#### **Clustering methods for Bioinformatics**



Marc A. Marti-Renom Adjunct Assistant Professor



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

## The Problem

	х	У	Z	
а	10	8	10	
b	10	0	9	
С	4	8.6	3	
d	7	8	3	
е	1	2	3	



## Good and Bad?





#### **Homogeneity and Separation Principles**

- Homogeneity: Elements within a cluster are close to each other
- Separation: Elements in different clusters are further apart from each other

...clustering is not an easy task!

Given these points a clustering algorithm might make two distinct clusters as follows





#### **Homogeneity & Separation**



- Microarray data are usually transformed into an intensity matrix (below)
- The intensity matrix allows to make correlations between different genes and to understand how genes functions might be related



#### Example microarray data



# **Clustering techniques**

- Hierarchical
- K-Means
- Clique Graphs
- Example of graph based clustering

#### Hierarchical clustering algorithm

- 1. Form **n** clusters each with one element
- 2. Construct a graph **T** by assigning one vertex to each cluster
- 3. while there is more than one cluster
- 4. Find the two closest clusters  $C_1$  and  $C_2$
- 5. Merge  $C_1$  and  $C_2$  into new cluster **C** with  $|C_1| + |C_2|$  elements
- 6. Compute distance **d** from **C** to all other clusters
- 7. Add a new vertex **C** to **T** and connect to vertices  $C_1$  and  $C_2$
- 8. Remove rows and columns of **d** corresponding to  $C_1$  and  $C_2$
- 9. Add a row and column to **d** corresponding to the new cluster **C**

10. return **T** 

The algorithm takes a *n*×*n* distance matrix *d* of pairwise distances between points as an input.

	$g_1$	$g_2$	$g_3$	$g_4$	$g_6$	$g_6$	$g_7$	$g_8$	$g_9$	$g_{10}$
$g_1$	0.0	8.1	9.2	7.7	9.3	2.3	5.1	10.2	6.1	7.0
$g_2$	8.1	0.0	12.0	0.9	12.0	9.5	10.1	12.8	2.0	1.0
$g_3$	9.2	12.0	0.0	11.2	0.7	11.1	8.1	1.1	10.5	11.5
$g_4$	7.7	0.9	11.2	0.0	11.2	9.2	9.5	12.0	1.6	1.1
$g_5$	9.3	12.0	0.7	11.2	0.0	11.2	8.5	1.0	10.6	11.6
<i>9</i> 6	2.3	9.5	11.1	9.2	11.2	0.0	5.6	12.1	7.7	8.5
$g_7$	5.1	10.1	8.1	9.5	8.5	5.6	0.0	9.1	8.3	9.3
$g_8$	10.2	12.8	1.1	12.0	1.0	12.1	9.1	0.0	11.4	12.4
<i>9</i> 9	6.1	2.0	10.5	1.6	10.6	7.7	8.3	11.4	0.0	1.1
$g_{10}$	7.0	1.0	11.5	1.1	11.6	8.5	9.3	12.4	1.1	0.0















#### Hierarchical clustering algorithm

- 1. Form **n** clusters each with one element
- 2. Construct a graph **T** by assigning one vertex to each cluster
- 3. while there is more than one cluster
- 4. Find the two closest clusters **C**<sub>1</sub> and **C**<sub>2</sub>
- 5. Merge  $C_1$  and  $C_2$  into new cluster **C** with  $|C_1| + |C_2|$  elements
- 6. Compute distance **d** from **C** to all other clusters
- 7. Add a new vertex **C** to **T** and connect to vertices  $C_1$  and  $C_2$
- 8. Remove rows and columns of **d** corresponding to  $C_1$  and  $C_2$
- 9. Add a row and column to **d** corresponding to the new cluster **C**
- 10. return **T**

#### Different ways to define distances between clusters may lead to different clusterings

#### Hierarchical clustering algorithm

Distance between two clusters is the **smallest** distance between any pair of their elements

$$d_{\min}(C^*, C) = \min_{x \in C^*, y \in C} d(x, y)$$

Distance between two clusters is the average distance between all pairs of their elements

$$d_{avg}(C^*, C) = \frac{1}{|C^*||C|} \sum_{x \in C^*, y \in C} d(x, y)$$

Hierarchical Clustering is often used to reveal evolutionary history



#### K-Means clustering algorithm



- Input: A set, V, consisting of n points
- Output: A K points x (cluster center) that minimizes the squared error distortion d(V,x) over all possible choices of x

1-Means Clustering problem is easy (K=1).

