# JMB



# Protein Fold Recognition by Mapping Predicted Secondary Structures

# Robert B. Russell, Richard R. Copley and Geoffrey J. Barton\*

University of Oxford Laboratory of Molecular Biophysics, The Rex Richards Building, South Parks Road Oxford, OX1 3QU, England

A strategy is presented for protein fold recognition from secondary structure assignments ( $\alpha$ -helix and  $\beta$ -strand). The method can detect similarities between protein folds in the absence of sequence similarity. Secondary structure mapping first identifies all possible matches (maps) between a query string of secondary structures and the secondary structures of protein domains of known three-dimensional structure. The maps are then passed through a series of structural filters to remove those that do not obey simple rules of protein structure. The surviving maps are ranked by scores from the alignment of predicted and experimental accessibilities. Searches made with secondary structure assignments for a test set of 11 fold-families put the correct sequence-dissimilar fold in the first rank 8/11 times. With cross-validated predictions of secondary structure this drops to 4/11 which compares favourably with the widely used THREADER program (1/11). The structural class is correctly predicted 10/11 times by the method in contrast to 5/11 for THREADER. The new technique obtains comparable accuracy in the alignment of amino acid residues and secondary structure elements. Searches are also performed with published secondary structure predictions for the von-Willebrand factor type A domain, the proteasome 20 S  $\alpha$  subunit and the phosphotyrosine interaction domain. These searches demonstrate how the method can find the correct fold for a protein from a carefully constructed secondary structure prediction, multiple sequence alignment and distance restraints. Scans with experimentally determined secondary structures and accessibility, recognise the correct fold with high alignment accuracy (86% on secondary structures). This suggests that the accuracy of mapping will improve alongside any improvements in the prediction of secondary structure or accessibility. Application to NMR structure determination is also discussed.

© 1996 Academic Press Limited

*Keywords:* structure prediction; fold recognition; threading; nuclear magnetic resonance; secondary structure mapping

\*Corresponding author

# Introduction

The flood of new protein sequences demands techniques to infer protein 3D structure from sequence alone. For  $\approx 30\%$  of protein sequences,

conventional alignment techniques (e.g. Lipman & Pearson, 1985; Altschul *et al.*, 1990; Smith & Waterman, 1981) or profile and pattern methods (e.g. Gribskov *et al.*, 1987; Barton & Sternberg, 1990) find similarities to a protein of known 3D structure (Chothia, 1992). The remaining 70% of protein sequences may adopt previously unseen protein folds. Alternatively, they may have topologies (folds) similar to known protein structures but share no detectable sequence similarity (e.g. Russell & Barton, 1994). Such fold similarities will normally not be found until both protein 3D structures have been determined experimentally (Orengo, 1994; Holm & Sander, 1994a). In an attempt to find fold similarities of this type in advance of 3D structure

Present address: R. B. Russell, Biomolecular Modelling Laboratory, Imperial Cancer Research Fund Laboratories, 44 Lincoln's Inn Fields, P.O. Box 123, London, WC2A 3PX, England.

Abbreviations used: 3D, three-dimensional; NMR, nuclear magnetic resonance; Ig, Immunoglobulin; SDM, site-directed mutagenesis; WWW, world wide web; the standard one- and three-letter abbreviations for the amino acids are also used throughout.

determination, several fold recognition techniques have been developed (see Bowie & Eisenberg, 1993; Wodak & Rooman, 1993; Jones & Thornton, 1993 and references therein). These techniques may locate some fold similarities that are undetectable by the comparison of sequence. However, the methods are often computationally intensive and many similarities still go undetected (Pickett *et al.*, 1992; Lemer *et al.*, 1996).

In parallel with the development of fold detection methods, the accuracy of secondary structure prediction has improved from  $\approx 65\%$  to  $\approx 72\%$  on average. Though this is only a small percentage increase, recent predictions are more useful, since the application of multiple sequence alignments improves the identification of the number, type and location of core secondary structure elements. Prediction from sequence alignments can also accurately identify the position of loops, and residues likely to be buried in the the protein core (Benner et al., 1994; Barton, 1995; Russell & Sternberg, 1995). Given a good secondary structure prediction, the next question to ask is how the secondary structures might be arranged into a tertiary fold. ab initio methods for folding secondary into tertiary structure search for possible arrangements of secondary structures that obey general packing rules (Cohen & Sternberg, 1980; Cohen et al., 1980, 1982; Smith-Brown et al., 1993; Sun et al., 1995). These methods have been applied in numerous blind predictions (Hurle et al., 1987; Cohen et al., 1986; Curtis et al., 1991; Jin et al., 1994; Huang et al., 1994) with varied results. A limitation is the number of packing combinations that must be considered. This can become unmanageable for >nine secondary structures (Cohen et al., 1982), though approaches to reduce the number of combinations have been described (Taylor, 1991; Clark et al., 1991).

The most successful predictions of protein tertiary structure in the absence of clear sequence similarity to a protein of known 3D structure, have been those where secondary structure predictions, and experimental information were combined to suggest resemblance to an already known fold. Correct folds have been predicted in this way for the  $\alpha$  subunit of tryptophan synthase (Crawford *et al.*, 1987), a family of cytokines (Bazan, 1990), and recently, for the von Willebrand factor type A domain (Edwards & Perkins, 1995), and the synaptotagmin C2 domain (Gerloff et al., 1995). Although the details of these studies differed, all used predicted secondary structures from multiple alignment, combined with the careful application of protein structural principles (often together with experimental data) to suggest a protein fold. Two automated methods for comparing predicted and experimental secondary structures have been described previously (Sheridan et al., 1985; Rost, 1995) with promising though limited preliminary results.

In this paper we show how secondary structure and accessibility prediction together with basic rules of protein structure may be used to find the correct fold within a database of protein structural domains. The method first generates all possible matches (referred to as maps) between query and database secondary structure patterns, allowing for insertions and deletions of whole secondary structure elements. Maps are filtered by a series of structural criteria to arrive at a collection of sensible template structures. The sequence of the query protein is then aligned to the template structures by matching predicted and observed patterns of residue accessibility. Finally, alignments are ranked by a score that combines accessibility matching with a penalty for differences in secondary structure length. The method is designed to cope with incorrect secondary structure assignments, insertions/deletions of whole secondary structure elements, and differences in the lengths and orientations of secondary structures.

# Theory and Algorithm

# Database of unique protein 3D structural domains

A database of protein 3D structural domains was derived from the Brookhaven Protein Databank (Bernstein et al., 1977). 930 non-identical chains were clustered by sequence comparison (Smith & Waterman, 1981; Barton, 1993) to leave 275 sequence families. One representative of each family was chosen to have the highest resolution and lowest *R*-factor. The representative structures were then split into 377 domains by eye. A sub-database of higher quality domains was created for analysis. This contained only those structures determined by X-ray crystallography, refined and of a resolution of 2.5 Å or better. Secondary structures for all domains were defined by the programs DSSP (definition of secondary structure in proteins; Kabsch & Sander, 1983) or by DEFINE (Richards & Kundrot, 1988) when only  $C^{\alpha}$  atoms were available. Axial coordinates were calculated for all secondary structures as described by Richards & Kundrot (1988). Extra axial coordinates were calculated at the N and C-terminal ends to allow for possible differences in secondary structure length. The domain database is available via the WWW (http://geoff.biop.ox.ac.uk/).

### Alignment of secondary structures

The secondary structure of the protein is represented as a sequence of H and B characters where each H represents an entire  $\alpha$  helix and each B a  $\beta$  strand. A fast method for generating all exact matching alignments between two strings that allows up to a maximum number of deletions from each string (Russell *et al.*, 1995) is used to find all maps between the query pattern of secondary structures and the domain database. The method is recursive, and reminiscent of regular expression matching. In this study up to two deletions were permitted from the query secondary structure

string, to allow for errors in the prediction. Up to five deletions were permitted from each database structure, to allow insertions or deletions of secondary structures typical of proteins having similar 3D structures in the absence of sequence similarity. Deletions from the database structure were only counted if they were contained within matched elements (overhanging deletions were ignored). Explicit mismatches were not allowed, but were treated as deletions from either the query or database structure. These values were chosen since they are typical of the expected accuracy of secondary structure prediction, and typical of insertions and deletions of secondary structure elements across members of a diverse structural family. In practice, the allowable deletions from query and database should be chosen on a case by case basis. For consistency, we kept the maximum numbers of deletions fixed during this study.

### Filters

The alignment method will find all maps between two strings of secondary structure elements, but due to the allowance for deletions, many of these will correspond to implausible topologies. Accordingly, seven filters are used to remove maps corresponding to nonsensical protein 3D structures and/or those not satisfying imposed experimental restraints.

#### Removing un-compact structures

Two filters exploit the radius of gyration,  $R_g$ , to remove non-compact maps. Analysis of the 275 high quality domains shows that  $R_g \leq 2.8L^{0.34} + 4.0$ , where *L* is the length of the structure in residues. For each map, a coarse  $R_g$  is first calculated by considering the centroids of secondary structures, and their C-terminal loops as point masses. A fine  $R_g$  is also calculated by considering all matched residues (plus C-terminal loops) as point masses. Maps are removed if either  $R_g$  value is greater than the maximum for compact domains of the same length.

# Loop length distance restraints

Analysis of the 275 high quality domains shows that the maximum distance  $D_{\text{max}}$  between axial coordinates that can be bridged by a loop of  $N_l$ residues is 11.621 ( $N_l$  + 0.25)<sup>0.359</sup> + 0.5 Å. Maps having any loop with distances larger than  $D_{\text{max}}$  + 4 Å are removed. 4 Å is added to allow for differences in the packing of database and query secondary structures, since similar structures with little sequence similarity can have shifts of up to 4 Å (Holm & Sander, 1995).

