

Computergestützte Strukturbiologie  
(Strukturelle Bioinformatik)

**Fold recognition**

Sommersemester 2009

Peter Güntert

**Sequence identity → Structural similarity**

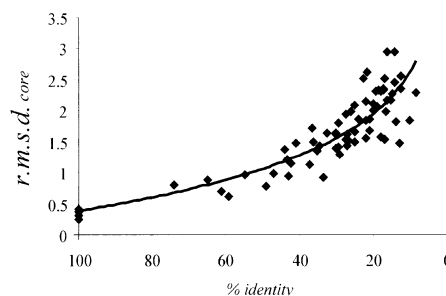


Figure 1.25 Relationships between sequence identity and structural similarity. The plot was obtained by using a larger set of proteins than in Figure 1.23, but the trend is essentially the same.

**Structural similarity ↗ Sequence identity**

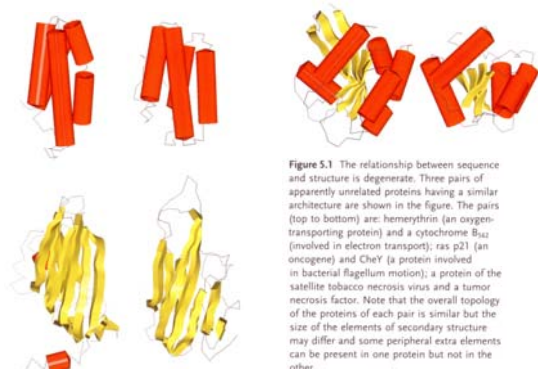


Figure 5.1 The relationship between sequence and structure is degenerate. Three pairs of apparently unrelated proteins having a similar architecture are shown in the figure. The pairs (top to bottom) are: hemerythrin (an oxygen-transporting protein) and a cytochrome B<sub>422</sub> (involved in electron transport); ras p21 (an oncogene) and CheY (a protein involved in bacterial flagellum motion); a protein of the satellite tobacco necrosis virus and a tumor necrosis factor. Note that the overall topology of the proteins of each pair is similar but the size of the elements of secondary structure may differ and some peripheral extra elements can be present in one protein but not in the other.

**Non-uniform distribution of folds**

- Few (~10) folds are shared by a large number (~30%) of known proteins
  - Large diversity in sequences and functions among members of these “superfolds”
- Examples:**
- Immunoglobulin fold
  - Rossmann fold
  - TIM barrel fold
  - Globin fold

**Methods for protein structure prediction**

Methods are distinguished according to the relationship between the target protein(s) and proteins of known structure:

- **Comparative modeling:** A clear evolutionary relationship between the target and a protein of known structure can be easily detected from the sequence.
- **Fold recognition:** The structure of the target turns out to be related to that of a protein of known structure although the relationship is difficult, or impossible, to detect from the sequences.
- **New fold prediction:** Neither the sequence nor the structure of the target protein are similar to that of a known protein.

**Scheme of protein structure prediction**

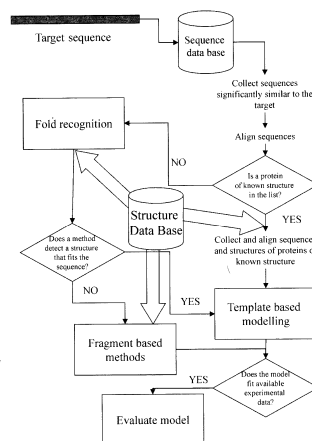


Figure 4.1 A guide to protein-structure prediction. The first step is always a search in the protein sequence database. Comparative modeling should be used when a protein of known structure sharing sequence similarity with the protein under examination is present in the database. If this is not so, fold-recognition methods should be applied and, should they fail, the user should resort to new fold or fragment-based methods. Note the central role played by the structure database in all these heuristic methods.

### Inverse protein folding problem

Which amino acid sequences fold into a known three-dimensional structure?

### Protein folding problem

Which three-dimensional structure is adopted by a given amino acid sequence?

### Fold recognition methods

- 3D profile methods
- Threading

### Profile-based fold recognition

- Physico-chemical properties of the amino acids of the target protein must “fit” with the environment in which they are placed in the modeled structure.

### Profile method for fold recognition

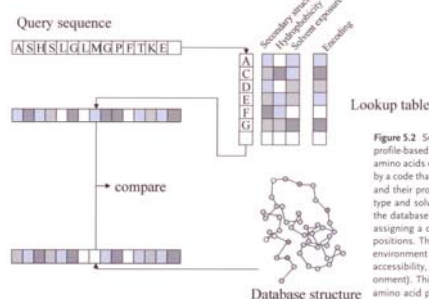


Figure 5.2 Schematic diagram of a possible profile-based method for fold recognition. The amino acids of the query sequence are replaced by a code that summarizes their hydrophobicity and their propensity for secondary structure type and solvent exposure. Each structure in the database is also encoded as a string by assigning a code to each of its amino acid positions. The code reflects their structural environment (secondary structure, solvent accessibility, and hydrophobicity of their environment). This does not depend on the actual amino acid present in the position analyzed. The string encoding the query sequence and each of the strings encoding the database structures are aligned and compared.

