

Abhinav Kumar,<sup>a,b</sup> Hsiu-Ju  
Chiu,<sup>a,b</sup> Herbert L. Axelrod,<sup>a,b</sup>  
Andrew Morse,<sup>a,c</sup> Marc-André  
Elslinger,<sup>a,d</sup> Ian A. Wilson<sup>a,d</sup> and  
Ashley Deacon<sup>a,b\*</sup>

<sup>a</sup>Joint Center for Structural Genomics,  
<http://www.jcsg.org>, USA, <sup>b</sup>Stanford  
Synchrotron Radiation Lightsource, SLAC  
National Accelerator Laboratory, Menlo Park,  
CA, USA, <sup>c</sup>Center for Research in  
Biological Systems, University of California, San  
Diego, La Jolla, CA, USA, and <sup>d</sup>Department of  
Molecular Biology, The Scripps Research  
Institute, La Jolla, CA, USA

Correspondence e-mail:  
[adeacon@slac.stanford.edu](mailto:adeacon@slac.stanford.edu)

Received 1 September 2009  
Accepted 3 March 2010

## Ligands in PSI structures

Approximately 65% of PSI structures report some type of ligand(s) that is bound in the crystal structure. Here, a description is given of how such ligands are handled and analyzed at the JCSG and a survey of the types, variety and frequency of ligands that are observed in the PSI structures is also compiled and analyzed, including illustrations of how these bound ligands have provided functional clues for annotation of proteins with little or no previous experimental characterization. Furthermore, a web server was developed as a tool to mine and analyze the PSI structures for bound ligands and other identifying features.

### 1. Introduction

International structural genomics initiatives, including the US-based Protein Structure Initiative (PSI; <http://www.nigms.nih.gov/Initiatives/PSI/>), have led to an unprecedented increase in the rate at which new protein structures are being solved and made available to the scientific community (Levitt, 2007). To date, these efforts have contributed over 7500 protein structures to the Protein Data Bank (PDB; Dutta *et al.*, 2009), more than half of which have come from the PSI. For the most part, the PSI effort has focused on determining unique structures from protein families that previously lacked any structural representative and on providing better structural coverage for large diverse protein families where more structures are needed to accurately model the entire family. Consequently, many of the proteins solved have little or no previous experimental characterization and have been classified as domains of unknown function (DUFs; Bateman *et al.*, 2010) or have only a tentative functional annotation based on amino-acid sequence homology. A variety of online tools and web-based search engines, such as *EBI-SSM* (Krissinel & Henrick, 2004), *DALI* (Holm *et al.*, 2008), *VAST* (Gibrat *et al.*, 1996) and *fast-SCOP* (Chi *et al.*, 2006), allow the inference of function based on structural similarity. However, these approaches have their limitations.

A significant number of the structures solved by structural genomics efforts can be assigned to superfolds (Orengo *et al.*, 1994), such as TIM-barrel and ferredoxin folds, whose members perform a wide diversity of biological functions. Thus, knowledge of the structure is often not sufficient to deduce the exact cellular function of a protein. To further aid in functional annotation, additional methods can be explored, such as catalytic residue matching and analysis of protein surface properties, although these methods usually only partially enhance the functional assignment (Binkowski *et al.*, 2005; Laskowski *et al.*, 2005a,b; Porter *et al.*, 2004).

Another challenge faced by large-scale structural biology efforts is to effectively disseminate the structural results to a broad scientific community. Although all of the PSI structures are deposited immediately into the PDB and rapidly released, only a small fraction of them have been described in publications in the scientific literature. Recently, efforts have been made to develop more streamlined web-based tools to rapidly disseminate key findings and new insights derived from these structures, as exemplified by the PSI Knowledgebase (<http://kb.psi-structuralgenomics.org>) and The Open Protein Structure Annotation Network (TOPSAN; <http://www.topsan.org/>; Krishna *et al.*, 2010). However, it is clear that complementary user-

friendly tools would be extremely beneficial to mine the latest structural data for functional and methodological insights. A rich source of functional data that is often overlooked in the PSI structures are the ligands that are identified during interpretation of the electron-density maps and subsequently built into the deposited

structures. More than half of the PSI structures (65%) contain bound ligands, such as metal ions, cofactors, substrates and effectors. Many of these ligands are acquired during protein production, whereas the remainder are incorporated into the protein at various steps during the purification and crystallization stages, as, for example, buffer

**Ligand Search Server**  
Send bug reports, comments to [Abhinav Kumar](mailto:Abhinav.Kumar@slac.stanford.edu)

(PSI Centers: ATCG3D BPSF BSGC CESH CHTSB CSGID CSMP ISFI JCSG MCSG NESG NYCOMPS NYSGXRC SECSG SGPP SSGCID TB)

Enter search queries in the boxes below:

Ligands:   
 Targets:   
 PDB Ids:   
 PFAM families:   
 Accession Ids:   
 Description:   
 Organisms:   
 PSI Centers:

Search tips:  
 You can enter partial search strings; search is case insensitive.  
 You can enter multiple queries in each box.  
 Entering queries in multiple boxes returns results common to all independent searches.  
 Entering a dot (.) in Ligands field and leaving the other fields blank shows the entire database.

Examples:  
 To search Ligands, enter ADP, UNL, etc. (Full ID) in the Ligands box.  
 To search Peptide ligands, enter Peptide, PHE, UNK etc. in the Ligands box.  
 To search Modified residues, enter ModRes, LLP, CSO etc. in the Ligands box.  
 To search Targets, enter TM, FB, CL, etc. in the Targets box.  
 To search PDB Ids, enter 2a6b, 2ur etc. in the PDB Ids box.  
 To search PFAM families, enter PF01234 etc. in the PFAM families box.  
 To search by Accession ids, enter YP\_012655.1 etc. in the Accession Ids box.  
 To search by Description, enter Hypothetical, Unknown, Hydrolase, etc. in the Description box.  
 To search by Organisms, enter Thermotoga, Coli, etc. (Full Word) in the Organisms box.

Note: Xtal ID information is restricted to JCSG members only. Email [Ashley Deacon](mailto:Ashley.Deacon@slac.stanford.edu) if you want this information.  
[Click here](#) to access the ligand visualization links.

