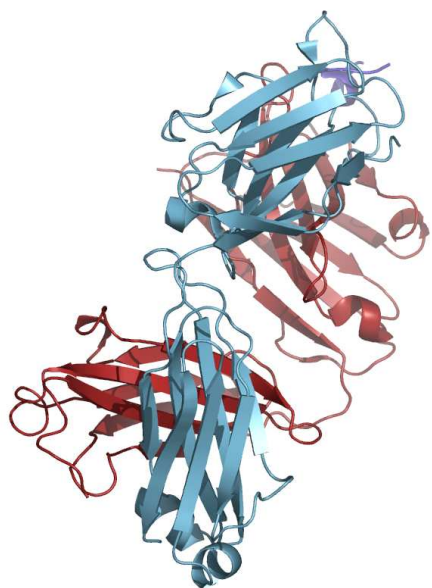


# lacy

Evolutionary trace report by **report\_maker**

July 26, 2010



## CONTENTS

<b>1</b>	<b>Introduction</b>
<b>2</b>	<b>Chain lacyH</b>
2.1	Q5M842 overview
2.2	Multiple sequence alignment for lacyH
2.3	Residue ranking in lacyH
2.4	Top ranking residues in lacyH 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.
<b>3</b>	<b>Chain lacyL</b>
3.1	Q66JS7 overview
3.2	Multiple sequence alignment for lacyL
3.3	Residue ranking in lacyL
3.4	Top ranking residues in lacyL and their position on the structure
3.4.1	Clustering of residues at 25% coverage.
3.4.2	Overlap with known functional surfaces at 25% coverage.
3.4.3	Possible novel functional surfaces at 25% coverage.

<b>4</b>	<b>Notes on using trace results</b>	<b>9</b>
4.1	Coverage	9
4.2	Known substitutions	9
4.3	Surface	9
4.4	Number of contacts	10
4.5	Annotation	10
4.6	Mutation suggestions	10
<b>5</b>	<b>Appendix</b>	<b>10</b>
5.1	File formats	10
5.2	Color schemes used	10
5.3	Credits	10
5.3.1	<b>Alistat</b>	10
5.3.2	<b>CE</b>	10
5.3.3	<b>DSSP</b>	11
5.3.4	<b>HSSP</b>	11
5.3.5	<b>LaTeX</b>	11
5.3.6	<b>Muscle</b>	11
5.3.7	<b>Pymol</b>	11
5.4	Note about ET Viewer	11
5.5	Citing this work	11
5.6	About report_maker	11
5.7	Attachments	11

1

## 1 INTRODUCTION

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

1 **Title:** Crystal structure of the principal neutralizing site of hiv- 1

1 **Compound:** Mol id: 1; molecule: igg1-kappa 59.1 fab (light chain);  
 1 chain: l; engineered: yes; mol id: 2; molecule: igg1-kappa 59.1 fab  
 1 (heavy chain); chain: h; engineered: yes; mol id: 3; molecule: hiv-  
 2 1 gp120 (mn isolate); chain: p; fragment: fragment (residues 308 -  
 2 332); engineered: yes

2 **Organism, scientific name:** Human Immunodeficiency Virus Type  
 1

3 lacy contains unique chains lacyH (221 residues) and lacyL (215  
 5 residues) Chain lacyP is too short (10 residues) to permit statistically  
 5 significant analysis, and was treated as a peptide ligand.

## 2 CHAIN 1ACYH

### 2.1 Q5M842 overview

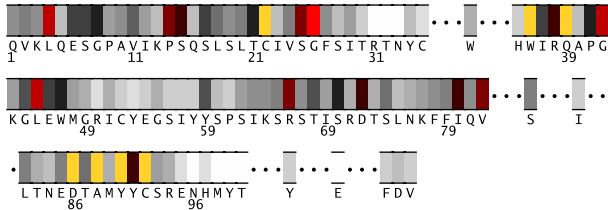
6 From SwissProt, id Q5M842, 63% identical to lacyH:

6 **Description:** Gamma-2a immunoglobulin heavy chain.

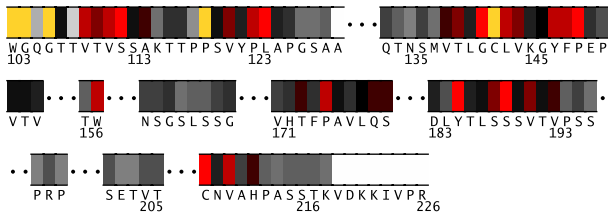
6 **Organism, scientific name:** Rattus norvegicus (Rat).

6 **Taxonomy:** Eukaryota; Metazoa; Chordata; Craniata; Verte-  
 7 brata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Glires;  
 7 Rodentia; Sciurognathi; Muridae; Murinae; Rattus.

