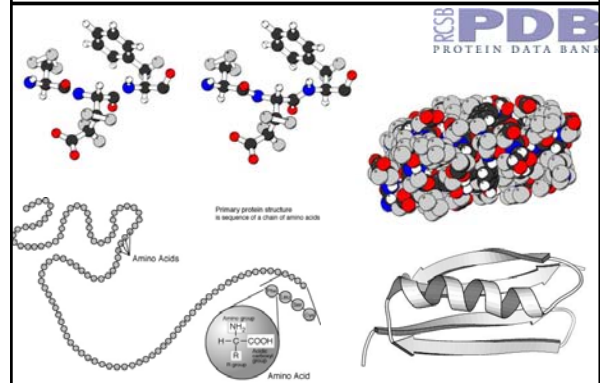


Protein Folding and Structure Prediction

Ingo Ruczinski

Department of Biostatistics, Johns Hopkins University

Proteins



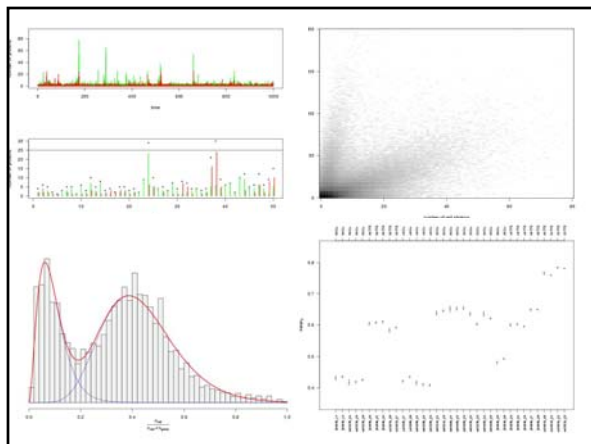
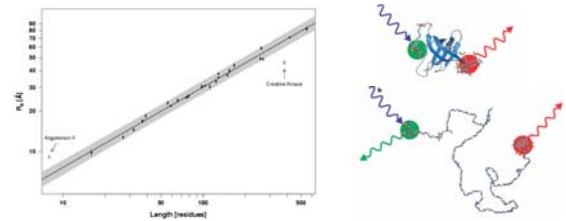
Protein Folding vs Structure Prediction

- Protein folding is concerned with the process of the protein taking its three dimensional shape.
- Protein structure prediction is solely concerned with the 3D structure of the protein, using theoretical and empirical means to get to the end result.

This presentation is bit about both. The role of statistics in protein folding and structure prediction is usually to support or discredit some hypothesis based on physical principles, and to build stochastic models.

Some Issues in Protein Folding

- Do chemically denatured proteins behave as random coils?
- What determines the folding rate of a protein?
- Are amino acids in proteins conserved because of folding kinetics?



Flavors of Structure Prediction

- Homology modeling,
- Fold recognition (threading),
- Ab initio (de novo, new folds) methods.

ROSETTA is mainly an ab initio structure prediction algorithm, although various parts of it can be used for other purposes as well (such as homology modeling).

Ab Initio Methods

- Ab initio: "From the beginning".
- Assumption 1: All the information about the structure of a protein is contained in its sequence of amino acids.
- Assumption 2: The structure that a (globular) protein folds into is the structure with the lowest free energy.
- Finding native-like conformations require:
 - A scoring function (potential).
 - A search strategy.

Rosetta

- The scoring function is a model generated using various contributions. It has a sequence dependent part (including for example a term for hydrophobic burial), and a sequence independent part (including for example a term for strand-strand packing).
- The search is carried out using simulated annealing. The move set is defined by a fragment library for each three and nine residue segment of the chain. The fragments are extracted from observed structures in the PDB.