Care is taken to allow a range of possible positions for the match of query and database structures. This allows for errors in secondary structure prediction, which may fail to predict the precise start or end of correctly identified elements, and allows for the observed differences between the lengths of secondary structure elements within proteins having similar topologies despite no significant sequence similarity. For a position x on a database secondary structure, and a minimum and maximum length for a query secondary structure,  $L_{mun}$  and  $L_{max}$ , the range of allowable positions of the query residue on the database structure (of length  $L_{obs}$ ) is given by:

 $x_{\min} = \min(L_{obs} - L_{max}, 0) - h + x$  $x_{max} = \max(L_{obs} - L_{min}, 0) + h + x$ 

where *h* is a leniency parameter, allowing for differences in the length of query and database secondary structures. h = 4 allows for differences typical of those found in proteins having similar 3D structures despite no sequence similarity.

# Poor $\beta$ sheets

The deletion of  $\beta$  strands from a  $\beta$  sheet can lead to maps corresponding to nonsensical 3D structures. Maps containing isolated  $\beta$  strands (i.e. those lacking hydrogen bonding partners) are removed. Maps are also removed if  $\beta$  strands are deleted from the centre of  $\beta$  sheets contained within the map.

Analysis of high quality domains shows that the number of  $C^{\alpha}$ – $C^{\alpha}$  contacts  $\leq 6$  Å made by a  $\beta$  strand ( $C_{\beta-\alpha\alpha}$ ) with any of its neighbouring  $\beta$  strands is always  $\geq N_{\beta} - 2$ , where  $N_{\beta}$  is the number of residues in the  $\beta$  strand. Thus maps are also removed if one or more  $\beta$  strands has  $C_{\beta-\alpha\alpha} < N_{\beta} - 2$ .

#### Adjacent parallel structures

Maps are removed if tandem secondary structures in the query are made to match parallel structures in the database by the deletion of intervening secondary structures. Genuine adjacent parallel structures within the database are allowed. This filter can be turned off in instances when there are long loops connecting query secondary structure elements, as in the phosphotyrosine interaction domain example (see Results).

# Distance restraints

Distance restraints may be imposed from the results of NMR experiments, knowledge of the disulphide linkages, or knowledge of residues involved in the active or binding site of the query. In this study, distance restraints are only included in the von Willebrand factor and proteasome examples (see Results). A tolerance value t = 4 Å is added to all distance restraints as for the loop length filtering.

# Consistency and redundancy

Maps are only kept if there is at least one placement of the query onto the database secondary

Table	1. Effect of	filters when	applied inde	pendently
Initial	Inital		Remaining	Remaining
maps	folds	Filter	maps (%)	folds (%)

maps	folds	Filter	maps (%)	folds (%)
204,783	212	R <sub>gc</sub>	182,696 (89.2)	210 (99.1)
		Loop	163,836 (80.0)	211 (99.5)
		Adj	161,470 (78.8)	212 (100)
		$R_{\rm gf}$	156,192 (76.3)	209 (98.5)
		Sheet	14,057 (7.0)	199 (93.9)
		Strand	13,336 (6.5)	192 (90.6)

 $R_{gc},$  coarse  $R_g;$  Loop, loop lengths; Adj, adjacent parallel;  $R_{gf},$  fine  $R_g;$  Sheet, poor  $\beta$  sheets; Strand,  $\beta$  strand with too few contacts.

structures where all distance restraints (loop length and/or experimental) are satisfied simultaneously.

After application of all the other filters, matches contained entirely within another match are considered redundant, and removed.

# Maps removed by each filter

It is illustrative to consider the fraction of maps removed by each of the filters described above. For example, scanning with a pattern derived from a DSSP assignment of secondary structure for thioredoxin that allows for two secondary structure element deletions from the query and five from the database, the initial alignment of secondary structure elements reduces the number of folds from  $377 \rightarrow 212$ . 165 folds have no match of secondary structures with the thioredoxin pattern. Table 1 illustrates the fractions of the initial 204,783 maps within 212 folds that are removed by each filter when applied independently. Table 2 shows for the same example, how the number of maps drops as the filters are applied in succession. The filters are independent of one another apart from consistency filtering, which must be applied after loop and distance restraint filtering, and redundancy filtering, which must be applied last. The order of filters shown in Table 2 was chosen so as to optimise speed.

The gradual elimination of maps and folds shows how the simple principles of protein structure are sufficient to reduce the number of possible alignments by two orders of magnitude. Interestingly, the number of folds drops very little after the generation of maps, suggesting that the filters are tending mostly to remove nonsensical maps associated with each identified fold rather than ruling out folds. Note that consistency filtering tends only to remove maps

Table	3.	Matrix for	scoring al	ignment of acc	essibilities
		b	e	u	gap
1.		0	0	0	1

b	2	-2	0	-1
e	-2	2	0	-1
u	0	0	0	-1
gap	-1	-1	-1	
b, burie the end.	ed; e, exposed;	u, unknown;	gap, residue t	hat overhangs

when tight loop lengths or distance restraints are included in the pattern.

# Fitting sequences on to 3D structures

Accessibilities for residues within each map are calculated quickly by exploiting the relationship between relative accessibility and the number of other  $C^{\beta}$  atoms within 7 Å ( $N_{C\beta7}$ ) of a residue's  $C^{\beta}$ atom.  $N_{CB7}$  is calculated by considering secondary structures and the C-terminal coils for the matched structures. Analysis of the high quality domains shows that helical residues are buried (b) when  $N_{C\beta7} \ge 3$ , exposed (e) when  $N_{C\beta7} = 0$  and intermediate/unknown (u) otherwise. Similarly, residues in  $\beta$ strands are b when  $N_{CB7} \ge 6$ , e when  $N_{CB7} \le 3$  and u otherwise. In the examples presented here, predicted accessibilities were taken from the SUB line within PHD (Rost & Sander, 1994) output, which highlights those regions predicted with confidence.

Given assignments of accessibility, the best alignment for each pair of secondary structures not permitting gaps within either secondary structure is found by applying the scoring matrix shown in Table 3. These values were chosen to prevent long overhanging gaps in the alignment of predicted and experimental secondary structures, and designed not to penalise mismatches too heavily. The total similarity score for the alignment is then defined as:

$$\left(\sum_{i=0}^{i=N} S_{\rm acc}\right) - L_{\rm diff}$$

where  $S_{acc}$  is the best score for a pair of matched secondary structures calculated by summing values from Table 3, *N* is the number of matched secondary structures, and  $L_{diff}$  is the total difference in the lengths of the two protein domains being compared. When calculating  $L_{diff}$  those secondary structures that have been equivalenced are ignored, since overhanging gaps are already penalised by the gap score in Table 3.

Table 2. Effect of filters when applied sequentially

									_
	Sheet	$R_{\rm gc}$	$R_{\rm gf}$	Loop	Strand	Adj	Cons.	Red.	
Maps 2	$204,783 \longrightarrow 1$	$4,057 \longrightarrow 1$	$3,575 \longrightarrow 1$	$2,534 \longrightarrow 1^{\circ}$	$1,435 \longrightarrow 702$	$74 \longrightarrow 5^{\circ}$	$108 \longrightarrow 5$	$108 \longrightarrow 25$	541
Folds	212 →	199 →	$195 \longrightarrow$	$194 \longrightarrow$	$192 \longrightarrow 12$	78 →	176 —→	176 → 1	176
R <sub>ac</sub> .	coarse R <sub>a</sub> : Lo	op. loop len	oths: Adi, a	diacent para	allel: R <sub>ef</sub> , fine	Ra: Shee	et poor ß	sheets: Str	rand

 $R_{gc}$ , coarse  $R_{g}$ ; Loop, loop lengths; Adj, adjacent parallel;  $R_{gf}$ , fine  $R_{g}$ ; Sheet, poor  $\beta$  sheets; Stranc  $\beta$  strand with too few contacts. Cons., consistency; Red., redundancy.

## Protein structure patterns for evaluation

Representatives (queries) from each of 11 structural families containing structural similarities despite no sequence similarity (Russell & Barton, 1994) were chosen to assess the method. The 11 queries are shown in Table 4 and represent a diversity of folds from all four protein folding classes. For all queries, there is at least one clear example of a similar fold in the database that does not show any detectable sequence similarity to the query. For reference, similar folds in the database were found by the STAMP (structural alignment of multiple proteins) structure comparison program (Russell & Barton, 1992) and with reference to the structural classification of proteins (**scop**) database (Murzin *et al.*, 1995).

Two patterns were defined for each of the 11 structures: (1) one taken directly from the DSSP secondary structure assignment and accessibility (i.e. perfect prediction) and (2) one from cross-validated secondary structure and accessibility prediction by the methods of Rost & Sander (1993, 1994). The PHD program and jack-knifed neural network architectures were kindly provided by Dr Burkhard Rost (EMBL). Experimental secondary structure summaries and accessibilities (a) were taken from DSSP (Kabsch & Sander, 1983). Predicted secondary structure summaries (b) were taken from the "PHD sec" entries and accessibilities from the "SUB acc" entries, since these most closely resembled the assignments from the  $N_{C\beta7}$  calculation of accessibility. PHD assignments of buried (b) and exposed (e) states were classified as buried (b) and exposed (e), with all other positions (i or no assignment) as unknown (u). Strands shorter than two residues, and helices shorter than four residues were ignored. The length of the secondary structure was given by the number of residues in each secondary structure (maximum = minimum), and the number of residues between the secondary structures was taken as the minimum loop length.

Patterns may also contain distance restraints, such as those available from NMR experiments, disulphide linkages, or SDM studies. Distance restraints were only added in the von-Willebrand factor and proteasome patterns (see Results).

# **Cross-validation**

Any predictive method that needs large numbers of parameters must be cross-validated to ensure that the method does not do artificially well on the examples used to derive the parameters. For cross-validation of the secondary structure and accessibility predictions, we used the jack-knifed neural-network architectures described by Rost & Sander (1993). Secondary structure and accessibility for each query protein were predicted by an architecture that did not include the query protein or any homologue.

