# BioinfoRmatika

October 6th 2018

# Content

- Sequences – string

  Library BioStrings

- Central dogma of molecular biology

- Transcription, translation, open reading frame (ORF)

_____

- Mutations

- Distance/difference between two DNA sequences

_____

- Genomic ranges

_____


- Break

_____
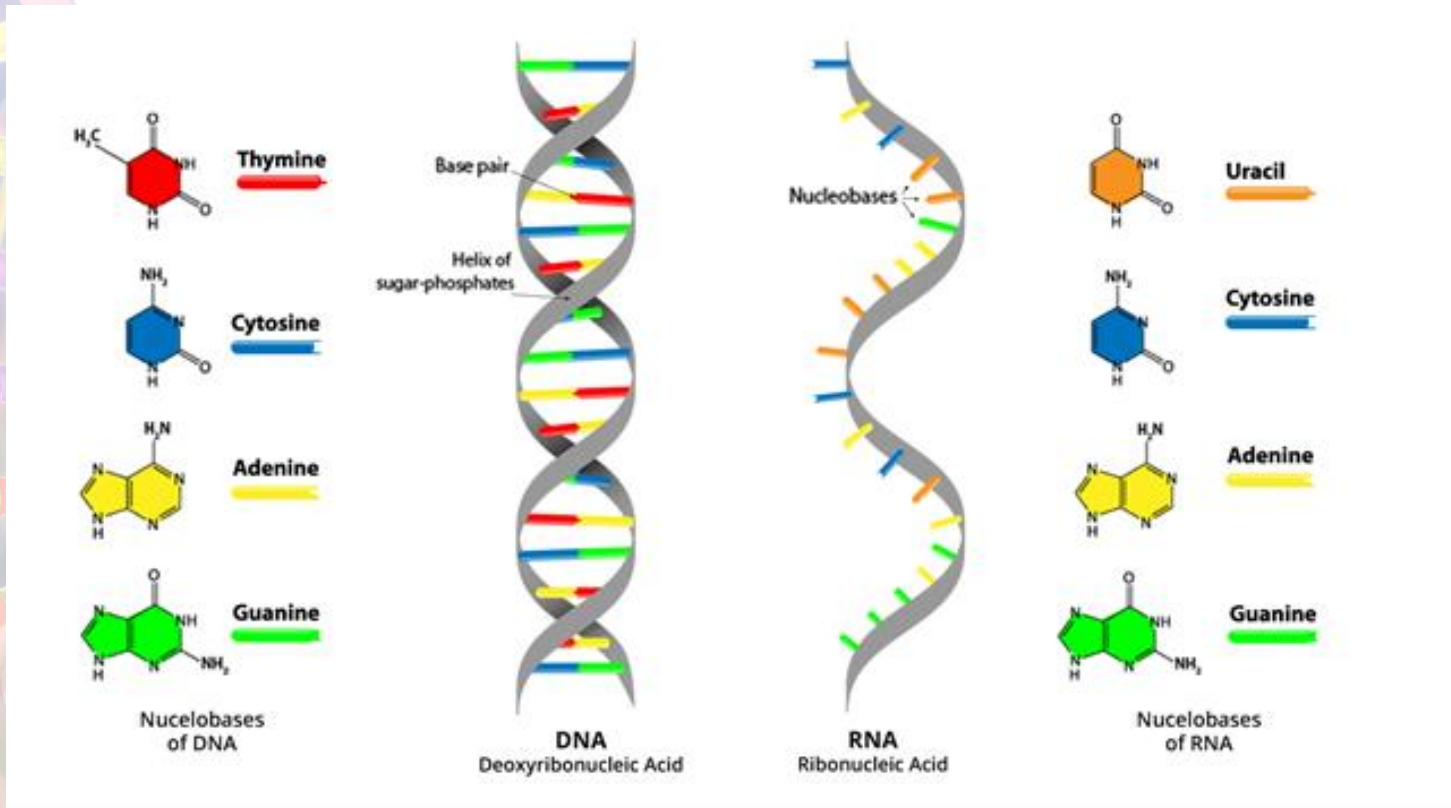
- Final assignment: finding a motif in a bacterial genome

# String

- A sequence of characters (char)
- s = "some text"

- substring = substr(s, start, end)

# DNA/RNA/proteins

Source: www.technologynetworks.com

# BioStrings

- **string = sequence**

- **DNAString: DNA sequence**

- **RNAString: RNA sequence**

- **AAString: aminoacid sequence (proteins)**

- s = "ACGT"

- DNAseq = DNAString(s)

Difference?

