# Comparative Protein Structure Prediction



**Marc A. Marti-Renom**
http://bioinfo.cipf.es/sgu/

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

PRINCIPE FELIPE
CENTRO DE INVESTIGACION

# 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.ncrf.edu/pdb_blast/ |
| PHD, TOPITS | S | http://www.predictprotein.org/ |

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

2

# Summary

- **INTRO**
- **MODELLER**
- **MOULDER**
- **MODEL(S) --> FUNCTION**
- **MODELLER example**

# Nomenclature

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

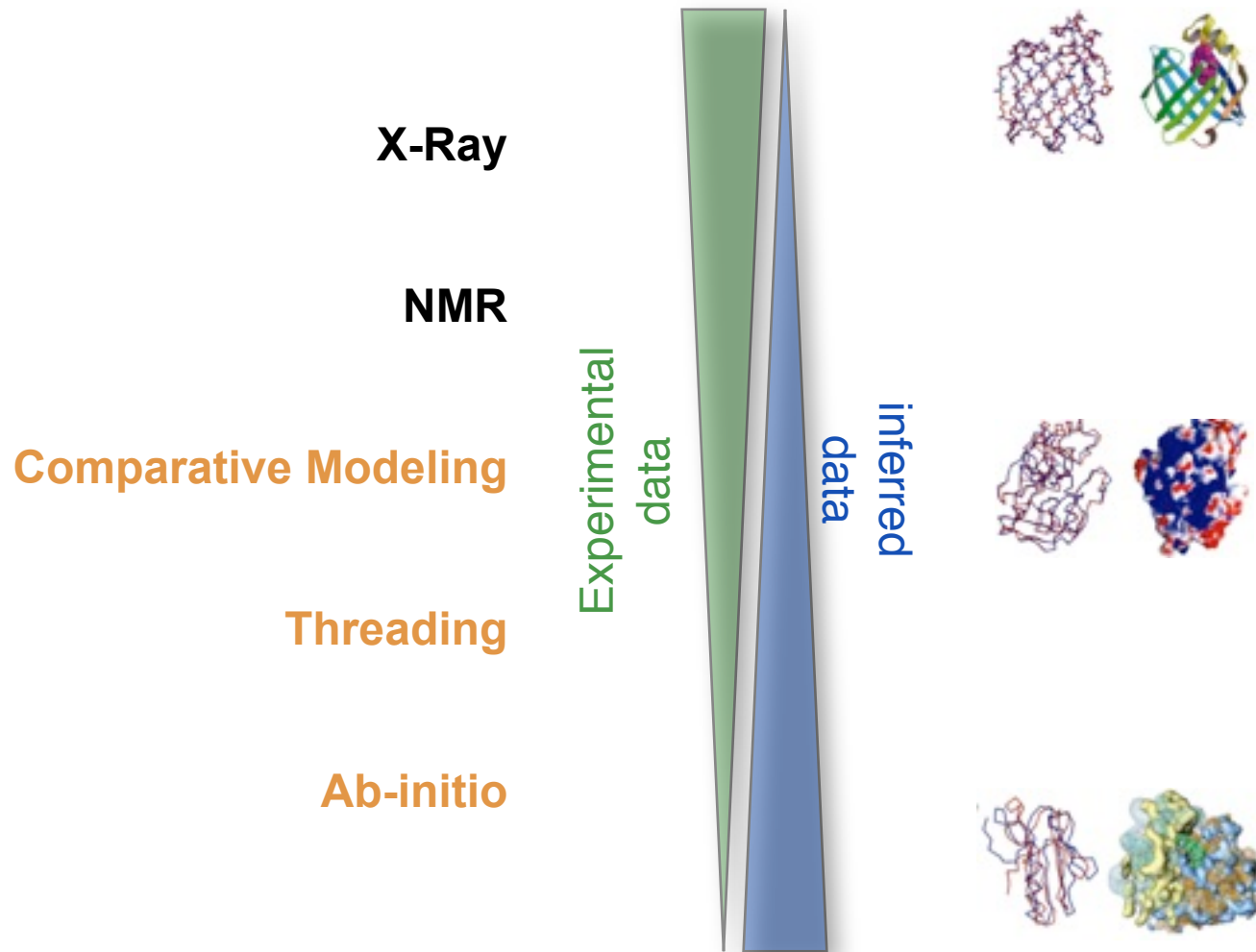**Similarity**: Score that quantifies the degree of relationship between two sequences.

