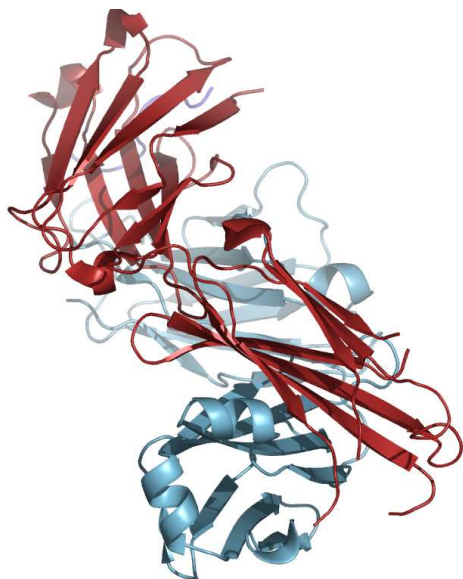


# 2hkf

Evolutionary trace report by **report\_maker**

March 17, 2010



## CONTENTS

<b>1</b>	<b>Introduction</b>	
<b>2</b>	<b>Chain 2hkfH</b>	
2.1	Q66K04 overview	
2.2	Multiple sequence alignment for 2hkfH	
2.3	Residue ranking in 2hkfH	
2.4	Top ranking residues in 2hkfH 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 2hkfL</b>	
3.1	Q65ZC0 overview	
3.2	Multiple sequence alignment for 2hkfL	
3.3	Residue ranking in 2hkfL	
3.4	Top ranking residues in 2hkfL 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>10</b>
4.1	Coverage	10
4.2	Known substitutions	10
4.3	Surface	10
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	11
5.3	Credits	11
5.3.1	<b>Alistat</b>	11
5.3.2	<b>CE</b>	11
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	12
5.7	Attachments	12

## 1 INTRODUCTION

1 From the original Protein Data Bank entry (PDB id 2hkf):  
 1 **Title:** Crystal structure of the complex fab m75- peptide  
 1 **Compound:** Mol id: 1; molecule: immunoglobulin light chain fab  
 1 fragment; chain: l; mol id: 2; molecule: immunoglobulin heavy chain  
 1 fab fragment; chain: h; mol id: 3; molecule: carbonic anhydrase 9;  
 1 chain: p; synonym: carbonic anhydrase ix, carbonate dehydratase ix,  
 2 ca-ix, caix, membrane antigen mn, p54/58n, renal cell carcinoma-  
 2 associated antigen g250, rcc-associated antigen g250, pmw1; ec:  
 2 4.2.1.1; engineered: yes  
**Organism, scientific name:** Mus Musculus;  
 4 2hkf contains unique chains 2hkfH (210 residues) and 2hkfL (219  
 residues) Chain 2hkfP is too short (9 residues) to permit statistically  
 6 significant analysis, and was treated as a peptide ligand.

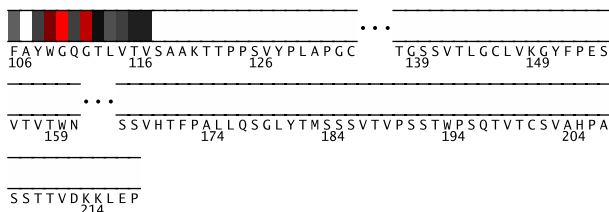
## 2 CHAIN 2HKFH

### 2.1 Q66K04 overview

6 From SwissProt, id Q66K04, 73% identical to 2hkfH:  
 7 **Description:** Hypothetical protein.  
**Organism, scientific name:** Mus musculus (Mouse).  
 7 **Taxonomy:** Eukaryota; Metazoa; Chordata; Craniata; Verte-  
 8 brata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Glires;  
 Rodentia; Sciurognathi; Muridae; Murinae; Mus.



**Fig. 1.** Residues 1-105 in 2hkfH colored by their relative importance. (See Appendix, Fig.14, for the coloring scheme.)



**Fig. 2.** Residues 106-218 in 2hkfH colored by their relative importance. (See Appendix, Fig.14, for the coloring scheme.)

