## Bioinformatics for molecular biology Structural bioinformatics tools, predictors, and 3D modeling – Structural Biology Review

Dr Jon K. Lærdahl, Research Scientist

Department of Microbiology, Oslo University Hospital - Rikshospitalet & Bioinformatics Core Facility/CLS initiative, University of Oslo

E-mail: jonkl@medisin.uio Phone: +47 22844784 Group: Torbjørn Rognes (http://www.ous-research.no/rognes) CF: Bioinformatics services (http://core.rr-research.no/bioinformatics) CLS: Bioinformatics education (http://www.mn.uio.no/ifi/english/research/networks/clsi) Main research area: Structural and Applied Bioinformatics







#### Now:

- Protein Structure Review
  - Amino acids, polypeptides, secondary structure elements, visualization, structure determination by X-ray crystallography and NMR methods, PDB

#### Later...

- Structure comparison and classification (CASP & SCOP)
- Predictors
- 3D structure modeling
  - Ab initio
  - Threading/fold recognition
  - Homology modeling
- Practical exercises
  - PyMOL & visualization
- Practical Exercises
  - Homology modeling of influenza neuraminidase (Tamiflu resistance?)
  - Other homology modeling
  - Threading
  - Your own project?

#### Stop me and ask questions!!



#### Amino acids – the building blocks Proteins are built from 20 naturally occurring amino Amino acids – the building blocks

acids. They have an

amino  $(-NH_2)$  and acidic

(-COOH) functional group

Jon K. Lærdahl, Structural Bioinformatics

The side chain group (R) determines the properties of the amino acid

Side chain OH. H<sub>2</sub>N H R HI  $^{+}H_{3}N$ Carboxyl <sup>+</sup>H<sub>3</sub>N group 0 Amino Zwitterionic form group found at physiological pH (Almost) invariably enantiomeric L-form (or S-form)

alpha carbon ( $C_{\alpha}$ )

## Amino acids

R-group properties:

• Large

Small

- Hydrophobic
   Aliphatic
  - Aromatic
- Polar
- Charged
  - Positive/negative charge

Increasing hydrophilicity/higher water (solvent) affinity

Structural Bioinformatics, Eds. P.E. Bourne & H. Weissig (Wiley, Hoboken, NJ, 2003)



#### Amino acids

Introduction to Protein Structure, C. Branden & J. Tooze (Garland, New York, 1998)



#### Amino acids

*Introduction to Protein Structure*, C. Branden & J. Tooze (Garland, New York, 1998)







## **Dihedral angles**

Proteins are polypeptides, *i.e.* many amino acids connected by peptide bonds



The peptide bond (light green) is a partial double bond and is fixed at  $\sim$ 180°, *i.e.* the green part is flat

Cis-form for peptide bond is extremely rare except for prolines (~25%).



The dihedral angles phi ( $\phi$ ) and psi ( $\psi$ ) determines the conformation of the peptide backbone

## Dihedral angles

Jon K. Lærdahl, **Structural Bioinformatics** 

connected by peptide bonds Н R Ø W Cα'

One  $(\phi, \psi)$  pair for each residue in a protein

Н

Structural Bioinformatics, Eds. P.E. Bourne & H. Weissig (Wiley, Hoboken, NJ, 2003)

R'

Н



### Ramachandran plot







#### Secondary structure – $\beta$ -sheets



#### Secondary structure – $\beta$ -sheets



#### Secondary structure – $\alpha$ -helices





Partial positive charge at Nterminus and negative charge at Cterminus, *i.e.* it is a *dipole* 

#### Secondary structure – 3 states

Three "states": α-helices (H) β-sheets (E) Loops/coils (C)





Loops/coils: • Loops may be hairpins or sharp turns • Random coils/irregular loops • Often "allowed" with insertions/deletions, *i.e.* evolutionary variable regions

Coil here: "Everything that is not helix or sheet"

Coil often means: "Everything that is not helix or sheet or some characteristic loops"

Often contains Gly (to give flexibility) or Pro (to "break up" secondary structure elements)

Left-handed helices

## Secondary structure – Gly & Pro



J. Richardson, Adv. Prot. Chem. 34, 167 (1981)

Glycine has no side chain and a more flexible backbone





Proline has very little flexibility in the backbone (disruptive to normal secondary structure)

#### Protein structure

Jon K. Lærdahl, Structural Bioinformatics

• Primary structure: Linear amino acid sequence

• Secondary structure: Local conformation of the peptide chain:

- α-helix
- β-sheet
- Tertiary structure: The full 3D structure
- Quaternary structure: Association of several proteins/peptide chains into protein complexes

Met-Ala-Leu-Asp-Asp-...