**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



**X-Ray**

**NMR**

**Comparative Modeling**

**Threading**

**Ab-initio**

Experimental data

inferred data

# 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.

YDL117W
(15-64)

KARYGWSGQTKGDLGFLEGDIMEVTRIAGVWYFYGKLLRNKKCSGYFPHF

Ser 30
Asn 49
Trp 31
Pro 47
Phe 50
Tyr 4

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**.

6

# Principles of protein structure

GFCHIKAYTRLIMVG...



Desulfovibrio vulgaris

Condrus crispus

Anabaena 7120

Anacystis nidulans

GFCHIKAYTRLIMVG...

Folding (physics)

*Ab initio* prediction

Evolution (rules)

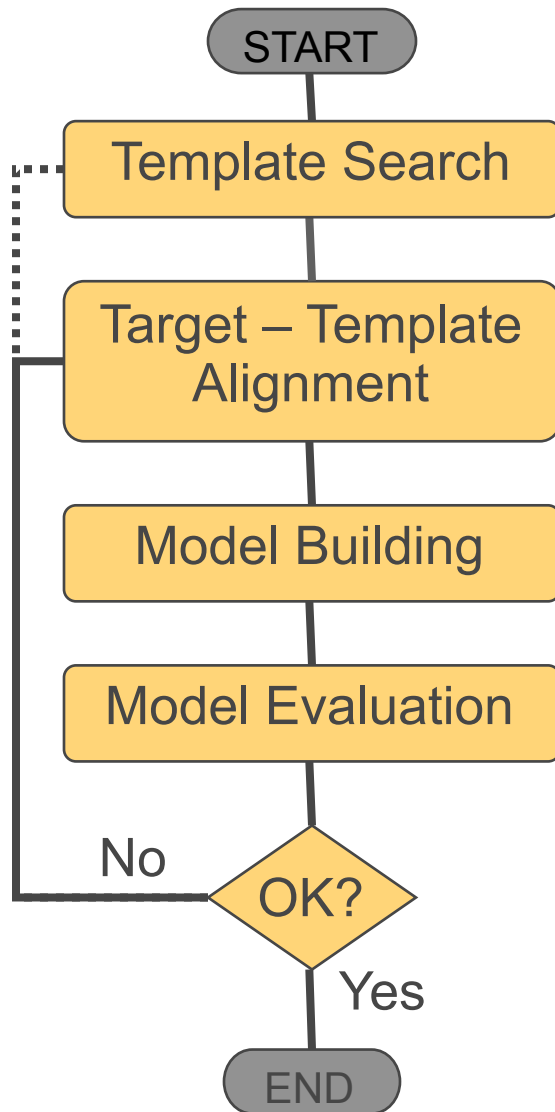Threading
Comparative Modeling

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

Thursday, April 23, 2009

# MODELLER

1. 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.*
2. M.A. Marti-Renom, et al.. *Comparative protein structure modeling of genes and genomes. Annu. Rev. Biophys. Biomol. Struct. 29, 291-325, 2000.*
3. A. Sali & T.L. Blundell. *Comparative protein modelling by satisfaction of spatial restraints. J. Mol. Biol. 234, 779-815, 1993.*
4. A. Fiser, R.K. Do, & A. Sali. *Modeling of loops in protein structures, Protein Science 9. 1753-1773, 2000.*

Thursday, April 23, 2009

# Steps in Comparative Protein Structure Modeling



START

**Template Search**

**Target – Template Alignment**

**Model Building**

**Model Evaluation**

No

OK?

Yes

END

TARGET   TEMPLATE

ASILPKRLFGNCEQTSDEG
LKIERTPLVPHISAQNVCLKI
DDVPERLIPERASFQWMN
DK

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

Start with a
Target Sequence

Template
Search

Target/Template
Alignment

Build model

Evaluate model

OK?

Output 3D Model

**Given an alignment...**

**extract spatial features
from the template(s)
and statistics from
known structures**

**apply these features
as restraints on your
target sequence**

**optimize to find the
best solution for the
restraints to produce
your 3D model**

MSVIPKR--GNCEQTSE
ASILPKRLFGNCEQTSD



+

FREQUENCY

4E6

3E6

2E6

1E6

15    17    19    21    23    25
$C_\alpha - C_\alpha$ DISTANCE          [Å]

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.