### Sequence-structure alignment

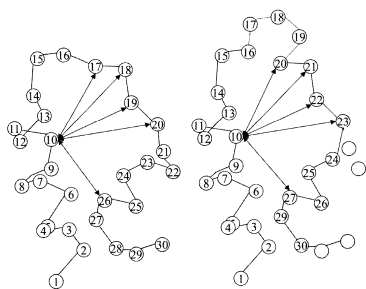


Figure 5.3 A query sequence can be positioned in a database structure in several ways, because there can be inserted and deleted residues, as shown in the right side of the figure. The interactions made by one amino acid, for example the one indicated with “10” in the figure, depend on the alignment of the rest of the sequence – the interactions of this amino acid (some of which are shown as arrows) are different in the two examples, reflecting two different alignments.

### Frozen approximation

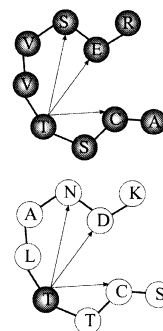


Figure 5.4 Schematic explanation of the frozen approximation. On the left, a database structure is shown with its original sequence (indicated by dark circles). In the right, the query sequence is positioned in the database structure in one of the many possible alignments. Calculation of the score should take into account which residues of the target sequence are in contact with, say, the threonine in the final alignment. In the frozen approximation, the interactions are computed by leaving the original sequence in every position of the database structure, except for the position occupied by the threonine. The procedure is repeated by substituting, in turn, each amino acid of the query sequence into a position of the target structure.

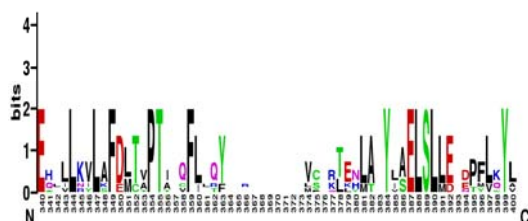
### Position Specific Iterated BLAST: PSI-BLAST

The purpose of PSI-BLAST is to look deeper into the database for matches to your query protein sequence by employing a scoring matrix that is customized to your query.

### PSI-BLAST is performed in five steps

1. Select a query sequence and search it against a protein sequence database: regular BLAST
2. PSI-BLAST constructs a multiple sequence alignment, then creates a "profile" or specialized position-specific scoring matrix (PSSM).
3. The PSSM is used as a query against the database.
4. PSI-BLAST estimates statistical significance ( $E$  values)
5. Repeat steps 2-4 iteratively, typically 5 times.  
At each new search, a new profile is used as the query.

### Position-specific scoring matrix (PSSM)



### PSI-BLAST: self-positives

- PSI-BLAST is useful to detect weak but biologically meaningful relationships between proteins.
- The main source of false positives is the erroneous amplification of sequences not related to the query. For instance, a query with a coiled-coil motif may detect thousands of other proteins with this motif that are not homologous.
- Once even a single non-related protein is included in a PSI-BLAST search above threshold, it will not go away.
- One way to check results: take newly found sequences and perform PSI-BLAST using them, then examine whether we 'fish' original sequence (reciprocal identification)

### A Method to Identify Protein Sequences That Fold into a Known Three-Dimensional Structure

JAMES U. BOWIE, ROLAND LÜTHY, DAVID EISENBERG

*Science* 253, 164-170 (1991)

