

# Protein Evolution

## Lecture 2

Lucy Colwell

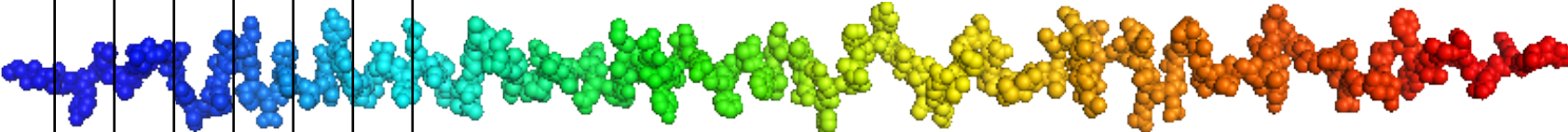
ICTS Winter School on Quantitative Systems Biology

Dec 13 2017

# How to extract useful information from protein sequences?

Each protein is constructed as a specific chain or string of amino acids. There are 20 different amino acids which are used with roughly equal frequencies across all proteins.

In a specific protein of interest, such as hemoglobin, the order of amino acids is highly important and contains all the information necessary to produce the folded, functional molecule.

1	2	3	4	5	6	7	8 .....
M	V	L	S	P	A	D	K .....
							

We would like to find a probability model for the sequence of amino acids that corresponds to each protein of interest.

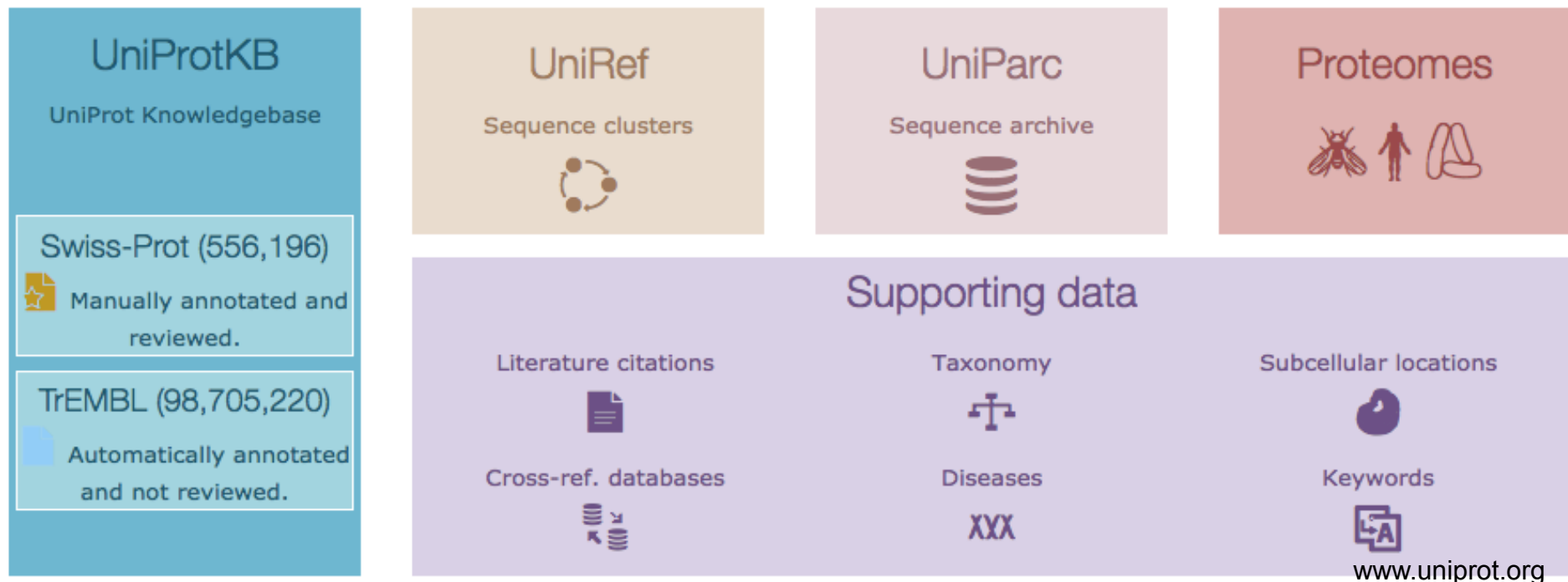
$$P(A_1, \dots, A_L) = \text{Probability that a sequence produces a folded, functional hemoglobin molecule.}$$

# Sequence Similarity

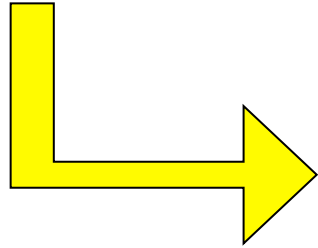
$P(A_1, \dots, A_L)$  = Probability that a sequence produces a folded, functional hemoglobin molecule.

To proceed, we need data – if we can find lots of examples of hemoglobin sequences, then we can treat this as an inverse problem, and look for a model that reproduces the statistics of our observed data.

How can we find this data? Over the last 50 years there has been an explosion of technological innovation, which has enabled us to determine the sequences of huge numbers of biological organisms, and thus of numerous proteins.




# Protein evolution

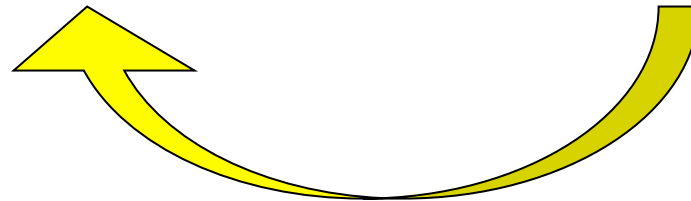


## Families of homologous proteins

- Similar/dissimilar sequence
- Common 3D structure
- Common function

Protein  Known family

VCVEVPSETEA...

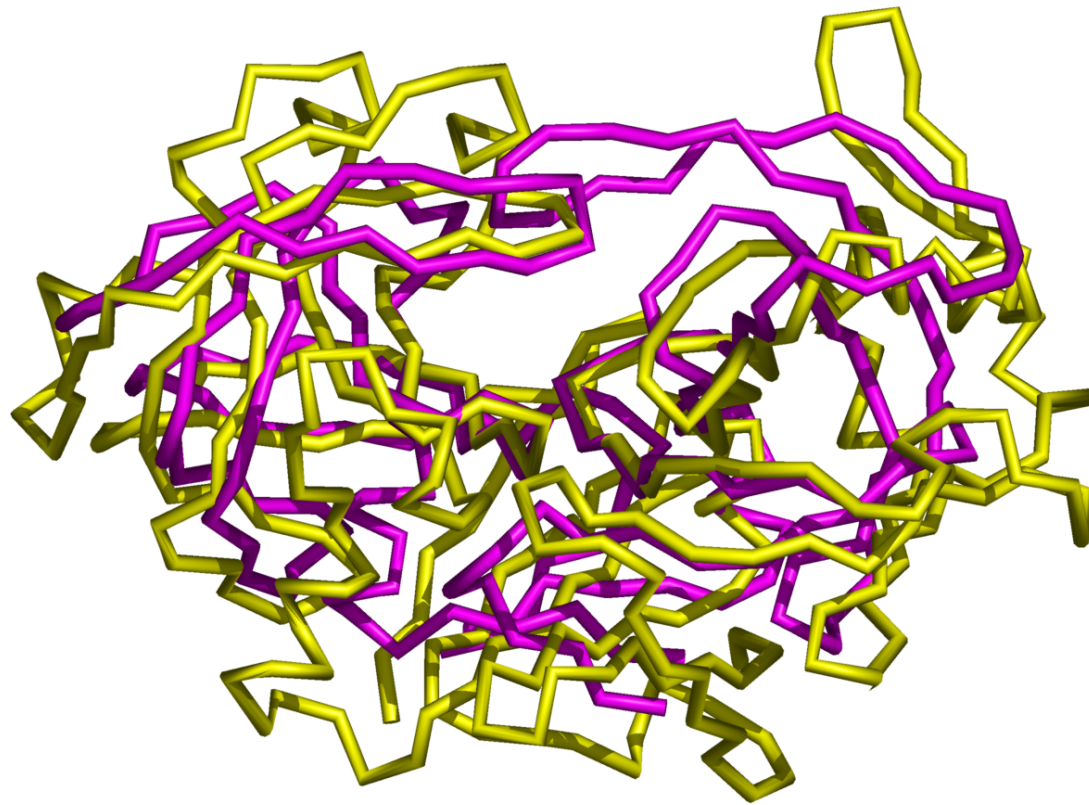


*3D structure, function*

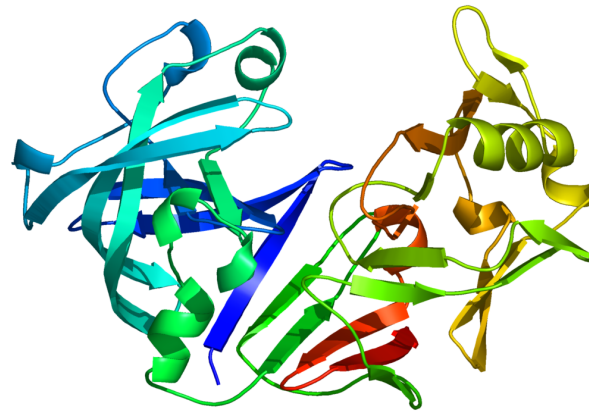




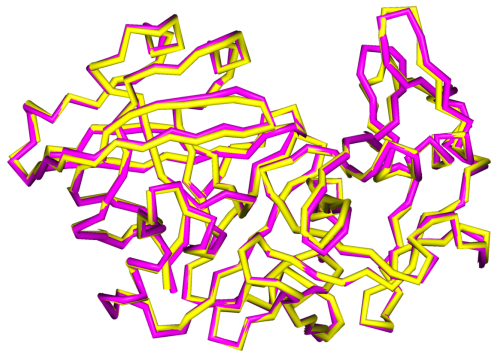
Similar or dissimilar?



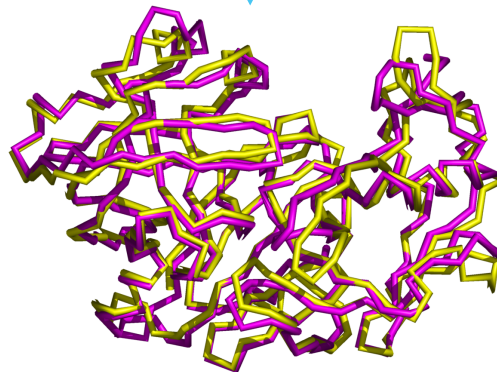
To compare proteins, we need to be able to measure similarity between their sequence and their 3D structures.



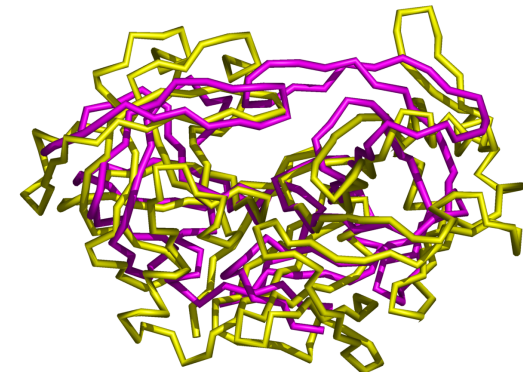
porcine pepsin



human pepsin 3A  
(PID 86%, RMSD 0.8 Å)



mouse renin  
(PID 43%, RMSD 1.6 Å)



HIV protease  
(PID 13%, RMSD 2.6 Å)

PID = percentage sequence identity, RMSD = root mean squared deviation

## Pairwise sequence alignment

To group sequences into families, we need to be able to answer the question: are two sequences related? This is the most basic sequence analysis task that we will come across. We approach this by first aligning the sequences, and then asking whether the alignment suggests that the sequences are related, or could have occurred by chance.

## Pairwise sequence alignment

To group sequences into families, we need to be able to answer the question: are two sequences related? This is the most basic sequence analysis task that we will come across. We approach this by first aligning the sequences, and then asking whether the alignment suggests that the sequences are related, or could have occurred by chance.

Consider this alignment between parts of the human haemoglobin alpha and beta chains:

HBA_HUMAN	GSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKL
	G+ +VK+HGKKV A++++AH+D++ +++++LS+LH KL
HBB_HUMAN	GNPKVKAHGKKVLGAFSDGLAHLNFKGTFTLSELHCDKL

Here identical positions are shown in the middle with letters, while similar positions are shown with plus signs. To define 'similar' we use a substitution matrix – more on this soon.

Note there are many positions at which the two residues are identical, and many others that we call 'functionally conservative' – for example the D-E pair towards the end (both are negatively charged amino acids).

## Pairwise sequence alignment

To group sequences into families, we need to be able to answer the question: are two sequences related? This is the most basic sequence analysis task that we will come across. We approach this by first aligning the sequences, and then asking whether the alignment suggests that the sequences are related, or could have occurred by chance.

Consider this alignment between parts of the human haemoglobin alpha and beta chains:

```
HBA_HUMAN  GSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKL
              G+ +VK+HGKKV  A+++++AH+D++ ++++++LS+LH  KL
HBB_HUMAN  GNPVKKAHGKKVLGAFSDGLAHLNKGTFATLSELHCDKL
```

Here identical positions are shown in the middle with letters, while similar positions are shown with plus signs. To define 'similar' we use a substitution matrix – more on this soon.

```
HBA_HUMAN  GSAQVKGHGKKVADALTNAVAHV---D---DMPNALSALSDLHAHKL
              ++ +++++H+ KV    + +A  ++                +L+ L+++H+ K
LGB2_LUPLU NNPELQAHAGKVFKLVYEAAIQLQVTGVVVTDATLKNLGSVHVS KG
```

In contrast the alignment with leghemoglobin from yellow lupin is much weaker, although these two proteins are evolutionarily related, and have the same 3D structures.

## Pairwise sequence alignment

To group sequences into families, we need to be able to answer the question: are two sequences related? This is the most basic sequence analysis task that we will come across. We approach this by first aligning the sequences, and then asking whether the alignment suggests that the sequences are related, or could have occurred by chance.

Consider this alignment between parts of the human haemoglobin alpha and beta chains:

```
HBA_HUMAN  GSAQVKGHGKKVADALTNAVAHVDDMPNALSASDLHAHKL
              G+ +VK+HGKKV  A+++++AH+D++ +++++LS+LH  KL
HBB_HUMAN  GNPVKVKAHGKKVLGAFSDGLAHLNKGTFATLSELHCDKL
```

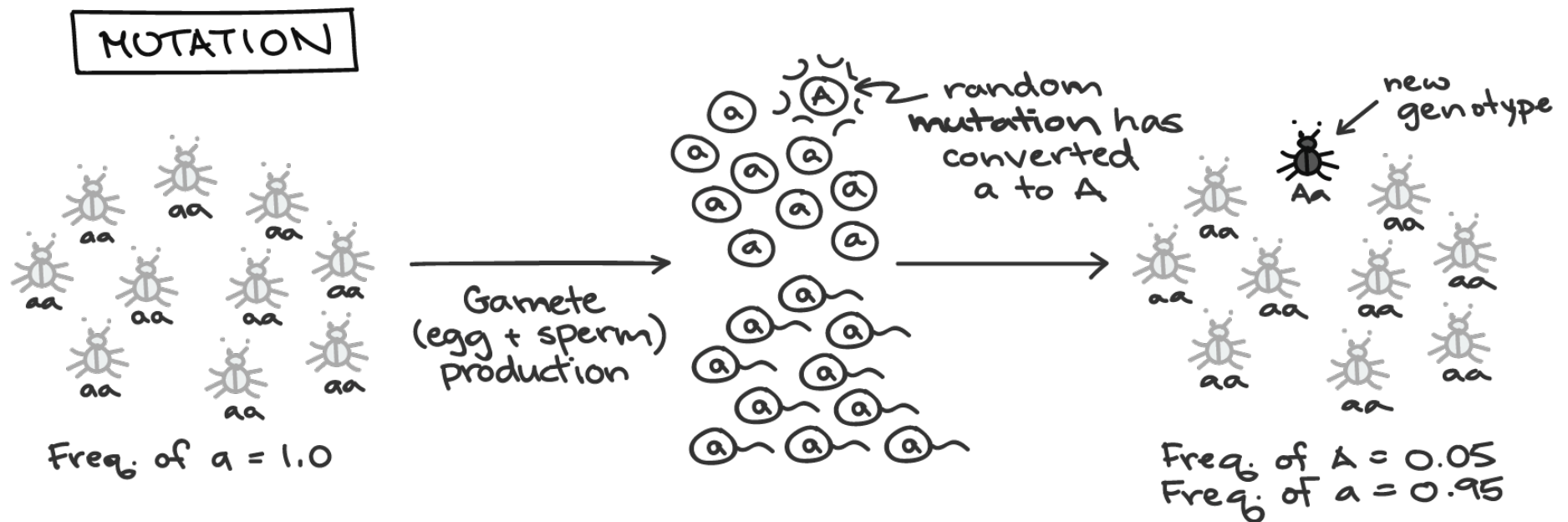
Here identical positions are shown in the middle with letters, while similar positions are shown with plus signs. To define 'similar' we use a substitution matrix – more on this soon.

```
HBA_HUMAN  GSAQVKGHGKKVADALTNAVAHVDDMPNALSASD----LHAHKL
              GS+ + G +   +D L  ++ H+ D+  A +AL D    ++AH+
F11G11.2   GSGYLVGDSLTFVDLL--VAQHTADLLAANAALLDEFPQFKAHQE
```

Finally this alignment with a nematode glutathione S-transferase homologue has a similar number of identities and similarities, but the structure and function are completely different.

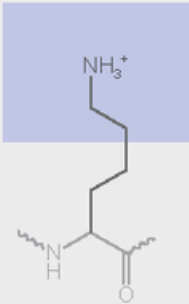
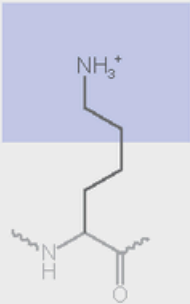
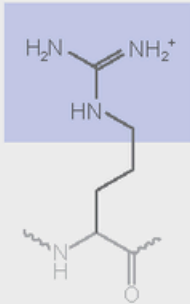
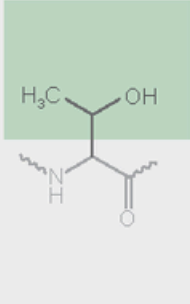
## Scoring models

The purpose of a scoring algorithm is to assess the evidence that two sequences have diverged from a common ancestor by a process of mutation and selection.



# Scoring models

The purpose of a scoring algorithm is to assess the evidence that two sequences have diverged from a common ancestor by a process of mutation and selection.

