

Shetti-Motif

Haitham Sobhy

Supplementary file – SI-1: This file, contains table S1 and figures S1-S9.

Supplementary file – SI-2: A compressed file (zip format) contains NCBI-BLASTp output files. The x-rich motifs differ from LCRs and may not be masked. For example, BLASTp was used to mask LCRs in four proteins encoded by *Acanthamoeba polyphaga* mimivirus (UniProt ID: 5757), as the following:

1. Collagen-like protein (L669: Q5UNS9), 1937 aa, is S-rich protein (450 Ser, ~23% of length), the LCR regions are masked by BLASTp.
2. Hypothetical protein (R661: E3VZK9), 218 aa, is enriched with Cys (~11% of length), but the region masked does not contain Cys residue.
3. Hypothetical protein (L725: Q5UNX5), 224 aa, is enriched with Cys and Lys (~17% of the length; 21 K and 18 C), which are not masked.
4. Putative TPR repeat-containing protein (R856: Q5UQQ7), 342 aa, KLN share ~27% of R856 length, (32, 30 and 30 Asn, Lys and Leu, respectively, but only 5 Cys and 1 Pro), but not masked.

Supplementary Table S2 (in a separate excel spreadsheet).

The tool is supplemented with built-in patterns/motifs obtained from experimentally validated literature. The table consists of the motif, description/function and references to PubMed IDs. See also, Fig. S3.

Supplementary Table S3 (in a separate excel spreadsheet).

The motif-containing proteins were counted and normalized (percentage %) to the total number of the proteins encoded by a virus. These normalized values were then used for further statistical analysis, see Fig. 2 in the manuscript.

Supplementary file – S1

Table S1. Shetti-Motif accepts PROSITE pattern syntax. Some reports use Latin symbols to refer to degenerate residues, e.g. θ , Ψ , Φ , etc., which refers to residues share a similar physicochemical property(ies), such as hydrophobicity, polarity, etc. For simplicity, Shetti-Motif accepts non-Latin symbols as listed in the table. These symbols are explicitly translated to the corresponding residues. Otherwise, users may insert alternative / degenerate residues between “[” and “]”.

Noteworthy, the physical properties of bases might be not the same within literatures. As an example, some reports refer to (F, I, L or V) residues as hydrophobic one. In Shetti-Motif, A, C, F, G, V, L, I, P, W, M or Y are considered as hydrophobic residues. This may lead to reporting additional motifs (e.g. contain A or Y, instead of those contain only F, I, L or V). To overcome this issue, users may write directly “[FILV]”, instead of the representative symbol.

Table S1.

Residues	Representative Motif	Symbol	Physical properties
Arg-Gly-Asp	RGD		
Pro-Pro-any amino acid-Tyr	PPxY		
P, and (E or D) residues \approx PE or PD	P[ED]		
P, and any residue but R	P{R}		
T is repeated 2, 2-5 or >3 times	PAST(2), PAST(2-5), or PAST(3,)		
E or D	[ED]	- “hyphen”	Negative - Acidic
H, K or R	[HKR]	+	Positive - Basic
S or T	[ST]	=	Alcohol
C or M	[CM]	*	Sulfur containing
I, V or L	[IVL]	?	Aliphatic
A, G or S	[AGS]	&	Tiny
F, H, W or Y	[FHWY]	@	Aromatic
D, E, H, K, N, Q, R, S or T	[DEHKNQRST]	%	Polar / hydrophilic
A, C, F, G, V, L, I, P, W, M or Y	[ACFGVLIPWMY]	!	Non-polar / hydrophobic
C, W, N, Q, S, T, Y, K, R, H, D or E	[CWNQSTYKRHDE]	#	H-bond
Any amino acid	Any amino acid	x	Any amino acid

Fig. S1. Result of sequences enriched with residues (x-rich motifs) appears in “Result - statistics” tab, which include the names of the proteins, protein length, number of motifs in each protein, and coverage of the residues to the full length of the protein (%). The results are tabulated and can be transferred to text editor or Excel software.

