

# BMI-206

## Structure-Structure comparisons Sequence-Structure comparisons

Marc A. Marti-Renom  
Assistant Adjunct Professor  
Department of Biopharmaceutical Sciences

February 3rd, 2005

# How to use this lectures

- Ask!
- Outline
  - Basic introduction
  - Theory (representation-scoring-optimization)
  - Available programs
  - Application
- Assignment
  - *The POM152 sequence. Modeling exercise.*

# Structure-Structure comparison

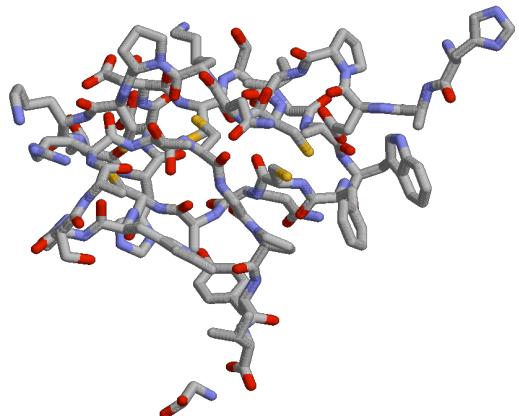
- Outline
  - Before we start...
    - Some theory
    - Coverage .vs. Accuracy
  - How can we compare structures...
    - SALIGN (properties comparison)
    - VAST (vector alignment)
    - CE (local heuristic comparison)
    - MAMMOTH (vector alignment)
  - How we classify the structural space...
    - SCOP (manual)
    - CATH (semi-automatic)
    - DBAli (fully automatic and comprehensive)

# Structure-Structure alignments

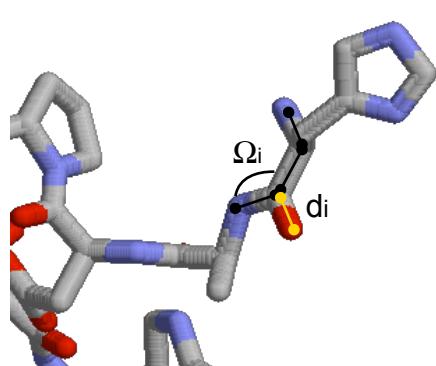
**As any other bioinformatics problem...**

- Representation
- Scoring
- Optimizer

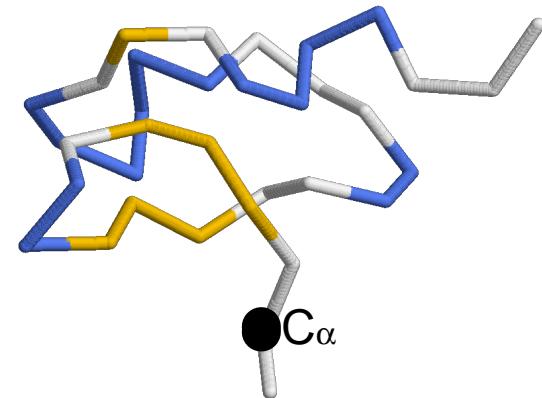
# Representation Structures



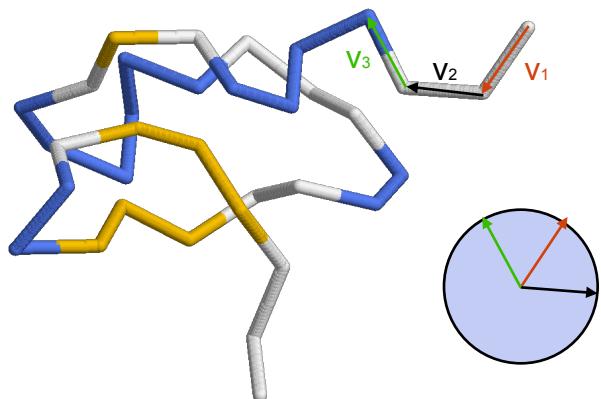
All atoms and coordinates



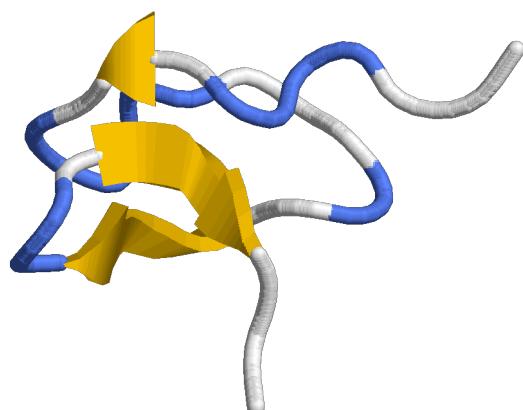
Dihedral space or distance space



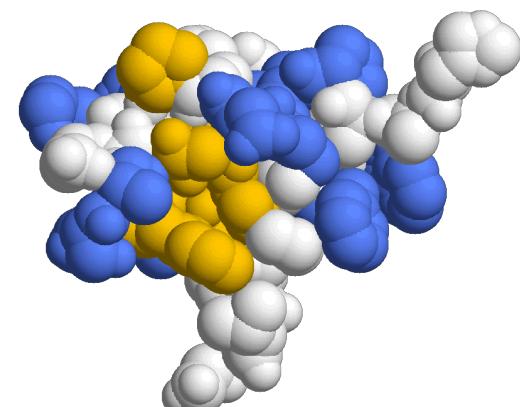
Reduced atom representation



Vector representation



Secondary Structure



Accessible surface (and others)

# Scoring Raw scores

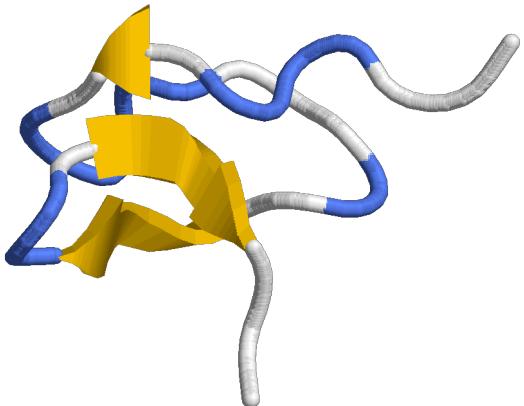
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	9	-1	-1	-3	0	-3	-3	-4	-3	-3	-3	-3	-3	-4	-4	-4	-1	-2	-2	-2
S	-1	4	1	-1	1	0	1	0	0	0	-1	-1	0	-1	-2	-2	-2	-2	-2	-3
T	-1	1	4	1	-1	1	0	1	0	0	0	-1	0	-1	-2	-2	-2	-2	-2	-3
P	-3	-1	1	7	-1	-2	-1	-1	-1	-1	-2	-2	-1	-2	-3	-3	-2	-4	-3	-4
A	0	1	-1	-1	4	0	-1	-2	-1	-1	-2	-1	-1	-1	-1	-1	-1	-2	-2	-3
G	-3	0	1	-2	0	6	-2	-1	-2	-2	-2	-2	-2	-3	-4	-4	0	-3	-3	-2
N	-3	1	0	-2	-2	0	6	1	0	0	-1	0	0	-2	-3	-3	-3	-3	-2	-4
D	-3	0	1	-1	-2	-1	1	6	2	0	-1	-2	-1	-3	-3	-4	-3	-3	-3	-4
E	-4	0	0	-1	-1	-2	0	2	5	2	0	0	1	-2	-3	-3	-3	-2	-3	-2
Q	-3	0	0	-1	-1	-2	0	0	2	5	0	1	1	0	-3	-2	-2	-3	-1	-2
H	-3	-1	0	-2	-2	-2	-1	1	0	0	8	0	-1	-2	-3	-3	-2	-1	2	-2
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5	2	-1	-3	-2	-3	-3	-2	-3
K	-3	0	0	-1	-1	-2	0	-1	1	1	-1	2	5	-1	-3	-2	-3	-3	-2	-3
M	-4	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5	1	2	-2	0	-1	-1
I	-1	-1	-2	-2	-3	-1	-4	-3	-3	-3	-3	-3	-3	1	4	2	1	0	-1	-3
L	-1	-2	-2	-3	-1	-1	-4	-4	-3	-3	-3	-2	-2	2	2	4	3	0	-1	-2
V	-1	-2	-2	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4	-1	-1	-3
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6	3	1
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	-2	-2	-1	-1	-1	-1	3	7	2	
W	-2	-3	-3	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11

2/

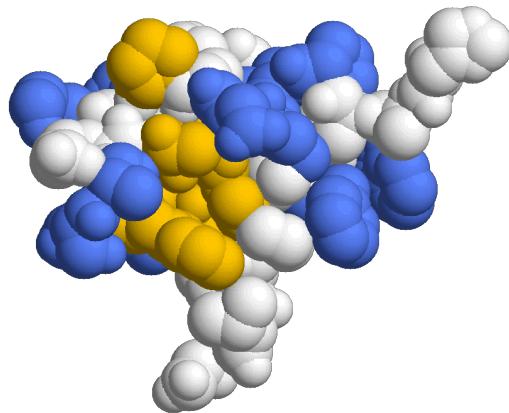
Aminoacid substitutions

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

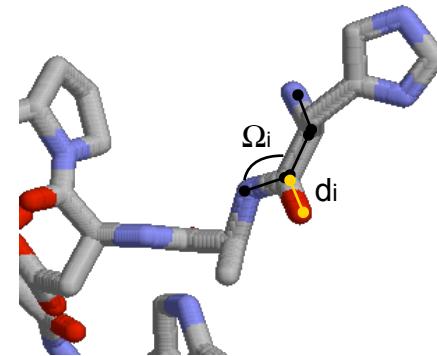
Root Mean Square Deviation



Secondary Structure (H,B,C)



Accessible surface (B,A [%])



Angles or distances

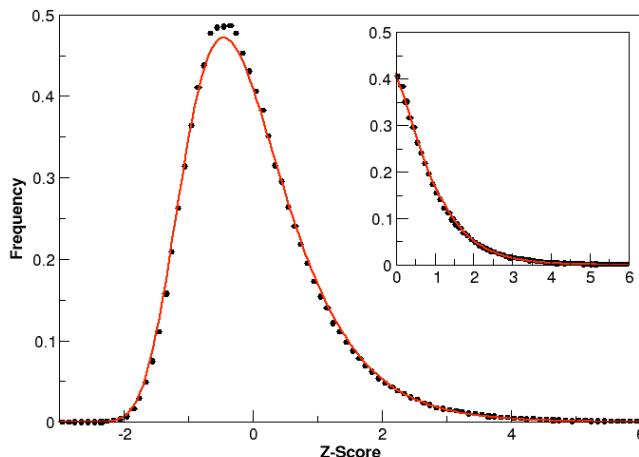
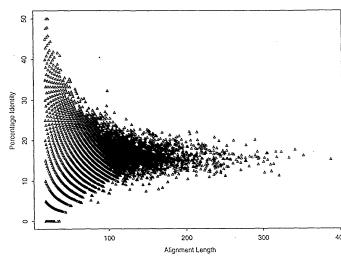
# Scoring

## Significance of an alignment (score)

*remember Patsy's class*

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

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

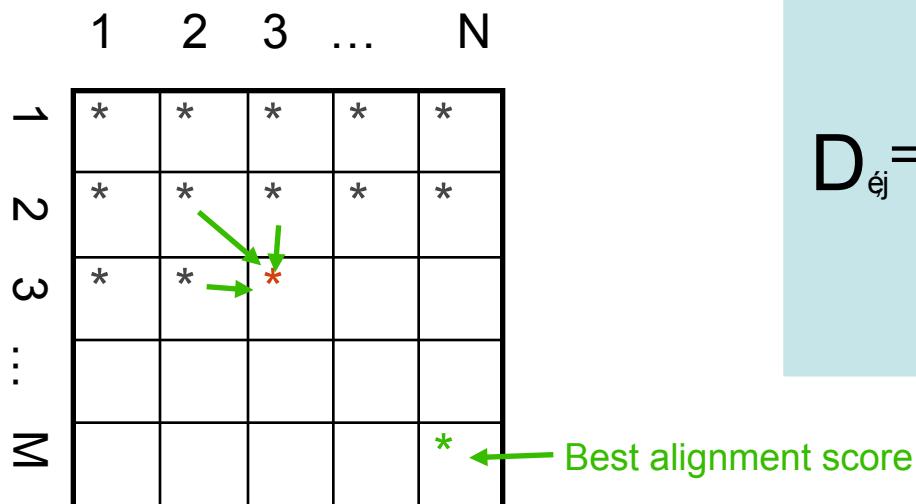
$$P(s \geq x) = 1 - \exp(e^{-\lambda (s-\mu)})$$

Karlin and Altschul, 1990 PNAS 87, pp2264

## Optimizer

# Global dynamic programming alignment

remember Patsy's class



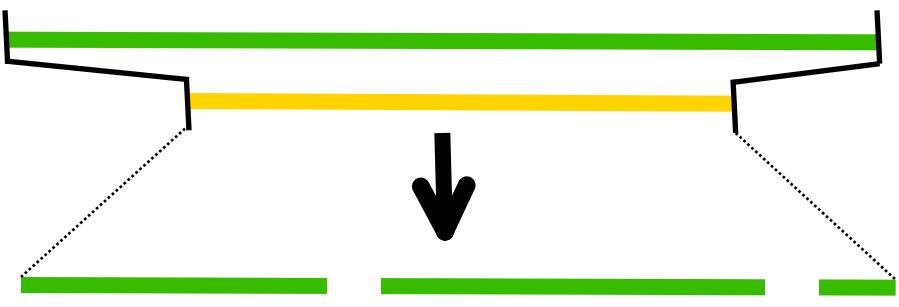
$$D_{ej} = \min \begin{cases} D_{i,j-1} + \text{Score}_{(\ddot{A}, rj)} \\ D_{i-1,j-1} + \text{Score}_{(ri, rj)} \\ D_{i-1,j} + \text{Score}_{(ri, \ddot{A})} \\ 0 \end{cases}$$

Backtracking to get the best alignment

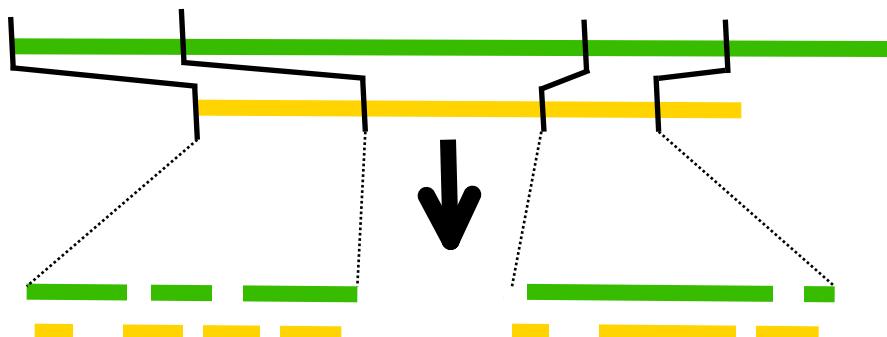
Optimizer

# Global .vs. local alignment

*remember Patsy's class*



Global alignment



Local alignment

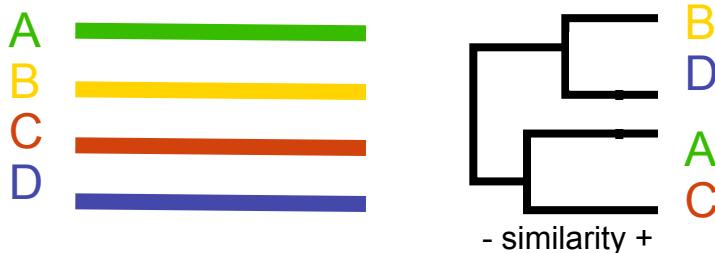
Optimizer

# Multiple alignment

*remember Patsy's class*

## Pairwise alignments

Example – 4 sequences A, B, C, D.



6 pairwise comparisons  
then cluster analysis

## Multiple alignments

Following the tree from step 1

B      yellow bar  
D      blue bar

Align the most similar pair

A      green bar  
C      orange bar

Align next most similar pair

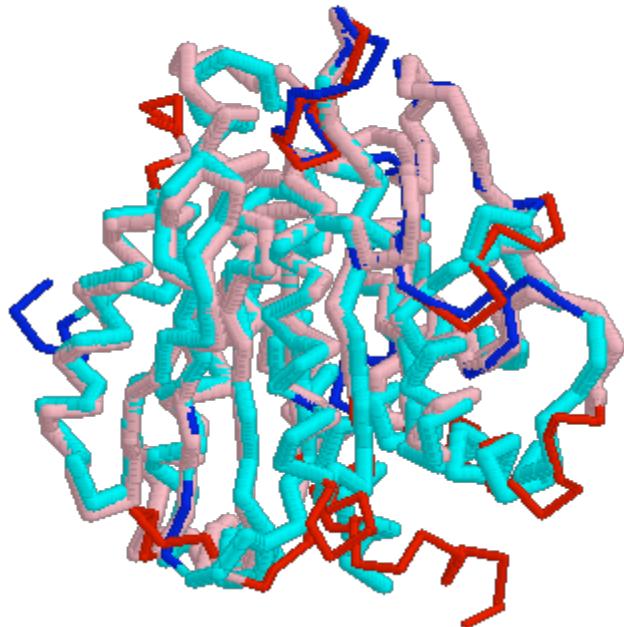
Align B-D with A-C

B      yellow bar  
D      blue bar  
A      green bar  
C      orange bar

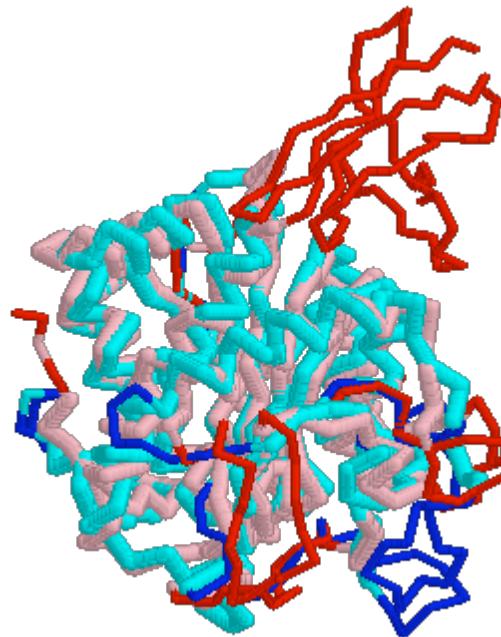


New gap in A-C to optimize  
its alignment with B-D

# Coverage .vs. Accuracy



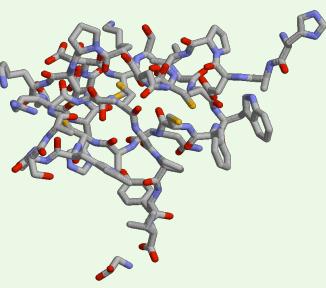
Coverage ~90% C $\alpha$



Coverage ~75% C $\alpha$

Same RMSD  $\sim 2.5\text{\AA}$

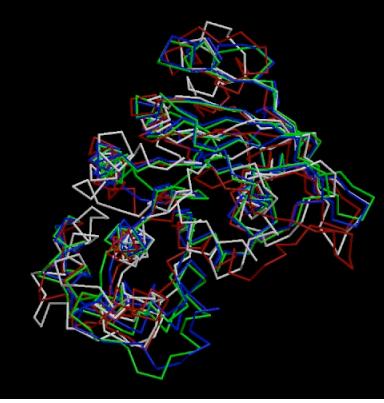
# Sequence-Structure alignment by properties conservation (SALIGN-MODELLER)





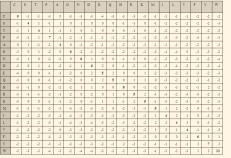
A 2D grid matrix with rows labeled 1, 2, 3, ..., M and columns labeled 1, 2, 3, ..., N. The matrix contains asterisks (\*) representing matches. A red arrow traces a path from the top-left to the bottom-right, indicating a local alignment. A green arrow points to the bottom-right corner, labeled "Best score".

Below the matrix are four horizontal bars labeled A (green), B (yellow), C (orange), and D (blue). To the right is a phylogenetic tree with nodes colored B (yellow), D (blue), A (green), and C (orange). Below the tree is the text "- similarity +".

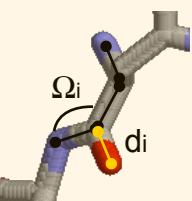


✓ Uses all available structural information  
✓ Provides the optimal alignment

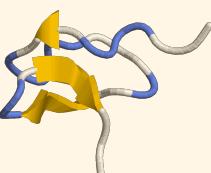
Computationally expensive



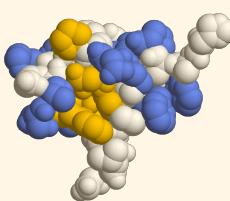
$R_{i,j}$



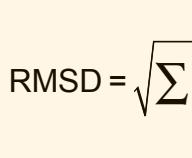
$D_{i(3),j(3)}$



$S_{i,j}$



$B_{i,j}$



$I_{i,j}$

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

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

BMI206

Madhusudhan et al. *in preparation*  
01/19/2005

# Structural alignment by properties conservation (SALIGN-MODELLER)

<http://alto.compbio.ucsf.edu/salign-cgi/index.cgi>

SALIGN Server

http://alto.compbio.ucsf.edu/salign-cgi/index.cgi Google

SALIGN Multiple Structure/Sequence Alignment Server

SALIGN is a general alignment module of the modeling program MODELLER

The alignments are computed using dynamic programming, making use of several features of the protein sequences and structures

Users can either upload their own sequences/structures to align or choose structures from the PDB

Sequences can either be pasted or uploaded as FASTA or PIR format alignment files

Paste sequence to align

Multiple sequences can be pasted by iteratively clicking 'upload' after every pasted sequence

Specify file to upload (PIR, FASTA, PDB, zip or .tar.gz)

Multiple files can be uploaded by iteratively clicking 'upload' after every file uploaded

Uploaded files:

No files uploaded

Enter 4 letter code(s) to choose PDB structures

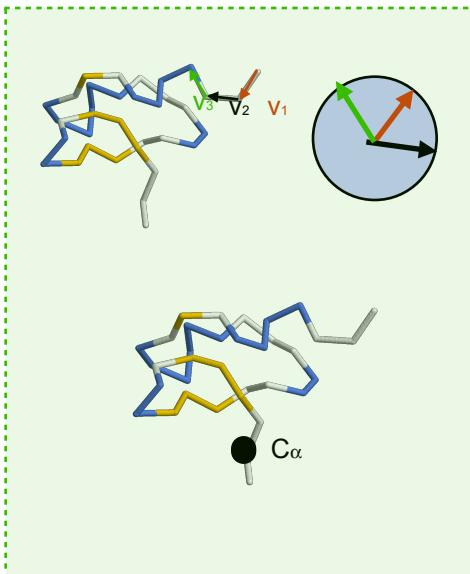
e-mail address, to receive results:

Reference:

Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A.

