DBAli tools

The AnnoLite and AnnoLyze programs for automatic annotation of protein structures

Structural Genomics Unit Bioinformatics Department Prince Felipe Resarch Center (CIPF), Valencia, Spain

Marc A. Marti-Renom

http://bioinfo.cipf.es/sgu/



DBAliv2.0 database

http://bioinfo.cipf.es/squ/DBAli/ http://www.salilab.org/DBAli/



Uses MAMMOTH for similarity detection

- ✓ VERY FAST!!!
- ✓ Good scoring system with significance

Ortiz AR, (2002) Protein Sci. 11 pp2606

- ✓ Fully-automatic
- ✓ Data is kept up-to-date with PDB releases
- ✓ Tools for "on the fly" classification of families.
- ✓ Easy to navigate
- Provides tools for structure analysis

Does not provide a stable classification similar to that of CATH or SCOP

Pairwise structure alignments
Last update:
Number of chains:
Number of structure-structure comparisons:*
Multiple structure alignments
Last update:
Number of representative chains:
Number of families:

Marti-Renom et al. 2001. Bioinformatics. 17, 746 Marti-Renom et al. 2006 Submitted.

DBAliv2.0 database

http://bioinfo.cipf.es/squ/DBAli/ http://www.salilab.org/DBAli/



DBAliv2.0 database

http://bioinfo.cipf.es/squ/DBAli/ http://www.salilab.org/DBAli/



For many protein structures function is *unknown*

	Structural Genomics*	Traditional methods
Annotated**	654	28,342
Not Annotated	506 (43.6%)	6,815 <mark>(19,4%)</mark>
Total deposited	1,160	35,157

* annotated as STRUCTURAL GENOMICS in the header of the PDB file **annotated with either CATH, SCOP, Pfam or GO terms in the MSD database 36,317 protein structures, as of August 8th, 2006

For 20% protein structures function is *unknown*

	Structural Genomics*	Traditional methods
Annotated**	654	28,342
Not Annotated	506 (43.6%)	6,815 <mark>(19,4%)</mark>
Total deposited	1,160	35,157

* annotated as STRUCTURAL GENOMICS in the header of the PDB file **annotated with either CATH, SCOP, Pfam or GO terms in the MSD database 36,317 protein structures, as of August 8th, 2006

AnnoLite results for chain <u>1 gp</u>: A based on <u>44</u> structural similar chains.

		7.5e-99		1,4-Beta-D-Glucan Cellobiohydrolase I, subunit A
SCOP:		0.00	<u>b.29.1.10</u>	Glycosyl hydrolase family 7 catalytic core
PFAM:	•	0.00	PF00840	Glycosyl hydrolase family 7
InterPro:	•	1.3e-99	IPR001722	Glycoside hydrolase, family 7
	•	6.0e-51	IPR008985	Concanavalin A-like lectin/glucanase
	۰	1.0e-42	IPR000254	Cellulose-binding region, fungal
EC Number:	•	1.2e-44	<u>3.2.1.91</u>	Cellulose 1,4-beta-celloblosidase.
	•	6.0e-41	3.2.1.4	Cellulase.
GO Molecular Function:	•	6.0e-36	0030248	cellulose binding 🟅
	•	8.4e-36	0016162	cellulose 1,4-beta-cellobiosidase activity 🟅
	•	1.0e-35	0004553	hydrolase activity, hydrolyzing O-glycosyl compounds 🕻
	•	1.4e-30	0008810	cellulase activity 🟅
	•	3.1e-20	0016798	hydrolase activity, acting on glycosyl bonds $\$
	٠	1.0e+0	0016787	hydrolase activity 🟅
GO Biological Process:	•	1.1e-63	0030245	cellulose catabolism 🗧
	•	1.2e-54	0000272	polysaccharide catabolism 🗧
	•	3.6e-20	0005975	carbohydrate metabolism 🗧
GO Cellular Component:	•	1.2e-23	0005576	extracellular region 🟅

Information annotated in the MSD database.

. High, . medium and . low confidence annotations not annotated in the MSD database.

High,
medium and
Medium and
Medium confidence annotations already annotated in the MSD database.

Benchmark set

	Number of chains
Initial set*	50,223
FULL annotation**	10,997
Non-redundant set***	1,879

*data from BioMart MSD.3 (release February 2005) **annotated with CATH, SCOP, Pfam, EC, InterPro, and GO terms in the MSD database **not two chains can be structurally aligned within 2A, superimposing more than 60% of their C atoms and have a length difference inferior to 30aa

Method



AnnoLite results for chain 1gpi:A based on 44 structural similar chains.

