

Characterization of Novel Proteins Based on Known Protein Structures

Walter A. Koppensteiner, Peter Lackner, Markus Wiederstein and Manfred J. Sippl*

Center for Applied Molecular Engineering, Institute for Chemistry and Biochemistry University of Salzburg Jakob-Haringer-Straße 3 A-5020 Salzburg, Austria

The genome sciences face the challenge to characterize structure and function of a vast number of novel genes. Sequence search techniques are used to infer functional and structural information from similarities to experimentally characterized genes or proteins. The persistent goal is to refine these techniques and to develop alternative and complementary methods to increase the range of reliable inference.

Here, we focus on the structural and functional assignments that can be inferred from the known three-dimensional structures of proteins. The study uses all structures in the Protein Data Bank that were known by the end of 1997. The protein structures released in 1998 were then characterized in terms of functional and structural similarity to the previously known structures, yielding an estimate of the maximum amount of information on novel protein sequences that can be obtained from inference techniques.

The 147 globular proteins corresponding to 196 domains released in 1998 have no clear sequence similarity to previously known structures. However, 75% of the domains have extensive structure similarity to previously known folds, and most importantly, in two out of three cases similarity in structure coincides with related function. In view of this analysis, full utilization of existing structure data bases would provide information for many new targets even if the relationship is not accessible from sequence information alone. Currently, the most sophisticated techniques detect of the order of one-third of these relationships.

© 2000 Academic Press

Keywords: structural genomics; functional genomics; structure prediction; structure comparison; side-chain orientation

*Corresponding author

Introduction

Several bacterial and eucaryotic genomes (TIGR†) have been released and the completion of the human genome project (Rowen *et al.*, 1997) is on the way. The challenge is to assign biological function to the novel sequences of these genomes. A full characterization of a protein contains its

molecular and cellular function, its three-dimensional structure and its interaction with other molecules. Frequently, the function and biological role of a hypothetical protein is inherited from a characterized protein using sequence comparison methods (Altschul *et al.*, 1990; Pearson, 1996).

The basis of sequence comparison is the conservation of structure and function among related proteins (Sander & Schneider, 1991). The limit of reliable inference using sequence similarity is the so-called twilight zone, where similarity becomes indistinguishable from random matches. On the other hand, proteins with insignificant sequence similarity can have similar tertiary structures (Pastore & Lesk, 1990). In fact, nature seems to be able to realize the enormously diverse biological functions by a limited number of folds (Chothia, 1992; Orengo *et al.*, 1994). Consequently, methods are desirable that can detect relationships beyond

Present address: W. A. Koppensteiner, ProCeryon Biosciences GmbH, Jakob-Haringer-Straße 3, A-5020 Salzburg, Austria.

Abbreviations used: FMN, flavin mononucleotide; PDB, Protein Data Bank; rms, root-mean-square; SCOP, structural classification of proteins.

† <http://www.tigr.org>

E-mail address of the corresponding author: sippl@came.sbg.ac.at

the twilight zone. Profile-based techniques like PSI-Blast (Altschul *et al.*, 1997) and Hidden-Markov models (Karplus *et al.*, 1998), and structure-based techniques like fold recognition (Bowie *et al.*, 1991; Sippl & Weitckus, 1992; Jones *et al.*, 1992; Bryant, 1996; Domingues *et al.*, 1999; Jones, 1999) have made progress in this direction (Jones, 1997; Koehl & Levitt, 1999). The success of these methods is limited by the information on structure and function contained in data bases. An analysis of this information content is the focus of this work.

The conservation of structure among distantly related proteins is the basis of structure-based functional annotations of uncharacterized proteins (Martin *et al.*, 1998; Russell *et al.*, 1998; Hegyi & Gerstein, 1999; Orengo *et al.*, 1999). The distinction of analogues and remote homologues has been perceived as a critical factor for the application of structures for function assignment (Flores *et al.*, 1993; Russell & Barton, 1994; Matsuo & Bryant, 1999). Analogous proteins are considered as a product of convergent evolution to a similar three-dimensional structure, while remote homologous originate from a common ancestor. There is agreement that a clear distinction is difficult to obtain because functional relatedness is often hard to prove (Holm & Sander, 1997; Murzin, 1998).

Here, we explore to what extent information on structure and function contained in current data bases can be used to characterize novel genes. The limits of inference of structure and function from data bases can be explored from a set of experimentally determined structures. We present an analysis of structures released by PDB (Bernstein *et al.*, 1977) in 1998. We derive a data set of proteins that do not have sequence similarity to previously determined structures, i.e. proteins that were made public prior to 1998 and we investigate structural and functional relationships to the previously known proteins. We address the following questions. (1) How many new proteins have structural similarity to previously known proteins? (2) How many new proteins are functionally related to previously known proteins? (3) To what extent do structural and functional similarity coincide?

The extent of structural similarity required to model a target from a template depends on the intended purpose and the desired accuracy of the derived model. The level of detail can range from an all-atom model to a rough arrangement of secondary-structure elements. Here, we assume that a suitable template needs to share at least all secondary-structure elements that build up the hydrophobic core of the target. If this condition is met then we consider the structural relationship to be on the fold level. We consider two additional levels of similarity that are relevant for structure prediction, called substructure and partial fold level; when a structure superimposes with a compact substructure of a larger protein, and when two folds have partial similarity where at least major parts of the secondary-structure elements are superimposable. All targets not assignable to one

of these levels are classified as novel folds. We regarded structures A and B to have the same fold when all core secondary-structure elements of A are superimposable with B, so that B can be an adequate template for A. This, however, does not always imply that A is also a suitable template for B. Here, the structures released in 1998 always represent the targets, and the previously determined structures represent the pool of possible templates.

The identification of functional relatedness among proteins requires a detailed and often elaborate analysis of the respective protein structures, especially when the similarity is weak. Here, we rely on the expert knowledge contained in the SCOP data base (Murzin *et al.*, 1995) and the reports of crystallographers and NMR spectroscopists. SCOP classifies protein domains in hierarchical levels called class, fold, superfamily, and family. This hierarchy reflects structural and functional similarities. Since SCOP is updated infrequently, new structures are often not found in the data base, so that we had to consult the reports on the respective structure determinations. Using these sources we determine the functional relationships of domains released in 1998 to previously known structures. We find that two-thirds of structurally similar pairs that are unrelated in sequence have a related function.

Superimposition of the C α traces is the standard approach to determine the extent of structural similarity of two proteins (Taylor & Orengo, 1989; Holm & Sander, 1993). Here, we incorporate side-chain orientation, represented by the C β positions, into structure comparison. Side-chain orientation is relevant for the selection of suitable templates, since any attempt to model a structure will run into enormous difficulties when side-chain orientation is not conserved. Also, side-chain orientation is more conserved in remote homologues than in analogues.

Here, we provide a detailed description of the data sets used, investigate the role of side-chain orientation in distinguishing significant and insignificant structural similarity, and analyze how much structural information can be derived from previously known structures. Finally, we investigate the correlation of structural and functional relationships and provide evidence that side-chain orientation facilitates the distinction of analogous and remote homologous proteins. In Conclusion we compare our results with a recent assessment of structure prediction methods and discuss the relevance of these findings for the structural genomics initiative and the annotation of whole genomes.

Results and Discussion

The data sets

In 1998 the PDB released 1792 new entries containing 3358 protein chains. PDB often contains multiple instances of the same protein. When

redundancies of greater than 95% sequence identity are removed, only 664 new chains remain. This set represents the total amount of structural information released in 1998. From this set, 490 sequences (74%) have significant sequence similarity to previously known structures. In principle, these structures could have been derived from the known data with a reasonable degree of accuracy.

We restrict the analysis to globular proteins and remove all non-globular chains, virus capsid and transmembrane proteins from the remaining 174 chains without significant sequence similarity to known structures. This data set contains 147 protein chains (22% of 664) consisting of 196 domains. There are 107 single-domain proteins in this data set, 32 contain two domains, seven contain three domains, and one contains four domains. The domains have no clear relationship to previously determined proteins and are used for the subsequent structural and functional analysis. The data sets used in this analysis are shown in Figure 1.

Significant versus insignificant similarity

The analysis presented here relies on the detection of structural relationships and, consequently, on the distinction of significant and insignificant similarity of protein structures. In particular, we find that side-chain orientation, as represented by the C^β atom positions, is needed for a proper defi-

nition of structurally equivalent residues. Often the C^α traces of two proteins show extensive similarity, while side-chains point in opposite directions. This is frequently observed in structures with a high content of β -sheet but less so with α -helices.

Figure 2 shows the superimposition of the ϵ -subunit of F_0F_1 -ATP synthase from *Escherichia coli* (1aqt) and the C-terminal subdomain of β -galactosidase from *E. coli* (1bgl). Based on C^α superimposition, the similarity is extensive (Figure 2(a)). Considering side-chain orientation, only one sheet of the β -sandwich can be superimposed (Figure 2(b)) and the number of equivalent residues drops from 61 (C^α) to 42 (C^α/C^β), corresponding to a difference of 31%. From the distinct function and different side-chain orientations it seems likely that the ATP synthase and β -galactosidase domains are unrelated, although the topology of the C^α traces appear similar. Therefore, we classify these domains as different folds. An attempt to derive a structural model for one of these domains using the other domain as a template will be very difficult, since the orientation of many side-chains has to be reversed.

