## Lecture 6

#### Introduction to bioinformatics (MVE510)

Autumn 2020

Additional reading: Mapping reads on a genomic sequence: an algorithmic overview and a practical comparative analysis. Schbath S, Martin V, Zytnicki M, Fayolle J, Loux V, and Gibrat JF. *Journal of Computational Biology*, 19(6) 2012.

### Repetition

- The Needleman-Wunsch algorithm can be used to find <u>global alignments</u>. The best alignment is identified by iteratively fill a alignment matrix and backtrack from the highest value.
- The Smith-Waterman algorithm can be find <u>local</u> <u>alignments</u>. The best alignment is found similarly to the NW algorithm but where all negative values has been replaced by zeros and the alignment can start and end at any positions.

### Repetition

- BLAST uses a 'seed-and-extend' algorithm to efficiently calculate local alignments. This makes it 50 times fast than Smith-Waterman.
- BLAST is however still to slow to efficiently map reads from next generation sequencing to a reference.

# The BLAST algorithm (Basic Local Alignment Search Tool)

#### Is BLAST fast enough?

No! Next generation DNA sequencing produce too many reads (50x faster than SW is not enough!).

#### Can we further improve the speed?

Yes, but we need some additional concepts from computer science!

### What is an index?

- An index is a <u>data structure</u> that enable fast lookup of <u>keys</u>
- In our case the data structure will built form the reference sequence and key our query sequence
- Many form of indices assume that the reference is long (e.g. a genome) and the query sequence is short (e.g. a NGS read)
- Once the index is built, we can use it to lookup to find where specific read match the reference.
- We will focus on exact matches but it is also possible to include mismatches and gaps.



#### Naive search



### Suffix trees

• Are highly efficient to find exact matches in a sequence. The complexity is

#### $O(L_R)$

for each read.

• Note that *L<sub>G</sub>* is gone (!!!). <u>The size of the reference</u> <u>does no longer matter!</u>

### Suffix trees

#### BUT.... (there is no free lunch!)

- Building the suffix trees takes time  $(O(L_G))$  but we only need to do this once.
- Suffix trees are very large. The suffix tree for the entire human genome is so large that it can not be stored in the memory of a standard computers.
- If we want to include mismatches and gaps, the complexity increases fast!

#### Reference: GAGAGGCAGC\$

Suffixes	Position
GAGAGGCAGC\$	1
AGAGGCAGC\$	2
GAGGCAGC\$	3
AGGCAGC\$	4
GGCAGC\$	5
GCAGC\$	6
CAGC\$	7
AGC\$	8
GC\$	9
C\$	10
\$	11

#### Reference: GAGAGGCAGC\$

Sorted suffixes	Sorted position
AGAGGCAGC\$	2
AGC\$	8
AGGCAGC\$	4
CAGC\$	7
C\$	10
GAGAGGCAGC\$	1
GAGGCAGC\$	3
GCAGC\$	6
GC\$	9
GGCAGC\$	5
\$	11

#### Reference: GAGAGGCAGC\$

Sorted suffixes	Sorted position	Cylinder suffix array	<b>BW-transform</b>
AGAGGCAGC\$	2	AGAGGCAGC\$G	G
AGC\$	8	AGC\$GAGAGGC	С
AGGCAGC\$	4	AGGCAGC\$GAG	G
CAGC\$	7	CAGC\$GAGAGG	G
С\$	10	C\$GAGAGGCAG	G
GAGAGGCAGC\$	1	GAGAGGCAGC\$	\$
GAGGCAGC\$	3	GAGGCAGC\$GA	А
GCAGC\$	6	GCAGC\$GAGAG	G
GC\$	9	GC\$GAGAGGCA	А
GGCAGC\$	5	GGCAGC\$GAGA	А
\$	11	\$GAGAGGCAGC	С

1	2	3	4	5	6
\$	А	\$A	AC	\$AC	ACG
G	А	GA	AC	GAC	ACT
А	С	AC	CG	ACG	CGA
А	С	AC	CT	ACT	CT\$
С	G	CG	GA	CGA	GAC
С	Т	СТ	Τ\$	CT\$	T\$A
Т	\$	Т\$	\$A	T\$A	\$AC
7	8	9	10	1:	1
\$ACG	ACGA	\$ACGA	A	ACGAC	\$ACGAC
GACT	ACT\$	GACT\$	A	ACT\$A	GACT\$A
ACGA	CGAC	ACGAC	(	CGACT	ACGACT
ACT\$	CT\$A	ACT\$A	(	CT\$AC	ACT\$AC
CGAC	GACT	CGACT	(	GACT\$	CGACT\$
CT\$A	T\$AC	CT\$AC	Т	<sup>-</sup> \$ACG	CT\$ACG
T\$AC	\$ACG	T\$ACG	ć	SACGA	T\$ACGA

12	13	14
ACGACT	\$ACGACT	ACGACT\$
ACT\$AC	GACT\$AC	ACT\$ACG
CGACT\$	ACGACT\$	CGACT\$A
CT\$ACG	ACT\$ACG	CT\$ACGA
GACT\$A	CGACT\$A	GACT\$AC
T\$ACGA	CT\$ACGA	T\$ACGAC
\$ACGAC	T\$ACGAC	\$ACGACT

<b>BW-transform</b>	Sorted BW-transform	Sorted position
G	A	2
С	A	8
G	A	4
G	С	7
G	С	10
\$	G	1
1 A ←	G	3
G <b>≁</b> <u>No</u>	match G	6
2 A	G	9
3 A	G	5
С	\$	11

BW-t	ransform	Sorted BW-tra	ansform	Sorted position
	G <b>←</b> No m	atchA	1	2
1	C ←	A	2	8
	G <b>←</b>	atch A	3	4
	G	С		7
	G	С		10
	\$	G		1
	A	G		3
	G	G		6
	A	G		9
	A	G		5
	С	\$		11

BW-transform	Sorted BW-transform	Sorted position
G	A	2
С	A	8
G	A	4
G ←	C 1	7
G	С	10
\$	G	1
A	G	3
G	G	6
A	G	9
А	G	5
С	\$	11



BW-transform	Sorted BW-transform	Sorted position
G	A	2
С	A	8
G	A	4
G	С	7
G	С	10
\$	G	1
A	G	3
G	G	6
A	G	9
А	G	5
С	\$	11



### The Borrows-Wheeler Aligner (BWA)

- Based on seed-and-extend approach.
- Seeds are identified using suffix arrays and the Burrows-Wheeler transform.
- Seeds are extended into full alignments using more accurate algorithms.
- Optimized for large volumes of reads from next generation DNA sequencing.
- You will use this aligner in the computer exercises.

### Choosing a suitable aligner



### Choosing a suitable aligner

#### **Smith-Waterman**

Slow but highly sensitive. Software: FASTA and EMBOSS.

#### Seed and extend

Faster but have reduced sensitivity.

Software: BLAST

### Choosing a suitable aligner

#### Suffix-array and Burrow-Wheeler transform

Very fast, low sensitivity. But often suitable for mapping sequence read to a reference.

Software: BWA

### Summary of this lecture

- An index is a data structure that enables fast lookup of exact matches in a DNA sequence.
- Suffix trees are crated from the suffixes of a DNA sequence. By traversing the tree we can find the position of any subsequence.
- The Burrows-Wheeler transform orders the genome in a specific way that makes identification of subsequences highly efficient.

### Summary of this lecture

 Read alignment using suffix arrays and the Burrows-Wheeler transform are as fast as a suffix tree but does not require the entire tree to be created and stored in memory.