

#### http://fraternalilab.kcl.ac.uk





## Sequence-structure gap

Despite the efforts from structural genomic projects the number of new structures per year has decreased

NEW sequencing techniques are becoming routinely available to scientists.....

Many genomes have been completely sequenced During the last 250 years, 1.2 million eukaryotic species have been identified and taxonomically classified.

Number of species estimated to exist on Earth: bacterial and archaea species, from 100,000 to 10 million; eukaryotic species, approximately 8.7 million.





Function is attributed to very few atoms absolutely conserved during the evolutionary process

## The Centrality of a 3D Structure



functional sites

motifs

electrostatics, clefts, patches





### Protein Structure Prediction: sequence vs structure

- The biochemical function (activity) of a protein is defined by its interactions with other molecules.
- The biological function is in large part a consequence of these interactions.
- The 3D structure is more informative than sequence because interactions are determined by residues that are close in space but are frequently distant in sequence.



As mentioned before structure is more conserved in evolution than sequence.

The net result is that patterns in space are frequently more recognizable than patterns in sequence.

### Protein Structure Prediction: the gap in numbers

#### Why Protein Structure Prediction?

	2018
Sequences (UniProt)	~556,000
Structures (PDB)	~14,000

UniProt: repository of protein sequences (www.uniprot.org) PDB: repository of protein structures (www.rcsb.org)

### Protein Structure Prediction: the gap in numbers

Chart 2

#### Single protein structures

#### Chart 1 9000 No structural data 8000 Structural domain Partial model 7000 Partial struct. Complete model 6000 Complete struct. Proteins 5000 4000 3000 2000 1000 0 Human Ecoli Fly Hpylori Worm Yeast Structural coverage of interactors Current Opinion in Structural Biology

#### Interaction protein pair structures



#### Structural coverage of interactors.

Structural coverage of interactions.

Stein, A., Mosca, R. & Aloy, P. Three-dimensional modeling of protein interactions and complexes is going 'omics. Curr. Opin. Struct. Biol. 21, 200–208 (2011).

## Protein Structure Prediction: Principles



## Folding

(physics)

Ab initio prediction

## **Evolution**

("statistical" rules)

Threading Comparative Modeling

### **Protein Structure Prediction: Principles**

During evolution protein structure is more conserved than sequence

The EMBO Journal vol.5 no.4 pp.823-826, 1986

The relation between the divergence of sequence and structure in proteins

Cyrus Chothia<sup>1</sup> and Arthur M.Lesk<sup>2</sup>



## Levitt conformational preferences of aa in globular proteins Biochemistry (1978) 17, 427

TABLE VI: Helix, $\beta$ Si	Conformational Pre heet and Reverse Tu	ferences <sup>a</sup> of Amin uns.	o Acids for α
type of secondary struct	favoring (h)	indifferent (i)	breaking (b)
$\alpha$ helix	Ala, Leu, Met, His. Glu, Gln. Lys, (Cys)	Val. Ile. Phe, Trp. Asp, Asn, Arg	Tyr. Thr. Gly, Ser. Pro
$\beta$ sheet	Val, Ile, Phe, (Trp), Tyr, Thr	Ala, Leu. Met, His. Gly, Ser. Arg	Glu. Gln. Lys, Asp, Asn, Pro. cys
reverse	Gly. Ser. Asp.	Gly. Gln, Lys,	Ala, Leu, Met.
turn	Asn. Pro	Tyr. Thr. (Arg)	His, Val, Ile. Phe, (Trp). (Cys), (Arg)

<sup>d</sup> These preferences are assigned with at least 75% confidence for the h and b classes (24 out of 40 are with at least 95% confidence). unless the amino acid is enclosed in parentheses when the confidence is as low as 56%. The confidence with which the i structure code is assigned is generally lower than for the h and b structure codes. In 1974 there were insufficient data (less than 2500 residues) to accurately determine the values.

The propensities were recalculated several times as more data became available, In 1998 using a dataset of over 33,000 residues leading to some notable differences in the two sets of propensities.

By 2004 the proteins structure datasets were large enough to derive residue propensities at different positions within alpha-helices.

One of the most accurate determinations of beta-turns residue propensities was by the group of Janet Thornton in 1999 based on almost 4000 beta-turns.

#### Conformational propensities



#### Conformational propensities



#### strong formers, formers, indifferent, breakers, strong breakers

## Ramachandran Plot



Note the clustering of low-energy states of single residues.

### Conformational propensities



## Secondary Structure Propensities

(Koehl-Levitt, 1999)

#### Alpha-helix propensity derived from designed sequences

A/L	R/K	N/M	D/F	C/P	Q/S	E/T	G/W	Н/Ү	I/V
-0.04	-0.30	0.25	0.27	0.57	-0.02	-0.33	1.24	-0.11	-0.26
-0.38	-0.18	-0.09	-0.01	0.	0.15	0.39	0.21	0.05	-0.06

#### Beta-sheet propensity derived from designed sequences

A/L	R/K	N/M	D/F	C/P	Q/S	E/T	G/W	Н/Ү	I/V
-0.12	0.34	1.05	1.12	-0.63	1.67	0.91	0.76	1.34	-0.77
0.15	0.29	-0.71	-0.67	0.	1.45	-0.70	-0.14	-0.49	-0.70

http://www.genome.jp/aaindex/

Anfinsen "Thermodynamic Hypothesis" states that the native conformation of a protein is adopted spontaneously.

In other words, there is sufficient information contained in the protein sequence to guarantee correct folding from any of a large number of unfolded states.

The Anfinsen experiment



A reasonable objection can be raised to the above result by suggesting that perhaps RNaseA was not completely unfolded in 8 M urea.

To address this class of objections, RNAseA was first reduced and denatured as above.

But in the second phase, the enzyme was first oxidized to form S-S bonds, and then the urea was removed, i.e. the order of steps in the second phase of the experiment was reversed.

The resulting activity was only about 1-2% of the untreated enzyme.

Sequence analysis showed a random assortment of S-S

## **Protein Structure Prediction**

- In theory, a protein structure can be solved computationally
- A protein folds into a 3D structure to minimizes its free potential energy
- The problem can be formulated as a search problem for minimum energy
- the search space is enormous
- the number of local minima increases exponentially

Computationally it is an exceedingly difficult problem



# Levinthal paradox

In 1969 Cyrus Levinthal noted that, because of the very large number of degrees of freedom in an unfolded polypeptide chain, the molecule has an astronomical number of possible conformations.

The estimate 3<sup>300</sup> or 10<sup>143</sup> appears in the original article. If the protein is to attain its corrected folded configuration by sequentially sampling all the possible conformations, it would require a time longer than the age of universe to arrive at its correct native conformation.

Levinthal himself was aware that proteins fold spontaneously and on short timescales, and that a random conformational search is therefore impossible.

Christian B. Anfinsen's 1971 Nobel Prize lecture revisits some of the same themes.

Protein folded states: 'explored' as contained in the PDB structures

- But how does a protein fold?
- According to Anfinsen and Levinthal a protein cannot visit all the possible φ and Ψ values before finding the native structure





It is likely that folding mechanisms vary significantly according to protein size, stability and structure.