The filters and matching algorithm described here use only a few geometric parameters all of

Table 4.	Proteins	used	to	assess	the	method	

0.1		~	
Code	Protein name	Class	Fold
hnf3	Hepatocyte nuclear	$\alpha + \beta$	Winged helix-turn-helix
	factor		DNA binding motif
1mba	Myoglobin	α	Globin
1plc	Plastocyanin	β	Greek-key $\beta$ sandwich
1rcb	Interleukin-4	α	Up-up-down-down 4-helix bundle
1shaa	<i>v-src</i> tyr kinase SH2 domain	α + β	SH2 domain fold
1ubq	Ubiquitin	α+β	β-Grasp
1wsya	α-Subunit of trp synthase	α/β	α/β-Barrel
2hmqa	Hemerythrin	α	Up-down-up-down 4-helix bundle
2pgd_I	6-Phosphogluconate dehydrogenase	α/β	Rossmann fold
2trxa	Thioredoxin	$\alpha/\beta$	Thioredoxin
4fgf	Basic fibroblast growth factor	β	β-Trefoil

Codes are the Brookhaven PDB code postfixed with the chain identifier code and domain number in Roman numerals where appropriate. hnf3 is assigned to a structure not in the PDB.

which are independent of the protein sequence. Accordingly, removal of query proteins and homologues from the set used to derive the equations above makes a negligible difference to the parameters.

## **Computational details**

Runs for the patterns shown in Table 4 take between 5 and 60 minutes on a Silicon Graphics Indigo 2 (150 MHZ IP22 Processor MIPS R4400). The MAP program is available from the authors. Contact GJB by e-mail: gjb@bioch.ox.ac.uk or see the WWW address http://geoff.biop.ox.ac.uk/ for details.

# Results

#### Assessing accuracy

Structural similarity is a continuum and for some fold types opinions differ as to what constitutes "similar". For example, thioredoxin has a  $\beta$ -sheet with helices packing on each side which superficially resembles a Rossmann fold domain. However, the topology of the sheet is different from a Rossmann fold: the connectivity is different, and it contains a mixture of parallel and antiparallel  $\beta$ hairpins rather than all parallel. To build a detailed model of thioredoxin based on a Rossmann fold would be incorrect, but recognising that thioredoxin has a "single sheet with helix on each side" is still useful. For some folds, e.g. the  $\beta$ -trefoils, there is no such ambiguity. We discuss the accuracy of our method using two grades of success "strict" and "loose", which are outlined in Table 5. Strict similarities are those where the topology of the structure in the database is nearly an exact match of that found in the query (e.g. plastocyanin and azurin). Loose similarities are those where the topologies are broadly similar, with additional

Query	Strict match	Loose match
hnf3	Winged DNA binding HTH domain	Any HTH DNA binding domain
1mba	Globins, phycocyanins colicin A	None
1plc	Cupredoxins	Any Greek key β sandwich
1rcb	Any up-up-down-down 4-helix bundle	Any four helix bundle
1shaa	BirA domain II	None
1ubq	Any β-grasp fold	None
1wsya	Any $\alpha/\beta$ barrel fold	None
2hmqa	Any up down up down 4-helix bundle	Any four helix bundle
2pgd_I	Rossmann folds	Any doubly wound α/β domain
2trxa	Any thioredoxin fold	Any doubly wound $\alpha/\beta$ domain
4fgf	Any β-trefoil	None

Table 5. Strict and loose matches with each query

secondary structures in one fold relative to another, and with some differences in topological ordering or orientation of equivalent secondary structure elements (e.g. plastocyanin and an Ig fold). Strict similarities tend to correspond with those specified by **scop** (Murzin *et al.*, 1995), whereas the loose similarities tend to correspond roughly with those identified by CATH (Orengo *et al.*, 1993) and by the assessors of the protein structure prediction challenge (Lemer *et al.*, 1996).

For comparison, we also scanned the same 11 queries against the database of domains using the fold recognition program THREADER (Jones *et al.*, 1992) with default parameters.

In addition to the recognition of the correct fold, it is important to consider how well the query is aligned onto the database structure. Two measures of alignment accuracy are given: (1) the fraction of correct residue equivalences found by each method % Res-Res, and (2) the fraction of correctly overlapping secondary structure elements found % Sec-Sec. Secondary structures were considered correctly matched if at least two residues from structurally equivalent secondary structures overlapped in the alignment generated by each method. % Res-Res is a strict definition, and broadly measures how accurate a 3D model would be if based on the alignment found. % Sec-Sec is a looser definition, and allows for slippages of secondary structures and thus indicates the accuracy of the predicted topology. The second measure is arguably a more reliable guide, since for many pairs of similar protein structures, alignments of sequence based on 3D structure are ambiguous. Problems arise when assessing the symmetrical  $\alpha/\beta$  barrel structures. Shifting the alignment of secondary structure elements by one  $\beta \alpha$  unit can lead to zero accuracy by these measures, though the resulting structure is largely correct. We thus report average accuracies with and without the  $\alpha/\beta$  barrels. To assess the overall alignment accuracies of each method, only those strict similarities that were not detectable by a sensitive sequence comparison algorithm (Barton, 1993) were considered. Similarities excluded were those with the globins, 1ECA, 1HBG and 1MYGA when scanning with sea hare myoglobin (1MBA), and that with 1PAZ when scanning with plastocyanin (1PLC). For all other examples, accuracies were included in the calculation of an average, regardless of whether the similarity was found at or near the top of the ranked lists. A total of 36 strict similarities were used in the calculation.

# Searches with 11 test proteins

The results of comparing the 11 protein structures to the database of domains using DSSP patterns, PHD patterns, and the THREADER program are shown in Table 6. The Table lists the top ten ranked domains for each query by each method. For each domain, the code, score, structural class and fold description are shown together with the alignment score and the percentage accuracies of the alignments at the residue (% Res-Res) and secondary structure (% Sec-Sec) level (see below). Within Table 6, domains classified as strict similarities (ignoring those detectable by sequence comparison) are shown in inverse text; loose similarities are shown as shaded. Table 7 summarises the rankings shown in Table 6 (see the legend).

Judging by the strict criteria shown in Table 5, 8/11 of the scans made with experimentally determined secondary structure (MAP(DSSP)) put the correct fold in the first rank. By the loose definition, the method located 10/11 folds in the first rank. Predictably, the scans based on patterns from secondary structure prediction fare worse. 4/11 folds were correctly ranked at position 1 by the strict criteria. However, this compares favourably with THREADER which placed one fold correctly in the first rank. When the loose definitions of fold similarity are used, our method placed 5/11 correct folds at the top of the list compared to 2/11 for THREADER. Expanding the definition of success to include any search that places a correct fold in the top ten, as described by Lemer et al. (1996), shows a similar trend (Table 7). The greater success of the DSSP derived patterns suggests that fold recognition by this method will improve alongside any improvements in secondary structure and accessibility prediction. The structural class of proteins (as identified using scop) in the top ten domains was more consistent by our method: MAP(PHD) scans lead to 10/11 correct protein class predictions for the first ranked protein, compared to 5/11 for THREADER. Although this improvement may be due mostly to the accuracy of the PHD predictions, the result suggests that other fold recognition methods could profit from the consideration of predicted secondary structures.

2

Our method (MAP) shows an improvement over THREADER with respect to detecting the correct fold. What of alignments of sequence to structure? Values for individual accuracies are given in Table 6. Reference alignments of 3D structures were found by the STAMP algorithm (Russell & Barton, 1992)

