# Master Bioinformatics for Health Sciences
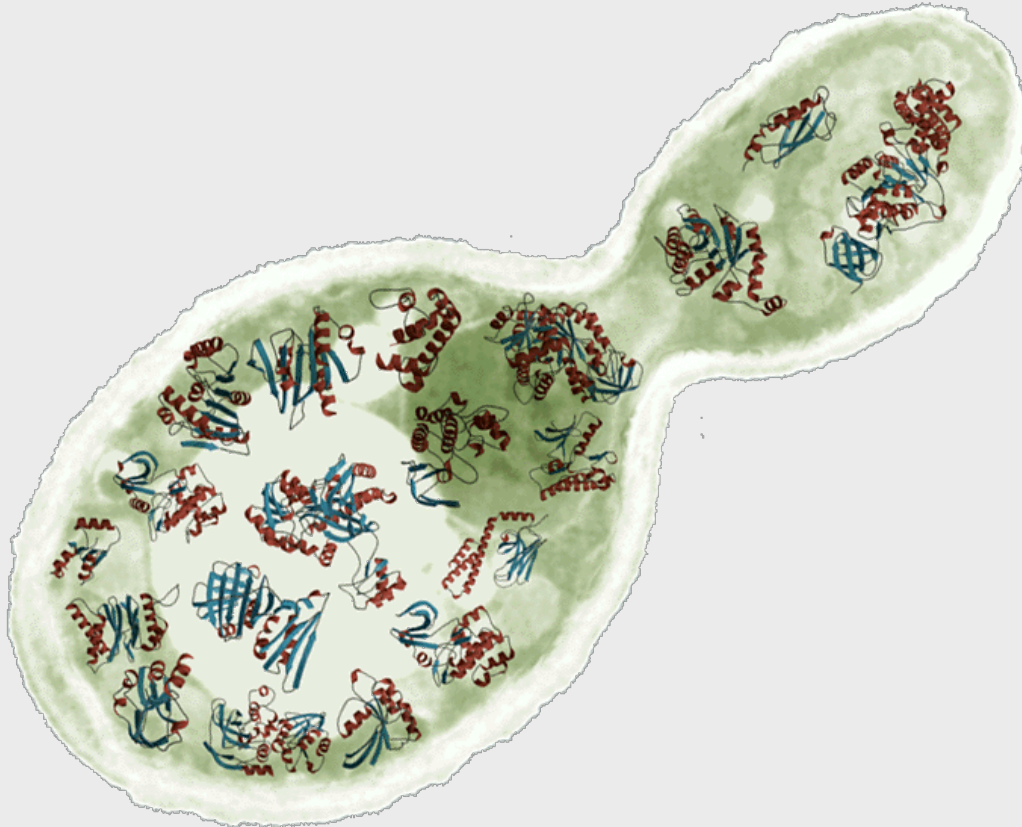## Comparative Protein Structure Prediction



**Marc A. Marti-Renom**

*Adjunct Assistant Professor*

http://salilab.org/~marcius

Depts. of Biopharmaceutical Sciences and Pharmaceutical Chemistry
California Institute for Quantitative Biomedical Research
University of California at San Francisco

UCSF
University of California
San Francisco

# DISCLAIMER!



| Name | Type[a] | World Wide Web address[b] |
|------|---------|---------------------------|
| **DATABASES** | | |
| CATH | S | http://www.biochem.ucl.ac.uk/bsm/cath/ |
| DBAli | S | http://www.salilab.org/DBAli/ |
| 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.nlm.nih.gov/ |
| PDB | S | http://www.rcsb.org/pdb/ |
| PSI | S | http://www.nigms.nih.gov/psi/ |
| 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 ASSIGNMENT** | | |
| 123D | S | http://123d.ncifcrf.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.nlm.nih.gov/BLAST/ |
| DALI | S | http://www2.ebi.ac.uk/dali/ |
| 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/ |
| FUGUE | S | http://www-cryst.bioc.cam.ac.uk/~fugue/ |
| LOOPP | S | http://ser-loopp.tc.cornell.edu/cbsu/loopp.htm |
| PDB-BLast/FASS | S | http://bioinformatics.ljcrf.edu/pdb_blast/ |
| PHD, TOPITS | S | http://www.predictprotein.org/ |

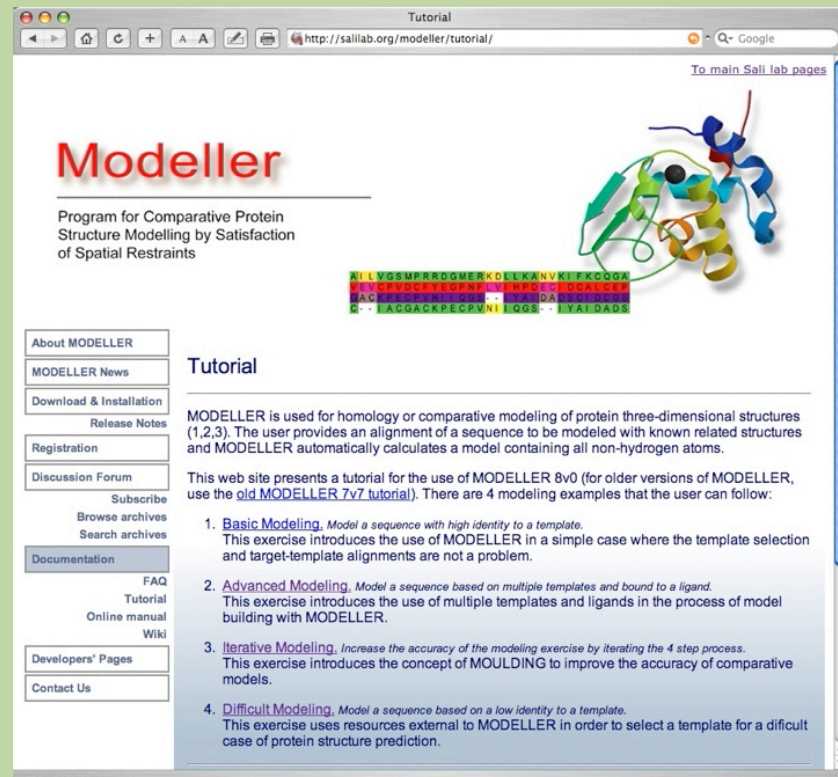http://salilab.org/bioinformatics_resources.shtml

2

# 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

# What are we going to do?

Ask!

Each day…

Basic introduction

Theory (representation-scoring-optimization)

Available programs

Application

# 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

**X-Ray**

**NMR**

**Comparative Modeling**

**Threading**

**Ab-initio**

Experimental data

inferred data

# Why protein structure prediction?

|  | Y 2005 | Y 2006 |
|---|---|---|
| Sequences | 1,700,000 | millions |
| Structures | 28,000 | 50,000 |

# Why protein structure prediction?

| | Y 2005 |
|---|---|
| Sequences | 1,700,000 |
| Structures | **900,000** |

http://salilab.org/modbase/

Theory

Experiment

# 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…



Desulfovibrio vulgaris

Condrus crispus

Anabaena 7120

Anacystis nidulans

GFCHIKAYTRLIMVG…

## Folding

Ab initio prediction

## Evolution

Threading
Comparative Modeling

# Steps in Comparative Protein Structure Modeling

START

Template Search

Target – Template Alignment

Model Building

Model Evaluation

OK?

No

Yes

END

TARGET

TEMPLATE

ASILPKRLFGNCEQTSDEGL
KIERTPLVPHISAQNVCLKID
DVPERLIPERASFQWMNDK

ASILPKRLFGNCEQTSDEGLKIERTPLVPHISAQNVCLKIDDVPERLIPE
MSVIPKRLYGNCEQTSEEAIRIEDSPIV---TADLVCLKIDEIPERLVGE

blbp.B99990001

PSEUDO ENERGY

0.8
0.0
-0.8
-1.6
-2.4

0    20    40    60    80    100   120
RESIDUE INDEX

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.