## 2.2 Multiple sequence alignment for 2hkfH

For the chain 2hkfH, the alignment 2hkfH.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 2hkfH.msf. Its statistics, from the *alostat* program are the following:

```

Format:                MSF
Number of sequences:   218
Total number of residues: 28883
Smallest:              74
Largest:               210
Average length:        132.5
Alignment length:      210
Average identity:       40%
Most related pair:     98%
Most unrelated pair:   15%
Most distant seq:      36%

```

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

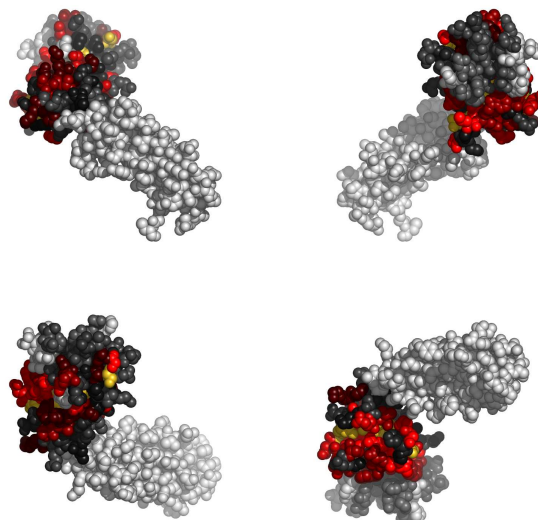
The alignment consists of 51% eukaryotic ( 51% vertebrata) sequences. (Descriptions of some sequences were not readily available.) The file containing the sequence descriptions can be found in the attachment, under the name 2hkfH.descr.

## 2.3 Residue ranking in 2hkfH

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

## 2.4 Top ranking residues in 2hkfH 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 2hkfH 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 2hkfH, colored by their relative importance. Clockwise: front, back, top and bottom views.

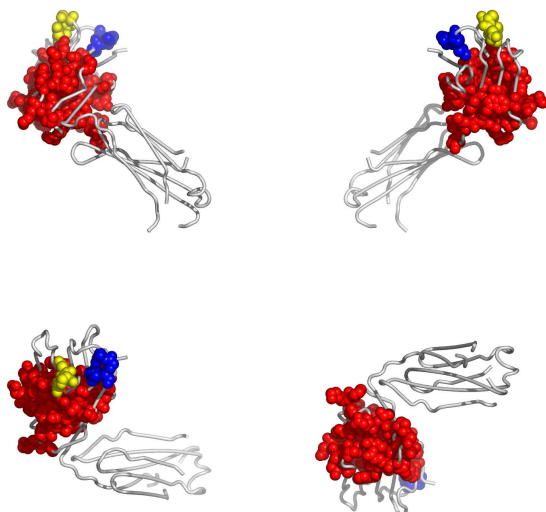
**2.4.1 Clustering of residues at 25% coverage.** Fig. 4 shows the top 25% of all residues, this time colored according to clusters they belong to. The clusters in Fig.4 are composed of the residues listed in Table 1.

cluster color	size	member residues
red	45	7, 14, 15, 17, 18, 19, 20, 21, 22, 36, 37, 38, 39, 41, 42, 45, 46, 47, 48, 49, 62, 65, 66, 67, 68, 69, 70, 71, 72, 73, 82, 83, 84, 85, 88, 91, 92, 93, 94, 96, 97, 98, 109, 110, 112
blue	4	24, 25, 26, 27
yellow	3	75, 76, 77

**Table 1.** Clusters of top ranking residues in 2hkfH.

**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 the peptide 2hkfP.** Table 2 lists the top 25% of residues at the interface with 2hkfP. The following table (Table



**Fig. 4.** Residues in 2hkfH, 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.

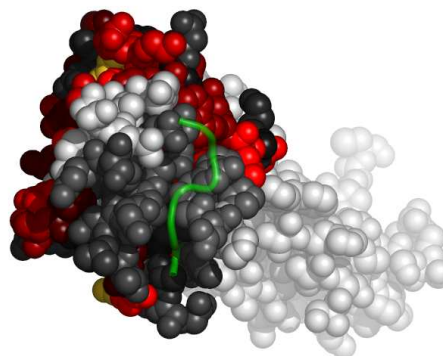
3) suggests possible disruptive replacements for these residues (see Section 4.6).

Table 2.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
26	G	G (54) S (9) Q (28) L (1) AT E (1) R. Y (1) KD	0.10	3/3	4.03
27	S	EF (32) S (29) Y (14) L (1) D (5) Q (1) HAP G (8) . IN T	0.24	13/8	3.31

**Table 2.** The top 25% of residues in 2hkfH at the interface with 2hkfP. (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
26	G	(R) (K) (FEWH) (M)
27	S	(R) (K) (H) (FQW)

**Table 3.** List of disruptive mutations for the top 25% of residues in 2hkfH, that are at the interface with 2hkfP.



**Fig. 5.** Residues in 2hkfH, at the interface with 2hkfP, colored by their relative importance. 2hkfP is shown in backbone representation (See Appendix for the coloring scheme for the protein chain 2hkfH.)