However, it becomes very difficult (NP-complete) for more than one center.

An efficient *heuristic* method for K-Means clustering is the Lloyd algorithm

K-Means clustering squared error distortion

- cost(P) = squared error distortion
- cost(P) = 0 for a clustering with all k points in the center of mass of each cluster.

$$d(V,X) = \frac{\sum_{i=1}^{n} d(v_i,X)^2}{n}$$

#### K-Means clustering Lloyd algorithm

- 1. Arbitrarily assign the **k** cluster centers
- 2. While the cluster **centers** keep changing
  - 3. Assign each data point to the cluster  $C_i$  corresponding to the closest cluster representative (center)  $(1 \le i \le k)$
  - 4. After the assignment of **all** data points, compute new cluster **representatives** according to the **center of gravity** of each cluster, that is, the new

cluster representative is 
$$\frac{\sum_{v \in C} v}{|C|}$$
 for all **v** in **C** for every

cluster **C** 

\*This may lead to merely a locally optimal clustering and leads to large rearrangements of nodes







This may lead to merely a locally optimal clustering and leads to large rearrangements of nodes



#### K-Means clustering clustering cost

- The clustering cost measures the quality of a particular partition **P**.
- **cost(P)** = squared error distortion
- cost(P) = 0 for a clustering with all k points in the center of mass of each cluster.

#### K-Means clustering "Greedy" algorithm

- 1. Select an arbitrary partition **P** into **k** clusters
- 2. while forever
- 3. **bestChange**  $\leftarrow 0$
- 4. for every cluster **C**
- 5. for every element i not in C
- 6. if moving i to cluster C reduces its clustering cost
- 7. if  $(cost(P) cost(P_{i \rightarrow C}) > bestChange$
- 8. **bestChange**  $\leftarrow$  cost(P) cost(P<sub>i  $\rightarrow$  c)</sub>
- 9. i\* ← i
- 10. C<sup>∗</sup> ← C
- 11. if **bestChange** > 0
- 12. Change partition **P** by moving **i**\* to **C**\*
- 13. else
- 14. return P

# Clique graph clustering

- A **clique** is a graph with every vertex connected to every other vertex
- A clique graph is a graph where each connected component is a clique





#### Clique graph clustering matrix to graph

- Turn the distance matrix into a distance graph
  - *ie*, genes are represented as vertices in the graph
  - Choose a distance threshold  $\theta$
  - If the distance between two vertices is below  $\theta$ , draw an edge between them
  - The resulting graph may contain cliques
  - These cliques represent clusters of closely located data points!

#### Clique graph clustering corrupted cliques problem

Input: A graph G

Output: The smallest number of additions and removals of edges that will transform *G* into a clique graph

# Clique graph clustering

 A graph can be transformed into a clique graph by (minimally) adding or removing edges



#### Clique graph clustering matrix to graph

The distance graph (threshold  $\theta$ =7) is transformed into a clique graph after removing the two highlighted edges

	$g_1$	$g_2$	$g_3$	$g_4$	$g_{5}$	$g_6$	$g_7$	$g_8$	$g_9$	$g_{10}$
$g_1$	0.0	8.1	9.2	7.7	9.3	2.3	5.1	10.2	6.1	7.0
$g_2$	8.1	0.0	12.0	0.9	12.0	9.5	10.1	12.8	2.0	1.0
$g_3$	9.2	12.0	0.0	11.2	0.7	11.1	8.1	1.1	10.5	11.5
$g_4$	7.7	0.9	11.2	0.0	11.2	9.2	9.5	12.0	1.6	1.1
$g_5$	9.3	12.0	0.7	11.2	0.0	11.2	8.5	1.0	10.6	11.6
$g_6$	2.3	9.5	11.1	9.2	11.2	0.0	5.6	12.1	7.7	8.5
$g_7$	5.1	10.1	8.1	9.5	8.5	5.6	0.0	9.1	8.3	9.3
$g_8$	10.2	12.8	1.1	12.0	1.0	12.1	9.1	0.0	11.4	12.4
<i>9</i> 9	6.1	2.0	10.5	1.6	10.6	7.7	8.3	11.4	0.0	1.1
$g_{10}$	-7.0	1.0	11.5	1.1	11.6	8.5	9.3	12.4	1.1	0.0

