

S T A M P

Structural Alignment of Multiple Proteins

Version 4.4

User Guide

Manual by: *Robert B. Russell, Tom Walsh and Geoff Barton*

When using this program please cite R. B. Russell & G. J. Barton, Multiple protein sequence alignment from tertiary structure comparison, *PROTEINS: Struct. Funct. Genet.*,**14**, 309–323, 1992, and this manual.

The STAMP authors contact addresses (19 January 2010) are:

Prof. Robert B. Russell (RBR)
Cell Networks, University of Heidelberg
Room 564, Bioquant
Im Neuenheimer Feld 267
69120 Heidelberg
Germany

Tel: +49 6221 54 513 62
Fax: +49 6221 54 514 86
Email: robert.russell@bioquant.uni-heidelberg.de
WWW: <http://www.russell.embl-heidelberg.de>

Prof. Geoffrey J. Barton (GJB)
College of Life Sciences
University of Dundee
Dow Street
Dundee DD1 5EH
UK

Tel: +44 1382 385860
FAX: +44 1382 385764
E-mail g.j.barton@dundee.ac.uk
WWW: <http://www.compbio.dundee.ac.uk>

Copyright (1997,1998,1999,2010) Robert B. Russell & Geoffrey J. Barton.
STAMP is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation. See the LICENSE file for more details.

STAMP is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.

Note that this work was originally developed in the Laboratory of
Molecular Biophysics, University of Oxford.

Contents

1	Introduction and Overview	4
1.1	Preface to Version 4.4 - Tom Walsh	4
1.2	Preface to Version 4.2 - Rob Russell	5
1.3	Overview	7
1.4	Background	7
1.5	A brief description of the package	10
1.5.1	Initial superimposition	10
1.5.2	Pairwise comparisons and alignments (PAIRWISE) . .	11
1.5.3	Multiple alignment (TREEWISE)	11
1.5.4	Structure database scanning (SCAN)	12
1.5.5	Displaying STAMP output (VER2HOR, DSTAMP, GSTAMP)	15
1.6	Some comments on interpreting structural similarities	16
1.7	The programs contained within the package	17
2	Worked examples	20
2.1	Setup of examples	20
2.2	Multiple alignment using an initial multiple sequence alignment	21
2.3	Database Scanning	25
2.4	Using SCAN mode as the starting point for multiple alignment	28
2.5	Generating a set of superimposed structures	30
2.6	Alignment without an initial multiple alignment using ROUGH-FIT	30
2.7	Protein domain databases	32
2.8	Generating transformed coordinates using TRANSFORM . . .	34
2.9	Generating averaged coordinates	35
2.10	Displaying/processing the output	37
2.10.1	POSTSTAMP	37

2.10.2	STAMP_CLEAN	38
2.10.3	Displaying text alignments	38
2.10.4	Pretty Alignments via ALSRIPT	40
2.10.5	Pretty Structures via MOLSCRIPT	41
3	Input and Output format for all programs	44
3.1	Describing domain structures	44
3.2	Transformations	47
3.3	Sequence format	47
3.4	Multiple alignment format	48
3.5	Output from STAMP database scanning mode	50
3.6	Output to standard output or log file	51
4	Summary of STAMP parameters	53
4.1	Main program (STAMP)	53
4.2	Summary of parameters for other programs	63
4.2.1	PDB checker (PDBC)	63
4.2.2	PDBSEQ	64
4.2.3	ALIGNFIT	65
4.2.4	VER2HOR	66
4.2.5	DSTAMP	66
4.2.6	SORTTRANS	68
4.2.7	TRANSFORM	68
4.2.8	PICKFRAME	69
4.2.9	MERGETRANS & EXTRANS	70
4.2.10	MERGESTAMP	71
4.2.11	AVESTRUC	71
4.2.12	GSTAMP	72
4.2.13	STAMP_CLEAN	74
4.2.14	Converting alignment formats using ACONVERT	74
5	Installation	75
5.1	Compiling/running	75
5.2	Configuring STAMP	76
5.3	Getting other programs	78
6	Some of our studies involving STAMP	79

Chapter 1

Introduction and Overview

1.1 Preface to Version 4.4 - Tom Walsh

The changes in this release are:

1. The biggest change is that STAMP 4.4 is now licensed under the GNU General Public License.
2. The manual has been updated to reflect some minor changes in the STAMP output format.
3. The source code has been modified to remove references to obsolete header files and allow STAMP to be compiled on Mac OS X.
4. Mac OS X has been added as a build target.
5. The bundled SCOP domain databases have been updated to SCOP release 1.75.

1.2 Preface to Version 4.2 - Rob Russell

First, a big acknowledgement to Steve Searle (European Bioinformatics Institute) for getting STAMP to run (finally) under OSF and 64 bit machines generally. Also thanks to Andrew Torda (Australian National University, Canberra), Dave Schuller (University of California-Irvine), Mike Tennant (SmithKline Beecham Pharmaceuticals, Harlow, UK), Asim Siddiqui (LMB, Oxford) G.P.S. Raghava (LMB, Oxford), and James Cuff (EBI) for their help and various painstaking trawls through my spaghetti code.

Apart from bugs, etc., the noticeable changes are:

1. STAMP now reads compressed PDB and DSSP files. It will also look for files that are stored in a Brookhaven-style directory structure (e.g. `distr/mb/pdb4mbn.ent`).
2. Output is now flushed (`fflush`) during scanning. Purely a cosmetic thing for those who want up to the minute output when the program is running.
3. AVESTRUC now has an option to calculate an average for all aligned positions. It also now outputs values to the temperature factor fields in the PDB output to denote those averaged positions corresponding to structurally equivalent regions (blue in RASMOL colouring by temperature) and those equivalenced fortuitously (red). It will also highlight positions showing identical or conserved residue character.
4. PDBSEQ has some new options, including the ability to output *separate* files for each domain, and now outputs a sensible description of the protein by considering the TITLE, COMPND and SOURCE entries in each PDB file. Note that the default format is now FASTA.
5. DSTAMP has been changed dramatically, and is now (I think) much more useful. The input files for ALSCRIPT are now much prettier, including cylinders/arrows for helices/strands and colouring/fonting according to residue property conservation within the sequence alignment. It can also be used on alignments not derived using STAMP (i.e. from GCG, AMPS etc.).
6. STAMP now appears to run smoothly under OSF. Once again thanks

to Steve Searle. Versions have also been compiled and tested on IRIX, Solaris and Linux.

7. CLUS2BLC and SMSF2BLC have been replaced by a general alignment conversion program written in Perl (ACONVERT). More details are given below. Note that this is a general alignment conversion utility that might be useful in contexts other than STAMP.

8. The significance of sequence identity following structural alignment is now estimated according to Murzin (1993), JMB, 230, 689-694.

9. Three new programs have been added to the package (see specific instructions below):

MERGETRANS allows one to combine transformations from a variety of different sources (i.e. ALIGNFIT and STAMP). It either uses a user-specified identifier to link the two files (i.e. one found in both files) or the first common identifier if none is specified.

MERGESTAMP. Like MERGETRANS this program permits one to merge various kinds of STAMP data. However, it considers more than merely the transformations, and attempts to combine the alignments as well. It can be used in exactly the same way as MERGETRANS (i.e. to combine files that contain only transformations), but will also attempt to merge alignments in the file, if they are present. The alignments must be in BLOCK format (see the depths of the manual for details, and for how to convert things like Clustal or MSF into BLOCK format). MERGESTAMP can combine files that do not contain transformations as well (i.e. those that contain only alignments), and can thus be used for sequence data handling as well.

EXTRANS allows one to select and extract particular domains from a transformation file.

1.3 Overview

STAMP is a package for the alignment of protein sequences based on three-dimensional (3D) structure. It provides not only multiple alignments and the corresponding ‘best-fit’ superimpositions, but also a systematic and reproducible method for assessing the quality of such alignments. It also provides a method for protein 3D structure data base scanning. In addition to structure comparison, the STAMP package provides input for programs to display and analyse protein sequence alignments and tertiary structures. Please note that, although STAMP outputs a sequence alignment, it is a program for 3D structures, and NOT sequences. If you are after a multiple sequence alignment for proteins of unknown 3D structure, stop reading now and contact GJB for information about AMPS, which can be used to perform multiple sequence alignments, or see www.jalview.org for GJB’s latest methods for this problem.

Comparison of 3D structures is a complicated business, particularly if one wants to do unusual things (i.e. reverse a strand direction, swap two segments of a structure around, only consider equivalent structures of greater than 10 residues, etc.). Complicated things are possible with STAMP but as a consequence, the method is very complex. Please be patient, and read this manual carefully.

Alternatively, if you only want to do fairly straightforward things, such as align a set of structures or search a database of structures for similarities, you can skip the remainder of this chapter and go straight to the next one (Chapter: 2), which contains a few worked examples that should demonstrate how to use STAMP in a black-box way.

1.4 Background

The aim of this work was to provide a set of multiple sequence alignments derived from structure alone. These alignments have obvious uses which have been described elsewhere [1, 2]. Numerous other means of deriving such alignments have been presented, but, at the time of the development of STAMP,

only one had been applied to alignments of more than two sequences, and no systematic method for assessing the quality of the alignments had been provided. These, then, were the goals of this work.

At the heart of the method is the Argos & Rossmann [3] equation for expressing the probability of equivalence of residue structural equivalence:

$$P_{ij} = \exp\frac{d_{ij}^2}{-2 \times E_1^2} \exp\frac{s_{ij}^2}{-2 \times E_2^2}$$

where d_{ij} is the distance between C_α atoms for residues i and j , and s_{ij} is a measure of the local main chain conformation. A detailed description of this equation, and how it has been applied to multiple structures is given in [1].

STAMP makes extensive use of the Smith-Waterman (SW) algorithm [4, 5, 6]. This is a widely used algorithm which allows fast determination of the best path through a matrix containing a numerical measure of the pairwise similarity of each position in one sequence to each position in another sequence. Within STAMP, these similarity values correspond to modified P_{ij} values (above).

The result of the SW algorithm applied to a matrix of modified P_{ij} values is a list of residue equivalences. From this list we may obtain a set of equivalenced C_α positions. These are used to obtain a best fit transformation and RMS deviation by a least squares method [7, 8]. This transformation can be applied in the relevant way to yield two new sets of coordinates for which calculation (and correction) of P_{ij} values, the SW path finding and the least squares fitting can be repeated in an iterative fashion until the two sets of coordinates, and the corresponding alignment, converge on a single solution.

This strategy has proved successful in the generation of tertiary structure-based multiple protein sequence alignment for a wide variety of diverse protein structural families [1, 9, 10, 11, 12]. The method can accurately superimpose and obtain alignments for families of proteins as structurally diverse as the greek key β sandwich folds (e.g. immunoglobulin domains, CD4, PapD chaperonin, azurin, superoxide dismutase, actinotaxin, prealbumin, etc.), the

aspartic proteinase N- and C-terminal lobes, the Rossmann fold domains, the globin folds (including phycoyanins and colicins), and many others.

It is important to remember that this method assumes overall topological similarity, and will not, without explicit intervention, be able to superimpose/align structures with common secondary structures in similar orientations, but different connectivity or topologies (such as the different types of four helix bundle proteins: up-down-up-down with up-up-down-down).

Two measures of alignment confidence are provided [1]

1. A structural similarity Score (S_c) is defined in order that overall alignment quality and structural similarity can be compared across a wide range of protein structural families. These are defined below.
2. A measure of individual residue accuracy P'_{ij} is defined in order that residue equivalences can be normalised with respect to both the number of structures in an alignment and the length of the structures being aligned.

Alignments having a structural similarity Score S_c between 5.5 and 9.8 imply a high degree of structural similarity and almost always suggest a functional and/or evolutionary relationship. Values between 2.5 and 5.5 correspond to more distantly related structures, and do not always imply a functional or evolutionary relationship. Values less than 2.0 generally indicate little overall structural similarity.

Stretches of three or more aligned positions with P'_{ij} values greater than 6.0 generally correspond to genuine topological equivalences, values between 4.0 and 6.0 are equivalent > 50% of the time, and values less than 4.0 are generally not equivalent. Stretches of residues having $P'_{ij} > 6.0$ generally correspond to regions of conserved secondary structure within a family of structures being compared. For multiple alignments, an alternative and more effective way of assessing residue-by-residue equivalence is provided in POST-STAMP (see below).

Both of these measures are referred to repeatedly below. For a more detailed description of their derivation please refer to [1]. In addition, RMSD

is used to refer to the root mean square deviation between atoms selected for a fit. The CUTOFF refers the lowest allowable P'_{ij} for the program to use a particular pair of residues in a fit (called 'C' in [1]).

1.5 A brief description of the package

What follows is a brief overview of each application of STAMP. A detailed description of each of these can be found in later sections.

1.5.1 Initial superimposition

The structure comparison algorithm of Argos & Rossmann [3], which is the method used by STAMP, requires that the protein structures being compared are approximately superimposed initially. If not then structural similarity may be undetected, and reliable superimpositions and alignments unattainable. This is a very important thing to remember about STAMP. If initial superimpositions do not yield high enough scores (i.e. $S_c < 2.0$) or if the structures are generally different, STAMP will warn you by printing 'LOW SCORE' warnings in its output.

The STAMP package provides three methods of arriving at an initial superimposition. The first of these is to make use of an alignment derived on the basis of sequence. The program ALIGNFIT requires that the sequences extracted from the PDB files (using the program PDBSEQ) are aligned vertically in AMPS block format (see format and examples below); one can use AMPS or another method of aligning sequences. The ACONVERT program is included in the distribution to facilitate converting alignments to AMPS format from other formats; it also possible to use Jalview (www.jalview.org) to perform the conversion. STAMP compares all possible pairs of structures by performing a least squares fit on all equivalenced C_α atoms. Once all pairwise comparisons are compared, the program makes use of a tree to superimpose multiply all coordinates following the tree. Thus the final superimposition output is the best possible fit of the structure given the alignment. For an example where ALIGNFIT is used to provide an initial superimposition, refer to the alignment of the serine proteinases in Chapter 2.

In instances where multiple sequence alignment is inaccurate, ALIGNFIT may still be used, though the initial superimpositions may not be accurate enough for STAMP to find structural similarity. In such cases, the best way to arrive at initial superimpositions is to use the SCAN option within STAMP. This option compares a query domain against a database of target domains and generates a set of superimpositions of the target domains onto the query domain. This set of superimpositions constitutes a multiple alignment that can be used by STAMP as an initial alignment. This works particularly well when structures are very diverse. For an example, see the alignment of the aspartyl proteinase N- and C- terminal lobes in Chapter 2.