**Search Results (7 hits)**  
 (FMN containing structures in PDB = 474 , PLP containing structures in PDB = 553 )

N	Target	Xtal ID	PDB ID	PFAM	Accession	Description	Organism	Ligands	PSI Center	Deposition Date
1	TM0831	72230	3csw	PF01063	BIG_0012.001391	Crystal Structure of a Putative Branched-chain Amino Acid Aminotransferase (tm0831) from Thermotoga Maritima at 2.15 Å Resolution	Thermotoga Maritima Msb8	CIT CL MPD PLP UNL	JCSG	4/10/2008
2	TM1785	37842	2ord	PF00202	MEGA_3.40.640.10	Crystal Structure of Acetylornithine Aminotransferase (ec 2.6.1.11) (acoat) (tm1785) from Thermotoga Maritima at 1.40 Å Resolution	Thermotoga Maritima Msb8	GOL PLP	JCSG	2/02/2007
3	TM1131	11504	1vp4	PF00155	MEGA_3.40.640.10	Crystal Structure of Aminotransferase, Putative (tm1131) from Thermotoga Maritima at 1.82 Å Resolution	Thermotoga Maritima Msb8	EDO FMT PLP UNL	JCSG	10/08/2004
4	SGT100		1tfv	PF01467 PF01687 PF06574	?	Crystal Structure of Adp, Amp, and Fmn Bound Tm379	Thermotoga Maritima	ADP AMP FMN	BSGC	5/07/2004
5	TM1255	3332	1ods	PF00155 PF00266	MEGA_3.40.640.10	Crystal Structure of Aspartate Aminotransferase (tm1255) from Thermotoga Maritima at 1.90 Å Resolution	Thermotoga Maritima	PLP SO4	JCSG	6/26/2003

Done

(a)

**Summary of Ligands (2246 structures)**  
 Number in the parentheses below indicates the number of structures in which that ligand occurs. Click on a ligand to see the details.

**Ligands** (348 structures; 198 different ligands)  
 UNL(101), UNX(23), SIN(8), DI(7), MAZ(6), GLC(6), MLT(4), GUN(4), PG6(3), SRT(3), SF4(3), SUC(3), PAF(3), NAG(3), NCA(3), APC(3), BAL(3), BU1(2), BGC(2), HHM(2), SAI(2), CLR(2), GNP(2), GAL(2), IBA(2), SAP(2), PCP(2), BIO(2), CEI(2), PG5(2), APR(2), DG(2), DA(2), C5P(2), G3H(2), G4P(2), PGR(2), MPO(2), ANP(2), 144(2), STE(2), PAJ(2), GLV(2), LMR(2), DHB(1), P4G(1), NIG(1), NRP(1), PRP(1), FBP(1), LGT(1), NIO(1), ABF(1), IPR(1), MTA(1), DIB(1), CP(1), MED(1), HEZ(1), A(1), 5GP(1), C(1), G(1), DIH(1), HED(1), G1P(1), GP2(1), DC(1), NBZ(1), ZMA(1), CEG(1), FRU(1), DIO(1), PLG(1), U(1), THF(1), B1M(1), ACM(1), ACP(1), DU(1), MMZ(1), OHA(1), NCN(1), 16A(1), M7P(1), 16G(1), XLF(1), CFS(1), 3GC(1), PEO(1), UMP(1), CTZ(1), CMK(1), ADE(1), KEG(1), RVF(1), XLS(1), BAM(1), ADI(1), NDG(1), PMP(1), TIM(1), B33(1), OXE(1), LUM(1), DD1(1), OXG(1), NDS(1), BEN(1), BMA(1), LFR(1), STH(1), FEO(1), G3P(1), OXN(1), BIU(1), FES(1), TYD(1), DGI(1), SEP(1), BET(1), GUD(1), STU(1), OXY(1), 8PP(1), CO2(1), SAL(1), 3PB(1), MP5(1), PGA(1), AES(1), CYN(1), 3PG(1), PNS(1), FUG(1), UUV(1), TRE(1), 3SL(1), 54D(1), LDA(1), TN4(1), GCS(1), TCL(1), SIB(1), PH2(1), NMN(1), MAL(1), SIC(1), MAN(1), RIP(1), MPR(1), RBE(1), ORO(1), DGT(1), PYR(1), DTP(1), DEP(1), UPG(1), HXA(1), 11X(1), GSP(1), CAA(1), GLA(1), NAI(1), AAT(1), 341(1), C2E(1), CHX(1), BDF(1), OSB(1), NPO(1), NGA(1), IHP(1), MYD(1), P2K(1), SOG(1), SVY(1), CAU(1), AGP(1), PYZ(1), AGS(1), TLP(1), 928(1), 12P(1), 1PS(1), OLC(1), DG2(1), DUT(1), CXS(1), GEO(1), PIM(1), PXP(1), TNE(1), G6P(1), G6Q(1)

**Peptides** (23 structures; 2 different peptides)  
 Peptide(22), UNK(1)

**Co-factors** (285 structures; 22 different co-factors)  
 FMN(46), NAD(38), FAD(24), COA(22), ADP(22), SAM(20), NAP(19), PLP(18), AMP(16), HEM(15), ATP(11), ACO(10), GDP(7), NDP(6), SAH(6), CDP(2), USP(2), GTP(2), COD(1), CNC(1), UTP(1), MLC(1)

**Metal ions** (899 structures; 24 different metal ions)  
 MG(271), ZN(227), NA(156), CA(121), NI(52), MN(37), FE(36), K(25), FE2(10), PI(8), CD(8), CO(8), HG(8), SM(2), PR(2), WO4(2), AU(2), ARS(1), YT3(1), BA(1), CS(1), SE(1), VO4(1), LI(1)

**Non-metal ions** (963 structures; 21 different non-metal ions)  
 SO4(443), CL(356), PO4(159), NO3(18), IOD(11), BR(11), SCN(11), CAC(10), POP(5), CO3(5), MLI(4), AZI(3), NH4(3), BCT(2), ALF(2), SO3(1), PO3(1), THU(1), PER(1), 1AL(1), OXL(1)

**Organics** (112 structures; 22 different organics)  
 BME(23), IPA(18), TLA(14), EOH(13), BEZ(6), PGO(5), ETX(5), AKG(5), DTT(4), DMS(3), TAR(2), OAA(2), ACE(2), GTT(2), AZI(1), XYL(1), MLA(1), SOH(1), PPI(1), MYR(1), MOH(1), LMT(1)

**Buffers** (368 structures; 14 different buffers)  
 ACT(146), FMT(60), ACY(59), CIT(39), TRS(24), EPE(23), IMD(18), MES(17), FLC(3), BTR(3), 10A(2), NHE(2), ICT(1), CPS(1)

**Precipitants** (173 structures; 13 different precipitants)  
 PEG(83), PG4(51), PGE(32), 1PE(10), P33(8), P6G(7), 2PE(6), PE4(4), PE8(3), PE5(2), BU3(1), 1PG(1), PEF(1)

**Salts** (2 structures; 2 different salts)  
 DPQ(1), AF3(1)

**Detergents** (2 structures)  
 BOG(2)

**Cryos** (761 structures; 4 different cryos)  
 EDO(379), GOL(370), MPD(41), MRD(9)

(b)

**Figure 1**

The Ligand Search Server and an example of its use. (a) The server's main page showing the search form and search example looking for PSI structures that contain either FMN or PLP bound to proteins from *Thermotoga maritima*. Tips on how to use the interface are displayed on the right and a partial list of structures is listed at the bottom. (b) A summary of all of the ligands bound to the PSI structures is displayed when the 'Summary' button is clicked.

