Comparative Protein Structure Prediction



Marc A. Marti-Renom





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

Program

Intro to comparative protein structure prediction

Template Search

Target – Template Alignment

Model Building

Model Evaluation

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



Objective

TO LEARN HOW-TO MODEL A 3D-STRUCTURE FROM A SEQUENCE AND A KNOWN STRUCTURE

DISCLAIMER!

Name	Type®	World Wide Web address ^b
DATABASES		
CATH	s	http://www.blochem.ucl.ac.uk/bsm/cath/
DBAII	s	http://www.salilab.org/DBAII/
GenBank	s	http://www.ncbi.nlm.nih.gov/Genbank/GenbankSearch.html
GeneCensus	s	http://bioinfo.mbb.yale.edu/genome
MODBASE	s	http://salilab.org/modbase/
MSD	s	http://www.rcsb.org/databases.html
NCBI	s	http://www.ncbi.nim.nih.gov/
PDB	s	http://www.rcsb.org/pdb/
PSI	s	http://www.nigms.nih.gov/psl/
Sacch3D	s	http://genome-www.stanford.edu/Sacch3D/
SCOP	s	http://scop.mrc-lmb.cam.ac.uk/scop/
TIGR	s	http://www.tigr.org/tdb/mdb/mdbcomplete.html
TrEMBL	s	http://srs.ebi.ac.uk/
FOLD ASSIGNM	ENT	
123D	s	http://123d.nciforf.gov/
3D-PSSM	s	http://www.sbg.bio.ic.ac.uk/~3dpssm/
BIOINBGU	s	http://www.cs.bgu.ac.il/~bioinbgu/
BLAST	s	http://www.ncbi.nim.nih.gov/BLAST/
DALI	s	http://www2.ebi.ac.uk/dall/
FASS	s	http://bioinformatics.burnham-inst.org/FFA5/index.html
FastA	s	http://www.ebi.ac.uk/fasta3/
FRSVR	s	http://foid.doe-mbi.ucia.edu/
PUGUE	s	http://www-cryst.bioc.cam.ac.uk/~fugua/

http://sgu.bioinfo.cipf.es/home/?page=resources

Programs, servers and databases

http://salilab.org



External Resources

PDB, Uniprot, GENBANK, NR, PIR, INTERPRO, Kinase Resource UCSC Genome Browser, CHIMERA, Pfam, SCOP, CATH

Nomenclature

Homology: Sharing a common ancestor, may have similar or dissimilar functions

Similarity: Score that quantifies the degree of relationship between two sequences.

Identity: Fraction of identical aminoacids between two aligned sequences (case of similarity).

Target: Sequence corresponding to the protein to be modeled.

Template: 3D structure/s to be used during protein structure prediction.

Model: Predicted 3D structure of the target sequence.

Nomenclature

Fold: Three dimensional conformation of a protein sequence (usually at domain level).

Domain: Structurally globular part of a protein, which may independently fold.

Secondary Structure: Regular subdomain structures composed by alphahelices, beta-sheets and coils (or loops).

Backbone: Protein structure skeleton composed by the carbon, nitrogen and oxygen atoms.

Side-Chain: Specific atoms identifying each of the 20 residues types.



General References

Protein Structure Prediction:

Marti-Renom et al. Annu. Rev. Biophys. Biomol. Struct. 29, 291-325, 2000. Baker & Sali. Science 294, 93-96, 2001.

Comparative Modeling:

Marti-Renom et al. Annu. Rev. Biophys. Biomol. Struct. 29, 291-325, 2000. Madhusudhan et al. The Proteomics Protocols Handbook. Ed. Walker. Humana Press Inc., Totowa, NJ. 831-860, 2005.

MODELLER:

Sali & Blundell. J. Mol. Biol. 234, 779-815, 1993.

Structural Genomics:

Sali. Nat. Struct. Biol. 5, 1029, 1998. Burley et al. Nat. Genet. 23, 151, 1999. Sali & Kuriyan. TIBS 22, M20, 1999. Sanchez et al. Nat. Str. Biol. 7, 986, 2000. Baker & Sali. Science 294, 93-96, 2001.

protein prediction .vs. protein determination



Why is it useful to know the structure of a protein, not only its sequence?

- 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.



In addition, since evolution tends to conserve function and function depends more directly on structure than on sequence, **structure is more conserved in evolution than sequence**.

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

Principles of protein structure

GFCHIKAYTRLIMVG...





Folding (physics)

Ab initio prediction

Evolution (rules) Threading Comparative Modeling



N. Eswar, et al. Comparative Protein Structure Modeling With MODELLER. Current Protocols in Bioinformatics, John Wiley & Sons, Inc., Supplement 15, 5.6.1-5.6.30, 2008.
 M.A. Marti-Renom, et al.. Comparative protein structure modeling of genes and genomes. Annu. Rev. Biophys. Biomol. Struct. 29, 291-325, 2000.
 A. Sali & T.L. Blundell. Comparative protein modelling by satisfaction of spatial restraints. J. Mol. Biol. 234, 779-815, 1993.
 A. Fiser, R.K. Do, & A. Sali. Modeling of loops in protein structures, Protein Science 9. 1753-1773, 2000.

MODELLER

Steps in Comparative Protein Structure Modeling





A. Šali, Curr. Opin. Biotech. 6, 437, 1995.
R. Sánchez & A. Šali, Curr. Opin. Str. Biol. 7, 206, 1997.
M. Marti et al. Ann. Rev. Biophys. Biomolec. Struct., 29, 291, 2000.