Hemoglobin, 1GZX



#### **Residue properties**



Proteins, T.E. Creighton (Freeman, New York, 1997)

#### **Residue properties**





Arg is "always" positively charged with pKa close to 12

His has pK<sub>a</sub> close to 7 and the local environment is often tuned to to give correct acid/base chemistry. Strong base at neutral pH/Strong nucleophile. Often a catalytic residue.

Table 1.2	Intrinsic pKa	Values	of Ionizable Groups	
Found in	Proteins			

Group	Observed $pK_a^a$			
α-Amino	6.8-8.0			
α-Carboxyl	3.5-4.3			
$\beta$ -Carboxyl (Asp)	3.9-4.0			
y-Carboxyl (Glu)	4.3-4.5			
$\delta$ -Guanido (Arg)	12.0			
ε-Amino (Lys)	10.4-11.1			
Imidazole (His)	6.0-7.0			
Thiol (Cys)	9.0-9.5			
Phenolic hydroxyl (Tyr)	10.0-10.3			

Proteins, T.E. Creighton (Freeman, New York, 1997)

#### Side chain conformations (Rotamers)



Analysis of many structures have shown that residues prefer one or a few conformations. These are called *rotamers* and are collected and distributed in *rotamer libraries* 

These libraries are used in computational modeling of protein 3D structure.

#### Very simply put:

- 1. Determine overall 3D structure of backbone
- 2. Add side chains
- Optimize side chains using conformations from rotamer libraries

## Stabilizing forces

Jon K. Lærdahl, Structural Bioinformatics



## Stabilizing forces

**IMPORTANT: Hydrophobic interaction forces** (minimizing the surface area of hydrophobic side chains exposed to solvent)





Introduction to Protein Structure, C. Branden & J. Tooze (Garland, New York, 1998)

the protein structure (e.g. in zinc fingers)

## Protein folding

What is making proteins fold and associate into a welldefined 3D structure?

• Proteins are often found in water and both protein-protein and protein-water interactions must be taken into account (*i.e.* interactions in folded vs. denatured state)

• *Dominant* forces responsible for tertiary structure are (believed to be) the hydrophobic interaction forces

- Residues with hydrophobic side chains are packed in the interior of the protein
- Charged and polar residues tend to be on the protein surface
- Polar backbone in the protein interior is "hidden" by building secondary structure elements

 Polar residue side chains in the core must be "neutralized" by interacting with other residues, e.g. in Hbond donor-acceptor pairs

 Charged residue side chains in the core must be "neutralized" by interacting with other residues through salt bridges Jon K. Lærdahl, Structural Bioinformatics



### Protein folding

Secondary structure elements ( $\alpha$ -helices &  $\beta$ -sheets) on the surfaces of proteins are often amphipathic (one hydrophilic and one hydrophobic side)



"Pattern" of every 3-4 residues hydrophobic

Patterns can be used for predictions by computational methods, *e.g.* predict secondary structure from primary sequence



http://cti.itc.virginia.edu/~cmg/Demo/wheel/wheelApp.html

Jon K. Lærdahl, Structural Bioinformatics

## Protein folding

#### Jon K. Lærdahl, Structural Bioinformatics

#### TLASTPALWASIPCPRSELRLDLV LPSGQS

Folding is spontaneous in the cell (but often with helper molecules, chaperones)



#### Put very simply:

- 1. Secondary structure forms transiently
- 2. Hydrophobic collapse, formation of stable secondary structure
- 3. Folding completes, formation of tertiary interactions

### Globular vs. membrane proteins

Jon K. Lærdahl, Structural Bioinformatics

Globular proteins • Soluble Surrounded by water Membrane proteins • In lipid bilayers • Hydrophobic surface facing membrane interior