12

# Utility of protein structure models, despite errors



D. Baker & A. Sali. Science 294, 93, 2001.

13

# General References

**Protein Structure Prediction:**

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

**Comparative Modeling:**

Marti-Renom el al. Annu. Rev. Biophys. Biomol. Struct. 29, 291-325, 2000.
Marti-Renom el al. Current Protocols in Protein Science 1, 2.9.1-2.9.22, 2002.

**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.
Vitkup et al. Nat. Struct. Biol. 8, 559, 2001.

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

# Programs, servers and databases
## http://salilab.org



**LS-SNP**
Web Server
http://salilab.org/LS-SNP
Predicts functional impact
of residue substitution

**PIBASE**
Database
http://salilab.org/pibase
Contains structurally defined
protein interfaces

**CCPR**
Center for Computational
Proteomics Research
http://www.ccpr.ucsf.edu

**MODLOOP**
Web Server
http://salilab.org/modloop
Models loops in protein
structures

**MODBASE**
Database
http://salilab.org/modbase
Fold assignments,alignments
models, model assessments
for all sequences related to a
known structure

**MODWEB**
Web Server
http://salilab.org/modweb
Provides a web interface to
MODPIPE

**MODELLER**
Program
http://salilab.org/modeller
Implements most operations
in comparative modeling

**DBALI**
Database
http://salilab.org/dbali
Contains a comprehensive
set of pairwise and multiple
structure-based alignments

**ICEDB**
Database/LIMS
http://nysgxrc.org
Tracks targets for structural
genomics by NYSGXRC

**MODPIPE**
Program
Automatically calculates
comparative models of many
protein sequences

**EVA**
Web Server
http://salilab.org/eva
Evaluates and ranks web
servers for protein structure
prediction

**LIGBASE**
Database
Ligand binding sites and
inheritance (accessible
through MODBASE)

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

16

# Acknowledgments

**COMPARATIVE MODELING**
**Andrej Sali**
M. S. Madhusudhan
Narayanan Eswar
David Eramian
Ursula Pieper
Ben Webb
Min-Yi Shen
Mark Peterson
Ash Stuart
Andras Fiser (AECOM)
Roberto Sanchez (MSSM)
Bino John (Pitsburg U.)
Eric Feyfant (GI)

**ASSEMBLIES**
Frank Alber
Damien Devos
Maya Topf
Dmitry Korkin
Fred Davis
Mike Kim

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

**BRCA1**
Alvaro Monteiro (Cornell)

**MODEL ASSESSMENT**
Francisco Melo (CU, Chile)
Alejandro Panjkovich (CU, Chile)

**FUNCTIONAL ANNOTATION**
Andrea Rossi
Rachel Karchin
Libusha Kelly
Nebojša Mirkovic

17

**Master Bioinformatics for Health Sciences**

# Comparative Protein Structure Prediction
## template selection & sequence-structure alignment*

**Marc A. Marti-Renom**
*Adjunct Assistant Professor*
http://salilab.org/~marcius

Depts. of Biopharmaceutical Sciences and Pharmaceutical Chemistry
California Institute for Quantitative Biomedical Research
University of California at San Francisco

UCSF
University of California
San Francisco

# Summary

◇ **Sequence and structure space (domains)**

◇ Domains from sequence

◇ Structure-Structure comparisons

◇ How can we compare structures

◇ How we classify the structural space

◇ Aligning sequences and structures

# Domain boundaries from sequence
## VERY DIFFICULT!!!!



MENFEIWVEKYRPRTLDEVVGQDEVIQRLKGYVERKNIPHLLFSGPPGTGKTATAIALARDLFGENWRDN
FIEMNASDERGIDVVRHKIKEFARTAPIGGAPFKIIFLDEADALTADAQAALRRTMEMYSKSCRFILSCN
YVSRIIEPIQSRCAVFRFKPVPKEAMKKRLLEICEKEGVKITEDGLEALIYISGGDFRKAINALQGAAAI
GEVVDADTIYQITATARPEEMTELIQTALKGNFMEARELLDRLMVEYGMSGEDIVAQLFREIISMPIKDS
LKVQLIDKLGEVDFRLTEGANERIQLDAYLAYLSTLAKK

# Domain boundaries from sequence (SnapDragon)



**Table 2.** Average accuracy percentages of linker prediction over 57 proteins

|  |  | Continuous set | Discontinuous set | Full set |
|---|---|---|---|---|
| Randomised background Z-score >2 | Coverage | 63.3 | 43.6 | 54.8 |
|  | Success | 27.2 | 31.1 | 28.9 |
| Self-normalised Z-score >1 | Coverage | 64.7 | 39.5 | 53.5 |
|  | Success | 26.6 | 31.7 | 28.9 |
| Self-normalised Z-score >2 | Coverage | 48.7 | 24.3 | 38.7 |
|  | Success | 41.3 | 28.3 | 29.9 |

# Domain boundaries from sequence and predicted SSE (DomSSEA)

|                                       | % Correctly assigned | |
|---------------------------------------|:-----------:|:----------:|
| Methods                               | All chains  | Multidomain chains |
| DomSSEA observed secondary structure  | 70.2        | 24.7       |
| DomSSEA predicted & consensus         | 68.6        | 24.0       |
| DomSSEA predicted & L/(N-1)           | 68.0        | 24.0       |
| DomSSEA predicted secondary structure | 68.7        | 23.6       |
| Absolute difference in length         | 62.0        | 8.4        |
| Average domain length & DGS-M         | 66.6        | 6.1        |
| FASTA alignment                       | 57.9        | 2.3        |
| Random (weighted)                     | 58.3        | 1.1        |
| DGS-M                                 | 76.6        | 0.0        |
| DGS-W                                 | 76.6        | 0.0        |

*Dersden et al. (2003) Prot. Science 11 pp2014*

# Prediction of Secondary Structure (PSI-PRED)

```
>gi42541361
MDIRSVSSLRGLLCLPPSWPRR
```

• Neural Network



✓ **Very simple idea**
✓ **Simple scoring**

**Obscure optimizer**



Raw profile from PSI-BLAST Log File

Window of 15 rows

15 x 20 scaled inputs to 1st network

1st Network
315 inputs
75 hidden units
3 outputs

Window of 15 x 3 outputs fed to 2nd network

2nd Network
60 inputs
60 hidden units
3 outputs

Final 3-state Prediction

*Jones DT. (1999) J. Mol. Biol. 292 pp195*

# Prediction of Secondary Structure (PSI-PRED)

**http://bioinf.cs.ucl.ac.uk/psipred/**

# **Template Selection**
## "Structural Space"

# Structure-Structure alignments

As any other bioinformatics problem…

- Representation
- Scoring
- Optimizer

# Structures



All atoms and coordinates

Dihedral space or distance space

Reduced atom representation

Vector representation

Secondary Structure

Accessible surface (and others)

# Raw scores



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

# 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.

Empirical



Sometimes approximated by Z-score (normal distribution).

Analytic

$$P(s) = e^{-\lambda \, (s-\mu)}$$

$$P(s \geq x) = 1 - \exp\left( e^{-\lambda \, (s-\mu)} \right)$$

*Karlin and Altschul, 1990 PNAS 87, pp2264*

# Global dynamic programming alignment



$$D_{i,j} = \min \begin{cases} D_{i,j-1} + Score_{(\Delta, rj)} \\ D_{i-1,j-1} + Score_{(ri, rj)} \\ D_{i-1,j} + Score_{(ri\Delta)} \end{cases}$$

Best alignment score

## Backtracking to get the best alignment

*Needleman and Wunsch (1970) J. Mol Biol, 3 pp443*

# Local dynamic programming alignment



$$D_{i,j}=\min \begin{cases} D_{i,j-1}+Score_{(\Delta, rj)} \\ D_{i-1,j-1}+Score_{(ri,rj)} \\ D_{i-1,j}+Score_{(ri,\Delta)} \\ 0 \end{cases}$$

Best score

Best local alignment

## Backtracking to get the best alignment

# Global .vs. local alignment



Global alignment

Local alignment

# Multiple alignment

## Pairwise alignments

Example – 4 sequences A, B, C, D.

