# JOY: protein sequence-structure representation and analysis

*Kenji Mizuguchi[1], Charlotte M. Deane[1], Tom L. Blundell[1],*
*Mark S. Johnson[2] and John P. Overington[3]*

[1]*Department of Biochemistry, University of Cambridge, Tennis Court Road, Cambridge CB2 1GA, UK,* [2]*Department of Biochemistry, University of Turku, Turku, Finland and* [3]*Computational Chemistry and Protein Crystallography, Pfizer Central Research, Sandwich, Kent CT13 9NJ, UK*

## Abstract

*Motivation: JOY is a program to annotate protein sequence alignments with three-dimensional (3D) structural features. It was developed to display 3D structural information in a sequence alignment and to help understand the conservation of amino acids in their specific local environments.*
*Results: The JOY representation now constitutes an essential part of the two databases of protein structure alignments: HOMSTRAD (http://www-cryst.bioc.cam.ac.uk/~homstrad) and CAMPASS (http://www-cryst.bioc.cam.ac.uk/~campass). It has also been successfully used for identifying distant evolutionary relationships.*
*Availability: The program can be obtained via anonymous ftp from torsa.bioc.cam.ac.uk from the directory /pub/joy/. The address for the JOY server is http://www-cryst.bioc.cam.ac.uk/cgi-bin/joy.cgi.*
*Contact: kenji@cryst.bioc.cam.ac.uk*

## Introduction

The comparative study of molecular sequences and structures has provided many insights into protein folding, stability and evolution. The display of sequence alignments (e.g. Parry-Smith and Attwood, 1992; Barton, 1993) can be very useful for encapsulating the essential features of a particular fold or family, and highlighting features of particular biological interest (e.g. active site residues, glycosylation sites). However, when the three-dimensional (3D) structure of a particular sequence is known, much important information is lost either in the representation as a single sequence, or as an alignment. There is, therefore, a need to marry the sequence and structural information within a common representation.

As part of our research into protein evolution, we have developed a program, JOY, to annotate automatically sequence alignments with particular structural features. The basic idea of the program was to display, in a highly compact notation, the local environment of a residue in terms of a number of

| solvent inaccessible | UPPER CASE | X |
| solvent accesible | lower case | x |
| positive φ | *italic* | *x* |
| *cis*−peptide | breve | x̆ |
| hydrogen bond to other sidechain | tilde | x̃ |
| hydrogen bond to mainchain amide | **bold** | **x** |
| hydrogen bond to mainchain carbonyl | underline | <u>x</u> |
| disulphide bond | cedilla | ç |
| α−helix | red | x |
| β−strand | blue | x |
| $3_{10}$−helix | maroon | x |

**Fig. 1.** Key to formatted alignments. Tilde and breve are not displayed in an HTML output. Colours are only for colour PostScript and HTML outputs.

structural descriptors that are considered to be among the most important to the fold. A typesetting code (Figure 1) was developed that can unambiguously represent the structural environment, with emphasis given to residues that have many interactions with other residues within the protein (and are thus likely to be important to the fold). JOY was initially developed (Overington *et al.*, 1990, 1992) as an aid in the interpretation of COMPARER alignments (Sali and Blundell, 1990; Zhu *et al.*, 1992) where many of the structural features displayed by JOY are used in the structural alignment procedure. The program has since been used in a number of other research projects (e.g. Hoffren *et al.*, 1991; Sali and Blundell, 1993; Emsley *et al.*, 1994; Srinivasan *et al.*, 1996) and the representation is particularly appropriate for discussions of sequence-structure features within a single family.

Another key design consideration for the original program was that it should not require the use of colours (for page charges can sometimes be prohibitive) and that it should use a 'standard' text processing language so that integration with manuscripts is easy. Initially, LaTeX (Lamport, 1994) was chosen as the output device, but the development of a wide variety of text and graphics processing tools has required a

more flexible system. It also became clear that colours could make a qualitative difference in the analysis by adding valuable information. Here we describe a new version of JOY, with significant enhancements, including PostScript and HTML output and a WWW server. The program now plays an essential role in the construction of the two databases of protein structure alignments: HOMSTRAD (HOMologous STRucture Alignment Database; Mizuguchi *et al.*,1998b; http://www-cryst.bioc.cam.ac.uk/~homstrad) and CAM-PASS (CAMbridge database of Protein Alignments organised as Structural Superfamilies; Sowdhamini *et al.*, submitted; http://www-cryst.bioc.cam.ac.uk/~campass). The new representation of alignments has also proved to be useful for analysing distantly related proteins.

