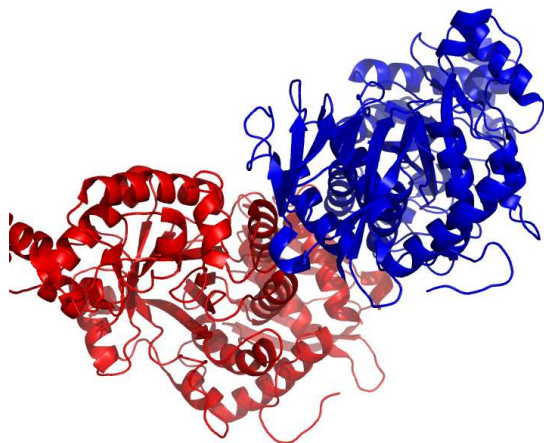


# Inou

Evolutionary trace report by **report\_maker**

September 22, 2008



## CONTENTS

### 1 Introduction

### 2 Chain InouA

- 2.1 P07686 overview
- 2.2 Multiple sequence alignment for InouA
- 2.3 Residue ranking in InouA
- 2.4 Top ranking residues in InouA and their position on the structure
  - 2.4.1 Clustering of residues at 25% coverage.
  - 2.4.2 Overlap with known functional surfaces at 25% coverage.
  - 2.4.3 Possible novel functional surfaces at 25% 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	<b>Alistat</b>	8
4.3.2	<b>CE</b>	9
4.3.3	<b>DSSP</b>	9
4.3.4	<b>HSSP</b>	9
4.3.5	<b>LaTex</b>	9
4.3.6	<b>Muscle</b>	9
4.3.7	<b>Pymol</b>	9
4.4	Note about ET Viewer	9
4.5	Citing this work	9
4.6	About report_maker	9
4.7	Attachments	9

## 1 INTRODUCTION

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

**Title:** Native human lysosomal beta-hexosaminidase isoform b

**Compound:** Mol id: 1; molecule: beta-hexosaminidase beta chain; chain: a, b; synonym: n-acetyl-beta-glucosaminidase, beta-n- acetyl-hexosaminidase, hexosaminidase b; ec: 3.2.1.52

**Organism, scientific name:** Homo Sapiens;

Inou contains a single unique chain InouA (480 residues long) and its homologue InouB.

## 1 2 CHAIN 1NOUA

### 2.1 P07686 overview

- 1 From SwissProt, id P07686, 96% identical to InouA:
- 1 **Description:** Beta-hexosaminidase beta chain precursor
- 2 (EC 3.2.1.52) (N-acetyl-beta- glucosaminidase) (Beta-N-
- 2 acetylhexosaminidase) (Hexosaminidase B) (Cervical cancer
- 2 proto-oncogene 7) (HCC-7).
- 2 **Organism, scientific name:** Homo sapiens (Human).
- 2 **Taxonomy:** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata;
- 3 Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates;
- 3 Catarrhini; Hominidae; Homo.
- 5 **Function:** Beta-hexosaminidase A is responsible for the degradation
- 7 of GM2 gangliosides, and a variety of other molecules containing
- 7 terminal N-acetyl hexosamines, in the brain and other tissues.
- 7 **Catalytic activity:** Hydrolysis of terminal non-reducing N-acetyl-
- 7 D-hexosamine residues in N-acetyl-beta-D-hexosaminides.
- 8 **Subunit:** There are 3 major forms of beta-hexosaminidase: hexo-
- 8 saminidase A is a trimer composed of one alpha chain, one beta-A
- 8 chain and one beta-B chain; hexosaminidase B is a tetramer of two
- 8 beta-A and two beta-B chains; hexosaminidase S is a homodimer of
- 8 two alpha chains. Some minor isozymes contain less processed forms
- 8 of the alpha or beta chains.
- 8 **Subcellular location:** Lysosomal.
- 8 **Ptm:** N-linked glycans at Asn-142 and Asn-190 consist of Man(3)-
- 8 GlcNAc(2) and Man(5 to 7)-GlcNAc(2), respectively.