The nucleation-condensation model has been supported by experimental evidence from several small proteins including chymotrysin inhibitor-II and barstar.

Bychkova and Ptitsyn have studied more than 20 proteins and found that nearly all adopted a molten globule state under mild denaturing conditions.

This points to the hydrophobic collapse model, a model favoured by many for the case of larger proteins.

But what if ... we explore the 'knowledge' of the structurally determined protein folded states: the 'explored' ones, contained in the PDB structures

How Can We Compare Sequences ? The Twilight Zone



### Some Basic Principles: Sequence identity

How is sequence identity defined?

It is the **fraction** of identical amino acids correctly aligned



14 residues 5 of which identical in the aligned positions:  $5/14 \times 100 = 35\%$  identity

# Structure Prediction: state of the art

#### Ab initio folding (force-field and simulation based)

1998 Duan and Kollman :36 residues, 1000 ns, 256 processors, 2 months

Recently examples of folding small proteins via computer simulations has been achieved V. PANDE (see movie)

## Ab initio folding (knowledge-based scoring functions)

Rosetta (BAKER) I-Tasser (ZHANG)

Deep Learning methods (from 2018)

#### Template-based (or knowledge-based) methods

- Homology modeling: sequence-sequence alignment, works if sequence identity > 30%
- Threading
- Protein threading: sequence-structure alignment, can go beyond the 25% limit

# Comparative modeling overview

Why build comparative models?

- Many more sequences available than structures (millions vs. tens of thousands)
- Many applications (e.g. determination of function) rely on structural information
- Structure is often more conserved than sequencesince evolution tends to preserve function,

# Comparative modeling overview

### How does it work?

- Extract information from known structures (one or more templates), and use to build the structure for the 'target' sequence
- Should also consider information from other sources: physical force fields, statistics (e.g. PDB mining)
- Classes of methods for comparative modeling
  - Assembly of rigid bodies (core, loops, sidechains)
  - Segment matching
  - Satisfaction of spatial restraints

# Comparative modeling by satisfaction of spatial restraints - MODELLER



2. Extract spatial restraints

3. Satisfy spatial restraints



A. Šali & T. Blundell. J. Mol. Biol. 234, 779, 1993. J.P. Overington & A. Šali. Prot. Sci. 3, 1582, 1994. A. Fiser, R. Do & A. Šali, Prot. Sci., 9, 1753, 2000.

Some Basic Principles: Homology

Pragmatic Definition of Homology:

The probability of two sequences to share more than 30% sequence identity by CHANCE is so low, one can safely assume that they share a common ancestor.

# 1. Align sequence with structures

First, must determine the template structures

- Simplistically, try to align the target sequence against every known structure's sequence
- In practice, this is too slow, so heuristics are used (e.g. BLAST)
- Profile or HMM searches are generally more sensitive in difficult cases (e.g. Modeller's profile.build method, or PSI-BLAST)
- Could also use threading or other web servers
- Alignment to templates generally uses global dynamic programming
  - Sequence-sequence: relies purely on a matrix of observed residue-residue mutation probabilities ('align')
  - Sequence-structure: gap insertion is penalized within secondary structure (helices etc.) ('align2d')
  - Other features and/or user-defined ('salign') or use an external program

# 2. Extract spatial restraints

Spatial restraints incorporate homology information, statistical preferences, and physical knowledge

- Template Cα- Cα internal distances
- Backbone dihedrals (φ/ψ)
- Sidechain dihedrals given residue type of both target and template
- Force field stereochemistry (bond, angle, dihedral)
- Statistical potentials
- Other experimental constraints
  etc.



## **Comparative Modelling**



(Adapted from slides in https://salilab.org/modeller/london.zip)
### 3. Satisfy spatial restraints

All information is combined into a single objective function

- Restraints and statistics are converted to an "energy" by taking the negative log
   Force field (CHARMM 22) simply added in
- Function is optimized by conjugate gradients and simulated annealing molecular dynamics, starting from the target sequence threaded onto template structure(s)
- Multiple models are generally recommended; 'best' model or cluster or models chosen by simply taking the lowest objective function score, or using a model assessment method such as Modeller's own DOPE or GA341, fit to EM density, or external programs such as PROSA or DFIRE

#### The UniProt database



The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.

UniProtKB	UniRef	UniParc	Proteomes	News 🔰 🚹 🗟	J
Swiss-Prot (546,439)	Sequence clusters	Sequence archive	** 1	Small is beautiful (and useful)   Evidences in the UniProtKB flat file format UniProt release 2014_09	
annotated and reviewed.	Supporting data			Ubiquitin caught at its own game	ľ
TrEMBL (83,955,074) Automatically	Literature citations	Taxonomy	Subcellular locations	the FTP site UniProt release 2014_08	l
annotated and not reviewed.	Cross-ref. databases ਛੋਮ ਨ ਛੋ	Diseases XXX	Keywords	Lark or owl? PER3 is the answer	*

#### Getting started

#### Q Text search

Our basic text search allows you to search all the resources available

#### SBLAST

Find regions of similarity between your sequences

#### Sequence alignments

Align two or more protein sequences using the Clustal Omega program

#### ⊥ Upload lists

This tool merges the "Retrieve" and "ID

#### UniProt data

You Tube

± Download latest release Get the UniProt data

#### 山 Statistics

View Swiss-Prot and TrEMBL statistics

#### Forthcoming changes Planned changes for the UniProt

knowledgebase

#### ESubmit your data

Submit your sequences and annotation updates

#### Protein spotlight



Moving Forward September 2014 Nature's imagination seems endless, and so is

Man's. For as long as humans have existed, they have twisted Nature to meet their own needs. Wood has been used to keep them warm. Whale oil has been used to make light. Water has been harnessed to make electricity. And when the era of bioengineering developed, it was not long before scientists found ways to tinker with an organism's genome for the benefits of mankind...

The UniProt database



#### **Comparative Modelling**



#### The FASTA format

		UniProtKB+	Advanced 🗸 🕻	2
••• < 8				1 and
BLAST Align Upload Lists			Help Con	itact
			∰ Baske	et 🗸
DUNHINI - CHET_	HAL53			
Protein	Chemotaxis protein CheY			
Gene	cheY			
Organism	Halobacterium salinarum (strain ATCC 29341 / DSM 6	571 / R1)		
Status	Reviewed - 👀 🍽 - Experimental evidence at	protein level <sup>i</sup>		
Display None	SBLAST ≅Align ⊡Format @Add to basket	O History	🖍 Comment (0)  Feedback 🗳 Help v	ideo
FUNCTION	Function <sup>i</sup> View format	×		
V NAMES & TAXONOMY	Involved in the transm FASTA (canonical)	emoreceptors and photore	ceptors to the flagellar motors. 💞 1 Publication 👻	
SUBCELLULAR LOCATION	GO - Molecular funct XML			
PATHOLOGY & BIOTECH	GFF			
PTM / PROCESSING	GO - Biological proce • chemotaxis & Source:			
EXPRESSION	Complete GO annotation			