1



**Fig. 1.** Residues 1-102 in lacyH colored by their relative importance. (See Appendix, Fig.14, for the coloring scheme.) Note that some residues in lacyH carry insertion code.



**Fig. 2.** Residues 103-226 in lacyH colored by their relative importance. (See Appendix, Fig.14, for the coloring scheme.) Note that some residues in lacyH carry insertion code.

## 2.2 Multiple sequence alignment for lacyH

For the chain lacyH, the alignment lacyH.msf (attached) with 218 sequences was used. The alignment was downloaded from the HSSP database, and fragments shorter than 75% of the query as well as duplicate sequences were removed. It can be found in the attachment to this report, under the name of lacyH.msf. Its statistics, from the *alistat* program are the following:

```

Format:                MSF
Number of sequences:   218
Total number of residues: 40014
Smallest:              87
Largest:               221
Average length:        183.6
Alignment length:      221
Average identity:      42%
Most related pair:     99%
Most unrelated pair:   0%
Most distant seq:      31%

```

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

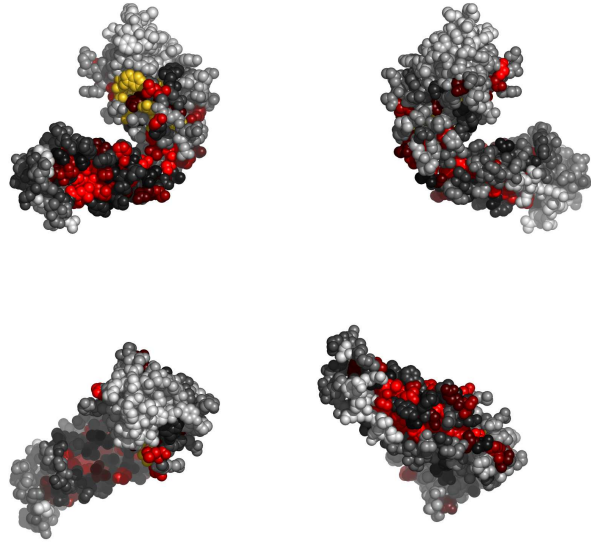
The alignment consists of 99% eukaryotic (99% vertebrata) sequences. (Descriptions of some sequences were not readily available.) The file containing the sequence descriptions can be found in the attachment, under the name lacyH.descr.

## 2.3 Residue ranking in lacyH

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

## 2.4 Top ranking residues in lacyH 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 lacyH 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.

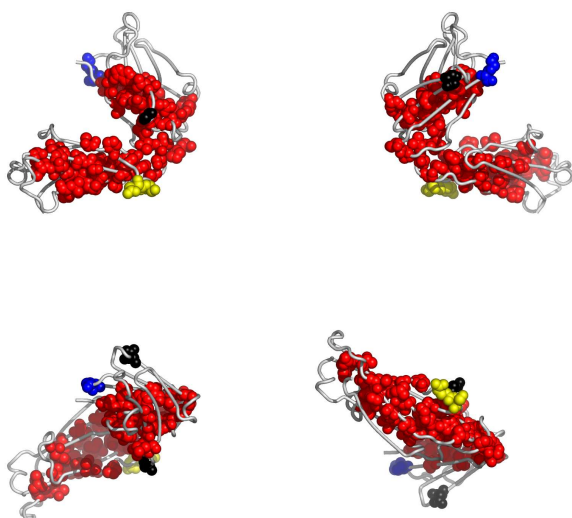


**Fig. 3.** Residues in lacyH, colored by their relative importance. Clockwise: front, back, top and bottom views.

**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	49	4, 14, 15, 22, 36, 38, 39, 45, 66, 80, 82, 86, 88, 90, 91, 92, 103, 104, 106, 109, 110, 111, 112, 114, 119, 121, 123, 124, 138, 139, 141, 142, 143, 144, 146, 147, 148, 149, 157, 173, 175, 185, 188, 189, 191, 193, 208, 210, 212
blue	2	25, 26
yellow	2	179, 180

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



**Fig. 4.** Residues in lacyH, 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.

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.

**Interface with lacyL.** Table 2 lists the top 25% of residues at the interface with lacyL. The following table (Table 3) suggests possible disruptive replacements for these residues (see Section 4.6).