DNAString is a special data structure with its own methods

# DNA String

- IUPAC (extended)

- "-" (the gap letter)
- "+" (the hard masking letter)
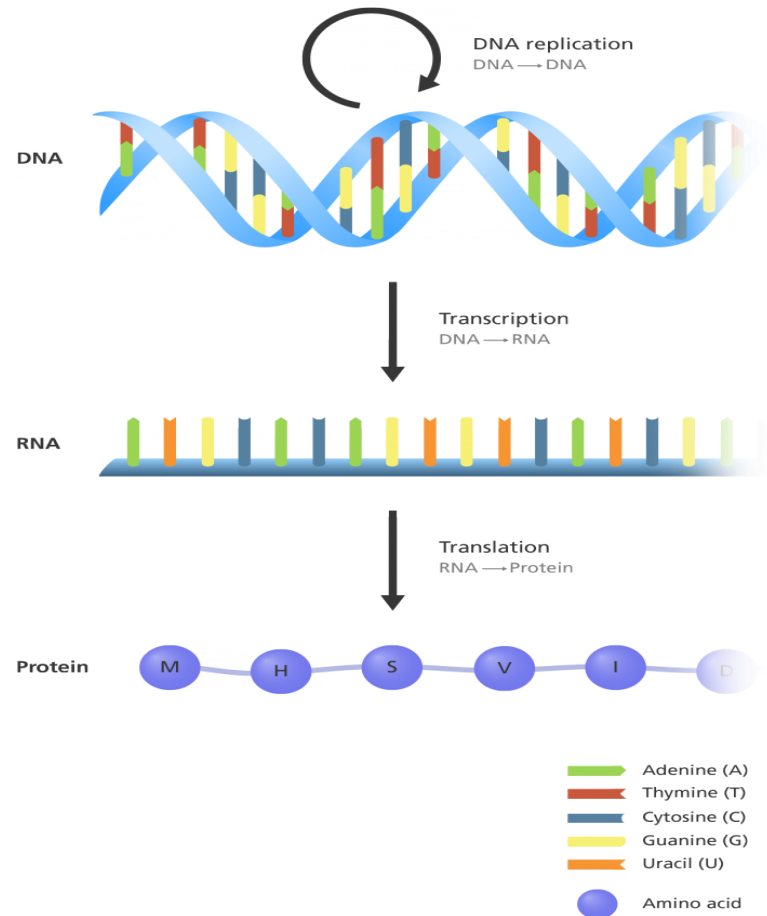- "." (the not a letter or not available letter)

- "*" (the stop letter)

Table 1.

IUPAC code for incomplete nucleic acid specification

| Symbol | Mnemonic | Translation |
| --- | --- | --- |
| A | | A (adenine) |
| C | | C (cytosine) |
| G | | G (guanine) |
| T | | T (thymine) |
| U | | U (uracil) |
| R | puRine | A or G (purines) |
| Y | pYrimidine | C or T/U (pyrimidines) |
| M | aMino group | A or C |
| K | Keto group | G or T/U |
| S | Strong interaction | C or G |
| W | Weak interaction | A or T/U |
| H | not G | A, C or T/U |
| B | not A | C, G or T/U |
| V | not T/U | A, C or G |

4

# The central dogma of molecular biology

**Useful functions:**
- **reverse**
- **complement**
- **reverseComplement**
- **translate**

# Transcription and translation



source: khanacademy

# Genetic code

- 20 amino acids
- 64 codons
- 3 stop codons



source: geneticcontrolandproteinfunction.wordpress.com

# Motifs

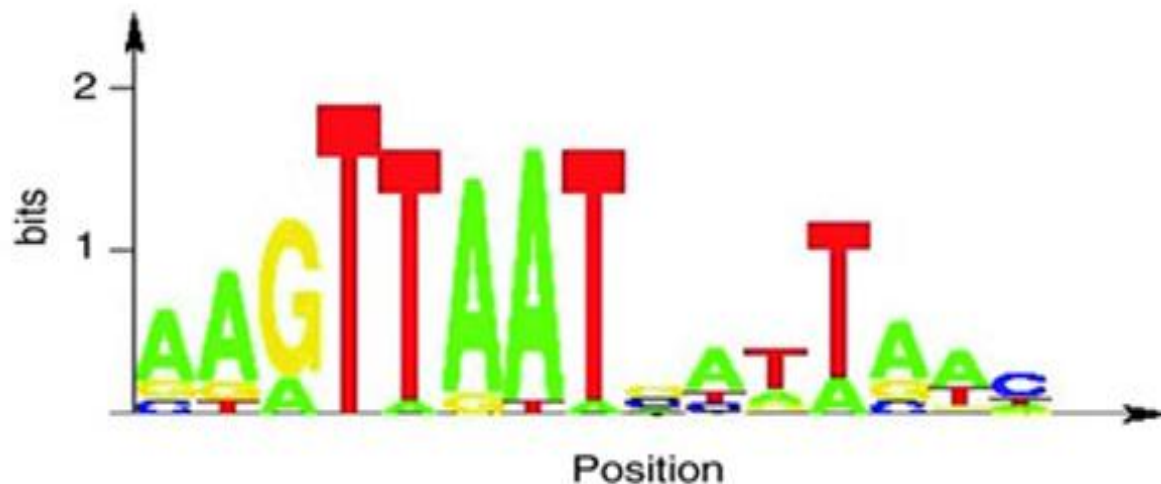The consensus sequence, position weight matrix and sequence logo of a motif



