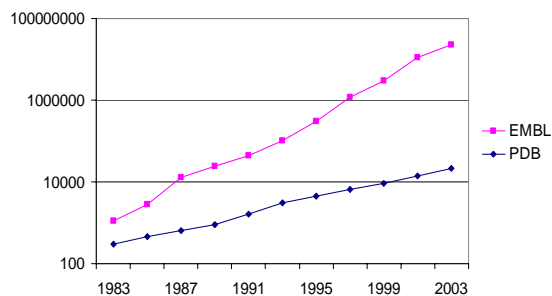


Structural Genomics

Iosif Vaisman

Email: ivaisman@gmu.edu

Dynamics of Database Growth



PDB Current Holdings

		Molecule Type				Total
		Proteins	Nucleic Acids	Protein/NA Compl	Other	
Exp. Meth.	X-ray	37254	1000	1726	24	40004
	NMR	5948	789	136	7	6880
	EM	109	11	40	0	160
	Other	83	4	4	2	93
	Total	43394	1804	1906	33	47137

PDB Redundancy

Method	Description	# of Clusters
blast	95% identity (one chain)	18106
blast	90% identity (one chain)	17367
blast	70% identity (one chain)	15507
blast	50% identity (one chain)	13201
blast	40% identity (one chain)	11457
blast	30% identity (one chain)	9458

Representative sets of known structures

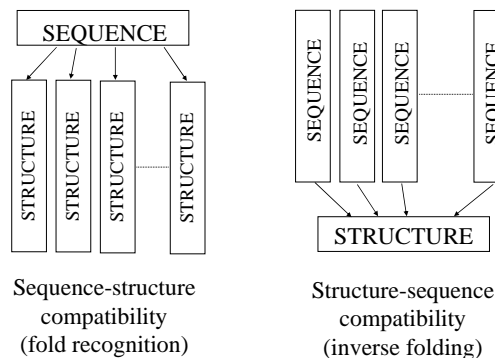
Filtering by parameters:

Resolution
R-factor
Sequence similarity
Experimental technique

PDB-REPRDB http://mbs.cbrc.jp/pdbreprdb/cgi/reprdb_menu.pl

PISCES <http://dunbrack.fccc.edu/PISCES.php>

Protein modeling



Sequence-structure correlations

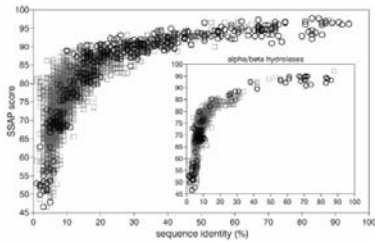
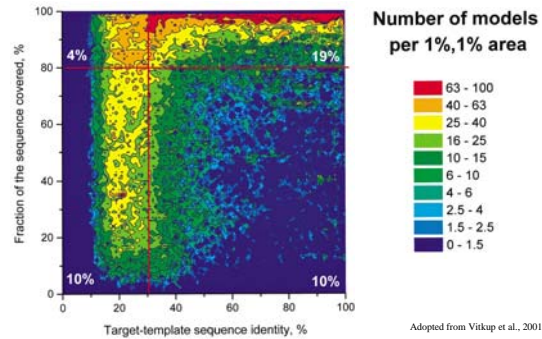


Fig. 1. Correlation between structure similarity (measured by the SSAP structure comparison algorithm, 0-100) and sequence similarity (measured by sequence identity) for all pairs of homologous domain structures in the CATH domain database.

Reifern and Orengo, 2005

Model structure coverage in sequence space



Adopted from Vitkup et al., 2001

Genome sequencing projects statistics

Organism	Complete	Draft assembly	In progress	total
Prokaryotes	590	402	460	1452
Archaea	47	4	31	82
Bacteria	543	398	429	1370
Eukaryotes	23	122	186	338
Animals	4	53	50	147
Mammals	2	21	26	49
Birds	1	3	3	3
Fish	1	2	4	9
Insects	1	19	20	40
Platworms	1	1	3	4
Roundworms	1	3	13	17
Amphibians			2	2
Reptiles			2	2
Other animals		6	19	25
Plants	3	3	24	40
Land plants	2	2	27	31
Green Algae	1	1	7	9
Fungi	10	52	31	93
Ascomycetes	8	46	21	75
Basidiomycetes	1	4	6	11
Other fungi	1	2	4	7
Protists	6	19	27	52
Alveolates	1	10	6	17
Kinetoplasts	1	2	6	9
Other protists	4	7	14	25
total	613	531	646	1790

Structural Genomics Project

- Organize known protein sequences into families.
- Select family representatives as targets.
- Solve the 3D structure of targets by X-ray crystallography or NMR spectroscopy.
- Build models for other proteins by homology to solved 3D structures.