Template Selection "Structural Space"

Structure-Structure alignments

As any other bioinformatics problem...

- Representation
 - Scoring
 - Optimizer

Representation

Structures



All atoms and coordinates







Reduced atom representation







Vector representation

Secondary Structure

Accessible surface (and others)

Scoring

Raw scores

	¢		Τ	7	А.		.8	3	1	9	я.	8	Ε.	м		5	7	1	T	π
٤.		-1	- 4				- 4	-4	-4	+	-1	-1	-1	-1	- 4	-4	-4	-1	-4	-
1	14		1		1	1.1	1	11		-1	-4	1	1	- 11	- 0	- 2	- 12	-4	a	
T.	1.1	1			-4							-4		- 4	-4	-4	- 4	-4	-0.	
F.	-4	-4	11	. 7	1	-1	-1	-	- 10	4	-4	- 4	- 18	- 2	- 2	-4	-1	-4	-11	
٩.	. +	1	- 41	-1			-1	. 4	-4	- 0	-4	1	-1	-1	-4	- 4	-1	-4	-4	
1	- 4		- 11	- 4			- 4	-11	- 2	2	-4	-1	- 4		-4	- 4	- 1	-1	-11	-
η	- 4	1		-1	10			- 1	.1		-1	. 0		. 4		8	-1	8	- 4	
2	-1	- 0	- 11	- 21	- 2	-4	1		- 2		-4	- 1	-1	- 3	- 21	- 4	-4	-1	- 3	
	- 4	- 1		10	- 4	-4	1	- 2	. 8	- 2		1	- 1	12	2	- 4	-1	-1	- 2	
٤.	- 4			-4	-4	-4			1			1	1		-4	. 4	-4	-4	-4	-
Ε.	- 4		- 8	- 4	0	- 4	1	- 1	. 8	. 0	٠		- 18	12	- 2	. 4	- 12	14	1.2	-
	- 4	-4	-4	- 4	-4	-4		-2		. 6			1	-4	. +	3	- 4	-1	1.1	-
E.	1.1	- 1		-18	- 4	- 4	- 0		- 6	1	-4	- 1		- 11	-2	4	- 4	-1	-4	
ĥĺ.	-4	-1	-0	- 4	- 4	-4	- 4	-0	् ने		-1	-4	- 4		. 1	1	-1	. 4	-1	-
	- 4	1	- 4	- 4	-4	-4	-4	- 2		-2	-4	-4	-1	1		- 1	1	. 4	-4	
	-1	-1	-8	-4	-0	-4	-14	-1	-0	-1	-4	-0	-1	1			- 4		-11	1
1	- 4	-1	21	2		-4	-4	- 10	2	-2	-4	-4	-1	- 11	. 8	- 1	- 4	1.4	- 4	-
٢.,	-1	-d	-0	-4	2	-1	-1	-0	18	0	-1	1.1	- 4		. 8			. 4	3	
1	-1	- 1	-2	- 11	-2	-4	- 4	-11	-2	4	- 1	-1	-1	- 41	.4	-4	4	- 1		
	1.1	1.1	-0	1.4		1.4	- 4	-4	- 0.	-0	-4	- 18	- 18	1	-0	- 4	1.1	1	- 2	11

Aminoacid substitutions

 $RMSD(x,y) = \sqrt{\left(\frac{1}{N}\right) \sum_{i=1}^{N} \left(\left\| \mathbf{x}(i) - \mathbf{y}(i) \right\|^{2} \right)}$

Root Mean Square Deviation







Secondary Structure (H,B,C)

Accessible surface (B,A [%])

Angles or distances

Scoring

Significance of an alignment (score)

Probability that the optimal alignment of two random sequences/structures of the same length and composition as the aligned sequences/structures have at least as good a score as the evaluated alignment.



Sometimes approximated by Z-score (normal distribution).



Optimizer Global dynamic programming alignment



Backtracking to get the best alignment

Optimizer

Local dynamic programming alignment



Backtracking to get the best alignment

Optimizer

Global .vs. local alignment



Optimizer

Multiple alignment

Pairwise alignments

Example – 4 sequences A, B, C, D.



6 pairwise comparisons then cluster analysis

Multiple alignments

Following the tree from step 1



Align B-D with A-C



Coverage .vs. Accuracy





Same RMSD ~ 2.5Å

Coverage ~90% C α

Coverage ~75% Ca

Structural alignment by properties conservation (SALIGN-MODELLER)



M. S. Madhusudhan, B. M. Webb, M. A. Marti-Renom, N. Eswar, A. Sali, Protein Eng Des Sel, (Jul 8, 2009).

Structural alignment by properties conservation (SALIGN-MODELLER)

http://salilab.org/DBAli



M. A. Marti-Renom et al., Nucleic Acids Res 35, W393 (Jul 1, 2007)

Vector Alignment Search Tool (VAST)



Gibrat JF et al. (1996) Curr Opin Struct Biol 3 pp377

Vector Alignment Search Tool (VAST)

http://www.ncbi.nlm.nih.gov/Structure/VAST/vast.shtml



Incremental combinatorial extension (CE)



Shindyalov IN, amd Bourne PE. (1998) Protein Eng. 9 pp739

Incremental combinatorial extension (CE)

http://cl.sdsc.edu/ce.html



University at Albany, NY)

Matching molecular models obtained from theory (MAMMOTH)