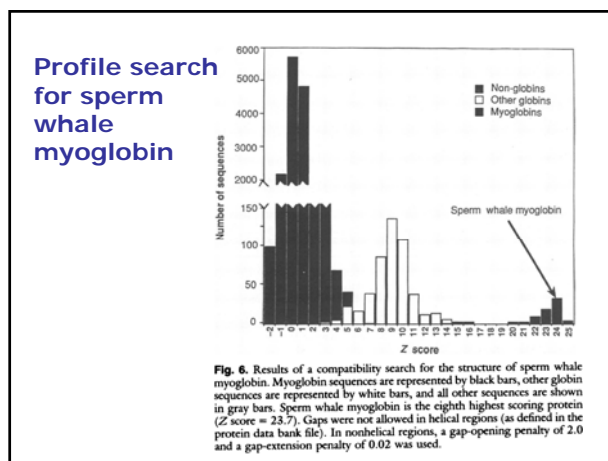
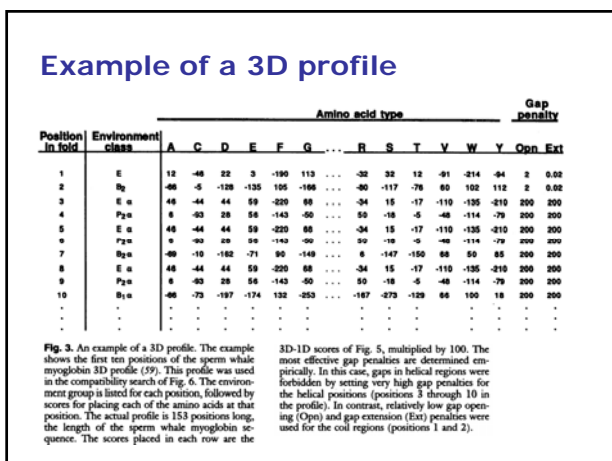
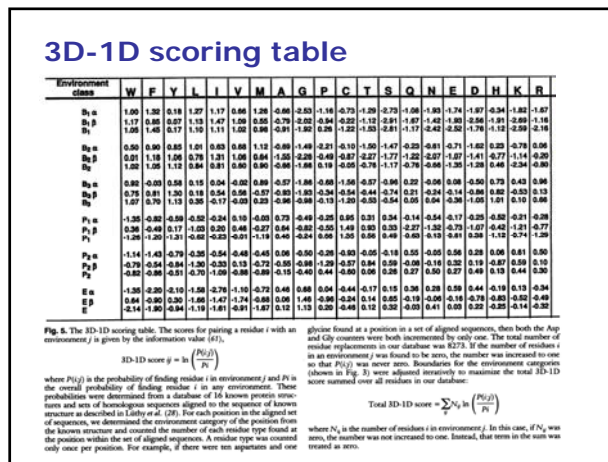
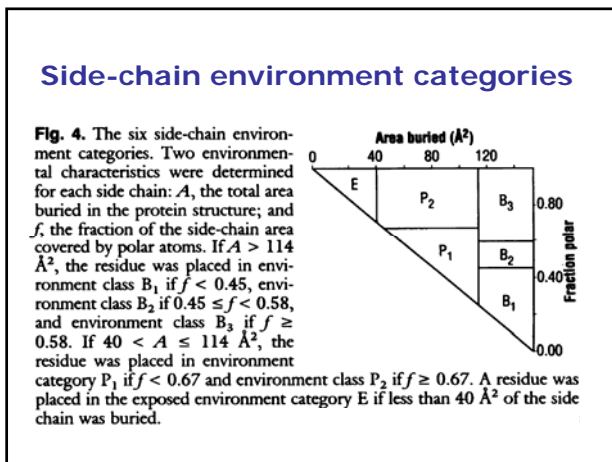
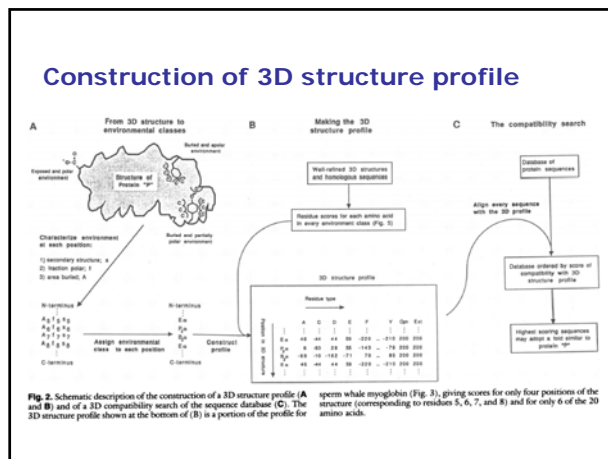
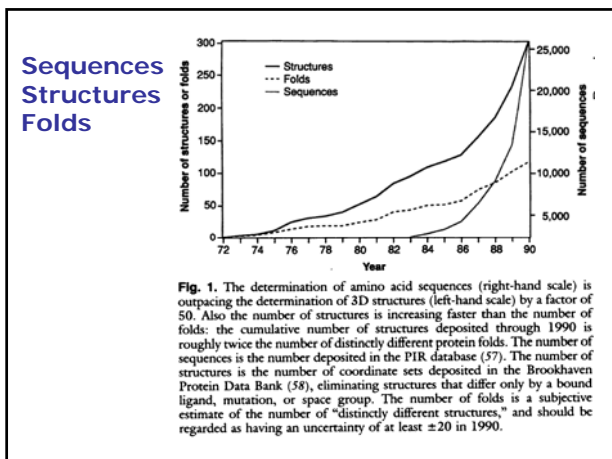
### 3D profile method

Find sequences that are most compatible with the environments of residues in the 3D structure.