Thursday, April 23, 2009

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

**Wrong fold**

**Miss alignments**

**Loop regions**

**Rigid body distortions**

**Side-chain packing**





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

Thursday, April 23, 2009

# Model Accuracy

## HIGH ACCURACY

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



Sidechains
Core backbone
Loops

X-RAY / MODEL

## 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.*

**MOULDER**

*John, Sali (2003). NAR pp31 3982*

13

# Moulding: iterative alignment, model building, model assessment

# 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 ($A_s$).

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

$$Z(score) = (score - \mu)/\sigma$$

$\mu$ … 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 | |
|---|---|---|---|---|---|---|---|---|
| | | | Cα RMSD [Å] | CE overlap [%] | Cα RMSD [Å] | CE overlap [%] | Cα 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 Cα RMSD 10.1Å

Final model Cα RMSD   3.6Å

a        b        c        d

# Can we use models to infer function?



*T. cruzi*

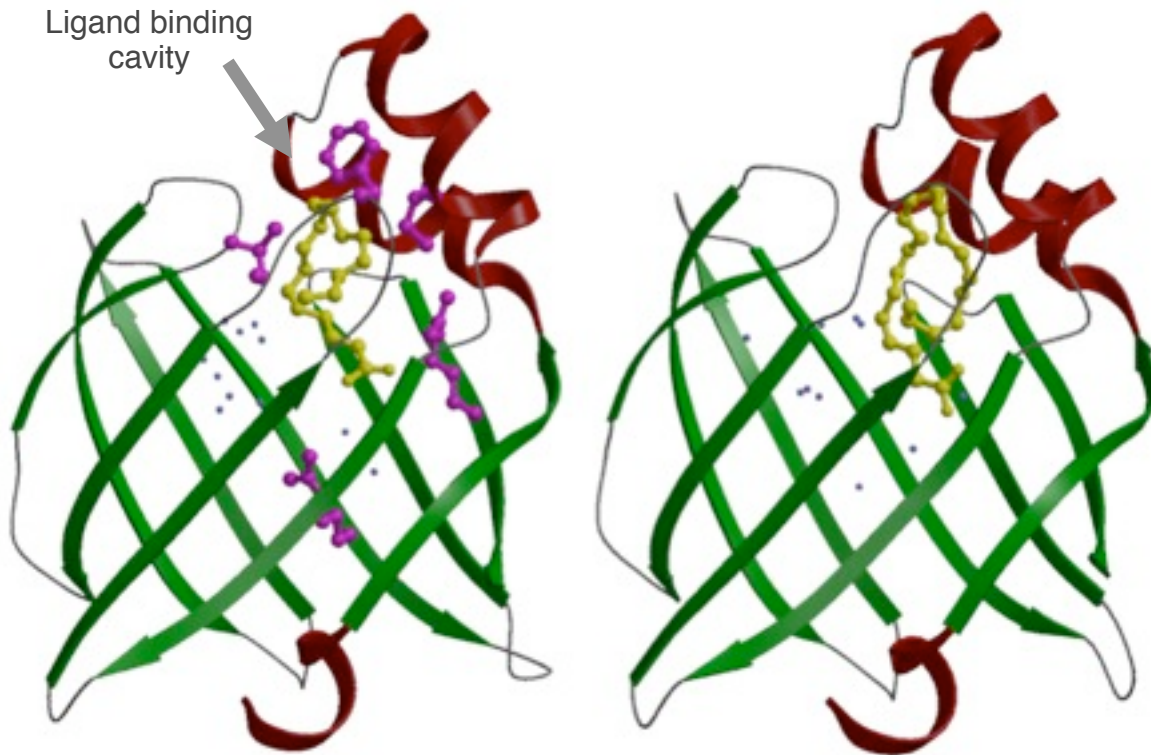# What is the physiological ligand of Brain Lipid-Binding Protein?

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

BLBP/oleic acid

Cavity is not filled

BLBP/docosahexaenoic acid

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.

# Structural analysis of missense mutations in human BRCA1 BRCT domains

Nebojsa Mirkovic, Marc A. Marti-Renom, Barbara L. Weber, Andrej Sali and Alvaro N.A. Monteiro

**Cancer Research (June 2004). 64:3790-97**

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.

# 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  L1705PS  F1761S<br>L1657P  1715NS1  M1775E<br>E1660G  722FF17  M1775K<br>H1686Q  34LG173  L1780P<br>R1699Q  8EG1743  I1807S<br>K1702E  RA1752  V1833E<br>Y1703HF  PF1761I  A1843T<br>1704S |
| **transcription activation** | | M1652I<br>A1669S | V1665M<br>D1692N<br>G1706A<br>D1733G<br>M1775V<br>P1806A |
| **?** | | | M1652T W1718S R1751P C1787S A1823T<br>V1653M T1720A R1751Q G1788D V1833M<br>L1664P W1730S R1758G G1788V W1837R<br>T1685A F1734S L1764P G1803A W1837G<br>T1685I E1735K I1766S V1804D S1841N<br>M1689R V1736A P1771L V1808A A1843P<br>D1692Y G1738R T1773S V1809A T1852S<br>F1695L D1739E P1776S V1809F P1856T<br>V1696L D1739G D1778N V1810G P1859R<br>R1699L D1739Y D1778G Q1811R<br>G1706E V1741G D1778H P1812S<br>W1718C H1746N M1783T N1819S |

# 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

# *S. cerevisiae* ribosome



Fitting of comparative models into 15Å cryo-electron 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.