## Description of JOY

JOY is not an alignment program. It acts as a post-processor to a structural (or optionally sequence) alignment program, taking an alignment file and producing annotated alignments. JOY takes an input alignment (or a single sequence) in a format similar to that of the NBRF/PIR format (Figure 2). Several extensions have been introduced to allow easy labelling of the alignment and mixing sequence and structural information. A character sequence can be mixed with amino acid sequences by using the delimiter line '>T1;text'. The line that follows '>P1;' or '>T1;text' must be one of 'sequence', 'structure' or 'text'. The 'sequence' means a sequence with no structural information. The 'structure' means that there is structural information for the sequence and thus the output alignment will include annotations. The 'text' must follow the line '>T1;text'. In the usual NBRF/PIR format, gaps are indicated by hyphens ('-'), but JOY also accepts slashes ('/'). The latter is translated to a blank in the output alignment (Figure 3) and is normally used to indicate a chain break (e.g. missing electron density, multi-chain alignments). The same alignment format is used in the modelling program MODELLER (Sali and Blundell, 1993).

JOY produces a number of output files. The 'featured' alignment can be output in four different formats. The black and white PostScript file is designed for inexpensive print out, while the colour PostScript file (Figure 3) contains more information and is easier to capture various features (see Figure 1). The LaTeX (Lamport, 1994) file is mainly for further text processing and integration with manuscripts. The HTML file is easily visualized with a WWW browser and is used for the HOMSTRAD and CAMPASS databases. In addition to the featured alignment, JOY produces a 'template' file containing data for a profile-based sequence search program PSLAVE (Johnson *et al.*, 1993). JOY also outputs a list of ungapped blocks derived from the input alignment, which can be used as initial equivalences for the

```
# family: ras
>N1;!5p21-0
>P1;5p21-0
structure:5p21-0
-------------MTEYKLVVVGAGGVGKSALTIQLIQN-----HFVDEY---DPTI------------EDSYR
KQVVIDGETCLLDILDTAGQEEYSA--MRDQYMRTGEGFLCVFAINNT----------KSFEDIHQYREQIKRVK
DSDDVPMVLVGNKCDLA-------------------ARTVESRQAQD----LARSY------GIPYIETSAK--
----------------TRQGVEDAFYTLVREIRQH-*
>P1;1eft-1
structure:1eft-1
--------------HVNVGTIGHVDHGKTTLTAALTFVTAAENPNVEVKDYGDIDKAPEERARGIT-INTAHVE
YET-----AKRHYSHVDCPGHADYIK--NMITGAAQMDGAILVVSAADGPM-----------PQTREHILLARQV-
GVPY--IVVFMNKVDMV-------------------DDPELLDLVEMEVRDLLNQYE-FPGDEVPVIRGSALLA
LEEMHKNPKTKRGENEWVDKIWELLDAIDEY-----*
>P1;1tada1
structure:1tada1
--------------RTVKLLLLGAGESGKSTIVKQMKIIH/-----------------------------IETQ
FSFK----DLNFRMFDVGGQRSERK--KWIHCFEGVTCIIFIAALSAYDMVLVEDDEVNRMHESLHLFNSICNHR
YFATTSIVLFLNKKDVFSEKIKKAHLSICFPDYNGPNTYEDAGNYIKVQ---FLELNMRRDVKEIYSHMTCAT--
----------------DTQNVKFVFDAVTDIIIK--*
>P1;1hura0
sequence
GNIFANLFKGLFGKKEMRILMVGLDAAGKTTILYKLKLG---------EI---VTTIPTI-------GFNVETVE
YK--------NISFTVWDVGGQDKIRPLWRHYFQNTQGLIFVVDSNDR----------ERVNEAREELMRMLAED
ELRDAVLLVFANKQDLP-------------------NA-MNAAEITDKL--GLHSLR---HRNWYIQATCAT--
----------------SGDGLYEGLDWLSNQLRNQK*
```

