

NGS Bioinformatics

Álvaro Cortés Calabuig
October 2019

Overview

- This workshop: why?
- Bioinformatics NGS terminology
- NGS bioinformatics pipelines

Typical RNA-Seq Experiment Produces:



Herenstraat 49 PO box 602,
3000 Leuven
tel: +32 16 33 08 21
mail: genomicscore@uzleuven.be

Project identification

Project Type RNA seq
Number Of Samples 16
Number Of Conditions 2
Condition 1 MMP-9_KO_water
Condition 2 MMP-9_KO_DSS

Used Read Information

Sample Code	Condition	Used Read Count	Size Factor
GCr24962	MMP-9_KO_water	4962926	0.481151532775072
GCr24963	MMP-9_KO_water	7160982	1.27716186751061
GCr24964	MMP-9_KO_water	6036545	1.06402993138704
GCr24965	MMP-9_KO_water	4594610	0.931302721027977
GCr24966	MMP-9_KO_water	6983397	0.951266712762562
GCr24971	MMP-9_KO_DSS	5571447	0.951144771034544
GCr24972	MMP-9_KO_DSS	7308085	1.129141569201
GCr24973	MMP-9_KO_DSS	6299628	1.11078799446147
GCr24978	MMP-9_KO_water	5922351	1.200961919252668
GCr24979	MMP-9_KO_water	3319976	0.643371024301831
GCr24980	MMP-9_KO_water	5768701	1.21066952669347
GCr24981	MMP-9_KO_water	6813443	1.45667411227976
GCr24986	MMP-9_KO_DSS	5754245	0.620228663906365
GCr24987	MMP-9_KO_DSS	10620259	1.43676485961185
GCr24988	MMP-9_KO_DSS	7968821	1.25322666394006
GCr24989	MMP-9_KO_DSS	6823721	1.28413830085057

Sample Relations

Data quality assessment and quality control are essential steps of any data analysis. Here we define the term quality a fitness for purpose. Our purpose is the detection of differentially expressed genes, and we are looking in particular in samples whose experimental treatment suffered from an abnormality that renders the data points obtained from them particular samples detrimental to our purpose.

Variance stabilized data is used to create sample-to-sample distances. With these distances sample clustering becomes possible. The clustering should reflect the experimental design correctly; samples are more similar when they have the same treatment. The heatmap with tree should show this effect (Figure ③). The Principal Component Analysis plot (PCA-plot) (Figure ④) is a 2 dimensional version of these analysis. Expected is that samples with a same treatment cluster together. Outliers and possible bias are easy to detect.

Top 30 highly expressed genes

By taking a look at the top 30 highly expressed genes, a first impression of the data can be made. The heatmap below shows this expression data (Figure ⑤). The data is normalized by using the variance stabilisation transformation. Sample should cluster together according the experimental factor. However if this is not the case, this doesn't imply that there is no difference. This plot shows only the highest expressed genes, not the differentially expressed genes.

10 report files: **differentially expressed genes**, reads quality report, counting report,

...



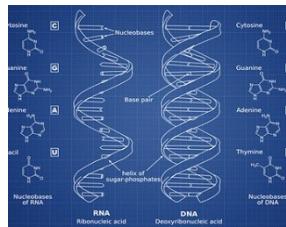
3GB of raw and analyzed data: Fastq, BAM, counts, normalized counts, gene expression files, heatmaps...