Thursday, April 23, 2009

# The Nucleopore complex Cell evolution (?)



Prokaryote  Early Eukaryote  Modern Eukaryote

Nup84
Nup85
Nup145C
Nup120
Nup133

Thursday, April 23, 2009

# Tropical Disease Initiative (TDI)
*Predicting binding sites in protein structure models.*



<http://www.tropicaldisease.org>



26

# Need is High in the Tail



DALY Burden Per Disease in Developed Countries

DALY Burden Per Disease in Developing Countries

Heart diseases

Rare diseases

DALY

Disease

Disease data taken from WHO, *World Health Report 2004*
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.*

Thursday, April 23, 2009

# Need is High in the Tail



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

Heart diseases

Rare diseases

DALY

Disease

Thursday, April 23, 2009

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

| Disease | DALY | | Disease | DALY |
|---|---|---|---|---|
| **Malaria*** | **46,486** | | Trichuriasis | 1,006 |
| 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.

28

Thursday, April 23, 2009

# DBAli$_{v2.0}$ database

http://www.dbali.org

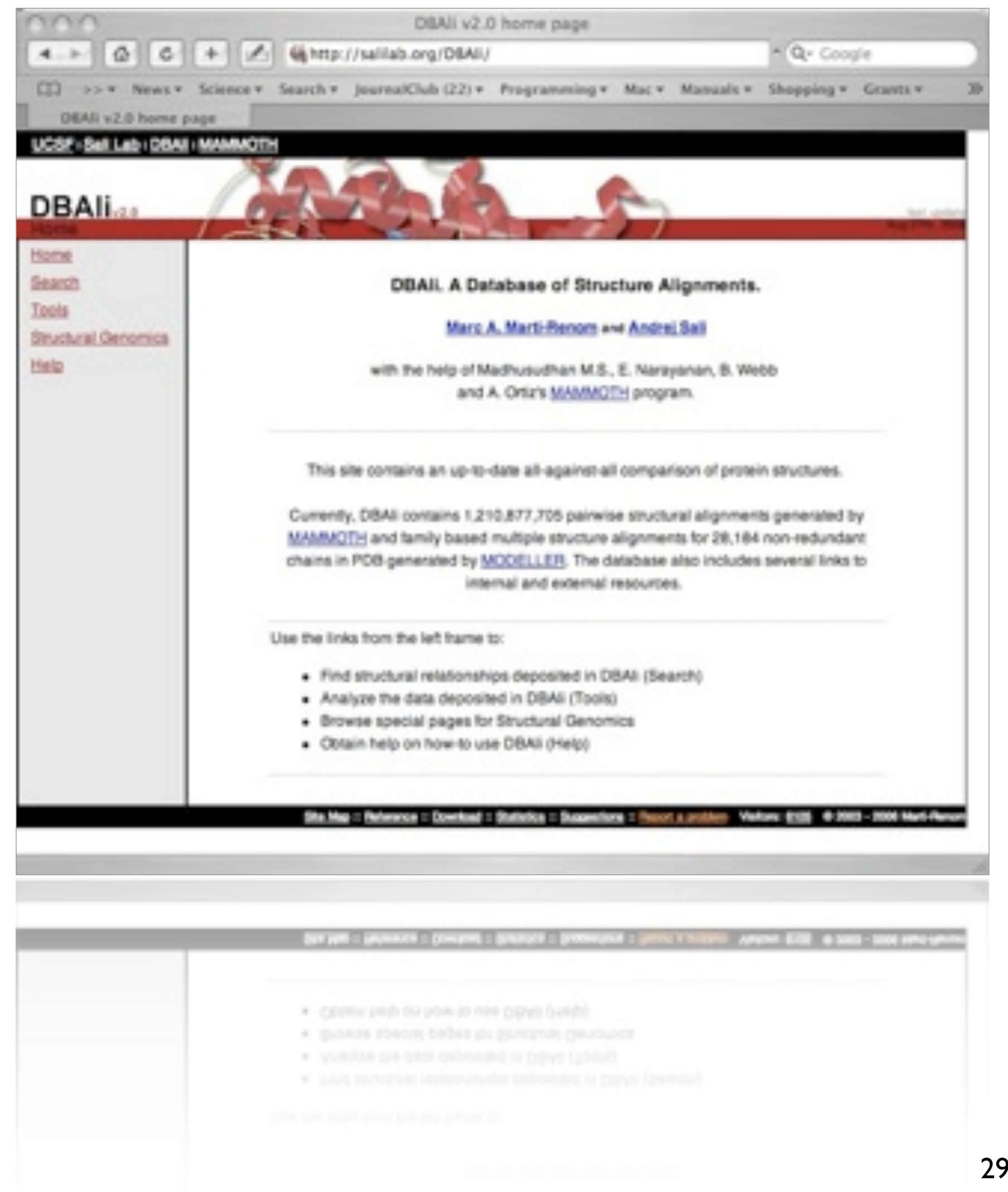


- ✓ Fully-automatic
- ✓ Data is kept up-to-date with PDB releases
- ✓ Tools for “on the fly” classification of families.
- ✓ Easy to navigate
- ✓ Provides tools for structure analysis

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

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

Uses MAMMOTH for similarity detection

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

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

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

# DBAli$_{v2.0}$ database

*Marti-Renom et al. BMC Bioinformatics (2007) Volume 8. Suppl S4*

# Method

# Sensitivity .vs. Precision

| | Optimal cut-off | Sensitivity (%)<br>Recall or TPR | Precision (%) |
|---|---|---|---|
| **Ligands** | 30% | 71.9 | 13.7 |

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

**~90-95% of residues correctly predicted**

*Marti-Renom et al. BMC Bioinformatics (2007) Volume 8. Suppl S4*

32

# Comparative docking

**Expansion**

co-crystalized protein/ligand

**2. Inheritance**

model

**1. Modeling**

crystalized protein

template

# Modeling Genomes

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



*A good model has MPQS of 1.0 or higher*

Thursday, April 23, 2009

# 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 | 1,409 | 478 | 43 | 20 |
| *M. leprae* | 1,605 | 1,178 | 893 | 310 | 25 | 6 |
| *M. tuberculosis* | 3,991 | 2,808 | 1,608 | 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** |

35

# *L. major* Histone deacetylase 2 + Vorinostat

*Template 1t64A a human HDAC8 protein.*



| PDB | ED | Template | SS | Model | G | 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.*

36