									Table 6								
			I. MAP(DSSP)						II. MAP(PHD)						III. Threader		
					H	NF-3 (w	vinge	d h	elix-turn-helix, HTH	H, Not	in Pl	DB)	1.25				
1BIA_I	27	α±β	W-HTH(HTH)	73.5	100.0	1r69	8	α	λ-rep(HTH)	Contraction of the	1	2fx2	-3.51	α/β	FLAVO(DWAB)		
3GAPA II	13	$\alpha + \beta$	W-HTH(HTH)	87.5	100.0	1avr_IV	0	α	annexin			1cd8	-2.63	β	IG(GKBS)		
1atna_III	8	α+β	actin			4blma_11	0	α+β	β-lactamase			1rn4	-2.38	α+β	ribo		
2end	2	α	T4.endonuclease			1avr_II	-3	α	annexin			1avr_IV	-2.13	α	annexin		
4icb	1	α	EF-HAND			2ts1_II	-4	α	Tyr tRNA synth		1.0	1hrha	-1.92	α/β	ribo		
4tnc_II	-3	α	EF-HAND			2c2c	-5	α	CYTOC			2trxa	-1.86	α/β	THIO(DWAB)		
3sdpa	-7	α	EF-HAND			1ycc	-6	α	CYTOC			2npx_II	-1.83	α+β	FAD-BIND		
1gmfa	-8	α	4HB-2(4HB)			1avr_I	-7	α	annexin		10	1aps	-1.81	α+β	acyl-phosphatase		
1ppn_I	-8	α+β	papain			1Imb3	-11	α	λ-rep(HTH)		The state	5ruba_l	-1.69	α+β	α/β-BARREL		
2aaia_II	-9	α+β	ricin a chain			1cc5	-13	α	CYTOC			1trb_II	-1.66	α+β	FAD-BIND		
						Sea	Hare	M	yoglobin (globin-fol	d, 1ME	BA)						
1HBG	68	α	GLOBIN	90.5	85.7	1HBG	33	α	GLOBIN	51.6	85.7	1ECA	-3.10	α	GLOBIN	72.2	100.0
1MYGA	65	α	GLOBIN	94.2	87.5	1MYGA	22	α	GLOBIN	31.4	87.5	1HBG	-2.94	α	GLOBIN	47.0	100.0
1ECA	49	α	GLOBIN	59.4	75.0	1ECA	18	α	GLOBIN	60.4	87.5	1pgd_l	-2.75	α/β	ROSS(DWAB)		
1COLA	22	α	GLOBIN	11.3	66.7	1COLA	-9	α	GLOBIN	0.0	33.3	1COLA	-2.70	Ο.	GLOBIN	18.1	66.7
1CPCA	17		GLOBIN	63.0	80.0	1CPCA	-12		GLOBIN	0.0	40.0	1bia_II	-2.34	α+β	SH2-like		
256ba	5	α	4HB-1			1ezm_II	-14	α	thermolysin			3chy	-2.32	α/β	FLAVO(DWAB)		
1gsta II	5	α	glut-S-trans.			1gsta_II	-16	α	glut-S-trans.			1llda_II	-2.09	α+β	ldh		
1avr II	2	α	annexin			4blma II	-17	α+β	β-lactamase			1llda I	-2.06	α/β	ROSS(DWAB)		
1ezm_II	2	α	thermolysin			2wrpr	-20	α	Trp rep(HTH)		- 0	1lfi_ll	-2.01	α/β	PBL(DWAB)		
1vsga_II	2	α	surface.glyc.			2scpa	-24	α	EF-hand		1	1pii_II	-1.95	α/β	α/β-BARREL		
			and the second second		Plas	tocvan	in (C	UP	. Greek-kev-beta-sa	andwid	ch. 1	PLC)			1. 11.14 (V2)		
2AZAA	18	В	CUP(GKBS)	66.7	100.0	1snwa I	-10	β	ser-prot			1PAZ	-3.16	ß	CUP(GKBS)	29.2	100.0
1hoe	12	ß	α-amylase.inh(GKBS)			1cd8	-12	B	IG(GKBS)	CT Years	- 2.43 (1981	6taa II	-2.83	ß	a-amylase.inh(GKBS)		A BALL
2bbkl	9	B	met.dehvdr(GKBS)	And Anna Andrews	anala in the	1bova	-13	ß	OB-FOLD	THURSDAY AND THE	0005666680.23	3cd4 I	-2.66	ß	IG(GKBS)		
2fbih II	7	ß	IG(GKBS)		Contraction of the	6taa II	-14	B	g-amylase inh(GKBS)	The Statement	AL REAL	1tie	-2.56	ß	β-TREFOIL	AND DESCRIPTION OF A DE	Carlo Specification (1)
3cd4 I	6	ß	IG(GKBS)		10 1 - 250 10 - 10 - 250	1hoe	-16	ß	a-amylase.inh(GKBS)		ALC: NOT	4faf	-2.28	ß	B-TREFOIL		
1snwa II	6	ß	ser-prot	and the second se	are and a second second	2sn3	-16	$\alpha + \beta$	scorpion.toxin	A CONTRACTOR OF ST	and a present of the	1tlk	-2.23	ß	IG(GKBS)	West Const Con	a contraction of the
2aaib II	4	ß	B-TREFOIL			1hcc	-18	ß	factor-H. 16th			2fbil I	-2.04	ß	IG(GKBS)		
2aaib	3	ß	B-TREFOIL			1tik	-18	ß	IG(GKBS)	ter and the start		2aaib I	-2.03	ß	ß-TREFOIL	Contraction of the Article	And Case of the Distance
1hsab	2	B	IG(GKBS)	and the state of the	ALC: NOT	4acr I	-18	ß	v-crystallin	acasees and a second second	Station and a lot	2mcm	-2.00	ß	macromycin(GKBS)	el istration	C.C.S. LAW
tooba		ß	sod(GKBS)	tille and a start of	AN SUR	24744	-20	ß	CUP(GKBS)	37.8	100.0	Inaha	-1 08	B	prealburnin/GKBS)		

Table 6. Evaluation of mapping

cont.

							Table 6				
			. MAP(DSSP)				II. MAP(PHD)				III. Threader
			Interle	eukin-4 (u	ıp-ι	ip-d	lown-down, four helix bundle,	1RCB)			
256ba	33	α	4HB-1	1avr_III	. 9	α	annexin	1pfka_II	-2.89	α/β	pfk(DWAB)
1bbha	32	α	4HB-1(4HB)	256ba	7	a	4HB-1	5p21	-2.65	α/β	G-PROT(DWAB)
20cya	29	α	4HB-1(4HB)	1GMFA	6	α	4HB-2(4HB) 0.0 100.0	1avr_l	-2.48	α	annexin
2asr	20	α	4HB-1(4HB)	7cata_III	4	α	catalase	1llda_l	-2.44	α/β	ROSS(DWAB)
1lpe	17	α	4HB-1(4HB)	2hmqa	2	α	4HB-1(4HB)	2mnr_II	-2.25	α/β	α/β-BARREL
1GMFA	16	α	4HB-2(4HB) 27.9 100.0	2end	-2	α	T4.endonuclease	2fx2	-2.09	α/β	FLAVO(DWAB)
1avr_III	16	α	annexin	1utg	-4	α	uteroglobin	6abp_II	-2.05	α/β	PBL(DWAB)
2hmqa	13	α	4HB-1(4HB)	2yhx_III	-4	α	hexokinase	2trxa	-1.99	α/β	THIO(DWAB)
7cata_III	11	α	catalase	1avr_II	-5	α	annexin	1ofv	-1.88	α/β	FLAVO(DWAB)
2end	7	α	T4.endonuclease	1avr_I	-6	α	annexin	1avr_II	-1.88	α	annexin
2.14				S	H2	dor	nain (SH2-fold, 1SHAA)				
1atna II	-10	α+β	ACTIN-ATPASE	3rubs	-10	α+β	RuBisCO small sub	1i1b	-3.12	β	β-TREFOIL
1ppn II	-11	$\alpha + \beta$	papain	4blma I	-14	α+B	B-lactamase	2fbjh_l	-2.99	β	IG(GKBS)
2npx III	-13	$\alpha + \beta$	nadp.perxidase	1atna II	-14	α+β	ACTIN-ATPASE	1tie	-2.80	β	β-TREFOIL
5ruba I	-14	$\alpha + \beta$	α/β-BARREL	1bova	-17	ß	OB-FOLD	1cd8	-2.32	β	IG(GKBS)
1sat I	-16	ß	ser-prot	1aps	-18	α+β	acyl-phosphatase	3dfr	-2.31	α/β	dhfr(DWAB)
1thm II	-17	α/β	subtilase	2fxb	-18	α+β	ferredoxin	1bbt3	-2.09	β	JELLYROLL
3b5c	-18	$\alpha + \beta$	cvtoc.b5	2hpr	-18	α+β	HPr	2fbjl_l	-2.00	β	IG(GKBS)
1aaqb	-18	ß	asp-prot	3hsc II	-18	α+β	ACTIN-ATPASE	6taa_II	-1.91	β	α-amylase.inh (GKBS)
2rspa	-18	β	asp-prot.	1bia III	-18	β	SH3-like	5fbpa_II	-1.80	α/β	sugar.phosph(DWAB)
2fbil II	-18	β	IG(GKBS)	5fd1	-19	α+β	ferredoxin	8adh_ll	-1.79	α/β	ROSS(DWAB)
					U	biaı	uitin (β-grasp, 1UBQ)				
1FXIA	12	α+B	β-GRASP 47.8 80.0	1PGX	4	α+β	β-GRASP 46.2 80.0	6taa_II	-2.30	β	α-amylase.inh(GKBS)
1PGX	11	α+B	β-GRASP 37.0 100.0	1alc_II	-8	α+β	a-lactalbumin	1avr_IV	-1.91	α	annexin
2sn3	6	α+β	scorpion.toxin	1rn4	-10	α+β	ribo	1bova	-1.79	β	OB-FOLD
1choi	6	$\alpha + \beta$	ovomucoid.inh	2cpke I	-10	α+β	ser/thr-kinase	1PGX	-1.73	$\alpha + \beta$	β-GRASP 43.8 80.0
1bova	5	ß	OB-FOLD	2sn3	-11	$\alpha + \beta$	scorpion.toxin	2msbb	-1.66	α/β	lectin
1m4	4	α+β	ribo	1choi	-11	α+β	Ovomucoid.inh	1cd8	-1.65	β	IG(GKBS)
1thm II	3	α/β	subtilase	3monb	-12	α+β	monellin	1gsta_I	-1.52	α/β	THIO(DWAB)
1fkb	0	α+β	FKBP-like	3sici	-12	α+β	subtilisin.inh	1trb_II	-1.45	α+β	FAD-BIND
1tasi	0	$\alpha + \beta$	ser-prot.inh.	3il8	-13	α+β	interleukin-8	1aba	-1.33	α/β	THIO(DWAB)
3sici	-2	α+β	subtilisin.inh	1ltsd	-14	β	OB-FOLD	3chy	-1.31	α/β	FLAVO(DWAB)

4

Table 6. cont.

.

cont.

Table	6.	cont.		