On average, 55% of the residues of structurally similar domains are equivalent. The respective number for unrelated pairs, including similarities of substructures and other partial similarities, is 36%. The averages differ considerably, but the corresponding distributions overlap, so that the number of equivalent residues is insufficient to separate

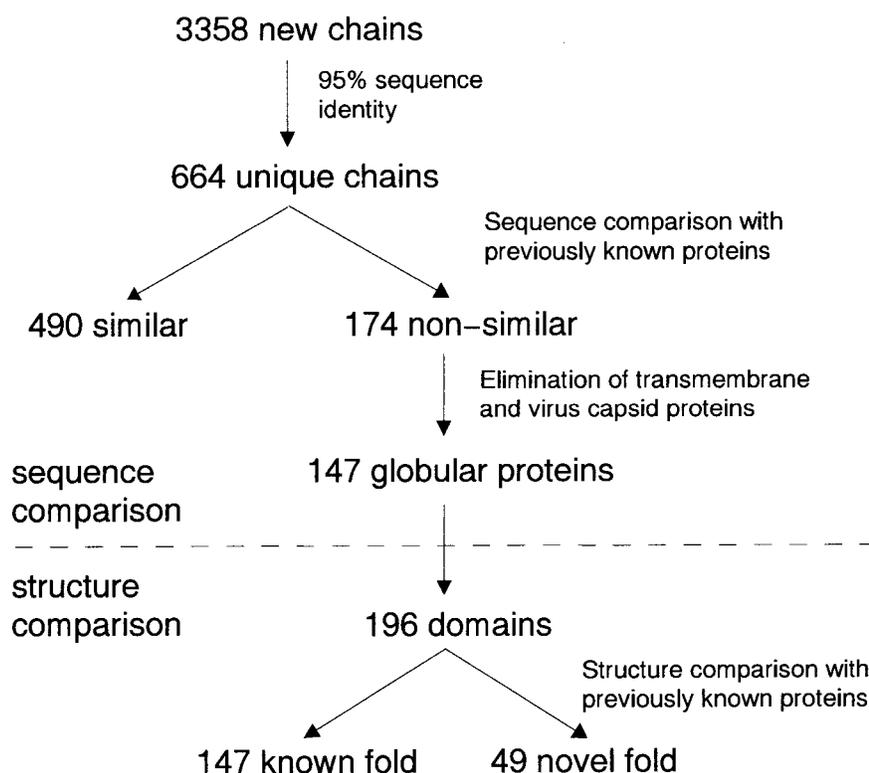


Figure 1. Preparation of data sets.

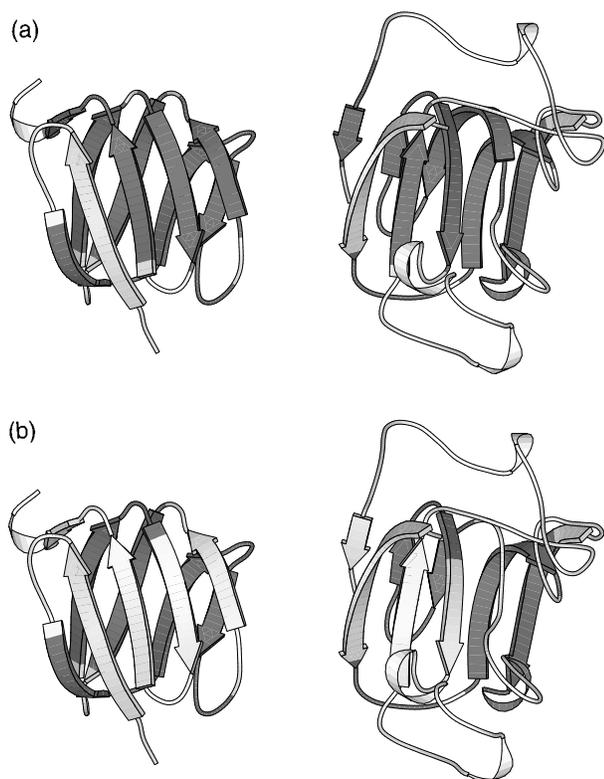


Figure 2. Comparison of the superimposition of the ϵ -subunit of F_0F_1 -ATP synthase from *E. coli* (1aqt, left) and the C-terminal subdomain of β -galactosidase from *E. coli* (1bgl, right). Structurally equivalent residues are dark gray. (a) C^α superimposition: strands of both sheets of the sandwich are equivalent. (b) C^α/C^β superimposition: only the strands of one sheet are equivalent because side-chain orientations are different. All ribbon representations were produced with the program Molscript (Kraulis, 1991).

cases of significant similarity from others. Nevertheless, many pairs having more than 50% equivalent residues, relative to the target size, share the same fold. This covers 66% of the significant similarities with an error rate of 4%. Extreme cases are lips.A and 1cau.A, two proteins having the same fold although they share only 23% equivalent residues. Another such case is 1kdx.A and 1vnc, which have different folds but 69% equivalent residues

In this analysis we observed eight pairs at the substructure level. Figure 3 depicts an example of a small protein that superimposes with a compact substructure of a larger protein. Ragweed pollen allergen (1bbg) is a small (40 amino acid residues) $\alpha + \beta$ fold and superimposes with the C-terminal end of profilin (1pne) with 31 equivalent residues (considering side-chain orientation). The respective alignment has only three small gaps. In principle, the larger protein could serve as a scaffold for modeling the smaller protein. Then, given the sequence of the smaller protein, the question is

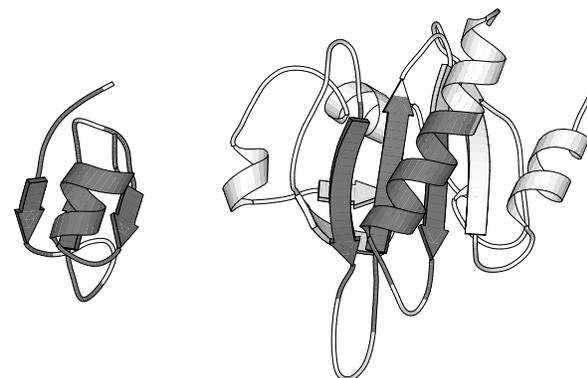


Figure 3. The small protein ragweed pollen allergen (1bbg, 40 amino acid residues) superimposed with a substructure of profilin (1pne, 140 residues). All secondary structure elements of 1bbg superimposes with the C-terminal part of 1pne (residues 84 to 129, 31 equivalent residues). Stretches of structurally equivalent residues are dark gray.

whether the structural similarity to the larger protein could be detected by prediction methods. The problem is that many structurally equivalent pairs are exposed in one protein but buried in the other structure. Therefore, methods that take the physicochemical environment into account may fail in such cases (F. S. Domingues, P.L. & M.J.S., unpublished results).

Partial fold similarity was observed in eight cases. An example is the similarity of oncogene product p14^{TCL1} (1jsg) and avidin (1avd). Both consist of an eight-stranded β -barrel but only five strands are structurally equivalent. Surprisingly, the proteins λ -exonuclease (1avq) and kinesin (2kin) are partially similar to functionally related proteins (Kovall & Matthews, 1997; Sack *et al.*, 1997).

Structural and functional similarities

The results of our classification are summarized in Table 1. Of the 196 domains, 75% have similarity on the fold level to previously known proteins. Only 49 domains of the structures released in 1998 have no clear structural similarity to a previously determined structure. The latter corresponds to 7.4% of the 664 non-redundant sequences with 92.6% being structurally similar to a previously known structure.

For comparison, Table 1 includes data based on C^α trace superimposition. Since this is less restrictive, more proteins having extensive similarity to a previously known structure are found. The effect is most pronounced for substructures, which increase from eight to 14 instances.

The SCOP data base was used to analyze functional relationships between the 196 domains and previously known structures (Table 2). However,

Table 1. Classification by C α and C α /C β superimposition

Level of similarity	C α superimposition		C α /C β superimposition	
	No.	(%)	No.	(%)
Fold	153	78.1	147	75.0
Substructure	14	7.1	8	4.1
Partial fold	9	4.6	8	4.1
Novel	20	10.2	33	16.8

The levels of similarity in the first column are explained in the text. Columns 2 to 5 indicate the absolute and relative numbers of the instances of the similarity levels, both for C α superimposition (columns 2 and 3) and C α /C β superimposition (columns 4 and 5).

only 120 of these domains are classified in the used release of SCOP (August 1998). For the remaining 76 domains, evidence regarding functional relationships was obtained from the literature. Of the 196 domains, 97 (49.5%) are functionally related to previously known proteins.

Structural similarity does not always coincide with functional relationship. On the other hand, in most cases functionally related pairs also have extensive structural similarity. Figure 4 depicts the concordance of structural and functional relationships of the domains. The functionally related domains contain three groups: for 45% the most similar structural template is functionally related, for 3% the functionally related template does not coincide with the "best" template (see below) and 1.5% are functionally related (same SCOP superfamily) but have weak structural similarity. From the functionally unrelated domains, 27% have structural similarity to known folds while 23.5%

are novel folds, substructures or partially similar folds. There are 94 protein pairs that are both structurally and functionally similar, corresponding to 64% of the structurally similar proteins.