Ortiz AR, (2002) Protein Sci. 11 pp2606

30

Matching molecular models obtained from theory (MAMMOTH)

http://ub.cbm.uam.es/mammoth/pair/index3.php



Classification of the structural space



SCOP_{1.75} database

http://scop.mrc-lmb.cam.ac.uk/scop/



Andreeva A., Howorth D., Brenner S.E., Hubbard T.J.P., Chothia C., Murzin A.G. (2004). SCOP database in 2004: refinements integrate structure and sequence family data. Nucl. Acid Res. 30(1), 264-265. [PDF]. Andreeva A., Howorth D., Brenner S.E., Hubbard T.J.P., Chothia C., Murzin A.G. (2004). SCOP database in 2004: refinements integrate structure and sequence family data. Nucl. Acid Res. 32:D225. D229. [PDF]. and Andreeva A., Howorth D., Chandonia J.-M., Brenner S.E., Hubbard T.J.P., Chothia C., Murzin A.G. (2008). Data growth and its impact on the SCOP database: new developments. Nucl. Acid Res. 35: D419-D425. [PDF].

Access methods

- · Enter scop at the top of the hierarchy
- Keyword search of SCOP entries
- SCOP parseable files
- All SCOP releases and reclassified entry history
- pre-SCOP preview of the next release
- SCOP domain sequences and pdb-style coordinate files (<u>ASTRAL</u>)
- Hidden Markov Model library for SCOP superfamilies (SLPERFAMILY)

scop help and information

- · Hidden Markov Model library for SCOP superfamilies (SUPERCAMILY)
- SCOP domain sequences and pdb-style coordinate files (ASTRAL)
 IEIdates Markow Model Breast for STOP superflowline (NURSEAMEN)
- · pre-score preview of the next release
- All SCOP releases and reclassified conty instery.
- SCOP parsonic files
- Velocity reaction
- · Enter scor at the top of the hierarchy

- Largely recognized as "standard of gold"
- ✓ Manually classification
- ✓ Clear classification of structures in:
 - CLASS FOLD SUPER-FAMILY FAMILY
- ✓ Some large number of tools already available

Manually classification Not 100% up-to-date Domain boundaries definition

Class	Number of folds	Number of superfamilies	Number of families		
All alpha proteins	284	507	871		
All beta proteins	174	354	742		
Alpha and beta proteins (a/b)	147	244	803		
Alpha and beta proteins (a+b)	376	552	1055		
Multi-domain proteins	66	66	89		
Membrane and cell surface proteins	58	110	123		
Small proteins	90	129	219		
Total	1195	1962	3902		

Murzin A. G., el at. (1995). J. Mol. Biol. 247, 536-540.

CATH_{3.2} database

http://www.cathdb.info



Uses FSSP for superimposition

- ✓ Recognized as "standard of gold"
- ✓ Semi-automatic classification
- Clear classification of structures in: CLASS ARCHITECTURE TOPOLOGY HOMOLOGOUS SUPERFAMILIES
- ✓ Some large number of tools already available
- ✓ Easy to navigate

Semi-automatic classification Domain boundaries definition

Class	Architecture	Topology	Homologous Superfamily	S35 Family	S60 Family	S95 Family	S100 Family	Domain
1	5	310	682	2078	2689	3540	6685	23491
2	20	196	438	2062	2902	4468	7656	29992
3	14	512	956	4558	6473	8135	16346	58967
4	1	92	102	173	217	301	445	1765
Total	40	1110	2178	8871	12281	16444	31132	114215

DBAliv2.0 database

http://salilab.org/DBAli/



- Analyze the data deposited in UDAR (Tools)

Uses MAMMOTH for superimposition

✓ Fully-automatic

- ✓ Data is kept up-to-date with PDB releases
- ✓ Tools for "on the fly" classification of families
- ✓ Up-to-date multiple structure alignments
- ✓ Easy to navigate
- ✓ Provides some tools for structure comparison

Does not provide a stable classification

Pairwise structure alignmen	ts
Last update:	October 6th, 2007
Number of chains:	96,804
Number of structure-structure comparisons:*	1,748,371,897
Multiple structure alignment	18
Last update:	August 1st, 2007
Number of representative chains:	34,637
Number of families:	12,732

Marti-Renom et al. 2001. Bioinformatics. 17, 746 Marti-Renom et al. 2007. BMC BMC Bioinformatics (2007) 8 (Suppl 4) S4 Marti-Renom et al. 2007. Nucleic Acid Research (2007) 35 W393-W397

Classification of the structural space Not an easy task!

Domain definition AND domain classification



Day, et al. (2003) Protein Sciences, 12 pp2150
template search and template-target alignment (pp_scan)

Marti-Renom, et al. (2004) Prot. Sci. 13 pp1071 Narayanan, et al. in prepration

PP_SCAN or profile-profile alignments





Seq.-Seq

Seq.-Str

Prof.-Seq

Prof.-Prof.





ALIGN: DP pairwise method

BLAST2SEQ: Local heuristic method

SEA: Local structure prediction method

SAM: HMM method
PSI-BLAST: Local search method that uses multiple sequence information for one of the sequences.
LOBSTER: HHM + Phylogeny Method

CLUSTALW: DP multiple sequence method. COMPASS: DP profile-profile method

PP_SCAN: DP pairwise method that uses multiple sequence information for both sequences.

