Comparative Protein Structure Prediction



Marc A. Marti-Renom





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

DISCLAIMER!

Name	Type#	World Wide Web address ^b
DATABASES		
CATH	s	http://www.blochem.ucl.ac.uk/bsm/cath/
DBAII	s	http://www.salifab.org/DBAII/
GenBank	s	http://www.ncbi.nlm.nih.gov/Genbank/GenbankSearch.htm
GeneCensus	s	http://bioinfo.mbb.yale.edu/genome
MODBASE	s	http://sailab.org/modbase/
MSD	s	http://www.rcsb.org/databases.html
NCBI	s	http://www.ncbi.nlm.nlh.gov/
PDB	s	http://www.rcsb.org/pdb/
PS1	5	http://www.nigms.nih.gov/psi/
Sacch3D	s	http://genome-www.stanford.edu/Sacch3D/
SCOP	5	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.nciferf.gov/
3D-PSSM	s	http://www.sbg.bio.ic.ac.uk/~3dpsam/
BIOINBGU	s	http://www.cs.bgu.ac.il/~bioinbgu/
BLAST	s	http://www.ncbi.nlm.nih.gov/BLAST/
DALI	s	http://www2.ebi.ac.uk/dall/
FASS	s	http://bioinformatics.burnham-inst.org/FFAS/index.html
FastA	s	http://www.ebi.ac.uk/fasta3/
FRSVR	s	http://fold.doe-mbi.ucla.edu/
PUGUE	s	http://www-cryst.bloc.cam.ac.uk/~fugue/

http://salilab.org/bioinformatics_resources.shtml

Summary

- INTRO
- Structural Space
- Profile-Profile alignment
- MOULDER
- MODELLER example

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.

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



Folding

Ab initio prediction

Evolution

Threading Comparative Modeling

Comparative Modeling by Satisfaction of Spatial Restraints (MODELLER)

3D GKITFYERGFQGHCYESDC-NLQP... SE GKITFYERG---RCYESDCPNLQP...



http://www.salilab.org/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.

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.

Typical errors in comparative models

MODEL X-RAY TEMPLATE

Region without a template



Incorrect template



Distortion/shifts in aligned regions



Misalignment



Sidechain packing



Marti-Renom et al. Annu.Rev.Biophys.Biomol.Struct. 29, 291-325, 2000.

Model Accuracy as a Function of Target-Template Sequence Identity



Sánchez, R., Šali, A. Proc Natl Acad Sci U S A. 95 pp13597-602. (1998).

Model Accuracy

HIGH ACCURACY

NM23 Seq id 77% Cα equiv 147/148 RMSD 0.41Å



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

Marti-Renom et al. Annu.Rev.Biophys.Biomol.Struct. 29, 291-325, 2000.

Classification of the structural space



SCOP_{1.71} database

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



- ✓ 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	226	392	645
All beta proteins	149	300	594
Alpha and beta proteins (a/b)	134	221	661
Alpha and beta proteins (a+b)	286	424	753
Multi-domain proteins	48	48	64
Membrane and cell surface proteins	49	90	101
Small proteins	79	114	186
Total	971	1589	3004

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

CATH_{3.1.0} 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



0	0	0	0	9	0	0	0	D
Mainly Alpha	5	305	652	1850	2329	3001	5587	1972
Mainly Beta	20	191	415	1860	2531	3846	6503	2553
Alpha Beta	14	496	922	3922	5303	6659	12998	4719
Few Secondary Structures	1	92	102	162	200	275	403	1426
Total	40	1084	2091	7794	10363	13781	25491	9388

Orengo, C.A., et al. (1997) Structure. 5. 1093-1108.

