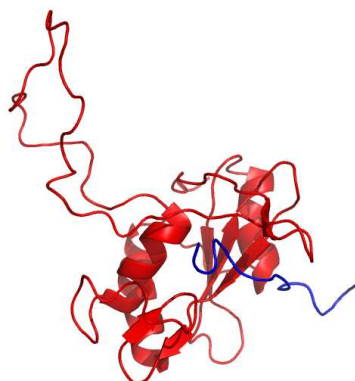


2bbu

Evolutionary trace report by **report_maker**

September 10, 2008



4.3.1	Alistat	5
4.3.2	CE	5
4.3.3	DSSP	5
4.3.4	HSSP	5
4.3.5	LaTex	5
4.3.6	Muscle	5
4.3.7	Pymol	6
4.4	Note about ET Viewer	6
4.5	Citing this work	6
4.6	About report_maker	6
4.7	Attachments	6

CONTENTS

1 Introduction

2 Chain 2bbuA

- 2.1 O35718 overview
- 2.2 Multiple sequence alignment for 2bbuA
- 2.3 Residue ranking in 2bbuA
- 2.4 Top ranking residues in 2bbuA and their position on the structure
 - 2.4.1 Clustering of residues at 28% coverage.
 - 2.4.2 Overlap with known functional surfaces at 28% coverage.
 - 2.4.3 Possible novel functional surfaces at 28% coverage.

3 Notes on using trace results

- 3.1 Coverage
- 3.2 Known substitutions
- 3.3 Surface
- 3.4 Number of contacts
- 3.5 Annotation
- 3.6 Mutation suggestions

4 Appendix

- 4.1 File formats
- 4.2 Color schemes used
- 4.3 Credits

1 INTRODUCTION

From the original Protein Data Bank entry (PDB id 2bbu):

Title: Solution structure of mouse socs3 in complex with a phosphopeptide from the gp130 receptor

Compound: Mol id: 1; molecule: suppressor of cytokine signaling 3; chain: a; fragment: kir/ess/sh2 domain/pest motif; synonym: socs-3, cytokine-inducible sh2 protein 3, cis-3, protein ef-10; engineered: yes; mol id: 2; molecule: gp130 phosphopeptide; chain: b; engineered: yes

Organism, scientific name: Mus Musculus;

- 2bbu contains a single unique chain 2bbuA (156 residues long).
- 1 Chain 2bbuB is too short (15 residues) to permit statistically significant analysis, and was treated as a peptide ligand. This is an
- 1 NMR-determined structure – in this report the first model in the file
- 1 was used.

2 CHAIN 2BBUA

2.1 O35718 overview

- 2 From SwissProt, id O35718, 91% identical to 2bbuA:
- 2 **Description:** Suppressor of cytokine signaling 3 (SOCS-3) (Cytokine-inducible SH2 protein 3) (CIS-3) (Protein EF-10).
- 3 **Organism, scientific name:** Mus musculus (Mouse).
- 3 **Taxonomy:** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Glires; Rodentia; Sciurognathi; Muroidea; Muridae; Murinae; Mus.
- 4 **Function:** SOCS family proteins form part of a classical negative feedback system that regulates cytokine signal transduction. SOCS3 is involved in negative regulation of cytokines that signal through the JAK/STAT pathway. Inhibits cytokine signal transduction by binding to tyrosine kinase receptors including gp130, LIF, erythropoietin, insulin, IL12, GCSF and leptin receptors. Binding to JAK2 inhibits its kinase activity. Suppresses fetal liver erythropoiesis. Regulates onset and maintenance of allergic responses mediated by T-helper
- 5 type 2 cells. Regulates IL-6 signaling in vivo.
- 5 **Subunit:** Interacts with multiple activated proteins of the tyrosine
- 5 kinase signaling pathway including IGF1 receptor, insulin receptor

and EPO receptor. Binding to JAK is mediated through the KIR and SH2 domains to a phosphorylated tyrosine residue within the JAK JH1 domain. Binds specific activated tyrosine residues of the leptin, EPOR and gp130 receptors.