**Fig. 2.** Sample input alignment for JOY. The alignment of the ras superfamily has been taken from the CAMPASS database and slightly modified for the purpose of illustration. The alignment is stored in a format similar to that of the NBRF/PIR format. A line beginning with the character 'C;' or '#' is a comment and ignored. In a special case where the word 'family:' follows a 'C;' or '#', the text after the 'family:' (in this example, ras) is used for the title of the output alignment. The line '>N1!5p21-0' indicates that the alignment will be labelled according to the PDB residue numbers of the structure corresponding to the code 5p21-0. The second line of the entry 1hura0 shows 'sequence', indicating that no structural information for this protein is used for annotations and that it will appear in the output alignment as a sequence only entry.

structure comparison program MNYFIT (Sutcliffe *et al.*, 1987).

Numbering of the alignment is possible in two different ways. First, a simple alignment-based numbering scheme is available. Second, a structure-based scheme labels the alignment according to the Protein Data Bank (PDB; Bernstein *et al.*, 1977; Abola *et al.*, 1987) residue numbers of a designated structure. Insertion codes are automatically placed within the alignment as required. In this way, it is possible to label an alignment in a family-based way (as is done in Figure 3 for the ras superfamily).

It is possible to display only a portion of the input alignment with the residues still in their proper structural environments. This option is useful when displaying aligned structural domains with each residue in its original environment within a whole protein. The segments to be displayed are read in from a separate input file.

A JOY server with a WWW interface has been developed (Figure 4). Users can submit either a PDB code or their own coordinate file. The output format can be specified by a menu and the formatted sequence can be retrieved in real time. Users can also submit a sequence alignment and retrieve the formatted alignment. In this case, only the atomic coordinates available from the PDB are used to format the user-supplied alignment.

