

# 1k4r

Evolutionary trace report by **report\_maker**

February 19, 2010



4.3.1	<b>Alistat</b>	9
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	10
4.5	Citing this work	10
4.6	About report_maker	10
4.7	Attachments	10

## CONTENTS

### 1 Introduction

### 2 Chain 1k4rA

- 2.1 P14336 overview
- 2.2 Multiple sequence alignment for 1k4rA
- 2.3 Residue ranking in 1k4rA
- 2.4 Top ranking residues in 1k4rA 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

## 1 INTRODUCTION

From the original Protein Data Bank entry (PDB id 1k4r):

**Title:** Structure of dengue virus

**Compound:** Mol id: 1; molecule: major envelope protein e; chain: a, b, c

**Organism, scientific name:** Tick-borne Encephalitis Virus;

1k4r contains a single unique chain 1k4rA (395 residues long) and its homologues 1k4rC and 1k4rB.

## 2 CHAIN 1K4RA

### 2.1 P14336 overview

- 1 From SwissProt, id P14336, 96% identical to 1k4rA:
- 1 **Description:** Genome polyprotein [Contains: Capsid protein C (Core protein); Envelope protein M (Matrix protein); Major envelope protein E; Nonstructural protein 1 (NS1); Nonstructural protein 2A (NS2A); Flavivirin protease NS2B regulatory subunit; Flavivirin protease NS3 catalytic subunit (EC 3.4.21.91); Nonstructural protein 4A (NS4A); Nonstructural protein 4B (NS4B); RNA-directed RNA polymerase (EC 2.7.7.48) (NS5)].
- 2 **Organism, scientific name:** Tick-borne encephalitis virus (Western subtype) (TBEV).
- 3 **Taxonomy:** Viruses; ssRNA positive-strand viruses, no DNA stage; Flaviviridae; Flavivirus; tick-borne encephalitis virus group.
- 4 **Function:** The small proteins NS2A, NS4A and NS4B are hydrophobic, suggesting a possible membrane-related function. NS5 may play a role in the viral RNA replication. The NS2B/NS3 protease complex processes the viral polyprotein.
- 8 **Catalytic activity:** Selective hydrolysis of -Xaa-Xaa—Yaa- bonds in which each of the Xaa can be either Arg or Lys and Yaa can be either Ser or Ala.
- 8 **Catalytic activity:** Nucleoside triphosphate + RNA(n) = diphosphate + RNA(n+1).
- 9 **Subunit:** NS3 and NS2B form a heterodimer. NS3 is the catalytic subunit, whereas NS2B strongly stimulates the latter (By similarity).
- 9 **Ptm:** Specific enzymatic cleavages in vivo yield mature proteins (By similarity).
- 9

**Miscellaneous:** The virion of this virus is a nucleocapsid covered by a lipoprotein envelope. The envelope contains two proteins: the protein M and glycoprotein E. The nucleocapsid is a complex of protein C and mRNA. In immature particles, there are 60 icosaedrically organized trimeric spikes on the surface. Each spike consists of three heterodimers of envelope protein M precursor (prM) and envelope protein E (By similarity).

**Similarity:** Contains 1 peptidase S7 domain.

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

## 2.2 Multiple sequence alignment for 1k4rA

For the chain 1k4rA, the alignment 1k4rA.msf (attached) with 158 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 1k4rA.msf. Its statistics, from the *alistat* program are the following:

```

Format:                MSF
Number of sequences:   158
Total number of residues: 57971
Smallest:              137
Largest:               395
Average length:        366.9
Alignment length:      395
Average identity:       54%
Most related pair:     99%
Most unrelated pair:   1%
Most distant seq:      54%

```

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

The alignment consists of 65% 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 1k4rA.descr.