Results Interpretation & Analysis

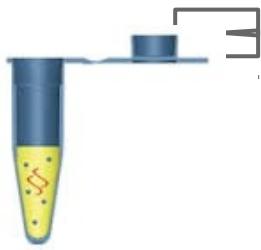


NGS Bioinformatics - This Workshop

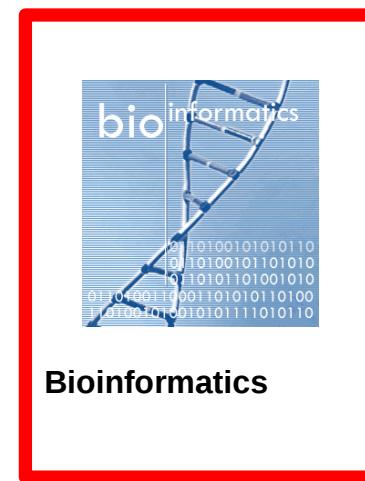
Experimental Design



Library Preparation

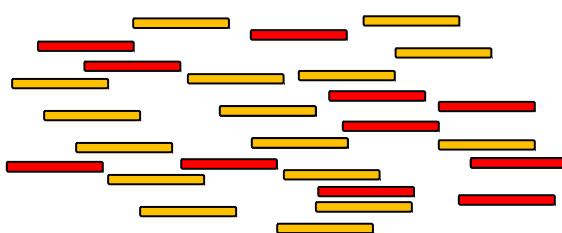


Sequencing



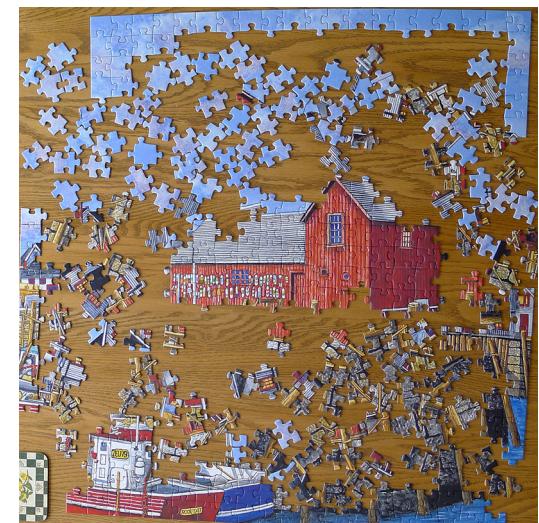
Follow up
and support

NGS Bioinformatics



NGS bioinformatics: interpretation and analysis
of NGS data using informatics tools

