# Patterns

Tore Samuelsson 2010

---

# Patterns -
## comparison to method of aligning evolutionary related sequences

Dynamic programming algorithm. Designed to deal with closely as well as distantly related sequences, taking into account gaps (indel mutations) and substitution matrix information



Applications in:

- Sequence assembly
- Classification
- Prediction of function
- Comparative genomics
- Phylogeny / Evolutionary history

---

# Patterns

## Regular expression
## "Regexp"
## "Regex"

Patterns / Regular expression matching

```
GRTKLPKLMKKWREKNRLYKMKWRAGGALKALK
```
Is there a match to **KWR** in this string?

Patterns / Regular expression matching

GRTKLPKLMKKWREKNRLYKMKWRAGGALKALK

**KWR**

---

Patterns / Regular expression matching

GRTKLPKLMKKWREKNRLYKMKWRAGGALKALK

**K.R**

---

Patterns / Regular expression matching

GRTKLPKLMKKWREKNRLYKMKWRAGGALKALK

**[KR].{0,2}[KR]**

## Regular expressions and NFA / DFA

Regular expressions are related to theory of *automata (machines)*:
Regular expressions can be translated into :

*Non-deterministic Finite Automaton* (NFA) &

*Deterministic Finite Automaton* (DFA)

(and this is what happens in Perl and other
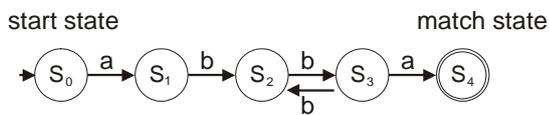utilities that make use of regexs)

---

## Regular expressions and finite state automata (state machines)

Regexp :      a(bb)+a

start state                                        match state

$S_0$ —a→ $S_1$ —b→ $S_2$ —b→ $S_3$ —a→ (($S_4$))
                        ←b—

---

start state                                   match state

$S_0$ —a→ $S_1$ —b→ $S_2$ —b→ $S_3$ —a→ (($S_4$))
                       ←b—

Regexp     a(bb)+a      Testing the string "abbbba"

STEP 0 → $S_0$ —a→ $S_1$ —b→ $S_2$ —b→ $S_3$ —a→ $S_4$      abbbba
                                              ←b—                       ^

STEP 1 → $S_0$ —a→ $S_1$ —b→ $S_2$ —b→ $S_3$ —a→ $S_4$      abbbba
                                              ←b—                       ^

STEP 2 → $S_0$ —a→ $S_1$ —b→ $S_2$ —b→ $S_3$ —a→ $S_4$      abbbba
                                              ←b—                       ^

STEP 3 → (S₀) —a→ (S₁) —b→ (S₂) ⇄b (S₃) —a→ (S₄)   abbbba
                                                         ^

STEP 4 → (S₀) —a→ (S₁) —b→ (S₂) ⇄b (S₃) —a→ (S₄)   abbbba
                                                         ^

STEP 5 → (S₀) —a→ (S₁) —b→ (S₂) ⇄b (S₃) —a→ (S₄)   abbbba
                                                          ^

STEP 6 → (S₀) —a→ (S₁) —b→ (S₂) ⇄b (S₃) —a→ (S₄)   abbbba
                                                          ^

---

## Deterministic and non-deterministic state machines

This machine (automaton) is **deterministic** because in any state each possible input character leads to at most *one* new state.

→ (S₀) —a→ (S₁) —b→ (S₂) →b← (S₃) —a→ ((S₄))
                              b

This machine is equivalent to the previous one, but is **non-deterministic** because in the state S2 it has multiple choices for b :

→ (S₀) —a→ (S₁) —b→ (S₂) —b→ (S₃) —a→ ((S₄))
                   b←

---

## Applications of patterns/regexps

When are patterns useful in bioinformatics?

* Finding patterns in DNA and protein sequences
  - Restriction enzyme sites (DNA)
  - Transcription factor binding sites (DNA)
  - PROSITE, using patterns to infer biological function
       from protein sequences
* Finding patterns to extract information from text files of
     bioinformatics applications,
         such as GenBank reports and BLAST output.