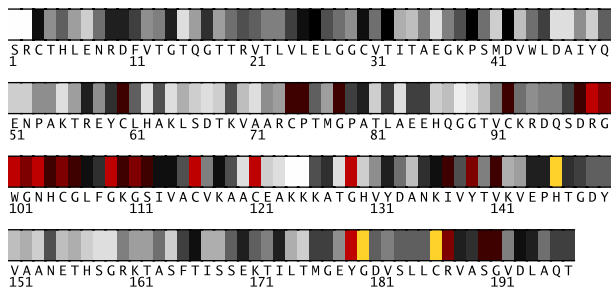
## 2.3 Residue ranking in 1k4rA

The 1k4rA sequence is shown in Figs. 1–2, with each residue colored according to its estimated importance. The full listing of residues in 1k4rA can be found in the file called 1k4rA.ranks.sorted in the attachment.

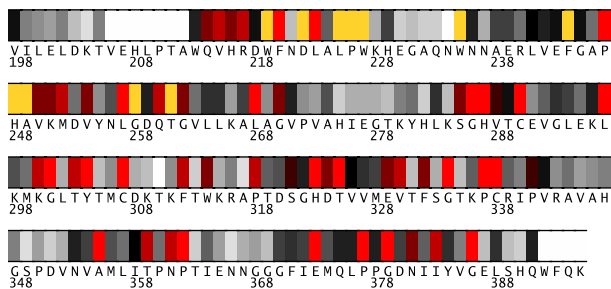
## 2.4 Top ranking residues in 1k4rA and their position on the structure

In the following we consider residues ranking among top 25% of residues in the protein. Figure 3 shows residues in 1k4rA colored by their importance: bright red and yellow indicate more conserved/important residues (see Appendix for the coloring scheme). A Pymol script for producing this figure can be found in the attachment.

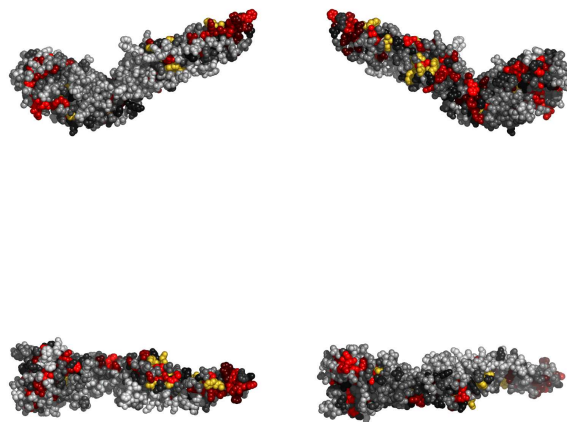
**2.4.1 Clustering of residues at 25% coverage.** Fig. 4 shows the top 25% of all residues, this time colored according to clusters they belong to. The clusters in Fig.4 are composed of the residues listed in Table 1.



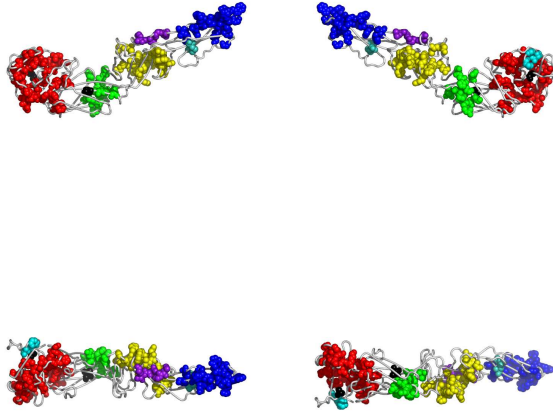
**Fig. 1.** Residues 1-197 in 1k4rA colored by their relative importance. (See Appendix, Fig.12, for the coloring scheme.)



**Fig. 2.** Residues 198-395 in 1k4rA colored by their relative importance. (See Appendix, Fig.12, for the coloring scheme.)



**Fig. 3.** Residues in 1k4rA, colored by their relative importance. Clockwise: front, back, top and bottom views.



**Fig. 4.** Residues in 1k4rA, colored according to the cluster they belong to: red, followed by blue and yellow are the largest clusters (see Appendix for the coloring scheme). Clockwise: front, back, top and bottom views. The corresponding Pymol script is attached.