**Table 1**  
Summary of ligands found in PSI and JCSG structures.

Type	% observed in PSI structures	% observed in JCSG structures	No. of unique compounds/entities
Ligands	12.3	15.6	285
Peptides	1.2	0.6	
Cofactors	9.4	10.9	22
Metals	24.7	26.4	24
Non-metals	27.2	40.7	21
Organics	3.0	3.9	23
Buffers	10.7	19.3	14
Precipitants	5.2	14.2	14
Cryoprotectants	21.3	51.6	3
Overall	65.0	85.2	

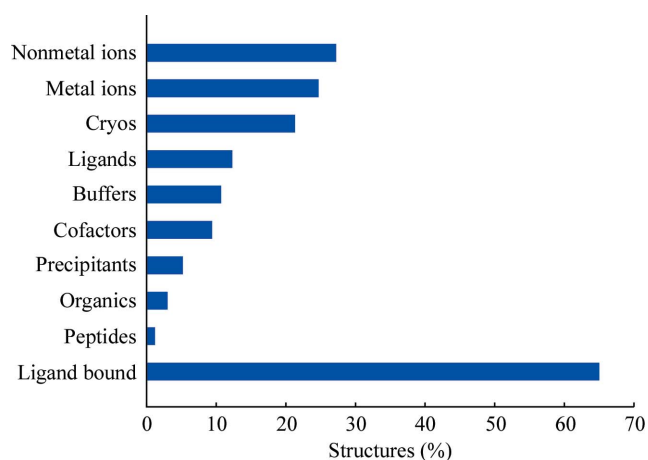
reagents, salts, precipitants and cryoprotectants. In many cases, these non-native ligands act as surrogates for the natural ligands owing to their similar biophysical properties. Their identification can often pinpoint favorable electrostatic regions or 'hot spots' on the protein and these surrogates often mimic the natural ligand–protein interactions, thus providing functional clues and insights.

The Joint Center for Structural Genomics (JCSG; <http://www.jcsg.org>) has designed the Ligand Search Server to be a fast and intuitive way to mine the PSI structures for detailed information regarding bound ligands. Searches can also be readily generated for entire families or for distinct classes of proteins or ligands, thus furthering collation and analysis of the functional knowledge derived from otherwise diverse sets of structures.

## 2. Methods

### 2.1. The Ligand Search Server

The JCSG Ligand Search Server ([http://smb.slac.stanford.edu/jcsg/Ligand\\_Search/](http://smb.slac.stanford.edu/jcsg/Ligand_Search/)) was created to mine PSI structures and to identify and classify the different types of bound ligands whether of functional relevance or not. The server also serves as a portal to complementary sites such as the Protein Data Bank (PDB; <http://www.pdb.org>), TOPSAN and Pfam (<http://pfam.sanger.ac.uk>; Finn *et al.*, 2008) which facilitate further exploration. The main user interface provides eight different search fields, including (i) the PDB ligand code, (ii) the PSI target name, (iii) the PDB code, (iv) the Pfam accession, (v) the protein/gene product accession ID, (vi) the structure description, as listed in the title of the PDB header, (vii) the source organism name



**Figure 2**  
Percentage of PSI structures that have any small-molecule ligand bound to them. The small molecules are categorized by their types.

**Table 2**  
Unique ligands found in PSI structures.

PDB code	Ligand name	Ligand ID	PSI center
1kph, 1kpi	Didecylidimethylammonium	10A	TBSGC
1z2l	Allantoate ion	1AL	NYSGXRC
1m33	3-Hydroxypropanoic acid	3OH	MCSG
1vr0	(2R)-3-Sulfolactic acid	3SL	JCSG
1y0g	2-[(2E,6E,10E,14E,18E,22E,26E)-3,7,11,15,19,23,27,31-Octamethyldotriacontan-2,6,10,14,18,22,26,30-octaenyl]phenol	8PP	NYSGXRC
1o8b	$\beta$ -D-Arabinofuranose-5'-phosphate	ABF	MCSG
1tuf	Azelaic acid	AZ1	NYSGXRC
1y80	Co-5-methoxybenzimidazolylcobamide	B1M	SECSG
2b4b	N-Ethyl-N-[3-(propylamino)propyl]propane-1,3-diamine	B33	NYSGXRC
2a3l	Coformycin 5'-phosphate	CF5	CECSG
2q09	3-[(4S)-2,5-Dioxoimidazolidin-4-yl]propanoic acid	DI6	NYSGXRC
2osu	6-Diazenyl-5-oxo-L-norleucine	DON	MCSG
2nw9	6-Fluoro-L-tryptophan	FT6	NESG
1p44	5-[[4-(9H-Fluoren-9-yl)piperazin-1-yl]carbonyl]-1H-indole	GEQ	TBSGC
2ou3	1H-Indole-3-carbaldehyde	I3A	JCSG
1x92	D-Glycero-D-mannopyranose-7-phosphate	M7P	MCSG
2gvc	1-Methyl-1,3-dihydro-2H-imidazole-2-thione	MMZ	NYSGXRC
1rtw	(4-Amino-2-methylpyrimidin-5-yl)methyl dihydrogen phosphate	MP5	NESG
2puz	N-(Iminomethyl)-L-glutamic acid	NIG	NYSGXRC
2od6	10-Oxohexadecanoic acid	OHA	JCSG
1n2h, 1n2i	Pantoyl adenylate	PAJ	TBSGC
1qpr	5-Phosphoribosyl-1-( $\beta$ -methylene) pyrophosphate	PPC	TBSGC
1xkl	2-Amino-4H-1,3-benzoxathiin-4-ol	STH	NESG
1bvr	Trans-2-hexadecenoyl-(N-acetyl-cysteamine)-thioester	THT	TBSGC
1lw4	3-Hydroxy-2-[(3-hydroxy-2-methyl-5-phosphono-oxymethylpyridin-4-ylmethyl)-amino]butyric acid	TLP	NYSGXRC

and (viii) the name of the PSI center. Each of these fields accepts multiple entries that are combined with a logical 'or' and entries in any of the eight search fields are then combined with a logical 'and' to generate the search query. A few search tips and examples are listed alongside the search form on the main page (see Fig. 1).