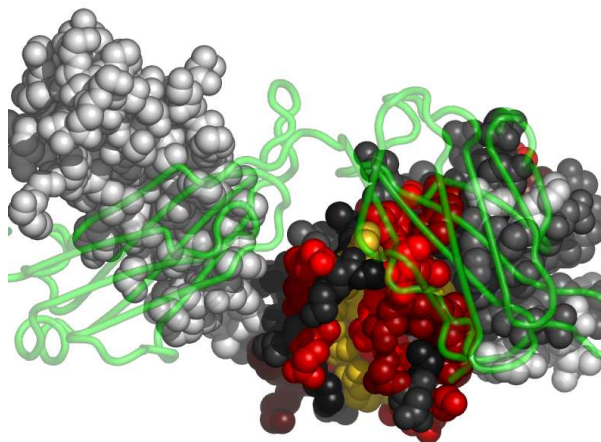
Figure 5 shows residues in 2hkfH colored by their importance, at the interface with 2hkfP.

**Interface with 2hkfL.** Table 4 lists the top 25% of residues at the interface with 2hkfL. The following table (Table 5) suggests possible disruptive replacements for these residues (see Section 4.6).

Table 4.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
39	Q	HQ (96) K (1) ER	0.03	12/0	2.78
45	L	P (44) L (50) F (1) AVM IX	0.05	57/14	3.58
110	G	G (84) . (14) A	0.09	5/5	3.65
97	Y	Y (83) SL F (12)	0.11	19/0	3.91

*continued in next column*

res	type	subst's (%)	cvg	noc/ bb	dist (Å)
37	V	H(1)W L(1) V(42) Y(39) I(10) F(4)HA	0.13	4/0	4.30
47	W	S(1) W(49) L(33) R(2) F(1)C Y(5) T(2)GMP NH	0.15	66/6	3.56
62	Y	.(25)H Y(48) R(11)A L(10)KP F(1)IC	0.16	1/1	4.96
109	W	F(39) W(44) . (15)R	0.17	54/3	3.45
46	E	Q(5) E(49) K(18) R(19) V(1)DZ T(2)INW	0.19	5/5	4.58



**Fig. 6.** Residues in 2hkfH, at the interface with 2hkfL, colored by their relative importance. 2hkfL is shown in backbone representation (See Appendix for the coloring scheme for the protein chain 2hkfH.)

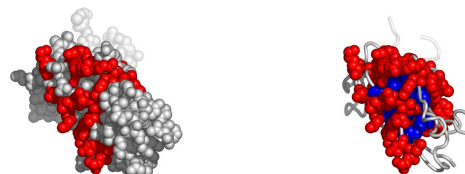
2.4.3 Possible novel functional surfaces at 25% coverage. One group of residues is conserved on the 2hkfH surface, away from (or substantially larger than) other functional sites and interfaces recognizable in PDB entry 2hkf. It is shown in Fig. 7. The right panel shows (in blue) the rest of the larger cluster this surface belongs to.

**Table 4.** The top 25% of residues in 2hkfH at the interface with 2hkfL. (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
39	Q	(Y)(T)(FW)(VCAG)
45	L	(YR)(TH)(KE)(SCG)
110	G	(KER)(HD)(Q)(FMW)
97	Y	(K)(Q)(E)(M)
37	V	(KE)(R)(Y)(Q)
47	W	(E)(K)(D)(Q)
62	Y	(K)(E)(Q)(M)
109	W	(E)(TKD)(Q)(SCG)
46	E	(H)(FW)(Y)(R)

**Table 5.** List of disruptive mutations for the top 25% of residues in 2hkfH, that are at the interface with 2hkfL.

Figure 6 shows residues in 2hkfH colored by their importance, at the interface with 2hkfL.



**Fig. 7.** A possible active surface on the chain 2hkfH. 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).

res	type	substitutions(%)	cvg
39	Q	HQ(96)K(1)ER	0.03
69	R	R(85)K(13)N	0.03
45	L	P(44)L(50)F(1)A	0.05