	Point mutations				
	No mutation	Silent	Nonsense	Missense	
				conservative	non-conservative
DNA level	TTC	TTT	ATC	TCC	TGC
mRNA level	AAG	AAA	UAG	AGG	ACG
protein level	Lys	Lys	STOP	Arg	Thr
					
	basic			polar	



# Scoring models

DNA level

tcatatgcaccccggt

S Y A P R

Process

Substitution

Indel

Type of mutation

Synonymous Substitution

Non-synonymous Substitution

Insertion

Protein sequence level

tcatatgca**cc**acgt  
S Y A **P** R

tcatatgca**tat**cgt  
S Y A **Y** R

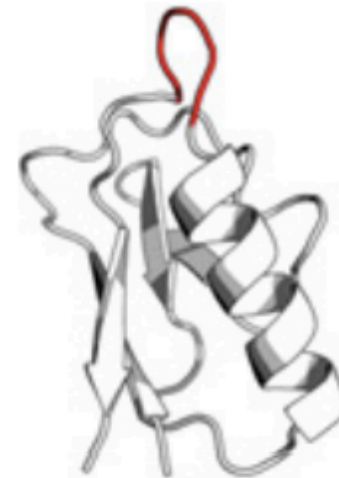
tcatat**acgtatgc**accacgt  
S Y **T Y A** P R

Residue conserved

Residue change

Embellishment

Protein structure level



Change on function

Possible, but very rare

Context dependent

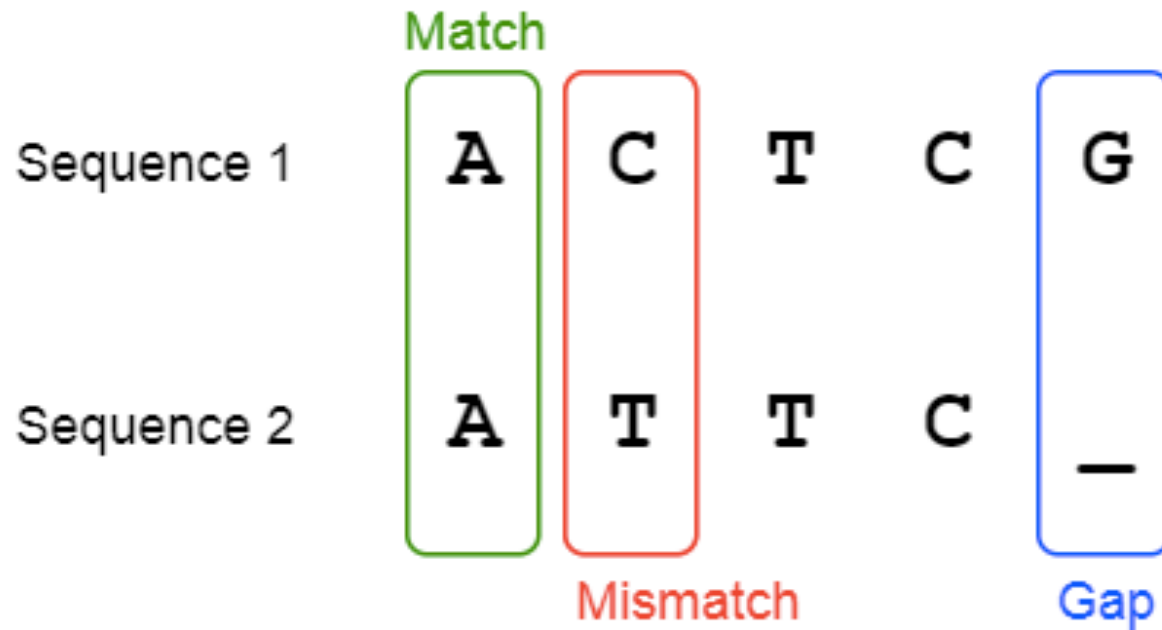
Context dependent

## Scoring models

Total score of an alignment will contain a term for each aligned pair of residues, plus terms for each gap introduced. This will correspond to the logarithm of the relative likelihood that the sequences are related, compared to being unrelated.

## Scoring models

Total score of an alignment will contain a term for each aligned pair of residues, plus terms for each gap introduced. This will correspond to the logarithm of the relative likelihood that the sequences are related, compared to being unrelated.



For example,  $S(A, A) = S(T, T) = S(C, C) = +5$

$S(C, T) = -3$

$S(\text{gap}) = ? -8?$

# Substitution Matrices

We need score terms for each aligned residue pair. The intuition for proteins alluded to in these slides could yield 210 scoring terms for all possible pairs of amino acids, but it is useful to have a guiding theory for what the scores mean.

# Substitution Matrices

We need score terms for each aligned residue pair. The intuition for proteins alluded to in these slides could yield 210 scoring terms for all possible pairs of amino acids, but it is useful to have a guiding theory for what the scores mean.

Consider two protein sequences  $x$  and  $y$ , we want to compare the random model:

$x_1, x_2, x_3, x_4, x_5$   
 $y_1, y_2, y_3, y_4, y_5$

$$P(x, y | R) = \prod_i q_{x_i} \prod_j q_{y_j}$$

with the match model:

$$P(x, y | M) = \prod_i p_{x_i y_i}$$

# Substitution Matrices

We need score terms for each aligned residue pair. The intuition for proteins alluded to in these slides could yield 210 scoring terms for all possible pairs of amino acids, but it is useful to have a guiding theory for what the scores mean.

Consider two protein sequences  $x$  and  $y$ , we want to compare the random model:

$$\begin{array}{l} x_1, x_2, x_3, x_4, x_5 \\ y_1, y_2, y_3, y_4, y_5 \end{array} \quad P(x, y | R) = \prod_i q_{x_i} \prod_j q_{y_j}$$

with the match model:

$$P(x, y | M) = \prod_i p_{x_i y_i}$$

So we want

$$\frac{P(x, y | M)}{P(x, y | R)} = \frac{\prod_i p_{x_i y_i}}{\prod_i q_{x_i} \prod_i q_{y_i}} = \prod_i \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}}$$

This leads to the score  $S = \sum_i s(x_i, y_i)$ , where  $s(a, b) = \log \left( \frac{p_{ab}}{q_a q_b} \right)$  is the log likelihood ratio that (a,b) is an aligned or match pair, vs occurred at random.

## Substitution Matrices – BLOSUM 50

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	5	-2	-1	-2	-1	-1	-1	0	-2	-1	-2	-1	-1	-3	-1	1	0	-3	-2	0
R	-2	7	-1	-2	-4	1	0	-3	0	-4	-3	3	-2	-3	-3	-1	-1	-3	-1	-3
N	-1	-1	7	2	-2	0	0	0	1	-3	-4	0	-2	-4	-2	1	0	-4	-2	-3
D	-2	-2	2	8	-4	0	2	-1	-1	-4	-4	-1	-4	-5	-1	0	-1	-5	-3	-4
C	-1	-4	-2	-4	13	-3	-3	-3	-3	-2	-2	-3	-2	-2	-4	-1	-1	-5	-3	-1
Q	-1	1	0	0	-3	7	2	-2	1	-3	-2	2	0	-4	-1	0	-1	-1	-1	-3
E	-1	0	0	2	-3	2	6	-3	0	-4	-3	1	-2	-3	-1	-1	-1	-3	-2	-3
G	0	-3	0	-1	-3	-2	-3	8	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4
H	-2	0	1	-1	-3	1	0	-2	10	-4	-3	0	-1	-1	-2	-1	-2	-3	2	-4
I	-1	-4	-3	-4	-2	-3	-4	-4	-4	5	2	-3	2	0	-3	-3	-1	-3	-1	4
L	-2	-3	-4	-4	-2	-2	-3	-4	-3	2	5	-3	3	1	-4	-3	-1	-2	-1	1
K	-1	3	0	-1	-3	2	1	-2	0	-3	-3	6	-2	-4	-1	0	-1	-3	-2	-3
M	-1	-2	-2	-4	-2	0	-2	-3	-1	2	3	-2	7	0	-3	-2	-1	-1	0	1
F	-3	-3	-4	-5	-2	-4	-3	-4	-1	0	1	-4	0	8	-4	-3	-2	1	4	-1
P	-1	-3	-2	-1	-4	-1	-1	-2	-2	-3	-4	-1	-3	-4	10	-1	-1	-4	-3	-3
S	1	-1	1	0	-1	0	-1	0	-1	-3	-3	0	-2	-3	-1	5	2	-4	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	2	5	-3	-2	0
W	-3	-3	-4	-5	-5	-1	-3	-3	-3	-3	-2	-3	-1	1	-4	-4	-3	15	2	-3
Y	-2	-1	-2	-3	-3	-1	-2	-3	2	-1	-1	-2	0	4	-3	-2	-2	2	8	-1
V	0	-3	-3	-4	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	0	-3	-1	5

Henikoff, Steven, and Jorja G. Henikoff. "Amino acid substitution matrices from protein blocks." PNAS 1992.

## **Alignment algorithms**

To construct the BLOSUM matrix a set of aligned, ungapped regions from protein families were assembled and clustered such that two sequences go in the same cluster if their percentage of identical residues exceeds 50%.



## Alignment algorithms

To construct the BLOSUM matrix a set of aligned, ungapped regions from protein families were assembled and clustered such that two sequences go in the same cluster if their percentage of identical residues exceeds 50%.

Once we have a scoring model, we then need a way to find the optimal alignment between two sequences. If we allow gaps, then there are

$$\binom{2n}{n} = \frac{(2n)!}{(n!)^2} \simeq \frac{2^{2n}}{\sqrt{\pi n}}$$

possible global alignments between two sequences of length  $n$ . This means it is not computationally feasible to enumerate all of these, even for moderate values of  $n$ .

# Alignment algorithms

To construct the BLOSUM matrix a set of aligned, ungapped regions from protein families were assembled and clustered such that two sequences go in the same cluster if their percentage of identical residues exceeds 50%.

Once we have a scoring model, we then need a way to find the optimal alignment between two sequences. If we allow gaps, then there are

$$\binom{2n}{n} = \frac{(2n)!}{(n!)^2} \simeq \frac{2^{2n}}{\sqrt{\pi n}}$$

possible global alignments between two sequences of length  $n$ . This means it is not computationally feasible to enumerate all of these, even for moderate values of  $n$ .

Luckily, dynamic programming algorithms provide a way of efficiently finding the optimal alignment. The idea is to start from an existing smaller alignment, and evaluate all possible next moves.

Global alignment: Needleman-Wunsch algorithm – align whole sequences.

Local alignment: Smith-Waterman algorithm – align subsequences.

## Next – multiple sequence alignment

```

Helix          AAAAAAAAAAAAAAAAAA  BBBBBBBBBBBBBBBBBBCCCCCCCCCCCC
HBA_HUMAN      -----VLSPADKTNVKAAWGKVGA--HAGEYGAEALERMFLSFPTTKTYFPHF
HBB_HUMAN      -----VHLTPEEKSAVTALWGKV---NVDEVGGEALGRLLVVYPWTQRFFESF
MYG_PHYCA      -----VLSEGEWQLVLHVWAKVEA--DVAGHGQDILIRLFKSHPETLEKFDRF
GLB3_CHITP     -----LSADQISTVQASFDKVKG-----DPVGILYAVFKADPSIMAKFTQF
GLB5_PETMA     PIVDTGSVAPLSAAEKTIRSAPVYS--TYETSGVDILVKFFTSTPAAQEFPKPF
LGB2_LUPLU     -----GALTESQAALVKSSWEEFNA--NIPKHTHRFFILVLEIAPAADKDLFS-F
GLB1_GLYDI     -----GLSAAQRQVIAATWKDIAGADNGAGVGKDCIKFLSAHPQMAAVFG-F
Consensus      Ls...  v a W kv . .   g . L.. f . P .   F F

```

```

Helix          DDDDDDDDEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE  FFFFFFFFFFFFFFFF
HBA_HUMAN      -DLS-----HGSAQVKGHGKKVADALTNVAHV--D--DMPNALSALSDLHAHKL-
HBB_HUMAN      GDLSTPDVAVMGNPKVKAHGGKVLGAFSDGLAHL---D--NLKGTFFATLSELHCDKL-
MYG_PHYCA      KHLKTEAEMKASEDLKKHGVTVLTALGAILKK---K-GHHEAELKPLAQSHATKH-
GLB3_CHITP     AG-KDLESIKGTAPFETHANRIVGFFSKIIGEL--P---NIEADVNTFVASHKPRG-
GLB5_PETMA     KGLTTADQLKKSADVRWHAERIINAVNDAVASM--DDTEKMSMKLRDLGKHAASF-
LGB2_LUPLU     LK-GTSEVPQNNPELQAHAGKVFKLVYEAQIQLQVTGVVVTDATLKNLGSVHVSKG-
GLB1_GLYDI     SG----AS---DPGVAALGAKVLAQIGVAVSHL--GDEGKMVAQMKAVGVRHKGYN
Consensus      .  t      .. . v..Hg kv. a   a...l   d   . a l. l   H .

```

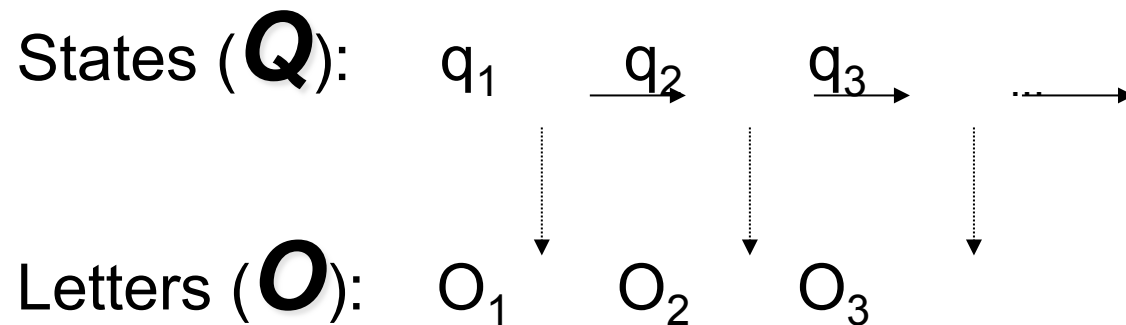
```

Helix          FFGGGGGGGGGGGGGGGGGGGGGGGGGGGGG  HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
HBA_HUMAN      -RVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR-----
HBB_HUMAN      -HVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVAVAGVANALAHKYH-----
MYG_PHYCA      -KIPIKYLEFISEAIIHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG
GLB3_CHITP     --VTHDQLNNFRAGFVSVMKAHT--DFA-GAEAAGATLDTFFGMIFSKM-----
GLB5_PETMA     -QVDPQYFKVLAAVIADTVAAG-----DAGFEKLSMICILLRSAY-----
LGB2_LUPLU     --VADAHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDAA---
GLB1_GLYDI     KHIKAQYFEPLGASLLSAMEHRIGGKMNAAKDAWAAAYADISGALISGLQS-----
Consensus      v.   f  l . . . . .   f   . aa. k. .   l sky

```

# Hidden Markov Models

A **Hidden Markov Model** (HMM) is a discrete-time finite-state Markov chain coupled with a sequence of letters emitted when the Markov chain visits its states.



The sequence **O** of emitted letters is called “the observed sequence” because we often know it while not knowing the state sequence **Q**, which we call “hidden”.

Used extensively throughout computational biology for ‘labeling’ data – and in particular for classifying and aligning protein sequence data.

# What is a hidden Markov model?

Sean R Eddy

**NATURE BIOTECHNOLOGY** VOLUME 22 NUMBER 10 OCTOBER 2004

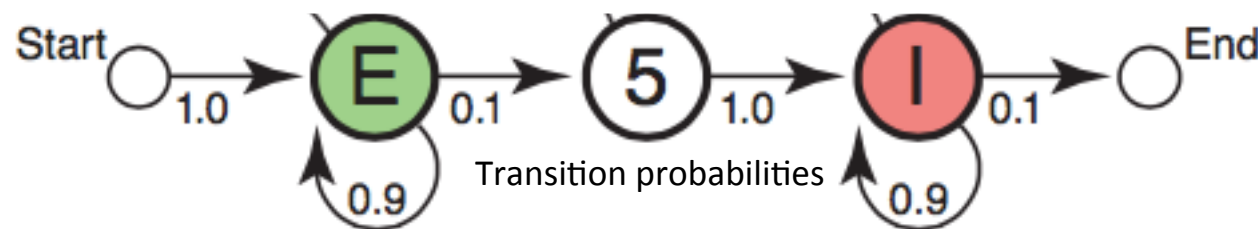
As a simple example, imagine the following 5' splice-site recognition problem. Assume we are given a DNA sequence that begins in an exon, contains one 5' splice site and ends in an intron. The problem is to identify where the switch from exon to intron occurred—where the 5' splice site (5'SS) is.

# What is a hidden Markov model?

Sean R Eddy

NATURE BIOTECHNOLOGY VOLUME 22 NUMBER 10 OCTOBER 2004

As a simple example, imagine the following 5' splice-site recognition problem. Assume we are given a DNA sequence that begins in an exon, contains one 5' splice site and ends in an intron. The problem is to identify where the switch from exon to intron occurred—where the 5' splice site (5'SS) is.



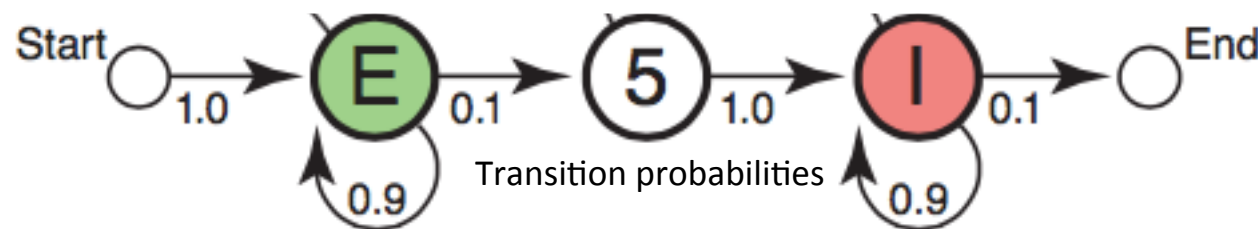
HMM with three states – exon, 5'SS and intron, and the probabilities of moving between these states.

# What is a hidden Markov model?

Sean R Eddy

NATURE BIOTECHNOLOGY VOLUME 22 NUMBER 10 OCTOBER 2004

As a simple example, imagine the following 5' splice-site recognition problem. Assume we are given a DNA sequence that begins in an exon, contains one 5' splice site and ends in an intron. The problem is to identify where the switch from exon to intron occurred—where the 5' splice site (5'SS) is.



HMM with three states – exon, 5'SS and intron, and the probabilities of moving between these states. To estimate which state each element of the DNA sequence is in, we need information about the statistics of exon, splice site and intron sequences – these are the emission probabilities.