Table 2.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
104	G	G(81) . (18)	0.02	4/4	4.05
103	W	W(77) . (18) F(3)S	0.03	60/2	3.40
39	Q	Q(91) . (5) K(2)HDR	0.04	13/0	2.90
189	S	S(76) . (20) T(1)IAW E	0.06	9/9	3.66
124	L	L(71) . (20) F(3) G(2)R I(1)Q	0.08	31/12	3.48
141	G	G(59)	0.08	6/6	4.22

*continued in next column*

Table 2. continued					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
143	L	. (22) A(12) V(4)PST L(73) . (21)TS VQIFRK	0.10	3/0	4.33
45	L	L(86)M . (4) P(6)ITH	0.11	60/15	3.50
123	P	P(69) . (20) L(2) S(3)NIE GTVA	0.12	19/8	3.29
175	P	P(70) . (20) R(1)L K(4)TES Q	0.13	23/9	2.86
139	T	T(25) A(21) . (22) I(15) V(12)Q S(1) C(4) S(59) . (20) T(11)QI HFNRYE	0.19	25/8	3.38
188	S	C(4) S(59) . (20) T(11)QI HFNRYE	0.19	8/0	3.39
179	Q	A(3) Q(32) . (20) S(28)V T(5) N(1)KP R(3) I(1)HD N(16) T(37) . (20) F(14)R S(3)V Q(1)H G(2)A	0.22	15/0	3.22
91	Y	F(20) Y(73) . (3)SW L(1) P(3) L(32) . (20)	0.24	21/0	3.61
178	L	L(1) P(3) L(32) . (20)	0.25	1/1	4.48

*continued in next column*

res	type	subst's (%)	cvg	noc/ bb	dist (Å)
		A(14) Q(15) Y(1) R(5) M(3)I G(1)D			

**Table 2.** The top 25% of residues in lacyH at the interface with lacyL. (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
104	G	(KER)(FQMWHD)(NLPI)(Y)
103	W	(K)(E)(Q)(D)
39	Q	(Y)(T)(FW)(VCAG)
189	S	(R)(K)(H)(Q)
124	L	(Y)(R)(T)(H)
141	G	(R)(K)(E)(H)
143	L	(Y)(R)(H)(T)
45	L	(R)(Y)(H)(T)
123	P	(R)(Y)(H)(K)
175	P	(Y)(R)(H)(T)
139	T	(R)(K)(H)(FW)
188	S	(KR)(FMWH)(Q)(E)
179	Q	(Y)(H)(FW)(T)
173	T	(K)(R)(QH)(M)
91	Y	(K)(Q)(R)(EM)
178	L	(Y)(R)(H)(T)

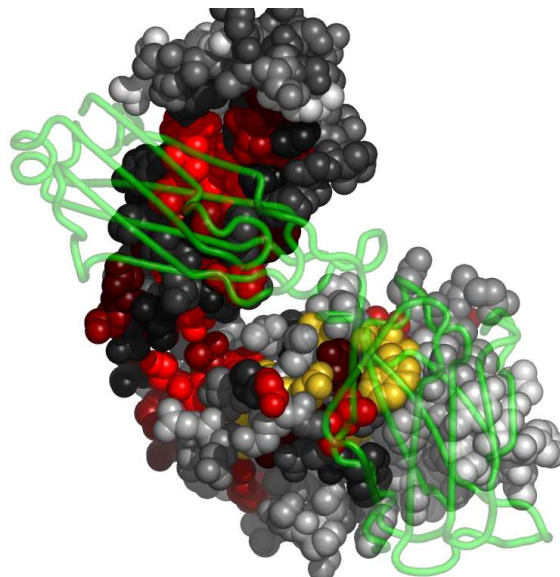
**Table 3.** List of disruptive mutations for the top 25% of residues in lacyH, that are at the interface with lacyL.

Figure 5 shows residues in lacyH colored by their importance, at the interface with lacyL.

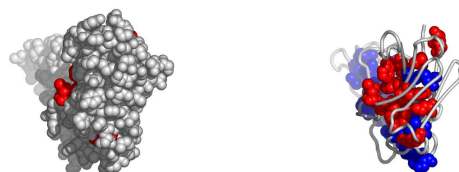
**2.4.3 Possible novel functional surfaces at 25% coverage.** One group of residues is conserved on the lacyH surface, away from (or substantially larger than) other functional sites and interfaces recognizable in PDB entry lacy. It is shown in Fig. 6. The right panel shows (in blue) the rest of the larger cluster this surface belongs to. The residues belonging to this surface "patch" are listed in Table 4, while Table 5 suggests possible disruptive replacements for these residues (see Section 4.6).

res	type	substitutions(%)	cvg
90	Y	Y(94).(3)FL	0.01
104	G	G(81).(18)	0.02

*continued in next column*



**Fig. 5.** Residues in lacyH, at the interface with lacyL, colored by their relative importance. lacyL is shown in backbone representation (See Appendix for the coloring scheme for the protein chain lacyH.)