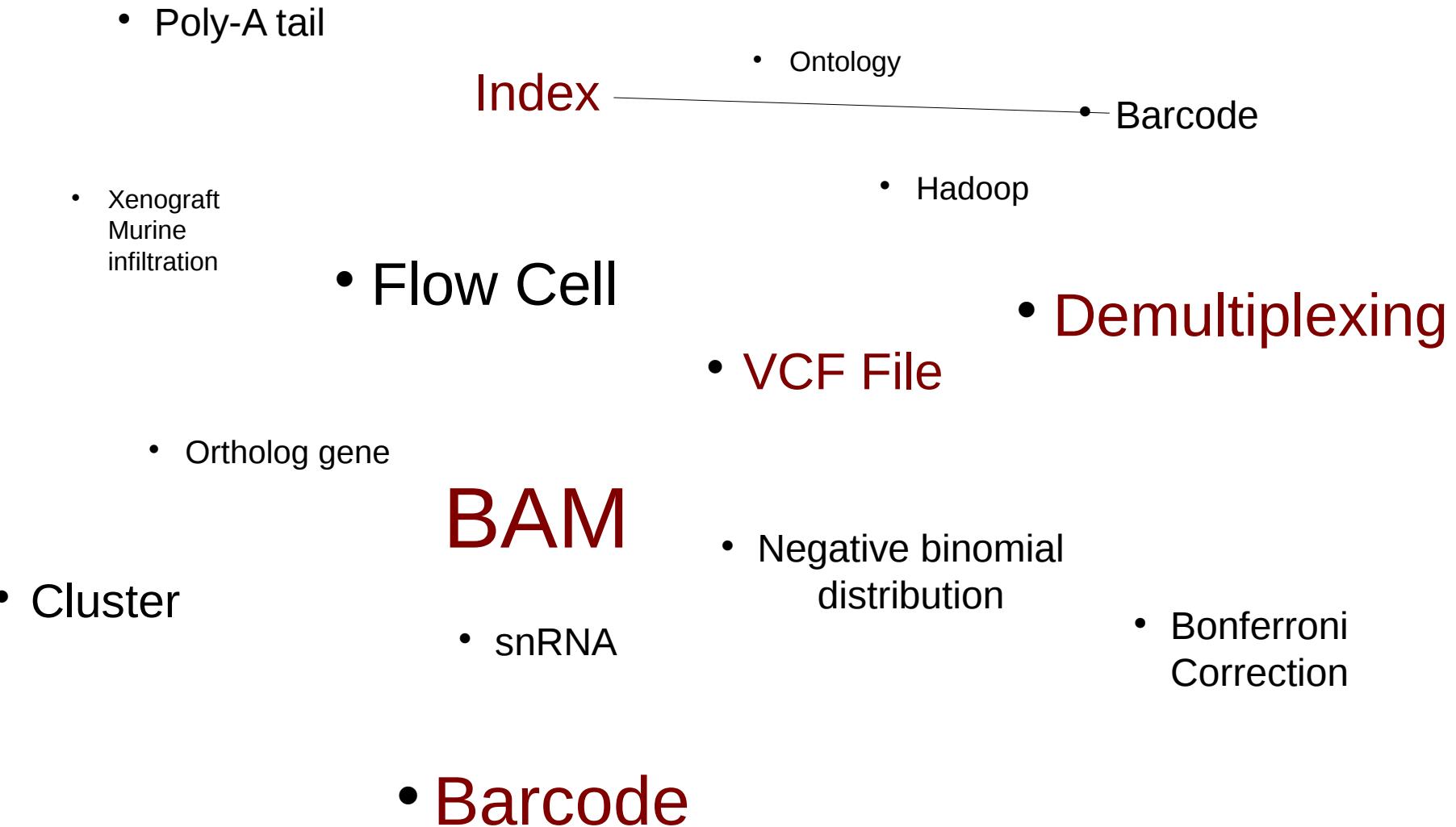
NGS Bioinformatics



Overview

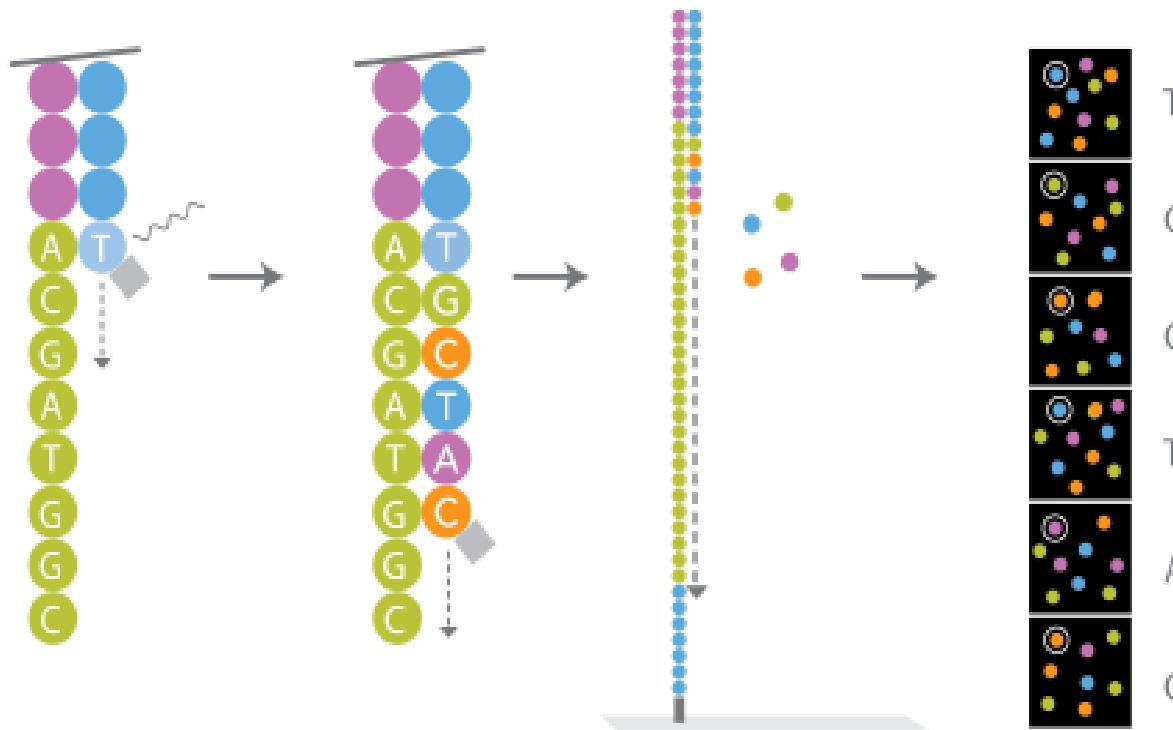
- This workshop: why?
- Bioinformatics NGS terminology
- NGS bioinformatics pipelines