Protein name	Length	Number of motifs	Coverage to protein length	Statistics
AY386371.1_prot_AA07364.1.1	11	1	[protein=11] [protein_id=AA07364.1] [location=complement(804..1808)]	334 2 C = 3.29; P = 1.5; V = 12.57;
AY386371.1_prot_AA07365.1.2	21	1	[protein=21] [protein_id=AA07365.1] [location=complement(1938..2963)]	341 1 C = 2.64; P = 4.4; V = 7.04;
AY386371.1_prot_AA07367.1.4	51	2	[protein=51] [protein_id=AA07367.1] [location=complement(3762..4232)]	156 2 C = 7.05; P = 0; V = 9.62;
AY386371.1_prot_AA07368.1.5	61	1	[protein=61] [protein_id=AA07368.1] [location=complement(4274..4732)]	152 1 C = 3.95; P = 2.63; V = 8.55;
AY386371.1_prot_AA07372.1.9	131	1	[protein=131] [protein_id=AA07372.1] [location=complement(8614..9474)]	286 1 C = 2.8; P = 3.15; V = 7.34;
AY386371.1_prot_AA07376.1.13	191	1	[protein=191] [protein_id=AA07376.1] [location=complement(11200..12774)]	524 1 C = 2.48; P = 2.67; V = 10.5;
AY386371.1_prot_AA07381.1.18	241	1	[protein=241] [protein_id=AA07381.1] [location=complement(14799..15443)]	214 1 C = 3.27; P = 3.27; V = 5.61;
AY386371.1_prot_AA07383.1.20	261	1	[protein=261] [protein_id=AA07383.1] [location=complement(16786..18711)]	641 1 C = 2.03; P = 3.12; V = 5.77;
AY386371.1_prot_AA07384.1.21	271	1	[protein=271] [protein_id=AA07384.1] [location=complement(18736..19845)]	369 1 C = 2.44; P = 2.98; V = 6.23;
AY386371.1_prot_AA07388.1.25	311	2	[protein=311] [protein_id=AA07388.1] [location=21335..21649]	104 2 C = 4.81; P = 6.73; V = 7.69;
AY386371.1_prot_AA07390.1.27	331	3	[protein=331] [protein_id=AA07390.1] [location=complement(23072..25120)]	682 3 C = 2.64; P = 2.79; V = 6.6;
AY386371.1_prot_AA07391.1.28	341	1	[protein=341] [protein_id=AA07391.1] [location=complement(25146..25700)]	184 1 C = 3.26; P = 4.89; V = 5.98;
AY386371.1_prot_AA07393.1.30	361	1	[protein=361] [protein_id=AA07393.1] [location=26444..27472]	342 1 C = 4.09; P = 2.63; V = 6.14;
AY386371.1_prot_AA07396.1.33	391	6	[protein=391] [protein_id=AA07396.1] [location=complement(30022..33042)]	1006 6 C = 2.49; P = 3.08; V = 6.76;
AY386371.1_prot_AA07397.1.34	401	1	[protein=401] [protein_id=AA07397.1] [location=33075..33559]	94 1 C = 3.19; P = 4.26; V = 4.26;
AY386371.1_prot_AA07403.1.40	471	1	[protein=471] [protein_id=AA07403.1] [location=complement(36245..37402)]	385 1 C = 2.08; P = 3.64; V = 8.57;
AY386371.1_prot_AA07408.1.45	521	1	[protein=521] [protein_id=AA07408.1] [location=42822..43490]	222 1 C = 3.6; P = 1.35; V = 5.41;
AY386371.1_prot_AA07423.1.60	671	1	[protein=671] [protein_id=AA07423.1] [location=52968..53471]	167 1 C = 2.4; P = 3.59; V = 10.18;
AY386371.1_prot_AA07424.1.61	681	1	[protein=681] [protein_id=AA07424.1] [location=53549..54550]	333 1 C = 0.9; P = 4.8; V = 4.8;
AY386371.1_prot_AA07426.1.63	701	1	[protein=701] [protein_id=AA07426.1] [location=complement(54999..55412)]	137 1 C = 5.84; P = 4.38; V = 7.3;
AY386371.1_prot_AA07427.1.64	711	1	[protein=711] [protein_id=AA07427.1] [location=55509..59366]	1285 1 C = 1.25; P = 3.5; V = 6.07;
AY386371.1_prot_AA07428.1.65	721	2	[protein=721] [protein_id=AA07428.1] [location=complement(59363..59872)]	169 2 C = 1.78; P = 3.55; V = 8.88;
AY386371.1_prot_AA07431.1.68	751	2	[protein=751] [protein_id=AA07431.1] [location=complement(61421..63814)]	797 2 C = 1; P = 2.76; V = 6.4;
AY386371.1_prot_AA07435.1.72	791	3	[protein=791] [protein_id=AA07435.1] [location=65982..68504]	840 3 C = 0.95; P = 2.98; V = 7.02;
AY386371.1_prot_AA07438.1.75	821	1	[protein=821] [protein_id=AA07438.1] [location=69660..70319]	219 1 C = 2.74; P = 5.94; V = 7.31;
AY386371.1_prot_AA07439.1.76	831	1	[protein=831] [protein_id=AA07439.1] [location=70393..72753]	786 1 C = 2.8; P = 4.2; V = 6.11;
AY386371.1_prot_AA07440.1.77	841	1	[protein=841] [protein_id=AA07440.1] [location=72750..74657]	635 1 C = 0.31; P = 3.15; V = 5.98;
AY386371.1_prot_AA07441.1.78	851	1	[protein=851] [protein_id=AA07441.1] [location=74690..75172]	160 1 C = 3.12; P = 2.5; V = 8.75;
AY386371.1_prot_AA07442.1.79	861	2	[protein=861] [protein_id=AA07442.1] [location=75194..75859]	221 2 C = 2.71; P = 1.81; V = 5.88;
AY386371.1_prot_AA07444.1.81	881	1	[protein=881] [protein_id=AA07444.1] [location=complement(76586..78481)]	631 1 C = 1.58; P = 2.69; V = 6.66;
AY386371.1_prot_AA07446.1.83	901	5	[protein=901] [protein_id=AA07446.1] [location=complement(79398..81059)]	553 5 C = 1.81; P = 4.16; V = 9.95;
AY386371.1_prot_AA07447.1.84	911	1	[protein=911] [protein_id=AA07447.1] [location=complement(81076..81531)]	151 1 C = 3.31; P = 3.97; V = 9.93;
AY386371.1_prot_AA07450.1.87	941	1	[protein=941] [protein_id=AA07450.1] [location=complement(82467..84440)]	657 1 C = 1.98; P = 3.96; V = 5.18;
AY386371.1_prot_AA07454.1.91	981	1	[protein=981] [protein_id=AA07454.1] [location=complement(86619..88760)]	713 1 C = 1.4; P = 2.52; V = 6.59;
AY386371.1_prot_AA07457.1.94	1011	1	[protein=1011] [protein_id=AA07457.1] [location=complement(89933..92641)]	902 1 C = 1.33; P = 3.44; V = 4.99;
AY386371.1_prot_AA07458.1.95	1021	1	[protein=1021] [protein_id=AA07458.1] [location=92656..93600]	314 1 C = 1.27; P = 3.5; V = 5.73;

