Bioinformatics for molecular biology

Structural bioinformatics tools, predictors, and 3D modeling – Structural Bioinformatics

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





Overview

Earlier...

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

Now

- 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!!



Structural bioinformatics

Jon K. Lærdahl, Structural Bioinformatics





To understand what is really going on in biology you need the 3D structure of the macromolecules, *i.e.* the proteins in particular!

Neuraminidase is a ^{CG} glycoside hydrolase enzyme found on the surface of the influenza virus



.OH



Structural bioinformatics

Experimental methods for determining protein 3D structure are very expensive in terms of money and time

Alternative: Use computational methods to determine protein structure

- Determine 3D structure with computers
- Determine secondary structure, structural disorder, domain boundaries, sites of post-translational modifications (PTMs), etc.
- Understand structure through computations
- Work with 3D structures, compare, classify, etc.
- Goal is to get biological insight!



Protein domains

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...





Protein domains

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

Comparing structures

Jon K. Lærdahl, Structural Bioinformatics

Superimposition/Alignment of 3D structures in space

The structures are superimposed in order to get corresponding main chain atoms as closely together as possible

If identical sequences – align all atoms

Non-identical sequences – align back-bone atoms only (usually *only* aligning Cα atoms!)

Structure is more conserved than sequence. A structural alignment can therefore be used to define the "correct" sequence alignment



In many cases we align domains, not full length proteins

Above, the two domains of NTH (blue) aligns nicely with the two C-terminal domains of OGG1 (green). The remaining domain of OGG1 is missing in NTH

Comparing structures

Jon K. Lærdahl, Structural Bioinformatics

Comparison of protein structures

Human NEIL1



E. coli endonuclease VIII

Brief demo in PyMOL!

>align molecule1, molecule2, object=matches

Aligned with RMSD = 1.41 Å



Root mean square deviation (RMSD) = square root of averaged sum of the squared differences of atomic (usually $C\alpha$) distances

$$RMSD = \sqrt{\frac{1}{N}\sum_{i=1}^{i=N}\delta_i^2}$$

Comparing structures

Root mean square deviation (RMSD) = square root of averaged sum of the squared differences of atomic (usually $C\alpha$) distances

Calculate RMSD by: Loop over equivalent positions *i* $_{P}$ Get coordinates for both Cas Calculate distance beetween Cas, δ_i Square δ_i and add to sum End loop Divide sum by number of pairs, *N*, and take square root

$$RMSD = \sqrt{\frac{1}{N}\sum_{i=1}^{i=N}\delta_i^2}$$

RMSD tells you how similar two structures are

RMSD of ~0.5 Å or less for "identical" structures



Comparing structures - Intermolecular method

Jon K. Lærdahl, Structural Bioinformatics



ARITGMVERLCQAFGPRLIQLDDVTYHGFPSLQALAGPEVEAHLR----PKVRRYQVLLSLMLSSQ--TKDQVTAGAMQRLRA-RGLTVDSILQTDD/

LEEQ<mark>GG</mark>LAWLQQLRESSYEEAHKALCILPGVGTKVADCICLMALDKPQ/ KQTS---AILQQHYGGDIPASVAELVALPGVGPKMAHLAMAVAWGTVS(



Comparing structures - Intermolecular method

Problems with Intermolecular method:

- RMSD depends on protein size
- Tricky to identify "equivalent residues" in the beginning
- Usually means that a sequence alignment is done first
 - Aligned residues are considered "equivalent"
 - Means the method is only useful for sequences that can be aligned by sequence comparison
- Several solutions suggested, but may give strange and non-optimal solutions

Important to check alignments visually!

- Iterative optimization:
 - First detect (often small) segments that can be aligned based on sequence

• Do 3D superimposition based on residues in these segments

• Based on 3D alignment, identify more residues that are close together and that are at "equivalent positions". Use this larger set of pairs to do a new 3D superimposition.

