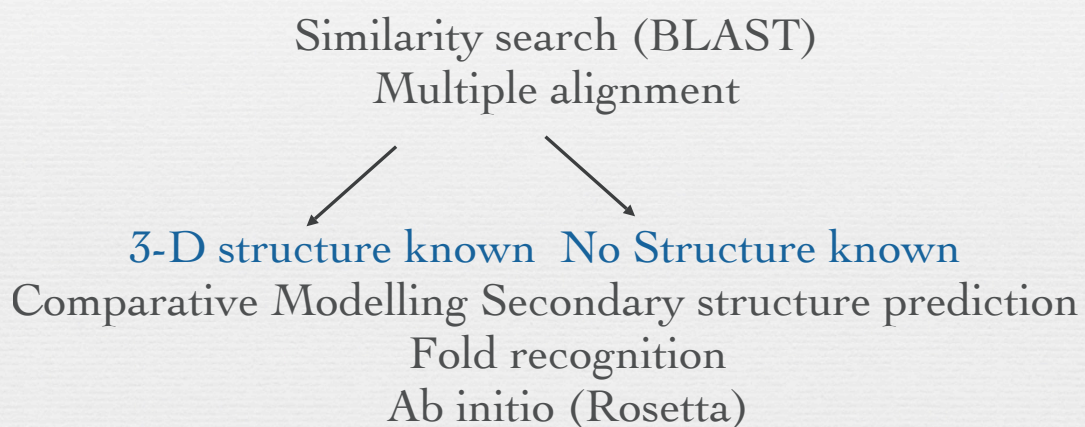


# Structure prediction, fold recognition and homology modelling

Marjolein Thunnissen  
Lund September 2009

## Steps in protein modelling



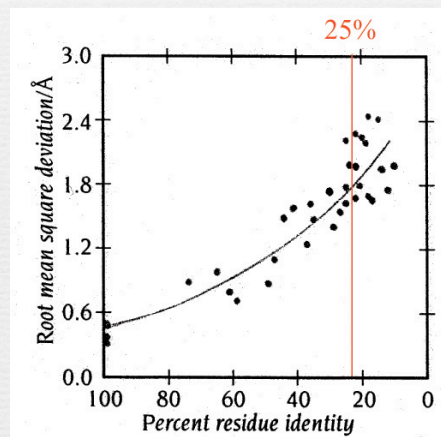
# Structure Prediction

- 1) **Prediction of secondary structure.**
  - a) Method of Chou and Fasman
  - b) Neural networks
  - c) hydrophobicity plots
- 2) **Prediction of tertiary structure.**
  - a) Ab initio structure prediction
    - b) Threading
      - 1D-3D profiles
    - c) Knowledge based potentials
  - c) Homology modelling

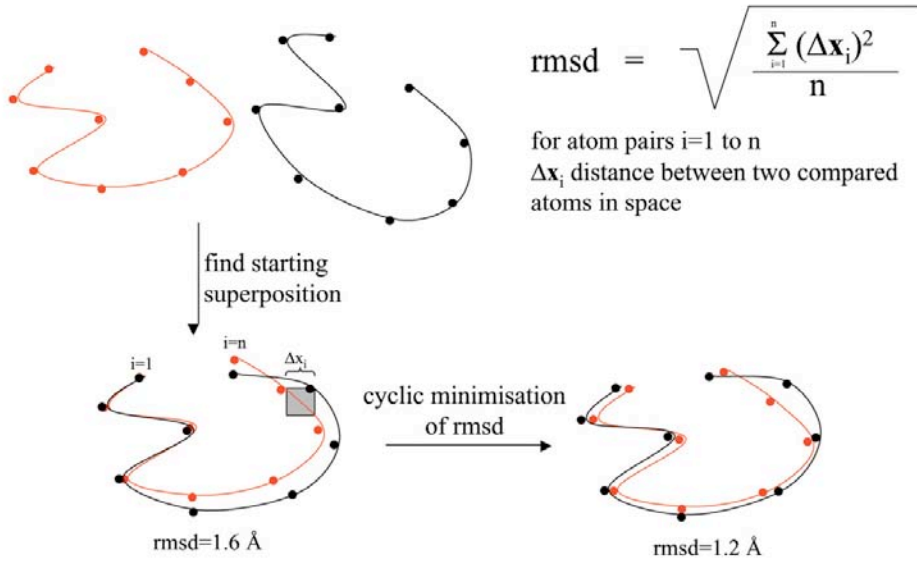
## How does sequence identity correlate with structural similarity

Analysis by Chotia and Lesk (89)

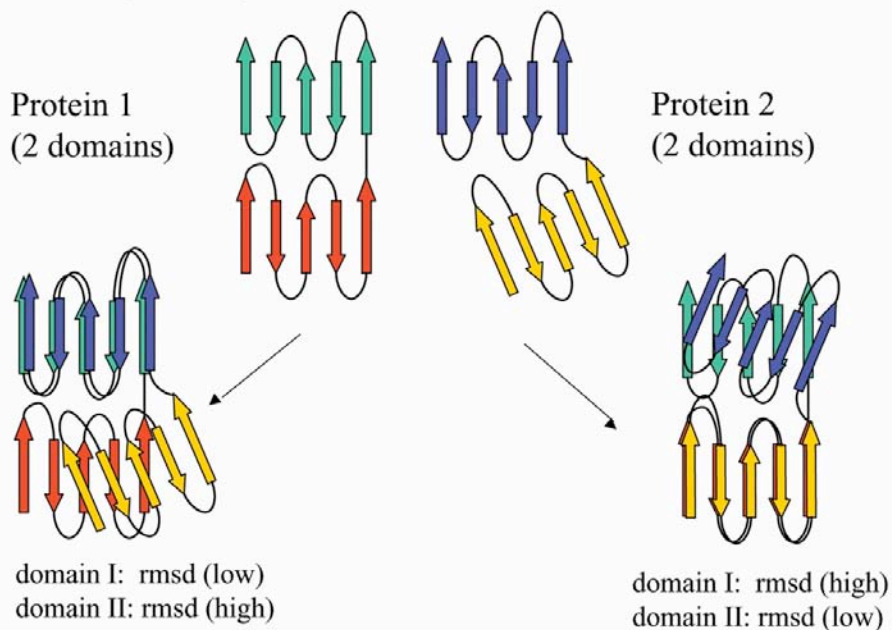
- 100% sequence identity: rmsd = experimental error
- <25% (twilight zone), structures might be similar but can also be different
- Rigid body movements make rmsd bigger