**Fig. 6.** A possible active surface on the chain lacyH. The larger cluster it belongs to is shown in blue.

res	type	substitutions(%)	cvg
106	G	G(81).(18)	0.02
36	W	W(93).(5)LR	0.03
103	W	W(77).(18)F(3)S	0.03
39	Q	Q(91).(5)K(2)HD	0.04
		R	
88	A	A(90)G(4).(4)CV	0.04
		T	
26	G	G(87).(7)ADTSKV	0.07
45	L	L(86)M.(4)P(6)I	0.11
		TH	
4	L	L(85).(11)V(2)I	0.13
		H	
25	S	S(84).(5)T(3)	0.18
		V(2)NHFYEP	

*continued in next column*

Table 4. continued			
res	type	substitutions(%)	cvg
66	R	R(72)K(18).(3) P(1)G(1)QSH	0.18
38	R	K(16)R(71).(5) Q(4)MPISW	0.20
80	I	L(71)M(22).(3)I F(1)V	0.24
91	Y	F(20)Y(73).(3)S WL(1)	0.24

Table 4. Residues forming surface "patch" in lacyH.

Table 5.		
res	type	disruptive mutations
90	Y	(K)(Q)(EMR)(N)
104	G	(KER)(FQMWHD)(NLPI)(Y)
106	G	(KER)(FQMWHD)(NLPI)(Y)
36	W	(E)(T)(KD)(SQCG)
103	W	(K)(E)(Q)(D)
39	Q	(Y)(T)(FW)(VCAG)
88	A	(KR)(E)(Y)(Q)
26	G	(R)(K)(E)(H)
45	L	(R)(Y)(H)(T)
4	L	(R)(Y)(T)(E)
25	S	(KR)(Q)(M)(H)
66	R	(TD)(Y)(E)(SVCAG)
38	R	(T)(Y)(D)(CG)
80	I	(YR)(T)(H)(KE)
91	Y	(K)(Q)(R)(EM)

Table 5. Disruptive mutations for the surface patch in lacyH.

Another group of surface residues is shown in Fig.7. The right panel shows (in blue) the rest of the larger cluster this surface belongs to.

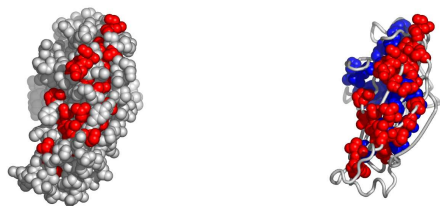


Fig. 7. Another possible active surface on the chain lacyH. The larger cluster it belongs to is shown in blue.

The residues belonging to this surface "patch" are listed in Table 6, while Table 7 suggests possible disruptive replacements for these residues (see Section 4.6).

Table 6.			
res	type	substitutions(%)	cvg
119	P	P(73).(21)G(3)L QD	0.05
112	S	S(70).(18)T(7) K(1)YLV	0.06
189	S	S(76).(20)T(1)I AWE	0.06
124	L	L(71).(20)F(3) G(2)RI(1)Q	0.08
141	G	G(59).(22)A(12) V(4)PST	0.08
185	Y	Y(72).(20)WHC S(1)F(2)T	0.09
111	V	V(78).(17)LI(1) MG	0.10
143	L	L(73).(21)TSVQI FRK	0.10
123	P	P(69).(20)L(2) S(3)NIEGTVA	0.12
157	W	W(74).(20)CQVHL RN	0.12
175	P	P(70).(20)R(1)L K(4)TESQ	0.13
148	F	F(56).(21)MY(3) IL(13)S(1)VC H(1)	0.14
121	V	I(18)V(50).(20) G(2)L(6)HK	0.15
110	T	T(74).(17)LE(1) DI(1)S(2)Q	0.16
14	P	P(85).(6)V(1) M(2)LTS(1)A	0.17
139	T	T(25)A(21).(22) I(15)V(12)QS(1)	0.19
188	S	C(4)S(59).(20) T(11)QIHFNRYE	0.19
180	S	S(36).(20)G(33) D(1)N(2)QER(1)Y PI	0.20
114	A	E(22)A(45).(22) VG(5)IPTD	0.21
138	V	V(57)A(17).(22) TIS	0.21
179	Q	A(3)Q(32).(20) S(28)VT(5)N(1)K PR(3)I(1)HD	0.22
15	S	G(71).(6)S(15) T(5)EAKD	0.23
173	T	N(16)T(37).(20) F(14)RS(3)VQ(1) HG(2)A	0.23

continued in next column

res	type	substitutions(%)	cvg
193	V	L(40)V(30).(21) I(2)RDPKQT(1)SG	0.23
146	G	N(4)D(41).(21) G(28)A(1)KITQ	0.25
178	L	P(3)L(32).(20) A(14)Q(15)Y(1) R(5)M(3)IG(1)D	0.25

Table 6. Residues forming surface "patch" in lacyH.