The environments are described by:

1. The area of the residue buried in the protein and inaccessible to solvent
2. The fraction of side-chain area that is covered by polar atoms (O and N)
3. The local secondary structure

Bowie, Lüthy & Eisenberg. *Science* 253, 164-170 (1991)

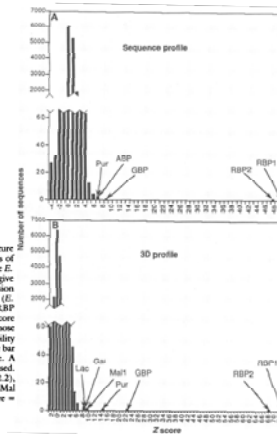


## Comparison of sequence homology and 3D profile search

**Table 1.** A comparison of a sequence homology search and a compatibility search with CRP. All proteins with Z scores greater than 6.0 in either the sequence homology search or the compatibility search are listed. Z score (1D) refers to the scores obtained from a sequence homology search with a sequence profile constructed with the Escherichia coli CRP sequence. Z score (3D) refers to the scores obtained from a structure compatibility search with a 3D profile constructed from the E. coli CRP structure (4F). Percent identity refers to the percentage of identical amino acids in the sequences aligned with the program BESTFIT (56). For the sequence homology search, a gap-opening penalty of 4.5 and a gap-extension penalty of 0.05 was used. For the structure compatibility search, a gap-opening penalty of 5.0 and a gap-extension penalty of 0.05 was used. In the sequence homology search, the next highest scoring protein after fur, Bam HI-DRF4 protein from *Foviplex virus*, had an insignificant Z score of 4.90.

Protein	Z score (3D)	Z score (1D)	Percent identity
cAMP receptor protein— <i>E. coli</i> (CRP)	46.53	73.99	100.0
cAMP receptor protein— <i>Salmonella typhimurium</i> (CRP)	44.13	72.45	99.5
Hypothetical 24.14D protein— <i>Lactobacillus casei</i>	11.84	12.74	25.6
Regulatory protein fixK— <i>Rhizobium meliloti</i>	10.65	9.26	21.1
Regulatory protein fixK— <i>E. coli</i>	9.20	7.03	21.2
Protein kinase, cAMP-dependent—bovine	8.24	—	22.0
Protein kinase type III regulatory chain—fruit fly	6.62	—	20.9
DNA polymerase accessory protein 44— <i>Bacteriophage T4</i>	6.58	—	19.7
Protein kinase type II regulatory chain—fruit fly	6.47	—	20.9
Protein kinase, cAMP-dependent, regulatory chain II- $\alpha$ —human	6.33	—	21.2
Protein kinase type I regulatory chain—fruit fly	6.15	—	20.9
Protein kinase, cAMP dependent, type II regulatory chain—bovine	6.06	—	20.9

## Comparison of a sequence homology search and a 3D profile search with ribose binding protein (RBP)



**Fig. 7.** Comparison of a sequence homology search and a structure compatibility search with ribose binding protein (RBP). (A) The results of a sequence homology search with a sequence profile constructed from the E. coli RBP sequence. The bar graph shows the number of sequences that give a particular Z score. A gap-opening penalty of 4.5 and a gap-extension penalty of 0.05 were used. The highest scoring proteins in (A) are RBP1 (*E. coli* RBP precursor, Z score = 49.0), RBP2 (*Salmonella typhimurium* RBP precursor, Z score = 47.9), GBP (*E. coli* galactose binding protein, Z score = 8.0), Pur (*E. coli* pur repressor, Z score = 6.1), and ABP (*E. coli* arabinose binding protein, Z score = 6.0). (B) The results of a structure compatibility search with a 3D profile constructed from the E. coli RBP structure. The bar graph shows the number of sequences that give a particular Z score. A gap-opening penalty of 5.0 and a gap-extension penalty of 0.2 were used. The highest scoring proteins labeled in (B) are RBP1 (Z score = 72.3), RBP2 (Z score = 68.9), GBP (Z score = 22.2), Pur (Z score = 14.2), Mal (*E. coli* Mal I protein, Z score = 9.0), Gal (*E. coli* gal repressor, Z score = 8.5), and Lac (*Klebsiella pneumoniae* lac repressor, Z score = 8.1).