<· ·

								Table 6					
11. 31	1 menters.	. MAP(DSSP)	harts-	1.1	1			II. MAP(PHD)			11010	Color.	III. Threader
	ARCH IN	and the second state of the second	Includes	. Constant	Tryptop	bhai	n sy	in thase ( $\alpha/\beta$ barrel,	1WSYA	A)†	an a	TANK .	
2TMDA I	<b>-19</b> αβ	α β-BARREL	0.0	0.0	1pfka_I	15	α/β	pfk(DWAB)	1.1	2liv_l	-3.15	α/β	PBL(DWAB)
1XIS	<b>-70</b> αβ	α β-BARREL	0.0	18.2	1pgd_l	14	α/β	ROSS(DWAB)		1pgd_l	-2.92	α/β	ROSS(DWAB)
ruct orn.	100				2MNR_II	6	αβ	α β-BARREL	62.8 8	1.8 1PII_I	-2.61	αβ	α β-BARREL 0.0
					2hsda	1	α/β	ROSS(DWAB)		2pmga_l	-2.53	α/β	P-glucomutase(DWAB)
					3pgk_II	-2	α/β	pgk(DWAB)		3chy	-2.51	α/β	FLAVO(DWAB)
				1.1	1akea	-9	α/β	p-loop.NTP.hydrolase		2fx2	-2.23	α/β	FLAVO(DWAB)
					8acn_I	-11	α/β	aconitase		1lap_II	-2.19	α/β	ROSS(DWAB)
				- 1	1wsyb_II	-22	α/β	Trp.synthase(DWAB)		5p21	-2.13	α/β	G-PROT(DWAB)
					1gky	-23	α/β	guanylate.kinase(DWAB)		1ofv	-2.10	α/β	FLAVO(DWAB)
					2pmga_II	-32	α/β	p-glucomutase(DWAB)		8adh_ll	-2.01	α/β	ROSS(DWAB)
	11			Hen	nerythrir	n (41	HB	I, 4HB, four helix bu	undle, 2	HMQA)			
256BA	<b>28</b> a	4HB-1	50.0	100.0	1avr_III	-2	α	annexin		8adh_II	-2.56	α/β	ROSS(DWAB)
prcm_II	19 α	prc memb. dom.			256BA	-7	α	4HB-1	0.0 2	5.0 3chy	-2.51	α/β	FLAVO(DWAB)
lutg	18 α	uteroglobin			7cata_III	-11	α	catalase		1pgd_l	-2.49	α/β	ROSS(DWAB)
1avr_III	18 α	annexin			1utg	-12	α	uteroglobin		2npx_IV	-2.37	α+β	nadp.peroxidase
7cata_III	16 α	catalase			1avr_II	-12	α	annexin		1ofv	-2.18	α/β	FLAVO(DWAB)
1avr_II	11 α	annexin			1avr_I	-13	α	annexin		2rsla	-2.17	α/β	γδ-resolvase
1LPE	10 α	4HB-1(4HB)	0.0	100.0	2ts1_II	-13	α	Tyr.tRNA.synth		2tmda_III	-2.12	α/β	tri-metam.dehydr.(DWAB)
Igsta_II	7α	glut-S-trans.	0-0-01		1prcc_I	-17	α	prc.cyt.dom.		2npx_II	-1.98	α+β	FAD-BIND
2yhx_III	7 α	hexokinase			1fiaa	-18	α	fis.protein		1avr_IV	-1.95	α	annexin
1BBHA	6 α	4HB-1(4HB)	30.4	100.0	1avr_IV	-19	α	annexin		1trb_II	-1.91	α+β	FAD-BIND
States 3	STR. N.E.	6-	phosp	hogli	uconate	deh	nvdi	rogenase (Rossmar	nn fold.	2PGD dor	nain	)	a second s
Schy	41 a/B	FLAVO(DWAB)			1LLDA I	1	αβ	ROSS(DWAB)	31.4 7	0.0 1pii_II	-2.80	α/β	α/β-BARREL
MDHA I	<b>40</b> α β	ROSS(DWAB)	60.4	87.5	3chy	-1	α/β	FLAVO(DWAB)		2liv_II	-2.73	α/β	PBL(DWAB)
ILLDA I	<b>37</b> α β	ROSS(DWAB)	40.4	80.0	4MDHA	-7	αβ	ROSS(DWAB)	21.4 6	2.5 1pfka_II	-2.70	α/β	pfk(DWAB)
3PGK I	<b>32</b> αβ	ROSS(DWAB)	38.9	100.0	5p21	-8	α/β	G-PROT(DWAB)		3chy	-2.33	α/β	FLAVO(DWAB)
BADH II	<b>29</b> αβ	ROSS(DWAB)	60.4	90.0	1pfka_II	-10	αβ	pfk(DWAB)		2pmga_II	-2.20	α/β	p-glucomutase(DWAB)
pfka_l	24 α/β	-pfk(DWAB)		mailane	2fx2	-12	α/β	FLAVO(DWAB)	Salar and a state of the	1gpr	-2.01	β	glucose.permease
116_11	24 a/B	PBL(DWAB)			2tmvp	-17	α	4HB-1(4HB)		2npx_II	-1.97	α+β	FAD-BIND
2tmda_II	20 α/β	tri-metam.dehydr(DWAB)			1rhd_II	-17	α/β	rhodanese		1trb_II	-1.79	α+β	FAD-BIND
letu	19 α/β	G-PROT(DWAB)			1etu	-19	α/β	G-PROT(DWAB)	The sector	1grca	-1.75	α/β	glyc.ribo.trans.(DWAB)
2liv H	17 g/B	PBL(DWAB)			2liv II	-23	ov/B	PBL(DWAB)		1drf	-1.69	α/β	dhfr(DWAB)

INVE (GHD)

×. 1

cont.

		I. MAP(DSSP)					Table 6 II. MAP(PHD)						III. Threader		
Constant Second	223.39	and the second		Th	iored	lox	in (THIO, DWAB.	2TRXA	.)		A STATE	1. 1. 1. 1		Pre Vice and	
1EGO	<b>9</b> α β	THIO(DWAB)	38.5 1	00.0 1GSTA I	3	αβ	THIO(DWAB)	0.0	60.0	IGSTA I	-3.36	αβ	THIO(DWAB)	0.0	40.0
1GSTA I	<b>5</b> α β	THIO(DWAB)	27.3 1	00.0 Schy	2	α/β	FLAVO(DWAB)	A CALL THE REAL PROPERTY OF		1snc	-2.75	ß	OB-FOLD		
3grs_III	4 α+f	3 glutathione.red		5fbpa_II	1	a/B	sugar.phosph(DWAB)	films the sheet	and the state	2fbih I	-2.64	ß	IG(GKBS)		
1ABA	-4 αβ	THIO(DWAB)	30.4 1	00.0 2fx2	-1	a/B	FLAVO(DWAB)	A Line of the local design	(YEL)	lofy	-2.24	α/β	FLAVO(DWAB)	A PROPERTY OF THE OWNER	-
4mdha_l	-4 a/B	ROSS(DWAB)		4mdha 1	-1	αß	BOSS(DWAB)	A State of the sta	Contraction of	1llda I	-2.03	ov/B	BOSS(DWAB)		
6abp_l	-8 a/β	PBL(DWAB)	And the second se	1pfka_II	-5	CVB	pfk(DWAB)	Carrie Carrie		2msbb	-1.96	α/β	lectin	Automatic Address of	And the owner of the owner
2fx2	-10 a/P	FLAVO(DWAB)	and the second	1EGO	-6	C B	THIO(DWAB)	0.0	20.0	4faf	-1.92	ß	B-TREFOIL		
1GP1A	<b>-10</b> α β	THIO(DWAB)	34.7 1	00.0 6abp 1	-7	αβ	PBL(DWAB)	and substantial line	and the second	2rsla	-1.88	$\alpha/\beta$	νδ-resolvase		
1ezm_l	-10 α+β	3 thermolysin		1ABA	-8	O B	THIO(DWAB)	57.9	60.0	1ltsd	-1.88	ß	OR-FOLD		
1lida_I	-11 a/β	ROSS(DWAB)	ACCOUNT AND AND A	1llda_I	-8	or/B	BOSS(DWAB)		and the set	4mdha I	-1.76	a/B	BOSS(DWAB)	and a subscript of	1
				Fibrobl	ast o	iro	wth factor (B-TRE	FOIL, 4	FGF	=)		and the second s		All destructions of	Sala and and a
1I1B	<b>23</b> β	β-TREFOIL	35.9 1	00.0 1cd8	-16	ß	IG(GKBS)	1 <b>0</b> 1 <u><u></u>,</u>		11trb	-4.18	α+β	FAD-RIND		
1TIE	<b>21</b> β	B-TREFOIL	35.4 1	00.0 1hoe	-18	ß	a-amylase inh(GKBS)			1cd8	-2 79	ß	IC(CKRS)		
2AAIB II	19 β	β-TREFOIL	0.0	10.0 1coba	-18	ß	sod(GKBS)		1	3chv	-2.34	ρ α/β			
2AAIB I	<b>12</b> β	β-TREFOIL	53.9	70.0 2mcm	-21	ß	macromycin(GKBS)			2nnx II	-2 20	a+B	FAD-RIND		
3cd4_I	8 β	IG(GKBS)		1snwa I	-22	ß	ser-prot		1	2rela	-2 12	a/B			
1coba	2 β	sod(GKBS)		2fbih I	-22	ß	IG(GKBS)		1	1TIE	-2.06	B	RTREEOII	0.0	10.0
1gpr	-3 β	glucose.permease		6taa II	-22	ß	g-amvlase.inh(GKBS)			2mshb	-1.99	α/β	lactin	0.0	10.0
2fbjh_I	-7 β	IG(GKBS)		2bbkl	-23	ß	met dehvdr (GKBS)			1hht3	-1.94	ß			
1tlk	-9 B	IG(GKBS)		4sbva	-23	ß	JELLYBOLL			1mba	-1.83	α+β	ribo		

2AAIB II -23 β β-TREFOIL

Table 6. cont.

2mcm

-11 B

macromycin(GKBS)

Results of running MAP using secondary structure assignments (I) and PHD secondary structure predictions (II) shown beside THREADER results (III) for 11 protein structures having type B and C similarities (Russell & Barton 1994) within the domain database. The first column for each method shows the top ten scoring domains, which are denoted by a PDB four letter code (Bernstein et al., 1977), a chain identifier as the fifth character (if any), followed by an underscore and a Roman numeral denoting the domain (if any). Bold inverted text denotes a correct match using the strict classification, grey backgrounds show loose classifications (see the text). The second column shows the score for each domain, the third the protein structure class, and the fourth the name of the fold/structure. Upper case denotes fold families under the strict definitions. Upper case names in parentheses (if present) denote the name of the loose family classification. The globins 1HBG, 1MYGA and 1ECA and the cupredoxin 1PAZ are sequence similar to the query so are not shown inverted and are not included in the evaluation statistics (see the text).