A day in bioinformatics: Terminology



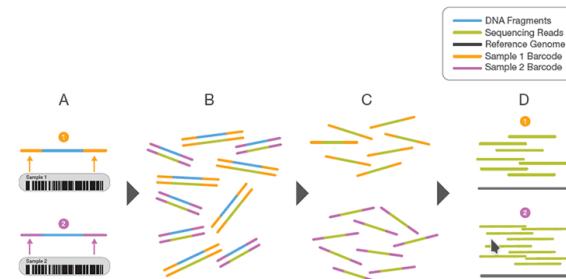
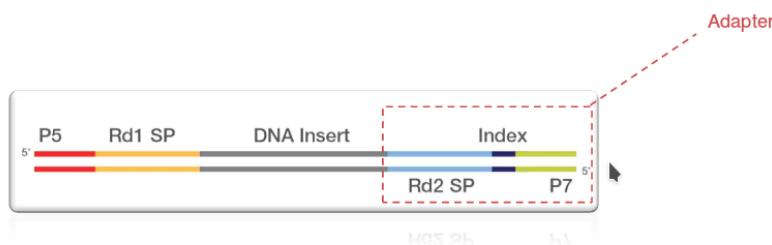
(De)Multiplexing (I)

Base Calling



(De)Multiplexing (II)

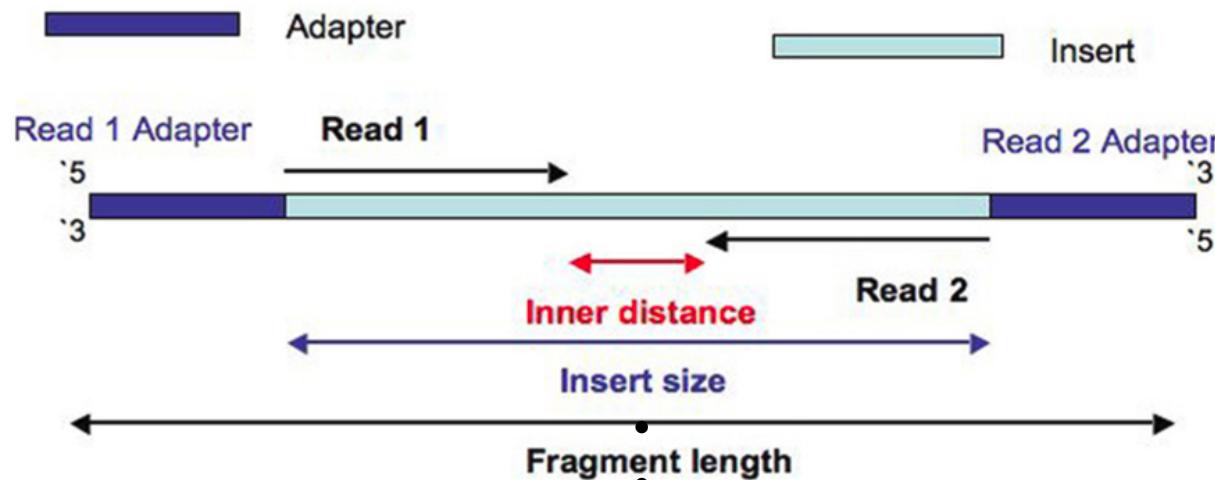
Multiple samples can be *pooled together or multiplexed* into one or more flowcells



From DNA Molecules to Digital Reads



Reads and Fragments



Fragment: biological entity
Read: bioinformatic concept

Fragment: the DNA template + adapters that were loaded on the sequencing machine (is not completely sequenced)

Read: a raw sequence originating from a sequencing machine

Single Read: Sequencing only from one end

Paired-end: Sequencing starting from both ends of the insert

Fastq format & Files

The result of demultiplexing is one (or two for paired-end reads) fastq files containing
raw reads