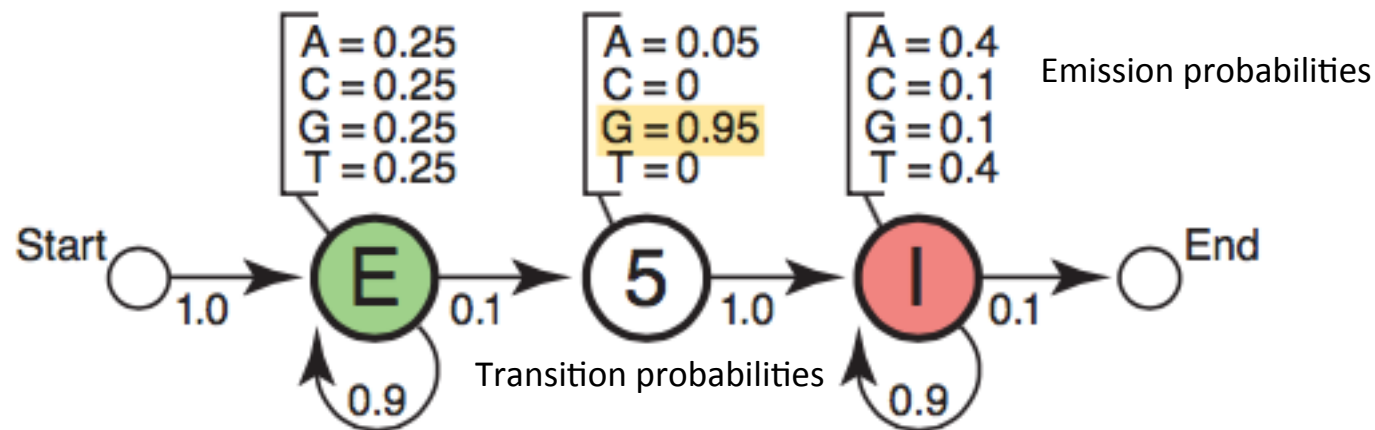


# What is a hidden Markov model?

Sean R Eddy

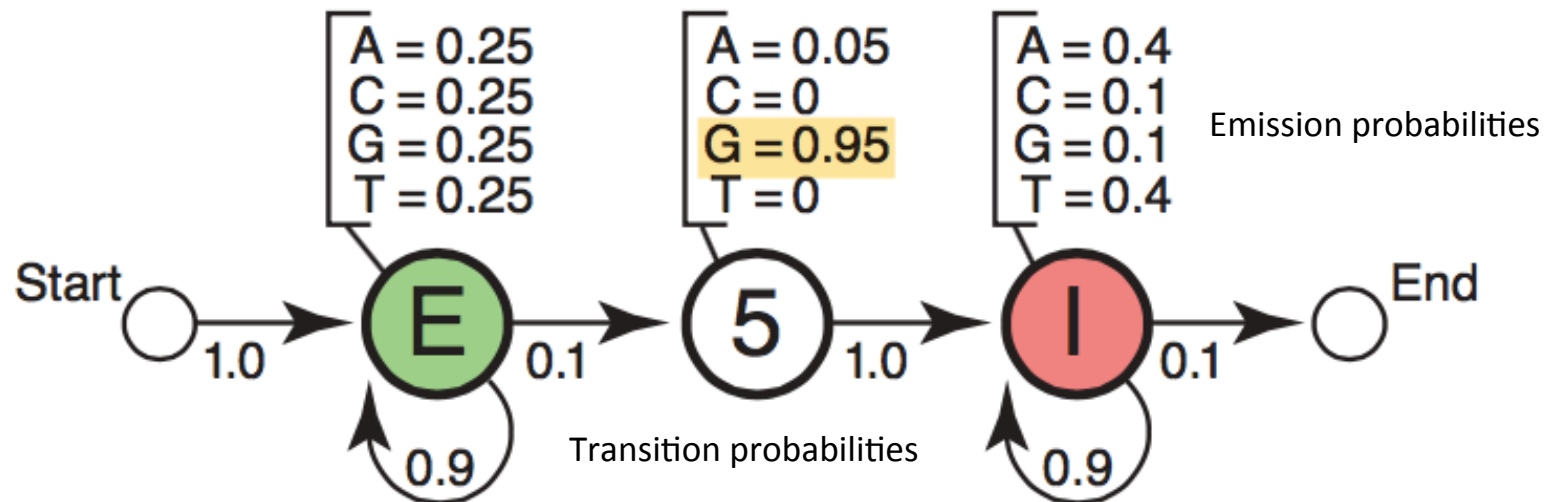
NATURE BIOTECHNOLOGY VOLUME 22 NUMBER 10 OCTOBER 2004

As a simple example, imagine the following 5' splice-site recognition problem. Assume we are given a DNA sequence that begins in an exon, contains one 5' splice site and ends in an intron. The problem is to identify where the switch from exon to intron occurred—where the 5' splice site (5'SS) is.



HMM with three states – exon, 5'SS and intron, and the probabilities of moving between these states. To estimate which state each element of the DNA sequence is in, we need information about the statistics of exon, splice site and intron sequences – these are the emission probabilities.

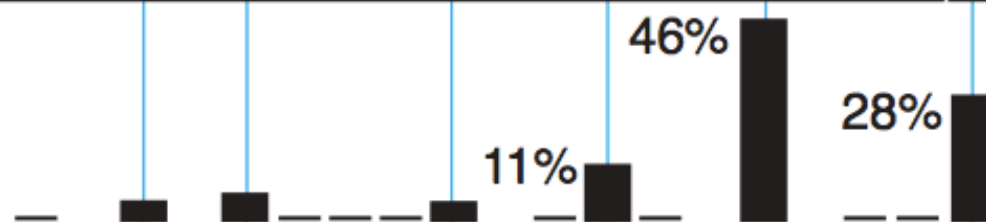




Sequence: **CTTCATGTGAAAGCAGACGTAAGTCA**

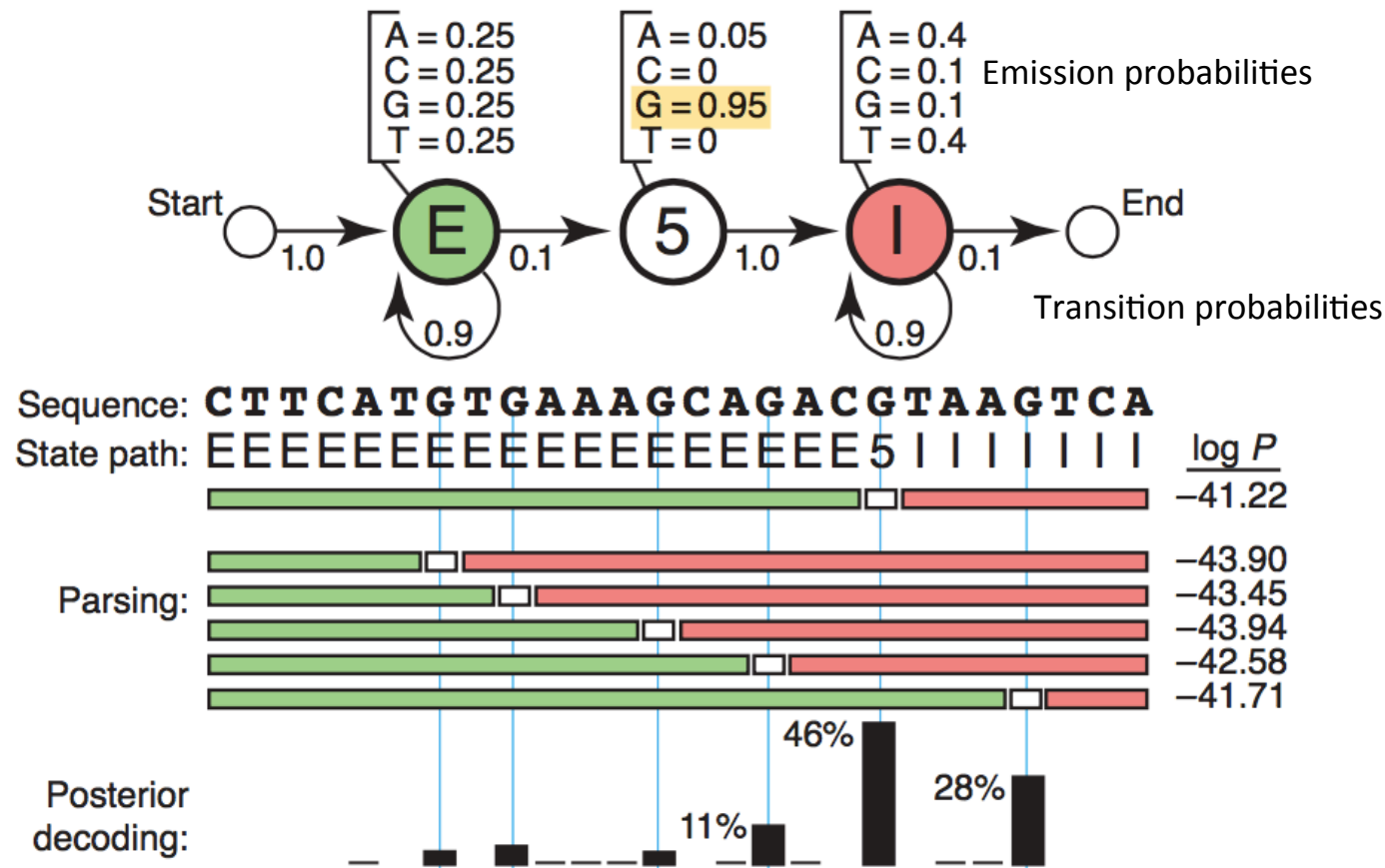


Posterior decoding:



A toy HMM for 5' splice site recognition.

Eddy, Nature Biotech, 2004



The probability  $P(S, \pi | \text{HMM}, \theta)$  that an HMM with parameters  $\theta$  generates a state path  $\pi$  and an observed sequence  $S$  is the product of all the emission probabilities and transition probabilities that were used.

Here for the 26 nucleotide sequence, there are 27 transitions and 26 emissions – the log of the product of all 53 probabilities yields  $\log P(S, \pi | \text{HMM}, \theta) = -41.22$ .

## Pfam 31.0 (March 2017, 16712 entries)

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. [More...](#)

### QUICK LINKS

[SEQUENCE SEARCH](#)[VIEW A PFAM ENTRY](#)[VIEW A CLAN](#)[VIEW A SEQUENCE](#)[VIEW A STRUCTURE](#)[KEYWORD SEARCH](#)[JUMP TO](#)

### YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...

Analyze your protein sequence for Pfam matches

View Pfam annotation and alignments

See groups of related entries

Look at the domain organisation of a protein sequence

Find the domains on a PDB structure

Query Pfam by keywords

Enter any type of accession or ID to jump to the page for a Pfam entry or clan, UniProt sequence, PDB structure, etc.

Or view the [help](#) pages for more information





**Family: *Globin* (PF00042)**

66 architectures

4699 sequences

12 interactions

1761 species

2557 structures

**Summary**

Domain organisation

Clan

Alignments

HMM logo

Trees

Curation &amp; model

Species

Interactions

Structures

**Jump to...**

enter ID/acc

**Go****Summary: Globin**

Pfam includes annotations and additional family information from a range of different sources. These sources can be accessed via the tabs below.

**Wikipedia: Globin****Pfam****InterPro**

This is the Wikipedia entry entitled "[Globin](#)". [More...](#)

**Globin**[Edit Wikipedia article](#)

Not to be confused with [globulin](#) or [globular protein](#).

The **globins** are a **superfamily** of **heme**-containing **globular proteins**, involved in **binding** and/or transporting **oxygen**. These proteins all incorporate the globin fold, a series of eight **alpha helical segments**. Two prominent members include **myoglobin** and **hemoglobin**. Both of these proteins reversibly bind oxygen via a **heme** prosthetic group. They are widely distributed in many **organisms**.<sup>[2]</sup>

**Contents** [\[hide\]](#)

- Structure
  - 1.1 Helix packing
- Evolution
  - 2.1 Sequence conservation
- Subfamilies
- Examples
- See also
- References

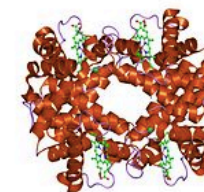
**Structure**

Globin superfamily members share a common **three-dimensional fold**.<sup>[3]</sup> This 'globin fold' typically consists of eight **alpha helices**, although some proteins have additional helix extensions at their termini.<sup>[4]</sup> Since the globin fold contains only helices, it is classified as an **all-alpha protein fold**.

The globin fold is found in its namesake globin **families** as well as in **phycocyanins**. The globin fold was thus the first protein fold discovered (myoglobin was the first protein whose structure was solved).

**Helix packing**

The eight helices of the globin fold core share significant nonlocal structure, unlike other **structural motifs** in which **amino acids** close to each other in **primary sequence** are also close in space. The helices pack together at an average angle of about 50 degrees, significantly steeper than other helical packings such as the **helix bundle**. The exact angle of helix packing depends on the sequence of the protein, because packing is mediated by the **sterics** and **hydrophobic** interactions of the amino acid **side chains** near the helix interfaces.

**Globin family**

the Structure of deoxyhemoglobin Rothschild 37 beta Trp----Arg: a mutation that creates an intersubunit chloride-binding site.<sup>[1]</sup>

**Identifiers**

Symbol	Globin
Pfam	PF00042 <a href="#">↗</a>
Pfam clan	CL0090 <a href="#">↗</a>
InterPro	IPR000971 <a href="#">↗</a>
PROSITE	PS01033 <a href="#">↗</a>
SCOP	1hba <a href="#">↗</a>
SUPERFAMILY	1hba <a href="#">↗</a>
CDD	cd01067 <a href="#">↗</a>



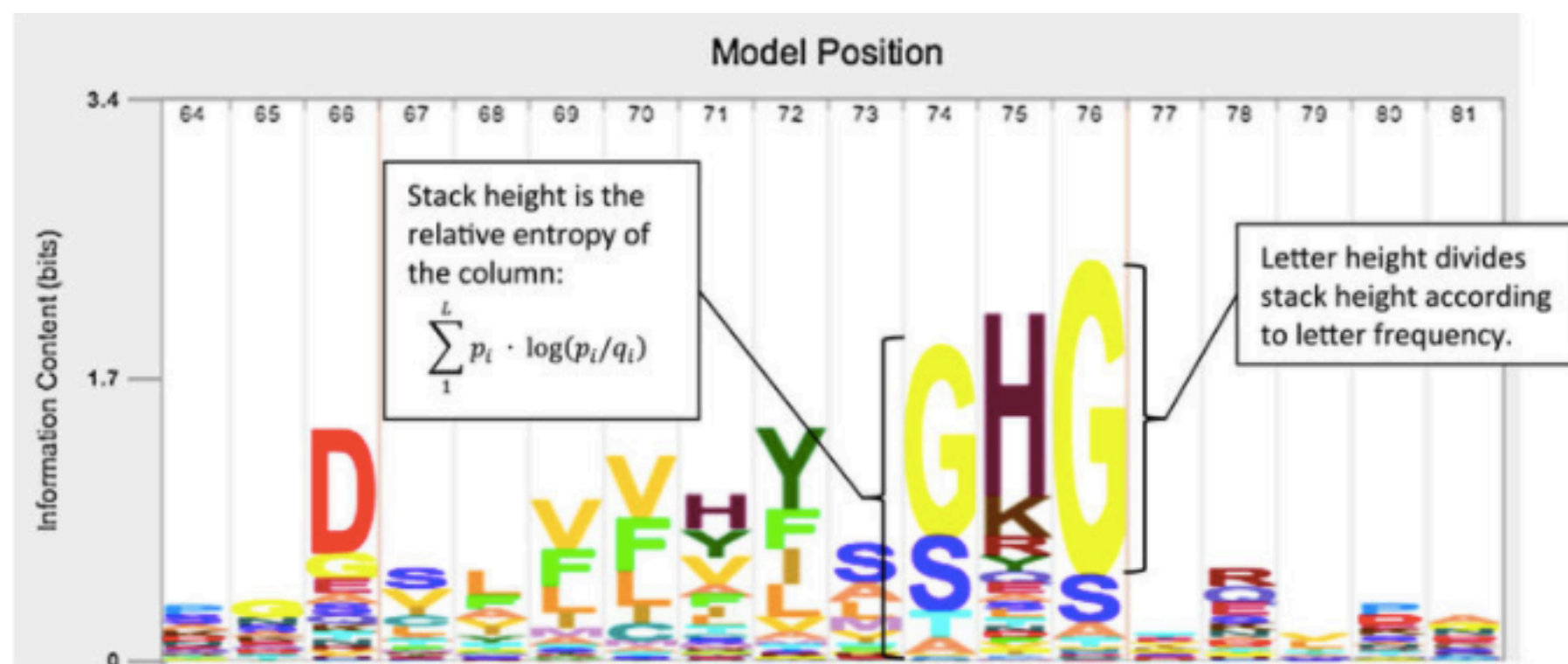


Figure 1

**Example profile logo.** This logo shows positions 64 to 81 of the Peptidase\_C14 profile HMM from Pfam (PF00656, Pfam 27.0), produced using Skyalign. The profile HMM was constructed using *hmmbuild* (default parameters) from HMMER 3.1 on the Pfam seed alignment. One of the active sites of this Caspase domain is found at position 75. This site is invariant in active peptidases, but not in this profile HMM. This is the result of two forces: (1) the Pfam alignment includes non-peptidase homologs, which do not contain a Histidine at this position, and (2) HMMER intentionally drives down the information content per position (using an approach called entropy weighting [12]) to increase sensitivity to remote homologs.



## Family: *Globin* (PF00042)

66 architectures

4699 sequences

12 interactions

1761 species

2557 structures

### Format an alignment

	Seed (73)	Full (4699)	Representative proteomes				UniProt (14977)	NCBI (21099)	Meta (34)
			RP15 (1041)	RP35 (2260)	RP55 (4063)	RP75 (5877)			
Alignment:	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Format:	Selex								
Order:	<input checked="" type="radio"/> Tree <input type="radio"/> Alphabetical								
Sequence:	<input checked="" type="radio"/> Inserts lower case <input type="radio"/> All upper case								
Gaps:	Gaps as "." or "-" (mixed)								
Download/view:	<input checked="" type="radio"/> Download <input type="radio"/> View								

Generate

20 July 1973, Volume 181, Number 4096

# SCIENCE

## Principles that Govern the Folding of Protein Chains

Christian B. Anfinsen

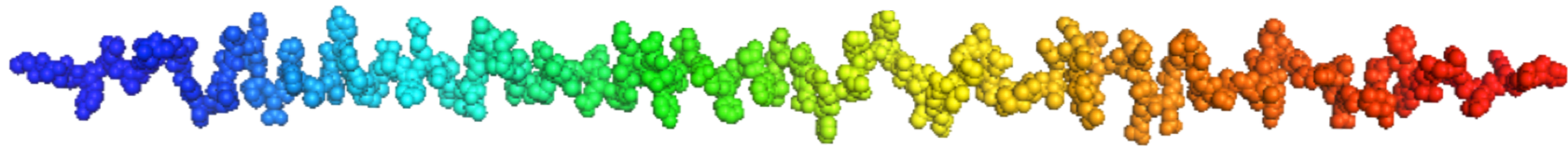
Empirical considerations of the large amount of data now available on correlations between sequence and three-dimensional structure (48), together with an increasing sophistication in the theoretical treatment of the energetics of polypeptide chain folding (49) are beginning to make more realistic the idea of the a priori prediction of protein conformation. It is certain that major advances in the understanding of cellular organization, and of the causes and control of abnormalities in such organization, will occur when we can predict, in advance, the three-dimensional phenotypic consequences of a genetic message.