Fig. S2. For sequences enriched with residues (x-rich motifs), the protein sequence headers, the sequences of the motifs and statistics are presented in tab “Result - sequences”. The results can be transferred to text editor software.

```

ShettiMotif
File Load patterns Help
Options Result - table Result - statistics Result - sequences
>|Id|AY386371.1_prot_AAR07364.1_1 [protein=1L] [protein_id=AAR07364.1] [location=complement(804..1808)]
;Motif_1 Motif length: 15; Residue coverage: 40% Residue(s) count: C = 1; P = 1; V = 4;
;CEVVAVGKMKRTPV
;Motif_2 Motif length: 15; Residue coverage: 40% Residue(s) count: C = 2; P = 0; V = 4;
;DVCVGDHLTVVKCFK
>|Id|AY386371.1_prot_AAR07365.1_2 [protein=2L] [protein_id=AAR07365.1] [location=complement(1938..2963)]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 0; P = 2; V = 3;
;IPSNVMPVTVTGEE
>|Id|AY386371.1_prot_AAR07367.1_4 [protein=5L] [protein_id=AAR07367.1] [location=complement(3762..4232)]
;Motif_1 Motif length: 25; Residue coverage: 32% Residue(s) count: C = 2; P = 0; V = 6;
;DDTFDFVFLTVYSMLVTVCLCV
;Motif_2 Motif length: 15; Residue coverage: 40% Residue(s) count: C = 2; P = 0; V = 4;
;MLVTVCVFLAL
>|Id|AY386371.1_prot_AAR07368.1_5 [protein=6L] [protein_id=AAR07368.1] [location=complement(4274..4732)]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 3; P = 1; V = 1;
;MDFCPGCLVDCLNR
>|Id|AY386371.1_prot_AAR07372.1_9 [protein=13L] [protein_id=AAR07372.1] [location=complement(8614..9474)]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 0; P = 2; V = 3;
;KRYVPMVPMFVLGHS
>|Id|AY386371.1_prot_AAR07376.1_13 [protein=19L] [protein_id=AAR07376.1] [location=complement(11200..12774)]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 0; P = 3; V = 2;
;PDMYPRKFGVNF
>|Id|AY386371.1_prot_AAR07381.1_18 [protein=24L] [protein_id=AAR07381.1] [location=complement(14799..15443)]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 1; P = 0; V = 4;
;IICFVDFVIVV
>|Id|AY386371.1_prot_AAR07383.1_20 [protein=26L] [protein_id=AAR07383.1] [location=complement(16786..18711)]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 0; P = 3; V = 2;
;VGITKYVEPSLPDK
>|Id|AY386371.1_prot_AAR07384.1_21 [protein=27L] [protein_id=AAR07384.1] [location=complement(18736..19845)]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 2; P = 2; V = 1;
;ACCLPVSTKYHYNF
>|Id|AY386371.1_prot_AAR07388.1_25 [protein=31R] [protein_id=AAR07388.1] [location=21335..21649]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 2; P = 0; V = 3;
;KVCDIRKVDVCDVSKA
;Motif_2 Motif length: 15; Residue coverage: 40% Residue(s) count: C = 2; P = 2; V = 2;
;SCLVKVEKPTPTCDR
>|Id|AY386371.1_prot_AAR07390.1_27 [protein=33L] [protein_id=AAR07390.1] [location=complement(23072..25120)]
;Motif_1 Motif length: 15; Residue coverage: 33% Residue(s) count: C = 2; P = 0; V = 3;
;KEVFTTEICVRVCD
  
```

Fig. S3. The built-in patterns are listed with check-boxes. The user may choose the pattern(s) of interest, then click on “Select these patterns” button. The pattern(s) are populated in the search text area. In “Options” tab, select or deselect all by ticking on “Select All”. The patterns were acquired from PROSITE databases (<http://prosite.expasy.org/>, <ftp://ftp.expasy.org/databases/prosite/>, last accessed August, 2015), or from literatures, see table S2 for references and PubMed IDs.