The Rosetta Scoring Function

$$P(\text{structure}|\text{sequence}) \propto P(\text{sequence}|\text{structure}) \times P(\text{structure})$$

Sequence dependent:

- hydrophobic burial
- residue pair interaction

Sequence independent:

- helix-strand packing
- strand-strand packing
- sheet configurations
- vdW interactions

The Sequence Dependent Term

$$P(\text{aa}_1, \dots, \text{aa}_n | X) =$$

$$\prod_i P(\text{aa}_i | X) \times \prod_{i < j} \frac{P(\text{aa}_i, \text{aa}_j | X)}{P(\text{aa}_i | X)P(\text{aa}_j | X)} \times \prod_{i < j < k} \frac{P(\text{aa}_i, \text{aa}_j, \text{aa}_k | X)P(\text{aa}_i | X)P(\text{aa}_j | X)P(\text{aa}_k | X)}{P(\text{aa}_i, \text{aa}_j | X)P(\text{aa}_i, \text{aa}_k | X)P(\text{aa}_j, \text{aa}_k | X)} \times \dots$$

The Sequence Dependent Term

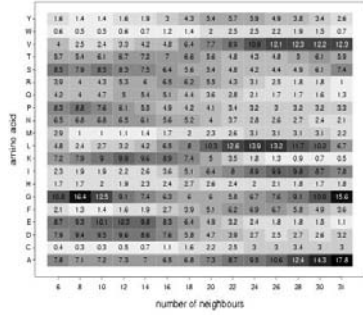
$$P(\text{sequence}|\text{structure}) \approx P_{\text{env}} \times P_{\text{pair}}$$

$$P_{\text{env}} = \prod_i P(\text{aa}_i | E_i)$$

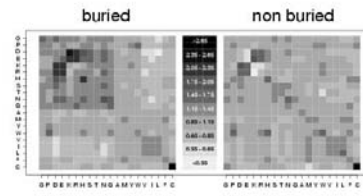
$$P_{\text{pair}} = \prod_{i < j} \frac{P(\text{aa}_i, \text{aa}_j | E_i, E_j, r_{ij})}{P(\text{aa}_i | E_i, r_{ij})P(\text{aa}_j | E_j, r_{ij})}$$

ID	length	Exptl.	resolution	B-factor	FreeRivision
1P2BL	52	EXPT	1.920	0.10	0.14
1P2GA	91	EXPT	0.870	0.14	0.15
1L3LA	74	EXPT	0.920	0.14	0.19
1Q2AL	23	EXPT	1.540	0.19	0.21
1Q2TA	32	EXPT	1.000	0.20	0.20
1Q2TA	81	EXPT	1.700	0.20	0.21
1Q2NB	95	EXPT	1.200	0.17	0.19
1Q2KA	20	EXPT	1.450	0.19	0.22
1Q2LA	44	EXPT	0.540	0.09	0.09
1E24A	72	EXPT	1.700	0.19	0.22
1Q2RA	99	EXPT	1.800	0.18	0.21
1Q2RZ	79	EXPT	1.700	0.20	0.23
1Q2RY	93	EXPT	1.700	0.20	0.23
1P2GA	95	EXPT	1.250	0.13	0.17
1Q2RA	62	EXPT	1.120	0.15	0.20
1P2LA	91	EXPT	0.000	0.00	1.00
1E2R0	83	EXPT	0.920	0.07	1.00
1E2R0	80	EXPT	1.800	0.19	1.00
1L2TA	82	EXPT	1.900	0.20	0.28
110GA	62	EXPT	1.700	0.20	0.23
1Q2LA	69	EXPT	1.800	0.20	0.22
1Q2EA	83	EXPT	1.800	0.17	1.00
1Q2GA	84	EXPT	1.400	0.14	0.20
1Q2EA	41	EXPT	1.850	0.19	0.22
1Q2EB	77	EXPT	1.800	0.17	1.00
1Q2GL	74	EXPT	0.950	0.10	0.13
1Q2PA	100	EXPT	1.350	0.14	0.18
1Q2EA	77	EXPT	1.900	0.19	0.22
1Q2RA	30	EXPT	1.850	0.17	1.00
1Q2CO	99	EXPT	1.250	0.14	1.00
1Q2QA	93	EXPT	1.240	0.17	0.19
1Q2TA	71	EXPT	0.970	0.11	1.00
1E2TA	61	EXPT	1.040	0.15	0.17
1E2B0	62	EXPT	1.400	0.18	1.00
1Q2FO	55	EXPT	1.700	0.17	0.23