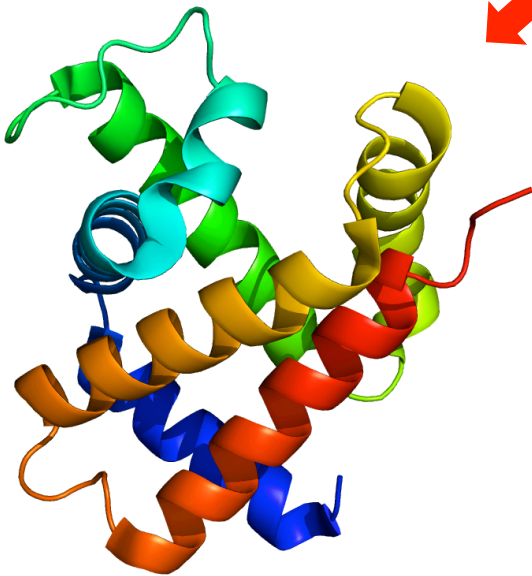
Anfinsen CB (1973). "Principles that govern the folding of protein chains". Science 181 (4096): 223–230.



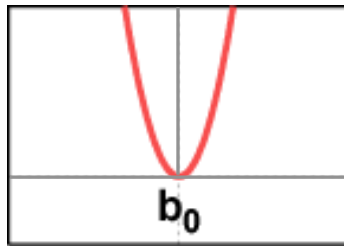
What makes this problem hard?



Polymer self-assembles into  
unique 3D structure

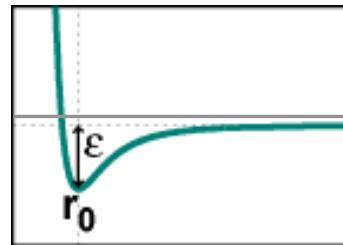


# Potential function for MD



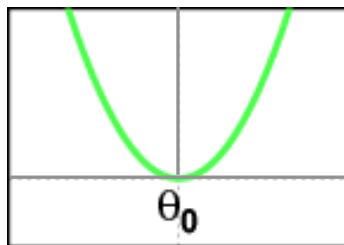
**Bond**

$$\sum_i^{bonds} K_{b,i} (b_i - b_{0,i})^2$$



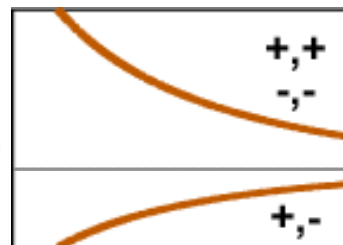
**van der Waals**

$$\sum_{pairs\cdot i,j} \left[ \epsilon_{ij} \left( \frac{r_{0,ij}}{r_{ij}} \right)^{12} - 2\epsilon_{ij} \left( \frac{r_{0,ij}}{r_{ij}} \right)^6 \right]$$



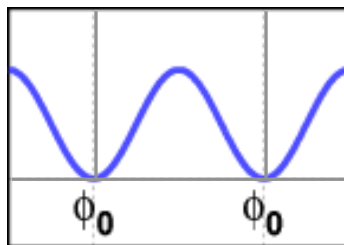
**Angle**

$$\sum_i^{bond\ angles} K_{\theta,i} (\theta_i - \theta_{0,i})^2$$



**Electrostatic**

$$\frac{1}{332} \sum_{pairs\cdot i,j} \left( \frac{q_i q_j}{r_{ij}} \right)$$



**Dihedral**

$$\sum_i^{torsion\ angles} K_{\phi,i} \{1 - \cos[n_i (\phi_i - \phi_{0,i})]\}$$

Evaluate forces and perform integration for every atom

Each picosecond of simulation time requires 500 iterations of cycle

E.g. w/ 50,000 atoms, each ps ( $10^{-12}$  s) involves 25,000,000 evaluations

$$F = -\frac{\partial U}{\partial x}$$

$$F = ma$$

$$a = \frac{v_3 - v_2}{\partial t}$$

$$v = \frac{x_3 - x_2}{\partial t}$$

$$E = U + K$$

$$\partial t = 2 \text{ fs}$$

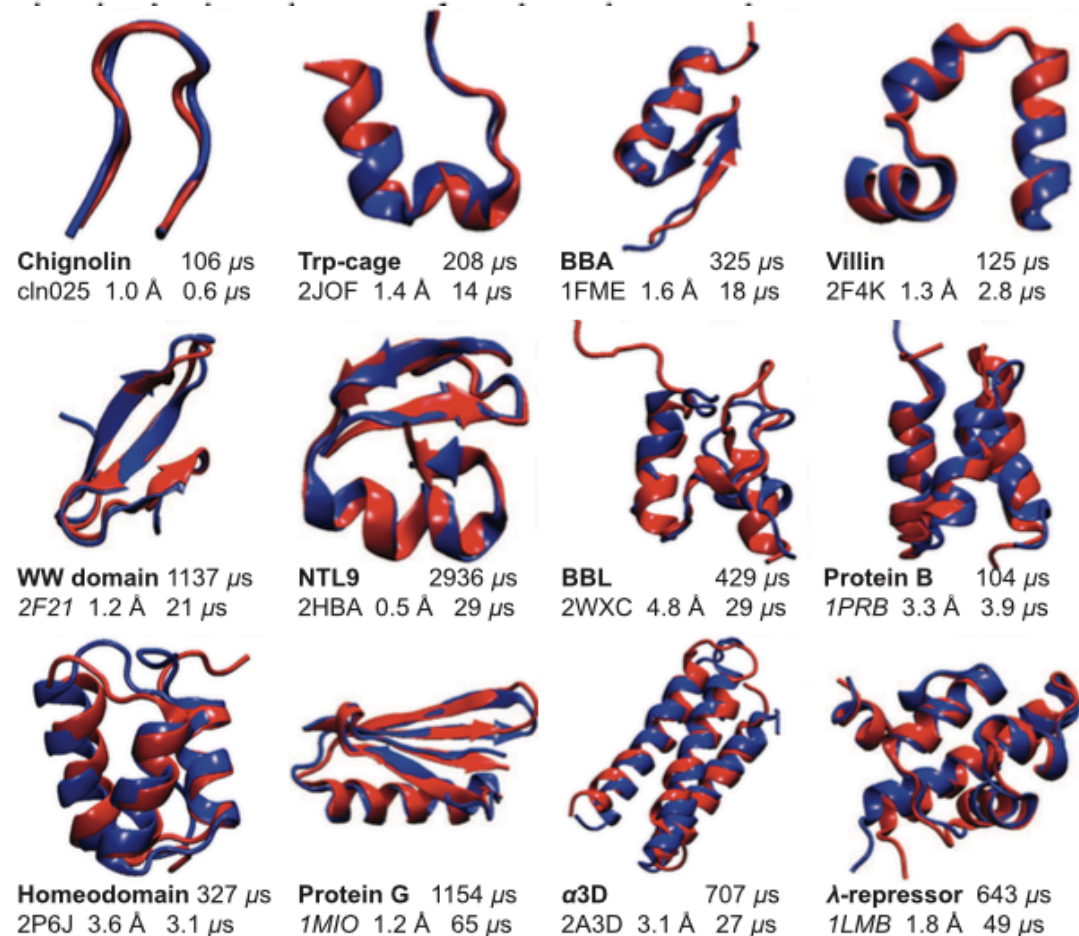
# How Fast-Folding Proteins Fold

Kresten Lindorff-Larsen,<sup>1\*†</sup> Stefano Piana,<sup>1\*†</sup> Ron O. Dror,<sup>1</sup> David E. Shaw<sup>1,2†</sup>

An outstanding challenge in the field of molecular biology has been to understand the process by which proteins fold into their characteristic three-dimensional structures. Here, we report the results of atomic-level molecular dynamics simulations, over periods ranging between 100  $\mu$ s and 1 ms, that reveal a set of common principles underlying the folding of 12 structurally diverse

Our research group has developed a specialized supercomputer, called Anton, which greatly accelerates the execution of atomistic molecular dynamics (MD) simulations. In addition, we recently modified the CHARMM force field in an effort to make it more easily transferable among different protein classes.

Lindorff-Larsen, Kresten, et al. "How fast-folding proteins fold." *Science* 334.6055 (2011): 517-520.



Rosetta – an algorithm that introduced assembly of sequence similar fragments was a major step forward in computational approaches to the protein folding problem.

Swiftly followed by ‘folding at home’ and the

## Crystal structure of a monomeric retroviral protease solved by protein folding game players

Firas Khatib<sup>1</sup>, Frank DiMaio<sup>1</sup>, Foldit Contenders Group, Foldit Void Crushers Group, Seth Cooper<sup>2</sup>, Maciej Kazmierczyk<sup>3</sup>, Mirosław Gilski<sup>3,4</sup>, Szymon Krzywda<sup>3</sup>, Helena Zabranska<sup>5</sup>, Iva Pichova<sup>5</sup>, James Thompson<sup>1</sup>, Zoran Popović<sup>2</sup>, Mariusz Jaskolski<sup>3,4</sup> & David Baker<sup>1,6</sup>

**Following the failure of a wide range of attempts to solve the crystal structure of M-PMV retroviral protease by molecular replacement, we challenged players of the protein folding game Foldit to produce accurate models of the protein. Remarkably, Foldit players were able to generate models of sufficient quality for successful molecular replacement and subsequent structure determination. The refined structure provides new insights for the design of antiretroviral drugs.**



**WHAT  
IF...**



**you  
could help  
find a cure?**



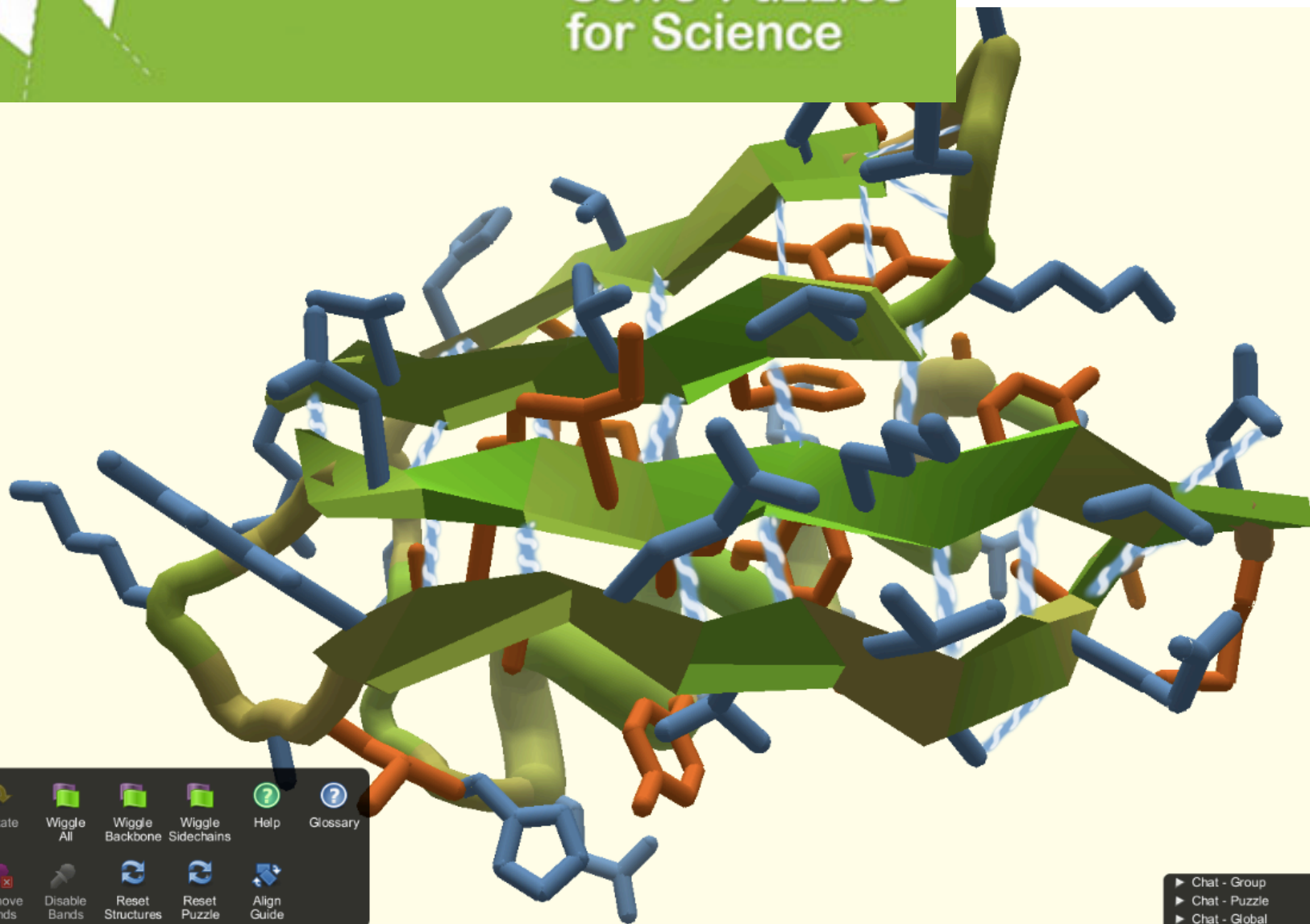
**PLAY VIDEO**

Help Stanford University scientists studying Alzheimer's, Huntington's, Parkinson's, and many cancers by simply running a piece of software on your computer.



# foldit

Solve Puzzles  
for Science



► Chat - Group auto show  
► Chat - Puzzle auto show  
► Chat - Global auto show  
► Notifications auto show



# How to extract useful information from protein sequences?

We can think of protein sequences from different species as samples from a probability distribution.....

	20	30	40	50	60	70	80	90	100																																																																									
Raccoon	V	E	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	A	D	A	I	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Ring-tailed coati	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	I	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Red fox	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	D	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I	
Weddell seal	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	N	A	I	M	S	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	D	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Harbor seal	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	A	D	A	I	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	D	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Eurasian badger	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I	
Red panda	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Ferret	V	D	E	V	G	G	E	T	I	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	S	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Fur seal	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	D	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
River otter	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Beach marten	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
European mink	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Ratel	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	K	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Colobus monkey	V	D	A	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	S	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I	
Lowland gorilla	V	D	A	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	S	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Human	V	D	A	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	S	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Spider monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I	
Brown spider monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	A	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Tamarin	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	A	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Moustached tamarin	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	A	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Capuchin	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
White Capuchin	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
white sapaïou	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
squirrel monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
red colobus	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	A	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Collared titi	V	X	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	S	N	X	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	S	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Common marmoset	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	H	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Howler monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	H	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E	I
Black spider monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G</																																	

## First order models

$P(A_1, \dots, A_L)$  = Probability that a sequence produces a folded, functional hemoglobin molecule.

To start, we might ask that the model matches the distribution of amino acids seen at each sequence position in the data (the set of available sequences).

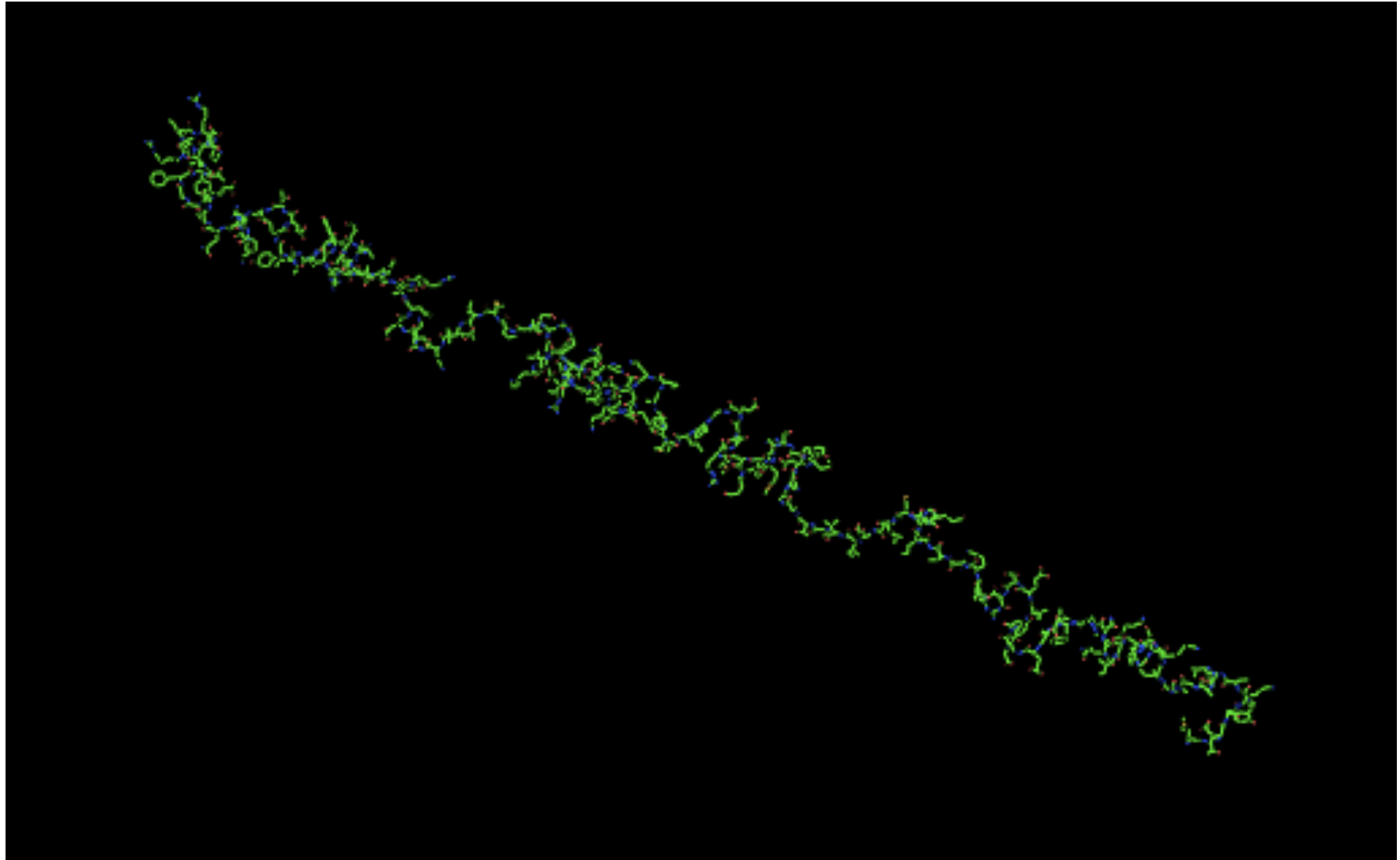
$$P_i(A_i) = \sum_{\{A_k | k \neq i\}} P(A_1, \dots, A_L) = f_i(A_i)$$

Hidden Markov Models are often used to model protein sequence evolution. Different sequence positions (variables) are assumed to evolve independently of one another – there are no interactions.

Does this model tell us anything about protein structure and function?



Those residues close in sequence will be close in 3D structure....



but their conservation level doesn't provide much information about the 3D structure.

## First order models

$P(A_1, \dots, A_L)$  = Probability that a sequence produces a folded, functional hemoglobin molecule.