1.5.2 Pairwise comparisons and alignments (PAIRWISE)

Given a suitable initial superimposition of structures, the best way to obtain a multiple alignment and superimposition of a diverse family of domains is to follow a hierarchy of similarity. This allows most similar domains to be compared/aligned first, and only makes comparisons/alignments between distantly related domains at a later time in the procedure.

Pairwise comparisons are an ideal way to obtain such a hierarchy. The PAIRWISE options in STAMP will result in all $N \times (N - 1)/2$ comparisons being performed and will output a matrix of pairwise similarities. This can then be used to produce a dendrogram, or *tree*, from which multiple alignments and superimpositions may be generated.

1.5.3 Multiple alignment (TREEWISE)

Given the initial set of superimpositions, and a set of PAIRWISE similarity scores, the TREEWISE option will perform all alignments that are possible given a dendrogram generated by considering the PAIRWISE scores. Statistics, transformations and alignments are output at each stage of the hierarchy so that a continuum of structure variation can be observed (i.e. the output will become more and more structurally varied as the search progresses).

Note that by default, STAMP performs both PAIRWISE and TREEWISE procedures together.

1.5.4 Structure database scanning (SCAN)

It is often desirable to compare a particular domain or protein structure to a database of known 3D structures in order that structurally similar proteins may be found.

Given a single protein domain (a *query*) and a list of domains to which it is to be compared (a *database*), STAMP can be used to perform all possible comparisons of the query to the database structures. The initial superimposition problem is solved by attempting more than one initial fit with each database structure. This can be done in one of two ways, which are named FAST and SLOW, for the obvious reasons.

In FAST mode, fits are performed by laying query sequence onto the database structure starting at every *i*th position, where *i* is an adjustable parameter usually set to five (i.e. the sequence is laid onto the 1st, 6th, 11th, etc. position). Diagrammatically, this looks like:

Q=query, D=database

```
Fit 1  Q  -----  
      D  -----  
Fit 2  Q   -----  
      D  -----  
Fit 3  Q    -----  
      D  -----  
<etc.>
```

This approach is fine if the query is a single domain, and there is a strong similarity in the database structure. However, if similarity is weaker, or if the query contains multiple domains (in which case it is advisable to split the query into multiple domains, if possible), then SLOW mode will perform more fits by sliding query and database sequences along each other like:

Q=query, D=database

```
Fit 1  Q  -----
```

```

D      -----
Fit 2  Q -----
D      -----
Fit 3  Q -----
D      -----
<etc.>
Fit N-2 Q -----
D -----
Fit N-1 Q -----
D -----
Fit N  Q -----
D -----

```

In this approach, initial superimpositions are calculated using many more fractions of query and database structure, making detection of weak similarities more likely.

The residues that are equivalenced by either FAST or SLOW procedures are used to perform an initial fit, which is refined by the conformation-based and distance-based fit used during PAIRWISE/TREEWISE comparison of distantly related structures. If a high enough similarity score (S_c) is found after these three steps, then the transformation is saved for further analysis. The output from SCAN mode is directly readable by STAMP so that once a list of domains similar to one's query is obtained, multiple alignment (ie. PAIRWISE and TREEWISE) can be performed.

The program PDBC can be used to generate a list of protein domains given a set of PDB identifier codes, and the program SORTTRANS can be used to sort the output from SCAN, and remove any redundancies.

The S_c values output in SCAN mode differ slightly from those output during a PAIRWISE comparison. The correction introduced to correct the SW Score according to the length of the sequence lengths is removed. During multiple alignment the start and end points of the domains to be superimposed should be known; thus one can penalise all positions which are not involved in the alignment. During a scan, however, it is desirable to detect sub alignments of the two structures being compared. Thus, the S_c for scanning may be defined

in one of three ways (a=query, b=database, p=path, i=insertion, L=length):

Scheme 1

$$S_c = \left(\frac{S_p}{L_p} \right) \left(\frac{L_p - i_a}{L_a} \right) \left(\frac{L_p - i_b}{L_b} \right)$$

As for multiple structure alignment. As discussed, this is generally not the best way to compare a query to the database, since one would not usually wish to penalise insertions or omitted missing segments within the database structure (due to truncation values, etc.). However, this scheme may be useful if one is scanning a database of structures known to exhibit a particular fold (i.e., if one is merely after accurate superimpositions for a family of known structures; see Chapter 2).

Scheme 2

$$S_c = \left(\frac{S_p}{L_p} \right) \left(\frac{L_p - i_a}{L_p} \right) \left(\frac{L_p - i_b}{L_p} \right)$$

L_a and L_b have been replaced by L_p to removed any dependence on query or database structure length. The second two terms lower the score if gaps in the path are placed in the query (a) or database structure (b). This avoids a consideration of length, but will allow short stretches of structural equivalences to score highly.

Scheme 3

$$S_c = \left(\frac{S_p}{L_p} \right) \left(\frac{L_p - i_a}{L_a} \right)$$

Only penalises insertions in the query sequence. If a small fraction of the query sequence is in the actual path, then S_c drops. This scheme is most useful if one wants only similarities to the entire protein under consideration, since it penalises any omissions from the query structure.

Scheme 4

$$S_c = \left(\frac{S_p}{L_p} \right) \left(\frac{L_p - i_b}{L_b} \right)$$

The opposite of 3. Only penalises insertions in the database sequence. If a small fraction of the database sequence is in the actual path, then S_c drops. This scheme may be useful if one is scanning with a collection of secondary structure elements, since gaps are to be expected within the query (i.e. since the loops have been omitted).

Scheme 5

$$S_c = \left(\frac{S_p}{L_p} \right)$$

Raw score, no length requirement, will report even short alignments between similar sub-structures. This scheme may be useful for the search for short stretches of structural similarity, such as supersecondary structures.

Scheme 6

$$S_c = \left(\frac{S_p}{L_a} \right) \left(\frac{L_a - i_a}{L_a} \right)$$

Vaguely similar to Scheme 3, but this only scores hits favourably if they involve a significant fraction of the query structure (i.e. similarities only containing part of the query will not stand out). This is useful when one is comparing a particular domain to a database and is not interested in local similarities. This is the default for scanning.

For the most part, all of these scoring schemes will yield similar numbers for very similar structures. However, when more distantly related structures are compared, it becomes more useful to use a scheme specific to the particular problem (i.e., whether one wishes to scan with secondary structures only, when one is after only very similar structures, etc.).

Schemes are specified by the STAMP parameter SCANSORE (see below). If you're not sure which scoring scheme to use then you should just use the default scheme.

1.5.5 Displaying STAMP output (VER2HOR, DSTAMP, GSTAMP)

There are several ways to view STAMP alignments in a more readable form than the BLOCK output format. The output files can be read by the

Jalview alignment editor (available from <http://www.jalview.org>). Alternatively, BLOCK format files can be converted to other alignment formats using the ACONVERT tool as discussed previously.

In addition, the STAMP package contains additional programs for converting STAMP output into a horizontal alignment format and creating input files suitable for creating figures of multiple sequence and structure alignments: VER2HOR reads STAMP output files and displays the alignments in a horizontal text format.

DSTAMP converts BLOCK format files into input for GJB's program ALSCRIPT, which can then be used to generate figures of alignments in Postscript format. Details of DSTAMP and how to obtain ALSCRIPT are given in CHAPTER VI. DSTAMP determines reliable regions given a set of criteria, and highlights sequence and secondary structure accordingly.

GSTAMP reads STAMP outputs and creates input suitable for creating molecular graphics figures using Per Kraulis' program MOLSCRIPT. One simply runs TRANSFORM on a STAMP alignment file, and then run GSTAMP on the same file to create a MOLSCRIPT input file which will produce separate Postscript files for each aligned structure, with structurally equivalent regions shown as ribbons, the rest as C_α trace.

1.6 Some comments on interpreting structural similarities

There is a wealth of literature on the nature of protein structural similarities, and this manual is not the place to review them. If you want to look into the subject, then I would refer you to some of my papers [11, 13, 14] and references therein.

An important aspect of assessing the meaning of structural similarity is discerning whether a similarity between proteins in the absence of obvious sequence identity implies a common evolutionary ancestor, and usually an associated similarity in molecular function. Some studies have found that it is possible to discern homology by the analysis of the sequence identity calcu-

lated following protein structure alignment. Note that this is a very different identity than that quoted during typical sequence comparison (e.g. BLAST, FASTA, SSEARCH, etc). During sequence comparison, the reported % identity is the result of optimising the alignment of two sequences, thus numbers as high as % 20-30 are possible for proteins that are definitely not homologous (i.e. those having different tertiary folds). However, if an alignment has been derived without consideration of the amino acid sequence, then lower % identities can still be significant. See Russell *et al.* , 1997, and Murzin, 1993 for examples, and more details.

STAMP reports both the % identity from structure comparison, defined as the percentage residue identities (m) within structurally equivalent residues (n), and an estimate of the statistical significance (reported as $P(m)$) of a given a particular combination of m and n . The latter is described in Murzin (1993). Values of $P(m)$ smaller than about 10^{-3} very often indicate that the pair of proteins belong to the same protein superfamily, which implies a common ancestor, and more importantly very often indicates a similarity in molecular function. Specifically, the $P(m)$ is calculated for a $p = 0.1$; please see Murzin (1993) for a more thorough explanation of this calculation.

1.7 The programs contained within the package

STAMP consists of the main program (usually referred to as STAMP) and several sub-programs. Briefly, the programs are:

STAMP	Main program, does PAIRWISE, tree construction, TREEWISE and SCAN modes.
ALIGNFIT	Given a list of domains and a multiple sequence alignment outputs an initial transformation.
PDBC	Finds and reports information about PDB files given a four (PDB) code and/or chain identifier.
TRANSFORM	Given a list of transformations, outputs the corresponding set of coordinates.
SORTTRANS	Sorts the output from SCAN, and removes repeated transformations.
PDBSEQ	Given a list of domains, extracts the corresponding sequences from the PDB files.
VER2HOR	Given a STAMP alignment file, outputs an easy to read text version of the alignment for quick analysis.
DSTAMP	Given a STAMP alignment file, outputs commands for GJB's program ALSCRIPT (alignment to Postscript).
GSTAMP	Given a STAMP alignment file, outputs commands for Per Kraulis' MOLSCRIPT program.

The programs contained within the package (continued):

AVESTRUC	Given a STAMP alignment file, generates an average set of main chain or C alpha coordinates for the structural family.
POSTSTAMP	Reanalyses a STAMP alignment file (provides a more accurate set of equivalences for alignments of more than one structure).
PICKFRAME	Given a transformation, transforms all other domains onto another (specified by the user).
MERGETRANS	Given two transformation files, merges them by centering on a common identifier, either the first common one found or one specified by the user.
MERGESTAMP	Given two files containing alignments or transformations or both merges them by centering on a common identifier, either the first common one found or one specified by the user.
EXTRANS	Given a transformation file and a list of domain identifiers it will output a new transformation file containing only the domains given
ACONVERT	General alignment format conversion utility.
STAMP_CLEAN	Tidies up STAMP alignments to remove nonsensical gaps

Chapter 2

Worked examples

The examples described below show how to apply STAMP to particular problems. However, by far the most common thing you might want to do with STAMP is to pairwise or multiply align a set of structures. Originally, STAMP was intended to align closely similar structures that you had already collected together - for this reason, we talk first about this way of using STAMP (See section:2.2). We recommend you work through this example as it introduces STAMP concepts and files. However after developing the SCAN option in STAMP we found that the best way to start a multiple structure alignment was to gather structures or structural domains firstly by carrying out a STAMP scan. So, we recommend that unless you already have a set of very similar structures to align, you follow the SCAN example (See section: 2.4) to gather and orientate domains for multiple structure alignment.

2.1 Setup of examples

The installation and configuration of STAMP are explained in Chapter 5. To run the worked examples you need only install STAMP and set the STAMPDIR environment variable as described below. It is not necessary to edit the configuration files.

All example output files may be found in the directory examples/ in the STAMP installation directory. There are four sub-directories in the exam-

ples directory corresponding to each of the four protein structure families discussed in the examples below (s_prot/, ac_prot/, ig/, globin/).

Before beginning you should set the STAMPDIR environment variable to the name of the directory containing the various STAMP defaults files. This directory is called defs/ within the STAMP installation directory. It contains files containing default parameters and rules for finding PDB and DSSP files. It is not necessary to edit these files to run the examples. The files are as follows:

1. STAMPDIR/stamp.defaults contains default values for STAMP command line parameters. Values in this file can be overridden by specifying options on the command line.
2. The file STAMPDIR/pdb.directories file contains a set of patterns which STAMP uses to find PDB files. This makes it possible for STAMP to use PDB files located in multiple directories and having various combinations of filename prefixes and suffixes.
3. STAMPDIR/dssp.directories is similar to pdb.directories but is used by STAMP to locate files containing secondary structure assignments for PDB structures generated using the DSSP program.

2.2 Multiple alignment using an initial multiple sequence alignment

Mammalian and Bacterial Serine Proteinases

(This example is discussed in Russell & Barton, (1992).)

Despite a pronounced functional similarity (a highly conserved catalytic triad), this family of proteins shows little overall sequence similarity. Indeed, sequence alignment methods generally fail to provide an accurate alignment of these protein sequences. In situations like these, STAMP can be used to provide an accurate alignment of protein sequences based on a comparison of 3D structure. This can often reveal regions of weak sequence similarity that are not detectable during a comparison of sequence. The files for this example are in the directory examples/s_prot in the STAMP installation directory.

The procedure in this example is to create a multiple sequence alignment which is fed into the ALIGNFIT program to create an initial rough multiple structure alignment which can then be refined by STAMP.

The list of the domains to be aligned is given in the file `s_prot.domains`. The sequences are extracted from the PDB files by using the domain file with PDBSEQ:

```
pdbseq -f s_prot.domains > s_prot.seqs
```

This produces the file `s_prot.seqs`. This file is used to generate a multiple alignment using AMPS, the alignment being stored in the file `s_prot_amps.align`. This file is in AMPS format, which is the only format that ALIGNFIT can read. However, alignments in other formats can be converted to AMPS format using the Jalview (www.jalview.org) or ACONVERT program. For example, if the alignment had been in Clustal W format it could have been converted by running:

```
aconvert -in c -out b < sprot_clw.aln > s_prot_amps.align
```

Running:

```
aconvert -h
```

will list the command-line arguments that ACONVERT accepts.

