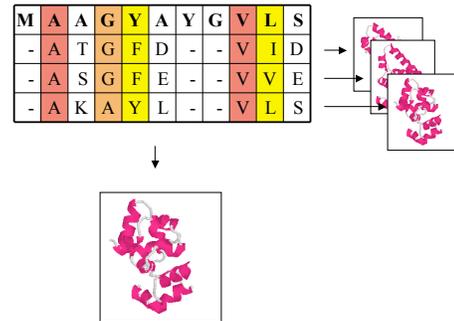


# Protein Structure Prediction

Ingo Ruczinski

Department of Biostatistics, Johns Hopkins University

## Homology Modeling

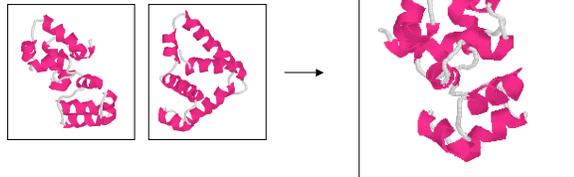


## Fold Recognition

Sequence:

M A A G Y A V L S

+  
Known folds



## Ab Initio Structure Prediction

M A A G Y A V L S

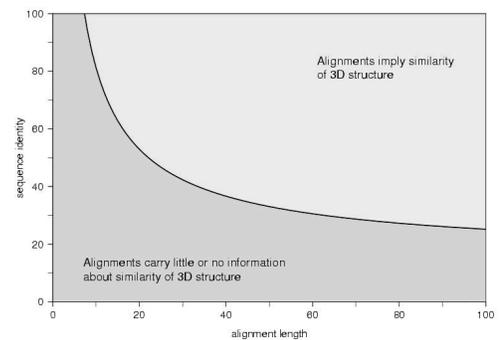


## Homology Modeling

- Align sequence to protein sequences with known structure.
- Construct and evaluate model of 3D structure from alignment.
- Requirement: Close match to template sequences with known 3D structure (sequence similarity of at least 25%).

Note: about 25% of the protein sequences in the Swiss-Prot database have templates for at least part of the sequence!

## Threshold for Structural Homology



Rost B, Protein Engineering 12 (1999).

## Homology Modeling Approach

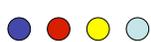
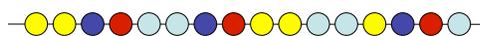
1. Find set of sequences related to target sequence.
2. Align target sequence to template sequences (key step).
3. Construct 3D model for core (backbone):
  - Conserved regions → conserved structure / coordinates.
  - Structure diverges → use sequence similarity, secondary structure prediction, manual prediction, etc. to fill in gaps.
4. Construct 3D models for loops:
  - Search loop conformation library, limited protein folding.
5. Model location of side chains
  - Search rotamer library, use molecular dynamics.
6. Optimize / verify the model
  - Improve likelihood / ensure legality of model.

## Quality Assessment

- Goal
  - Ensure predicted 3D structure is possible / probable in practice
  - Based on general knowledge of protein structures
- Criteria
  - Carbon backbone conformations allowed (Ramachandran map)
  - Legal bond lengths, angles, dihedrals
  - Peptide bonds are planar
  - Side chain conformations correspond to ones in rotamer library
  - Hydrogen-bonding of polar atoms if buried
  - Proper environments for hydrophobic / hydrophilic residues
  - No bad atom-atom contacts
  - No holes inside 3D structure
  - Solvent accessibility

## Fold Recognition

- The input sequence is threaded on different folds from a library of known folds.
- Using scoring functions, we get a score for the compatibility between the sequence and the structures.