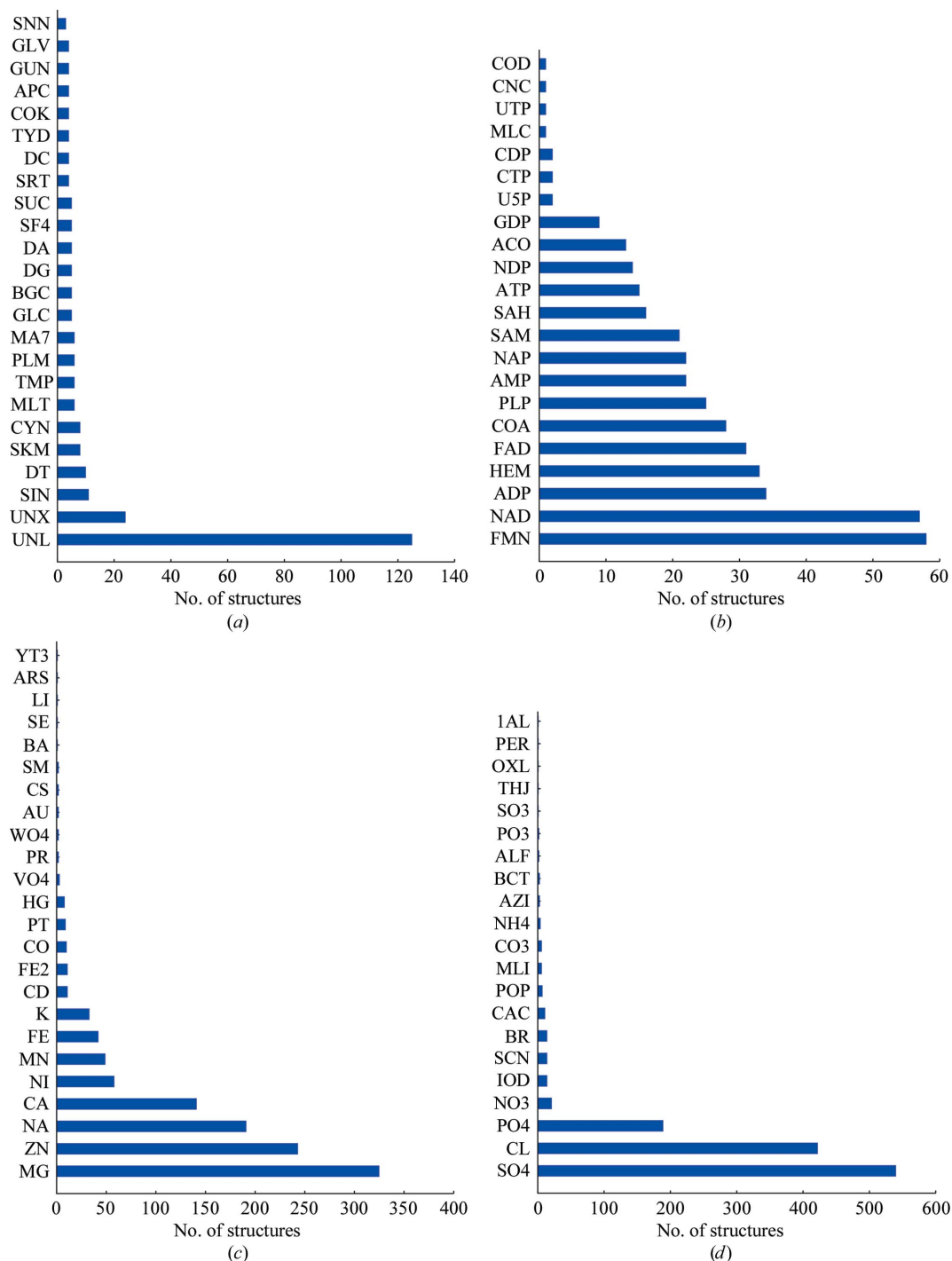
The 'Search' button submits the query against a locally maintained database which contains information on all of the PSI structures deposited in the PDB. The query results are returned as a single page that contains a concise tabular report at the top, which contains a row for every PDB structure that matches the query, lists the protein identifier used by the individual PSI center, the PDB code, the Pfam family name, the gene accession ID, the structure description, the source organism name, the bound ligands, the contributing PSI center and the deposition date. An additional column, 'Xtal ID', is included for JCSG structures which provides a link to specific information on the crystal used for structure solution, including all of the data and log files produced at various stages of structure solution and refinement. Most of the report fields are linked to other web resources to explore the structures further. This tabular report can also be exported to an Excel spreadsheet. Next, a ligand-visualization section provides links to HIC-Up (<http://alpha2.bmc.uu.se/hicup/>; Kleywegt, 2007) and Ligand Expo (<http://ligand-depot.rcsb.org>; Feng *et al.*, 2004) for each of the ligands found. Several summary sections that include information on the nature of the ligands found, the associated Pfam families and the source organisms follow. A 'Summary' button is also provided, which if used without any search query will generate an overall statistical report on all of the PSI structures in the database. More concise 'Summary' reports can be produced by including query values in the form fields.

## 2.2. Treatment of ligands at JCSG

During structure determination at the JCSG, attempts are made to account for all significant electron density observed during refinement. In addition to solvent molecules and chemical reagents used during protein production and crystallization, potential biological molecules, such as enzyme cofactors, substrates, products or their derivatives, which are presumably relevant to the protein's function and clearly supported by the electron density and chemical environment, are modeled into the structures, even if these molecules were not explicitly present in the reagents used during the protein preparation and crystallization stages.

The JCSG routinely uses X-ray fluorescence to identify metals that are bound in the structures. This technique allows the identification of most metals in the sample with a single experimental spectrum. When multiple metals are detected, X-ray diffraction data sets are then collected above and below the relevant X-ray absorption edges of the

comment, are modeled into the structures, even if these molecules were not explicitly present in the reagents used during the protein preparation and crystallization stages.



**Figure 3** Distribution of various ligands by category and relative frequency. Only the most common of these small molecules are shown. The names of the ligands follow the IDs used in the PDB and their full names can be obtained from the Ligand Expo Server (<http://ligand-depot.rcsb.org/ld-search.html>). (a) The 'Ligands' category includes biological ligands, such as substrates/products or their analogs. (b) The 'Cofactors' category includes various cofactors of enzymes but excludes ions, which are shown in the 'Metal ions' (c) and 'Non-metal ions' (d) categories.

metals and anomalous difference Fourier maps are calculated in order to unequivocally locate and confirm the identity of the bound metals. Lighter metals, such as Mg and Na, cannot be determined by X-ray fluorescence owing to limitations in our experimental setup; therefore, these are usually identified based on their binding geometries and environment.