Tissue specificity: Low expression in lung, spleen and thymus. Expressed in Th2 but not TH1 cells.

Developmental stage: In the developing brain, expressed at low levels from E10 stages to young adulthood (P25) with peak levels from E14 to P8. In the cortex, first expressed uniformly in all cells at E14. Not expressed in the retina. Highly expressed in fetal liver progenitors at E12.5.

Induction: By a subset of cytokines including EPO, leptin, LIF, IL-2, IL-3, IL-4, IGF1, growth hormone and prolactin.

Domain: The ESS and SH2 domains are required for JAK phosphotyrosine binding. Further interaction with the KIR domain is necessary for signal and kinase inhibition.

Ptm: Phosphorylated on tyrosine residues after stimulation by the cytokines, IL-2, EPO or IGF1 (By similarity).

Similarity: Contains 1 SH2 domain.

Similarity: Contains 1 SOCS box domain.

About: This Swiss-Prot entry is copyright. It is produced through a collaboration between the Swiss Institute of Bioinformatics and the EMBL outstation - the European Bioinformatics Institute. There are no restrictions on its use as long as its content is in no way modified and this statement is not removed.

2.2 Multiple sequence alignment for 2bbuA

For the chain 2bbuA, the alignment 2bbuA.msf (attached) with 8 sequences was used. The alignment was assembled through combination of BLAST searching on the UniProt database and alignment using Muscle program. It can be found in the attachment to this report, under the name of 2bbuA.msf. Its statistics, from the *alistat* program are the following:

```

Format:                MSF
Number of sequences:   8
Total number of residues: 1129
Smallest:              131
Largest:               156
Average length:       141.1
Alignment length:     156
Average identity:     49%
Most related pair:    97%
Most unrelated pair:  32%
Most distant seq:    45%
```

Furthermore, 21% of residues show as conserved in this alignment.

The alignment consists of 87% eukaryotic (87% vertebrata) sequences. (Descriptions of some sequences were not readily available.) The file containing the sequence descriptions can be found in the attachment, under the name 2bbuA.descr.

2.3 Residue ranking in 2bbuA

The 2bbuA sequence is shown in Fig. 1, with each residue colored according to its estimated importance. The full listing of residues in 2bbuA can be found in the file called 2bbuA.ranks_sorted in the attachment.

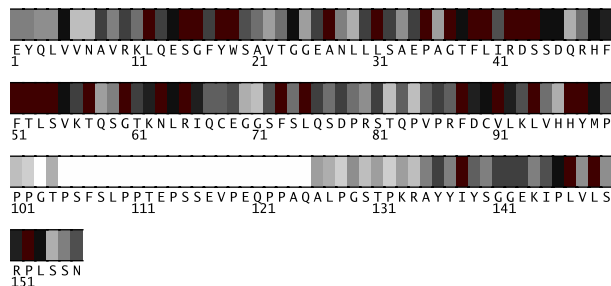


Fig. 1. Residues 1-156 in 2bbuA colored by their relative importance. (See Appendix, Fig.6, for the coloring scheme.)

2.4 Top ranking residues in 2bbuA and their position on the structure

In the following we consider residues ranking among top 28% of residues in the protein (the closest this analysis allows us to get to 25%). Figure 2 shows residues in 2bbuA colored by their importance: bright red and yellow indicate more conserved/important residues (see Appendix for the coloring scheme). A Pymol script for producing this figure can be found in the attachment.