- similarity +

6 pairwise comparisons
then cluster analysis

## Multiple alignments

Following the tree from step 1

Align the most similar pair

Align next most similar pair

Align B-D with A-C

New gap in A-C to optimize
its alignment with B-D

# Coverage .vs. Accuracy



Same RMSD ~ 2.5Å

Coverage ~90% Cα                    Coverage ~75% Cα

# Structural alignment by properties conservation (SALIGN-MODELLER)



✓ Uses all available structural information
  ✓ Provides the optimal alignment

Computationally expensive

$R_{i,j}$   $D_{,i(3),j(3)}$   $S_{i,j}$   $B_{i,j}$   $I_{i,j}$

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

*Madhusudhan et al. in preparation*

# Structural alignment by properties conservation (SALIGN-MODELLER)

**http://salilab.org/DBAli**
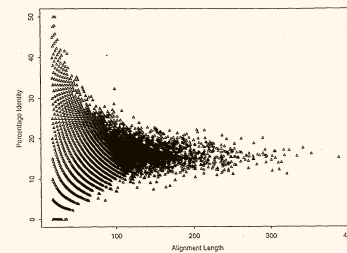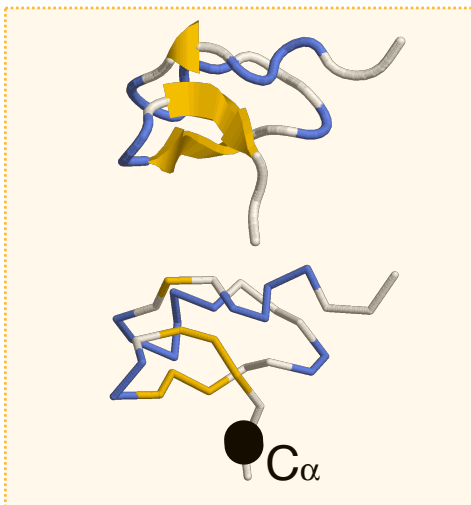


*Madhusudhan, in preparation*
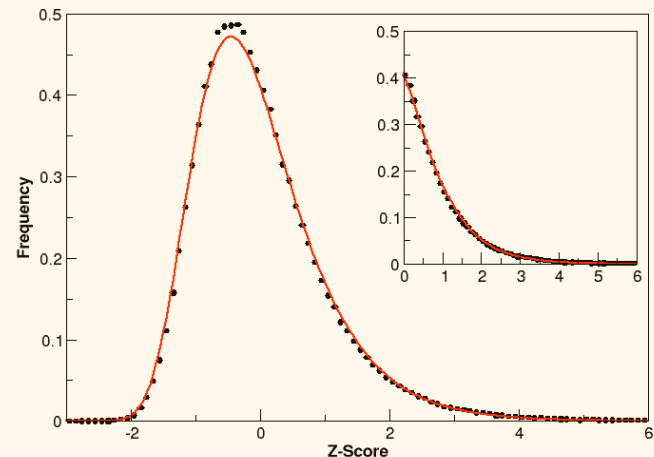
# Vector Alignment Search Tool (VAST)



Graph theory search
of similar SSE
Refining by Monte Carlo
at all atom resolution



✓ Good scoring system with significance

Reduces the protein representation



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

*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)

Exhaustive combination
of fragments

Longest combination of
AFPs

Heuristic similar to
PSI-BLAST

✓ FAST!
✓ Good quality of local alignments

Complicated scoring and heuristics

8 residues peptides

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

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

# Incremental combinatorial extension (CE)

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

# Matching molecular models obtained from theory (MAMMOTH)
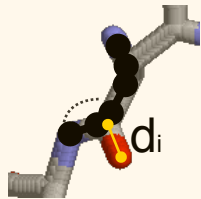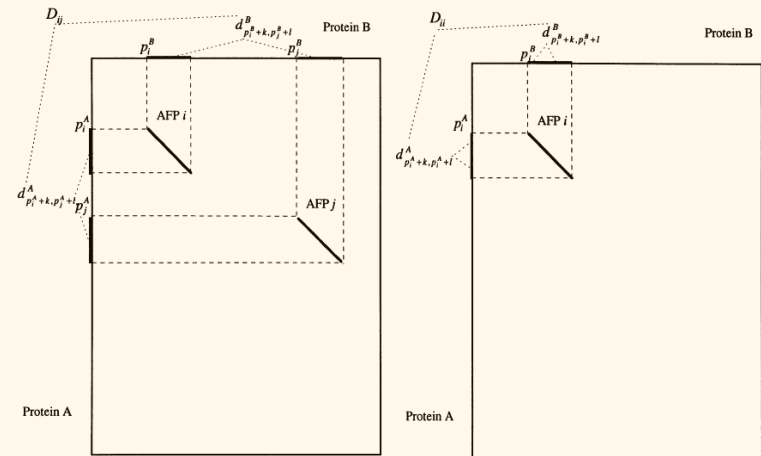


✓ VERY FAST!
✓ Good scoring system with significance

Reduces the protein representation

$$URMS^R = \sqrt{2.0 - \frac{2.84}{\sqrt{n}}}$$

$$S_{AB} = \frac{\left(URMS^R - URMS^{AB}\right)D}{URMS^R}$$

*Ortiz AR, (2002) Protein Sci. 11 pp2606*

# Matching molecular models obtained from theory (MAMMOTH)

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

# Classification of the structural space

**SCOP classification**



**Large Graph Layout**

*Adai AT, Date SV, Wieland S, Marcotte EM. J Mol Biol. 2004 Jun 25;340(1):179-90*

# SCOP₁.₆₅ 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 | 179 | 299 | 480 |
| All beta proteins | 126 | 248 | 462 |
| Alpha and beta proteins (a/b) | 121 | 199 | 542 |
| Alpha and beta proteins (a+b) | 234 | 349 | 567 |
| Multi-domain proteins | 38 | 38 | 53 |
| Membrane and cell surface proteins | 36 | 66 | 73 |
| Small proteins | 66 | 95 | 150 |
| Total | 800 | 1294 | 2327 |

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

# CATH₂.₆.₀ database

`http://www.biochem.ucl.ac.uk/bsm/cath/`



## 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**



| | Version | 2.6.0 | | | | | |
|---|---|---|---|---|---|---|---|
| | Date | 11-04-2005 | | | | | |

| | A | T | H | S | N | O | D |
|---|---|---|---|---|---|---|---|
| Mainly Alpha | 5 | 251 | 465 | 1402 | 2189 | 3705 | 14105 |
| Mainly Beta | 19 | 160 | 311 | 1443 | 2961 | 4329 | 18771 |
| Alpha Beta | 14 | 414 | 706 | 3014 | 4781 | 7660 | 33080 |
| Few Secondary Structures | 1 | 82 | 90 | 144 | 232 | 285 | 1098 |
| Preliminary single domain assigments | 10 | 808 | 809 | 906 | 967 | 1090 | 3012 |
| Multi-domain domains | 1 | 12 | 12 | 16 | 25 | 36 | 109 |
| CATH-35 Sequence families | 1 | 4707 | 4707 | 4719 | 4768 | 4862 | 6168 |
| | 1 | 22 | 22 | 27 | 33 | 38 | 198 |

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

# DBAli<sub>v2.0</sub> database
## http://salilab.org/DBAli/



## 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 (yet) a stable classification**

| Pairwise structure alignments | |
|---|---|
| Last update: | June 6th, 2005 |
| Number of chains: | 65,286 |
| Number of structure-structure comparisons:* | 791,171,210 |
| **Multiple structure alignments** | |
| Last update: | May 14th, 2005 |
| Number of representative chains: | 22,324 |
| Number of families: | 8,737 |