#### Beta-barrel porin (1PRN)



#### Membrane proteins

Jon K. Lærdahl, Structural Bioinformatics

Rhodopsin (1QHJ)

Co-factor/prosthetic group retinal:



Covalent (Schiff bond) linkage to protein Lys residue

Many apo-proteins need cofactors/prosthetic groups to become functional

## PTMs

Jon K. Lærdahl, Structural Bioinformatics

Post-translational modifications (PTMs), *i.e.* chemical modification after translation, *e.g.* 

- Glycosylation (addition of sugar groups to *e.g.* Asn, Ser, or Thr)
- Phosphorylation of Ser/Thr by kinases
- Methylation of Lys in histones
- Ubiquitination (addition of the protein ubiquitin to Lys)
- Methionine aminopeptidases may remove N-terminal Met
- Many, many more!!

Bhaumik et al., Nat. Struct. Mol. Biol. 14, 1008 (2007)



PTMs of human histones include acetylation (ac), methylation (me), phosphorylation (ph) and ubiquitination (ub1)

Even if you know the complete 3D structure of the apo-protein you may be unable to understand the function of the protein if you have no information about the PTMs!

#### Visualization of protein structure

Jon K. Lærdahl, Structural Bioinformatics



Human OGG1, a DNA repair enzyme that recognizes and excises oxidized **DNA** bases





**Ribbons/Cartoon** 

Software (advanced graphics rendering):

- RasMol
- Swiss-PDBViewer (freeware; also homology) modeling)
- Molscript (command-line-based)
- Jmol (open-source Java viewer)
- PyMOL (open-source, user-sponsored)
- Many more both free and very expensive

We will use some of these at the Exercises!

#### Visualization of protein structure

Jon K. Lærdahl, Structural Bioinformatics





#### Visualization of protein structure

Jon K. Lærdahl, Structural Bioinformatics



Publication quality graphics from PyMOL

#### Movies, interactivity etc.



The structure of Bacillus stearothermophilus Fpg protein borohydride-trapped with DNA oligo as determined by Fromme and Verdine, Nat. Struct. Biol. **9**, 544 (2002), PDB: 1L1Z.

The graphics was generated with PyMOL

## Structural disorder in proteins

Jon K. Lærdahl, Structural Bioinformatics

Not all proteins have a regular 3D structure for the full sequence
The full protein, segments or small parts may be structurally disordered/intrinsically unstructured

Predicted 20% of human proteins have disordered segments of length >50 residues (1% in *E. coli*) (J.J. Ward *et al.*, *J. Mol. Biol.* **337**, 635 (2004))



Increasing content of stable three-dimensional structure

H.J. Dyson & P.E. Wright, Nat. Rev. Mol. Cell Biol. 6, 197 (2005)

#### Experimental determination of protein structure – X-ray Crystallography

- Necessary to grow protein crystals
  - Often (extremely) difficult
- Diffraction in X-ray beam
- Must solve "phase problem" (due to unknown timing of diffraction waves hitting the detector):
  - Molecular replacement (use the known structure of similar protein)
  - Multiple isomorphus replacement (generate crystals with heavy atoms, *e.g.* by soaking)
- Strong X-ray source needed to get high accuracy (Synchrotron)



Jon K. Lærdahl, Structural Bioinformatics

Li *et al.*, *Acta Cryst.* **D55**, 1023 (1999)

Proteins are located in a lattice, in a repeated and oriented fashion







#### **Experimental determination of protein structure – X-ray** Crystallography

Jon K. Lærdahl, Structural Bioinformatics

Diffraction pattern & solved phases: Electron density map ("electron cloud"):

• Model protein primary sequence into electron density map

- Resolution:
  - Low ~5.0 Å
  - Intermediate ~2.0-2.5 Å
  - High ~1.2 Å (Only at this very high, and rare, resolution it is possible to locate hydrogen atoms. H-atoms are therefore usually not visible in the structures.