Homologous versus analogous proteins

Typically, a protein has structural similarity to several proteins and expert knowledge is essential to distinguish homologous from analogous proteins. A measure that facilitates this distinction would be advantageous. The nine cases where an analogous protein has more equivalent residues in common with the target than a homologous protein are reduced to five when side-chain orientation is taken into account. Although not perfect, side-chain orientation helps to distinguish analogous from homologous pairs. We also encountered some instances where the consideration of side-chain orientation reveals a closer homologue (SCOP family instead of superfamily), e.g. for the RNA-binding domain of the transcriptional terminator protein ρ (1a62).

An example is G:T/U mismatch-specific DNA glycosylase from *E. coli* (1mug) (Barrett *et al.*, 1998), which is in the same SCOP superfamily as human uracil-DNA glycosylase (1akz). The C α trace of 1mug shares higher similarity with cutinase from *Fusarium solani* (1cex, 97 equivalent residues, different folds in SCOP) than with 1akz (89 residues). When side-chain orientation is taken into account, the ranking is reversed: 1mug shares 77 equivalent residues with 1akz and 73 with 1cex. A similar situation is observed for the proteins interleukin-1 receptor (1itb.B), the T-domain of the brachyury transcription factor (1xbr.A), and robustoxin (1qdp).

For five proteins, an analogue has greater similarity than the homologue also when side-chain orientation is taken into account (Table 3). We describe two examples where side-chain orientation fails to correctly discriminate homologous from analogous pairs.

The β -subunit of protein farnesyltransferase (1ft1.B) has the fold of an α - α toroid consisting of six α -helical hairpins arranged in a closed circular array. This type of fold tolerates variations in size. The functionally unrelated protein glucoamylase (1gai) has the same fold with an identical number of α -helical hairpins, while the functionally related

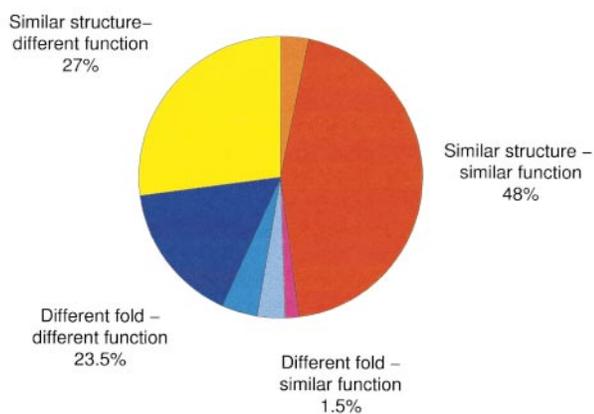


Figure 4. Coincidence of structural and functional similarities of the data set domains. Red sector, domains which have extensive structural and functional relationship to a known structure; orange sector, domains where a template with functional relationship exists but an analogous protein has highest structural similarity; yellow sector, domains having a structural template with different function; magenta sector, domains that are classified in the same SCOP superfamily without extensive structural similarity; blue sectors, domains with a novel fold, including substructures (light blue) and partial fold similarity (very light blue).

Table 2. Structural and functional relationships of the data set domains

Target	Domain	Length	Template	Eq	(%)	Level	Homologue	Target name
1a0b.-.		125	1cgo.-.	75	(60.0)	Fold		C-terminal HPT domain of ArcB
1a0i.-.	2-240	239	1ckm.A.-	126	(52.7)	Fold	Superfamily	ATP-dependent DNA ligase
1a0i.-.	241-349	109	1ah9.-.1	52	(47.7)	Fold	Superfamily	
1a17.-.		166	1sly.-.	93	(56.0)	Fold		Protein phosphatase 5
1a2z.A.-		220	1ecp.A.-	128	(58.2)	Fold		Pyrrolidone carboxyl peptidase
1a3c.-.		181	1hgx.A.-	136	(75.1)	Fold	Tomchick <i>et al.</i> (1998)	Transcriptional attenuation protein PyrR
1a5r.-.1		103	1ubi.-.	63	(61.2)	Fold	Bayer <i>et al.</i> (1998)	SUMO-1
1a5t.-.	1-167	167	2reb.-.	104	(62.3)	Fold	Superfamily	δ' subunit of DNA polymerase III
1a5t.-.	168-330	163	1wer.-.	71	(43.6)	Novel		
1a62.-.		130	1mjc.-.	49	(37.7)	Fold	Family	Transcription termination factor ρ
1a68.-.		95	1tfr.-.	39	(41.1)	Novel		Shaker potassium channel
1a6q.-.	2-296	295	1pma.1.-	83	(28.1)	Partial		Protein Ser/Thr phosphatase 2C
1a6q.-.	297-368	72	1bgw.-.	53	(73.6)	Fold		
1a74.A.-		163	1kit.-.	39	(23.9)	Novel		Homing endonuclease I-PpoI
1a7j.-.		290	2ak3.A.-	105	(36.2)	Fold	Superfamily	Phosphoribulokinase
1aa2.-.		108	1aoa.-.	96	(88.9)	Fold	Family	Calponin homology domain
1acc.-.	14-258	245	1msp.A.-	62	(25.3)	Fold		Anthrax protective antigen
1acc.-.	259-485	227	1rho.A.-	62	(27.3)	Fold		
1acc.-.	486-604	119	1gua.B.-	49	(41.2)	Fold		
1acc.-.	605-735	131	1rho.A.-	56	(42.7)	Fold		
1ad1.A.-		266	1fdy.D.-	161	(60.5)	Fold		Dihydropteroate synthetase
1ad6.-.		185	1vin.-.	88	(47.6)	Fold	Kim & Cho (1997)	Retinoblastoma tumor suppressor
1af7.-.	11-90	90	1pnk.B.-	54	(60.0)	Fold		Methyltransferase CheR
1af7.-.	91-284	194	1vid.-.	107	(55.2)	Fold	Superfamily	
1af9.-.	875-1109	235	1sac.A.-	118	(50.2)	Fold	Superfamily	Tetanus neurotoxin
1af9.-.	1110-1315	206	1tie.-.	128	(62.1)	Fold	Superfamily	
1ah1.-.1		129	1tcr.A.-	89	(69.0)	Fold	Family	Immunoreceptor CTLA-4
1ahj.A.-		207	1rom.-.	55	(26.6)	Novel		Nitrile hydratase
1ahj.B.-	1-112	112	1cpc.A.-	48	(42.9)	Novel		
1ahj.B.-	113-212	100	1pse.-.1	49	(49.0)	Fold		
1ahk.-.		129	1fie.A.-	50	(38.8)	Fold	Family	Major mite allergen Der f 2
1aiw.-.1		62	1fgp.-.1	31	(50.0)	Partial		Endoglucanase Z
1aj6.-.		219	1yer.-.	119	(54.3)	Fold	Superfamily	DNA gyrase B, N-terminal domain
1al0.1.-		152	2cae.-.	51	(33.6)	Novel		Scaffolding protein gpD
1am2.-.		199	1at0.-.	124	(62.3)	Fold	Superfamily	GyrA intein
1amx.-.		180	1anu.-.	71	(39.4)	Fold		Collagen-binding domain of adhesin
1ao6.A.-	1-190	190	1ar1.A.-	62	(32.6)	Novel		Serum albumin
1ao6.A.-	191-382	192	1rom.-.	69	(35.9)	Novel		
1ao6.A.-	383-582	200	1ddt.-.	63	(31.5)	Novel		
1aos.A.-	1-112	112	1fur.A.-	60	(53.6)	Fold	Superfamily	Argininosuccinate lyase
1aos.A.-	113-362	250	1fur.A.-	232	(92.8)	Fold	Superfamily	
1aos.A.-	363-462	100	2abk.-.	55	(55.0)	Fold		
1ap0.-.1		73	1sap.-.	42	(57.5)	Fold	Superfamily	Chromatin modifier protein 1
1ap8.-.1		213	1vao.A.-	71	(33.3)	Partial		Translation initiation factor eIF4e
1apj.-.1		74	1uby.-.	28	(37.8)	Novel		Fibrillin
1aqt.-.		138	1xsm.-.	52	(37.7)	Novel		ϵ subunit of F ₀ F ₁ -ATP synthase
1aqu.A.-		297	1gky.-.	97	(32.7)	Fold	Kakuta <i>et al.</i> (1997)	Estrogen sulfotransferase
1auk.-.		489	1alk.A.-	168	(34.4)	Fold	Superfamily	Arylsulfatase A
1auo.A.-		218	1tht.A.-	147	(67.4)	Fold	Superfamily	Carboxylesterase
1auv.A.-		311	1gsa.-.	195	(62.7)	Fold	Superfamily	Synapsin I
1auz.-.1		116	1wab.-.	70	(60.3)	Fold		Spollaa
1avq.A.-		228	1cfr.-.	64	(28.1)	Partial	Superfamily	λ -Exonuclease
1aw5.-.		326	1ad1.A.-	165	(50.6)	Fold		5-Aminolevulinate dehydratase
1aw8.B.-		91	1cxs.A.-	53	(58.2)	Fold	Superfamily	Aspartate decarboxylase
1axj.-.1		122	1hav.A.-	67	(54.9)	Fold	(Liepinsh <i>et al.</i> , (1997))	FMN-binding protein
1ay2.-.		158	1bcp.A.-	47	(29.7)	Novel		Pilin
1azo.-.		232	1pvu.A.-	67	(28.9)	Fold	Superfamily	DNA mismatch repair protein MutH
1bak.-.1		119	1btk.B.-	82	(68.9)	Fold	Family	G-protein coupled receptor kinase 2
1baq.-.1		139	1ivh.A.-	60	(43.2)	Sub.		Antitermination factor NusB
1bbg.-.		40	1pne.-.	31	(77.5)	Sub.		Ragweed allergen Amb V
1bcc.A.-	4-234	231	1bia.-.	76	(32.9)	Fold		Cytochrome <i>b</i> _{c1}
1bcc.A.-	235-445	211	1azs.B.-	69	(32.7)	Fold		
1bcc.D.-		241	1cxc.-.	69	(28.6)	Fold	Zhang <i>et al.</i> (1998)	
1bd9.A.-		187	1rlw.-.	58	(31.0)	Fold		Phosphatidylethanolamine-binding protein
1bdy.A.-		123	1rlw.-.	85	(69.1)	Fold	Pappa <i>et al.</i> (1998)	C2 domain from protein kinase C- δ
1bf5.A.-	136-317	182	1dlc.-.	103	(56.6)	Fold		STAT-1
1bf5.A.-	318-490	173	1a02.N.-	87	(50.3)	Fold	Chen <i>et al.</i> (1998)	
1bf5.A.-	491-710	220	1sha.A.-	69	(31.4)	Fold	Chen <i>et al.</i> (1998)	
1bg8.A.-		89	1ao6.A.-	56	(62.9)	Fold		HdeA
1bgf.-.		124	2brd.-.	63	(50.8)	Novel		Transcription factor STAT-4
1bhe.-.		376	1rmg.-.	274	(72.9)	Fold	Pickersgill <i>et al.</i> (1998)	Polygalacturonase
1bja.A.-		95	1smt.A.-	72	(75.8)	Fold	Finnin <i>et al.</i> (1997)	Transcription factor MotA