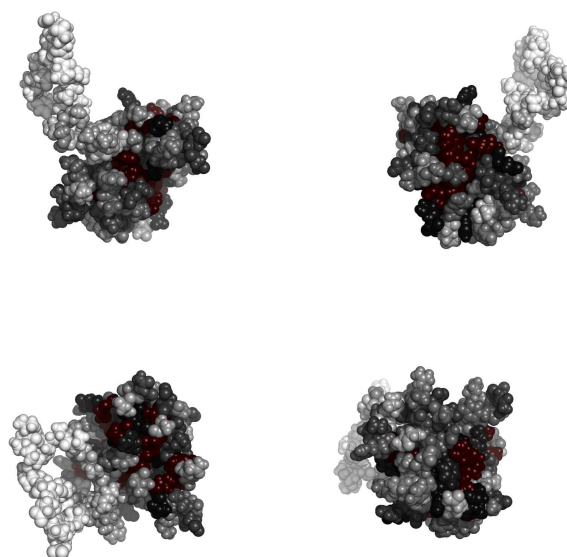


Fig. 2. Residues in 2bbuA, colored by their relative importance. Clockwise: front, back, top and bottom views.

2.4.1 Clustering of residues at 28% coverage. Fig. 3 shows the top 28% of all residues, this time colored according to clusters they belong to. The clusters in Fig.3 are composed of the residues listed in Table 1.

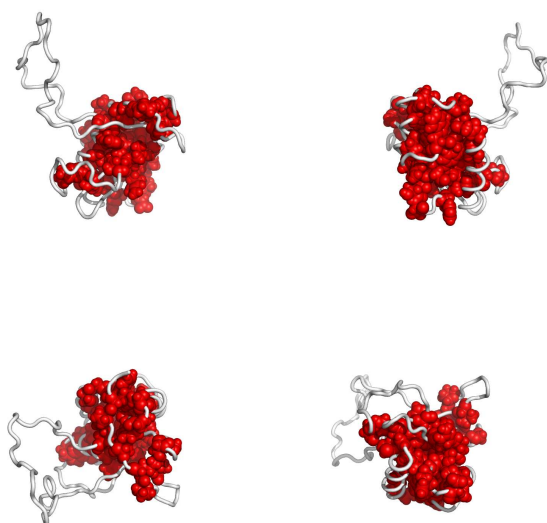


Fig. 3. Residues in 2bbuA, colored according to the cluster they belong to: red, followed by blue and yellow are the largest clusters (see Appendix for the coloring scheme). Clockwise: front, back, top and bottom views. The corresponding Pymol script is attached.

Table 1.		
cluster color	size	member residues
red	43	5, 12, 15, 16, 18, 19, 24, 27, 31, 35 37, 39, 40, 42, 43, 44, 45, 46, 49 51, 52, 53, 54, 57, 60, 63, 64, 65 73, 75, 88, 90, 91, 94, 97, 98, 99 138, 146, 147, 149, 152, 153

Table 1. Clusters of top ranking residues in 2bbuA.

2.4.2 *Overlap with known functional surfaces at 28% coverage.* The name of the ligand is composed of the source PDB identifier and the heteroatom name used in that file.

Interface with the peptide 2bbuB. Table 2 lists the top 28% of residues at the interface with 2bbuB. The following table (Table 3) suggests possible disruptive replacements for these residues (see Section 3.6).

Table 2.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
52	T	T(100)	0.21	55/0	3.24
63	N	N(100)	0.21	133/50	1.90
65	R	R(100)	0.21	66/3	2.16
75	L	L(100)	0.21	37/3	3.42
97	H	H(100)	0.21	93/2	1.94
138	I	I(100)	0.21	13/4	3.68

continued in next column

Table 2. continued					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
149	L	L(100)	0.21	2/0	4.59

Table 2. The top 28% of residues in 2bbuA at the interface with 2bbuB. (Field names: res: residue number in the PDB entry; type: amino acid type; subst's: substitutions seen in the alignment; with the percentage of each type in the bracket; noc/bb: number of contacts with the ligand, with the number of contacts realized through backbone atoms given in the bracket; dist: distance of closest approach to the ligand.)