**ras**

```
                                    10              20         A B C D E      30
5p21-0(  1 )              m t e Ỹ k L V V V G a g g V g K s̃ a L T̲ i q̃ L i q n - - - - - h f v d e y
left-1 ( 11 )               h̃ v ñ V G T̃ I G h v d H g K ĩ t L T̲̃ A A L T̲ f V T̲ a A e ñ p n V ẽ v k̃
1tada1 ( 28 )               r̲ t V k̲ L L L L G a g ẽ S̲ g k s t I v k Q̲̃ M k i i h̃
1hura        G N I F A N L F K G L F G K K E M R I L M V G L D A A G K T T I L Y K L K L G - - - - - - - - - E I
                         β β β β β β             α α α α α α α α α
```

```
               A B C        A B C D E F G H I J K L M     40              50              60
5p21-0( 33 )  - - - d p ĩ i - - - - - - - - - - - - - ẽ d s y r̃ k q̃ v v I d̃ g ē ĩ C̲ l L d̃ l l D̃ T̲ A G q ē e ỹ s a
left-1 ( 46 )  d̃ y g d I d̃ k a p e E r̲ a r̃ g i ĩ - i n t a h v ẽ Ỹ ẽ T̲̃ - - - - a k̃ r̲ h̃ Ỹ s H̃ V D̲ C̲ P G h̃ a d̃ ỹ i k
1tada1 ( 181 )                    i ẽ t q̃ f s̃ f k̃ - - - - d L ñ F r̲ M f D̲ v g G q r̃ s e r̲ k̲
1hura        - - - V T T I P T I - - - - - - - G F N V E T V E Y K - - - - - - - - N I S F T V W D V G G Q D K I
                         β β β β β                         β β β β β β
```

```
               A B     70            80         A B C D E F G H I J     90              100
5p21-0( 67 )  - - m R̃ d̃ q y M r t̲ G e̲ G F L C̃ V F A I n n ĩ - - - - - - - - - - k S̲ f e d̲ I h̃ q v̲ R̃ ē q̃ I k̃ r̃ V k̲
left-1 ( 91 )  - - n̲ M i t G A a q̲ M d̃ g A I L V V S̃ A a d̃ G p m - - - - - - - - - - - p q̃ T̲ r̃ ē H̲ I I I A r̲ q v -
1tada1 ( 206 )  - - k W̲ i h c̲ F e g V t C̲ I I F I A a L S a Y̲ d̃ m̃ v L v e d d e v ñ̲ r̲ M h̃ ē s I h̲ I F ñ s̲ I C̃ ñ̲ h̲ r̲
1hura        R P L W R H Y F Q N T Q G L I F V V D S N D R - - - - - - - - - - E R V N E A R E E L M R M L A E D
                         β β β β β β                         α   α α α α α α α α
```

```
                          110                120  A B C D E F G H I J K L M N O P Q R S T       130     A B
5p21-0( 105 ) d̃ s̲ d d V p M̃ V L V G Ñ̲ k̲ c̃ d̃ l a - - - - - - - - - - - - - - - - - - - - - a r̲ ĩ V e s̃ r q A q d̲ - -
left-1 ( 127 ) g V p y - - I V V F M N̲ k̲ v d̃ m v - - - - - - - - - - - - - - - - - - - - d d̃ p ē I I d̲ I V ẽ m e V
1tada1 ( 254 ) y F a t T s I V L F L Ñ̲ k̲ k d̃ v F s̲ ẽ k̃ I k̃ k a h L s̲ i c̲̃ f p d Y̲ n g p n̲ ĩ ỹ ē d̃ A G ñ y I k̃ v q̲ -
1hura        E L R D A V L L V F A N K Q D L P - - - - - - - - - - - - - - - - - N A - M N A A E I T D K L
                         β β β β β β                                α α α α α α
```

```
               C D         A B C D E F    140         A B C D E F G H I J K L M N O P Q R    150
5p21-0( 133 ) - - l A r s̲ ỹ - - - - - - g i p Ỹ i ē T̃ S̃ A k - - - - - - - - - - - - - - - - - - t r̃ q̃ g V ē d̃ A F
left-1 ( 155 ) r̲ d̃ l L ñ̲ q̃ ỹ e - F p G d̃ e V p V i r̲ G S̃ A l l A l ē ē M h̲ k̃ ñ p k̃ T̲̃ k̃ r̃ g e n e w V d k̃ I w ẽ L L
1tada1 ( 303 ) - - F l e l Ñ̲ m̃ r̲ r d v k e l y S h m̃ T̃ c̃ A ĩ - - - - - - - - - - - - - - - - - d ĩ q n V k f V F
1hura        - - G L H S L R - - - H R N W Y I Q A T C A T - - - - - - - - - - - - - - - - S G D G L Y E G L
                         α α             β β β                         α α α α α
```

```
                       160
5p21-0( 157 ) y ĩ L V r̃ e i r̃ q h̲
left-1 ( 204 ) d̲̃ a i d̲̃ ē y
1tada1 ( 333 ) d̃ a V T̲ d̃ i I i k
1hura        D W L S N Q L R N Q K
                         α α α α α
```