Repeat until RMSD is converged



Comparing structures - Intramolecular method

- May be used for any two or more structures
- Does not depend on sequence similarity
- Does not necessarily generate physical superimposition
- Instead structural similarity measure based on internal structural statistic for each protein chain
- Based on building and comparing distance matrices for the structures
 - For example matrix A of all C distances in protein A and matrix B for protein B
 - "Align" matrices to get best overlap
- Used in the most popular structure comparison tools, for example DALI
- Used for example to find which protein in the PDB is most similar to a new structure
- Intermolecular method:
 - Similar structures
 - Gives physical superimposition
- Intramolecular method:
 - Can be used for any two or more structures





Jon K. Lærdahl, Structural Bioinformatics

Comparing structures – Some tools:

Jon K. Lærdahl, Structural Bioinformatics

STAMP (http://www.compbio.dundee.ac.uk/downloads/stamp): Unix program for iterative intermolecular alignment

Similar algorithms are often included in Viewers (*e.g.* DeepView & PyMOL)

Dali server				Institute of Biotechnology	Intramole				
SERVICES & TOOLS	GROUP MEMBERS	NEWS & VACANCIES	RESEARCH	PUBLICATIONS	method				
					method				
Protein Stru	cture Databas	e Searching b	y DaliLite v. 3						
The Dali server is a netw compares them against cases, comparing 3D str	vork service for comparing pro those in the Protein Data Bar uctures may reveal biologically	otein structures in 3D. You s nk (PDB). You receive an em y interesting similarities that	ubmit the coordinates of a qu ail notification when the sear are not detectable by compar	uery protein structure and Dali rch has finished. In favourable ring sequences.	(Finland)				
Requests can also be s in PDB format.	ubmitted by e-mail to dali-sen	ver at helsinki dot fi. The boo	y of the e-mail message mu	st contain atomic coordinates					
If you want to know the st	If you want to know the structural neighbours of a protein already in the Protein Data Bank (PDB), you can find them in the Dali Database.								
If you want to superimpos									
Upload a structu	re:(Browse							
Or enter PDB ide (Keyword search for PDB	ntifier: chain:	(optional)							
Job name:									
		(optional)							
Enter email addr	ess for notification:	(
		(recommended)							

Results

Comparing structures – Some tools:

- Parseable data
- Matches to PDB90

The match list is truncated at 500 hits.

Query: 1ebmA

MOLECULE: 8-OXOGUANINE DNA GLYCOSYLASE;

Select neighbours (check boxes) for viewing as multiple structural alignment or 3D superimposition. The list of neighbours is sorted by Z-score. Similarities with a Z-score lower than 2 are spurious. Each neighbour has links to pairwise structural alignment with the query structure, to pre-computed structural neighbours in the Dali Database, and to the PDB format coordinate file where the neighbour is superimposed onto the query structure.



Summary

No:	Chain	Z	rmsd	lali	nres	%id PDB	Description
<u>1</u> :	<u>lebm-A</u>	99.9	0.0	0	314	0 PDB	MOLECULE: 8-OXOGUANINE DNA GLYCOSYLASE;
<u>2</u> :	<u>1n3a-A</u>	51.9	0.0	0	314	0 PDB	MOLECULE: N-GLYCOSYLASE/DNA LYASE;
<u>3</u> :	<u>1m3q-A</u>	51.8	0.0	0	314	0 <u>PDB</u>	MOLECULE: 8-OXOGUANINE DNA GLYCOSYLASE;
<u>4</u> :	<u>11ww-A</u>	51.8	0.0	0	314	0 <u>PDB</u>	MOLECULE: 8-OXOGUANINE DNA GLYCOSYLASE;

<u>13</u> :	<u>2jhj-B</u> 24.3	0.0	0	291	0 PDB	MOLECULE:	3-METHYLADENINE DNA-	GLYCOSYLASE;	
<u>14</u> :	<u>2jhn-A</u> 23.9	0.0	0	293	0 <u>PDB</u>	MOLECULE:	3-METHYLADENINE DNA-	GLYCOSYLASE;	
<u>15</u> :	<u>3d4v-D</u> 22.9	0.0	0	281	0 <u>PDB</u>	MOLECULE:	DNA-3-METHYLADENINE	GLYCOSYLASE	2;
<u>16</u> :	<u>3d4v-B</u> 22.8	0.0	0	282	0 <u>PDB</u>	MOLECULE:	DNA-3-METHYLADENINE	GLYCOSYLASE	2;
<u>17</u> :	<u>3cwu-D</u> 22.8	0.0	0	282	0 <u>PDB</u>	MOLECULE:	DNA-3-METHYLADENINE	GLYCOSYLASE	2;
18:	<u>3cwt-D</u> 22.8	0.0	0	281	0 <u>PDB</u>	MOLECULE:	DNA-3-METHYLADENINE	GLYCOSYLASE	2;
19:	3cvs-D 22.7	0.0	0	282	0 PDB	MOLECULE:	DNA-3-METHYLADENINE	GLYCOSYLASE	2;

