

## Protein Structure Prediction using ROSETTA

Ingo Ruczinski

Department of Biostatistics, Johns Hopkins University

## Protein Folding vs Structure Prediction

- Protein folding is concerned with the process of the protein taking its three dimensional shape. The role of statistics is usually to support or discredit some hypothesis based on physical principles.
- 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 about the latter.

## 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 Humble Beginnings

- Kim Simons and David Baker tackle ab initio structure prediction (1995/96).
- A bit later, Charles Kooperberg and Ingo Ruczinski join the project.
- Two publications appear:
  - Simons et al (1997): Assembly of protein tertiary structures from fragments with similar local sequences using simulated annealing and Bayesian scoring functions, JMB 268, pp 209-25.
  - Simons et al (1999): Improved recognition of native-like protein structures using a combination of sequence-dependent and sequence-independent features of proteins, Proteins 34, pp 82-95.
- With the help of Richard Bonneau and Chris Bystroff, Rosetta is used for the first time on unknown targets in CASP3 (1998).

## The Rosetta Scoring Function

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

- |  |   |
|--|---|
| Sequence dependent:  | Sequence independent:   |
| <ul style="list-style-type: none"> <li>hydrophobic burial</li> <li>residue pair interaction</li> </ul> | <ul style="list-style-type: none"> <li>helix-strand packing</li> <li>strand-strand packing</li> <li>sheet configurations</li> <li>vdW interactions</li> </ul> |

## The Sequence Dependent Term

$$P(aa_1, \dots, aa_n | X) = \prod_i P(aa_i | X) \times \prod_{i < j} \frac{P(aa_i, aa_j | X)}{P(aa_i | X)P(aa_j | X)} \times \prod_{i < j < k} \frac{P(aa_i, aa_j, aa_k | X)P(aa_i | X)P(aa_j | X)P(aa_k | X)}{P(aa_i, aa_j | X)P(aa_i, aa_k | X)P(aa_j, aa_k | X)} \times \dots$$

## The Sequence Dependent Term

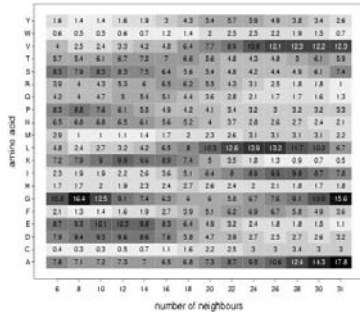
$$P(\text{sequence}|\text{structure}) \approx P_{\text{env}} \times P_{\text{pair}}$$

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

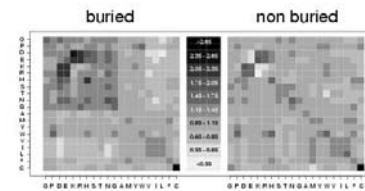
$$P_{\text{pair}} = \prod_{i < j} \frac{P(aa_i, aa_j | E_i, E_j, r_{ij})}{P(aa_i | E_i)P(aa_j | E_j, r_{ij})}$$

ID	length	Depth	resolution	B-factor	RMSD
1F2EL	82	DEPT	1.920	0.20	0.18
1F2GA	91	DEPT	0.970	0.14	0.15
1L5LA	74	DEPT	0.850	0.14	0.19
1Q2AA	23	DEPT	1.580	0.19	0.21
1Q2TA	32	DEPT	1.000	0.20	0.20
1P7IA	81	DEPT	1.700	0.20	0.21
1Q2NS	95	DEPT	1.200	0.17	0.19
1Q2CA	20	DEPT	1.450	0.19	0.22
1E2GA	44	DEPT	0.540	0.09	0.09
1E2LA	72	DEPT	1.700	0.19	0.22
1Q2RA	99	DEPT	1.900	0.18	0.21
1Q2RE	79	DEPT	1.700	0.20	0.25
1Q2FF	90	DEPT	1.700	0.20	0.15
1P2GA	95	DEPT	1.250	0.13	0.17
1Q2BA	62	DEPT	1.120	0.15	0.20
1F2LA	91	DEPT	2.000	0.20	1.00
1P2NO	53	DEPT	0.920	0.07	1.00
1P2FO	80	DEPT	1.800	0.19	1.00
1L2TA	82	DEPT	1.900	0.20	0.28
1Q2GA	42	DEPT	1.700	0.20	0.23
1Q2LA	89	DEPT	1.800	0.20	0.22
1Q2EA	83	DEPT	1.800	0.17	1.00
1E2DA	84	DEPT	1.400	0.16	0.20
1Q2EA	84	DEPT	1.850	0.19	0.22
1E2ED	77	DEPT	1.800	0.17	1.00
1Q2EA	74	DEPT	0.950	0.10	0.19
1Q2FA	100	DEPT	1.250	0.16	0.20
1L2KA	77	DEPT	2.000	0.19	0.22
1Q2GA	30	DEPT	1.450	0.17	1.00
1P2CO	99	DEPT	1.220	0.15	1.00
1P2GA	93	DEPT	1.240	0.17	0.19
1Q2IA	71	DEPT	0.970	0.12	1.00
1L2TA	81	DEPT	1.040	0.15	0.17
1E2DO	62	DEPT	1.400	0.18	1.00
1P2FO	55	DEPT	1.700	0.17	0.23