## Sequence compatibility search with a 3D structure profile for actin

Protein	Z score
69 of 71 Actin Sequences	88.11
Kinase-related transforming protein (fgr)—feline sarcoma virus	17.47
Actin SC—fruit fly	9.29
68-KD Heat shock protein—mouse	8.12
70-KD Heat shock protein—frog	7.96
70-KD Major heat shock—fruit fly	7.03
70-KD Heat shock cognate protein-bovine	6.99
HNRNP complex, protein C—frog	6.74
70-KD Heat shock cognate protein—human	6.31

**Fig. 8.** Sequence compatibility search with a 3D structure profile for actin (47). All sequences that received a Z score of 6.0 or greater are listed. A gap-opening penalty of 5.0 and a gap-extension penalty of 0.2 were used. The fgr protein is the result of a gene fusion between actin and a tyrosine-specific protein kinase (63). The bovine HSC70 protein, known to have a similar structure to actin, received a Z score of 6.99 and is shown in bold type.

## Threading

- Sequences are fitted directly onto the backbone coordinates of known protein structures.
- Matching of sequences to backbone coordinates is performed in 3D space, incorporating specific pair interactions explicitly.

## A new approach to protein fold recognition

D. T. Jones\*†, W. R. Taylor† & J. M. Thornton\*

\* Biomolecular Structure and Modelling Unit, Department of Biochemistry and Molecular Biology, University College, Gower Street, London WC1E 6BT, UK  
 † Laboratory of Mathematical Biology, National Institute for Medical Research, The Ridgeway, Mill Hill, London, NW7 1AA, UK

The prediction of protein tertiary structure from sequence using molecular energy calculations has not yet been successful; an alternative strategy of recognizing known motifs<sup>1</sup> or folds<sup>2-4</sup> in sequences looks more promising. We present here a new approach to fold recognition, whereby sequences are fitted directly onto the backbone coordinates of known protein structures. Our method for protein fold recognition involves automatic modelling of protein structures using a given sequence, and is based on the frameworks of known protein folds. The plausibility of each model, and hence the degree of compatibility between the sequence and the proposed structure, is evaluated by means of a set of empirical potentials derived from proteins of known structure. The novel aspect of our approach is that the matching of sequences to backbone coordinates is performed in full three-dimensional space, incorporating specific pair interactions explicitly.

*Nature* 358, 86-89 (1992)

## Threading

- A library of different protein folds is derived from the database of protein structures.
- Each fold is considered as a chain tracing through space; the original sequence being ignored completely.
- The test sequence is then optimally fitted to each library fold, allowing for relative insertions and deletions in loop regions.
- The 'energy' of each possible fit (or threading) is calculated by summing the proposed pairwise interactions and the solvation energy.
- The library of folds is then ranked in ascending order of total energy, with the lowest energy fold being taken as the most probable match.

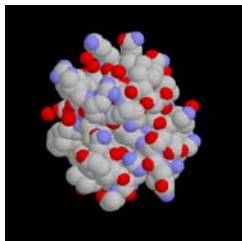
## Knowledge-based (pair) potentials

$$E(r) = -kT \ln[f(r)]$$

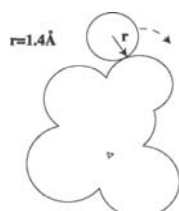
- $r$  distance between two atoms (or some other parameter, like dihedral angles or solvent accessible surface)
- $E(r)$  is the energy at  $r$
- $f(r)$  is the probability density at  $r$
- $k$  is the Boltzmann constant
- $T$  is the absolute temperature

## Solvent accessible surface

Represent atoms as spheres with appropriate radii and eliminate overlapping parts.



Mathematically roll a sphere all around that surface.



The sphere's center traces out a surface as it rolls.

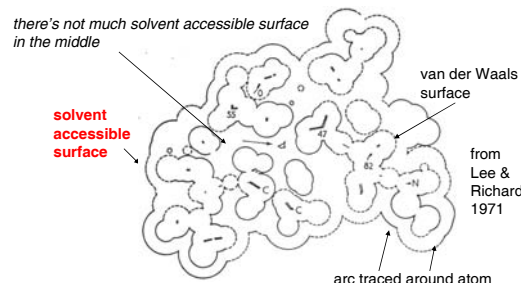
Lee & Richards, 1971  
Shrake & Rupley, 1973

## Cross-section (slice) of a protein structure:

Inner surfaces here are van der Waals. Outer surface is that traced out by the center of the sphere as it rolls around the van der Waals' surface. If any part of the arc around a given atom is traced out, that atom is accessible to solvent. The solvent accessible surface of the atom is defined as the sum of the arcs traced around an atom.

there's not much solvent accessible surface in the middle

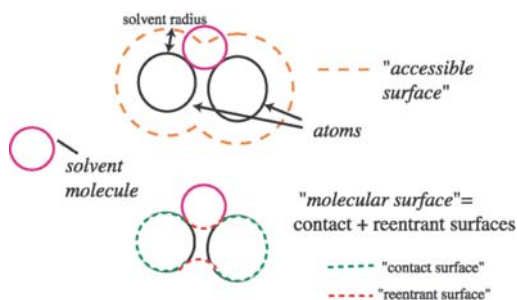
solvent accessible surface



van der Waals surface  
from Lee & Richards, 1971

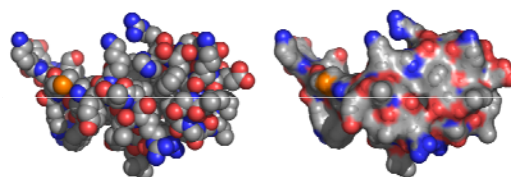
arc traced around atom

## Accessible surface/Molecular surface



These are *alternative* ways of representing the surface which is essentially in contact with solvent.