res	type	disruptive mutations
119	P	(Y)(R)(H)(T)
112	S	(R)(K)(H)(FW)
189	S	(R)(K)(H)(Q)
124	L	(Y)(R)(T)(H)
141	G	(R)(K)(E)(H)
185	Y	(K)(Q)(M)(E)
111	V	(Y)(R)(KE)(H)
143	L	(Y)(R)(H)(T)
123	P	(R)(Y)(H)(K)
157	W	(E)(K)(TD)(Q)
175	P	(Y)(R)(H)(T)
148	F	(K)(E)(Q)(D)
121	V	(Y)(E)(R)(K)
110	T	(R)(H)(K)(FW)
14	P	(R)(Y)(H)(K)
139	T	(R)(K)(H)(FW)
188	S	(KR)(FMWH)(Q)(E)
180	S	(R)(K)(H)(FW)
114	A	(R)(K)(Y)(H)
138	V	(R)(K)(Y)(E)
179	Q	(Y)(H)(FW)(T)
15	S	(R)(K)(FWH)(QM)
173	T	(K)(R)(QH)(M)
193	V	(Y)(R)(E)(H)
146	G	(R)(H)(FW)(E)
178	L	(Y)(R)(H)(T)

Table 7. Disruptive mutations for the surface patch in lacyH.

### 3 CHAIN 1ACYL

#### 3.1 Q66JS7 overview

From SwissProt, id Q66JS7, 81% identical to lacyL:

**Description:** LOC243469 protein.

**Organism, scientific name:** *Mus musculus* (Mouse).

**Taxonomy:** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Glires; Rodentia; Sciurognathi; Muridae; Murinae; Mus.

#### 3.2 Multiple sequence alignment for lacyL

For the chain lacyL, the alignment lacyL.msf (attached) with 149 sequences was used. The alignment was downloaded from the HSSP

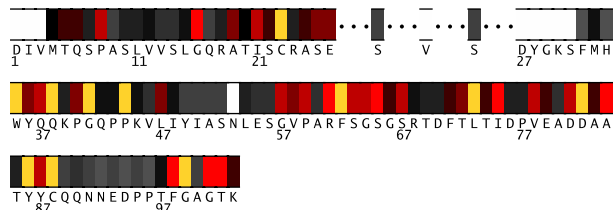


Fig. 8. Residues 1-103 in lacyL colored by their relative importance. (See Appendix, Fig.14, for the coloring scheme.) Note that some residues in lacyL carry insertion code.

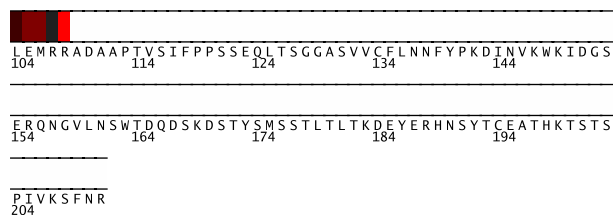


Fig. 9. Residues 104-211 in lacyL colored by their relative importance. (See Appendix, Fig.14, for the coloring scheme.) Note that some residues in lacyL carry insertion code.

database, and fragments shorter than 75% of the query as well as duplicate sequences were removed. It can be found in the attachment to this report, under the name of lacyL.msf. Its statistics, from the *alistat* program are the following:

```

Format:                               MSF
Number of sequences: 149
Total number of residues: 23785
Smallest: 70
Largest: 215
Average length: 159.6
Alignment length: 215
Average identity: 39%
Most related pair: 98%
Most unrelated pair: 0%
Most distant seq: 33%

```

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

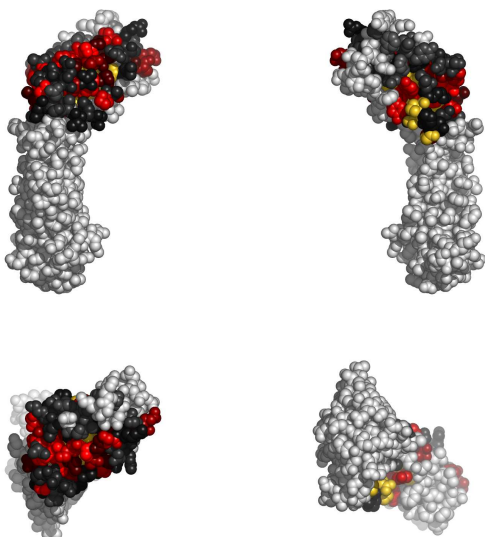
The alignment consists of 98% eukaryotic (98% vertebrata), and <1% viral sequences. (Descriptions of some sequences were not readily available.) The file containing the sequence descriptions can be found in the attachment, under the name lacyL.descr.

#### 3.3 Residue ranking in lacyL

The lacyL sequence is shown in Figs. 8–9, with each residue colored according to its estimated importance. The full listing of residues in lacyL can be found in the file called lacyL.ranks\_sorted in the attachment.