<u>38</u> :	<u>1pu8-A</u> 14.7	0.0	0	215	0 PDB	MOLECULE:	3-METHYLADENINE DNA GLYCOSYLASE;
<u>39</u> :	<u>2h56-B</u> 14.1	0.0	0	217	0 <u>PDB</u>	MOLECULE:	DNA-3-METHYLADENINE GLYCOSIDASE;
<u>40</u> :	2abk 14.0	0.0	0	211	0 <u>PDB</u>	MOLECULE:	ENDONUCLEASE III;
<u>41</u> :	<u>1rrt-A</u> 14.0	0.0	0	346	0 PDB	MOLECULE:	MUTY;
<u>42</u> :	<u>1vrl-A</u> 14.0	0.0	0	346	0 <u>PDB</u>	MOLECULE:	5'-D(*AP*AP*GP*AP*CP*(80G)P*TP*GP*GP*AP
<u>43</u> :	<u>1kea-A</u> 13.5	0.0	0	217	0 PDB	MOLECULE:	POSSIBLE G-T MISMATCHES REPAIR ENZYME;
44:	<u>1kg6-A</u> 13.0	0.0	0	224	0 <u>PDB</u>	MOLECULE:	A/G-SPECIFIC ADENINE GLYCOSYLASE;
<u>45</u> :	<u>1kg5-A</u> 12.9	0.0	0	225	0 PDB	MOLECULE:	A/G-SPECIFIC ADENINE GLYCOSYLASE;
46:	<u>1kg3-A</u> 12.8	0.0	0	224	0 PDB	MOLECULE:	A/G-SPECIFIC ADENINE GLYCOSYLASE;
<u>47</u> :	<u>1kqj-A</u> 12.7	0.0	0	225	0 PDB	MOLECULE:	A/G-SPECIFIC ADENINE GLYCOSYLASE;
48:	1muy-A 12.6	0.0	0	225	0 PDB	MOLECULE:	ADENINE GLYCOSYLASE;

(http://ekhidna.biocenter.helsinki.fi/dali_server)

Dali

- Compare 2 structures
- Compare multiple structures
- Search a database of structures for the most similar structures with a pdb-file query
- Search database with PDB id query
- Z-score > 4 usually indicates significant level of similarity

CEAlign

demo

Alternatives:

CE

SSAP

More...

VAST+ (at NCBI)

Protein structure evolution

The origin of this gene/protein is (very likely) before the last common

ancestor of S. cerevisiae (yeast), human, mouse, rat, and fruit fly

- Some of the amino acids have not mutated in >1 billion years
- Neutral mutation rate in mammals is ~0.01 base pair/5 million yr



Jon K. Lærdahl.

