

## Introduction to Bioinformatics

Iosif Vaisman

Email: [ivaisman@gmu.edu](mailto:ivaisman@gmu.edu)

## Sequence patterns

```
KKFAQSTNLKSHILT
KQFShSAQLRAHIST
GKFSDSNQLKSHMLV
KDISSESRLRTHMFK
KRFSHSGSYSSHIS
KRFSHSGSFSSHMTS
KTLSDRLEYQQHMLK
```

## Sequence patterns

```
KKFAQSTNLKSHILT
KQFShSAQLRAHIST
GKFSDSNQLKSHMLV
KDISSESRLRTHMFK
KRFSHSGSYSSHIS
KRFSHSGSFSSHMTS
KTLSDRLEYQQHMLK
```

## Sequence patterns

```
KKFAQSTNLKSHILT
KQFShSAQLRAHIST
GKFSDSNQLKSHMLV
KDISSESRLRTHMFK
KRFSHSGSYSSHIS
KRFSHSGSFSSHMTS
KTLSDRLEYQQHMLK
```

## Regular Expressions

Operation	Regular Expression	Example
Concatenation	DLIV	DLIV
Alternation	D[LIV]K	DLK
Replication	DL(2,5)K	DLLK

## Regular Expressions

Patterns described in a standard way are known as *regular expressions*

<b>x</b>	ANY		
<b>[ ]</b>	OR	[ILV]	I or L or V
<b>{ }</b>	NOT	{DE}	not D or E
<b>( )</b>	repetitions	x(2,3)	x-x or x-x-x
<b>-</b>	separator		
<b>&lt;</b>	N-terminal		
<b>&gt;</b>	C-terminal		
<b>.</b>	END		

## Regular Expressions

[AC]-x-V-x(4)-{ED}.

[Ala or Cys]-any-Val-any-any-any-any-{any but Glu or Asp}

```

... LKHVAYVFQALIYWIK...
... AVEMAGVKYLQVQHGS...
... LYTGAIVTNNDGPYMA...
... KEYKCKVEKELTDICN...
    
```

## PROSITE Database

The screenshot shows the PROSITE database interface. It includes a search bar, navigation links, and a list of motifs. The motifs are listed in a table with columns for motif name, accession number, and sequence. The motifs are: DMA\_VIBCH|Q08318 (85) SCTQWMPFF 77, HEMK\_MYCLE|P45832 (181) DLFVAQPTL 100, MT57\_ECOLI|P25240 (111) DGALGNPPF 13, MTC1\_CHVN1|Q01511 (172) NFFVLDPPY 8, MTC1\_COREQ|P42828 (71) QLSFSCPPF 49, MTH2\_HAEHA|P00473 (32) KIAFFDPQY 52, MTH3\_HAEIN|P43871 (23) HAIISDIPY 73, MTM1\_MICAM|P50190 (306) AAVLTNPPF 14, MTM2\_MORBO|P23192 (25) QLAVIDPPY 10, MTM3\_MYCSP|P43641 (37) QVIYADPPW 13, MTR1\_RHOSH|P14751 (60) DLIICDPY 8.

## PROSITE Database

Current version contains 1079 documentation entries that describe 1459 different patterns, rules and profiles/matrices

[ST]-x(2)-[DE]

Casein kinase II phosphorylation site

[AG]-x(4)-G-K-[ST]

ATP/GTP-binding site motif A (P-loop)

Y-x-[NQH]-K-[DE]-[IVA]-F-[LM]-R-[ED]

Heat shock hsp90 proteins family signature

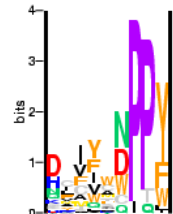
<http://www.expasy.ch/prosite>

## Blocks Database

Blocks are multiply aligned ungapped segments corresponding to the most highly conserved regions of proteins

N-6 Adenine-specific DNA methylases proteins  
width=9 seqs=78

DMA_VIBCH Q08318	(85)	SCTQWMPFF	77
HEMK_MYCLE P45832	(181)	DLFVAQPTL	100
MT57_ECOLI P25240	(111)	DGALGNPPF	13
MTC1_CHVN1 Q01511	(172)	NFFVLDPPY	8
MTC1_COREQ P42828	(71)	QLSFSCPPF	49
MTH2_HAEHA P00473	(32)	KIAFFDPQY	52
MTH3_HAEIN P43871	(23)	HAIISDIPY	73
MTM1_MICAM P50190	(306)	AAVLTNPPF	14
MTM2_MORBO P23192	(25)	QLAVIDPPY	10
MTM3_MYCSP P43641	(37)	QVIYADPPW	13
MTR1_RHOSH P14751	(60)	DLIICDPY	8



<http://www.blocks.fhcr.org/>

## Pfam Database

Pfam is a large collection of multiple sequence alignments and hidden Markov models covering many common protein domains

Zinc finger, C2H2 type

```

TYY1 HUMAN/383-407 YVCPF.DGCN...KKFAQSTNLKSHILT...H
ZG52 XENLA/61-83 YTCT...QCN...KQFSSHAQLRAHIST...H
KRUP DROME/306-328 YTCE...ICD...KQFSSDNLKSHMLV...H
YKQ8 CAEEL/78-102 YKCT...VCR...KDISSSESRLTHMFRQ.HH
DEFI CHICK/268-292 YECP...NCK...KRFSSHSYSSHISSK.KC
ZFH1 DROME/389-413 FGCD...NGG...KRFSSHSYSSHISSK.KC
YLS7 CAEEL/42-65 YLCY...YCG...KTLSDRLEYQQHMLK...VH
ZFA MOUSE/542-564 FKCD...ICL...LTFSDTEVQQHALV...H
BASO HUMAN/719-742 FQCD...ICK...KTFKACSVKIHKN...MH
HUNB DROME/297-319 FQCD...KCS...YTCVNKSMNLNHRKS...H
SFP1 YEAST/598-623 FKCPV.IGCE...KTYKNQNLKYHRLH...GH
ZG29 XENLA/62-84 FVCT...VCG...KTYKYHGLNTHLHS...H
    
```

<http://pfam.wustl.edu/>

## Other Motif Databases

**PRINTS** : a compendium of protein fingerprints.

A fingerprint is a group of conserved motifs used to characterise a protein family

<http://bioinf.man.ac.uk/dbbrowser/PRINTS/>

**DOMO** : a protein domain database

<http://www.infobiogen.fr/~gracy/domo/home.htm>

**ProDom** : a protein domain database

<http://protein.toulouse.inra.fr/prodom.html>