(a) Distance matrix, d (distances shorter than 7 are shown in bold). After transforming the distance graph into the clique graph, the dataset is partitioned into three clusters



#### Clique graph clustering corrupted cliques problem

- Corrupted Cliques problem is NP-Hard, some heuristics exist to approximately solve it:
- **CAST** (Cluster Affinity Search Technique): a practical and fast algorithm:
  - **CAST** is based on the notion of genes *close* to cluster *C* or *distant* from cluster *C*
  - Distance between gene *i* and cluster *C*:

d(i, C) = average distance between gene *i* and all genes in C

 Gene i is close to cluster C if d(i,C) < θ and distant otherwise

# Cast algorithm

- 1.  $\mathbf{P} \leftarrow \emptyset$
- 2. while  $S \neq Ø$
- 3.  $\mathbf{v} \leftarrow$  vertex of maximal degree in the distance graph **G**
- 4. **C** ← {**v**}
- 5. while a close gene i not in C or distant gene i in C exists
- 6. Find the nearest close gene i not in C and add it to C
- 7. Remove the farthest distant gene i in C
- 8. Add cluster **C** to partition **P**
- 9. S ← S \ C
- 10. Remove vertices of cluster **C** from the distance graph **G**
- 11. return P

S – set of elements, G – distance graph,  $\theta$  – distance threshold

#### Clustering application defining domains using a graph based clustering method



MCL - a cluster algorithm for graphs http://micans.org/mcl/


More significant

## **ModParts algorithm**

















#### 

# **Repetitions as recurrent fragments**

**Ribosomal protein S6 (1ris)**  $\alpha+\beta \rightarrow Ferrodoxin Like domain$ 



# **Domains as recurrent fragments**

2163 chains from Islam et al. 1995 → 569 Non-redundant <2Å && <30aa diff.

Divide randomly into two sets Remove of incomplete or obsolete entries.

## Training set $\rightarrow$ 242 chains Testing set $\rightarrow$ 234 chains



 $1dhr_{-}$  (dihydropteridine reductase )

# **Domains as recurrent fragments**









# http://www.bioalgorithms.info

### AN INTRODUCTION TO BIOINFORMATICS ALGORITHMS

NEIL C. JONES AND PAVEL A. PEVZNER



### **BMI-206**

## Structure-Structure comparisons Sequence-Structure comparisons

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

February 2nd, 2006

## How to use this lectures

### •Ask!

### Outline

- Basic introduction
- Theory (representation-scoring-optimization)
- Available programs
- Application

## **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



Aminoacid substitutions

 $\mathsf{RMSD} = \sqrt{\sum \left( \mathbf{X}_i - \overline{\mathbf{X}} \right)^2}$ 

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.



Sometimes approximated by Z-score (normal distribution).

### Optimizer Global dynamic programming alignment

remember Patsy's class





### Backtracking to get the best alignment

## Optimizer Global .vs. local alignment

remember Patsy's class



# 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 Align the most similar pair Align next most similar pair Align B-D with A-C New gap in A-C to optimize its alignment with B-D