Structural Bioinformatics

OGG1_YEAS7/1-376	1 - MSYKFGKLAINKSELCLANVLQAGQSFRVIWDEKLNQYSTTMKIGQQEKYSVVILRQDEENEILEFVAVGDCGNQ75
OGG1_MOUSE/1-345	1 - MLFRSWLPSSMRHRTLSSSPALWASIPCPRSELRLDLVLASGOSFRWKEQSPAHWSGVLADQVWTLTQTEDQLYCTVYRGDDSQVSRPTLEEL93
OGG1 RAT/1-345	1 - MLFSSSLSSSMRHRTLTSSPALWASIPCPRSELRLDLVLASGQSFRMREQSPAHMSGVLADQWVTLTQTEDQLYCTVYRGDKGQVGRPTLEEL93
OGG1 HUMAN/1-345	1 - MPARALL PRRMGHRTLASTPALWASIPCPRSELRLDLVLPSGQSFRMR EQSPAHMSGVLADQWVTLTQTEEQLHCTVYRGDKSQASRPTPDEL 93
0GG1 FLY/1-343	1 MLAHNLGFHKKRLFSNMKAVLQDRGVIGLSLEECDLERTLLGGQSFRARSICDGNRTKYGGVVFNTYWVLQDEESFITYEAY-GTSSPLATKDYSSL96
OGG1 YEAST/1-376	76 DALKTHLMKYFRLDVSLKHLFDNW I PSDKAFAKLSP OG I RILAQEPVETLISFICSSNNN ISRITRMCNSLCSNFGNLITTIDGVAYHSFPTS EELT 173
OGG1 MOUSE/1-345	94 ETLHNYFOLDVSLAOLYSH-WASVDSHFORVAOKFOGVRLLRODPTECLFSFICSSNNNIARITGMVERLCOAFGPRLIOLDDVTYHGFPNL HALA 188
OGG1 RAT/1-345	94 ETLHKYFOLDVSLTOLYSH-WASVDSHFOSVAOKFOGVRLLRODPTECLESFICSSNNNIARITGMVERLCOAFGPRLVOLDDVTYHGFPNLHALA 188
OGG1 HUMAN/1-345	94 EAVRY FOLDYTLAOLYTH-WGSYDSHFOEVAOKFOGYRLLRODP LECLESFLCSSNNN LAR LTGMVERLCOAFGPRL LOLDDYTYHGFPSLOALA 188
0GG1 FLY/1-343	97 ISDYLRVDEDLKVNOKD-WISKDDNEVKELS KEVRLLSOEPEENLESELCSONNNIKRLSSMIEWECATEGTKIGHENGADAYTEPTINREHDLP 190
0001_12#1010	
OGG1 YEAST/1-376	174 SRATEAKLEELGEGYRAKYILLETARKI VNDKAEAN I TSDITTYLOSI CKDAOYEDVREHLMSYNGVGPKVADGVGLMGLHMDG I VPVDVHVSR LAKRDYG I SAN 276
0GG1_MOUSE/1-345	189 GPEAETHLEKLGLGYRARYVRASAKALLEFOGGP
0GG1 RAT/1-345	189 GREVETHERKI GLOVEARAWYCASAKALLEEOGGP
0661 HUMAN/1-345	189 CREVEAU BKLCLCKAR AND CLCKAR
0661 ELV/1 343	
0661_FEI/1=345	
0.0.01 YEAST/1 376	
OGGI_TEAST/1-376	
OGG1_MOUSE/1-345	285 - ISOAKGPS PLANKELG NFFRNL WGPYAGWAQAVLFSADLROPS - LSREPPAKRK
OGG1_RAT/1-345	285 - TSOTKOPS PLANKELG NFFRNL WGPYAGWAQAVLFSADLROON - LSR EPPAKRK
OGG1_HUMAN/1-345	285 - TSQAKGPS PQTNKELG NFFRSL WGPYAGWAQAVLFSADLRQSR - HAQEPPAKRR KGSKGPEG 345
OGG1_FLY/1-343	286 GQKNVTKKIYEEVSKHFQKLHGKYAGAAAILFSADLSQFQ-NTSTVACKKK

Protein structure evolution

96.50 OGG1_MOUSE 86.00 96.50 OGG1_RAT **Common ancestor** 182.50 - OGG1 HUMAN 146.71 OGG1_FLY **Ancestor B** OGG1_YEAST **Ancestor C** The overall structure **Ancestor D** of this protein is the same in all these organisms – i.e. many mutations does not change the structure and/or

Jon K. Lærdahl, Structural Bioinformatics

function

Protein structure evolution ^s

Jon K. Lærdahl, Structural Bioinformatics

Proteins that fold in the same way, i.e. "have the same fold" are often homologs. Structure evolves slower than sequence Sequence is less conserved than structure



EndoIII (2ABK):

OGG1 (1EBM):