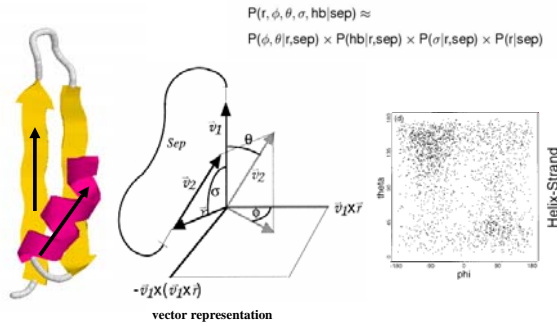
Hydrophobic Burial



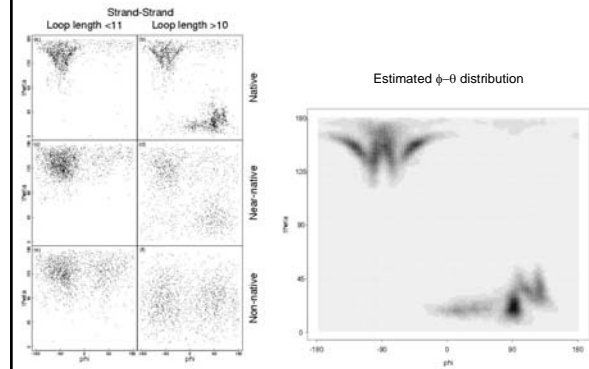
Residue Pair Interaction



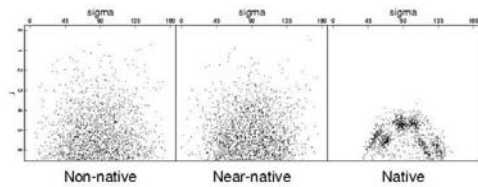
The Sequence Independent Term



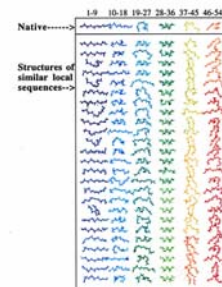
Strand Packing – Helps!