Amino acids with different chemical properties

Library of known folds:



## Homology Modeling Web Pages

MODELLER

<http://salilab.org/modeller/modeller.html>

SWISS-MODEL

<http://www.expasy.org/swissmod/SWISS-MODEL.html>

## Quality Assessment Programs

VERIFY3D

[http://shannon.mbi.ucla.edu/DOE/Services/Verify\\_3D](http://shannon.mbi.ucla.edu/DOE/Services/Verify_3D)

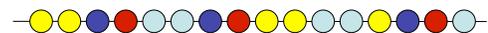
PROCHECK

<http://www.biochem.ucl.ac.uk/~roman/procheck/procheck.html>

WHATIF

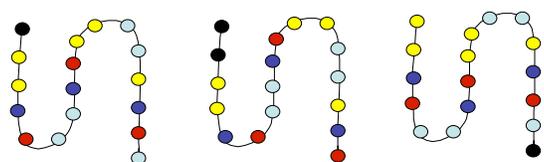
<http://www.cmbi.kun.nl/whatif/>

## Fold Recognition



- Hydrogen donor (blue circle)
- Hydrogen acceptor (red circle)
- Hydrophobic (yellow circle)
- Glycin (light blue circle)

Good score!



## Fold Recognition

---

- This method is less accurate than homology modeling, but can be applied in more cases.
- When the real fold of the input sequence is not represented in the structural database, we do not get a good solution (duh).
- The most important part is the accuracy of the scoring function. The scoring function is the major difference between the approaches used for fold recognition.

## Profile Based Scoring Functions

---

- In methods based on structural profiles, for every fold a profile is built based on structural features of the fold and the compatibility of every amino acid to the features.
- The structural features of each position are based on the combination of secondary structure, solvent accessibility, and the properties of the local environment (such as hydrophobicity, etc).

## Contact Potentials

---

- This method is based on predefined tables which include (pseudo-energetic) scores for each interaction of two amino acids.
- This method makes use of a distance matrix for the representation of different folds.
- For each pair of amino acids which are close in space, the interaction energy is summed up. The total sum is the indication for the "fitness" of the sequence for the given structure .

## Web Sites for Fold Recognition

---

3D-PSSM  
<http://www.bmm.icnet.uk/~3dpssm>

LIBRA I  
[http://www.ddbj.nig.ac.jp/htmls/Email/libra/LIBRA\\_I.html](http://www.ddbj.nig.ac.jp/htmls/Email/libra/LIBRA_I.html)

UCLA DOE  
<http://www.doe-mpi.ucla.edu/people/frsvr/frsvr.html>

123D  
<http://www-Immb.ncifcrf.gov/~nicka/123D.html>

PROFIT  
<http://lore.came.sbg.ac.at/home.html>

## Ab Initio Methods

---

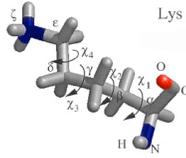
- Ab initio: "From the beginning".
- Assumption 1: All the information about the structure of a protein is contained in its sequence of amino acids.
- Assumption 2: The structure that a (globular) protein folds into is the structure with the lowest free energy.
- Finding native-like conformations require: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

## Representations of the Protein

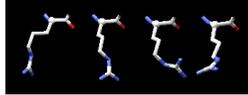
---

- Sidechain: represented as all atoms, rotamers, carbon  $\alpha$  or  $\beta$ , centroids.
- Backbone: torsion angles restricted to discrete values commonly seen in known structures (using a small set of pre-selected  $\phi$ - $\psi$  angles, angles chosen from secondary structure elements, selection of fragments of known structures), secondary structure rigid bodies, lattice models.

## Rotamer Libraries



Some members of the rotamer library:



## Potential Functions

- So-called “molecular mechanics” potentials model the force that determine protein conformation using physically based functional forms (van der Waals, Coulomb).
- Potentials empirically derived from known structures in the Protein Data Bank.

## Search Strategies

- Molecular dynamics. Not really feasible for ab initio prediction per se.
- Probabilistic search algorithms (simulated annealing, genetic algorithms) generate ensembles of candidate structures. Additional methods to discriminate between those are needed.