Table 1.		
cluster color	size	member residues
red	33	32, 39, 42, 146, 179, 180, 297, 300 301, 303, 304, 307, 312, 314, 318 321, 323, 324, 325, 326, 329, 330 332, 334, 337, 338, 341, 359, 361 362, 373, 383, 386
blue	24	74, 75, 78, 98, 99, 100, 101, 102 103, 104, 105, 106, 109, 110, 111 112, 244, 247, 248, 249, 250, 251 252, 254
yellow	17	60, 121, 129, 214, 215, 216, 217 219, 220, 222, 223, 225, 226, 227 235, 268, 270
green	12	26, 137, 139, 186, 187, 190, 191 285, 286, 287, 288, 290
purple	5	257, 258, 260, 261, 262
azure	3	377, 379, 381
turquoise	2	92, 116

**Table 1.** Clusters of top ranking residues in 1k4rA.

**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.

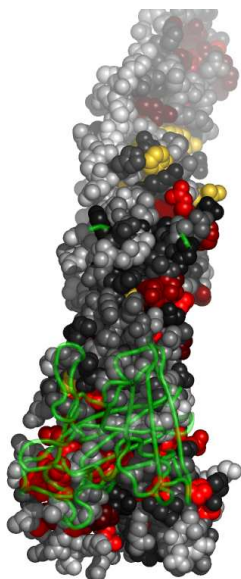
**Interface with 1k4rC4.** By analogy with 1k4rC – 1k4rC4 interface. Table 2 lists the top 25% of residues at the interface with 1k4rC4. The following table (Table 3) suggests possible disruptive replacements for these residues (see Section 3.6).

Table 2.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
303	T	T (77) S (15) . (6)	0.13	3/2	4.63

**Table 2.** The top 25% of residues in 1k4rA at the interface with 1k4rC4. (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 3.		
res	type	disruptive mutations
303	T	(KR) (FMWH) (Q) (LPI)

**Table 3.** List of disruptive mutations for the top 25% of residues in 1k4rA, that are at the interface with 1k4rC4.



**Fig. 5.** Residues in 1k4rA, at the interface with 1k4rC4, colored by their relative importance. 1k4rC4 is shown in backbone representation (See Appendix for the coloring scheme for the protein chain 1k4rA.)

Figure 5 shows residues in 1k4rA colored by their importance, at the interface with 1k4rC4.

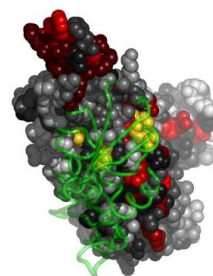
**Interface with 1k4rC.** Table 4 lists the top 25% of residues at the interface with 1k4rC. The following table (Table 5) suggests possible disruptive replacements for these residues (see Section 3.6).

Table 4.					
res	type	subst's (%)	cvg	noc/ bb	dist (Å)
227	W	W (98) S . (1)	0.04	2/0	4.54

**Table 4.** The top 25% of residues in 1k4rA at the interface with 1k4rC. (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
227	W	(K) (E) (Q) (D)

**Table 5.** List of disruptive mutations for the top 25% of residues in 1k4rA, that are at the interface with 1k4rC.



**Fig. 6.** Residues in 1k4rA, at the interface with 1k4rC, colored by their relative importance. 1k4rC is shown in backbone representation (See Appendix for the coloring scheme for the protein chain 1k4rA.)

Figure 6 shows residues in 1k4rA colored by their importance, at the interface with 1k4rC.

**Interface with 1k4rB.** By analogy with 1k4rC – 1k4rB interface. Table 6 lists the top 25% of residues at the interface with 1k4rB. 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)
258	G	G(97) . (1)AR	0.03	4/4	4.13
261	T	E(55) T(40) . (1) R(2)M	0.03	1/1	4.77
257	L	L(95) . (1)C Y(1)M	0.08	11/11	3.56
260	Q	Q(94) H(1) . (1)RK E(1)	0.12	10/0	3.38
101	W	W(96) . (3)	0.14	122/14	3.25
254	V	V(89) . (1) I(5) A(3)	0.17	22/18	3.14
329	E	E(76) K(13) . (6)D Q(1)	0.17	1/1	4.53
102	G	G(95)N . (3)	0.18	32/32	2.94
262	G	G(89) . (1) P(1) A(6)TRS	0.18	21/21	3.41
321	S	T(55) S(37) . (6)	0.20	21/21	3.30
106	G	G(95)A . (3)	0.21	4/4	4.00
98	D	D(95) . (3)A	0.22	6/0	4.61
104	H	G(57) H(37) . (3)N	0.23	17/0	3.23