## Molecular surface of proteins



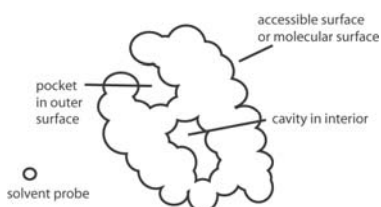
depiction of heavy atoms (O, N, C, S) in a protein as van der Waals spheres

depiction of the corresponding "molecular surface"--volume contained by this surface is vdW volume plus "interstitial volume"--spaces in between

## The irregular surface of proteins: pockets and cavities

• A **pocket** is an empty concavity on a protein surface which is accessible to solvent from the outside.

• A **cavity** or **void** in a protein is a pocket which has no opening to the outside. It is an interior empty space inside the protein.

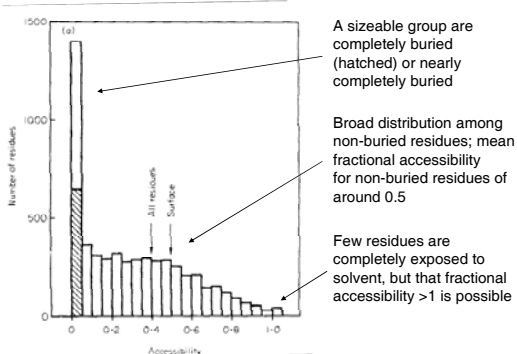


Pockets and cavities can be critical features of proteins in terms of their binding behavior, and identifying them is usually a first step in structure-based ligand design etc.

## Fractional accessibility

- Calculate total solvent accessible surface of protein structure (also can calculate solvent accessible surface for individual residues/sidechains within the protein).
- Model the accessible surface area in a *disordered* or *unfolded* protein using accessible surface area calculations on *model tripeptides* such as Ala-X-Ala or Gly-X-Gly.
- From these we can calculate *what fraction* of the surface is buried (inaccessible to solvent) by virtue of being within the folded, native structure of the protein.
- *Fractional accessibility* is computed by dividing the accessible surface area in the native protein structure by the accessible surface in the modelled unfolded protein. The *residue fractional accessibility* and *side-chain fractional accessibility* are calculated for individual residues or side-chains in the structure.

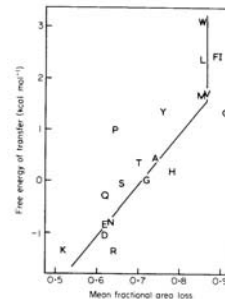
### Distribution of residue fractional accessibilities



### Residue fractional accessibility correlates with free energies of transfer for amino acids between water and organic solvents

The interior of a protein is akin to a nonpolar solvent in which the nonpolar sidechains are buried. Polar sidechains, on the other hand, are usually on the surface.

However, some polar side chains do get buried, and it must also be remembered that the backbone for every residue is polar, including those with nonpolar side chains. So a lot of polar moieties do get buried in proteins.



Miller, Janin, Lesk & Chothia (1987)  
Fauchere & Pliska (1983)

### Solvation potential

Similarly, the solvation potential for an amino-acid residue *a* is defined as

$$\Delta E_{solv}^a(r) = -RT \ln \left[ \frac{f^a(r)}{f(r)} \right]$$

where *r* is the % residue accessibility (relative to residue accessibility in GGXGG fully extended pentapeptide), *f<sup>a</sup>(r)* is the frequency of occurrence of residue *a* with accessibility *r*, and *f(r)* is the frequency of occurrence of all residues with accessibility *r*.