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

46

# Classification of the structural space
## *Not an easy task!*

Domain definition AND domain classification

47

# Application (ModDom)



**Assigning domains from structure**

Less significant

{1,2}{3,4}{4,5}
{5,6}{6,7}{7,8}{8,9}

Lower MAMMOTH P-value

{1,2,3,4}

{6,7,8,9}

{5,6,7,8,9}

{3,4,5,6,7,8,9}

Upper MAMMOTH P-value

{all}

More significant

{1,2,3,4}

{6,7,8,9}

{5,6,7,8,9}

{3,4,5,6,7,8,9}

| # | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 3 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 4 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

53

| # | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 3 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 4 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |
| 6 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 2 |
| 7 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 2 |
| 8 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 2 |
| 9 | 0 | 0 | 0 | 0 | 1 | 2 | 2 | 2 | 2 |

55

| # | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 3 | 1 | 1 | 2 | 2 | 1 | 1 | 1 | 1 | 1 |
| 4 | 1 | 1 | 2 | 2 | 1 | 1 | 1 | 1 | 1 |
| 5 | 0 | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 2 |
| 6 | 0 | 0 | 1 | 1 | 2 | 3 | 3 | 3 | 3 |
| 7 | 0 | 0 | 1 | 1 | 2 | 3 | 3 | 3 | 3 |
| 8 | 0 | 0 | 1 | 1 | 2 | 3 | 3 | 3 | 3 |
| 9 | 0 | 0 | 1 | 1 | 2 | 3 | 3 | 3 | 3 |

Threshold #3 MCL Cluster level (-I)

Stijn van Dongen (http://micans.org/mcl/)

| 1phh | 290-329 | | 2.7Å | 3.1 |
|------|---------|--|------|-----|
| 1hadB | 72-111 | | | |



| 1phh | 279-373 | | 3.9Å | 4.7 |
|------|---------|--|------|-----|
| 1bke | 310-410 | | | |





Conservation

Residue number

1phh (Oxydoreductase from Pseudomonas fluorescens)

| 1phh | 1-213 | 3.0Å | 8.1 |
|------|-------|------|-----|
| 1qjdA | 125-379 |  | |

| 1phh | 1-319 | 3.6Å | 9.8 |
|------|-------|------|-----|
| 1gerA | 3-327 |  | |



Conservation

Residue number

1phh (Oxydoreductase from Pseudomonas fluorescens)

| 1phh | 1-378 | | 3.8Å | 10.3 |
|------|-------|---|------|------|
| 1feaC | 2-464 | | | |

| 1phh | 1-316 | | 3.8Å | 17.2 |
|------|-------|---|------|------|
| 1l9dB | 2-364 | | | |



Conservation

Residue number

**1phh** (Oxydoreductase from Pseudomonas fluorescens)

# 1phh (Oxydoreductase from Pseudomonas fluorescens)

# Benchmark dataset

2163 chains from Islam et al. 1995 → 569 Non-redundant

<2Å && <30aa diff.

Divide randomly into two sets

Remove of incomplete or obsolete entries.

FINAL:

Training set → 242 chains

Testing set → 234 chains

# Scoring function    R = Volume/ASA



Domain → max(<dist f(R)>)

-0.11    -0.10    -0.08    -0.09

5-46

47-84

85-192

193-239

1-84

85-239

1-239

1dhr_ (dihydropteridine reductase )

65

# Domain assignment for a non-redundant set of 234 protein chains



DP (233)
DAD (234)
DALI (229)
CATH (224)
PUU (233)
3DEE (225)
DOMAK (231)
DIAL (233)
X-Ray (234)
PDP (234)
ModDom (234)
SCOP (228)

% correct predictions

100 — 80 — 60 — 40 — 20 — 0

X-Ray  CATH  DALI  3Dee  DAD  DIAL  DOMAK  DP  ModDom  PDP  PUU

- All chains (228)
- 1 domain (157)
- 2 domains (60)
- 3 domains (9)
- 4 domains (2)

# What are domains?



# Structural recurrent fragments

# G-protein (1gotB) *all-β* → *7 bladed beta propeller domain*

# Ribosomal protein S6 (1ris) α+β → *Ferrodoxin Like domain*



**1ee9A** 17.9% id. 2.3Å

**6timB** 11.1% id. 2.6Å

Legend (contour values):
- 0.6-0.7
- 0.5-0.6
- 0.4-0.5
- 0.3-0.4
- 0.2-0.3
- 0.1-0.2
- 0-0.1

# Cytochrome C Peroxidase (2cyp) *all-α → CCP-like domain*



27/35  3.2Å

29/34  2.9Å

# Barnase Domain-Swapping



Barnase (1brn:L)
conservation profile

# chymotrypsin inhibitor 2



1-37 | 38-64          1-40 | 41-64

*Neira JL, Davis B, Ladurner AG, Buckle AM, Gay GP, Fersht AR. 1996. Fold Des 1:189-208.*
*Ladurner AG, Itzhaki LS, de Prat GG, Fersht AR. 1997. J Mol Biol 273:317-329.*

# Sequence space .vs. Structure space

**The PDB is a covering set of small protein structures.**

# Sequence space .vs. Structure space

## Structure map @ >20% sequence identity



~22,000 nodes
186 clusters
Larger cluster contains 20,231 chains

# Sequence space .vs. Structure space

# SALIGN

## aligning profiles (PP_SCAN in MODELLER)

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

A) Training Set　　B) Testing Set

# SALIGN 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)

# SALIGN protocols accuracy

| SALIGN protocol | CE overlap [%] | Shift score |
|:---:|:---:|:---:|
| CC$_{PBP}$ | 55 ± 23 | 0.61 ± 0.24 |
| CC$_{HH}$ | **56 ± 23** | **0.61 ± 0.24** |
| CC$_{HS}$ | **56 ± 24** | **0.62 ± 0.23** |
| CC$_{MAT}$ | 51 ± 25 | 0.55 ± 0.27 |
| ED$_{PBP}$ | 54 ± 24 | 0.60 ± 0.25 |
| ED$_{HH}$ | 54 ± 24 | 0.59 ± 0.26 |
| ED$_{HS}$ | 55 ± 24 | 0.59 ± 0.26 |
| DP$_{PBP}$ | 55 ± 23 | 0.61 ± 0.24 |
| DP$_{HH}$ | 56 ± 23 | 0.60 ± 0.25 |
| DP$_{HS}$ | 55 ± 24 | 0.61 ± 0.24 |
| JS$_{HH}$ | 53 ± 24 | 0.60 ± 0.24 |
| JS$_{HS}$ | 54 ± 24 | 0.60 ± 0.24 |
| Ave$_{MAT}$ | 49 ± 26 | 0.52 ± 0.29 |
| TOP | 62 ± 20 | 0.67 ± 0.20 |

# SALIGN 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 |
| $CC_{HH}$ | 56 ± 23 | 0.61 ± 0.24 |
| $CC_{HS}$ | 56 ± 24 | 0.62 ± 0.24 |
| TOP | 62 ± 20 | 0.67 ± 0.20 |

# SALIGN success

# Alignment accuracy (CE overlap)

*200 pairwise DBAli alignments*

PSI-BLAST (sequence-profile alignment)    43%

SEA (local structure alignment)    49%

SALIGN (profile-profile alignment)    56%

# Program

Intro to comparative
protein structure prediction

Template Search*

Target – Template
Alignment*

Model Building

Model Evaluation



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

# Protein Structure Prediction
## model building & model assessment

**Marc A. Marti-Renom**

*Adjunct Assistant Professor*

http://salilab.org/~marcius

Depts. of Biopharmaceutical Sciences and Pharmaceutical Chemistry

California Institute for Quantitative Biomedical Research

University of California at San Francisco

UCSF
University of California
San Francisco

# Summary

◈ **Model building with MODELLER**

◇ Points and restraints

◇ Model accuracy

◇ Modeling loops

◇ Evaluating models

◇ MOULDER

◇ Modeing genes (examples)

◇ Modeling genomes (large-scale modeling)