**Recognition sites of restriction enzymes**

```
EcoRI  -GAATTC-
       -CTTAAG-

BamHI  -GGATCC-
       -CCTAGG-

XhoII  -RGATCY-    R = A or G
       -YCTAGR-    Y = C or T

PpuMI  -RGGWCCY-   W = A or T
       -YCCWGGR-

FokI   -GGATGNNNNNNNNNNNNNNNNN-
       -CCTACNNNNNNNNNNNNNNNNN-
```
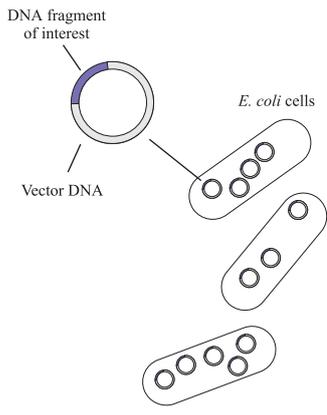
**Cloning**

DNA fragment
of interest

*E. coli* cells

Vector DNA

**Restriction enzymes are used in cloning**

DNA fragment of interest to be cloned

```
G AATTC——— G AATTC
CTTAA G ——— CTTAA G
```

Plasmid vector DNA

## Slide 1

### Cutting in silico

http://rna.lundberg.gu.se/cutter2/

**w e b**
**c u t**
**t e r**

Please enter a title for this sequence:
Untitled sequence

Paste the DNA sequence into the box below

Please select the type of analysis you would like
- Linear sequence analysis
- Circular sequence analysis
- Find sites which may be introduced by silent mutagenesis

Please indicate how you would like the restriction sites displayed
- Map of restriction sites
- Table of sites, sorted alphabetically by enzyme name
- Table of sites, sorted sequentially by base pair number

## Slide 2

### Web-cutter output

```
                                        Eco31I   PshBI
              BseRI                      BspHI    AsnI
tcaaggctcctctcttgattcaggctgactcacagtgttggtctcatgacaaattaatttatatttcacatgagg base pairs
agttccgaggagagaactaagtccgactgagtgtcacaaccagagtactgtttaattaaatataaagtgtactcc 1 to 75
                                        RcaI     VspI
                                        BsaI     AseI



              Eco57I                                          NspI
gacaggatgatttcttgaaaaaccttcagaaatacttagccatcatctgccttattctcactggcaagacatgca base pairs
ctgtcctactaaagaactttttggaagtctttatgaatcggtagtagacggaataagagtgaccgttctgtacgt 76 to 150
                                                              MslI
```

## Slide 3

### Transcription factor databases contain information about DNA sequences recognized by transcription factors

**Regulation of transcription**

DNA

Repressor

Activator

Initiation factors

RNA polymerase

Promoter

Transcription start site

Protein recognition elements of DNA

Transcription factor binding sites - example:

Estrogen response element

5' **AGGTCA**NNN**TGACCT**

## Transcription factor databases contain information about DNA sequences recognized by transcription factors

```
Name                         Sequence                            Trans. Factor
-----------------------      -----------------------------       ---------------
UAS(G)-pMH100            0    CGGAGTACTGTCCTCCG              0 :  GAL4
TFIIIC-Xls-5S.1         0    TGGATGGGAG                     0 :  TFIIIC
HSE_CS_inverted_repeat  0    CTNGAANNTTCNAG                 0 :  HSTF
ZDNA_CS                 0    GCGTGTGCA                      0 :  unknown
GCN4-his3-189           0    ATGACTCAT                      0 :  GCN4
Ad5_EIA_element_I       0    AGGAAGTGAAA                    0 :  unknown
Ad5_EIA_element_II      0    GGGCGTAACCGAGTAAGATTTGGCCATTTTC 0 : unknown
BPV-E2_CS1              0    ACCGNNNNCGGT                   0 :  (BPV-E2)
Alb_DEI                 0    GATTTTGTAATGG                  0 :  C/EBP
Alb_DEII                0    TTTTTGGCAAAGAT                 0 :  CTF/NF-1
Alb_DEIII               0    GCAAGGGATTTAGTT                0 :  unknown
Alb_PEI                 0    TGGTTAATGATCTACAGT             0 :  APF/HNF1
BPV-E2_CS2              0    ACCNNNNNNNGGT                  0 :  BPV-E2
CACA                    0    CACACACACA                     0 :  unknown
dc_box                  0    TSATTTGCAT                     0 :  unknown
GALV_E                  0    AGAAATAGATGAGTCAACAG           0 :  unknown
GCRE                    0    TGACTC                         0 :  GCN4
Pit-1_CS1               0    WTATYCAT                       0 :  Pit-1
```