**Fig. 3.** Sample output alignment of JOY. The colour PostScript output for the alignment shown in Figure 2 is displayed. The aligned proteins are: H-ras p21 (5p21-0), first domain of elongation factor Tu (1eft-1), first domain of transducin-α (1tada1) and ADP-ribosylation factor 1 (1hura). The number in parentheses after a protein code indicates the PDB residue number at the beginning of each alignment block. The top line shows alignment positions in 'ras' numbering.
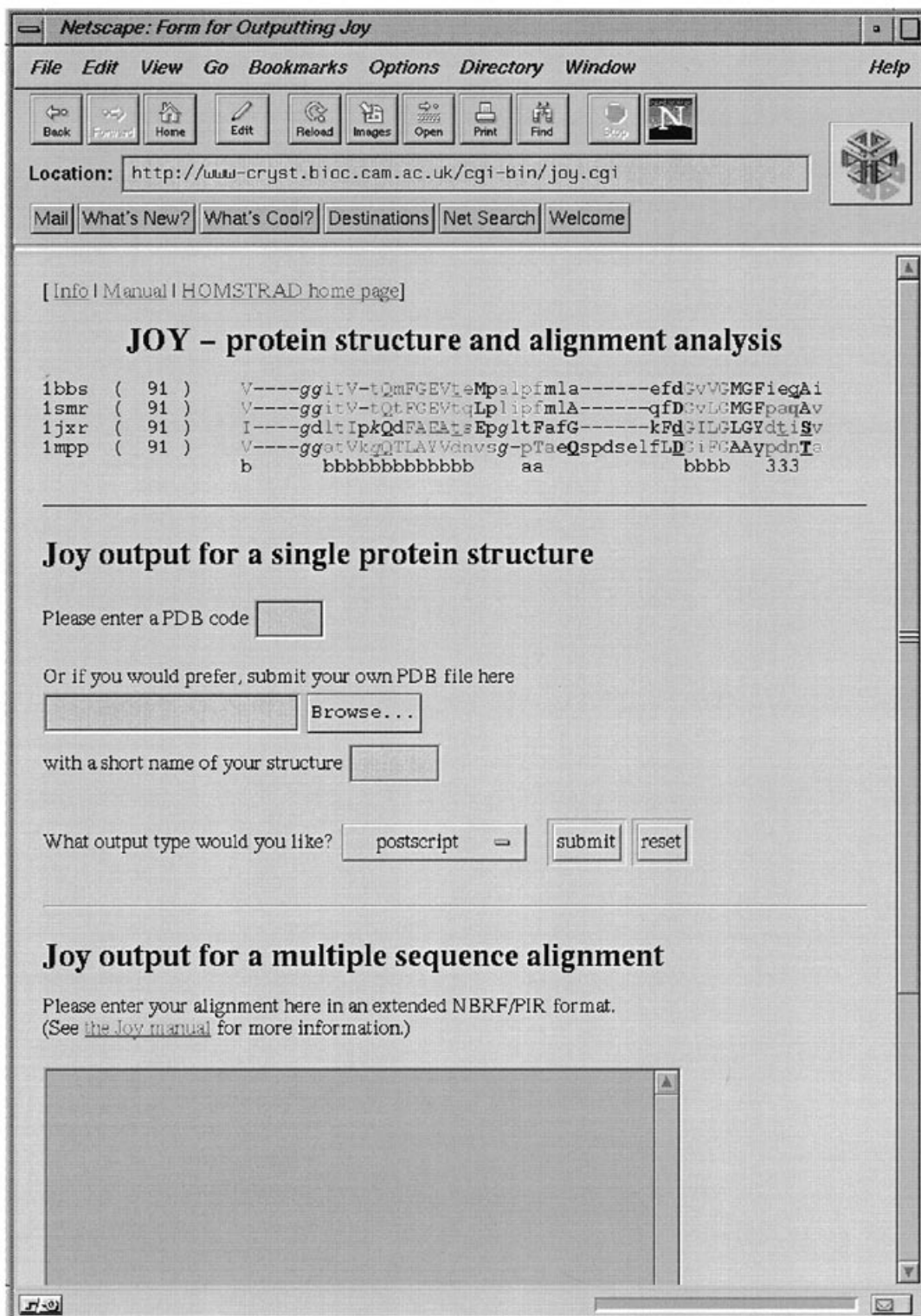
**Fig. 4.** WWW interface to the JOY server. The top part shows textports and an option menu for the formatting of a single protein structure. Below is a textport for the submission of a multiple alignment.

JOY requires a series of datafiles containing information about secondary structures, solvent accessibility and hydrogen bonding. These are produced automatically from a PDB file by supporting programs. These programs and the choice of structural features displayed are described below.

## Structural features represented by JOY

### Solvent accessibility

The partitioning of residues between a polar aqueous phase and a generally hydrophobic phase (the core of a globular protein) is established to be a major determinant in the process of protein folding. Residues in the solvent-inaccessible core of a protein are more conserved and are thus more useful for identifying distant evolutionary relationships. The program PSA is used to calculate the relative solvent-accessible surface area of all residues in a protein. The program uses an implementation of the algorithm of Lee and Richards (1971). Residues are defined as inaccessible by comparison to an extended conformation, and by default a 7% relative accessibility cut-off is applied (Hubbard and Blundell, 1987). The cut-off value can be set as a command line argument of JOY. Since solvent accessibility is known to be one of the most important attributes of a residue, a striking distinction is required in the typeset alignment. For this reason, we display buried residues in upper case and accessible residues in lower case.