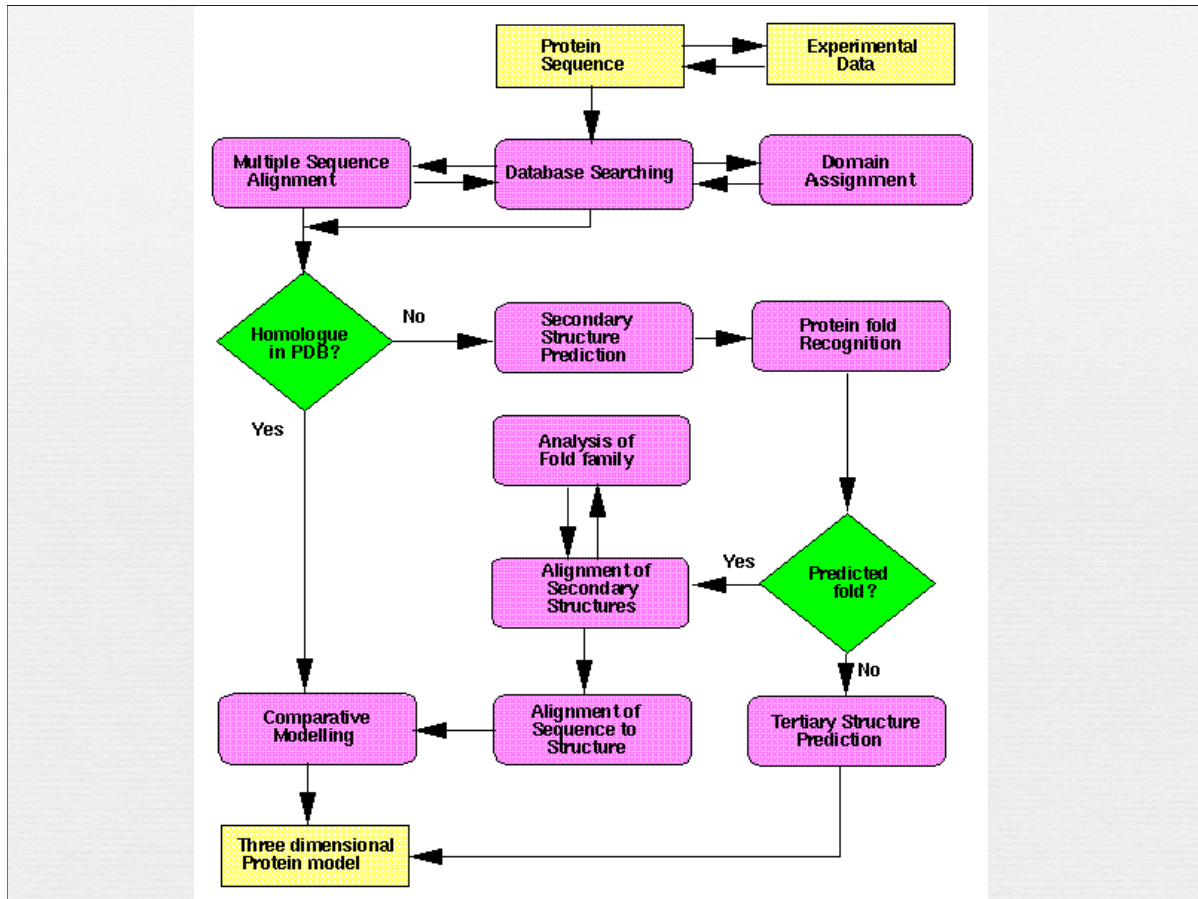


## Root mean square deviation (rmsd)



## Rigid body movements increase the rmsd





## Secondary structure prediction:

Take the sequence and, using rules derived from known structures, predict the secondary structure that is most likely to be adopted by each residue

## Why secondary structure prediction ?

- A major part of the general folding prediction problem.
- The first method of obtaining some structural information from a newly determined sequence. Rules governing  $\alpha$ -helix and  $\beta$ -sheet structures provide guidelines for selecting specific mutations.
- Assignment of sec. str. can help to confirm structural and functional relationship between proteins when sequences homology is weak (used in threading experiments).
- Important in establishing alignments during model building by homology; the first step in attempts to generate 3D models

## *Some interesting facts 2nd structure predictions*

- based on primary sequence only
- accuracy 64% -75%
- higher accuracy for  $\alpha$ -helices than  $\beta$ - strands
- accuracy is dependent on protein family
- predictions of engineered proteins are less accurate

## Methods:

- **Statistical methods** based on studies of databases of known protein structures from which structural propensities for all amino acids are calculated. However, these methods do not take into account physico-chemical knowledge about proteins.
- **Physico-chemical methods** (helical wheels, hydrophobicity profiles etc.).
- **Hybrid methods** combines the first two.