Now that we have the multiple alignment, we can run ALIGNFIT on it:

```
alignfit -f s_prot_amps.align -d s_prot.domains -out s_prot_alignfit.trans
```

giving the output:

```
ALIGNFIT R.B. Russell 1995
Reading in block file...
Blocfile read: Length: 261
Reading in coordinate descriptions...
Reading coordinates...
Checking for inconsistencies...
Doing pairwise comparisons...
Doing treewise comparisons...
ALIGNFIT done.
Look in the file s_prot_alignfit.trans for output and details
```

The final transformation is in the file s_prot_alignfit.trans.

This provides an initial set of transformations for use by STAMP. To run STAMP type:

```
stamp -l s_prot_alignfit.trans -prefix s_prot
```

This should produce the following output on the terminal:

STAMP Structural Alignment of Multiple Proteins

Version 4.4 (May 2010)

by Robert B. Russell & Geoffrey J. Barton
Please cite PROTEINS, v14, 309-323, 1992

Sc = STAMP score, RMS = RMS deviation, Align = alignment length
Len1, Len2 = length of domain, Nfit = residues fitted
Secs = no. equivalent sec. strucs. Eq = no. equivalent residues
%I = seq. identity, %S = sec. str. identity
P(m) = P value (p=1/10) calculated after Murzin (1993), JMB, 230, 689-694
(NC = P value not calculated - potential FP overflow)

	No.	Domain1	Domain2	Sc	RMS	Len1	Len2	Align	NFit	Eq.	Secs.	%I	%S	P(m)
Pair	1	4cha	3est	7.74	1.15	239	240	242	217	214	20	41.59	84.58	1.32e-33
Pair	2	4cha	2ptn	7.81	0.97	239	223	234	206	203	20	47.78	91.63	8.32e-43
Pair	3	4cha	1ton	6.92	1.19	239	227	241	191	189	19	40.21	89.42	8.18e-28
Pair	4	4cha	3rp2a	7.46	1.09	239	224	235	203	199	18	37.19	85.43	6.38e-24
Pair	5	4cha	2pkaab	7.28	1.14	239	232	241	203	202	20	37.13	90.10	6.62e-25
Pair	6	4cha	1sgt	7.09	1.26	239	223	239	197	191	20	36.13	90.05	2.86e-22
Pair	7	4cha	2sga	3.64	1.66	239	181	240	109	101	15	28.71	84.16	8.84e-08
Pair	8	4cha	3sgbe	3.62	1.56	239	185	240	105	95	15	26.32	85.26	3.48e-06

<etc.>

Reading in matrix file s_prot.mat...

Doing cluster analysis...

Cluster: 1 (2ptn & 2pkaab) Sc 8.50 RMS 1.07 Len 232 nfit 216

See file s_prot.1 for the alignment and transformations

Cluster: 2 (2sga & 3sgbe) Sc 8.36 RMS 0.65 Len 191 nfit 166

See file s_prot.2 for the alignment and transformations

Cluster: 3 (1ton & 2ptn 2pkaab) Sc 9.03 RMS 0.73 Len 239 nfit 205

See file s_prot.3 for the alignment and transformations

Cluster: 4 (3rp2a & 1ton 2ptn 2pkaab) Sc 8.73 RMS 0.93 Len 242 nfit 206

See file s_prot.4 for the alignment and transformations

Cluster: 5 (3est & 3rp2a 1ton 2ptn 2pkaab) Sc 8.51 RMS 1.13 Len 258 nfit 208

See file s_prot.5 for the alignment and transformations

Cluster: 6 (4cha & 3est 3rp2a 1ton 2ptn 2pkaab) Sc 8.18 RMS 1.01 Len 260 nfit 201

See file s_prot.6 for the alignment and transformations

Cluster: 7 (2alp & 2sga 3sgbe) Sc 8.35 RMS 1.06 Len 203 nfit 168

See file s_prot.7 for the alignment and transformations


```

Cluster: 8 ( 1sgt & 4chaa 3est 3rp2a 1ton 2ptn 2pkaab ) Sc 7.70 RMS 1.11 Len 267 nfit
See file s_prot.8 for the alignment and transformations
Cluster: 9 ( 1sgt 4chaa 3est 3rp2a 1ton 2ptn 2pkaab & 2alp 2sga 3sgbe ) Sc 4.7
See file s_prot.9 for the alignment and transformations

```

The various fields describe details of the pairwise and treewise comparisons: S_c , RMS deviation, the alignment length (Align), the length of each structure in residues (Len1, Len2), the number of atoms used in the RMS fit (Nfit), the number of equivalent secondary structure elements (Secs), and the number of equivalent residues (see above, Eq.).

STAMP will also produce several files:

s_prot.mat – a file containing the information used to derive the structural similarity tree (i.e. the output from the PAIRWISE) mode. This is an upper diagonal matrix containing the pairwise S_c values.

s_prot.N – a series of files containing transformations and alignments created by running the TREEWISE mode in STAMP. Each file corresponds to a *node* in the similarity tree (i.e. a *cluster*), where two groups of one or more structures have been combined to form an alignment and transformations. The higher the value of N the more structurally dissimilar the proteins contained in the file are. Highly similar structures are clustered (aligned/superimposed) at an early stage in the program's run, with more distantly related structures being clustered towards the end.

The top of each s_prot.N file contains the information needed to generate superimposed coordinates using TRANSFORM. For example, running:

```
transform -f s_prot.9 -g -o s_prot.pdb
```

will create a PDB file containing all of the structures from the alignment in s_prot.8.

After these details, various details of the similarity (RMS deviation, S_c value, etc) are given. The bottom portion of the file contains the structural alignment in STAMP format. Positions that do not include gaps contain information as to the degree of local structural similarity, such as the distance between (averaged) C_α atoms, and the P'_{ij} value.

Methods for displaying sequence alignments and structures are described below.

2.3 Database Scanning

Database scanning within STAMP is unpublished, apart from a brief description in a figure legend [16], but it has been fairly well tested since version 2.0. Indeed, two novel similarities have resulted in publications [9, 16].

Immunoglobulin domain

One example of a scan is given. The light chain variable domain of the immunoglobulin 2FB4 is used to scan a small database of other protein domains containing both a diverse collection of related folds (greek key folds, including azurin, superoxide dismutase, CD4, etc.), and completely unrelated folds (such as globins). See the directory examples/ig for this example.

The 2FB4 domain is described in 2fb4lv.domain. To scan this against the database type:

```
stamp -l 2fb4lv.domain -s -n 2 -slide 5 -prefix 2fb4lv_stamp -d some.domains -cut
```

‘-s’ specifies the SCAN mode ‘-slide’ describes how many residues to slide the query sequence (2fb4lv) along each sequence in the file some.domains to provide each initial fit (i.e. the sequence of 2fb4lv is layed on top of each database sequence at postions 1, 6, 11, etc.). ‘-cut’ tells the program to cut down each domain read in from some.domains according to where the similarity is found. If it is not specified, the output will contain domain descriptors identical to those found in ‘some.domains’. When one is comparing a single-domain query to a database structure having multiple domains, it is desirable to do this. Try running it both ways (with and without -cut) and look at the output to see the difference. (e.g. CHAIN A is converted to A 1 _ to A 60 _ in one descriptor in the SCAN output and A 120 _ to A 175 _ in another, since there are two repeats of the query domain in the database structure).

The above run should write the following to the standard output (again, ignoring the header):

STAMP Structural Alignment of Multiple Proteins

Version 4.4 (May 2010)

by Robert B. Russell & Geoffrey J. Barton
Please cite PROTEINS, v14, 309-323, 1992

Results of scan will be written to file 2fb4lv_stamp.scan

Fits = no. of fits performed, Sc = STAMP score, RMS = RMS deviation

Align = alignment length, Nfit = residues fitted, Eq. = equivalent residues

Secs = no. equiv. secondary structures, %I = seq. identity, %S = sec. str. identity

P(m) = P value (p=1/10) calculated after Murzin (1993), JMB, 230, 689-694

(NC = P value not calculated - potential FP overflow)

	Domain1	Domain2	Fits	Sc	RMS	Len1	Len2	Align	Fit	Eq.	Secs	%I	%S	P(m)
Scan	2fb4lv	2fb4lc	1	4.317	2.120	111	105	127	55	46	8	10.87	78.26	1.00e+00
Scan	2fb4lv	2fb4l	1	9.799	0.001	111	166	111	111	111	11	100.00	97.30	0.00e+00
Scan	2fb4lv	1mcplv	1	7.848	1.165	111	113	116	96	95	0	49.47	40.00	2.05e-22
Scan	2fb4lv	1mcphv	1	6.921	1.500	111	122	126	85	81	0	30.86	34.57	1.44e-07
Scan	2fb4lv	1cmsC	1	2.507	1.639	111	148	157	28	24	4	4.17	62.50	1.00e+00
Scan	2fb4lv	3cd4	1	5.939	1.334	111	166	114	78	75	12	20.00	76.00	4.10e-03
Scan	2fb4lv	2hhbb	0	0.000	100.000	111	146	0	0	75	0	0.00	0.00	1.00e+00
Scan	2fb4lv	3dpa	0	0.000	100.000	111	166	0	0	75	0	0.00	0.00	1.00e+00
Scan	2fb4lv	3sgbe	0	1.940	2.313	111	166	204	25	17	3	5.88	88.24	1.00e+00
Scan	2fb4lv	1acx	1	4.152	2.454	111	108	133	57	43	4	16.28	72.09	7.26e-02
Scan	2fb4lv	2abxa	0	0.000	100.000	111	74	0	0	43	0	0.00	0.00	1.00e+00
Scan	2fb4lv	1101	0	0.000	100.000	111	164	0	0	43	0	0.00	0.00	1.00e+00
Scan	2fb4lv	2azaa	1	4.063	2.463	111	129	134	49	35	5	14.29	82.86	1.00e+00
Scan	2fb4lv	1rnt	0	1.503	2.545	111	104	148	17	13	3	15.38	69.23	1.00e+00
Scan	2fb4lv	2sodo	1	3.611	2.365	111	151	158	42	32	8	9.38	71.88	1.00e+00
Scan	2fb4lv	2pcy	1	3.788	2.052	111	99	125	47	39	6	30.77	79.49	2.27e-04
Scan	2fb4lv	8atca	0	0.000	100.000	111	166	0	0	39	0	0.00	0.00	1.00e+00

See the file 2fb4lv_stamp.scan

where all of the fields are as for the PAIRWISE mode, save for Fits, which indicates the number of fits that were saved to the file '2fb4lv_stamp.scan'. Note that for domain descriptors (see some.domains) containing two Ig type folds (e.g. 2fb4l, 1cd4, etc.) that more than one fit has been saved, since the search found both of the Ig type folds in each of these two proteins. Not also that 'Fits' is zero for several of the examples, indicating that the no similarity was found within these proteins. Where more than one Fit is output for a domain in the database, the best S_c , RMS etc. are reported.

2fbjlv_stamp.scan will contain all the transformations output during the scan.

Several of these will be redundant, since it is possible for a particular match to be found twice. To remove repeated transformations, or those not considered interesting, run the program SORTTRANS on the output.

```
sorttrans -f 2fb4lv_stamp.scan -s Sc 2.0 > 2fb4lv_stamp.sorted
```

This sorts the input file by S_c values, and leaves only those non-redundant domain descriptions having an $S_c \geq 2.0$. A cutoff of 2.0 is generally a good choice pairwise comparisons with a score lower than this tend to produce poor quality alignments.

```
sorttrans -f 2fb4lv_stamp.scan -s rms 1.5 > 2fb4lv_stamp.sorted
```

sorts the input file by RMSD values, and leaves only those domain descriptions having an RMSD ≤ 1.5 Å. Despite its predominance in the literature, RMSD is not a very good means of measuring structural similarity, since low RMSDs can usually be obtained for any two structures if one considers a small enough set of residues.

```
sorttrans -f 2fb4lv_stamp.scan -s nfit 40 > 2fb4lv_stamp.sorted
```

sorts the input file by the number of atoms used in the final fitting, and leaves only those domain descriptions where $nfit \geq 40$.

```
sorttrans -f 2fb4lv_scan -s n_sec 6 > 2fb4lv_stamp.sorted
```

sorts the input file by the number of equivalent secondary structures, and leaves only those having 6 or more secondary structures equivalent.

Combinations of these can be used to select out interesting domains from a scan output. Probably the best combination involves S_c and $nfit$ (ie. score and $nfit$), since large structures can give fortuitously large S_c values with very few fitted atoms.

The final output is in the file 2fb4lv_stamp.sorted. This is the result of

the first example (i.e. -s Sc 2.0). Note that several structures similar to the Ig type domain have been detected, and appear (according to S_c) in the order one might expect from knowledge of the 3D structures, sequences and functions of these proteins.

The output from scanning can be used as input for other modes of the program. Once you have performed a scan, and have sorted the 'hits' down to an interesting set, you can then use the output from scan as the input for a multiple alignment. This is discussed in the next section.

2.4 Using SCAN mode as the starting point for multiple alignment

In certain instances initial fits based on multiple sequence alignment will be far from accurate, such that even an initial conformation based fit will not be able to correct the poor initial superposition, and even genuine structural homology will be missed. In these instances it is possible to make use of the SCAN mode to provide a more accurate initial superimposition.

To do this one need only select one representative of the domains to be superimposed and use this domain in a sensitive scan *of the other domains*. By applying the same techniques as used for the scan with the Ig light variable domain (see the previous section) one can create a set of transformations of the searched domains onto the query domain. This set of transformations constitutes a rough multiple structure alignment which can be used by STAMP as the starting point for an accurate alignment.