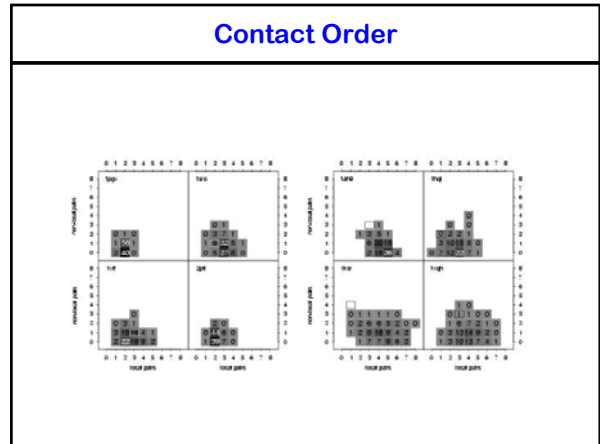
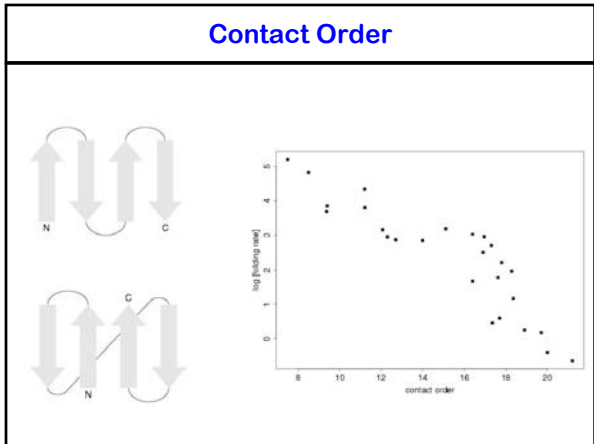
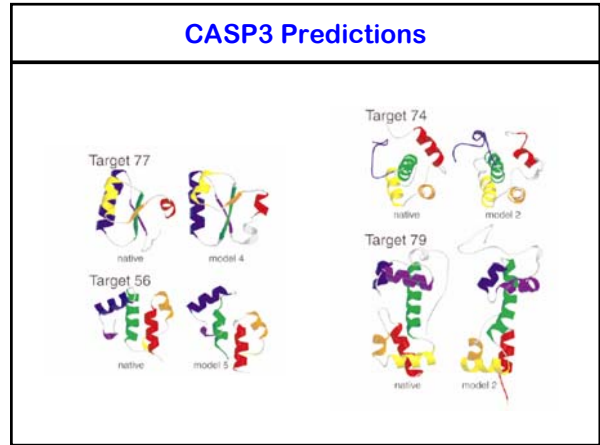
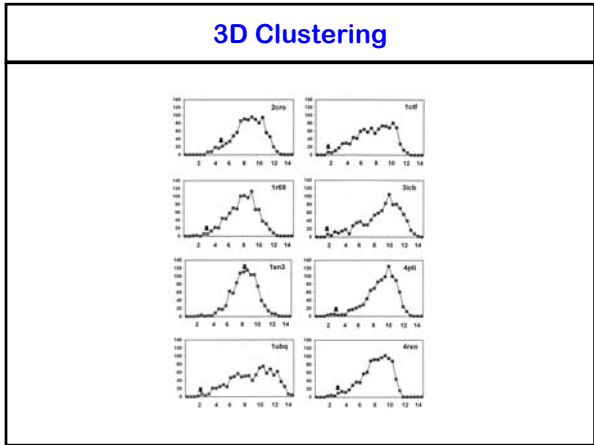
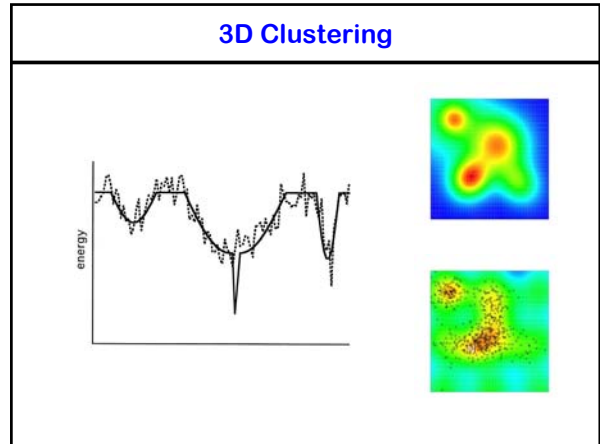
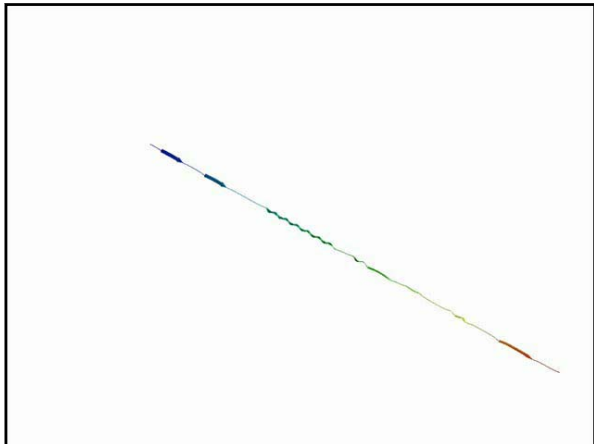


Shear Angles – Help not!