# Protein patterns and the PROSITE database

## Proteins that bind the nucleotides ATP or GTP share a short sequence motif

### Entry in PROSITE for the ATP/GTP binding site motif

```
ID   ATP_GTP_A; PATTERN.
AC   PS00017;
DT   APR-1990 (CREATED); APR-1990 (DATA UPDATE); NOV-1990 (INFO
UPDATE).
DE   ATP/GTP-binding site motif A (P-loop).
PA   [AG]-x(4)-G-K-[ST].
CC   /TAXO-RANGE=ABEPV;
3D   1EFM; 1ETU; 1Q21; 2Q21; 4Q21; 5Q21; 6Q21;
DO   PDOC00017;
```



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

---

### More PROSITE: a category of zinc finger proteins



---

### PROSITE entry of C2H2 type of zinc finger protein

```
ID   ZINC_FINGER_C2H2; PATTERN.
AC   PS00028;
DT   APR-1990 (CREATED); JUN-1994 (DATA UPDATE); NOV-1997 (INFO UPDATE).
DE   Zinc finger, C2H2 type, domain.
PA   C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H.
NR   /RELEASE=35,69113;
```

## Recognition of DNA by zinc-finger binding protein



Zinc finger binding motifs are extremely common.
Among 23000 human genes, 845 have C2H2 zinc finger domains !!

---

## More PROSITE: Disulfide bonds in EGF domain

- First identified in epidermial growth factor, but present in many other proteins
- Function unclear, but found in extracellular part of membrane proteins
- Includes six cysteins involved in disulfide bonds

```
       2                4       5                      6
       +----------------+       +------------------+
       |                |       |                  |
C-x(0,48)-C-x(3,12)-C-x(1,70)-C-x(1,6)-C-x2-G-x(0,21)-G-x2-C
|                        |
+----------------+
1                3
```



---

## PROSITE



http://www.expasy.org/prosite/ (Swiss Institute of Bioinformatics)

**Release 20.67, of 05-Oct-2010 (contains**
**1308 *patterns* and 909 *profiles/matrices*).**

A - x(4,7) - G -x (5,6) - D

"Non-probabilistic"

| A | 10 | 0  | 0  | 10 | 5 | 4 |
|---|----|----|----|----|---|---|
| K | 0  | 0  | 0  | 0  | 0 | 0 |
| G | 0  | 0  | 0  | 0  | 1 | 4 |
| L | 0  | 10 | 10 | 0  | 4 | 2 |
| .. etc. | | | | | | |

"Probabilistic"

## PROSITE syntax

| | | Perl |
|---|---|---|
| D | matches D | |
| [AS] | matches A or S | |
| x | any symbol (amino acid) | **.** |
| x(3) or x3 | any three symbols (amino acids) | **.{3}** |
| {AT} | matches any symbol *except* A and T | [^AT] |
| [AS]2 | Matches two positions with either A or S (i.e. AA, AS , SS, or SA) | [AS]{2} |
| x(3,7) | a sequence of symbols between 3 and 7 in length (for instance GAT, GCRE, PPLKM, GTTREC or PPPPPPP) | **.{3,7}** |

## Variation in classification accuracy

ZINC_FINGER_C2H2_1, PS00028; **Zinc finger C2H2 type domain signature** (PATTERN)

Consensus pattern:  **C - x(2,4) - C - x(3) - [LIVMFYWC] - x(8) - H - x(3,5) - H**
*The 2 C's and the 2 H's are zinc ligands*

Sequences known to belong to this class detected by the pattern:  ALL

Other sequence(s) detected in Swiss-Prot:  42.