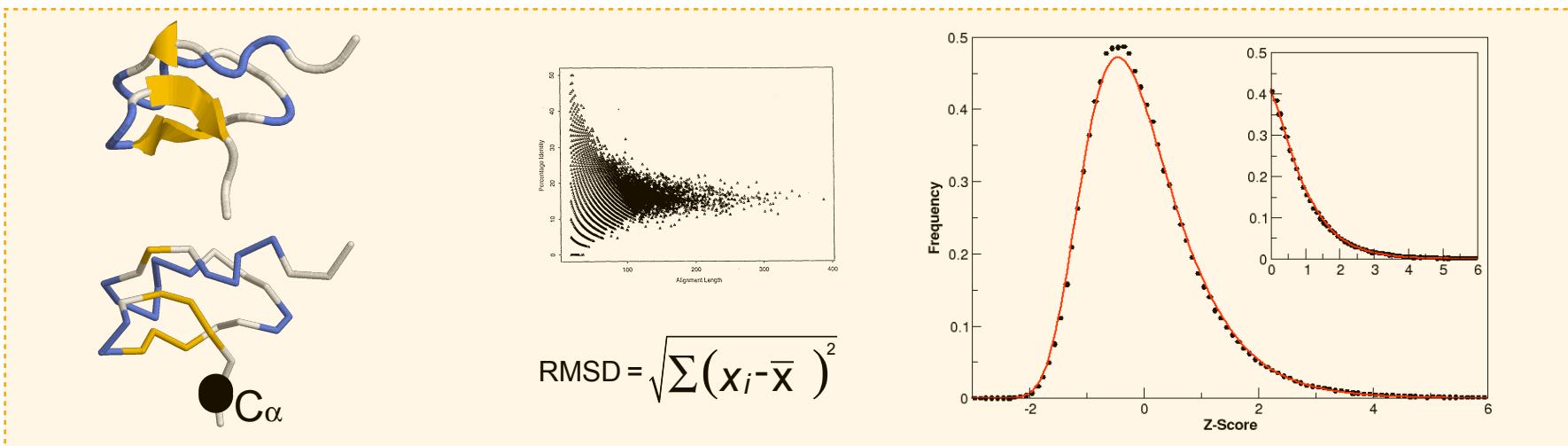
SALIGN - a multiple structure/sequence alignment tool, under preparation

# 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



# Vector Alignment Search Tool (VAST)

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

NCBI VAST Home Page - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Search Favorites Media Go Links

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

NCBI Structure PubMed Entrez BLAST OMIM Books TaxBrowser Entrez Structure

Search Entrez Structure for Go

VAST Help Vector Alignment Search Tool try:

Comprehensive help and frequently asked questions

VAST Search Submit structure database searches

VAST Search Help Help on submitting VAST Searches

VAST Search FAQ More help on VAST Search

Linking to VAST direct WWW access to the VAST server

nr-PDB non-redundant protein structure subsets

MMDB

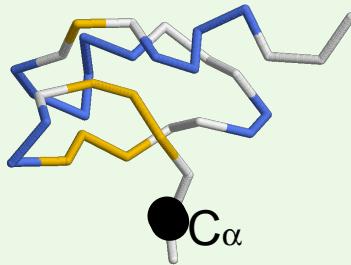
Protein structure neighbors in Entrez are determined by direct comparison of 3-dimensional protein structures with the VAST algorithm. Each of the more than 87,804 domains in MMDB is compared to every other one. From the MMDB Structure summary pages, retrieved via Entrez, structure neighbors are available for protein chains and individual structural domains. If you already know a PDB/MMDB-Id you can try this at once, using the input form in the right column.

On the Structure summary page, use "3d Domains" or "Protein" to retrieve a list of similar structures. For example, click on a bar with a chain identifier such as "B", or the bar below the Chain B with a domain identifier such as "1", to get a list of neighbors. The results of the precompiled VAST search will then present structural neighbor graphically. Using the check boxes in the leftmost column of this graph, select those structures you would like to see superimposed and click on "View 3D Structure" to view these with the mime-typed helper application you have installed (e.g., Cn3D).

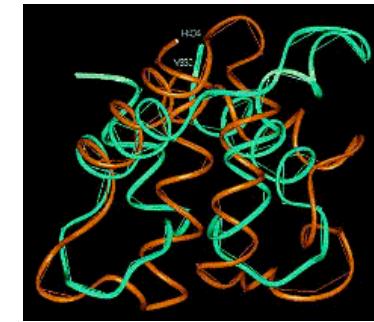
VAST Search is a service that allows searching for structural neighbors starting with a set of 3D-coordinates specified by the user. This service is meant to be used with newly determined protein structures that are not yet part of MMDB. Structure neighbors for proteins already in MMDB have been pre-computed and can simply be looked up from MMDB's Structure Summary via PDB/MMDB Code Get

Install and test structure alignment viewers:  
Get Cn3D v4.1 and look at this example to test! Read a bit more about VAST...

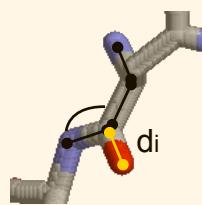
# Incremental combinatorial extension (CE)



- Exhaustive combination of fragments
- Longest combination of AFPs
- Heuristic similar to PSI-BLAST

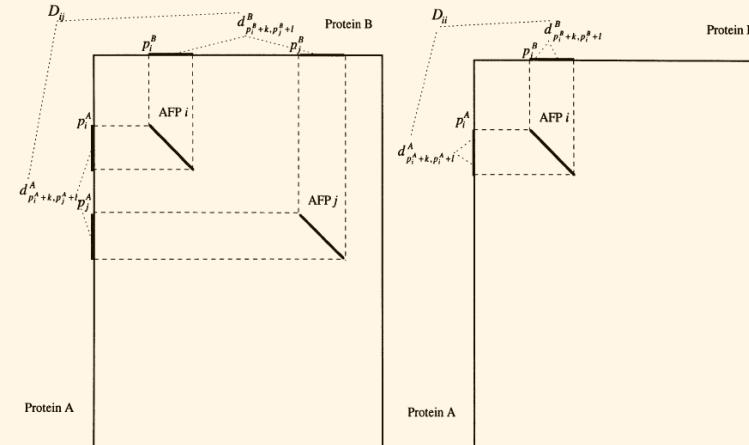


✓ FAST!  
✓ Good quality of local alignments  
**Complicated scoring and heuristics**



8 residues peptides

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

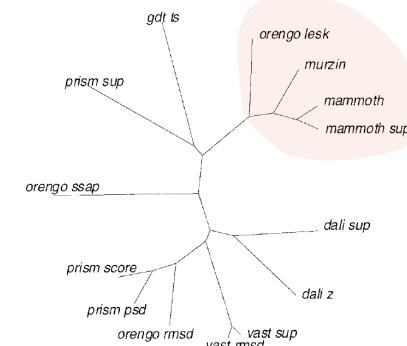
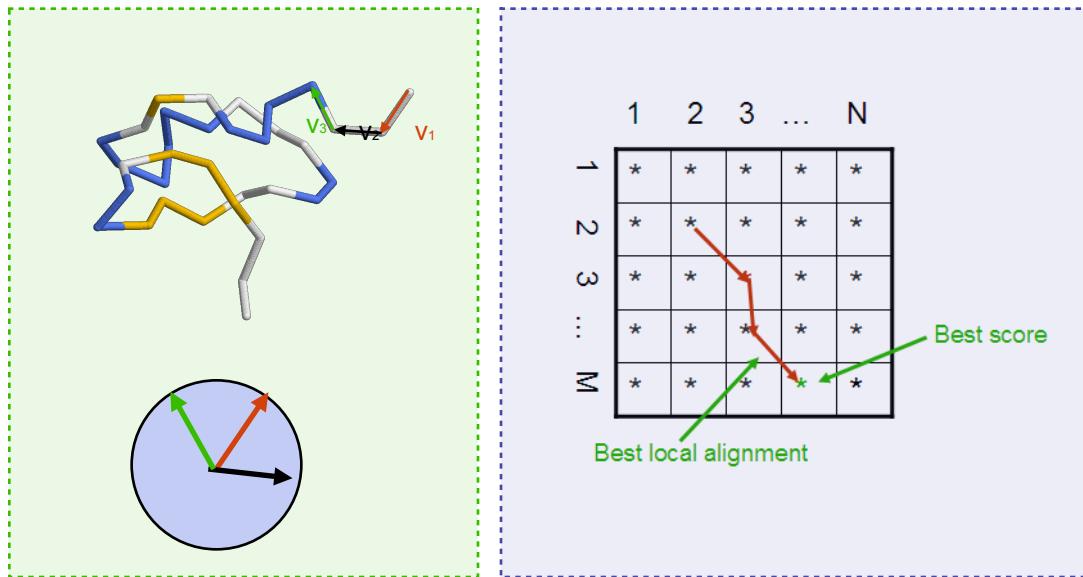


# Incremental combinatorial extension (CE)

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

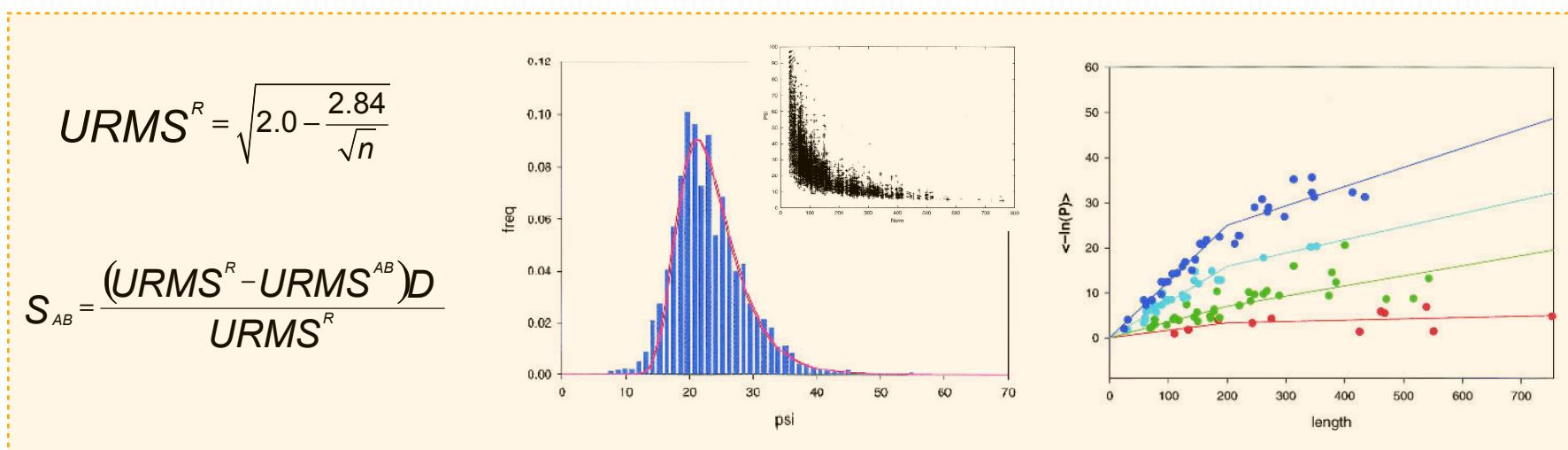
The screenshot shows a Microsoft Internet Explorer window with the title "CE Home Page - Combinatorial Extension - Microsoft Internet Explorer". The address bar contains the URL "http://cl.sdsc.edu/ce.html". The main content area features a protein structure visualization with colored sticks representing atoms. A specific residue, Y335, is highlighted in yellow and labeled "HO". Below the visualization, text reads: "Structural similarity between Acetylcholinesterase and Calmodulin found using CE (Tsigelny et al., *Proc Sci*, 2000, 9:180)". To the right, the text "Databases and Tools for 3-D Protein Structure Comparison and Alignment" and "Using the Combinatorial Extension (CE) Method" is displayed. At the bottom left, there are two large buttons: "FIND" and "CALCULATE". The "FIND" button is associated with text: "Find structural alignments by selecting from [ALL](#) or [REPRESENTATIVES](#) from the PDB.". The "CALCULATE" button is associated with text: "Calculate structural alignment for [TWO CHAINS](#) either from the PDB or uploaded by the user. Calculate structural neighbors for one protein [UPLOADED BY THE USER AGAINST THE PDB.](#)". There is also a link to "Calculate [MULTIPLE STRUCTURE ALIGNMENT](#)". A "More Info" link with a question mark icon is located near the bottom right. The browser interface includes standard menu bars (File, Edit, View, Favorites, Tools, Help), toolbars, and status bars at the bottom.

# Matching molecular models obtained from theory (MAMMOTH)



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

Reduces the protein representation



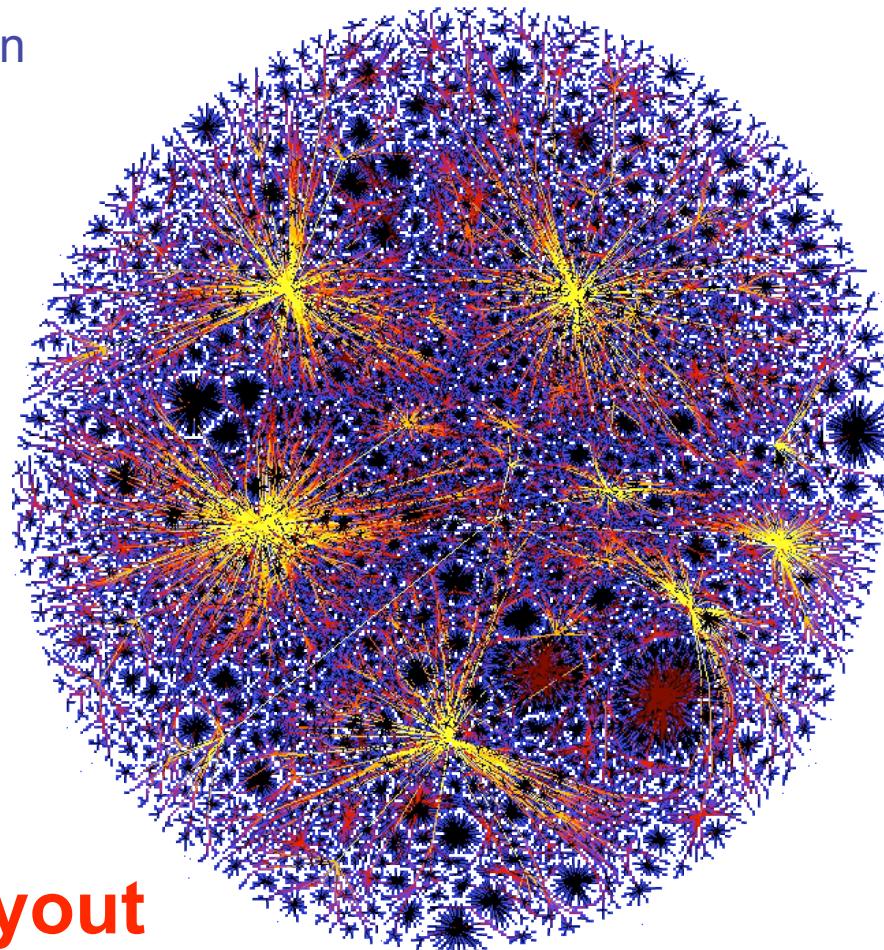
# Matching molecular models obtained from theory (MAMMOTH)

<http://fulcrum.physbio.mssm.edu:8083/>

The screenshot shows a Microsoft Internet Explorer window displaying the MAMMOTH web application. The title bar reads "Protein Structure Alignment Server - Microsoft Internet Explorer". The address bar contains the URL "http://fulcrum.physbio.mssm.edu:8083/mammoth/". The main content area features a logo of a mammoth on the left and the text "MAMMOTH" and "MAatching Molecular Models Obtained from THeory" in the center. Below this, there are two input fields: one for "PREDICTION" coordinates (PDB format) and one for "EXPERIMENT" coordinates (PDB format), each with a "Browse..." button. The bottom status bar shows "Done" and "Internet".

# Classification of the structural space

SCOP classification



## Large Graph Layout

Alex Adai

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

<http://bioinformatics.icmb.utexas.edu/lgl/>

# SCOP 1.65 database

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

The screenshot shows the SCOP homepage. At the top, there's a navigation bar with links for File, Edit, View, Favorites, Tools, and Help. The address bar shows the URL. Below the bar, there's a logo with three icons: a square, a triangle, and a question mark. The main content area has a title "Structural Classification of Proteins". Below it, a welcome message says: "Welcome to SCOP: Structural Classification of Proteins. 1.65 release (December 2003). 20619 PDB Entries. 1 Literature Reference. 54745 Domains (excluding nucleic acids and theoretical models). Folds, superfamilies, and families [statistics here](#). New folds [superfamilies](#) [families](#). [List of obsolete entries and their replacements](#)". A section for "Authors" lists Alexey G. Murzin, Loredana Lo Conte, Antonina Andreeva, Dave Howorth, Bartlett G. Ailey, Steven E. Brenner, Tim J. P. Hubbard, and Cyrus Chothia, with an email address [scop@mrc-lmb.cam.ac.uk](mailto:scop@mrc-lmb.cam.ac.uk). A "Reference" section cites Murzin et al. (1995) and Murzin et al. (2002). A "Major changes" section discusses updates from version 1.63 to 1.65. A "Access methods" section lists various ways to search and browse the database. A note at the bottom says: "SCOP [mirrors](#) around the world may speed your access."

## News

- SCOP has been updated to include all PDB entries released up to 1 August 2003. See [folds, superfamilies, and families statistics](#).
- Several parts of the SCOP classification have been restructured, especially in this release and in the previous one. You can browse the subset of the classification affected by these changes in a SCOP-view form for modifications occurred between [1.63 and 1.65](#), or [previous releases](#). Changes appear as comments associated to [domain entries](#), with links to the revised classification. You can use the SCOP navigation buttons to move up in the hierarchy and to expand or collapse entries. The list of [obsolete entries and their replacements](#) is also available online.
- SCOP identifiers now appear explicitly in the web pages (in [squared brackets](#)).
- Links from a SCOP domain to the corresponding SWISSPROT and EC entries have been added (see the [l icon](#)). Thanks to Sameer Velankar and Phil McNeil from the EBI-MSD group and to Virginie Mittard from the EBI sequence database group for providing the most up-to-date map between PDB chains and SWISSPROT, EC identifiers.
- It is now possible to use SSM to search the up-to-date PDB archive using a SCOP domain entry (via the [l icon](#)) or to

- ✓ 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

# CATH<sub>2.5.1</sub> database

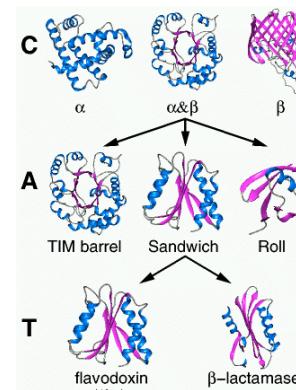
<http://www.biochem.ucl.ac.uk/bsm/cath/>

The screenshot shows the CATH Protein Structure Classification Database (UCL) interface. The top menu includes File, Edit, View, Favorites, Tools, and Help. The address bar shows the URL. The main content area displays the CATH Protein Structure Classification, version 2.5.1, released in January 2004. It features a search bar, a 'Goto' section with links to SSAP Server, GRATH Server, DHS, and Gene3D, and a 'Navigation' section with links to Home and Top of hierarchy. The central text discusses the hierarchical classification of protein domain structures at four major levels: Class(C), Architecture(A), Topology(T), and Homologous superfamily(H). It also provides a detailed explanation of the architecture level and lists references by Orengo et al. (1997) and Pearl et al. (2000).

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

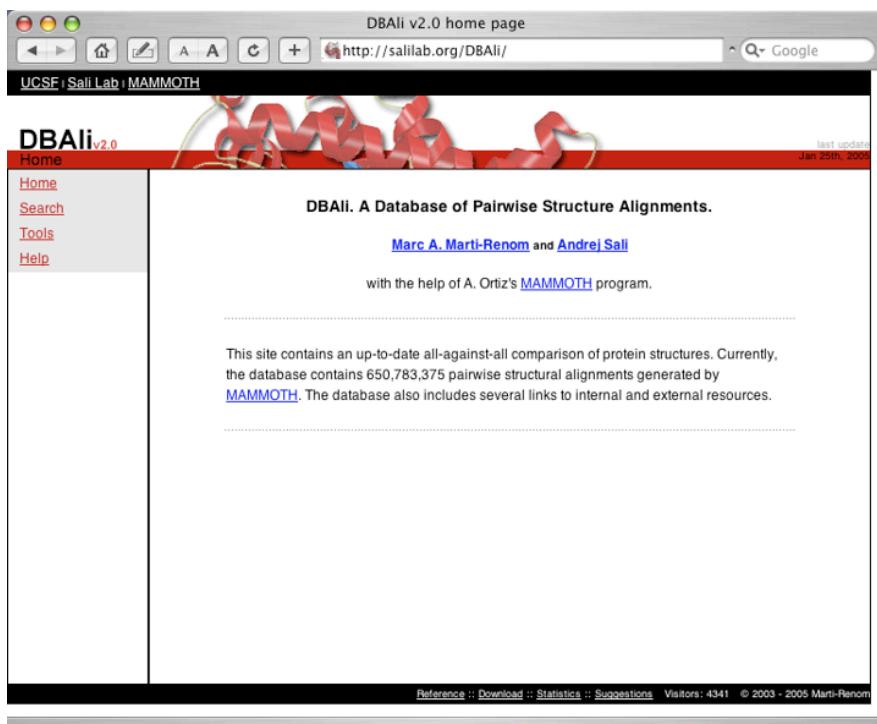


Version	2.5.1						
Date	28-01-2004						
	A	T	H	S	N	I	D
Mainly Alpha	5	227	428	948	1713	3946	10155
Mainly Beta	19	139	292	951	2344	5011	14259
Alpha Beta	12	368	648	2010	3631	8639	23025
Few Secondary Structures	1	86	91	114	225	378	952
Multi-domain chains	1	1053	1057	1071	2186	5801	12471
Preliminary single domain assignments	1	371	374	422	479	789	1663
Multi-domain domains	2	31	31	49	67	139	287
CATH-35 Sequence families	1	997	997	997	1108	2154	3431
Fragments from multi-chain domains	1	28	28	30	33	56	106

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

# DBAli v2.0 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.
- ✓ Easy to navigate
- ✓ Provides some tools for structure comparison

Does not provide (yet) a stable classification

Last updated:

January 25th, 2005

Number of chains in database:

60,656

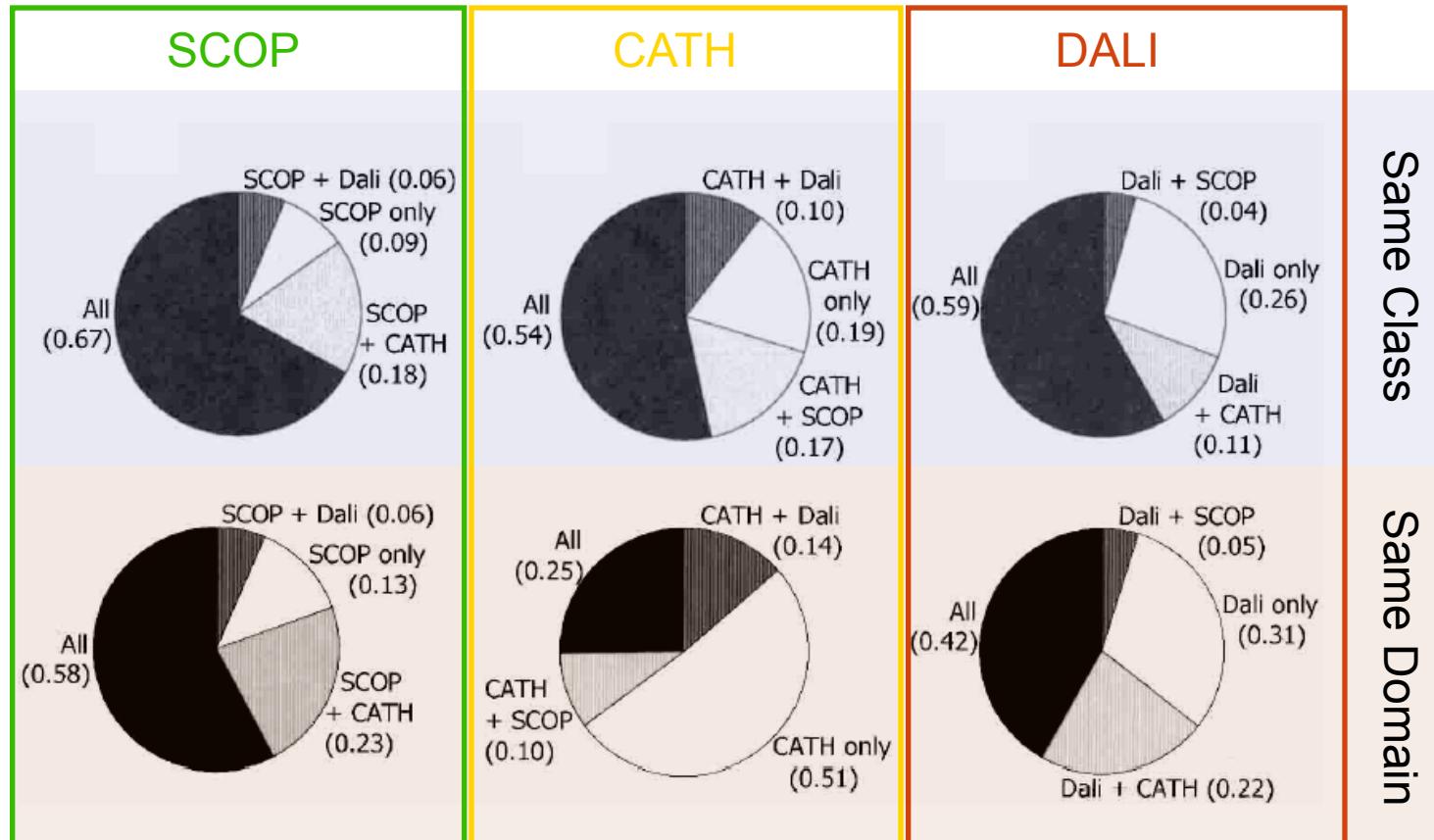
Number of structure-structure comparisons:

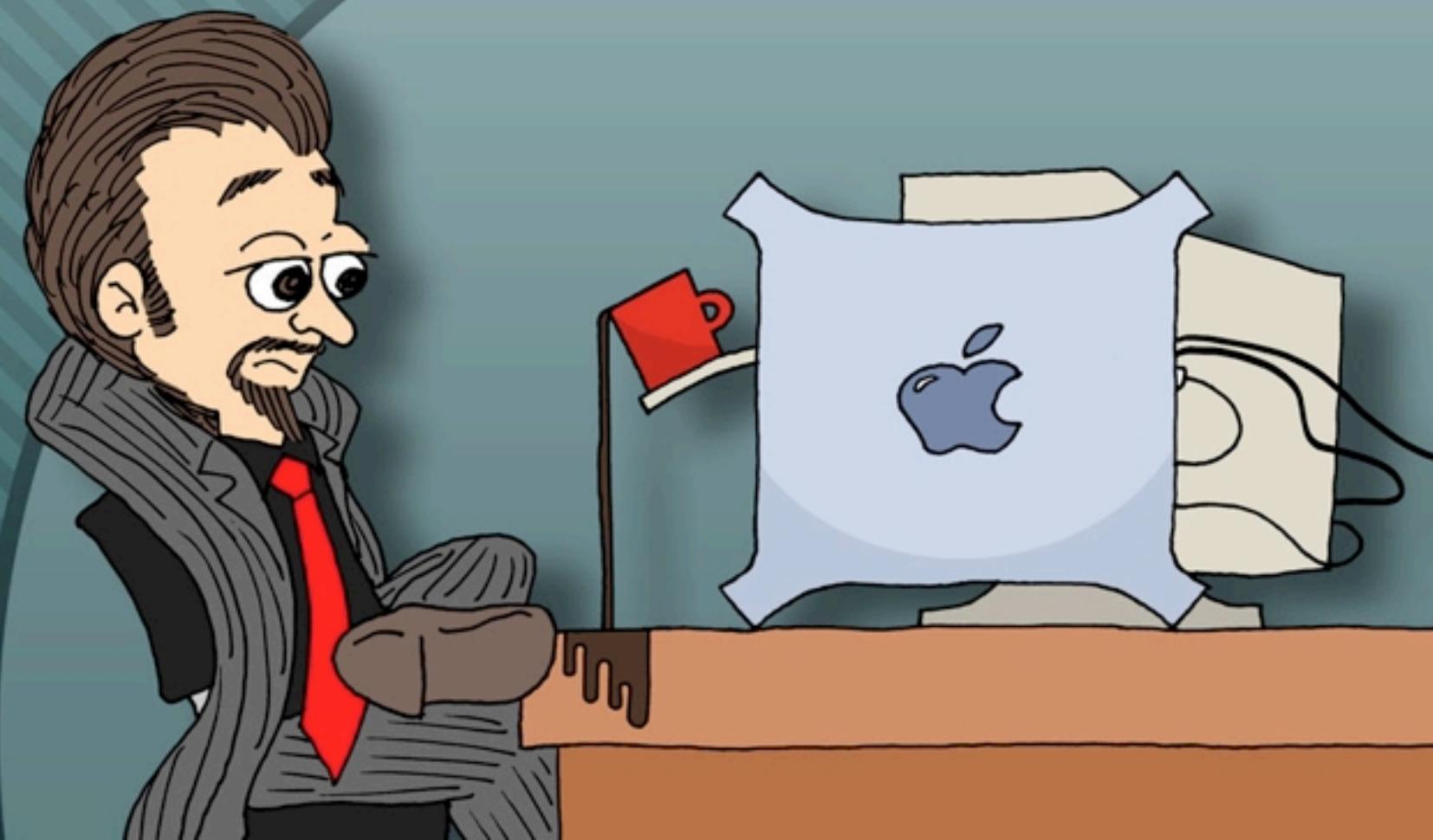
650,783,375

# Classification of the structural space

## *Not an easy task!*

Domain definition AND domain classification





# Sequence-Structure comparison

- Outline
  - Before we start...
    - Some theory...
    - Domain boundaries
  - Structural predictions from sequence...
    - SALIGN (gap penalties and substitution matrices)
    - mGenThreader (SSE prediction and alignment/potential scores)
    - Fugue (gap penalties and substitution matrices)
    - 3D-Jury (as a meta server example)

# General overview (Threading)

- Matches sequences to 3D structures
  - Requires a scoring function to asses the fit of a sequence to a given fold
  - Scoring functions deried from known structures and include atom contact and solvation terms evaluated in a pairwise fashion
  - May include secondary structure terms, multiple alignments...
- Threading servers available using several different approaches
  - Fold recognition server at Imperial College, UK  
<http://www.sbg.bio.ic.ac.uk/~3dpssm/>
  - PredictProtein server at EMBL  
<http://www.embl-heidelberg.de/predictprotein/predictprotein.html>
  - Protein sequence-structure threading at NCBI  
<http://www.ncbi.nlm.nih.gov/Structure/RESEARCH/threading.shtml>

# Template comparison methods

- Uses 3D “templates” for searching structural databases
  - active site or binding site templates generated to reflect functionally important structural signatures
- Available software/servers
  - Template Search and Superposition (TESS), Thornton Group  
<http://www.biochem.ucl.ac.uk/bsm/PROCAT/PROCAT.html>  
Wallace AC; Borkakoti N; Thornton JM. (1997) *Protein Science* **6** pp2308
  - “Fuzzy Functional Forms”, Skolnick - commercial availability  
Fetrow, JS and Skolnick, J (1998) *J. Mo. Biol.* **281** pp949
  - Spatial Arrangements of Side-chain and Main-chain (SPASM), Kleywegt, Univ. of Uppsala  
<http://portray.bmc.uu.se/cgi-bin/dennis/spasm.pl>  
Kleywegt GJ (1999). *J. Mol. Biol.* **285** pp1887

# 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.

# Sequence-Structure alignments

**As any other bioinformatics problem...**

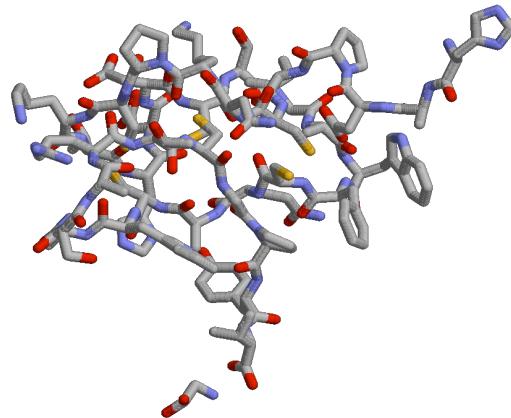
- Representation
- Scoring
- Optimizer

Representation

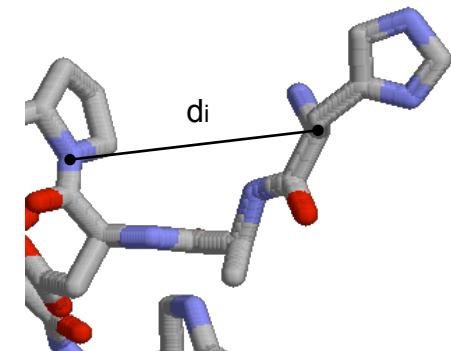
# Sequence/Structures

>gi42541361  
MDIRSVSSLRGILLCLPPSWPRR

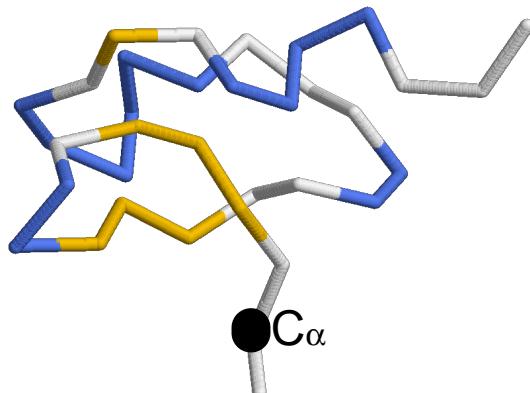
Primary sequence



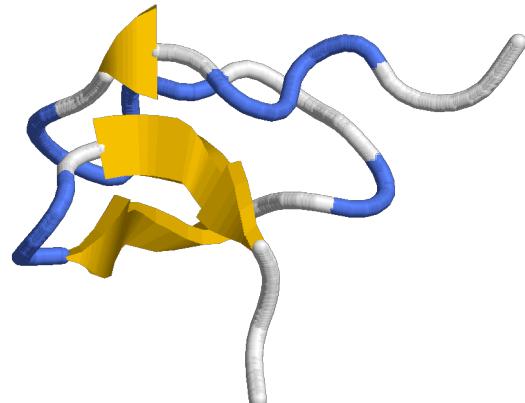
All atoms and coordinates



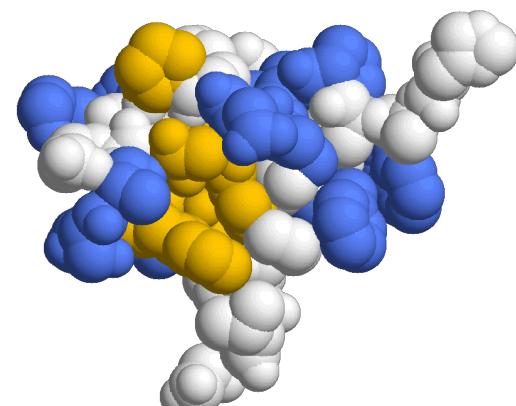
Distance space



Reduced atoms representation



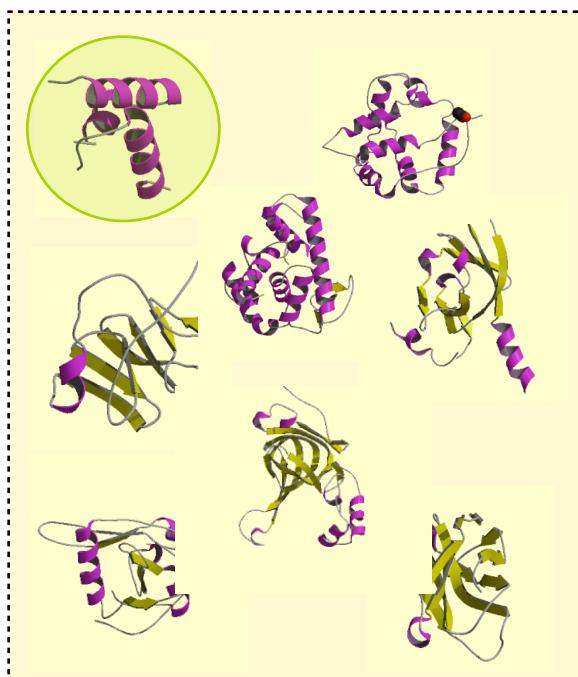
Secondary Structure



Accessible surface

# Scoring Statistical Potentials (background)

Structural space



Sequence space

MKLLIVLTCISLCSCICTVVQRCASNKPHVLEDPCKVQH  
HLSVNQCVLLPQCCPKSCKICTHLISIEVVLTCRAVDKM  
MHVNCVEQCSLQDCIKIAPRVLKTCILCVLKPCLTSH  
VHLVQPTSCCCKNCICHVEIRSLDILTKSVQLACLVPM  
⋮  
MQCCRVQKICDLLAVELCKLHISTPSCKILCVVTSVPHN

# Scoring Statistical Potential (inspiration)

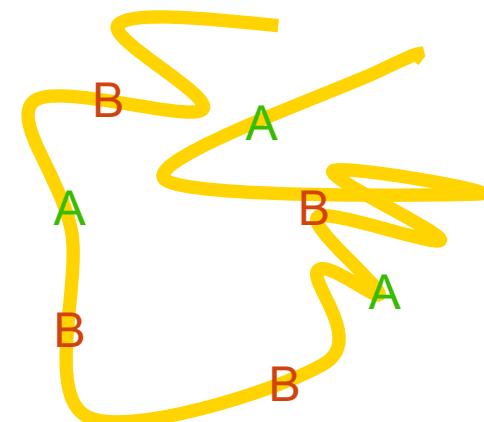
$$K = \frac{[AB]}{[A] \cdot [B]}$$

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

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]:

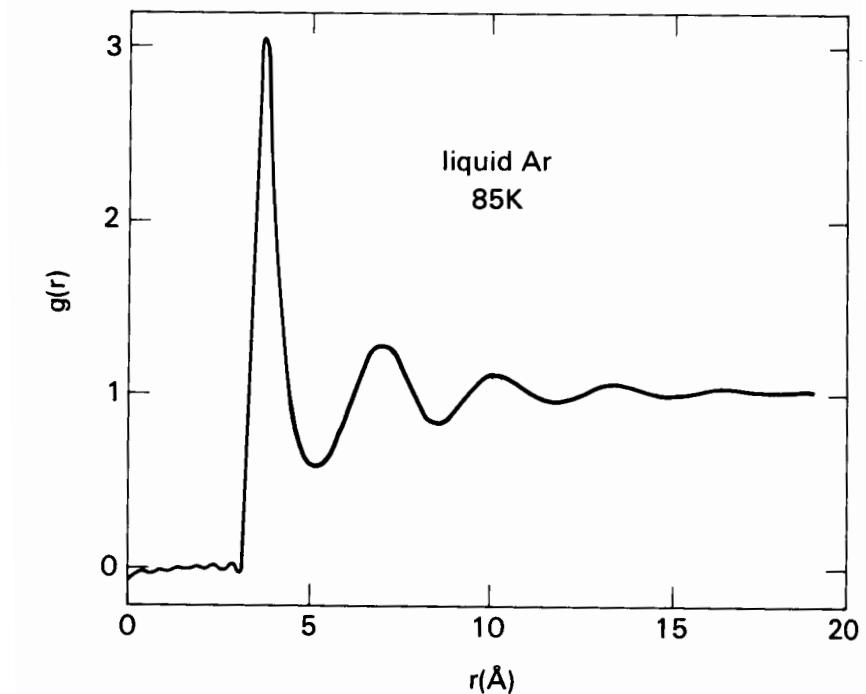
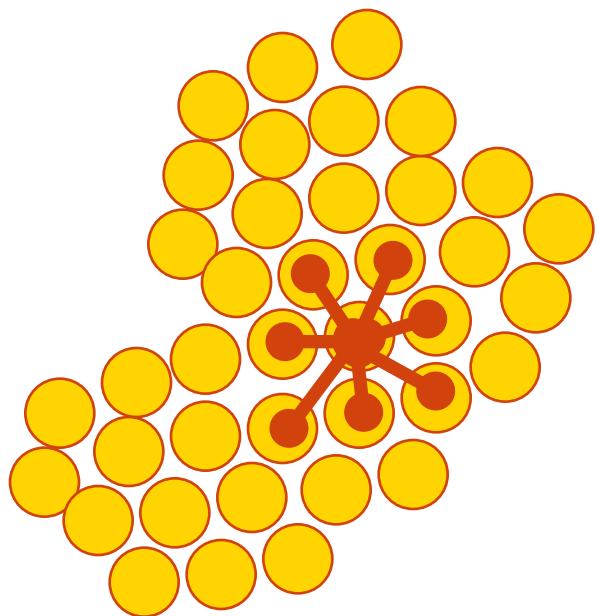
$$\Delta E = -kT \ln \left( \frac{N_1}{N_2} \right) \quad (1)$$

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.



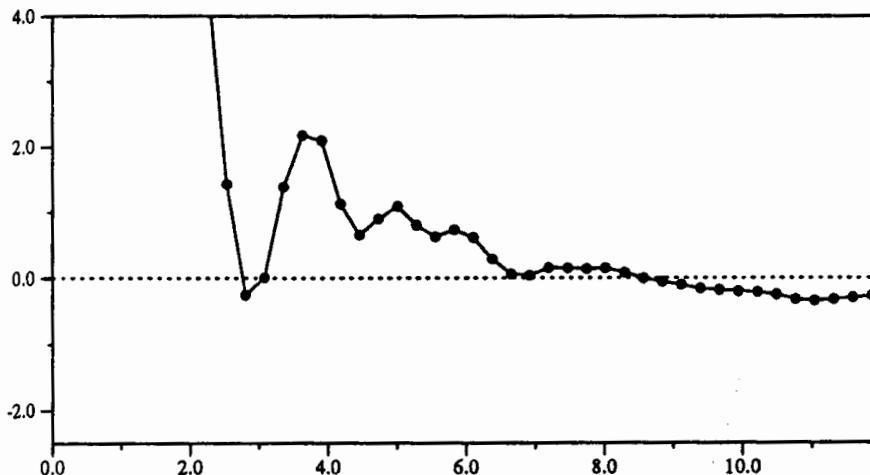
Tanaka and Sheraga (1975) PNAS, 72 pp3802  
 Sippl, (1990) J.Mo.Biol. 213 pp859  
 Godzik, (1996) Structure 15 pp363

# Scoring Statistical Potential (reference state)

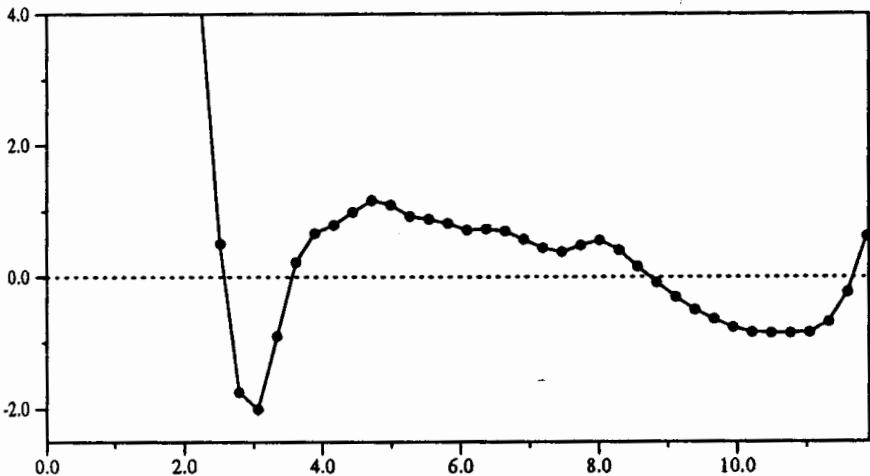


# Scoring Statistical Potential... Hydrogen Bonds

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

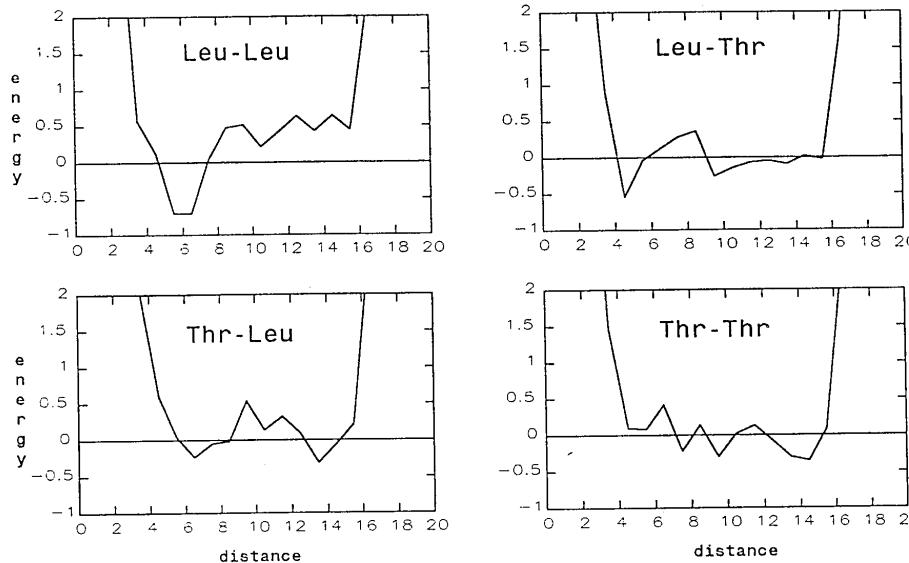
$$\rho_{NO}(r) = \sum_{ij} \delta(\mathbf{r} - \mathbf{r}_{ij})$$