>sp|B0R4K1|CHEY\_HALS3 Chemotaxis protein CheY OS=Halobacterium salinarum (strain ATCC 29341 / DSM 671 / R1) GN=cheY PE=1 SV=1 MAKQVLLVDDSEFMRNLLREILEEEFEIADEAENGVEAVEMYKEYDPDLVMMDIVMPIRD GIEATSEIKEFDAGAHIIMCTSIGQEEKMKKAVKAGADGYITKPFQKPSVMDAISDVLTA

### **Comparative Modelling**



Task:

- Query: new sequence (300 aa)
- Database (searching space): very many sequences
- Goal: find sequences related to query

We want:

- fast tool
- primarily a filter: most sequences will be unrelated to the query
- fine-tune the alignment later



BLAST is a program designed for rapidly comparing your sequence with every sequence in a database and report the most similar sequences

A good general reference is in wikipedia http://en.wikipedia.org/wiki/BLAST

BLAST

(http://blast.ncbi.nlm.nih.gov/Blast.cgi)

BLAST®	Basic Local Alignment Search Tool					
Home Recen	t Results Saved Strategies Help					
NCBI/ BLAST/ blastp	suite Standard Protein BLAST					
blastr blastr blastr tblastr tblastr						
	BLASTP programs search protein databases using a protein query. more					
Enter Query S	aquence					
Enter accession	number(s), gi(s), or FASTA sequence(s) 😣 Clear Query subrange 😣					
>sp B0R4K1 CH (strain ATCC 29)	EY_HALS3 Chemotaxis protein CheY OS=Halobacterium salinarum S41 / DSM 671 / R1) GN=cheY PF=1 SV=1					
MAKQVLLVDDSE	FMRNLLREILEEEFEIADEAENGVEAVEMYKEYDPDLVMMDIVMPIRD					
GIEATSEIKEFDAG	AHIIMCTSIGQEEKMKKAVKAGADGYITKPFQKPSVMDAISDVLTA					
Or, upload file	(Chance Site) no file selected					
Job Title						
	Enter a descriptive title for your BLAST search 😣					
□ Align two or m	ore sequences 😡					
Choose Search	n Set					
Database	Protein Data Bank proteins(pdb)					
Organism	Enter organism name or id-completions will be suggested Decklored Exclude					
optional	Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown. 🛞					
Exclude	□ Models (XM/XP) □ Uncultured/environmental sample sequences					
Optional Entropy Owners						
Optional	Fotor on Entror query to limit search Q					
Program Selection						
Algorithm  e blastp (protein-protein BLAST)						
	O PSI-BLAST (Position-Specific Iterated BLAST)					
	O PHI-BLAST (Pattern Hit Initiated BLAST)					
	O DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)					
	Choose a BLAST algorithm 😣					
BLAST	Search database Protein Data Bank proteins(odb) using Blasto (protein-protein BLAST)					
DENST	□ Show results in a new window					

#### Understanding the BLAST output



#### Descriptions

Sequences producing significant alignments:	
Select: All None Selected:0	
🛱 Alignments 🐷 Download 👻 <u>GenPept</u> <u>Graphica</u> <u>Distance tree of results</u> <u>Multiple alignment</u>	0
Description	Max Total Query E score score cover value Ident Accession
Chain A, The Crystal Structure Of An Activated Thermotoga Maritima Chevy With N- Terminal Region Of Flim [Thermotoga maritima MS88]	115 115 97% 2e-33 49% 4IGA A
Chain A, Chey From Thermotoga Maritima (Mn-Iii) [Thermotoga maritima]	115 115 97% 3e-33 49% 3TMY A
Chain Y, Chemotaxis Kinase Chea P2 Domain In Complex With Response Regulator Chey From The Thermophile Thermotoga Maritima [Thermotoga maritima]	114 114 96% 1e-32 49% 1UOS Y
Chain A, Crystal Structure Of The Chex-Chey-Bef3-Mg+2 Complex From Borrelia Burgdorferi [Borrelia burgdorferi]	81.6 81.6 95% 7e-20 40% 3HZH A
Chain A, Computational Design Of An Eight-Stranded (BetaALPHA)-Barrel From Fragments Of Different Folds [Thermotoga maritima]	82.4 82.4 75% 1e-19 45% 2LLE A
Chain A, Crystal Structure Of Rv1626 From Mycobacterium Tuberculosis [Mycobacterium tuberculosis H37Rv]	80.9 80.9 93% 3e-19 42% <u>1S8N A</u>
Chain X, A BetaALPHA-Barrel Built By The Combination Of Fragments From Different Folds [Thermotoga maritima]	79.7 79.7 75% 9e-19 45% 3CWO X
Chain A, Crystal Structure Of Two-component Response Regulator, Luxr Family, From Aurantimonas Sp. Si85-9a1 [Aurantimonas manganoxydans Si85-9A1]	78.2 78.2 97% 2e-18 34% 3CZ5 A
Chain A, Domain-Swapping In The Sporulation Response Regulator Spo0a [Geobacillus stearothermophilus]	72.4 72.4 97% 1e-16 34% 1DZ3 A
Chain C, Structure Of A Histidine Kinase-response Regulator Complex Reveals Insights Into Two-component Signaling And A Novel Cis- Autophosphorylation Mechanism [Thermotoga maritima]	71.6 71.6 99% 2e-16 36% <u>3DGE C</u>

#### Alignments

■Download → <u>GenPept</u> <u>Graphics</u>		▼ Next ▲ Previous 🎍 Descriptions
Chain A, The Crystal Structure Of An Activated Thermotoga Maritima Chey With N- Terminal Region Of Flim Sequence ID: <u>adb/4IGAIA</u> Length: 123 Number of Matches: 1	Alignments	Related Information
Next Match & Previous Match         Score       Expect Method       Identities       Positives       Gaps         115 bits(289) 2e-33       Compositional matrix adjust. 58/18(49%) 85/118(72%) 1/118(0%)       Identities       Positives       Gaps         0uerry       1       MakovLUvbSIFMMLLRILLET-FIENDEARMGVEAVETMXETDPLW+T01MPE IK 59       MikerVLVbBIFMMLLRILLET-FIENDEARMGVEAVETMXETDPLW+T01MPE IK 59       MikerVLVbBIFMMLLRILLET-FIENDEARMGVEAVETMXETDPLW+T01MPE IK 59         %bjet       4       MSKVLVDDARFMMHLKDITKAGYTVAGEATMGREAVEKYKELKEDIVTMDITMPIM 63       0         Querry       60       DEIEARDETHM-TRIGGERKMKAWAGABGETTMKFTGVCS3MD 61130V 117		<u>Structure</u> -3D structure displays

#### BLAST Scoring (E-value)

To assess whether a given alignment constitutes evidence for homology, it helps to know how strong an alignment can be expected from chance alone.

In this context, "chance" can mean the comparison of (i) real but non-homologous sequences; (ii) real sequences that are shuffled to preserve compositional properties; or (iii) sequences that are generated randomly based upon a DNA or protein sequence model.