Table 3.		
res	type	disruptive mutations
52	T	(KR) (FQMWH) (NELPI) (D)
63	N	(Y) (FTWH) (SEVCARG) (MD)
65	R	(TD) (SYEVCLAPIG) (FMW) (N)
75	L	(YR) (TH) (SKECG) (FQWD)
97	H	(E) (TQMD) (SNKVCLAPIG) (YR)
138	I	(YR) (TH) (SKECG) (FQWD)
149	L	(YR) (TH) (SKECG) (FQWD)

Table 3. List of disruptive mutations for the top 28% of residues in 2bbuA, that are at the interface with 2bbuB.

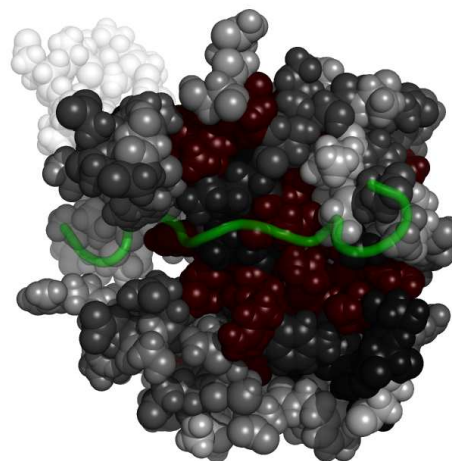


Fig. 4. Residues in 2bbuA, at the interface with 2bbuB, colored by their relative importance. 2bbuB is shown in backbone representation (See Appendix for the coloring scheme for the protein chain 2bbuA.)

Figure 4 shows residues in 2bbuA colored by their importance, at the interface with 2bbuB.

2.4.3 Possible novel functional surfaces at 28% coverage. One group of residues is conserved on the 2bbuA surface, away from (or substantially larger than) other functional sites and interfaces recognizable in PDB entry 2bbu. It is shown in Fig. 5. The right panel shows (in blue) the rest of the larger cluster this surface belongs to.

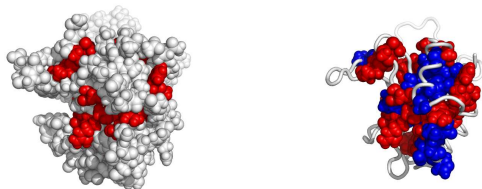


Fig. 5. A possible active surface on the chain 2bbuA. The larger cluster it belongs to is shown in blue.

The residues belonging to this surface "patch" are listed in Table 4, while Table 5 suggests possible disruptive replacements for these residues (see Section 3.6).

res	type	substitutions(%)	cvg
12	L	L(100)	0.21
15	S	S(100)	0.21
16	G	G(100)	0.21
18	Y	Y(100)	0.21
19	W	W(100)	0.21
35	P	P(100)	0.21
37	G	G(100)	0.21
39	F	F(100)	0.21
42	R	R(100)	0.21
44	S	S(100)	0.21
52	T	T(100)	0.21
57	T	T(100)	0.21
60	G	G(100)	0.21
63	N	N(100)	0.21
65	R	R(100)	0.21
73	F	F(100)	0.21
75	L	L(100)	0.21
88	F	F(100)	0.21
97	H	H(100)	0.21
138	I	I(100)	0.21
147	L	L(100)	0.21
149	L	L(100)	0.21

Table 4. Residues forming surface "patch" in 2bbuA.