### Secondary structure and main chain conformation

Secondary structure assignments are calculated with the program SSTRUC [D.K.Smith, unpublished; now also part of the PROCHECK suite of programs (Laskowski *et al.*, 1993)]. This program calculates the secondary structural state according to the definition of Kabsch and Sander (1983). It also provides information about the main chain dihedral angles $\phi$, $\psi$ and $\omega$. In a colour PostScript and an HTML output of JOY, repeating elements of secondary structure ($\alpha$, $3_{10}$ and $\pi$ helices, and $\beta$ strands) are shown in different colours. In a black and white PostScript file, secondary structures can be displayed, by specifying a command line option, underneath each sequence. By default, the consensus secondary structure is shown underneath the output alignment. The definition of consensus is that a fraction of >0.7 is in a particular conformational state at a given position.

The $\phi$ and $\psi$ angles are used to assign an 'Efimov conformational assignment' for a residue (Efimov, 1986). Efimov assignments can be displayed underneath each sequence. This annotation is useful for the display and identification of particular structural motifs such as $\beta$ turns (Wilmot and Thornton, 1990).

Our studies into the structural constraints on residue substitution during divergent evolution identified another important feature related to main chain conformation (Overington *et al.*, 1990). The presence of a positive main chain $\phi$ angle places severe restriction on what residue types may be accommodated at a particular position. Residues with a positive $\phi$ angle are shown in italics in the formatted alignment (Figure 1). These positions (often occupied by glycine residues, the only residue that lacks a side chain and can therefore readily adopt conformations on either side of the $\phi$ angle) are often found as the C-terminal cap to an $\alpha$ helix (Richardson and Richardson, 1988), at key positions with $\beta$ hairpins (Sibanda and Thornton, 1985) or within $\beta$ bulges (Richardson, 1981). We also show, for structural interest, the positions of *cis* peptide bonds with a breve above the residue (not in an HTML output; see Figure 1).

### Hydrogen bonds

Hydrogen bonds are calculated with the program HBOND (J.Overington, unpublished). HBOND identifies pairs of a heavy atom donor and acceptor, and calculates the associated energy and angles for each interaction. As we describe below, our main purpose for using the program is to identify side chain hydrogen bonds, but in a structure defined by X-ray crystallography side chain positions are often less accurate than the main chain positions, and it is often difficult to differentiate nitrogen and oxygen atoms in the side chain of asparagine and glutamine. Considering these ambiguities, we currently adopt a lenient definition of hydrogen bonds: a single distance cut-off of 3.5 Å. However, energy and angles printed out by HBOND can be used for refining the list of hydrogen bonds by further applying an energy-based or angular cut-off.

For the formatted alignment, the hydrogen bond data are split into four classes: (i) main chain to main chain (i.e. those responsible for secondary structure); (ii) side chain to main chain amide; (iii) side chain to main chain carbonyl; (iv) side chain to side chain/hetero atom group. Since the main chain-to-main chain hydrogen bond information is largely contained in the secondary structure representation, we do not explicitly display these data. It has been known that a buried side chain oxygen that is hydrogen bonded to a main chain amide proton plays a larger role in residue conservation than hydrogen bonds to a main chain carbonyl or to another side chain (Overington *et al.*, 1990). We therefore chose bold face to denote residues whose side chain is hydrogen bonded to a main chain amide. Residues whose side chain is hydrogen bonded to a main chain carbonyl, or to another side chain, are displayed as underlined and with a tilde, respectively (Figure 1). The deconvolution of these interactions can aid in the rapid identification of features such as helix N-terminal caps and potential salt bridges. Buried aspartic acid and threonine

in the active site of aspartic proteinase exemplify the importance of hydrogen bonds to main chain groups (see the HOMSTRAD entry of aspartic proteinase). Alternatively, unusual features can be easily identified; for example, a buried but non-hydrogen bonded polar residue might suggest local errors within a model.