- Gives a *static* picture of the protein in the crystal which might not correspond closely to situation in solution
- Bottleneck: Crystallization (and phase problem)
- No electron density for structurally disordered regions



(~3.5 Å resolution)

28, 677 (2000) (1.9 Å resolution)

## Experimental determination of protein structure – NMR Spectroscopy

Jon K. Lærdahl, Structural Bioinformatics

Nuclear Magnetic Resonance (NMR) Spectroscopy:

• Based in energy levels of magnetic nuclei (e.g. <sup>13</sup>C and <sup>15</sup>N) in a very strong external magnetic field probed my a radio frequency signal

• Determines distances between all labeled atoms in a protein

- Structure model built from distances
- Structure solved in solution
  - No need to grow crystals

Can be used to study proteins dynamics & behavior in solution

• Can currently only be employed for proteins of limited size (a few hundred residues)





#### **Experimental** determination of protein

structure

Jon K. Lærdahl, Structural Bioinformatics

#### X-ray Crystallography:

Pros:

- Can be used for huge protein complexes
  - 10.000s of atoms in *e.g.* complete ribosomes



B.S. Schuwirth, *Science* **310**, 827 (2005)

- Can in fortunate cases give very high resolution (Atom position uncertainty ~0.2 Å or less) *Cons*:
- Usually (extremely!) tricky to grow crystals
  - Membrane proteins are particularly difficult
  - Proteins with disordered segments are difficult
- Need to solve phase problem
- Does not give insight into dynamics and protein disorder
- Large amounts of protein needed
- Usually missing H-atoms
- Disordered loops/regions are not visible

#### NMR Spectroscopy:

Pros:

- Can be used directly on proteins in solution
- No need for crystallization
- Dynamics studies
- Both ordered and disordered proteins (usually an ensemble of 20-40 models)



#### Cons:

- Only applicable for small proteins (<200 residues?)
- Huge amounts of protein needed

All experimental methods: Labor intensive and requiring (very) expensive instruments Membrane proteins *extremely tricky The experimental structures are also models!* 

# Modeling of atoms into electron density

Jon K. Lærdahl, Structural Bioinformatics



X-ray crystallography



NMR

# Modeling of atoms into electron density

Jon K. Lærdahl, Structural Bioinformatics

1PRN



The experimental structures are also "models"!

And heavily depends on computers/software



Remember, when looking at an *experimental structure* (X-ray):

- Resolution and R-factor gives you an idea about the quality of the experimental model
  - Resolution ~ 3 Å: side chains may be wrong rotamer or missing, main chain normally ok
  - Resolution ~ 2 Å: most side chains should be ok
  - Resolution < 1.5 Å: high accuracy structure
  - Resolution < 1.2 Å: may even be possible to determine positions for hydrogen atoms

• Due to structural flexibility or "problems" in crystals, some regions, typically loops or N-/C-terminus may have little visible electron density.

- In some cases this gives gaps in the sequences or missing side chains
- In other cases people put in residues/atoms anyway, in reasonable positions
- The Uppsala Electron Density Server can be useful

#### **Protein Structure Database**

Protein Data Bank (PDB) <u>www.pdb.org</u>:

The home of all experimental proteins structures



Jon K. Lærdahl, Structural Bioinformatics

>113,000 structures
 Not all are unique

Some few 1000 unique protein folds

126,551,501,141 bases in 135,440,924 sequence records in the traditional GenBank divisions as of April 2011

PDB identifiers are on the form 1LYZ, 2B6C, 1T06 (and does not "mean" anything)