Aspartic Proteinase Domains

The output files for this example are in the directory examples/ac_prot. In this example the aspartyl proteinase N- and C-terminal lobes are aligned. The N-terminal domain of 1CMS (in the file 1cmsN.domain) is used as the query domain to scan a list of aspartyl proteinase N- and C-terminal do-

mains (ac_prot.domains). Running:

```
stamp -l 1cmsN.domain -n 2 -s -slide 5 -d ac_prot.domains -prefix ac_prot
```

should produce:

STAMP Structural Alignment of Multiple Proteins

Version 4.4 (May 2010)

by Robert B. Russell & Geoffrey J. Barton
Please cite PROTEINS, v14, 309-323, 1992

Results of scan will be written to file ac_prot.scan

Fits = no. of fits performed, Sc = STAMP score, RMS = RMS deviation

Align = alignment length, Nfit = residues fitted, Eq. = equivalent residues

Secs = no. equiv. secondary structures, %I = seq. identity, %S = sec. str. identity

P(m) = P value (p=1/10) calculated after Murzin (1993), JMB, 230, 689-694

(NC = P value not calculated - potential FP overflow)

	Domain1	Domain2	Fits	Sc	RMS	Len1	Len2	Align	Fit	Eq.	Secs	%I	%S	P(m)
Scan	1cmsN	1cmsN	1	9.800	0.000	175	175	175	175	175	18	100.00	94.86	0.00e+00
Scan	1cmsN	1cmsC	1	3.214	2.065	175	148	204	68	62	11	19.35	83.87	1.11e-02
Scan	1cmsN	4apeN	1	8.209	1.300	175	178	182	162	159	15	30.82	87.42	2.82e-13
Scan	1cmsN	4apeC	1	3.420	1.948	175	152	205	70	67	13	14.93	79.10	6.11e-02
Scan	1cmsN	3appN	1	7.976	1.280	175	174	183	157	157	19	29.30	90.45	1.01e-11
Scan	1cmsN	3appC	1	3.269	2.068	175	149	205	68	59	12	13.56	81.36	1.00e+00
Scan	1cmsN	2aprN	1	8.435	1.075	175	178	178	164	162	15	33.33	85.80	4.60e-16
Scan	1cmsN	2aprC	1	3.304	1.973	175	147	200	67	64	13	18.75	76.56	1.37e-02
Scan	1cmsN	4pepN	1	8.836	0.930	175	173	174	169	169	15	57.99	83.43	3.00e-53
Scan	1cmsN	4pepC	1	3.223	2.105	175	152	206	68	57	10	21.05	85.96	6.17e-03

See the file ac_prot.scan

The file ac_prot.scan will contain all 10 domains superimposed onto 1cmsN. Note that we haven't run the program with the '-cut' option, since the file ac_prot.domains contains an assignment of domains. Running SORTTRANS removes any redundancies:

```
sorttrans -f ac_prot.scan -s Sc 2.5 > ac_prot.sorted
```

and running STAMP will generate the multiple alignment as described for the examples above.

```
stamp -l ac_prot.sorted -prefix ac_prot
```

2.5 Generating a set of superimposed structures

The TRANSFORM program included in the STAMP package can be used to generate sets of superimposed structures from the output of a STAMP multiple alignment run. For example, the output from the multiple alignment in the previous section can be fed into TRANSFORM as follows:

```
transform -f ac_prot.9 -g -o ac_prot.pdb
```

This will read in the files, transform the coordinates and save them to the file ac_prot.pdb (with each chain labelled starting with a different letter).

2.6 Alignment without an initial multiple alignment using ROUGHFIT

This method described in this section, where the ROUGHFIT mode is used to create an initial alignment, was the one originally used in STAMP for cases where an initial multiple sequence alignment was not available. The SCAN mode (see the aligning section on aligning protease domains) now makes it possible to create a reasonable starting alignment even for cases where an accurate alignment based on sequence is impossible. Apart from cases where the structures are homologous or of very similar length, using the SCAN mode generally produces better results than using ROUGHFIT. Accordingly, ROUGHFIT is deprecated in favour of using SCAN mode as a starting point. It is documented here for the sake of completeness.

This method avoids having to create an initial multiple sequence alignment and tends to work for homologous proteins, or those having very similar lengths despite no sequence similarity.

Globins

Since the globin sequences are of similar length an initial superimposition accurate enough to proceed with STAMP can be obtained by merely aligning the N-terminal ends of the sequences and using whatever equivalences result to obtain an initial superimposition. The command ROUGH (ROUGHFIT

procedure) is used. In addition, an initial conformation based fit is performed in order that any inaccuracies in this initial superimposition may be corrected. See the directory examples/globins.

To run STAMP in this example, type:

```
stamp -l globin.domains -rough -n 2 -prefix globin
```

This should produce the following on the standard output (ignoring the header):

STAMP Structural Alignment of Multiple Proteins

Version 4.4 (May 2010)

by Robert B. Russell & Geoffrey J. Barton
Please cite PROTEINS, v14, 309-323, 1992

Running roughfit.

Sc = STAMP score, RMS = RMS deviation, Align = alignment length
Len1, Len2 = length of domain, Nfit = residues fitted
Secs = no. equivalent sec. strucs. Eq = no. equivalent residues
%I = seq. identity, %S = sec. str. identity
P(m) = P value (p=1/10) calculated after Murzin (1993), JMB, 230, 689-694
(NC = P value not calculated - potential FP overflow)

No.	Domain1	Domain2	Sc	RMS	Len1	Len2	Align	NFit	Eq.	Secs.	%I	%S	P(m)
Pair 1	2hhbb	2hhba	8.19	1.38	146	141	147	136	135	7	44.44	82.96	4.82e-25
Pair 2	2hhbb	2lhb	7.11	1.39	146	149	151	127	127	7	26.77	86.61	4.90e-08
Pair 3	2hhbb	4mbn	8.06	1.38	146	153	151	141	139	8	25.18	87.05	1.57e-07
Pair 4	2hhbb	1ecd	6.89	2.04	146	136	144	127	119	7	20.17	86.55	3.91e-04
Pair 5	2hhbb	1lh1	5.89	2.39	146	153	155	120	110	6	17.27	80.91	6.62e-03
Pair 6	2hhba	2lhb	6.54	1.66	141	149	150	122	119	7	34.45	88.24	3.97e-13
Pair 7	2hhba	4mbn	7.78	1.39	141	153	148	136	133	8	27.07	87.97	1.51e-08
Pair 8	2hhba	1ecd	6.61	2.18	141	136	145	124	118	8	17.80	87.29	3.47e-03
Pair 9	2hhba	1lh1	5.95	2.20	141	153	153	117	106	6	14.15	82.08	4.45e-02
Pair 10	2lhb	4mbn	7.13	1.23	149	153	149	131	130	8	25.38	90.77	2.82e-07
Pair 11	2lhb	1ecd	6.42	1.93	149	136	145	124	123	8	18.70	87.80	1.34e-03
Pair 12	2lhb	1lh1	5.74	2.11	149	153	155	117	105	6	19.05	85.71	2.04e-03
Pair 13	4mbn	1ecd	7.46	1.64	153	136	145	134	132	8	21.21	87.88	6.21e-05
Pair 14	4mbn	1lh1	6.76	2.35	153	153	155	135	133	6	17.29	84.21	3.35e-03
Pair 15	1ecd	1lh1	5.96	2.59	136	153	149	121	114	6	15.79	87.72	1.62e-02

Reading in matrix file globin.mat...

Doing cluster analysis...

Cluster: 1 (2hhbb & 2hhba) Sc 8.19 RMS 1.38 Len 147 nfit 136

See file globin.1 for the alignment and transformations

Cluster: 2 (4mbn & 2hhbb 2hhba) Sc 8.96 RMS 1.31 Len 151 nfit 138

See file globin.2 for the alignment and transformations

Cluster: 3 (1ecd & 4mbn 2hhbb 2hhba) Sc 8.35 RMS 1.81 Len 146 nfit 128

See file globin.3 for the alignment and transformations


```

Cluster: 4 ( 2lhb & 1ecd 4mbn 2hhbb 2hhba ) Sc 8.24 RMS 1.23 Len 152 nfit 120
See file globin.4 for the alignment and transformations
Cluster: 5 ( 1lh1 & 2lhb 1ecd 4mbn 2hhbb 2hhba ) Sc 7.70 RMS 2.46 Len 160 nfit 121
See file globin.5 for the alignment and transformations

```

where the output and files are as described for the serine proteinase example above, with ‘s_prot’ replaced with ‘globin’.

-rough performs the initial superimpositions (ROUGHFIT) and -n 2 means that the conformation biased fit will be performed before the final fit. This conformation biased fit is usually necessary when the initial superimpositions are approximate.

ROUGHFIT will not always work. Note that in this example all the pairwise S_c values are above 5.6, suggesting strong structural similarity. If when using the ROUGHFIT option you find low S_c values (the program will flag the values with the message ‘LOW SCORE’), this usually means that ROUGHFIT hasn’t managed to generate a good enough starting superimposition, and you should try using SCAN mode to generate an initial alignment, as described in the previous section.

2.7 Protein domain databases

The program PDBC may be used to output a set of STAMP readable domain descriptions. Given a list of four letter brookhaven codes and an optional set of chains. This will only work if you have a suitable ‘pdb.directories’ file. See the chapter on installation for details on how to do this.

```

pdbc -d 2hhba >! globin_fold.domains
pdbc -d 2hhbb >> globin_fold.domains
pdbc -d 4mbn >> globin_fold.domains
pdbc -d 1lh1 >> globin_fold.domains
pdbc -d 1cola >> globin_fold.domains
pdbc -d 1cpca >> globin_fold.domains

```

will produce the following output (ignoring comments, which are specified by a ‘%’ in column 0):

```

/(PDB PATH)/pdb2hhb.ent 2hhba { CHAIN A }
/(PDB PATH)/pdb2hhb.ent 2hhbb { CHAIN B }

```

```
/(PDB PATH)/pdb4mbn.ent 4mbn { ALL }  
/(PDB PATH)/pdb1lh1.ent 1lh1 { ALL }  
/(PDB PATH)/pdb1col.ent 1cola { CHAIN A }
```

Where (PDB PATH) denotes the location of the relevant PDB file on your system. Note that your PDB files may be called (code).pdb instead, or may follow some other convention. This is OK, see Chapter 5 (installation) for details as to setting this up.

Note that there doesn't need to be a filename in the domain file. One can merely leave it as 'Unknown' or some other string (i.e. not empty spaces), and the programs will try and find where the file corresponding to the four letter code is on your system. In other words, the files given in this distribution should work on your system, provided that you have all the PDB files.

Note that PDBC can be used to probe information about a PDB entry by using the '-q' option. Try it and see. This is a good test of whether STAMP has been set up properly on your system. If you just want to test where STAMP is looking for PDB and DSSP files, then use the '-m' (minimal) options. This just reports PDB/DSSP files if found and exits.

STAMP database comparisons are computationally intensive, so it is prudent to avoid comparisons that are redundant (e.g. multiple mutants or binding studies of the same protein, T4 lysozyme for example).

The STAMP distribution contains a series of non-redundant databases derived by a parsing of the SCOP database. These are located in the 'STAMPDIR' directory. The files are derived from SCOP release 1.75. The files were created using the scop2stamp program, which can be found in the 'bin/' directory of the STAMP installation. Running 'scop2stamp' without arguments will list the options that this program accepts.

Domain database	N	Description
scop.dom	109747	All PDB domains classified in SCOP
scop_species.dom	13816	One representative per species of each SCOP protein.
scop_prot.dom	9621	One representative of each SCOP protein
scop_fam.dom	3883	One representative of each SCOP family
scop_supf.dom	1950	One representative of each SCOP superfamily
scop_fold.dom	1190	One representative of each SCOP fold

The complete set of SCOP domains contains a high degree of redundancy. The amount of time required to search it will depend on your particular system but you should expect it to take on the order of 20 hours of CPU time on the current generation of processors. If you have access to multiple CPUs, it is possible to divide the database into subsets, search the individual subsets in parallel on multiple CPUs, and then aggregate the search outputs into a single results file which can be filtered using SORTTRANS in the usual way.

2.8 Generating transformed coordinates using TRANSFORM

The program TRANSFORM can be used with any STAMP alignment file containing domain descriptions to output a set of PDB format files for display or further analysis. For example, running:

```
transform -f globin.5
```

should write the following to the standard output:

```
TRANSFORM R.B. Russell, 1995
Using PDB files
Files will not include heteroatoms
Files will not include waters
Domain 1, 1lh1 => to 1lh1.pdb
Domain 2, 2lhb => to 2lhb.pdb
Domain 3, 1ecd => to 1ecd.pdb
Domain 4, 4mbn => to 4mbn.pdb
Domain 5, 2hhbb => to 2hhbb.pdb
Domain 6, 2hhba => to 2hhba.pdb
```

A set of PDB format files containing the superimposed coordinates is generated. Running the program as shown above will produce one PDB file for

each domain identifier. If one wishes to look at the superimposed structures *together* (in the same file), then the option `-g` (i.e. graphics) can be used:

```
transform -f globin.5 -g -o globins.pdb
```

which should output the following:

```
TRANSFORM R.B. Russell, 1995
Using PDB files
Files will not include heteroatoms
Files will not include waters
All coordinates will be in file globins.pdb
Domain  1,  1lh1 => to globins.pdb (chain A)
Domain  2,  2lhb => to globins.pdb (chain B)
Domain  3,  1ecd => to globins.pdb (chain C)
Domain  4,  4mbn => to globins.pdb (chain D)
Domain  5,  2hhbb => to globins.pdb (chain E)
Domain  6,  2hhba => to globins.pdb (chain F)
```

This options puts transformed coordinates for each domain into one file (specified by `-o`, in this example it is 'globins.pdb'). Each domain will be labelled sequentially with a different chain identifier (i.e. A, B, C, etc.). Note that only 'globins.pdb' is included in the example directory.

By default, TRANSFORM does not include heteroatoms in the output. If you wish heteroatoms to be included, then add the `-het` option to the transform command. If you wish waters to be included in the file add the `-hoh` option. Note that heteroatoms/waters are sometimes included that fall outside the range of your domain descriptor. This may seem silly, but it is difficult to determine which heteroatoms are associated with which residues given PDB format.

2.9 Generating averaged coordinates

It may also be useful to have a set of *averaged* coordinates derived from a protein structural family. This makes it possible to see what portions of the structure are common to all members of the family (i.e. the common core). The program AVESTRUC takes the output from STAMP (i.e. an aligned family of protein structures), and generates a PDB file containing averaged

coordinates for the common core as identified by STAMP. For example, to generate the averaged coordinates for the aspartic proteinase domains one needs to type:

```
avestruc -f ac_prot.8 -o ac_prot_ave.pdb
```

The file `ac_prot_ave.pdb` will contain a set of averaged C_α atoms taken by averaging the coordinates for those positions within the file `ac_prot.8` that are found to be structurally equivalent. To obtain a poly Alanine set of coordinates (i.e. including main chain and C_β coordinates), type:

```
avestruc -f ac_prot.8 -o ac_prot_ave.pdb -polyA
```

Note that this will only work if all main chain atoms are found in the file (i.e. it won't work if the PDB files contain only C_α atoms).

A useful feature in AVESTRUC that was added in STAMP version 4.1 is the use of the `-ident` and `-cons` options. The program now labels all residues in the averaged model as 'UNK'. If positions are totally conserved across all structures in the averaged model, the '`-ident`' option will name residues accordingly. The `-cons` option will label residues additionally as conserved in character if all amino acids in the set have the following properties:

SMA	small
TIN	tiny
POL	polar
HYD	hydrophobic
POS	positive
NEG	negative
CHA	charged
ARO	aromatic
ALI	aliphatic
BRA	C_β branched

See Taylor (1986) for a description of amino acid properties.