- Fastq files are **human readable** (not binary) text files.
- Referred as **raw data**.
- Fastq files are the **raw diamonds** of any NGS bioinformatics project
- Fastq files are often compressed using zip or gzip



```
-rwxr-Xr-x 1 vsc31439 lp_biogenomics 32G 17 janv. 15:05 GC036462.R1.fastq.gz
-rwxr-Xr-x 1 vsc31439 lp_biogenomics 39G 17 janv. 15:49 GC036462.R2.fastq.gz
-rwxr-Xr-x 1 vsc31439 lp_biogenomics 31G 17 janv. 16:25 GC036463.R1.fastq.gz
-rwxr-Xr-x 1 vsc31439 lp_biogenomics 37G 17 janv. 17:07 GC036463.R2.fastq.gz
-rwxr-Xr-x 1 vsc31439 lp_biogenomics 31G 17 janv. 17:41 GC036464.R1.fastq.gz
-rwxr-Xr-x 1 vsc31439 vsc31439 37G 17 janv. 18:21 GC036464.R2.fastq.gz
```



```
: vsc31420@hpc-p-login-1 /storage/leuven/stg_00019/full_genomes/test_ws 11:48 $ ls -lha GC036463.R1.fastq
-rwxr-Xr-x 1 vsc31420 vsc31420 134G 29 mars 11:24 GC036463.R1.fastq
```

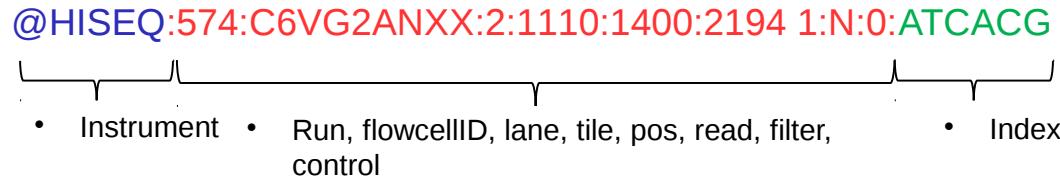
4x



Raw Reads

```
1 @HISEQ:574:C6VG2ANXX:2:1110:1400:2194 1:N:0:ATCACG
2 GGGGGATTCTCACTAGGTCTCAAGGTCTCTCACTCTCGTAGTGTCCCAG
3 +
4 CCCCCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG
```

- Line 1: Read identifier and is followed by a sequence that
 - Unique, platform dependent
 - Begins with a '@' character



- Line 2: Raw sequence of nucleotides: **read**
- Line 3: begins with a '+' character and is optionally followed by the same sequence identifier.
- Line 4: Quality values for the sequence in Line 2

Reads and Fastq Format

- **Fastq format?**

- Plain-text file, where each read and complementary information occupies 4 consecutive lines
- Typical size 500M compressed, 2000M unzipped

```
@HISEQ:574:C6VG2ANXX:2:1110:1400:2194 1:N:0:ATCACG  
GGGGGATTCTCACTAGGTCTCAAGGTCTCTCACTCTCGTAGTGTCCCCAG  
+  
CCCCCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG  
@HISEQ:574:C6VG2ANXX:2:1110:1560:2177 1:N:0:ATCACG  
ATGGTCCAGCAAGGGGTATGCTGAGAAGGGGAGCAGTTCAGAACCCATCAG  
+  
CCCCCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG  
@HISEQ:574:C6VG2ANXX:2:1110:1583:2223 1:N:0:ATCACG  
CTACCTTCACTATCAACATAGCAAACACACACCTTAGCTCCAGCTATTAACA  
+  
CCCCCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG  
@HISEQ:574:C6VG2ANXX:2:1110:1609:2245 1:N:0:ATCACG  
AGCTTAAGAGGCAGTACAGACACAGCCAGCTTCTCAGGTGATCCATGAACAC
```

- **Sequencing depth:** The total number of sequences generated for a sample.

- Usually expressed in fragments or reads
- How **deep** to sequence? Experiment!