### 3.4 Top ranking residues in lacyL and their position on the structure

In the following we consider residues ranking among top 25% of residues in the protein. Figure 10 shows residues in lacyL 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.



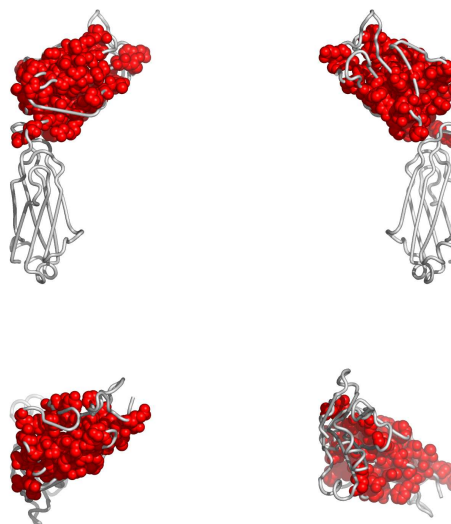
**Fig. 10.** Residues in lacyL, colored by their relative importance. Clockwise: front, back, top and bottom views.

**3.4.1 Clustering of residues at 25% coverage.** Fig. 11 shows the top 25% of all residues, this time colored according to clusters they belong to. The clusters in Fig.11 are composed of the residues listed in Table 8.

Table 8.		
cluster color	size	member residues
red	54	4, 5, 6, 8, 16, 19, 20, 21, 22, 23, 25 26, 27, 35, 36, 37, 38, 40, 41, 44 47, 57, 58, 59, 61, 62, 63, 64, 65 66, 67, 71, 72, 73, 74, 75, 78, 79 81, 82, 83, 84, 86, 87, 88, 98, 99 101, 102, 103, 104, 105, 106, 108

**Table 8.** Clusters of top ranking residues in lacyL.

**3.4.2 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.



**Fig. 11.** Residues in lacyL, 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.

**Interface with lacyH.** Table 9 lists the top 25% of residues at the interface with lacyH. The following table (Table 10) suggests possible disruptive replacements for these residues (see Section 4.6).

Table 9.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
44	P	L(25) P(67) . (6)S	0.03	49/19	3.40
38	Q	Q(89) . (4) H(2)K W(1) L(1)	0.04	17/0	2.90
98	F	. (28) F(67) W(4)Y	0.06	57/0	3.59
87	Y	Y(70) . (4) F(20) H(1)S L(2)	0.10	34/0	3.50
36	Y	V(22) Y(57) F(8) . (4) L(2)	0.15	38/0	3.00

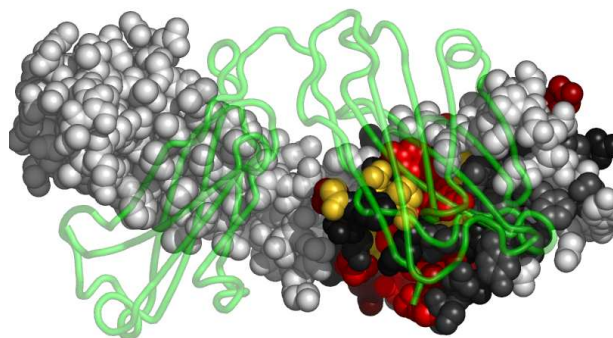
*continued in next column*

Table 9. continued					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
		I ( 4 ) H ( 1 )			

**Table 9.** The top 25% of residues in lacyL at the interface with lacyH. (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 10.		
res	type	disruptive mutations
44	P	( R ) ( Y ) ( H ) ( T )
38	Q	( Y ) ( T ) ( S ) ( C G )
98	F	( K ) ( E ) ( Q ) ( D )
87	Y	( K ) ( Q ) ( M ) ( E )
36	Y	( K ) ( Q ) ( E ) ( R )

**Table 10.** List of disruptive mutations for the top 25% of residues in lacyL, that are at the interface with lacyH.



**Fig. 12.** Residues in lacyL, at the interface with lacyH, colored by their relative importance. lacyH is shown in backbone representation (See Appendix for the coloring scheme for the protein chain lacyL.)

Figure 12 shows residues in lacyL colored by their importance, at the interface with lacyH.

**3.4.3 Possible novel functional surfaces at 25% coverage.** One group of residues is conserved on the lacyL surface, away from (or

substantially larger than) other functional sites and interfaces recognizable in PDB entry 1acy. It is shown in Fig. 13. The right panel shows (in blue) the rest of the larger cluster this surface belongs to.



**Fig. 13.** A possible active surface on the chain lacyL. The larger cluster it belongs to is shown in blue.

The residues belonging to this surface "patch" are listed in Table 11, while Table 12 suggests possible disruptive replacements for these residues (see Section 4.6).