**Ptm:** The beta-A and beta-B chains are produced by proteolytic processing of the precursor beta chain.

**Disease:** Defects in HEXB are the cause of Sandhoff disease (SD) [MIM:268800]; also known as GM2-gangliosidosis type II. SD is a progressive neurodegenerative disorder characterized by an accumulation of GM2 gangliosides, particularly in neurons. It is clinically indistinguishable from Tay-Sachs disease.

**Similarity:** Belongs to the glycosyl hydrolase 20 family.

**Caution:** Ref.3 sequence differs from that shown due to a frameshift in position 21.

**Database:** NAME=HEXBdb; NOTE=HEXB mutation database; WWW="http://www.hexdb.mcgill.ca/?Topic=HEXBdb".

**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 InouA

For the chain InouA, the alignment InouA.msf (attached) with 86 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 InouA.msf. Its statistics, from the *alistat* program are the following:

```

Format:                MSF
Number of sequences:   86
Total number of residues: 38878
Smallest:              363
Largest:               480
Average length:        452.1
Alignment length:      480
Average identity:      30%
Most related pair:     99%
Most unrelated pair:   19%
Most distant seq:     36%

```

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

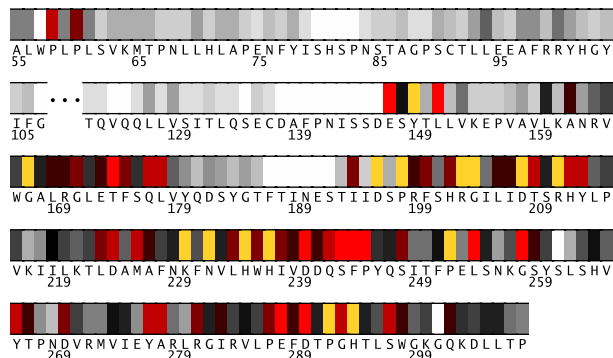
The alignment consists of 50% eukaryotic ( 12% vertebrata, 4% arthropoda, 18% fungi, 6% plantae), and 48% prokaryotic sequences. (Descriptions of some sequences were not readily available.) The file containing the sequence descriptions can be found in the attachment, under the name InouA.descr.

## 2.3 Residue ranking in InouA

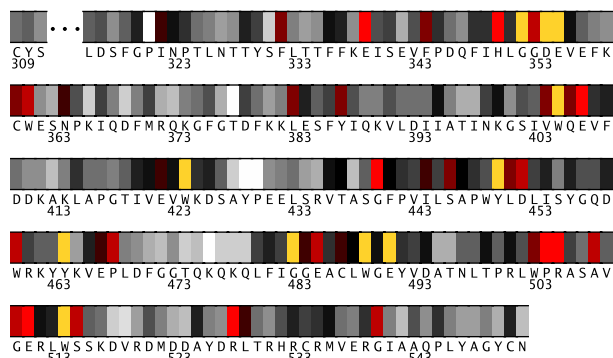
The InouA sequence is shown in Figs. 1–2, with each residue colored according to its estimated importance. The full listing of residues in InouA can be found in the file called InouA.ranks\_sorted in the attachment.

## 2.4 Top ranking residues in InouA and their position on the structure

In the following we consider residues ranking among top 25% of residues in the protein . Figure 3 shows residues in InouA colored by their importance: bright red and yellow indicate more conserved/important residues (see Appendix for the coloring scheme). A



**Fig. 1.** Residues 55-308 in InouA colored by their relative importance. (See Appendix, Fig.11, for the coloring scheme.)