## 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

$$P(aa_1, \dots, aa_n|X) =$$

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

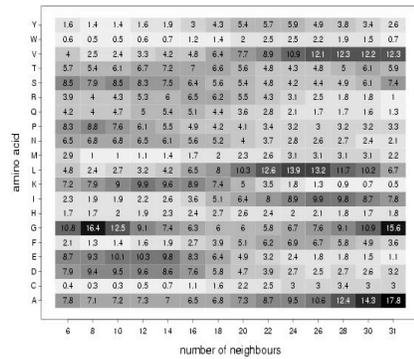
# The Sequence Dependent Term

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

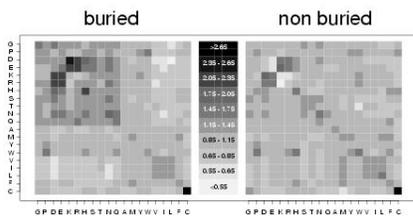
$$P_{\text{env}} = \prod_i P(\text{aa}_i|E_i)$$

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

# Hydrophobic Burial



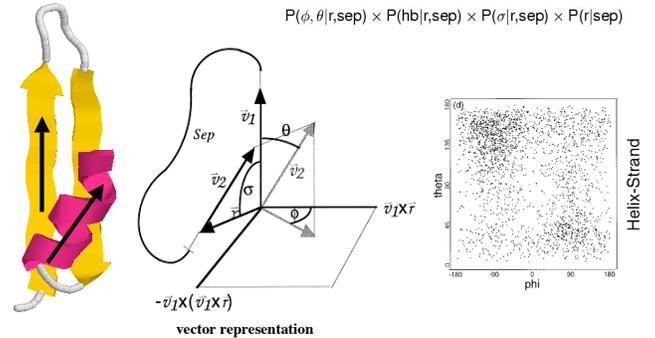
# Residue Pair Interaction



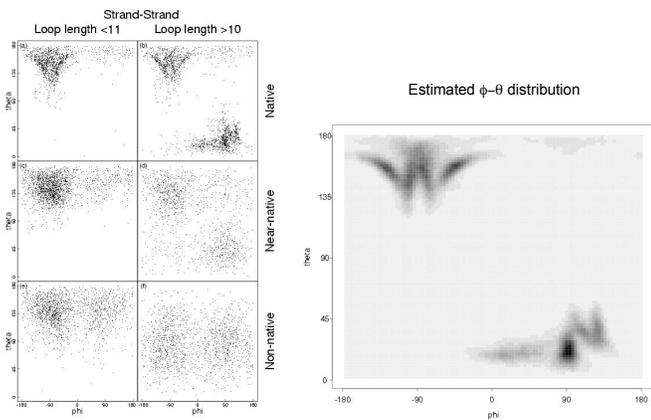
# The Sequence Independent Term

$$P(r, \phi, \theta, \sigma, \text{hb}|\text{sep}) \approx$$

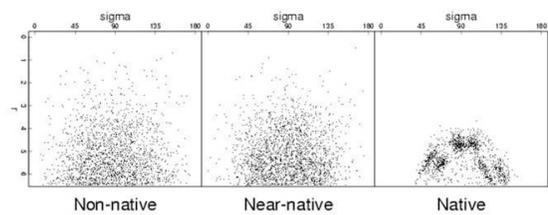
$$P(\phi, \theta|r, \text{sep}) \times P(\text{hb}|r, \text{sep}) \times P(\sigma|r, \text{sep}) \times P(r|\text{sep})$$