Nevertheless, in many cases a suitable ligand cannot be unambiguously assigned to the electron density and the true identity of the ligand is inconclusive without further experimentation. The JCSG has adopted the policy of including these ligands as 'unknown ligands' and they are identified in the PDB file as UNL. The density is modeled by positioning a group of connected atoms that match the overall shape and a relevant description is included in the 'REMARK 3' field of the PDB header. To date, this strategy has surprisingly not been widely adopted by other PSI centers as it provides extremely valuable information that can be searched by a simple query; thus, the majority of these UNL-bound structures have been deposited by the JCSG (90%). Furthermore, all structures, including bound ligands, are internally peer-reviewed by at least one other scientist as a quality-control step prior to deposition in the PDB.

### 3. Overall statistics

A preliminary analysis of the 4200 currently available PSI structures shows that more than 2700 structures (~65%) contain small-molecule ligands of some kind. These ligands can be loosely classified as biological ligands (substrate, products, cofactors, inhibitors and their analogs) or surrogates, as well as peptides, ions, buffer molecules, crystallization reagents and cryoprotectants. This classification scheme is described in more detail in Table 1 and Fig. 1. Most of the functionally relevant biological ligands, including cofactors, were not explicitly added to the crystallization experiments. Hence, these ligands are endogenous to the expression systems and were acquired during protein production.

The overall distribution of the various types of ligands bound to PSI structures is shown in Fig. 2. It is of note that the JCSG reports more ligands in their structures compared with other PSI centers, particularly for the various ligands used as crystallization agents (buffers, precipitants and cryoprotectants); however, we also report more ligands in other categories. One possible explanation for this increased reporting of ligands comes from the standardized refinement and structure validation procedures implemented at the JCSG, in which specific steps (manual inspection and modeling of appropriate ligands) are undertaken to verify that all unmodeled electron density is properly accounted for. Indeed, a significant number of JCSG structures also contain 'unknown ligands' (UNLs), which refer to bound ligands that could not be unambiguously identified based on the electron density. The majority of these UNLs appear to be of biological importance since they are often located in crevices or cavities that resemble known active-site pockets or are identified based on comparison to structural homologs or other biochemical evidence. A survey of the number of biological ligands bound to PSI structures (Fig. 3) indicates that succinic acid (SIN), thymidine-5'-monophosphate (DT) and palmitic acid (PLM) are the most frequently observed and are likely to originate from the expression system. Similarly, flavin mononucleotide (FMN), nicotinamide adenine dinucleotide (NAD) and flavin adenine dinucleotide (FAD) are the most common cofactors. Magnesium ( $Mg^{2+}$ ), zinc ( $Zn^{2+}$ ) and sodium ( $Na^+$ ) are the most common metal ions and sulfate ( $SO_4^{2-}$ ), chloride ( $Cl^-$ ) and phosphate ( $PO_4^{3-}$ ) are the most common non-metal ions that are found in PSI structures. These particular ions are

**Table 3**

Ligands bound to proteins of unknown function, excluding common crystallization reagents and cryoprotectants.

Ligand†	Count	PDB codes
UNL	31	1vk9 1vpy 2aam 2g8l 2i8d 2opk 2pnk 2q9k 2qdr 2qe8 3cnx 3d82 3e80 3ebt 3ejv 3ez0 3ezu 3f7s 3ff0 3fgv 3fgy 3fh1 3fka 3flj 3fkd 3g16 3gi7 3giw 3gzz 3h3h 3hrq
ZN	29	1q9u 1sed 1su0 1t8h 1vk9 1vpy 1xaf 1xv2 1y7p 1ylo 2az4 2g7z 2gnr 2hek 2i9w 2oh3 2pg3 2pjs 2r8c 2rjb 3chv 3cjp 3di4 3dza 3e02 3e49 3feq 3fm2 3h0n
NA	29	1nnh 1nnw 1q8c 1sed 1vk1 1vmf 1vmh 1vmj 1yx1 1z67 1zl0 2asf 2fbl 2gkp 2hhg 2idl 2il5 2okq 2p0o 2pnk 2q3l 2qsv 2qzi 2ra9 3dnx 3f7c 3frm 3grd 3h0n
MG	28	1tzz 1z6n 1zd0 1zke 2a5z 2f4i 2fdr 2g80 2gfu 2h5n 2hx0 2i3d 2i71 2iec 2nn5 2o35 2oy9 2p3p 2p97 3bpd 3c5p 3cnx 3cu3 3e2v 3e06 3etk 3fa5 3hdg
CA	19	1sum 1vly 2arh 2esh 2g42 2gvj 2i6h 2pr7 2qng 2rld 3bdv 3bfm 3bvc 3db7 3dt5 3en8 3fyb 3g0k 3h36
UNX	11	1xrg 1xx7 1y81 1y82 1yb3 1ybx 1yby 1ybz 1yd7 1yem 1zd0
NI	10	1sum 1xx7 2aj7 2o8q 2ou6 2qe9 2qjv 3bvc 3d82 3h0n
FE	6	1sum 2rg4 3bv6 3bww 3dby 3hcl
K	5	1vph 1zl0 2aj7 2rgq 3hcl
NO3	4	1t6a 1t6s 3dde 3fof
MN	4	2p0n 3ck2 3fij 3gg7
COA	4	1q6y 1y81 1yre 2hqy
PT	2	1nnw 1yem
PLM	2	1mgp 1pzx
HG	2	1pvm 1qz4
FMN	2	2i51 2iml
SNN	1	3esm
SIN	1	3cqy
SE	1	2arh
SAM	1	2qe6
SAH	1	3go4
RIP	1	1y7p
NDP	1	1xkq
NBZ	1	3bgu
NAP	1	1i36
NAD	1	2o2z
HXA	1	2g7z
GLC	1	2esr
GDP	1	2hek
CO3	1	3c9q
CO	1	2h9f
BR	1	2hek
BEZ	1	2q9r
AU	1	1she
ATP	1	3gbu

† The names of the ligands follow the IDs used in the PDB; their full names can be obtained from the Ligand Expo server (<http://ligand-depot.rcsb.org/id-search.html>).

often present in the expression, purification and crystallization solutions, which may account for their frequent observation. A further analysis of the biological ligands reveals that 25 are unique to PSI structures and have not been observed previously in other structures deposited in the PDB, again indicating the richness and diversity of the information that is being derived from such structure determinations of proteins of unknown function (Table 2).