1bjn.A.-	1-258	258	1cl1.A.-	137	(53.1)	Fold	Hester <i>et al.</i> (1999)	Phosphoserine aminotransferase
1bjn.A.-	259-362	104	2dkb.-.	84	(80.8)	Fold	Hester <i>et al.</i> (1999)	
1bkb.-.	4-75	72	1vie.-.	46	(63.9)	Fold		Transcription initiation factor 5A
1bkb.-.	76-139	64	1ah9.-.1	55	(85.9)	Fold	Peat <i>et al.</i> (1998)	
1bnl.A.-		178	1ixx.A.-	56	(31.5)	Novel		Endostatin
1bqv.-.1		110	1rlr.-.	51	(46.4)	Fold		Transcription factor Ets-1
1brz.-.1		54	2sn3.-.	39	(72.2)	Fold	Superfamily	Brazzein
1c25.-.		161	1rhs.-.	84	(52.2)	Fold		Phosphatase CDC25A
1c3d.-.		294	1sqc.-.	188	(63.9)	Fold		C3d
1cl1.A.-	1-256	256	1ajs.A.-	155	(60.5)	Fold	Clausen <i>et al.</i> (1996)	Cystathionine β -lyase
1cl1.A.-	257-395	139	2dkb.-.	80	(57.6)	Fold	Clausen <i>et al.</i> (1996)	
1crx.A.-	20-130	111	1a0p.-.	65	(58.6)	Fold	Guo <i>et al.</i> (1997)	Cre recombinase
1crx.A.-	131-341	211	1ae9.A.-	121	(57.3)	Fold	Guo <i>et al.</i> (1997)	
1cv8.-.	2-79	78	1ppn.-.	52	(66.7)	Fold	Hofmann <i>et al.</i> (1993)	Staphopain
1cv8.-.	80-174	95	1the.A.-	65	(68.4)	Fold	Hofmann <i>et al.</i> (1993)	
1d2n.A.-	505-676	172	1a5t.-.	120	(69.8)	Fold	Lenzen <i>et al.</i> (1998)	N-ethylmaleimide-sensitive factor
1d2n.A.-	677-750	74	1maz.-.	48	(64.9)	Sub.		
1dfx.-.		125	1f13.A.-	59	(47.2)	Fold		Desulfoferrodoxin
1dhs.-.		344	2mas.A.-	98	(28.5)	Fold	(Liao <i>et al.</i> (1998))	Deoxyhypusine synthase
1dps.A.-		167	1aew.-.	128	(76.6)	Fold	Grant <i>et al.</i> (1998)	DNA-binding protein Dps
1e2a.A.-		105	1aj3.-.1	74	(70.5)	Fold		Enzyme IIA ^{lactose}
1fgj.A.-		546	1occ.A.-	92	(16.8)	Novel		Hydroxylamine oxidoreductase
1fsz.-.	23-231	209	1hdc.D.-	119	(56.9)	Fold		Cell-division protein FtsZ
1fsz.-.	232-356	125	1com.K.-	67	(53.6)	Fold		
1ft1.A.-		377	1lrv.-.	80	(21.2)	Fold		Protein farnesyltransferase
1ft1.B.-		437	1gai.-.	152	(34.8)	Fold	(Park <i>et al.</i> (1997))	
1ftr.A.-	1-17, 163-296	151	1cg2.A.-	65	(43.0)	Fold		Formyltransferase
1ftr.A.-	18-162	145	1cg2.A.-	70	(48.3)	Fold		
1g31.A.-		111	1aon.O.-	61	(55.0)	Fold	Hunt <i>et al.</i> (1997)	Co-chaperonin Gp31
1gc1.G.-		321	1ord.A.-	47	(14.6)	Novel		Envelope protein gp120
1hei.A.-	187-323	136	1pjr.-.	88	(64.7)	Fold	Yao <i>et al.</i> (1997)	RNA helicase
1hei.A.-	324-485	161	1pjr.-.	70	(43.5)	Fold	Yao <i>et al.</i> (1997)	
1hei.A.-	486-629	144	1ah7.-.	50	(34.7)	Novel		
1hkb.A.-	16-476	461	1gla.G.-	175	(38.0)	Fold	Superfamily	D-Glucose 6-phosphotransferase
1hkb.A.-	477-914	438	1gla.G.-	174	(39.7)	Fold	Superfamily	
1hp8.-.		68	1sly.-.	47	(69.1)	Sub.		p8-MTSCP1
1hus.-.		155	1rlr.-.	52	(33.5)	Novel		Ribosomal protein S7
1if1.A.-		113	1pue.E.-	67	(59.3)	Fold	Superfamily	Interferon regulatory factor 1
1ihp.-.		438	1rpa.-.	202	(46.1)	Fold	Superfamily	Phytase
1ips.A.-		331	1cau.A.-	76	(23.0)	Fold		Isopenicillin N synthase
1itb.B.-	1-100	100	2ncm.-.1	71	(71.0)	Fold	Family	Type-1 interleukin-1 receptor
1itb.B.-	101-204	104	2ncm.-.1	73	(70.2)	Fold	Family	
1itb.B.-	205-315	111	1tlk.-.	91	(82.0)	Fold	Family	
1jdw.-.		423	1esc.-.	57	(13.5)	Novel		L-Arg:Gly amidinotransferase
1jfr.A.-		262	1bro.A.-	160	(61.1)	Fold	Superfamily	Lipase
1jsg.-.		114	1avd.A.-	46	(40.4)	Partial		p14 ^{TCL1}
1kdx.A.1		81	1vnc.-.	52	(64.2)	Sub.		KIX domain of CBP
1kwa.A.-		88	1pdr.-.	76	(86.4)	Fold	Family	hCASK-2 PDZ domain
1ltm.-.	42-264	223	154l.-.	74	(33.2)	Fold	van Asselt <i>et al.</i> (1998)	Transglycosylase
1ltm.-.	265-361	97	1csh.-.	44	(45.4)	Novel		
1lxd.A.-		100	1gua.-.	65	(65.0)	Fold	Family	Ras-interacting domain of RaiGDS
1mb1.-.		130	1apm.E.-	46	(35.4)	Fold		Transcription factor Mbp1
1moq.-.	243-450	208	1fsz.-.	96	(46.2)	Fold		Glucosamine-6-phosphate synthase
1moq.-.	451-608	158	1bmt.A.-	82	(51.9)	Fold		
1mpg.A.-	1-99	99	1ytb.A.-	56	(56.6)	Fold	Superfamily	3-Methyladenine DNA glycosylase II
1mpg.A.-	100-282	183	2abk.-.	130	(71.0)	Fold	Superfamily	
1mro.A.-	102-276	175	1bdp.-.	57	(32.6)	Novel		Methyl-coenzyme M reductase
1mro.A.-	277-549	273	1rom.-.	76	(27.8)	Novel		
1mro.B.-	1-43, 181-442	305	1aj8.A.-	69	(22.6)	Novel		
1mro.B.-	44-180	137	4mp.A.-	64	(46.7)	Sub.		
1mro.C.-		247	1vhi.A.-	59	(23.9)	Novel		
1mrp.-.	1-101, 226-277	153	1dmb.-.	123	(80.4)	Fold	Family	Ferric iron binding protein
1mrp.-.	102-235, 278-309	166	1sbp.-.	110	(66.3)	Fold	Family	
1mug.A.-		168	1akz.-.	77	(45.8)	Fold	Superfamily	G:T/U-specific DNA glycosylase
1ned.A.-		183	1pma.1.-	156	(85.2)	Fold	Family	HslV (ClpQ) protease
1nkr.-.	6-101	96	1zxq.-.	73	(76.0)	Fold	Fan <i>et al.</i> (1997)	Inhibitory receptor p58-cl42
1nkr.-.	102-200	99	1ac6.A.-	76	(76.8)	Fold	Fan <i>et al.</i> (1997)	
1nlr.-.		234	1xyo.A.-	133	(56.8)	Fold	Sulzenbacher <i>et al.</i> (1997)	Endo- β -1,4-glucanase
1noc.A.-		388	2min.B.-	54	(13.9)	Novel		Inducible nitric oxide synthase
1np1.A.-		184	1aqb.-.	124	(67.4)	Fold	Family	Nitrophorin 1
1pfo.-.	30-388	359	1auk.-.	55	(15.3)	Novel		Perfringolysin O
1pfo.-.	389-500	112	1tsh.A.-	52	(46.4)	Fold		