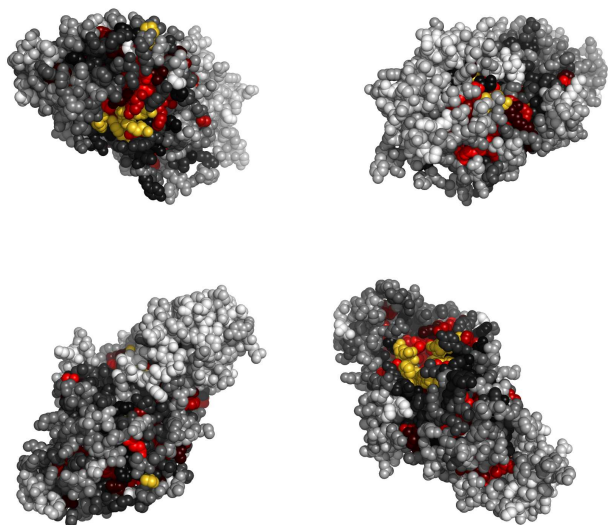
**Fig. 2.** Residues 309-552 in InouA colored by their relative importance. (See Appendix, Fig.11, for the coloring scheme.)

Pymol script for producing this figure can be found in the attachment.

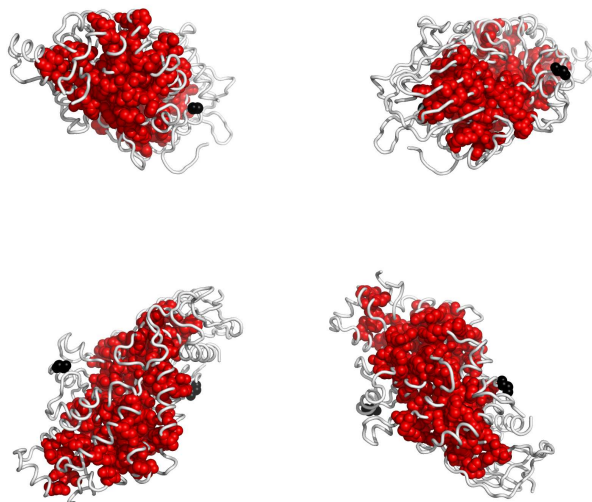
2.4.1 *Clustering of residues at 25% coverage.* Fig. 4 shows the top 25% of all residues, this time colored according to clusters they belong to. The clusters in Fig.4 are composed of the residues listed in Table 1.

Table 1.		
cluster color	size	member residues
red	116	58, 60, 147, 149, 151, 162, 167 169, 170, 171, 173, 174, 175, 177 178, 194, 196, 198, 199, 200, 202 203, 204, 206, 207, 208, 209, 211 212, 213, 219, 223, 224, 226, 227 230, 232, 234, 235, 236, 237, 238 239, 240, 241, 242, 243, 244, 245 248, 252, 254, 258, 259, 266, 267 270, 277, 278, 281, 283, 287, 288 289, 290, 292, 293, 294, 297, 298 302, 322, 332, 339, 344, 350, 352 353, 354, 355, 360, 361, 364, 383

*continued in next column*



**Fig. 3.** Residues in InouA, colored by their relative importance. Clockwise: front, back, top and bottom views.



**Fig. 4.** Residues in InouA, 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
		387, 404, 405, 406, 407, 422, 424 444, 446, 450, 451, 452, 460, 464 467, 468, 483, 484, 485, 487, 489
<i>continued in next column</i>		

Table 1. continued		
cluster color	size	member residues
		491, 503, 504, 505, 508, 510, 511 514, 515, 528, 529

**Table 1.** Clusters of top ranking residues in InouA.

2.4.2 *Overlap with known functional surfaces at 25% coverage.*  
The name of the ligand is composed of the source PDB identifier and the heteroatom name used in that file.

**Glycerol binding site.** Table 2 lists the top 25% of residues at the interface with InouGOL300 (glycerol). 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)
149	Y	Y (97) H (1) F (1)	0.05	9/9	4.15
162	A	G (13) A (69) S (13) V (1) . (1)	0.22	14/13	3.67

**Table 2.** The top 25% of residues in InouA at the interface with glycerol. (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.)