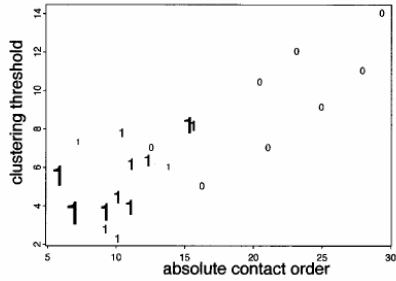


Fragment Selection

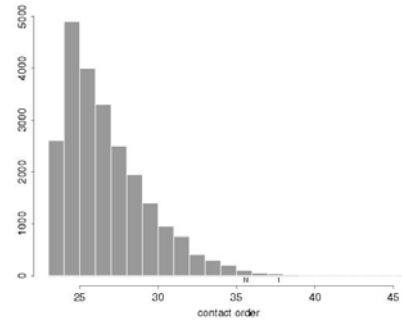




Clustering and Contact Order



Decoy Enrichment in CASP4

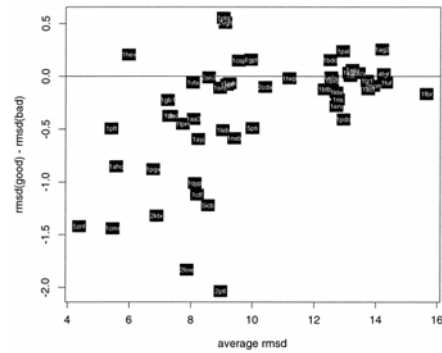


A Filter for Bad β -Sheets

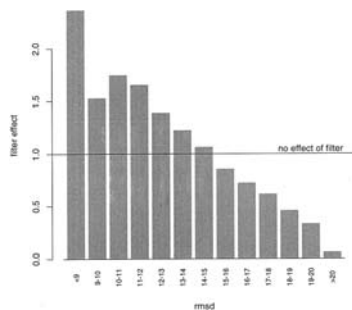
Many decoys do not have proper sheets. Filtering those out seems to enhance the rmsd distribution in the decoy set. Bad features we see in decoys include:

- No strands,
- Single strands,
- Too many neighbours,
- Single strand in sheets,
- Bad dot-product,
- False handedness,
- False sheet type (barrel),
- ...

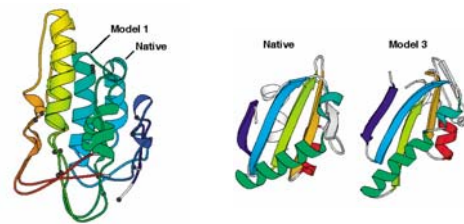
A Filter for Bad β -Sheets

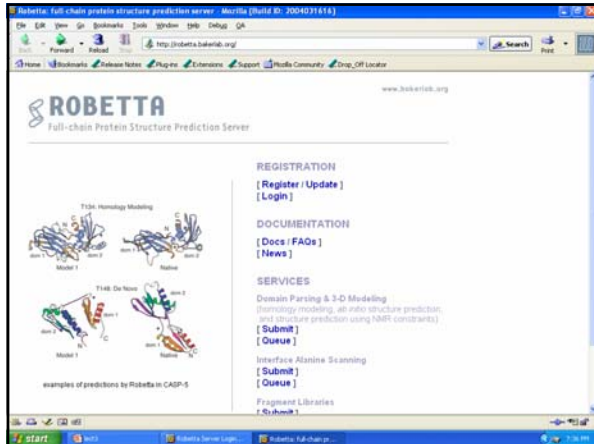


A Filter for Bad β -Sheets



Rosetta in CASP4





Applications and Other Uses of Rosetta

- Other uses of Rosetta:
 - Homology modeling.
 - Rosetta NMR.
 - Protein interactions (docking).
- Applications of Rosetta:
 - Functional annotation of genes.
 - Novel protein design.

Functional Inference

Collaborators

Collaborators = People who I troubled way more than I should have.

David Baker	University of Washington
Kevin Plaxco	UC Santa Barbara
Richard Bonneau	Institute for Systems Biology
Chris Bystroff	Rensselaer Polytechnic Institute
Dylan Chivian	University of Washington
Charles Kooperberg	Fred Hutchinson Cancer Research Center
Carol Rohl	UC Santa Cruz
Kim Simons	Harvard University
Charlie Strauss	Los Alamos National Laboratory
Jerry Tsai	Texas A&M