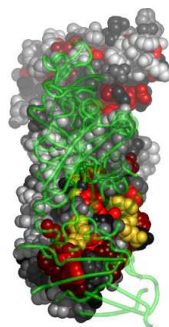
**Table 6.** The top 25% of residues in 1k4rA at the interface with 1k4rB. (Field names: res: residue number in the PDB entry; type: amino acid type; substs: substitutions seen in the alignment; with the percentage of each type in the bracket; noc/bb: number of contacts with the ligand, with the number of contacts realized through backbone atoms given in the bracket; dist: distance of closest approach to the ligand. )

Table 7.		
res	type	disruptive mutations
258	G	(E)(KD)(R)(FQMWH)
261	T	(FWH)(R)(K)(Q)
257	L	(R)(Y)(H)(K)

*continued in next column*

Table 7. continued		
res	type	disruptive mutations
260	Q	(Y)(T)(FW)(VCAG)
101	W	(KE)(TQD)(SNCG)(R)
254	V	(YR)(KE)(H)(QD)
329	E	(FW)(H)(Y)(VCAG)
102	G	(R)(E)(K)(FWH)
262	G	(KER)(H)(FWD)(Q)
321	S	(KR)(FWH)(QM)(LPI)
106	G	(KER)(HD)(Q)(FMW)
98	D	(R)(H)(FW)(Y)
104	H	(E)(MD)(TQ)(SKVLAPI)

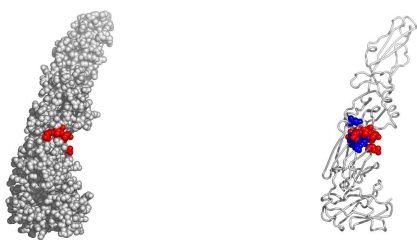
**Table 7.** List of disruptive mutations for the top 25% of residues in 1k4rA, that are at the interface with 1k4rB.



**Fig. 7.** Residues in 1k4rA, at the interface with 1k4rB, colored by their relative importance. 1k4rB is shown in backbone representation (See Appendix for the coloring scheme for the protein chain 1k4rA.)

Figure 7 shows residues in 1k4rA colored by their importance, at the interface with 1k4rB.

**2.4.3 Possible novel functional surfaces at 25% coverage.** One group of residues is conserved on the 1k4rA surface, away from (or substantially larger than) other functional sites and interfaces recognizable in PDB entry 1k4r. It is shown in Fig. 8. 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 8, while Table 9 suggests possible disruptive replacements for these residues (see Section 3.6).



**Fig. 8.** A possible active surface on the chain 1k4rA. The larger cluster it belongs to is shown in blue.

Table 8.			
res	type	substitutions(%)	cvg
286	G	G(93).(6)	0.06
287	H	H(92).(6)Y	0.06
139	Y	Y(93)F(5).(1)	0.15
187	R	E(51)R(37)K(5) Q(1).(1)N(1)H	0.16
285	S	S(71)A(12).(6) G(8)	0.17
190	S	S(90)T(3)A(3)N . (1)M	0.20
191	G	G(93)S(2)A(1) . (1)TM	0.20
26	E	E(90).(9)	0.25

**Table 8.** Residues forming surface "patch" in 1k4rA.

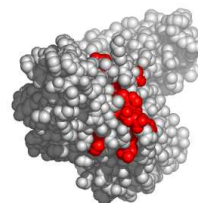
Table 9.		
res	type	disruptive mutations
286	G	(KER)(FQMWHD)(NLPI)(Y)
287	H	(E)(M)(Q)(D)
139	Y	(K)(Q)(M)(E)
187	R	(T)(Y)(VA)(CG)
285	S	(KR)(QH)(FMW)(E)
190	S	(R)(K)(H)(FW)
191	G	(R)(K)(E)(H)
26	E	(FWH)(VCAG)(YR)(T)