*continued in next column*

Table 6. continued			
res	type	substitutions(%)	cvg
		VMIX	
94	A	G(8)A(90).T	0.05
42	G	G(90)E(3)D(2) K(1)SH	0.06
73	S	S(82)T(15)ELRAF	0.06
38	R	Q(40)R(42)K(11) L(4)HDY	0.07
72	I	G(45)I(33)L(14) V(1)F(1)A(1)P M(1).	0.07
88	L	V(14)L(79)A(2)T IM(1)FP	0.08
110	G	G(84).(14)A	0.09
21	S	S(69)C(3)T(22) I(1)REP.N	0.10
91	E	E(82)D(5)A(9)N G(1)XK	0.10
97	Y	Y(83)SLF(12) H(1)W	0.11
112	G	G(84).(15)N	0.11
71	T	T(45)S(47)VRNE A(1)M.I(1)K	0.12
75	D	.(9)D(53)E(2) S(28)T(1)BIY(1) ANF	0.12
37	V	L(1)V(42)Y(39) I(10)F(4)HA	0.13
41	P	P(84)H(3)LA(3) S(5)RT(1)Q	0.14
67	K	P(43)K(37)N(5)G IS(4)R(1)EQ(4)T M	0.14
47	W	S(1)W(49)L(33) R(2)F(1)CY(5) T(2)GMPNH	0.15
68	D	S(19)G(38)D(29) A(6)K(1)PTN(3)E V	0.15
62	Y	.(25)HY(48) R(11)AL(10)KP F(1)IC	0.16
17	S	P(5)S(50)R(21) T(12)D(3)K(3) G(1)A.(1)	0.17
109	W	F(39)W(44).(15) R	0.17
65	S	G(46)S(24)K(10) W(5)FA(4)T(2)L N(1)DEP.R	0.18
77	S	S(49)A(9)T(30) D(3)G(1)N(1) P(1)QWYK	0.18

continued in next column

Table 6. continued			
res	type	substitutions(%)	cvg
46	E	Q(5)E(49)K(18) R(19)V(1)DZT(2) INW	0.19
93	T	A(3)T(40)E(10) S(11)L(2)F(20) D(1)V(5)CI(2) M(1)KG	0.19
18	L	L(36)V(35)A(17) M(1)R(1)HQI(3).	0.20
14	P	T(4)P(73)HL(2) S(11)EA(1)G.(2) QV	0.21
76	D	G(38)N(22).(4) K(9)I(3)D(4) T(12)AS(4)FR	0.21
19	K	S(14)R(16)T(36) K(23)E(3)V(1)QA N.IG	0.22
82	Y	T(7)Y(33)S(12) .(29)C(3)F(9)HR LD	0.22
15	K	Q(3)G(69)IS(9) V(6)K(1)A.(2) L(3)RDTE	0.23
84	Q	T(36)Q(33)E(4) K(11)S(2)ZGA I(1)H(1)D(2) N(2)RYF	0.23
7	S	S(72)P(13)A(3) T(3).(5)ED	0.25
89	K	Q(27)R(19)T(18) E(19)K(9)LZSG D(2)A	0.25

Table 6. Residues forming surface "patch" in 2hkhH.

Table 7.		
res	type	disruptive mutations
39	Q	(Y)(T)(FW)(VCAG)
69	R	(T)(Y)(D)(SVCAG)
45	L	(YR)(TH)(KE)(SCG)
94	A	(KR)(E)(YQH)(D)
42	G	(R)(FKW)(EH)(M)
73	S	(KR)(H)(Q)(FW)
38	R	(T)(D)(VCAG)(Y)
72	I	(R)(Y)(H)(T)
88	L	(R)(Y)(H)(TK)
110	G	(KER)(HD)(Q)(FMW)
21	S	(R)(K)(H)(FW)
91	E	(FWH)(Y)(R)(CG)
97	Y	(K)(Q)(E)(M)

continued in next column

Table 7. continued		
res	type	disruptive mutations
112	G	(R) (E) (K) (FWH)
71	T	(R) (H) (K) (FW)
75	D	(R) (H) (FKW) (Y)
37	V	(KE) (R) (Y) (Q)
41	P	(Y) (R) (H) (T)
67	K	(Y) (FW) (T) (H)
47	W	(E) (K) (D) (Q)
68	D	(R) (H) (FW) (Y)
62	Y	(K) (E) (Q) (M)
17	S	(R) (K) (FWH) (YQM)
109	W	(E) (TKD) (Q) (SCG)
65	S	(R) (K) (H) (FW)
77	S	(R) (K) (H) (FW)
46	E	(H) (FW) (Y) (R)
93	T	(R) (K) (H) (FW)
18	L	(Y) (R) (T) (H)
14	P	(R) (Y) (H) (K)
76	D	(R) (H) (FW) (Y)
19	K	(Y) (FW) (T) (H)
82	Y	(K) (Q) (M) (R)
15	K	(Y) (FW) (T) (H)
84	Q	(Y) (FW) (H) (T)
7	S	(R) (K) (H) (FW)
89	K	(Y) (FW) (T) (H)

