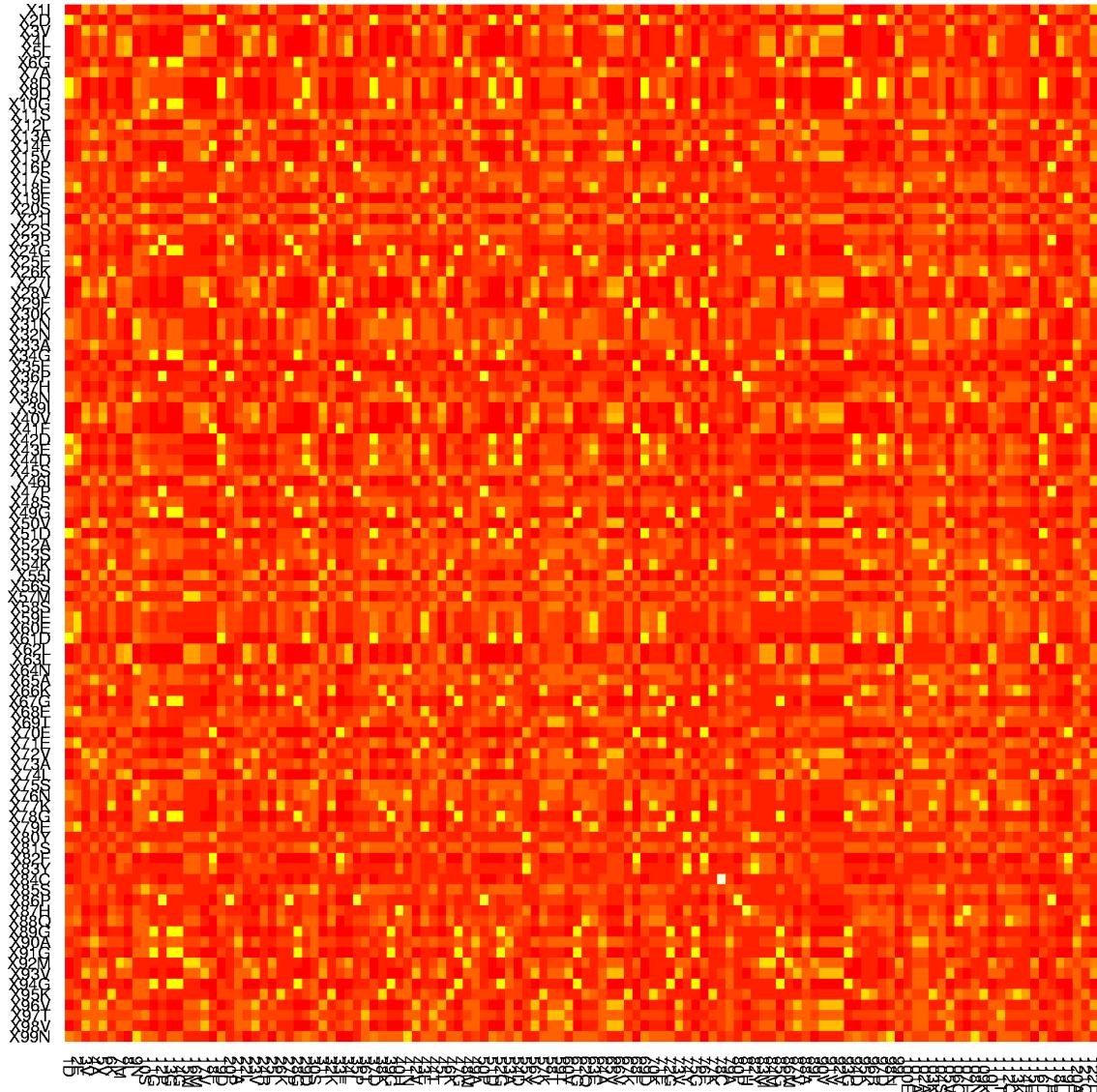
1G66 Acetylxylan esterase



1WAB Acetylhydrolase

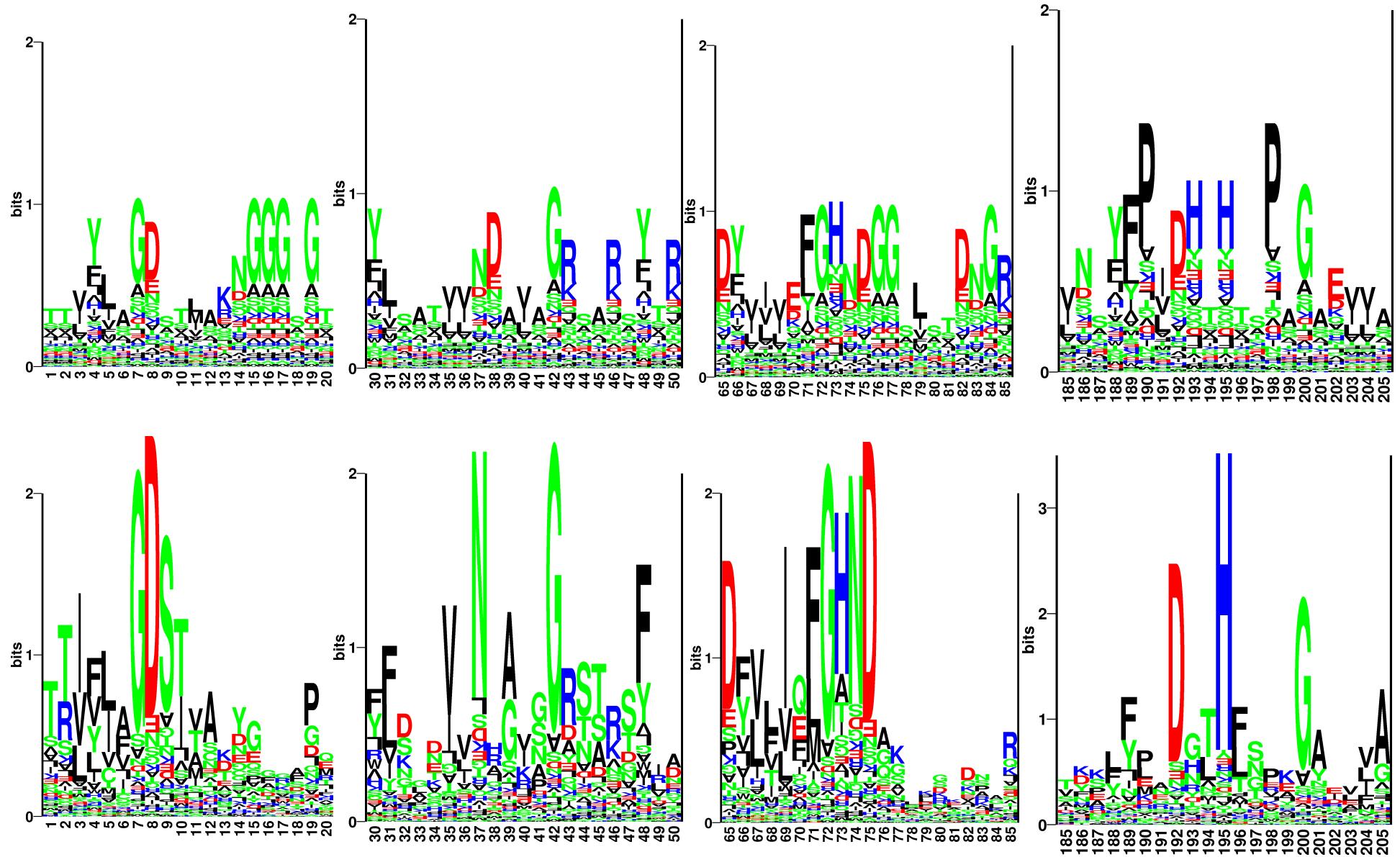
When Blast fails!

1K7A.A



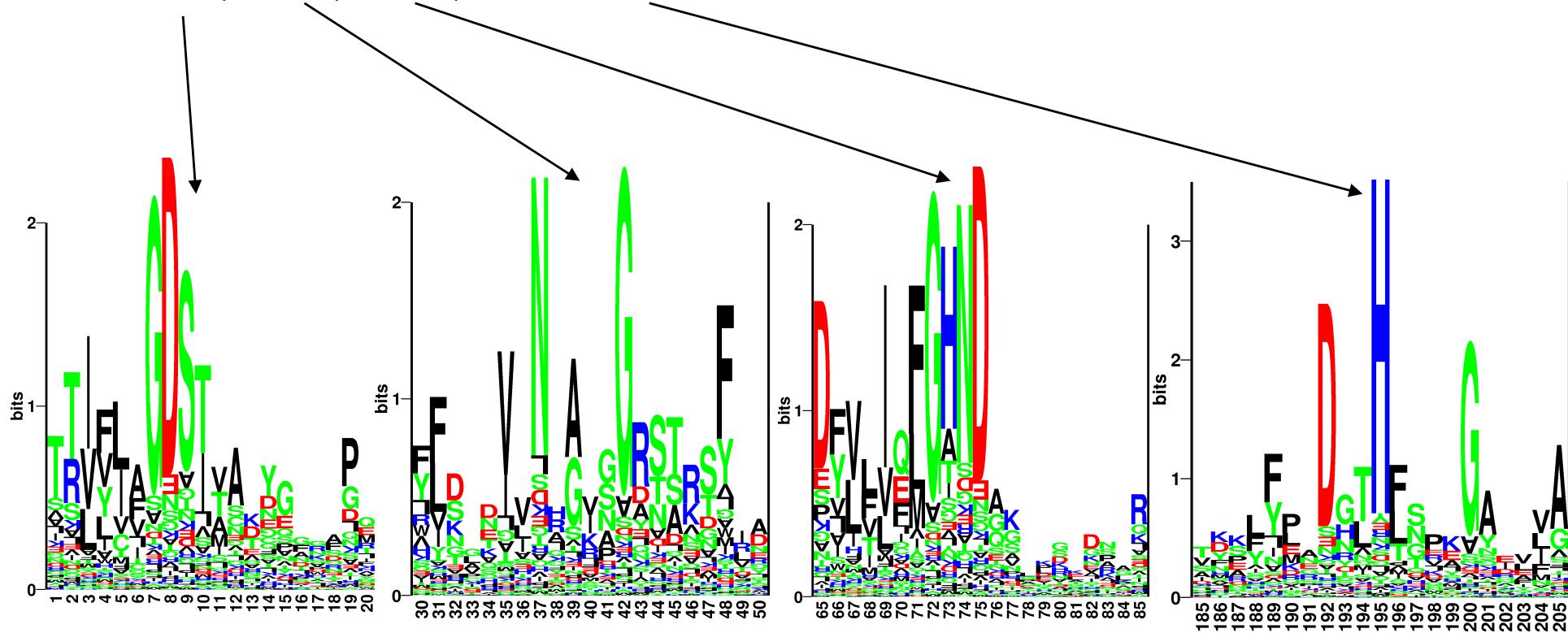
1WAB._

Example. (SGNH active site)

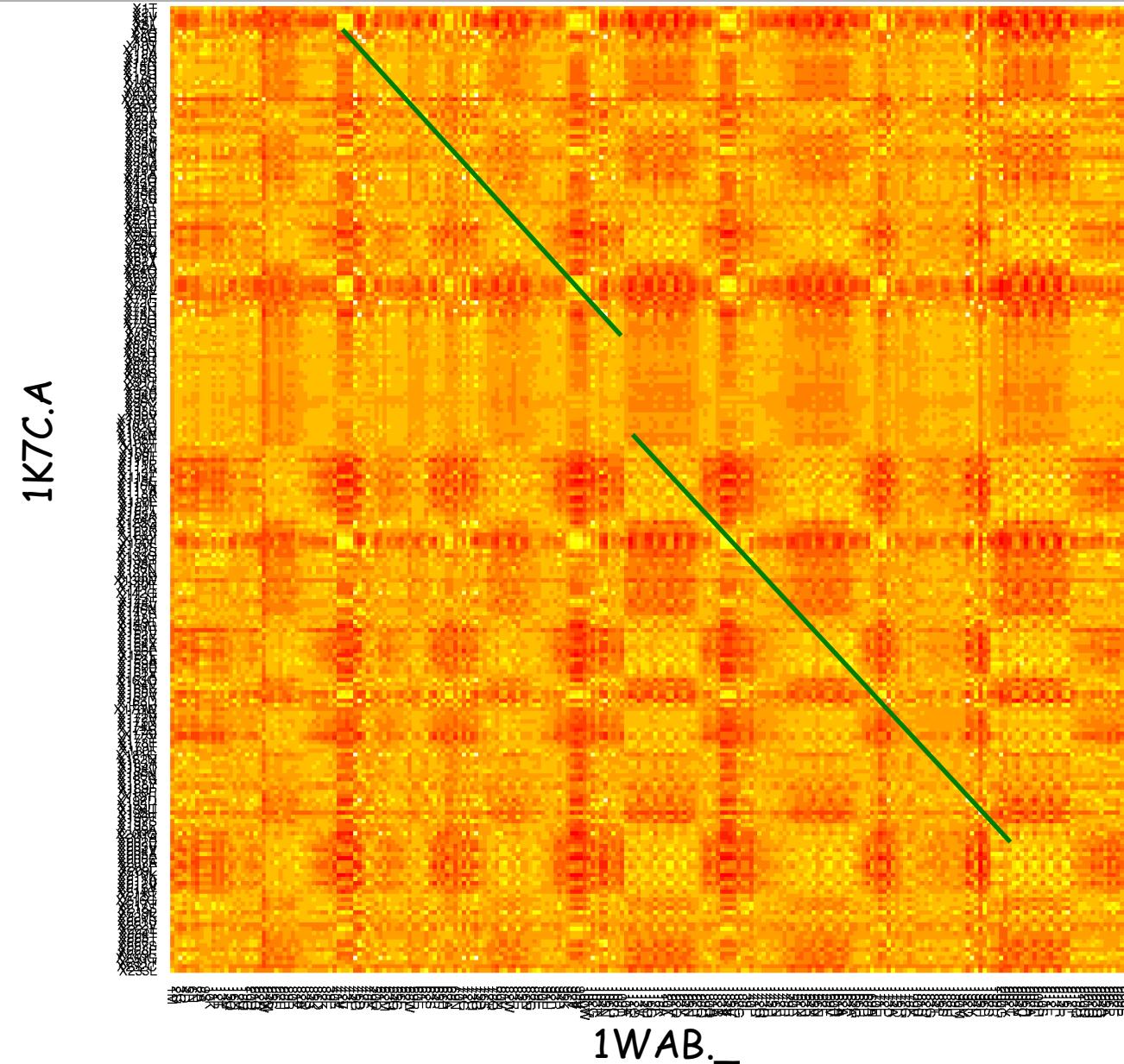


Example. Where is the active site?