# *L. major* Histone deacetylase 2 + Vorinostat

## *Literature*

## 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

### 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 | 🔗 | Template | 🔗 | Model | ☀ | 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:**

*For the treatment of human immunovirus (HIV) infections.*

38

# *P. falciparum* thymidylate kinase + zidovudine

## NMR *Water-LOGSY* and *STD* experiments



ATM

Zidovudine

dTMP

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

39

# TDI's kernel

## http://tropicaldisease.org/kernel

*Ortí et al . "A kernel for open source drug discovery in tropical diseases". PLoS Neglected Tropical Diseases. (2009) **3**:e18*
*Ortí et al . "A Kernel for the Tropical Disease Initiative". Nature Biotechnology. (2009) **27**:320-321*



40

# Modeling Genomes

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



*A good model has MPQS of 1.1 or higher*

41

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

Thursday, April 23, 2009

43

# Comparative Protein Structure Prediction
## MODELLER tutorial

$$\$>mod9v6 \ model.py$$

**Marc A. Marti-Renom**
http://bioinfo.cipf.es/sgu/

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

PRINCIPE FELIPE
CENTRO DE INVESTIGACION

Thursday, April 23, 2009

# Obtaining **MODELLER** and related information

◇ MODELLER 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 ☺☺

# MODELLER 9v6
## 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)

47

# 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*
```

# 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 `mod9v6 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 `mod9v6 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 `mod9v6 align.py` in the directory where you have the python file.
MODELLER will produce a `align.log` file

52

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

```
# 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 `mod9v6 align.py` in the directory where you have the python file.
MODELLER will produce a `align.log` file

53

# 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 `mod9v6 model.py` in the directory where you have the python file.
MODELLER will produce a `model.log` file

Thursday, April 23, 2009

# 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 `mod9v6 model.py` in the directory where you have the python file.
MODELLER will produce a `model.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 `mod9v6 model.py` in the directory where you have the python file.
MODELLER will produce a `model.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

PyMol

http://pymol.sourceforge.net/



Model file →
`blbp.B99990001.pdb`

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

# MODWEB

# MODBASE

http://salilab.org/modbase

*Search Page*



*Model Details*



*Sequence Overview*



*Model Overview*



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

# Acknowledgments

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

Thursday, April 23, 2009