#### **Protein Structure Database**

Jon K. Lærdahl, Structural Bioinformatics

Search for	Showing 1 - 2	5 of 31 Results				Results :	25 ᅌ	Page: 1 of 2
"OGG1"	Filter: Check	All 📀 View:	Detailed	💿 Download Results 🛎		Reports: Select one	Sort:	Relevance
	Ø 1KO9 ≛ ₿ ♣	Native Struct Authors: Release:	Bjoras, M.,	e Human 8-oxoguanine DNA Glycosyla Seeberg, E., Luna, L., Pearl, L.H., Barrett, T.E. 9	se hOGG1			
		Experiment: Compound:	X-RAY DIFI 1 Polymer 1 Ligand (	FRACTION with resolution of 2.15 Å [ Display Full Polymer Details   Display for All Results ] Display Full Ligand Details   Display for All Results ]	Residue Count 345			
		Citation: Search Hit:	Reciproca (2002) J.M _entity_src	al "flipping" underlies substrate recognition a lol.Biol. <b>317:</b> 171-177 ( <i>Display Full Abstract</i>   <i>Display</i> c_gen.pdbx_gene_src_gene <b>OGG1</b>	nd catalytic activation by (for All Results ]	the human 8-oxo-guanine l	DNA giyo	osylase.
	© 1HU0 ≛ ⊯ ♣	CRYSTAL STR Authors:	Fromme, J.C	OF AN HOGG1-DNA BOROHYDRIDE TR C., Bruner, S.D., Yang, W., Karplus, M., Verdine, G.L.	APPED INTERMEDIA	TE COMPLEX		
	State .	Experiment: Compound:	X-RAY DIFI 3 Polymers 2 Ligands /	FRACTION with resolution of 2.35 Å 5 [ Display Full Polymer Details   Display for All Results ] [ Display Full Ligand Details   Display for All Results ]	Residue Count 354			
		Citation:	Product-A (2003) Nat	Assisted Catalysis in base-excision DNA Repai t.Struct.Biol. 10: 204-211 ( <i>Display Full Abstract</i>   <i>Dis</i> c gen pdby gene src gene QGG1	r splay for All Results }			
	<b>□ 1LWV</b> ★ ■ ■	Borohydride- Authors: Release:	Fromme, J.C	Dogg1 Intermediate Structure Co-Crys	stallized with 8-aming	oguanine		
	States	Experiment: Compound:	X-RAY DIFI 3 Polymers 2 Ligands (	- FRACTION with resolution of 2.30 Å 5 [ <i>Display Full Polymer Details</i>   <i>Display for All Results</i> ] [ <i>Display Full Ligand Details</i>   <i>Display for All Results</i> ]	Residue Count 354			
		Citation: Search Hit:	Product-A (2003) Nat	Assisted Catalysis in Base Excision DNA Repair t.Struct.Biol. 10: 204-211 ( <i>Display Full Abstract</i>   <i>Dis</i> c_gen.pdbx_gene_src_gene ogg1	r splay for All Results ]			
	Ø 1LWW ≟ 🖹 🖡	Borohydride- Authors: Release:	trapped h Fromme, J.C 2003-02-2	Cogg1 Intermediate Structure Co-Crys C., Bruner, S.D., Yang, W., Karplus, M., Verdine, G.L. S	stallized with 8-brom	oguanine		
	State Bar	Experiment: Compound:	X-RAY DIF 3 Polymers 2 Ligands (	FRACTION with resolution of 2.10 Å 5 ( <i>Display Full Polymer Details</i>   <i>Display for All Results</i> ) [ <i>Display Full Ligand Details</i>   <i>Display for All Results</i> ]	Residue Count 354			
		Citation:	Product-A (2003) Nat	Assisted Catalysis in Base Excision DNA Repai t.Struct.Biol. 10: 204-211 [ Display Full Abstract   Dis	r splay for All Results }			