PP_SCAN protocols

Profile generation

- PSI-Blast (PBP)
- Henikoff & Henikoff (HH)
- Henikoff & Henikoff + Similarity (HS)
- Henikoff & Henikoff substitution matrix (MAT)

Profile comparison

- Correlation coefficient (CC)
- Euclidean distance (ED)
- Dot product (DP)
- Jensen-Shannon distance (JS)
- Average value (Ave)

PP_SCAN protocols accuracy

SALIGN protocol	CE overlap [%]	Shift score
ССрвр	55 ± 23	0.61 ± 0.24
ССнн	56 ± 23	0.61 ± 0.24
ССнѕ	56 ± 24	0.62 ± 0.23
ССмат	51 ± 25	0.55 ± 0.27
ЕДрвр	54 ± 24	0.60 ± 0.25
ЕДнн	54 ± 24	0.59 ± 0.26
EDHs	55 ± 24	0.59 ± 0.26
DРрвр	55 ± 23	0.61 ± 0.24
DРнн	56 ± 23	0.60 ± 0.25
DPHs	55 ± 24	0.61 ± 0.24
JSнн	53 ± 24	0.60 ± 0.24
JSнs	54 ± 24	0.60 ± 0.24
Ауемат	49 ± 26	0.52 ± 0.29
ТОР	62 ± 20	0.67 ± 0.20

PP_SCAN accuracy

Method	CE overlap	Shift score			
CE	100 ± 0	1.00 ± 0.00			
BLAST	26 ± 29	0.32 ± 0.33			
PSI-BLAST	43 ± 31	0.48 ± 0.35			
SAM	48 ± 26	0.50 ± 0.34			
LOBSTER	50 ± 27	0.51 ± 0.32			
SEA	49 ± 27	0.53 ± 0.29			
ALIGN	42 ± 25	0.44 ±0.28			
CLUSTALW	43 ± 27	0.44 ± 0.31			
COMPASS	43 ± 32	0.49 ± 0.35			
ССнн	56 ± 23	0.61 ± 0.24			
ССнѕ	56 ± 24	0.62 ± 0.24			
ТОР	62 ± 20	0.67 ± 0.20			



PP_SCAN success



Alignment accuracy (CE overlap) 200 pairwise DBAli alignments

PSI-BLAST (sequence-profile alignment) 43%

SEA (local structure alignment) 49%

PP_SCAN (profile-profile alignment) 56%



model building and model assessment

Information about a protein can come from three distinct sources



Experimental observations





Statistical rules



Laws of physics

Classes of methods for comparative protein structure modeling

- Model building by assembly of rigid bodies core, loops, sidechains.
- Model building by segment matching.
- Model building by satisfaction of spatial restraints.

Comparative modeling by satisfaction of spatial restraints MODELLER



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.

Multiple Templates

Local similarity extracted from closest template



Templates Target

KSINPIHGDNCEQTSDEGLKIERTPL----QWLKSSICDMRGLIPE ASILPKRLFGNCEQTSDEGLKIERTPLVPHISAQNVCLKIDDVPERLIPE MSVIPKRLYGNCEQTSEEAIRIEDSPIVRWISAQLVCLKIDEIPERLVGE

Modeling ligands and using external restraints



Homology derived restraint

External Restraint



Accuracy and applicability of comparative models

Comparative modeling by satisfaction of spatial restraints Types of errors and their impact



Marti-Renom etal. Ann Rev Biophys Biomol Struct (2000) 29, 291

"Biological" significance of modeling errors



NMR – X-RAY Erabutoxin 3ebx Erabutoxin 1era

NMR Ileal lipid-binding protein 1eal



CRABPII1opbBFABP1ftpAALBP1lib40% seq. id.

X-RAY Interleukin 1 β 41bi (2.9Å) Interleukin 1 β 2mib (2.8Å)



Model Accuracy

HIGH ACCURACY

NM23 Seq id 77%

Cα equiv 147/148 RMSD 0.41Å



Sidechains Core backbone Loops

MEDIUM ACCURACY

CRABP Seq id 41%

Cα equiv 122/137 RMSD 1.34Å



Sidechains Core backbone Loops Alignment

LOW ACCURACY

EDN Seq id 33%

Cα equiv 90/134 RMSD 1.17Å



Sidechains Core backbone Loops Alignment Fold assignment

X-RAY / MODEL

Utility of protein structure models, despite errors



Model Assessment (PMF)

Scoring Statistical Potential (inspiration)

$$K = \frac{\begin{bmatrix} AB \end{bmatrix}}{\begin{bmatrix} A \end{bmatrix} \cdot \begin{bmatrix} B \end{bmatrix}}$$
$$\Delta G = -RT \ln(K) = -RT \ln \frac{\begin{bmatrix} AB \end{bmatrix}}{\begin{bmatrix} A \end{bmatrix} \cdot \begin{bmatrix} B \end{bmatrix}}$$

From statistical physics, we know that energy difference between two states (ΔE) and the ratio of their occupancies (N₁:N₂) are related [9]:

$$\Delta E = -kT \ln \left(\frac{N_1}{N_2}\right) \qquad (1)$$

in which T is the absolute temperature and k is the Boltzmann's constant. As we are interested in an interaction energy between two amino acid side chains, it would seem natural to define N_1 as the number of interactions between these two residues types in a group of real protein structures, a number which is readily available from simple database analysis. But this number must be compared with the number of interactions in some other system, N_2 , to obtain the energy difference between them.

$+ B \rightleftharpoons AB$



Tanaka and Sheraga (1975) PNAS, **72** pp3802 Sippl, (1990) J.Mo.Biol. **213** pp859 Godzik, (1996) Structure **15** pp363

Scoring Significance of an alignment (score)

Energy Z-score the model with respect the energy of random models (or rest of decoys).