Another feature is that the temperature factors reflect whether positions are

structurally conserved, or simply fortuitously aligned. If you add the option ‘-aligned’ to the command line, all positions that are not matched to a gap will be considered in the generation of the averaged model. If you then colour your model according to temperature in a structure viewer, the blue regions will correspond to those that are structurally equivalent (as you have defined or by default) whereas the red regions will show those that are simply in the same position in the sequence alignment.

2.10 Displaying/processing the output

2.10.1 POSTSTAMP

There is something inherently wrong with the way STAMP assigns equivalences within multiple alignments. It considers an average set of C_α coordinates and uses an average set of probabilities to derive equivalences when more than two structures are involved, and as a consequence, it sometimes appears to go wrong during this process. Usually this is only when very distantly related proteins are being considered. A fix to this problem is to consider each pair of structures within the alignment separately, and to recalculate the *raw* Rossmann and Argos probabilities. One need then define positions as structurally equivalent when *all* pairs of structures have a P_{ij} value larger than a cutoff at a particular residue position.

For example, for ten structures, there are $(10 \times 9/2) = 45$ pairs. For a position to be structurally equivalent across all members of the family, P_{ij} should be ≥ 0.5 for all 45 pairs.

POSTSTAMP does just this. It adds two new STAMP format fields to a STAMP alignment file: one tells whether the above is true (1) or false (0) for each position (i.e. is each position structurally equivalent across all members of the family); the second tells how many pairwise comparison have P_{ij} greater than or equal to the cutoff (e.g. 0.5).

For example,

```
poststamp -f globin.5 -min 0.5
```

Creates a file globin.5.post, containing the above data for a P_{ij} value of 0.5.

2.10.2 STAMP_CLEAN

When aligning more than one structure, STAMP will usually create alignments that are fairly meaningless within regions that are not structurally equivalent across all structures. Such regions may have meaning for particular sub-families of structures, but for the purposes of display, are nonsensical. STAMP_CLEAN is a useful program that takes a STAMP alignment file and ‘cleans up’ such gaps. To run the program, for example (using the POST-STAMP output file generated above):

```
stamp_clean globin.5.post 3 > globin.5.clean
```

will create a file globin.5.clean where all gaps not lying within structurally equivalent regions, and having fewer than 3 aligned residues in a row (i.e blocks where all sequences are not aligned with gap) are shortened to their minimum length.

2.10.3 Displaying text alignments

There are two ways to display STAMP alignments in a horizontal format. The first is simply to use ACONVERT to change the STAMP block file format into another format such as MSF or Clustal. The format would be:

```
aconvert -in b -out c < <stamp alignment file>
```

Where ‘-out c’ denotes Clustal format.

ACONVERT does not use any of the STAMP specific parts of the alignment (i.e reliable structural equivalences, etc.). There is a program specifically designed for displaying these data in a vertical format. VER2HOR takes a STAMP alignment file and outputs a horizontal text format. For example, to display the globin alignment, one needs to type:

```
ver2hor -f globin.5.clean
```


2.10.4 Pretty Alignments via ALSCRIPT

DSTAMP generates input files for GJBs ALSCRIPT program. Given a STAMP alignment file, DSTAMP can be run to create a fairly pretty alignment. Detailed descriptions of the parameters are given below. As a quick example, using the globin example,

```
dstamp -f globin.5.clean -prefix globin_align
```

will create a file called:

```
globin_align.als
```

Which contains a set of ALSCRIPT commands. To get a pretty Postscript alignment, one needs to run alscript:

```
alscript globin_align.als
```

The file globin_align.ps will be created, and is previewable or printable on a Postscript printer. And is shown in Figure 1.

By default, residues occurring within structurally equivalent regions will be boxed in the sequence alignment. Helices and strands will appear as cylinders and arrows (coil/turn regions are not shown). Conserved residues will be in inverse text, positions showing a conservation of polar character will be in bold, those showing conservation of hydrophobic character will be shaded and those showing a conservation of small size will be shown in a smaller font. It is possible to modify the output format (parameters are described in a Chapter 4). I would also recommend only using the DSTAMP output as a starting point, and refine the ALSCRIPT file yourself to give the best

alignment. The automated procedure can give some ugly results.

Figure 1 Globin alignment as discussed in the text.

2.10.5 Pretty Structures via MOLSCRIPT

GSTAMP can be used to display the structurally equivalences found by STAMP. It works by creating an input file for MOLSCRIPT [18] (contact Per Kraulis to obtain a copy).

As for DSTAMP, a detailed description of parameters is given later. Here is a quick example, using the first globin alignment (i.e. containing only two structures).

First one needs to generate transformed PDB coordinates using the program TRANSFORM:

```
transform -f globin.5.clean
```

This will create 2 PDB files with coordinates superimposed: 2hhbb.pdb and 2hhba.pdb.

```
gstamp -f globin.5.clean
```

This reads in the six structures and the alignment and outputs six molscript files called (domain identifier).molscript.

One must then run molscript on each of these files that one wants to display. For illustration, we will run two very distantly related globins:

```
molscript < 1lh1.molscript > 1lh1.ps  
molscript < 2hhba.molscript > 2hhba.ps
```

To give the two postscript files are shown in Figure 2.

Figure 2 Superimpositions of globin 1lh1 (left) and 2hhba (right).

By default, GSTAMP will show equivalent helix, strand and coil residues as MOLSCRIPT α helix, β strand and coil, with un-equivalent regions being shown as C_α trace.

At best, GSTAMP will give only a starting point for further refinement. Invariably, one will need to modify the orientation of the image for the best

view, and probably need to tweak the assignments of helix and strand to look clear; MOLSCRIPT will not work, for example, if one has very short β strands.

Chapter 3

Input and Output format for all programs

3.1 Describing domain structures

Every entry in a STAMP input file is called a ‘domain’. This term is a bit of a misnomer, since ‘domains’ needn’t be single domains (though it is usually best to do structure comparisons at the domain level).

The problem of defining domains such that a wide variety of possibilities may be used (e.g. all the coordinates in a PDB file, one chain, bits of one chain, two chains, one chain and bits of another, etc) is solved by defining a domain by: 1) a file, 2) an identifier, and 3) a list of ‘objects’, from the file, to be included in the domain. An object is defined as a run of C_{α} coordinates, and a domain may contain more than one object.

Domains are stored in STAMP in files which may contain one or more of such domain definitions.

The format of these files must be as follows:

```
<file name> <identifier label> { <objects> }
```

or,

```

<file name> <identifier label> { <objects> [RETURN]
R11 R12 R13   V1
R21 R22 R23   V2
R31 R32 R33   V3 }

```

<file name> is the full name (including path) of the PDB file in which the coordinate information is to be found. If you don't know the precise location of the file, then just call it UNK or something (i.e. not a blank), and the programs should be able to find the appropriate PDB file using the domain identifier field. Note that finding the PDB file using the identifier relies on a set of rules defined in the `pdb.directories` configuration file (see Chapter 5 for details of the location of the file, its format and how STAMP uses it to locate PDB files).

<identifier label> is a short name to be used by the program. eg. 4mbn1. The domain identifiers in a STAMP input file must be unique.

DSSP secondary structure files can only be found by STAMP by using the domain identifier in a similar fashion as described for PDB files above. Again, see Chapter 5 for details of how this works.

<objects> are coordinate descriptions, and may be one of three types:

1. ALL

all C_{α} 's from the file.

2. CHAIN X

only C_{α} 's labeled as chain X.

3. <chain1> <number1> <insert1> to <chain2> <number2> <insert2>
e.g. B 20 _ to B 67 P

only C_{α} 's between (and including) the two full brookhaven residue names (chain, number, insertion code; the '_' character denotes a space)

N.B. THERE MUST BE AT LEAST ONE SPACE BETWEEN THE VARIOUS FIELDS. Combinations of these are allowed within one domain, e.g. '

CHAIN A B 1 _ to B 65 _ ‘

$R_{11} \rightarrow R_{33}$ and $V_1 \rightarrow V_3$ are a rotation matrix and translation vector, respectively.

Thus, a full description of three domains might look something like this:

```
/data/newpdb/pdb/pdb1ton.ent 1ton { ALL
0.9876 0.34 0.543 19.23
1.0 2.34 0.98473332 1.0
0.023 0.94 4.345 20.0 }
/data/newpdb/pdb/pdb2kai.ent 2kai_Kallikrien { CHAIN X CHAIN Y }
/data/newpdb/pdb/pdb3sgb.ent 3sgbe_SGprotease { E 20 _ to E 160 P
1.0 0.0 0.0 0.0
0.0 1.0 0.0 0.0
0.0 0.0 1.0 0.0 }
```

Note the spaces. There must be spaces separating each keyword or datum to be read, even between the braces. For example:

```
/data/newpdb/pdb/pdb3sgb.ent 3sgb_protease{E 20 _ to E 160P}
```

would not be allowed.

In the second domain (Kallikrein) the transformation will be set equal to the identity matrix with a translation of zero, since none has been supplied.

The domains must be listed at the start of a file (ie. nothing must come before them in a file), but anything may come afterwards, provided that it contains no braces (ie. { or }) unless they are on lines containing ‘%’ in the first column.

It is possible to *reverse* the direction of an object in a domain description. For example, if one has two objects, one can reverse the direction of one or more of these by placing the word "REVERSE" in front of the object, e.g.:

```
/data/newpdb/pdb/pdb4mbn.ent { REVERSE _ 1 _ to _ 20 _ _ 21 _ to _ 120 _ }
```

3.2 Transformations

Transformations, which may or may not be included in the domain definition given above are in the sense:

$$\begin{array}{r} X_{\text{new}} \\ Y_{\text{new}} \\ Y_{\text{new}} \end{array} = \begin{array}{c} | R_{11} R_{12} R_{13} | X_{\text{old}} \\ | R_{21} R_{22} R_{23} | Y_{\text{old}} + V_2 \\ | R_{31} R_{32} R_{33} | Z_{\text{old}} \\ + V_3 \end{array}$$

OR

$$\begin{array}{l} X_{\text{new}} = (R_{11}*X_{\text{old}} + R_{12}*Y_{\text{old}} + R_{13}*Z_{\text{old}}) + V_1 \\ Y_{\text{new}} = (R_{21}*X_{\text{old}} + R_{22}*Y_{\text{old}} + R_{23}*Z_{\text{old}}) + V_2 \\ Z_{\text{new}} = (R_{31}*X_{\text{old}} + R_{32}*Y_{\text{old}} + R_{33}*Z_{\text{old}}) + V_3 \end{array}$$

If initial transformations are obtained in some other way (eg. those taken from a PDB file) they may be passed to STAMP if they are in the above format. As far as I can make out, this is the standard used in the PDB, but one can never be sure.

If no transformation is given, then the domain is assigned a unity rotation matrix and zero translation vector.

3.3 Sequence format

When necessary, STAMP programs read sequence information in NBRF (PIR) format. For example, user defined secondary structure assignment might be supplied in a file that looks like:

```
>Tonin
Tonin secondary structure Author's assignments
----EEEE-----EEEE-- <etc.> --HHHH---*
>Kallikrien
Kallikrien secondary structure -- visual inspection
----EEEEEEE---E-EEEE-- <etc.> --GGHHH---*
>SGprotease
S. Griseus protease secondary structure.
----EEEE---EEEE-EEEEEEE-- <etc.> --GGHGHG---*
```


This is essentially NBRF (PIR) format. Note the position of the asterix. Comments must be limited to the single line between the >identifier and the start of the sequence string.

3.4 Multiple alignment format

STAMP alignment output consists first of a list of domain descriptions and relevant transformations. After this an alignment may or may not be output.

Multiple alignments are displayed as follows (see STAMPDIR/examples/globin_stamp_trans.6):

```

/data/newpdb/pdb/pdb1lh1.ent 1lh1 { ALL
1.00000  0.00000  0.00000  0.00000
0.00000  1.00000  0.00000  0.00000
0.00000  0.00000  1.00000  0.00000  }
/data/newpdb/pdb/pdb2hhb.ent 2hhba { CHAIN A
0.71639  0.34414  0.60691  19.45435
<Etc.>
-0.31092  -0.94263  0.12159  68.85890  }

Alignment Score Sc = 7.665619
Alignment length Lp = 156
RMS deviation after fitting on 116 atoms = 2.434597
Secondary structures are from DSSP

>1lh1 (cluster A) sequence
>2hhba (cluster B) sequence
>2hhbb (cluster B) sequence
>4mbn (cluster B) sequence
>1ecd (cluster B) sequence
>2lhb (cluster B) sequence
>space
>1lh1_dssp (cluster A) secondary structure from DSSP
>2hhba_dssp (cluster B) secondary structure from DSSP
>2hhbb_dssp (cluster B) secondary structure from DSSP
>4mbn_dssp (cluster B) secondary structure from DSSP
>1ecd_dssp (cluster B) secondary structure from DSSP
>2lhb_dssp (cluster B) secondary structure from DSSP
#T -- '1' = used in the final fit
#P -- averaged Pij
#A -- distance between averaged CA atoms in angstroms
#G -- $P_{ij}{\prime}$ value
BBBBB ABBBBB use Pij Distance $P_{ij}{\prime}$
* iteration 1
P
I
V

```

```

D
T
G
S
V
G V A - - -
AVHV P ---- -
LLLLLL ----- 1 0.50337 1.90006 6.98400
TSTSS ----- 1 0.49631 2.00483 6.88900
EPPEAA HHHHHH 1 0.55533 1.89926 7.68300
SAEGDA HHHHHH 1 0.60834 1.80863 8.39600
QDEEQE HHHHHH 1 0.70134 1.64212 9.64700
AKKWIK HHHHHH 1 0.75434 1.52204 10.36000
ATSQST HHHHHH 1 0.75137 1.51092 10.32000
LNALTK HHHHHH 1 0.80831 1.36142 11.08600
VVVVI HHHHHH 1 0.85737 1.21626 11.74600
KKTLQR HHHHHH 1 0.83537 1.27448 11.45000

<Etc.>

ITNKG I HHHHHH 1 0.85737 1.04393 11.74600
VVADML HHHHHH 1 0.84332 1.11847 11.55700
ILLIIL HHHHHH 1 0.81232 1.20349 11.14000
KTAAFR HTHHHH 1 0.80035 1.22529 10.97900
KSHASS HTTHHT 1 0.73137 1.29476 10.05100
EKKKKA HTTHHT 1 0.60031 1.66495 8.28800
M Y H H
D K H H
D E H H
A L H H
*
```

The ‘>’ and ‘#’ characters tell the routines that read alignments what is to be contained in each field. A ‘>’ character denotes a character string which is to be displayed vertically, and a ‘#’ character denotes a string of numbers to be displayed separated by spaces. Thus in the above example we have 13 character strings vertically (6 amino acid sequences, 1 string of spaces and 6 DSSP assignments) and 6 numeric fields (corresponding to various details from STAMP) specified. The actual alignment will be contained within ‘*’ characters as shown. Accordingly, no occurrence of ‘>’, ‘#’ and ‘*’ characters should occur outside of these contexts.