Select ALL

- 6-cysteine motif, degradation of chitin and chitotriose ::: Cx(13,20)Cx(5,6)Cx(9,19)Cx(10,14)Cx(4,14)C
- Adenovirus-2/5 nuclear localization signal (NLS) ::: KRAR
- Adenovirus-C - putative heparan sulfate-binding site/motif, important for post-internalization steps of virus infection ::: KKTK
- Adenovirus-D fiber flexibility motif ::: KLGXGLxF[DN]
- Adenovirus-D fiber flexibility motif ::: KxGGLxF[DN]
- Adhesion protein motif ::: KxGFFKR
- Adhesion protein motif ::: SVSVGMKPSRPR
- Binding of Vif to human APOBEC3G, ElonginB and C, and cullin 5, suppression of APOBEC3 proteins ::: PPLP
- Binding of Vif to human APOBEC3G, ElonginB and C, and cullin5, suppression of APOBEC3 proteins ::: SLxYLA
- Binding to ESCRT, Paramyxoviruses budding ::: IPxV
- binding to integrins (Adenovirus and foot-and-mouth disease virus) ::: RGD
- Binding to integrins, and viral attachment to cellular receptors ::: DLxxL
- Binding to Rb (LxCxE motif) ::: [L]x[Cx][DE]
- Budded virions production and nucleocapsid assembly ::: Cx5CxnHx6C (C2HC zinc finger)
- Caveolin scaffolding domain (CSD) ::: !x!xxxx!
- Caveolin scaffolding domain (CSD) ::: !x!xxxx!xx!
- Caveolin scaffolding domain (CSD) ::: !xxxx!xx!
- Cell surface of Plasmodium falciparum, Thrombospondin type-1 (TSP1) repeat profile, adhesive, Immunodominant surface antigen on the sporozoite (the infe
- Clathrin-binding motifs, clathrin-box ::: L!x! [DE]
- Clathrin-binding motifs, clathrin-box ::: L[L] [DEN] [LF] [DE]
- Clathrin-binding motifs, W-box ::: PW!xxW
- cleavage of NS1 from the NS1-NS2A region of flavivirus ::: [LM]VxSxVxVx
- cleavage site for Influenza A ::: [QE] [ST] RGLF
- cleavage site of HA of H9N2 avian influenza KSS[RG]LF motif ::: KSS[RG]LF
- cleavage site of HA of H9N2 avian influenza vaccine strain ::: RSS[RG]LF
- could allow SUMO to bind to substrate ::: [VI]x[VI][VI]
- dynein binding motifs, protein transport, viruses transport inside the cell (Adenovirus, ASF virus, Papilloma virus, Rabies virus, Mokola Virus and Ebolavirus) :
- dynein binding motifs, protein transport, viruses transport inside the cell (Adenovirus, ASF virus, Papilloma virus, Rabies virus, Mokola Virus and Ebolavirus) :
- enhance virion-release, anti-tetherin activity ::: DSGxxS
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: [ENT] [HNPS] [ILV] Y [ADEG]
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EDLY
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EHIYD
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EHLYA
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: ENIYE
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EPIYA
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EPIYG
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EPIYD
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EPIYD
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EPLYA
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EPLYA
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: EPVYA
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: ESIYE
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: NPLYE
- EPIYA-related motif, interactions between bacterial effectors and host SH2 domain-containing proteins ::: TPLYA
- FIL-rich domain, agnoprotein function, productive viral infection ::: L[F] [VI] F [VIL] [LE] [LF] LLxF

Fig. S4. After clicking “Search pattern” button, the protein containing the motifs, and the motifs and their positions on the sequences are listed in a table, in “Result - table” tab. The results are tabulated and can be moved to text editor or Excel software.

Note: for multiple pattern, the results are printed to text file.