	Con	f. P-value	Link	Description
CATH:	•	7.5e-99	2.70.100.10	1,4-Beta-D-Glucan Cellobiohydrolase I, subunit A
SCOP:	•	0.00	<u>b.29.1.10</u>	Glycosyl hydrolase family 7 catalytic core
PFAM:	•	0.00	PF00840	Glycosyl hydrolase family 7
InterPro:	•	1.3e-99	IPR001722	Glycoside hydrolase, family 7
	•	6.0e-51	IPR008985	Concanavalin A-like lectin/glucanase
	۰	1.0e-42	IPR000254	Cellulose-binding region, fungal
EC Number:	٠	1.2e-44	<u>3.2.1.91</u>	Cellulose 1,4-beta-cellobiosidase.
	•	6.0e-41	3.2.1.4	Cellulase.
GO Molecular Function:	•	6.0e-36	0030248	cellulose binding 🟅
	•	8.4e-36	0016162	cellulose 1,4-beta-cellobiosidase activity 🗧
	•	1.0e-35	0004553	hydrolase activity, hydrolyzing O-glycosyl compounds ζ
	•	1.4e-30	0008810	cellulase activity 🗧
	٠	3.1e-20	0016798	hydrolase activity, acting on glycosyl bonds $\ \zeta$
	٠	1.0e+0	<u>0016787</u>	hydrolase activity 💪
GO Biological Process:	•	1.1e-63	0030245	cellulose catabolism 🟅
	•	1.2e-54	0000272	polysaccharide catabolism 🟅
	•	3.6e-20	0005975	carbohydrate metabolism 🐛
GO Cellular Component:	•	1.2e-23	0005576	extracellular region 🟅

Information annotated in the MSD database.

. High, . medium and . low confidence annotations not annotated in the MSD database.

High, O medium and O low confidence annotations already annotated in the MSD database.

Scoring function

Fisher's 2x2 contingency test

	Non- similar	Similar	Total
Annotated	а	b	a+b
Not Annotated	с	d	c+d
Total	a+c	b+d	n

1b78A SCOP c.51.4.1	Similar	Not similar	Total
Annotated	4	2	6
Not Annotated	0	71,096	71,096
Total	4	71,098	71,102

$$p = \binom{a+b}{a} \binom{c+d}{c} / \binom{n}{a+c}$$
$$= \frac{(a+b)!(c+d)!(a+c)!(b+d)!}{n!a!b!c!d!}$$

$$p = 1.78e^{-19}$$

Sensitivity .vs. Precision

	Optimal cut-off	Sensitivity (%) Recall or TPR	Precision (%)
SCOP fold	1e-6	92.7	88.4
CATH fold	1e-3	95.7	90.1
InterPro	1e-3	88.4	78.2
PFam family	1e-4	90.5	82.8
EC number	1e-4	93.3	79.7
GO Molecular Function	1e-1	84.3	80.9
GO Biological Process	1e-3	85.5	74.8
GO Cellular Component	1e-2	77.6	58.6

Sensitivity =
$$\frac{TP}{TP + FN}$$
 Precision = $\frac{TP}{TP + FP}$

LigandAv. binding ste seq. kd.Av. residue conservationResidues in predicted binding ste (size proportional to the local conservation)MO259.030.18548 49 52 62 63 66 67 113 116CRY20.000.11123 29 31 37 44 48 49 83 85 94 96 103 121BOG20.000.11119 20 21 48 49 51 96 98 136ACY15.870.16323 29 31 37 44 45 81 83 85 94 96 98 103 121 135Inherited partners:PartnerÅv. ste sog. id.	<u>d.113.1.1</u>	23.68	<u>0.948</u>	1920 8182	50 51 52 53 54 55 56 57 58 77 78 79 80 83 84 85 93 95 97 99 134 135 138 142 145	Real Providence
Ligand Av. binding ste seq. id. Av. resdue conservation Residues in predicted binding ste (size proportional to the local conservation) MO2 59.03 0.185 48 49 52 62 63 66 67 113 116 CRY 20.00 0.111 23 29 31 37 44 48 49 83 85 94 96 103 121 BOG 20.00 0.111 19 20 21 48 49 51 96 98 136 ACY 15.87 0.163 23 29 31 37 44 45 81 83 85 94 96 98 103 121 135	Partner	Av. binding site seq. id.	Av. residue conservatior	1	Residues in predicted binding site (size proportional to the local conservation)	
Ligand Av. binding site seq. id. Av. residue conservation Residues in predicted binding site (size proportional to the local conservation) MO2 59.03 0.185 48 49 52 62 63 66 67 113 116 CRY 20.00 0.111 23 29 31 37 44 48 49 83 85 94 96 103 121 BOG 20.00 0.111 19 20 21 48 49 51 96 98 136	ACY	15. artners:1	87	<u>0.163</u>	23 29 31 37 44 45 81 83 85 94 96 98 103 121 135	0
LigandAv. binding site seq. kd.Av. residue conservationResidues in predicted binding site (size proportional to the local conservation)MO259.030.18548 49 52 62 63 66 67 113 116CRY20.000.11123 29 31 37 44 48 49 83 85 94 96 103 121	<u>80G</u>	20.	.00	<u>0.111</u>	19 20 21 48 49 51 96 98 136	CH CH
Ligand Av. binding site seq. id. Av. residue conservation (size proportional to the local conservation) MO2 59.03 0.185 48 49 52 62 63 66 67 113 116		20.			23 29 31 37 44 48 49 83 85 94 96 103 121	1017
Ligand Av. binding site Av. residue Residues in predicted binding site conservation (size proportional to the local conservation)					48 49 52 62 63 66 67 113 116	4