**Table 9.** Disruptive mutations for the surface patch in 1k4rA.

Another group of surface residues is shown in Fig.9. 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	antn
219	W	W(98).(1)	0.02	
222	D	D(98).(1)	0.02	

*continued in next column*



**Fig. 9.** Another possible active surface on the chain 1k4rA.

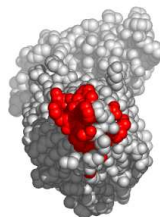
Table 10. continued				
res	type	substitutions(%)	cvg	antn
226	P	P(98).(1)	0.02	
225	L	L(95)Y(2).(1)F	0.03	
235	W	W(98)R.(1)	0.03	
258	G	G(97).(1)AR	0.03	
261	T	E(55)T(40).(1) R(2)M	0.03	
227	W	W(98)S.(1)	0.04	
223	L	L(96)I(1)V.(1)	0.08	
257	L	L(95).(1)CY(1)M	0.08	
268	L	L(93).(5)M(1)	0.09	
215	V	V(97).(1)L(1)	0.10	
121	C	C(96).(3)	0.11	S-S
260	Q	Q(94)H(1).(1)RK E(1)	0.12	
217	R	R(79)K(16)E(2) . (1)	0.13	
216	H	H(87)D(1)N(9) . (1)	0.16	
214	Q	L(52)Q(36)A(1) N(3)S(2)VI(1) . (1)M	0.18	
262	G	G(89).(1)P(1) A(6)TRS	0.18	
270	G	G(89).(6)K(2)E D(1)	0.18	
60	C	C(92).(7)	0.23	S-S
241	L	L(83)M(13)V(1) . (1)	0.25	

**Table 10.** Residues forming surface "patch" in 1k4rA.

Table 11.		
res	type	disruptive mutations
219	W	(KE) (TQD) (SNCG) (R)
222	D	(R) (FWH) (VCAG) (KY)
226	P	(YR) (TH) (SCG) (KE)
225	L	(R) (TK) (YE) (SCHG)
235	W	(E) (TD) (K) (Q)
258	G	(E) (KD) (R) (FQMWH)
261	T	(FWH) (R) (K) (Q)
227	W	(K) (E) (Q) (D)
223	L	(YR) (H) (T) (KE)
257	L	(R) (Y) (H) (K)
268	L	(Y) (R) (T) (H)
215	V	(Y) (R) (KE) (H)
121	C	(KER) (FQMWH) (NLPI) (Y)
260	Q	(Y) (T) (FW) (VCAG)
217	R	(T) (YVCAG) (SD) (FLWPI)
216	H	(TE) (M) (VQCAG) (SD)
214	Q	(Y) (H) (T) (FW)
262	G	(KER) (H) (FWD) (Q)
270	G	(FWR) (H) (K) (Y)
60	C	(KER) (FQMWH) (NLPI) (Y)
241	L	(Y) (R) (H) (T)

**Table 11.** Disruptive mutations for the surface patch in 1k4rA.

Another group of surface residues is shown in Fig.10. The residues



**Fig. 10.** Another possible active surface on the chain 1k4rA.

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	antn
249	A	A(99) .	0.00	
244	F	F(99) .	0.01	
248	H	H(98) .Y	0.01	
247	P	A(22)P(75) .T	0.09	
252	M	Q(51)M(43) .(1)	0.13	
		I(1)RWV(1)		
99	R	R(96) .(3)	0.14	
101	W	W(96) .(3)	0.14	
103	N	N(96) .(3)	0.14	
109	G	G(96) .(3)	0.14	
250	V	T(38)V(44)K(13)	0.16	
		.(1)A(1)X		
100	G	G(95)S .(3)	0.17	
254	V	V(89) .(1)I(5)	0.17	
		A(3)		
102	G	G(95)N .(3)	0.18	
105	C	C(95) .(3)X	0.19	S-S
251	K	K(80)R(15) .(1)	0.19	
		T(2)S		
74	C	C(93) .(6)	0.21	S-S
75	P	P(93) .(6)	0.21	
78	G	G(93) .(6)	0.21	
92	C	C(93) .(6)	0.21	S-S
106	G	G(95)A .(3)	0.21	
98	D	D(95) .(3)A	0.22	
104	H	G(57)H(37) .(3)N	0.23	
110	K	K(95)E .(3)	0.23	
112	S	S(84)G(12) .(3)R	0.23	

Table 12. Residues forming surface "patch" in 1k4rA.