### 4. Unknown ligands (UNLs)

Examples of some UNL structures are shown in Fig. 4. About 75% of the UNL-bound structures now have some functional annotation and, therefore, biophysical and biochemical experiments can be designed to confirm the identity of the unknown ligands based on size and shape of the electron density as well as the nature of the environment surrounding the bound ligand. For example, in several instances the UNL resembles benzoic acid or nitrobenzene (PDB codes 2f4p, 2ig6, 2pbl, 3d82, 3ecf, 3ejv and 3ff0). However, these compounds were not modeled as such since neither was present in any of the reagents used nor was there any correlation with the

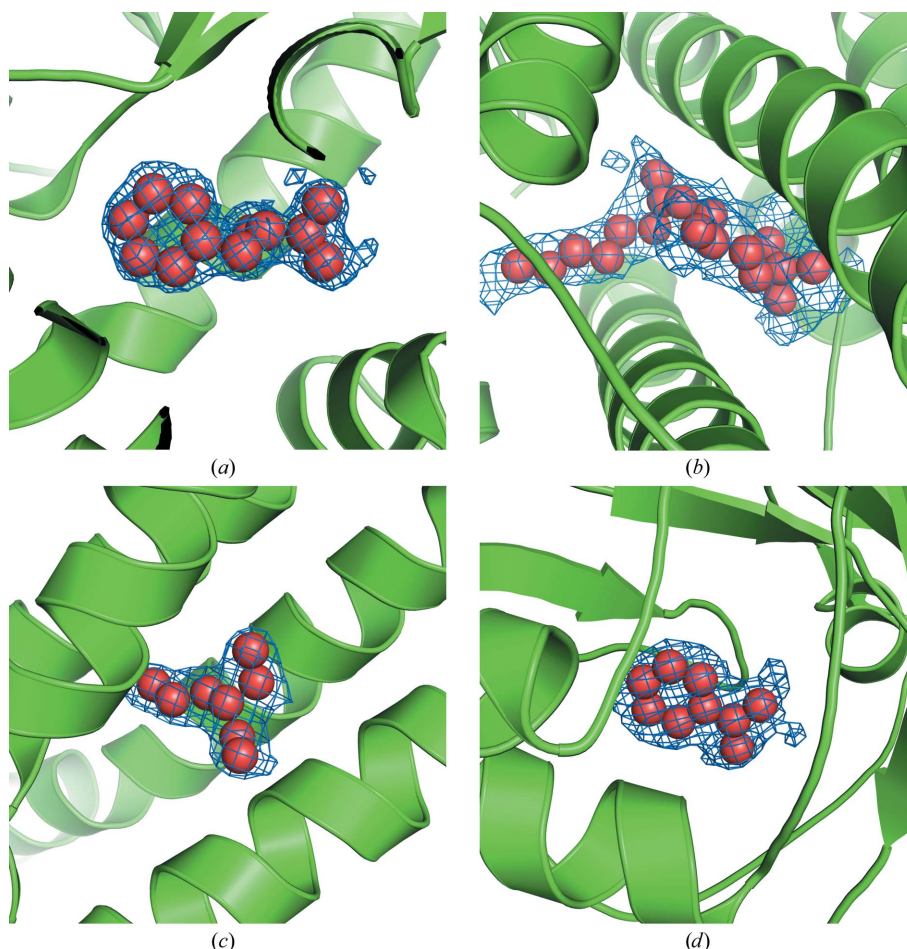
protein function. Uptake of endogenous molecules by proteins during the expression/purification stages is more common than is often appreciated, as exemplified by the occurrence of benzoic acid in 59 other structures in the PDB. However, in other cases, the UNL can provide functional clues about the protein. For instance, protein NP\_823353.1 (PDB code 3giw) is annotated as a protein of unknown function (Pfam DUF574) with an unknown ligand bound (<http://www.topsan.org/Proteins/JCSG/3giw>). The UNL resembles phenylalanine and the protein is structurally similar to SAM-dependent methyltransferases (Martin & McMillan, 2002; Fig. 4a), suggesting the possibility that it could be a phenylethanolamine *N*-methyltransferase (PNMT; Wong *et al.*, 1992), histamine *N*-methyltransferase (HNMT; Rutherford *et al.*, 2008) or catechol-*O*-methyl transferase (COMT; Weinshilboum *et al.*, 1999).

## 5. Metals bound to PSI structures

Approximately 25% of PSI structures and 27% of the JCSG structures contain metal ions (Table 1). Zn<sup>2+</sup> and Mg<sup>2+</sup> ions are among the most prevalent ligands in PSI structures, with 5.7 and 7.6% occurrence, respectively. Of the 857 structures determined by the JCSG as of July 2009, 226 contained metal ions (50 Zn<sup>2+</sup>, 19 Fe<sup>3+</sup>, 28

Ni<sup>2+</sup>, 43 Mg<sup>2+</sup>, 41 Ca<sup>2+</sup>, 47 Na<sup>+</sup>, 11 K<sup>+</sup>, three Mn<sup>2+</sup>, two Co<sup>2+</sup> and one Li<sup>+</sup>). The majority of Fe<sup>3+</sup> and Zn<sup>2+</sup> ions in PSI structures have a higher probability of being biologically relevant, since they are less frequently present in the crystallization buffers. For example, only 20% of the structures containing Zn<sup>2+</sup> ions report a zinc salt in the crystallization conditions. Other metals are potentially less biologically relevant as they are more frequently used during protein purification or crystallization. PSI structures containing Ca<sup>2+</sup>, Mg<sup>2+</sup> and Na<sup>+</sup> ions were obtained when such salts were used in 77, 64 and 61% of their crystallization conditions, respectively.

The identification of a bound metal can often aid in identification of the active site in a protein. For example, the crystal structures of three proteins of unknown function, YP\_164873.1 from *Silicibacter pomeroyi* DSS-3 (PDB code 3chv), YP\_556190.1 from *Burkholderia xenovorans* LB400 (PDB code 3e49) and YP\_555544.1 from *B. xenovorans* LB400 (PDB code 3e02), revealed structural similarity to 3-keto-5-aminohexanoate cleavage protein (YP\_293392.1) from *Ralstonia eutropha* Jmp123 (PDB code 3c6c), although their sequence identity (27–32%) is relatively low. Pairwise structural alignments gave a C<sup>α</sup> r.m.s.d. of 1.6 Å for 264 aligned residues between 3chv and 3c6c, a C<sup>α</sup> r.m.s.d. of 1.6 Å for 275 aligned residues between 3e49 and 3c6c and a C<sup>α</sup> r.m.s.d. of 1.7 Å for 259 aligned residues between 3e02 and 3c6c. All four structures share a conserved



**Figure 4**

Unknown ligands (UNL) in a few PSI structures. The UNL atoms are represented as red spheres enveloped by electron-density mesh ( $2F_o - F_c$  density contoured at  $1\sigma$  level above the mean) and surrounded by the protein rendered in cartoon representation. In many cases, the ligand could have been assigned as one or a few potential compounds, but is still annotated as a UNL since we have no definitive proof of its identity. (a) A protein of unknown function, NP\_823353.1 from *Streptomyces avermitilis*, at 1.45 Å resolution. (b) A protein of unknown function possessing a ferritin-like fold (YP\_832262.1; PDB code 3ez0) from *Arthrobacter* sp. Fb24 at 2.33 Å resolution. (c) A protein of unknown function from *Geobacter sulfurreducens* possessing a GGDEF domain (NP\_951600.1; PDB code 3ezu) at 1.95 Å resolution. (d) Phzb2 (NP\_250591.1; PDB code 3ff0) with a cystatin-like fold and an unknown function in phenazine biosynthesis from *Pseudomonas aeruginosa* at 1.90 Å resolution.

