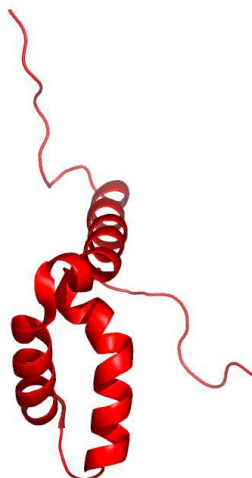


2yul

Evolutionary trace report by **report_maker**

November 15, 2009



4.3.3	DSSP	4
4.3.4	HSSP	4
4.3.5	LaTex	5
4.3.6	Muscle	5
4.3.7	Pymol	5
4.4	Note about ET Viewer	5
4.5	Citing this work	5
4.6	About report_maker	5
4.7	Attachments	5

CONTENTS

1 Introduction

2 Chain 2yulA

- 2.1 Q9H6I2 overview
- 2.2 Multiple sequence alignment for 2yulA
- 2.3 Residue ranking in 2yulA
- 2.4 Top ranking residues in 2yulA and their position on the structure
 - 2.4.1 Clustering of residues at 26% coverage.
 - 2.4.2 Possible novel functional surfaces at 26% 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
 - 4.3.1 **Alistat**
 - 4.3.2 **CE**

1 INTRODUCTION

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

Title: Solution structure of the hmg box of human transcription factor sox-17

Compound: Mol id: 1; molecule: transcription factor sox-17; chain: a; fragment: hmg (high mobility group) box; engineered: yes

Organism, scientific name: Homo Sapiens;

2yul contains a single unique chain 2yulA (82 residues long). This is an NMR-determined structure – in this report the first model in the file was used.

2 CHAIN 2YULA

2.1 Q9H6I2 overview

- 1 From SwissProt, id Q9H6I2, 92% identical to 2yulA:
- 1 **Description:** Transcription factor SOX-17.
- 1 **Organism, scientific name:** Homo sapiens (Human).
- 1 **Taxonomy:** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Catarrhini; Hominidae; Homo.
- 2 **Function:** Probable transcriptional activator in the premeiotic germ cells. It binds to the sequences 5'-AACAAAT-3' or 5'-AACAAAG-3' (By similarity).
- 2 **Subcellular location:** Nuclear (Potential).
- 3 **Tissue specificity:** Expressed in adult heart, lung, spleen, testis, ovary, placenta, fetal lung, and kidney. In normal gastrointestinal tract, it is preferentially expressed in esophagus, stomach and small intestine than in colon and rectum.
- 3 **Similarity:** Contains 1 HMG box DNA-binding domain.
- 3 **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.
- 4 **2.2 Multiple sequence alignment for 2yulA**
- 4 For the chain 2yulA, the alignment 2yulA.msf (attached) with 562 sequences was used. The alignment was downloaded from the HSSP database, and fragments shorter than 75% of the query as well as

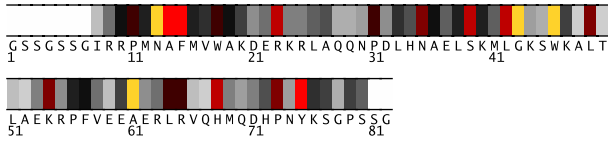


Fig. 1. Residues 1-82 in 2yulA colored by their relative importance. (See Appendix, Fig.5, for the coloring scheme.)

duplicate sequences were removed. It can be found in the attachment to this report, under the name of 2yulA.msff. Its statistics, from the *alistat* program are the following:

```

Format:                MSF
Number of sequences:   562
Total number of residues: 43056
Smallest:              62
Largest:               82
Average length:        76.6
Alignment length:      82
Average identity:       52%
Most related pair:     99%
Most unrelated pair:   5%
Most distant seq:      39%

```

Furthermore, <1% of residues show as conserved in this alignment.

The alignment consists of 47% eukaryotic (40% vertebrata, <1% arthropoda, 2% fungi) sequences. (Descriptions of some sequences were not readily available.) The file containing the sequence descriptions can be found in the attachment, under the name 2yulA.descr.

2.3 Residue ranking in 2yulA

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

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

In the following we consider residues ranking among top 26% of residues in the protein (the closest this analysis allows us to get to 25%). Figure 2 shows residues in 2yulA 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.

2.4.1 Clustering of residues at 26% coverage. Fig. 3 shows the top 26% 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.

Table 1.		
cluster color	size	member residues
red	18	11, 13, 14, 15, 18, 23, 35, 36, 39
<i>continued in next column</i>		

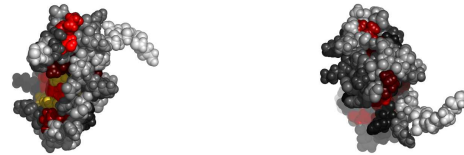
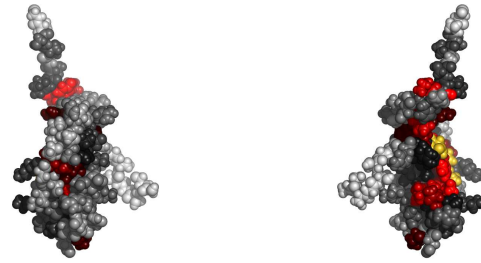


Fig. 2. Residues in 2yulA, colored by their relative importance. Clockwise: front, back, top and bottom views.

