

Lab 2: From Genomics to Sequence Alignment

Intro to Genomics and
Sequence Alignment
Techniques

As. Univ. Drd. Ing. Bozdog Alexandru



Learning Objectives

- ✓ Define genus, genome, sequences
- ✓ ADN → ARN → protein
- ✓ Recognize biological file formats
- ✓ Understanding sequence alignment
- ✓ Practical application with BLAST, Clustal and Biopython



Genomics

- The study of an organism's complete set of genetic information.
- The genome includes both genes (coding) and non-coding DNA.
- 'Genome': the complete genetic information of an organism.

vs

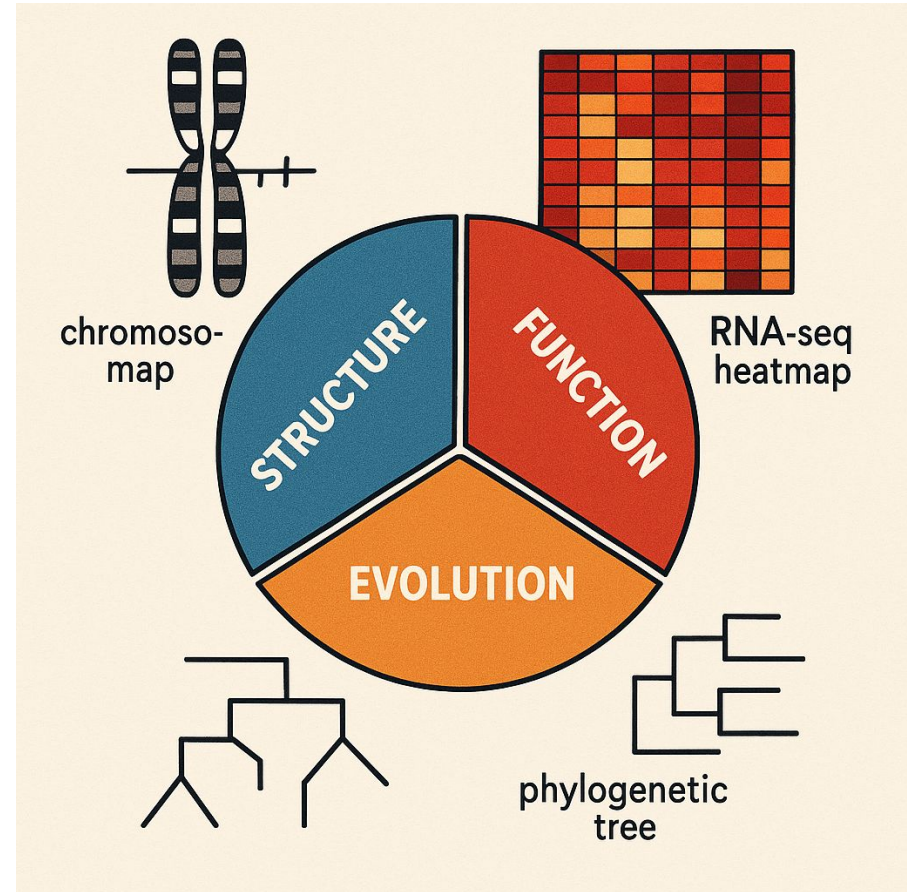


Genetics

- The study of heredity
- The study of the function and composition of single genes.
- 'Gene': specific sequence of DNA that codes for a functional molecule.

What is Genomics?

- Study of the whole DNA of an organism
- Fields of interest: Structural, Functional, Comparative, Translational
- Data obtained from high-throughput sequencing
- Relevance → medicine, evolution, biotechnology