1phm.-.	200-354	155	1ahs.A.-	65	(41.9)	Fold		Peptidylglycine α -hydroxylating monooxygenase
1phm.-.	45-199	155	1cwp.A.-	65	(41.9)	Fold		
1pin.A.-		163	1fkb.-.	41	(25.2)	Fold	Ranganathan <i>et al.</i> (1997)	Peptidyl-prolyl <i>cis-trans</i> isomerase
1poi.A.-		317	1dea.B.-	75	(23.7)	Fold		Glutaconate CoA-transferase
1poi.B.-		260	1qor.A.-	83	(31.9)	Fold		
1prx.A.-		224	1gp1.A.-	111	(49.6)	Fold	Superfamily	HorF6 peroxidase
1ps1.A.-		337	5eas.-.	215	(63.8)	Fold	Superfamily	Pentalene synthase
1qdp.-1		42	1agg.-1	30	(71.4)	Fold	Family	Robustoxin
1rdr.-.	67-367	301	1mml.-.	65	(21.6)	Fold	Hansen <i>et al.</i> (1997)	RNA polymerase
1rdr.-.	12-37, 368-461	120	2lbd.-.	50	(41.7)	Novel		
1rkd.-.		309	1nah.-.	100	(32.4)	Partial		Ribokinase
1rmg.-.		422	1tyu.-.	216	(51.2)	Fold		Rhamnogalacturonase A
1ryp.2.-		233	1pma.1.-	184	(79.0)	Fold	Family	Proteasome
1sfp.-.		114	2bpa.2.-	70	(61.4)	Fold		Acidic seminal protein
1shk.A.-		173	1ukz.-.	108	(62.4)	Fold	Krell <i>et al.</i> (1996)	Shikimate kinase
1skn.P.-		92	1cnt.1.-	50	(54.3)	Novel		Transcription factor Skn-1
1tdj.-.	5-335	331	2tys.B.-	247	(74.6)	Fold	Gallagher <i>et al.</i> (1998)	Threonine deaminase
1tdj.-.	336-497	162	1psd.A.-	63	(38.9)	Fold		
1tmk.A.-		216	1kin.A.-	132	(61.1)	Fold	Family	Thymidylate kinase
1toh.-.		343	1jet.A.-	60	(17.5)	Novel		Tyrosine hydroxylase
1tub.A.-	1-205	205	1fsz.-.	127	(62.0)	Fold	Nogales <i>et al.</i> (1998)	Tubulin
1tub.A.-	206-381	176	1fsz.-.	109	(61.9)	Fold	Nogales <i>et al.</i> (1998)	
1tyf.A.-		193	1nzy.A.-	131	(67.9)	Fold	Superfamily	Clp protease
1uag.-.	1-93	93	1mio.B.-	75	(80.6)	Fold		UDP-N-acetylmuramoyl-L-Ala:D-Glu ligase
1uag.-.	94-298	205	1rkd.-.	105	(51.2)	Fold		
1uag.-.	299-437	139	1iso.-.	74	(53.2)	Fold		
1uch.-.		230	1the.A.-	55	(23.9)	Novel	Superfamily	Ubiquitin C-terminal hydrolase
1uea.B.-	1-107	107	1tii.D.-	54	(50.5)	Fold		Metalloproteinase inhibitor 1
1uea.B.-	108-181	74	3grs.-.	36	(48.6)	Novel		
1uox.-.		295	1gtq.A.-	84	(28.5)	Fold		Urate oxidase
1vde.A.-	1-180, 416-454	219	1at0.-.	117	(53.4)	Fold	Superfamily	Homing endonuclease PI-SceI
1vde.A.-	181-298	118	1af5.-.	66	(55.9)	Fold	Superfamily	
1vde.A.-	299-415	117	1af5.-.	69	(59.0)	Fold	Superfamily	
1vgh.-1		55	2mbr.-.	32	(58.2)	Novel		Heparin-binding domain
1vub.A.-		101	1vie.-.	38	(37.6)	Fold		Topoisomerase poison CcdB
1wja.A.-		47	2cgp.A.-	35	(74.5)	Fold		HIV-1 integrase
1xat.-.	3-163	161	1tdt.A.-	89	(55.3)	Fold	Superfamily	Xenobiotic acetyltransferase
1xat.-.	164-210	47	1dkx.A.-	42	(89.4)	Fold		
1xbr.A.-		184	1svc.P.-	71	(38.6)	Fold	Superfamily	Brachyury transcription factor
1ygs.-.		234	1a25.A.-	52	(22.2)	Partial		Smad4 tumor suppressor
1yub.-.	1-181	181	1vpt.-.	103	(56.9)	Fold	Family	rRNA methyltransferase
1yub.-.	182-245	64	1ddf.-.	42	(65.6)	Fold		
2ezk.-.		99	1jhg.A.-	49	(49.5)	Fold	(Schumacher <i>et al.</i> (1997))	I β domain of Mu transposase
2fmr.-1		65	2pii.-.	49	(75.4)	Fold	(Musco <i>et al.</i> (1997))	Fragile X protein
2hfh.-1		109	1hst.A.-	51	(46.8)	Fold	Superfamily	Genesis
2hgf.-.		97	1gup.A.-	46	(47.4)	Sub.		Hepatocyte growth factor
2kin.A.-		238	1vom.-.	90	(37.8)	Partial	Family	Kinesin
2kin.B.-		100	1aa6.-.	43	(43.0)		Substructure	
2pth.-.		193	1ecp.A.-	110	(57.0)	Fold		Peptidyl-tRNA hydrolase
2sak.-.		121	2qil.A.-	51	(42.1)	Fold		Staphylokinase
4rnp.A.-	7-294	294	2dtr.-.	61	(20.7)	Novel		RNA polymerase
4rnp.A.-	295-883	589	1bdp.-.	204	(34.6)	Fold	Superfamily	

Structural and functional relationships of all data set domains. The columns indicate the target identifier (PDB code with chain and model identifier) with eventual domain boundaries, the length of the domain, the name of the "best" template, the absolute number of equivalent residues and relative to the target domain length, the assigned level of structural similarity (fold, partial, substructure, novel; see the text for explanation), the functional relationship of target and template, and the target name. Homologue specifies the level and source of evidence of functional relationship. This could be SCOP superfamily or family or information obtained from the structure publications. A reference in parentheses means that a homologue exists but an analogue with higher similarity is the best template.

protein 5-epi-aristolochene (5eas) consists of only four α -helical hairpins. This is an extreme example where a functionally unrelated protein has extensive structure similarity, whereas a functionally related protein is much smaller and, therefore, has

fewer equivalent residues (152 with 1gai *versus* 105 with 5eas).

The second example concerns the FMN-binding protein from *Desulfovibrio vulgaris* (1axj), which forms a small β -barrel and binds the cofactor flavin

mononucleotide (FMN). This protein has high similarity (67 equivalent residues) to the C-terminal domain of hepatitis A virus 3C proteinase (1hav.A) (Figure 5), which is an analogous relationship (Liepinsh *et al.*, 1997). The structural relationship of 1axj with the N-terminal domain of phthalate dioxygenase reductase from *Pseudomonas cepacia* (2pia) is a typical example of partial fold similarity (Liepinsh *et al.*, 1998). 1axj shares only 41 equivalent residues with 2pia, which binds FMN in the same region as 1axj (Figure 6(a)). Both proteins are in the same SCOP ferredoxin reductase-like superfamily. The sequence of FMN-binding protein is circularly permuted relative to 2pia so that the N-terminal β -hairpin of 2pia corresponds to the C-terminal hairpin in 1axj (Figure 6(b)). The superimposition algorithm preserves sequence order so that these substructures cannot be superimposed.

Homologous proteins having weak structural similarity

There are a few examples of proteins having related function but rather dissimilar structures. Such similarities can have several origins. An example that most likely diverged by extensive permutations in the sequence is human deubiquitinating enzyme UCH-L3 (1uch). This protein belongs to the SCOP superfamily of cysteine proteinases. Its catalytic triad superimposes quite well with that of other members of this superfamily (Johnston *et al.*, 1997). However, the sequential order of the catalytic residues is different in deubiquitinating enzyme and cysteine proteinases of the papain family. The structural similarity is confined to a five-stranded β -sheet, while the flanking helices are not superimposable. The main difference in topology is that the helix that accommodates the catalytic cysteine residue is located between strands 2 and 3 in deubiquitinating enzyme and at the N terminus in papain-like cysteine proteinases. This is an example where the functional relatedness of two proteins, which most

likely evolved from a common ancestor, is most difficult to detect by automated, structure-based methods.