# Strand Packing – Helps!



# Shear Angles – Help not!



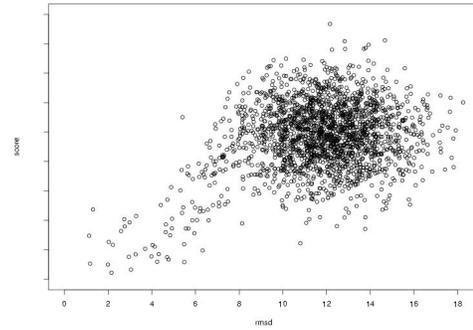
## The Model

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

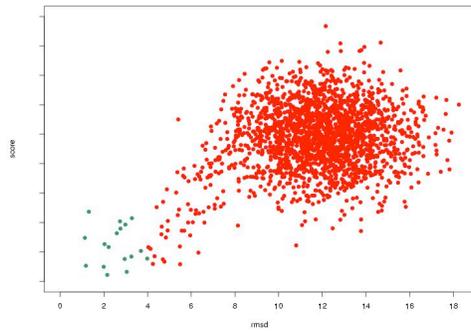
- $\log P(\text{structure}|\text{sequence}) \propto$
- $\log P(\text{sequence}|\text{structure}) - \log P(\text{structure})$

$$g(\text{rmsd}) = w_{\text{protein}} + w_{\text{HS}} \log P_{\text{HS}} + w_{\text{SS}} \log P_{\text{SS}} + w_{\text{VdW}} \text{VdW} + w_{\text{sheet}} \log P_{\text{sheet}} + w_{\text{seq}} (\log P_{\text{env}} + \log P_{\text{par}})$$

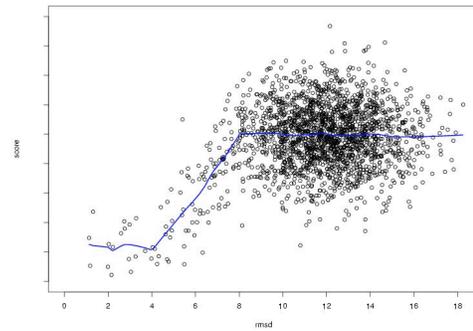
## Parameter Estimation



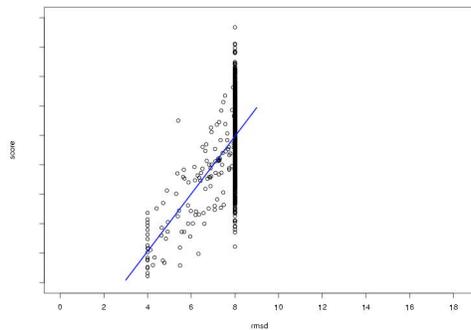
## Parameter Estimation



## Parameter Estimation



## Parameter Estimation

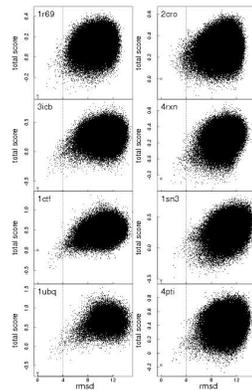


A screenshot of a Mozilla browser window displaying the 'Decoys 'R' Us' website. The page title is 'Decoys 'R' Us - Mozilla (Build ID: 2004031616)'. The address bar shows 'http://www3.oup.co.uk/jwr/database/summary/365'. The website content includes:

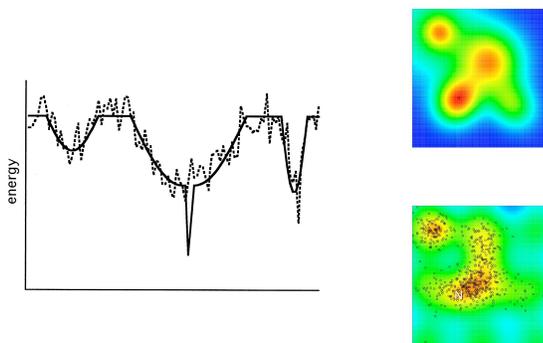
- Navigation links: HOME, HELP, FEEDBACK, SUBSCRIPTIONS, ARCHIVE, SEARCH ARTICLES, TABLE OF CONTENTS
- Section: **Decoys 'R' Us**
- URL: <http://d4.stanford.edu/>
- Contact: [ram@cslb.stanford.edu](mailto:ram@cslb.stanford.edu)
- Database Description: Computer-generated protein conformations based on sequence data
- Category: Structure Databases; Subcategory: Protein structure
- Utility links: Completion Paper, Category List, Alphabetical List, Category/Paper List, Search Summary Pages

The browser's status bar at the bottom shows 'Start', 'Int3', 'Decoys 'R' Us - Mozill...', and the time '7:10 PM'.