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

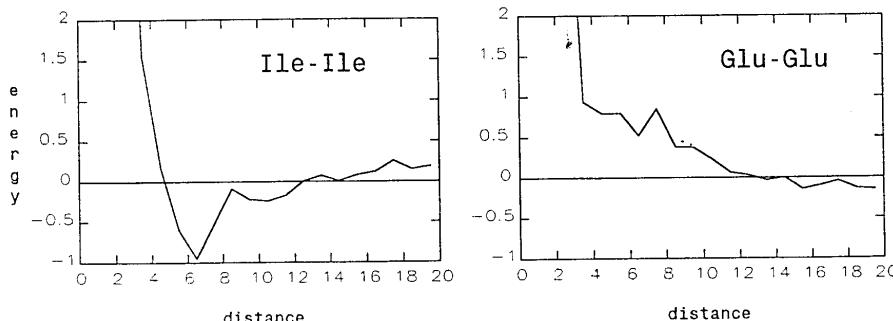
$$W_{NO}(r) = -kT \ln(g_{NO}(r))$$

# Scoring Statistical Potential... Distance Potentials

Long range free energy



Short range free energy



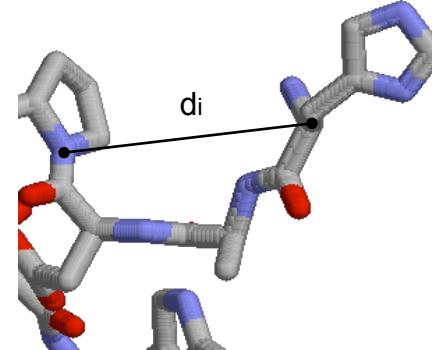
# Scoring

## Raw scores of an alignment

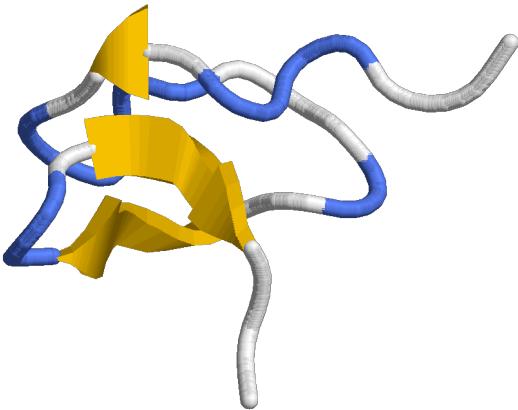
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
C	9	-1	-1	-3	0	-3	-3	-3	-4	-3	-3	-3	-3	-1	-1	-1	-1	-2	-2	-2	
S	-1	4	1	-1	1	0	1	0	0	0	-1	-1	0	-1	-2	-2	-2	-2	-3		
T	-1	1	4	1	-1	1	0	1	-1	-1	-1	-1	-1	-1	-2	-3	-3	-2	-3		
P	-3	-1	1	7	-1	-1	-2	-1	-1	-1	-1	-2	-2	-1	-2	-3	-3	-2	-4	-3	
A	0	1	-1	-1	4	0	-1	-2	-1	-1	-1	-2	-1	-1	-1	-1	-1	-2	-2	-3	
G	-3	0	1	-2	0	6	-2	-1	-2	-2	-2	-2	-2	-3	-4	-4	0	-3	-3	-2	
N	-3	1	0	-2	-2	0	6	1	0	0	-1	0	0	-2	-3	-3	-3	-2	-4		
D	-3	0	1	-1	-2	-1	1	6	2	0	-1	-2	-1	-3	-3	-4	-3	-3	-4		
E	-4	0	0	-1	-1	-1	-2	0	2	5	2	0	0	1	-2	-3	-3	-3	-2	-3	
Q	-3	0	0	-1	-1	-2	0	0	2	5	0	1	1	0	-3	-2	-2	-3	-1	-2	
H	-3	-1	0	-2	-2	-2	1	1	0	0	8	0	-1	-2	-3	-3	-2	-1	2	-2	
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5	2	-1	-3	-2	-3	-3	-2	-3	
K	-3	0	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5	-1	-3	-2	-3	-3	-2	-3
M	-1	-1	-1	-2	-1	-3	-2	-3	0	-2	-1	-1	-1	5	1	2	-2	0	-1	-1	
I	-1	-2	-2	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4	2	1	0	-1	-3	
L	-1	-2	-2	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4	3	0	-1	-2	
V	-1	-2	-2	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4	-1	-1	-3	
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6	3	1	
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	-2	-2	-2	-1	-1	-1	-1	3	7	2	
W	-2	-3	-3	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11	

2/

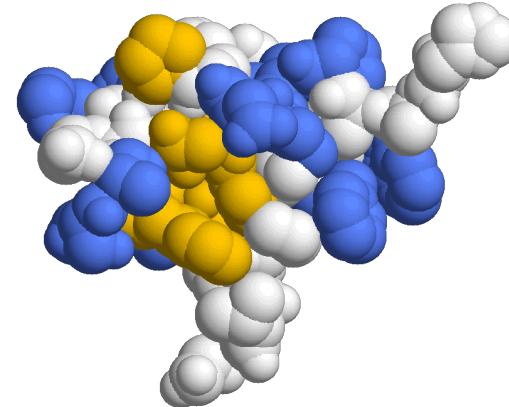
Aminoacid substitutions



Distance space



Secondary Structure (H,B,C)



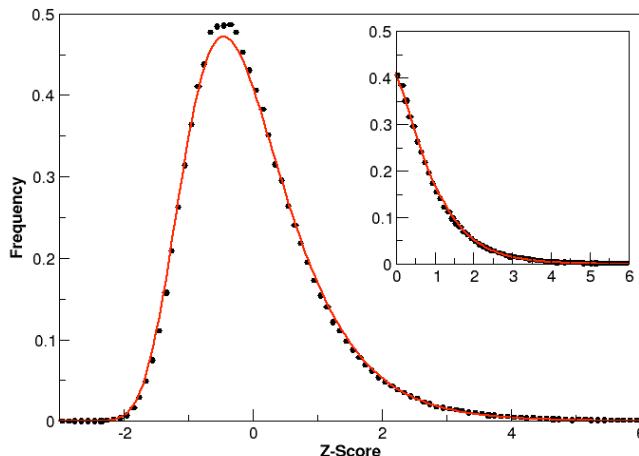
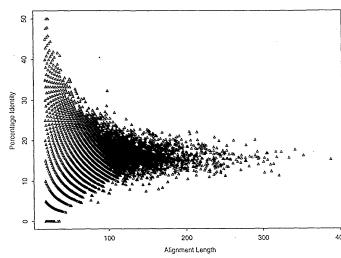
Accessible surface (B,A [%])

# Scoring

## Significance of an alignment (score)

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

Empirical



Sometimes approximated by Z-score (normal distribution).

Analytic

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

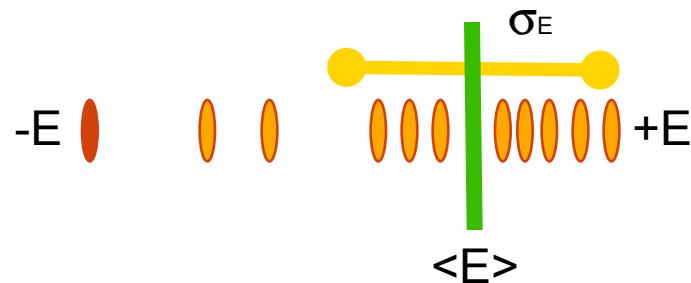
$$P(s \geq x) = 1 - \exp(e^{-\lambda (s-\mu)})$$

Karlin and Altschul, 1990 PNAS 87, pp2264

## Scoring

# Significance of an alignment (score)

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

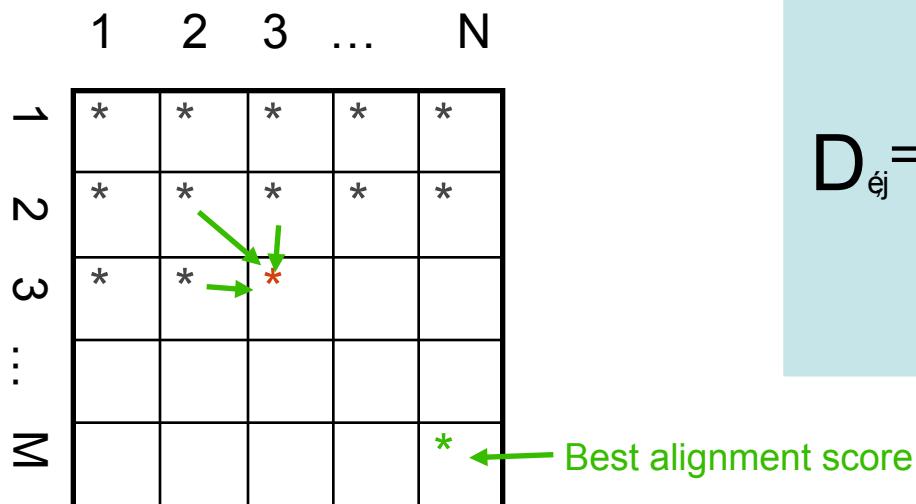


$$Zscore = \frac{(\langle E \rangle - E_m)}{\sigma_E}$$

## Optimizer

# Global dynamic programming alignment

remember Patsy's class



$$D_{ej} = \min \begin{cases} D_{i,j-1} + \text{Score}_{(\ddot{A}, rj)} \\ D_{i-1,j-1} + \text{Score}_{(ri, rj)} \\ D_{i-1,j} + \text{Score}_{(ri, \ddot{A})} \\ 0 \end{cases}$$

Backtracking to get the best alignment

# Applications of PMFs

Model assessment.

*Ab initio* folding simulations.

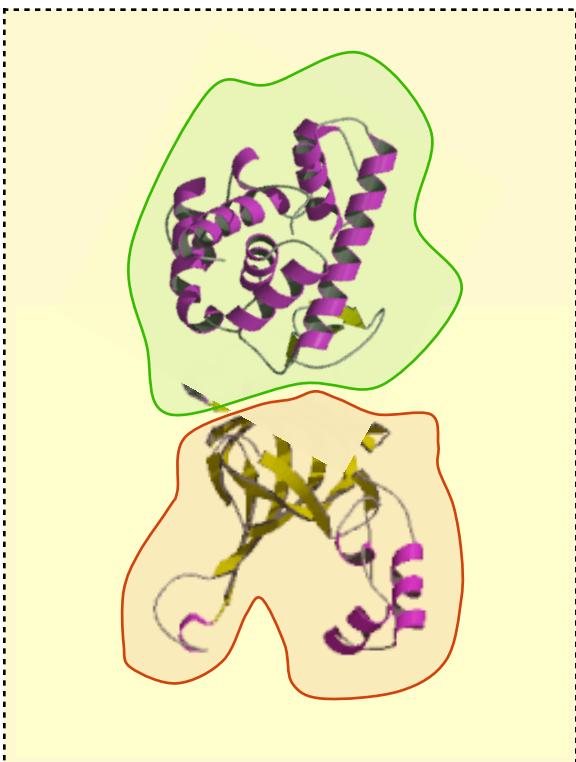
Sequence-structure matching (threading).

Comparative protein structure modeling  
(loops, sidechains, ...).

Secondary structure prediction, *etc.*

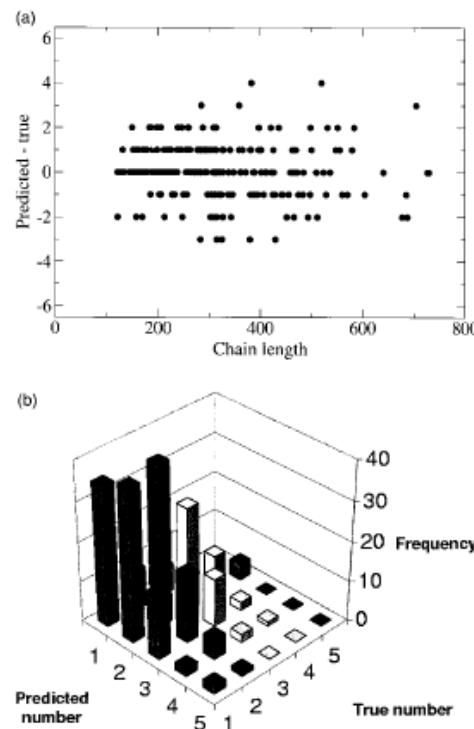
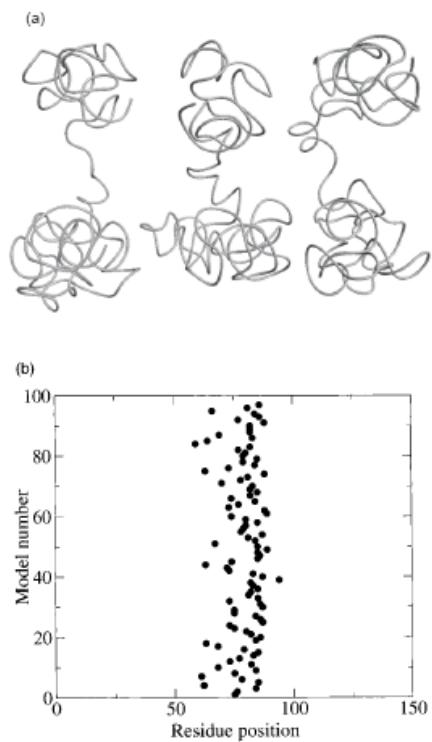
# Domain boundaries from sequence

**VERY DIFFICULT!!!!**



MENFEIWVEKYRPTLDEVVGQDEVIQRLKGYVERKNIPHLLFSGPPGTGKTATAIALARDLFGENWRDN  
FIEMNASDERGIDVVRHKIKEFARTAPIGGAPFKIIFLDEADALTADAQAAALRTMEMYSKSCRFLSCN  
YVSRIIEPIQSRCAVFRFKPVPKEAMKKRLLEICEKEGVKITEDGLEALIYISGGDFRKAINALQGAAAI  
GEVVVDADTIYQITATARPEEMTELIQTALKGNFMEARELLDRLMVEYGMMSGEDIVAQLFREIISMPIKDS  
LKVQLIDKLGEVDFRLTEGANERIQLDAYLAYLSTLAKK

# Domain boundaries from sequence (SnapDragon)

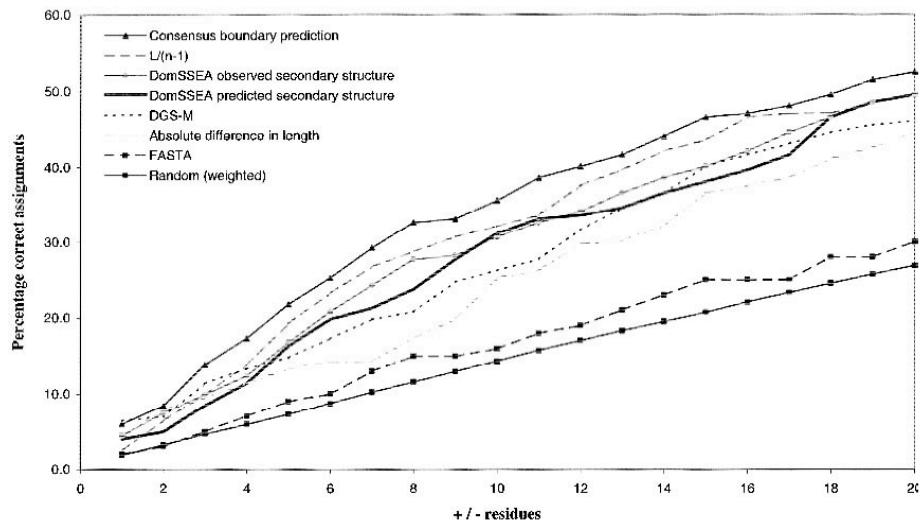
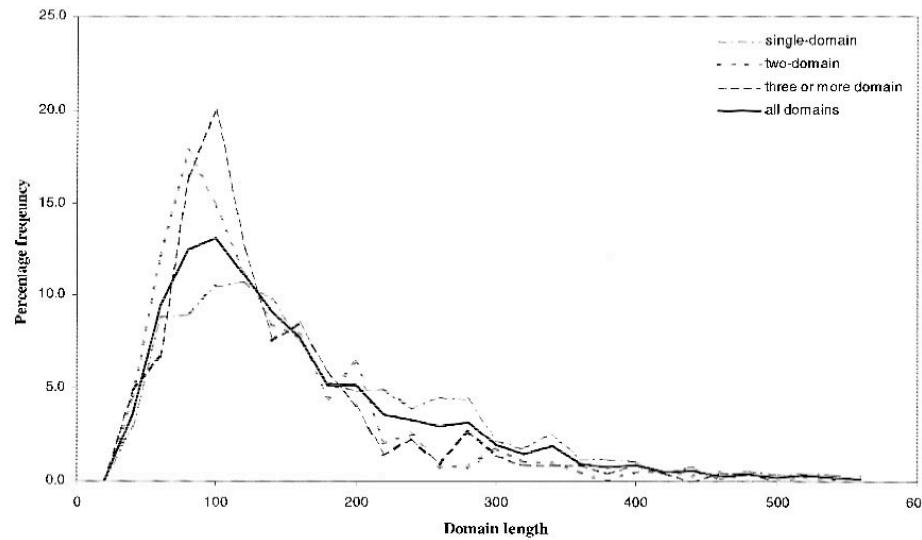


**Table 2.** Average accuracy percentages of linker prediction over 57 proteins

		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)

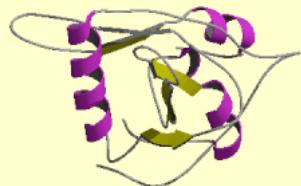
Methods	% Correctly assigned	
	All chains	Multidomain chains
DomSSEA observed secondary structure	70.2	24.7
DomSSEA predicted & consensus	68.6	24.0
DomSSEA predicted & $L/(N-1)$	68.0	24.0
DomSSEA predicted secondary structure	68.7	23.6
Absolute difference in length	62.0	8.4
Average domain length & DGS-M	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



# Prediction of Secondary Structure (PSI-PRED)

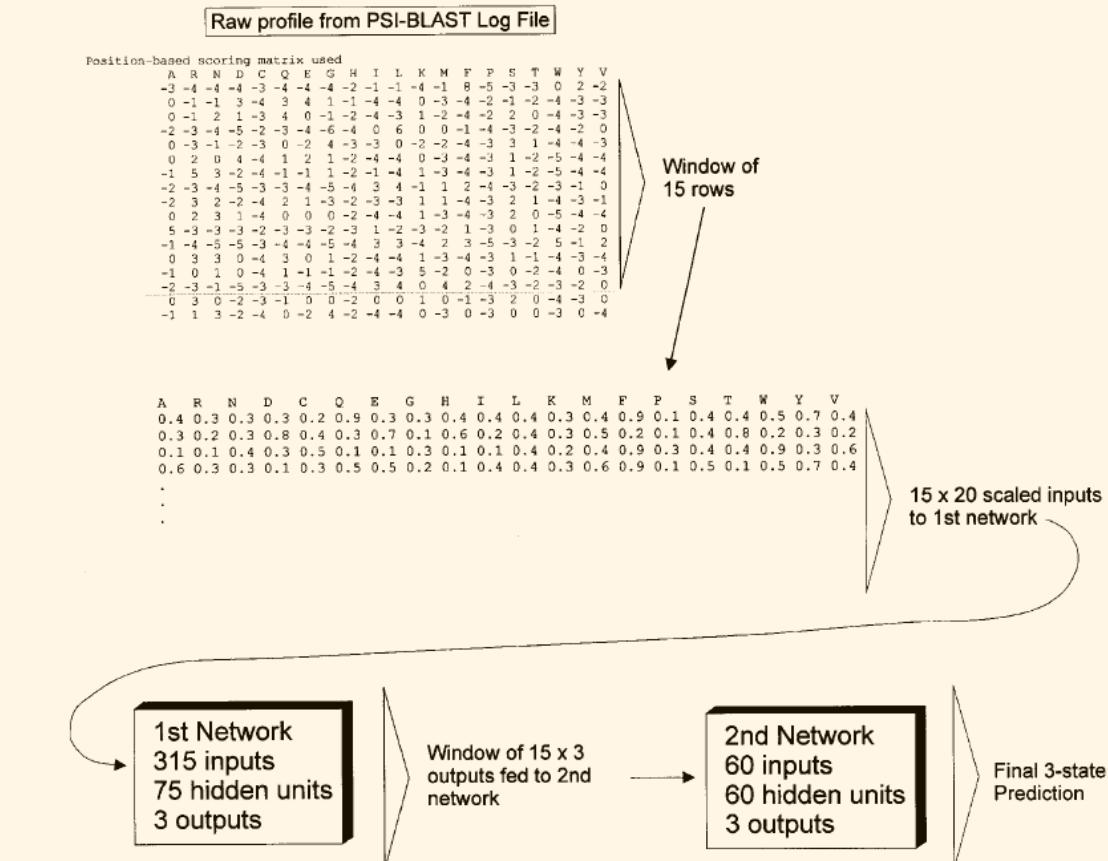
>gi42541361  
MDIRSVSSLRGLLCLPPSWPRR

- Neural Network



- ✓ Very simple idea
- ✓ Simple scoring

Obscure optimizer



# Prediction of Secondary Structure (PSI-PRED)

<http://bioinf.cs.ucl.ac.uk/psiform.html>

PSIPRED Protein Structure Prediction Server - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Search Favorites Media Mail Links

Address http://bioinf.cs.ucl.ac.uk/psiform.html

**Bioinformatics Unit**

**PSIPRED home>**

**The PSIPRED Protein Structure Prediction Server**

**Info** We suggest that you do not bookmark this page as it is liable to move. It is best to access the server via the [PSIPRED home page](#), which has more information about the methods and a full reference list.