Conclusion

Here, we investigated the use of protein structure in the structural and functional characterization of protein sequences. We used the 3358 protein chains released in 1998 by the PDB to simulate a situation where a set of sequences has to be annotated. A large part of these entries is redundant, in the sense that the sequences have a high percentage of identities to previously known structures. Removing the sequence redundancies and all non-globular structures, 147 proteins corresponding to 196 domains remain. This set of 196 domains is most interesting, since it contains the information that is not accessible by sequence-based methods. The goal of this analysis was to document the information that could be gained from these proteins if proper search tools are available.

The analysis reveals that 75% of these 196 domains have extensive similarity to previously determined structures. These 75% correspond to the maximum amount of structural information that could be retrieved from a data base if a perfect technique is available. Most likely this amount will increase with the growing number of available protein structures. Two-thirds of structurally similar proteins are also functionally related. With respect to the 196 domains, this corresponds to 50% of functional coincidences.

In a recent study, Orengo *et al.* (1999) obtained a much higher percentage, 83%, of functional coincidences. Our analysis differs in several aspects. Their cutoff for significant sequence similarity is 30% sequence identity. However, as was discussed by Brenner *et al.* (1998), the use of extreme value statistics reveals significant similarities way below this threshold. Therefore, their estimate of 83% contains cases that are considered to be homologous by these techniques. Since we used FastaA E-values in this analysis, the number of homologous proteins appear to be significantly smaller.

The question of what extent current techniques are able to use the available information is most interesting. A reasonable estimate requires a suitable benchmark. Fortunately, the recent CASP experiment (Koehl & Levitt, 1999; Sippl, 1999) provides an invaluable estimate of the current state of the art in structure prediction. In particular, the results obtained in the fold recognition category are relevant in view of the current analysis. These techniques predict structures using fold data bases. In the last CASP experiment there were 23 targets in the fold recognition category, i.e. sequences having no significant sequence similarity to proteins of known structure. Hence, the CASP experiment yields an estimate for the amount of structural

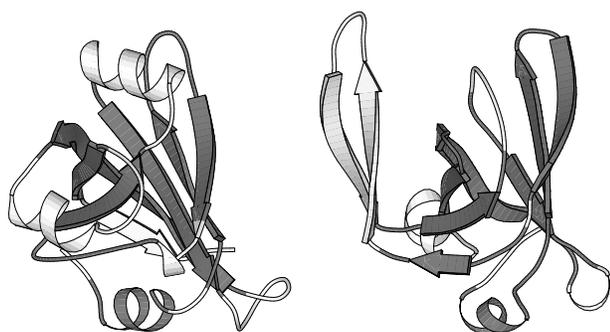


Figure 5. Superimposition of FMN-binding protein from *Desulfovibrio vulgaris* (1axj) with the C-terminal domain of hepatitis A virus 3C proteinase (1hav.A). All strands of the core have structurally equivalent counterparts (67 equivalent residues, dark gray).

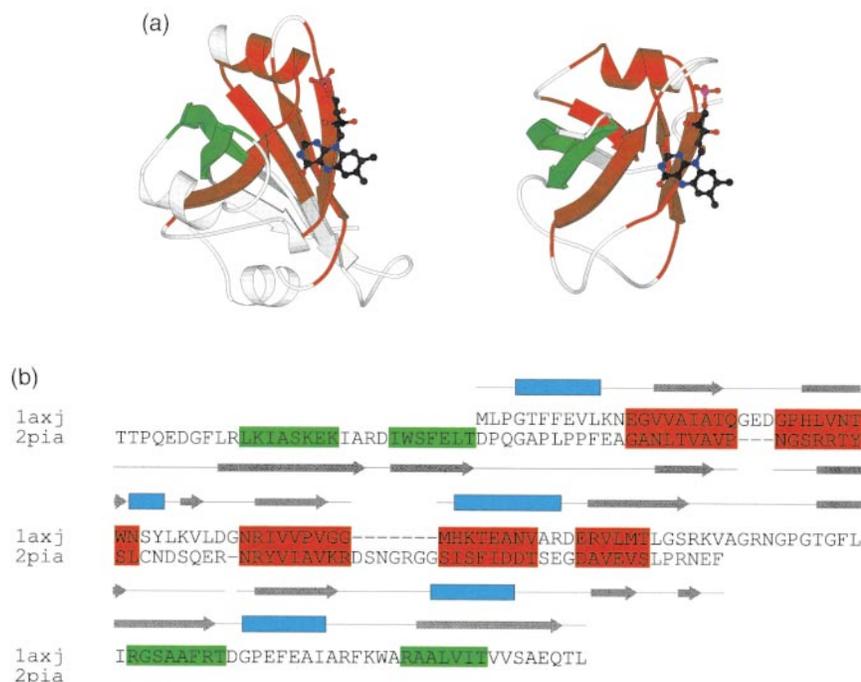


Figure 6. (a) Superimposition of FMN-binding protein from *Desulfovibrio vulgaris* (1axj) with the N-terminal domain of phthalate dioxygenase reductase from *Pseudomonas cepacia* (2pia). Both structures are ligated with the cofactor FMN in a similar orientation. Four of the six strands are shared by both structures (41 equivalent residues, colored red). The green colored hairpins are in the same structural position but not structurally equivalent because the sequences are circularly permuted. (b) Alignment of the structure superimposition. Helices (blue boxes), strands (gray arrows).

information that can be obtained from prediction methods that employ fold data bases.

Twenty of the 23 fold recognition targets are similar to previously known structures and 48% of these are remote homologues (Murzin, 1999), a number that is similar to the 50% homologues obtained here. Therefore, the prediction success

observed in CASP3 can be used to estimate how many of the structural coincidences observed for the 1998 proteins could be obtained from current structure prediction methods. The success rate is in the order of 30 to 40% (Murzin, 1999). Hence, there is an enormous potential to increase the amount of information that can be used for the

Table 3. Highly similar analogous structure pairs

Target name	PDB-Id	Length	Template name	PDB-ID	Equiv
FMN-binding protein (Liepinsh <i>et al.</i> , 1997)	1axj	122	Hepatitis A virus 3C proteinase	1hav.A	67
			Phthalate dioxygenase reductase	2pia	41
Deoxyhypusine synthase (Liao <i>et al.</i> , 1998)	1dhs	344	Purine nucleoside hydrolase	2mas	98
			Pyruvate decarboxylase	1pvd	89
Protein farnesyltransferase (Park <i>et al.</i> , 1997)	1ft1.B	437	Glucoamylase	1gai	152
			5-Epi-aristolochene synthase	5eas	105
Transposase (Schumacher <i>et al.</i> , 1997)	2ezk	99	Trp repressor	1jhg.A	49
			TC3 transposase	1tc3.C	43
KH domain of FMR1 (Musco <i>et al.</i> , 1997)	2fmr	65	Signal transduction protein PII	2pii	49
			Vigilin	1vig	44

Proteins that have an analogous structure with higher similarity than a potential homologue. The analogue is always specified before the homologue. In three cases (1dhs, 2ezk, and 2fmr) the difference in equivalent residues is only marginal ($\approx 10\%$), while in the other cases the analogues have a substantially higher number of equivalent residues (see the text for a discussion).

structural and functional characterization of novel proteins by using and improving these techniques.

These results are relevant in light of the structural genomics initiatives whose goal is to determine at least one member of all biologically relevant protein families. On the one hand, the structure data bases created in these projects will enable template-based prediction methods to compute approximate structures for virtually every protein of interest, at least in principle. On the other hand, the prediction methods will be an asset for structure determination, since they are instrumental in identifying proteins that most likely have a novel fold.

Large-scale annotations of whole proteomes indicate that currently sufficiently accurate models can be built for 8% to 17% of the sequences based on sequence comparison and comparative modeling (Andrade *et al.*, 1999). This limit can be further extended when all remote sequence-structure relationships are exploited. To estimate this range, we assume that the fold types in whole proteomes are distributed like the domains in PDB. Then the fraction of transmembrane and other non-globular proteins is approximately 20% (Frishman & Mewes, 1999). The estimated number of proteins unrelated in sequence but having significant similarity to a known structure is $(100 - 17 - 20) \times 0.75 = 47\%$. Two-thirds of these will be functionally related according to the results obtained here (Figure 7).

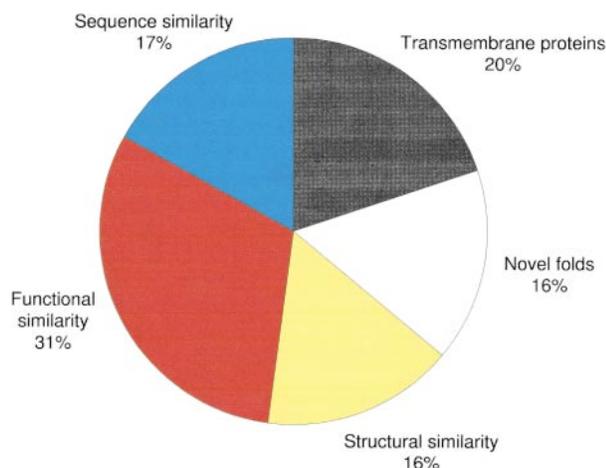


Figure 7. Current limits of structural and functional inference for proteomes inferred from structure data bases. Blue sector, the fraction of structures that can be inferred from pairwise sequence comparison (17%); gray sector, corresponds to the estimated percentage of transmembrane and non-globular proteins (20%). White sector, 16% of proteomes are most probably novel folds. Red sector, the fraction of proteins that are both structurally and functionally similar to an already determined protein structure (31%). Yellow sector, proteins with structural similarity to a known structure (16%).