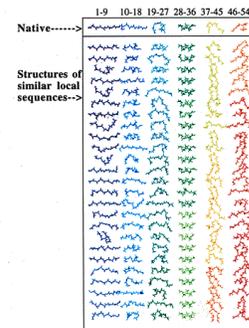
## Validation Data Set



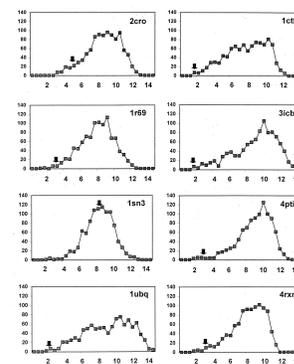
## 3D Clustering



## Fragment Selection



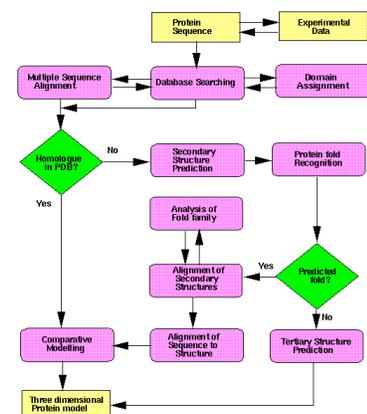
## 3D Clustering



## Assessing Structure Prediction

- CASP (Critical Assessment of Protein Structure Prediction)
  - Competitions measuring current state of the art in protein structure prediction.
  - Researchers predict structure of actual protein sequences.
  - Compare with laboratory determination of structure.
  - Held in 1994, 1996, 1998, 2000, 2002, 2004.
- CAFASP (Critical Assessment of Fully Automated Protein Structure Prediction).

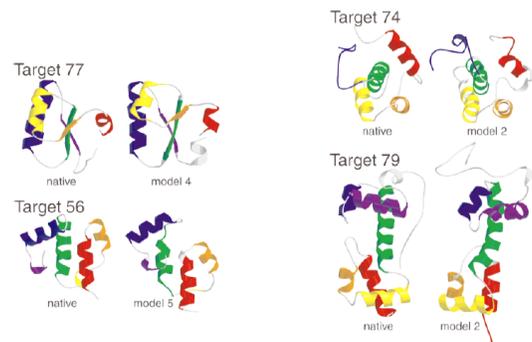
## Protein Structure Prediction



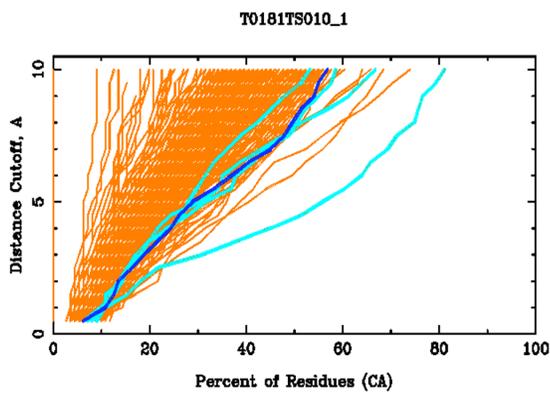
## CASP3 Protocol

- 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.