#### **Protein Structure Database**

Jon K. Lærdahl, Structural Bioinformatics



## PDB entry – an example in PDB format

Standard since early 1970s
FORTRAN compatible format
Some limitations

Number of atoms
Number of chains
Length of fields

Not good for parsing by computers

HEADER	$\mathbf{L}$	YASE/DNA 24-JAN-00 1EBM
TITLE	C	RYSTAL STRUCTURE OF THE HUMAN 8-OXOGUANINE GLYCOSYLASE
TITLE	2	(HOGG1) BOUND TO A SUBSTRATE OLIGONUCLEOTIDE
COMPND	М	OL_ID: 1;
COMPND	2	MOLECULE: 8-OXOGUANINE DNA GLYCOSYLASE;
COMPND	3	CHAIN: A;
COMPND	4	FRAGMENT: CORE FRAGMENT (RESIDUES 12 TO 325);
COMPND	5	SYNONYM: AP LYASE;
COMPND	б	ENGINEERED: YES;
COMPND	7	MUTATION: YES;
COMPND	8	MOL_ID: 2;
COMPND	9	MOLECULE: DNA (5'-D(*GP*CP*GP*TP*CP*CP*AP*(OXO)
COMPND	10	GP*GP*TP*CP*TP*AP*CP*C)-3');
COMPND	11	CHAIN: C;
COMPND	12	ENGINEERED: YES;
COMPND	13	MOL_ID: 3;
COMPND	14	MOLECULE: DNA (5'-
COMPND	15	D(*GP*GP*TP*AP*GP*AP*CP*CP*TP*GP*GP*AP*CP*GP*C)-3');
COMPND	16	CHAIN: D;
COMPND	17	ENGINEERED: YES
SOURCE	М	OL_ID: 1;
SOURCE	2	ORGANISM_SCIENTIFIC: HOMO SAPIENS;
SOURCE	3	EXPRESSION_SYSTEM: ESCHERICHIA COLI;
SOURCE	4	EXPRESSION_SYSTEM_COMMON: BACTERIA;
SOURCE	5	EXPRESSION_SYSTEM_VECTOR_TYPE: PLASMID;
SOURCE	6	EXPRESSION_SYSTEM_PLASMID: PET30A-HOGG1;
SOURCE	7	MOL_ID: 2;
SOURCE	8	SYNTHETIC: YES;
SOURCE	9	MOL_ID: 3;
SOURCE	10	SYNTHETIC: YES
KEYWDS	D	NA REPAIR, DNA GLYCOSYLASE, PROTEIN/DNA
EXPDTA	Х	-RAY DIFFRACTION
AUTHOR	S	.D.BRUNER, D.P.NORMAN, G.L.VERDINE
REVDAT	1	20-MAR-00 1EBM 0
JRNL		AUTH S.D.BRUNER, D.P.NORMAN, G.L.VERDINE
JRNL		TITL STRUCTURAL BASIS FOR RECOGNITION AND REPAIR OF THE
JRNL		TITL 2 ENDOGENOUS MUTAGEN 8-OXOGUANINE IN DNA
JRNL		REF NATURE V. 403 859 2000
JRNL		REFN ASTM NATUAS UK ISSN 0028-0836
REMARK	1	
REMARK	2	RESOLUTION. 2.10 ANGSTROMS.
REMARK	3	



```
PDB entry –
an example
in mmCIF
format
```

Newer data format and alternative to "PDB format"

No limitations in number of atoms, chains, fields etc.
Better suited for automatic parsing/processing

```
data 1EBM
entry.id
            1EBM
audit conform.dict name
                               mmcif pdbx.dic
audit conform.dict version
                               1.044
audit conform.dict location
                               http://mmcif.pdb.org/dictionaries/ascii/mmcif pdbx.
database 2.database code
PDB
   1EBM
NDB
    PD0117
RCSB RCSB010437
#
database PDB rev.num
                                  1
database PDB rev.date
                                  2000-03-20
database PDB rev.date original
                                  2000-01-24
database PDB rev.status
                                  ?
database PDB rev.replaces
                                  1EBM
                                  0
database PDB rev.mod type
#
pdbx database status.status code
                                     REL
pdbx database status.entry id
                                     1EBM
pdbx database status.deposit site
                                     RCSB
pdbx database status.process site
                                     RCSB
pdbx database status.SG entry
                                      .
#
loop
audit author.name
'Bruner, S.D.'
'Norman, D.P.'
'Verdine, G.L.'
#
citation.id
                                    primary
                                    'Structural basis for recognition
citation.title
citation.journal abbrev
                                    Nature
_citation.journal_volume
                                    403
citation.page first
                                    859
citation.page last
                                    866
```

## Structural bioinformatics

Experimental structure is hard to get

The 3D structure on a protein is determined by the amino acid sequence (primary structure)

There are many orders of magnitude more sequences available than there are structures



How do we get information about structure from sequence?



Jon K. Lærdahl, Structural Bioinformatics

Jon K. Lærdahl, Structural Bioinformatics

Domain: Compact part of a protein that represents a structurally independent region

Domains are often separate functional units that may be studied separately

Domains fold independently? Not always...