DBAliv2.0 database

http://bioinfo.cipf.es/squ/services/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 alignm	sents
Last update:	February 15th, 2007
Number of chains:	88,276
Number of structure-structure comparisons:"	1,425,479,365
Multiple structure alignme	ents
Last update:	January 23rd, 2007
Number of representative chains:	30,900
Number of families:	11,615

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

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 (build_profile & pp_scan)

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

Preparation of Sequence Database

Generation of Alignment Scores

Construction of PSSM

Position-Specific Scoring Matrix

Data-dependent Pseudocounts

Position-Based Sequence Weights

Assessment of Statistical Significance

Select Sequences Based on E-value

Create Multiple Alignment

T	AS	TRI	TKI			VENI	RDFU	01 NN	VPGAGE	l pagpfaqmil
K	İΒ		2	þ	υ	r	^	TP	ILNVIG	Y SVEEIQDIFL
N	PF.	þ.		ĽΥ	KN.	RL	22	1 2	TVGHAH	I AGSKFAPNPN
QS	-	8L ,			#SN	MRS	-	6 3 8	LRREEE	A ENDEAQXQXM
t	he	U I	niv	er	sa	I p	10	teir	re	source

1,803,406

1,774,668

LENGTH FILTER (≤30aa / ≥3000aa)

SEG FILTER

(≥40aa / ≥40% of length)

1,460,796

	90%	799,201
SEQID FILTER	80%	688,726
	70%	609,238

60% 532,251

Preparation of Sequence Database

Generation of Alignment Scores

Construction of PSSM

Position-Specific Scoring Matrix

Data-dependent Pseudocounts

Position-Based Sequence Weights

Assessment of Statistical Significance

Select Sequences Based on E-value

Create Multiple Alignment

		s	м	L	к	Р
	0	0	0	0	0	0
т	0	S11	S12	S13	S14	S15
с	0	S21	S22	S23	S24	S25
I	0	S31	S32	S33	S34	S 35
R	0	S41	S42	S43	S44	S45

Score-only Implementation of Smith-Waterman Dynamic Programing Algorithm

Miller & Myers, 1988

Preparation of Sequence Database Generation of Alignment Scores Construction of PSSM Position-Specific Scoring Matrix Data-dependent Pseudocounts

Position-Based Sequence Weights

Assessment of Statistical Significance

Select Sequences Based on E-value

Create Multiple Alignment

G ... 450 400 500 400 450 А 400 400 250 300 С 400 300 450 550 440 300 450 700 350 400 450 350 300 300 \boldsymbol{a} G 300 350 450 :750 300 300 450 450 300 н . . . $w_{ia} = \frac{1}{\lambda_{ia}} \ln \left(\frac{p_{ia}}{P_{a}} \right)$

where:

 λ_u is a scaling factor

 p_{ia} is the estimated probability of residue *a* to be found at position *i*

 P_a is the background probability of residue a

Henikoff & Henikoff, 1994

Preparation of Sequence Database Generation of Alignment Scores Construction of PSSM **Position-Specific Scoring Matrix Data-dependent Pseudo-counts Position-Based Sequence Weights Assessment of Statistical Significance Select Sequences Based on E-value Create Multiple Alignment**

$$p_{ia} = \frac{\alpha_i}{\alpha_i + \beta} f_{ia} + \frac{\beta}{\alpha_i + \beta} \sum_{b=1}^{20} f_{ib} \frac{q_{ab}}{P_b}$$

where:

 f_{ia} , f_{ib} are the observed weighted counts of residues a, b at position i

 q_{ab} are the target frequencies implicit in the substitution matrix (BLOSUM62)

$$\alpha_i = N_{diff}^i - 1 \qquad \beta = 10$$

where:

 N^{i}_{diff} is the number of different residues at *i*

Tatusov et.al., 1994; Altschul et.al., 1997

Preparation of Sequence Database Generation of Alignment Scores Construction of PSSM **Position-Specific Scoring Matrix Estimation of Target Frequencies Position-Based Sequence Weights Assessment of Statistical Significance Select Sequences Based on E-value Create Multiple Alignment**



where:

 $n^{j_{m}}$ is the number of times the residue in sequence *m* occurs in the column

> Henikoff & Henikoff, 1994; Wang & Dunbrack, 2004

Preparation of Sequence Database Generation of Alignment Scores Construction of PSSM Position-Specific Scoring Matrix Estimation of Target Frequencies Position-Based Sequence Weights Assessment of Statistical Significance Select Sequences Based on E-value

Create Multiple Alignment



Pearson, 1998

Preparation of Sequence Database

Generation of Alignment Scores

Construction of PSSM

Position-Specific Scoring Matrix

Estimation of Target Frequencies

Position-Based Sequence Weights

Assessment of Statistical Significance

Re-align Significant Alignments

Create Multiple Alignment

		S	М	L	к	P
	0	0	0	0	0	0
т	0	S 11	S12	S13	S14	S15
с	0	S21	S22	S23	S24	S25
I	0	S31	S32	s33	5 34	<mark>S</mark> 35
R	0	S41	S42	S43	S44	S45