Table 11.			
res	type	substitutions(%)	cvg
41	G	G ( 89 ) . ( 4 ) HP ( 2 ) E ( 1 ) D	0.03
44	P	L ( 25 ) P ( 67 ) . ( 6 ) S	0.03
38	Q	Q ( 89 ) . ( 4 ) H ( 2 ) K W ( 1 ) L ( 1 )	0.04
99	G	. ( 27 ) G ( 72 )	0.04
62	F	F ( 81 ) . ( 4 ) W ( 8 ) L ( 3 ) I ( 2 ) A	0.05
98	F	. ( 28 ) F ( 67 ) W ( 4 ) Y	0.06
101	G	. ( 27 ) G ( 71 ) Y	0.06
16	G	G ( 87 ) . ( 8 ) QT ( 2 ) R SD	0.07
61	R	R ( 80 ) . ( 4 ) S ( 6 ) K ( 4 ) TA ( 1 ) GYH	0.07
84	A	A ( 83 ) . ( 4 ) G ( 11 ) V	0.07
65	S	S ( 85 ) . ( 4 ) GT ( 2 ) L ( 2 ) A ( 2 ) DQR ( 1 )	0.08
37	Q	R ( 26 ) Q ( 60 ) . ( 4 ) L ( 4 ) SPH ( 1 ) IYK	0.10
87	Y	Y ( 70 ) . ( 4 ) F ( 20 ) H ( 1 ) SL ( 2 )	0.10
63	S	I ( 5 ) S ( 69 ) . ( 4 ) T ( 10 ) A ( 8 ) VML	0.11
8	P	G ( 4 ) P ( 69 ) . ( 12 ) S ( 9 ) E ( 2 ) HIA	0.12
64	G	V ( 2 ) G ( 66 ) . ( 4 ) L ( 1 ) I ( 11 ) K ( 7 ) W ( 1 ) A ( 2 ) FSRP	0.12
57	G	S ( 10 ) G ( 67 ) . ( 4 ) R	0.13

*continued in next column*



Table 11. continued			
res	type	substitutions(%)	cvg
59	P	W(6)N(1)A(4) T(4)ED K(8)P(64).(4) N(3)S(4)RHY(7)D Q(4)T(1)E(1)	0.14
67	S	N(4)S(73)D(9) . (4)T(2)F(2)E A(2)L(1)	0.14
81	D	E(66).(5)D(6) A(12)G(6)KBTN	0.14
36	Y	V(22)Y(57)F(8) . (4)L(2)I(4) H(1)	0.15
40	P	P(73)A(4).(4) S(8)L(4)V(1) T(2)QM	0.15
26	S	S(72)T(6).(6) D(2)Q(1)G(2) R(2)N(6)I(1)	0.16
27	E	G(26)S(18)Q(26) . (6)E(5)D(1) N(6)R(2)T(2) Y(3)K	0.17
58	V	V(52).(4)A(2) I(18)T(6)Y(7)DR F(3)WL(4)	0.17
47	L	I(23)M(4)L(52) . (4)W(2)TV(11)A R	0.18
72	T	N(3)S(22)Y(8) T(52).(4)A(1) E(2)V(2)I(1) F(1)D	0.19
105	E	. (28)T(32)D(4) E(30)L(1)ISNG	0.19
6	Q	Q(60).(22)T(8) E(4)K(1)IV(1)SR A	0.20
79	E	K(6)Q(43)L(2) . (4)E(23)T(12) R(4)VZD	0.20
83	A	S(6)E(24)D(4) F(16).(4)L(4) T(22)M(2)V(6) A(6)I(2)	0.20
66	G	R(16)K(16)G(37) . (4)A(1)S(5) N(6)T(2)MQ(2) I(2)L(4)	0.21
104	L	. (28)L(42)V(28) F	0.21
74	T	Q(10)T(47).(4)	0.22

continued in next column

Table 11. continued			
res	type	substitutions(%)	cvg
22	S	K(17)S(6)EA(5)D V(1)N(4)LIR T(53)S(34).(4) C(2)YA(1)R(1) P(1)E	0.23
103	K	. (28)K(45)Q(7)T R(11)LM(1)N(1) S(1)PAG	0.23
5	T	T(55).(27)L(9) I(1)S(1)DV(1)K A(2)	0.24
25	A	A(46)G(18)L(4) . (6)S(10)T(3) V(8)EP(1)F	0.24
4	M	L(31).(27)S(2) M(24)V(11)Q(1) T(1)	0.25
20	T	T(62)R(10).(4) S(10)EI(1)PK(6) L(1)A(1)	0.25

Table 11. Residues forming surface "patch" in lacyL.