ATP_GTP_A, PS00017; **ATP/GTP-binding site motif A (P-loop)** (PATTERN with a high probability of occurrence!)

Consensus pattern:  **[AG] - x(4) - G - K - [ST]**

Sequences known to belong to this class detected by the pattern:  a majority

Other sequence(s) detected in Swiss-Prot:  in addition to the proteins listed above, the 'A' motif is also found in a number of other proteins. Most of these proteins probably bind a nucleotide, but others are definitively not ATP- or GTP-binding (as for example chymotrypsin, or human ferritin light chain).

ASN_GLYCOSYLATION, PS00001; **N-glycosylation site** (PATTERN with a high probability of occurrence!)

Consensus pattern:  **N - {P} - [ST] - {P}**
*N is the glycosylation site*

• Scan Swiss-Prot/TrEMBL entries against PS00001
• view ligand binding statistics

---



* True positive - member of the protein family that is matched by the pattern

* False positive - a protein which is *not* a member of the protein family is matched by the pattern

* True negative - a protein which is *not* a member of the protein family is *not* matched by the pattern

* False negative - member of the protein family is *not* matched by pattern

## Classification accuracy

$$Specificity \quad : \quad \frac{TP}{(TP+FP)}$$

$$Sensitivity \quad : \quad \frac{TP}{(TP+FN)}$$

- Ex: The flavodoxin family (PS00201)

  True positives:     51

  False positives:    4

  False negatives:   9

$$\frac{51}{(51+4)} \approx 0.93$$

$$\frac{51}{(51+9)} \approx 0.85$$

---

## Important application of databases like PROSITE

*Biological information may be extracted from a protein sequence.*

A 'new' protein sequence has been identified using
bioinformatics methods (like in a genome project).
A scan of PROSITE using this sequence (regular
expression matching) can give important
clues as to the biological function of the protein.

---

## Common applications of PROSITE  (ScanProsite/Motifscan)
### * Search a sequence for all PROSITE patterns/profiles
### * Search database with a pattern/profile

## More applications of patterns ...

**Finding minor introds**

**Finding minor introns**



A class of U12 (minor) introns are highly conserved in sequence:

RT|**AT**CCTYT................TTCCTTTR......**AC**|Y
exon | 5'        intron              3'| exon

[AG]TATCCT[CT]T.{30,1500}TTCCTT[AG]......AC[TC]

---

## More applications of patterns ...

**Finding the iron responsive element**

IRE binding
protein

CAGUG[CU].{10,200}AUG

A G
C   U
    G

IRE        ferritin mRNA

5'          AUG

**... and finding ORFs**
`AUG(...){0,200}?(UAA|UAG|UGA)`

match as few characters as possible
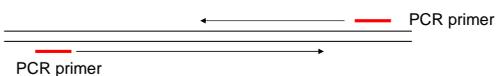
---

## More applications of patterns ...

**When BLAST will fail to identify regions of identity in a genome**

For instance: Find occurrences of 'AGCTGCAAAAA'

To find such short regions of identity BLAST is often difficult to use.
There are however applications when you want to find such
matches, like when searching for matches to oligonucleotides used in PCR or
probes used in microarrays. Or you may want to search for matches to very short
peptide sequences potentially encoded by the genome.

PCR primer

PCR primer

## More applications of patterns ...

**Repeats ...**

like in the Huntingtin protein

---CAGCAGCAGCAGCAGCAGCAGCAGCAGCAGCAGCAGCAG---

```
(...)\1{8,20}
```

stored
in memory

recall what is
in memory

---

## Regular expressions are used in ...

* Patterns as used in PROSITE

* programming and scripting languages, for instance *perl*

* operating system tools , for instance *grep*

* text editors

* EMBOSS program *fuzzpro*

---

## UNIX utility grep / egrep

```
egrep "[AG].{4}GK[TS]" /dbs/swissprot
```

will print all lines containing matching the pattern
corresponding to the ATP/GTP binding site motif in PROSITE