(http://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-1.html)

Useful rules of thumb:

The E-value indicates the probability of encountering a false positive (expected by chance) in the sequence hit list. Therefore if one expects many hits from the database (e.g. 100) the E-values should be low enough to exclude the chance of non related hit.

E-values between 10<sup>-4</sup> and 10<sup>-6</sup> should be sufficient to exclude this chance.

If one expects few hits from the database (e.g. 1 to 10) the E-values should be high enough to allow distant sequences to be included in the hit list.

E-values between 10<sup>-1</sup> and 10<sup>-4</sup> should be high enough for this purpose.

Understanding the BLAST output

#### Hit List

Select: All None Selected:0							
Alignments EDownload <u>GenPept</u> Graphics Distance tree of results Multiple alignment						<	
Description	Max score	Total score	Query cover	E value	Ident	Accessio	
Chain A, The Crystal Structure Of An Activated Thermotoga Maritima Chey With N- Terminal Region Of Flim [Thermotoga maritima MSB8]	115	115	97%	2e-33	49%	4IGA_A	
Chain A, Chey From Thermotoga Maritima (Mn-Iii) [Thermotoga maritima]	115	115	97%	3e-33	49%	<u>3TMY_A</u>	
Chain Y, Chemotaxis Kinase Chea P2 Domain In Complex With Response Regulator Chey From The Thermophile Thermotoga Maritima [Thermotoga mar	<u>i</u> 114	114	96%	1e-32	49%	<u>1U0S_Y</u>	
Chain A, Crystal Structure Of The Chex-Chey-Bef3-Mg+2 Complex From Borrelia Burgdorferi [Borrelia burgdorferi]	81.6	81.6	95%	7e-20	40%	<u>3HZH_A</u>	
Chain A, Computational Design Of An Eight-Stranded (BetaALPHA)-Barrel From Fragments Of Different Folds [Thermotoga maritima]	82.4	82.4	75%	1e-19	45%	<u>2LLE_A</u>	
Chain A, Crystal Structure Of Rv1626 From Mycobacterium Tuberculosis [Mycobacterium tuberculosis H37Rv]	80.9	80.9	93%	3e-19	42%	<u>1S8N_A</u>	
Chain X, A BetaALPHA-Barrel Built By The Combination Of Fragments From Different Folds [Thermotoga maritima]	79.7	79.7	75%	9e-19	45%	<u>3CWO_X</u>	
Chain A, Crystal Structure Of Two-component Response Regulator, Luxr Family, From Aurantimonas Sp. Si85-9a1 [Aurantimonas manganoxydans Si85-9	78.2	78.2	97%	2e-18	34%	<u>3CZ5_A</u>	
Chain A, Domain-Swapping In The Sporulation Response Regulator Spo0a [Geobacillus stearothermophilus]	72.4	72.4	97%	1e-16	34%	<u>1DZ3_A</u>	
Chain C, Structure Of A Histidine Kinase-response Regulator Complex Reveals Insights Into Two-component Signaling And A Novel Cis- Autophosphoryla	t 71.6	71.6	99%	2e-16	36%	3DGE_C	

### **Comparative Modelling**



### **Comparative Modelling: Template Alignment**

Multiple Sequence Alignment

	A
target	QVQLQESGGDLVQPGGSLKLSCAASGFTFSSYGMSWVRQTPDKRL
1st-hit	EVQLVESGGDLVKPGGSLKLSCAASGFSFSSYGMSWVRQTPDKRL
2st-hit	EVKLVESGGGLVQPGGSLKLSCATSGFTFSDYYMYWVRQNSEKRL
3st-hit	EVQLVESGGGLVQPGGSLRLSCAASGFNIKDTYIHWVRQAPGKGL
4st-hit	. VQLQQSGTELMKPGRSLKISCKTTGYIFSNYWIEWVKQRPGHGL
5st-hit	QGQLQQSGAELVRPGSSVKISCKASGYAFSSFWVNWVKQRPGQGL
6st-hit	QVQLQQSGAELVRPGASVKLSCKASGYTFISYWINWVKQRPGQGL
7st-hit	QIQLVQSGPELKKPGETVKISCKASGYAFTNYGVNWVKEAPGKEL
8st-hit	IVLTQSPASLAVSLGQRATISCRASKSVSTSGYSHIHWYQQKPGQPP
9st-hit	IVLTQSPGSLAVSLGQRATISCRASESVDDDGNSFLHWYQQKPGQPP
10st-hit	IVLTQSPGSLAVSLGQRATISCRASESVDDDGNSFLHWYQQKPGQPP
11st-hit	IVLTQSPGSLAVSLGQRATISCRASESVDDDGNSFLHWYQQKPGQPP
12st-hit	IVLTQSPGSLAVSLGQRATISCRASESVDDDGNSFLHWYQQKPGQPP
13st-hit	.VMTQSPASLVVSLGQRATISCRASESVDSYCKSFMHWYQQKPGQPP
14st-hit	. VMTQSPSSLSVSAG <mark>ER</mark> VTMSCKSSQSLLNSGNQKNFLAWYQQKPGQPP
15st-hit	. VMTQSPSSLSVSAG <mark>ER</mark> VTMSCKSSQSLLNS <mark>GNQK</mark> NFLAWYQQKPGQPP
16st-hit	. VMTQSPSSLSVSAG <mark>ER</mark> VTMSCKSSQSLLYKDG. KNFLAWYQQKPGQPP
17st-hit	. VMTQSPSSLTVTAGEKVTMSCTSSQSLFNSCKQKNYLTWYQQKPGQPP
18st-hit	. VMTQSPSSLTVTAGEKVTMSCTSSQSLFNSCKQKNYLTWYQQKPGQPP

A widely used software for sequence alignments is

T-Coffee <u>http://tcoffee.vital-it.ch/apps/tcoffee/do:regular</u>

### **Comparative Modelling**



### **Comparative Modelling: Model Building**

Comparative Modeling by Satisfaction of Spatial Restraints MODELLER MODBASE Swissmodel

> **3D** GKITFYERGFQGHCYESDC-NLQP... SE GKITFYERG---RCYESDCPNLQP...

1. Extract spatial restraints 2. Satisfy spatial restraints

#### http://www.salilab.org/modeller

(Adapted from slides in https://salilab.org/modeller/london.zip)

A. Šali & T. Blundell. J. Mol. Biol. 234, 779, 1993.
J.P. Overington & A. Šali. Prot. Sci. 3, 1582, 1994.
A. Fiser, R. Do & A. Šali, Prot. Sci., 9, 1753, 2000.

### Loop modeling

- Often, there are parts of the sequence which have no detectable templates (usually loops)
- "Mini folding problem" these loops must be sampled to get improved conformations
- Database searches only complete for 4-6 residue loops
- Modeller uses conformational search with a custom energy function optimized for loop modeling (statistical potential derived from PDB)
  - Fiser/Melo protocol ('loopmodel')
  - Newer DOPE + GB/SA protocol ('dope\_loopmodel')

## Accuracy of loop models as a function of amount of optimization



# Fraction of loops modeled with medium accuracy (<2Å)



### **Comparative Modelling**



### **Comparative Modelling**



### Typical errors in comparative models



Marti-Renom et al. Annu.Rev.Biophys.Biomol.Struct. 29, 291-325, 2000.