# **Coverage .vs. Accuracy**





Same RMSD ~ 2.5Å

Coverage ~90% Cα

Coverage ~75% Cα

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



# Structural alignment by properties conservation (SALIGN-MODELLER)

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

A A     A     A     A     A     A     A     A     A     A     A     A     A     A     A     A     A           A   <	0	SALICN Server
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 far.gz) Multiple files can be uploaded by iteratively clicking 'upload' after every file uploaded (rocalized string not found) localized string not found (upload) e-mail address, to receive results: e-mail address, to receive results: cutomit 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	▶ @ ∠	A A C + Mttp://alto.compbio.ucsf.edu/salign-cgi/index.cgi S ^ Q 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 tocalized string not found tocalized string not found topoad Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures summi tocalized string not found current summi tocalized string not found Reference: Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation		
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 tocalized string not found tupload Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures submit tocalized string not found Reference: Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation	SALIGN M	ultiple Structure/Sequence Alianment 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 (rocalized string not found) Lupload Upload Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit calized string not found Reference: Nathusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation		
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 tocalized string not found localized string not found Upload Uploaded files: No files uploaded by tochoose PDB structures e-mail address, to receive results: Submit localized string not found Reference: Mathusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation	SALIGN is a ger	eral alignment module of the modeling program MODELLER
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 (ocalized string not found) localized string not found Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: (Submit) [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	The alignments	are computed using dynamic programming, making use of several features of the protein sequences and structures
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 tocalized string not found localized string not found tupload Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit 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	Users can eithe	er upload their own sequences/structures to align or choose structures from the PDB
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.qz) Multiple files can be uploaded by iteratively clicking 'upload' after every file uploaded tocalized string not found tocalized string not found Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit tocalized string not found Reference: Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation	sequences can	either be pasted or uploaded as FASTA or PIR format alignment files
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 Iocalized string not found Iocalized string not found Upload Uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit Iocalized string not found Reference: Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation	Paste sequence	to align
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 localized string not found Upload Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit 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	Multiple sequen	ces 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 ileratively clicking 'upload' after every file uploaded tocalized string not found localized string not found Upload Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit 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		
Specify file to upload (PIR, FASTA, PDB, zip or .tar.gz) Multiple files can be uploaded by ileratively clicking 'upload' after every file uploaded tocalized string not found localized string not found Upload Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit 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		
Specify file to upload (PIR, FASTA, PDB, zip or.tar.gz) Multiple files can be uploaded by iteratively clicking 'upload' after every file uploaded (tocalized string not found) localized string not found (upload) Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit localized string not found (upload) Submit 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		
Multiple files can be uploaded by iteratively clicking 'upload' after every file uploaded	Specify file to up	load (PIR, FASTA, PDB, .zip or .tar.gz)
Iocalized string not found       Iocalized string not found         Uploaded files:       No files uploaded         Enter 4 letter code(s) to choose PDB structures       Iocalized string not found         e-mail address, to receive results:       Iocalized string not found         Submit       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	Multiple files car	a 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:         Submit       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	localized string	g not found) localized string not found
Uploaded files: No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: Submit 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	Upload	
No files uploaded Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results: submit 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	Uploaded files:	
Enter 4 letter code(s) to choose PDB structures e-mail address, to receive results:  submit 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	No files uploade	ad the second
e-mail address, to receive results: submit 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		
e-mail address, to receive results: Submit 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	Enter 4 letter co	de(s) to choose PDB structures
e-mail address, to receive results:		
e-mail address, to receive results: Submit Cocalized string not found Reference: Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation		
Submit Discalized string not found Submit Incalized string not found Reference: Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation	a mail address	
Submit) 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	e-mail address,	
Submit 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		
Reference: Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation	Submit loca	alized string not found
Madhusudhan, M.S., Marti-Renom, M.A., Eswar, N., Sali, A. SALIGN - a multiple structure/sequence alignment tool, under preparation	Reference:	
SALIGN - a multiple structure/sequence alignment tool, under preparation	Madhusudhan, I	M.S., Marti-Renom, M.A., Eswar, N., Sali, A.
	SALIGN - a mult	iple 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