- 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**RSARSYTREGRFENIADVVTAGDYVIVEFGH**N**DGGSISTDN
S G N

1WAB . _ EVVFIGDSLVQILMHQCE---IWRELFS---PLHALNFGIG**G**DSTQHVLW--RLENGELEHIRPKIVVVVWGT**NN**HG-----

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

1K7C . A GVEYVDHWSYVDSIYETLGNATVNSYFPIDHT**H**TSPAGAEVVAEAFLKAVVCTGTSL
H

1WAB . _ RAHFLLADPG---FVHSDG--TISHHDMDYDYL**H**LSRLGYTPVCRALHSLLLRL---L

Handout exercise

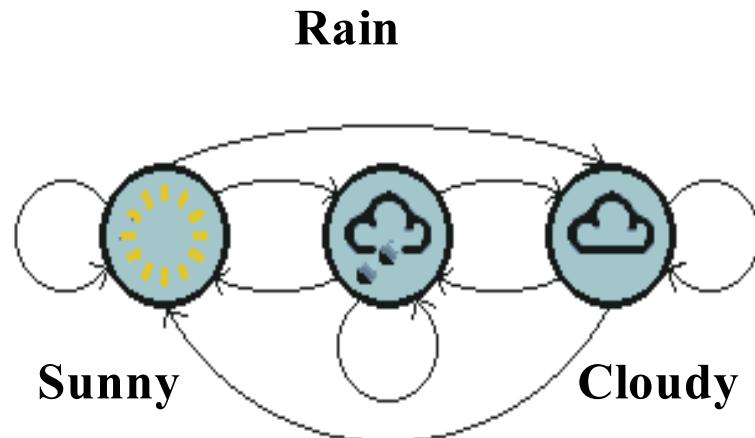
Using Psi-Blast Profiles

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

Hidden Markov Models, HMM's

Markov Chains



States : Three states - sunny, cloudy, rainy.

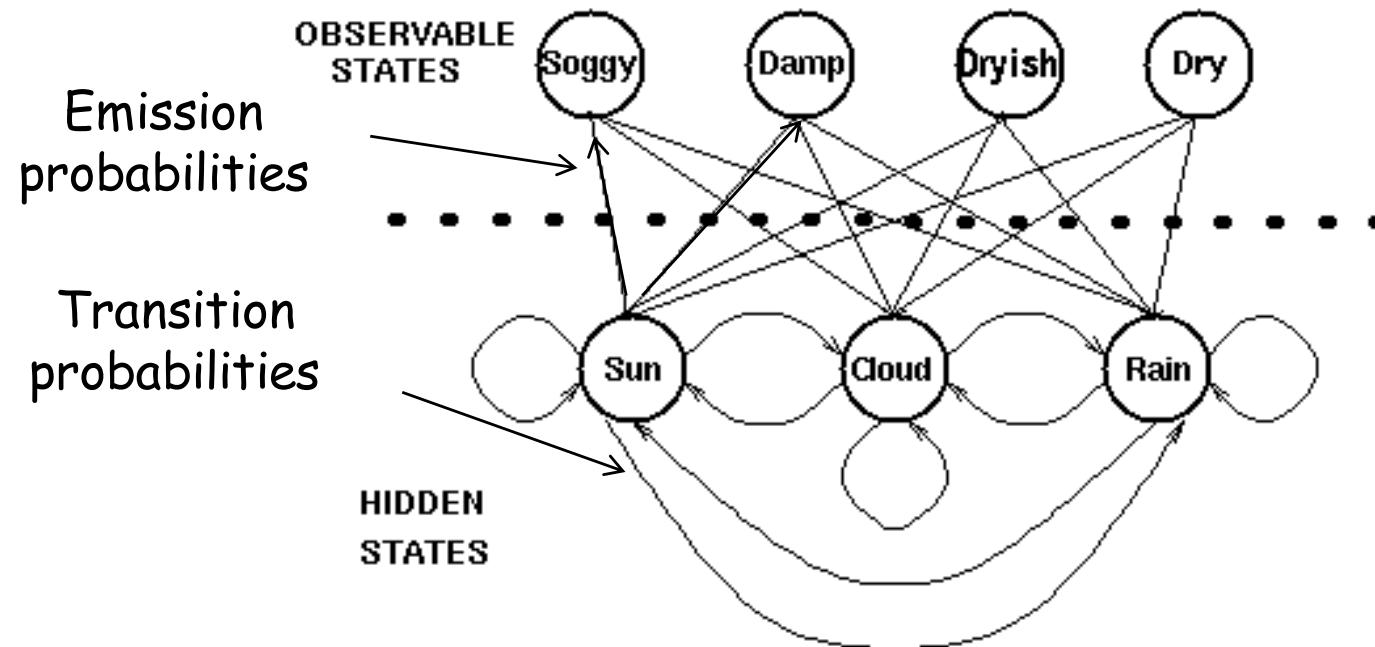
		weather today			
		Sun	Cloud	Rain	
weather yesterday	Sun	0.5	0.25	0.25	
	Cloud	0.375	0.125	0.375	
		Rain	0.125	0.625	0.375

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
R BIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

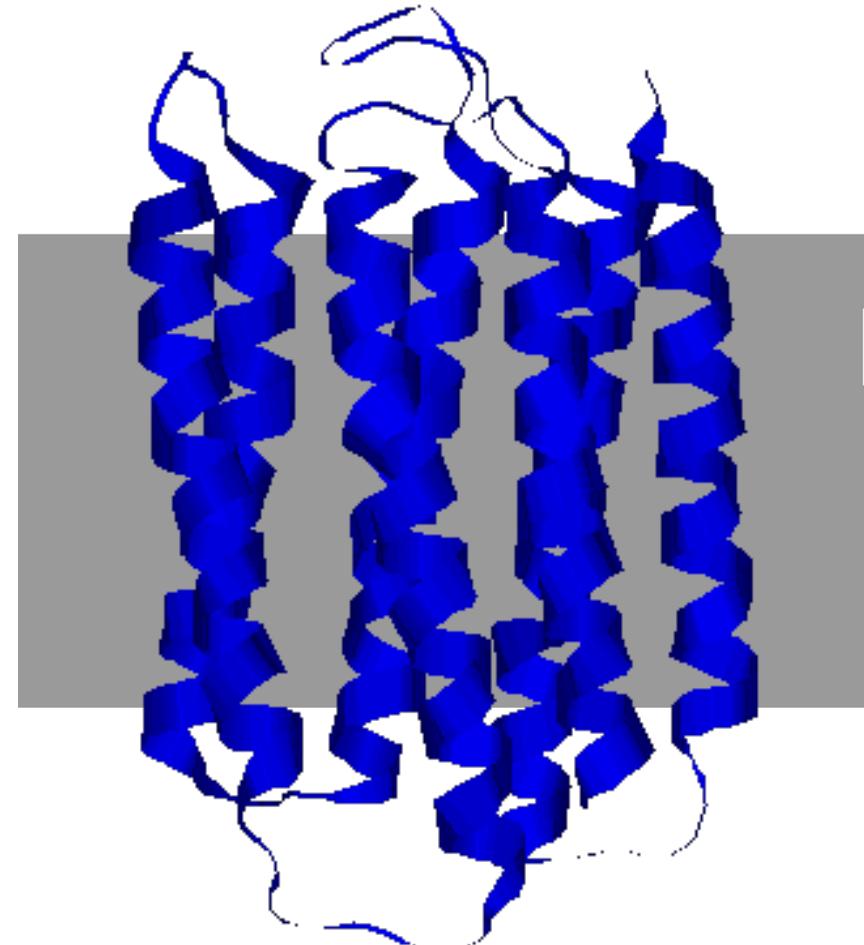
>3UTX:A|PDBID|CHAIN|SEQUENCE

QAQITGRPEWIWLALGTALMGLGTLYFLVKGMGVSDPDAKKFYAIATLVPAlAFTMYLSMLLGYGLTMVP
FGGEQNPIYWARYADWLFTTPLLLDLALLVDADQGTILALVGADGIMIGTGLVGALTkvSYRFVWWAIST
AAMLYILYVLFFGFTSKAESMRPEVASTFKVLRNVTVVLWSAYPVVWLIGSEGAGIVPLNIETLLFMVLDVS
AKVGFGLILLRSRAIFGEAEAPEPSAGDGAAATSD

Extra cellular

Trans membrane

Intra cellular



Weight matrix (PSSM) construction

SLLPAIVEL YLLPAIVHI TLWVDPYEV GLVPFLVSV KLLEPVLLL LLDVPTAAV LLDVPTAAV LLDVPTAAV
LLDVPTAAV VLFRGGPRG MVDGTLLLL YMNGTMSQV MLLSVPLLL SLLGLLVEV ALLPPINIL TLIKIQHTL
HLIDYLVTS ILAPPVKL ALFPQLVIL GILGFVFTL STNRQSGRQ GLDVLTAKV RILGAVAKV QVCERIPTI
ILFGHENRV ILMEHIHKL ILDQKINEV SLAGGIIGV LLIENVASL FLLWATAEA SLPDFGISY KKREEAPSL
LERPGGNEI ALSNLEVKL ALNELLQHV DLERKVESL FLGENISNF ALSDHHIYL GLSEFTEYL STAPPAGV
PLDGEYFTL GVLVGVALI RTLDKVLEV HLSTAFARV RLDSYVRSL YMNGTMSQV GILGFVFTL ILKEPVHGV
ILGFVFTLT LLFGYPVYV GLSPTVWLS WLSLLVPFV FLPSDFFPS CLGGILTMV FIAGNSAYE KLGEFYNQM
KLVALGINA DLMGYIPLV RLVTLKDIV MLLAVLYCL AAGIGILTV YLEPGPVTA LLDGTATLR ITDQVPFSV
KTWGQYWQV TITDQVPFS AFHHVAREL YLNKIQNSL MMRKLAILS AIMDKNIIL IMDKNIILK SMVGNWAKV
SLLAPGAKQ KIFGSLAFL ELVSEFSRM KLTPLCVTL VLYRYGSFS YIGEVLVSV CINGVCWTW VMNILLQYV
ILTVIDGVL KVLEYVIKV FLWGPRALV GLSRYVARL FLLTRILTI HLGNVKYLV GIAGGLALL GLQDCTMLV
TGAPVTYST VIYQYMDL VLPDVFIRC VLPDVFIRC AVGIGIAVV LVVLGLLAV ALGLGLLPV GIGIGVLA
GAGIGVAVL IAGIGILAI LIVIGILIL LAGIGLIAA VDGIGILTI GAGIGVLTA AAGIGIIQI QAGIGILLA
KARDPHSGH KACDPHSGH ACDPHSGHF SLYNTVATL RGPGRAFVT NLVPMVATV GLHCYEQLV PLKQHFQIV
AVFDRKSDA LLDFVRFMG VLVKSPNHV GLAPPQHLL LLGRNSFEV PLTFGWCYK VLEWRFDSR TLNAWVKVV
GLCTLVAML FIDSYICQV IISAVVGIL VMAGVGSPY LLWTLVVLL SVRDRLARL LLMDCSGSI CLTSTVQLV
VLHDDILLEA LMWITQCFL SLLMWITQC QLSLLMWIT LLGATCMFV RLTRFLSRV YMDGTMSQV FLTPKKLQC
ISNDVCAQV VKTDGNPPE SVYDFFVWL FLYGALLA VLFSSDFRI LMWAKIGPV SLLLELEEV SLSRFWSWA
YTAFTIPSI RLMKQDFSV RLPRIFCSC FLWGPRAYA RLLQETELV SLFEGIDFY SLDQSVVEL RLNMFTPYI
NMFTPYIGV LMIPIINV TLFIGSHVV SLVIVTTFV VLQWASLAV ILAKFLHWL STAPPHVNV LLLLTVLTV
VVLGVVFGI ILHNGAYSL MIMVKCWMY MLGHTHTMEV SLADTNSLA LLWAARPRL GVALQTMQ
GLYDGMEHL KMVELVHFL YLQLVFGIE MLMAQEALA LMAQEALAF VYDGREHTV YLSGANLNL RMFPNAPYL
EAAGIGILT TLDSQVMSL STPPPGRTRV KVAELVHFL IMIGVLGVV ALCRWGLLL LLFAGVQCQ VLLCESTAV
YLTAFARV YLLEMILWRL SLDDYNHIV RTLDKVLEV GLPVEYLQV KLIANNTRV FIYAGSLSA KLVANNTRL
FLDEFMEGV ALQPGTALL VLDGLDVLL SLYSFPEPE ALYVDSLFF SLLQHLLIGL ELTLGEFLK MINAYLDKL
AAGIGILTV FLPSDFFPS SVRDRLARL SLREWLLRI LLSAWILTA AAGIGILTV AVPDEIPPL FAYDGKDYI
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 acid 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
-

HMM (a simple example)