### Model Accuracy as a Function of Target-Template Sequence Identity



Sánchez, R., Šali, A. Proc Natl Acad Sci U S A. 95 pp13597-602. (1998).

### Model accuracy



Marti-Renom et al. Annu.Rev.Biophys.Biomol.Struct. 29, 291-325, 2000.

### Applications of protein structure models



D. Baker & A. Sali. Science 294, 93, 2001.

### **Comparative Modelling: Model Evaluation**

### Ramachandran plot

- B. Beta strand
- A. Right handed helix
- L. Left handed helix
- Color coding
  - White. Disallowed
  - Red. Most favorable
  - Yellow Allowed rogion

unit

• Glycine 1



Comparative Modelling: Model Evaluation Let's find the wrong structure!





### Procheck: Bond Lenghts









**Comparative Modelling: Model Evaluation** 

Let's find the wrong structure!



#### **Comparative Modelling: Model Evaluation**



About

#### Citation

Please cite the following paper when referring to MolProbity in print or during a presentation:

Simon C. Lovell, Ian W. Davis, W. Bryan Arendall III, Paul I. W. de Bakker, J. Michael Word, Michael G. Prisant, Jane S. Richardson, David C. Richardson (2003) Structure validation by C-alpha geometry: phi, psi, and C-beta deviation. Proteins: Structure, Function, and Genetics. 50: 437-450.

The following paper is also relevant when using MolProbity with nucleic acids:

Ian W. Davis, Laura Weston Murray, Jane S. Richardson, David C. Richardson (2004) MolProbity: structure validation and all-atom contact analysis for nucleic acids and their complexes. Nucleic Acids Research. 32: W615-W619.

#### http://molprobity.biochem.duke.edu/

#### Rotamer outliers 2.90% Goal: <1% Ramachandran outliers 6.00% Goal: <0.2% Ramachandran favored 77.33% Goal: >98% Protein Geometry CB deviations >0.25Å Goal: 0 5 Residues with bad bonds: 0.00% Goal: <1% Residues with bad angles: Goal: <0.5% 2.63%

#### 1sddb

#### 1czta

	Rotamer outliers	1.46%	Goal: <1%
	Ramachandran outliers	2.01%	Goal: <0.2%
Protein	Ramachandran favored	89.26%	Goal: >98%
Geometry	Cβ deviations >0.25Å	0	Goal: 0
	Residues with bad bonds:	0.66%	Goal: <1%
	Residues with bad angles:	1.99%	Goal: <0.5%

### Calmodulin

### How



### Secondary Structure Assignment

(given a 3D structure, assign secondary structural elements)

The DSSP program defines 7 secondary structure states

H : alpha helix

B : residue in isolated beta-bridge

- E : extended strand, participates in beta ladder
- G: 3-helix (3/10 helix)
- I:5 helix (pi helix)
- T : hydrogen bonded turn
- S:bend

The secondary structure assignment with DSSP over a database of structures can be used as 'standard of truth' for secondary structure prediction methods.

Kabsch & Sander, Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. Biopolymers, 22(12), 2577-2637 (1983). <u>http://swift.cmbi.ru.nl/gv/dssp/</u> <u>http://www.cmbi.ru.nl/hsspsoap/</u>

web-server

#### Secondary Structure Assignment of a PDB file

Output

The output from DSSP contains secondary structure assignments and other information, one line per residue. Extract from 1est.dssp (simplified):

HEADER 17-MAY-76 1EST HYDROLASE (SERINE PROTEINASE) . . . 240 0 TOTAL NUMBER OF RESIDUES, NUMBER OF CHAINS, - 1 NUMBER OF SS-BRIDGES (TOTAL, INTRACHAIN, INTERCHAIN) ACCESSIBLE SURFACE OF PROTEIN (ANGSTROM\*\*2) 10891.0 162 67.5 TOTAL NUMBER OF HYDROGEN BONDS OF TYPE O(I)-->H-N(J) ; PER 100 RESIDUES 0 0.0 TOTAL NUMBER OF HYDROGEN BONDS IN PARALLEL BRIDGES; PER 100 RESIDUES 84 35.0 TOTAL NUMBER OF HYDROGEN BONDS IN ANTIPARALLEL BRIDGES; PER 100 RESIDUES 26 10.8 TOTAL NUMBER OF HYDROGEN BONDS OF TYPE O(I) -->H-N(I+2) 30 12.5 TOTAL NUMBER OF HYDROGEN BONDS OF TYPE O(I) -->H-N(I+3) 10 4.2 TOTAL NUMBER OF HYDROGEN BONDS OF TYPE O(1)-->H-N(1+4) . . . # RESIDUE AA STRUCTURE BP1 BP2 ACC N-H-->O O-->H-N N-H-->O O-->H-N 2 17 V B 3 +A 182 0A 8 180,-2.5 180,-1.9 1,-0.2 134,-0.1 ... Next two lines wrapped as a pair ... TCO KAPPA ALPHA PHI PSI X-CA 2-CA Y-CA -0.776 360.0 8.1 -84.5 125.5 -14.7 34.4 34.8 .-- sequential resnumber, including chain breaks as extra residues .-- original PDB resname, not nec. sequential, may contain letters .-- amino acid sequence in one letter code .-- secondary structure summary based on columns 19-38 xxxxxxxxxxxxxxxxxxx recommend columns for secstruc details -- 3-turns/helix .-- 4-turns/helix .-- 5-turns/helix .-- geometrical bend .-- chirality .-- beta bridge label .-- beta bridge label .-- beta bridge partner resnum .-- beta bridge partner resnum .-- beta sheet label -- solvent accessibility # RESIDUE AA STRUCTURE BP1 BP2 ACC 35 36 37 38 50 34 0 39 51 -KL 36 980

The secondary structure assignment with DSSP over a database of structures can be used as 'standard of truth' for secondary structure prediction methods.

## DSSP uses mainly hydrogen bond assessments to discriminate secondary structure elements

What is a Hydrogen Bond?

A hydrogen bond is formed when a proton (H) covalently attached to one electronegative donor atom (D) is shared with another electronegative acceptor atom (A).

One of the widely used schemes was proposed by Morokuma (1977) in which abinitio calculations describe the interaction energy of a hydrogen bond in terms of electrostatic, charge transfer, polarization, exchange repulsion and coupling.


### DSSP: H-bond energy cutoff

Hydrogen bonds in proteins have little wave-function overlap and are well described by an electrostatic model.<sup>13</sup> We calculate the electrostatic interaction energy between two H-bonding groups by placing partial charges on the C,O  $(+q_1, -q_1)$  and N,H  $(-q_2, +q_2)$  atoms, i.e.,

 $E = q_1 q_2 (1/r(\text{ON}) + 1/r(\text{CH}) - 1/r(\text{OH}) - 1/r(\text{CN}))*f$ 

with  $q_1 = 0.42e$  and  $q_2 = 0.20e$ , e being the unit electron charge and r(AB) the interatomic distance from A to B. In chemical units, r is in angstroms, the dimensional factor f = 332, and E is in kcal/mol. A good H bond has about -3 kcal/mol binding energy. We choose a generous cutoff to allow for bifurcated H bonds and errors in coordinates and assign an H bond between C=O of residue i and N-H of residue j if E is less than the cutoff, i.e., "Hbond(i,j)=: [E < -0.5kcal/mole]."