Half-cystine residues are considered as a separate class of side chain-to-side chain interaction. These are displayed with a cedilla (in an attempt to convey their attachment to another half-cystine residue). This information is also calculated by the HBOND program.

HBOND also calculates potential hydrogen bonds and covalent bonds between protein and ligand. This is useful for identifying ligand-binding sites and for recognizing different cofactors within a homologous family. This information, however, is not used for JOY formatting at the moment.

## Implementation

JOY and the supporting programs are written in FORTRAN77 and ANSI C, and should compile and run on most UNIX platforms. Control of the program is through the use of UNIX command line arguments. These allow control of output formats, accessibility cut-off and other options. JOY assumes default filename extensions (e.g. .pdb for a PDB file, .sst for a secondary structure data file) to read in sequence and structural data. If JOY cannot find these files, it automatically runs appropriate supporting programs to generate them. The CGI program for the JOY server is written in PERL. The server runs on a Silicon Graphics Power Challenge with six 180 MHz R10000 processors.

## Discussion and conclusion

Taken together, these modifications of the one-letter amino acid code allow rapid correlation of structural with sequence features and aid the visual identification of residues likely to be essential to the fold (Figures 2 and 3). Particularly visible in this representation are the buried polar residues known to be amongst the most conserved in evolution (Overington *et al.*, 1990, 1992, 1993). One example is a buried glutamine conserved in a family of nerve growth factor (NGF) and related proteins. A JOY alignment of NGF, coagulogen and a *Drosophila* protein Spätzle enabled us to identify several key residues including this glutamine, which are conserved in these proteins, and this led to a hypothesis that Spätzle is also a member of this family (Mizuguchi *et al.*, 1998a). This example shows that JOY is useful for analysing distantly related proteins, where the overall sequence identity is low but several key residues are conserved that play important roles in stabilizing the fold.

The JOY representation is now used in two databases of protein structure alignments. In HOMSTRAD, protein alignments are compiled at the level of homologous families, while CAMPASS presents alignments of more distantly related proteins. In both databases, JOY representation plays an essential role in characterizing each family or superfamily.

## References

Abola,E.E., Bernstein,F.C., Bryant,S.H., Koetzle,T.F. and Weng,J. (1987) Protein Data Bank. In *Crystallographic Databases—Information Content, Software Systems, Scientific Applications*, Data Commission of the International Union of Crystallography, Bonn/Cambridge/Chester.

Barton,G.J. (1993) ALSCRIPT a tool to format multiple sequence alignments. *Protein Eng.*, **6**, 37–40.

Bernstein,F.C., Koetzle,T.F., Williams,G.J.B., Meyer,E.F., Brice,M.D., Rodgers,J.R., Kennard,O., Shimanouchi,T. and Tasumi,M. (1977) The Protein Data Bank: A computer-based archival file for macromolecular structures. *J. Mol. Biol.*, **112**, 535–542.

Efimov,A.V. (1986) Standard conformations of a polypeptide chain in irregular protein regions. *Mol. Biol. (Mosk)*, **20**, 250–260.

Emsley,J., White,H.E., O'Hara,B.P., Oliva,G., Srinivasan,N., Tickle,I.J., Blundell,T.L., Pepys,M.B. and Wood,S.P. (1994) Structure of pentameric human serum amyloid P component. *Nature*, **367**, 338–345.

Hoffren,A.-M., Saloheimo,M., Thomas,P., Overington,J.P., Johnson,M.S. and Blundell,T.L. (1991) Modelling the lignin peroxidase LIII of *Phlebia radiata* using a knowledge-based approach. *J. Chim. Phys.*, **88**, 2659–2662.

Hubbard,T.J. and Blundell,T.L. (1987) Comparison of solvent-inaccessible cores of homologous proteins: definitions useful for protein modelling. *Protein Eng.*, **1**, 159–171.

Johnson,M.S., Overington,J.P. and Blundell,T.L. (1993) Alignment and searching for common protein folds using a data bank of structural templates. *J. Mol. Biol.*, **231**, 735–752.