ACA---**ATG**

TCAACT**ATC**

ACAC--**AGC**

AGA---**ATC**

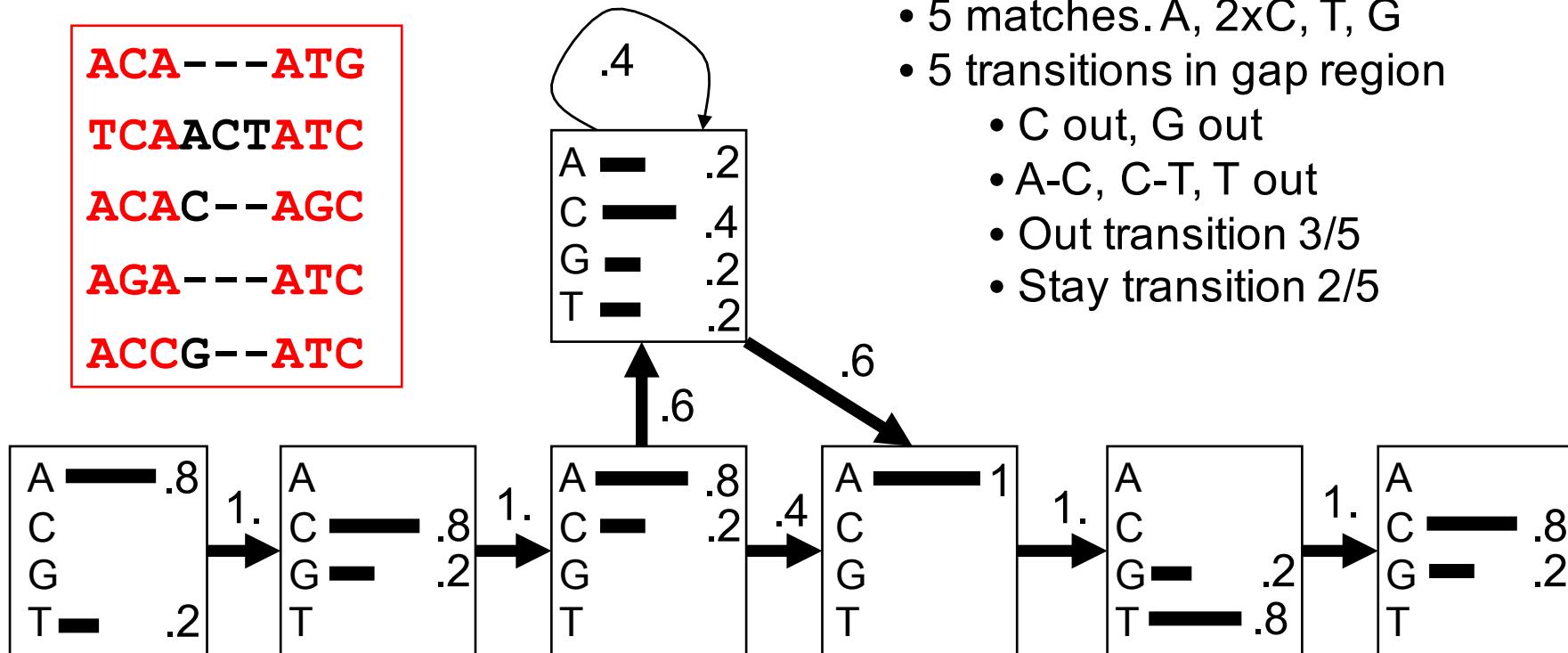
ACCG--**ATC**



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

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
 - $\text{Log}(\text{P}(S)/0.25^L)$
 - Positive score means more likely than Null model
- This is just like we did for PSSM $\text{log}(p/q)$!

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

Note!

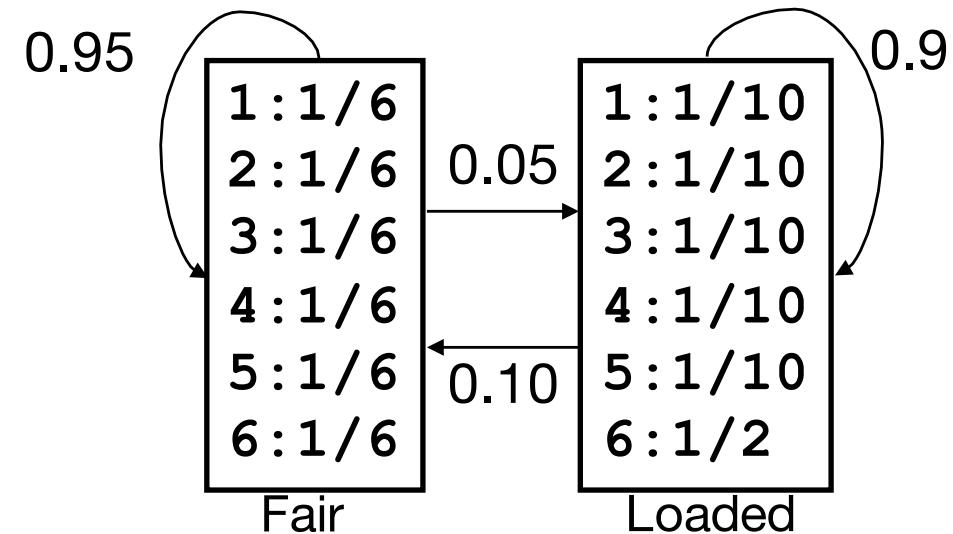
Exceptional:

TGCT--AGG = -0.97

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



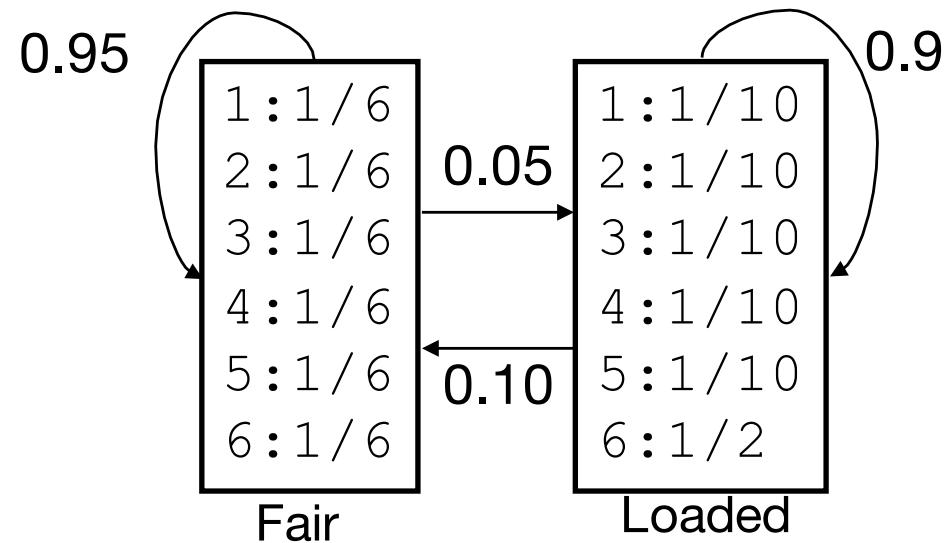
Aligning a sequence to an HMM

- 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

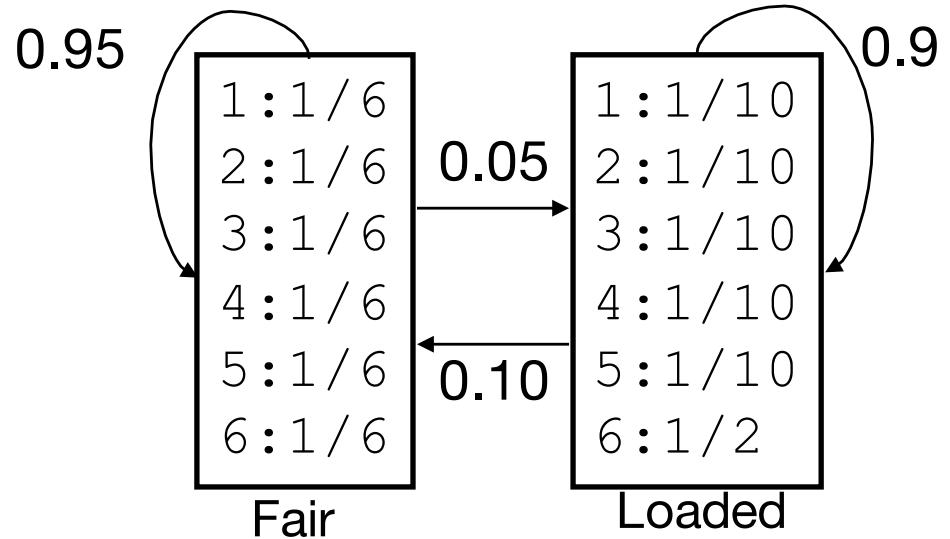
The unfair casino: Loaded dice $p(6) = 0.5$; switch fair to load:0.05; switch load to fair: 0.1

- 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



$$\text{FFF} = 0.5 * 0.167 * 0.95 * 0.167 * 0.95 * 0.167 = 0.0021$$

$$\text{FFL} = 0.5 * 0.167 * 0.95 * 0.167 * 0.05 * 0.5 = 0.00333$$

$$\text{FLF} = 0.5 * 0.167 * 0.05 * 0.5 * 0.1 * 0.167 = 0.000035$$

$$\text{FLL} = 0.5 * 0.167 * 0.05 * 0.5 * 0.9 * 0.5 = 0.00094$$

$$\text{LFF} = 0.5 * 0.1 * 0.1 * 0.167 * 0.95 * 0.167 = 0.00013$$

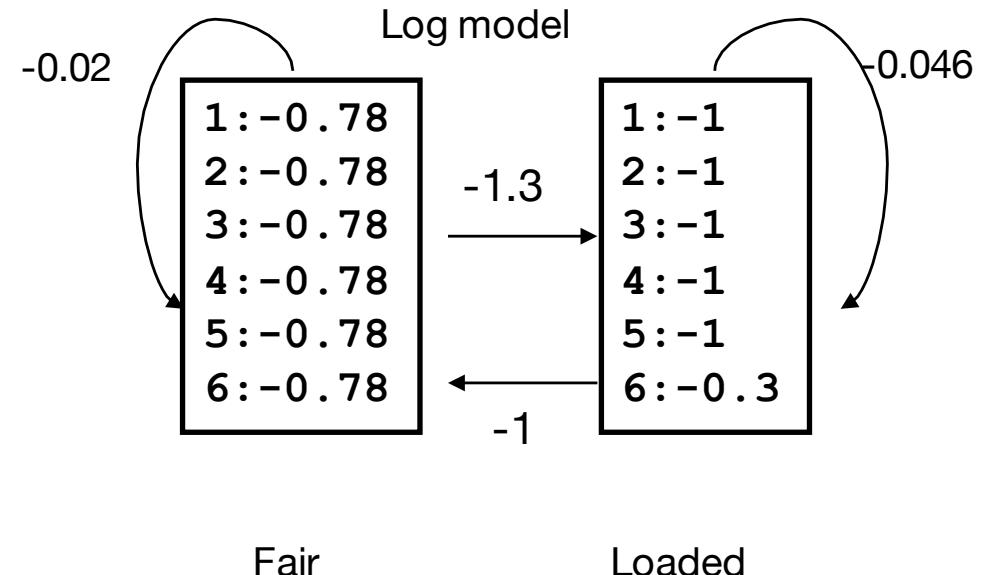
$$\text{LFL} = 0.5 * 0.1 * 0.1 * 0.167 * 0.05 * 0.5 = 0.000021$$

$$\text{LLF} = 0.5 * 0.1 * 0.9 * 0.5 * 0.1 * 0.167 = 0.00038$$

$$\text{LLL} = 0.5 * 0.1 * 0.9 * 0.5 * 0.9 * 0.5 = 0.0101$$

Or in log space

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



$$\log(P(LLL|M)) = \log(0.5 * 0.1 * 0.9 * 0.5 * 0.9 * 0.5) = \log(0.0101)$$

or

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

Model decoding (Viterbi)