Table 13.			
res	type	disruptive mutations	
249	A	(KYER) (QHD) (N) (FTMW)	
244	F	(KE) (TQD) (SNCG) (R)	
248	H	(E) (M) (Q) (D)	
247	P	(R) (Y) (H) (K)	
252	M	(Y) (T) (H) (S)	
99	R	(TD) (SVCLAPIG) (YE) (FMW)	
101	W	(KE) (TQD) (SNCG) (R)	
103	N	(Y) (FTWH) (SVCAG) (ER)	
109	G	(KER) (FQMWH) (NLPI) (Y)	
250	V	(Y) (R) (E) (K)	
100	G	(KR) (E) (FMWH) (Q)	
254	V	(YR) (KE) (H) (QD)	
102	G	(R) (E) (K) (FWH)	
105	C	(KER) (FWH) (MD) (Q)	
251	K	(Y) (FW) (T) (VA)	
74	C	(KER) (FQMWH) (NLPI) (Y)	
75	P	(YR) (TH) (SCG) (KE)	
78	G	(KER) (FQMWH) (NLPI) (Y)	

*continued in next column*

Table 13. continued		
res	type	disruptive mutations
92	C	(KER) (FQMWH) (NLPI) (Y)
106	G	(KER) (HD) (Q) (FMW)
98	D	(R) (H) (FW) (Y)
104	H	(E) (MD) (TQ) (SKVLAPI)
110	K	(Y) (FW) (T) (VCAG)
112	S	(K) (FMWR) (QH) (ELPI)

Table 13. Disruptive mutations for the surface patch in 1k4rA.

Another group of surface residues is shown in Fig. 11. The right panel shows (in blue) the rest of the larger cluster this surface belongs to.



Fig. 11. Another possible active surface on the chain 1k4rA. The larger cluster it belongs to is shown in blue.

The residues belonging to this surface "patch" are listed in Table 14, while Table 15 suggests possible disruptive replacements for these residues (see Section 3.6).

Table 14.				
res	type	substitutions(%)	cvg	antn
146	H	H(98) .(1)	0.02	
180	G	G(98) .(1)	0.02	
301	G	G(93) .(6)	0.06	
307	C	C(92) .(6)R	0.06	S-S
323	H	H(93) .(6)	0.06	
325	T	T(93) .(6)	0.06	
337	P	P(93) .(6)	0.06	
338	C	C(93) .(6)	0.06	S-S
362	P	P(93) .(6)	0.06	
304	Y	EY(92) .(6)	0.07	
386	G	G(92) .(6)V	0.07	
297	L	L(87) .(6)V(5)	0.08	
334	G	G(89) .(6)K(2)N	0.09	
318	P	P(81)M(12) .(6)	0.10	
312	F	F(89) .(6)LM(1)	0.11	
		Y(1)		
359	T	T(80)S(12) .(6)	0.11	
179	Y	Y(93)F(4) .(1)HD	0.12	

*continued in next column*



res	type	substitutions(%)	cvg	antn
300	K	K(88).(6)R(3) V(1)	0.12	
361	N	N(81)T(11).(6)I	0.12	
303	T	T(77)S(15).(6)	0.13	
330	V	L(38)V(53).(6)I	0.13	
324	D	G(55)D(36).(6) Q(1)	0.15	
332	F	Y(67)F(22).(6) V(3)	0.15	
314	W	F(41)MW(37) I(12).(6)L	0.17	
329	E	E(76)K(13).(6)D Q(1)	0.17	
321	S	T(55)S(37).(6)	0.20	
341	P	P(89)Q(1).(6)AT SL	0.22	
39	P	P(90).(9)	0.25	
358	I	V(41)I(51).(6)	0.25	

Table 14. Residues forming surface "patch" in 1k4rA.