The As and Bs just above the ‘*’ symbol refer to the members of the two cluster (branches) which are brought together during this alignment.

Briefly, the numeric fields are:

#T 1 or 0, 1 shows those residues used to determine the fit of the two sets of structures.

#P averaged Rossmann and Argos P_{ij} value

#A distance between averaged C_α atoms

#G corrected P_{ij} value (P_{ij}')

Note that the program POSTSTAMP adds two new fields:

#B 1 if all pairwise $P_{ij} \geq$ the user defined minimum, 0 otherwise

#R the total number of pairwise comparisons having $P_{ij} \geq$ the cutoff out of $N \times (N - 1)/2$

3.5 Output from STAMP database scanning mode

Output from STAMP scans consists of a list of domains and a corresponding set of scores, lengths and other numbers that can be used to sort and understand the output.

The format is as follows (see examples/ig/stamp_scan.trans):

```
% Output from STAMP scanning routine
%
% Domain 2fb4lv was used to scan the domain database:
% some.domains
% 1 fits were performed
% Fit 1 E1= 20.000, E2= 3.800, CUT= 1.000
% Approximate fits (alignment from N-termini) were performed
% at every 10 residue of the database sequences
% Transformations were output for Sc= 2.000
%
% Domain used to scan
# Sc= 10.000 RMS= 0.01 Len= 999 nfit= 999 Seqid= 100.00 Secid= 100.00 q_len= 111 d_len= 111 n_sec= 100 n_e
./pdb2fb4.ent 2fb4lv { L 1 _ to L 109 _ }
# Sc= 4.332 RMS= 1.556 len= 123 nfit= 57 seq_id= 21.82 sec_id= 74.55 q_len= 111 d_len= 105 n_sec= 1
/db/pdb_all/2fb4.pdb 2fb4lc_1 { L 110 _ to L 214 _ }
```

```

-0.69582  -0.69016   0.19878   132.36090
-0.57870   0.37482  -0.72430    67.11302
0.42538  -0.61902  -0.66021   109.94639  }
# Sc=   9.799 RMS=   0.001 len=  111 nfit=  111 seq_id= 100.00 sec_id= 97.30 q_len=  111 d_len=  216 n_sec=
/db/pdb_all//2fb4.pdb 2fb4l_1 { CHAIN L
1.00000   0.00001  -0.00002    0.00194
-0.00001   1.00000   0.00000    0.00193
0.00002  -0.00000   1.00000   -0.00027  }
# Sc=   7.842 RMS=   1.128 len=  116 nfit=   95 seq_id= 48.94 sec_id= 86.17 q_len=  111 d_len=  113 n_sec=
/db/pdb_all//1mcp.pdb 1mcp1v_1 { L 1 _ to L 113 _
-0.54068   0.77937   0.31662    37.59489
0.79918   0.35840   0.48255    0.56629
0.26261   0.51394  -0.81664   -33.14809  }
<etc.>

```

(note that the lines beginning with ‘#’ symbols have been wrapped here) ‘%’ denotes a comment, and ‘#’ denotes numbers corresponding to the domain description described below (both will be ignored by all programs except for SORTTRANS, which uses the ‘#’ fields to sort and interpret the data.

‘Sc’ is the STAMP Score for the comparison of the query to each database sequence. ‘RMS’ is the RMS difference between equivalenced atoms, ‘len’ is the alignment length, ‘nfit’ is the number of atoms used during the final fit of the two domains, ‘seq_id’ and ‘sec_id’ are the sequence and secondary structure identities, ‘q_len’ and ‘d_len’ are the lengths of the query and database structure (in residues), ‘n_sec’ is the number of equivalenced secondary structures, and ‘n_equiv’ are the number of residues found within stretches of 3 or more having $P'_{ij} \geq 6$. These fields are used during any run of SORTTRANS to sort and remove redundant/poor superimpositions. ‘fit_pos’ is the Brookhaven numbering of the position in the database sequence to which the query’s N-terminal end was aligned for the initial fit. The transformation supplied is that for the superimposition of the database structures onto the query.

3.6 Output to standard output or log file

STAMP now keeps fairly quiet during its running, updating the user only after a pairwise/treewise/scan comparison has been completed. You can get lots of other output by using the -V (verbose) option. If you want a lot of output to be written to a file instead of the standard output, you can use -V

in conjunction with `-logfile <file name>`.

Chapter 4

Summary of STAMP parameters

4.1 Main program (STAMP)

The format for running STAMP is:

```
stamp -l <starting domain file> -s -o <output file> -P <parameter file>
-n <1 or 2 fits> -d <database file for scans>
-slide <slide value>
-pen1 <gap penatly 1> -pen2 <gap pentalty 2>
-prefix <output file prefix>
-v
-rough
-cut
-<parameter> <value>
```

If you have old STAMP parameter files, they can be read by using the command `stamp -P <parameter file>`. This means that the old file can be read in exactly the same way as for version 2.0.

In general, all commands can be specified by `-<parameter> <value>`. For example, `'-first_pairpen 0.5'`. However, I have made some abbreviations for frequently used commands, these are:

```
-l <starting domain file>           same as -listfile <list file>
-o <output file>                   same as -logfile <output file>
-n <1 or 2>                         same as -npass <1 or 2>
-pen1 <gap penalty 1>             same as -first_pairpen <gap penalty 1>
```

```

-pen2 <gap penalty 2>    same as -second_pairpen <gap penalty 2>
-prefix <output prefix>  same as -transprefix <output prefix>
-s                        same as -scan true
-d <database file>       same as -database <database file>
-slide <slide parameter> same as -scanslide <slide parameter>
-cut                      same as -co true
-rough                   same as -roughfit true

```

Default parameters are always looked for in the file STAMPDIR/stamp.defaults. You can personalise this as you like, but I would recommend using the defaults, unless you have a thorough understanding of the method. The values described below were essentially chosen to mimic the successful and well-tested parameters [1].

I would recommend using the command line parameters. The commands, and their arguments are given below. The command line parameters are case insensitive. To use a parameter one need only type ‘-<parameter> <value>’ or use one of the short forms listed above.

STAMP can also be supplied with a parameter file. Parameters in a parameter file can be supplied in the format:

```
<Parameter> <Value> <Optional Comments> [return]
```

eg.

```

PAIRWISE Yes    Perform pairwise calculations
E1 3.8
E2 3.8
CUTOFF 4.5

```

The input is read in an open format. Generally, data are expected to be separated by spaces or return characters. The number and position of spaces, tabs and returns generally should not matter with the exception of PDB format, which is read as the fixed format described in the brookhaven documentation.

The possible parameters are listed below. Strings, characters, floats and integers are as expected (though strings may not contain spaces). Boolean variables may be set by any of the following:

TRUE == TRUE, True, T, true, Yes, YES, yes, Y, 1
FALSE == FALSE, False, F, false, No, NO, n, 0

LOGFILE <string>

This is the file into which the log is to be written. If 'stdout' is supplied then the information is written to the standard output.

Default LOGFILE = stdout

LISTFILE <string> (or '-l <string>' or '-f <string>')

This is the name of a file that contains the location and description of the domains to be analysed and, if desired, an initial transformation.

Default LISTFILE = domain.list

SECTYPE <integer>

This must be set to 0 (no secondary structure assignment) or to one of the following values:

SECTYPE = 1 Output from Kabsch and Sander's DSSP program [19].

SECTYPE = 2 Secondary structure summary format. A string of residue by residue secondary structure assignments for each domain is to be read in from SECFILE in the format specified in the previous chapter.

Note that it is not possible to mix assignments. This is probably not a very realistic thing to do anyway, since assignments can differ substantially. If you really want to do this, then the only possible way is to set SECTYPE = 3, and define each secondary structure independently in SECFILE.

Default SECTYPE = 1 (for DSSP).

SECFILE <string>

The file from which user specified secondary structure assignments are to be read (ie. SECTYPE = 2 only).

Default SECFILE = stamp.sec

PAIRWISE <boolean>

If TRUE, then pairwise comparisons are to be performed for each possible pair of domains described in LISTFILE. A matrix of pairwise (S_c) scores will be output (to MATFILE).

Default PAIRWISE = TRUE

N.B. Many of the following parameters also apply to TREEWISE and SCAN comparisons. For clarity they are discussed here in the PAIRWISE comparison context.

NPASS <1 or 2> (or '-n <1 or 2>')

Whether one or two fits are to be performed. The idea is that the initial fit can be used with a conformation biased set of parameters to improve the initial fit prior to fitting using distance and conformation parameters. The parameters described below are called 'first_' and 'second_' accordingly. When NPASS = 1, then only the 'second_' (or unprefixd) parameters are used. Default NPASS = 1

SW <0 or 1>

If set to 0, then the entire M x N matrix will be calculated and used during the Smith Waterman path finding routine. If set to 1, then a corner cutting routine will be used (to save time). Note that corner cutting will nullify many of the parameters specified in [1], and recommended only for SCAN mode. Accordingly, corner cutting parameters are specified below (after SCAN).

PAIRPEN <float> (or '-pen1 <float>'/ '-pen2 <float>')

(first_PAIRPEN)

(second_PAIRPEN)

Smith-Waterman gap penalty to be used during the fitting. second_PAIRPEN and PAIRPEN are equivalent. (PAIRPEN is also relevant to treewise fitting) Defaults PAIRPEN = second_PAIRPEN = 0.0 first_PAIRPEN = 0.0

E1 <float>

E2 <float>

(first_E1,first_E2)

(second_E2,second_E2)

Rossmann and Argos parameters to be used during the fitting. Rossmann and Argos suggested that E1 = E2 = 3.8 lead to good superimpositions, and further suggested that E1 = 20.0 and E2 = 3.8 would relax the distance requirement, and allow poor initial superimpositions to be improved. The

defaults are defined accordingly.

Defaults:

E1 = second_E1 = 3.8

E2 = second_E2 = 3.8

first_E1 = 20.0

first_E2 = 3.8

I would not recommend modifying these parameters, since I really don't know what changing them will do. If it ain't broke, don't fix it as my father would say.

NA <float>

NB <float>

NASD <float>

NBSD <float>

NSD <float>

NMEAN <float>

Parameters used to define P_{ij}' and S_c values. These are defined in [1]. I wouldn't change these.

Defaults:

NA = -0.9497

NB = 0.6859

NASD = -0.4743

NBSD = 0.01522

NMEAN = 0.02

NSD = 0.1

CUTOFF <float>

(first_CUTOFF)

(second_CUTOFF)

This is the minimum P_{ij}' value allowed for atoms to be used for a least squares fit. Equivalences above this value will be used to determine a transformation and RMS deviation.

Defaults:

CUTOFF = second_CUTOFF = 4.5
first_CUTOFF = 1.0

PAIRALIGN <boolean>
If true, then each final pairwise alignment will be output to the log file.
Default PAIRALIGN = FALSE

COLUMNS <integer>
Number of sequence positions to be displayed per line when either PAIRALIGN,
SCANALIN or TREEALIGN is set to TRUE.
Default COLUMNS = 80

SCORETOL <float>
This is the percent S_c difference that will result in convergence being reached.
In other words, if $100 \times \text{abs}|S_c - S_{c,old}|/S_{c,old} \leq \text{SCORETOL}$ then the fitting
will be considered done.
Default SCORETOL = 1.0

MAXPITER <integer>
The maximum number of iterations allowed during the pairwise comparisons.
This prevents a particular fit, which jumps between two values rather than
converging, from lasting indefinitely.
Default MAXPITER = 10

MATFILE <string>
This is the file which contains an upper diagonal matrix consisting of the
pairwise Scores (either 1/RMS, or S_c) for each comparison. It may then be
used to derive a tree, if desired, for treewise analysis.
Default MATFILE = <stamp_prefix>.mat

ROUGHFIT <boolean> (or '-rough' to set to TRUE)
If set to TRUE, then an initial rough superimposition will be performed by
aligning the N-terminal ends of the sequences and fitting on whatever atoms
this process equivalences. Probably this is too crude for structures that differ
quite a bit, but if they are very similar, one can use this to avoid having to

perform a multiple sequence alignment.

TREEWISE <boolean>

If TRUE, then a treewise comparison is performed by following a derived hierarchy. Reads in the matrix file specified (either created by PAIRWISE or some other method), derives a tree (dendrogram), and does a tree-based alignment.

Default TREEWISE = TRUE

TREEPEN <float>

(first_TREEPEN)

(second_TREEPEN)

Value subtracted from the P_{ij} matrix at positions where a residue is to be aligned with a gap. For details see [1].

Defaults TREEPEN = second_TREEPEN = 0.0 first_TREEPEN = 0.0

MAXTITER <int>

As for MAXPITER, but applied to the treewise case.

Default MAXPITER = 10

TREEALIGN <boolean>

As for PAIRALIGN, only for treewise comparisons.

Default TREEALIGN = TRUE

STAMPPREFIX <string> (or '-prefix <string>')

This is the name of the family of files that will be produced from a multiple alignment. The files will be named STAMPPREFIX.<N>, where N is the number of the cluster after which the alignment has been derived. There are always one fewer clusters than their are domains being compared.

Default STAMPPREFIX = 'stamp_trans'

SCAN <boolean> (or simply '-s' to set true)

If TRUE, then SCAN mode is selected. TREEWISE and PAIRWISE are set to FALSE. The first domain described in LISTFILE (the query) is used to scan all the domains listed in DATABASE. The parameters for scanning are described below. The output of a SCAN run appears in the file called STAMPPREFIX.scan.

Default SCAN = FALSE

DATABASE <string> (or -d <string>)

The list of domains to be compared with the query during a scan.

Default DATABASE = domain.database

MAXSITER <int>

As for MAXPITER and MAXTITER, but for scanning. Equivalent within the program to MAXPITER.

Default MAXSITER = 10

SCANALIGN <boolean>

As for PAIRALIGN and TREEALIGN, but for scanning. Equivalent within the program to MAXPITER.

Default SCANALIGN = FALSE

SCANSORE <integer>

Specifies how the Sc value is to be calculated. This depends on the particular application. The values are described in the first chapter.