## Structural Propensities

- Due to the size, shape and charge of the side chain, each amino acid may “fit” better in one type of secondary structure than another.
- Classic example: The rigidity and main chain angle of proline cannot be accommodated in an  $\alpha$ -helical structure.

## Examples of statistical methods:

- Two classical methods that use the statistical approach (previously determined propensities):
  - Chou-Fasman
  - Garnier-Osguthorpe-Robson

## Chou-Fasman method

Has been one of the most popular methods. Based on calculation of the propensity of each residue to form  $\alpha$ -helix or  $\beta$ -strand.

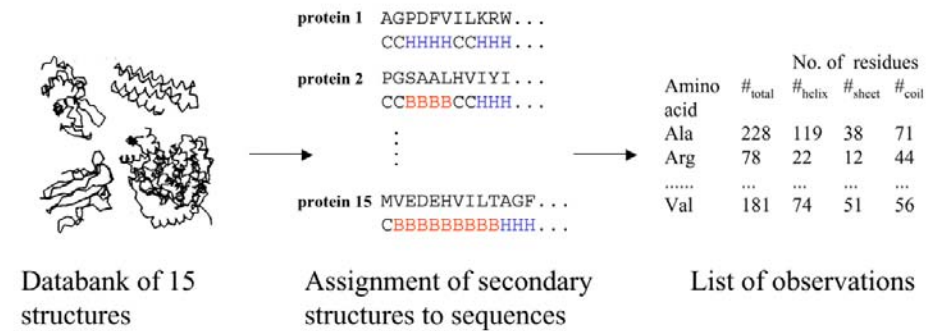
Uses table of conformational parameters (propensities) determined primarily from measurements of secondary structure by CD spectroscopy.

Can result in ambiguity if a region has high propensities for both helix and sheet.

## The Chou & Fassman method for secondary structure prediction

Chou & Fassman (1974) *Biochemistry* **13**, 211-221, 222-245

1. Probabilities for all amino acids to be either in  $\alpha$ -helix,  $\beta$ -sheet or coil.  
20 \* 3 probabilities
2. Set knowledge-based rules to apply probability tables for prediction



## From observations to probabilities

Aminoacid	observations				frequency			probability		
	# <sub>total</sub>	# <sub>helix</sub>	# <sub>sheet</sub>	# <sub>coil</sub>	f <sub>helix</sub>	f <sub>sheet</sub>	f <sub>coil</sub>	P <sub>helix</sub>	P <sub>sheet</sub>	P <sub>coil</sub>
Ala	228	119	38	71	0.52	0.97	0.31	1.45	0.97	0.66
Arg	78	22	12	44	0.28	0.15	0.56	0.79	0.90	1.20
.....	...	...	...	...	...	...	...	...	...	...
Val	181	74	51	56	0.41	0.28	0.31	1.14	1.65	0.66
average <f>					0.36	0.17	0.47	1.0	1.0	1.0

$$\frac{119}{228} = 0.52$$

$$\frac{\#Ala_{helix}}{\#Ala_{total}} = f_{Ala_{helix}}$$

$$\frac{0.52}{0.36} = 1.45$$

$$f_{Ala_{helix}} / \langle f_{helix} \rangle = P_{Ala_{helix}}$$



## Chou-Fasman propensities (partial table)

Chou-Fasman Parameters

Residue	$P_{\alpha}$	Residue	$P_{\beta}$	Residue	$P_t$
<b>Glu</b>	<b>1.51</b>	<b>Val</b>	<b>1.70</b>	<b>Asn</b>	<b>1.56</b>
<b>Met</b>	<b>1.45</b>	<b>Ile</b>	<b>1.60</b>	<b>Gly</b>	<b>1.56</b>
<b>Ala</b>	<b>1.42</b>	<b>Tyr</b>	<b>1.47</b>	<b>Pro</b>	<b>1.52</b>
<b>Leu</b>	<b>1.21</b>	<b>Phe</b>	<b>1.38</b>	<b>Asp</b>	<b>1.46</b>
<b>Lys</b>	<b>1.16</b>	<b>Trp</b>	<b>1.37</b>	<b>Ser</b>	<b>1.43</b>
<b>Phe</b>	<b>1.13</b>	<b>Leu</b>	<b>1.30</b>	<b>Cys</b>	<b>1.19</b>
<b>Gln</b>	<b>1.11</b>	<b>Cys</b>	<b>1.19</b>	<b>Tyr</b>	<b>1.14</b>
<b>Trp</b>	<b>1.08</b>	<b>Thr</b>	<b>1.19</b>	<b>Lys</b>	<b>1.01</b>
<b>Ile</b>	<b>1.08</b>	<b>Gln</b>	<b>1.10</b>	<b>Gln</b>	<b>0.98</b>
<b>Val</b>	<b>1.06</b>	<b>Met</b>	<b>1.05</b>	<b>Thr</b>	<b>0.96</b>
<b>Asp</b>	<b>1.01</b>	<b>Arg</b>	<b>0.93</b>	<b>Trp</b>	<b>0.96</b>
<b>His</b>	<b>1.00</b>	<b>Asn</b>	<b>0.89</b>	<b>Arg</b>	<b>0.95</b>
<b>Arg</b>	<b>0.98</b>	<b>His</b>	<b>0.87</b>	<b>His</b>	<b>0.95</b>
<b>Thr</b>	<b>0.83</b>	<b>Ala</b>	<b>0.83</b>	<b>Glu</b>	<b>0.74</b>
<b>Ser</b>	<b>0.77</b>	<b>Ser</b>	<b>0.75</b>	<b>Ala</b>	<b>0.66</b>
<b>Cys</b>	<b>0.70</b>	<b>Gly</b>	<b>0.75</b>	<b>Met</b>	<b>0.60</b>
<b>Tyr</b>	<b>0.69</b>	<b>Lys</b>	<b>0.74</b>	<b>Phe</b>	<b>0.60</b>
<b>Asn</b>	<b>0.67</b>	<b>Pro</b>	<b>0.55</b>	<b>Leu</b>	<b>0.59</b>
<b>Pro</b>	<b>0.57</b>	<b>Asp</b>	<b>0.54</b>	<b>Val</b>	<b>0.50</b>
<b>Gly</b>	<b>0.57</b>	<b>Glu</b>	<b>0.37</b>	<b>Ile</b>	<b>0.47</b>