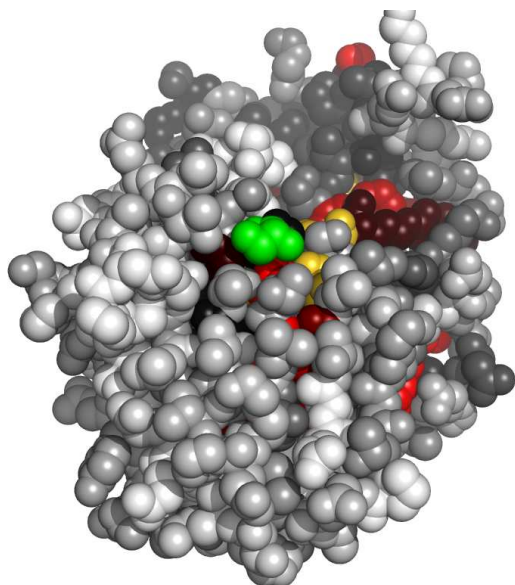
Table 3.		
res	type	disruptive mutations
149	Y	(K) (Q) (EM) (N)
162	A	(KR) (E) (Y) (Q)

**Table 3.** List of disruptive mutations for the top 25% of residues in InouA, that are at the interface with glycerol.

Figure 5 shows residues in InouA colored by their importance, at the interface with InouGOL300.

**Glycerol binding site.** Table 4 lists the top 25% of residues at the interface with InouGOL301 (glycerol). The following table (Table 5) suggests possible disruptive replacements for these residues (see Section 3.6).

Table 4.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
504	P	P (95)	0.06	8/3	3.37
<i>continued in next column</i>					



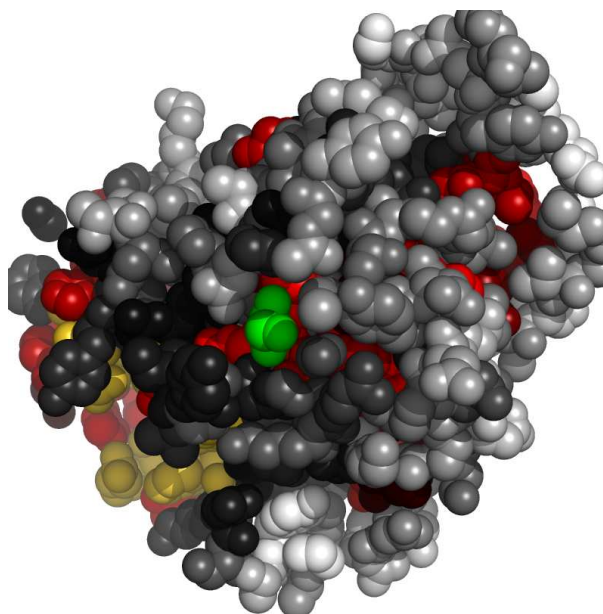
**Fig. 5.** Residues in 1nouA, at the interface with glycerol, colored by their relative importance. The ligand (glycerol) is colored green. Atoms further than 30Å away from the geometric center of the ligand, as well as on the line of sight to the ligand were removed. (See Appendix for the coloring scheme for the protein chain 1nouA.)

Table 4. continued					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
505	R	Q(3)	0.06	13/3	3.28
		.(1)			
		R(97)			
		T(1)			
460	W	.(1)	0.14	17/4	3.12
		L(17)			
		W(68)			
		V(4)			
		T(1)			
		K(1)			
		.(2)			
		G(1)			
		M(2)			
I(1)					

**Table 4.** The top 25% of residues in 1nouA at the interface with glycerol. (Field names: res: residue number in the PDB entry; type: amino acid type; substs: 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 5.		
res	type	disruptive mutations
504	P	(Y)(R)(TH)(CG)
505	R	(D)(LPI)(E)(VA)
460	W	(E)(K)(D)(TR)

**Table 5.** List of disruptive mutations for the top 25% of residues in 1nouA, that are at the interface with glycerol.



**Fig. 6.** Residues in 1nouA, at the interface with glycerol, colored by their relative importance. The ligand (glycerol) is colored green. Atoms further than 30Å away from the geometric center of the ligand, as well as on the line of sight to the ligand were removed. (See Appendix for the coloring scheme for the protein chain 1nouA.)

