## Lecture 5

### 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.

# What we did we do in the previous lecture?

- SNPs, indels and structural variants can be identified by whole genome sequencing
- Three main steps: read mapping, quality refinement and SNP calling
- SNP calling is done through statistical models describing the errors in the data
- Analysis of whole genome sequencing results in a massive data reduction: from billions of observations to a handful of mutations



## Today's agenda

- Pair-wise alignment of DNA sequences
  - The Needleman-Wunsch algorithm
  - The Smith-Waterman algorithm
- Computational complexity
- BLAST: The Basic Local Alignment Search Tool

## Our main problem this lecture

### The sequence reads – many and short

 Hundreds of million reads from sequencing, each only a few hundred nucleotides long

### The reference – few and large

- The human genome consisting of 23 chromosomes (total 3.2×10<sup>9</sup> nucleotides)
- A bacterial genome (one chromosome, 4.6×10<sup>6</sup> nucleotides for *E. coli*)

## Our main problem this lecture

How do we find where our sequence read match ('align') to the reference sequences?

## Example of alignments

Alignment 1

Alignment 2

Alignment 3





## Two main forms of alignment

### **Global alignment**

Two sequences are aligned over their full length.

### Local alignment

Two sequences are aligned based on their best matching subsequences.

<u>Local alignments</u> are used to match short sequence reads against long reference sequences.

## Example

Let  $S(a, b) = \begin{cases} 5, a = b \\ -4, a \neq b \end{cases}$  and d = -7.

Score the following alignments.

Alignment 1

	AGTCTAGT	<b>S-5+5+5-1+5+5+5-31</b>
	AGTATAGT	3-3+3+3+3+3+3+3+3=31
Alignment 2		
	TTGA-TGA	<b>S-5 7+5+5 7+5+5+5-16</b>
	T-GACTGA	3-3-7+3+3-7+3+3+3-10
Alignment 3		
	TGTAACCT	S=-7-7+5+5+5-4+5+5-7-7=-9
	TAAGCTAG	

## Example

			1	2	3	4	5	6
			Т	С	G	G	А	Т
		0	-5	-10	-15	-20	-25	-30
1	Т	-5						
2	С	-10						
3	С	-15						
4	A	-20						
5	Т	-25						

## Example

			1	2	3	4	5	6
			Т	С	G	G	A	Т
		0	-5	-10	-15	-20	-25	-30
1	Т	-5	5 ←	04	∎ -5♦	-10	-15∢	-20
2	С	-10	0	10	5 🗲	0 <	-5	-10
3	С	-15	-5	5	6 🕇	1	-4	<b>-</b> 9
4	А	-20	<b>∎</b> -10			2	6 4	<b>1</b>
5	Т	-25	-15	-5	-4	-3	1	11

## Computational complexity

- Measures how many computations that are needed to execute an algorithm
- Fewer computations means faster algorithm

Computational complexity is measured using *O* notation (*O* stands for *ordo* which means order in latin).

## Computational complexity

If an algorithm is O(n) this means that the number of computations grows linearly with respect to n (which can, for example, be the number of input data).

Other examples:  $O(n^2)$ ,  $O(\log n)$ , O(nm), etc.



## Computational complexity

In our application, Smith-Waterman will have a complexity of

### $O(L_R L_G)$

for each read. Here,

 $L_R$  = the length of the reads,

 $L_G$  = the length of the genome (reference).

This can become <u>very</u> slow if  $L_G$  is large (which it is in our case).

### Futhermore....

Alignment of a 100 bases long sequence read against the first human chromosome requires a matrix of dimension 100×247,199,720. This which will take 100 gigabytes of RAM if stored in memory.

Smith-Waterman are too slow and resource heavy! Can we make the alignment faster?

BLAST can be used to match any sequence ('query') against a reference ('database'). The algorithm work as follows

- 1. Create a table of words (subsequences) of size W and their location in the reference (index, hash).
- 2. Find the position of the words present in the sequence. These positions are called 'seeds'.
- 3. Extend the alignment around the seeds.

BLAST is around 50 times faster than Smith-Waterman.