OGG1/1-345	1 M <mark>P</mark> ARALL <mark>PRRMGH</mark> RTLASTPALWASIPCPRSELRLDLVLPSGQSFRAREQSPAHASGVLADQVATLTQ6	58
NTH1/1-312	1 MCSPQESGMTALSARMLTRSRSLGPGAGPRGCREEPGPLRREAAAE	47
OGG1/1-345	69 TEEQLHCTVYRGDKSQASRPTPDELEAVRKYFQLDVTLAQL <mark>YHHMG</mark> SVDSHFQEVAQKFQGVRLLRQD 1	36
NTH1/1-312	48 - ARKSHSPVKRPRKAQRLRVAYEGSDSEKGEGAEPLKVPVWEPQDWQQQLVNIRAMRNKKDA 1	108
OGG1/1-345	137 PIECLF <mark>S</mark> FICSSNNNIA <mark>R</mark> ITGMVERLCQAFGPRLIQLDDVTYHGFPSLQALAGPEVEAHLRKLG 2	200
NTH1/1-312	109 PVDHLGTEHCYDSSAPPKVRRYQVLLSLMLSSQTKDQVTAGAMQRLRA-RGLTVDSILQTDDATLG 1	173
OGG1/1-345	201 - L <mark>GYRARY</mark> VSASA <mark>R</mark> A ILEEQ <mark>GG</mark> LAWLQQLRESSYEEAHKALCILPGVGTKVADCICLMALDKPQAVPV	267
NTH1/1-312	174 KLIYPVGFWRSKVKYIKQTSAILQQHYGGDIPASVAELVALPGVGPKMAHLAMAVAWGTVSGIAV 2	238
OGG1/1-345	268 DVHMAH I AQ <mark>R</mark> DYSWHPTTSQAKGPSPQTNKELGNFF-RSLWGPYAGAAQAVLFSADLRQSRHAQ3	330
NTH1/1-312	239 DTHVHRIANR-LRWTKKATKSPEETRAALEEWLPRELWHEINGLLVGFGQQTCLPVHPRCHACLN3	302
OGG1/1-345	331 EPPAKRRKGSKGPEG 33	345
NTH1/1-312	303 GALCPAAQGL 33	312

Protein structure evolution

Pyrobaculum aerophilum AGOG

G.M. Lingaraju *et al. Structure* **13**, 87 (2005)

Hardly any detectable sequence similarity to human OGG1, and *E. coli* EndollI and MutY

Still clearly the same protein fold (overall structure)

Evolution has "eroded away" sequence similarity but left the structure intact



Protein structure evolution ^s

Jon K. Lærdahl, Structural Bioinformatics



Evolution has "eroded away" sequence similarity but left the structure intact

Protein structure alignments

Proteins that fold in the same way, i.e. "have the same

fold" are often homologs.

Structure evolves slower than sequence

Sequence is less conserved than structure

If BLAST gives no homologs (*i.e.* sequence based)

Instead: Search with protein *structure* (pdb-file) in *structure database* (e.g. PDB) to find more remote homologs

- For example using DALI
- Much more sensitive than sequence search
- Problems
 - Much smaller database (PDB vs. Genbank)
 - Need 3D structure of protein

Use structure comparisons to classify, group and cluster proteins. Build protein structure families and hierarchies





Jon K. Lærdahl, Structural Bioinformatics

Protein structure classification

- Based on taking all structures of PDB
- Remove redundancy (*i.e.* keep only one copy of "identical" structures)
- Split structures into domains
- Group domains/proteins based on similarity
- Two main classification schemes: SCOP & CATH



Structural Classification of Proteins

Scop Classification Statistics

SCOP: Structural Classification of Proteins. 1.73 release 34494 PDB Entries (26 Sep 2007). 97178 Domains. 1 Literature Reference (excluding nucleic acids and theoretical models)

Almost 100% manually generated
Proteins grouped into hierarchy of classes, folds, superfamilies and families

Class	Number of folds	Number of superfamilies	Number of families
All alpha proteins	259	459	772
All beta proteins	165	331	679
Alpha and beta proteins (a/b)	141	232	736
Alpha and beta proteins (a+b)	334	488	897
Multi-domain proteins	53	53	74
Membrane and cell surface proteins	50	92	104
Small proteins	85	122	202
Total	1086	1777	3464



SCOP

- Families
 - Sequence identity ~30% or higher
 - Very similar structures
 - Clearly homologous proteins
- Superfamilies
 - Contains families
 - May have no or little sequence similarity
 - Common fold
 - Are probably evolutionary related
- Folds
 - Contains superfamilies
 - Difficult level of classification
 - Same major secondary structure elements (α -helices and β -sheets) with same connections
 - Not always homologs

- Classes
 - Upper level of classification (4 major, 3 minor)
 - Contains folds
 - Based on secondary structure composition and "general features"
 - *e.g.* all- α , all- β , "membrane and cell surface" and "small proteins"
 - α/β : One β -sheet with strands connected by single α -helices
 - α + β : α -helical and β -sheet part separated in sequence

Class	Number of folds	Number of superfamilies	Number of families
All alpha proteins	259	459	772
All beta proteins	165	331	679
Alpha and beta proteins (a/b)	141	232	736
Alpha and beta proteins (a+b)	334	488	897
Multi-domain proteins	53	53	74
Membrane and cell surface proteins	50	92	104
Small proteins	85	122	202
Total	1086	1777	3464