	Sequence header	Motifs (Position: Motif)	Number of motifs in the protein
1	kl AY386371.1_prot_AAR07375.1_12 [protein=17L] [protein_id=AAR07375.1] [location=complement(10635..11066)]	'105: RGD'; '133: RGD';	2
2	kl AY386371.1_prot_AAR07418.1_55 [protein=62L] [protein_id=AAR07418.1] [location=complement(49816..50763)]	'257: RGD';	1
3	kl AY386371.1_prot_AAR07429.1_66 [protein=73R] [protein_id=AAR07429.1] [location=59887..60456]	'126: RGD';	1
4	kl AY386371.1_prot_AAR07439.1_76 [protein=83R] [protein_id=AAR07439.1] [location=70393..72753]	'344: RGD';	1
5	kl AY386371.1_prot_AAR07494.1_131 [protein=142R] [protein_id=AAR07494.1] [location=123296..124225]	'223: RGD';	1
6	kl X69198.1_prot_CAA48949.1_8 [gene=D8L] [protein=D8L] [protein_id=CAA48949.1] [location=complement(8602..9054)]	'24: RGD';	1
7	kl X69198.1_prot_CAA48967.1_26 [gene=C6L] [protein=C6L] [protein_id=CAA48967.1] [location=complement(21874..22317)]	'136: RGD';	1
8	kl X69198.1_prot_CAA48980.1_39 [gene=C19L] [protein=C19L] [protein_id=CAA48980.1] [location=complement(33321..33806)]	'123: RGD';	1
9	kl X69198.1_prot_CAA49028.1_87 [gene=I4L] [protein=I4L] [protein_id=CAA49028.1] [location=complement(80101..82488)]	'137: RGD';	1
10	kl X69198.1_prot_CAA49036.1_95 [gene=F5R] [protein=F5R] [protein_id=CAA49036.1] [location=89132..91489]	'344: RGD';	1
11	kl X69198.1_prot_CAA49110.1_169 [gene=B1R] [protein=B1R] [protein_id=CAA49110.1] [location=152700..153602]	'211: RGD';	1
12	kl AY243312.1_prot_AAO89320.1_41 [gene=VACWR041] [protein=dUTPase] [protein_id=AAO89320.1] [location=complement(31038..31481)]	'136: RGD';	1
13	kl AY243312.1_prot_AAO89333.1_54 [gene=VACWR054] [protein=unknown] [protein_id=AAO89333.1] [location=complement(42460..42903)]	'109: RGD';	1
14	kl AY243312.1_prot_AAO89369.1_90 [gene=VACWR090] [protein=unknown] [protein_id=AAO89369.1] [location=complement(78062..79114)]	'19: RGD';	1
15	kl AY243312.1_prot_AAO89381.1_102 [gene=VACWR102] [protein=RAP94] [protein_id=AAO89381.1] [location=complement(89298..91685)]	'137: RGD';	1
16	kl AY243312.1_prot_AAO89389.1_110 [gene=VACWR110] [protein=NTPase interacts with A20R] [protein_id=AAO89389.1]	'344: RGD';	1
17	kl AY243312.1_prot_AAO89462.1_183 [gene=VACWR183] [protein=ser/thr kinase] [protein_id=AAO89462.1] [location=163878..164780]	'211: RGD';	1
18	kl AY243312.1_prot_AAO89469.1_190 [gene=VACWR190] [protein=soluble interferon-gamma receptor-like protein] [protein_id=AAO89469.1]	'137: RGD';	1
19	kl AY243312.1_prot_AAO89473.1_194 [gene=VACWR194] [protein=ser/thr protein kinase-like protein] [protein_id=AAO89473.1] [location=17...	'199: RGD';	1
20	kl AY243312.1_prot_AAO89485.1_206 [gene=VACWR206] [protein=unknown] [protein_id=AAO89485.1] [location=183734..184306]	'149: RGD';	1
21	kl AY386264.1_prot_AAR98232.1_4 [protein=ORF007 dUTPase] [protein_id=AAR98232.1] [location=complement(5191..5700)]	'145: RGD';	1

Fig. S5. The percentage (in %) of the protein harbouring these motifs to the entire protein dataset (proteome) appears in a text area, in “Result - statistics” tab. The results are tabulated and can be moved to text editor or Excel software.

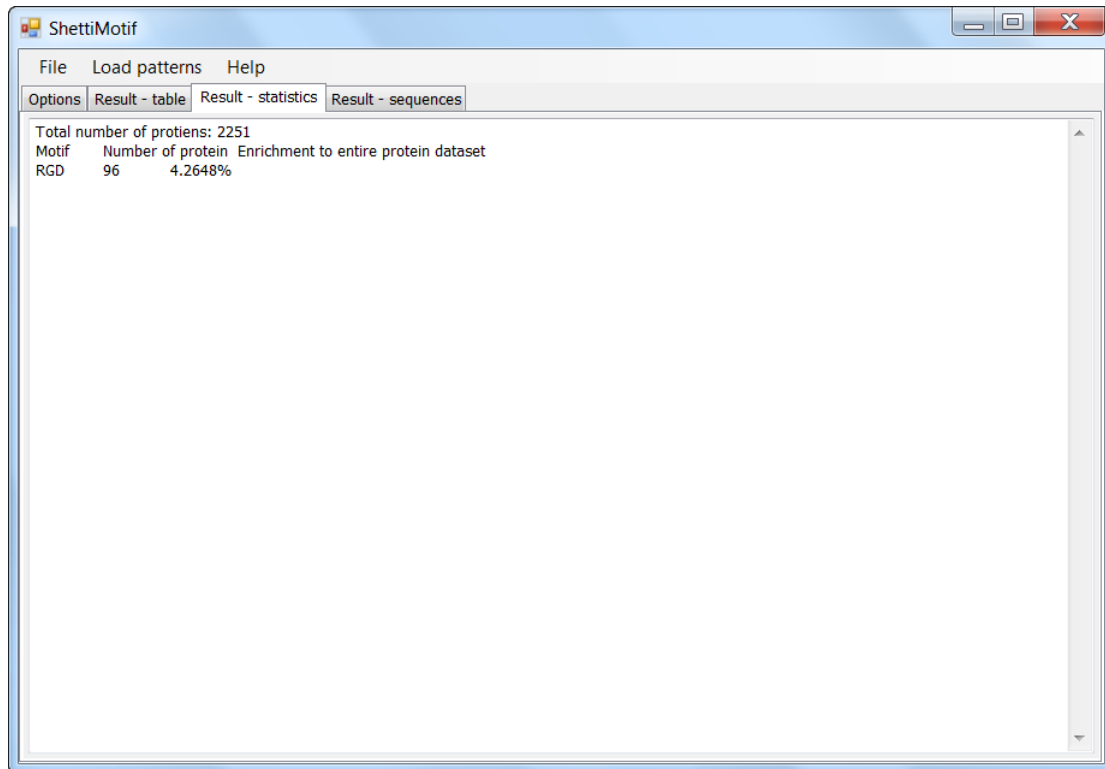


Fig. S6. The sequences of the motif-containing proteins appear in a text area, in “Result - sequences” tab, which can be copied to text editor software.