Zn<sup>2+</sup>-binding site in which almost all of the active-site residues are identical. Other examples of using structural knowledge about a bound metal to enhance the functional annotation are presented elsewhere in this issue. Bakolitsa and coworkers provide an example of the identification of Zn and Ni bound to the structure of the DUF1470 protein (Bakolitsa *et al.*, 2010). Axelrod and coworkers provide another good example where binding of Zn<sup>2+</sup> in the zinc-finger domain combined with structural comparisons suggest that two of the PF02663 Pfam family members in this study may bind nucleic acids and possibly function as transcriptional regulators (Axelrod *et al.*, 2010). These results have revealed functional and structural diversity within the PF02663 family.

## 6. Functional clues

### 6.1. Proteins of unknown function

Submitting the query 'Unknown', 'Uncharacterized', 'Hypothetical' or 'DUF' in the Description field of the Ligand Search Server finds 593 PSI structures (~14% of the total) that lack any functional annotation. The vast majority (474 structures) have been assigned to families in Pfam based on their amino-acid sequence.

About 66% of these 600 or so functionally unannotated proteins have one or more bound ligands. A closer examination of those ligands that are most likely to be biologically relevant (excluding crystallization and cryogenic reagents, although in some cases these may also provide clues to function) indicates that the most frequently found are either metal ions (22% of all ligands) or ligands with unknown identity (UNL; 5%), as shown in Table 3. Further analysis is necessary to determine their functional relevance. In a few cases, analysis of these ion-binding sites has already yielded definitive functional insights (see §5).

### 6.2. PSI contribution to new Pfam families

One of the key goals of PSI has been to increase the structural coverage of protein family space. Pfam coverage by the current set of PSI structures now extends to 1630 families; for approximately 700 (~43%) of these the PSI has provided the first structural representative. Over 150 of these Pfam families are populated by a single structure. Analysis of these first structural representatives representing 700 families indicates that over 175 of these structures contain

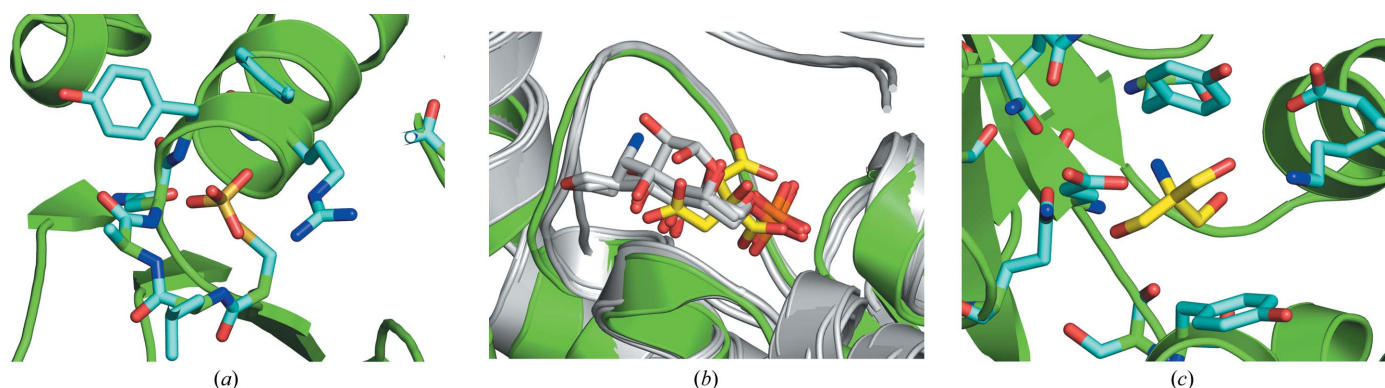
some biologically relevant ligands. Of these, Zn<sup>2+</sup> tops the list as the most frequently observed ligand in about 38 structures, followed by Mg<sup>2+</sup> in 35 structures, Na<sup>+</sup> in 23 structures, UNL in 16 structures and Ni<sup>2+</sup> in 12 structures.

### 6.3. Biological relevance of common molecules bound to proteins

Common reagents used during purification and crystallization, such as SO<sub>4</sub><sup>2-</sup>, Cl<sup>-</sup> or PO<sub>4</sub><sup>3-</sup> ions, buffer molecules such as Tris (2-amino-2-hydroxymethyl-propane-1,3-diol) or citrate, and precipitants such as polyethylene glycols *etc.*, often bind to proteins and are identified during structure refinement. In some cases, these bound reagents improve our understanding of putative binding sites on proteins and help to identify functionally relevant interactions by mimicking substrates. Here, we discuss three such examples (Fig. 5). A SO<sub>4</sub><sup>2-</sup> ion bound in the active site of YP\_001181608.1 (PDB code 3gxx; <http://www.topsan.org/Proteins/JCSG/3gxx>) mimics a substrate phosphate moiety and lends support to its annotation as a phosphatase. Similarly, a citrate molecule helped to identify the active site in YP\_001089791.1 (PDB code 3g68; <http://www.topsan.org/Proteins/JCSG/3g68>), where comparison of structurally similar proteins with a substrate bound in a similar location to the citrate led to the identification of likely active-site residues. In another example, the buffer molecule Tris is bound in the active site of the protein and emulates a sugar moiety in YP\_001304206.1 (PDB code 3h3l; <http://www.topsan.org/Proteins/JCSG/3h3l>).

## 7. Data mining of ligands in crystal structures for improving methodology

In addition to being a rich source of functional clues, the ligands bound to PSI structures can also serve as a source of data to improve crystallographic methods and map interpretation. As an example, we examined the frequency with which various cryoprotectant reagents are observed in crystal structures. We limited our analysis to JCSG structures, since we also had the precise crystallization and cryoprotective conditions used for each structure. Analysis of about 800 structures indicates that the most frequently observed cryoprotectant is ethylene glycol (EDO), with a probability of ~82% of being found in the structure if used in the crystallization/cryoprotective conditions, as shown in Table 4. The next on the list are polyethylene glycol 200 (PEG 200) and glycerol (GOL), with around a 56 or 55% chance,



**Figure 5**

Common reagents bound in the active sites of proteins. The protein structures are shown in cartoon representation and colored green or gray. The bound ligands are drawn as sticks and are colored yellow (carbon), red (oxygen) and blue (nitrogen). The interacting residues are also drawn as sticks with their C atoms colored cyan. (a) An SO<sub>4</sub><sup>2-</sup> ion bound in the active site of protein YP\_001181608.1 (PDB code 3gxx). (b) A citrate molecule bound to YP\_001089791.1 (PDB code 3g68) helped to identify the potential active site and was supported by substrates (gray) bound to the same location in structurally similar proteins (gray; PDB codes 1mos, 2bpl, 2poc and 2v4m). (c) A Tris molecule bound in the active site of YP\_001304206.1 (PDB code 3h3l).