res	type	disruptive mutations
12	L	(YR)(TH)(SKECG)(FQWD)
15	S	(KR)(FQMWH)(NYELPI)(D)
16	G	(KER)(FQMWH)(NYLPI)(SVA)
18	Y	(K)(QM)(NEVLAPIR)(D)
19	W	(KE)(TQD)(SNCRG)(M)
35	P	(YR)(TH)(SKECG)(FQWD)
37	G	(KER)(FQMWH)(NYLPI)(SVA)
39	F	(KE)(TQD)(SNCRG)(M)
42	R	(TD)(SYEVCLAPIG)(FMW)(N)
44	S	(KR)(FQMWH)(NYELPI)(D)
52	T	(KR)(FQMWH)(NELPI)(D)
57	T	(KR)(FQMWH)(NELPI)(D)
60	G	(KER)(FQMWH)(NYLPI)(SVA)
63	N	(Y)(FTWH)(SEVCARG)(MD)
65	R	(TD)(SYEVCLAPIG)(FMW)(N)
73	F	(KE)(TQD)(SNCRG)(M)
75	L	(YR)(TH)(SKECG)(FQWD)
88	F	(KE)(TQD)(SNCRG)(M)
97	H	(E)(TQMD)(SNKVCLAPIG)(YR)
138	I	(YR)(TH)(SKECG)(FQWD)
147	L	(YR)(TH)(SKECG)(FQWD)
149	L	(YR)(TH)(SKECG)(FQWD)

Table 5. Disruptive mutations for the surface patch in 2bbuA.

3 NOTES ON USING TRACE RESULTS

3.1 Coverage

Trace results are commonly expressed in terms of coverage: the residue is important if its "coverage" is small - that is if it belongs to some small top percentage of residues [100% is all of the residues in a chain], according to trace. The ET results are presented in the form of a table, usually limited to top 25% percent of residues (or to some nearby percentage), sorted by the strength of the presumed evolutionary pressure. (I.e., the smaller the coverage, the stronger the pressure on the residue.) Starting from the top of that list, mutating a couple of residues should affect the protein somehow, with the exact effects to be determined experimentally.

3.2 Known substitutions

One of the table columns is "substitutions" - other amino acid types seen at the same position in the alignment. These amino acid types may be interchangeable at that position in the protein, so if one wants to affect the protein by a point mutation, they should be avoided. For example if the substitutions are "RVK" and the original protein has an R at that position, it is advisable to try anything, but RVK. Conversely, when looking for substitutions which will *not* affect the protein, one may try replacing, R with K, or (perhaps more surprisingly), with V. The percentage of times the substitution appears in the alignment is given in the immediately following bracket. No percentage is given in the cases when it is smaller than 1%. This is meant to be a rough guide - due to rounding errors these percentages often do not add up to 100%.

3.3 Surface

To detect candidates for novel functional interfaces, first we look for residues that are solvent accessible (according to DSSP program) by at least 10\AA^2 , which is roughly the area needed for one water molecule to come in the contact with the residue. Furthermore, we require that these residues form a “cluster” of residues which have neighbor within 5\AA from any of their heavy atoms.

Note, however, that, if our picture of protein evolution is correct, the neighboring residues which *are not* surface accessible might be equally important in maintaining the interaction specificity - they should not be automatically dropped from consideration when choosing the set for mutagenesis. (Especially if they form a cluster with the surface residues.)

3.4 Number of contacts

Another column worth noting is denoted “noc/bb”; it tells the number of contacts heavy atoms of the residue in question make across the interface, as well as how many of them are realized through the backbone atoms (if all or most contacts are through the backbone, mutation presumably won’t have strong impact). Two heavy atoms are considered to be “in contact” if their centers are closer than 5\AA .

3.5 Annotation

If the residue annotation is available (either from the pdb file or from other sources), another column, with the header “annotation” appears. Annotations carried over from PDB are the following: site (indicating existence of related site record in PDB), S-S (disulfide bond forming residue), hb (hydrogen bond forming residue, jb (james bond forming residue), and sb (for salt bridge forming residue).