## Chou-Fasman method:

- Calculation rules are somewhat arbitrary
- Example: Method for helix
- Search for nucleating region where 4 out of 6 a.a. have  $P_{\alpha} > 1.03$
- Extend until 4 consecutive a.a. have an average  $P_{\alpha} < 1.00$
- If region is at least 6 a.a. long, has an average  $P_{\alpha} > 1.03$ , and average  $P_{\alpha} > \text{average } P_{\beta}$ , consider region to be helix

## Garnier-Osguthorpe-Robson (GOR):

Build on Chou-Fasman  $P_{ij}$  values

- Probability of an amino-acid to be in a specific structural element depends on amino acid type of residue itself and neighboring atoms
- evaluate each residue PLUS adjacent 8 N-terminal and 8 carboxyl-terminal residues
- sliding window of 17
- underpredicts  $\beta$ -strand regions
- GOR III method accuracy ~64%

## Accuracy of predictions

- Both methods are only about 55-65% accurate
- A major reason is that while they consider the local context of each sequence element, they do not consider the global context of the sequence - the type of protein
  - The same amino acids may adopt a different configuration in a cytoplasmic protein than in a membrane protein

## “Adaptive” methods

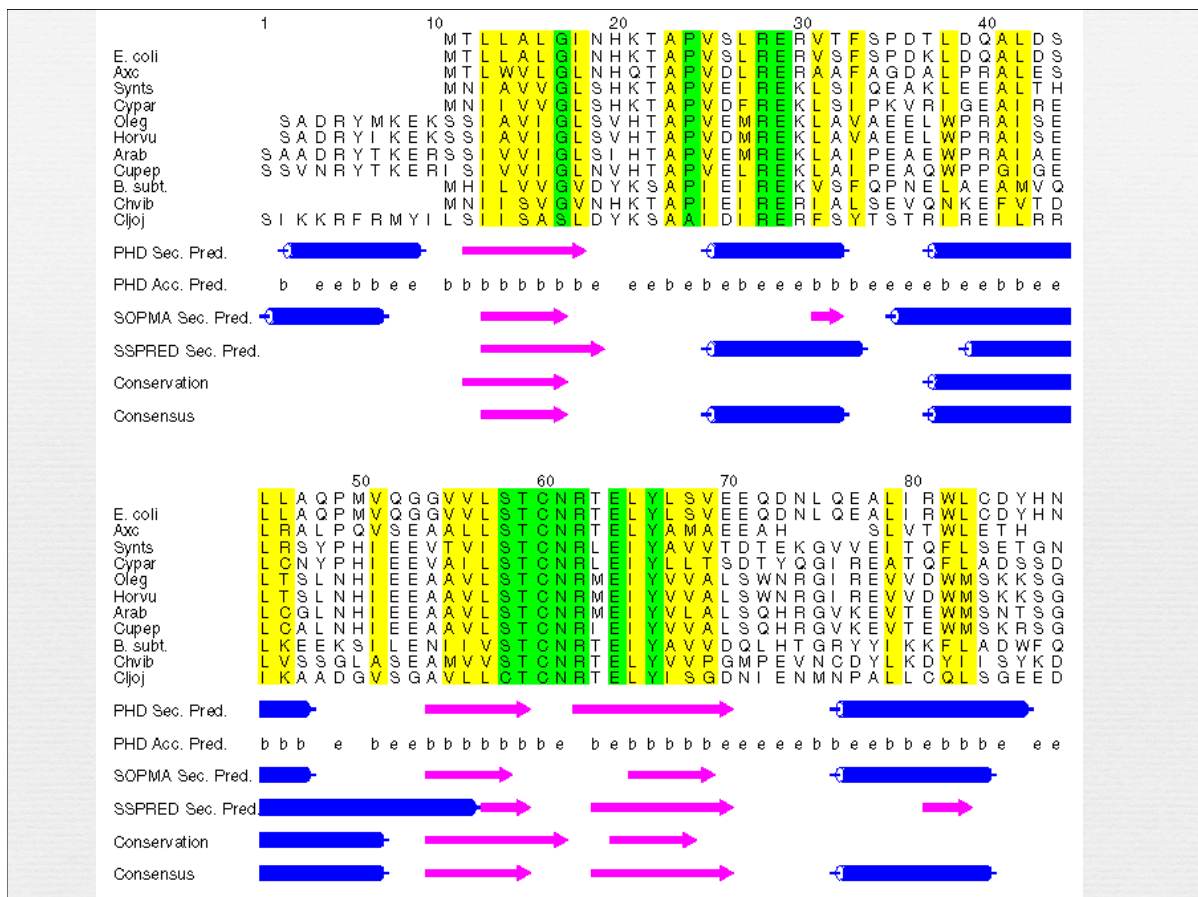
- Neural network methods - train network using sets of known proteins then use to predict for query sequence
  - **Nnpredict**
- Homology-based methods - predict structure using rules derived from proteins homologous to query sequence (multiple seq. alignment):
  - **SOPM**
  - **PHD**

## Information from multiple sequence alignment increases accuracy:

- Positions of insertions and deletions suggest regions with surface loops.
- Conserved Gly or Pro suggest a  $\beta$ -turn.
- Hydrophobic residues conserved at  $i$ ,  $i+2$ ,  $i+4$ , and separated by unconserved hydrophilic residues suggest a surface  $\beta$ -strand.
- A short run of hydrophobic a.a. (4 residues) may suggest a buried  $\beta$ -strand, a longer stretch (20 residues) - a membrane spanning helix.
- Pairs of conserved hydrophobic a.a. separated by pairs of unconserved or hydrophilic residues suggest a helix with one face packed against the protein core. Likewise an  $i$ ,  $i+3$ ,  $i+4$ ,  $i+7$  pattern of conserved hydrophobic residues.