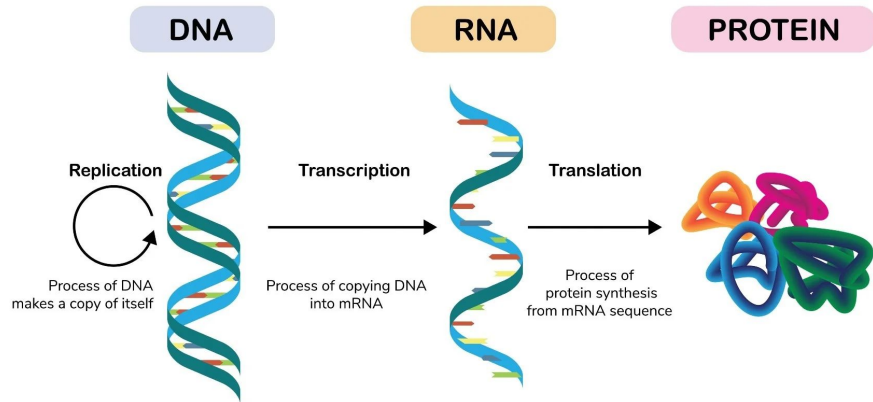
Central Dogma

BIOLOGY ● ● ●

Central Dogma

DNA → RNA → Protein

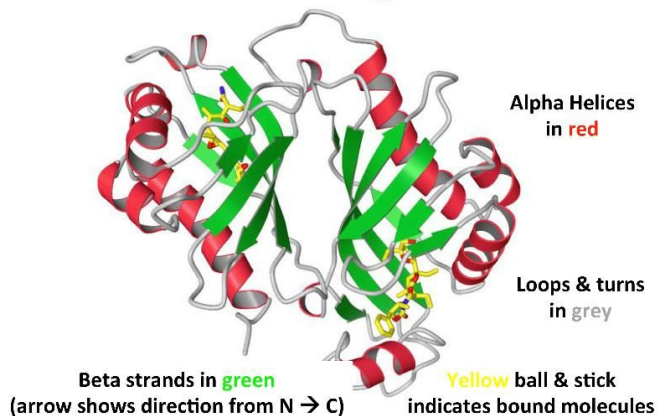
- Transcription and Translation
- Each level - different alphabet
- Flow of information = flow of data



Sequences and Biological formats

Type	Alphabet	Format	Example
DNA	A T C G	FASTA / FASTQ	>TP53 ATGCGTAAC
RNA	A U C G	FASTA	AUGCGAU
Protein	20 aa	FASTA / PDB	MKTAYIAKQ

Ribbon Diagrams

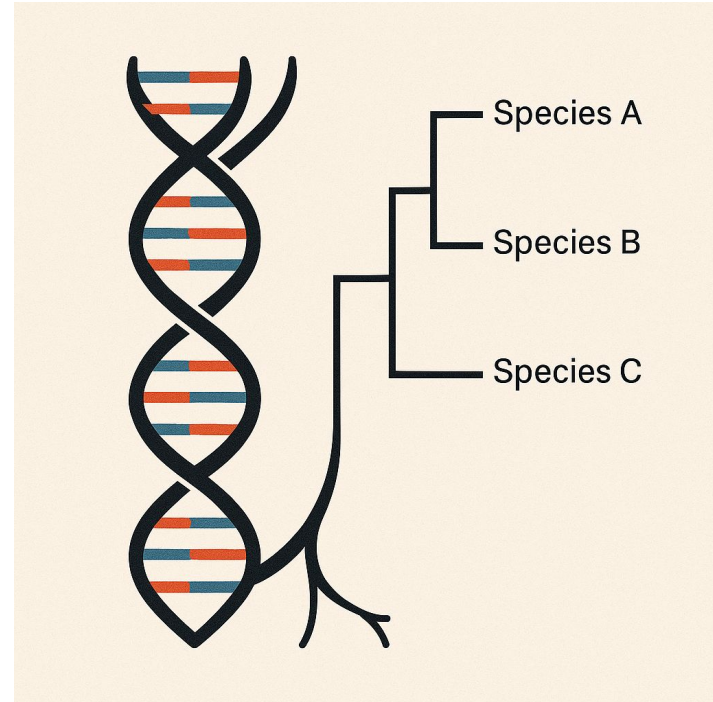


```

>NM_011040.3 Mus musculus transformation related protein 53 (Trp53), transcript variant 1, mRNA
TTTCCCTCCCACGTGCTCACCCTGCTAAAGTTCTGTAGCTTCAGTTTCATTG6GACCATCTG6CTGTA
GGTAGCACTACAGTTAGGGGGCACCTAGCATTCAAGCCCTCATCTCCTCTCCACAGAGGAGTGTAC
GCTTCCGAGAGCTGATGACTGCCATGAGGAGTCCAGTGGGATATCAGCTCGAGCTCCTCTGAG
CCAGGAGACATTTTCAGGCTTATGGAAGTACTCTCCAGAGATATCTGCCATCACTCATGATG
GACGATCTGTTGCTCCCAAGATGTTGAGGAGTTTTTGAAGGCCCAAGTGAAGCCCTCCGAGTGTCA
GAGCTCCTGCAAGCAGAGACCTGTACCCGAGACCCCTG6G6CCAGTGGCCCTGCCCCGACCTCCATG
GCCCTGTGATCTTTTGTCCCTTCTCAAAAAGTTACCAAGGCACTATG8CTTCCACTG6G6CTTCTG
CAGCTG6GACAGCCCAAGTCTGTTATGTGACGCTACTCTCTCCCTCAATAAGATTTCTGCCAGCTGG
CGAAGACGTGCCCTGTGAGTTGTGGTCAAGCCACACCTCCAGCTGGGAGCCGTGTCCGCGCATG6G
CATCTACAAGAGTCCAGCACATGAGGAGGTCGTGAGAGCTG6G6CCCACTGAGCGCTGCTCCGAT
G6TGTG6GCTG6CTCCTCCCAAGCATCTATCCG6GTGGAAGGAAATTTGATCCCGAGTATCTGGAAG
ACAGGCAAGCTTTTCCGACAGAGCTG6TGTGATCTTATGAGCCACCGAGGCGGCTGTGAGTATACCAC
CATCCACTACAAGTACATGTGTAATAGCTCTGATG6G6G6CATGAGCCGCACTATCCTTACCAC
ATCACACTGGAAGACTCCAGTGGGAACCTTCTG6GACG6GACAGCTTTGAGGTTCTGTTTGTGCTG6C
CTG6GAGAGACCCGCTACAGAGAGAGAAATTTCCBAAAAAGAGTCTTTG6CCTGAACTG6CCCC
AGGAGCGCAAGAGAGCGCTGCCACCTGCACAGCGCTCTCCCCGCAAAAGAAAAACCACTTGTAT
GAGAGATTTACCTTCAAGATCCG6G6GCTAAAGCTTCCAGATGTTCCG6GAGCTGAATGAGGCT
TAGAGTTAAGGATGCCATGCTACAGAGAGTCTG6GAGCAGCAGGCTCACTGAGCTACCTGAAGAC
CAAGAGG6GCAAGTCTACTTCCG6CATAAAAAACAATG6TCAAGAAAGTGG6G6CTGACTCAGACTGA
CTG6CTCTGATCCGCTCCCATCAGGCTCCCTCTCTGCTGCTTATGACTTCAGG6CTGAGAG
GACATCTCCG6GCTCTGCTGCTTTTTTACCTTGTAGCTAG6GCTCAGGCTCTCTGAGTATG
  