Materials and Methods

Preparation of dataset

The data sets for the presented analysis were extracted from the structures released by the PDB in 1998, where we regarded the time-stamp of the file at the PDB server as the release date. (This is the date where a structure becomes available for public use. In some cases it might happen that there is a long delay between publication and release.) These sequences were filtered with a threshold of 95% sequence identity to remove identical and highly similar sequences (mutants) to derive a set of unique sequences. This is necessary because PDB files may contain several instances of the same sequence and because experimentalist often deposit slightly different structures of the same protein, like ligated and non-ligated forms. These unique sequences were compared with the sequences in the PDB at the end of 1997 using FastA (Pearson, 1996, 1998) and all sequences with significant similarity (E -value < 0.01) to an existing protein structure were removed from the set. From the resulting set, all chains with non-compact structure, transmembrane proteins and virus capsid proteins were also eliminated. All structures of the final data set were split into domains according to information from SCOP (Murzin *et al.*, 1995), CATH (Orengo *et al.*, 1997) or publications by the experimentalists (Table 2).

Structure comparison

All structure comparisons were performed with the program ProSup, which implements rigid-body superimposition of two proteins (Feng & Sippl, 1996). Parameters were set as described by P.L., W.A.K., M.J.S. & F.S. Domingues (unpublished). ProSup measures the extent of similarity by the number of structurally equivalent residues. Two residues are considered as equivalent if their C^α atoms are closer than 5 Å after superimposition. The rms deviation of equivalent residues is held approximately constant by ProSup in the range of 2 to 3 Å. Side-chain orientation is represented by the C^β atom positions. When side-chain orientation is incorporated into structure superimposition, both C^α and C^β atoms have to be closer than 5 Å after optimal superimposition of the C^α trace. One feature of ProSup is the ability to generate a list of alternative alignments. For the present analysis, only the alignments with the highest number of equivalent residues are considered for evaluation.

Each domain of the data set was compared with the structures of the PDB at the time of its release. To save computing time and to avoid too many structure libraries, we chose the following strategy: three snapshots (January 27, May 15 and September 9) of the PDB with less than 40% sequence identity were extracted and the domains of the data set were compared with the entries of the current snapshot at the time of release. The results of the ProSup data base searches (a list of protein chains sorted by the number of equivalent residues) were inspected by eye. The selected best template was normally the chain with highest number of equivalent residues. Only in cases where another chain gave a more reasonable alignment (shorter gaps, more compact in three dimensions, etc.) was this one considered as the best template.

Acknowledgments

We are grateful to Alessandro Monge and Francisco S. Domingues for valuable suggestions to improve the manuscript. This work was supported by grants P11601-GEN, P11205-MOB and P13710-MOB of the Austrian Fonds zur Förderung der wissenschaftlichen Forschung.

References

- Altschul, S. F., Gisch, W., Miller, W., Myers, E. W. & Lipman, D. J. (1990). Basic local alignment search tool. *J. Mol. Biol.* **215**, 403-410.
- Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang, J., Zhang, Z., Miller, W. & Lipman, D. J. (1997). Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucl. Acids Res.* **25**, 3389-3402.
- Andrade, M. A., Brown, N. P., Leroy, C., Hoersch, S., de Daruvar, A., Reich, C., Franchini, A., Tamames, J., Valencia, A., Ouzounis, C. & Sander, C. (1999). Automated genome sequence analysis and annotation. *Bioinformatics*, **15**, 391-412.
- Barrett, T. E., Savva, R., Panayotou, G., Barlow, T., Brown, T., Jiricny, J. & Pearl, L. H. (1998). Crystal structure of a G:T/U mismatch-specific DNA glycosylase: mismatch recognition by complementary-strand interactions. *Cell*, **92**, 117-129.
- Bayer, P., Arndt, A., Metzger, S., Mahajan, R., Melchior, F., Jaenicke, R. & Becker, J. (1998). Structure determination of the small ubiquitin-related modifier SUMO-1. *J. Mol. Biol.* **280**, 275-286.
- Bernstein, F. C., Koetzle, T. F., Williams, G. J. B., Meyer, E. F., Jr, Brice, M. D., Rogers, J. R., Kennard, O., Shimanouchi, T. & Tasumi, M. (1977). The Protein Data Bank: a computer-based archival file for molecular structures. *J. Mol. Biol.* **112**, 535-542.
- Bowie, J. U., Lüthy, R. & Eisenberg, D. (1991). A method to identify protein sequences that fold into a known three-dimensional structure. *Science*, **253**, 164-170.
- Brenner, S. E., Chothia, C. & Hubbard, T. J. (1998). Assessing sequence comparison methods with reliable structurally identified distant evolutionary relationships. *Proc. Natl Acad. Sci. USA*, **95**, 6073-6078.
- Bryant, S. H. (1996). Evaluation of threading specificity and accuracy. *Proteins: Struct. Funct. Genet.* **26**, 172-185.
- Chen, X., Vinkemeier, U., Zhao, Y., Jeruzalmi, D., Darnell, J. E., Jr & Kuriyan, J. (1998). Crystal structure of a tyrosine phosphorylated STAT-1 dimer bound to DNA. *Cell*, **93**, 827-839.
- Chothia, C. (1992). One thousand families for the molecular biologist. *Nature*, **357**, 543-544.
- Clausen, T., Huber, R., Laber, B., Pohlenz, H. D. & Messerschmidt, A. (1996). Crystal structure of the pyridoxal-5'-phosphate dependent cystathionine β -lyase from *Escherichia coli* at 1.83 Å. *J. Mol. Biol.* **262**, 202-224.
- Domingues, F. S., Koppensteiner, W. A., Jaritz, M., Prlic, A., Weichenberger, C., Wiederstein, M., Flöckner, H., Lackner, P. & Sippl, M. J. (1999). Sustained performance of knowledge-based potentials in fold recognition. *Proteins: Struct. Funct. Genet. Suppl.* **3**, 112-120.
- Fan, Q. R., Mosyak, L., Winter, C. C., Wagtmann, N., Long, E. O. & Wiley, D. C. (1997). Structure of the inhibitory receptor for human natural killer cells resembles haematopoietic receptors. *Nature*, **389**, 96-100.
- Feng, Z.-K. & Sippl, M. J. (1996). Optimum superimposition of protein structures: ambiguities and implications. *Fold. Des.* **1**, 123-132.
- Finnin, M. S., Cicero, M. P., Davies, C., Porter, S. J., White, S. W. & Kreuzer, K. N. (1997). The activation domain of the MotA transcription factor from bacteriophage T4. *EMBO J.* **16**, 1992-2003.
- Flores, T. P., Orengo, C. A., Moss, D. S. & Thornton, J. M. (1993). Comparison of conformational characteristics in structurally similar protein pairs. *Protein Sci.* **2**, 1811-1826.
- Frishman, D. & Mewes, H.-W. (1999). Genome-based structural biology. *Prog. Biophys. Mol. Biol.* **72**, 1-17.
- Gallagher, D. T., Gilliland, G. L., Xiao, G., Zondlo, J., Fisher, K. E., Chinchilla, D. & Eisenstein, E. (1998). Structure and control of pyridoxal phosphate dependent allosteric threonine deaminase. *Structure*, **6**, 465-475.
- Grant, R. A., Filman, D. J., Finkel, S. E., Kolter, R. & Hogle, J. M. (1998). The crystal structure of Dps, a ferritin homolog that binds and protects DNA. *Nature Struct. Biol.* **5**, 294-303.
- Guo, F., Gopaul, D. N. & van Duyne, G. D. (1997). Structure of Cre recombinase complexed with DNA in a site-specific recombination synapse. *Nature*, **389**, 40-46.
- Hansen, J. L., Long, A. M. & Schultz, S. C. (1997). Structure of the RNA-dependent RNA polymerase of poliovirus. *Structure*, **15**, 1109-1122.
- Hegyvi, H. & Gerstein, M. (1999). The relationship between protein structure and function: a comprehensive survey with application to the yeast genome. *J. Mol. Biol.* **288**, 147-164.
- Hester, G., Stark, W., Moser, M., Kallen, J., Markovic-Housley, Z. & Jansonius, J. N. (1999). Crystal structure of phosphoserine aminotransferase from *Escherichia coli* at 2.3 Å resolution: comparison of the unligated enzyme and a complex with α -methyl-L-glutamate. *J. Mol. Biol.* **286**, 829-850.
- Hofmann, B., Schomburg, D. & Hecht, H. J. (1993). Crystal structure of a thiol proteinase from *Staphylococcus aureus* V-8 in the E-64 inhibitor complex. *Acta Crystallog. sect. A*, **49**, 102.
- Holm, L. & Sander, C. (1993). Protein structure comparison by alignment of distance matrices. *J. Mol. Biol.* **233**, 123-138.
- Holm, L. & Sander, C. (1997). New structure - novel fold? *Structure*, **5**, 165-171.
- Hunt, J. F., van der Vies, S. M., Henry, L. & Deisenhofer, J. (1997). Structural adaptations in the specialized bacteriophage T4 co-chaperonin Gp31 expand the size of the Anfinsen cage. *Cell*, **90**, 361-370.
- Johnston, S. C., Larsen, C. N., Cook, W. J., Wilkinson, K. D. & Hill, C. P. (1997). Crystal structure of a deubiquitinating enzyme (human UCH-L3) at 1.8 Å resolution. *EMBO J.* **16**, 3787-3796.
- Jones, D. T. (1997). Progress in protein structure prediction. *Curr. Opin. Struct. Biol.* **7**, 377-387.
- Jones, D. T. (1999). GenTHREADER: an efficient and reliable protein fold recognition method for genomic sequences. *J. Mol. Biol.* **287**, 797-815.
- Jones, D. T., Taylor, W. R. & Thornton, J. M. (1992). A new approach to protein fold recognition. *Nature*, **358**, 86-89.