# Guidelines:

- Since no method is the best, it is sensible to try different methods and compare the results - region where the methods agree are likely to be correctly predicted.
- Build a consensus !



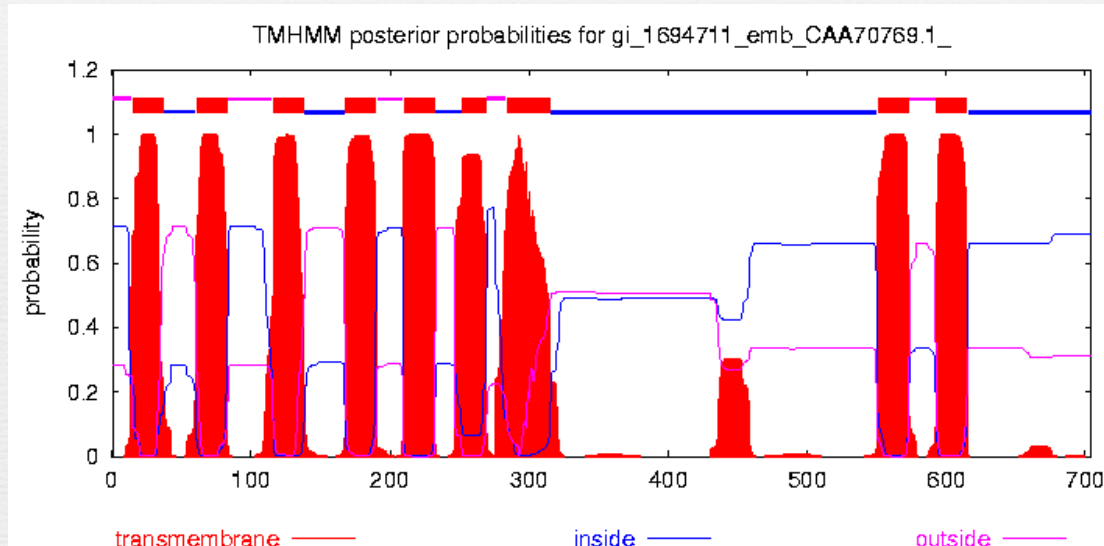
## Hydrophobicity profiles

- Tendency of a residue to occur at the surface or the interior can be described by a partition coefficient between these two phases (hydrophobicity scale).
- The profile is computed by averaging the hydrophobicity within a moving window.
- The window size depends on the size of the structural element needed to be resolved. For secondary structure elements the window must be larger than a single turn (more than 4 residues), but smaller than a large segment (a helix, less than approx. 12 residues). For a membrane spanning region the window should have the size of the expected segment (approx.. 20 residues). Small windows are noisy (too many details).

## Hydrophobicity profiles

- Can be used to predict turns, exterior and interior regions of a molecule.
- Can be applied to distinguish between membrane and soluble proteins.
- Mostly used to identify transmembrane helices in membrane proteins.

Prediction of transmembrane helices for Arabidopsis ferric reductase (FRO1) by a TMHMM (v. 2.0) at: <http://www.cbs.dtu.dk/services/TMHMM-2.0/>



## Tertiary structure prediction

- Ab initio prediction of protein 3D structures is still problematic at present (Rosetta). However, proteins often adopt similar folds despite no significant sequence or functional similarity. Nature has apparently restricted the number of protein folds.

## How to recognise a fold?

Even with no homologue of known 3D structure, it may be possible to find a suitable fold for your protein among known 3D structures using fold recognition methods.

Methods of protein fold recognition attempt to detect similarities between protein 3D structure that are not accompanied by significant sequence similarity - **find a fold that is compatible with a particular sequence, or rather than predicting how a sequence will fold, predict how well a fold will fit a sequence.**

## Fold recognition on the Internet:

- Guide to predicting protein 3D structure (**highly recommended**): <http://www.sbg.bio.ic.ac.uk/people/rob/CCP11BBS>. The guide is from 1996, theory and flowchart information is still applicable today.
- JPRED (secondary structure prediction) <http://www.compbio.dundee.ac.uk/www-jpred/>
- Phyre <http://www.sbg.bio.ic.ac.uk/~phyre/>
- TOPITS/predictprotein (<http://www.predictprotein.org>)
- UCLA-DOE Structure Prediction Server (<http://fold.doe-mbi.ucla.edu/>) Only verification.