Fig. S7. The PROSITE flat file can be converted to a table format. In “Result - table” tab, the results are tabulated and can be moved to text editor or Excel software.

	Identification	Accession number	Description	Pattern
1	ASN_GLYCOSYLATION; PATTERN	PS00001	N-glycosylation site	N-(P)-[ST]-[P]
2	CAMP_PHOSPHO_SITE; PATTERN	PS00004	cAMP- and cGMP-dependent protein kinase phosphorylation site	[RK](2)-x-[ST]
3	PKC_PHOSPHO_SITE; PATTERN	PS00005	Protein kinase C phosphorylation site	[ST]-x-[RK]
4	CK2_PHOSPHO_SITE; PATTERN	PS00006	Casein kinase II phosphorylation site	[ST]-x(2)-[DE]
5	TYR_PHOSPHO_SITE; PATTERN	PS00007	Tyrosine kinase phosphorylation site	[RK]-x(2,3)-[DE]-x(2,3)-Y
6	MYRISTYL; PATTERN	PS00008	N-myristoylation site	G-(EDRKHPFFW)-x(2)-[STAGCN]-[P]
7	AMIDATION; PATTERN	PS00009	Amidation site	x-G-[RK]-[RK]
8	ASX_HYDROXYL; PATTERN	PS00010	Aspartic acid and asparagine hydroxylation site	C-x-[DN]-x(4)-[FY]-x-C-x-C
9	GLA_1; PATTERN	PS00011	Vitamin K-dependent carboxylation domain	E-x(2)-[ERK]-E-x-C-x(6)-[EDR]-x(10,11)-[FYA]-[YW]
10	PHOSPHOPANTHETHEINE; PATTERN	PS00012	Phosphopantetheine attachment site	[DEQGSTALMKRH]-[LIVFYSTAC]-[GNQ]-[LIVFYAG]-[DNEKHS]-S-[LIVMST]-[PCF]
11	ER_TARGET; PATTERN	PS00014	Endoplasmic reticulum targeting sequence	[KRHQSA]-[DENQ]-E-L>
12	RGD; PATTERN	PS00016	Cell attachment sequence	R-G-D
13	ATP_GTP_A; PATTERN	PS00017	ATP/GTP-binding site motif A (P-loop)	[AG]-x(4)-G-K-[ST]
14	EF_HAND_1; PATTERN	PS00018	EF-hand calcium-binding domain	D-[W]-[DNS]-[LIVFYW]-[DENSTG]-[DNQGHK]-[GP]-[LIVMC]-[DENQSTAGC]-x(2)
15	ACTININ_1; PATTERN	PS00019	Actinin-type actin-binding domain signature 1	[EQ]-[LIVY]-x-[ATV]-[FY]-[LDAM]-[T]-W-[PG]-N
16	ACTININ_2; PATTERN	PS00020	Actinin-type actin-binding domain signature 2	[LIVM]-x-[SGNL]-[LIVMN]-[DAGENRS]-[SAGPNVT]-x-[DNEAG]-[LIVM]-x-[DEAGQ]
17	KRINGLE_1; PATTERN	PS00021	Kringle domain signature	[FY]-C-[RH]-[NS]-x(7,8)-[IWI]-C
18	EGF_1; PATTERN	PS00022	EGF-like domain signature 1	C-x-C-x(2)-[V]-x(2)-G-[C]-x-C
19	FN2_1; PATTERN	PS00023	Fibronectin type-II collagen-binding domain signature	C-x(2)-P-F-x-[FYWIV]-x(7)-C-x(8,10)-W-C-x(4)-[DNSR]-[FYW]-x(3,5)-[FYW]-x-[F]
20	HEMOPEXIN; PATTERN	PS00024	Hemopexin domain signature	[LIFAT]-[IL]-x(2)-W-x(2,3)-[PE]-x-[VF]-[LIVMFY]-[DENQS]-[STA]-[AV]-[LIVMFY]
21	P_TREFOIL_1; PATTERN	PS00025	P-type Trefoil domain signature	[RRH]-x(2)-C-x-[FVPSTV]-x(3,4)-[ST]-x(3)-C-x(4)-C-C-[FYWH]
22	CHIT_BIND_1_1; PATTERN	PS00026	Chitin recognition or binding domain signature	C-x(4,5)-C-C-S-x(2)-G-x-C-G-x(3,4)-[FYW]-C
23	HOMEBOX_1; PATTERN	PS00027	'Homeobox' domain signature	[LIVMFYGG]-[ASLVR]-x(2)-[LIVMSTACN]-x-[LIVM]-[Y]-x(2)-[L]-[LIV]-[RKNQSTAIY]
24	ZINC_FINGER_C2H2_1; PATTERN	PS00028	Zinc finger C2H2 type domain signature	C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H
25	LEUCINE_ZIPPER; PATTERN	PS00029	Leucine zipper pattern	L-x(6)-L-x(6)-L-x(6)-L
26	NUCLEAR_REC_DBD_1; PATTERN	PS00031	Nuclear hormones receptors DNA-binding region signature	C-x(2)-C-x(1,2)-[DENAVSPHKQT]-x(5,6)-[HNY]-[FY]-x(4)-C-x(2)-C-x(2)-F(2)-x-R
27	ANTENNAPEIDIA; PATTERN	PS00032	'Homeobox' antennapedia-type protein signature	[LIVMFE]-[FY]-P-W-M-[RQTA]
28	ENGRAILED; PATTERN	PS00033	'Homeobox' engrailed-type protein signature	L-M-A-[EQ]-G-L-Y-N

Fig. S8. The UniProt flat file is converted to a table format. The table includes UniProt ID, name of organism, protein name, taxonomy, sequence, PROSITE pattern IDs, InterPro and Pfam IDs. The results are tabulated, and blank cells denotes absence of those IDs.