Full Implementation of Smith-Waterman Dynamic Programing Algorithm

Gotoh, 1987

Generation of Alignment Scores Construction of PSSM Position-Specific Scoring Matrix Estimation of Target Frequencies Position-Based Sequence Weights Assessment of Statistical Significance Re-align Significant Alignments

Preparation of Sequence Database

Create Multiple Alignment

VLSEGEWQLVIWMQLC -LSEGEWQLVTFLNLC TLAEGEYQLI--LNLC T--IAADGEYNLVALC



Only 26 (out of 6600) profiles showed corruption



PP_SCAN or profile-profile alignments





% of equivalent positions

eq.-Seq

(n

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
DPнs	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
ССнз	56 ± 24	0.62 ± 0.24
ТОР	62 ± 20	0.67 ± 0.20



PP_SCAN success



Alignment accuracy (CE overlap) 200 pairwise DBAli alignments





John, Sali (2003). NAR pp31 3982

Moulding: iterative alignment, model building, model assessment



Genetic algorithm operators





...TSSONMKLGVFWGY... ...VSSCNGDLHMKVGV...



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) .

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

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

Benchmark with the "very difficult" test set

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

			Initial pr	ediction	Final p	rediction	Best pi	rediction
Target -template	Sequence identity [%]	Coverage [% aa]	Cα RMSD [Å]	CE overlap [%]	RMSD [Å]	CE overlap [%]	RMSD [A]	CE 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



4.4%



Comparative Protein Structure Prediction MODELLER tutorial

\$>mod9v1 model.py

Marc A. Marti-Renom



CRINCIPE FELIPE

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

Obtaining MODELLER and related information

MODELLER (9v1) web page

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

- Download Software (Linux/Windows/Mac/Solaris)
- HTML Manual
- ♦ Join Mailing List







Using MODELLER

No GUI! 😕

- Controlled by command file 88
- Script is written in PYTHON language ③
- You may know Python language is simple <a>©©

MODELLER 9v1 Python interface

- Modeller Python interface uses classes, e.g.:
 - 'alignment' holds and manipulates aligned sequences
 - 'model' holds and manipulates protein models
 - 'environ' keeps the configuration of the environment
 - 'profile' holds and manipulates sequence profiles
 - 'sequence_db' is for sequence databases
- These behave just like ordinary Python classes, but Modeller Fortran code is linked to them
- The Modeller data is automatically freed when the Python object is deleted (explicitly or implicitly)

MODELLER 8 class hierarchy object modobject model automodel loopmodel alignment environ 'object' is a standard Python class density · 'modobject' provides

- modobject provides basic functions for most Modeller classes
- Not all classes are shown in this diagram

Using MODELLER

INPUT:

- Target Sequence (FASTA/PIR format)
- Template Structure (PDB format)
- Python file

OUTPUT:

- Target-Template Alignment
- Model in PDB format
- Other data

Modeling of BLBP Input

Target: Brain lipid-binding protein (BLBP)
 BLBP sequence in PIR (MODELLER) format:

>P1;blbp

sequence:blbp::::::::

VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSID DRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFGDIVAVRCYEKA*

```
# Example for: alignment.align()
# This will read two sequences, align them, and write the alignment
# to a file:
loq.verbose()
env = environ()
aln = alignment(env)
mdl = model(env, file='1hms')
aln.append model(mdl, align codes='lhms')
aln.append(file='blbp.seq', align codes=('blbp'))
# The as1.sim.mat similarity matrix is used by default:
aln.align(gap penalties 1d=(-600, -400))
aln.write(file='blbp-1hms.ali', alignment format='PIR')
aln.write(file='blbp-1hms.pap', alignment format='PAP')
```

```
# Example for: alignment.align()
# This will read two sequences, align them, and write the alignment
# to a file:
log.verbose()
env = environ()
aln = alignment(env)
mdl = model(env, file='1hms')
aln.append model(mdl, align codes='1hms')
aln.append(file='blbp.seq', align_codes=('blbp'))
# The as1.sim.mat similarity matrix is used by default:
aln.align(gap penalties 1d=(-600, -400))
aln.write(file='blbp-1hms.ali', alignment format='PIR')
aln.write(file='blbp-1hms.pap', alignment format='PAP')
```

```
# Example for: alignment.align()
# This will read two sequences, align them, and write the alignment
# to a file:
loq.verbose()
env = environ()
aln = alignment(env)
mdl = model(env, file='1hms')
aln.append model(mdl, align_codes='1hms')
aln.append(file='blbp.seq', align codes=('blbp'))
# The as1.sim.mat similarity matrix is used by default:
aln.align(gap penalties 1d=(-600, -400))
aln.write(file='blbp-1hms.ali', alignment_format='PIR')
aln.write(file='blbp-1hms.pap', alignment format='PAP')
```

```
# Example for: alignment.align()
# This will read two sequences, align them, and write the alignment
# to a file:
log.verbose()
env = environ()
aln = alignment(env)
mdl = model(env, file '1hms')
aln.append model(mdl, align codes='1hms')
aln.append(file='blbp.seq', align codes=('blbp'))
# The as1.sim.mat similarity matrix is used by default:
aln.align(gap penalties 1d=(-600, -400))
aln.write(file='blbp-1hms.al,', alignment format='PIR')
aln.write(file='blbp-1hms.pap', alignment format='PAP')
```

Modeling of BLBP STEP 1: Align blbp and 1hms sequences *Output*

>P1;1hms
<pre>structureX:1hms: 1 : : 131 : :undefined:undefined:-1.00:-1.00</pre>
VDAFLGTWKLVDSKNFDDYMKSLGVGFATRQVASMTKPTTIIEKNGDILTLKTHSTFKNTEISFKLGVEFDETTA
DDRKVKSIVTLDGGKLVHLQKWDGQETTLVRELIDGKLILTLTHGTAVCTRTYEKE*
>P1;blbp
sequence:blbp: :::::::0.00:0.00
VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSI
DDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFGDIVAVRCYEKA*

Modeling of BLBP STEP 1: Align blbp and 1hms sequences *Output*

>P1;1hms
<pre>structureX:1hms: 1 : : 131 : :undefined:undefined:-1.00:-1.00</pre>
VDAFLGTWKLVDSKNFDDYMKSLGVGFATRQVASMTKPTTIIEKNGDILTLKTHSTFKNTEISFKLGVEFDETTA
DDRKVKSIVTLDGGKLVHLQKWDGQETTLVRELIDGKLILTLTHGTAVCTRTYEKE*
>P1;blbp
sequence:blbp: ::::::0.00:0.00
VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSI
DDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFGDIVAVRCYEKA*

Modeling of BLBP STEP 1: Align blbp and 1hms sequences *Output*

_aln.pos 1hms blbp _consrvd	1(VDAFLGTWKI VDAFCATWKI **** ***) LVDSKNFDD LTDSQNFDE * ** ***	20 YMKSLGVGFA YMKALGVGFA	30 ATRQVASMTKE ATRQVGNVTKE **** ***	40 PTTIIEKNGDI PTVIISQEGGR	50 LTLKTHSTFK VVIRTQCTFK * ***	60 INTEISFKLGV INTEINFQLGE
_aln.p 1hms blbp _consrvd	70 EFDETTADDF EFEETSIDDF ** ** ***	80 RKVKSIVTL RNCKSVVRL	90 DGGKLVHLQK DGDKLIHVQK ** ** * **	100 CWDGQETTLVF CWDGKETNCTF	110 ELIDGKLILT EIKDGKMVVT	120 LTHGTAVCTF LTFGDIVAVF	130 RTYEKE RCYEKA * * *

```
# Homology modelling by the automodel class
from modeller.automodel import *  # Load the automodel class
log.verbose()
                                  # request verbose output
env = environ()
                                   # create a new MODELLER environment
# directories for input atom files
env.io.atom files directory = './:../atom files'
a = automodel(env,
             alnfile = 'blbp-1hms.ali', # alignment filename
             knowns = '1hms',
                                          # codes of the templates
             sequence = 'blbp')
                                            # code of the target
a.starting model= 1
                                  # index of the first model
a.ending model = 1
                                   # index of the last model
                                   # (determines how many models to calculate)
                                    do the actual homology modelling
a.make()
```

```
# Homology modelling by the automodel class
from modeller.automodel import *  # Load the automodel class
log.verbose()
                            # request verbose output
                                  # create a new MODELLER environment
env = environ()
# directories for input atom files
env.io.atom files directory = './:../atom files'
a = automodel(env,
             alnfile = 'blbp-1hms.ali', # alignment filename
             knowns = '1hms',
                                         # codes of the templates
             sequence = 'blbp')
                                           # code of the target
a.starting model= 1
                                 # index of the first model
                                  # index of the last model
a.ending model = 1
                                  # (determines how many models to calculate)
                                   # do the actual homology modelling
a.make()
```

```
# Homology modelling by the automodel class
from modeller.automodel import *  # Load the automodel class
log.verbose()
                      # request verbose output
env = environ()
                                  # create a new MODELLER environment
# directories for input atom files
env.io.atom files directory = './:../atom files'
a = automodel(env,
             alnfile = 'blbp-1hms.ali', # alignment filename
             knowns = '1hms', # codes of the templates
              sequence = 'blbp')
                                         # code of the target
                               # index of the first model
a.starting model= 1
a.ending model = 1
                                 # index of the last model
                                  # (determines how many models to calculate)
                                   do the actual homology modelling
a.make()
```