# Threading or fold recognition

Task: detection of remote homologues behind the twilight zone

Method

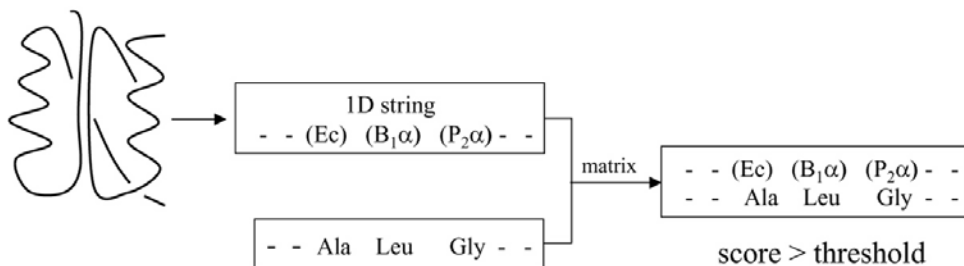
- 1) Database of folds
- 2) Dock side chains according to query sequence onto a structural scaffold from a fold library
- 3) Validate structure/sequence match by energy calculation
  - 1D-3D profiles
  - knowledge based potentials
  - atomic force fields

Applications: Threading: Fold database + sequence  
Reverse Threading: Sequence database + fold

## 1D - 3D Profiles

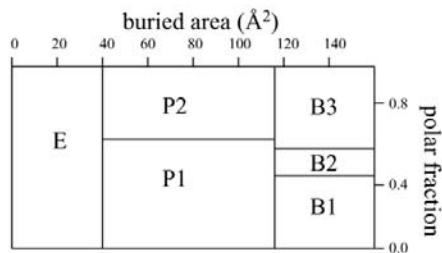
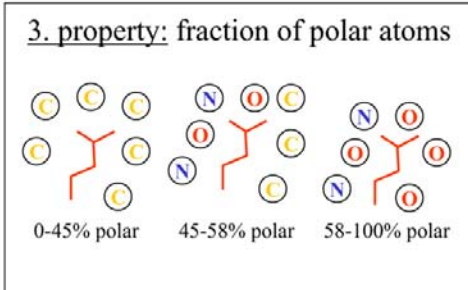
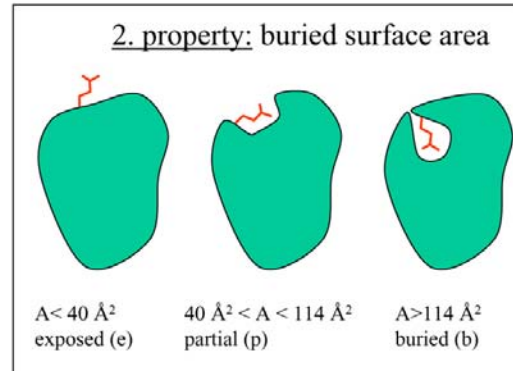
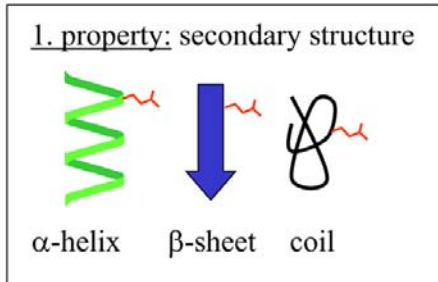
Bowie, Lüthy & Eisenberg, *Science* **253**, 164 (1991)

- compare sequence (1D) with structure (3D)
- convert 3D-structure into a 1D-string of environment classes, compare 1D-sequence with 1D-string using a scoring matrix
- if (score > threshold) then the sequence will fold into this structure



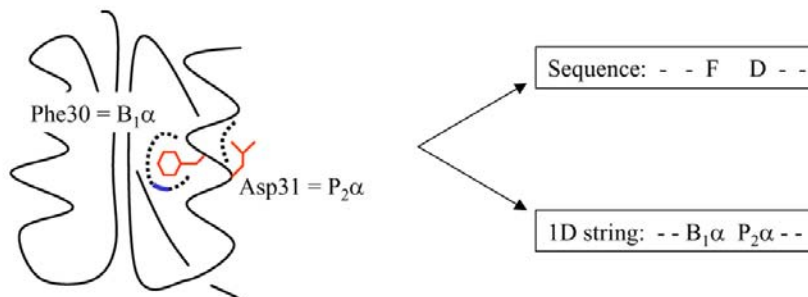


## Three properties define the class of a given residue (??)



- Combination of properties 2 & 3 gives 6 classes (E, P1, P2, B1, B2, B3)
- Combination of 6 classes with property 1 gives 18 classes ( $E\alpha$ ,  $E\beta$ ,  $E\gamma$ ,  $P1\alpha$ ,  $P1\beta$ , ...,  $B3\beta$ ,  $B3\gamma$ )

Any protein structure can be converted into a 1D-string and a 1D-sequence

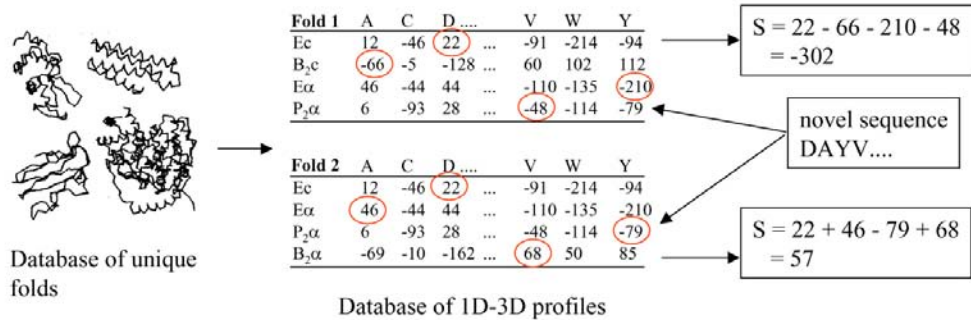




## Applications for 1D-3D profiles

### 1. Threading: What is the fold for a give sequence ?

- convert all structures into 1D-3D profiles
- compare the query sequence with all 1D-3D profiles
- calculate the Z-score for every sequence/profile pair



### 2. Inverse threading: Which sequence matches a given fold?

- convert the query fold into a 1D-3D profile
- compare the query profile against all protein sequences