As a general rule of thumb, use SCANSORE=6 for large database scans, when you are scanning with a small domain, and wishing to find all examples of this domain – even within large structures. Use SCANSORE=1 when you wish to obtain a set of transformations for a set of domains which you know are similar (and have defined fairly precisely as domains rather than the larger structure that they may be a part of).

Default SCANSORE = 6

SKIPAHEAD <boolean>

If set to TRUE, then the program will skip over all hits. In other words, if a similarity is found with a particular starting fit position, then the next fit position will be the last residue of the similar region. This is not always desirable, since there can be more than one hit within repetitive structures, such as α/β barrels.

Default SKIPAHED = TRUE

OPD <boolean>

Means “One Per Domain”. When the first hit for a domain is found during a SCAN (i.e. with S_c above SCANCUT), the rest of the comparisons involving that domain are skipped. Means that multiple matches involving the probe and database structures will be missed.

Default OPD = FALSE

SCANCUT <float>

If SCANMODE = 1, then S_c must be \geq SCANCUT in order for a transformation to be output.

Default SCANCUT = 2.0

SCANSLIDE <integer> (or ‘-slide <integer>’) This is the number of residues that a query sequence is ‘slid’ along a database sequence to derive each initial superimposition. Initially, the N-terminus of the query is aligned to the 1st residue of the database, once this fit has been performed and refined, and tested for good structural similarity, the N-terminus is aligned with the $1 + \langle \text{SCANSLIDE} \rangle$ th position, and the process repeated until the end of the database sequence has been reached.

Default SCANSLIDE = 5

SCANTRUNC <boolean>

If TRUE, then sequences from DATABASE that are more than SCANTRUNCFACTOR x the length of the query sequence are truncated to this size. This saves a lot of CPU time, as comparisons between things that are vastly different in size are largely meaningless. Moreover, since most scans will be done with discrete domains, then this allows separate domains in large proteins to be compared to the query separately.

Default SCANTRUNC = TRUE

SCANTRUNCFACTOR <float>

The largest size of sequence which may be compared to the query sequence (expressed as SCANTRUNCFACTOR x query sequence length). Structures in the DATABASE that are larger than this will be truncated to this size if SCANTRUNC = TRUE.

Default SCANTRUNCFACTOR = 2.0

SLOWSCAN <boolean>

If set to TRUE, then the SLOW method of getting the initial fits for scanning will be used (See chapter 1).

Default SLOWSCAN = FALSE

MIN_FRAC <float>

This is the minimum ratio of database length/query length to be allowed. In other words, if a database structure is too small (ie. if database length/query length < MIN_FRAC), then the comparison will be skipped. Whether to use this or not depends on whether or not one is interested in sub alignments where only a part of the query structure is used. The default implies that all comparisons will be performed.

Default MIN_FRAC = 0.001

SECSCREEN <boolean>

If TRUE, then an initial comparison between query and DATABASE secondary structure assignments (if available) is performed. A secondary structure distance is defined by:

$$D_{sec} = \sqrt{(\|Q_h - D_h\|^2 + \|Q_b - D_b\|^2)}$$

where Q_h and Q_b are the percent of Helix and Beta structure in the query, and D_h and D_b are the same for the database sequence. If Dist is larger than a threshold (SECSCREENMAX) then the comparison will be ignored.

Default SECSCREEN = true

SECSCREENMAX <float>

This is the maximum value of Dist (above) tolerated. If Dist is larger than SECSCREENMAX then the comparison is ignored. For screening to be effective, it is important that secondary structure assignments are accurate (preferably done using the same program).

Default SECSCREENMAX = 60.0 (this is very lenient; 40 is usually safe)

CCFACTOR <float>

Corner cutting factor. This is approximately the maximum number of gaps to be tolerated in any pairwise comparison. Only used if SW = 1. For a more detailed explanation, refer to [6] (pp 279 – 281).

Default CCFACTOR = 30.0

CCADD <boolean>

If TRUE, then the difference between query and database sequence lengths will be added to CCFACTOR. Probably this is only realistic when SCANT-RUNC is set TRUE.

Default CCADD = FALSE

PRECISION <integer>

Since STAMP works as much as possible with integers, this is what all floating point values are multiplied by during conversion. A value of 1000 has never presented us with any problems.

Default PRECISION = 1000

MAX_SEQ_LEN <integer>

The maximum length of alignment tolerated. The program ought to inform you when this value is surpassed.

Default MAX_SEQ_LEN = 1500

4.2 Summary of parameters for other programs

4.2.1 PDB checker (PDBC)

This is a simple program which looks for the location of a four letter PDB code (using the list of directories, prefixes and suffixes supplied in the file ./pdb.directories or if this does not exist STAMPDIR/pdb.directories) There are several options:

```
pdhc -q <four letter code>
```

will merely report useful information (number of atoms, the occurrence of HETATM, resolution, etc.) about each chain found in the PDB file which corresponds to the four letter code supplied.


```
pdbc -d <four letter code>[<chains to be considered>]
```

this outputs a domain description (or more than one if more than one chain is given. Sequential use of this program can be used to create a list of domains for use in scanning.

```
pdbc -m <four letter code>
```

this will just report the location of PDB and DSSP files. Good for a quick test of whether PDB codes can be found in the files specified in STAMPDIR.

Output is to standard output.

4.2.2 PDBSEQ

This program takes a list of protein domains (ie. a LISTFILE) and outputs a series of sequences derived from the described PDB files. The format is:

```
pdbseq -f <domain file> [-min <val> -max <val> -separate  
-format <fasta> -v -tl <max title length>]
```

‘-min/max <val>’ specify the minimum/maximum sequence length to be output. If the length of a sequence is less than min or greater than max, the sequence will be skipped (useful particularly if one wants to ignore very short PDB sequence, such as peptide inhibitors, etc.).

The output is in NBRF (PIR) format, and is written to the standard output. Using ‘-format <fasta>’ will make the output as FASTA format.

The option ‘-separate’ will produce files for each domain in the input file. These files are named ‘ID’.seq.

The program outputs a title line that attempts to describe the protein sequence according to the definitions given in the PDB file. The TITLE, COMPND and SOURCE lines are strung together (in that order). The

option `-tl ;number;` (`tl` = title limit) specifies the maximum length of this string. This description will always be postfixed (after a “:”) by the range of residues considered (i.e. All, Chain a, etc.).

4.2.3 ALIGNFIT

ALIGNFIT takes a multiple sequence alignment of proteins of known 3D structures and uses it to superimpose them. It requires two files: an AMPS multiple sequence alignment (block format), and a domain description file. An optional parameter file may be supplied; if none is given the program simply uses default parameters.

The format is:

```
alignfit -a <AMPS file> -d <domain file>
(-P <optional parm file> -<parameter> <value>)
```

`-P` can be used to read in an old ALIGNFIT parameter file (version 3.0 and earlier) The possible parameters, and their defaults are (names are case insensitive):

PAIRWISE <boolean>

If TRUE, then pairwise comparisons will be performed to derive a matrix (MATFILE).

Default PAIRWISE = TRUE

TREewise <boolean>

If TRUE, then treewise comparisons will be performed to derive a final transformation.

Default TREewise = TRUE

MATFILE <string>

The file into which the results of PAIRWISE are output.

Default MATFILE = alignfit.mat

MAX_SEQ_LEN <integer>

The maximum length of alignment to be tolerated.
Default is 3000

For most purposes, the default parameters should suffice. Note that one can use ACONVERT to convert CLUSTAL and MSF formats to block format, so that one can use alignments created using other programs (e.g. PILEUP, CLUSTAL, etc.) as a starting point for superimposition.

4.2.4 VER2HOR

This program provides a horizontal alignment given a STAMP alignment file (i.e. a text alignment written to the standard output). The format is:

```
ver2hor -f <stamp alignment file> [ -columns <width> ]
```

'-columns' specifies the number of columns to be used in the alignment output. This program is explained by example in the Worked Examples chapter. It also accepts most DSTAMP (see below) commands (i.e. those that are relevant to text output) from the command line.

4.2.5 DSTAMP

This program provides input for ALSCRIPT [20], GJB's program for the display of multiple sequence alignments. To get a copy of this program, contact GJB at the address at the front this manual.

The format is:

```
dstamp -f <STAMP alignment file> -prefix <output prefix>  
(-<parameter> [<value>])
```

where <parameter> is one of the many parameters described below. The new command line argument -P reads in parameter files, so if you have old DSTAMP files, they can still be read in this way.

The parameters for DSTAMP, and their defaults, are (parameter names are

case insensitive):

prefix <string>

Prefix specifying the name of the alscript (.als) and postscript files (.ps) to be generated.

Default prefix = 'alscript'

t <character>

The type of STAMP data to be used (ie. the first letter that occurs after the '#' characters in STAMP multiple alignment output). Default t = 'G'

c <float>

The minimum (or maximum in the case of RMS deviation) value to make a position considered as reliably aligned.

Default c = 6.0

w <integer>

The minimum length of a stretch of reliable regions to be allowed.

Default w = 3

ignore <integer>

This is the number of sequences that can be ignored during the calculation of residue or residue-property conservation (i.e. if ignore = 1 you allow one 'error' in one sequence during the calculation of conserved positions).

colour

Boolean parameter. If specified, the output will be in colour (via alscript).

motif

Boolean parameter. If specified, then a motif is written in the space between the sequence alignment and the aligned secondary structures.

The output is an ALSRIPT command file.

4.2.6 SORTTRANS

This program takes the output from a scan, and cleans and sorts the output. It removes repeated transformations by a simple least squares comparison of the matrices and vectors for those transformations which have the same identifier.

The format is:

```
sorttrans -f <scan output file> -s <keyword> <cutoff> [-t -i]
```

-f reads output from STAMP scanning, -s tells the program how to sort the output. The keyword tells which method to use. There are 8 possible keywords:

Sc	sort by Sc
rms	RMS deviation
nfit	number of fitted atoms
len	alignment length
frac	nfit/len
q_frac	nfit/q_len (q_len = length of query structure)
d_frac	nfit/d_len (d_len = length of database structure)
n_sec	number of equivalent secondary structure elements
seq_id	percent sequence identity
sec_id	percent secondary structure identity

sorted transformations are written to the standard output.

The option -i ==> identifiers only. Consider only the best transformation per identifier.

The option -n ==> ignore domain descriptors. This means that only the filename and the transformations are used. This is useful if you have different domain names attributed to the same region of the structure.

4.2.7 TRANSFORM

This program takes a transformation file, either from ALIGNFIT, STAMP, or SORTTRANS and outputs a series of PDB format files containing the specified coordinates transformed as specified in the given file.

The format is:

```
transform -f <transformation file> [ -g -het -hoh -o <output file> ]
```

options:

‘-het’ Include hetero atoms. Hetero atoms are normally not included in the output.

‘-hoh’ Include waters.

‘-g’ Graphics output. This mode puts all transformed coordinates into a single PDB file, and labels the chains for domains sequentially (after their order in the transformation file) with A, B, C.. etc. This allows fast analysis of the structures graphically (i.e. using Rasmol) since one need only colour each chain a different colour to see the superimposition. The default file for writing the coordinates using this mode is ‘all.pdb’, but this can be changed (see below).

‘-o <output file>’ When using ‘-g’, this option allows the specification of a file to contain the transformed coordinates. The default is ‘all.pdb’

The PDB files will be named <identifier>.pdb (except when running using the ‘-g’ option).

4.2.8 PICKFRAME

It is often the case that one wishes a particular protein structure to be the ‘parent’ of the superimposition, i.e. the structure that is un-transformed. Accordingly, the program PICKFRAME allows one to select a particular reference frame for a particular domain identifier. Given a transformation file and an identifier, the program will set the selected identifier’s transformation to the unit matrix and zero vector, and transform the other structures accordingly. This is useful if one wishes to combine different transformation files (i.e. if a multidomain protein has two domains, with each being similar to a separate domain).

The format is:

```
pickframe -f <transformation file> -i <domain identifier>
```

The output will be to the standard output (i.e. one need just pipe the results into a file).

This program is very useful if one wishes to superimpose STAMP results for two different domains from the same protein. Since one can just make all transformations relative to the PDB file containing the two domains, and then combine the output into one transformation file.

4.2.9 MERGETRANS & EXTRANS

Sometimes one has several transformations and wants to combine them. For example, one may have transformations from an ALIGNFIT run (i.e. taken from a multiple alignment) and those from a STAMP run and want to combine them, since they have at least one domain in common. This would avoid having to run the more time-consuming STAMP program on things where similarity was obvious (i.e. clear sequence homologues). MERGETRANS allows this to be done.

The format is:

```
mergetrans -f1 <transformation file1> -f2 <transformation file2> [-i <domain identifier>]
```

If an identifier is given, then that identifier will be used to link the two files (provided it can be found in both). Otherwise the program will simply search for the first identifier that is exactly in common across the two transformation files.

One may also wish to extract particular transformations from a file. To do this, use EXTRANS as follows:

```
extrans -f <transformation file> -i <id1> <id2> <id3> ... <idN>
```

A new transformation file will be output to the standard output containing only those domains that have been input on the command line.

4.2.10 MERGESTAMP

Sometimes one has several files containing transformations or alignments or both and wants to combine them. Alignments/transformations from STAMP may need to be combined with (for example) an alignment of a single PDB sequence with its homologues from a sequence database search, etc. MERGESTAMP does just this. It is essentially an extension of MERGETRANS.

The format is:

```
mergestamp -f1 <transformation file1> -f2 <transformation file2> [-i <domain identifier>]
```

If an identifier is given, then that identifier will be used to link the two files (provided it can be found in both). Otherwise the program will simply search for the first identifier that is exactly in common across the two transformation/alignment files.

4.2.11 AVESTRUC

For various reasons, it is often useful to derive ‘average’ structures (i.e. for homology modelling, molecular replacement search objects, etc.). STAMP output provides an obvious starting point for obtaining an average structure. AVESTRUC reads in a STAMP alignment file, and generates another PDB file containing averaged coordinates (either as C alpha or as a polyalanine structure).

The format is:

```
avestruc -f <STAMP alignment file>  
[ -polyA -c <STAMP char> -t <threshold> -w <window> -aligned ]
```


‘-f’ specifies the file to be considered. Note that this MUST BE a STAMP alignment file, containing both transformations and a sequence alignment. It will not work on transformation files lacking sequence alignment data or STAMP data.

‘-polyA’ generate polyalanine model, the default is a C alpha model

‘-c <STAMP char>’ ‘-t <threshold>’ ‘-w <window>’

these three parameters tell the program how to define structurally equivalent residues. ‘STAMP char’ is the label of the STAMP field specified by the ‘#’ character in the alignment file. ‘threshold’ is the minimum (or maximum in the case of RMS deviation) value of the specified STAMP parameter tolerated, and ‘window’ tells the minimum number of residues over which this must be true for structural equivalence. This is less complicated than it sounds.