History of Structural Genomics

1995	SG project proposed in Japan	2000 Sep.	NIGMS Protein Structure Initiative starts in US with 7 Centers
1997 Apr.	SG pilot project starts at RIKEN Inst.	2000 Nov.	International Conference on SG (ICSG 2000), Yokohama, Japan / International SG Task Force Meeting / OECD/GSF Meeting
1997	SG studies initiated through DOE, NIGMS in US	2001 Jan.	OECD/CSTP/GSF, Paris, France – Further Study on SG
1998/99	Initial SG projects start in Canada, Germany, US	2001 Apr.	2nd International SG Meeting, Airlie House, US – Start of ISGO
1999 June	Call for SG pilot projects issued by NIGMS/NIH	2001 Sep.	NIGMS Protein Structure Initiative adds 2 new centers
2000 Jan.	OECD Committee for Scientific and Technological Policy (CSTP) proposes to initiate SG studies	2002 Mar.	European Commission announces funding of Structural Proteomics in Europe (SPINE)
2000 Apr.	1st International SG Meeting, Hinxton, UK	2002 Apr.	National project on Protein Structural and Functional Analyses starts in Japan
2000 June	OECD/Global Science Forum (GSF) and SG Workshop, Florence, Italy	2002 Oct.	ISGO International Conference on SG (ICSG 2002), Berlin, Germany
2000 Sep.	SG: From Gene to Structure to Function, Cambridge, UK		

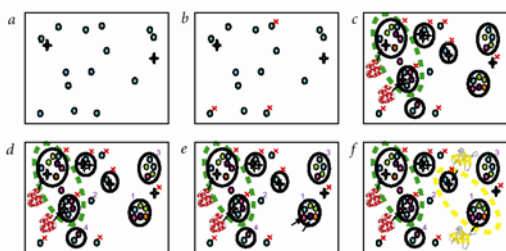
Heinemann, 2002

Goals of structural genomics

- Provision of enough structural templates to facilitate homology modeling of most proteins
- Structures of all proteins in a complete proteome
- Structural elucidation of a complete biological pathway
- Structural elucidation of a complete disease

Phil Bourne, 2005

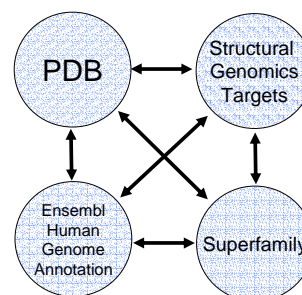
Target selection



- a) realm of interest
 b) family exclusion - impossible
 c) family exclusion - known
 d) prioritization
 e) selection
 f) analysis and interpretation

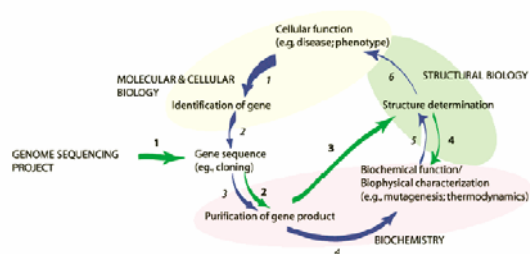
S. Brenner, 2000

Coverage of the Human Genome By Structure



Xie and Bourne, 2005

Structural genomics shortcuts



Yee et al., *Acc. Chem. Res.* **2003**, *36*, 183-189

NIGMS Protein Structure Initiative

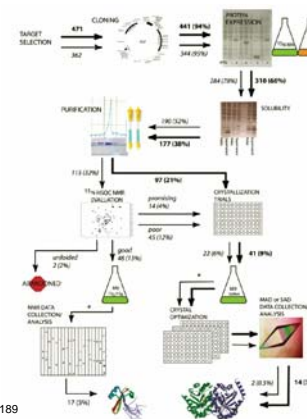
	12/2001	12/2002	12/2003	12/2004	11/2007
Selected	11214	21872	42726	74637	148735
Cloned	5465	11277	23237	45353	104366
Expressed	2860	6115	13602	25536	66458
Purified	1505	2823	5291	8398	25989
Crystallized	336	1161	1876	3199	9578
Diffraction	96	438	767	1651	4919
Crystal structure	87	314	545	1260	3781
PDB	76	247	569	1488	5235

Targets by genome

Organism	Number of targets	% Of all targets
<i>Caenorhabditis elegans</i>	4674	17.4
<i>Arabidopsis thaliana</i>	3900	14.5
<i>Homo sapiens</i>	3257	12.1
<i>Pyrococcus furiosus</i>	2179	8.1
<i>Thermotoga maritima</i>	1860	6.9
<i>Mycobacterium tuberculosis</i>	1476	5.5
<i>Escherichia coli</i>	1272	4.7
<i>Saccharomyces cerevisiae</i>	1254	4.7
<i>Bacillus subtilis</i>	1220	4.5
<i>Bacillus stearothermophilus</i>	764	2.8