## Inverse threading example:

### Which sequences fit the sperm whale myoglobin fold?

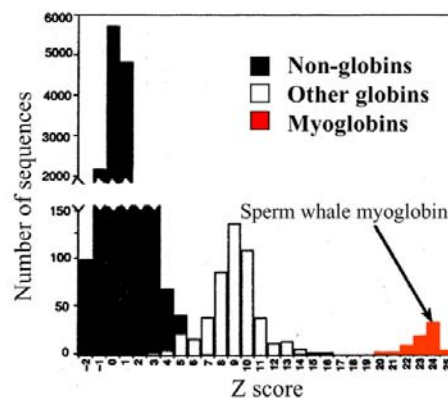
Input:

sperm whale myoglobin 1D-3D profile  
protein sequence database (n: #entries)

$$\langle S \rangle = \frac{\sum_{i=1}^n S_i}{n} \quad \text{average S-score}$$

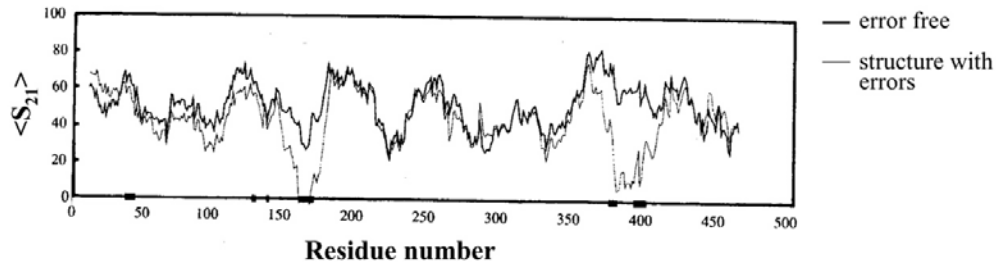
$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (\langle S \rangle - S_i)^2} \quad \text{standard deviation}$$

$$Z_i = (S_i - \langle S \rangle) / \sigma \quad \text{Z-score}$$



## Example 2: Verification of experimentally determined protein structures

$\langle S_{21} \rangle$  Average S-score for a window of 21 residues



$\langle S_{21} \rangle$  high

residues are in their preferred environment

$\langle S_{21} \rangle$  low

residues are in a non-natural environment  
(Glu in hydrophobic core)

## 2D to 3D: Comparative (homology) modelling

- A prediction of 3D structure is most successful when a structures of one or more homologues are known.
- Homologous proteins always contain a core region where the general fold of the chain is very similar.

However:

Even in core regions side-chain conformations may vary.

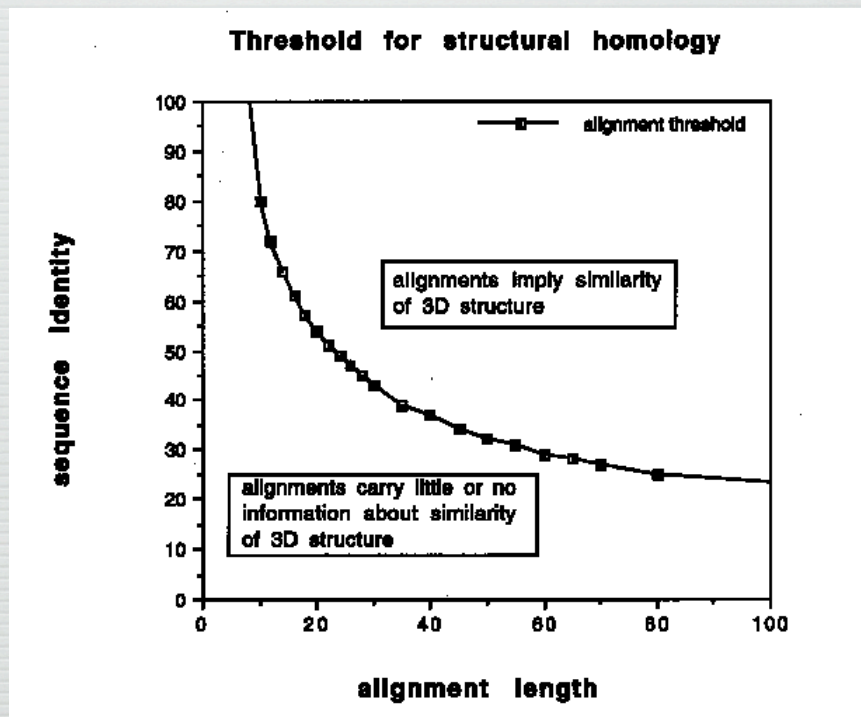
## Modelling:

The modelling process can be subdivided into 9 stages:

1. template recognition;
2. alignment;
3. alignment correction;
4. backbone generation;
5. generation of canonical loops (data based);
6. side chain generation plus optimisation;
7. *ab initio* loop building (energy based);
8. overall model optimisation (energy minimisation);
9. model verification with optional repeat of previous steps.

### What can be modelled ?

Homology threshold for structurally reliable alignments as a function of alignment length.