/dbs/swissprot is the Swissprot protein sequence database
in FASTA format

## Sample session with 'fuzzpro'

```
% fuzzpro
Protein pattern search
Input sequence(s): tsw:*
Search pattern: [FY]-[LIV]-G-[DE]-E-A-Q-x-[RKQ](2)-G
Number of mismatches [0]:
Output report [100k_rat.fuzzpro]:
```

---

## Perl

Example code that uses a regular expression:

variable

```
$str = 'GACACAGGGATCGGGGATC';
```

binding operator  regexp

```
if ( $str =~ /[AG][CT]CG/ )  {
```

regexp delimiter

Is there a match of [AG][CT]CG in the variable $str?

```
print "Found match!";

}
```

regexp will match here

GACACAGGG<u>ATCG</u>GGGATC

---

## The substitution and transliteration operators

*Transcribing DNA into RNA.*

```
$dna = 'GCAATGG';
print "The DNA sequence is $dna\n";
$rna = $dna;
$rna =~ s/T/U/g;
print "and the RNA sequence is $rna\n";
```

global modifier;
do the substitution operation
for all matches

substitution operator

$$s/T/U/g;$$

REPLACEMENT; text to replace PATTERN

PATTERN; regexp to be replaced by REPLACEMENT

## The substitution and transliteration operators

*Count the bases in a DNA sequence using tr*

```
$dna = 'GCAATGNGATTACTTCG';

$basecount = ($dna =~ tr/ACGT//;)
$nonbase = length($dna)-$basecount;
print "There are $basecount As,Cs,Ts,and Gs \n";
print "There is/are $nonbase other symbol(s) \n";
```

**tr/ACGT//;**
operation does not change the string!

---

## Metacharacters

| | |
|---|---|
| . | match any character |
| ^ | match beginning of string |
| $ | match end of string |
| ? | optional match |

Quantifiers

| | |
|---|---|
| * | Any number of characters, including zero |
| + | One or more characters |
| {m,n} | minimum m , maximum n characters |

---

## Reading a file in PERL and checking for regular expresssions

Program will read the file with the name 'myfile' and that has a sequence in EMBL format. The program will print to the screen all lines starting with 'FT', i.e lines with feature table information.

```
open IN, 'myfile';

while (<IN>) {

if (/^FT/) {print ;}

}
close IN;
```

## Example of how regexps may be used to parse the output from a BLAST search

'bl2seq' is a program to blast two sequences against each other.
This is an output using two sequences seq1 and seq2:

```
Query= seq1
        (31 letters)
>seq2
        Length = 29
 Score = 42.1 bits (21), Expect = 1e-10
 Identities = 27/29 (93%)
 Strand = Plus / Plus

Query: 3  acgacgtacacgactagtcaggcggagct 31
          |||||||  |||||||||||| ||||||||
Sbjct: 1  acgacgttcacgactagtcacgcggagct 29

.... (and some more text ) ......
```

---

## Example of how regexps may be used to parse the output from a BLAST search

Let's say you have many output files like this and you want to make a script to present the names of the two sequences as well as the Expect value listed. Here's an example of code: ( we assume that the blast output is in a file called 'bl2seq.out')

```perl
open IN, 'bl2seq.out';
while (<IN>) {  # we read one line at a time, each
            # line is stored in the default variable $_)
chomp; # remove the end of line character in $_
if (/Query= (.*)/) {print "$1 "; }
if (/^>(.*)/) {print "$1 ";}
if (/Expect = (.*)/) {print "$1\n";}
}
close IN;
```

.* means any number (including zero) of any characters; the regexp algorithm is *greedy* so it will try to find the longest substring that matches

---

## Regular expression match basics

1. Quantifiers * + ? {m,n} are *greedy*

Perl example:

```perl
$str = "GGAAGG";
$str =~ s/(G.*G)//;
```

$str is now empty

2. The match that begins earliest (leftmost) wins

Perl example:

```perl
$str = "GGAAGG";
$str =~ s/GG/XX/;
print $str;
```

$str is now XXAAGG