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

$$FF = 0.5 * 0.167 * 0.95 * 0.167$$

$$\text{Log}(FF) = -0.30 - 0.78 - 0.02 - 0.78 = -1.88$$

$$LF = 0.5 * 0.1 * 0.1 * 0.167$$

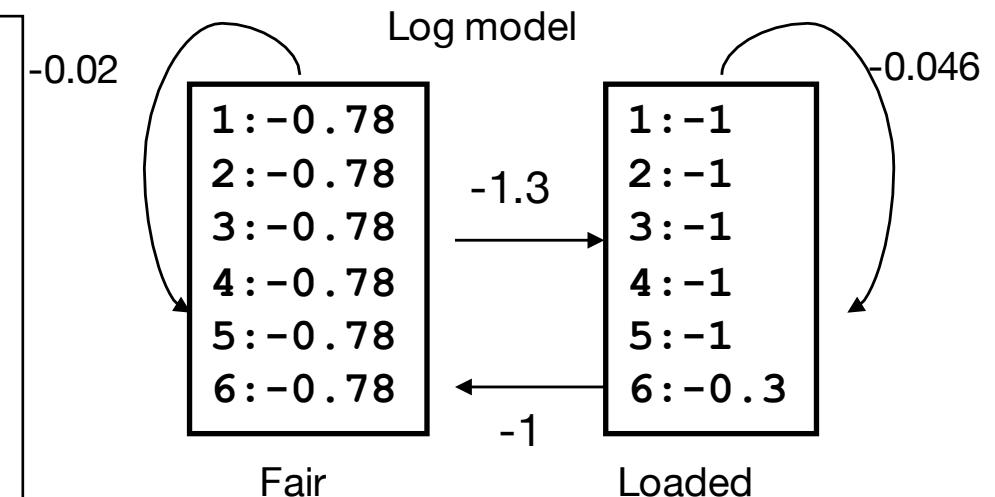
$$\text{Log}(LF) = -0.30 - 1 - 1 - 0.78 = -3.08$$

$$FL = 0.5 * 0.167 * 0.05 * 0.5$$

$$\text{Log}(FL) = -0.30 - 0.78 - 1.30 - 0.30 = -2.68$$

$$LL = 0.5 * 0.1 * 0.9 * 0.5$$

$$\text{Log}(LL) = -0.30 - 1 - 0.046 - 0.3 = -1.65$$

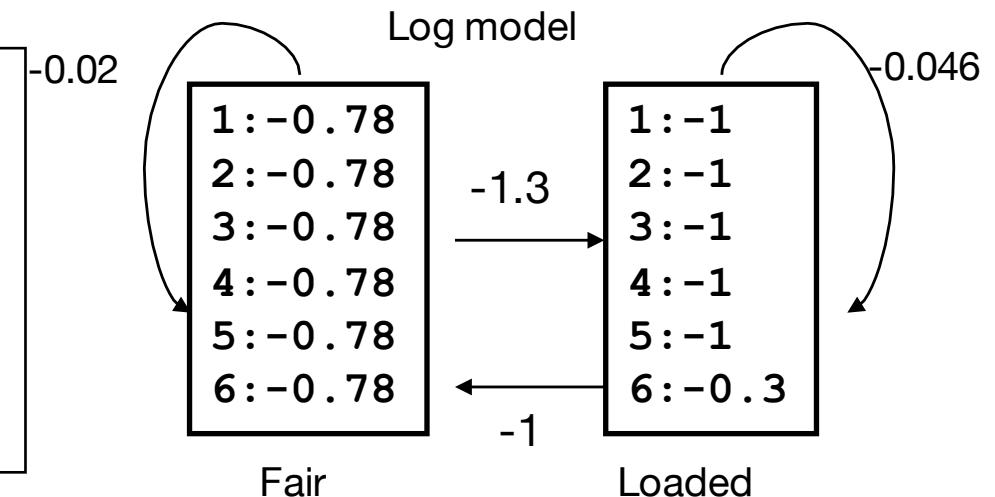


	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88							
L	-1.30	-1.65							

Model decoding (Viterbi)

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

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

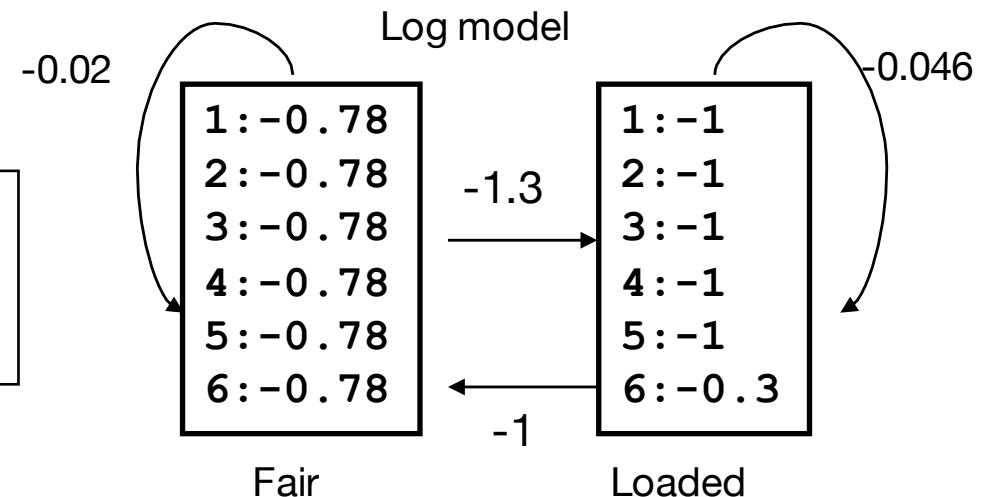


	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88							
L	-1.30	-1.65							

Model decoding (Viterbi)

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

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



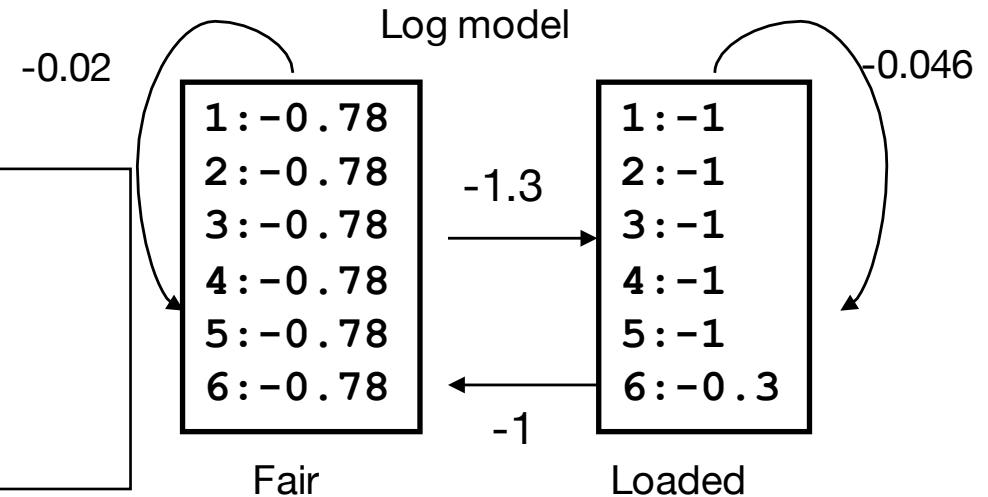
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68						
L	-1.30	-1.65	-1.99						

Model decoding (Viterbi) - Or doing it recursively

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

```

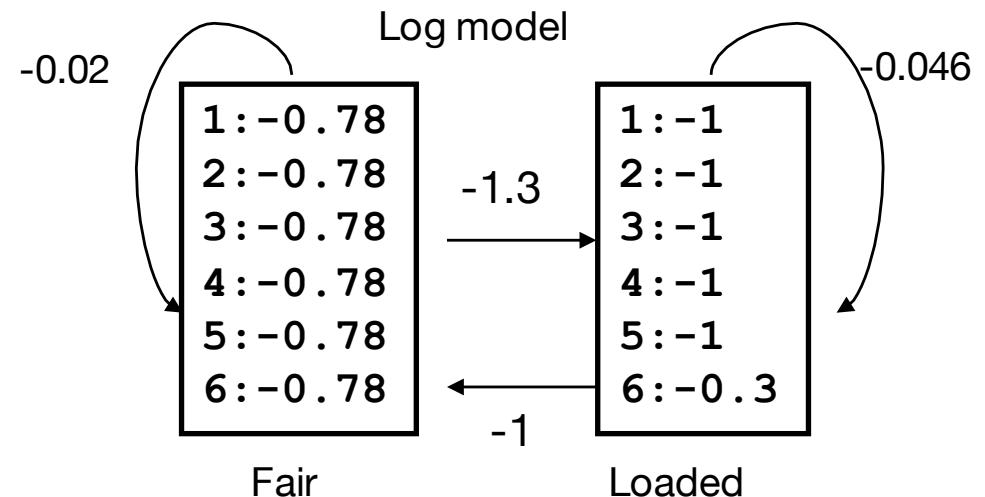
FFF = -0.78 + max(-1.88-0.02, -1.99-1)
      = -0.78 - 1.88 - 0.02
Log(P(FFF)) = -2.68
LLL = -0.3 + max(-1.88-1.3, -1.65-0.046)
      = -0.3 - 1.65 - 0.046
Log(P(LLL)) = -1.99
    
```



	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?



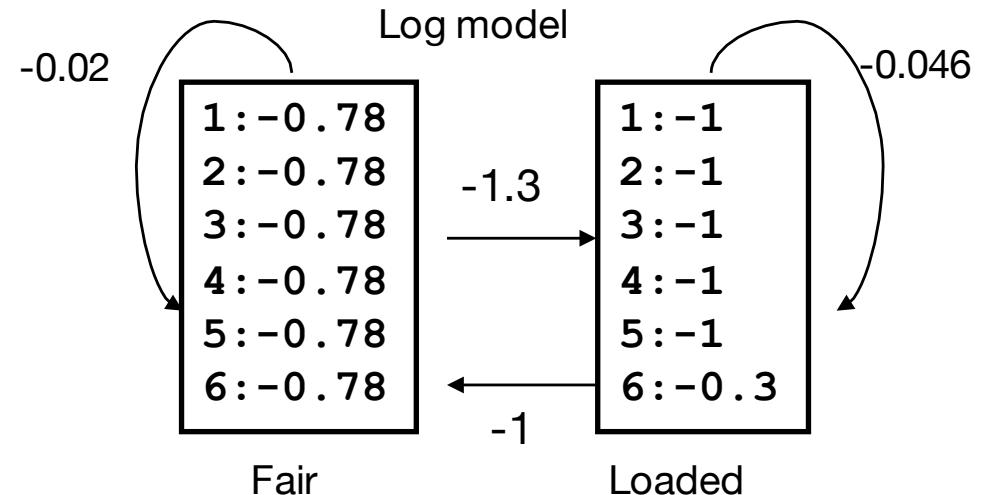
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68						
L	-1.30	-1.65	-1.99						

Model decoding (Viterbi)

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

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

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



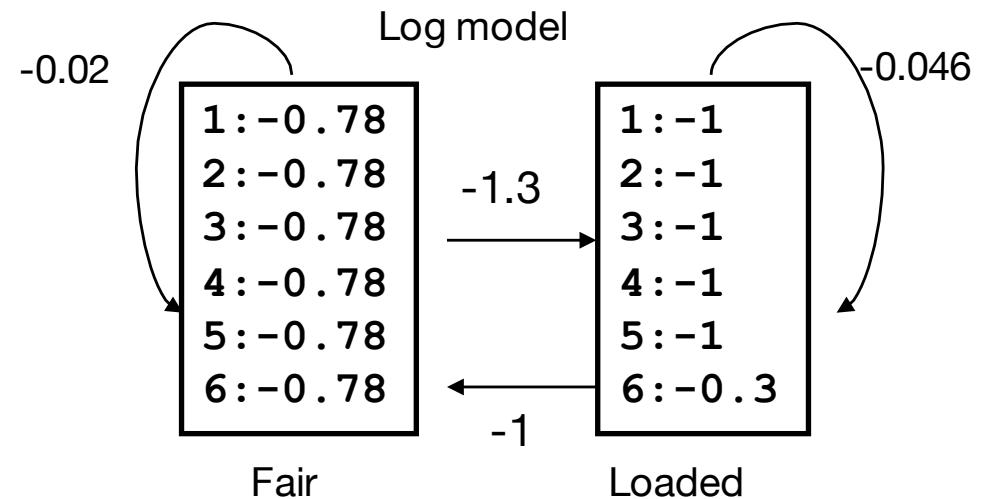
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-					
L	-1.30	-1.65	-1.99	-					

Model decoding (Viterbi)

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

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

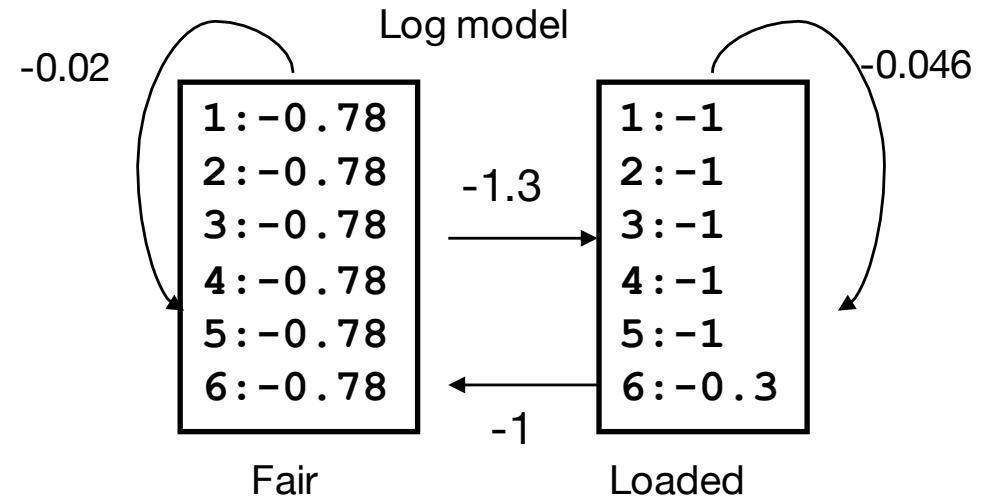
$$\boxed{-0.78 - 1 - 1.99 = -3.77}$$



	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

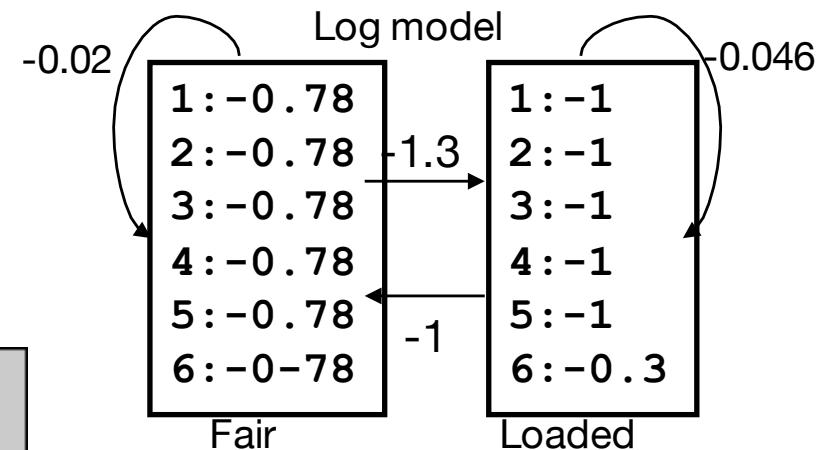
Transition

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

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

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



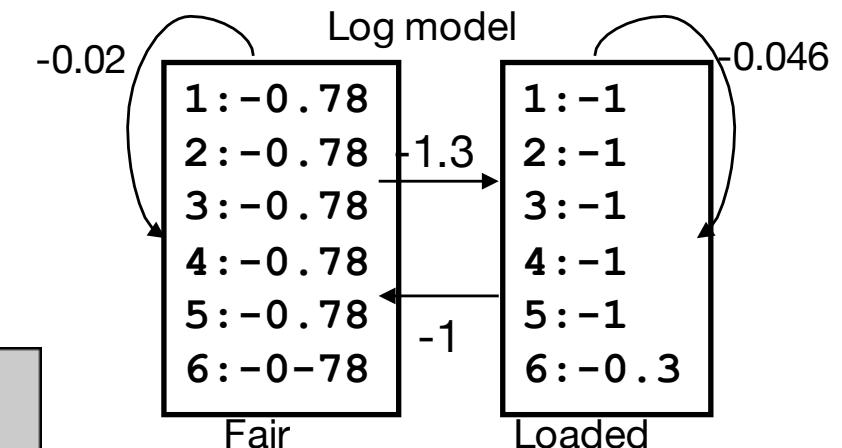
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-3.48					
L	-1.30	-1.65	-1.99						