**Table 7.** Disruptive mutations for the surface patch in 2hkfL.

### 3 CHAIN 2HKFL

#### 3.1 Q65ZC0 overview

From SwissProt, id Q65ZC0, 89% identical to 2hkfL:

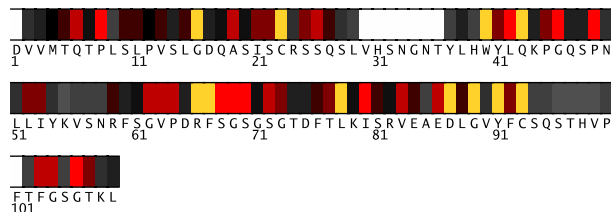
**Description:** Kappa light chain C region (Fragment).

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

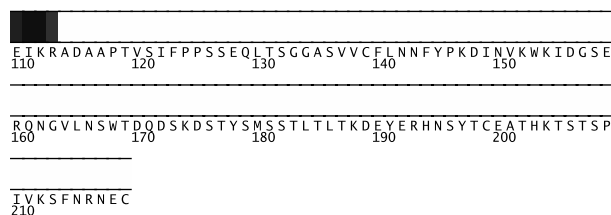
**Taxonomy:** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Glires; Rodentia; Sciurognathi; Muridae; Murinae; Mus.

#### 3.2 Multiple sequence alignment for 2hkfL

For the chain 2hkfL, the alignment 2hkfL.msf (attached) with 247 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 2hkfL.msf. Its statistics, from the *alstat* program are the following:



**Fig. 8.** Residues 1-109 in 2hkfL colored by their relative importance. (See Appendix, Fig.14, for the coloring scheme.)



**Fig. 9.** Residues 110-219 in 2hkfL colored by their relative importance. (See Appendix, Fig.14, for the coloring scheme.)

```

Format: MSF
Number of sequences: 247
Total number of residues: 38614
Smallest: 73
Largest: 219
Average length: 156.3
Alignment length: 219
Average identity: 47%
Most related pair: 99%
Most unrelated pair: 11%
Most distant seq: 35%

```

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

The alignment consists of 47% eukaryotic ( 47% 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 2hkfL.descr.

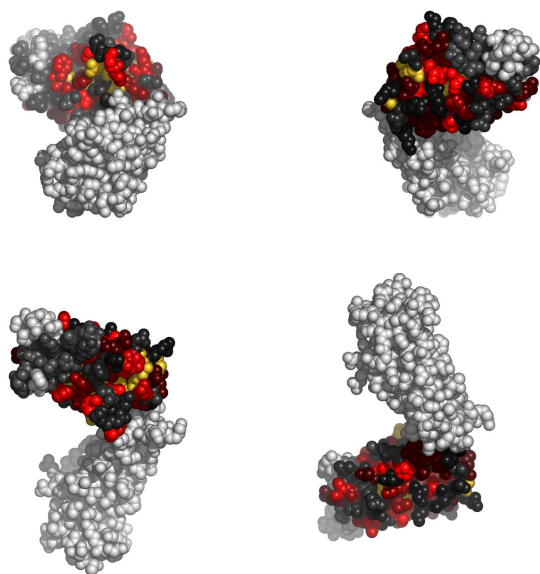
#### 3.3 Residue ranking in 2hkfL

The 2hkfL sequence is shown in Figs. 8–9, with each residue colored according to its estimated importance. The full listing of residues in 2hkfL can be found in the file called 2hkfL.ranks\_sorted in the attachment.

#### 3.4 Top ranking residues in 2hkfL and their position on the structure

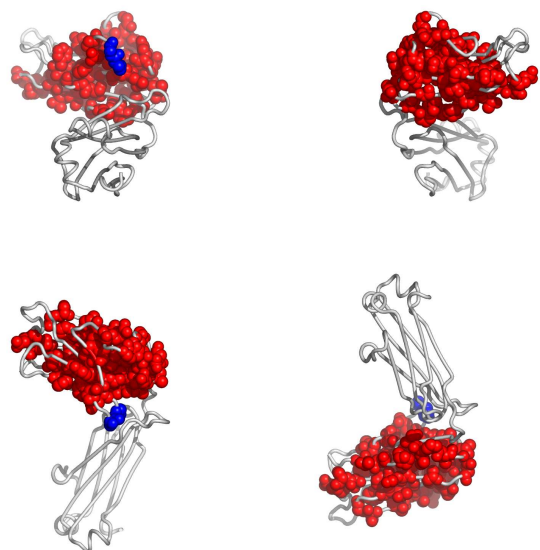
In the following we consider residues ranking among top 25% of residues in the protein . Figure 10 shows residues in 2hkfL colored by their importance: bright red and yellow indicate more conserved/important residues (see Appendix for the coloring scheme). A Pymol script for producing this figure can be found in the attachment.