Kabsch W, Sander C., Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features, Biopolymers (1983) 22:2577-637.

### DSSP: H-bond energy cutoff



Fig. 1. H bond between peptide units is described here by the dominant electrostatic part E (see text) of the H-bond energy, drawn in contours of constant E at 0.5 kcal/mol intervals as a function of the distance, d, and the alignment angle  $\theta$ . Dotted lines, E positive or zero; broken lines, E negative. An ideal H bond has d = 2.9 Å,  $\theta = 0$ , and E = -3.0 kcal/mol. We assume an H bond for E up to -0.5 kcal/mol (solid line). Thus, misalignment of up to 63° is allowed at the ideal length; an N-O distance of up to d = 5.2 Å is allowed for perfect alignment. This definition of H bonds is particularly simple and physically meaningful. It is more general than the historical definition of hydrogen "bond" and could be called polar interaction.

# Model accuracy

## If we know the answer

$$RMSD = \sqrt{\frac{1}{N} \bullet \sum_{i,j} d_{i,j}^2}$$

Fraction correct =  $N_c/N$  $N_c$  = number correct (dij<4Å)



### Databases of multiple alignments ----> Domains

Very early in the days of protein sequence analysis, it was observed that some protein sequences contained long segments that were very similar to other proteins, while the rest of the sequence in that protein had no detectable similarity.

Today, we take more or less for granted that proteins are composed of domains, segments of sequence which have been joined together by genetic events during evolution so that the new protein has a function that is based on the activities of the domains it contains.

Often the domains detectable by sequence analysis correspond to structural domains in the 3D structure as well. There are now many well-documented cases where it has been shown that domains can exists perfectly well in isolation, when excised from the original protein. Surprisingly often, a domain can be expressed and folded all on its own.

There are today several databases that keep track of which domains have been discovered, which proteins are involved, and that store the multiple sequence alignments of the relevant segments of the protein sequences. We will discuss

one such of databases, Pfam.

Also, several of the primary sequence databases now contain information about the domains in the sequence entries.

The idea behind Pfam is twofold:

1. Create and maintain good-quality multiple sequence alignments of well-defined protein sequence domains from proteins in SWISS-PROT.

2. Use these multiple alignments for creating so-called HMMs(Hidden Markov Models) machine learning algorithms, that can be used in profile searches of sequence databases. https://en.wikipedia.org/wiki/Hidden\_Markov\_model

The multiple alignment used to define a domain (protein family) in Pfam are called the seed alignment. It is created by a curator, or taken from the literature. It is used to generate a profile HMM for identifying other sequences in the databases (SWISS-PROT and TREMBL) that contain the domain. The search results are inspected to decide which cutoff should be used for that particular Pfam entry. The search hits are then aligned automatically into a so-called full alignment.

There are today several databases that keep track of which domains have been discovered, which proteins are involved, and that store the multiple sequence alignments of the relevant segments of the protein sequences. We will discuss one such of databases, Pfam. Also, several of the primary sequence databases now contain information about the domains in the sequence entries.

Pfam

http://pfam.sanger.ac.uk/

The alignments can be converted into hidden Markov models (HMM), which can be used to search for domains in a query protein sequence. The software HMMER (by Sean Eddy) is the computational foundation for Pfam.

http://nar.oxfordjournals.org/content/early/2015/05/04/nar.gkv397.full http://hmmer.org/

HMMER is often used together with a profile database, such as <u>Pfam</u> or many of the databases that participate in <u>Interpro</u>. But HMMER can also work with query *sequences*, not just profiles, just like BLAST. For example, you can search a protein query sequence against a database with **phmmer**, or do an iterative search with **jackhmmer**.

The domain structure of protein sequences in SWISS-PROT and TrEMBL are available directly from the Pfam web sites, and it is also possible to search for domains in other sequences using servers at the web sites.

### Prosite

PROSITE www.expasy.ch/prosite/

PROSITE is a database of protein families and domains. It consists of biologically significant sites, patterns and profiles that help to reliably identify to which known protein family (if any) a new sequence belongs.

It was started by Amos Bairoch, is part of SWISS-PROT and is maintained in the same way as SWISS-PROT. The basis of it are regular expressions describing characteristic subsequences of specific protein families or domains. PROSITE has been extended to contain also some profiles, which can be described as probability patterns for specific protein sequence families.

The site above can be used to search by keyword or other text in the entries, to search for a pattern in a sequence, or to search for proteins in SWISS-PROT that match a pattern.



PROSITE

Home I Contact

Home | ScanProsite | ProRule | Documents | Downloads | Links | Funding



#### Database of protein domains, families and functional sites

PROSITE consists of documentation entries describing protein domains, families and functional sites as well as associated patterns and profiles to identify them [More... / References / Commercial users].

PROSITE is complemented by ProRule, a collection of rules based on profiles and patterns, which increases the discriminatory power of profiles and patterns by providing additional information about functionally and/or structurally critical amino acids [More...].

#### Forthcoming changes to the profile format

#### Release 20.111 of 04-Feb-2015 contains 1716 documentation entries, 1308 patterns, 1107 profiles and 1105 ProRule.



#### Domain prediction by Smart

SMARTEE	SMART MODE: NORMAL GENOMIC	Simple Modular Architecture	keywords Search SMART
Schultz et al. (1998) Proc. Natl. Acad. Sci. USA 95, 5857-5864		Research	
Letunic et al. (2012) Nucleic Acids Res, doi:10.1093/nar/gkr931		Tool	
HOME SETUP FAQ ABOUT GLOSSARY WHAT'S NEW FEEDBACK			

#### Select your default SMART mode

You can use SMART in two different modes: normal or genomic. The main difference is in the underlying protein database used. In Normal SMART, the database contains Swiss-Prot, SP-TrEMBL and stable Ensembl proteomes. In Genomic SMART only the proteomes of completely sequenced genomes are used; Ensembl for metazoans and Swiss-Prot for the rest. The complete list of genomes in Genomic SMART is available here.

The protein database in Normal SMART has significant redundancy, even though identical proteins are removed. If you use SMART to explore domain architectures, or want to find exact domain counts in various genomes, consider switching to **Genomic** mode. The numbers in the domain annotation pages will be more accurate, and there will not be many protein fragments corresponding to the same gene in the architecture query results. Remember you are exploring a limited set of genomes, though.

#### Different color schemes are used to easily identify the mode you're in.

Normal mode	Genomic mode	Genomic mode				
SMART MODE: NORMAL GENOMIC Res Too	smart mode: Normal GENOMIC	Simple Modular Architecture Research Tool				

#### Click on the images above to select your default mode.

Information about your selected mode is stored in a browser cookie. If you for whatever reason don't want/can't use cookies, access SMART through this page.