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

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

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



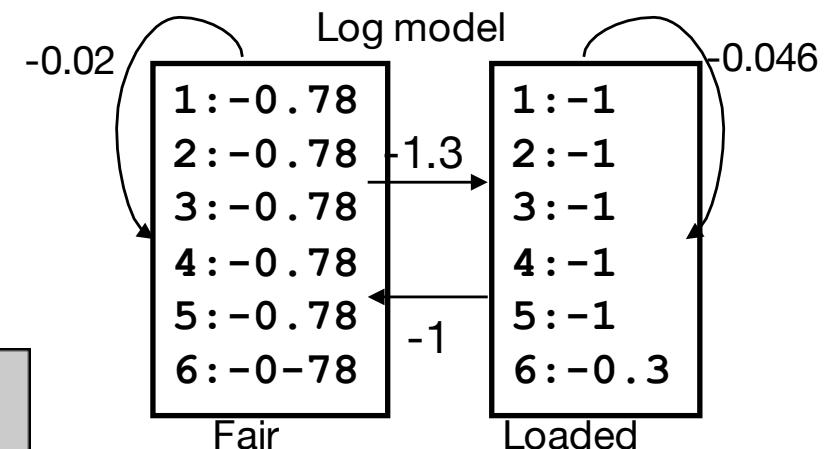
	5	6	6	6	1	1	2	3	4
F	-1.08	-1.88	-2.68	-3.48			-4.92	-5.73	-6.53
L	-1.30	-1.65	-1.99		-3.39	-4.44		-6.52	-7.57

Model decoding (Viterbi). Can you do it?

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

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

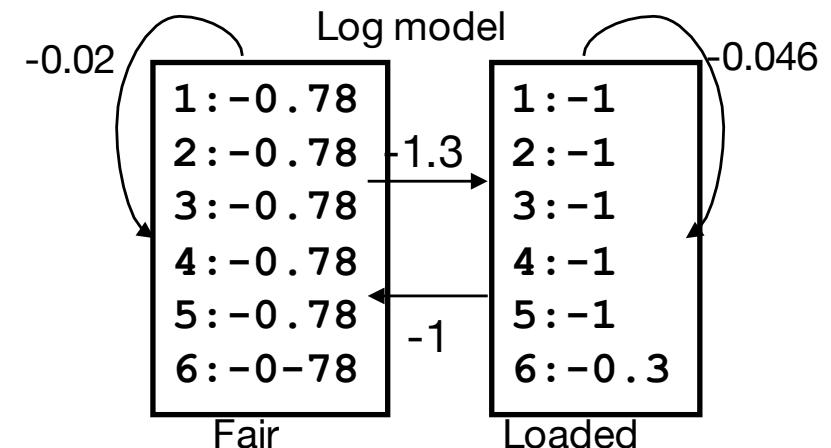
$$\log(P_l(i+1)) = \log(p_l(i+1)) + \max_k (\log(P_k(i)) + \log(a_{kl}))$$



	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	-4.44	-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								

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

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

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

Conserved deletion Non-conserved Insertion

The diagram shows two protein sequences aligned vertically. Regions of conservation are highlighted with light blue boxes. A red annotation 'Conserved' points to a box around the first sequence's 'SISPGE' motif. A red annotation 'deletion' points to a gap in the second sequence's alignment. A red annotation 'Non-conserved Insertion' points to a box around the second sequence's 'QCP' motif. The sequences are:

ADDGSLAFVPSEF--SISPGEKIVFKNNAGFPNIVFDED S IPSGV DASK I SMSEE DLLN
TVNGAI --PGPLIAERLKE G Q NVRVTNTLDEDTS I HWHGLLVPFGMDGVPGVS FPG--- I
-TSMAPAFGVQE FYRTVKQ G D E VT VTI T ----- N IDQIED -V SHGFVVVN HGV SME --- I
IE -- KMKYL TPEVF YTI KAGE TVYWVN G E VMPHN V AFKKGIV -- GEDAF R GEMMTK D ---
-TSVAPSFSQPSF -LTVKE G D E VT VIV TNLDE ----- IDDLTHGFTMGN HGV AME --- V
ASAETMVFE PDFLV LIEIGPGDRVRFV PTHK -SHNAATIDGMVPEGVE GFKSR INDE ---
TVNGQ -- FPGPRLAGVARE GD QVLV KV VN HVAEN ITIHW HG VQLGTGWADGPAYVT QCP I
TKAVVLT FNTSVE ICLVMQ G TSIV --- AAESHP LHLHG FNP S FN LVD P M ERNTAGVP

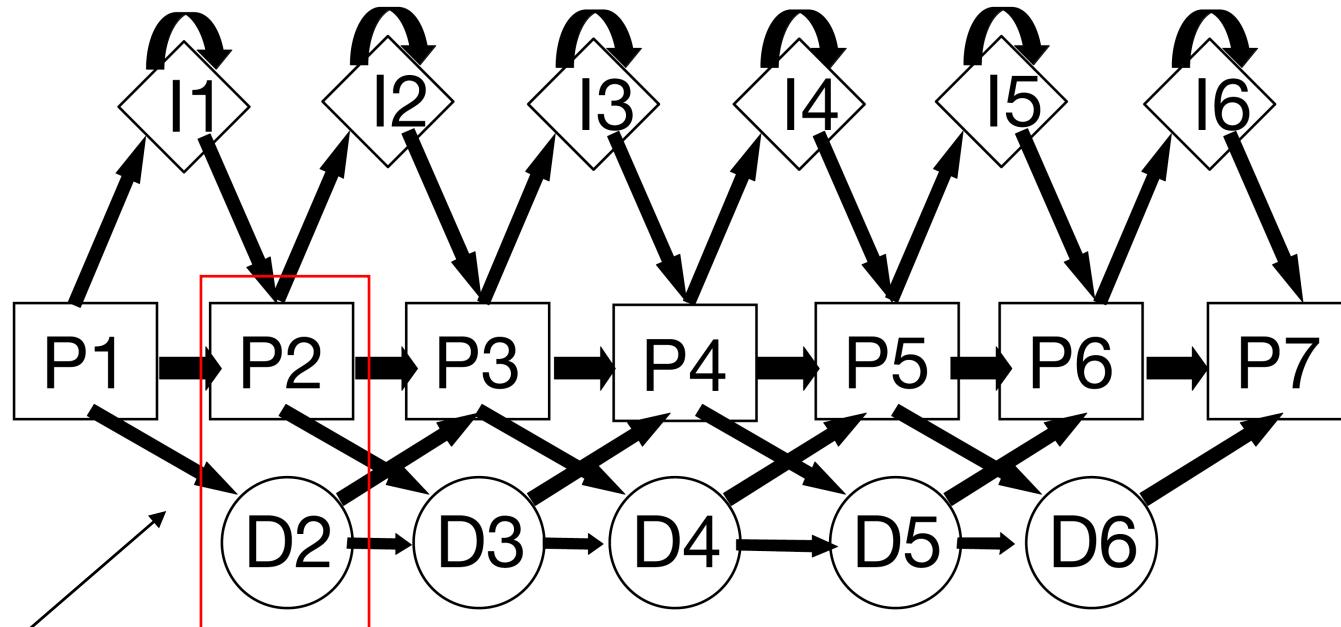
Matching any thing
but *G* => large
negative score