**Reference ('database')** 

ACGAGTGAGTGCCGAGTACGTAGCGTAGGAGTGAGTTGGAGTGAGACGTGAGT

#### Table (hash) of words of size 7

Maximum 4<sup>7</sup>=2<sup>14</sup>=16384 rows

Word	Position
ACGAGTG	1
CGAGTGA	2
GAGTGAG	3, 38
AGTGAGT	4, 30





BLASTN 2.2.26 [Sep-21-2011]			
Query= read1 (50 letters)			
Database: CCDS_nucleotide.20160908.fna 32,554 sequences; 55,631,471 total letters			
Searchingdone			
Sequences producing significant alignments:	Score (bits)	E Valı	ıe
CCDS73971.1 Hs108 chr17		84	1e-16
CCDS73970.1 Hs108 chr17		84	1e-16
CCDS73969.1 Hs108 chr17		84	1e-16
>CCDS73971.1 Hs108 chr17 Length = 909			
Score = 83.8 bits (42), Expect = 1e-16 Identities = 48/50 (96%) Strand = Plus / Plus			
Query: 1 aggctgctcccccggtggcctctgcaccagcagctcctacaccggcggcc 50 	)		

- The word size changes the sensitivity
  - Large word size means faster algorithm but lower sensitivity
  - Smaller word size means slower algorithms but higher sensitivity.
- BLAST can compare
  - Nucleotide vs nucleotide ('blastn')
  - Protein vs protein ('blastp')
  - Nucleotide vs protein ('blastx')
  - Protein vs nucleotide ('tblastn')



https://www.ncbi.nlm.nih.gov/genbank/

$\leftrightarrow$ $\rightarrow$ C $\blacksquare$ Säke	rt https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastn&PAGE_TYPE=BlastSearch&LINK_LOC=blasthome			☆	: 1 0
NIH U.S. Nationa	Library of Medicine NCBI National Center for Biotechnology Information			Sign in 1	o NCBI
BLAST <sup>®</sup> » bla	stn suite	Home	Recent Results	Saved Strategies	Help
	Standard Nucleotide BLAST				
blastn blastp blas	tx (blastr blastr				
Enter Query S	equence BLASTN programs search nucleotide databases using a nucleotide query. more		Reset page	Bookmark	
Enter accession n	umber(s), gi(s), or FASTA sequence(s) 🥹 Clear Query subrange 🐌				
aggotgotococcegteg	stttg:sc:sg:sg:tcttg:sc:gg:sg:sc!     From       To				
Or, upload file	Valij fil Ingen fil har valts				
Job Title					
	Enter a descriptive title for your BLAST search 🥹				
Align two or me	re sequences 🕑				
Choose Searc	h Set				
Database	⊖ Human genomic + transcript ⊙Mouse genomic + transcript ⊛Others (nr etc.):				
	Nucleotide collection (nr/nt)				
Organism Optional	Enter organism name or id-completions will be suggested  Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown @				
Exclude Optional	Models (XM/XP) Uncultured/environmental sample sequences				
Limit to	Sequences from type material				
Entrez Query	Yeu 1013 Create custom database				
Optional	Enter an Entrez query to limit search 🤪				
Program Sele	tion				
Optimize for	Highly similar sequences (megablast)				
	More dissimilar sequences (discontiguous megablast)				
	Somewhat similar sequences (blastn)				
	Choose a BLAST aigontinm 🥪				
BLAST	Search database Nucleotide collection (nr/nt) using Megablast (Optimize for highly similar sequences)				
Algorithm parame	ters				
					You
BLAST is a registered tr	ademark of the National Library of Medicine		S	upport center Mailing	list 🛄
4					•