*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



**Fig. 10.** Residues in 2hkfL, colored by their relative importance. Clockwise: front, back, top and bottom views.

belong to. The clusters in Fig.11 are composed of the residues listed



**Fig. 11.** Residues in 2hkfL, 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.

in Table 8.

Table 8.		
cluster color	size	member residues
red	53	4, 5, 6, 8, 10, 11, 12, 13, 15, 16, 19, 21, 22, 23, 25, 26, 27, 40, 41, 42, 43, 49, 52, 53, 59, 62, 63, 64, 66, 67, 68, 69, 70, 72, 73, 76, 77, 78, 80, 81, 83, 84, 86, 87, 88, 89, 91, 92, 93, 103, 104, 106, 107
blue	2	45, 46

**Table 8.** Clusters of top ranking residues in 2hkfL.

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.

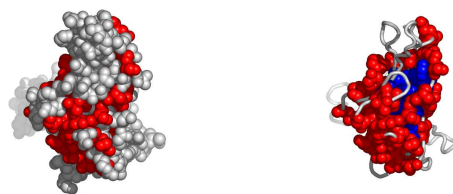
**Interface with 2hkfH.** Table 9 lists the top 25% of residues at the interface with 2hkfH. The following table (Table 10) suggests possible disruptive replacements for these residues (see Section 4.6).

Table 9.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
43	Q	Q(93) E(1) L(1)	0.05	17/0	2.78
46	G	H(2)WK G(90)H P(2) D(3)K E(1)S	0.06	1/1	4.48
49	P	P(82) L(10) F(1)A V(3)E.I	0.06	48/20	3.45
103	F	F(83) . (12)L W(3)V	0.14	49/1	3.58
92	F	Y(79) F(15) H(2) L(1)S.	0.16	27/0	4.02
41	Y	Y(72) L(2) V(10) F(9) I(2)A H(1)S	0.18	42/0	2.79

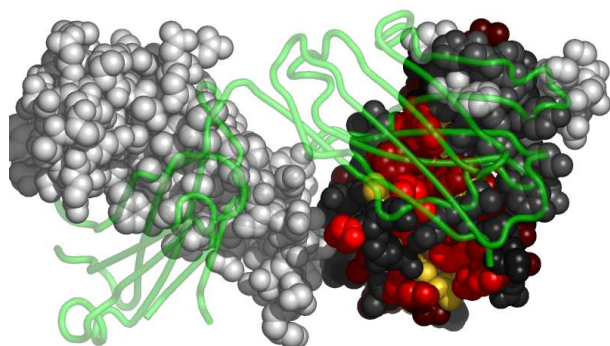
**Table 9.** The top 25% of residues in 2hkfL at the interface with 2hkfH. (Field names: res: residue number in the PDB entry; type: amino acid type; substs: substitutions seen in the alignment; with the percentage of each type in the bracket; noc/bb: number of contacts with the ligand, with the number of contacts realized through backbone atoms given in the bracket; dist: distance of closest approach to the ligand.)

Table 10.		
res	type	disruptive mutations
43	Q	(Y)(T)(FWH)(CG)
46	G	(R)(K)(FW)(H)
49	P	(R)(Y)(H)(T)
103	F	(KE)(Q)(TD)(R)
92	F	(K)(E)(Q)(D)
41	Y	(K)(Q)(ER)(M)

**Table 10.** List of disruptive mutations for the top 25% of residues in 2hkfL, that are at the interface with 2hkfH.



**Fig. 13.** A possible active surface on the chain 2hkfL. The larger cluster it belongs to is shown in blue.



**Fig. 12.** Residues in 2hkfL, at the interface with 2hkfH, colored by their relative importance. 2hkfH is shown in backbone representation (See Appendix for the coloring scheme for the protein chain 2hkfL.)

Figure 12 shows residues in 2hkfL colored by their importance, at the interface with 2hkfH.