# Information about a protein can come from three distinct sources



Experimental observations



Statistical rules



Laws of physics

# Modeling by optimization

**There is nothing but points and restraints on them.**

P(r/I) feature

P(R/I) molecule

# 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.

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

# Comparative modeling by satisfaction of spatial restraints
## MODELLER

```
3D   GKITFYERGFQGHCYESDC-NLQP…
SEQ  GKITFYERG---RCYESDCPNLQP…
```

1. Extract spatial restraints

2. Satisfy spatial restraints

$$F(R) = \prod_i p_i (f_i / I)$$

*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.*

# Accuracy and applicability of comparative models

# "Biological" significance of modeling errors



**NMR – X-RAY**
Erabutoxin 3ebx
Erabutoxin 1era

**NMR**
Ileal lipid-binding protein
1eal

**CRABPII** 1opbB
**FABP** 1ftpA
**ALBP** 1lib
40% seq. id.

**X-RAY**
Interleukin 1β 41bi (2.9Å)
Interleukin 1β 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**

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

9

# Typical errors in comparative models

MODEL
X-RAY
TEMPLATE

Incorrect template

Misalignment

Region without a template

Distortion/shifts in aligned regions

Sidechain packing



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

10

# Utility of protein structure models, despite errors

D. Baker & A. Sali. Science 294, 93, 2001.

# Modeling of loops in protein structures

# Loop Modeling in Protein Structures



α+β barrel: flavodoxin



IG fold: immunoglobulin



antiparallel β-barrel

13

# Loop modeling strategies



- database is complete only up to 4-6 residues
- even in DB search, the different conformations must be ranked
- loops longer than 4 residues need extensive optimization
- DB method is efficient for specific families (eg, canonical loops in Ig's, β–hairpins)

# Loop Modeling by Conformational Search



1. Protein representation.

2. Energy (scoring) function.

3. Optimization algorithm.

# Energy Function for Loop Modeling

The energy function is a sum of many terms:

1. Stereochemistry (CHARMM).

2. Mainchain conformation ($\Phi$, $\Psi$).

3. Non-bonded contacts.

# Energy Function for Loop Modeling

1) Statistical preferences for dihedral angles:



2) Restraints from the CHARMM-22 force field:



3) Statistical potential for non-bonded contacts:

# Mainchain Terms for Loop Modeling

# Optimization of Objective Function

- Test set: 40 randomly selected loops of known structures, for each length from 1 to 14 residues.

- Starting conformation: Loop atoms were spaced evenly on a line spanning the two anchor regions, then randomized by ± 5 Å.

- To simulate real comparative modeling situations, performance of the loop modeling problem was determined by predicting loops in only approximately correct environment.



Native loop region (99-106) in 1nba

Starting loop conformation

Distorted environment

X-ray structure of 1nba

# Optimization of Objective Function

# Calculating an Ensemble of Loop Models



5p21 4552, 0.25Å

(a)

Objective functon

RMSD$_{mnch,local}$

1alc, 3441, 3.17Å

(b)

Objective functon

RMSD$_{mnch,local}$

# Accuracy of loop models as a function of amount of optimization



8 residue long loops

Legend:
- Global fit, $E_{min}$
- Local fit, $E_{min}$
- Global fit, $RMSD_{min}$
- Local fit, $RMSD_{min}$

Y-axis: $<RMSD_{loop,mnch}>$ [Å]

X-axis: Number of independent optimizations

# Accuracy of Loop Modeling



RMSD=0.6Å                    RMSD=1.1Å                    RMSD=2.8Å

HIGH ACCURACY (<1Å)     MEDIUM ACCURACY (<2Å)     LOW ACCURACY (>2Å)

50% (30%) of             40% (48%) of             10% (22%) of
8-residue loops          8-residue loops          8-residue loops

*A. Fiser, R. Do & A. Šali, Prot. Sci. 9, 1753, 2000.*

# Fraction of Loops Modeled With at Least Medium Accuracy

# Problems in Practical Loop Modeling

1. Decide which regions to model as loops.
2. Correct alignment of anchor regions & environment.
3. Modeling of a loop.



T0058: 80-85
$RMSD_{mnch}$ loop = 1.09 Å
$RMSD_{mnch}$ anchors = 0.29 Å

T0076: 46-53
$RMSD_{mnch}$ loop = 1.37 Å
$RMSD_{mnch}$ anchors = 1.52 Å

# Potentials of Mean Force (PMF)

# Empirical energy functions (PMF)

Idea: **energy leads to structure, thus it should be possible to infer energy from many known structures**

To be used in: **model refinement and assessment**

Properties needed:

Deep minimum at correct state (native)
Smooth (energy landscape)
Simple (CPU calculation)

Types:

Contact potential
Distance potentials
Surface potentials

# Approximations/Limitations in PMFs

**Database size.**

**PMF versus Energy (additive/higher order terms).**

**Reference state.**

**Physical origin.**

# Potentials of Mean Force

**As any other bioinformatics problem…**

- **Representation**
- **Scoring**
- **Optimizer**

# Sequence/Structures

>gi42541361
MDIRSVSSLRGLLCLPPSWPRR

Primary sequence

All atoms and coordinates

$d_i$

Distance space

$\bullet C_\alpha$

Reduced atoms representation

Secondary Structure

Accessible surface

# Statistical Potentials (background)

Structural space



Sequence space

MKLLIVLTCISLCSCICTVVQRCASNKPHVLEDPCKVQH

HLSVNQCVLLPQCCPKSCKICTHLISIEVVLTCRAVDKM

MHVNCVEQCSLQDCIKIAPRVLKTCILCVLKPCLTSVSH

VHLVQPTSCCCKKNCICHVEIRSLDILTKSVQLACLVPM

⋮

MQCCRVQKICDLLAVELCKLHISTPSCKILCVVTSVPHN

# Statistical Potential (inspiration)

$$K = \frac{[AB]}{[A] \cdot [B]}$$

$$\Delta G = -RT \ln(K) = -RT \ln \frac{[AB]}{[A] \cdot [B]}$$

A + B ⇌ AB

From statistical physics, we know that energy difference between two states ($\Delta E$) and the ratio of their occupancies ($N_1 : N_2$) are related [9]:

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

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.

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

# Scoring
# Statistical Potential (reference state)



liquid Ar
85K

*Theory of simple liquids 2nd edition JP Hansen and IR McDonald, Academic Press.*

# Statistical Potential… Hydrogen Bonds

Long range free energy



Short range free energy



Free energy of the protein backbone hydrogen bond N · · · O compiled from a database of 289 X-ray structures

$$\rho_{NO}(r) = \sum_{ij} \delta(r - r_{ij})$$

$$g_{NO}(r) = \frac{\rho_{NO}(r)}{\rho^2}$$

$$w_{NO}(r) = -kT \ln\left(g_{NO}(r)\right)$$

# Statistical Potential… Distance Potentials

## Long range free energy



## Short range free energy

36

# Significance of an alignment (score)

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

$$Zscore = \frac{\left(\langle E \rangle - E_m\right)}{\sigma_E}$$

# ProsaII

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

## Deriving

Structural space



## Scoring



$-E$   $+E$

$\sigma_E$

$<E>$

$$Zscore = \frac{\left(\langle E \rangle - E_m\right)}{\sigma_E}$$

# ANOLEA

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

## Deriving

Structural space



## Scoring



$$Zscore = \frac{\left(\langle E \rangle - E_m\right)}{\sigma_E}$$

all atom potential

# Verify3D

## Deriving

Structural space



## Scoring

# DFIRE

## Deriving

Structural space



## Scoring

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

41

# DOPE (MODELLER)

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



## Deriving

Structural space

## Scoring

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

MOULDER

# Moulding: iterative alignment, model building, model assessment

# Iterative process… MOULDER

45

# Genetic algorithm operators

## Single point cross-over

...TSSQ—NMKLGVFWGY——...
...V—SSCN——GDLHMKVGV...

...TSSQNMK——LGVFWGY...
...VSSCNGDLHMKV——GV...