Adopted from O'Toole et al., 2004

M. thermoautotrophicum structural genomics project



Yee et al., *Acc. Chem. Res.* **2003**, *36*, 183-189

Current results

Group or SG center	Targets and nonidentical chains	New PDB families total family size	Novel structures (30% ID)	New SCOP folds	New SCOP fold or superfamily
SG centers					
Berkeley Structural Genomics Center (BSGC)	57 (57 chains)	32 (5757)	41	4	6
Center for Eukaryotic Structural Genomics (CESG)	48 (48 chains)	7 (281)	28	0	0
Joint Center for Structural Genomics (JCSG)	186 (187 chains)	32 (4875)	92	3	4
Albion Center for Structural Genomics (ACSJG)	224 (229 chains)	55 (5512)	163	18	25
Northeast Structural Genomics Consortium (NESGC)	159 (159 chains)	52 (4812)	108	15	26
New York Structural Genomics Research Consortium (NYSGRC)	166 (171 chains)	27 (3922)	90	6	9
Southeast Collaboratory for Structural Genomics (SECSG)	67 (67 chains)	6 (2079)	25	0	1
Structural Genomics of Pathogenic Proteins Consortium (SGPP)	26 (26 chains)	1 (129)	8	2	2
TB Structural Genomics Consortium (TBSG)	99 (99 chains)	9 (5938)	42	0	1
PSI centers total of 9 centers above	1032 (1043 chains)	211 (30,360)	597	48	74
Japanese center (RIKEN)	466 (718 chains)	50 (6860)	289	10	20
Other International SG total, excluding all centers above	169 (183 chains)	33 (5877)	69	6	9
Non-SG groups (since 2000)					
Non-SG structural biology total	17,096 (23,747 chains)	928 (249,171)	2,521	269	478
Steffi group	44 (159 chains)	23 (4190)	31	7	12
Huber group	185 (273 chains)	8 (6279)	38	5	10
Heuts group	14 (54 chains)	14 (9760)	20	2	3

Structural genomics target database

Welcome to TargetDB - Mozilla Firefox

Home | About | Contact | Privacy Policy | Terms of Service | Sitemap

TargetDB | Target Search for Structural Genomics

TargetDB, a target registration database, provides status and tracking information on the progress of the production and solutions of structures. Search sequences from the NIH P50 Structural Genomics Center and other structural projects below.

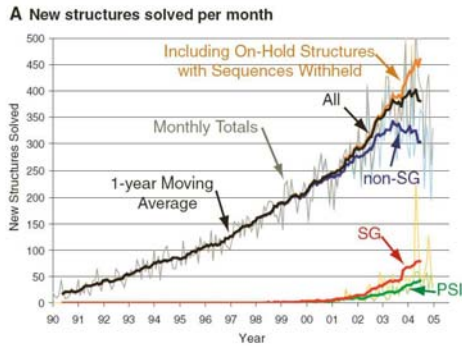
- Download all targets (XML¹ and tab delimited formats)
- Download target data document type definition.
- The XML document is organized according to the recommendations of the International Task Force on Target Tracking.
- Guidelines for preparing TargetDB data file. Last update: May 25, 2005
- Target Status Query Feature
- Summary Reports

Target sequence lists are also maintained at the following contributing sites:

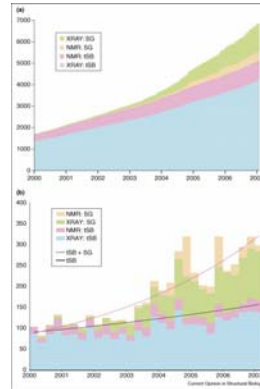
[BSGC | BSGC | BSGC | ISPI | JCSG | MCSG | MSGP | NESG | NYSGRC | OPPF | RIKEN | S2F | SECSG | SGPP | SGPP | SPNE | TB | XMTB | YSG]

Using the Target Search Form:

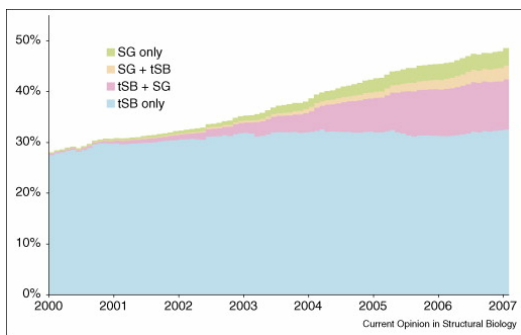
Current results



Current results



Structural coverage of the Swiss-Prot database



Grabowski et al., 2007

Practical applications of structural genomics

Fig. 3 Mis-annotation of the Rv3853 gene from *M. tuberculosis*. Originally annotated as the terminal SAM-dependent methyltransferase of menaquinone biosynthesis (MenG), Rv3853 has a monomer fold (left) that is completely different from that of a typical SAM-dependent methyltransferase such as CnaA1 (right).