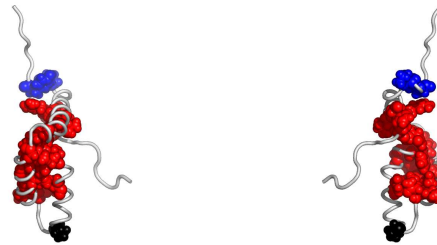


Fig. 3. Residues in 2yulA, colored according to the cluster they belong to: red, followed by blue and yellow are the largest clusters (see Appendix for the coloring scheme). Clockwise: front, back, top and bottom views. The corresponding Pymol script is attached.

Table 1. continued		
cluster color	size	member residues
		42, 43, 46, 49, 54, 61, 64, 65, 68
blue	2	73, 75

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

2.4.2 *Possible novel functional surfaces at 26% coverage.* One group of residues is conserved on the 2yulA surface, away from (or substantially larger than) other functional sites and interfaces recognizable in PDB entry 2yul. It is shown in Fig. 4. The residues

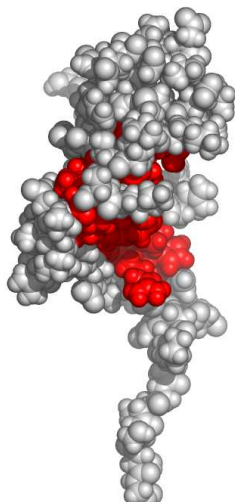


Fig. 4. A possible active surface on the chain 2yulA.

belonging to this surface "patch" are listed in Table 2, while Table 3 suggests possible disruptive replacements for these residues (see Section 3.6).

Table 2.			
res	type	substitutions(%)	cvg
46	W	W(99)CGQS	0.01
43	G	AG(97)S(2)	0.02
61	A	A(91)S(1)Q(5)YR PVG	0.04
13	N	N(96)S(1)L.PHDT	0.05
14	A	A(93)C(1)S(1) P(1)G.T	0.06
15	F	F(95)W(1)Y(2)L. S	0.07
39	S	TS(92)X(1)MA(1) CN(2)GE.L	0.10
42	L	L(91)I(4)C(1)VS A.PRMTG	0.11
68	H	F(2)LH(93)KVA Y(1).SQPM	0.12
23	R	R(90)H(5)Q(1)VS GAL.YK	0.13

continued in next column

Table 2. continued			
res	type	substitutions(%)	cvg
35	N	N(90)S(2)TV(1)K CHFPIGQDMX	0.16
49	L	L(80)M(12)E(5)A IQVSTFG	0.17
54	K	K(92)Q(1)R(3)VI SPM	0.18
11	P	P(94). (3)SANRGL	0.20
18	W	F(9)W(81)Y(8)ER IG	0.22
64	L	YL(83)E(4)V(2) I(2)A(2)D(1) M(1)FKN	0.23
65	R	K(17)R(58)Q(12) L(4)S(5)HNATD	0.24

Table 2. Residues forming surface "patch" in 2yulA.

Table 3.		
res	type	disruptive mutations
46	W	(KE)(D)(QR)(T)
43	G	(KR)(E)(QH)(FMW)
61	A	(E)(KYR)(HD)(Q)
13	N	(Y)(H)(FWR)(T)
14	A	(R)(K)(E)(Y)
15	F	(K)(E)(Q)(D)
39	S	(R)(K)(H)(FW)
42	L	(R)(Y)(H)(KE)
68	H	(E)(T)(D)(Q)
23	R	(D)(T)(E)(Y)
35	N	(Y)(H)(FW)(R)
49	L	(R)(Y)(H)(K)
54	K	(Y)(T)(FW)(CG)
11	P	(Y)(R)(H)(T)
18	W	(K)(E)(Q)(T)
64	L	(YR)(H)(T)(K)
65	R	(TY)(D)(E)(CG)

Table 3. Disruptive mutations for the surface patch in 2yulA.

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 [LPNQDEMIK], large [WFYHR], hydrophobic [LPVAMWFI], polar [GTCY]; positively [KHR], or negatively [DE] charged, aromatic [WFYH], 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

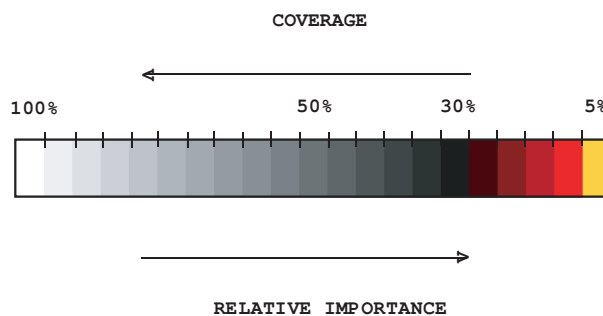


Fig. 5. Coloring scheme used to color residues by their relative importance.

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
- 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. 5.

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:

- 2yulA.complex.pdb - coordinates of 2yulA with all of its interacting partners
- 2yulA.etvx - ET viewer input file for 2yulA
- 2yulA.cluster_report.summary - Cluster report summary for 2yulA
- 2yulA.ranks - Ranks file in sequence order for 2yulA
- 2yulA.clusters - Cluster descriptions for 2yulA
- 2yulA.msf - the multiple sequence alignment used for the chain 2yulA
- 2yulA.descr - description of sequences used in 2yulA msf
- 2yulA.ranks_sorted - full listing of residues and their ranking for 2yulA