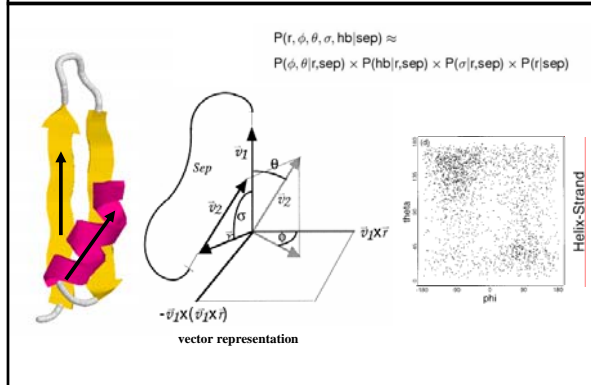
## Hydrophobic Burial



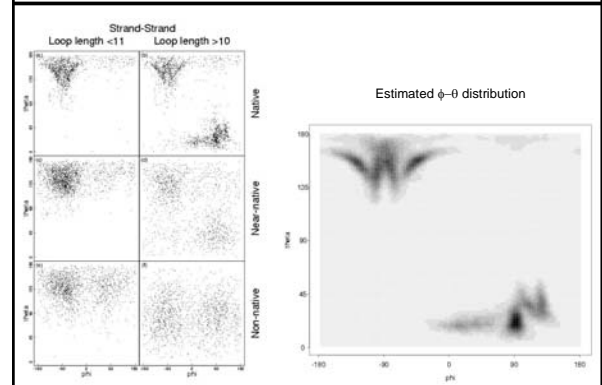
## Residue Pair Interaction



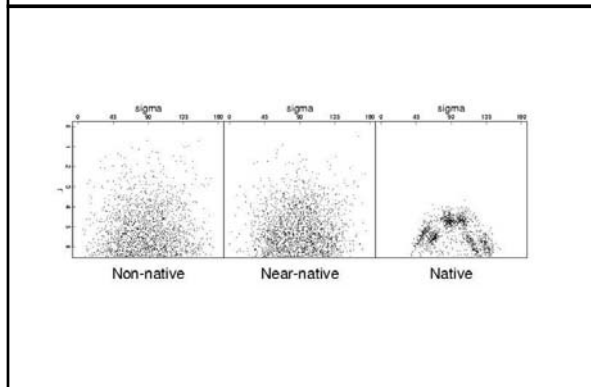
## The Sequence Independent Term



## Strand Packing – Helps!



## Shear Angles – Help not!



## The Model

$$P(\text{structure}) = P_A^{w_A} P_B^{w_B} P_C^{w_C}, \quad w_X > 0.$$

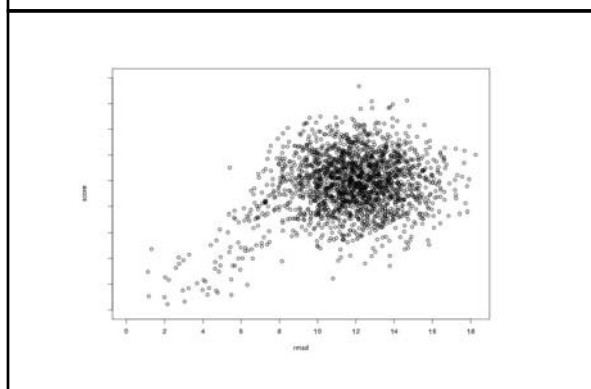
$$-\log P(\text{structure}|\text{sequence}) \propto$$

$$-\log P(\text{sequence}|\text{structure}) - \log P(\text{structure})$$

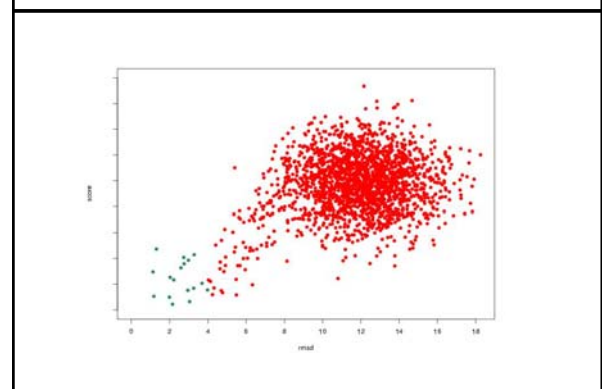
$$g(\text{rmsd}) = w_{\text{pssm}} + w_{\text{hb}} \log P_{\text{hb}} + w_{\text{ss}} \log P_{\text{ss}} + w_{\text{vdw}} \text{VdW} +$$

$$w_{\text{shear}} \log P_{\text{shear}} + w_{\text{sep}} (\log P_{\text{sep}} + \log P_{\text{sep}})$$