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

Profile HMM's

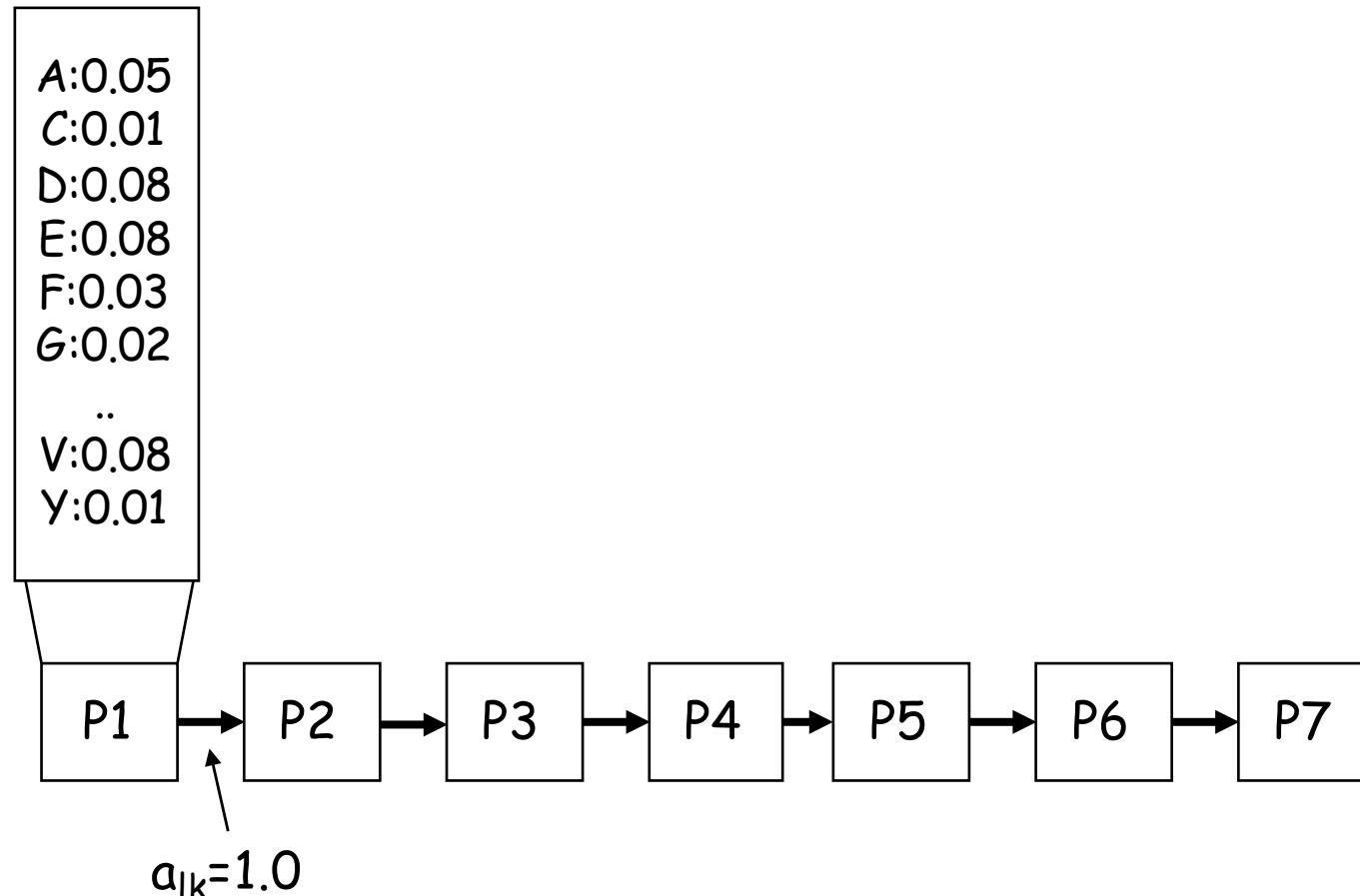


All P/D pairs must
be visited once

$L_1 - Y_2 A_3 V_4 R_5 - I_6$
 $P_1 D_2 P_3 P_4 I_4 P_5 D_6 P_7$

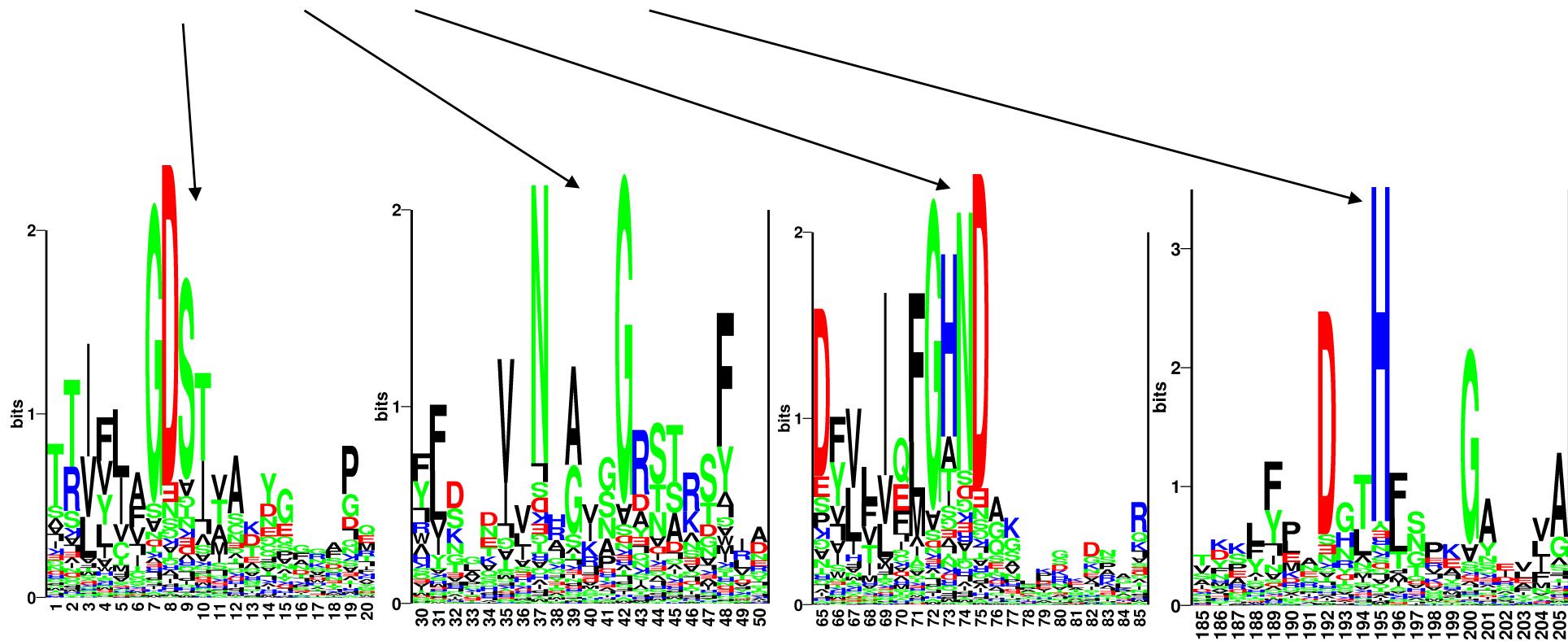
Profile HMM

- Un-gapped profile HMM is just a sequence profile



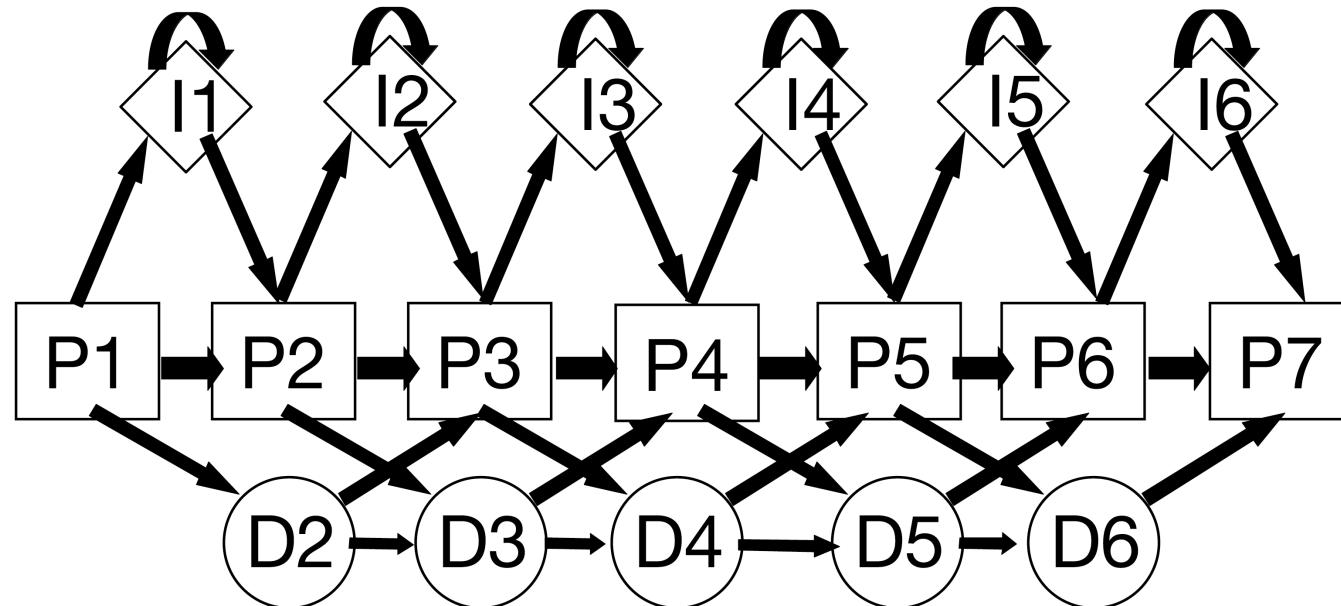
Example. Where is the active site?

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



Profile HMM

- Profile HMM (deletions and insertions)



Profile HMM (deletions and insertions)

QUERY	HAMDIRCYHSGG	PLHL	GE	EDFNGQ	SCIVCPWHKYKITLATGE	GLYQSINPKDPS
Q8K2P6	HAMDIRCYHSGG	PLHL	GE	EDFNGQ	SCIVCPWHKYKITLATGE	GLYQSINPKDPS
Q8TAC1	HAMDIRCYHSGG	PLHL	GD	EDFDGRPCIVCPWHKYKITLATGE	GLYQSINPKDPS	
Q07947	FAVQDTCTHGDW	ALSE	GYI	DGD	VVECTLHFGKFCVRTGK	VKAL-----PA
P0A185	YATDNLCTHGSA	RMSD	GYI	EGRE	IECPLHQGRFDVCTGK	ALC-----APV
P0A186	YATDNLCTHGSA	RMSD	GYI	EGRE	IECPLHQGRFDVCTGK	ALC-----APV
Q51493	YATDNLCTHGAA	RMSD	GFI	EGRE	IECPLHQGRFDVCTGR	ALC-----APV
A5W4F0	FAVQDTCTHGDW	ALSD	GYI	DGD	VVECTLHFGKFCVRTGK	VKAL-----PA
P0C620	FAVQDTCTHGDW	ALSD	GYI	DGD	VVECTLHFGKFCVRTGK	VKAL-----PA
P08086	FAVQDTCTHGDW	ALSD	GYI	DGD	VVECTLHFGKFCVRTGK	VKAL-----PA
Q52440	FATQDQCTHGEW	SLSE	GGY	LDGD	VVECSLHMGKFCVRTGK	-----V
Q7N4V8	FAVDDRCSHGNA	SISE	GYI	ED	NATVECPLHTASFCLRTGK	ALCL-----PA
P37332	FATQDRCTHGDW	SLSDGGYI	EGD	-----	VVECSLHMGKFCVRTGK	-----V
A7ZPY3	YAINDRCSHGNA	SMSE	GYI	EDD	ATVECPLHAASFCLKTGK	ALCL-----PA
P0ABW1	YAINDRCSHGNA	SMSE	GYI	EDD	ATVECPLHAASFCLKTGK	ALCL-----PA
A8A346	YAINDRCSHGNA	SMSE	GYI	EDD	ATVECPLHAASFCLKTGK	ALCL-----PA
P0ABW0	YAINDRCSHGNA	SMSE	GYI	EDD	ATVECPLHAASFCLKTGK	ALCL-----PA
P0ABW2	YAINDRCSHGNA	SMSE	GYI	EDD	ATVECPLHAASFCLKTGK	ALCL-----PA
Q3YZ13	YAINDRCSHGNA	SMSE	GYI	EDD	ATVECPLHAASFCLKTGK	ALCL-----PA
Q06458	YALDNLEPGSEAN	VLSR	GLI	GDAGGEP	IVISPLYKQRIRLRDG	-----



Core



Insertion



Deletion

The HMMer program

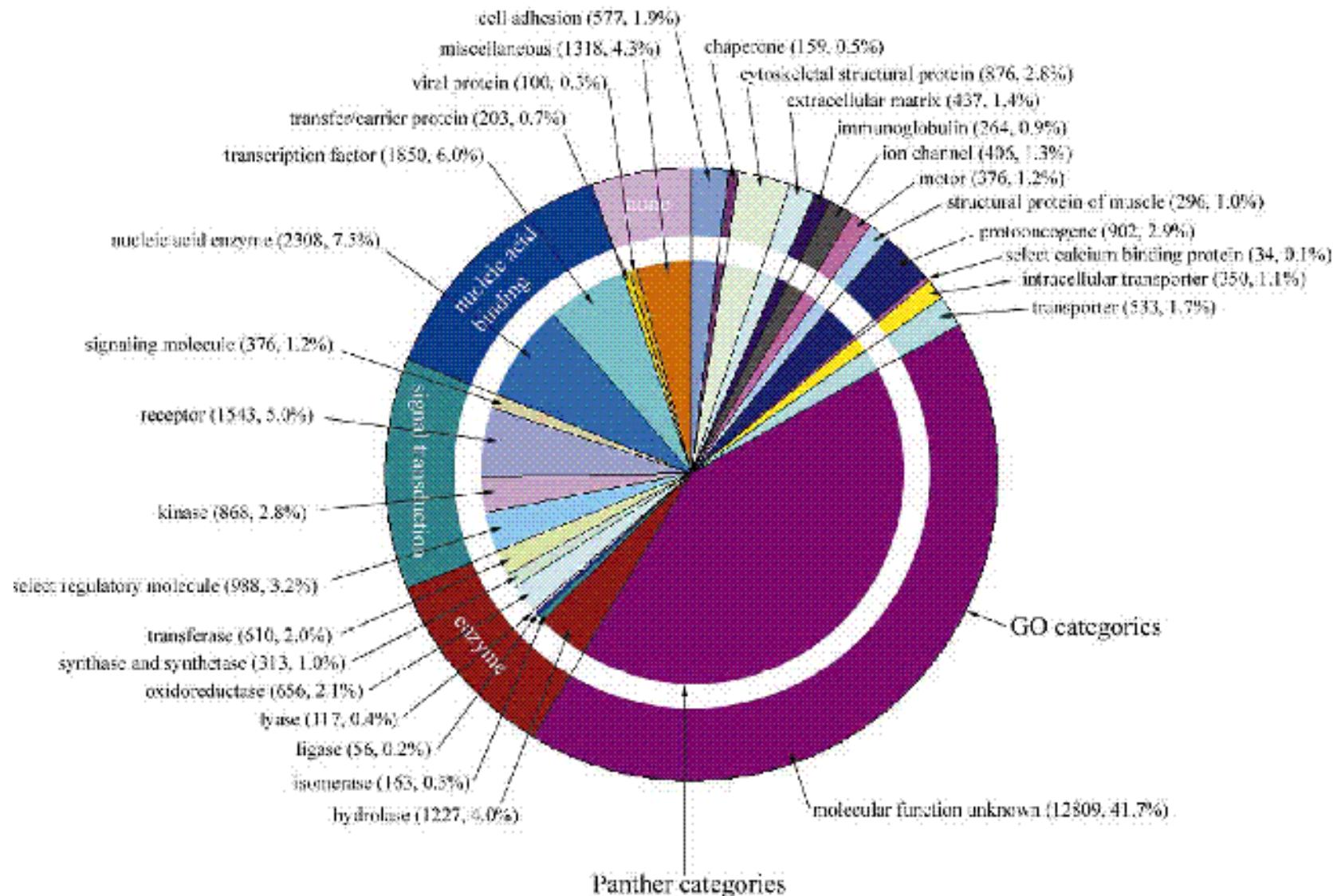
- HMMer is a open source program suite for profile HMM for biological sequence analysis
- Used to make the Pfam database of protein families
 - <http://pfam.sanger.ac.uk/>

Protein homology modeling and sequence profiles

Background. Why protein modeling?

- Because it works!
 - Close to 50% of all new sequences can be homology modeled
- Experimental effort to determine protein structure is very large and costly
- The gap between the size of the protein sequence data and protein structure data is large and increasing

Homology modeling and the human genome



Is it really impossible?

Protein homology modeling is only possible
if %id greater than 30-50%

WRONG!!!

Why %id is so bad!!

1200 models sharing 25-95% sequence identity with the submitted sequences (www.expasy.ch/swissmod)

Probabilities of SWISS-MODEL accuracy for target-template identity classes.