Reference Sequence in Fasta Format

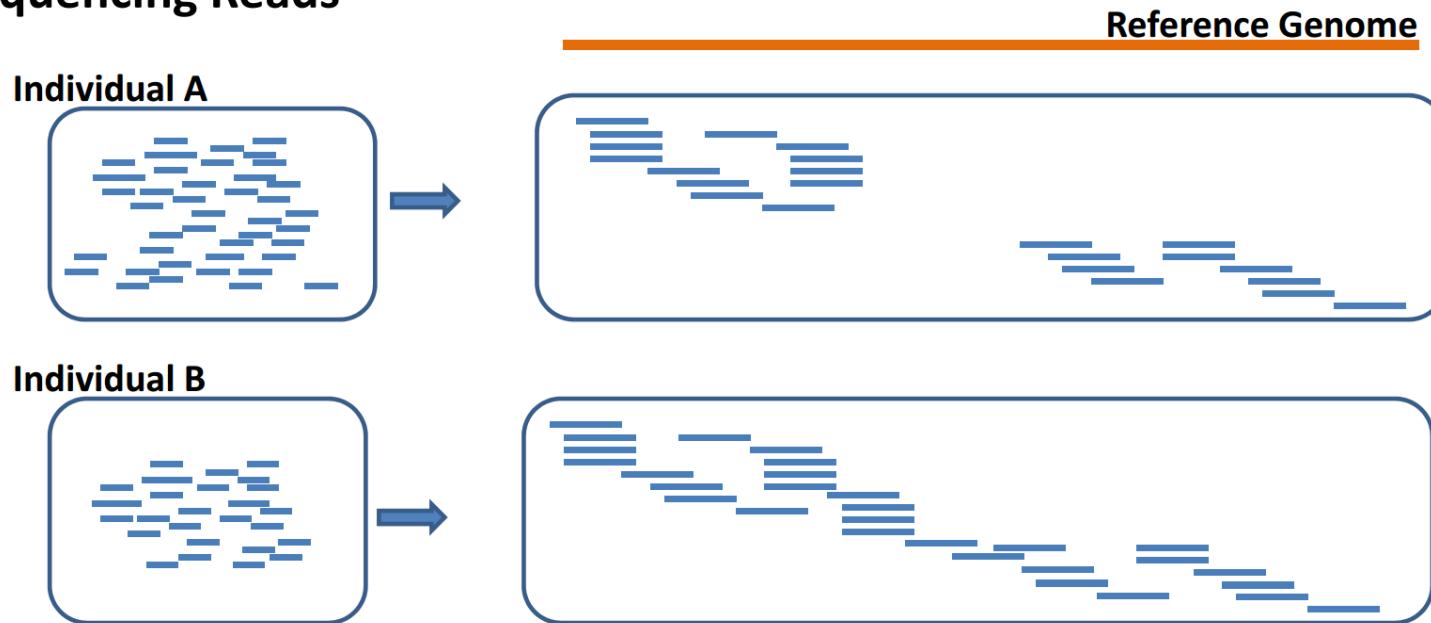
- genome.fa **human-readable** nucleotide sequence
- Species dependent
- Mouse genome: 2.6GB
- **Evolves**

AATAAGTCAATGGCCTTCTACACAAAGAATAAACAGGCTGAGAAAGAAATTAGGGAA
ACAAACACCCTCTCAATAGTCACAAATAATATAACATATCTCGGCGTGAECTAACTAAG
GAAGTGAAAGATCTGTATGATAAAAACCTCAAGTCTGAAGAAAGAAATTAAAGAAGAT
CTCAGAAGATGGAAAGATCTCCATGCTCATGGATTGGCAGGATCAATATTGAAAAATG
GCTATCTTGCCAAAAGCAATCTACAGATTCAATGCAATCCCATCAAATTCCAACCTCAA
TTCTTCAACGAATTAGAAGGGCAATTGCAAATTCTCTGTAATAACAAAAACCTAGG
ATAGCAAAAAGTCTTCTCAAGGATAAAAGAACCTCTGGTGAATCACCAGCCTGACCTA
AAGCTTACTACAGAGCAATTGGTAAAAACTGCATGGTACTGGTATAGAGACAGACAA
GTAGACCAATGGAATAGAATTGAAGACCCAGAAATGAACCCACACACCTATGGCACTTG
ATCTTCGACAAGGGAGCTAAACCATCCAGTGGAAAGAAAGACAGCATTTCACAAATGG
TGCTGGCACAACTGGTTGTATCATGTTAGAAGAATGCGAATCGATCCATACCTATCTCCT
TGTACTAAGGTCAAATCTAAATGGATCAAAGAACTTCACATAAAACCAGAGACACTGAAA
CTTATAGAGGAGAAAGTGGGAAAAGCCTGAAAGATATGGGCACAGGGAAAATTCTG
AACAGAACAGCAATGGCTGTGCTGTAAGATTGAGAATTGACAATGGACCTAATGAAA
CTCCAAAGTTCTGCAAGGCAAAAGACACCGTCATAAAGAGAAAGAGACACCACAGAT
TGGGAAAGGATCTTACCTATCTAAATCAGATAGGGACTAATATCCAACATATATAAA
GAACCTAAGAAGGTGGACTTCAGAAAATCAAACACCCATTAAAAATGGGCTCAGAA
CTGAACAAAGAATTCTCACCTGAGTTATACCGAATGGCAGAGAAGCACCTGAAAAATGC
TCAACATCCTTAATCATCAGGGAAATGCAAATCAAACACCCCTGAGATTCCACCTCACA
CCAGTCAGAATGTCTAAGATCAAATTCAAGGTGACAGCAGATGCTGGCGAGGATGTGGA
GAAAGAAGAACACTCCTCCATTGTTGGGGATTGCAAGGCTGTACAACCACTCTGGAAA
TCCGTCTGGGGTTCTCAGAAAATTGGACATAGTACTACCGGAGGATCCAGCAATACCT
CTCCTGGGATATATCCAGAAGATGCCCAACTGGTAAGAAGGACACATGCTCCACTATG
TTCATAGCAGCCTTATTATAATAGCCAGAAGCTGAAAGAACCCAGATGCCCTCAACA
GAGGAATGGATACAGAAAATGTGGTACATCTACACAATGGAGTACTACTCAGCTATTAAA
AAGAATGAATTATGAAATTCTAGCCAAATGGATGGACCTGGAGGGCATCATCCTGAGT