Prosall

http://www.came.sbg.ac.at



ANOLEA

http://protein.bio.puc.cl/cardex/servers/anolea/



Verify3D

http://nihserver.mbi.ucla.edu/Verify_3D/



DFIRE

http://sparks.informatics.iupui.edu/

Deriving

Structural space



Scoring

Pseudo-Energy with respect a ideal gas-phase reference state

DOPE (MODELLER)

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

Deriving

Structural space



Scoring

Pseudo-Energy with respect a ideal spherical protein as a reference state





John, Sali (2003). NAR pp31 3982

Moulding: iterative alignment, model building, model assessment



Genetic algorithm operators







Also, "two point crossover" and "gap deletion".

Composite model assessment score

Weighted linear combination of several scores:

- Pair (Pp) and surface (Ps) statistical potentials;
- Structural compactness (S_C);
- Harmonic average distance score (H_a);
- Alignment score (A_S) .

$Z = 0.17 Z(P_P) + 0.02 Z(P_s) + 0.10 Z(S_c) + 0.26 Z(H_a) + 0.45 (A_s)$

 $Z(\text{score}) = (\text{score-}\mu)/\sigma$ μ ... average score of all models σ ... standard deviation of the scores

Benchmark with the "very difficult" test set

D. Fischer threading test set of 68 structural pairs (a subset of 19)

Target -template Sequence identity [%]		Initial prediction		Final prediction		Best prediction		
	Sequence identity [%]	Coverage [% aa]	Cα RMSD [Å]	overlap [%]	RMSD [Å]	overlap [%]	RMSD [Å]	overlap [%]
1ATR-1ATN	13.8	94.3	19.2	20.2	18.8	20.2	17.1	24.6
1BOV-1LTS	4.4	83.5	10.1	29.4	3.6	79.4	3.1	92.6
1CAU-1CAU	18.8	96.7	11.7	15.6	10.0	27.4	7.6	47.4
1COL-1CPC	11.2	81.4	8.6	44.0	5.6	58.6	4.8	59.3
1LFB-1HOM	17.6	75.0	1.2	100.0	1.2	100.0	1.1	100.0
1NSB-2SIM	10.1	89.2	13.2	20.2	13.2	20.1	12.3	26.8
1RNH-1HRH	26.6	91.2	13.0	21.2	4.8	35.4	3.5	57.5
1YCC-2MTA	14.5	55.1	3.4	72.4	5.3	58.4	3.1	75.0
2AYH-1SAC	8.8	78.4	5.8	33.8	5.5	48.0	4.8	64.9
2CCY-1BBH	21.3	97.0	4.1	52.4	3.1	73.0	2.6	77.0
2PLV-1BBT	20.2	91.4	7.3	58.9	7.3	58.9	6.2	60.7
2POR-2OMF	13.2	97.3	18.3	11.3	11.4	14.7	10.5	25.9
2RHE-1CID	21.2	61.6	9.2	33.7	7.5	51.1	4.4	71.1
2RHE-3HLA	2.4	96.0	8.1	16.5	7.6	9.4	6.7	43.5
3ADK-1GKY	19.5	100.0	13.8	26.6	11.5	37.7	7.7	48.1
3HHR-1TEN	18.4	98.9	7.3	60.9	6.0	66.7	4.9	79.3
4FGF-81IB	14.1	98.6	11.3	24.0	9.3	30.6	5.4	41.2
6XIA-3RUB	8.7	44.1	10.5	14.5	10.1	11.0	9.0	34.3
9RNT-2SAR	13.1	88.5	5.8	41.7	5.1	51.2	4.8	69.0
AVERAGE	14.2	85.2	9.6	36.7	7.7	44.8	6.3	57.8

Application to a difficult modeling case 1BOV-1LTS



Sequence identity 4.4%

Initial model C α RMSD 10.1Å

Final model C α RMSD 3.6Å

d

Can we use models to infer function?







Modeling genes

What is the physiological ligand of Brain Lipid-Binding Protein?

Predicting features of a model that are not present in the template

BLBP/oleic acid BLBP/docosahexaenoic acid Cavity is not filled Cavity is filled Ligand binding cavity 1. BLBP binds fatty acids. 2. Build a 3D model. 3. Find the fatty acid that fits most snuggly into the ligand binding cavity.

L. Xu, R. Sánchez, A. Šali, N. Heintz, J. Biol. Chem. 271, 24711, 1996.
Do mast cell proteases bind proteoglycans? Where? When?

Predicting features of a model that are not present in the template

- 1. mMCPs bind negatively charged proteoglycans through electrostatic interactions
- 2. Comparative models used to find clusters of positively charged surface residues.
- 3. Tested by site-directed mutagenesis.







Huang *et al. J. Clin. Immunol.* **18**,169,1998. Matsumoto *et al. J.Biol.Chem.* **270**,19524,1995. Šali *et al. J. Biol. Chem.* **268**, 9023, 1993.



S. cerevisiae ribosome



Fitting of comparative models into 15Å cryoelectron density map.

43 proteins could be modeled on 20-56% seq.id. to a known structure.

The modeled fraction of the proteins ranges from 34-99%.