To start, we might ask that the model matches the distribution of amino acids seen at each sequence position in the data (the set of available sequences).

$$P_i(A_i) = \sum_{\{A_k | k \neq i\}} P(A_1, \dots, A_L) = f_i(A_i)$$

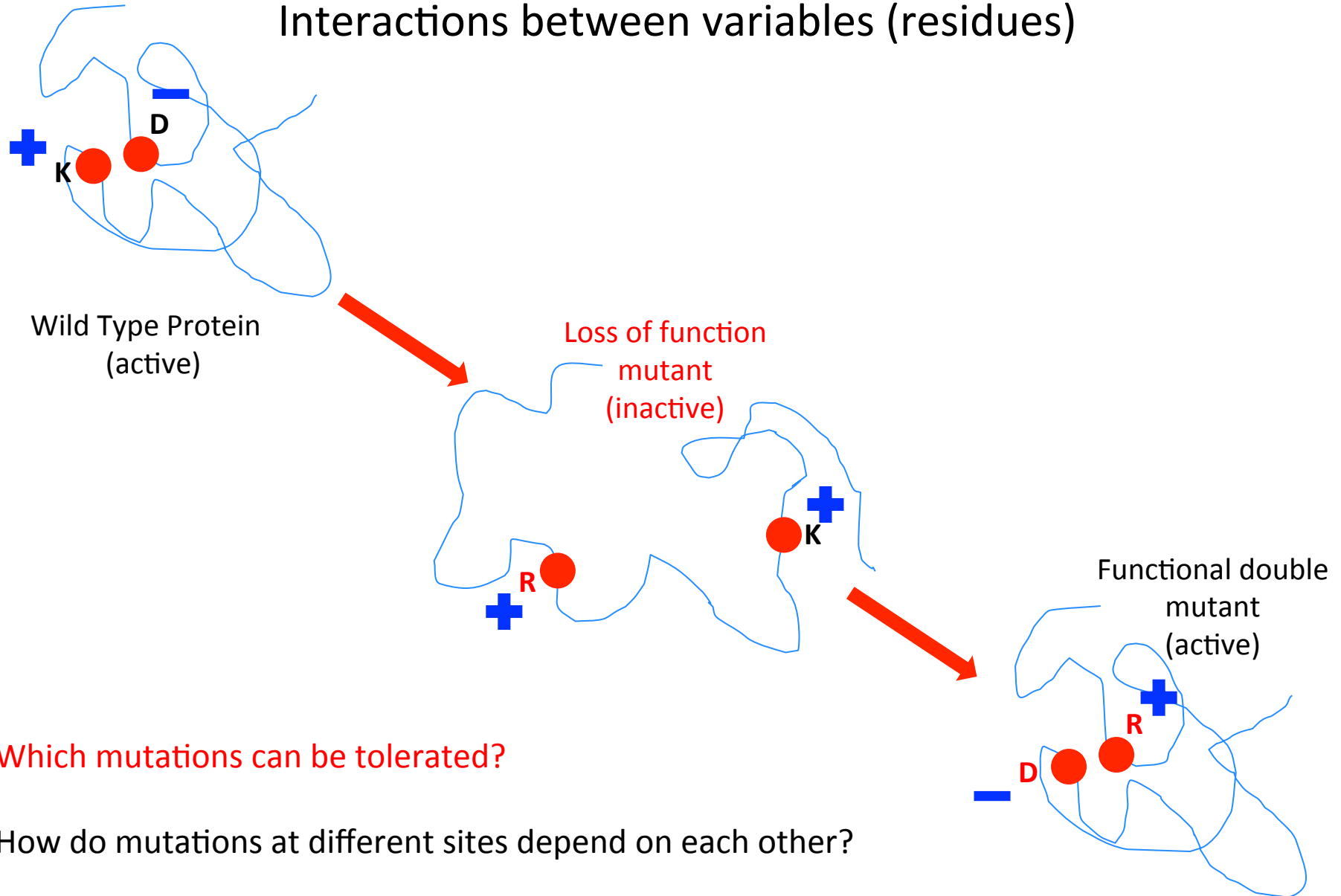
Hidden Markov Models are often used to model protein sequence evolution. Different sequence positions (variables) are assumed to evolve independently of one another – there are no interactions.

Does it tell us anything about the 3D protein structure.....NOT MUCH

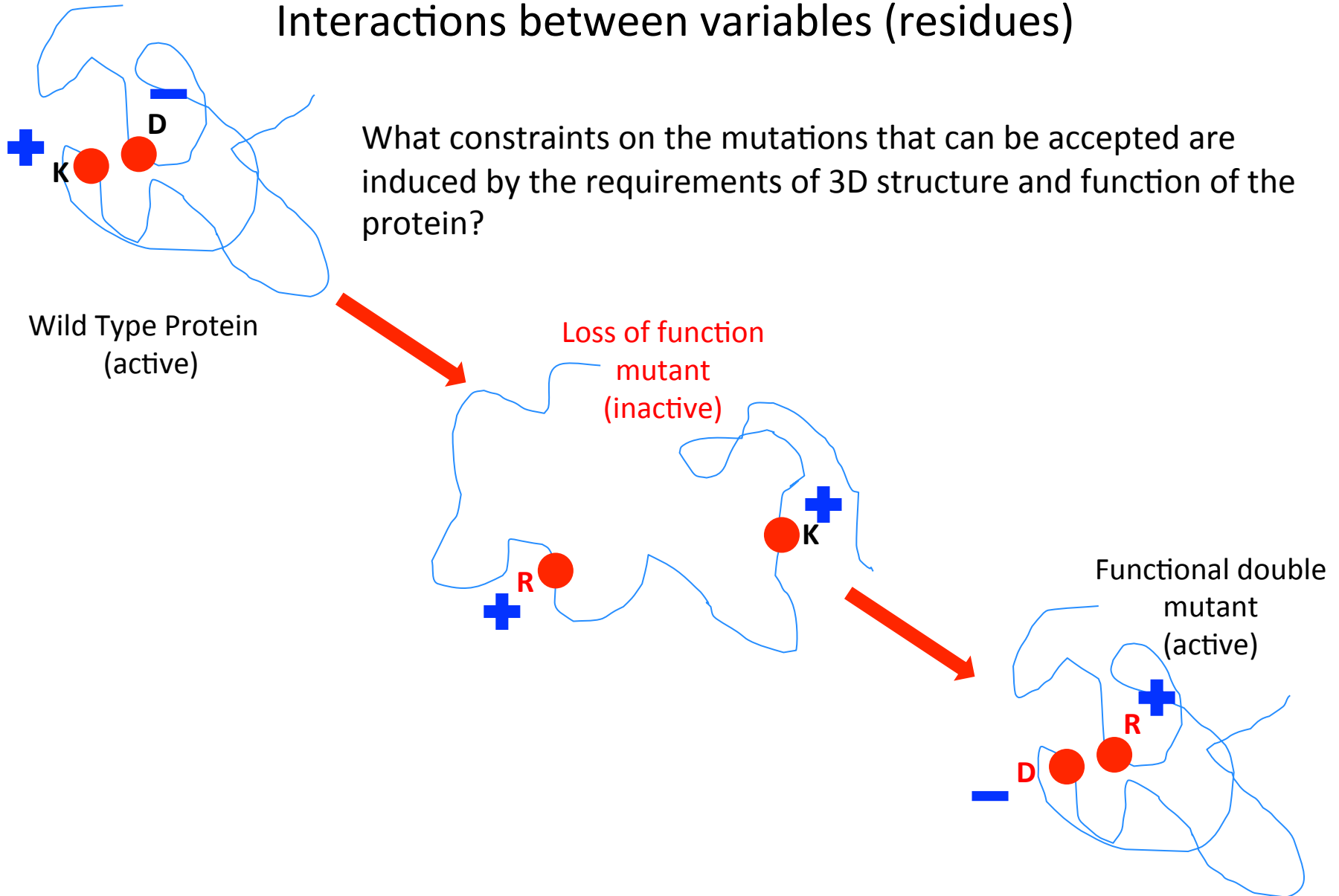
How can we get to 3D protein structure and/or function?

**INTERACTIONS!!!**

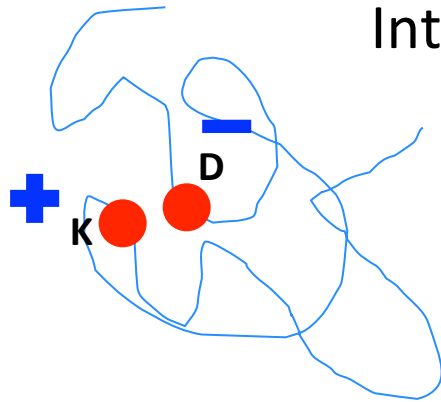
## Interactions between variables (residues)



## Interactions between variables (residues)



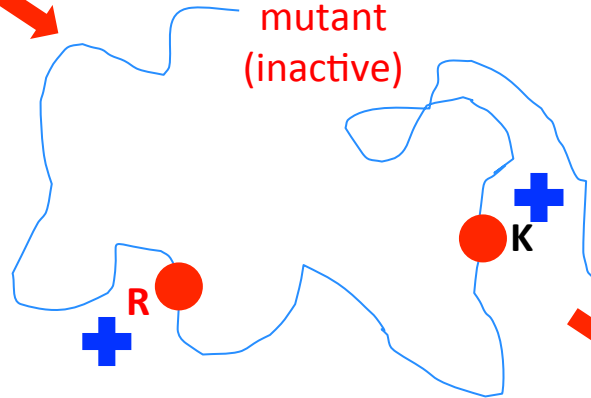
## Interactions between variables (residues)



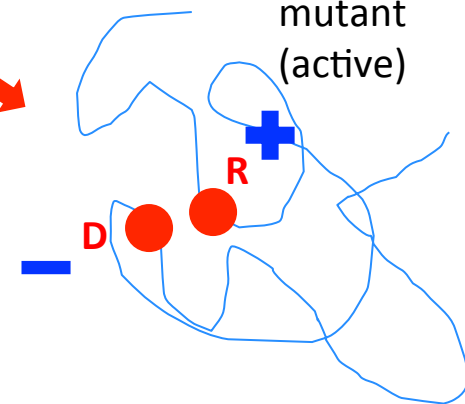
Wild Type Protein  
(active)

What constraints on the mutations that can be accepted are induced by the requirements of 3D structure and function of the protein?

Loss of function  
mutant  
(inactive)



Functional double  
mutant  
(active)



	1	2	3	4	5	6
<i>H.s</i>	D	A	I	L	V	K
<i>M.m</i>	D	A	I	L	A	K
<i>B.t</i>	R	A	I	M	V	D
<i>G.g</i>	R	V	I	L	V	D

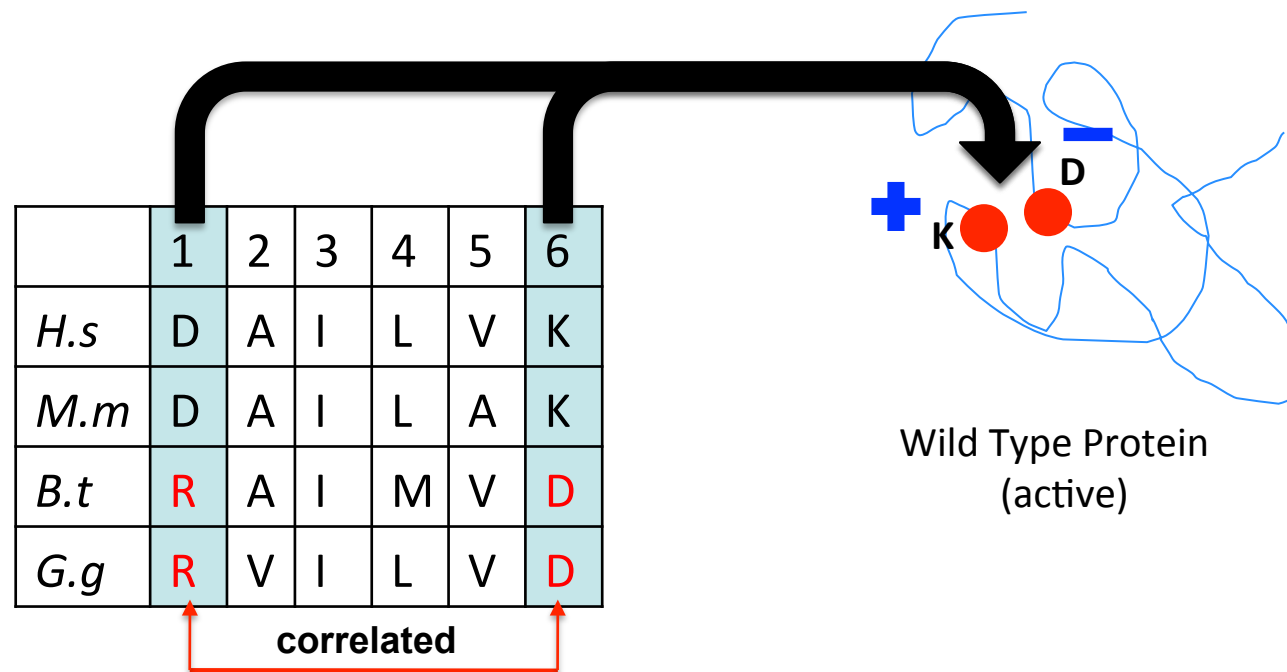
correlation

# Interactions lead to covariance of amino acids

	20	30	40	50	60	70	80	90	100																																																																								
Raccoon	V	E	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	A	D	A	I	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Ring-tailed coati	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	I	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Red fox	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	T	P	D	A	V	M	G	N	A	K	V	K	A	H	G	K	K	V	L	N	S	F	S	D	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Weddell seal	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	P	N	A	I	M	S	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	D	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	Q	L	H	V	D	P	E
Harbor seal	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	A	D	A	I	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	D	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Eurasian badger	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	Y	F	D	S	F	G	D	L	S	T	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Red panda	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Ferret	V	D	E	V	G	G	E	T	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	P	D	A	V	M	S	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Fur seal	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	D	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
River otter	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Beach marten	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
European mink	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Ratel	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	N	S	F	S	E	G	L	R	N	L	D	N	L	K	G	T	F	A	K	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Colobus monkey	V	D	A	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	A	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	S	L	K	G	T	F	S	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Lowland gorilla	V	D	A	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	S	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Human	V	D	A	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	S	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Spider monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	A	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Brown spider monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	A	L	S	T	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Tamarin	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	A	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Moustached tamarin	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	A	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Capuchin	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
White Capuchin	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
white sapajou	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
squirrel monkey	V	E	D	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
red colobus	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	D	S	F	G	D	L	S	T	A	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Collared titi	V	X	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	S	N	X	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	S	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Common marmoset	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	N	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	T	H	L	D	N	L	K	G	T	F	A	H	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Howler monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	T	P	D	A	V	M	H	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Black spider monkey	V	D	E	V	G	G	E	A	L	G	R	L	L	V	V	P	W	T	Q	R	F	F	E	S	F	G	D	L	S	S	P	D	A	V	M	G	N	P	K	V	K	A	H	G	K	K	V	L	G	A	F	S	D	G	L	A	H	L	D	N	L	K	G	T	F	A	Q	L	S	E	L	H	C	D	K	L	H	V	D	P	E
Go																																																																																	

# Exploit correlation structure of protein sequences

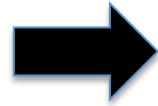
What constraints on the mutations that can be accepted are induced by the requirements of 3D structure and function of the protein?



Can we **invert correlations** in the sequence data to provide information about **the 3D protein structure**?

## Measure pair correlations in the sequence alignment

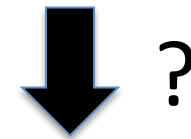
	1	2	3	4	5	6
<i>H.s</i>	K	D	I	L	V	D
<i>M.m</i>	K	D	I	L	V	D
<i>S.c</i>	D	D	I	K	V	H
<i>S.p</i>	D	E	I	L	V	H



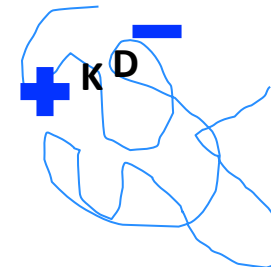
	K	D	I	L	V	D
K	X					I
D		X				
I			X			
L				X		
V					X	
D	I					X

Compute (for example) the mutual information for each pair of columns

$$MI_{ij} = \sum_{A_i, A_j=1}^q f_{ij}(A_i, A_j) \ln \left( \frac{f_{ij}(A_i, A_j)}{f_i(A_i)f_j(A_j)} \right)$$



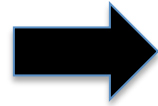
Wild Type Protein





## Measure pair correlations in the sequence alignment

	1	2	3	4	5	6
<i>H.s</i>	K	D	I	L	V	D
<i>M.m</i>	K	D	I	L	V	D
<i>S.c</i>	D	D	I	K	V	H
<i>S.p</i>	D	E	I	L	V	H

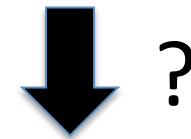


	K	D	I	L	V	D
K	X					I
D		X				
I			X			
L				X		
V					X	
D	I					X

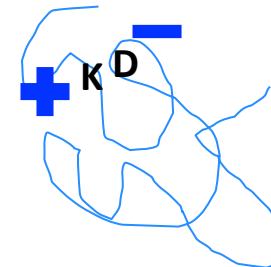
Compute (for example) the mutual information for each pair of columns

$$MI_{ij} = \sum_{A_i, A_j=1}^q f_{ij}(A_i, A_j) \ln \left( \frac{f_{ij}(A_i, A_j)}{f_i(A_i)f_j(A_j)} \right)$$