13.3 40.0 1gpr

-1.79 β glucose.permease

Strict fold classifications: 4HB-1, up-down-up-down four helix bundle (4HB); 4HB-2, up-up-down-down (interleukin-4 type) 4HB; GLOBIN, globin-type folds; W-HTH, winged helix-turn-helix (HTH) folds; EF-HAND, calcium binding EF hands; CYTOC, cytochromes c; THIO, thioredoxin-like folds; FLAVO, flavodoxin-like folds; ROSS, Rossman folds; PBL, periplasmic binding protein-like folds; ACTIN-ATPASE, actin/HSC-70/hexokinse like folds; G-PROT, G-protein (ras) like folds; FAD-BIND, FAD/NAD binding protein-like folds; α β-BARREL, α β (TIM) barrels; β-GRASP, β-grasp (ferredoxin) like folds; IG, immunoglobulin superfamily; CUP, cuppredoxins (plastocyanin-like);  $\beta$ -TREFOIL,  $\beta$ -trefoils (interleukin-1- $\beta$ -like); OB-FOLD, oligonucleotide/oligosaccharride binding folds.

Loose fold classifications: 4HB: 4HB-1, 4HB-2, ferritin; HTH: W-HTH, λ-rep., trp-rep.; DWAB (doubly-wound-α β): ROSS, FLAVO, THIO, PBL, G-PROT: sugar phosphatase, pfk, pgk, dhfr; GKBS (greek key  $\beta$  sandwich), IG, CUP,  $\alpha$ -amylase inhibitor, sod, macromycin, prealbumin.

Other abbreviations: sod, superoxide dismutase; pfk, phosphrofuctokinase; pgk, phosphoglcerate kinase; dhfr, dihdrofolate reducatse; ldh, lactate dehydrogenase; ser-prot, serine proteinase; asp-prot, aspartic proteinase; inh., inhibitor; rep., repressor; glut, glutathione; red., reductase; thym. phosph., thymidine phosphorylase; ribo., ribonuclease; glyc., glycoprotein; P-glucomutase, phosphoglucomutase; glyc. ribo trans., glycinamide ribotransferase.

Method	Strict (1st)	Loose (1st)	Class (1st)	Strict (Top 10)	Loose (Top 10)
MAP (DSSP)	8/11	10/11	11/11	10/11	10/11
MAP(PHD)	4/11	5/11	10/11	9/11	10/11
THREADER	1/11	2/11	5/11	6/11	7/11

Table 7.	Summary	of	fold	recognition	success	rates

Summary of fold recognition success rates. Strict and Loose refer to the criteria for structural similarity discussed in the text. Class refers to structural class success as discussed in the text. (1st) refers to success measured as a correct fold at rank 1, (Top 10) as a correct fold in the top ten ranked structures.

for all strict similarities with the 11 protein families. The averaged values for % Res-Res and % Sec-Sec are shown in Table 8. MAP(DSSP), MAP(PHD) and THREADER give % Res-Res of 35, 15 and 11%, respectively and % Sec-Sec of 75, 43 and 37%. If one ignores the repetitive  $\alpha/\beta$  barrel alignments, accuracies improve slightly with % Res-Res 39, 15 and 13% and % Sec-Sec of 86, 49 and 50 % for MAP(DSSP), MAP(PHD) and THREADER. None of the methods perform well by the % Res-Res criterion, though % Sec-Sec suggests that the correct topology is achieved about 50% of the time. The high % Sec-Sec for MAP(DSSP) scans suggests that alignment accuracy, like fold recognition, will improve with developments in secondary structure and accessibility prediction.

How useful are the detected loose similarities? For some examples, loose similarities imply only a broadly similar architecture, and may not immediately be used for homology modelling studies. However, for others the loose similarity genuinely represents a feasible modelling template. For example, the PHD prediction of hepatocyte nuclear factor 3 (HNF-3) failed to predict two short  $\beta$ strands found in the native structure, and thus the MAP search did not detect BirA domain I (PDB code 1BIA) or GAP domain II (3GAPA) as possible templates. However, the search with the predominantly helical prediction did rank another helixturn-helix motif first, as shown in Figure 1. The core three helices have been aligned correctly at the secondary structure level and a prediction of this type could be useful in the absence of experimental 3D structure information.

#### Fold recognition from published predictions

In the tests above only the type and length of secondary structures, the loop length observed in the query structure, and the pattern of burial and exposure, observed or predicted for each secondary structure segment were used in the search. Many published predictions are augmented by human insight, contain detailed predictions of loop lengths,

 Table 8. Average alignment accuracies for 36 strict similarities

	All alig	nments	Ignoring $\alpha/\beta$ barrels		
Method	% Res-Res	% Sec-Sec	% Res-Res	% Sec-Sec	
MAP (DSSP)	35	75	39	86	
MAP (PHD)	15	43	15	49	
THREADER	11	37	13	50	

and consider experimental distance restraints. All of this information can be used with the MAP method described here. To test the method under these circumstances, we considered three predictions: (1) the von Willebrand factor (vWf) prediction by Edwards & Perkins (1995), (2) the proteasome prediction by Lupas et al. (1994) and (3) a prediction for the phosphotyrosine interaction domain (PID) by Bork & Margolis (1995). All of these predictions were made from very diverse sequences, which is likely to improve prediction accuracy (Russell & Sternberg, 1995). The predictions also comprise carefully constructed sequence alignments, that can provide tight loop-length distance restraints. For the three searches, a larger and more up-to-date database of 780 protein domains was scanned (A. S. Siddiqui, personal communication). Subsequent 3D structure determination has shown all three of these proteins to resemble previously observed folds (Lee et al., 1995; Brannigan et al., 1995; Zhou et al., 1995).

# The vWF domain

Perkins and co-workers (Perkins *et al.*, 1994; Edwards & Perkins, 1995) used an alignment of 92 sequences together with spectroscopic data, and prediction algorithms to predict that the vWf domain would comprise a repeating arrangement of  $\beta$  strands and  $\alpha$  helices. Edwards & Perkins combined a THREADER scan with analysis of the location of active site residues, a putative disulphide bridge, and the principles of protein 3D structure. They suggested that the vWf domain would be most likely to resemble ras p21. The subsequently determined 3D structures (Lee *et al.*, 1995) showed this prediction of secondary structure and fold to be largely correct (Russell & Sternberg, 1995).

Our mapping technique allows many of the features exploited by Perkins et al. to be combined in a search. Figure 2 shows a vWf pattern based on the prediction of Perkins & co-workers (Perkins et al., 1994; Edwards & Perkins, 1995). In addition to a pattern of predicted secondary structures, the pattern also contains detailed information as to the loop lengths, and details of two distance restraints: one from a pair of aspartic acid residues thought to be involved in a metal binding site (constrained to have their axial coordinates within 15 Å), and a putative disulphide bond (constrained to have their axial coordinates within 9.5 Å). A tolerance of t = 4 Å was added to each of these restraints to allow for changes in secondary structure packing across similar protein 3D structures.



**Figure 1.** An example of a useful "loose" similarity between 3D structure detected using the MAP method and a secondary structure prediction. (a) The alignment found by the method between the predicted pattern for HNF-3 and the helical DNA binding motif within phage 434 repressor. Boxed, bold-faced, upper-case regions indicate aligned predicted and experimental secondary structures. Sec denotes the PHD prediction for HNF-3, and a three-state DSSP secondary structure assignment for 434 repressor. Bur shows predicted and experimental states of burial for HNF-3 and 434 repressor: b, buried; e, exposed; u, intermediate/unknown. (b) The equivalent alignment found using the STAMP (Russell & Barton, 1992) structure comparison algorithm. Boxed, bold-faced, upper-case regions dicate structural equivalences. Sec denotes DSSP three-state secondary structures for both proteins. (c) and (d) The crystallographic structures of the matched regions of HNF-3 and 434 repressor, with structurally equivalent residues shown in ribbon/coil format, and unequivalent regions shown as C trace. The N and C-termini of the structures are labelled.

A comparison of the vWf pattern to the database of 780 domains finds elongation factor Tu (PDB code 1ETU), ras P21 (821P) and Che-Y (3CHY) as the three top scoring folds, with other double-wound,  $\alpha/\beta$ , Rossmann-type folds following in the top 20 scoring folds. The top three scoring proteins are highly similar to the recently solved structures of the vWf, with ras P21/elongation factor Tu being the most similar (Lee *et al.*, 1995).

# The proteasome

Lupas *et al.* (1994) predicted the secondary structure for the 20 S proteasome  $\alpha$  subunits by a variety of algorithms. We took their predicted pattern of secondary structure elements and accessibility and searched the database of 780 non-redundant protein domains. Without imposing any experimental distance restraints, the method

finds seven folds (173 maps). The top scoring fold, according to the amphipathicity scoring scheme, is that of glutamine amidotransferase (PDB code 1GPH), which is structurally and functionally similar to the proteasome (Lowe *et al.*, 1995; Brannigan *et al.*, 1995).

5

A small number of weak distance restraints can make a significant difference to the results of this search. If alignment positions identified as putative active site residues by Lupas *et al.*, by the method of Benner and co-workers (Benner *et al.*, 1993), are required to have axial coordinates within 15 Å (tolerance of 4 Å) of each other, only four folds (19 maps) remain, with the correct fold still at the first rank. Although distance restraints are not always available prior to 3D structure determination, our results suggest that they should be used to aid fold recognition whenever possible.



Figure 2. Search pattern for the von-Willebrand factor type A domain (derived from Edwards & Perkins, 1995) as discussed in the text.  $\alpha$  Helices are indicated by cylinders,  $\beta$  strands by arrows. The range of numbers given beside each secondary structure or loop are the range of predicted lengths. Bullets ( $\bullet$ ) show those secondary structures that are required for any possible map (i.e. those involved in distance restraints). Two distance restraints, one from a putative disulphide bond (9.5 Å) and the other from knowledge of two residues thought to be involved in metal coordination (15 Å) are shown to the left of the Figure.