C. Spahn, R. Beckmann, N. Eswar, P. Penczek, A. Sali, G. Blobel, J. Frank. Cell 107, 361-372, 2001.

Common Evolutionary Origin of Coated Vesicles and Nuclear Pore Complexes

mGenThreader + *SALIGN* + *MOULDER*

D. Devos, S. Dokudovskaya, F. Alber, R. Williams, B.T. Chait, A. Sali, M.P. Rout. Components of Coated Vesicles and Nuclear Pore Complexes Share a Common Molecular Architecture. *PLOS Biology* **2(12)**:e380, 2004

yNup84 complex proteins



All Nucleoporins in the Nup84 Complex are Predicted to Contain β -Propeller and/or α -Solenoid Folds



NPC and Coated Vesicles Share the β -Propeller and α -Solenoid Folds and Associate with Membranes



NPC and Coated Vesicles Both Associate with Membranes



Alber et al. The molecular architecture of the nuclear pore complex. Nature (2007) vol. 450 (7170) pp. 695-701

A Common Evolutionary Origin for Nuclear Pore Complexes and Coated Vesicles? The proto-coatomer hypothesis



Tropical Disease Initiative (TDI)

Predicting binding sites in protein structure models.



http://www.tropicaldisease.org



Need is High in the Tail

DALY Burden Per Disease in Developed Countries DALY Burden Per Disease in Developing Countries



Disease data taken from WHO, <u>World Health Report 2004</u> DALY - Disability adjusted life years

DALY is not a perfect measure of market size, but is certainly a good measure for importance.

DALYs for a disease are the sum of the years of life lost due to premature mortality (YLL) in the population and the years lost due to disability (YLD) for incident cases of the health condition. The DALY is a health gap measure that extends the concept of potential years of life lost due to premature death (PYLL) to include equivalent years of 'healthy' life lost in states of less than full health, broadly termed disability. One DALY represents the loss of one year of equivalent full health.

"Unprofitable" Diseases and Global DALY (in 1000's)

Malaria*	46,486	Trichuriasis	I,006	
Tetanus	7,074	Japanese encephalitis	709	
Lymphatic filariasis*	5,777	Chagas Disease*	667	
Syphilis	4,200	Dengue*	616	
Trachoma	2,329	Onchocerciasis*	484	
Leishmaniasis*	2,090	Leprosy*	199	
Ascariasis	1,817	Diphtheria	185	
Schistosomiasis*	1,702	Poliomyelitise	151	
Trypanosomiasis*	1,525	Hookworm disease	59	

Disease data taken from WHO, World Health Report 2004

DALY - Disability adjusted life year in 1000's.

* Officially listed in the WHO Tropical Disease Research disease portfolio.

Comparative docking



DBAliv2.0 database

http://www.dbali.org



Modeling Genomes

data from models generated by ModPipe (Eswar, Pieper & Sali)



A good model has MPQS of 1.0 or higher

Summary table

models with inherited ligands

29,271 targets with good models, 297 inherited a ligand/substance similar to a known drug in DrugBank

	Transcripts	Modeled targets	Selected models	Inherited ligands	Similar to a drug	Drugs
C. hominis	3,886	1,614	666	197	20	13
C. parvum	3,806	1,918	742	232	24	13
L. major	8,274	3,975	١,409	478	43	20
M. leprae	١,605	1,178	893	310	25	6
M. tuberculosis	3,991	2,808	I,608	365	30	10
P. falciparum	5,363	2,599	818	284	28	13
P. vivax	5,342	2,359	822	268	24	13
T. brucei	7,793	1,530	300	138	13	6
T. cruzi	19,607	7,390	3,070	769	51	28
T. gondii	9,210	3,900	1,386	458	39	21
TOTAL	68,877	29,271	11,714	3,499	297	143

L. major Histone deacetylase 2 + Vorinostat

Template 1t64A a human HDAC8 protein.



PDB	ED	Template	653	Model	C+	Ligand	Exact	SupStr	SubStr	Similar
1c3sA	83.33/80.00	1t64A	36.00/1.47	LmjF21.0680.1.pdb	90.91/100.00	SHH	DB02546	DB02546	DB02546	DB02546



DB02546 Vorinostat

Small Molecule; Approved; Investigational

Drug categories:

Anti-Inflammatory Agents, Non-Steroidal Anticarcinogenic Agents Antineoplastic Agents Enzyme Inhibitors

Drug indication:

For the treatment of cutaneous manifestations in patients with cutaneous T-cell lymphoma who have progressive, persistent or recurrent disease on or following two systemic therapies.



L. major Histone deacetylase 2 + Vorinostat

Literature

Proc. Natl. Acad. Sci. USA Vol. 93, pp. 13143–13147, November 1996 Medical Sciences

Apicidin: A novel antiprotozoal agent that inhibits parasite histone deacetylase

(cyclic tetrapeptide/Apicomplexa/antiparasitic/malaria/coccidiosis)

Sandra J. Darkin-Rattray^{*†}, Anne M. Gurnett^{*}, Robert W. Myers^{*}, Paula M. Dulski^{*}, Tami M. Crumley^{*}, John J. Allocco^{*}, Christine Cannova^{*}, Peter T. Meinke[‡], Steven L. Colletti[‡], Maria A. Bednarek[‡], Sheo B. Singh[§], Michael A. Goetz[§], Anne W. Dombrowski[§], Jon D. Polishook[§], and Dennis M. Schmatz^{*}

Departments of *Parasite Biochemistry and Cell Biology, [‡]Medicinal Chemistry, and [§]Natural Products Drug Discovery, Merck Research Laboratories, P.O. Box 2000, Rahway, NJ 07065

ANTIMICROBIAL AGENTS AND CHEMOTHERAPY, Apr. 2004, p. 1435–1436 0066-4804/04/\$08.00+0 DOI: 10.1128/AAC.48.4.1435–1436.2004 Copyright © 2004, American Society for Microbiology. All Rights Reserved. Vol. 48, No. 4

Antimalarial and Antileishmanial Activities of Aroyl-Pyrrolyl-Hydroxyamides, a New Class of Histone Deacetylase Inhibitors

P. falciparum tymidylate kinase + zidovudine

Template 3tmkA a yeast tymidylate kinase.