Figure 6 shows residues in 1nouA colored by their importance, at the interface with 1nouGOL301.

**Interface with 1nouB.** Table 6 lists the top 25% of residues at the interface with 1nouB. The following table (Table 7) suggests possible disruptive replacements for these residues (see Section 3.6).

Table 6.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
252	P	P(98) G(1)	0.03	3/0	4.16
213	Y	F(63) P(3) Y(30) A(1) W(1)	0.13	1/0	4.62

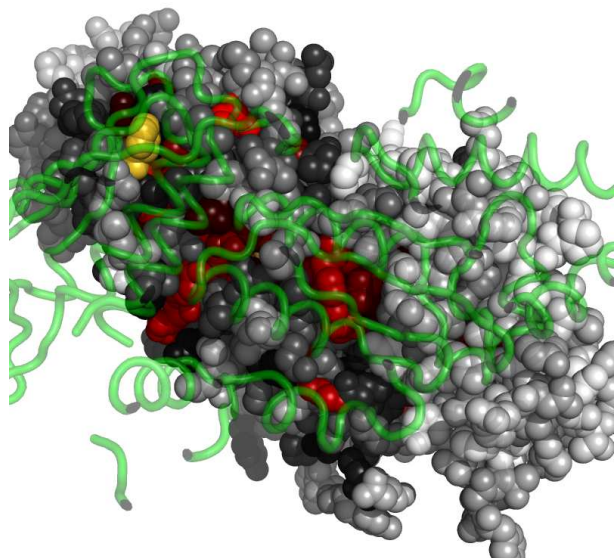
*continued in next column*

res	type	subst's (%)	cvg	noc/ bb	dist (Å)
270	D	D(66) Q(17) F(4) E(10) A(1)	0.18	15/1	2.65
267	T	T(72) S(25) N(1) R(1)	0.23	22/2	2.75

**Table 6.** The top 25% of residues in InouA at the interface with InouB. (Field names: res: residue number in the PDB entry; type: amino acid type; substs: 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.)

res	type	disruptive mutations
252	P	(R)(Y)(H)(KE)
213	Y	(K)(Q)(E)(R)
270	D	(R)(H)(FYW)(K)
267	T	(FKWR)(MH)(EQ)(LPI)

**Table 7.** List of disruptive mutations for the top 25% of residues in InouA, that are at the interface with InouB.



**Fig. 7.** Residues in InouA, at the interface with InouB, colored by their relative importance. InouB is shown in backbone representation (See Appendix for the coloring scheme for the protein chain InouA.)

Figure 7 shows residues in InouA colored by their importance, at the interface with InouB.

**Interface with InouB1.** By analogy with InouB – InouB1 interface. Table 8 lists the top 25% of residues at the interface with InouB1. The following table (Table 9) suggests possible disruptive replacements for these residues (see Section 3.6).

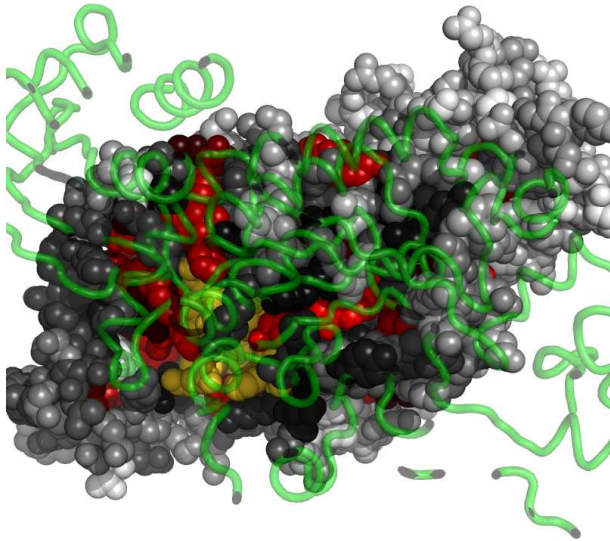
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
211	R	R(100)	0.02	19/16	3.13
491	E	E(98) D(1)	0.02	36/24	3.00
452	D	D(84) N(11) E(1) T(1) C(1)	0.10	15/8	3.30
213	Y	F(63) P(3) Y(30) A(1) W(1)	0.13	1/1	4.96
212	H	H(62) N(30) Q(2) T(4)	0.14	48/2	2.91
242	Q	Q(67) L(1) E(16) S(3) N(2) P(6) A(2)	0.14	7/0	3.81