# The phosphotyrosine interaction domain

Bork & Margolis (1995) recently identified a new phosphotyrosine interaction domain (PID) involved in the cytoplasmic signalling cascade. They constructed an alignment of several diverse members of this sequence family, and performed a prediction of secondary structure. We ran the PHD program on a slightly more up-to-date alignment of PID proteins (P. Bork, personal communication) to predict the secondary structure and accessibility. A search pattern was made from the prediction, and the loop length ranges taken from the multiple alignment. The pattern of nine secondary structures was BBHBBBBBH and these elements are numbered

sequentially from one to nine below. Since there were two long loops connecting the predicted secondary structures, the adjacent parallel filter was not used during the search. Structures corresponding to the best alignment with each of the top six scoring folds are shown in Figure 3. Recent structure determination has shown the PID (PTB domain) to resemble the plekstrin homology (PH) domain in structure and function (Zhou et al., 1995). By the accessibility scoring scheme, the top ranked fold is not a PH domain, although a PH domain (from dynamin) is ranked at position 2. The top six folds are illustrative in that they show how the method can suggest alternative plausible folds that satisfy a pattern of predicted secondary structures and accessibilities.

The best scoring fold (Figure 3(a)) is that of profilin (PDB code 2BFPP), and the best scoring map gives an anti-parallel  $\beta$  sheet with the strand order 218754 (predicted strand 6 is deleted) with one helix packing against each face. The second best scoring fold (Figure 3(b)) is a correct match with the PH domain from human dynamin (1DYNB), having deleted the first predicted  $\alpha$  helix from the PID pattern. The third best scoring fold (Figure 3(c)) comes from *Staphylococcus aureus*  $\beta$  lactamase (1BLH, domain 1), with an anti-parallel  $\beta$  sheet of order 54876, with both helices packing against one face. The fourth and fifth best scoring folds (Figure 3(d) and (e)) come from members of the Ig superfamily, and comprise alternative arrangements of  $\beta$  strands to form a Greek key  $\beta$  sandwich. Both of the predicted  $\alpha$  helices from the PID pattern have been deleted in these matches. Finally, the sixth best scoring fold (Figure 3(f)) comes from the tryptic core of Escherichia coli lac repressor (1TLFD domain 4), and comprises a parallel  $\beta$  sheet (42576) with both helices packing against one face. This fold is perhaps the least plausible, since it would require three crossover connections between adjacent and parallel  $\beta$  strands.

The method has suggested plausible alternative structures that can be scrutinised, in the absence of 3D structural data, by way of further experiments, secondary structure predictions, or even other methods of fold recognition. The results show how the predicted secondary structure elements can be accommodated into a compact, plausible protein fold, and encouragingly, the method has identified the correct fold high in the list of alternatives.

# **Discussion and Conclusions**

In this paper we have presented a new method for protein fold recognition which exploits recent improvements in protein secondary structure prediction, and can use other information such as predictions of accessibility, loop lengths and experimental data to restrict possible folds. When applied to predicted secondary structures and accessibilities, the method has been shown to be slightly better than one widely used fold recognition method (Jones *et al.*, 1992) at detecting the correct

.

3



e) Rat CD4 domain 2 (1CID\_II)

f) Tryptic core of lac repressor (1TLFD\_IV)

Figure 3. Maps from the top six scoring folds found during a search with the PID pattern. Details are given in the text.

fold for 11 test examples. The alignments generated by the method are of comparable accuracy at the residue-residue and secondary structure alignment level. When the query is defined by experimental secondary structures and accessibilities, the method is highly successful at recognising the correct fold.

This suggests that the mapping method will improve alongside any future improvement in secondary structure and accessibility prediction. The method also has the advantage of being computationally inexpensive, and so allows for multiple searches to be performed quickly.

The simplicity of the technique suggests several enhancements that could improve accuracy even further. The method of aligning sequences onto 3D structures might be developed by the use of empirically derived pair-potentials or accessibility preferences (e.g. Sippl, 1990; Jones *et al.*, 1992), or by the identification of favourable interaction sites between secondary structures (Cohen & Sternberg, 1980; Cohen *et al.*, 1980, 1982). A more sophisticated alignment and ranking procedure is under development.

The initial alignment and filtering procedures are perhaps the most unique feature of this technique. Other techniques for fold-recognition tend only to provide a single sequence alignment of query and database structures. The use of a secondary structure element alignment method has the advantage that exhaustive comparisons of two proteins can be performed; most folds identified have an ensemble of alternative alignments that can be explored further.

Since most protein structure similarities occur at the domain level, it is advantageous, whenever possible to split both query and database structures into domains. The problem of assigning domains for protein 3D structures has been the subject of revived interest (Holm & Sander, 1994b; Siddiqui & Barton, 1995; Sowdhamini & Blundell, 1995; Islam et al., 1995) and is likely to lead to accessible databases of protein structural domains. Assigning domains within proteins of unknown 3D structure is more problematic, though approaches based on sequence homology (Pongor et al., 1994; Sonnhammer & Kahn, 1994) are undoubtedly the most promising; the vWf and PID proteins above are both examples of domains that occur in a variety of multi-domain contexts.

The method described here has applications in protein structure determination by NMR. During NMR structure determination, a preliminary secondary structure assignment (equivalent to a very accurate prediction) and a small number of distance restraints may be available early in the study. However, these data are usually insufficient to determine a unique structure by distance geometry or molecular dynamics (Smith-Brown et al., 1993). Our results for the vWF and proteasome domains suggest that the data may be sufficient to locate a similar fold in the database if one is present. Folds predicted from distance restraints and secondary structure assignment may be used to guide the assignment of cross-peaks and thus speed up the structure determination process. Clearly, the alternative consistent topologies may also give clues as to possible structural/functional/evolutionary relationships that are generally not known until after 3D structure determination (such as that described by Matthews et al., 1994).

We have shown that secondary structure predictions of typical accuracy, together with simple principles of protein 3D structures and/or experimental data can be used to recognise correct protein folds in a library of domains. These results and others (Edwards & Perkins, 1995; Russell & Sternberg, 1995; Gerloff *et al.*, 1995) suggest that secondary structure prediction, experimental data, and protein structural principles should be used to augment protein fold recognition whenever possible.

#### Acknowledgements

We thank Professor L. N. Johnson for encouragement and support. We are indebted to Dr D. T. Jones (University of Warwick) for giving advice on the THREADER program and its database, Dr B. Rost (EMBL, Heidelberg) for providing the PHD program, Drs P. Bork (EMBL, Heidelberg) and S. J. Perkins (Royal Free Hospital, London) for providing prediction data for the PID and vWF domains, Dr S. K. Burley (Rockefeller University, New York) for providing the coordinates of the HNF-3 structure and Mr A. S. Siddiqui (LMB, Oxford) for providing a database of protein structural domains. R.B.R. thanks Dr C. P. Ponting (Fibrinolysis Research Unit, Oxford) for helpful discussions. R.B.R. and G.J.B. thank the Royal Commission for the Exhibition of 1851 and the Royal Society for support. R.R.C. is funded by an MRC studentship. This research was funded in part by a grant from the BBSRC (UK).

# References

- Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. (1990). Basic local alignment search tool. J. Mol. Biol. 215, 403–410.
- Barton, G. J. (1993). An efficient algorithm to locate all locally optimal alignments. *Comp. App. Biosci.* 9, 729–734.
- Barton, G. J. (1995). Protein secondary structure prediction. *Curr. Opin. Struct. Biol.* **5**, 372–376.
- Barton, G. J. & Sternberg, M. J. E. (1990). Flexible protein sequence patterns: a sensitive method to detect weak structural similarities. *J. Mol. Biol.* 212, 389–402.
- Bazan, J. F. (1990). Structural design and molecular evolution of a cytokine receptor superfamily. *Proc. Natl Acad. Sci. USA*, **87**, 6934–6938.
- Benner, S. A., Badcoe, I., Cohen, M. A. & Gerloff, D. L. (1993). Predicted secondary structure for the src homology 3 domain. J. Mol. Biol. 229, 295–305.
- Benner, S. A., Gerloff, D. L. & Jenny, T. F. (1994). Predicting protein crystal structures. *Science*, 265, 1642–1644.
- Bernstein, F. C., Koetzle, T. F., Williams, G. J. B., Meyer, E. F., Brice, M. D., Rodgers, J. R., Kennard, O., Shimanovichi, T. & Tasumi, M. (1977). The Protein Data Bank: a computer-based archival file for macromolecular structures. J. Mol. Biol. 112, 535–542.
- Bork, P. Margolis, B. (1995). A phosphotyrosine interaction domain. Cell, 80, 693-694.
- Bowie, J. U. & Eisenberg, D. (1993). Inverted protein structure prediction. Curr. Opin. Struct. Biol. 3, 437–444.
- Brannigan, J. A., Dodson, G., Duggleby, H. J., Moody,