→

...TSSQ—NMK——LGVFWGY...
...V—SSCNGDLHMKV——GV...

...TSSQNMKLGVFWGY——...
...VSSCN——GDLHMKVGV...

## Gap insertion

...TSSQNMKLGVFWGY...
...VSSCNGDLHMKVGV...

→

...TSSQN——MKLGVFWGY...
...VSSCNGDLHMKVG——V...

## Gap shift

...T——SSQNMKLGVFWGY...
...VSSCNGDLHMKVGV——...

→

...—T—SSQNMKLGVFWGY...
...VSSCNGDLHMKVGV——...

...T—S—SQNMKLGVFWGY...
...VSSCNGDLHMKVGV——...

...——TSSQNMKLGVFWGY...
...VSSCNGDLHMKVGV——...

...TS——SQNMKLGVFWGY...
...VSSCNGDLHMKVGV——...

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

# Composite model assessment score

Weighted linear combination of several scores:

- Pair ($P_p$) and surface ($P_s$) statistical potentials;

- Structural compactness ($S_c$);

- Harmonic average distance score ($H_a$);

- Alignment score ($A_s$).

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

$$Z(\text{score}) = (\text{score} - \mu)/\sigma$$

$\mu$ … average score of all models

$\sigma$ … 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 [%] | Coverage [% aa] | Initial prediction | | Final prediction | | Best prediction | |
|---|---|---|---|---|---|---|---|---|
| | | | $C\alpha$ RMSD [Å] | CE overlap [%] | $C\alpha$ RMSD [Å] | CE overlap [%] | $C\alpha$ RMSD [Å] | 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



Sequence identity         4.4%

Initial model C$\alpha$ RMSD 10.1Å

Final model C$\alpha$ RMSD   3.6Å

# **Modeling genes**

# Structural analysis of missense mutations in human BRCA1 BRCT domains

Cannot measure the functional impact of every possible SNP at all positions in each protein! Thus, prediction based on general principles of protein structure is needed.



*Mirkovic et al. (2004) Cancer Research 64 pp3790*

# Human BRCA1 and its two BRCT domains



RING NLS BRCT

Globular regions
Nonglobular regions
200 aa

BRCA1 BRCT repeats, 1jnx

**MYRIAD**

*BRACAnalysis*™
Comprehensive BRCA1-BRCA2 Gene Sequence Analysis Result

| | | |
|---|---|---|
| Niece Singer, MS<br>Strang Cancer Prevention Center<br>428 E 72nd St<br>New York, NY 10021 | **SPECIMEN**<br>Specimen Type: Blood<br>Draw Date: n/a<br>Accession Date: Oct 27, 2000<br>Report Date: Nov 17, 2000 | **PATIENT**<br>Name:<br>Date of Birth: Feb 02, 1953<br>Patient ID:<br>Gender: Female<br>Accession #: 00019998<br>Requisition #: 58694 |

Physician: Fred Gilbert, MD

**Test Result**

| Gene Analyzed | Specific Genetic Variant |
|---|---|
| BRCA2 | H2116R |
| BRCA1 | None Detected |

## Interpretation

### GENETIC VARIANT OF UNCERTAIN SIGNIFICANCE

The BRCA2 variant H2116R results in the substitution of arginine for histidine at amino acid position 2116 of the BRCA2 protein. Variants of this type may or may not affect BRCA2 protein function. Therefore, the contribution of this variant to the relative risk of breast or ovarian cancer cannot be established solely from this analysis. The observation by Myriad Genetic Laboratories of this particular variant in an individual with a deleterious truncating mutation in BRCA2, however, reduces the likelihood that H2116R is itself deleterious.

Authorized Signature:

Brian E. Ward, Ph.D.
Laboratory Director

Thomas S. Frank, M.D.
Medical Director

53

# Missense mutations in BRCT domains by function



|  | cancer associated | not cancer associated | ? | | |
|---|---|---|---|---|---|
| **no transcription activation** | C1697R<br>R1699W<br>A1708E<br>S1715R<br>P1749R<br>M1775R | | M1652K<br>L1657P<br>E1660G<br>H1686Q<br>R1699Q<br>K1702E<br>Y1703HF<br>1704S | L1705PS<br>1715NS1<br>722FF17<br>34LG173<br>8EG174<br>3RA175<br>2PF1761<br>I | F1761S<br>M1775E<br>M1775K<br>L1780P<br>I1807S<br>V1833E<br>A1843T |
| **transcription activation** | | M1652I<br>A1669S | V1665M<br>D1692N<br>G1706A<br>D1733G<br>M1775V<br>P1806A | | |
| **?** | | | M1652T W1718S<br>V1653M T1720A<br>L1664P W1730S<br>T1685A F1734S<br>T1685I E1735K<br>M1689R V1736A<br>D1692Y G1738R<br>F1695L D1739E<br>V1696L D1739G<br>R1699L D1739Y<br>G1706E V1741G<br>W1718C H1746N | R1751P<br>R1751Q<br>R1758G<br>L1764P<br>I1766S<br>P1771L<br>T1773S<br>P1776S<br>D1778N<br>D1778G<br>D1778H<br>M1783T | C1787S A1823T<br>G1788D V1833M<br>G1788V W1837R<br>G1803A W1837G<br>V1804D S1841N<br>V1808A A1843P<br>V1809A T1852S<br>V1809F P1856T<br>V1810G P1859R<br>Q1811R<br>P1812S<br>N1819S |

54

"Decision" tree for predicting functional impact of genetic variants

55

# Putative binding site on BRCA1



Putative binding site predicted in 2003 and accepted for publication on March 2004.

Williams *et al.* 2004 Nature Structure Biology. **June 2004 11**:519

Mirkovic *et al.* 2004 Cancer Research. **June 2004 64**:3790

# 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

Cavity is not filled

Ligand binding cavity

BLBP/docosahexaenoic acid

Cavity is filled



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.*

# Does RuvB have the same fold as δ' of E.coli DNA polymerase III?

```
Ec d'   MRWYPWLRPDFEKLVASYQAGRG----HHALLIQALPGMGDDALIYALSRYLLCQQPQGHKSCGHCRG
RUVB    LEEYVGQPQVRSQMEIFIKAAKLRGDALDHLLIFGPPGLGKTTLANIVANEMG--------------

Ec d'   CQLMQAGTHPDYYTLAPEKGKATLGVDAVREVTEKLNEAARLGGAKVVWVTDAALLTDAAANALLKTL
RUVB    ----------VNLRTT-------SGPVLEKAGDLAAMLTNLEPHDVLFIDEIHRLSPVVEEVLYPAM

Ec d'   ----------------EEPPAETWFFLATREPERL---LATLRSRCRLHYLAPPPEQYAVTWLSRE
Ppdp    EDYQLDIMIGEGPAARSIKIDLPPFTLIGATTRAGSLTSPLRDRFGIVQRLEFY--QVPDLQYIVSRS

Ec d'   VTM-----SQDALLAALRLSAGSPGAALALFQ------------GDNWQARETLCQALAYSVPSGD--
RUVB    ARFMGLEMSDDGALEVARRARGTPRIANRLLRRVRDFAEVKHDGTISADIAAQALDMLNVDAEGFDYM

Ec d'   -WYSLLAALN---HEQAPARLHWLATLLMDALKR/VTNVDVPGLVAELANHL---SPSRLQAILGDVC
RUVB    DRKLLLAVIDKFF-GGPVGLDNLAAAIGEERETIE--DVLEPYLIQQGFLQRTPRGRMATTRAWNHFG

        Ec d'   HIREQLMSVAGANRELLITDLLLRIEHYLQPGVVLP
        RUVB    ITPPEMP-----------------------------
```



Energy profiles (ProsaII by M. Sippl)

*B. Guenther, et al. Cell 91, 335, 1997.*
*Yamada, K., et al. Proc.Nat.Acad.Sci.USA 98,1442, 2001.*