The default is as described in [1]:

STAMP char = ‘G’ (i.e. P_{ij}^t)

threshold = 6.0

window = 3

(i.e. stretches of three or more residues having $P_{ij}^t > 6.0$ are considered equivalent)

‘-aligned’ this flag will generate an averaged position for all positions structures are present at a position (i.e. positions not containing any gaps are deemed equivalent). The temperature factor will then distinguish between genuine structural equivalences and fortitously aligned residues.

‘-ident’ ‘-cons’ these flags will name residues either as a single amino acid type (ident) or a conserved type (cons) according to the sequence alignment. See the appropriate sections in the preceding chapter for a further explanation.

4.2.12 GSTAMP

Like DSTAMP, this program takes STAMP output and translates it in to input for another program, namely Per Kraulis’ program MOLSCRIPT. The

program allows one to create multiple molscript files (i.e. one for each structure in the STAMP alignment file), or a single molscript file for an average structure. Appropriate PDB files for these alternatives must be generated by using TRANSFORM and AVESTRUC, respectively, prior to running MOLSCRIPT.

When multiple structure are considered, structurally equivalent regions (specified as for AVESTRUC) are shown as MOLSCRIPT helix, strand or coil. Non-structurally equivalent regions are shown as C_{α} trace. For an example of how this looks, see Figure 1 [11] or Figure 1 in [12].

The rest is up to you. Once MOLSCRIPT input files have been generated, they can be modified to suit your particular display needs (i.e. using colour, etc.).

The format is:

```
gstamp -f <STAMP alignment file>  
[ -c <STAMP char> -t <threshold> -w <window> -aligned -a -cons ]
```

-f, -c, -t, -w and -aligned is as for AVESTRUC and DSTAMP.

-a specifies that an average structure is to be used.

-cons specifies how the secondary structures are to be define in the MOLSCRIPT files. By default, structures are displayed as helix or strand only if *all* structures are helix or strand at the positions. '- cons' means that structures are displayed as helix or strand if the majority of structures are helix or strand at the positions. In both cases, the remaining structures are drawn as 'coil'.

BUG: sometimes GSTAMP will output single residue strands for Molscript input. It is therefore necessary to modify the Molscript output to correct the odd mistake (single residue strands produce funny pictures in my version of MOLSCRIPT — try it and see).

4.2.13 STAMP_CLEAN

This program allows you to tidy up gaps that are not meaningful in the context of a multiple sequence alignment derived from structure. In other words, regions that are not similar across all members of a structural family can be ‘cleaned’ to remove isolated residues aligned in the middle of nowhere. Note that one doesn’t always want to do this (since the sub-alignments can be meaningful).

The format is:

```
stamp_clean <stamp alignment file> <minimum segment length> > <output file>
```

The <minimum segment length> is the minimum number of residues that is to be considered significant. I always use 3, since this means that short stretches of 1-2 residues that are surrounded by gaps (i.e. in any sequence) are ‘cleaned’. Try it and see what I mean.

4.2.14 Converting alignment formats using ACONVERT

ACONVERT is a utility for interconverting alignment formats. It can be found installed as bin/aconvert in the STAMP installation directory. The typical usage is:

```
aconvert [-in <type> -out <type>] < <input file> > <output file>
```

where ‘type’ is one of ‘c’, ‘m’, ‘b’, ‘f’, ‘p’, which denote CLUSTAL, MSF, AMPS/BLOCK, FASTA and PIR format respectively. If no ‘-in’ argument is given, the program tries to guess the format, though note that this can sometimes fail (the program will usually issue an error in this case). For example to convert a STAMP alignment into CLUSTAL format, one would run :

```
aconvert -in b -out c < stamp_trans.10 > stamp_trans.10.aln
```

Chapter 5

Installation

5.1 Compiling/running

STAMP requires an ANSI C compiler (e.g. GCC) for installation.

STAMP is distributed as a gzipped tar file, which must be uncompressed and untarred to create the installation directory.

On most Unix and Unix-like systems, one can install STAMP with:

```
gunzip STAMP.tarfile.gz
tar -xvf STAMP.tarfile
cd stamp
./BUILD <system type> (e.g. BUILD sgi)
```

should work.

The systems that are available are:

- linux (should also work on Cygwin and MinGW)
- osx (Mac OS X)
- sgi (IRIX64 version 6.2)
- mips4-sgi (IRIX64 R10K version 6.2)
- dec (OSF1 version 4.0)
- sun (SunOS sol4 5.5.1)

All of these are specified by a makefile in the `src/` sub-directory. If your system isn't one of the above, then you can probably just use the one that is nearest, and edit the makefile accordingly.

Note that there are several precompiled executables in the distribution. Files found in the directories `bin/linux`, `bin/osx`, `bin/sgi`, `bin/sun` and `bin/dec`. You will overwrite these if you attempt a 'BUILD' as discussed above. Only the Linux and OS X binaries are current.

The built executables are copied into the directory `bin/<system-type>`, which should be added to your `PATH` environment variable, or linked/copied to some central directory, such as `/usr/local/bin`.

5.2 Configuring STAMP

To use STAMP, the user must set the environment variable `STAMPDIR` to the full path of the subdirectory `'/defs'` in which the installation was made. The directory containing the STAMP binaries, which is `STAMPDIR/bin` should also be included in the user's `PATH` environment variable.

STAMP reads PDB coordinate information and DSSP secondary structure assignments. Thus, you should have copies of the PDB and DSSP files for the structures in which you are interested (although DSSP files are not strictly required).

STAMP input files do not require that the full paths of the PDB / DSSP files being loaded should be specified. As an alternative to a full path, STAMP can find PDB and DSSP files for a domain by using only the domain identifier and sets of patterns defined in the files `'STAMPDIR/pdb.directories'` and `'STAMPDIR/dssp.directories'`. The format of each line in these files is:

```
<directory> <prefix> <suffix> [RETURN]
```

For example, the default `pdb.directories` file looks like this:

```
- - -  
./ _ _  
./ _ .pdb  
./ _ .pdb.Z  
./ _ .pdb.gz  
./ pdb .ent
```

STAMP searches for the PDB/DSSP files corresponding to a domain by taking the first four characters of the domain identifier as a PDB code and combining the code with each of patterns in turn to construct a test file name. If the test file exists, then that is used as the source of PDB coordinates or DSSP records. The fields in the pattern on each line are:

1. Directory path
2. File prefix
3. File suffix

If a field has the value ‘_’, then it is ignored when creating a test filename. For example, suppose STAMP is searching for the PDB file for a domain with the identifier 4cha. Using the default `pdb.directories` file, STAMP will attempt to open the following sequence of files:

```
4cha ./4cha ./4cha.pdb.Z ./4cha.pdb.gz ./pdb4cha.ent
```

The first file which it finds will be loaded to find coordinates for the domain. If you specify the full path to the PDB files in a STAMP input file, or the PDB files are in the directory in which you run STAMP, then the default `pdb.directories` file will be sufficient and you need not modify it.

A recent modification (version 4.2) is to look in each of the ‘distr’ type sub-directories for filenames. Some people store PDB files in a format, e.g.

```
<directory>/ab/pdb1abc.ent
```

Where the two letter sub-directory name corresponds to the second two characters in the four letter PDB code (i.e. ignoring the leading number). STAMP now handles these file types. If you just specify the top directory, the program will explore suitable two-letter sub-directories corresponding to each file it is looking for.

dssp.directories contains a description as to where possible DSSP files may be found. The format is as for pdb.directories, e.g.

```
-- --  
./ - .dssp  
./ - .dssp.Z
```

For example, the DSSP file for 4mbn might be found in the file ./4mbn.dssp

STAMP now reads compressed files (.Z or .gz suffixes). In order for this to work properly, you must have the programs zcat (.Z) and gunzip (.gz) installed on your system.

5.3 Getting other programs

There are several other programs that are useful to have when using STAMP:

DSSP – Definition of Secondary Structure in Proteins, Kabsch & Sander.
Contact

<http://swift.cmbi.kun.nl/gv/dssp/>

Note that this is the WWW page for both the program and a database of precomputed DSSP files corresponding to PDB entries.

Jalview – a cross-platform multiple alignment editor.

WWW page: <http://www.jalview.org>

ALSCRIPT – displays alignments in PostScript format, contact GJB (see address above)

WWW page: <http://www.compbio.dundee.ac.uk/>

MOLSCRIPT – displays PDB structures in PostScript format, contact:

<http://www.avatar.se/molscript/>

Chapter 6

Some of our studies involving STAMP

STAMP has been used in numerous published studies. Several novel similarities uncovered by STAMP have appeared in the literature: the similarity between the SH2 domain and domain II of *E. coli* biotin operon protein [9]; the similarity between HIV matrix protein p17 and Interferon gamma [16] and numerous others [12, 21, 22].

STAMP has also aided several other investigations into protein structure. STAMP alignments have been used to determine the best accuracy of secondary structure prediction from multiple sequence alignment [10]. It has been used to investigate the conservation of various protein structural features across structural similar (but apparently non-homologous) proteins [11, 13] and has been used for several investigations into protein domain structure [12, 23, 24].

STAMP has also proved extremely useful when assessing the results of protein structure prediction by fold recognition [25, 26, 27].

Most recently, STAMP has been used to investigate various aspects of protein function and evolution, in addition to doing large scale superimpositions of the entire protein database according to SCOP [13, 14], and problems associated with alignments for protein comparative modelling [28].

Bibliography

- [1] R. B. Russell and G. J. Barton. Multiple protein sequence alignment from tertiary structure comparison: assignment of global and residue confidence levels. *Proteins*, 14:309–323, 1992.
- [2] A. Sali and T. L. Blundell. Definition of general topological equivalence in protein structures, a procedure involving comparison of properties and relationships thorough simulated annealing and dynamic programming. *J. Mol. Biol.*, 212:403–428, 1990.
- [3] P. Argos and M. Rossmann. Exploring structural homology of proteins. *J. Mol. Biol.*, 105:75–95, 1976.
- [4] T. F. Smith and M. S. Waterman. Identification of common molecular subsequences. *J. Mol. Biol.*, 147:195–197, 1981.
- [5] G. J. Barton. An efficient algorithm to locate all locally optimal alignments. *Comp. App. Biosci.*, 9:729–734, 1993.
- [6] D. Sankoff and J. B. Kruskal, editors. *Time warps, string edits, and macromolecules: The theory and practice of sequence comparison*. Addison-Wesley, Inc., Reading, Mass., USA, 1983.
- [7] W. Kabsch. *Acta crystallographica*, A34:827, 1978.
- [8] A.D. McLachlan. Gene duplication in the structural evolution of chymotrypsin. *J. Mol. Biol.*, 128:49–79, 1979.
- [9] R. B. Russell and G. J. Barton. An SH2–SH3 domain hybrid. *Nature*, 364:765, 1993.

- [10] R. B. Russell and G. J. Barton. The limits of protein secondary structure prediction accuracy from multiple sequence alignment. *J. Mol. Biol.*, 234:951–957, 1993.
- [11] R. B. Russell and G. J. Barton. 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, 1994.
- [12] R. B. Russell. Domain insertion. *Prot. Eng.*, 7:1407–1410, 1994.
- [13] R. B. Russell, M. A. Saqi, R. A. Sayle, P. A. Bates, and M. J. E. Sternberg. Recognition of analogous and homologous protein folds: Analysis of sequence and structure conservation. *J. Mol. Biol.*, 269:423–439, 1997.
- [14] R. B. Russell, M. A. S. Saqi, P. A. Bates, R. A. Sayle, and M. J. E. Sternberg. Recognition of homologous and analogous protein folds: assessment of prediction success and associated alignment accuracy using empirical substitution matrices. *Prot. Eng.*, 11:1–9, 1998.
- [15] A. G. Murzin. Sweet-tasting protein monellin is related to the cystatin family of thiol proteinase inhibitors. *J. Mol. Biol.*, 230:689–694, 1993.
- [16] S. Matthews, P. Barlow, J. Boyd, G. Barton, R. Russell, H. Mills, M. Cunningham, N. Meyers, N. Burns, N. Clark, S. Kingsman, A. Kingsman, and I. Campbell. Structural similarity between the p17 matrix protein of HIV-1 and interferon- γ . *Nature*, 370:666–668, 1994.
- [17] W. R. Taylor. Classification of amino acid conservation. *J. Theor. Biol.*, 119:205–218, 1986.
- [18] P. J. Kraulis. Molscript: a program to produce both detailed and schematic plots of protein structures. *J. App. Cryst.*, 24:964–950, 1991.
- [19] W. Kabsch and C. Sander. A dictionary of protein secondary structure. *Biopolymers*, 22:2577–2637, 1983.
- [20] G. J. Barton. Alscript: A tool to format multiple sequence alignments. *Prot. Eng.*, 6:37–40, 1993.

- [21] R. B. Russell and M. J. E. Sternberg. A novel binding site in catalase is suggested by similarity to the calycin superfamily. *Prot. Eng.*, 9:107–111, 1996.
- [22] R. B. Russell and M. J. E. Sternberg. Two new examples of protein structural similarities within the structure-function twilight zone. *Prot. Eng.*, 10:333–338, 1997.
- [23] M. J. E. Sternberg, H. Hegyi, S. A. Islam, J. Luo, and R. B. Russell. Towards an intelligent system for the automatic assignment of domains in globular proteins. *Proceedings of the 3rd Annual Conference on Intelligent Systems for Molecular Biology*, pages 376–383, 1995.
- [24] A. S. Siddiqui and G. J. Barton. Continuous and discontinuous domains: an algorithm for the automatic generation of reliable protein domain definitions. *Prot. Sci.*, 4:872–874, 1995.
- [25] R. B. Russell, R. R. Copley, and G. J. Barton. Protein fold recognition from secondary structure assignments. *Proc. 28th Hawaii. Int. Conf. Sys. Sci. IEEE Press*, 5:302–311, 1995.
- [26] R. B. Russell, R. R. Copley, and G. J. Barton. Protein fold recognition by mapping predicted secondary structures. *J. Mol. Biol.*, 259:349–365, 1996.
- [27] R. B. Russell, P. D. Sasieni, and M. J. E. Sternberg. Supersites within superfolds. binding site similarity in the absence of homology. *J. Mol. Biol.*, 282:903–918, 1998.
- [28] M. A. S. Saqi, R. B. Russell, and M. J. E. Sternberg. Misleading local sequence alignments: implications for comparative protein modelling. *Prot. Eng.*, 11:627–630, 1998.
- [29] Anonymous.