Dividing a protein structure into domains: no "right way to do it" or "correct algorithm", *i.e.* **a lot of subjectivity involved** 



Most people would agree there are two domains here

Three domains? One domain? Two?

SCOP vs. CATH?

Very often we model, compare, classify *domains* – not full-length proteins

Jon K. Lærdahl, Structural Bioinformatics



Instead of working with full length proteins that may be

- very large
- contain one or many separate modules (*i.e.* domains)
- have both structured and unstructured parts

We often instead work with protein domains that are

- more compact
- can be studied separately
  - function
  - structure by X-ray crystallography/NMR
  - bioinformatics modeling
- may be viewed as the "spare parts" building up full-length proteins

Many proteins are structured domains, "spare parts", connected by short loops or long disordered regions

Far from trivial to detect boundaries between domains from sequence only:





Jon K. Lærdahl, Structural Bioinformatics

Domains have a "signature sequence" that can be described as a HMM Logo Domains can be "switched". They can be viewed as "spare parts" that can be used to build new proteins through evolution



Pfam HMM-logo for the GRF zinc finger domain

#### Protein domains Nature of the protein universe

Jon K. Lærdahl, Structural Bioinformatics

#### PNAS 106, 11079 (2009)

#### Michael Levitt<sup>1</sup>

Department of Structural Biology, Stanford University, Stanford, CA 94305-5126

Contributed by Michael Levitt, May 9, 2009 (sent for review April 20, 2009)

The protein universe is the set of all proteins of all organisms. Here, all currently known sequences are analyzed in terms of families that have single-domain or multidomain architectures and whether they have a known three-dimensional structure. Growth of new single-domain families is very slow: Almost all growth comes from new multidomain architectures that are combinations of domains characterized by  $\approx$ 15,000 sequence profiles. Single-domain families are mostly shared by the major groups of organisms, whereas multidomain architectures are specific and account for species diversity. There are known structures for a quarter of the single-domain families, and >70% of all sequences can be partially modeled thanks to their membership in these families.

An obvious way to cluster sequences into families is by pairwise comparison (4) of all sequences preceded by indexing (5) to eliminate close pairs. Such a combination led to massive clustering of millions of protein sequences from both known species and environmental samples by Yooseph et al. (6). Their remarkable conclusion was that the number of protein families as measured by the number of sequence clusters showed no sign of saturation. Indeed, the cluster count was increasing at the same rate as new sequences were being determined. This result

www.pnas.org/cgi/doi/10.1073/pnas.0905029106

featured in a recent report on the Protein Structure Initiative (7) that expressed concern that because the number of new families is expanding rapidly determining three-dimensional structures for a representative of each family may not be possible (8).

Here, we approach the problem differently. Instead of clustering entire protein sequences (6), we rely on the occurrence of protein sequence patterns termed "sequence profiles." These patterns can be derived from a few members of the family and then used to add new members that match the same pattern.

> (6) Yooseph D, *et al.* (2007) The Sorcerer II global ocean sampling expedition: Expanding the universe of protein families. PLoS Biol **5**:e16.





Jon K. Lærdahl, Structural Bioinformatics

PNAS 106, 11079 (2009)



Fig. 1. As the NR database grows, the number of different multidomain architecture (MDA) families found by CDART is increasing rapidly with year (*Left*) or added sequence (*Right*). In contrast, the number of single-domain architecture (SDA) families is increasing much more slowly. Because the number of sequences is growing exponentially, fractional sequence coverage (number of sequences in a SDA or MDA family divided by the total number of NR sequences) has dropped slightly from 0.88 to 0.76; more than three-quarters of current sequences still contain a domain recognized by a known sequence profile. Merged CDART sequence profiles are used here. Corresponding results with unmerged CDART sequence profiles are given in Fig. S1. The solid curves marked "2008" were made with a release of CDART from February 9, 2008, which contained fewer sequence profiles (24,083 compared with 27,036). This gave rise to smaller numbers of SDA and MDA families and lower coverage. During this time, the number of sequences in the NR database increased by 2 million.

There are known structures for a quarter of the single-domain families, and >70% of all sequences can be partially modeled thanks to their membership in these families.

## End