res	type	disruptive mutations
146	H	(E)(TQMD)(SNVCLAPIG)(K)
180	G	(KER)(FQMWH)(NLPI)(Y)
301	G	(KER)(FQMWH)(NLPI)(Y)
307	C	(E)(D)(KM)(FW)
323	H	(E)(TQMD)(SNVCLAPIG)(K)
325	T	(KR)(FQMW)(NLPI)(E)
337	P	(YR)(TH)(SCG)(KE)
338	C	(KER)(FQMWH)(NLPI)(Y)
362	P	(YR)(TH)(SCG)(KE)
304	Y	(K)(M)(VQA)(LPIR)
386	G	(KER)(HD)(Q)(FMW)
297	L	(Y)(R)(H)(T)
334	G	(FEW)(HR)(YD)(KM)
318	P	(Y)(R)(T)(H)
312	F	(K)(E)(T)(Q)
359	T	(KR)(FMWH)(Q)(LPI)
179	Y	(K)(Q)(M)(NE)
300	K	(Y)(T)(FW)(SCG)
361	N	(Y)(H)(FW)(R)
303	T	(KR)(FMWH)(Q)(LPI)
330	V	(YR)(KE)(H)(QD)
324	D	(R)(FWH)(Y)(VA)
332	F	(K)(E)(Q)(D)
314	W	(KE)(T)(D)(QR)
329	E	(FW)(H)(Y)(VCAG)
321	S	(KR)(FWH)(QM)(LPI)
341	P	(R)(Y)(H)(TK)

*continued in next column*

res	type	disruptive mutations
39	P	(YR)(TH)(SCG)(KE)
358	I	(Y)(R)(H)(T)

Table 15. Disruptive mutations for the surface patch in 1k4rA.

### 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 [EK RQM], OH-group possession [SDETY], and NH2 group possession [NQRK]. The suggestions are listed according to how different they appear to be from the original amino acid, and they are grouped in round brackets if they appear equally disruptive. From left to right, each bracketed group of amino acid types resembles more strongly the original (i.e. is, presumably, less disruptive) These suggestions are tentative - they might prove disruptive to the fold rather than to the interaction. Many researcher will choose, however, the straightforward alanine mutations, especially in the beginning stages of their investigation.

## 4 APPENDIX

### 4.1 File formats

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

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

### 4.2 Color schemes used

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

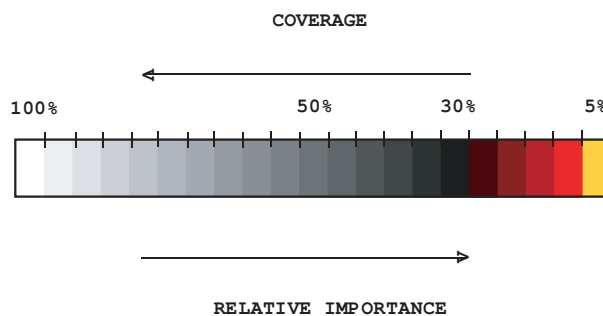


Fig. 12. Coloring scheme used to color residues by their relative importance.

The colors used to distinguish the residues by the estimated evolutionary pressure they experience can be seen in Fig. 12.

### 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:

- 1k4rA.complex.pdb - coordinates of 1k4rA with all of its interacting partners
- 1k4rA.etvx - ET viewer input file for 1k4rA
- 1k4rA.cluster\_report.summary - Cluster report summary for 1k4rA
- 1k4rA.ranks - Ranks file in sequence order for 1k4rA
- 1k4rA.clusters - Cluster descriptions for 1k4rA
- 1k4rA.msf - the multiple sequence alignment used for the chain 1k4rA
- 1k4rA.descr - description of sequences used in 1k4rA msf
- 1k4rA.ranks\_sorted - full listing of residues and their ranking for 1k4rA
- 1k4rA.1k4rC4.if.pml - Pymol script for Figure 5
- 1k4rA.cbvcg - used by other 1k4rA – related pymol scripts
- 1k4rA.1k4rC.if.pml - Pymol script for Figure 6
- 1k4rA.1k4rB.if.pml - Pymol script for Figure 7