You can easily change modes later, by clicking on the links in the 'SMART MODE' header box, or in your personal preference settings ('SETUP' link in the menu):



http://www.cs.wright.edu/~mraymer/cs790/Homology\_Modeling.ppt

Economist.com

### Moore's law turns 50 Ever more from Moore

#### A microchip pioneer's prediction has a bit more life left in it

Apr 18th 2015 From the print edition

The Economist

NEWS of the death of Moore's law has always been greatly exaggerated. People started to pronounce it deceased not long after Gordon Moore, co-founder of Intel, a chipmaker, published on April 19th 1965 a paper arguing that the number of transistors that can be etched on a given surface area of silicon would double every year. In a later paper he corrected his forecast to every two years, which has come to be stated as his "law". Regularly proving sceptics wrong, however, the exponential growth kept going (see chart), driving the digital revolution.



ulations of n based on years over

Energy landscapes of proteins



http://www.youtube.com/watch?v=YANAso8Jxrk

Molecular Mechanics Force Fields provide a 'recipe' (equations and parameters) to calculate the potential energy of a protein from its atomic coordinates



the atomic Cartesian coordinates give forces on every atom. Each term is discussed more fully in the text.

**MM Force Fields** 

### Assignment of Atom Types (based on the chemical environment)

Table 1. List of Atom Types<sup>a</sup>

atom	type	description					
carbon	СТ	any sp <sup>3</sup> carbon	0.00	OW	en <sup>3</sup> ovugan in TID3D water		
	C	any carbonyl sp <sup>2</sup> carbon	oxygen	04	sp oxygen in rir or water		
	CA	any aromatic sp <sup>2</sup> carbon and (Ce of Arg)		OH	sp <sup>2</sup> oxygen in alcohols, tyrosine, and		HC
	CM	any sp <sup>2</sup> carbon, double bonded		05	protonated carboxyne acids		11
	CC	sp <sup>2</sup> aromatic in 5-membered ring with one		03	sp <sup>2</sup> oxygen in ethers		H
		substituent + next to nitrogen (Cy in His)		0	sp <sup>2</sup> oxygen in amides	- 🚺 LIA 👗	
	CV	sp <sup>2</sup> aromatic in 5-membered ring next to carbon		02	sp <sup>2</sup> oxygen in anionic acids	С 🚺 ПІЛ	
		and lone pair nitrogen (e.g. Cô in His (ô))	sulfur	S	sulfur in methionine and cysteine		
	CW	sp <sup>2</sup> aromatic in 5-membered ring next to carbon		SH	sulfur in cysteine		HC
		and NH (e.g. Co in His ( $\epsilon$ ) and in Trp)	phosphorus	Р	phosphorus in phosphates		
	CR	sp <sup>2</sup> aromatic in 5-membered ring next to	hydrogen	н	H attached to N CT	UHC U	
	an	two nitrogens (Cy and Ce in His)		HW	H in TIP3P water	H H	
	CB	sp <sup>2</sup> aromatic at junction of 5- and 6-membered		HO	H in alcohols and acids	нс ''	0 10
		rings (Co in 1rp) and both junction atoms		HS	H attached to sulfur		•
	<b>C</b> *	In Ade and Gua		HA	H attached to aromatic carbon		
	C.	two carbons (e.g. Cy in Trp)		HC	H attached to alighatic carbon with		
	CN	sp <sup>2</sup> junction between 5- and 6-membered rings		ne	no electron-withdrawing substituents		
	C. (	and bonded to CH and NH (Ce in Trp)		HI	H attached to alightic carbon with		
	CK	sp <sup>2</sup> carbon in 5-membered aromatic between N			one electron-withdrawing substituent		
		and N-R (C8 in purines)		112	H attached to alighatic carbon with		
	CO	sp <sup>2</sup> carbon in 6-membered ring between		<b>H</b> 2	two electron withdrawing substituents		
	-	lone pair nitrogens (e.g. C2 in purines)		112	We dectron-withdrawing substituents		
nitrogen	N	sp <sup>2</sup> nitrogen in amides		пэ	H attached to anphatic carbon with		
-	NA	sp2 nitrogen in aromatic rings with hydrogen			Infee electron-withdrawing substituents		
		attached (e.g. protonated His, Gua, Trp)		HP	H attached to carbon directly bonded to		
	NB	sp2 nitrogen in 5-membered ring with lone pair			formally positive atoms (e.g. C next to		
		(e.g. N7 in purines)			NH <sub>3</sub> . of lysine)		
	NC	sp <sup>2</sup> nitrogen in 6-membered ring with lone pair		H4	H attached to aromatic carbon with one		
		(e.g. N3 in purines)			electronegative neighbor (e.g. hydrogen on		
	N*	sp <sup>2</sup> nitrogen in 5-membered ring with carbon			C5 of 1rp, C6 of 1ny)		
		substituent (in purine nucleosides)		H5	H attached to aromatic carbon with two		
	N2	sp <sup>e</sup> nitrogen of aromatic amines and			electronegative neighbors (e.g. H8 of Ade and		
		guanidinium ions			Gua and H2 of Ade)		
	N3	sp <sup>o</sup> nitrogen			-		

MM Force Fields											
						• • • • •			CT H	IC	
				ŀ	IC V	Н					
Assignment of parameters							_		H1 🖌 HC		HC
Abolg milent of parameters						HC	K				
									N C	N	HC
							СТ	HC	Ų. 📕	V	
$U^{X-1}$	Y = K	$X^{X-Y} \times$	$(r - r^{X-1})$	Y) <sup>-</sup>				НС	H O	HC	
• stre	tch	r	(* * eq	)				110		7	
Bond Parameters											
bond	K <sup>b</sup>	req	bond	Kb	r <sub>eq</sub> c	bond	$K_r^b$	$r_{eq}^{c}$	bond	Kr <sup>b</sup>	$r_{eq}^{c}$
C-CA	469.0	1.409	CA-HA	367.0	1.080	CM-HA	367.0	1.080	CT-S	227.0	1.810
C-CB	447.0	1.419	CA-N2	481.0	1.340	CM-N*	448.0	1.365	CT-SH	237.0	1.810
C-CM	410.0	1.444	CA-NA	427.0	1.381	CN-NA	428.0	1.380	CV-H4	367.0	1.080
C-CT	317.0	1.522	CA-NC	483.0	1.339	CQ-H5	367.0	1.080	CV-NB	410.0	1.394
C-N	490.0	1.335	CB-CB	520.0	1.370	CQ-NC	502.0	1.324	CW-H4	367.0	1.080
C-N*	424.0	1.383	CB-CN	447.0	1.419	CR-H5	367.0	1.080	CW-NA	427.0	1.381
C-NA	418.0	1.388	CB-N*	436.0	1.374	CR-NA	477.0	1 343	HN	434.0	1.010
C-NC	457.0	1.358	CB-NB	414.0	1.391	CR-NB	488.0	1.335	H-N*	434.0	1.010
C-0	570.0	1.229	CB-NC	461.0	1.354	CT-CT	310.0	1.526	H-N2	434.0	1.010
C-02	656.0	1.250	CC-CT	317.0	1.504	CT-F	367.0	1.380	H-N3	434.0	1.010
C-OH	450.0	1.364	CC-CV	512.0	1.375	CT-H1	340.0	1.090	H-NA	434.0	1.010
C*-CB	388.0	1.459	CC-CW	518.0	1.371	CT-H2	340.0	1.090	HO-OH	553.0	0.960
C*-CT	317.0	1.495	CC-NA	422.0	1.385	CT-H3	340.0	1.090	HO-OS	553.0	0.960
C*-CW	546.0	1.352	CC-NB	410.0	1.394	CT-HC	340.0	1.090	HS-SH	274.0	1.336
C*-HC	367.0	1.080	CK-H5	367.0	1.080	CT-HP	340.0	1.090	O2-P	525.0	1.480
CA-CA	469.0	1.400	CK-N*	440.0	1.371	CT-N	337.0	1.449	OH-P	230.0	1.610
CA-CB	469.0	1.404	CK-NB	529.0	1.304	CT-N*	337.0	1.475	OS-P	230.0	1.610
CA-CM	427.0	1.433	CM-CM	549.0	1.350	CT-N2	337.0	1.463	OW-HW	553.0	0.9572
CA-CN	469.0	1.400	CM-CT	317.0	1.510	CT-N3	367.0	1.471	S-S	166.0	2.038
CA-CT	317.0	1.510	CM-H4	367.0	1.080	CT-OH	320.0	1.410			
CA-H4	367.0	1.080	CM-H5	367.0	1.080	CT-OS	320.0	1.410			