**3.4.3 Possible novel functional surfaces at 25% coverage.** One group of residues is conserved on the 2hkfL surface, away from (or substantially larger than) other functional sites and interfaces recognizable in PDB entry 2hkf. It is shown in Fig. 13. 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
87	D	D(98)TE.	0.02
66	R	R(94)K(2)YLNUGH	0.03

*continued in next column*

Table 11. continued			
res	type	substitutions(%)	cvg
16	G	G(95)RD.S(1)EQ	0.04
89	G	A(84)G(14)S.	0.04
43	Q	Q(93)E(1)L(1) H(2)WK	0.05
67	F	F(93)L(3)NAT I(1)W	0.05
46	G	G(90)HP(2)D(3)K E(1)S	0.06
49	P	P(82)L(10)F(1)A V(3)E.I	0.06
69	G	G(84)L(3)A(2)K I(6)WDVP	0.06
70	S	S(89)T(3)LR(1)I A(1)NGFDEK	0.07
8	P	P(80).(8)S(4)T E(1)A(1)KIGH	0.08
68	S	S(82)T(12)K(1)V IERAML	0.08
42	L	Q(74)L(7)R(14) H(1)KDYW	0.09
106	G	G(86).(13)	0.09
6	Q	Q(80).(14)T(2)Z SIV	0.10
72	S	S(83)D(6)I(1) Y(2)A(1)LT(2)E. PF	0.10
86	E	E(82)A(5)D(4) G(6)XK.	0.10
19	A	V(51)A(31)L(10) I(4).S	0.11
26	S	T(7)S(80)AD(4)R N(3)IF.EQ	0.11
62	G	G(85)W(5)REPA S(4)KNM	0.12
104	G	G(85).(12)WAF	0.12
63	V	V(58)I(29)AT(4) L(3)NK(1)YFGR	0.13

*continued in next column*



Table 11. continued			
res	type	substitutions(%)	cvg
83	V	L(61)V(22)M(5) A(9)ITDS	0.13
64	P	P(80)N(4)R(1) S(4)Q(1)GTA(1)M K(3)ED	0.14
103	F	F(83).(12)LW(3) V	0.14
21	I	I(52)L(39)M(6)Q V.F	0.15
45	P	P(83)S(6)A(3) L(2)T(1)GRQM	0.15
22	S	S(55)C(2)T(36)E N(1)P(1)ARY.KH	0.16
92	F	Y(79)F(15)H(2) L(1)S.	0.16
25	S	G(20)S(13)A(50) V(4)L(4)R(1) T(1)EQP(1)WI.M	0.17
77	T	S(23)T(61)Y(6) A(2)I(1)V(1)HDF ER	0.17
41	Y	Y(72)L(2)V(10) F(9)I(2)AH(1)S	0.18
73	G	G(76)I(2)S(6) E(2)A(1)N(2)R K(4)DV(1).	0.18
52	L	M(4)L(66)I(14) W(2)V(9)SGFA	0.19
53	I	I(79)M(3)L(2) V(4)YENG(3)F S(1)W(1)TH	0.19
27	Q	S(21)Q(43)L(1) G(13)E(3)D(1) T(2)N(4)R(2) K(1)Y(1)A(1).	0.20
88	L	E(27)L(6)F(26) T(11)V(7)D(4) A(7)I(3).S(2) M(2)C	0.20
13	V	G(12)V(39)L(14) A(22).(2)T(2) E(1)RKIM(1)FS	0.21
59	R	R(56)L(23)G(3)F S(3)NK(4)Y(2)D T(1)PMIEQ.	0.21
5	T	T(74).(14)I(2) L(4)NEA(1)S(1)V	0.22
81	S	S(71)N(12)T(7) K(1)H(1)V(1)Y G(1)AQR	0.22
10	S	S(54)T(20).(8) L(1)A(3)F(1)	0.23

continued in next column

Table 11. continued			
res	type	substitutions(%)	cvg
15	L	I(3)PY(1)VEN G(1)H P(54)I(2)L(20) T(1)V(14).(1) A(2)RQKS	0.23
11	L	A(4)L(48).(8) F(1)M(6)E(1) V(25)SG(1)KQTDN	0.24
84	E	Q(45)E(36)T(5) L(1)ZK(4)VGR(5) D	0.24
4	M	L(44).(14)M(29) S(1)V(7)Q(1)KI	0.25
12	P	S(67).(8)P(7) A(5)T(4)L(2)IYK V(1)E	0.25

Table 11. Residues forming surface "patch" in 2hkfL.