**Table 8.** The top 25% of residues in InouA at the interface with InouB1. (Field names: res: residue number in the PDB entry; type: amino acid type; substs: 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.)

res	type	disruptive mutations
211	R	(TD)(SYEVCLAPIG)(FMW)(N)
491	E	(FWH)(R)(YVCAG)(T)
452	D	(R)(FWH)(K)(Y)
213	Y	(K)(Q)(E)(R)
212	H	(E)(TMD)(VQA)(SKCLPIG)
242	Q	(Y)(H)(FW)(T)

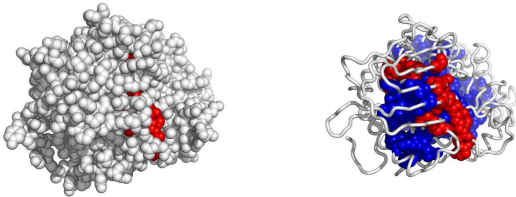
**Table 9.** List of disruptive mutations for the top 25% of residues in InouA, that are at the interface with InouB1.

Figure 8 shows residues in InouA colored by their importance, at the interface with InouB1.



**Fig. 8.** Residues in InouA, at the interface with InouB1, colored by their relative importance. InouB1 is shown in backbone representation (See Appendix for the coloring scheme for the protein chain InouA.)

2.4.3 *Possible novel functional surfaces at 25% coverage.* One group of residues is conserved on the InouA surface, away from (or substantially larger than) other functional sites and interfaces recognizable in PDB entry Inou. It is shown in Fig. 9. The right panel shows (in blue) the rest of the larger cluster this surface belongs to.



**Fig. 9.** A possible active surface on the chain InouA. The larger cluster it belongs to is shown in blue.

The residues belonging to this surface "patch" are listed in Table 10, while Table 11 suggests possible disruptive replacements for these residues (see Section 3.6).

Table 10.			
res	type	substitutions(%)	cvg
196	D	D(98)S(1)	0.02
198	P	P(100)	0.02
232	N	N(100)	0.02
<i>continued in next column</i>			

Table 10. <i>continued</i>			
res	type	substitutions(%)	cvg
149	Y	Y(97)H(1)F(1)	0.05
147	E	E(82)G(16) . (1)	0.08
224	D	D(86)E(4)K(1) N(5)V(2)	0.12
277	Y	Y(88)H(1)F(8) A(1)R(1)	0.12
468	P	P(82)I(2)F(9) L(2)D(1)V(2)	0.12
515	S	S(62)F(1)T(22) H(11)D(1)A(1)	0.12
200	F	F(66)Y(29)T(1) L(1)H(2)	0.15
60	P	P(83)E(3)D(5) A(2) . (4)	0.16
281	R	R(87)Y(2)L(2) K(3)H(3)N(1)	0.16
199	R	R(70)A(2)I(1) L(10)S(2)K(5) N(4)H(1)V(1)	0.20
170	R	W(6)Y(25)H(39) R(22)N(1)A(2) Q(1) . (1)	0.21
227	A	A(75)S(18)E(1) T(3)G(1)	0.21
162	A	G(13)A(69)S(13) V(1) . (1)	0.22
467	E	E(32)D(38)N(9) Q(1)V(2)A(8) M(2)P(2)R(2) S(1)	0.22
169	L	F(37)R(9)I(4) L(36)Y(1)T(3) M(3)V(2)K(1) . (1)	0.24
283	I	I(69)V(23)L(4) . (2)	0.24