Mapping to Reference Genome

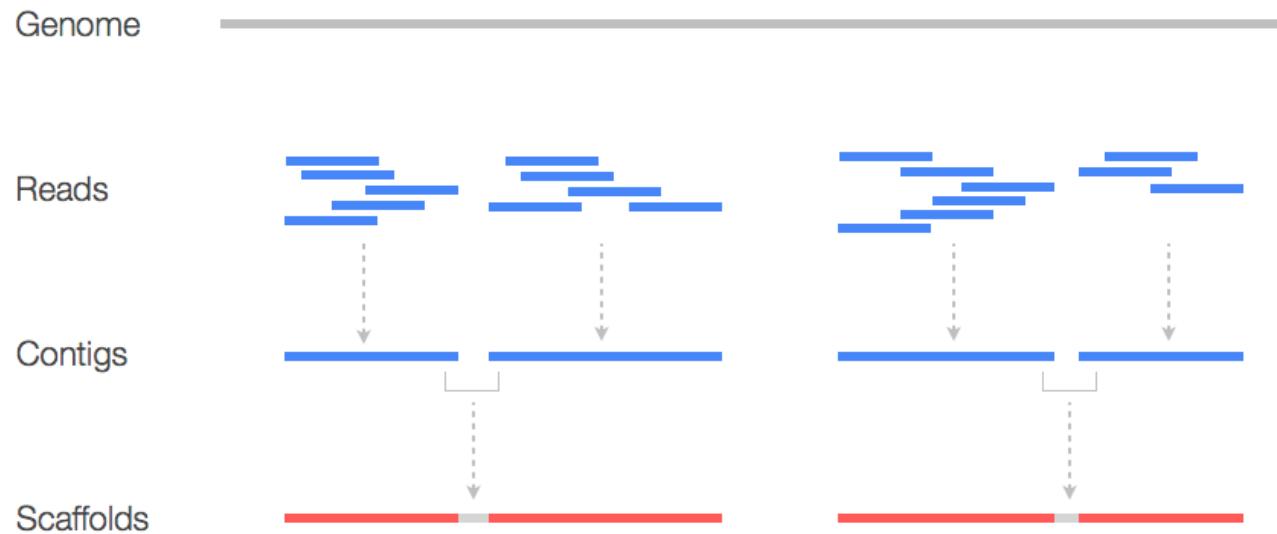
Mapping refers to the process of aligning short reads to a reference sequence

Sequencing Reads



Assembly

- The **generation of a reference**, from scratch (*de novo*) or reference assisted.
- Overlapping reads are merged to **contigs** (smallest unitable unit without unknown bases)
- Contigs that belong together, but where the connecting sequence is unknown, can be connected to **scaffolds**, inserting N's for the unknown bases