### Pairwise pseudo-energy terms

For specified atoms (*Cβ* → *Cβ* for example) in a pair of residues *ab*, topological level (sequence separation) *k* and distance interval *s*, the potential is given by the following expression

$$\Delta E_k^{ab} = RT \ln [1 + m_{ab}\sigma] - RT \ln \left[ 1 + m_{ab}\sigma \frac{f_k^{ab}(s)}{f_k(s)} \right]$$

where *m<sub>ab</sub>* is the number of pairs *ab* observed at topological level *k*, *σ* is the weight given to each observation, *f<sub>k</sub>(s)* is the frequency of occurrence of all residue pairs at topological level *k* and separation distance *s*, and *f<sub>k</sub><sup>ab</sup>(s)* is the equivalent frequency of occurrence of residue pair *ab*. *RT* is taken as 0.582 kcal mol<sup>-1</sup>. Short- (sequence separation, *k* ≤ 10), medium- (11 ≤ *k* ≤ 30) and long- (*k* > 30) range potentials have been calculated between the following atom pairs: *Cβ* → *Cβ*, *Cβ* → *N*, *Cβ* → *O*, *N* → *Cβ*, *N* → *O*, *O* → *Cβ* and *O* → *N*.

### Statistically derived potentials

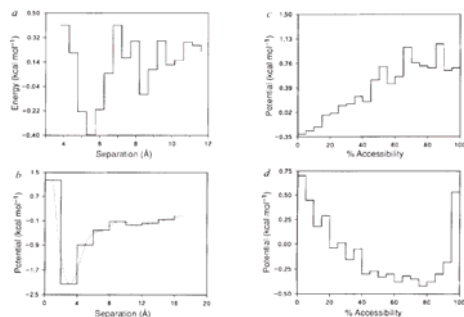
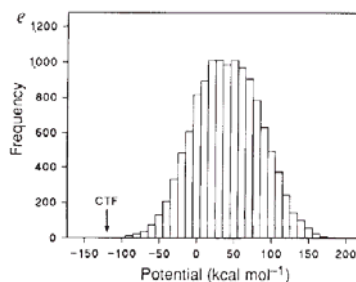


FIG. 3 Samples of the statistically derived potentials are shown. a Short-range (*k* = 3) Ala-Ala Cβ → Cβ interaction. Low-energy states are observed for distances around 6 Å, corresponding mainly to α-structure, and 9 Å, corresponding mainly to β-structure. b Long-range (*k* > 30) Cys-Cys Cβ → Cβ interaction. The most significant energy minimum around 4 Å corresponds to disulphide bridge formation. c Solvation potential for leucine, and d solvation potential for glutamic acid.

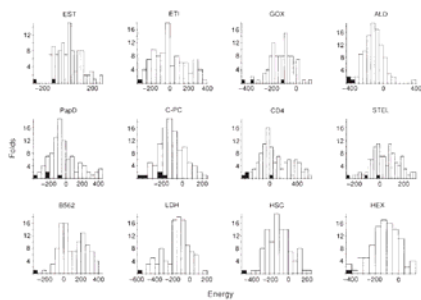
### Threading histogram for 1CTF



Threading histogram for the C-terminal ribosomal protein fragment, 1CTF. All possible threadings of the CTF sequence on the CTF structure were computed (secondary structure gaps disallowed) and the energies of each threading calculated. The native threading is indicated, and was found to be the lowest energy threading.

### Threading results

FIG. 3 For a number of test cases (see Table 3) the histogram of energies for randomly threading onto each of the 102 folds is given. In each histogram, the positions of folds expected to match the given sequence (that is, those folds similar to the known fold of the test sequence) are shown as filled bars. For example, in the case of LDH (lactate dehydrogenase), the expected match is the 10th fold (lactate dehydrogenase). This match is shown as a single filled bar representing an energy of  $-577$  kcal mol<sup>-1</sup>, an energy which is lower than that achieved by any other fold. As noted in the text, in some cases expected folds are apparently not detected. This occurs for two reasons: either the expected structures are not sufficiently similar to the native fold or the optimisation method does not succeed in producing a satisfactory alignment.



### Summary of trial fold-recognition searches

TABLE 1 Summary of trial fold-recognition searches

Test protein	Source	Fold	Best match	JC	% Sequence identity	Matches
C-phycoerythrin $\beta$ (C-PC)	Red algae	Oblon	1MBR	101	7	1, 2, 9, 16, 25
GroEL-like protein (GOL)	Sonch	TM barrel	1R5V(A)	52	10	1, 3, 49
Muscle aldolase (ALD)	Human	TM barrel	4U(A)	80	6	1, 2, 3
Lactate dehydrogenase (LDH)	Dogfish	Rossmann	4M2H(A)	87	15	1*
Elastase (EST)	Pig	Trypsin	4TTP	110	35	1, 14
CD4	Human	$\beta$	2F84(B)	87	10	1, 2, 31
Stathmin (STEL)	Yeast	Cu binding	2A24(A)	18	14	1, 6, 20
Cytochrome B562 (B562)	E. coli	4-helix bundle	2MR	78	6	1
Trypsin inhibitor DE-3 (TI)	Rat	Interleukin 1 $\beta$	1IB	14	5	1
PaO-chaperonin	E. coli	$\beta$	2F84(L)	64	15	1, 5, 9, 35
70k, Heat-shock cognate (HSC)	Cow	Actin	1ATN(A)	94	9	1
Hexokinase B (HXK)	Yeast	Actin	1ATN(A)	0	12	1