3.6 Mutation suggestions

Mutation suggestions are completely heuristic and based on complementarity with the substitutions found in the alignment. Note that they are meant to be **disruptive** to the interaction of the protein with its ligand. The attempt is made to complement the following properties: small [AVGSTC], medium [LPNQDEMILK], large [WPHYHR], hydrophobic [LPVAMWFI], polar [GTCY]; positively [KHR], or negatively [DE] charged, aromatic [WPHYH], long aliphatic chain [EK RQM], OH-group possession [SDETY], and NH2 group possession [NQRK]. The suggestions are listed according to how different they appear to be from the original amino acid, and they are grouped in round brackets if they appear equally disruptive. From left to right, each bracketed group of amino acid types resembles more strongly the original (i.e. is, presumably, less disruptive) These suggestions are tentative - they might prove disruptive to the fold rather than to the interaction. Many researcher will choose, however, the straightforward alanine mutations, especially in the beginning stages of their investigation.

4 APPENDIX

4.1 File formats

Files with extension “ranks_sorted” are the actual trace results. The fields in the table in this file:

- alignment# number of the position in the alignment
- residue# residue number in the PDB file
- type amino acid type
- rank rank of the position according to older version of ET

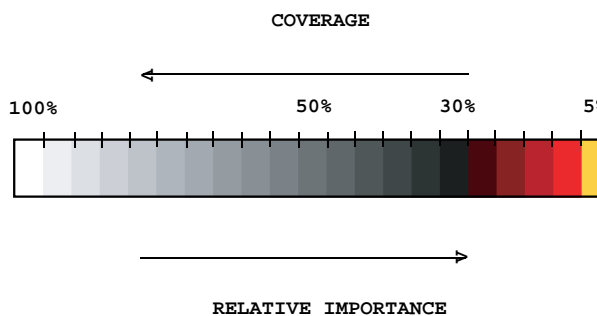


Fig. 6. Coloring scheme used to color residues by their relative importance.

- variability has two subfields:
 1. number of different amino acids appearing in in this column of the alignment
 2. their type
- rho ET score - the smaller this value, the lesser variability of this position across the branches of the tree (and, presumably, the greater the importance for the protein)
- cvg coverage - percentage of the residues on the structure which have this rho or smaller
- gaps percentage of gaps in this column

4.2 Color schemes used

The following color scheme is used in figures with residues colored by cluster size: black is a single-residue cluster; clusters composed of more than one residue colored according to this hierarchy (ordered by descending size): red, blue, yellow, green, purple, azure, turquoise, brown, coral, magenta, LightSalmon, SkyBlue, violet, gold, bisque, LightSlateBlue, orchid, RosyBrown, MediumAquamarine, DarkOliveGreen, CornflowerBlue, grey55, burlywood, LimeGreen, tan, DarkOrange, DeepPink, maroon, BlanchedAlmond.

The colors used to distinguish the residues by the estimated evolutionary pressure they experience can be seen in Fig. 6.

4.3 Credits

4.3.1 Alistat *alistat* reads a multiple sequence alignment from the file and shows a number of simple statistics about it. These statistics include the format, the number of sequences, the total number of residues, the average and range of the sequence lengths, and the alignment length (e.g. including gap characters). Also shown are some percent identities. A percent pairwise alignment identity is defined as $(\text{idents} / \text{MIN}(\text{len1}, \text{len2}))$ where *idents* is the number of exact identities and *len1*, *len2* are the unaligned lengths of the two sequences. The “average percent identity”, “most related pair”, and “most unrelated pair” of the alignment are the average, maximum, and minimum of all $(N)(N-1)/2$ pairs, respectively. The “most distant seq” is calculated by finding the maximum pairwise identity (best relative) for all *N* sequences, then finding the minimum of these *N* numbers (hence, the most outlying sequence). *alistat* is copyrighted

by HHMI/Washington University School of Medicine, 1992-2001, and freely distributed under the GNU General Public License.

4.3.2 CE To map ligand binding sites from different source structures, `report_maker` uses the CE program: <http://cl.sdsc.edu/>. Shindyalov IN, Bourne PE (1998) "Protein structure alignment by incremental combinatorial extension (CE) of the optimal path". *Protein Engineering* 11(9) 739-747.