**Table 4**

Frequency of cryoprotectant reagents found as bound ligands in JCSG structures.

Values in parentheses correspond to occurrences in all of the other structures in the PDB. These numbers are obtained from structures that report the use of these compounds in the crystal-growth conditions in their headers.

Cryo reagent†	No. of times used in crystallization/cryoprotective conditions	No. of times observed in structures	% observed
EDO	348 (888)	284 (184)	81.6 (20.7)
GOL	302 (3079)	167 (723)	55.3 (23.5)
MPD	106 (2558)	47 (373)	44.3 (14.6)
PEG 200	78	44	56.4
PEG 400	70	21	30.0

† Three-letter codes: EDO, ethylene glycol; GOL, glycerol; MPD, 2-methyl-2,4-pentanediol.

respectively, of being observed in the structure. A comparative analysis performed with all of the structures in the PDB, although limited because the crystal growth and cryoprotective conditions are often missing from the deposition record, indicates a much smaller frequency of observation of these compounds in crystal structures. For example, of the 888 structures that list ethylene glycol as a crystallization/cryoprotective component in the PDB header, it is observed in only 184 (20.7%) of these structures. Similarly, only 723 (23.5%) structures indicate the presence of bound glycerol out of 3079 structures that report its use during crystallization. The high frequency of occurrence of these cryoprotectants in our structures suggests that more care should be taken in general to identify these molecules during model building and refinement if present in the crystallization/cryoprotective conditions and to include cryoprotectants in addition to the crystallization conditions in the PDB header.

## 8. Conclusions

We have provided an overview of the various types of ligands bound in PSI structures and have tabulated their relative frequencies. Furthermore, we have described how ligands are identified and modeled into the structures at JCSG. The sheer number and diversity of ligands found in JCSG structures, based on a rigorous and systematic interpretation of the electron-density maps, suggests that for many structures in the PDB, ligands may have been overlooked or not adequately characterized. The observation of bound ligands, including unknown ligands and common chemical reagents mimicking potential biological ligands, often enhances the functional annotation of novel, uncharacterized proteins and generates hypotheses which can be validated experimentally. The JCSG Ligand Search Server provides an easy tool to survey the large collection of novel PSI structures for their bound ligands.

This work was supported by National Institutes of Health Protein Structure Initiative grant No. U54 GM074898 from the National Institute of General Medical Sciences (<http://www.nigms.nih.gov>). Most of this research was carried out at the Stanford Synchrotron Radiation Lightsource (SSRL). The SSRL is a national user facility operated by Stanford University on behalf of the US Department of Energy, Office of Basic Energy Sciences. The SSRL Structural Molecular Biology Program is supported by the Department of Energy, Office of Biological and Environmental Research and by the National Institutes of Health (National Center for Research Resources, Biomedical Technology Program and the National Institute of General Medical Sciences). Rachel K. Green from the RCSB PDB helped in the analysis of cryoprotectants in PDB structures. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institute of General Medical Sciences or the National Institutes of Health.

## References

- Axelrod, H. L. *et al.* (2010). *Acta Cryst.* **F66**, 1335–1346.
- Bakolitsa, C. *et al.* (2010). *Acta Cryst.* **F66**, 1198–1204.
- Bateman, A., Coggill, P. & Finn, R. D. (2010). *Acta Cryst.* **F66**, 1148–1152.
- Binkowski, T. A., Joachimiak, A. & Liang, J. (2005). *Protein Sci.* **14**, 2972–2981.
- Chi, P.-H., Shyu, C.-R. & Xu, D. (2006). *BMC Bioinformatics*, **7**, 362.
- Dutta, S., Burkhardt, K., Young, J., Swaminathan, G. J., Matsuura, T., Henrick, K., Nakamura, H. & Berman, H. M. (2009). *Mol. Biotechnol.* **42**, 1–13.
- Feng, Z., Chen, L., Maddula, H., Akcan, O., Oughtred, R., Berman, H. M. & Westbrook, J. (2004). *Bioinformatics*, **20**, 2153–2155.
- Finn, R. D., Tate, J., Mistry, J., Coggill, P. C., Sammut, S. J., Hotz, H. R., Ceric, G., Forslund, K., Eddy, S. R., Sonnhammer, E. L. & Bateman, A. (2008). *Nucleic Acids Res.* **36**, D281–D288.
- Gibrat, J. F., Madej, T. & Bryant, S. H. (1996). *Curr. Opin. Struct. Biol.* **6**, 377–385.
- Holm, L., Kaariainen, S., Rosenstrom, P. & Schenkel, A. (2008). *Bioinformatics*, **24**, 2780–2781.
- Kleywegt, G. J. (2007). *Acta Cryst.* **D63**, 94–100.
- Krishna, S. S., Weekes, D., Bakolitsa, C., Elsliger, M.-A., Wilson, I. A., Godzik, A. & Wooley, J. (2010). *Acta Cryst.* **F66**, 1143–1147.
- Krissinel, E. & Henrick, K. (2004). *Acta Cryst.* **D60**, 2256–2268.
- Laskowski, R. A., Watson, J. D. & Thornton, J. M. (2005a). *J. Mol. Biol.* **351**, 614–626.
- Laskowski, R. A., Watson, J. D. & Thornton, J. M. (2005b). *Nucleic Acids Res.* **33**, W89–W93.
- Levitt, M. (2007). *Proc. Natl Acad. Sci. USA*, **104**, 3183–3188.
- Martin, J. L. & McMillan, F. M. (2002). *Curr. Opin. Struct. Biol.* **12**, 783–793.
- Orengo, C. A., Jones, D. T. & Thornton, J. M. (1994). *Nature (London)*, **372**, 631–634.
- Porter, C. T., Bartlett, G. J. & Thornton, J. M. (2004). *Nucleic Acids Res.* **32**, D129–D133.
- Rutherford, K., Parson, W. W. & Daggett, V. (2008). *Biochemistry*, **47**, 893–901.
- Weinshilboum, R. M., Otterness, D. M. & Szumlanski, C. L. (1999). *Annu. Rev. Pharmacol. Toxicol.* **39**, 19–52.
- Wong, D. L., Lesage, A., Siddall, B. & Funder, J. W. (1992). *FASEB J.* **6**, 3310–3315.