Table 12.		
res	type	disruptive mutations
41	G	(R)(K)(FEWH)(Q)
44	P	(R)(Y)(H)(T)
38	Q	(Y)(T)(S)(CG)
99	G	(KER)(FQMWHD)(NLPI)(Y)
62	F	(KE)(T)(Q)(D)
98	F	(K)(E)(Q)(D)
101	G	(K)(E)(R)(M)
16	G	(R)(FW)(K)(H)
61	R	(D)(E)(LPI)(T)
84	A	(KER)(Y)(QHD)(N)
65	S	(R)(K)(H)(FW)
37	Q	(Y)(T)(FW)(H)
87	Y	(K)(Q)(M)(E)
63	S	(R)(K)(H)(YQ)
8	P	(R)(Y)(H)(TK)
64	G	(E)(R)(K)(D)
57	G	(R)(K)(E)(H)
59	P	(Y)(R)(H)(T)
67	S	(R)(K)(H)(FW)
81	D	(R)(FW)(H)(Y)
36	Y	(K)(Q)(E)(R)
40	P	(R)(Y)(H)(T)
26	S	(R)(FKWH)(YM)(Q)
27	E	(FW)(H)(Y)(VAR)
58	V	(K)(R)(E)(Y)

continued in next column

Table 12. <i>continued</i>		
res	type	disruptive mutations
47	L	(Y)(R)(H)(T)
72	T	(R)(K)(H)(Q)
105	E	(H)(FW)(R)(Y)
6	Q	(Y)(H)(FW)(T)
79	E	(FWH)(Y)(R)(CG)
83	A	(R)(Y)(K)(H)
66	G	(R)(E)(H)(K)
104	L	(R)(Y)(T)(KEH)
74	T	(R)(H)(FW)(K)
22	S	(K)(R)(H)(Q)
103	K	(Y)(FW)(T)(H)
5	T	(R)(K)(H)(FW)
25	A	(R)(K)(Y)(E)
4	M	(Y)(H)(R)(T)
20	T	(R)(H)(K)(FW)

**Table 12.** Disruptive mutations for the surface patch in lacyL.

## 4 NOTES ON USING TRACE RESULTS

### 4.1 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.

### 4.2 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%.

### 4.3 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.)

### 4.4 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  $5\text{\AA}$ .

### 4.5 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).

### 4.6 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.

## 5 APPENDIX

### 5.1 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

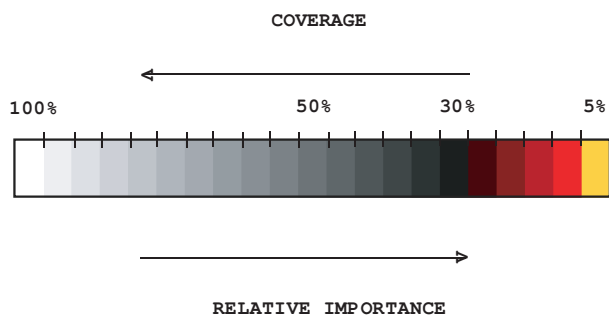


Fig. 14. Coloring scheme used to color residues by their relative importance.

- 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

## 5.2 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. 14.

## 5.3 Credits

**5.3.1 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.

**5.3.2 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.

**5.3.3 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>.

**5.3.4 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/>

**5.3.5 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).

**5.3.6 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/>

**5.3.7 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.)

## 5.4 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.

## 5.5 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.

## 5.6 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.

## 5.7 Attachments

The following files should accompany this report:

- lacyH.complex.pdb - coordinates of lacyH with all of its interacting partners
- lacyH.etvx - ET viewer input file for lacyH
- lacyH.cluster\_report.summary - Cluster report summary for lacyH
- lacyH.ranks - Ranks file in sequence order for lacyH
- lacyH.clusters - Cluster descriptions for lacyH
- lacyH.msf - the multiple sequence alignment used for the chain lacyH
- lacyH.descr - description of sequences used in lacyH msf
- lacyH.ranks\_sorted - full listing of residues and their ranking for lacyH
- lacyH.lacyL.if.pml - Pymol script for Figure 5
- lacyH.cbvvg - used by other lacyH – related pymol scripts
- lacyL.complex.pdb - coordinates of lacyL with all of its interacting partners
- lacyL.etvx - ET viewer input file for lacyL
- lacyL.cluster\_report.summary - Cluster report summary for lacyL
- lacyL.ranks - Ranks file in sequence order for lacyL
- lacyL.clusters - Cluster descriptions for lacyL
- lacyL.msf - the multiple sequence alignment used for the chain lacyL
- lacyL.descr - description of sequences used in lacyL msf
- lacyL.ranks\_sorted - full listing of residues and their ranking for lacyL
- lacyL.lacyH.if.pml - Pymol script for Figure 12
- lacyL.cbvvg - used by other lacyL – related pymol scripts