- Chothia, C. (1992). One thousand families for the molecular biologist. *Nature*, **357**, 543–544.
- Clark, D. A., Shirazi, J. & Rawlings, C. J. (1991). Protein topology prediction through constraint-based search and the evaluation of topological folding rules. *Protein Eng.* 4, 751–760.
- Cohen, F. E. & Sternberg, M. J. E. (1980). On the use of chemically derived distance constraints in the prediction of protein structure with myoglobin as an example. J. Mol. Biol. 137, 9–22.
- Cohen, F. E., Sternberg, M. J. E. & Taylor, W. R. (1980). Analysis and prediction of protein beta-sheet structures by a combinatorial approach. *Nature*, 285, 378–382.
- Cohen, F. E., Sternberg, M. J. E. & Taylor, W. R. (1982). The analysis and prediction of the tertiary structure of globular proteins involving the packing of  $\alpha$  helices against a  $\beta$ -sheet: a combinatorial appraoch. *J. Mol. Biol.* **156**, 821–862.
- Cohen, F. E., Kosen, P. A., Kuntz, I. D., Epstein, L. B., Ciardelli, T. L. & Smith, K. A. (1986). Structure-activity studies of interleukin-2. *Science*, 234, 349–352.
- Crawford, I. P., Niermann, T. & Kirschner, K. (1987). Predictions of secondary structure by evolutionary comparison: application to the  $\alpha$  subunit of tryptophan synthase. *Proteins: Struct. Funct. Genet.* 1, 118–129.
- Curtis, B. M., Presnell, S. R., Srinivasan, S., Sassenfeld, H., Klinke, R., Jeffery, E., Cosman, D., March, C. J. & Cohen, F. E. (1991). Experimental and theoretical studies of the three-dimensional structure of human interleukin-4. *Proteins: Struct. Funct. Genet.* 11, 111–119.
- Edwards, Y. J. K. & Perkins, S. J. (1995). The protein fold of the von Willebrand factor type A is predicted to be similar to the open twisted  $\beta$ -sheet flanked by  $\alpha$ -helices found in human ras-p21. *FEBS Letters*, **358**, 283–286.
- Gerloff, D. L., Chelvanayagam, G. & Benner, S. A. (1995). A predicted consensus structure for the proteinkinase c2 homology (c2h) domain, the repeating unit of synaptotagmin. *Proteins: Struct. Funct. Genet.* 22, 299–310.
- Gribskov, M., McLachlan, A. D. & Eisenberg, D. (1987). Profile analysis: detection of distantly related proteins. Proc. Natl Acad. Sci. USA, 84, 4355–4358.
- Holm, L. & Sander, C. (1994a). Searching protein structure databases has come of age. *Proteins: Struct. Funct. Genet.* 19, 165–173.
- Holm, L. & Sander, C. (1994b). Parser for protein folding units. Proteins: Struct. Funct. Genet. 19, 256–268.
- Holm, L. & Sander, C. (1995). 3-d lookup: fast protein structure database searches at 90% reliability. Proc. 3rd Int. Conf. Intel. Sys. Mol. Biol. 179–187.
- Huang, Z., Gabriel, J., Baldwin, M. A., Fletterick, R. J., Prusiner, S. B. & Cohen, F. E. (1994). Proposed three-dimensional structure for the cellular prion protein. *Proc. Natl Acad. Sci. USA*, **91**, 7139–7143.
- Hurle, M. R., Matthews, C. R., Cohen, F. E., Kuntz, I. D., Toumadje, A. & Johnson, W. C. (1987). Prediction of the tertiary structure of the α-subunit of tryptophan synthase. *Proteins: Struct. Funct. Genet.* **2**, 221–224.
- Islam, S. A., Luo, J. C. & Sternberg, M. J. E. (1995). Identification and analysis of domains in proteins. *Protein Eng.* 8, 513–525.

- Jin, L., Cohen, F. E. & Wells, J. A. (1994). Structure from function: screening structural models with functional data. Proc. Natl Acad. Sci. USA, 91, 113–117.
- Jones, D. & Thornton, J. (1993). Protein fold recognition. J. Comp. Aid. Mol. Des. 7, 439–456.
- Jones, D. T., Taylor, W. R. & Thornton, J. M. (1992). A new approach to protein fold recognition. *Nature*, 358, 86–89.
- Kabsch, W. & Sander, C. (1983). A dictionary of protein secondary structure. *Biopolymers*, 22, 2577–2637.
- Lee, J., Rieu, P., Arnaout, M. A. & Liddington, R. (1995). Crystal structure of the A domain from the α subunit of integrin CR3 (CD11b/CD18). *Cell*, **80**, 631–638.
- Lemer, C., Rooman, M. J. & Wodak, S. J. (1996). Protein structure prediction by threading methods: evaluation of current techniques. *Proteins: Struct. Funct. Genet.* 23, 337–355.
- Lipman, D. J. & Pearson, W. R. (1985). Rapid and sensitive protein similarity searches. *Science*, 227, 1435– 1441.
- Lowe, J., Stock, D., Jap, B., Zwickl, P., Baumeister, W. & Huber, R. (1995). Crystal structure of the 20s proteasome from archaeon *T. acidophilum* at 3.4 Å resolution. *Science*, **268**, 533–539.

.

5

- Lupas, A., Koster, A. J., Walz, J. & Baumeister, W. (1994). Predicted secondary structure of the 20s proteasome and model structure of the putative peptide channel. *FEBS Letters*, **354**, 45–49.
- Matthews, S., Barlow, P., Boyd, J., Barton, G., Russell, R., Mills, H., Cunningham, M., Meyers, N., Burns, N., Clark, N., Kingsman, S., Kingsman, A. & Campbell, I. (1994). Structural similarity between the p17 matrix protein of HIV-1 and interferon-γ. *Nature*, **370**, 666–668.
- Murzin, A., Brenner, S. E., Hubbard, T. & Chothia, C. (1995). scop: a structural classification of proteins database for the investigation of sequences and structures. J. Mol. Biol. 247, 536–540.
- Orengo, C. (1994). Classification of protein folds. Curr Opin. Struct. Biol. 4, 429-440.
- Orengo, C. A., Flores, T. P., Taylor, W. R. & Thornton, J. M. (1993). Identification and classification of protein fold families. *Protein Eng.* **6**, 485–500.
- Perkins, S. J., F., Smith, K. F., Williams, S. C., Haris, P. I., Chapman, D. & Sim, R. B. (1994). The secondary structure of the von Willebrand factor type a domain in factor b of human complement by Fourier transform infrared spectroscopy. J. Mol. Biol. 238, 104–119.
- Pickett, S. D., Saqi, M. A. S. & Sternberg, M. J. E. (1992). Evaluation of the sequence template method for protein structure prediction. Discrimination of the  $\beta/\alpha$ -barrel fold. *J. Mol. Biol.* **228**, 170–187.
- Pongor, S., Hatsagi, Z., Degtyarenko, K., Fabian, P., Skerl, V., Hegyi, H., Murvai, J. & Bevilacqua, V. (1994). The SBASE protein domain library, release 3.0—a collection of annotated protein-sequence segments. *Nucl. Acids Res.* 22, 3610–3615.
- Richards, F. M. & Kundrot, C. E. (1988). Identification of structural motifs from protein coordinate data: secondary structure and first-level supersecondary structure. *Proteins: Struct. Funct. Genet.* 3, 71–84.
- Rost, B. (1995). TOPITS: threading one-dimensional predictions into three-dimensional structures. Proc. 3rd. Int. Conf. Intel. Sys. Mol. Biol. 314–321.
- Rost, B. & Sander, C. (1993). Prediction of protein secondary structure at better than 70% accuracy. J. Mol. Biol. 232, 584–599.
- Rost, B. & Sander, C. (1994). Conservation and prediction

of solvent accessibility in protein families. *Proteins: Struct. Funct. Genet.* **20**, 216–226.

- Russell, R. B. & Barton, G. J. (1992). Multiple protein sequence alignment from tertiary structure comparison: assignment of global and residue confidence levels. *Proteins: Struct. Funct. Genet.* 14, 309–323.
- Russell, R. B. & Barton, G. J. (1994). Structural features can be unconserved in proteins with similar folds: an analysis of side-chain to side-chain contacts, secondary structure and accessibility. J. Mol. Biol. 244, 332–350.
- Russell, R. B. & Sternberg, M. J. E. (1995). How good are we? *Curr. Biol.* **5**, 488–490.
- Russell, R. B., Copley, R. R. & Barton, G. J. (1995). Protein fold recognition from secondary structure assignments. Proc. 28th Hawaii. Int. Conf. Sys. Sci. IEEE Press, 5, 302–311.
- Sheridan, R. P., Dixon, J. & Venkataraghavan, R. (1985). Generating plausible protein folds by secondary structure similarity. *Int. J. Pept. Protein Res.* 25, 132–143.
- Siddiqui, A. S. & Barton, G. J. (1995). Continuous and discontinuous domains—an algorithm for the automatic-generation of reliable protein domain definitions. *Protein Sci.* 4, 872–884.
- Sippl, M. J. (1990). Calculation of conformational ensembles from potentials of mean force. An approach to the knowledge-based prediction of local structures in globular proteins. *J. Mol. Biol.* **213**, 859–883.

- Smith, T. F. & Waterman, M. S. (1981). Identification of common molecular subsequences. J. Mol. Biol. 147, 195–197.
- Smith-Brown, M. J., Kominos, D. & Levy, R. M. (1993). Global folding of proteins using a limited number of distance constraints. *Protein Eng.* 6, 605–614.
- Sonnhammer, F. L. L. & Kahn, D. (1994). Modular arrangement of proteins as inferred from analysis of homology. *Protein Sci.* 3, 482–492.
  Sowdhamini, R. & Blundell, T. L. (1995). An automatic
- Sowdhamini, R. & Blundell, T. L. (1995). An automatic method involving cluster-analysis of secondary structures for the identification of domains in proteins. *Protein Sci.* 4, 506–520.
- Sun, S., Thomas, P. D. & Dill, K. A. (1995). A simple protein folding algorithm using a binary code and secondary structure constraints. *Protein Eng.* 8, 769–778.
- Taylor, W. R. (1991). Towards protein tertiary fold prediction using distance and motif constraints. *Protein Eng.* 4, 853–870.
- Wodak, S. J. & Rooman, M. J. (1993). Generating and testing protein folds. Curr. Opin. Struct. Biol. 3, 247–259.
- Zhou, M., Ravichandran, K. S., Olejniczak, E. T., Petros, A. M., Meadows, R. P., Sattler, M., Harlan, J. E., Wade, W. S., Burakoff, S. J. & Fesik, S. W. (1995). Structure and ligand recognition of the phosphotyrosine binding domain of shc. *Nature*, **378**, 584–592.

Edited by F. E. Cohen

(Received 19 December 1995; received in revised form 21 March 1996; accepted 1 April 1996)