PDB file

Can be viewed with Chimera http://www.cql.ucsf.edu/chimera/

Rasmol

http://www.openrasmol.org

PyMol

http://pymol.sourceforge.net/



Model file \rightarrow blbp.B99990001.pdb X Viewing

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



MODWEB

http://salilab.org/modweb



MODBASE

http://salilab.org/modbase

Search Page User Login ModBase Search Page ModWeb Modelling Server **Current Logins** Home Help: MOD/ **Database of Comparative Protein Structure Models** Welcome to ModBase, a database of three-dimensional protein models calculated by comparative modeling. (Old ModBase Interface) General Information Search Statistics ModBase search form **Project Pages** Search type Model Defaulty • Documentation Authors and All available datasets are selected in Select specific dataset(s) Acknowledgements Publications Todo List Search by properties **Related Resources** Property 2 ALL -Organism @ ALL • 97 Note: MODBASE contains theoretically calculated models, not experimentally Advanced search determined structures. The models may contain significant errors.

Model Details

Mod Home	User Login	ModBase Se	arch Page	ModWeb Modelling Server + Current Lo				
Sequence Information Primary Database Link I Organism II Annotation	P 43632 (KI254 Homo satients killer cell immun associated trans	HUMAN) oglobulin-like recep cript 8) inkat-flide	tor 2ds4 precurso	r (mho class ide nk cell receptor) (na cell receptor clore cl-39) (p58 nk	tural killer			
Sequence Length	304		10.00000000000					
Model Information								
Perform action on this		equence Model Cr equence Identify -Value Kodel Score arget Region trotein Length emplate PDB Jode	werage werage werag					
8	T	emplate Region lataset	6-200 snp-human2					
Filtered models for our	ment sequence (S	how all models)						
1 in 1								
Cross-references					_			

Sequence Overview

SegId Fold MScore		hypothetical protein	Pseudomonas aeruginosa	3738
Seg Id Fold HScore	□ <u> Q8G9W1</u>	hypothetical protein	Escherichia coli	1140
SegId Fold MScore	□ <u>Q8CY62</u>	hypothetical protein spr1965	Streptococcus pneumoniae, Streptococcus pneumoniae	1038

Model Overview

£96	• Г	Q8G8C7	hypothetical protein	<u>Pseudomonas</u> aeruginosa	4996	2089-2158	70	37.00	7e-14	1.00	1dnyA	8-78
教育	• -	Q8G8C7	hypothetical protein	<u>Pseudomonas</u> aeruginosa	4996	492-1017	526	36.00	1e-82	1.00	<u>1amuA</u>	19-529
	• □	Q8G9W1	hypothetical protein	Escherichia coli	1140	349-1135	787	35.00	0	1.00	1r9dA	6-783

Utility of protein structure models, despite errors



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

FUNDING

Prince Felipe Research Center Marie Curie Reintegration Grant STREP EU Grant MODEL ASSESSMENT Francisco Melo (CU) Alejandro Panjkovich (CU)

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

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)

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 L Alfonso Valencia (CNB/LIAI

CAMP

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

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