Table 12.		
res	type	disruptive mutations
87	D	(R)(FWH)(KVA)(YCG)
66	R	(D)(T)(E)(Y)
16	G	(R)(FW)(H)(K)
89	G	(KR)(E)(QH)(FMWD)
43	Q	(Y)(T)(FWH)(CG)
67	F	(KE)(R)(TD)(Q)
46	G	(R)(K)(FW)(H)
49	P	(R)(Y)(H)(T)
69	G	(R)(KE)(H)(Y)
70	S	(R)(K)(H)(FW)
8	P	(R)(Y)(H)(T)
68	S	(R)(K)(H)(Y)
42	L	(Y)(T)(R)(CG)
106	G	(KER)(FQMWHD)(NLPI)(Y)
6	Q	(Y)(H)(FW)(T)
72	S	(R)(K)(H)(Q)
86	E	(FWH)(Y)(R)(VCAG)
19	A	(R)(Y)(K)(E)
26	S	(R)(K)(H)(FW)
62	G	(R)(E)(K)(H)
104	G	(K)(E)(R)(Q)
63	V	(E)(Y)(R)(K)
83	V	(R)(Y)(K)(H)
64	P	(Y)(R)(H)(T)
103	F	(KE)(Q)(TD)(R)
21	I	(Y)(R)(T)(H)
45	P	(Y)(R)(H)(T)
22	S	(R)(K)(FW)(H)
92	F	(K)(E)(Q)(D)

continued in next column

Table 12. continued		
res	type	disruptive mutations
25	S	(R) (K) (H) (Y)
77	T	(K) (R) (Q) (MH)
41	Y	(K) (Q) (ER) (M)
73	G	(R) (KE) (H) (FW)
52	L	(R) (Y) (K) (H)
53	I	(R) (Y) (KH) (T)
27	Q	(Y) (FWH) (T) (VA)
88	L	(R) (Y) (H) (K)
13	V	(Y) (R) (K) (E)
59	R	(T) (Y) (D) (CG)
5	T	(R) (K) (H) (FW)
81	S	(R) (K) (FWH) (EM)
10	S	(R) (K) (Q) (H)
15	L	(Y) (R) (H) (T)
11	L	(Y) (R) (H) (T)
84	E	(FWH) (Y) (R) (VA)
4	M	(Y) (H) (T) (R)
12	P	(R) (Y) (H) (K)

**Table 12.** Disruptive mutations for the surface patch in 2hkfL.

## 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

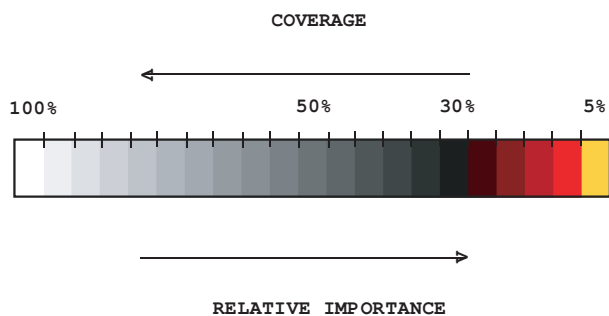


Fig. 14. Coloring scheme used to color residues by their relative importance.

- 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:

- 2hkfH.complex.pdb - coordinates of 2hkfH with all of its interacting partners
- 2hkfH.etvx - ET viewer input file for 2hkfH
- 2hkfH.cluster\_report.summary - Cluster report summary for 2hkfH
- 2hkfH.ranks - Ranks file in sequence order for 2hkfH
- 2hkfH.clusters - Cluster descriptions for 2hkfH
- 2hkfH.msf - the multiple sequence alignment used for the chain 2hkfH
- 2hkfH.descr - description of sequences used in 2hkfH msf
- 2hkfH.ranks\_sorted - full listing of residues and their ranking for 2hkfH
- 2hkfH.2hkfP.if.pml - Pymol script for Figure 5
- 2hkfH.cbv - used by other 2hkfH – related pymol scripts
- 2hkfH.2hkfL.if.pml - Pymol script for Figure 6
- 2hkfL.complex.pdb - coordinates of 2hkfL with all of its interacting partners
- 2hkfL.etvx - ET viewer input file for 2hkfL
- 2hkfL.cluster\_report.summary - Cluster report summary for 2hkfL
- 2hkfL.ranks - Ranks file in sequence order for 2hkfL
- 2hkfL.clusters - Cluster descriptions for 2hkfL
- 2hkfL.msf - the multiple sequence alignment used for the chain 2hkfL
- 2hkfL.descr - description of sequences used in 2hkfL msf
- 2hkfL.ranks\_sorted - full listing of residues and their ranking for 2hkfL
- 2hkfL.2hkfH.if.pml - Pymol script for Figure 12
- 2hkfL.cbv - used by other 2hkfL – related pymol scripts