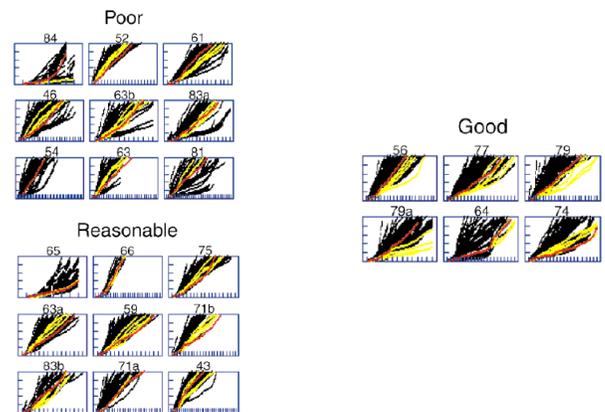
## CASP3 Predictions



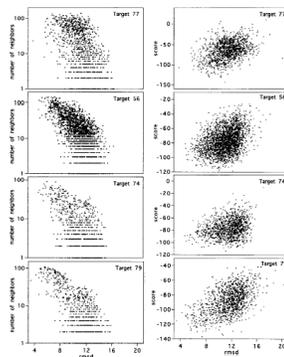
## Hubbard Plot



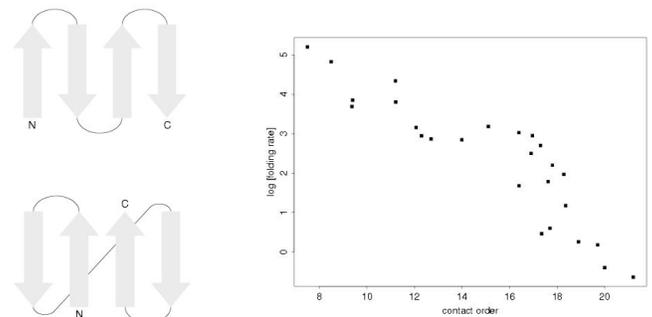
## CASP3 Results



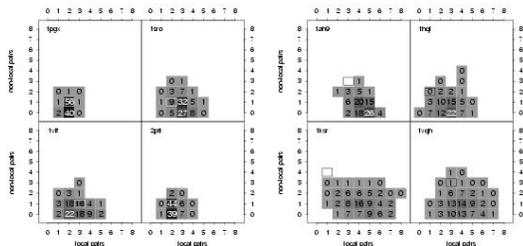
## 3D Clustering in CASP3



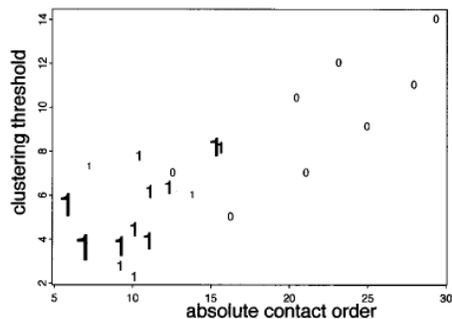
## Contact Order



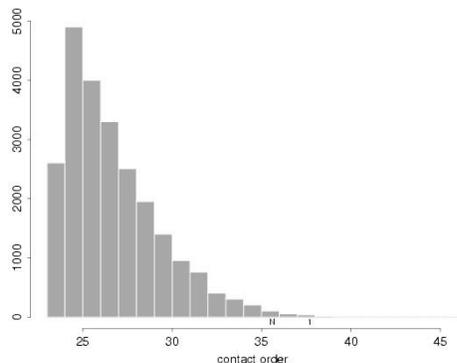
## Contact Order



## Clustering and Contact Order



## Decoy Enrichment in CASP4

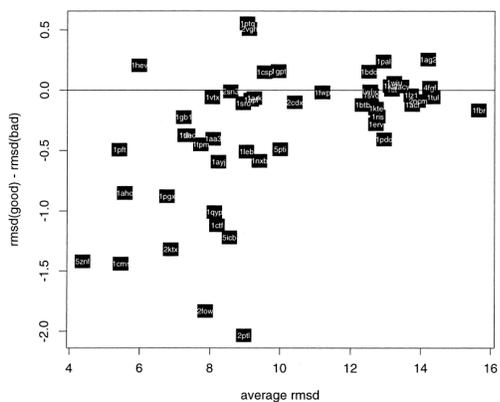


## A Filter for Bad $\beta$ -Sheets

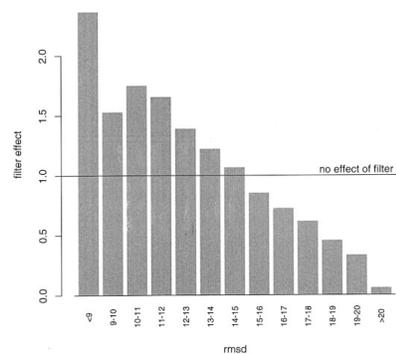
Many decoys do not have proper sheets. Filtering those out seems to enhance the rmsd distribution in the decoy set. Bad features we see in decoys include:

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

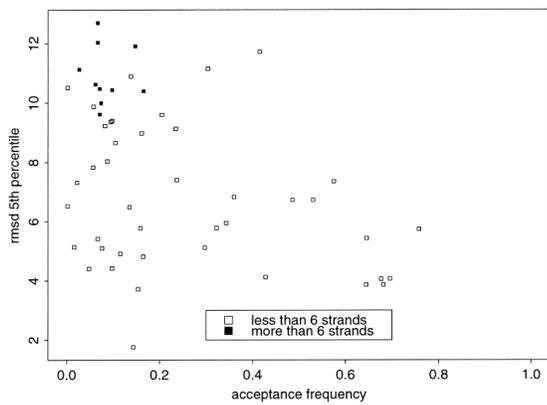
## A Filter for Bad $\beta$ -Sheets



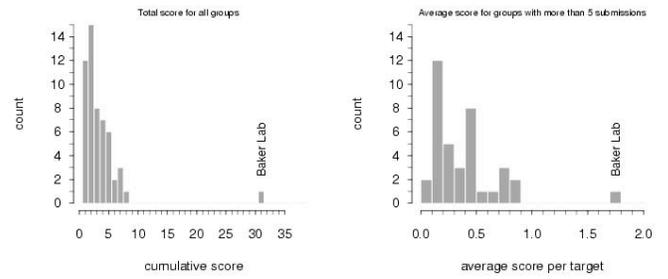