**Table 10.** Residues forming surface "patch" in InouA.

Table 11.		
res	type	disruptive mutations
196	D	(R) (FWH) (K) (Y)
198	P	(YR) (TH) (SKECG) (FQWD)
232	N	(Y) (FTWH) (SEVCARG) (MD)
149	Y	(K) (Q) (EM) (N)
147	E	(FWH) (R) (VA) (Y)
224	D	(R) (FWH) (Y) (CG)
277	Y	(K) (Q) (E) (M)
468	P	(R) (Y) (H) (T)
515	S	(K) (R) (Q) (M)
200	F	(K) (E) (Q) (D)
60	P	(R) (Y) (H) (T)
<i>continued in next column</i>		

Table 11. continued		
res	type	disruptive mutations
281	R	(T)(D)(E)(SY)
199	R	(TY)(D)(E)(CG)
170	R	(TD)(E)(SCG)(YVA)
227	A	(R)(K)(YH)(E)
162	A	(KR)(E)(Y)(Q)
467	E	(H)(Y)(FW)(R)
169	L	(Y)(R)(T)(H)
283	I	(YR)(H)(T)(KE)

Table 11. Disruptive mutations for the surface patch in InouA.

Another group of surface residues is shown in Fig.10. The right panel shows (in blue) the rest of the larger cluster this surface belongs to.

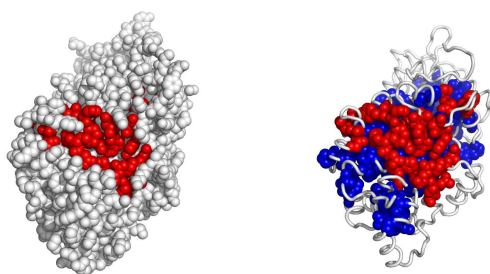


Fig. 10. Another possible active surface on the chain InouA. The larger cluster it belongs to is shown in blue.

The residues belonging to this surface "patch" are listed in Table 12, while Table 13 suggests possible disruptive replacements for these residues (see Section 3.6).

Table 12.			
res	type	substitutions(%)	cvg
208	D	D(100)	0.02
211	R	R(100)	0.02
450	Y	Y(100)	0.02
489	W	W(100)	0.02
491	E	E(98)D(1)	0.02
424	W	W(97)Y(2)	0.03
294	H	H(98).(1)	0.05
354	D	D(98)R(1)	0.05
355	E	E(98)I(1)	0.05
504	P	P(95)Q(3).(1)	0.06
505	R	R(97)T(1).(1)	0.06
243	S	G(47)T(3)S(47) A(1)	0.07
258	G	G(83)D(9)S(1) A(5)	0.08
407	E	E(65)D(25)Q(4)	0.09

*continued in next column*

Table 12. continued			
res	type	substitutions(%)	cvg
452	D	S(4) D(84)N(11)E(1) T(1)C(1)	0.10
213	Y	F(63)P(3)Y(30) A(1)W(1)	0.13
212	H	H(62)N(30)Q(2) T(4)	0.14
242	Q	Q(67)L(1)E(16) S(3)N(2)P(6) A(2)	0.14
293	G	G(76)A(19).(1) S(2)	0.14
460	W	L(17)W(68)V(4) T(1)K(1).(2) G(1)M(2)I(1)	0.14
270	D	D(66)Q(17)F(4) E(10)A(1)	0.18
406	Q	D(20)E(37)N(4) G(1)Q(23)H(5) T(2)S(3)A(1)	0.18
241	D	D(56)A(9)S(20) G(2)T(8)E(2)	0.23
259	S	S(31)A(53)M(3) G(11)	0.23
267	T	T(72)S(25)N(1) R(1)	0.23
207	I	L(61)V(17)I(15) T(1)M(2)F(2)	0.24
219	I	V(26)F(3)L(10) I(58)M(1)	0.25
447	A	P(33)A(19)T(2) S(13)N(17)H(3) D(6)I(1)E(1)	0.25
488	L	L(59)S(4)M(17) F(2)Q(1)I(10) V(3)T(1)	0.25

Table 12. Residues forming surface "patch" in InouA.