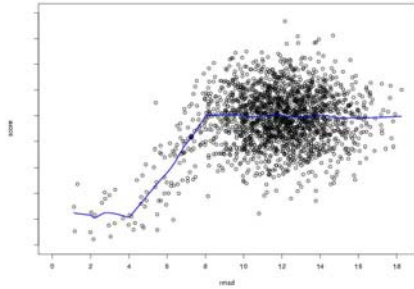
## Parameter Estimation



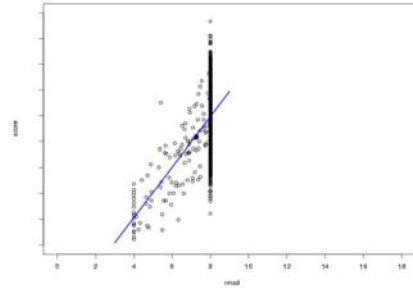
## Parameter Estimation



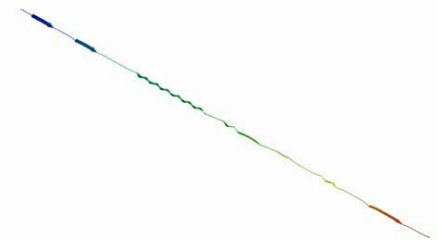
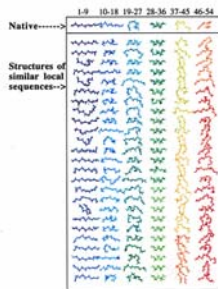
### Parameter Estimation



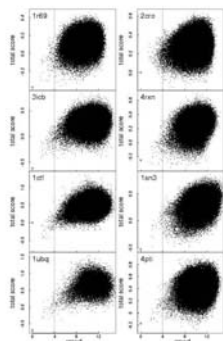
### Parameter Estimation



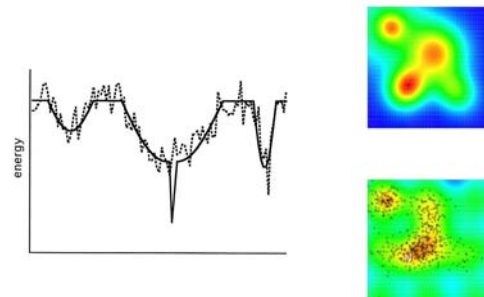
### Fragment Selection



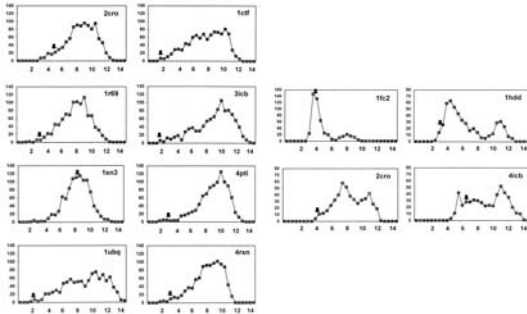
### Validation Data Set



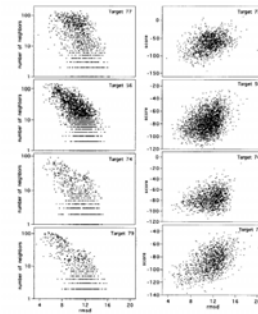
### 3D Clustering



### 3D Clustering



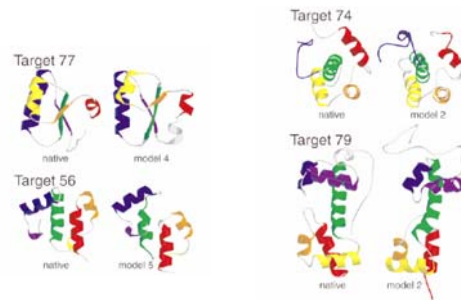
### 3D Clustering in CASP3



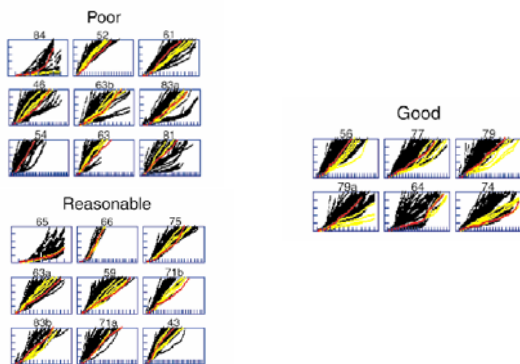
### CASP3 Protocol

- Construct a multiple sequence alignment from  $\phi$ -blast.
- Edit the multiple sequence alignment.
- Identify the ab initio targets from the sequence.
- Search the literature for biological and functional information.
- Generate 1200 structures, each the result of 100,000 cycles.
- Analyze the top 50 or so structures by an all-atom scoring function (also using clustering data).
- Rank the top 5 structures according to protein-like appearance and/or expectations from the literature.