8

# ORF

= open reading frame (a sequence of codons uninterrupted by a STOP codon)

ATGCAATGGGGAAATGTTACCAGGTCCGAACTTATTGAGGTAAGACAGATTTAA

```
1.  ATG CAA TGG GGA AAT GTT ACC AGG TCC GAA CTT ATT GAG GTA AGA CAG ATT TAA
2.  A TGC AAT GGG GAA ATG TTA CCA GGT CCG AAC TTA TTG AGG TAA GAC AGA TTT AA
3.  AT GCA ATG GGG AAA TGT TAC CAG GTC CGA ACT TAT TGA GGT AAG ACA GAT TTA A
```

There are 3 reading frames

# The distance between two sequences

# Hamming distance

For two sequences of equal length the number of diferring positions

GAGCCTACTAACGGGAT
CATCGTAATGACGGCCT

# Point mutations

# Transitions/transversions

- *Transition:* between A and G, or. C and T
- G → A or  A → G
- C → T or T → C

- *Transversions*
- G → T or C
- A → T or C
- C → A or G
- T → A or G

# Mutations



(1) Insertion

(2) Substitution

**9 bases**

GAGACTTAC

(3) Deletion

A → GAGAC**A**TTAC → **10 bases**
GAGAC**A**TTAC
CTCTG**T**AATG

A
C → GAGA**A**TTAC → **9 bases**
GAGA**A**TTAC
CTCT**T**AATG

C → GAGA**-**TTAC → **8 bases**
GAGATTAC
CTCTAATG

14

# Edit distance

- transforming sequences
- transformation with the minimal number of changes
- operations:  **insertions**

  **deletions**

  **supstitutions**
- operations get penalties

```
I N T E * N T I O N
| | | | | | | | | |
* E X E C U T I O N
d s s     i s
```

*Levenshtein edit distance*

# Used to align sequences



**Local Alignment**

Target Sequence
5' ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA 3'
|||| ||||||| |||||||||||||||
Query Sequence 5' TACTCACGGATGAGGTACTTTAGAGGC 3'

**Global Alignment**

Target Sequence
5' ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA 3'
||||||||||| ||||||| |||||||||||||| |||||||
5' ACTACTAGATT----ACGGATC--GTACTTTAGAGGCTAGCAACCA 3'
Query Sequence

16

# IRanges



reduce

disjoin

# The *E. coli* genome

- 4.6 million nucleotides
- 4288 protein-coding genes
- 2584 operons

Binding of the transcription initiation complex to the promoter



Figure 8.1 *E. coli* RNA polymerase

THE CELL: A MOLECULAR APPROACH 7e, Figure 8.1
© 2016 Sinauer Associates, Inc.

# Final assignment

Discovering a motif:
left - a multiple sequence alignment
right - a motif logo
down - logarithmic position weight matrix

```
          aaTTGCGTCAtttc
gccgtcatactgTGACGTCTttcag
        actgaTGACGTCCatg
        gctcgtTGACGTCAccaaga
gagcggagcccgTGACGCGGccgagcggc
    tctctctttCCAGGTATctc
                . . .
        ggcttTGACGTCAgcctggc
tggaatctctgcTGACGTCAcgacactccgca
  cggcgggcatTGACGTCAaacggcagc
acccctccccgcTGACCTCActcgagccgccg
```

$$M(i, x) = \log_2 \frac{\text{frequency of letter } x \text{ at position } i}{\text{background frequency of letter } x} \qquad (1)$$

$$
\begin{array}{l}
A \\ C \\ G \\ T
\end{array}
\begin{bmatrix}
-3.219 & -3.219 & 3.785 & -3.219 & 1.396 & -3.219 & 2.084 & 3.467 \\
1.396 & 1.396 & -3.219 & 3.585 & -3.219 & 2.488 & 3.334 & -3.219 \\
1.396 & 3.690 & -3.219 & 2.084 & 3.690 & -3.219 & 1.396 & 1.396 \\
3.585 & -3.219 & -3.219 & -3.219 & -3.219 & 3.467 & 1.396 & 2.084
\end{bmatrix}
$$

20

# Final assignment

The consensus sequences (TTAGACA and TATAAT) and the position weight matrices for the two motifs in the sigma70 binding site

22