# Modeling genomes

# Structural Genomics

**Characterize most protein <span style="color:green">sequences</span> based on related known <span style="color:darkred">structures</span>**

1. The number of "**families**" is much **smaller** than the number of proteins.
2. **Any one** of the members of a family is **fine**.



protein space

There are **~16,000** families (90%) @ 30% sequence identity cutoff

*Sali. Nat. Struct. Biol. **5**, 1029, 1998.*
*Sali et al. Nat. Struct. Biol., **7**, 986, 2000.*
*Sali. Nat. Struct. Biol. **7**, 484, 2001.*
*Baker & Sali. Science **294**, 93, 2001.*
*Vitkup et al. Nat. Struct. Biol. **8**, 559, 2001*

# Structure Space & Structural Genomics

**Structure map @ >30% sequence identity**



Isolated cluster of two structural genomics entries in PDB corresponding to *1l6r* (chains A and B) and *1kyt* (chains A and B)

MshbS, Lmbe-Related Proteins (including chains from *1q74*, and *1q7t* PDB entries) structuraly joined by *1uan* chain A from structural genomics

Thymidylate Synthase Complementing Proteins (including chains from *1o2a*, *1o2b*, *1o24*, *1o25*, *1o26*, *1o27*, and *1o28* PDB entries) structurally joined by *1kq4* chains from structural genomics

# MODPIPE2.0
## Large-Scale Protein Structure Modeling



*Eswar et.al., (2003) Nucl.Acids.Res. 31(13)*

# ModBase Statistics

Large-scale modeling of the TrEMBL-SWISSPROT databases

| | |
|---|---:|
| **Sequences (total)** | 1,679,742 |
| **Sequences (modeled)** | 964,442 |
| **Models** | 2,947,461 |

# MODBASE

http://salilab.org/modbase

*Search Page*



*Model Details*



*Sequence Overview*



*Model Overview*



*Pieper et al. (2004) Nucleic Acids Research 32, D217-D222*

**Protein Structure Prediction**
**SUMMARY**

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

# A suite of programs, servers and databases for comparative protein structure modeling

## http://salilab.org

**LS-SNP**
**Web Server**
http://salilab.org/LS-SNP
Predicts functional impact of residue substitution

**PIBASE**
**Database**
http://salilab.org/pibase
Contains structurally defined protein interfaces

**CCPR**
**Center for Computational Proteomics Research**
http://www.ccpr.ucsf.edu

**MODLOOP**
**Web Server**
http://salilab.org/modloop
Models loops in protein structures

**MODBASE**
**Database**
http://salilab.org/modbase
Fold assignments,alignments models, model assessments for all sequences related to a known structure

**MODWEB**
**Web Server**
http://salilab.org/modweb
Provides a web interface to MODPIPE

**MODELLER**
**Program**
http://salilab.org/modeller
Implements most operations in comparative modeling

**DBALI**
**Database**
http://salilab.org/dbali
Contains a comprehensive set of pairwise and multiple structure-based alignments

**ICEDB**
**Database/LIMS**
http://nysgxrc.org
Tracks targets for structural genomics by NYSGXRC

**MODPIPE**
**Program**
Automatically calculates comparative models of many protein sequences

**EVA**
**Web Server**
http://salilab.org/eva
Evaluates and ranks web servers for protein structure prediction

**LIGBASE**
**Database**
Ligand binding sites and inheritance (accessible through MODBASE)

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

| Name | Type[a] | World Wide Web address[b] |
|---|---|---|
| **DATABASES** | | |
| CATH | S | http://www.biochem.ucl.ac.uk/bsm/cath/ |
| DBAli | S | http://www.salilab.org/DBAli/ |
| 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.nlm.nih.gov/ |
| PDB | S | http://www.rcsb.org/pdb/ |
| PSI | S | http://www.nigms.nih.gov/psi/ |
| 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 ASSIGNMENT** | | |
| 123D | S | http://123d.ncifcrf.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.nlm.nih.gov/BLAST/ |
| DALI | S | http://www2.ebi.ac.uk/dali/ |
| 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/ |
| FUGUE | S | http://www-cryst.bioc.cam.ac.uk/~fugue/ |
| LOOPP | S | http://ser-loopp.tc.cornell.edu/cbsu/loopp.htm |
| PDB-BLast/FASS | S | http://bioinformatics.ljcrf.edu/pdb_blast/ |
| PHD, TOPITS | S | http://www.predictprotein.org/ |

**http://salilab.org/bioinformatics_resources.shtml**

# Happy Modeling!

# Master Bioinformatics for Health Sciences
## MODELLER tutorial

$>mod8v1 model.py

**Marc A. Marti-Renom**
*Adjunct Assistant Professor*
http://salilab.org/~marcius

Depts. of Biopharmaceutical Sciences and Pharmaceutical Chemistry
California Institute for Quantitative Biomedical Research
University of California at San Francisco

UCSF
University of California
San Francisco

# Steps in Comparative Protein Structure Modeling



**TARGET**

**TEMPLATE**

ASILPKRLFGNCEQTSDEGL
KIERTPLVPHISAQNVCLKID
DVPERLIPERASFQWMNDK

ASILPKRLFGNCEQTSDEGLKIERTPLVPHISAQNVCLKIDDVPERLIPE
MSVIPKRLYGNCEQTSEEAIRIEDSPIV---TADLVCLKIDEIPERLVGE

*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.*

# Comparative modeling by satisfaction of spatial restraints
## MODELLER



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

1. Extract spatial restraints

2. Satisfy spatial restraints

$$F(R) = \prod_i p_i (f_i / I)$$

*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.*

3

# Typical errors in comparative models

MODEL
X-RAY
TEMPLATE

Incorrect template

Misalignment

```
              10        20        30        40        50        60
EDN   ---KPPQFTWAQWFETQHINMTSQQCTNAMQVINNYQRRCKNQNTFLLTTFANVVNVCGNPNMTCPSN
7RSA  KETAAAKFERQHMDSSTSAASSSNYCNQMMKSRNLTKDRCKPVNTFVHESLADVQAVCSQKNVAC-KN
           aaaaaaaaaa       aaaaaaaaaaaaaa    bbbbbbb aaaaaaaaa

         70        80        90       100       110       120       130
EDN   KTRKNCHHSGSQVPLIHCNLTTPSPQNISNCRYAQTPANMFYIVACDNRDQRRDPPQYPVVPVHLDRII
7RSA  -GQTNCYQSYSTMSITDCRETGSS--KYPNCAYKTTQANKHIIVACEGN---------PYVPVHFDASV
          bbbbb   bbbbbbbbbbb    aaaaabbbbbbbbbbbbbbbbbbbbb     bbbbbbbbbbb
```

Region without a
template

Distortion/shifts in
aligned regions

Sidechain packing

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

4

# Utility of protein structure models, despite errors



D. Baker & A. Sali. Science 294, 93, 2001.

# **New features in MODELLER 8**

# Overview

- SALIGN: versatile alignment module

- MODPIPE2.0: Large-scale protein structure modeling

- profile.build(): Iterative database searching & profile construction

- profile.scan(): fold-assignment & profile-profile alignments

- New objective functions

- Mod-EM: Fitting a model into an EM map

- Additional features and bugs fixed

# SALIGN: Versatile alignment command in MODELLER



Best score

Best local alignment

A
B
C
D

B
D
A
C

- similarity +

✓ **Uses all available structural information**
✓ **Provides the optimal alignment**

• **Computationally expensive**

$$\Omega_i$$
$$d_i$$

$$RMSD = \sqrt{\sum (x_i - \overline{x})^2}$$

$$R_{i,j} \qquad D_{i(3),j(3)} \qquad S_{i,j} \qquad B_{i,j} \qquad I_{i,j}$$

$$Score_{i,j} = w_1 * R_{i,j} + w_2 * D_{i(a),j(a)} + w_3 * S_{i,j} + w_4 * B_{i,j} + w_5 * I_{i,j} + w_6 * X_{i,j}$$

