

Swiss Institute of Bioinformatics

#### INTRODUCTION TO SEQUENCING DATA ANALYSIS

## Read alignment

Deepak Tanwar Frédéric Burdet

April 23-25, 2025 Adapted from previous year courses





#### Learning Objectives

Understand the concept and purpose of sequence alignment

Distinguish between global and local alignment strategies and their algorithms

Describe how short and long sequencing reads are aligned to reference genomes

Understand indexing strategies (e.g., BWT, suffix arrays) used in fast read alignment tools



In bioinformatics, **alignment** refers to the process of arranging sequences of DNA, RNA, or proteins to identify regions of similarity.



## Alignment types

#### Pairwise alignment

- A. Global Alignment Needleman-Wunsch Algorithm
  - DNA, RNA, or protein sequences of similar length
- B. Local Alignment Smith-Waterman Algorithm
  - Protein or RNA/DNA domains

#### Multiple Sequence Alignment (MSA)

Aligns three or more sequences to detect conserved regions.

#### Read Alignment (e.g. BWA, Bowtie2, minimap2)

• Mapping short reads (DNA/RNA) to a reference genome



## Pairwise alignment

#### **Global Alignment**

- Aligns entire sequences from end to end
- Introduces gaps to optimize alignment length and score

#### **Local Alignment**

- Aligns the best-matching subsequences
- Useful for finding **domains** or **conserved motifs**



Sequence1 : A T A C C G G A T A T T

 $Sequence 2: A \ A \ C \ G \ G \ A \ C \ C \ T$ 



Sequence1 : A T A C C G G A T A T T

 $Sequence 2: A \ A \ C \ G \ G \ A \ C \ C \ T$ 





Sequence1: A T A C C G G A T A T T

 $Sequence 2: A \ A \ C \ G \ G \ A \ C \ C \ T$ 



Alignment scoring



Sequence1: A T A C C G G A T A T T

 $Sequence 2: A \ A \ C \ G \ G \ A \ C \ C \ T$ 



Alignment scoring

Sı	S2	Score
Α	-	-1
Т	-	-1
А	А	+1
С	А	-1
С	С	+1
G	G	+1
G	G	+1
А	А	+1
Т	С	-1
А	С	-1
Т	С	-1
Т	Т	+1
Total		0



Sequence1 : A T A C C G G A T A T T

 $Sequence 2: A \ A \ C \ G \ G \ A \ C \ C \ T$ 

Alignment scoring



Sequence1 : A T A C C G G A T A T T

Sequence2 : A A C G G A C C C T





Sequence1 : A T A C C G G A T A T T

Sequence2 : A A C G G A C C C T



Sı	S2	Score	
Α	-	0	
Т	-	0	
Α	А	+1	
С	А	-1	
С	С	+1	
G	G	+1	
G	G	+1	
Α	А	+1	
Т	-	0	
Α	-	0	
Т	-	0	
Т	-	0	
Tot	al	4	



#### Quiz 1

Which of the following algorithms is used for global sequence alignment?

- A. Smith-Waterman
- B. BLAST
- C. Needleman-Wunsch
- D. Clustal Omega



#### Quiz 2

In local alignment using the Smith-Waterman algorithm, which of the following best describes the approach?

A. Aligns sequences end-to-end, including gaps.
B. Identifies the best matching subsequences within two sequences.
C. Aligns three or more sequences for motif detection.
D. Only aligns protein sequences.



For phylogeny, function prediction, or motif discovery.



For phylogeny, function prediction, or motif discovery.

Α	Т	С	G
Α	Т	G	G
A	С	G	



For phylogeny, function prediction, or motif discovery.





For phylogeny, function prediction, or motif discovery.





For phylogeny, function prediction, or motif discovery.





For phylogeny, function prediction, or motif discovery.





For phylogeny, function prediction, or motif discovery.

Tool: Clustal Omega, MAFFT Command-line example with Clustal Omega



**MSA score**: 3 - 1 - 1 - 1 = 0



## **Read Alignment**

Mapping millions of short reads to a reference genome. Read Aligners: BWA, Bowtie2, STAR, etc.





Aim: generate a 'phonebook' for fast searches





Aim: generate a 'phonebook' for fast searches

# Reference: TAATA\$



Aim: generate a 'phonebook' for fast searches

# Reference: TAATA\$

0	Т	Α	Α	Т	Α	\$
1	A	А	Т	А	\$	
2	A	Т	А	\$		
3	Т	А	\$			
4	A	\$				
5	\$					



Aim: generate a 'phonebook' for fast searches





Aim: generate a 'phonebook' for fast searches







## Reference: TAATA\$ Query: ATA





#### Indexing and querying

Suffix array: large, same sequence stored multiple times BWT: only first and last columns are stored -> still enables fast querying



suffix array



## Indexing and querying

Suffix array: large, same sequence stored multiple times BWT: only first and last columns are stored -> still enables fast querying



suffix array

Burrows-Wheeler Transformation

\$	Т	Α	Α	Т	Α
Α	\$	Т	Α	Α	Т
Α	Α	Т	Α	\$	Т
Α	Т	Α	\$	Т	Α
Т	A	\$	Т	Α	Α
Т	Α	Α	Т	Α	\$



#### Global vs local

#### Global (end-to-end)

Read:GACTGGGCGATCTCGACTTCG||||||||||Reference:GACTG--CGATCTCGACATCG



#### Quiz 3

## What is the primary role of the Burrows-Wheeler Transform (BWT) in read alignment?

- A. To assemble the genome de novo.
- B. To visualize phylogenetic trees.
- C. To enable fast searching by storing only first and last columns.
- D. To simulate RNA splicing events.



#### Software

**Basic alignment:** 

- » bowtie2 (BWT; default = global)
- >>> bwa-mem (BWT; default = local )

Splice-aware (RNA-seq):

- » hisat2
- » STAR

Long reads + short reads + splice-aware:

» minimap2





## Mapping quality





## Mapping quality



 $\begin{array}{l} MAPQ \\ = -10 log_{10} \Pr\{mapping \ position \ is \ wrong\} \end{array}$ 

$$-10log_{10}(0.01) = 20$$
  
 $-10log_{10}(0.5) = 3$ 





#### Quiz 4

Which tool is best suited for aligning RNA-seq reads with splice awareness?

A. BWA B. Bowtie2 C. STAR D. Clustal Omega



#### Quiz 5

#### What distinguishes global from local alignment in the context of Bowtie2?

- A. Global allows for splicing, local does not.
- B. Global requires fewer reads than local.
- C. Global aligns full sequences end-to-end; local allows for partial matches.
- D. Global is only for RNA, local is for DNA.



#### Summary

Global alignment aligns full sequences end-to-end

Local alignment is suited for partial matches like conserved domains

**MSA** aligns multiple sequences to reveal evolutionary or functional patterns

Read alignment maps millions of sequencing reads to a reference genome

Suffix arrays and BWT are used to enable fast searching within large genome indexes

Tools like **BWA**, **Bowtie2**, **STAR**, and **minimap2** support various alignment types, with specific strengths for short reads, long reads, or spliced transcripts





## Thank you

DATA SCIENTISTS FOR LIFE





sib.swiss