**Input Sequence** [Help](#)  
Input sequence (single letter code)

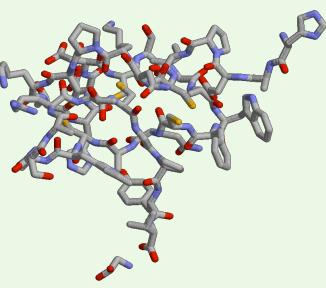
**Choose Prediction Method** [Help](#)  
 Predict Secondary Structure (PSIPRED v2.4)  
 Predict Transmembrane Topology (MEMSAT)  
 Fold Recognition(GenTHREADER - quick)  
 Fold Recognition (mGenTHREADER - with profiles and predicted secondary structure)

**Filtering Options** [Help](#)  
 Mask low complexity regions  
 Mask transmembrane helices  
 Mask coiled-coil regions  
Warning: Turn off all filtering if you are running MEMSAT

**Submit Sequence** E-mail address [Help](#)  
Password (only required for commercial e-mail addresses) [Help](#)  
Short name for sequence [Help](#)  
Predict Clear form

Internet

# Sequence-Structure alignment by properties conservation (SALIGN-MODELLER)

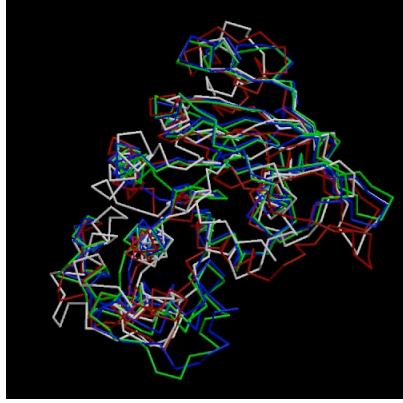


1	2	3	...	N
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*

Best score  
Best local alignment

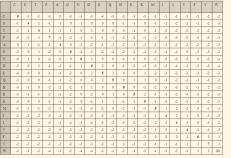
A (green)  
B (yellow)  
C (orange)  
D (blue)

- similarity +

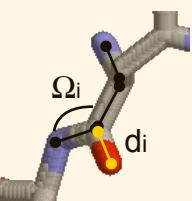


✓ Uses all available structural information  
 ✓ Provides the optimal alignment

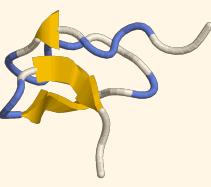
Computationally expensive



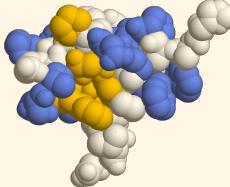
$R_{i,j}$



$D_{i(3),j(3)}$



$S_{i,j}$



$B_{i,j}$

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

$I_{i,j}$

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

# Sequence-Structure alignment by properties conservation (SALIGN-MODELLER)

<http://alto.compbio.ucsf.edu/salign-cgi/index.cgi>

SALIGN Server

http://alto.compbio.ucsf.edu/salign-cgi/index.cgi

Google

### SALIGN Multiple Structure/Sequence Alignment Server

SALIGN is a general alignment module of the modeling program MODELLER

The alignments are computed using dynamic programming, making use of several features of the protein sequences and structures

Users can either upload their own sequences/structures to align or choose structures from the PDB

Sequences can either be pasted or uploaded as FASTA or PIR format alignment files

Paste sequence to align

Multiple sequences can be pasted by iteratively clicking 'upload' after every pasted sequence

Specify file to upload (PIR, FASTA, PDB, zip or .tar.gz)

Multiple files can be uploaded by iteratively clicking 'upload' after every file uploaded

Localized string not found

Uploaded files:

No files uploaded

Enter 4 letter code(s) to choose PDB structures

e-mail address, to receive results:

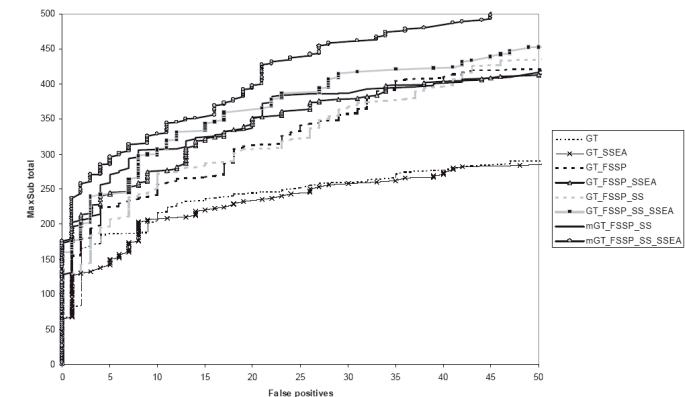
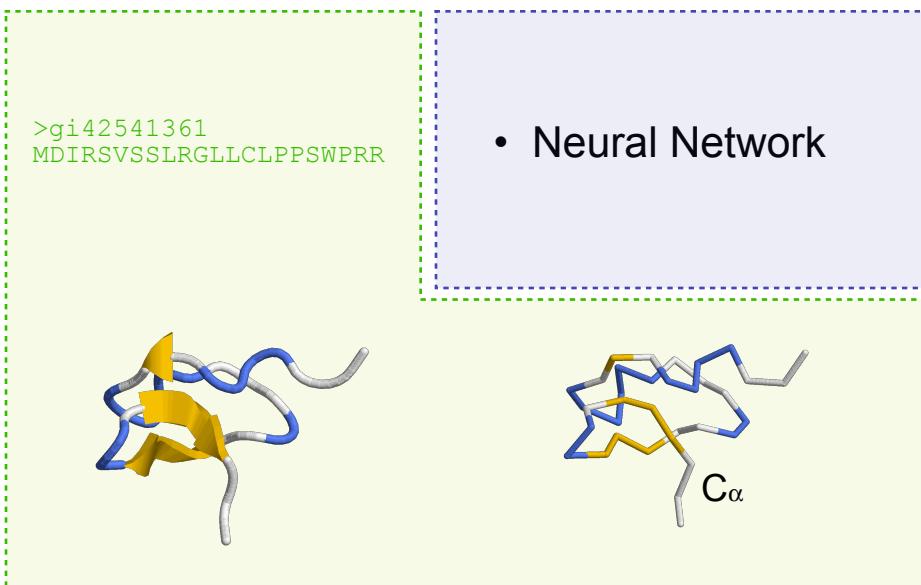
Localized string not found

Reference:

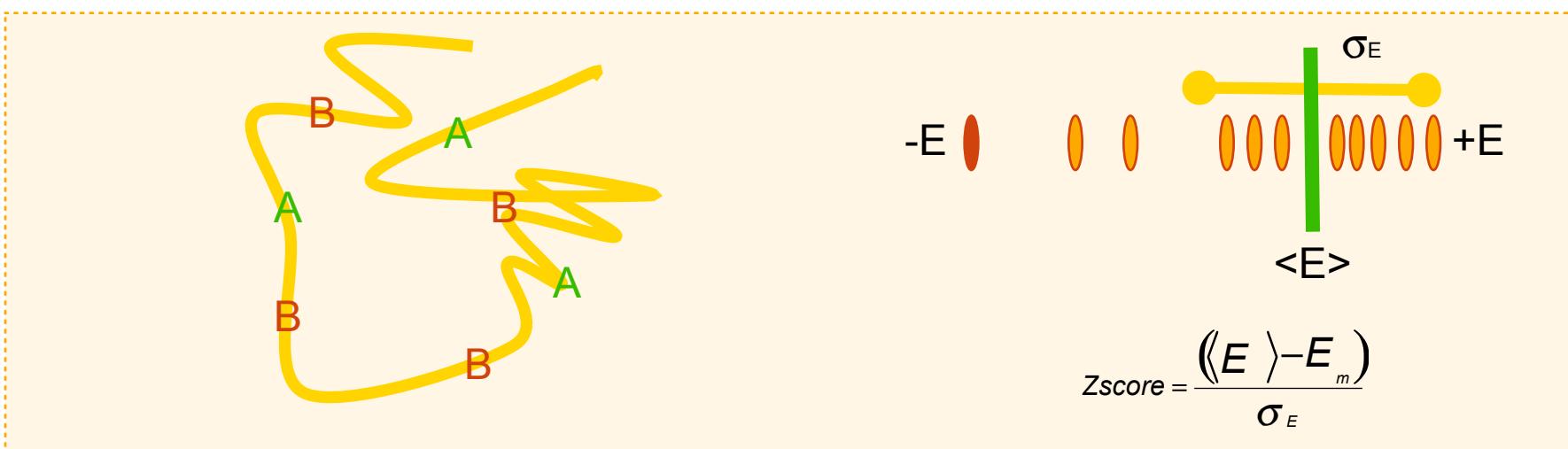
Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A.

SALIGN - a multiple structure/sequence alignment tool, under preparation

# Threading (mGenThreader)



✓ Good row and significance scoring  
Obscure optimizer



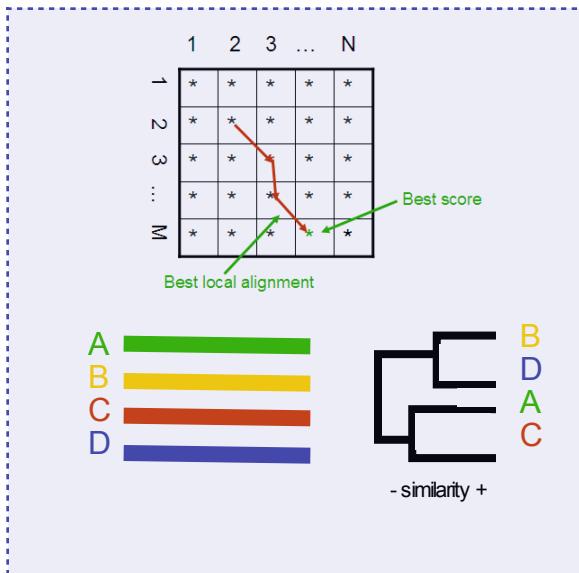
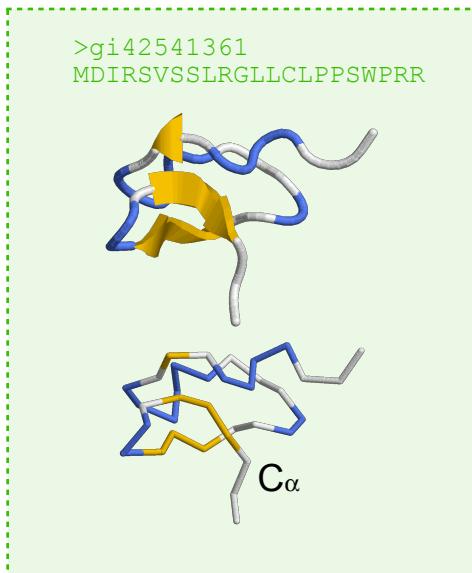
# Threading (mGenThreader)

<http://bioinf.cs.ucl.ac.uk/psiform.html>

The screenshot shows the PSIPRED Protein Structure Prediction Server interface. At the top, there's a banner for the Bioinformatics Unit at UCL. The main content area has several sections:

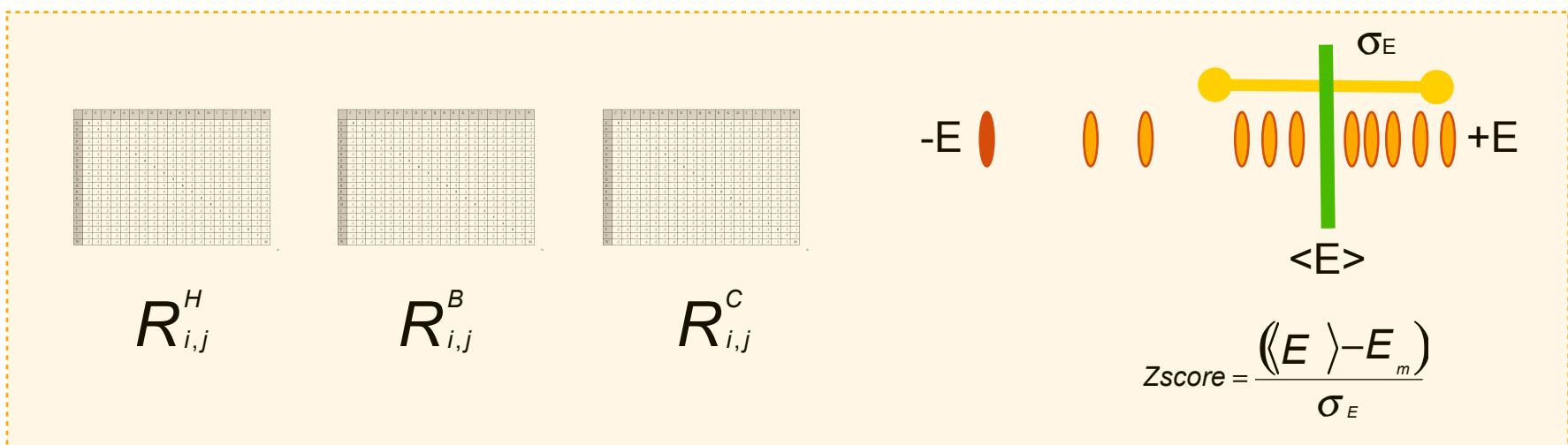
- PSIPRED home**: A link to the main PSIPRED page.
- Info**: A note suggesting not to bookmark this page as it is liable to move, and directing users to the [PSIPRED home page](#) for more information.
- Input Sequence**: A text input field labeled "Input sequence (single letter code)" with a help link.
- Choose Prediction Method**: A section with a red arrow pointing to it. It contains a "Help" link and four radio button options:
  - Predict Secondary Structure (PSIPRED v2.4)
  - Predict Transmembrane Topology (MEMSAT)
  - Fold Recognition(GenTHREADER - quick)
  - Fold Recognition (mGenTHREADER - with profiles and predicted secondary structure)
- Filtering Options**: A section with a "Help" link and three checkboxes:
  - Mask low complexity regions
  - Mask transmembrane helices
  - Mask coiled-coil regionsA warning message at the bottom of this section says: "Warning: Turn off all filtering if you are running MEMSAT".

# Remote homology detection (FUGUE)



- ✓ Uses most of the structural information
- ✓ Easy to access either locally and on the web
- ✓ Good row and significance scoring

Does not uses multiple sequence information



$$R_{i,j}^H$$

$$R_{i,j}^B$$

$$R_{i,j}^C$$

# Remote homology detection (FUGUE)

<http://www-cryst.bioc.cam.ac.uk/fugue/>

The screenshot shows a Microsoft Internet Explorer window displaying the FUGUE homepage. The title bar reads "FUGUE: sequence-structure homology recognition and alignment engine - Microsoft Internet Explorer". The address bar shows the URL "http://www-cryst.bioc.cam.ac.uk/fugue/". The page content includes the FUGUE logo, the text "Crystallography and Biocomputing Unit Department of Biochemistry, University of Cambridge", and a crest. Below this, a main heading states "Sequence-structure homology recognition using environment-specific substitution tables and structure-dependent gap penalties". A horizontal line separates this from a "Submit your protein sequence" section. This section contains links for "SEARCH STRUCTURAL DATABASE", "ALIGN SEQUENCE WITH STRUCTURE", "DOWNLOAD", and "DOCUMENTATION". Another horizontal line separates this from a "Methods" section. The "Methods" section contains text about the FUGUE program's functionality and a link to a summary of how it works. It also mentions the original paper by Shi et al. (2001) and some practical information from a book by Núñez Miguel et al. (2001). At the bottom, there is a note about the HOMSTRAD database.

FUGUE: sequence-structure homology recognition and alignment engine - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Media Mail Links

Address http://www-cryst.bioc.cam.ac.uk/fugue/ Go Links

**FUGUE**

Crystallography and Biocomputing Unit  
Department of Biochemistry, University of Cambridge

Sequence-structure homology recognition using environment-specific substitution tables and structure-dependent gap penalties

---

**Submit your protein sequence**

[SEARCH STRUCTURAL DATABASE](#)

[ALIGN SEQUENCE WITH STRUCTURE](#)

[DOWNLOAD](#)

[DOCUMENTATION](#)

---

**Methods**

FUGUE is a program for recognizing distant homologues by sequence-structure comparison. It utilizes environment-specific substitution tables and structure-dependent gap penalties, where scores for amino acid matching and insertions/deletions are evaluated depending on the local environment of each amino acid residue in a known structure. Given a query sequence (or a sequence alignment), FUGUE scans a database of structural profiles, calculates the sequence-structure compatibility scores and produces a list of potential homologues and alignments.

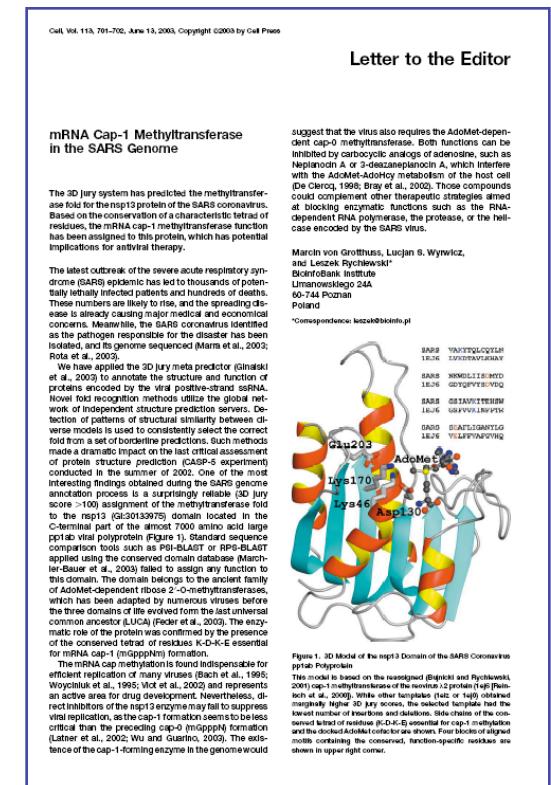
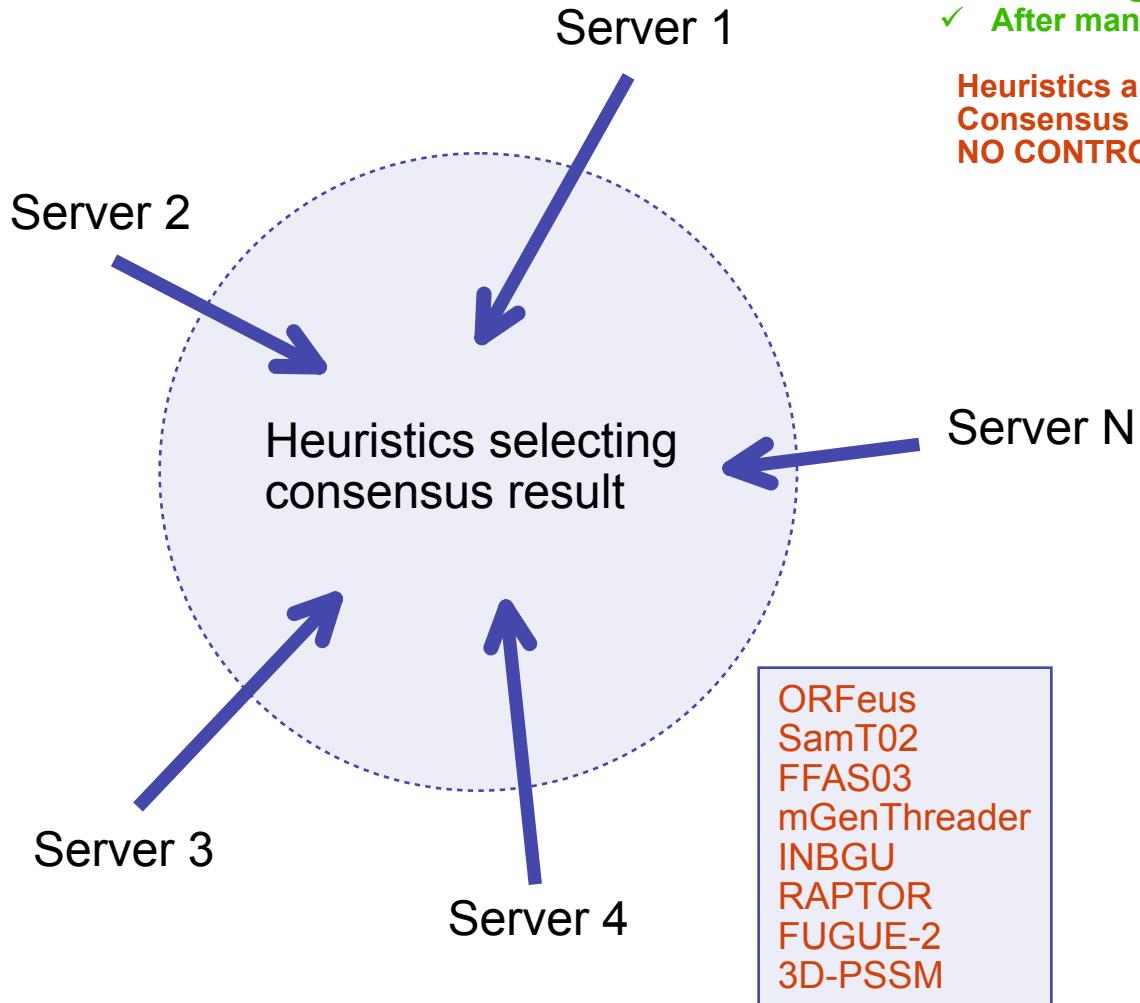
[Here](#) is a summary of how it works.

Read the original paper for more details:  
[J. Shi, T. L. Blundell, and K. Mizuguchi \(2001\). FUGUE: sequence-structure homology recognition using environment-specific substitution tables and structure- dependent gap penalties. \*J. Mol. Biol.\*, 310, 243-257.](#)  
[Medline](#), [Article on-line](#), [PDF \(local only\)](#).

Some practical information can be found in:  
[R. Núñez Miguel, J. Shi and K. Mizuguchi \(2001\). Protein Fold Recognition and Comparative Modeling using HOMSTRAD, JOY and FUGUE. In \*Protein Structure Prediction: Bioinformatic Approach\*. International University Line publishers, La Jolla, 143-169.](#)  
[PDF \(local only\)](#)

Click [here](#) for information about the [HOMSTRAD](#) database.