*Madhusudhan et al. in preparation*

# MODPIPE2.0
## Large-Scale Protein Structure Modeling



*Eswar et.al., (2003) Nucl.Acids.Res. 31(13)*

# profile.build(): Iterative database searching & profile construction

- Rigorous Smith-Waterman local alignments
- Statistical significance using z-scores
- Null model for statistics from actual sequences
- Built-in mechanism to detect profile divergence

# profile.scan(): Large-scale fold-assignment & profile-profile alignments

- Correlation Coefficient as the scoring scheme
- Smith-Waterman local alignments
- Statistical significance using z-scores
- Null model for statistics from actual sequences



| 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 |
| $CC_{HH}$ | 56 ± 23 | 0.61 ± 0.24 |
| $CC_{HS}$ | 56 ± 24 | 0.62 ± 0.24 |
| TOP | 62 ± 20 | 0.67 ± 0.20 |

# New objective functions in Modeller

◇ Discrete Optimized Protein Energy (DOPE)

  ◇ A highly accurate distance dependent statistical potential

  ◇ The best performer among 29 tested scoring functions

  ◇ Has been released with Modeller 8v0

  ◇ Applications to loop modeling and model assessment

◇ Solvation model in Modeller

  ◇ The GB/SA solvation model is implemented in Modeller

  ◇ A newer scheme for point charges that concurs with CHARMM22

  ◇ Will be included in Modeller 8v1

◇ Other modifications

  ◇ Minor modifications to the atomic van der Waals radii (radii.lib)

# Mod-EM: Fitting a model into an EM map

**Allows protein models to be docked into lower-resolution electron microscopy maps for scoring**

**Representation**:

*map:* density in voxels, $\rho^{EM}(r)$
*protein:* based on atomic mass, $\rho^{probe}(r)$ - can be represented by different functions: Gaussian, uniform sphere model, hybrid Gaussian/sphere model
*filters:* cutoff filters, sqr filter, Laplacian

**Scoring function:**

Cross-correlation function (C) between the map density and the probe 'density':

$$C = \frac{\sum_{i=1}^{M} \rho_i^{EM} \left( \sum_{j=1}^{N} \rho_{i,j}^{probe} \div \right)}{\sqrt{\sum_{i=1}^{M} \left( \rho_i^{EM} \right)^2 \sum_{i=1}^{M} \left( \sum_{j=1}^{N} \rho_{i,j}^{probe} \div \right)^2}}$$

*N* - number of atoms in the probe

*M* - number of cryoEM grid points covered by the probe density



**Optimization methods (density.grid_search()):**

Exhaustive rigid rotation and translation
Monte Carlo optimization

*Topf et al., J. Struct. Biol.* **149**, *191-203 (2005)*

# Additional features

- **New methods for model assessment:**

  - DOPE

  - GA341

  - ModEM scoring

- **Improved user interface:**

  - Python interface allows for more powerful and flexible scripting, and easier integration with other applications

  - Legacy TOP scripts are still parsed

# Obtaining **MODELLER** and related information

- MODELLER (8v0) 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 ☹☹
- Script is written in PYTHON language ☺
- You may know Python language is simple ☺☺

# Using MODELLER

- INPUT:
  - Target Sequence (FASTA/PIR format)
  - Template Structure (PDB format)
  - TOP command 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*
```

# Modeling of BLBP
## STEP 1: Align blbp and 1hms sequences
### *Python script for target-template alignment*

```python
# 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')
```

Run by typing `mod8v0 align.py` in the directory where you have the python file.
MODELLER will produce a `align.log` file

# Modeling of BLBP
## STEP 1: Align blbp and 1hms sequences
### *Python script for target-template alignment*

```python
# 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')
```

Run by typing `mod8v0 align.py` in the directory where you have the python file.
MODELLER will produce a `align.log` file

# Modeling of BLBP
## STEP 1: Align blbp and 1hms sequences
### *Python script for target-template alignment*

```python
# 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')
```

Run by typing `mod8v0 align.py` in the directory where you have the python file.
MODELLER will produce a `align.log` file

# Modeling of BLBP
## STEP 1: Align blbp and 1hms sequences
### Python script for target-template alignment

```python
# 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')
```

Run by typing `mod8v0 align.py` in the directory where you have the python file.
MODELLER will produce a `align.log` file

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

```
>P1;1hms
structureX:1hms:   1 : : 131 : :undefined:undefined:-1.00:-1.00
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
structureX:1hms:    1 : : 131 : :undefined:undefined:-1.00:-1.00
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              10        20        30        40        50        60
1hms      VDAFLGTWKLVDSKNFDDYMKSLGVGFATRQVASMTKPTTIIEKNGDILTLKTHSTFKNTEISFKLGV
blbp      VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGE
_consrvd  ****    **** ** *** *** **********  **** **    *      *  ****** * **


_aln.p    70        80        90        100       110       120       130
1hms      EFDETTADDRKVKSIVTLDGGKLVHLQKWDGQETTLVRELIDGKLILTLTHGTAVCTRTYEKE
blbp      EFEETSIDDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFGDIVAVRCYEKA
_consrvd  ** **   ***   ** * *** ** * **** **    **   ***   ***  *    * * * ***
```

# Modeling of BLBP
## STEP 2: Model the blbp structure using the alignment from step 1.
### Python script for model building

```python
# 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)
a.make()                             # do the actual homology modelling
```

Run by typing `mod8v0 model.py` in the directory where you have the python file. MODELLER will produce a `align.log` file

# Modeling of BLBP
## STEP 2: Model the blbp structure using the alignment from step 1.
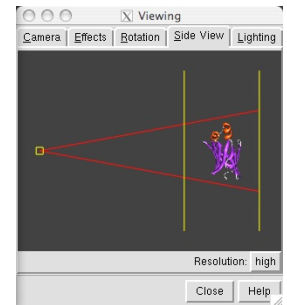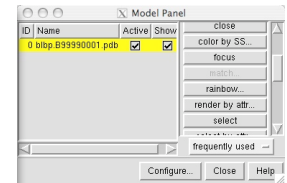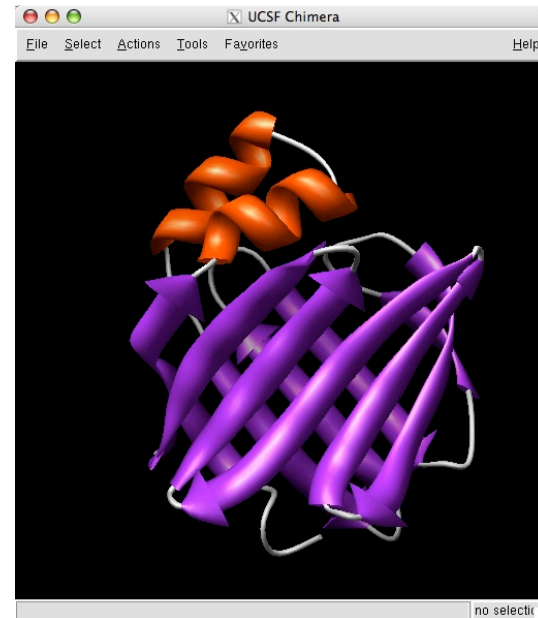### *Python script for model building*

```python
# 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)
a.make()                              # do the actual homology modelling
```

Run by typing `mod8v0 model.py` in the directory where you have the python file.
MODELLER will produce a `align.log` file

# Modeling of BLBP
## STEP 2: Model the blbp structure using the alignment from step 1.
### *Python script for model building*

```python
# 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)
a.make()                            # do the actual homology modelling
```

Run by typing `mod8v0 model.py` in the directory where you have the python file.
MODELLER will produce a `align.log` file

# Modeling of BLBP
## STEP 2: Model the blbp structure using the alignment from step 1.
### *Python script for model building*

PDB file

Can be viewed with Chimera

http://www.cgl.ucsf.edu/chimera/

Rasmol

http://www.openrasmol.org



Model file → blbp.B99990001

# http://www.salilab.org/bioinformatics_resources.shtml

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