NCBI CDD pfam00959	with Query Sequence added - Mi	crosoft Internet Explorer		
Ble Edit View Fav	vorites Ipols Help			4
🔇 Back + 🜍 + 🗷	👔 🐔 🔎 Search 👷 Pavorit	tes 🜒 Media 🤬 🍰 👿	- 🖵 🚳	
Iddress 🗿 http://www	nchi nim nih nov /Structure/rdd/rdder	v mižarchin = 28 mavaln = 108 račkona = 38	ed=15128ab=5.0.23 💌 🌄 Go 🛛 Unics	» 🖷
		-		
SNON	Conco	rued Domain Databa	<b>50</b>	
> NCBI	Conse	iveu Domain Databa	30	
PubMed	Nucleotide Protein	Structure CDD	Taxonomy Help?	
CD: pfam0	0959.11. Phage_lysozyme.	Query added PSSM-Id: 1512	Source: Pfam[US]. F	Pfam[Uk
Description: Phage I	ysozyme. This family includes larr	ibda phage lysozyme and E. coli en	dolysin.	
Ftatues Alianmo	of from course	Created: 12 Do	Related: COG3772. Cl	064678
Aligned: 10 rows	antirom source	DSSN: 110 cr	0-2003 Representative: Conserve	
Proteine: IClick by	are for CDAPT summary of Protein	PSSME 110 C	numms representative: Consensus	
Subset Rows	up to 10 💌 sequences m	ost similar to the query 💌		
	10 20	30 40	50 60	
	1 XNTIGIG		*	
21zm (query)	24 YYTIGIGHiltKSPSL	NAAKseldkaigrnongV	ITKD 61	
1PDL_A 1	97 YPTIGIGhl-iMKQPVB	DMAQInkvlsk	gvgreitgnpgsITME 239	
1L47 at 500552	24 YYTIGIGhl-1TKSPSL	NAARSeldkai	grncngvITKD 61	
gi 126600	48 VNTVCHGHTGKDI	MLGK	YTKA 69	
g1 138699	29 IPTIGVGhtgkVDGNSV	ASGNT	ITAE 54	
g1 126602	30 HYTIGYGHYGSDV	SPRQVV	ITAK 51	
g1 6016519	30 YYTIGYGHYGSDV	MPCQV	ISEE 51	
91 2010440	17 HF11010INLOV	1012	RI 36	
				>

✓ Good scoring system with significance

### **Reduces the protein representation**



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

# **Vector Alignment Search Tool (VAST)**

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



# **Incremental combinatorial extension (CE)**





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

## **Incremental combinatorial extension (CE)**

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



# Matching molecular models obtained from theory (MAMMOTH)



Ortiz AR, (2002) Protein Sci. 11 pp2606

# Matching molecular models obtained from theory (MAMMOTH)

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



# **Classification of the structural space**



#### Alex Adai

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

# **SCOP**<sub>1.65</sub> database

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



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

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

Class	Number of folds	Number of superfamilies	Number of families	
All alpha proteins	218	376	608	
All beta proteins	144	290	560	
Alpha and beta proteins (a/b)	136	222	629	
Alpha and beta proteins (a+b)	279	409	717	
Multi-domain proteins	46	46	61	
Membrane and cell surface	47	88	99	
Small proteins	75	108	171	
Total	945	1539	2845	

# CATH<sub>2.5.1</sub> database

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



Uses FSSP for superimposition

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

### Semi-automatic classification Domain boundaries definition

۲	۵	0		8		0	D
Mainly Alpha	5	251	465	1402	2189	3705	14105
Mainly Beta	19	160	311	1443	2961	4329	18771
Alpha Beta	14	414	706	3014	4781	7660	33080
Few Secondary Structures	1	82	90	144	232	285	1098
Preliminary single domain assigments	10	808	809	906	967	1090	3012
Multi-domain domains	1	12	12	16	25	36	109
CATH-35 Sequence families	1	4707	4707	4719	4768	4862	6168
	1	22	22	27	33	38	198

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

# DBAliv2.0 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

#### **Pairwise structure alignments**

Last update:	January 30th, 2006
Number of chains:	73,817
Number of structure-structure comparisons:*	1,013,210,249
Multiple structure alignm	ients
Last update:	November 16th, 2005
Number of representative chains:	24,828
Number of families:	9.588
## 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 derived 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
<u>http://www.sbg.bio.ic.ac.uk/~3dpssm/</u>
PredictProtein server at EMBL
<u>http://www.embl-heidelberg.de/predictprotein/predictprotein.html</u>
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),
 <a href="http://portray.bmc.uu.se/cgi-bin/dennis/spasm.pl">http://portray.bmc.uu.se/cgi-bin/dennis/spasm.pl</a>

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.