In each case the database included 102 protein chains, except where the test protein was itself in the database, in which case it was excluded. Templates for each chain were constructed as described in the text, with residues not in helices or strands (as calculated by DSSP<sup>17</sup>), assigned as loop residues. For the 70k heat-shock cognate protein and hexokinase searches, the coordinates for actin were also included (coordinates deposited under the code 1ATN, but not yet released). Proteins with >25% sequence identity to the test protein were also excluded from the calculation of potentials. The pairwise and solvation terms were summed and stored separately, and standard deviations ( $s_{d_{ij}}$  and  $s_{d_{sol}}$ ) for the two contributing factors calculated over the set of 102 folds. To balance the contributions of the pairwise and solvation terms, the final energy was taken as  $E = E_{ij} + W_{sol}$ , where  $W = (s_{d_{ij}}/s_{d_{sol}})^2$ . The confidence of the match (JC) is given in terms of the absolute energy difference between the top scoring fold and the next highest scoring, different, fold. The 'best match' column gives the Brookhaven ID of the best matching chain fold (including chain identity where appropriate), along with the sequence identity between the best matching chain and the test protein. Residues in the sorted list of threading energies of similar folds are also shown. A constant set of alignment parameters (gap penalty for example) was used for all database searches shown. Typical execution times for a single search of 102 chains are around 100 minutes on a Unix workstation (Silbourne S-602). The 102 chains used were as follows: 351C, A256B, 2AK1, 1ABP, A530N, B40P, 3APD, N4ATC, B6ATC, A2AZA, N8LA, S8P2, 2CA1, ATCA1, 1C25, 1C26, A2C7, 1C24, 2C2V, 3C2A, 2C2A, HOPR, 5C2A, 2C2P, 4C2P, 1C2R, 2C2R, 1C2E, 1C2E, 1C2F, 1C2Y, 2C2Y, 3C2P, A4C2P, A12C2F, 1C2CA, 1C2C7, HOPB4, L2P4, 1FD2, 1FK1, 3FK1, 4FK1, A3GAP, 3GAP, 1C2C, 015D1, 3GAP, A3H4, 1AP, A3P4, B3P4, 14E, 11B, 3K8, 3C0, 110, 2UR, 4E3, 1LH, 111R, A2L7, 1121, 1MBA, 14ED, 4M2H, 2MR, 2D0, A2P4, 5P4, 1P2, 1P7, A1P4, 3P4, 4P4, 1PH, 5P1, 4P1, 1R0, 2R, 2R2, 7R5, 4R0, 25G, 45G, 11K3, 29K, 0250, 25S, 25T, 1T05, 1T05, 1T05, 4T0, A1T0, 11T0, 11T0, A2W0, B2W0, A1W1, B1W1, A4A1, A1V1.

\* Other topologically similar (yet structurally different) parallel  $\beta$ -folds were positioned at 3, 7, 11, 12, 13, 17, 19, 31, 34, 82.

### Critical components of fold recognition techniques

- Techniques producing useful alignments between sequences and structures
- Criteria for identifying native-like sequence/structure combinations
- Energy functions or parameter sets providing a reasonable description of protein-solvent systems

### Fold recognition results from CASP

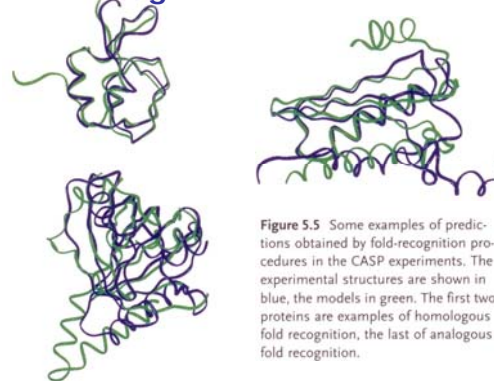


Figure 5.5 Some examples of predictions obtained by fold-recognition procedures in the CASP experiments. The experimental structures are shown in blue, the models in green. The first two proteins are examples of homologous fold recognition, the last of analogous fold recognition.

### Literatur

- Anna Tramontano: *Protein Structure Prediction, Wiley-VCH, 2006.*
- J. U. Bowie, R. Lüthy & D. Eisenberg. *Science* 253, 164-170 (1991) [3D profiles]
- D. T. Jones, W. R. Taylor & J. M. Thornton. *Nature* 358, 86-89 (1992) [Threading]
- M. J. Sippl. *Current Opinion in Structural Biology* 5, 229-235 (1995) [Database-derived potentials]