PDB	C0	Template	655	Model	0	Ligand	Exact	SupStr	SubStr	Similar
2tmkB	100.00/100.00	3tmkA	41.00/1.49	PFL2465c.2.pdb	82.61/100.00	ATM		DB00495		DB00495



DB00495 Zidovudine Small Molecule; Approved Drug categories: Anti-HIV Agents Antimetabolites Nucleoside and Nucleotide Reverse Transcriptase Inhibitors Drug indication:

P. falciparum tymydilate kinase + zidovudine

NMR Water-LOGSY and STD experiments



Leticia Ortí, Rodrigo J. Carbajo, and Antonio Pineda-Lucena

TDI's kernel

http://tropicaldisease.org/kernel



TDI's kernel

http://tropicaldisease.org/kernel

L. Orti et al., Nat Biotechnol 27, 320 (Apr, 2009).

CORRESPONDENCE

A kernel for the Tropical Disease Initiative

To the Editor: Io the Editor: Identifying proteins that are good drug targets and finding drug leads that bind to them is generally a challenging problem. It is particularly difficult for neglected tropical diseases, such as malaria and tuberculosis. TDI website above. where research resources are relatively scarce1. Fortunately, several developments with our software pipeline6.7 for scarde. Fortunately, several developments improve our ability to add with drug sequencing of many complete genomes sequencing of many complete genomes of organism that case tropical diseases: second, the determinism of a large number of composal mixtures, thind, the creation of composal mixtures, including already terms, backing large data that blud predicting ligands that blud predicting l approved drugs; and fourth, the availability of linked 297 proteins from approved drags and fourth, the availability of improved bioinformatic analysis, including to the pathogan genomes with methods for comparative protein instrumer. In pathogan genomes with a second second second second second legand seconding and fund design. Therefore, the second second detarility in play-second second second second second second drags leads for neglected topolar diseases. Here we encourage a collobariation amount sections to engage in drag diseases. Here we encourage a collobariation amount sections to engage in drag diseases. Here is neglected topolar diseases. Here is neglected to the Tropical Disease Initiative (TDI, http:// were tested for their binding www.tropicaldisease.org/)². As the Linux to a known drug by NMR the Toojcal Disease Initiative (TD).http:// were tested for their binding worktropicalificases(r) As the Linux kernel did for open source cold development. sectoropy, vikidning one weagent that the TD leared may hop or our predictions of Gig L and support the sector of our compactions absence of a critical mass of precessing work absence of a critical mass of precessing work and was vibrary of our companying of our companying predictions based on this initiated that volumers can bound on interioritania); Predictions to decide of this semet comparements several of the initiatives on neglected tropical diseases^{1–5}, including collaborative web portals (e.g., and facilities to test additional predictions. http://www.thesynapticleap.org/), public-ber of the set of the set

Table 1 TDI kernel genomes

3,806

5 363 818 822

(Å



Science Commons protocol for implementin open access data (http://sciencecommons. org/projects/publishing/open-access-data-protocol/), which prescribes standard academic attribution and facilitates tracking of the standard state of the state of work but imposes no other restrictions. W do not seek intellectual property rights in the actual discoveries based on the TDI kernel, in actual discoveries based on the 119 kernel, in the hope of reinvigorating drug discovery for neglected tropical discase?. By minimizing restrictions on the data, including viral terms that would be inherited by all derivative works, we hope to attract as many cychalls as negregable beinger to see and immerse the hord we possibly can to use and improve the kernel. Although many of the drugs in the kernel are proprietary under diverse types of rights, we believe that the existence of public domain pairs of targets and compounds will reduce the royalties that patent owners can charge and sponsors must pay. This should decrease the large sums of money governments and