Computer Cluster

NGS data means big data...means big computing power

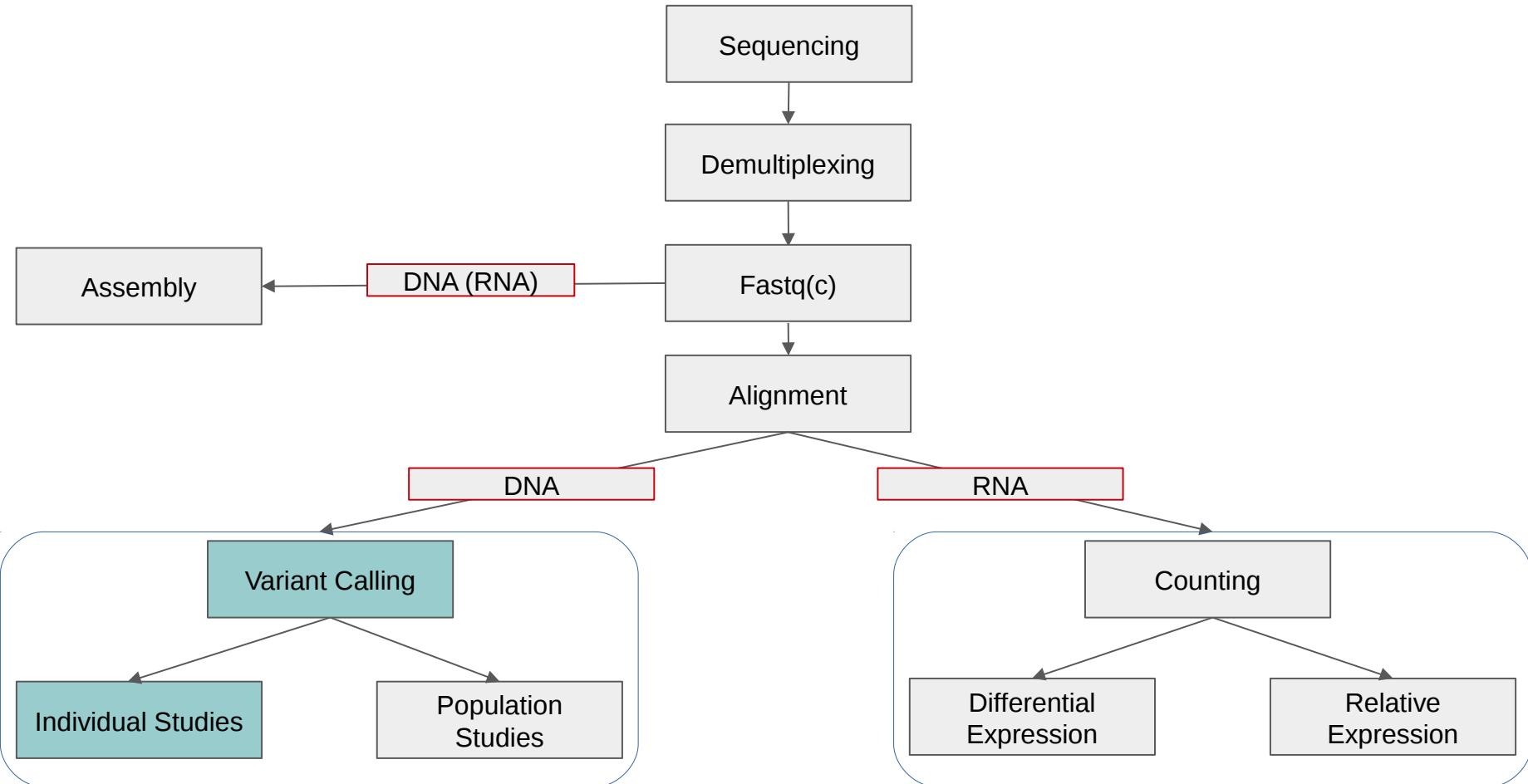


- Whole Human Genome:
- 300Gb, gzip file
- Exome data:
- 6GB
- RNA-Seq
- 1GB



- NGS data is usually analyzed on a **supercomputer or cluster**.
 - UZ Leuven: **Avalok/Hydra**
 - KU Leuven: **VSC Flemish Super Computer**
 - **Google genomics, etc.**

NGS Common Pipelines



Digging deeper

- Introduction to RNA-Seq pipelines and bioinformatics
 - Friday the 8th of November
- High performance computing for genomics
 - Hands-on Workshop (Speakers to determine)
 - Friday the 6th of December
 - Organized in conjunction with VSC, Lab. Computational Biology, GC
- Introduction to Single-Cell omics
 - January

Thanks!