	Identification	Accession number	Description (recommended name)	Gene name	Organism species	Organelle	Length	MW	PROSITE	InterPro	Pfam
21	DHCR7_MIMIV	QSUQ14	Probable 7-dehydroc...	MIMI_R807	Acanthamoeba polyphaga mimivirus (APMV)		447	51692		IPR001171;	PF01222;
22	DIPP_MIMIV	QSUQW2	Putative	MIMI_L375	Acanthamoeba polyphaga mimivirus (APMV)		360	43662	PS51462;PS50158;	IPR000086;IPR0157...	PF00293;
23	DNLI_MIMIV	QSUP20	DNA ligase	MIMI_R303	Acanthamoeba polyphaga mimivirus (APMV)		636	72089	PS50172;	IPR001357;IPR0...	PF00533;PF016...
24	DNMK_MIMIV	QSUQ70	Putative	MIMI_R512	Acanthamoeba polyphaga mimivirus (APMV)		193	22106		IPR027417;	
25	DPOLX_MIMIV	Q7T6Y4	Probable DNA polyme...	MIMI_L318	Acanthamoeba polyphaga mimivirus (APMV)		354	40612	PS00522;	IPR002054;IPR0...	PF14792;PF147...
26	DPOL_MIMIV	QSUQR0	DNA polymerase	MIMI_R322	Acanthamoeba polyphaga mimivirus (APMV)		1740	201945	PS50818;PS50819;	IPR006172;IPR0061...	PF00136;PF03104
27	DRTS_MIMIV	QSUQG3	Bifunctional dihydrofo...	MIMI_R497	Acanthamoeba polyphaga mimivirus (APMV)		563	65063	PS51330;PS00091;	IPR024072;IPR0...	PF00186;PF00303
28	END4_MIMIV	QSUPY4	Putative endonuclease 4	MIMI_R296	Acanthamoeba polyphaga mimivirus (APMV)		312	34857	PS00730;PS51432;	IPR001719;IPR0182...	PF01261;
29	FPG_MIMIV	QSUQ00	Probable formamidop...	MIMI_L315	Acanthamoeba polyphaga mimivirus (APMV)		287	33463	PS51068;	IPR015886;IPR0...	PF01149;PF06831
30	GFAT_MIMIV	Q7T6X6	Probable glutamine--fructose-6-phosphoaminotransferase	MIMI_L619	Acanthamoeba polyphaga mimivirus (APMV)		606	68570	PS51278;PS51464;	IPR017932;IPR0058...	PF01380;
31	GLNA_MIMIV	QSUR44	Putative glutamine sy...	MIMI_R565	Acanthamoeba polyphaga mimivirus (APMV)		353	40079	PS00180;PS00181;	IPR008147;IPR0...	PF00120;PF03951
32	GLRX_MIMIV	QSUQ14	Probable glutaredoxin	MIMI_R195	Acanthamoeba polyphaga mimivirus (APMV)		106	12084	PS00195;PS51354;	IPR011767;IPR0021...	PF00462;
33	GNA1_MIMIV	QSUP29	Probable glucosamine...	MIMI_L316	Acanthamoeba polyphaga mimivirus (APMV)		148	17017	PS51186;	IPR016181;IPR0...	PF00583;
34	HSP70_MIMIV	QSUQ49	Heat shock 70 kDa protein homolog	MIMI_L393	Acanthamoeba polyphaga mimivirus (APMV)		634	70514	PS00297;PS00329;PS01...	IPR018181;IPR0290...	PF00012;
35	HSP7L_MIMIV	QSUPU0	Heat shock protein 7...	MIMI_L254	Acanthamoeba polyphaga mimivirus (APMV)		941	107204	PS01036;	IPR018181;IPR0...	PF00012;
36	IF4EH_MIMIV	QSUQG4	Eukaryotic translation initiation factor 4E homolog	MIMI_L496	Acanthamoeba polyphaga mimivirus (APMV)		272	30949		IPR023398;IPR0010...	PF01652;
37	KITH_MIMIV	QSUP25	Thymidine kinase	MIMI_L258	Acanthamoeba polyphaga mimivirus (APMV)		225	25909	PS00603;	IPR027417;IPR0...	PF00265;
38	LONH_MIMIV	QSUPT0	Lon protease homolog	MIMI_L251	Acanthamoeba polyphaga mimivirus (APMV)		1023	116827		IPR003593;IPR0039...	PF00004;PF05362
39	MCAR_MIMIV	QSUPV8	Mitochondrial carrier-4...	MIMI_L276	Acanthamoeba polyphaga mimivirus (APMV)		237	27319	PS50920;	IPR018108;IPR0...	PF00153;
40	MCE_MIMIV	QSUQX1	Probable mRNA-capping enzyme	MIMI_R382	Acanthamoeba polyphaga mimivirus (APMV)		1170	136508	PS51562;	IPR023577;IPR0013...	PF01331;PF02940
41	MGMT_MIMIV	QSUINU9	Probable methylated-...	MIMI_R693	Acanthamoeba polyphaga mimivirus (APMV)		149	16851	PS00374;	IPR001497;IPR0...	PF01035;
42	MUTSL_MIMIV	QSUQU6	Putative DNA mismatch repair protein mutS homolog L359	MIMI_L359	Acanthamoeba polyphaga mimivirus (APMV)		1124	130330	PS00486;	IPR007695;IPR0004...	PF01624;PF05192