# Meta-Servers (3D-Jury)



# Meta-Servers (3D-Jury)

<http://bioinfo.pl/Meta/>

Meta Server Job List, BioInfo.PL - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Refresh Search Favorites Media Mail Home Links

Address http://bioinfo.pl/Meta/ Go Links

BIOINFO.PL: META Meta Server Job List [ABOUT] [SERVERS] [BENCHMARKS] [STATUS]

Structure Prediction Meta Server Input Page  
0 jobs from 64.54.249. in the last week

Your E-mail: \_\_\_\_\_

Target Name: \_\_\_\_\_

Amino Acid Sequence only (in one letter code):

Skip: Queue:  
 PDB-Blast 1  
 3D-Jigsaw 43  
 ESyPred3D 1  
 GRDB 1  
 FFAS03 1  
 Sam-T99 1  
 SUPERFAMILY 1  
 INBGU 39  
 FUGUE2 1  
 3D-PSSM 1  
 mGenTHREADER  
 psipred  
 profsec 1  
Pcons2 1  
3D-ShotGun 11  
3D-Jury

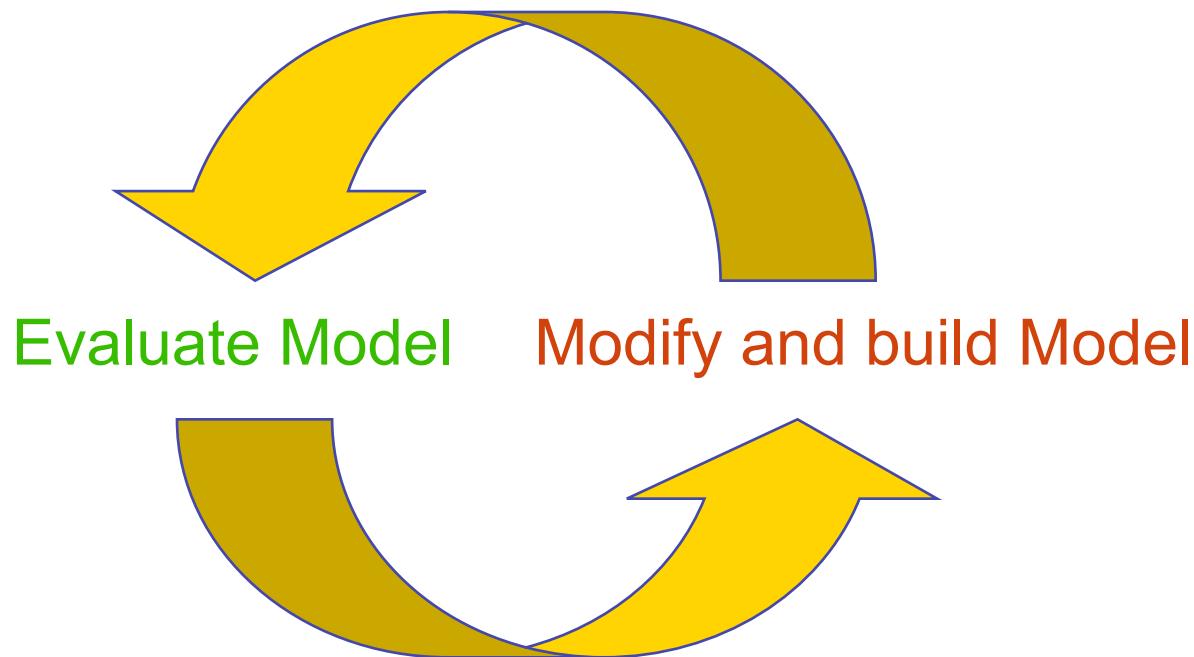
Reset Clear Format Submit

Please submit domains separately  
Please remove coiled coil regions  
Check LiveBench for evaluation of the reliability of the servers  
Results are stored only for 1 month  
Jobs queued for more than 7 days for servers with queue>30 are skipped  
Use is limited to 10 jobs per week per domain  
Please contact us in case of problems with interpretation of results  
Please contact us if You plan larger analysis projects  
Some servers return only models, no alignments (target sequence is shown)

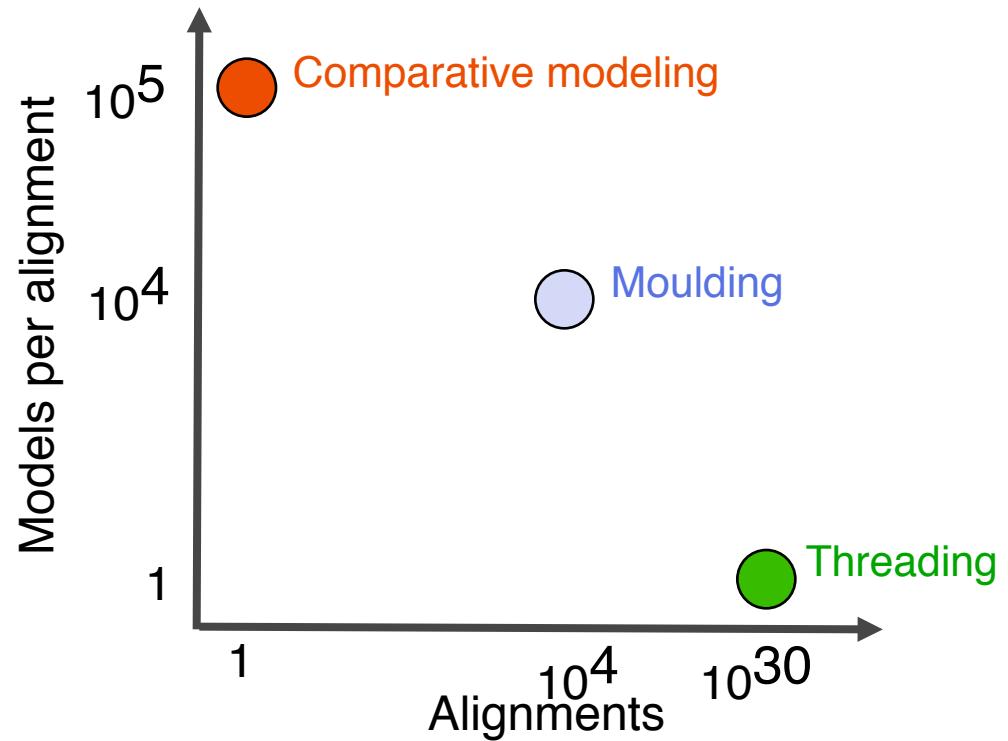
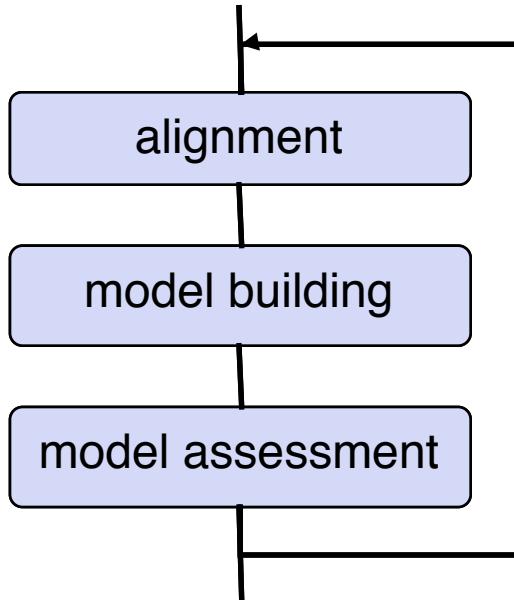
Please cite the prediction servers and 3D-Jury:  
Ginalski K, Elofsson A, Fischer D, Rychlewski L.  
"3D-Jury: a simple approach to improve protein structure predictions."  
Bioinformatics. 2003 May 22;19(8):1015-8. [PubMed]

Internet

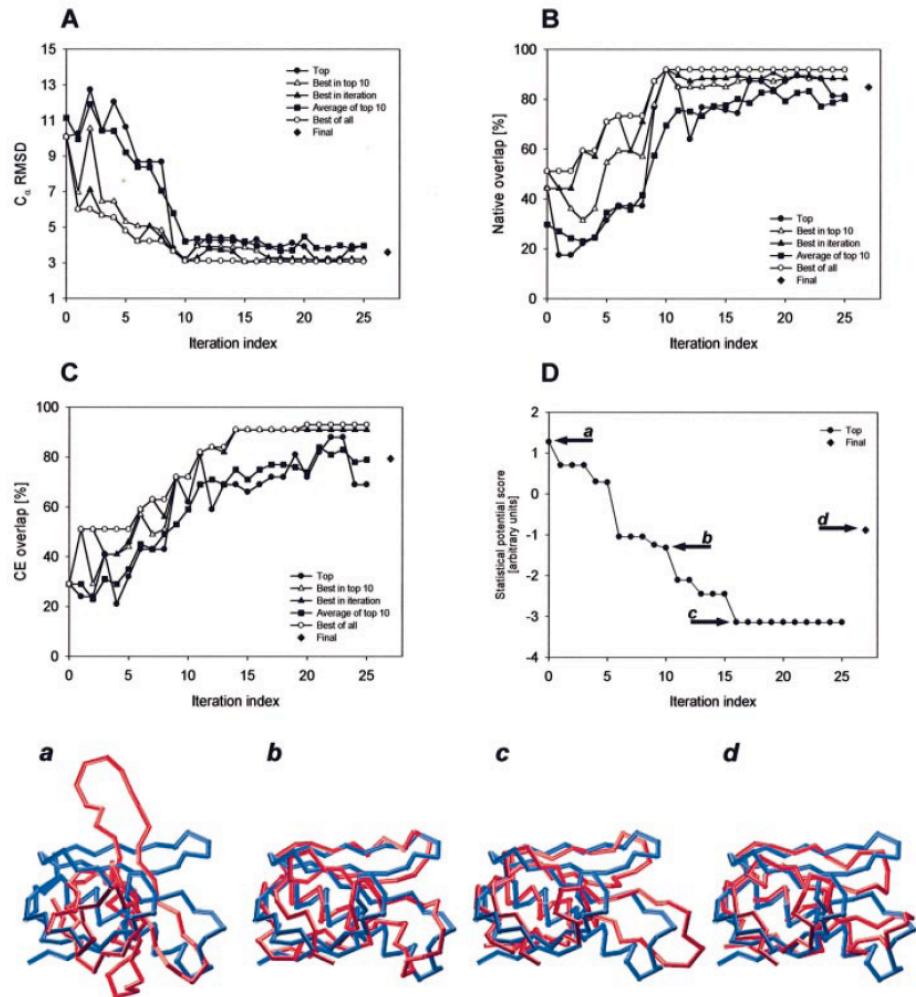
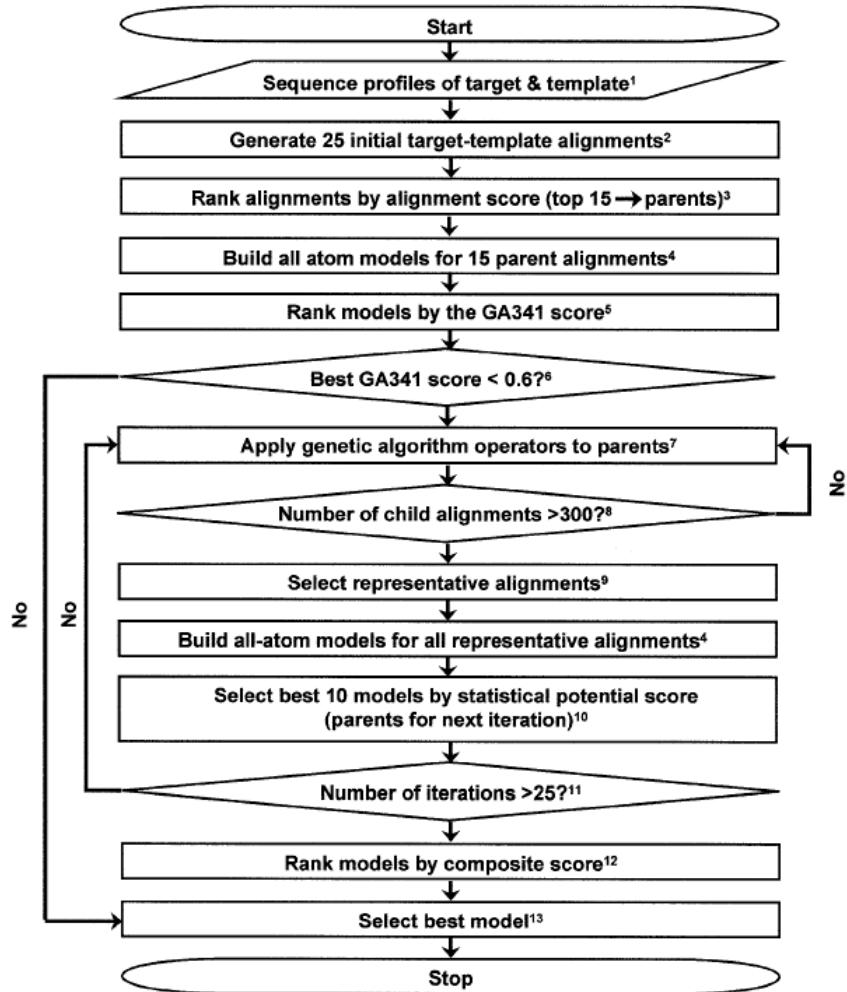
# Iterative process... better models(?)



# Moulding: iterative alignment, model building, model assessment



# Iterative process... MOULDER



# Genetic algorithm operators

## Single point cross-over

...TSSQ-N**MKLGVFWGY**...  
...V-S**SCN**GDLHMKVGV...



...TSSQ-N**MK**LGVFWGY...  
...V-S**SCN**GDLHMKV—GV...

...TSSQ**N**M**KLGVFWGY**...  
...V**S**SCN—GDLHMKVGV...

## Gap insertion

...TSSQ**N**M**KLGVFWGY**...  
...V**S**SCN**GDLHMKVGV**...



...TSSQ**N**M**KLGVFWGY**...  
...V**S**SCN**GDLHMKVG**V...

## Gap shift

...**T**—**S**SQ**N**M**KLGVFWGY**...  
...V**S**SCN**GDLHMKVGV**—...



...—**T**—SQ**N**M**KLGVFWGY**...  
...V**S**SCN**GDLHMKVGV**—...

...—**T**—**S**—SQ**N**M**KLGVFWGY**...  
...V**S**SCN**GDLHMKVGV**—...

...—**T**—SQ**N**M**KLGVFWGY**...  
...V**S**SCN**GDLHMKVGV**—...

...**TS**—SQ**N**M**KLGVFWGY**...  
...V**S**SCN**GDLHMKVGV**—...

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 ( $A_s$ ).

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

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

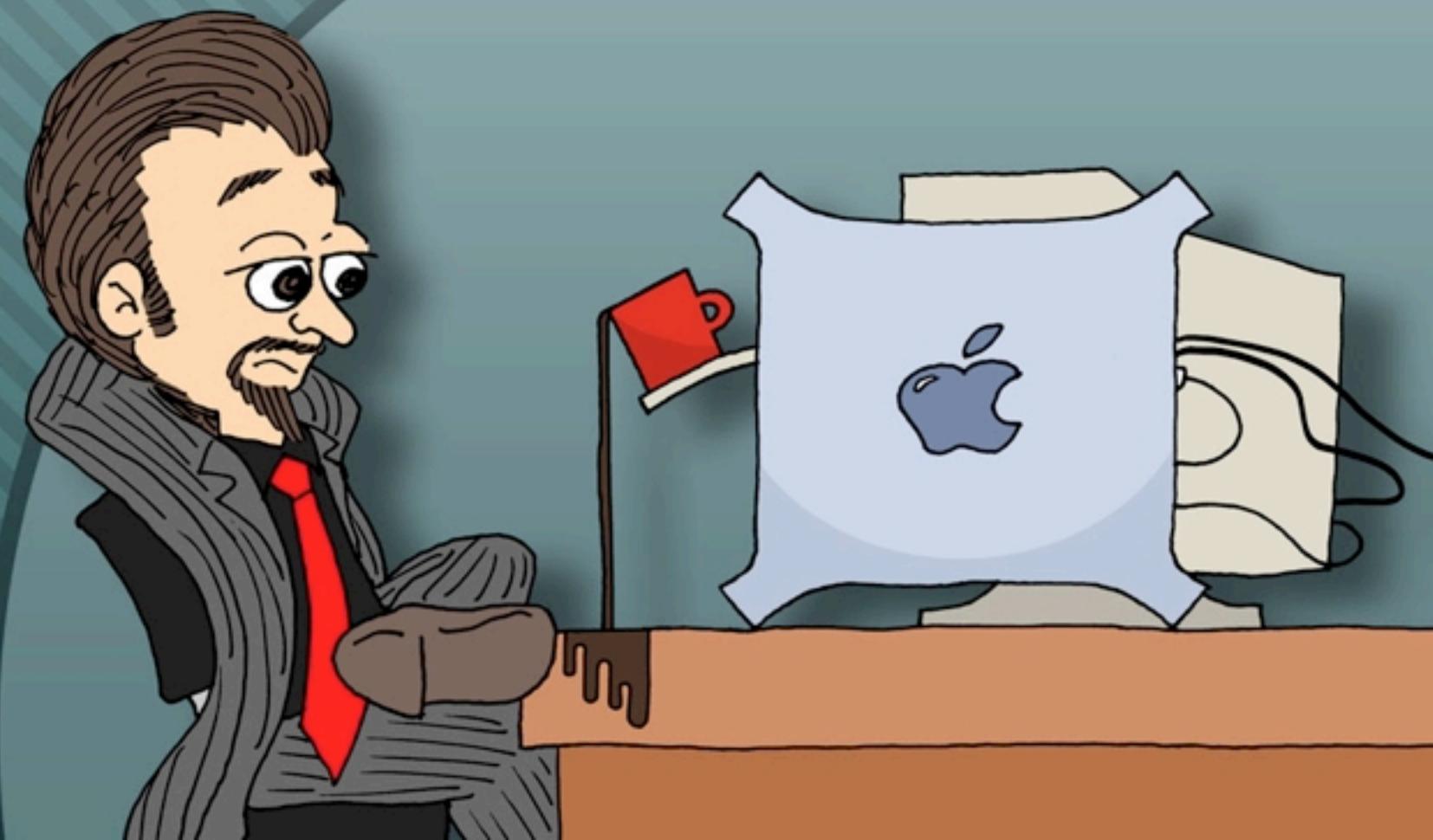
$\mu$  ... average score of all models

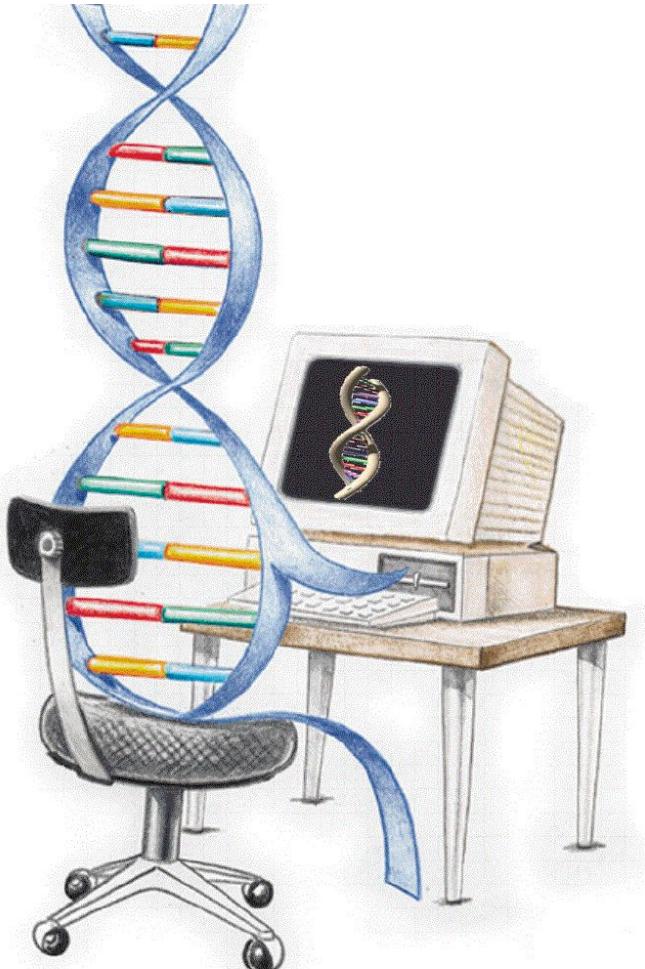
$\sigma$  ... standard deviation of the scores

# 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 <sub>α</sub> RMSD [Å]	CE overlap [%]	C <sub>α</sub> RMSD [Å]	CE overlap [%]	C <sub>α</sub> 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
<b>AVERAGE</b>	<b>14.2</b>	<b>85.2</b>	<b>9.6</b>	<b>36.7</b>	<b>7.7</b>	<b>44.8</b>	<b>6.3</b>	<b>57.8</b>





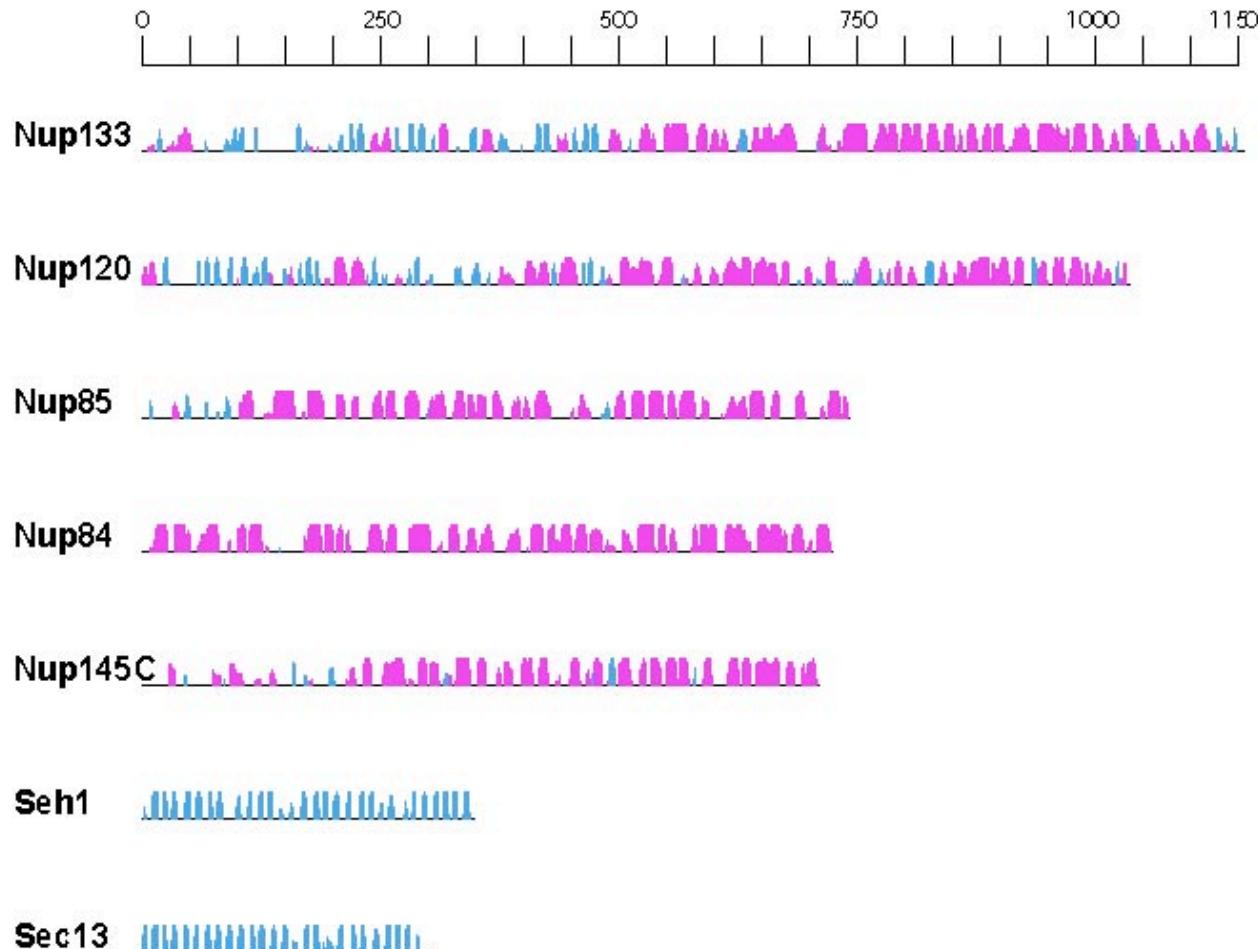
some biology?  
please...