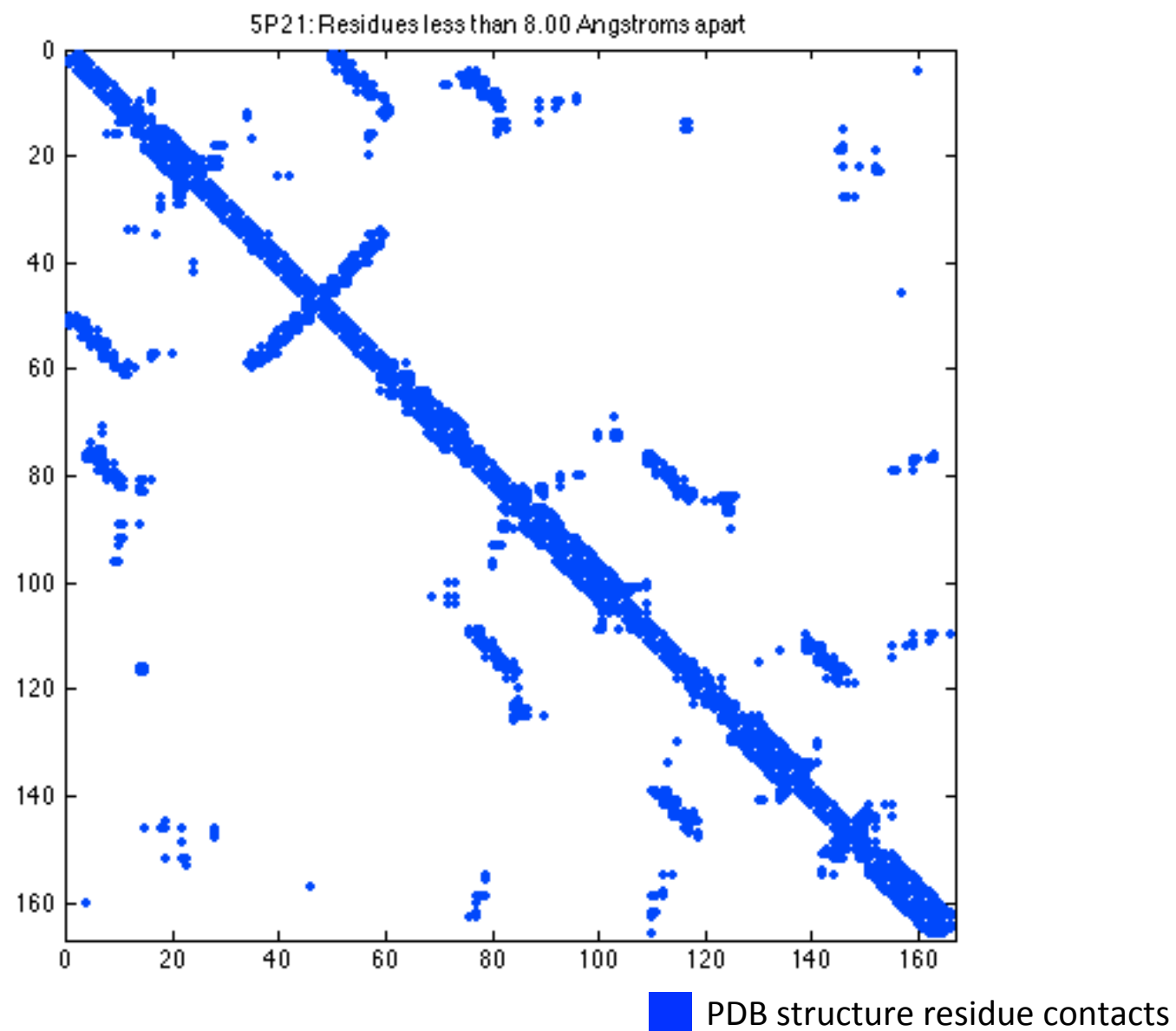
**Are the most correlated residue pairs close in tertiary protein structure?**



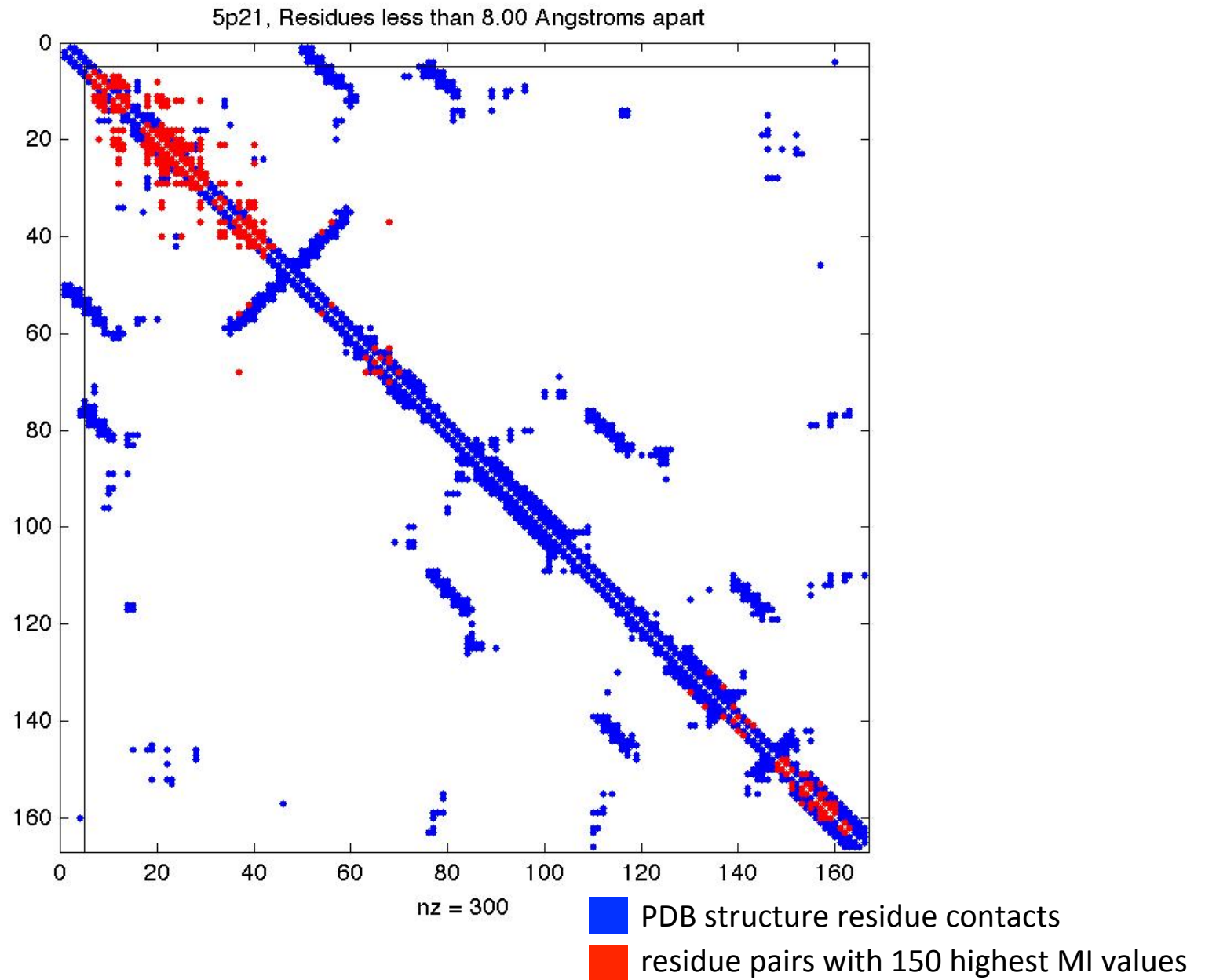
Wild Type Protein



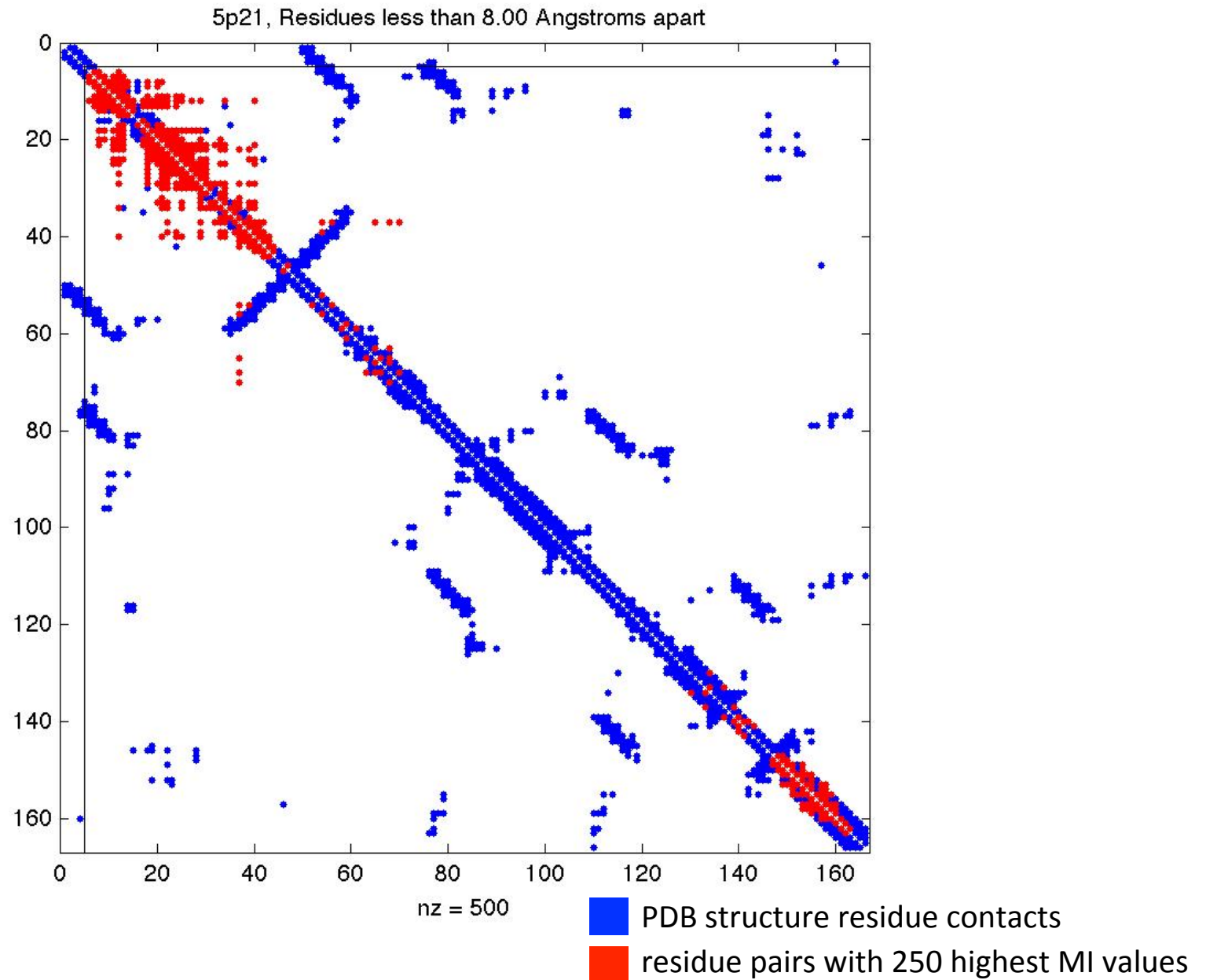
Testing: Are highly correlated pairs close in structure?



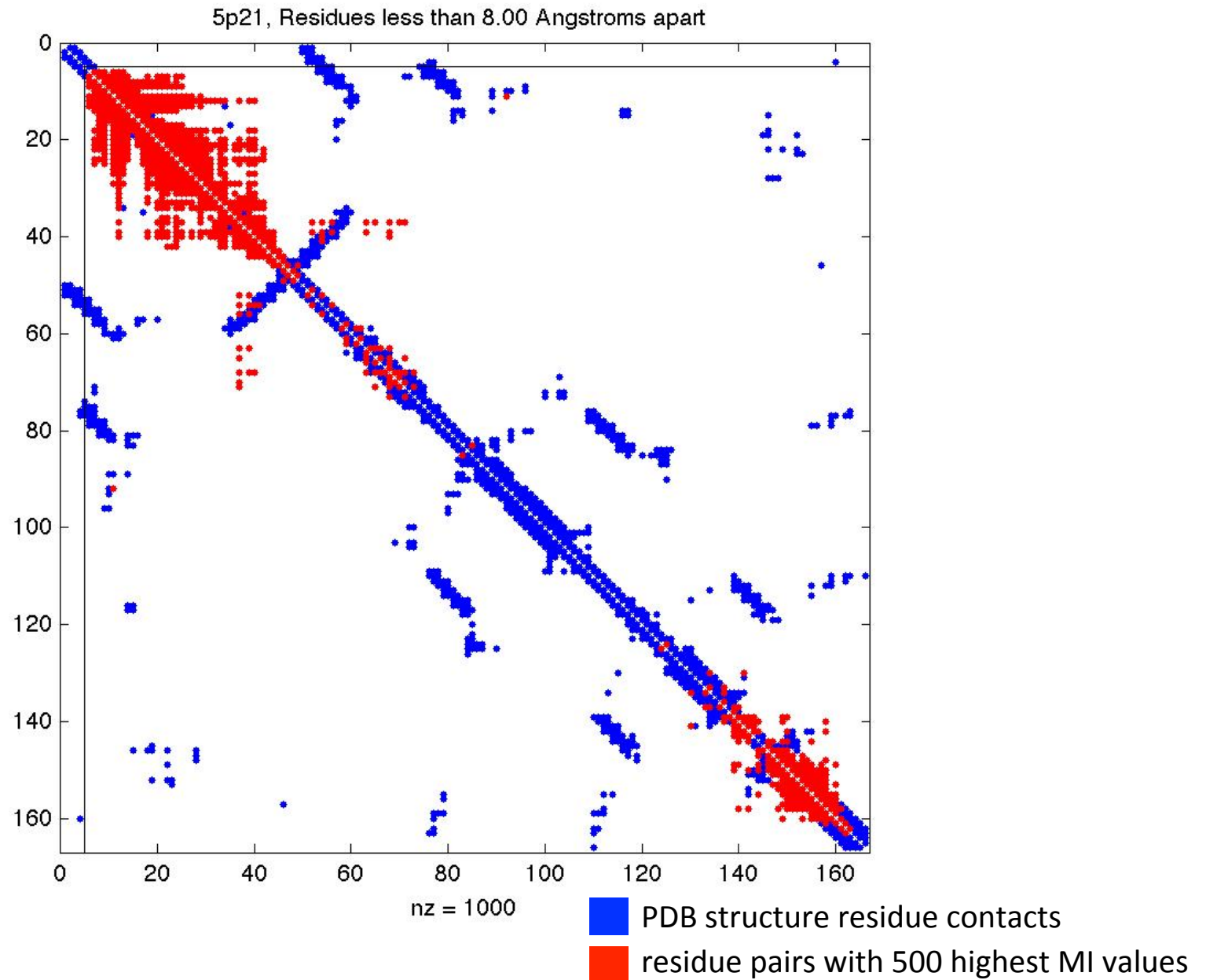
## Testing MI: Are high-scoring pairs close in structure?



## Testing MI: Are high-scoring pairs close in structure?

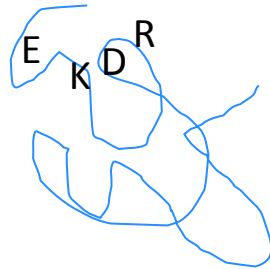


## Testing MI: Are high-scoring pairs close in structure?

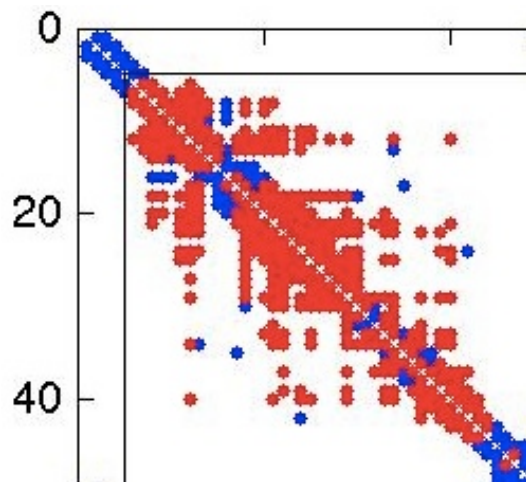


## Why does this fail? Construct and analyze a model

Wild Type Protein



K ↔ D

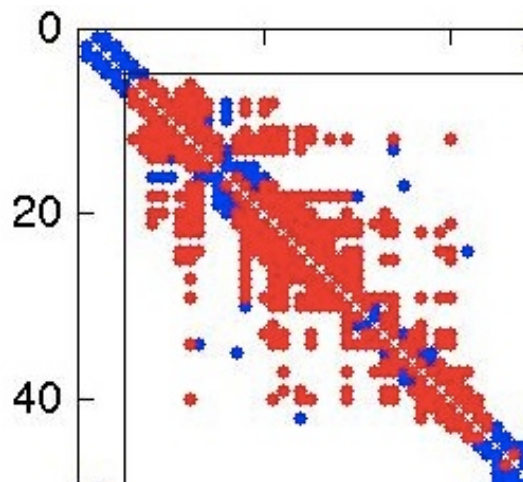
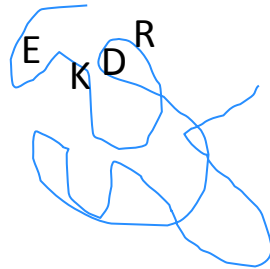


	K	D	E	R
K	X	I		
D	I	X		
E			X	
R				X

- PDB structure residue contacts
- residue pairs with 100 highest MI values

# Construct and analyze a model

Wild Type Protein

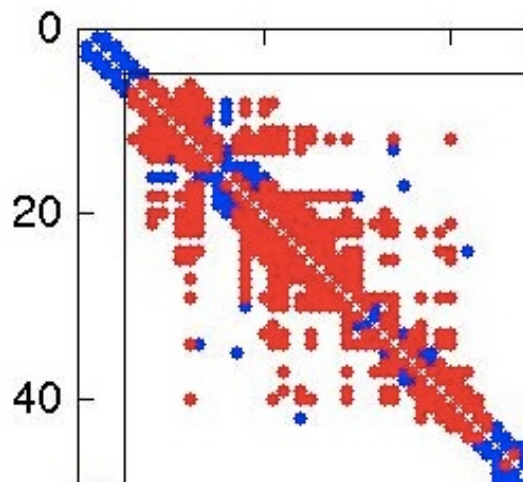
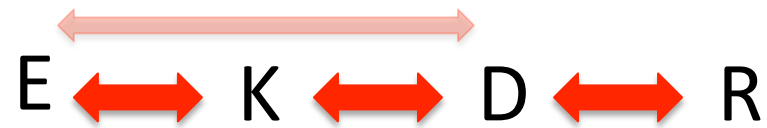
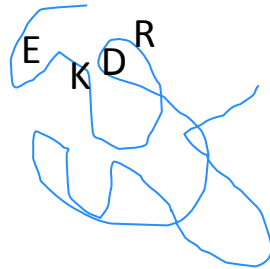


	K	D	E	R
K	X	I	I	
D	I	X		I
E	I		X	
R		I		X



- PDB structure residue contacts
- residue pairs with 100 highest MI values

# Construct and analyze a model

Wild Type Protein



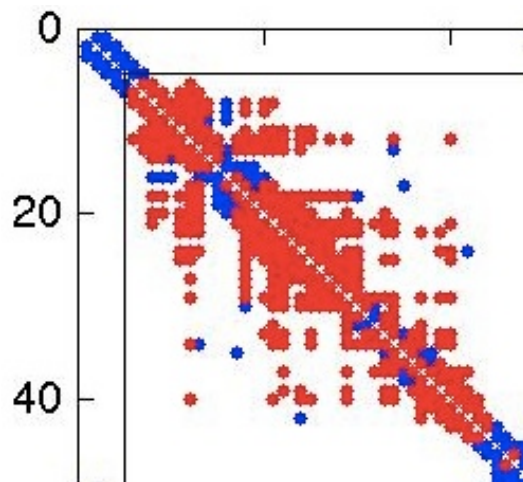
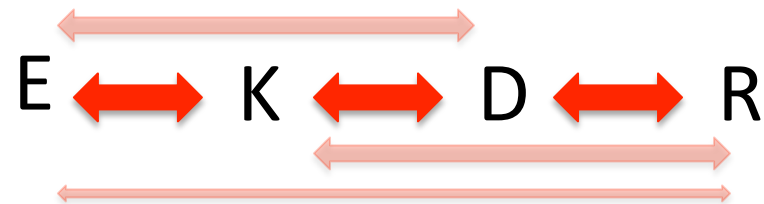
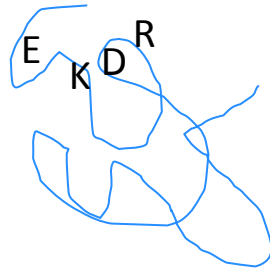
	K	D	E	R
K	X	I	I	
D	I	X	C	I
E	I	C	X	
R		I		X