Table 13.		
res	type	disruptive mutations
208	D	(R)(FWH)(KYVCAG)(TQM)
211	R	(TD)(SYEVCLAPIG)(FMW)(N)
450	Y	(K)(QM)(NEVLAPIR)(D)
489	W	(KE)(TQD)(SNCRG)(M)
491	E	(FWH)(R)(YVCAG)(T)
424	W	(K)(E)(Q)(D)
294	H	(E)(TQMD)(SNVCLAPIG)(K)
354	D	(FW)(YVCAHRG)(T)(KM)
355	E	(H)(FYWR)(CG)(TVA)
504	P	(Y)(R)(TH)(CG)
505	R	(D)(LPI)(E)(VA)

*continued in next column*

Table 13. continued		
res	type	disruptive mutations
243	S	(KR) (QH) (FMW) (E)
258	G	(R) (K) (EH) (FQW)
407	E	(FWH) (YR) (VCAG) (T)
452	D	(R) (FWH) (K) (Y)
213	Y	(K) (Q) (E) (R)
212	H	(E) (TMD) (VQA) (SKCLPIG)
242	Q	(Y) (H) (FW) (T)
293	G	(KR) (E) (QH) (FMWD)
460	W	(E) (K) (D) (TR)
270	D	(R) (H) (FYW) (K)
406	Q	(Y) (FWH) (T) (VA)
241	D	(R) (H) (FW) (K)
259	S	(R) (K) (H) (FYQW)
267	T	(FKWR) (MH) (EQ) (LPI)
207	I	(R) (Y) (H) (TK)
219	I	(R) (Y) (T) (H)
447	A	(R) (Y) (K) (EH)
488	L	(R) (Y) (H) (T)

**Table 13.** Disruptive mutations for the surface patch in 1nouA.

### 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  $10\text{\AA}^2$ , 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  $5\text{\AA}$  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  $5\text{\AA}$ .

#### 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 [*EKRQM*], OH-group possession [*SDETY*], and NH<sub>2</sub> 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.

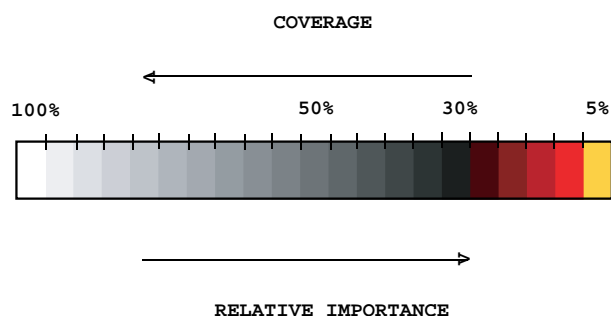
### 4 APPENDIX

#### 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
- `variability` has two subfields:
  1. number of different amino acids appearing in in this column of the alignment
  2. their type





**Fig. 11.** Coloring scheme used to color residues by their relative importance.

- `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. 11.

## 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  $10\text{\AA}^2$ , 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<sup>A</sup>T<sub>E</sub>X; 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:

- InouA.complex.pdb - coordinates of InouA with all of its interacting partners
- InouA.etvx - ET viewer input file for InouA
- InouA.cluster\_report.summary - Cluster report summary for InouA
- InouA.ranks - Ranks file in sequence order for InouA
- InouA.clusters - Cluster descriptions for InouA
- InouA.msf - the multiple sequence alignment used for the chain InouA
- InouA.descr - description of sequences used in InouA msf
- InouA.ranks\_sorted - full listing of residues and their ranking for InouA
- InouA.InouGOL300.if.pml - Pymol script for Figure 5
- InouA.cbcvg - used by other InouA – related pymol scripts
- InouA.InouGOL301.if.pml - Pymol script for Figure 6
- InouA.InouB.if.pml - Pymol script for Figure 7
- InouA.InouB1.if.pml - Pymol script for Figure 8