```

Why sequence alignment?

- Detect homology (common origin)
- Identify mutations, insertions, deletions
- Deduce function and structure
- Build phylogenetic trees



What is Sequence Alignment

- Rearranging sequences to maximize similarity
- Introducing gaps (–) for insertions/deletions
- Calculating an optimal score

Ancestral sequence:

A A T G C G A T G T C C
A A T G C G A T G T C C

Sequence derived from ancestral sequence:

A A T G A C G A – T G T G C C
A G T G – C G A G T T T – A C

mismatches

indels

Alignment:

A A T G A C G A – T G T G C C
| | | | | |
A G T G – C G A G T T T – A C

Types of Alignment

Type	Purpose	Algorithm	Use case
Global	Full length	Needleman–Wunsch	Similar genes
Local	Best fit region	Smith–Waterman	Motifs, domains
Semi-global	No edge loss	Hybrid	NGS Reads

Scores and Substitution Matrices

$$S = \sum_{i=1}^L s(x_i, y_i) - \sum \text{penalizări pentru gap-uri}$$

unde:

- $s(x_i, y_i)$ este scorul din matrice (ex. BLOSUM);
- penalizările provin din evenimentele de inserție/deleție.

- DNA: +1 match / -1 mismatch
- Protein: Matrices (Blossum / PAL)
- Gap penalties: linear, affine
- Goal: Maximize score

A	C	G	T
-1	-2	-2	1
-2	-2	-2	-1
-2	-2	-2	0

Gap Opening
vs. Extension

A C T G G



Gap opening

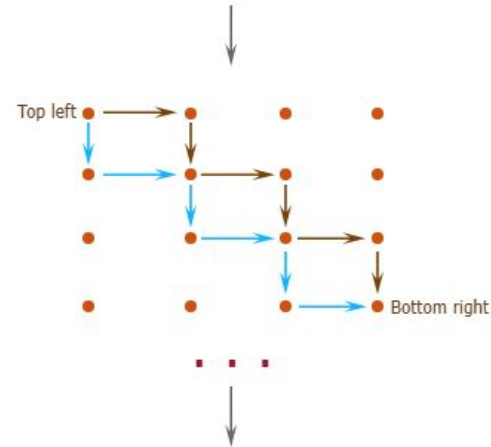
A C T G G

|

Dynamic programming

- 1 Initiate matrix
- 2 Recurently : $S(i, j) = \max(\text{diag} + \text{match}, \text{stânga} + \text{gap}, \text{sus} + \text{gap})$
- 3 Backtracking → Optimal Alignment
- 4 $O(n \times m)$ Complexity