Fig. S9. Combining (mapping) both UniProt and PROSITE files results in a new table, which includes PROSITE ID, accession numbers of the protein containing the pattern, description of the pattern and the pattern consensus. The results are tabulated, and blank cells denotes absence of consensus pattern in PROSITE flat file.

	PROSITE	Accession numbers	Description	Pattern
1	PS51278	QSUQE1;Q7T6X6;	Glutamine amidotransferase type 2 domain profile	
2	PS00122	QSUR02;	Carboxylesterases type-B serine active site	F[GR]Gx(4)[LIVM]x[LIV]xGxS[STAG]G
3	PS00941	QSUR02;	Carboxylesterases type-B signature 2	[EDA][D]G[CL][YTF][LIVT][DNS][LIV][LIVFYW]x[PQR]
4	PS00086	QSUQI3;	Cytochrome P450 cysteine heme-iron ligand signature	[FW][SGNH]x[GD][F][RKH#T][P]C[LIVMFAP][GAD]
5	PS00191	QSUR80;	Cytochrome b5 family, heme-binding domain signature	[FY][LIVMK]{}{}[Q]HP[GAA]G
6	PS50255	QSUR80;	Cytochrome b5 family, heme-binding domain profile	
7	PS51462	QSUQW2;	Nudix hydrolase domain profile	
8	PS50158	QSUQW2;	Zinc finger CCHC-type profile	
9	PS50172	QSUP20;	BRCT domain profile	
10	PS00522	Q7T6Y4;	DNA polymerase family X signature	G[SG][LFY]xR[GE]x(3)[SGCL]x[D][LIVM]D[LIVMFY](3)x(2)[SAP]
11	PS50818	QSUQR0;	Intein C-terminal splicing motif profile	
12	PS50819	QSUQR0;	Intein DOD-type homing endonuclease domain profile	
13	PS51330	QSUQG3;	Dihydrofolate reductase (DHFR) domain profile	
14	PS00091	QSUQG3;	Thymidylate synthase active site	Rx(2)[LIVMT]x(2,3)[FWY][QNYDI]x(8,13)[LVESI]xPC[HAVMLC]x(3)[QMTLHD][FYWL]x(0,1)[LV]
15	PS00730	QSUPY4;	AP endonucleases family 2 signature 2	[GSARY][LIVMF][CT][LIVMFY]DTCH
16	PS51432	QSUPY4;	AP endonucleases family 2 profile	
17	PS51068	QSUQ00;	Formamidopyrimidine-DNA glycosylase catalytic domain profile	
18	PS51464	Q7T6X6;	SIS domain profile	
19	PS00180	QSUR44;	Glutamine synthetase signature 1	[FYWL]DGSSx(6,8)[DENQSTAK][SA][DE]x(2)[LIVMFY]
20	PS00181	QSUR44;	Glutamine synthetase putative ATP-binding region signature	KP[LIVMFYA]x(3,5)[NPAT][GA][GSTAN][GA]xHx(3)S
21	PS00195	QSUQ14;	Glutaredoxin active site	[LIVMD][FYSA]x(4)C[PV][FYWF]Cx(2)[TAV]x(2,3)[LV]
22	PS51354	QSUQ14;	Glutaredoxin domain profile	
23	PS51186	QSUP29;QSUR52;	Gcn5-related N-acetyltransferase (GNAT) domain profile	
24	PS00297	QSUQ49;	Heat shock hsp70 proteins family signature 1	[IV]DLGT[ST]x[SC]
25	PS00329	QSUQ49;	Heat shock hsp70 proteins family signature 2	[LIVMF][LIVMFY][DN][LIVMF]G[GS][AST]x(3)[ST][LIVM][LIVMFC]
26	PS01036	QSUQ49;QSUPU0;	Heat shock hsp70 proteins family signature 3	[LIVMY]x[LIVMF]xGGx[ST](L,S)[LIVM]Px[LIVM]x[DEQKRSTA]
27	PS00603	QSUP25;	Thymidine kinase cellular-type signature	[GA]x(1,2)[DE]Yx[STAPV]xC[NKR]x[CH][LIVMFYWH]
28	PS50920	QSUPV8;	Solute carrier (Solcar) repeat profile	
29	PS51557	QSUQY1;	RNA (guanine-N7)-methyltransferase (EC 2.1.1.55) domain profile	