https://www.ncbi.nlm.nih.gov/genbank/

☆

S

С

Säkert https://blast.ncbi.nlm.nih.gov/Blast.cgi

1	Jignments 📳 Download 🐱 GenBank Graphics Distance tree of results						
	Description	Max score	Total score	Query cover	E value	Ident	Accession
	PREDICTED: Gorilla gorilla gorilla tumor protein p53 (TP53), transcript variant X3, mRNA	82.4	82.4	100%	4e-13	96%	XM 004058511.2
	PREDICTED: Gorilla gorilla gorilla tumor protein p53 (TP53), transcript variant X2, mRNA	82.4	82.4	100%	4e-13	96%	XM 019013137.1
	PREDICTED: Gorilla gorilla gorilla tumor protein p53 (TP53), transcript variant X1, mRNA	82.4	82.4	100%	<mark>4</mark> e-13	96%	XM 019013136.1
	Homo sapiens isolate PG/1/2016 TP53 (TP53) gene, partial cds	82.4	82.4	100%	4e-13	96%	KU902450.1
	Homo sapiens phosphoprotein p53 (p53) gene, complete cds	82.4	82.4	100%	4e-13	96%	AH002919.2
	PREDICTED: Pan troglodytes tumor protein p53 (TP53), transcript variant X3, mRNA	82.4	82.4	100%	4e-13	96%	XM 001172077.4
	PREDICTED: Pan troglodytes tumor protein p53 (TP53), transcript variant X2, mRNA	82.4	82.4	100%	4e-13	96%	XM 016931471.1
	PREDICTED: Pan troglodytes tumor protein p53 (TP53), transcript variant X1, mRNA	82.4	82.4	100%	4e-13	96%	XM 016931470.1
	PREDICTED: Pan paniscus tumor protein p53 (TP53), transcript variant X1, mRNA	82.4	82.4	100%	4e-13	96%	XM 014341943.1
	PREDICTED: Pan paniscus tumor protein p53 (TP53), transcript variant X2, mRNA	82.4	82.4	100%	4e-13	96%	XM 003810066.2
	TPA: Homo sapiens Processed transcript p53-mRNA (p53 gene)	82.4	82.4	100%	4e-13	96%	HG975427.1
	Tupaia chinensis tumor protein p53 (TP53), mRNA	82.4	82.4	100%	4e-13	96%	NM 001287369.1
	Homo sapiens isolate C138 nonfunctional tumor suppressor p53 (TP53) gene, partial sequence	82.4	82.4	100%	4e-13	96%	JQ752242.1

Download + OE	nBank Graphics			Vext 🛦 Previous 🛕 Descriptions
PREDICTED: Go	rilla gorilla gorilla tur	or protein p53 (TP53), tra	anscript variant X3, mRNA	
sequence ID: <u>XIVI_U</u>	14050511.2 Length: 25	of Number of Matches: 1		<b>Related Information</b>
Range 1: 405 to 454	GenBank Graphics	tities Cans	V Next Match A Previous Match	
82.4 bits(90)	4e-13 48/5	0(96%) 0/50(0%	) Plus/Plus	
Query 1 AGGCT       Sbjct 405 AGGCT			50 454	
Download v <u>Ge</u>	nBank Graphics			▼ Next ▲ Previous ▲ Descriptions
PREDICTED: Go Sequence ID: <u>XM_0</u> Range 1: 456 to 50!	rilla gorilla gorilla tum 1 <u>9013137.1</u> Length: 26 i <u>GenBank Graphics</u>	Ior protein p53 (TP53), tra 18 Number of Matches: 1	Next Match A Previous Match	Related Information
Score 82.4 bits(90)	Expect Iden 4e-13 48/5	tities Gaps 0(96%) 0/50(0%	) Plus/Plus	
Query 1 AGGCT       Sbjct 456 AGGCT Download ~ Ge	CTCCCCCGGTGGCCTCTGCA		50 505	▼ Next ▲ Previous ▲ Descriptions
	rilla gorilla gorilla turr	ior protein p53 (TP53), tra	anscript variant X1, mRNA	_
PREDICTED: Go	19013136.1 Length: 26	91 Number of Matches: 1		Deleted Televentine
PREDICTED: Go Sequence ID: <u>XM_0</u>			🔻 Next Match 🔺 Previous Match	Related Information
PREDICTED: Go Sequence ID: <u>XM_0</u> Range 1: 529 to 578	GenBank Graphics			
PREDICTED: Go Sequence ID: <u>XM_0</u> Range 1: 529 to 578 Score 82.4 bits(90)	<u>GenBank</u> Graphics Expect Iden 4e-13 48/5	tities Gaps 0(96%) 0/50(0%	Strand ) Plus/Plus	

### 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 more concepts from computer science!

## Summary of lecture 5

- 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.

## Summary of lecture 5

- 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.