# Common Evolutionary Origin of Coated Vesicles and Nuclear Pore Complexes

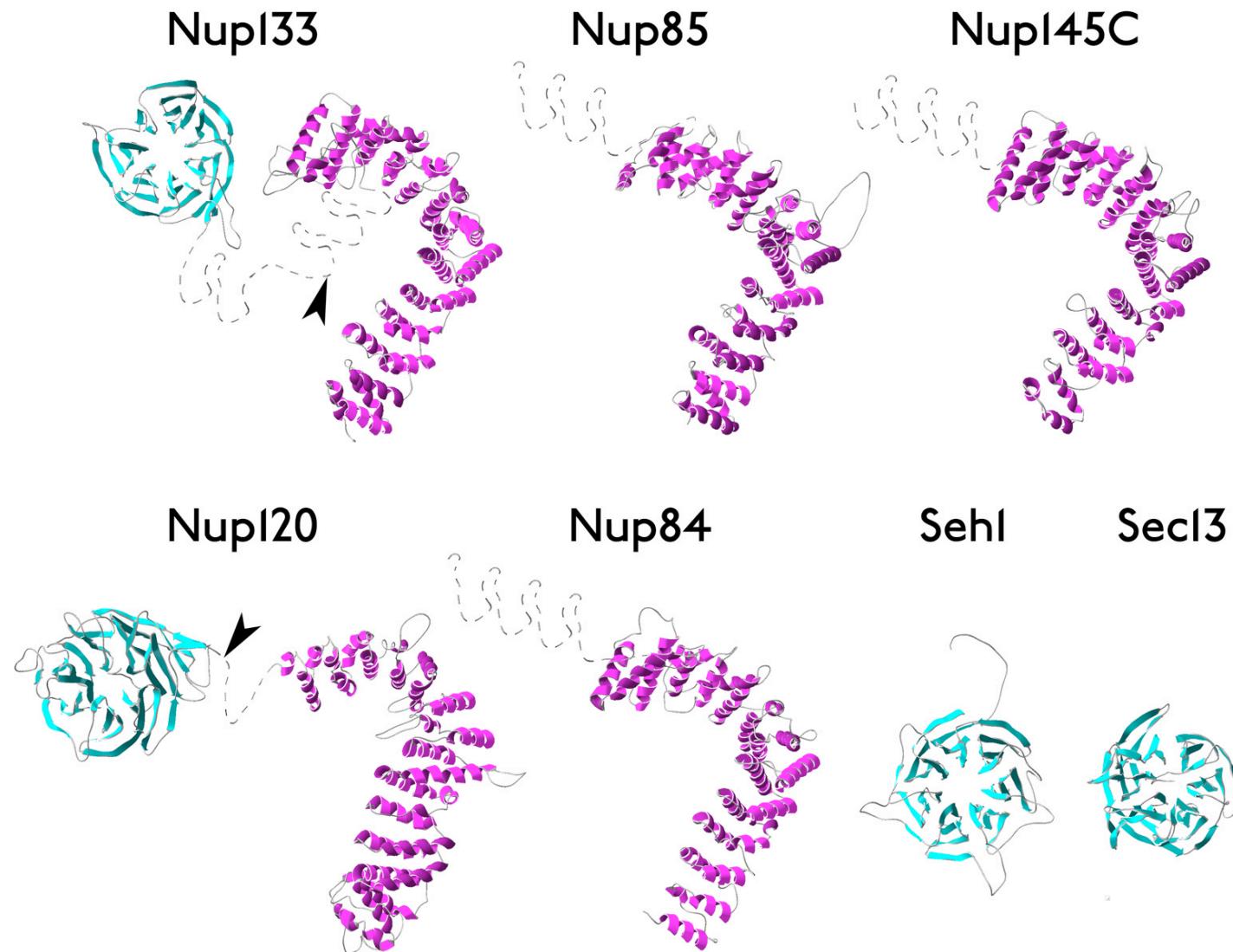
*mGenThreader + SALIGN + MOULDER*

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

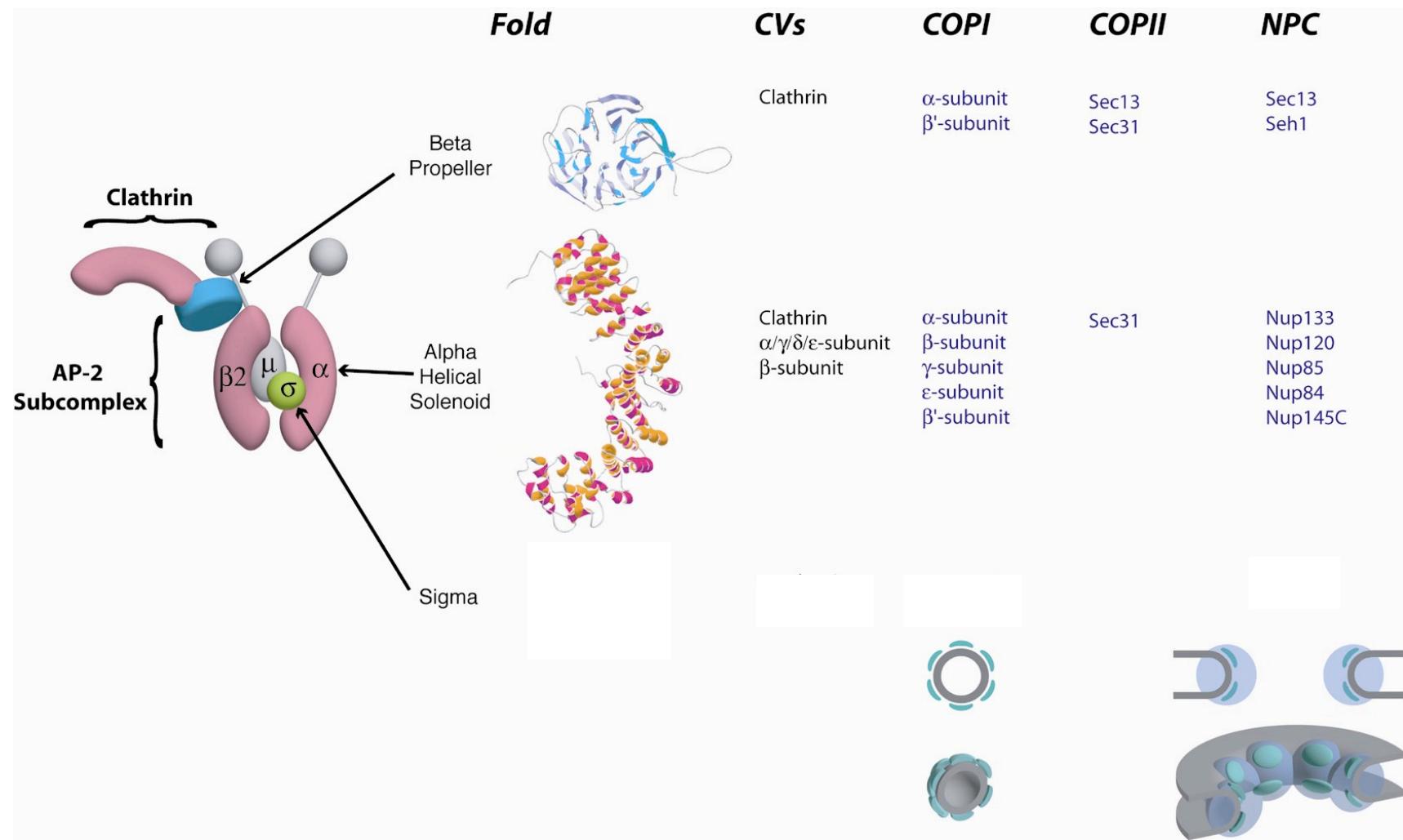
# yNup84 complex proteins



# All Nucleoporins in the Nup84 Complex are Predicted to Contain $\beta$ -Propeller and/or $\alpha$ -Solenoid Folds



# NPC and Coated Vesicles Share the $\beta$ -Propeller and $\alpha$ -Solenoid Folds and Associate with Membranes

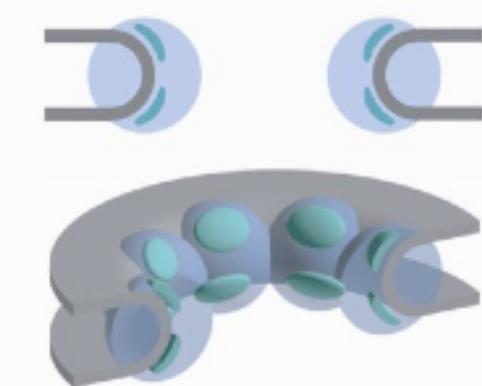


# NPC and Coated Vesicles Both Associate with Membranes

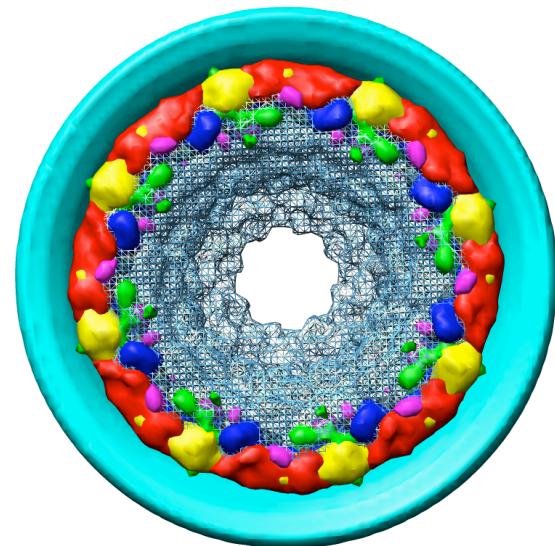
Coated Vesicle



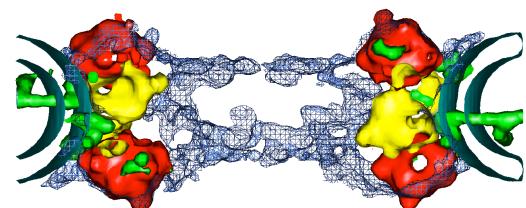
NPC model



top view



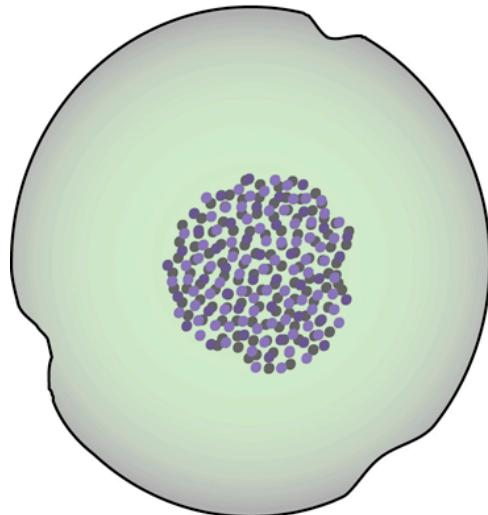
side view



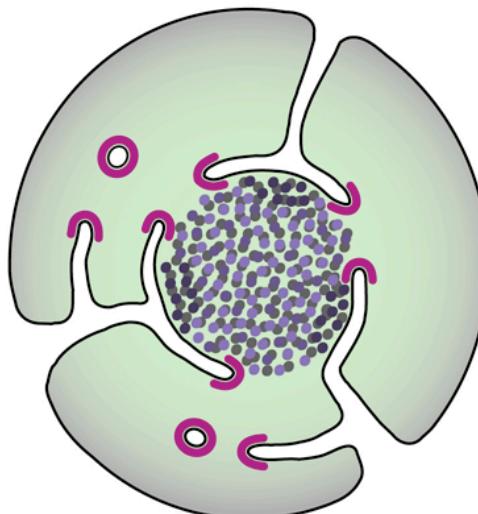
Nup 84 complex

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

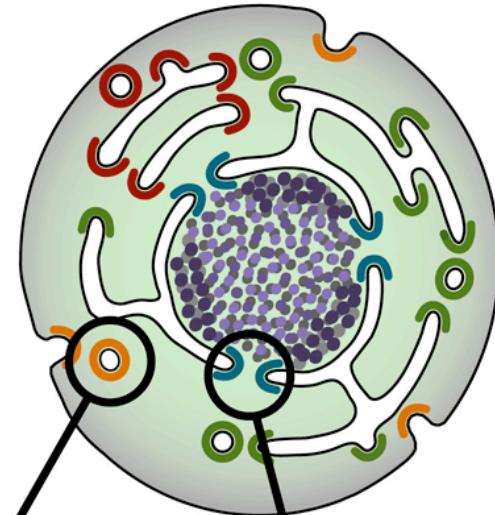
Prokaryote



Early Eukaryote

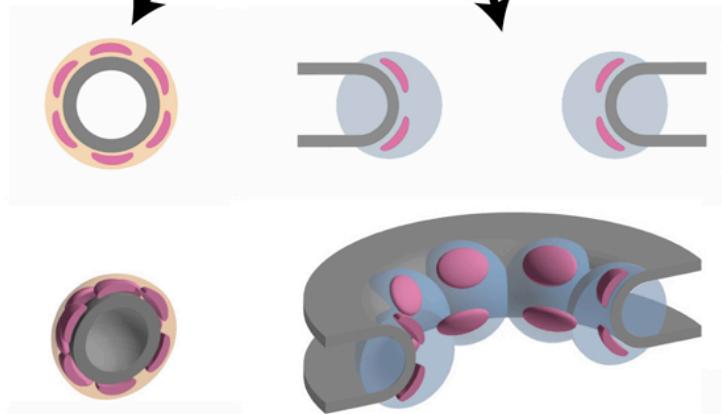


Modern Eukaryote



A simple coating module containing minimal copies of the two conserved folds evolved in proto-eukaryotes to bend membranes.

The progenitor of the NPC arose from a membrane-coating module that wrapped extensions of an early ER around the cell's chromatin.



# **Course assignment**

## **The POM152 nucleoporin protein**

# Course assignment

## The POM152 nucleoporin protein

### Introduction

The POM152 protein functions as a component of the nuclear pore complex (NPC). NPC components, collectively referred to as nucleoporins (NUPs), can play the role of both NPC structural components and of docking or interaction partners for transiently associated nuclear transport factors. POM152 is important for the de novo assembly of NPCs.

The nuclear pore complex (NPC) constitutes the exclusive means of nucleocytoplasmic transport. NPCs allow the passive diffusion of ions and small molecules and the active, nuclear transport receptor-mediated bidirectional transport of macromolecules such as proteins, RNAs, ribonucleoparticles (RNPs), and ribosomal subunits across the nuclear envelope. The 55-60 MDa NPC is composed of at least 31 different subunits. Due to its 8-fold rotational symmetry, all subunits are present with 8 copies or multiples thereof. POM152 is known to interact with NUP188.

### Assignment

1. Predict the domain boundaries for the POM152 sequence
2. Search for a suitable template/s for the POM152 domains
3. Align the sequences of POM152 domains against the template/s sequences
4. Build a 3D-models of the POM152 domains
5. Evaluate the models
6. Indicate possible applications of the models

**GRADING:** The entire assignment is worth 20 points.

<http://www.ncbi.nlm.nih.gov/entrez/viewer.fcgi?db=protein&val=730249>

The screenshot shows the NCBI Entrez Protein search results for the protein POM152. The search bar at the top contains "POM152". Below the search bar, there are various filters and options: "Display" set to "GenPept", "Send" button, "all to file" dropdown, "Range: from begin to end", "Features" checkboxes (SNP, CDD, MGCG, HPRD), and a checked checkbox for "MGCG". The results table has columns for Locus, Definition, Accession, Version, DBSource, Keywords, Source, Organism, Reference, Authors, Title, Journal, PubMed, and Remark. The first result for POM152 is shown, detailing its 1337 amino acid linear sequence (PLN 25-JAN-2005), its definition as the Nuclear pore protein POM152 (Pore membrane protein POM152) (P150), and its source as Saccharomyces cerevisiae (baker's yeast). It also lists its Eukaryote, Fungi, Ascomycota, Saccharomycotina, Saccharomycetes, Saccharomycetales, Saccharomycetaceae, and Saccharomyces taxonomic levels. The reference is to Wozniak et al. (1994) in J. Cell. Biol., 125(1), 31-42, which states POM152 is an integral protein of the pore membrane domain of the yeast nuclear envelope. The journal entry is dated Feb 1, 1995, and the PMID is 730249. The remark section includes information about the nucleotide sequence, partial protein sequence, glycosylation, and repeats, with STRAIN=W303. The next reference is to Bowman et al. (1997) in Nature, 387 (6632 SUPPL), 90-93, which discusses the nucleotide sequence of Saccharomyces cerevisiae chromosome XIII. The journal entry is dated Jun 12, 1997, and the PMID is 9169873. The remark section includes information about the nucleotide sequence (LARGE SCALE GENOMIC DNA) and STRAIN=S288c / AB972. The final reference is to Nehrbass et al. (1996) in J. Cell. Biol., 133 (6), 1153-1162, which discusses the yeast nucleoporin Nup188p interacts genetically and physically with the core structures of the nuclear pore complex. The journal entry is dated Jul 1, 1996, and the PMID is 8682855.

BMI 206

# MODELLER TUTORIAL

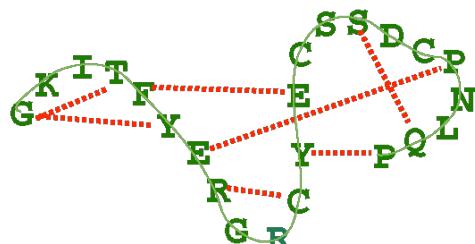
<http://www.salilab.org/modeller/tutorial/>

Marc A. Marti-Renom  
Assistant Adjunct Professor  
Department of Biopharmaceutical Sciences

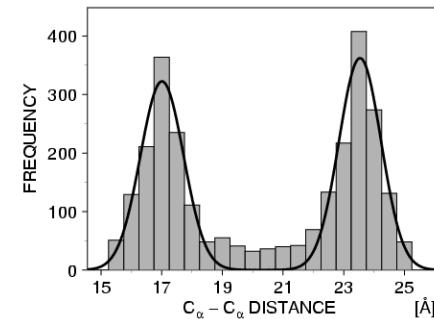
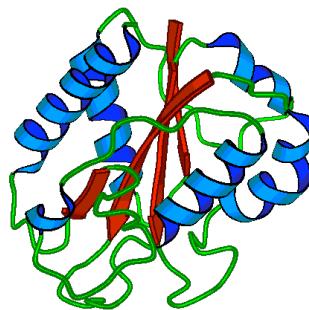
# Comparative Modeling by Satisfaction of Spatial Restraints (MODELLER)

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

## 1. Extract spatial restraints

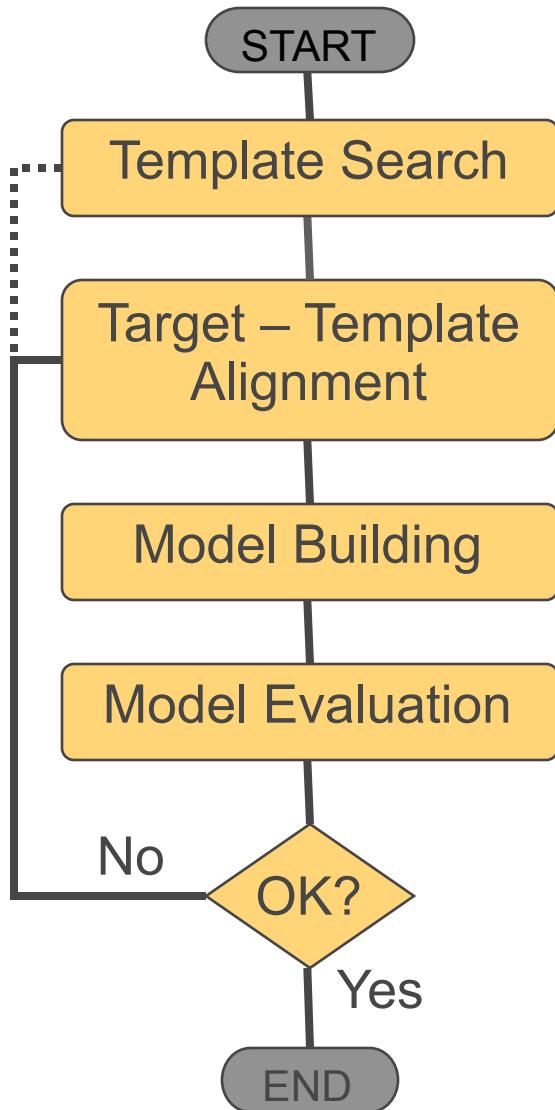


## 2. Satisfy spatial restraints



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

# Steps in Comparative Protein Structure Modeling



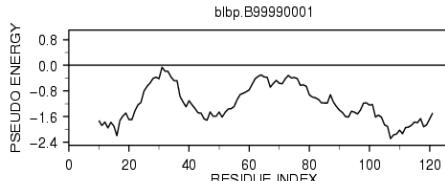
## TARGET

ASILPKRLFGNCEQTSDEG  
LKIERTPLVPHISAQNVLCKI  
DDVPERLIPERASFQWMN  
DK

## TEMPLATE



ASILPKRLFGNCEQTSDEGLKIERPLVPHISAQNVLCKIDDVPERLIP  
MSVIPKRLYGNCEQTSEEAIRIEDSPIV---TADLVCLKIDEIPELVGE



A. Šali, *Curr. Opin. Biotech.* 6, 437, 1995.

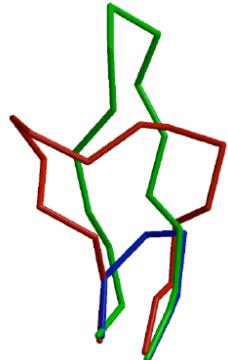
R. Sánchez & A. Šali, *Curr. Opin. Str. Biol.* 7, 206, 1997.