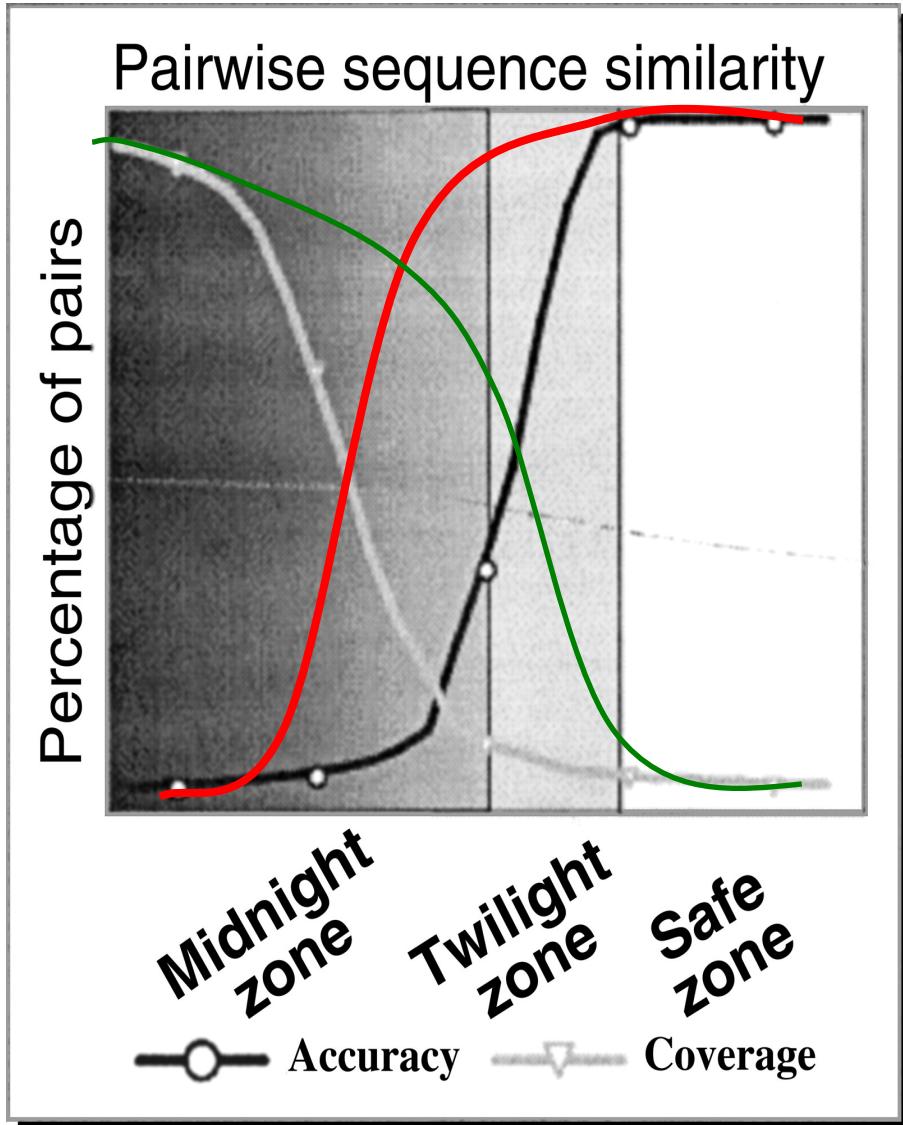
Percent sequence identity ^a	Total number of models ^b	Percent models with rmsd lower than 1 Å	Percent models with rmsd lower than 2 Å	Percent models with rmsd lower than 3 Å	Percent models with rmsd lower than 4 Å	Percent models with rmsd lower than 5 Å	Percent models with rmsd higher than 5 Å
25-29	125	0	10	30	46	67	33
30-39	222	0	18	45	66	77	23
40-49	156	9	44	63	78	91	9
50-59	155	18	55	79	86	91	9
60-69	145	38	72	85	91	92	8
70-79	137	42	71	82	85	88	12
80-89	173	45	79	86	94	95	5
90-95	88	59	78	83	86	91	9

a: Range of sequence identity between target and template sequence.

b: Total number of models in any given class of sequence identity. The table summarises 1201 model – control structure pairs.

c: Probability in percent that a model, sharing X% sequence identity with its template, deviates by 1 Å or less from the corresponding experimental control structure. The following columns provide these probabilities for other rms deviations.

Identification of fold



- If sequence similarity is high, proteins share structure (Safe zone)
- If sequence similarity is low, proteins may share structure (Twilight zone)
- Most proteins do not have a high sequence homologous partner

How can we do it?

- Identify template(s) - initial alignment
 - Can give you the protein function
- Improve alignment
 - Can give you the active site
- Backbone generation
- Loop modeling
 - Most difficult part
- Side chains
- Refinement
- Validation

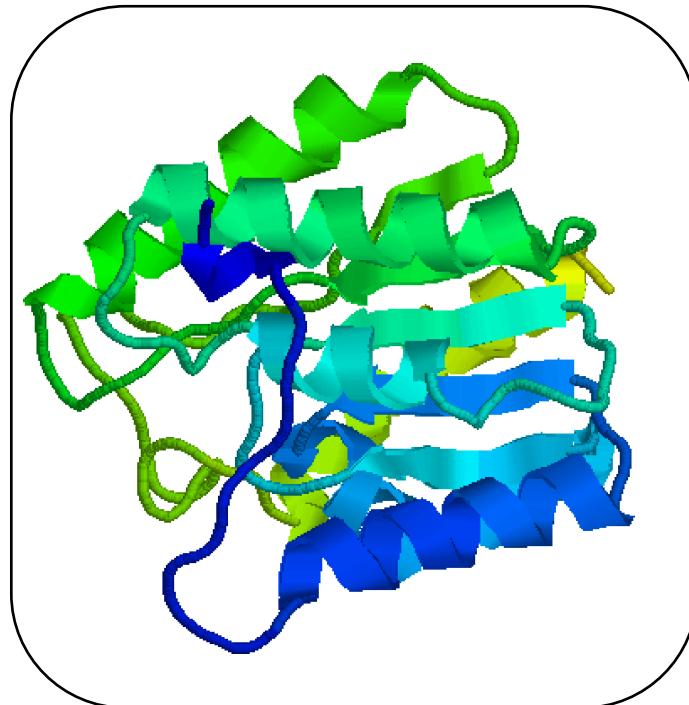
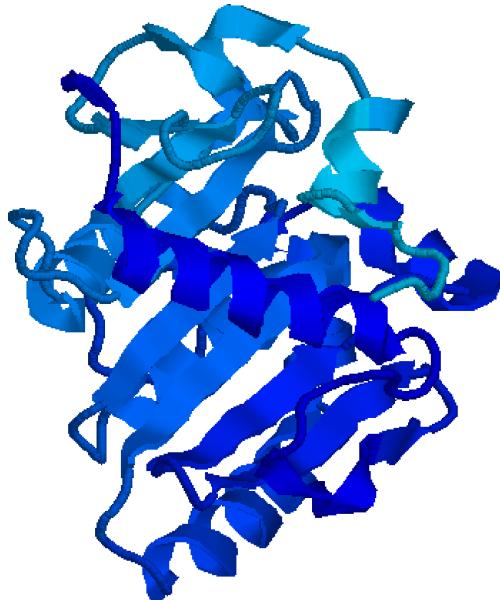
How to do it?

- Identify fold (template) for modeling
 - Find the structure in the PDB database that resembles your query protein the most
- Can be used to predict function
 - and maybe active sites
- Align protein sequence to template
 - Simple alignment methods
 - Sequence profiles
- Model side chains and loops

How to do it?

>1K7C.A

TTVYLAGDSTMAGNGGGSGTNGWGEYLASYLSATVVNDAVAGRSARSYTREGRFENIADV
VTAGDYVIVEFGHNDGGSLSTDNGRTDCSGTGAEVCYSVYDGVNTELTFPAYLENAAKL
FTAKGAKVILSSQTPNNPWETGTFVNSPTRFVEYAAEVAGVEYVDHWSYVDSIYETL
GNATVNSYFPIDHTHTSPAGAEVVAEAFLKAVVCTGTSVLTTSFEGTCL

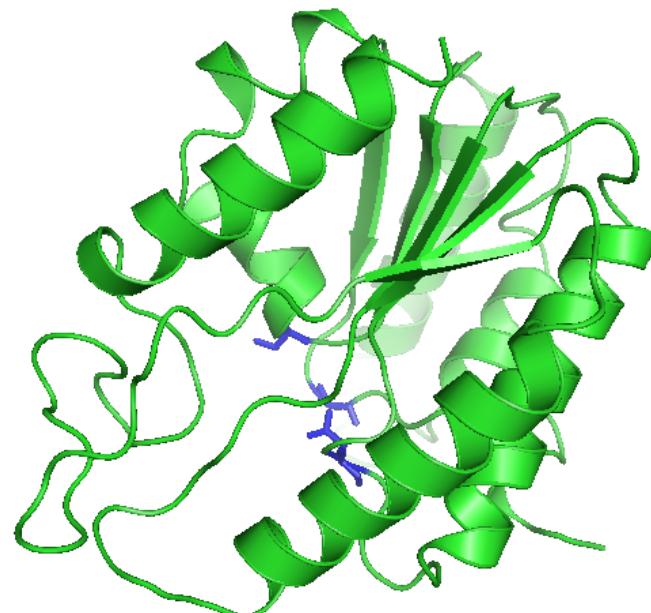


How to do it?

>1K7C.A

TTVYLAGDSTMAKNNGGSGTNGWGEYLASYLSATVVNDAVAGRSARSYTREGRFENIADV
VTAGDYVIVEFGHNDGGSLSTDNGRTDCSGTGAEVCYSVYDGVNTELTFPAYLENAAKL
FTAKGAKVILSSQTPNNPWETGTFVNSPTRFVEYAAEVAGVEYVDHWSYVDSIYETL
GNATVNSYFPIDHTHTSPAGAEVVAEAFLKAVVCTGTSKSVLTTSFEGTCL

1K7C.A TTVYLAGD**S**TMAKNNGGSGTNGWGEYLASYLSATVVNDAV**G**RSARSYTREGRFENI**N**ADVVTAGDYVIVEFGH**N**
1WAB._ EVVFIGD**S**LVQLMHQCE---IWRELFS---PLHALNFGIG**G**DSTQHVLW--RLЕНGELEHIRPKIVVVWVG**T****N**

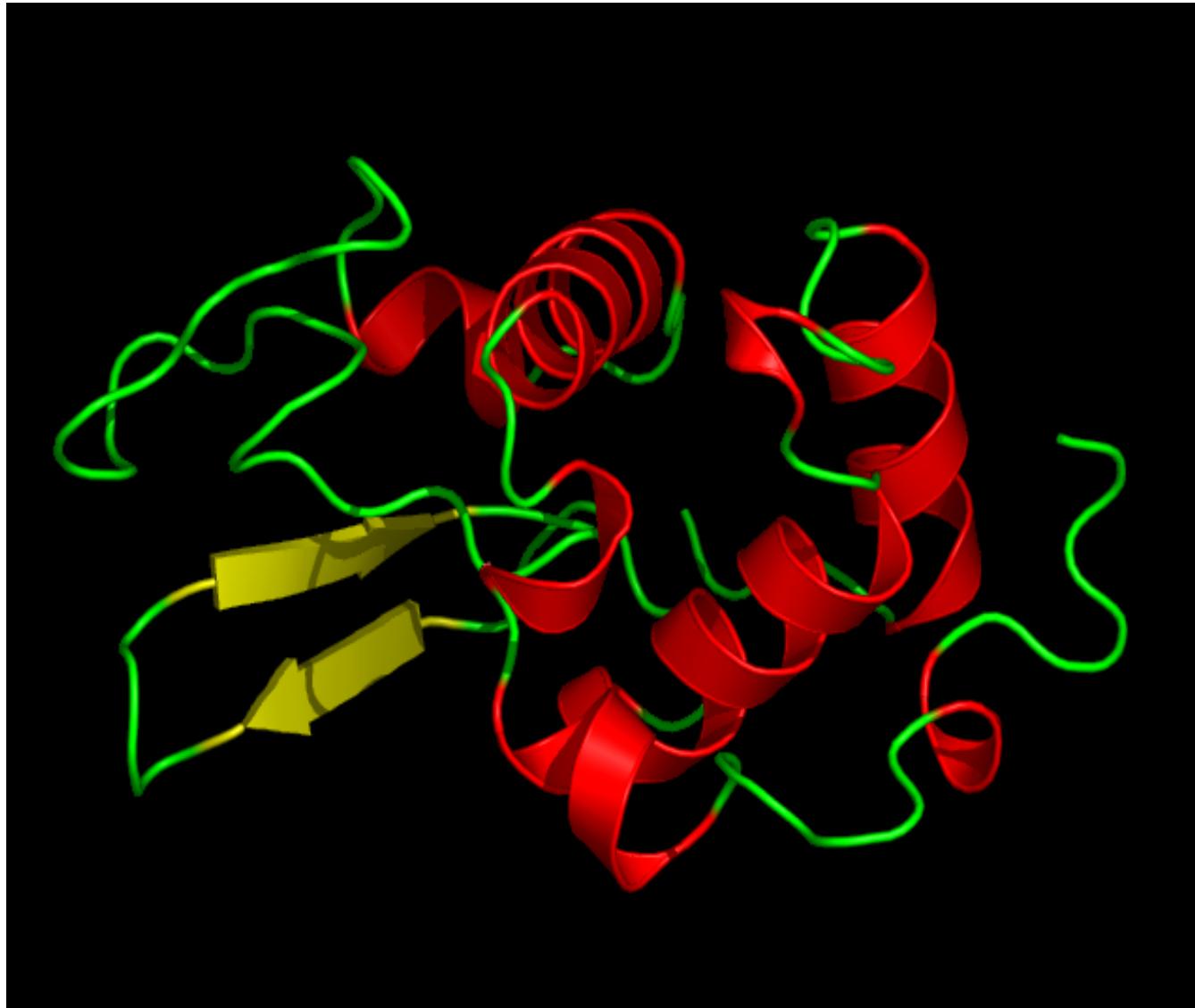


Including structure

- Sequence within in a protein superfamily share remote sequence similarity
- , but they share high structural similarity
- Structure is known for template
- Predict structural properties for query
 - Secondary structure
 - Surface exposure
- Position specific gap penalties derived from secondary structure and surface exposure

Predicting local Protein Structure

$\alpha+\beta$: Lysozyme (1jsf)



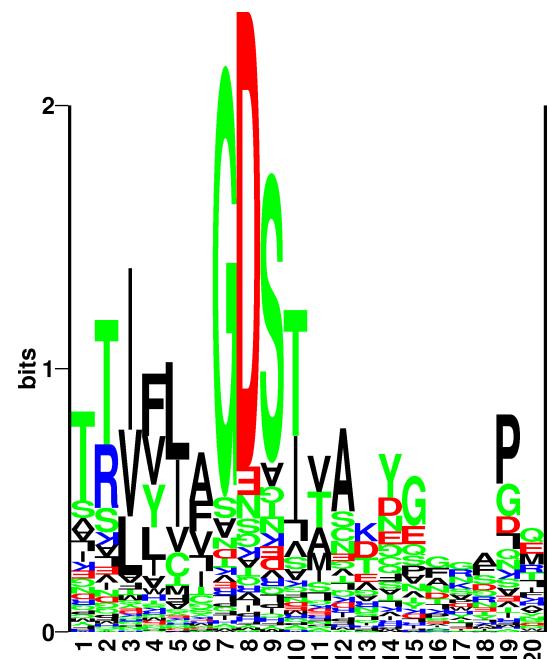
Improvement of accuracy

1974 Chou & Fasman	~50-53%
1978 Garnier	63%
1987 Zvelebil	66%
1988 Quian & Sejnowski	64.3%
1993 Rost & Sander	70.8-72.0%
1997 Frishman & Argos	<75%
1999 Cuff & Barton	72.9%
1999 Jones	76.5%
2000 Petersen et al.	77.9%