**Force Field Parametrisation** 

Force Fields parameters are derived from experimental data or calculations performed at a higher level of theory (Quantum Mechanics)

Equilibrium bond distances and angles: X-ray crystallography

**Bond and angle force constants**: vibrational spectra, normal mode calculations with Quantum Mechanics (QM)

**Dihedral angle parameters**: difficult to measure directly with experiments; fit to QM calculations for rotations around a bond with other motions fixed

**Atom charges**: fit to experimental liquid properties, ESP charge fitting to reproduce electrostatic potentials of high level QM, X-ray crystallographic electron density

van der Waals parameters: often most difficult to determine, fit to experimental liquid properties, intermolecular energy fitting

Force Fields for Biomolecules

FFs commonly used for biomolecules:

AMBER CHARMM GROMOS OPLS

Improved over time and validated against experimental data, including:

- experimental structures
- secondary structure propensities
- NMR data (e.g. order parameters, chemical shifts, NOEs etc...)



**Figure 3. Improvement of force fields over time.** For each force field, we assigned a score depending on the agreement with experiments in the tests presented here. Low scores indicate good agreement with experiments. These scores are plotted against the year in which the force field was published. For the force fields that involve multiple corrections (e.g., ff99SB\*-ILDN), we use the year of the most recently published correction.

Lindorff-Larsen, K. *et al. PLoS ONE* **7**, e32131 (2012).

# Ab-initio folding from first principles Force Fields and MD

### Study of folding mechanisms





http://www.youtube.com/watch?v=gFcp2Xpd29I

Voelz, V. A. et al. J. Am. Chem. Soc. 132, 1526–1528 (2010).

Simulation by Computer of a protein fast folder

### Ab initio folding (force-field and simulation based)

http://www.youtube.com/watch?v=gFcp2Xpd29I

https://www.youtube.com/watch?v=sD6vyfTtE4U

## Learning outcomes of this lecture:

Why do we need to predict protein structures Importance of 3D structure Knowledge Folding: Anfinsen theorem Levinthal paradox Structure Prediction: Comparative Modelling: which are the main steps involved? Main problems encountered in Comparative Modelling What is DSSP and which information gives? Assessing protein structures **Analysing Structures** Analysing Dynamics of Structures

# Suggested readings

Secondary structure predictions:

 DSSP: Kabsch W, Sander C.Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical

features. Biopolymers. 1983 22(12):2577-637.

- 2) Chou PY, Fasman GD. Prediction of the secondary structure of proteins from their amino acid sequence. Adv Enzymol Relat Areas Mol Biol. 1978;47:45-148.
- 3) Chen H, Gu F, Huang Z. Improved Chou-Fasman method for protein secondary structure prediction. BMC Bioinformatics. 2006 12:7 Suppl 4:S14

Tertiary structure predictions:

- I) Forster MJ (2002) Molecular Modelling in structural biology. Micron 33, 365-384
- 2) Sutcliffe MJ, Hayes FR, Blundell TL. Knowledge based modelling of homologous proteins, Part II: Rules for the conformations of substituted sidechains. Protein Eng. 1987 1(5):385-92.
- 3) Sutcliffe MJ, Haneef I, Carney D, Blundell TL. Knowledge based modelling of homologous proteins, Part I: Three-dimensional frameworks derived from the simultaneous superposition of multiple structures. Protein Eng. 1987 1(5):377-84.

### Very recent methods

Protein structure determination using metagenome sequence data. Ovchinnikov S, Park H, Varghese N, Huang PS, Pavlopoulos GA, Kim DE, Kamisetty H, Kyrpides NC, Baker D. Science. 2017 Jan 20;355(6322):294-298. doi: 10.1126/science.aah4043.

De novo structure prediction with deeplearning based scoring R.Evans \*,1 , J.Jumper \*,1 , J.Kirkpatrick \*,1 , L.Sifre \*,1 , T.F.G.Green 1 , C.Qin ,1 , A.Zidek 1 , A.Nelson 1 , A.Bridgland 1 , H.Penedones 1 , S.Petersen 1 , K.Simonyan 1 , S.Crossan 1 , D.T.Jones 2 , D.Silver 1 , K.Kavukcuoglu 1 , D.Hassabis 1 , A.W.Senior (Goodle Deepmind)

Deep learning extends de novo protein modelling coverage of genomes using iteratively predicted structural constraints. Greener JG, Kandathil SM, Jones DT. Nat Commun. 2019 Sep 4;10(1):3977. doi: 10.1038/s41467-019-11994-0.

## Web resources

### Introductory MD simulation tutorials and How-tos

### **GROMACS** (freely available)

http://www.gromacs.org/Documentation/Tutorials#General\_GROMACS\_Use http://www.gromacs.org/Documentation/How-tos

#### **NAMD** (freely available)

http://www.ks.uiuc.edu/Training/Tutorials/index-all.html#namd

#### **AMBER**

http://ambermd.org/tutorials/#basic\_tut

**VMD** (visualisation of MD trajectories, freely available) http://www.ks.uiuc.edu/Training/Tutorials/vmd/tutorial-html/

### YASARA

http://www.yasara.org/movies.htm

Servers **MDWeb** 

#### MDWeb

http://mmb.irbbarcelona.org/MDWeb/ (System set-up)

#### We-NMR

https://www.wenmr.eu/wenmr/molecular-dynamics-software (run simulations on the GRID, requires registration)

### Databases

**MolMovDB** (database of macromolecular motions, contains morphs between different conformations) http://www.molmovdb.org/