M. Martí et al. *Ann. Rev. Biophys. Biomol. Struct.*, 29, 291, 2000.

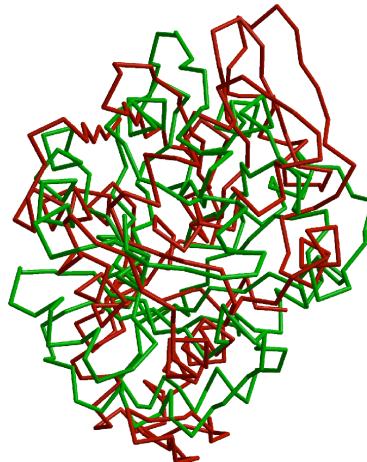
# Typical errors in comparative models

MODEL  
X-RAY  
TEMPLATE

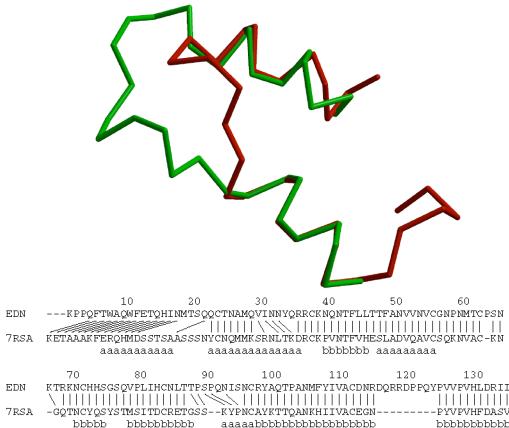
Region without a template



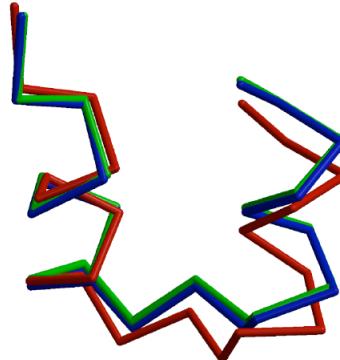
Incorrect template



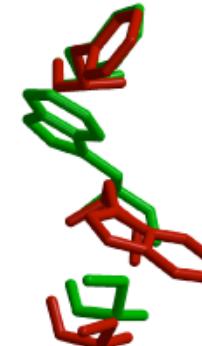
Misalignment



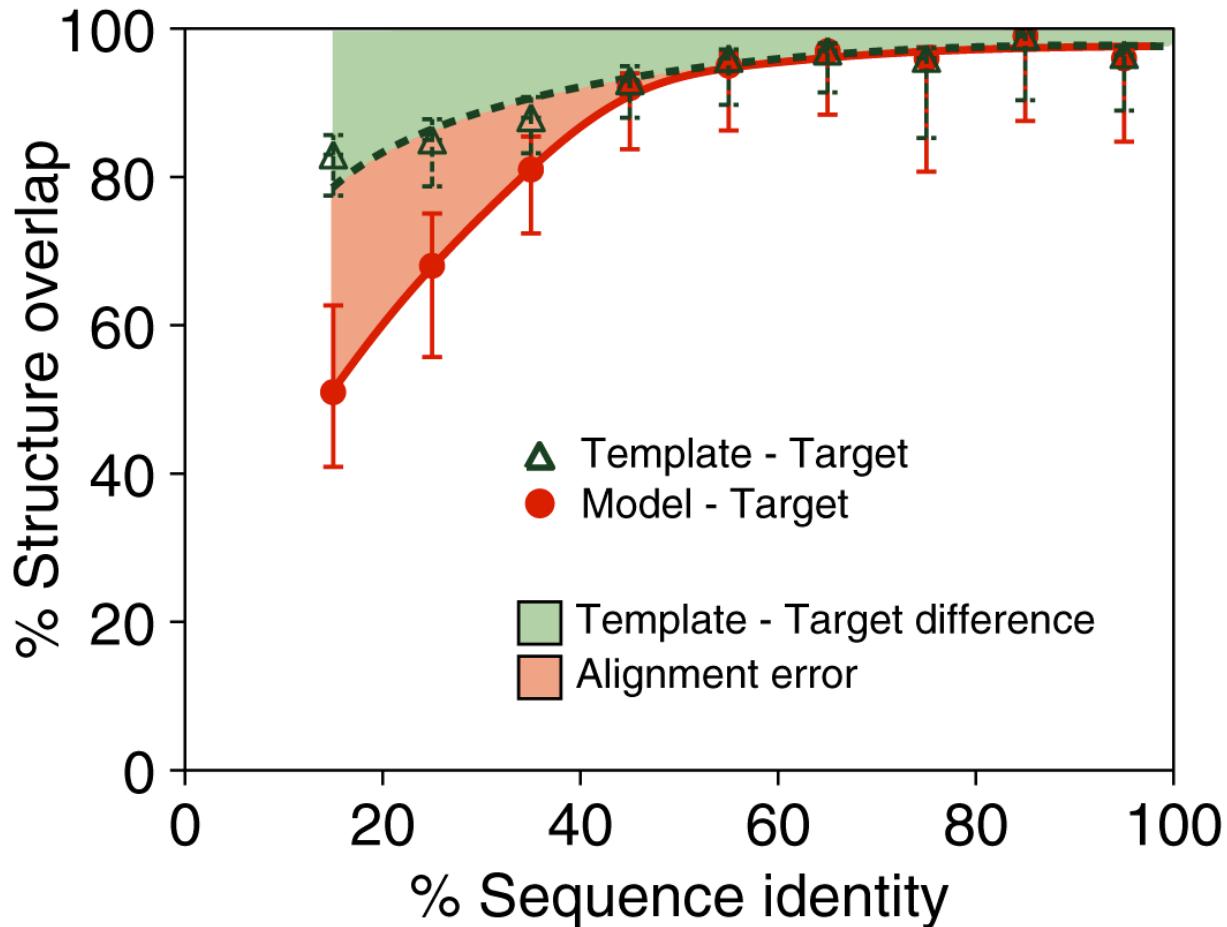
Distortion/shifts in aligned regions



Sidechain packing



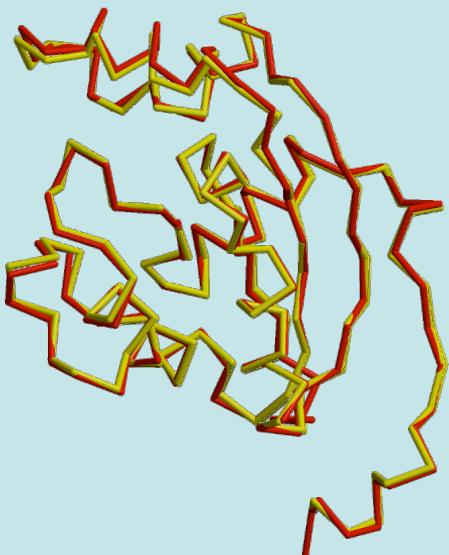
## Model Accuracy as a Function of Target-Template Sequence Identity



# Model Accuracy

## HIGH ACCURACY

NM23  
Seq id 77%  
 $C\alpha$  equiv 147/148  
RMSD 0.41Å

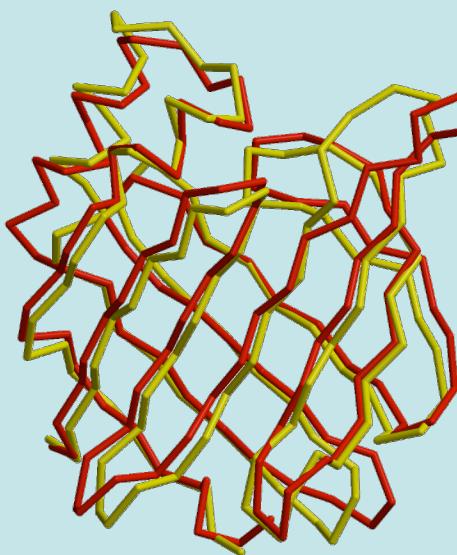


Sidechains  
Core backbone  
Loops

X-RAY / MODEL

## MEDIUM ACCURACY

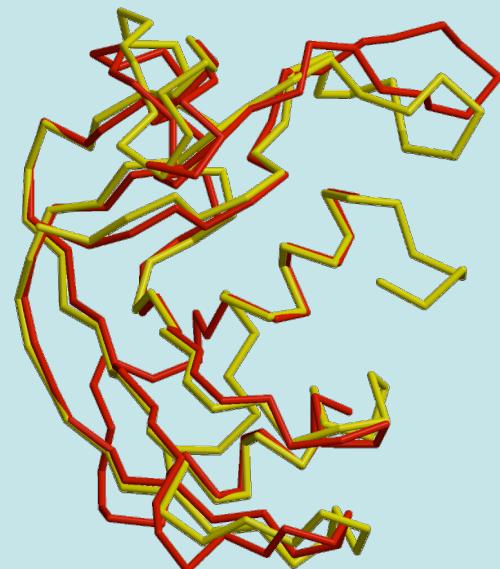
CRABP  
Seq id 41%  
 $C\alpha$  equiv 122/137  
RMSD 1.34Å



Sidechains  
Core backbone  
Loops  
Alignment

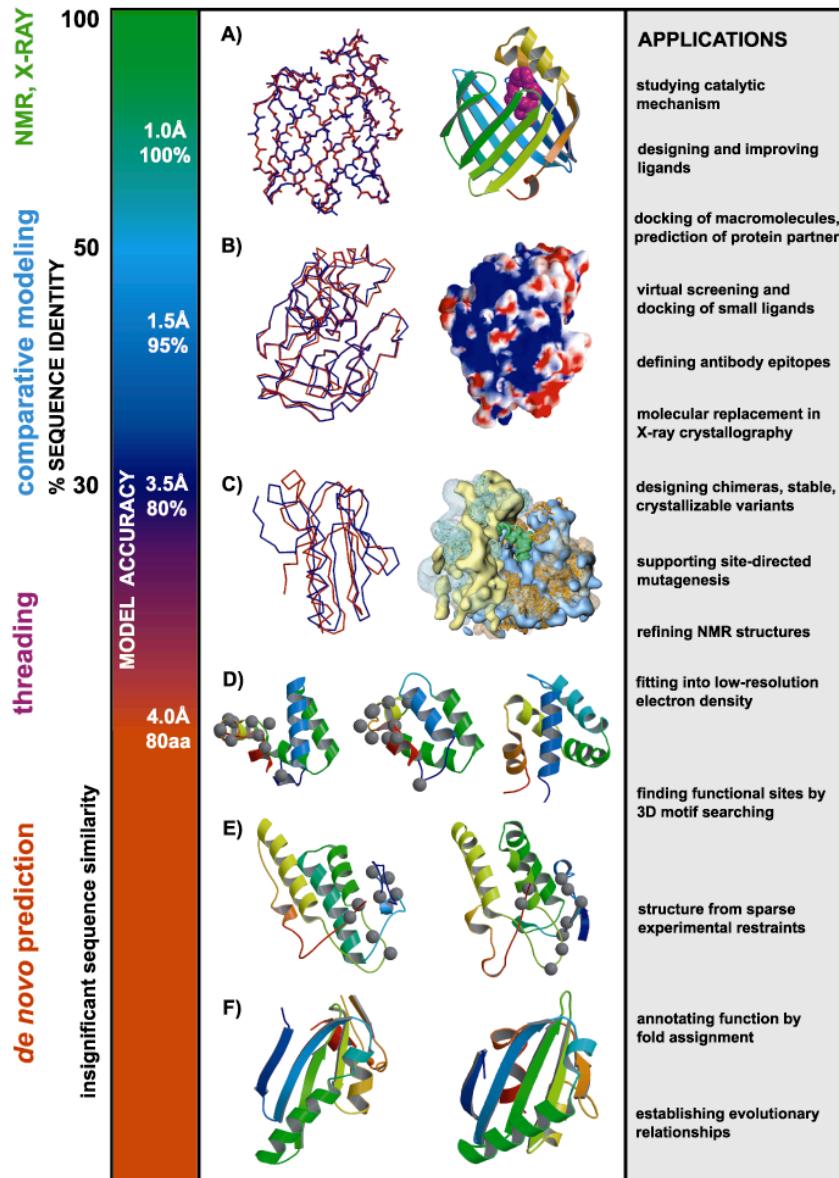
## LOW ACCURACY

EDN  
Seq id 33%  
 $C\alpha$  equiv 90/134  
RMSD 1.17Å



Sidechains  
Core backbone  
Loops  
Alignment  
Fold assignment

# Applications of Protein Structure Models



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

# Obtaining MODELLER and related information

- MODELLER (7v7) 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) ☹☹
- Script is written in TOP language ☹☹☹
- TOP 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:::::::::::
VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIIISQEGGKVVI
RTQCTFKNTEINFQLGEFEETSI
DDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMWVTLT
FGDIVAVRCYEKA*
```

- PSI-BLAST template search: Template: PDB file 1HMS:\_\_

# Modeling of BLBP

STEP 1: Align **blbp** and **1hms** sequences  
*TOP script for target-template alignment*

```
READ_MODEL FILE = '1hms.pdb'
SEQUENCE_TO_ALI ALIGN_CODES = '1hms'
READ_ALIGNMENT FILE = 'blbp.seq', ALIGN_CODES = 'blbp', ADD_SEQUENCE = on
ALIGN
WRITE_ALIGNMENT FILE='blbp-1hms.ali', ALIGNMENT_FORMAT = 'PIR'
WRITE_ALIGNMENT FILE='blbp-1hms.pap', ALIGNMENT_FORMAT = 'PAP'
```

Run by typing `mod align.top` directory where you have the TOP file.  
MODELLER will produce a `align.log` file

# Modeling of BLBP

STEP 1: Align **blbp** and **1hms** sequences  
*TOP script for target-template alignment*

```
READ_MODEL FILE = '1hms.pdb'
SEQUENCE_TO_ALI ALIGN_CODES = '1hms'
READ_ALIGNMENT FILE = 'blbp.seq', ALIGN_CODES = 'blbp', ADD_SEQUENCE = on
ALIGN
WRITE_ALIGNMENT FILE='blbp-1hms.ali', ALIGNMENT_FORMAT = 'PIR'
WRITE_ALIGNMENT FILE='blbp-1hms.pap', ALIGNMENT_FORMAT = 'PAP'
```

Run by typing `mod align.top` directory where you have the TOP file.  
MODELLER will produce a `align.log` file

# Modeling of BLBP

STEP 1: Align blbp and 1hms sequences

*TOP script for target-template alignment*

```
READ_MODEL FILE = '1hms.pdb'
SEQUENCE_TO_ALI ALIGN_CODES = '1hms'
READ_ALIGNMENT FILE = 'blbp.seq', ALIGN_CODES = 'blbp', ADD_SEQUENCE = on
ALIGN
WRITE_ALIGNMENT FILE 'blbp-1hms.ali', ALIGNMENT_FORMAT = 'PIR'
WRITE_ALIGNMENT FILE 'blbp-1hms.pap', ALIGNMENT_FORMAT = 'PAP'
```

Run by typing `mod align.top` directory where you have the TOP file.  
MODELLER will produce a `align.log` file

# Modeling of BLBP

STEP 1: Align blbp and 1hms sequences

*TOP script for target-template alignment*

```
READ_MODEL FILE = '1hms.pdb'
SEQUENCE_TO_ALI ALIGN_CODES = '1hms'
READ_ALIGNMENT FILE = 'blbp.seq', ALIGN_CODES = 'blbp', ADD_SEQUENCE = on
ALIGN
WRITE_ALIGNMENT FILE='blbp-1hms.ali', ALIGNMENT_FORMAT = 'PIR'
WRITE_ALIGNMENT FILE='blbp-1hms.pap', ALIGNMENT_FORMAT = 'PAP'
```

Run by typing mod align.top directory where you have the TOP 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
VDAFLGTWKLVDSKNFDDYMKSLGVGFATRQVASMTKPTTIEKNGDILTLKTHSTFKNTEISFKLGVEFDETTA
DDRKVKSIVTLDGGKLVHLQKWDGQETTLVRELIIDGKLILTLHGTAVCTRTRYEKE*

>P1;blbp
sequence:blbp: : : : : : 0.00: 0.00
VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSI
DDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMWVTLTFGDIVAVRCYEKA*
```

# Modeling of BLBP

## STEP 1: Align **blbp** and **1hms** sequences

### *Output*

```
>P1;1hms
structureX:1hms: 1 : : 131 : :undefined:undefined:-1.00:-1.00
VDAFLGTWKLVDSKNFDDYMKSLGVGFATRQVASMTKPTTIEKNGDILTLKTHSTFKNTEISFKLGVEFDETTA
DDRKVKSIVTLDGGKLVHLQKWDGQETTLVRELIIDGKLILTLHGTAVCTRTRYEKE*
>P1;blbp
sequence:blbp: : : : : : 0.00: 0.00
VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSI
DDRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMWVTLTFGDIVAVRCYEKA*
```

# Modeling of BLBP

## STEP 1: Align **blbp** and **1hms** sequences

### *Output*

_aln.pos	10	20	30	40	50	60
1hms	VDAFLGTWKLVD SKNFDD YMKS LGVG FATRQ VASMT KPTT IIEKNGD ILTLKTH STFKNT					
blbp	VDAFCATWKL TD SQNF D EYMK ALGV GFATRQ VGNV TKPTV IIS QEGG KV VIRT QCT FKNT					
_consrvd	***** *					
_aln.pos	70	80	90	100	110	120
1hms	EISFKLGVEFDETTADDRKVKSIVTLGGKLVHLQKWDGQETTLVRELIDGKLILTLHG					
blbp	EINFQLGEFEETSIDDRNCKSVVR LDGDKL IHVQKWDGKETNCTREIKDGKMVVTLFG					
_consrvd	* * * * * * * *** *					
_aln.pos	130					
1hms	TAVCTR TYEKE					
blbp	DIVAVRCYEKA					
_consrvd	* * ***					

# Modeling of BLBP

STEP 2: Model the blbp structure using the alignment from step 1.  
*TOP script for model building*

```
INCLUDE  
  
SET ALNFILE = 'blbp-1hms.ali'  
SET KNOWNS = '1hms'  
SET SEQUENCE = 'blbp'  
SET STARTING_MODEL = 1  
SET ENDING_MODEL = 1  
CALL ROUTINE = 'model'
```

Run by typing mod model.top.  
Check file model.log

# Modeling of BLBP

STEP 2: Model the blbp structure using the alignment from step 1.  
*TOP script for model building*

```
INCLUDE  
  
SET ALNFILE = 'blbp-1hms.ali'  
  
SET KNOWNS = '1hms'  
  
SET SEQUENCE = 'blbp'  
  
SET STARTING_MODEL = 1  
  
SET ENDING_MODEL = 1  
  
CALL ROUTINE = 'model'
```

Run by typing mod model.top.  
Check file model.log

# Modeling of BLBP

STEP 2: Model the blbp structure using the alignment from step 1.  
*TOP script for model building*

```
INCLUDE  
  
SET ALNFILE = 'blbp-1hms.ali'  
SET KNOWNS = '1hms'  
SET SEQUENCE = 'blbp'  
SET STARTING_MODEL = 1  
SET ENDING_MODEL = 1  
CALL ROUTINE = 'model'
```

Run by typing mod model.top  
Check file model.log

## Modeling of BLBP

STEP 2: Model the blbp structure using the alignment from step 1.  
*Output coordinates file*

Model file → blbp.B99990001

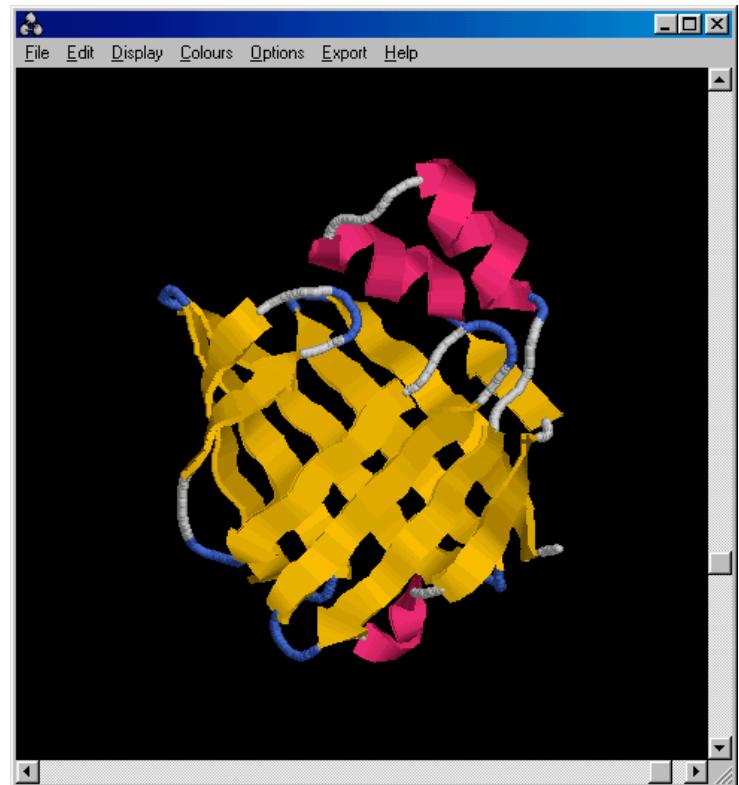
PDB file

Can be viewed with Chimera

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

Rasmol

<http://www.bernstein-plus-sons.com/software/rasmol/>



[http://www.salilab.org/bioinformatics\\_resources.shtml](http://www.salilab.org/bioinformatics_resources.shtml)

Andrej Sali Lab  
http://salilab.org/bioinformatics\_resources.shtml

Andrej Sali Lab

University of California, San Francisco  
Departments of Biopharmaceutical Sciences and Pharmaceutical Chemistry, and  
California Institute for Quantitative Biomedical Research

UCSF

Modeler

Bioinformatics Resources

Programs and World Wide Web servers useful in comparative modeling

Name	Type <sup>a</sup>	World Wide Web address <sup>b</sup>
<b>DATABASES</b>		
CATH	S	<a href="http://www.biochem.ucl.ac.uk/bsm/cath/">http://www.biochem.ucl.ac.uk/bsm/cath/</a>
GenBank	S	<a href="http://www.ncbi.nlm.nih.gov/Genbank/GenbankSearch.html">http://www.ncbi.nlm.nih.gov/Genbank/GenbankSearch.html</a>
GeneCensus	S	<a href="http://bioinfo.mbb.yale.edu/genome">http://bioinfo.mbb.yale.edu/genome</a>
ModBASE	S	<a href="http://salilab.org/modbase/">http://salilab.org/modbase/</a>
MSD	S	<a href="http://www.rcsb.org/databases.html">http://www.rcsb.org/databases.html</a>
NCBI	S	<a href="http://www.ncbi.nlm.nih.gov/">http://www.ncbi.nlm.nih.gov/</a>
PDB	S	<a href="http://www.rcsb.org/pdb/">http://www.rcsb.org/pdb/</a>
PRESAGE	S	<a href="http://presage.berkeley.edu/">http://presage.berkeley.edu/</a>
PSI	S	<a href="http://www.structuralgenomics.org/">http://www.structuralgenomics.org/</a>
Sacch3D	S	<a href="http://genome-www.stanford.edu/Sacch3D/">http://genome-www.stanford.edu/Sacch3D/</a>
SCOP	S	<a href="http://scop.mrc-lmb.cam.ac.uk/scop/">http://scop.mrc-lmb.cam.ac.uk/scop/</a>
TIGR	S	<a href="http://www.tigr.org/tigr/mdb/mdbcomplete.html">http://www.tigr.org/tigr/mdb/mdbcomplete.html</a>
TrEMBL	S	<a href="http://era.ebi.ac.uk/">http://era.ebi.ac.uk/</a>
<b>FOLD ASSIGNMENT</b>		
123D	S	<a href="http://www.lecb.ncifcrf.gov/~nicka/123D.html">http://www.lecb.ncifcrf.gov/~nicka/123D.html</a>
3D-PSSM	S	<a href="http://www.bmn.icnet.uk/~3dpssm/html/ffrecog.html">http://www.bmn.icnet.uk/~3dpssm/html/ffrecog.html</a>
BIOINBGU	S	<a href="http://www.cs.bgu.ac.il/~bioinbgu/">http://www.cs.bgu.ac.il/~bioinbgu/</a>
BLAST	S	<a href="http://www.ncbi.nlm.nih.gov/BLAST/">http://www.ncbi.nlm.nih.gov/BLAST/</a>
DALI	S	<a href="http://www2.ebi.ac.uk/dali/">http://www2.ebi.ac.uk/dali/</a>
FASS	S	<a href="http://bioinformatics.burnham-inst.org/FFAS/index.html">http://bioinformatics.burnham-inst.org/FFAS/index.html</a>
FastA	S	<a href="http://www.ebi.ac.uk/fasta3/">http://www.ebi.ac.uk/fasta3/</a>
FRSVR	S	<a href="http://www.doe-mbi.ucla.edu/~frsyr/preds/MG/MG.html">http://www.doe-mbi.ucla.edu/~frsyr/preds/MG/MG.html</a>
FUGUE	S	<a href="http://www-crypt.bioc.cam.ac.uk/~fugue">http://www-crypt.bioc.cam.ac.uk/~fugue</a>
Koonin Group	S	<a href="ftp://ncbi.nlm.nih.gov/pub/koonin/FOLDS/index.html">ftp://ncbi.nlm.nih.gov/pub/koonin/FOLDS/index.html</a>
LOOPP	S	<a href="http://www.tc.cornell.edu/reports/NIH/resource/CompBiologyTools/looppp/">http://www.tc.cornell.edu/reports/NIH/resource/CompBiologyTools/looppp/</a>
PDB-Blast/FASS	S	<a href="http://bioinformatics.ljcrf.edu/pdb_blast/">http://bioinformatics.ljcrf.edu/pdb_blast/</a>
PHD_TOPIST	S	<a href="http://www.embl-heidelberg.de/predictprotein/predictprotein.html">http://www.embl-heidelberg.de/predictprotein/predictprotein.html</a>
PROFIT	P	<a href="http://www.came.sbg.ac.at">http://www.came.sbg.ac.at</a>
SAM-T99/T98	S	<a href="http://www.cse.ucsc.edu/research/compbio/HMM-apps/">http://www.cse.ucsc.edu/research/compbio/HMM-apps/</a>
THREADER	S	<a href="http://insulin.brunel.ac.uk/threader/threader.html">http://insulin.brunel.ac.uk/threader/threader.html</a>
ToPLign/123D	S	<a href="http://cartan.gmd.de/Genome/">http://cartan.gmd.de/Genome/</a>

Go to "http://salilab.org/our\_resources.shtml"

# 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.