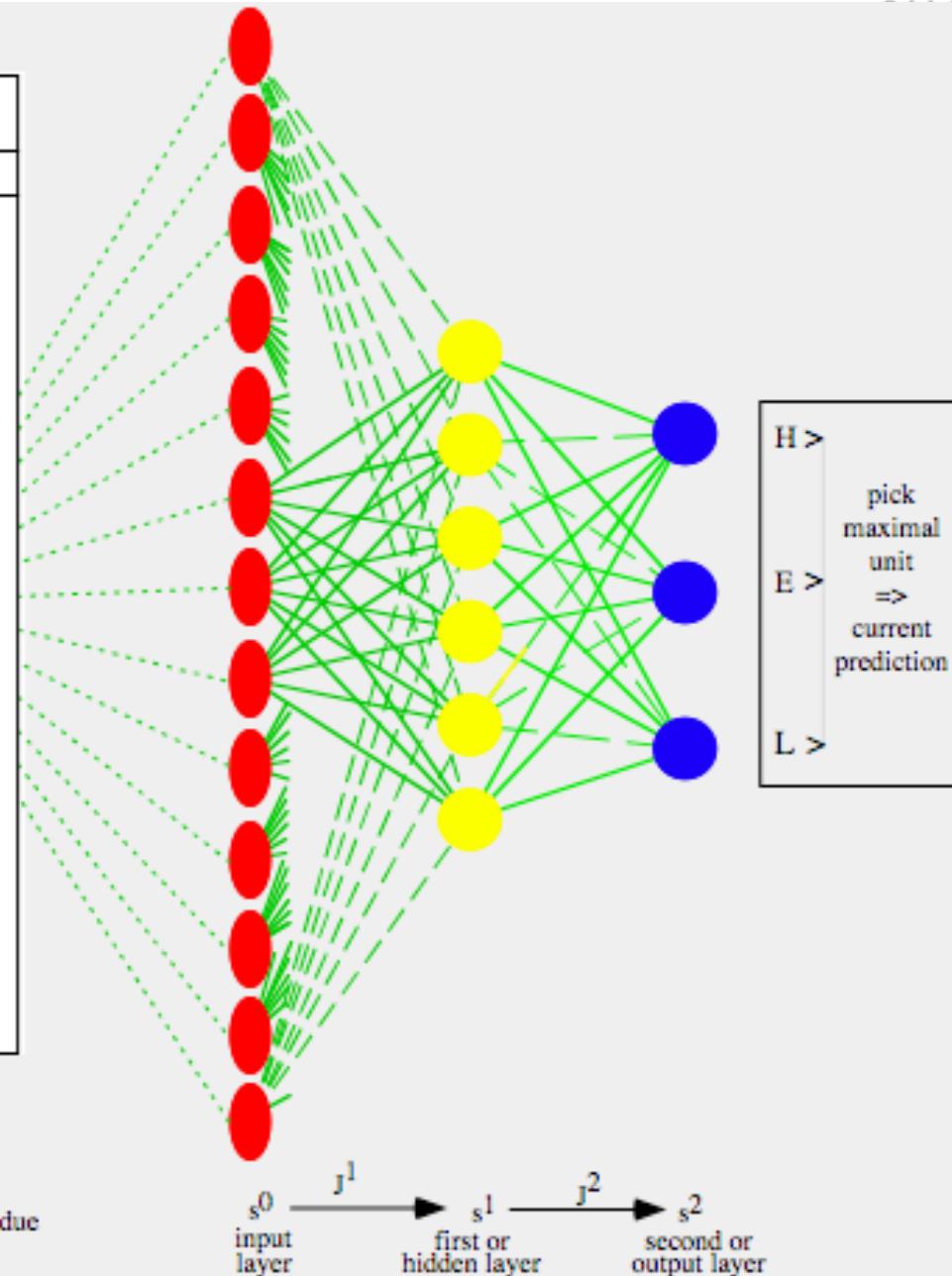
Sequence profiles



1	50
fyn_human	VTLFVALYDY EARTEDDLSF HKGEKFQILN SSEGDWWEAR SLTTGETQI
yrk_chick	VTLFIALYDY EARTEDDLSF QKGEKFHIIN NTEGDWWEAR SLSSGATQI
fgr_human	VTLFIALYDY EARTEDDLTF TKGEKFHILN NTEGDWWEAR SLSSGKTQI
yes_chick	VTVFVALYDY EARTTDDLSF KKGERFQIIN NTEGDWWEAR SIATGKTQI
src_avis2	VTTFVALYDY ESRTETDLSF KKGERLQIVN NTEGDWWLAH SLTTGQTQI
src_aviss	VTTFVALYDY ESRTETDLSF KKGERLQIVN NTEGDWWLAH SLTTGQTQI
src_avisr	VTTFVALYDY ESRTETDLSF KKGERLQIVN NTEGDWWLAH SLTTGQTQI
src_chick	VTTFVALYDY ESRTETDLSF KKGERLQIVN NTEGDWWLAH SLTTGQTQI
stk_hydat	VTIFVALYDY EARISEDLSF KKGERLQIIN TADGDWWYAR SLITNSEQI
src_rsvpa ESRIETDLSF KKGERLQIVN NTEGDWWLAH SLTTGQTQI
hck_human	..IVVALYDY EAIHHEDLSF QKGDQMVVLE ES.GEWWKAR SLATRKEQI
blk_mouse	..FVVALFDY AAVNDRDLQV LKGEKLQVLR .STGDDWWLAR SLVTGREGV
hck_mouse	.TIVVALYDY EAIHREDLSF QKGDQMVVLE .EAGEWWKAR SLATKKEQI
lyn_human	..IVVALYPY DGIHPDDLSF KKGEKMKVLE .EHGEWWKAK SLLTKKEQFI
lck_human	..LVIALHSY EPSHDGDLGF EKGEQIRILE QS.GEWWKAQ SLTTGQEQUI
ss81_yeast ALYPY DADDDdeISF EQNEILQVSD .IEGRWWKAR R.ANGETGII
abl_mouse	..LFVALYDF VASGDNTLSI TKGEKIRVLG YnnGEWCEAQ ..TKNGQQNV
abl1_human	..LFVALYDF VASGDNTLSI TKGEKIRVLG YnnGEWCEAQ ..TKNGQQNV
src1_drome	..VVSLYDY KSRDESDSL SF MKGDRMEVID DTESDWWRVV NLTRRQEQLI
mysd_dicdi ALYDF DAESSMELSF KEGDIITVLD QSSGDWWDAE L..KGRCKV
yfj4_yeast VALYSF AGEESGDPF RKGDVITILK ksQNDWWTGR V..NGREGIF
abl2_human	..LFVALYDF VASGDNTLSI TKGEKIRVLG YNQNGEWSEV RSKNG.QQNV
tec_human	.EIVVAMYDF QAAEGHDLRL ERGQEYLILE KNDVHWWRAR D.KYGNEGQI
abl1_caeel	..LFVALYDF HGVGEEQLSL RKGDQVRILG YNKNNEWCEA RlrLGEIGWV
txk_human ALYDF LPREPCNLAL RRAEYYLILE KYNPHWWKAR D.RLGNEQLI
yha2_yeast	VRRVVALYDL TTNEPDELSF RKGDVITVLE QVYRDWWKGA L..RGNMGIF
abp1_sacex AEYDY EAGEDNELTF AENDKIINIE FVDDDWLGE LETTGQKGLF

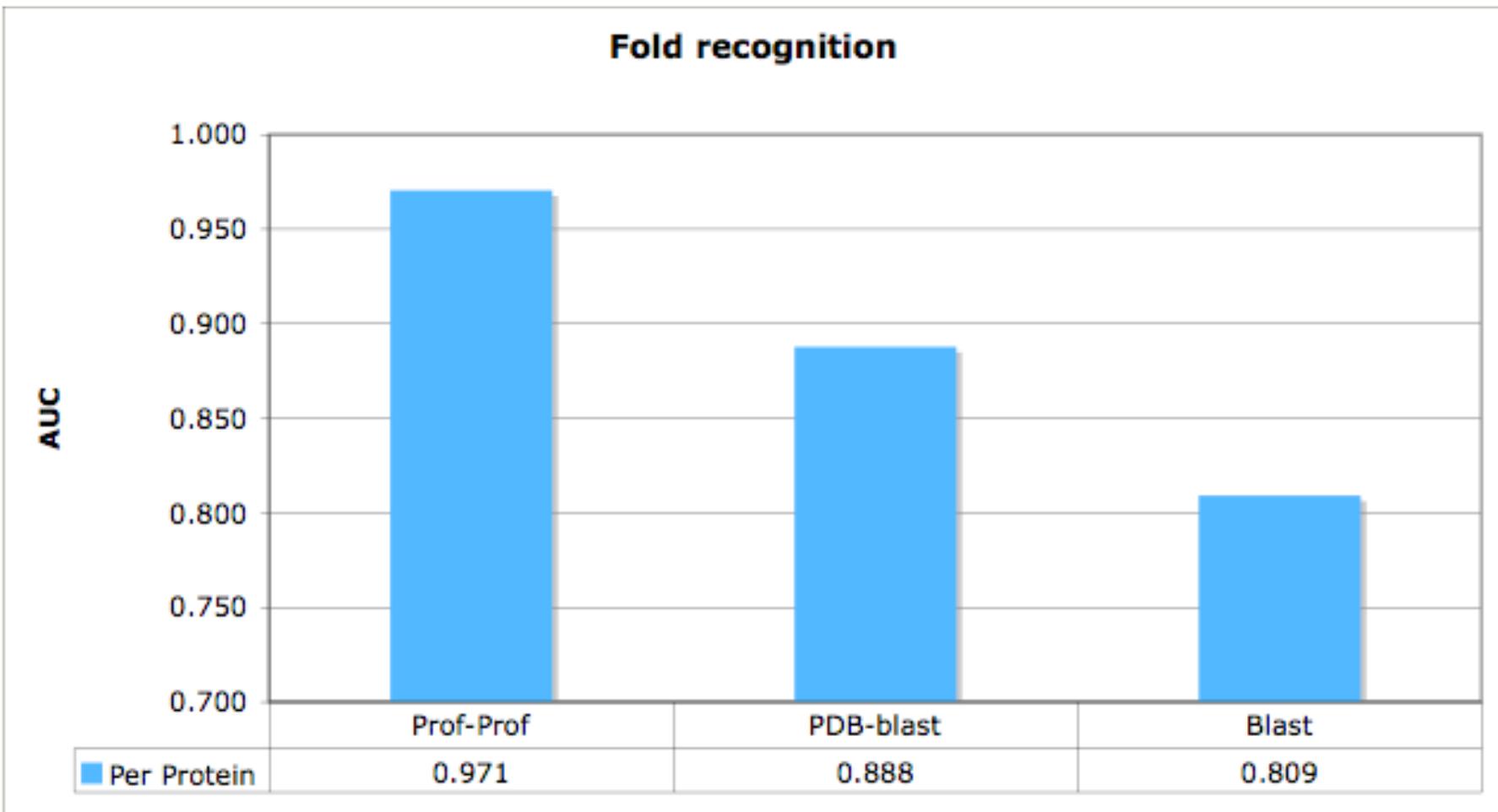


Protein	Alignments	profile table
:	:	GSAPD NTEKQ CVHIR LMYFW
G	GGGG	5.....
Y	YYYY5..
I	IIEE2...3..
Y	YYYY5..
D	DDDD5.....
P	PPPP5..
E	AEAA	..3..2..
D	VVEE1..2..2..
G	GGGG	5.....
D	DDDD5.....
P	PPPP5..
D	DTDD4..1..
D	NQNN1..3...1
G	GNGG	4....1....
V	VIVV4..1..
N	EPKK1..1..12..
P	PPPP5..
G	GGGG	5.....
T	TTTT5..
D	EKSA	..1..1..11..
F	FFFF5..
:	:



corresponds to the the 21×3 bits coding for the profile of one residue

CpHModel - Fold recognition performance



CpHModels 3.2

CENTERFO RBIOLOGI CALSEQU ENCEANA LYSIS CBS	EVENTS	NEWS	RESEARCH GROUPS	CBS PREDICTION SERVERS	CBS DATA SETS	PUBLICATIONS	EDUCATION
	STAFF	CONTACT	ABOUT CBS	INTERNAL	CBS BIOINFORMATICS TOOLS	CBS COURSES	OTHER BIOINFORMATICS LINKS

[CBS](#) >> [CBS Prediction Servers](#) >> CPHmodels

CPHmodels 3.2 Server

CPHmodels 3.2 is a protein homology modeling server. The template recognition is based on profile-profile alignment guided by secondary structure and exposure predictions.

New in version 3.2: An improvement in the alignment algorithm in case of remote homology modeling where a structure dependant gap penalty has been introduced. Also, there are changes in the output format. A summary line has been included to make parsing easier.

View the [version history](#) of this server.

[Instructions](#)

[Output format](#)

[Article abstract](#)

SUBMISSION

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

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

no file selected

Restrictions:

Only one sequence per submission with not less than 15 and not more than 4,000 amino acids.

Confidentiality:

The sequences are kept confidential and will be deleted after processing.

Take home message

- Identifying the correct fold is only a small step towards successful homology modeling
 - You can do reliable fold recognition AND homology modeling even for low sequence homology
 - Use sequence profiles and local protein structure (predictions) to align sequences
-