Finkelstein et al. (1995) Proteins 23, pp142

# **Sequence-Structure alignments**

#### As any other bioinformatics problem...

- Representation
- Scoring
- Optimizer

# Representation Sequence/Structures



Reduced atoms representation

Secondary Structure

Accessible surface

#### 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]:

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

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 Potentials (background)

#### Structural space



Sequence space

MKLLIVLTCISLCSCICTVVQRCASNKPHVLEDPCKVQH HLSVNQCVLLPQCCPKSCKICTHLISIEVVLTCRAVDKM MHVNCVEQCSLQDCIKIAPRVLKTCILCVLKPCLTSVSH VHLVQPTSCCCKKNCICHVEIRSLDILTKSVQLACLVPM

MQCCRVQKICDLLAVELCKLHISTPSCKILCVVTSVPHN

# Scoring Statistical Potential... Hydrogen Bonds

Long range free energy



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

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

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

$$\boldsymbol{W}_{NO}(r) = -kT\ln\left(\boldsymbol{g}_{NO}(r)\right)$$

Sippl (1996). JMB 260 pp644

# Scoring Statistical Potential... Distance Potentials



Long range free energy

Sippl (1993). JCAM 7 pp473

#### Scoring Raw scores of an alignment



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.







#### Scoring Significance of an alignment (score)

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



#### Optimizer Global dynamic programming alignment

remember Patsy's class





#### 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!!!!**



MENFEIWVEKYRPRTLDEVVGQDEVIQRLKGYVERKNIPHLLFSGPPGTGKTATAIALARDLFGENWRDN FIEMNASDERGIDVVRHKIKEFARTAPIGGAPFKIIFLDEADALTADAQAALRRTMEMYSKSCRFILSCN YVSRIIEPIQSRCAVFRFKPVPKEAMKKRLLEICEKEGVKITEDGLEALIYISGGDFRKAINALQGAAAI GEVVDADTIYQITATARPEEMTELIQTALKGNFMEARELLDRLMVEYGMSGEDIVAQLFREIISMPIKDS LKVQLIDKLGEVDFRLTEGANERIQLDAYLAYLSTLAKK

## Domain boundaries from sequence (SnapDragon)



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

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

## Domain boundaries from sequence and predicted SSE (DomSSEA)

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



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

#### **Prediction of Secondary Structure (PSI-PRED)**



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

#### **Prediction of Secondary Structure (PSI-PRED)**

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

PSIPRED Protein Structure Predictio File Edit <u>V</u> iew F <u>a</u> vorites <u>T</u> ools	i Server - Microsoft Internet Expl <u>H</u> elp	orer	
3 Back 🔹 🕥 🕤 🗷 😰 🏠 🔎	Search 📌 Favorites 🔏 Media	🛛 🖉 · 📚 🔟 · 🖵 📖 🥸	
goress (a) http://bomhcs.ud.ac.uk/pat		Bioinformatics Unit	
PSIPRE home>	D The	PSIPRED Protein Structure Prediction Server	
Info	We sug the <u>PSI</u>	gest that you do not bookmark this page as it is liable to move. It is best to access the server via <u>PRED home page</u> , which has more information about the methods and a full reference list.	
Input Sequ	Help Input se ence	equence (single letter code)	
Choos Predia Metho	e OPred tion OFold d OFold	lict Secondary Structure (PSIPRED v2.4) lict Transmembrane Topology (MEMSAT) Recognition(GenTHREADER - quick) Recognition (mGenTHREADER - with profiles and predicted secondary structure)	
Filteri Option	ng Pasi Masi NS Masi Warning:	k low complexity regions k transmembrane helices k coiled-coil regions Turn off all filtering if you are running MEMSAT	
Subm Seque	E-mail a ence Passwo Short na Predict	Iddress <u>Help</u> ord (only required for commercial e-mail addresses) <u>Help</u> ame for sequence <u>Help</u> Clear form	
			🔮 Internet