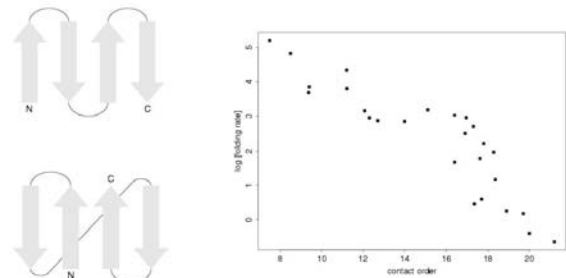
### CASP3 Predictions



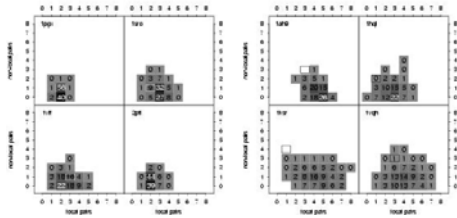
### CASP3 Results



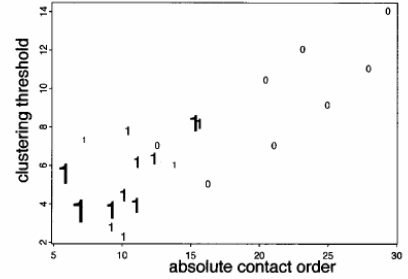
### Contact Order



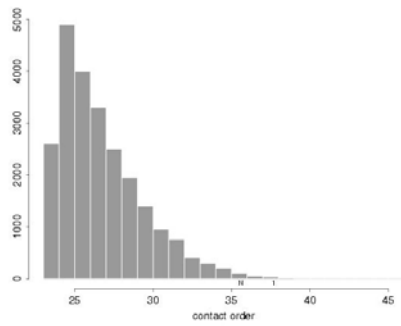
### Contact Order



### Clustering and Contact Order



### Decoy Enrichment in CASP4

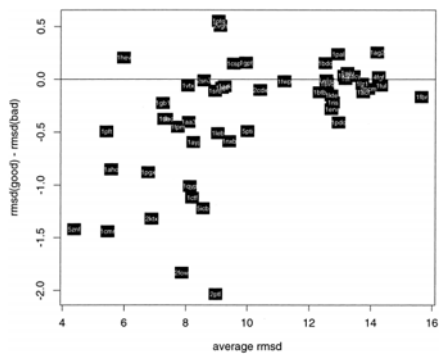


### A Filter for Bad $\beta$ -Sheets

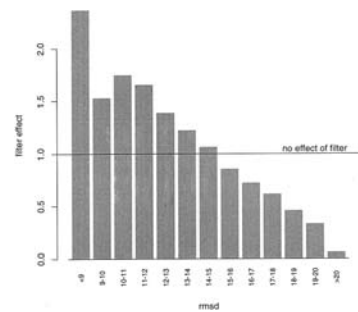
Many decoys do not have proper sheets. Filtering those out seems to enhance the rmsd distribution in the decoy set. Bad features we see in decoys include:

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

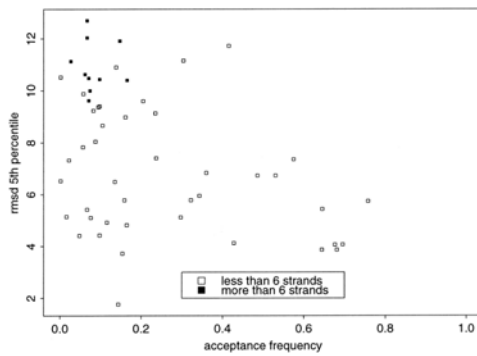
### A Filter for Bad $\beta$ -Sheets



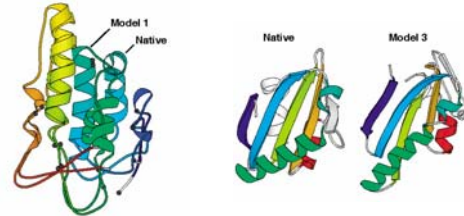
### A Filter for Bad $\beta$ -Sheets



## A Filter for Bad $\beta$ -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.

## Collaborators

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

David Baker	University of Washington
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

## Rosetta Developers