-  PDB structure residue contacts
-  residue pairs with 100 highest MI values





# Construct and analyze a model

Wild Type Protein



	K	D	E	R
K	X	I	I	C
D	I	X	C	I
E	I	C	X	C
R	C	I	C	X

-  PDB structure residue contacts
-  residue pairs with 100 highest MI values

## Probability model

$P(A_1, \dots, A_L)$  = Probability a sequence is part of protein family

$$P_i(A_i) = \sum_{\{A_k | k \neq i\}} P(A_1, \dots, A_L) = f_i(A_i)$$

$$P_{ij}(A_i, A_j) = \sum_{\{A_k | k \neq i, j\}} P(A_1, \dots, A_L) = f_{ij}(A_i, A_j)$$

Satisfy these constraints and choose the maximum entropy or least constrained model:

$$S = - \sum_{\{A_i | i=1, \dots, L\}} P(A_1, \dots, A_L) \ln P(A_1, \dots, A_L)$$

## Probability model

$P(A_1, \dots, A_L)$  = Probability a sequence is part of protein family

$$= \frac{1}{Z} \exp \left\{ - \sum_i h_i A_i - \sum_{(i,j)} e_{ij} (A_i, A_j) \right\}$$

## Probability model

$P(A_1, \dots, A_L)$  = Probability a sequence is part of protein family

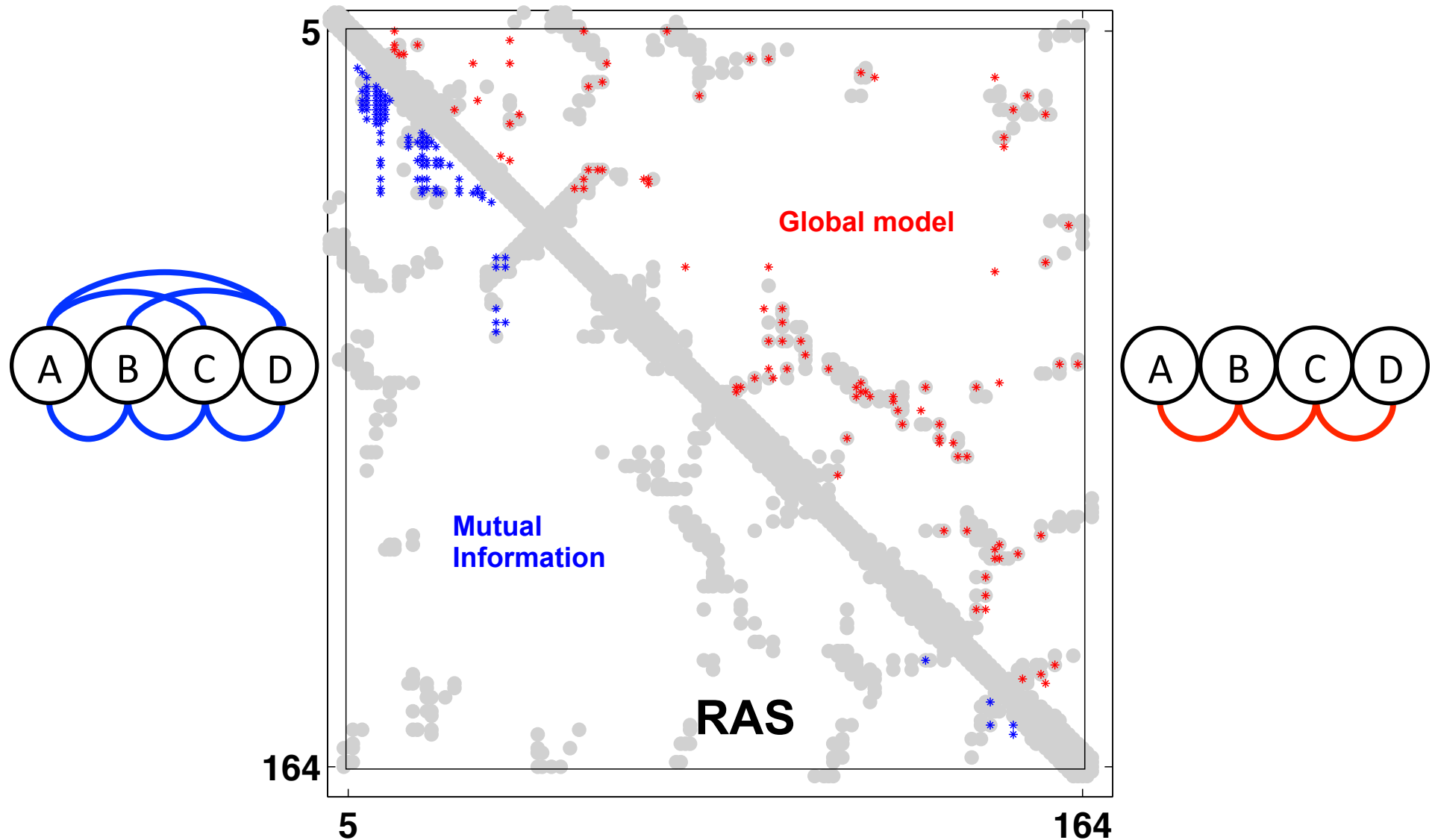
$$= \frac{1}{Z} \exp \left\{ - \sum_i h_i A_i - \sum_{(i,j)} e_{ij} (A_i, A_j) \right\}$$

naïve mean field approximation, assume small couplings

$$e_{ij} (A_i, A_j) = (C^{-1})_{ij} (A_i, A_j)$$

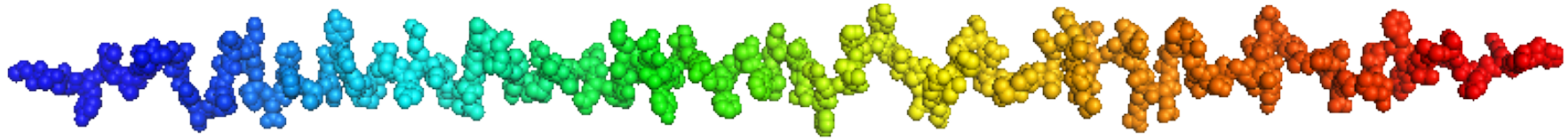
## Sequences to structure:

Solution: A global probability model – high scoring pairs are close in structure!



## Sequences to structure:

Is this enough information to fold the protein?

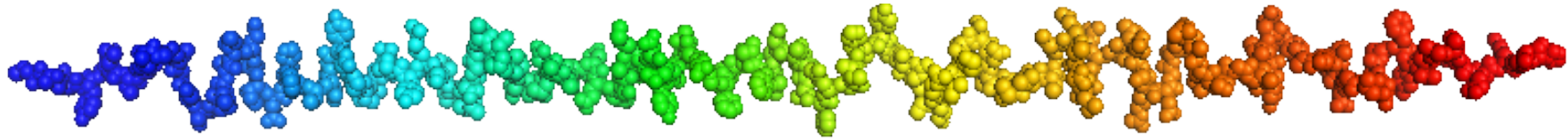


Start with an extended polypeptide, add predicted secondary structure

Our hypothesis is the high scoring pairs will be close in the structure – so **constrain the distance between the residues in each pair.**

## Sequences to structure:

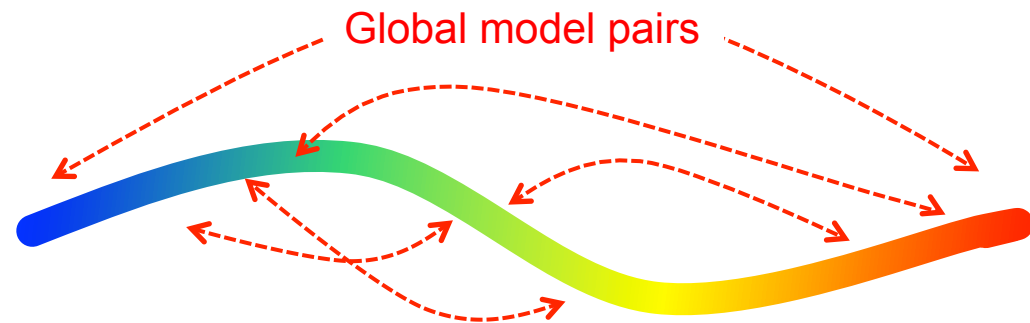
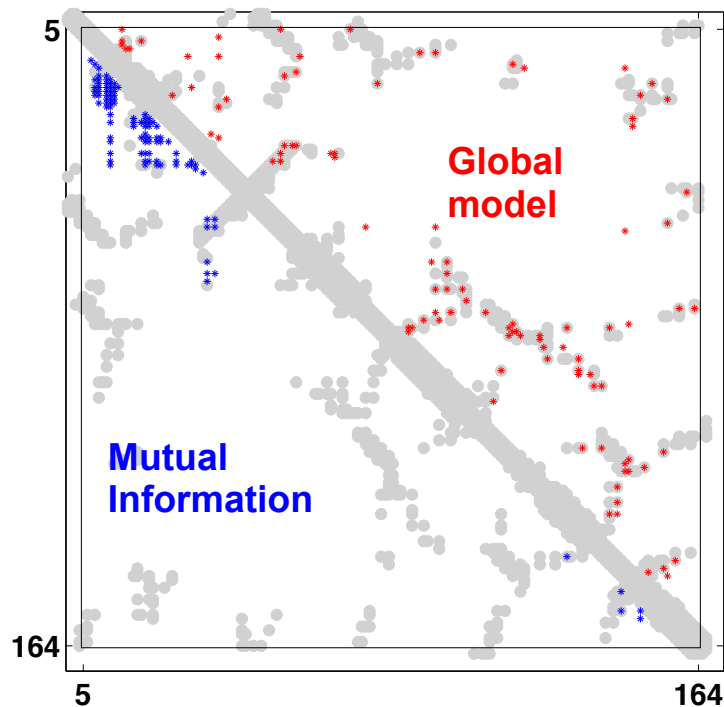
Is this enough information to fold the protein?



Start with an extended polypeptide, add predicted secondary structure

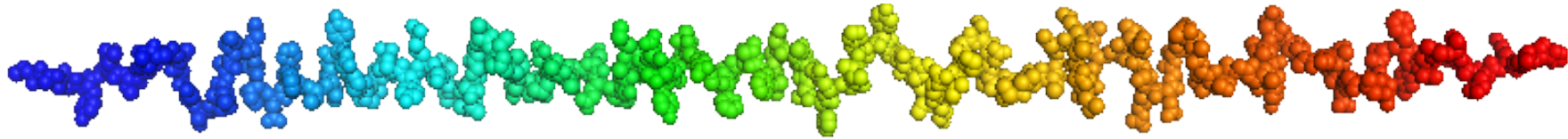
Our hypothesis is the high scoring pairs will be close in the structure – so **constrain the distance between the residues in each pair**.

**This massively reduces the space of possible 3D conformations of the protein**





# Folding the protein

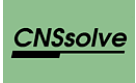


Start with an extended polypeptide

Our hypothesis is the high scoring pairs will be close in the structure – so we **constrain the distance between the two residues in each pair.**

Crystallography & NMR System

<http://cns-online.org/v1.3/>



**Crystallography & NMR System**

**Main Menu**

- ▶ About CNS
- ▶ Download
- ▶ Installation
- ▶ Getting started
- ▶ Input files
- ▶ Modules
- ▶ Libraries
- ▶ Utilities
- ▶ Tutorial
- ▶ Syntax Manual
- ▶ CNS wiki

*Version:* 1.3  
*Patch level:* 0  
*Status:* general release

*Authors:*

A.T.Brunger, P.D.Adams, G.M.Clore, W.L.Delano, P.Gros,  
R.W.Grosse-Kunstleve, J.-S.Jiang, J.M.Krahn,  
J.Kuszewski, M.Nilges, N.S.Pannu, R.J.Read, L.M.Rice,  
G.F.Schroeder, T.Simonson, G.L.Warren

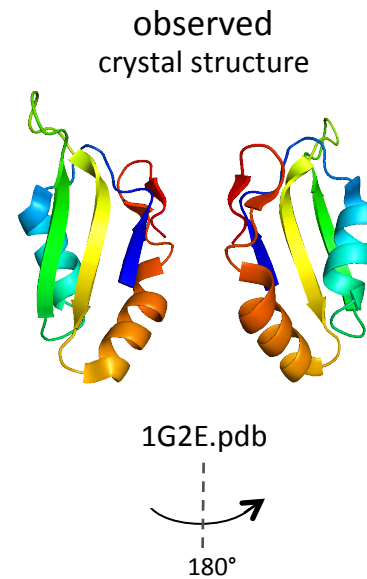
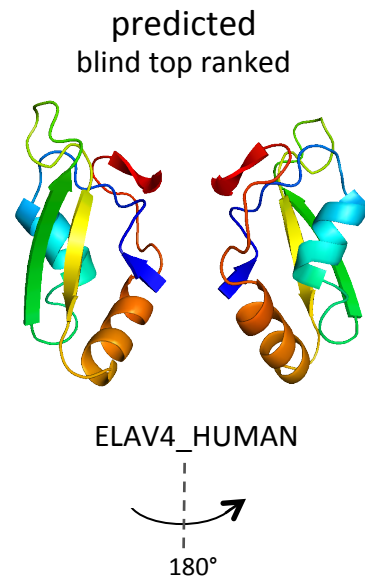
Copyright (c) 1997-2010 Yale University [License](#)

Use these distance constraints together with predicted secondary structure in the standard distance geometry and simulated annealing protocol from CNS to generate structures

## Sequences to structure:

Correlations in the mutation patterns of amino acids within a protein used to predict tertiary structure.

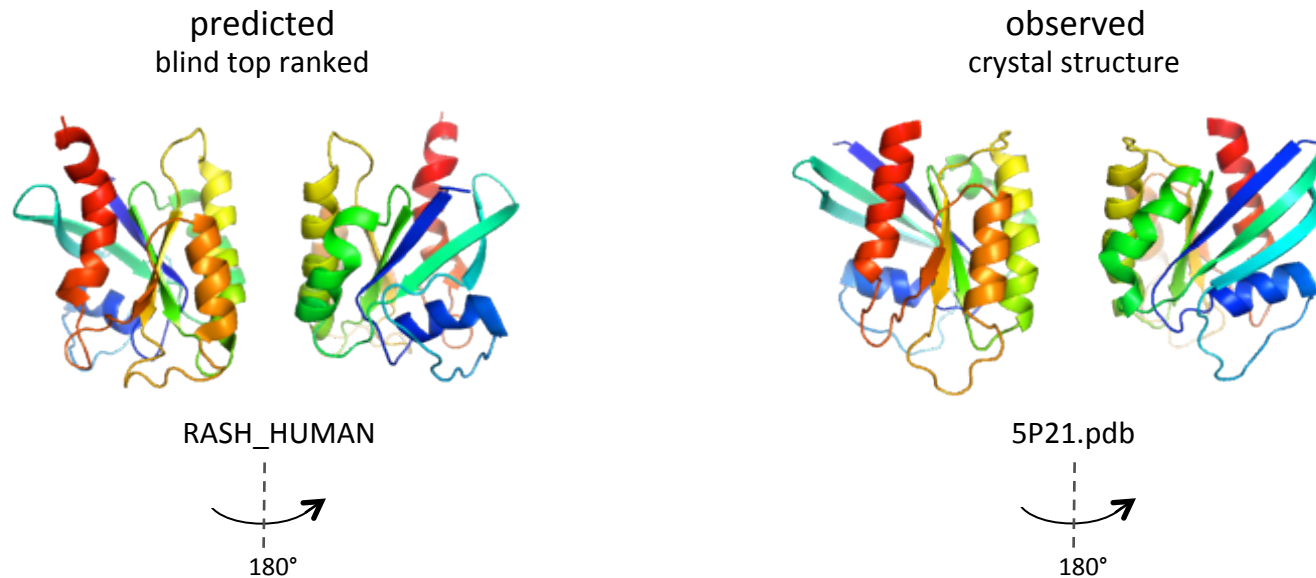
An RNA binding domain: 2.9Å C-alpha rmsd over 67 residues



## Sequences to structure:

Correlations in the mutation patterns of amino acids within a protein used to predict tertiary structure.

RAS: 3.5Å C-alpha rmsd over 161 residues



Beta strands positioned well enough to predict correct registration with the exception of the beta one strand, which was very recently shown to exist in both the conformations that we find.