Benchmark

	Number of chains
Initial set*	78,167
LigBase**	30,126
Non-redundant set***	4,948 (8,846 ligands)

*all PDB chains larger than 30 aminoacids in length (8th of August, 2006) **annotated with at least one ligand in the LigBase database

***not two chains can be structurally aligned within 3A, superimposing more than 75% of their C atoms, result in a sequence alignment with more than 30% identity, and have a length difference inferior to 50aa

	Number of chains
Initial set*	78,167
πBase **	30,425
Non-redundant set***	4,613 (11,641 partnerships)

*all PDB chains larger than 30 aminoacids in length (8th of August, 2006)

**annotated with at least one partner in the Base database

***not two chains can be structurally aligned within 3A, superimposing more than 75% of their C atoms, result in a sequence alignment with more than 30% identity, and have a length difference inferior to 50aa

Method



Inherited ligands: 4						
Ligand	Av. binding site seq. id.	Av. residue conservation	Residues in predicted binding site (size proportional to the local conservation)			
<u>MO2</u>	59.03	<u>0.185</u>	48 49 52 62 63 66 67 113 116			
CRY	20.00	<u>0.111</u>	23 29 31 37 44 48 49 83 85 94 96 103 121			
<u>80G</u>	20.00	<u>0.111</u>	19 20 21 48 49 51 96 98 136			
<u>ACY</u>	15.87	<u>0.163</u>	23 29 31 37 44 45 81 83 85 94 96 98 103 121 135			







Scoring function

Ligands

Partners





Sensitivity .vs. Precision

	Optimal cut-off	Sensitivity (%) Recall or TPR	Precision (%)
Ligands	30%	71.9	13.7
Partners	40%	72.9	55.7

Sensitivity =
$$\frac{TP}{TP + FN}$$
 Precision = $\frac{TP}{TP + FP}$

Example (2azwA) Structural Genomics Unknown Function

Molecule: MutT/nudix family protein



Acknowledgments

COMPARATIVE MODELING Andrej Sali

M. S. Madhusudhan Narayanan Eswar Min-Yi Shen **Ursula Pieper** Ben Webb Maya Topf

MODEL ASSESSMENT David Eramian Min-Yi Shen Damien Devos

FUNCTIONAL ANNOTATION Andrea Rossi Fred Davis

Prince Felipe Research Center Marie Curie Reintegration Grant

MODEL ASSESSMENT

Francisco Melo (CU) Alejandro Panjkovich (CU)

STRUCTURAL GENOMICS Stephen Burley (SGX) John Kuriyan (UCB) NY-SGXRC

FUNCTIONAL ANNOTATION Fatima Al-Shahrour Joaquin Dopazo

BIOLOGY

Jeff Friedman (RU) James Hudsped (RU) Partho Ghosh (UCSD) Alvaro Monteiro (Cornell U) Stephen Krilli (St.George H)

Tropical Disease Initiative

Stephen Maurer (UC Berkeley) Arti Rai (Duke U) Andrej Sali (UCSF) Thomas Kepler (Duke U) Ginger Taylor (TSL)

CCPR Functional Proteomics Patsy Babbitt (UCSF) Fred Cohen (UCSF) Ken Dill (UCSF) Tom Ferrin (UCSF) John Irwin (UCSF) Matt Jacobson (UCSF) Tack Kuntz (UCSF) Andrej Sali (UCSF) Brian Shoichet (UCSF) Chris Voigt (UCSF)

EVA

Burkhard Rost (Columbia U) Alfonso Valencia (CNB/UAM)

CAMP

Xavier Aviles (UAB) Hans-Peter Nester (SANOFI) Ernst Meinjohanns (ARPIDA) Boris Turk (IJS) Markus Gruetter (UE) Matthias Wilmanns (EMBL) Wolfram Bode (MPG)