VOLUME 27 NUMBER 4 APRIL 2009 NATURE BIOTECHNOLOGY

Transcripts^b Modeled targets^c Similar^d

1,386

L. Orti et al., PLoS Negl Trop Dis 3, e418 (2009).

PLOS REGLECTED TROPKAL DISEASES OPEN CACCESS Freely available on line A Kernel for Open Source Drug Discovery in Tropical Diseases Leticia Orti^{1,2}, Rodrigo J. Carbajo², Ursula Pieper³, Narayanan Eswar³", Stephen M. Maurer⁴, Arti K. Rai⁵, Ginger Taylor⁶, Matthew H. Todd⁷, Antonio Pineda-Lucena², Andrej Sali³, Marc A. Marti-Renom¹ 1 Structural Genomics Unit, Bioinformatics and Genomics Department, Centro de Investigación Principe Felipe, Valencia, Spain, 2 Structural Biology Laboratory, Medicinal Chemistry Department, Centro de Investigación Principe Felipe, Valencia, Spain, 3 Department of Bioengineering and Therapeutic Sciences, Department of Bennamedia Chemistry, and California Intellio. To example researce and the segment of the seg Abstract Background Conventional nature based drug development incentive work lady for the developing world, where commonical markets are usually and its one variations. For the reason the part dicade has meet noticitive experimentation with alternative R8D institutions ranging from prister-public patternity to development prists. Despite extensis discussion, however, one of the most promoting inservance—quest source drug discurge—hyst market R8D. We argue that the similarity B4Oc has been the absence of a critical mass of prescriptions work that volunteers can improve through a source that the similarity back has been presented by source drug discurge—hyst market discussion. Nevere 4.1. Attehnology Phincipal Tradings lines, we use a compactional pilerine for III comparative structure modeling at supports, III predicting the localization of light addings stare, on their uncertex, and III assessing be limitality of the predictive light of the localization of light addings stare, on their uncertex, and III assessing be limitality of the predictive light of the localization of light addings stare, on their uncertex, and III assessing be limitality of the predictive light of the localization of light addings stare, on their uncertex, and III assessing the localization of light and light adding stare to their uncertex and III assessing uncertex. The limit are lawnow days are proteinstally interproved as source of potential days targets and dug candidates around which an online gene source commission can unclusive. Limits and an uncertex adding stare addings and uncertex and the limits are lawnow days can uncetter. Limits are addentiated as the adding stare adding stare to the stare stares. Commiting on eard on unclusive star adding stare adding stare adding stare adding stare stares and which an online gene source committing and and unclusive stares around which an online gene source committing and and unclusive stares around which an online gene source committing and and unclusive stares around which an online gene source committing and and unclusive stares around which an online gene source committing and and unclusive stares are adding stares around which an online gene source committing and and unclusive stares are adding stares are add Conclusions/Significance: The TDI kernel, which is being offered under the Creative Commons attribution share-alike license for free and unrestricted use, can be accessed on the World Wide Web at http://www.tropicaldisase.org. We hope that the kernel will facilitate collaborative efforts towards the discovery of new drugs against paralites that cause tropical disease: Citation: Onl L, Carbajo RJ, Pieper U, Eswar N, Maurer SM, et al. (2009) A Kernel for Open Source Drug Discovery in Tropical Diseases. PLoS Negl Trop Dis 3(4) e118. doi:10.1371/journalpentd.0000418 Editor: Timothy G. Geary, McGill University, Canada Received December 29, 2008; Accepted March 23, 2009; Published April 21, 2009 Copyright: © 2009 Orti et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permit use, distribution, and reproduction in any medium, provided the original author and source are credited. Funding: WMAR advantedges the separat from a Sparsh Ministerio de Educación y Carricia para (BICD20710027). As advantedegas the support from testingent from Sparsh Ministerio de Carrieria a terrestario para (SUCIDAD 1994). Del Carrieria para (BICD2071027). A VI. advantedegas the support from testingent historia de Carrieria a terrestario para (SUCIDAD 1994). Del Carrieria para (BICD207102). A VI. advantedegas the support from testingent historia de Carrieria a terrestario para (SUCIDAD 1994). Del Carrieria para (BICD207102). A VI. advantedegas the support from testingent historia de Carrieria a terrestario para (SUCIDAD 1994). Del Carrieria para (BICD207102). A VI. advantedegas the support from testingent historia de Carrieria a terrestario para (SUCIDAD 1994). Del Carrieria de VI. advantedegas the support from testingent historia de Carrieria a terrestario para (SUCIDAD 1994). Del Carrieria de VI. advantedegas the support from testingent historia de Carrieria a terrestario de Carrieria a terrestario de Carrieria de VI. advantedegas the support from testingent historia de Carrieria de VI. advantedegas the support from testingent historia de Carrieria de VI. advantedegas the support from testingent historia de Carrieria de VI. advantedegas the support from testingent historia de Carrieria de VI. advantedegas the support from testingent historia de Carrieria de VI. advantedegas the support from testingent historia de Carrieria de VI. advantedegas the support from testingent historia de Carrieria de VI. advantedegas the support from testingent historia de Carrieria de VI. advantedea de Carrieria de VI. advantedea de Carrieria de VI. advantedea Competing Interests: The authors have declared that no competing interests exist. * E-mail: sali@salilab.org (AS); mmarti@cipf.es (MAM-R) R Current address: DuPont Knowledge Center, Hyderabad, India Introduction Three is the Arthliphenality previous trapers and the second seco www.nlosntds.org 1 Anril 2009 | Volume 3 | Issue 4 | e418

93

Acknowledgments

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

COMPARATIVE MODELING Andrej Sali M. S. Madhusudhan Narayanan Eswar Min-Yi Shen Ursula Pieper Ben Webb Maya Topf (Birbeck College)

MODEL ASSESSMENT David Eramian Min-Yi Shen Damien Devos

FUNCTIONAL ANNOTATION Andrea Rossi (Rinat-Pfizer) Fred Davis (Janelia Fram)

FUNDING

Prince Felipe Research Center **Ministerio de Educación y Ciencia** STREP UE Grant Marie Curie Reintegration Grant MODEL ASSESSMENT Francisco Melo (CU) Alejandro Panjkovich (CU)

NMR Antonio Pineda-Lucena Leticia Ortí Rodrigo J. Carbajo

MAMMOTH Angel R. Ortiz

FUNCTIONAL ANNOTATION Fatima Al-Shahrour Joaquin Dopazo

BIOLOGY

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



Tropical Disease Initiative Stephen Maurer (UC Berkeley) Arti Rai (Duke U) Andrej Sali (UCSF) Ginger Taylor (TSL) Matthew Todd (U Sydney)

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)