## A Filter for Bad $\beta$ -Sheets



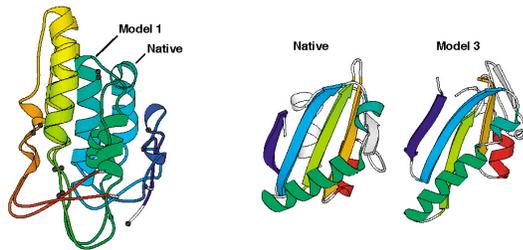
## A Filter for Bad $\beta$ -Sheets



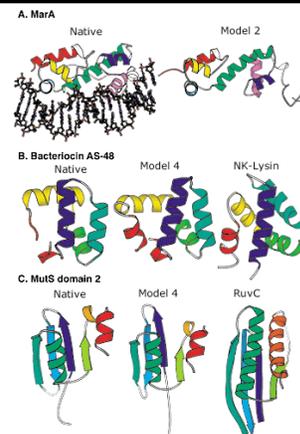
## CASP 4



## Rosetta in CASP4



## CASP 4



## 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.

A screenshot of the Rosetta Full-chain Protein Structure Prediction Server website. The page features the Rosetta logo and navigation links. It includes sections for 'REGISTRATION' (with links for Register/Update and Login), 'DOCUMENTATION' (with links for Docs/FAQs and News), and 'SERVICES' (listing Domain Parsing & 3-D Modeling, Interface Alanine Scanning, and Fragment Libraries). Examples of predictions by Rosetta in CASP-5 are shown, including T134: Homology Modeling and T148: De Novo. The browser window shows the URL http://rosetta.bakerlab.org/ and the time 7:36 PM.