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



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

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

) 🖯		SALIGN Server		
	A A C + Mhttp	://alto.compbio.ucsf.edu/sal	ign–cgi/index.cgi	😡 ^ 🔍 Google
SALIGN M	ultiple Structure/Segu	ence Alignment Ser	ver	
OALIONI		enec Angiment oer	VCI	
SALIGN is a ger	eral alignment module of the mod	eling program MODELLER		
The alignments	re computed using dynamic prog	ramming, making use of severa	I features of the protein s	equences and structures
Users can eithe	upload their own sequences/st	ructures to align or choose str	uctures from the PDB	
sequences can	ither be pasted or uploaded as F	ASTA or PIR format alignment fil	es	
Paste sequence	o align			
Multiple sequen	es can be pasted by iteratively cli	cking 'upload' after every pasted	d sequence	
			7	
Specify file to up	oad (PIR, FASTA, PDB, .zip or .ta	.gz)		
Multiple files car	be uploaded by iteratively clickin	g 'upload' after every file upload	ed	
localized string	not found localized string not	ound		
Indexed				
opioad				
Uploaded files:				
No files uploade	I.			
Enter 4 letter co	e(s) to choose PDB structures			
e-mail address	receive results:			
o	in control for and in the second seco			
Submit loca	ized string not found			
Reference:				
Madhusudhan,	I.S., Marti-Renom, M.A., Eswar, N	, Sali, A.		
SALIGN - a mult	ble structure/sequence alignment	tool, under preparation		

## **Threading (mGenThreader)**



McGuffin LJ, Jones DT. (2003) Bioinformatics, 19, pp874

## **Threading (mGenThreader)**

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

🗿 PSIPR	ED Protein Structure Prediction	Server - Microsoft Internet Explorer	- 🗆 🛛
<u>File E</u> o	dit <u>V</u> iew F <u>a</u> vorites <u>T</u> ools <u>t</u>	Help	
G Back	- 🕄 🛛 🖻 🔹 🖉	Search 📌 Favorites 🜒 Media 🕢 😥 💀 💌 🔻 🦵 🏭 🦄	
A <u>d</u> dress	http://bioinf.cs.ucl.ac.uk/psipr	ed/psiform.html 👽 🔁 Go Links *	' 🔁 -
		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 <u>PSIPRED home page</u> , 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)	
2	Filtering Options	Help Mask low complexity regions Mask transmembrane helices Mask coiled-coil regions Marning: Turn off all filtering if you are running MEMSAT	~
Done Done		Unternet 🔮	

## **Remote homology detection (FUGUE)**



## **Remote homology detection (FUGUE)**

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



#### **Meta-Servers (3D-Jury)**



#### **Meta-Servers (3D-Jury)**

#### http://bioinfo.pl/Meta/

🚰 Meta Server Job L	ist, Biolnfo.PL - Microsoft Internet Explorer				
<u>File E</u> dit <u>V</u> iew F	avorites Tools Help			A.	
🔇 Back 🔹 🕥 🐇 😰 🏠 🔎 Search 👷 Favorites 🜒 Media 🕢 🍰 🐨 🍃 🏭 🐨 🖵 🏭 🦓					
Address 🙆 http://bio	info.pl/Meta/		🖌 🄁 Co	Links » 📆 🗸	
BIOINFO.PL:MEI	A Meta Server Job List	[ABOUT] [SERVERS]	[ [BENCHMARKS]	[STATUS]	
		Skip:	Queue:		
	O jobs from 64.54.249. in the last week	PDB-Blast	1		
		✓ 3D-Jigsaw	43		
	Your E-mail:	ESyPred3D	1		
	Target Name:	GRDB	1		
	Amino Acid Sequence only (in one letter code):	FFAS03	1		
		Sam-T99	1		
		SUPERFAMILY	1		
	×	INBGU	39		
	Reset Clear Format Submit	FUGUE2	1		
		3D-PSSM	1		
	Please submit domains separately Please remove coiled coil regions	mGenTHREADER			
	Check LiveBench for evaluation of the reliability of the servers Results are stored only for 1 month	psipred			
	Jobs queued for more than 7 days for servers with queue>30 are skipped	profsec	1		
	Please contact us in case of problems with interpretation of results	Pcons2	1		
	Please contact us if You plan larger analysis projects Some servers return only models, no alignments (target sequence is shown)	3D-ShotGun 3D-Jury	11		
	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]				
				×	
e			🥝 Interne	et 💡	

# **Iterative process... better models(?)**



#### Evaluate Model Modify and build Model



## Moulding: iterative alignment, model building, model assessment



# Iterative process... MOULDER





John & Sali (2003). NAR 31 pp3982

st in iteratio

20

10

Iteration index

eration inde

Average of top 10
 Best of all
 Final

25

Top
 Final

25

### Genetic algorithm operators



Also, "two point crossover" and "gap deletion".
### Composite model assessment score

Weighted linear combination of several scores:

Pair (Pp) and surface (Ps) statistical potentials;

Structural compactness (S<sub>C</sub>);

Harmonic average distance score (H<sub>a</sub>);

Alignment score (A<sub>S</sub>).