- Kakuta, Y., Pedersen, L. G., Carter, C. W., Negishi, M. & Pedersen, L. (1997). Crystal structure of estrogen sulphotransferase. *Nature Struct. Biol.* **4**, 904-908.
- Karplus, K., Barrett, C. & Hughey, R. (1998). Hidden Markov models for detecting remote protein homologies. *Bioinformatics*, **14**, 846-856.
- Kim, H. Y. & Cho, Y. (1997). Structural similarity between the pocket region of retinoblastoma tumour suppressor and the cyclin-box. *Nature Struct. Biol.* **4**, 390-395.
- Koehl, P. & Levitt, M. (1999). A brighter future for protein structure prediction. *Nature Struct. Biol.* **6**, 108-111.
- Kovall, R. & Matthews, B. W. (1997). Toroidal structure of λ -exonuclease. *Science*, **277**, 1824-1827.
- Kraulis, P. J. (1991). Molscript: a program to produce both detailed and schematic plots of protein structures. *J. Appl. Crystallog.* **24**, 946-950.
- Krell, T., Horsburgh, M. J., Cooper, A., Kelly, S. M. & Coggins, J. R. (1996). Localization of the active site of type II dehydroquinases. Identification of a common arginine-containing motif in the two classes of dehydroquinases. *J. Biol. Chem.* **271**, 24492-24497.
- Lenzen, C. U., Steinmann, D., Whiteheart, S. W. & Weis, W. I. (1998). Crystal structure of the hexamerization domain of N-ethylmaleimide-sensitive fusion protein. *Cell*, **94**, 525-536.
- Liao, D. I., Wolff, E. C., Park, M. H. & Davies, D. R. (1998). Crystal structure of the NAD complex of human deoxyhypusine synthase: an enzyme with a ball-and-chain mechanism for blocking the active site. *Structure*, **6**, 23-32.
- Liepinsh, E., Kitamura, M., Murakami, T., Nakaya, T. & Otting, G. (1997). Pathway of chymotrypsin evolution suggested by the structure of the FMN-binding protein from *Desulfovibrio vulgaris*. *Nature Struct. Biol.* **4**, 975-979.
- Liepinsh, E., Kitamura, M., Murakami, T., Nakaya, T. & Otting, G. (1998). Common ancestor of serine proteases and flavin-binding domains. *Nature Struct. Biol.* **5**, 102-103.
- Martin, A. C., Orengo, C. A., Hutchinson, E. G., Jones, S., Karmirantzou, M., Laskowski, R. A., Mitchell, J. B., Taroni, C. & Thornton, J. M. (1998). Protein folds and functions. *Structure*, **15**, 875-884.
- Matsuo, Y. & Bryant, S. H. (1999). Identification of homologous core structures. *Proteins: Struct. Funct. Genet.* **35**, 70-79.
- Murzin, A. G. (1998). How far divergent evolution goes in proteins. *Curr. Opin. Struct. Biol.* **8**, 380-387.
- Murzin, A. G. (1999). Structure classification-based assessment of CASP3 predictions for the fold recognition targets. *Proteins: Struct. Funct. Genet. Suppl.* **3**, 88-103.
- Murzin, A. G., 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.
- Musco, G., Kharrat, A., Stier, G., Fraternali, F., Gibson, T. J., Nilges, M. & Pastore, A. (1997). The solution structure of the first KH domain of FMR1, the protein responsible for the fragile X syndrome. *Nature Struct. Biol.* **4**, 712-716.
- Nogales, E., Wolf, S. G. & Downing, K. H. (1998). Structure of the $\alpha\beta$ tubulin dimer by electron crystallography. *Nature*, **391**, 199-203.
- Orengo, C. A., Jones, D. T. & Thornton, J. M. (1994). Protein superfamilies and domain superfolds. *Nature*, **372**, 631-634.
- Orengo, C. A., Michie, A. D., Jones, S., Jones, D. T., Swindells, M. B. & Thornton, J. M. (1997). CATH: a hierarchical classification of protein domain structures. *Structure*, **5**, 1093-1108.
- Orengo, C. A., Pearl, F. M., Bray, J. E., Todd, A. E., Martin, A. C., Lo Conte, L. & Thornton, J. M. (1999). The CATH database provides insights into protein structure/function relationships. *Nucl. Acids Res.* **27**, 275-279.
- Pappa, H., Murray-Rust, J., Dekker, L. V., Parker, P. J. & McDonald, N. Q. (1998). Crystal structure of the C2 domain from protein kinase C- δ . *Structure*, **6**, 885-894.
- Park, H. W., Boduluri, S. R., Moomaw, J. F., Casey, P. J. & Beese, L. S. (1997). Crystal structure of protein farnesyltransferase at 2.25 Angstrom resolution. *Science*, **275**, 1800-1804.
- Pastore, A. & Lesk, A. M. (1990). Comparison of the structures of globins and phycocyanins: evidence for evolutionary relationship. *Proteins: Struct. Funct. Genet.* **8**, 133-155.
- Pearson, W. R. (1996). Effective protein sequence comparison. *Methods Enzymol.* **266**, 227-258.
- Pearson, W. R. (1998). Empirical statistical estimates for sequence similarity searches. *J. Mol. Biol.* **276**, 71-84.
- Peat, T., Newman, J., Waldo, G. S., Berendzen, J. & Terwilliger, T. C. (1998). Structure of translation initiation factor 5A from *Pyrobaculum aerophilum* at 1.75 Å resolution. *Structure*, **6**, 1207-1214.
- Pickersgill, R., Smith, D., Worboys, K. & Jenkins, J. (1998). Crystal structure of polygalacturonase from *Erwinia carotovora ssp. carotovora*. *J. Biol. Chem.* **273**, 24660-24664.
- Ranganathan, R., Lu, K. P., Hunter, T. & Noel, J. P. (1997). Structural and functional analysis of the mitotic rotamase Pin1 suggests substrate recognition is phosphorylation dependent. *Cell*, **89**, 875-886.
- Rowen, L., Mahairas, G. & Hood, L. (1997). Sequencing the human genome. *Science*, **278**, 605-607.
- 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., Sasieni, P. D. & Sternberg, M. J. E. (1998). Supersites within superfolds. Binding site similarity in the absence of homology. *J. Mol. Biol.* **282**, 903-918.
- Sack, S., Muller, J., Marx, A., Thormahlen, M., Mandelkow, E. M., Brady, S. T. & Mandelkow, E. (1997). X-ray structure of motor and neck domains from rat brain kinesin. *Biochemistry*, **36**, 16155-16165.
- Sander, C. & Schneider, R. (1991). Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins: Struct. Funct. Genet.* **9**, 56-68.
- Schumacher, S., Clubb, R. T., Cai, M., Mizuuchi, K., Clore, G. M. & Gronenborn, A. M. (1997). Solution structure of the Mu end DNA-binding I β subdomain of phage Mu transposase: modular DNA recognition by two tethered domains. *EMBO J.* **16**, 7532-7541.
- Sippl, M. J. (1999). Who solved the protein folding problem? *Structure*, **7**, R81-R83.
- Sippl, M. J. & Weitckus, S. (1992). Detection of native like models for amino acid sequences of unknown three dimensional structure in a data base of known protein conformations. *Proteins: Struct. Funct. Genet.* **13**, 258-271.

- Sulzenbacher, G., Shareck, F., Morosoli, R., Dupont, C. & Davies, G. J. (1997). The *Streptomyces lividans* family 12 endoglucanase: construction of the catalytic core, expression, and X-ray structure at 1.75 Å resolution. *Biochemistry*, **36**, 16032-16039.
- Taylor, W. R. & Orengo, C. (1989). Protein structure alignment. *J. Mol. Biol.* **208**, 1-22.
- Tomchick, D. R., Turner, R. J., Switzer, R. L. & Smith, J. L. (1998). Adaptation of an enzyme to regulatory function: structure of *Bacillus subtilis* PyrR, a *pyr* RNA-binding attenuation protein and uracil phosphoribosyltransferase. *Structure*, **6**, 337-350.
- Van Asselt, E. J., Perrakis, A., Kalk, K. H., Lamzin, V. S. & Dijkstra, B. W. (1998). Accelerated X-ray structure elucidation of a 36 kDa muramidase/transglycosylase using wARP. *Acta Crystallog. sect. D*, **54**, 58-73.
- Yao, N., Hesson, T., Cable, M., Hong, Z., Kwong, A. D., Le, H. V. & Weber, P. C. (1997). Structure of the hepatitis C virus RNA helicase domain. *Nature Struct. Biol.* **4**, 463-467.
- Zhang, Z., Huang, L., Shulmeister, V. M., Chi, Y. L., Kim, K. K., Hung, L. W., Crofts, A. R., Berry, E. A. & Kim, S. (1998). Electron transfer by domain movement in cytochrome *bc₁*. *Nature*, **392**, 677-684.

Edited by R. Huber

(Received 9 September 1999; received in revised form 28 December 1999; accepted 28 December 1999)