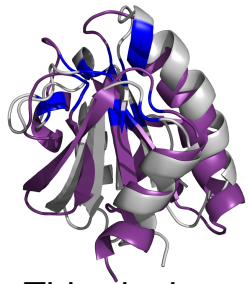
# Comparison of predicted and observed structures



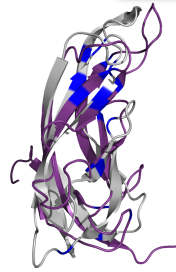
Crystal structure\*



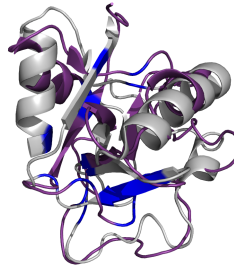
Predicted structure



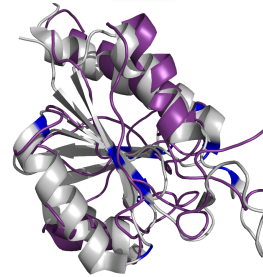
Thioredoxin:  
3.5 Å, 97 res



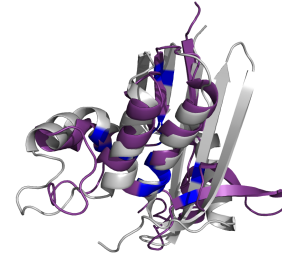
Cadherin: 3.8 Å,  
88 residues



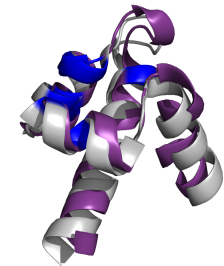
Lectin: 4.0 Å  
100 res



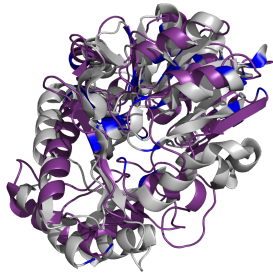
MOEA\_ECOLI:  
4.0 Å 128 res



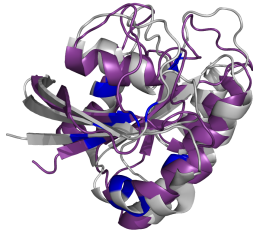
RnaseH: 3.5 Å  
114 res



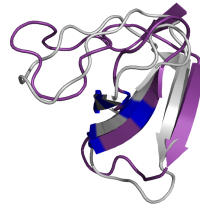
Q9KZ96\_STRCO:  
1.46 Å 49 res



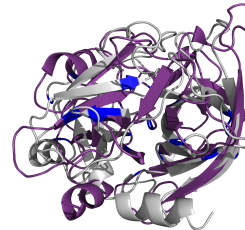
O85354\_CAUCR:  
3.74 Å, 278 res



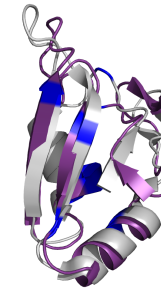
Ras: 3.2 Å,  
156 res



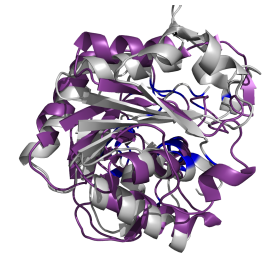
SH3: 3.6 Å,  
47 res



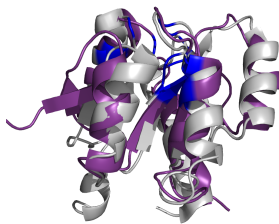
Trypsin: 4.27 Å,  
186 res



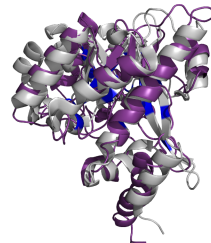
RRM\_1: 3.16 Å,  
71 residues



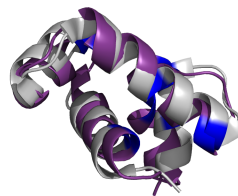
ISPD\_ECOLI:  
4.0 Å, 152 res



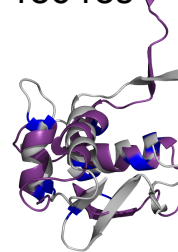
CheY: 2.98 Å,  
107 residues



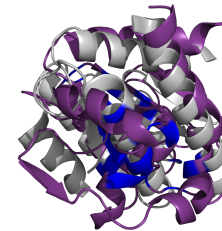
Q3SJE6\_THIDA:  
3.6 Å, 206 res



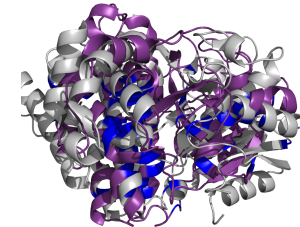
GERE\_BACSU:  
1.8 Å, 49 res



OmpR: 3.9 Å  
62 res



GMHA\_VIBCH :  
4.0 Å, 133 res



Q7WRJ3\_KLEPN  
4.5 Å, 232 res

Challenge : transmembrane proteins

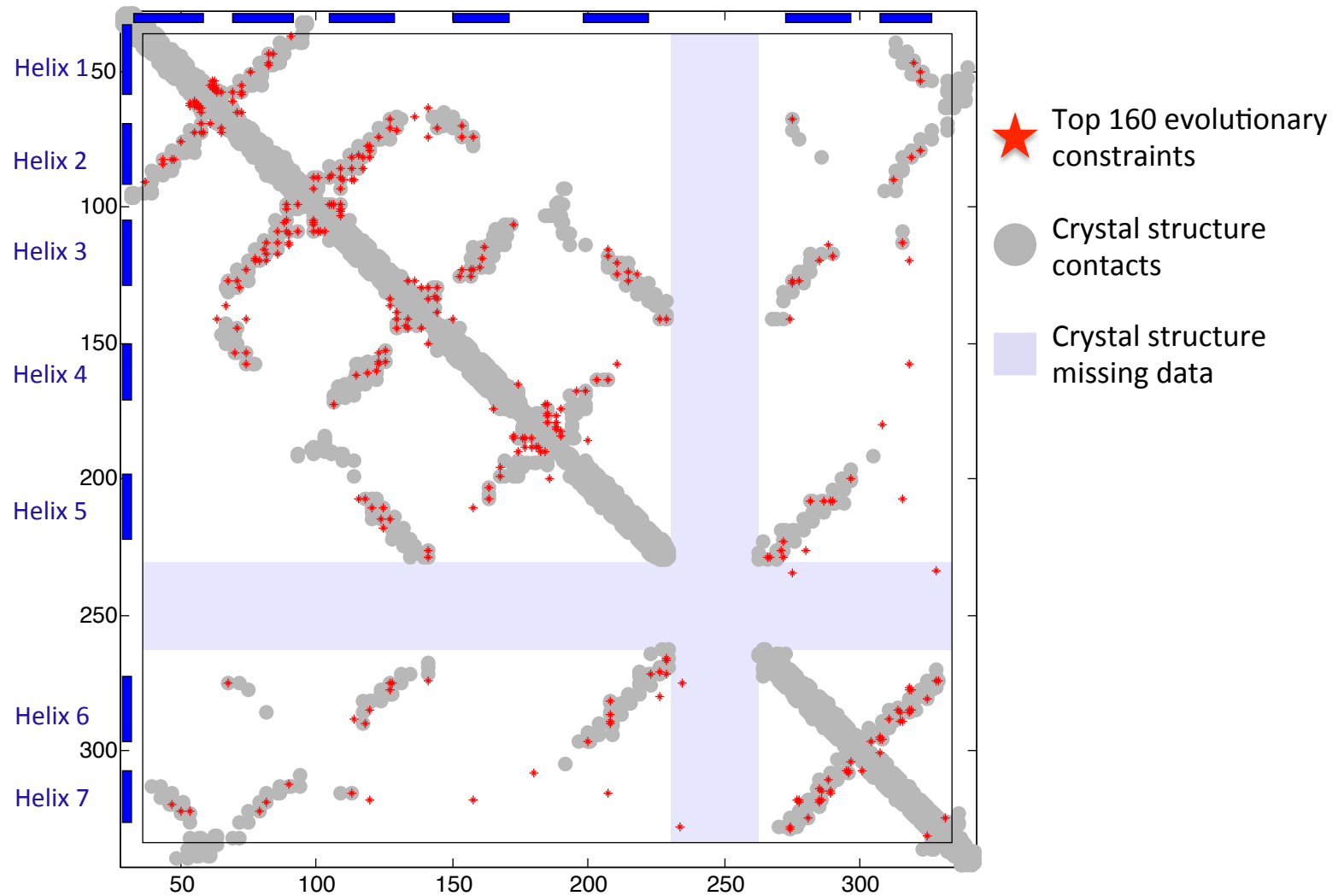
Hard experimentally

No high throughput

Drug targets - > 40%

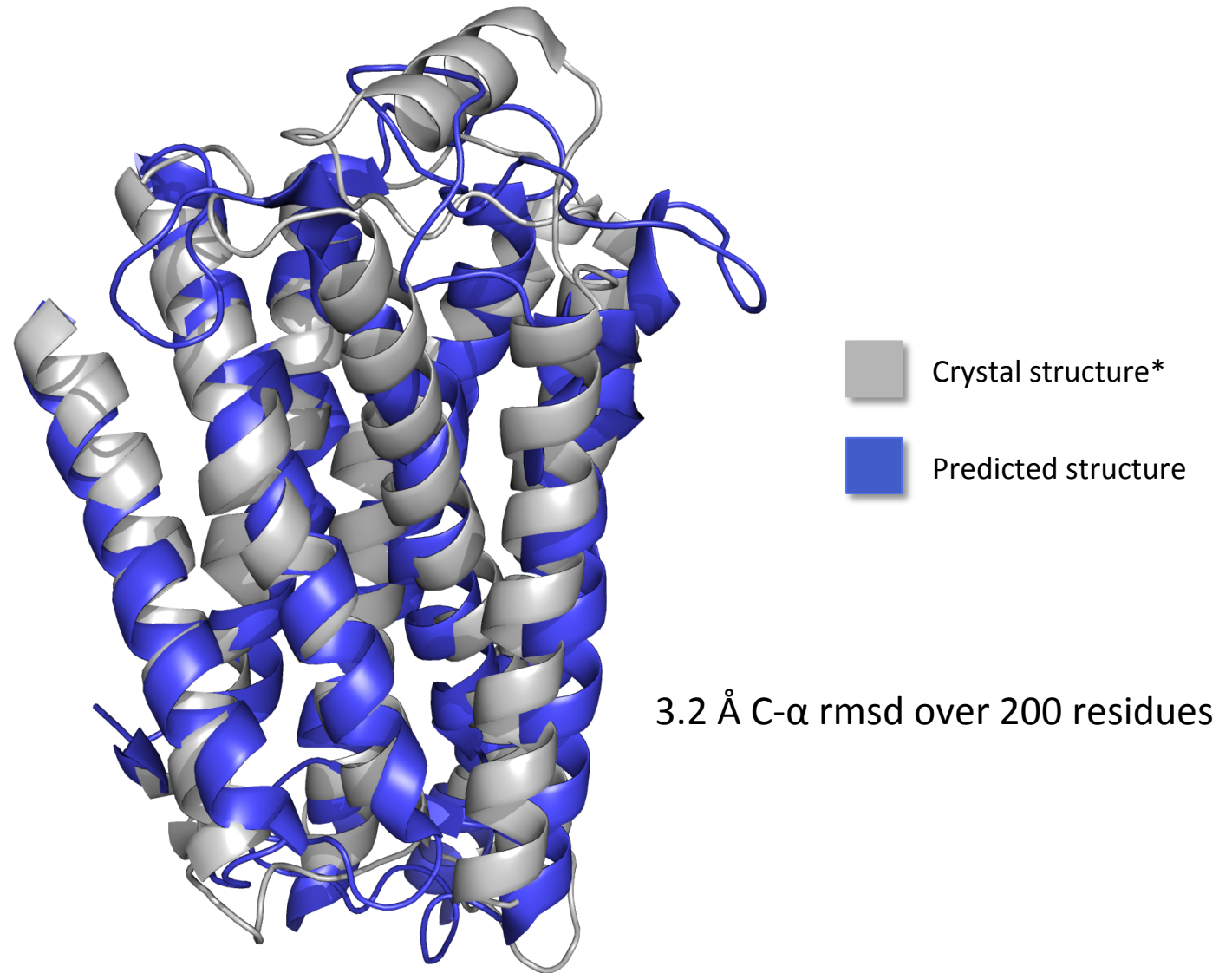
*De novo* prediction

## $\beta 2$ adrenergic receptor : evolutionary constraints



The  $\beta 2$  adrenergic receptor, a seven transmembrane helix GPCR family member.

## Blind prediction for $\beta 2$ adrenergic receptor



\*Crystal Structure: Rasmussen SG, DeVree BT, Zou Y, Kruse AC, Chung KY, Kobilka TS, Thian FS, Chae PS, Pardon E, Calinski D, Mathiesen JM, Shah ST, Lyons JA, Caffrey M, Gellman SH, Steyaert J, Skinotis G, Weis WI, Sunahara RK, Kobilka BK. Nature. 2011 Jul 19.

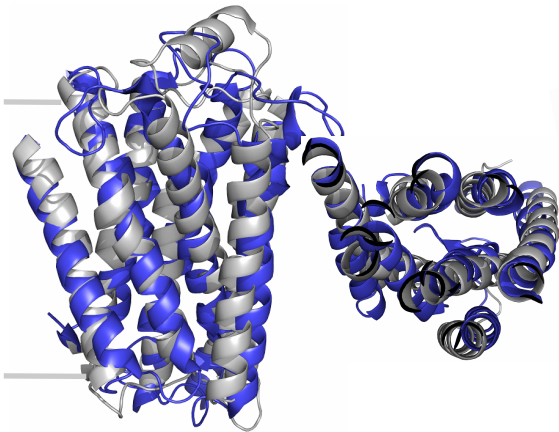


## Sequences to structure:

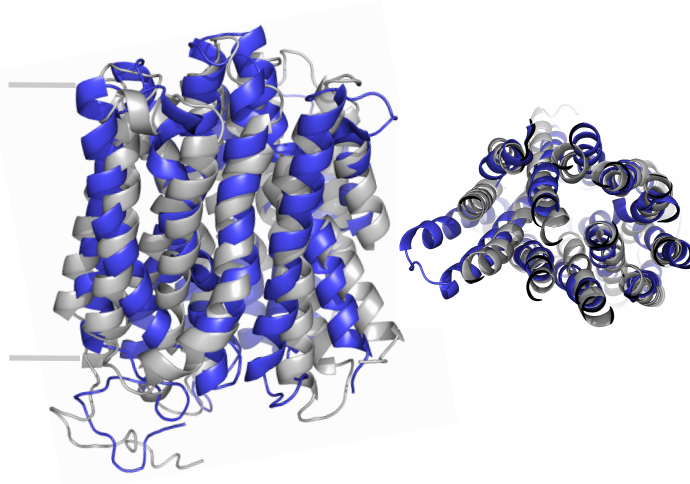
### Blind structure prediction.

■ observed    ■ predicted

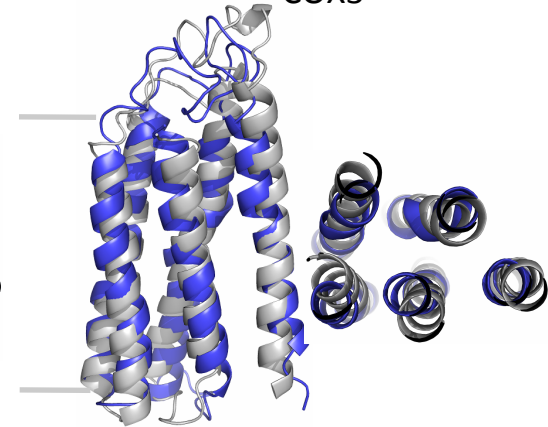
β2 adrenergic receptor



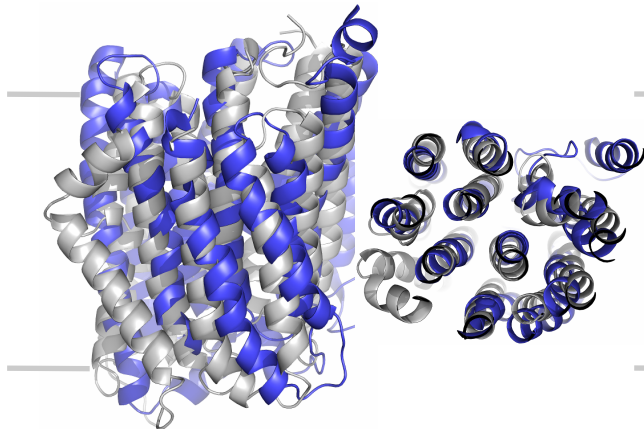
G-3-P transporter GlpT



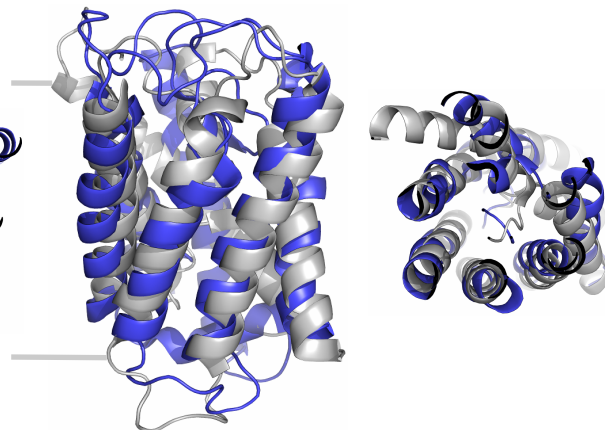
COX3



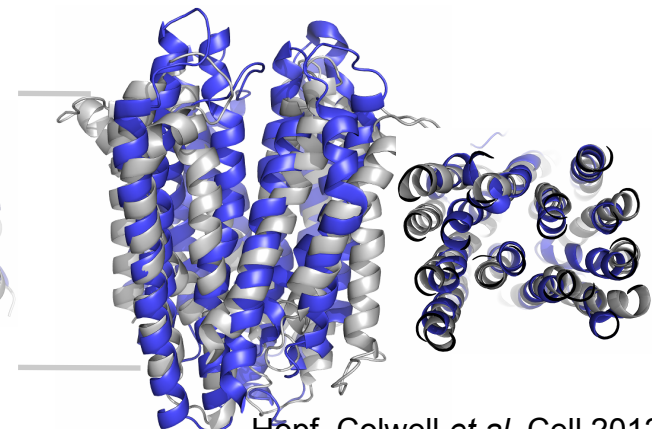
NADH-quinone oxidoreductase  
subunit N 1



Aquaporin-0



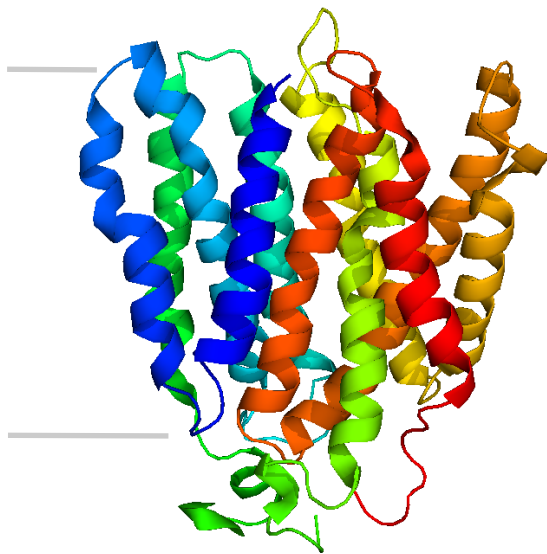
Multidrug resistance protein norM



Hopf, Colwell *et al*, Cell 2012

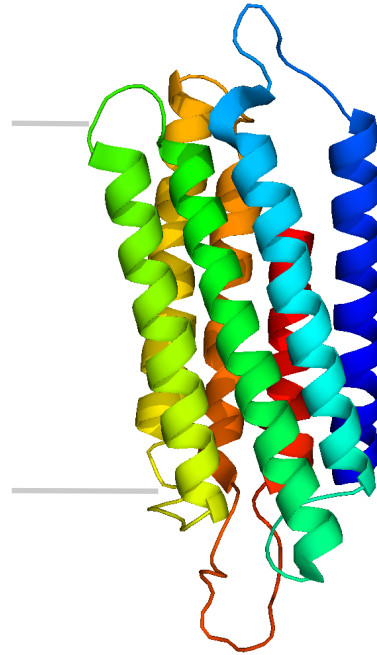
# 11 medically important membrane proteins of unknown structure predicted

OCTN1



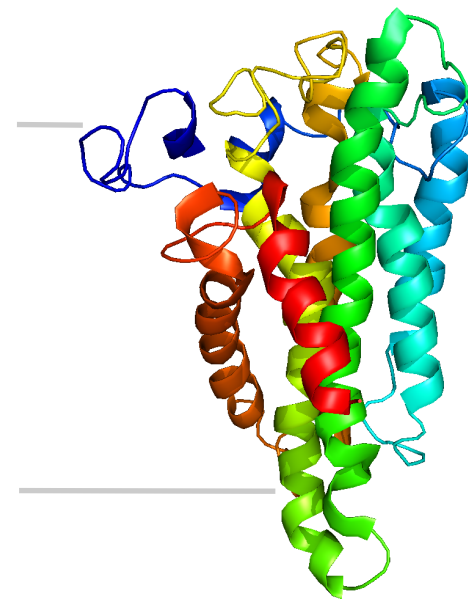
Crohn's disease,  
rheumatoid arthritis

Adiponectin receptor 1



diabetes, obesity,  
cancer

MT-ND1



LHON, MELAS,  
Alzheimer, Parkinson