 $Z = 0.17 Z(P_P) + 0.02 Z(P_S) + 0.10 Z(S_C) + 0.26 Z(H_a) + 0.45 (A_S)$ 

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

- $\mu$  ... average score of all models
- $\sigma \dots$  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)

	Sequence identity [%]	Coverage [% aa]	Initial prediction		Final prediction		Best prediction	
Target -template			Ca RMSD [Å]	CE overlap [%]	Ca RMSD [Å]	CE overlap [%]	Ca 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



# 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



Nup85 Jakes & Bar (10,4 Autor) A JARDIG MELLA

Nup84 Abd. td. Abic 84 (0.6 td. 4.8 5560. ) (1) (0.6) 4.8

Nup145C h. sharle cal al least a state that a share

Seh1 JAAAA JALJAAA JAAA

Sec13 (III)(III)(III)(III)

# 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



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





# 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...



http://www.salilab.org/modeller

A. Šali & T. Blundell. J. Mol. Biol. 234, 779, 1993.
J.P. Overington & A. Šali. Prot. Sci. 3, 1582, 1994.
A. Fiser, R. Do & A. Šali, Prot. Sci., 9, 1753, 2000.

### **Steps in Comparative Protein Structure Modeling**



### TARGET **TEMPLATE** ASILPKRLFGNCEQTSDEG I KIERTPI VPHISAONVCI KI **DDVPERLIPERASFQWMN** DK ASILPKRLFGNCEQTSDEGLKIERTPLVPHISAQNVCLKIDDVPERLIPE MSVIPKRLYGNCEOTSEEAIRIEDSPIV---TADLVCLKIDEIPERLVGE blbp.B99990001 -SEUDO ENERGY 0.8 0.0 -0.8 -1.6

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.

60 80 RESIDUE INDEX

100

120

-2.4

20

40

### Typical errors in comparative models

MODEL X-RAY TEMPLATE

Region without a template



Incorrect template



Distortion/shifts in aligned regions



Misalignment



Sidechain packing



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

Model Accuracy as a Function of Target-Template Sequence Identity



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

### Model Accuracy

#### **HIGH ACCURACY**

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



#### MEDIUM ACCURACY

CRABP Seq id 41% Cα equiv 122/137 RMSD 1.34Å



Sidechains Core backbone Loops Alignment LOW ACCURACY

EDN Seq id 33% Cα equiv 90/134 RMSD 1.17Å



Sidechains Core backbone Loops Alignment Fold assignment

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

#### **Applications of Protein Structure Models**



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

# Obtaining MODELLER and related information

- MODELLER (8v0) web page
- http://www.salilab.org/modeller/
  - Download Software (Linux/Windows/Mac/Solaris)
  - HTML Manual
  - ♦ Join Mailing List







# **Using MODELLER**

# No GUI! 😕

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

# Using MODELLER

### INPUT:

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

### OUTPUT:

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

# Modeling of BLBP Input

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

>P1;blbp

sequence:blbp::::::::

VDAFCATWKLTDSQNFDEYMKALGVGFATRQVGNVTKPTVIISQEGGKVVIRTQCTFKNTEINFQLGEEFEETSID DRNCKSVVRLDGDKLIHVQKWDGKETNCTREIKDGKMVVTLTFGDIVAVRCYEKA\*

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

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

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

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

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

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

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

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

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

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

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

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

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

PDB file Can be viewed with Chimera http://www.cgl.ucsf.edu/chimera/ Rasmol http://www.openrasmol.org



### Model file → blbp.B9990001

#### http://www.salilab.org/bioinformatics\_resources.shtml



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



### 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.