Kabsch,W. and Sander,C. (1983) Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers*, **22**, 2577–2637.

Lamport,L. (1994) *LATEX: A Document Preparation System*, 2nd edn. Addison Wesley, Reading, Mass.

Laskowski,R.A., MacArthur,M.W., Moss,D.S. and Thornton,J.M. (1993) PROCHECK: a program to check the stereochemical quality of protein structures. *J. Appl. Crystallogr.*, **26**, 283–291.

Lee,B. and Richards,F.M. (1971) The interpretation of protein structures: estimation of static accessibility. *J. Mol. Biol.*, **14**, 379–400.

Mizuguchi,K., Parker,J.S., Blundell,T.L. and Gay,N.J. (1998a) Getting knotted: a model for the structure and activation of Spätzle. *Trends Biochem. Sci.*, **23**, 239–242.

Mizuguchi,K., Deane,C.M., Blundell,T.L. and Overington,J.P. (1998b) HOMSTRAD: a database of protein structure alignments for homologous families. *Protein Sci.,* (in press).

Overington,J.P., Johnson,M.S., Sali,A. and Blundell,T.L. (1990) Tertiary structural constraints on protein evolutionary diversity: templates, key residues and structure prediction. *Proc. R. Soc. London Ser. B Biol. Sci.*, **241**, 132–145.

Overington,J.P., Donnelly,D., Johnson,M.S., Sali,A. and Blundell,T.L. (1992) Environment specific amino acid substitution tables: Tertiary templates and prediction of protein folds. *Protein Sci.*, **1**, 216–226.

Overington,J.P., Zhu,Z.-Y., Sali,A., Johnson,M.S., Sowdhamini,R., Louie,G.V. and Blundell,T.L. (1993) Molecular recognition in protein families: A database of aligned three-dimensional structures of related proteins. *Biochem. Soc. Trans.*, **21**, 597–604.

Parry-Smith,D.J. and Attwood,T.K. (1992) ADSP—a new package for computational sequence analysis. *Comput. Applic. Biosci.*, **8**, 451–459.

Richardson,J.S. (1981) The anatomy and taxonomy of protein structure. *Adv. Protein Chem.*, **34**, 167–339.

Richardson,J.S. and Richardson,D.C. (1988) Amino acid preferences for specific locations at the ends of α-helices. *Science*, **240**, 1648–1652.

Sali,A. and Blundell,T.L. (1990) The definition of topological equivalence in homologous and analogous structures: A procedure involving comparison of local properties and structural relationships through dynamic programming and simulated annealing. *J. Mol. Biol.*, **212**, 403–428.

Sali,A. and Blundell,T.L. (1993) Comparative protein modelling by satisfaction of spatial restraints. *J. Mol. Biol.*, **234**, 779–815.

Sibanda,B.L. and Thornton,J.M. (1985) Beta-hairpin families in globular proteins. *Nature*, **316**, 170–174.

Sowdhamini,R., Burke,D.F., Huang,J.-F., Mizuguchi,K., Nagarajaram,H.A, Srinivasan,N., Steward,R.E. and Blundell,T.L. CAMPASS: A database of structurally aligned protein superfamilies. Submitted.

Srinivasan,N., Bax,B., Blundell,T.L. and Parker,P.J. (1996) Structural aspects of the functional modules in human protein kinase-C alpha deduced from comparative analyses. *Proteins*, **26**, 217–235.

Sutcliffe,M.J., Haneef,I., Carney,D. and Blundell,T.L. (1987) Knowledge based modelling of homologous proteins, Part I: Three-dimensional frameworks derived from the simultaneous superposition of multiple structures. *Protein Eng.*, **1**, 377–384.

Wilmot,C.M. and Thornton,J.M. (1990) β-turns and their distortions: a proposed new nomenclature. *Protein Eng.*, **3**, 479–493.

Zhu,Z.Y., Sali,A. and Blundell,T.L. (1992) A variable gap penalty function and feature weights for protein 3-D structure comparisons. *Protein Eng.*, **5**, 43–51.