Counting all possible paths from top left to bottom right of a m X n matrix



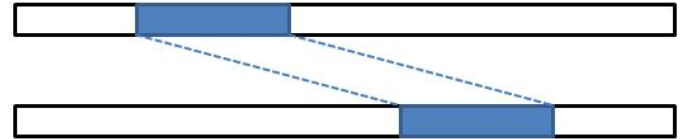
The all possible paths from top left to bottom right is : 20

Local vs Global Alignment

- **Global:** Covers complete sequence
- **Local:** covers only maximum score region
- **BLAST** = quick Smith-Waterman



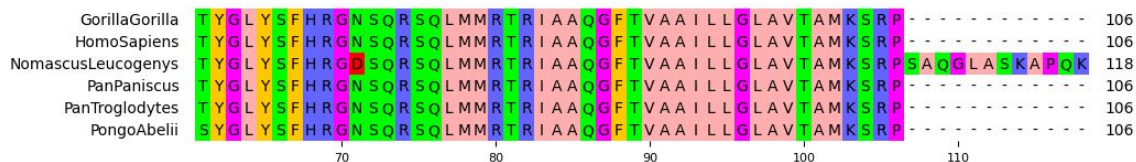
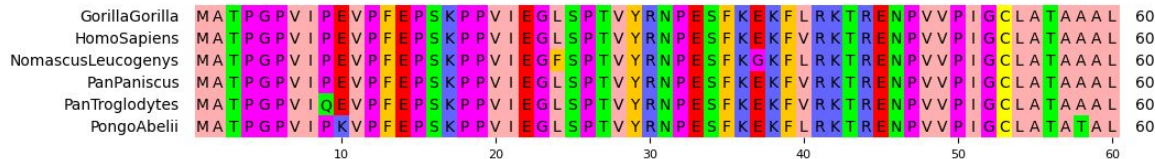
Global Alignment



Local Alignment

Multiple Sequence Alignment

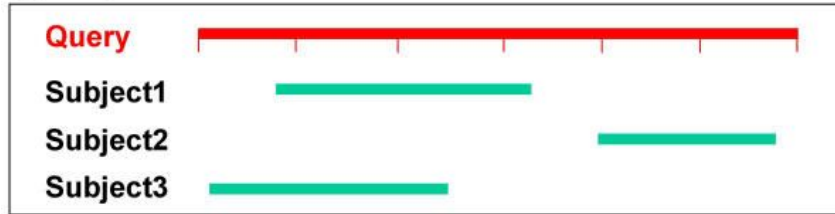
- 2 or more sequences aligned simultaneously
- Heuristics: progressive (Clustal Ω), consistency (T-Coffee)
- Identify conserved motifs, domains, evolutionary relationships



Tools and Libraries

```
from Bio import pairwise2
```

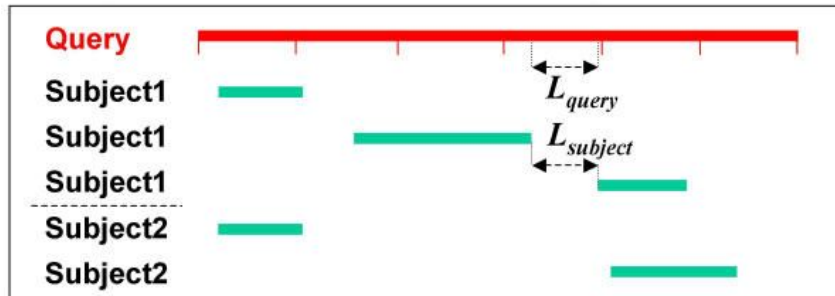
Type1 alignment: continuous match



CLI: `needle`, `water`, `blastn`, `clustalo`

Python (Biopython): `pairwise2`, `AlignIO`,
`Bio.Blast.NCBIWWW`

Type2 alignment: discontinuous matches in the same subject



Results :Alignment, Score, E-value

Alignment Evaluation

Metric	Semnification	Interpretation
Score	Sum of matches - penalties	Algorithm specific
% Identity	Exact matches - length	Similarity Degree
E-value	Probability Score appears by chance	Smaller = More significant
Consensus	Dominant character per column	MSA/ conservation

Use Cases & Conclusions

- Identify mutations and SNP-s
- Construct phylogenetic trees
- Predict proteic domains
- Drug repositioning
- Metagenomic analysis

✓ Genomics = data source

✓ Sequence alignment = First step of analysis

✓ Instruments: BLAST, Clustal, Biopython



Next lab : Aligning NGS reads and variant analysis