4.3.3 DSSP In this work a residue is considered solvent accessible if the DSSP program finds it exposed to water by at least 10\AA^2 , which is roughly the area needed for one water molecule to come in the contact with the residue. DSSP is copyrighted by W. Kabsch, C. Sander and MPI-MF, 1983, 1985, 1988, 1994 1995, CMBI version by Elmar.Krieger@cmbi.kun.nl November 18,2002,

<http://www.cmbi.kun.nl/gv/dssp/descrip.html>.

4.3.4 HSSP Whenever available, `report_maker` uses HSSP alignment as a starting point for the analysis (sequences shorter than 75% of the query are taken out, however); R. Schneider, A. de Daruvar, and C. Sander. "The HSSP database of protein structure-sequence alignments." *Nucleic Acids Res.*, 25:226-230, 1997.

<http://swift.cmbi.kun.nl/swift/hssp/>

4.3.5 LaTeX The text for this report was processed using L^AT_EX; Leslie Lamport, "LaTeX: A Document Preparation System Addison-Wesley," Reading, Mass. (1986).

4.3.6 Muscle When making alignments "from scratch", `report_maker` uses Muscle alignment program: Edgar, Robert C. (2004), "MUSCLE: multiple sequence alignment with high accuracy and high throughput." *Nucleic Acids Research* 32(5), 1792-97.

<http://www.drive5.com/muscle/>

4.3.7 Pymol The figures in this report were produced using Pymol. The scripts can be found in the attachment. Pymol is an open-source application copyrighted by DeLano Scientific LLC (2005). For more information about Pymol see <http://pymol.sourceforge.net/>. (Note for Windows users: the attached package needs to be unzipped for Pymol to read the scripts and launch the viewer.)

4.4 Note about ET Viewer

Dan Morgan from the Lichtarge lab has developed a visualization tool specifically for viewing trace results. If you are interested, please visit:

<http://mammoth.bcm.tmc.edu/traceview/>

The viewer is self-unpacking and self-installing. Input files to be used with ETV (extension .etvx) can be found in the attachment to the main report.

4.5 Citing this work

The method used to rank residues and make predictions in this report can be found in Mihalek, I., I. Reš, O. Lichtarge. (2004). "A Family of Evolution-Entropy Hybrid Methods for Ranking of Protein Residues by Importance" *J. Mol. Bio.* **336**: 1265-82. For the original version of ET see O. Lichtarge, H.Bourne and F. Cohen (1996). "An Evolutionary Trace Method Defines Binding Surfaces Common to Protein Families" *J. Mol. Bio.* **257**: 342-358.

`report_maker` itself is described in Mihalek I., I. Res and O. Lichtarge (2006). "Evolutionary Trace Report Maker: a new type of service for comparative analysis of proteins." *Bioinformatics* **22**:1656-7.

4.6 About report_maker

`report_maker` was written in 2006 by Ivana Mihalek. The 1D ranking visualization program was written by Ivica Reš. `report_maker` is copyrighted by Lichtarge Lab, Baylor College of Medicine, Houston.

4.7 Attachments

The following files should accompany this report:

- 2bbuA.complex.pdb - coordinates of 2bbuA with all of its interacting partners
- 2bbuA.etvx - ET viewer input file for 2bbuA
- 2bbuA.cluster_report.summary - Cluster report summary for 2bbuA
- 2bbuA.ranks - Ranks file in sequence order for 2bbuA
- 2bbuA.clusters - Cluster descriptions for 2bbuA
- 2bbuA.msf - the multiple sequence alignment used for the chain 2bbuA
- 2bbuA.descr - description of sequences used in 2bbuA msf
- 2bbuA.ranks_sorted - full listing of residues and their ranking for 2bbuA
- 2bbuA.2bbuB.if.pml - Pymol script for Figure 4
- 2bbuA.cbvcg - used by other 2bbuA – related pymol scripts