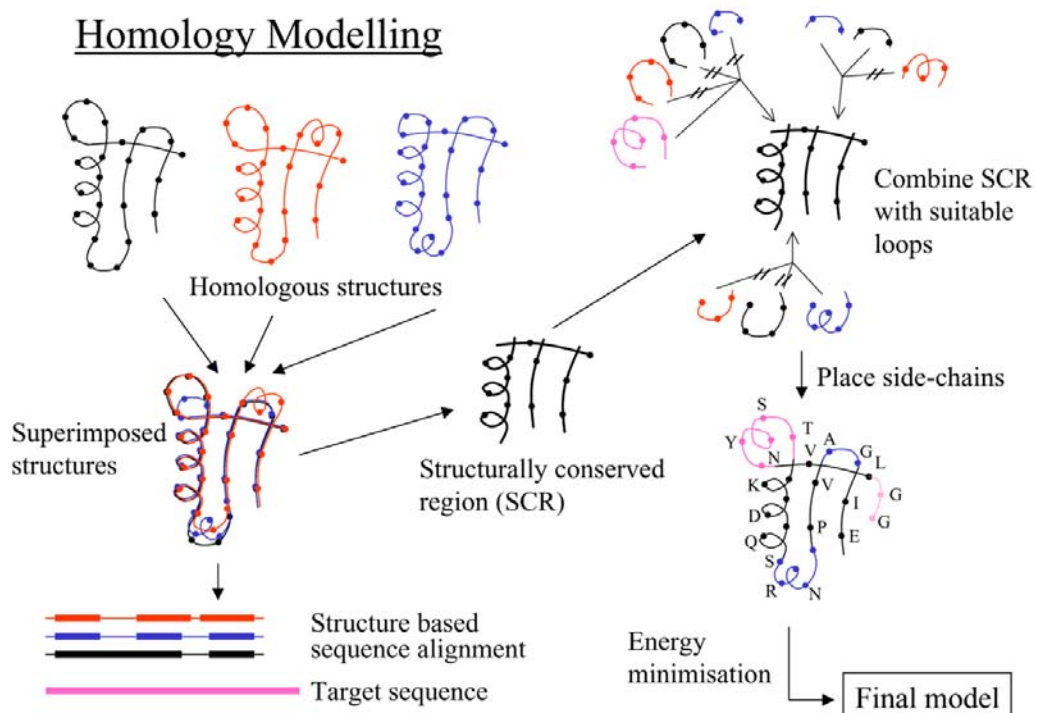
## Some general rules:

- Sequence alignments, particularly those involving proteins having low percent sequence identities can be inaccurate. Thus, a model built using the alignment will be wrong in some places. Look over the alignment carefully before building a model.
- The quality of protein models built using homology to a template protein structure is normally determined by the RMS errors in models of proteins of which the structure is known.
- Visual checking of the model is important: **check the Ramachandran plot and the energy of your model in SwissPdbviewer**, hydrophobic residues should be buried, polar and charged exposed, charged residues avoid having hydrophobic neighbours (Asp-Leu), might help to build a model from another homologous sequence and compare the results, check against secondary structure prediction,

## More rules:

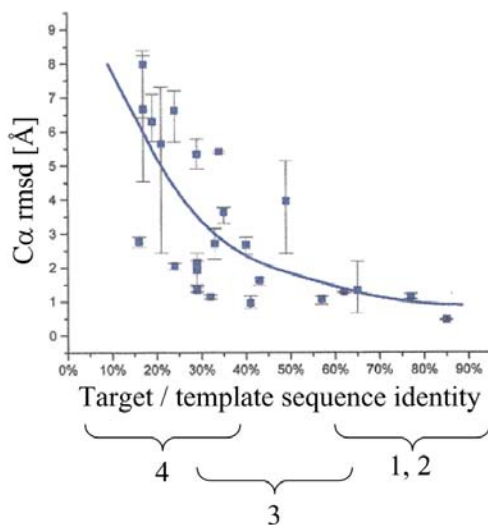
- The observed residue burial or exposure should be compared to residue burial or exposure in the model.
- The conservation of residue properties in experimental structure and model.
- Whether or not the side chains on the core beta-strands pointed in towards the barrel or out towards the helices
- The hydrogen bonding pattern of the beta-strands and helices should be checked.

## Homology Modelling



## Accuracy of homology modelling

Critical assessment of structure prediction (CASP) *Proteins* Supp. 3, (1999)



### Errors in homology modelling

1. side-chain packing
2. main-chain distortions of fragments (loops, helices, etc.)
3. main-chain distortions of newly modelled fragments
4. miss-alignments

## Some results:

- 63% of sequences sharing 40-49% identity with template yield a model deviating by less than 3 Å from the control structure.
- The number increases to 79% for seq. Identities ranging from 50 to 59%.
- Below 30 % the accuracy rapidly degrades.
- The most reliable part of the model is the portion it shares with the template, while loop and other non-

## Some links

- List with several different software packages:

[http://en.wikipedia.org/wiki/  
Protein\\_structure\\_prediction\\_software](http://en.wikipedia.org/wiki/Protein_structure_prediction_software)

- 3D JIGSAW

- <http://www.bmm.icnet.uk/servers/3djigsaw/>

- SwissModel

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

CASP does checks on reliability of the different software packages.



# Model Validation

- Does it look like a protein
  - Is chemistry ok (eg. Hydrophobic in core)
  - Geometry ok (eg Phi-Psi angles, bond lengths)
  - Amino acid environment A “correct” model can be completely wrong
  - Atom packing correct
- Accuracy (if we know the answer)
  - Rmsd
  - Fraction of correct no of modelled residues
- Use validation programs

## Evaluation of model accuracy

- Many programs come from validation programs for experimentally determined structures
- Check for proper protein stereochemistry
  - Procheck (<http://biotech.embl-ebi.ac.uk:8400/>)
    - Ramachandran plot, bond length etc
  - Whatif/Whatcheck (<http://swift.cmbi.kun.nl/swift/whatcheck/>)
    - Packing quality
  - MolProbity (<http://molprobity.biochem.duke.edu/>)
- Check sequence vs structure
  - Verify3D ([http://www.doe-mbi.ucla.edu/Services/Verify\\_3D/](http://www.doe-mbi.ucla.edu/Services/Verify_3D/))