# RNA-Seq data analysis

**17-18 octobre 2019**

**Céline Noirot et Matthias Zytnicki**

GENOTOUL Bioinfo · MIA TOULOUSE

---

# Material

- **Slides:**
  - pdf : one per page
    http://genoweb.toulouse.inra.fr/~formation/19_Rnaseq_Cli/doc/Rnaseq_training_012019.pdf
  - pdf : three per page with comment lines
    http://genoweb.toulouse.inra.fr/~formation/19_Rnaseq_Cli/doc/Rnaseq_training_012019_3p.pdf

- **Hands on:**
  - Exercises:
    http://genoweb.toulouse.inra.fr/~formation/19_Rnaseq_Cli/doc/Tps/RNAseq_TP_ligne_cmd_ennonce-Octobre2019.pdf
  - Data files: http://genoweb.toulouse.inra.fr/~formation/19_Rnaseq_Cli/data
  - Results files: http://genoweb.toulouse.inra.fr/~formation/19_Rnaseq_Cli/doc/Tps/Correction.txt

GENOTOUL Bioinfo · MIA TOULOUSE

---

# Session organisation

**Day 1**

**Morning (9h00 -12h30) :**
- Biological reminds
- Sequence quality
  Theory & exercises
- Spliced read mapping
  Theory & Exercises & Visualisation

**Afternoon (14h-17h) :**
- Expression quantification
  Theory + exercises

- mRNA calling
  Theory & exercises & Visualisation

**Day 2**

**Morning (9h00 -12h30) :**
- Models comparison
  Theory & exercises

- Hovering differential gene expression analyse
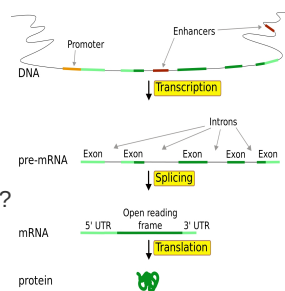
GENOTOUL Bioinfo · MIA TOULOUSE

## Summary – Biological reminds

✔ Transcriptome specificity
✔ High throughput sequencers
✔ Illumina protocol, paired-end library, directional library
✔ Experimental protocol
✔ RNAseq specific bias
✔ How to retrieve public data

## Context

Prerequis :
- Reference genome available
- RNAseq sequencing
  (sequence of transcript)

Try to answer to :
- How to map transcript to the genome ?
- How to discover new transcript ?
- What are the alternative transcript ?



Source : en.wikipedia.org/wiki/User:Forluvoft/sandbox

## Transcriptome variability

- Many types of transcripts (mRNA, ncRNA, cis-natural antisense, fusion gene ...)

- Many isoform (non canonical splice sites, intron retention …)

- Number of transcripts
  - possible variation factor between transcripts: $10^6$ or more,
  - expression variation between samples.

- Allele specific expression

2

## Transcriptome variability *(ENCODE)*

GENCODE    Data    Stats

Statistics about the current Human GENCODE Release (version 28)

\* The statistics derive from the **gtf file** @ that contains only the annotation of the main chromosomes

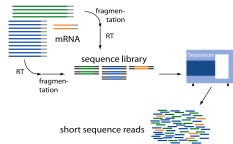For details about the calculation of these statistics please see the **README_stats.txt** @ file.

**Compare with the previous release (GENCODE 27) »**

Version 28 (November 2017 freeze, GRCh38) - Ensembl 92, 93

General stats

| | | | |
|---|---|---|---|
| Total No of Genes | 58381 | Total No of Transcripts | 203835 |
| Protein-coding genes | 19901 | Protein-coding transcripts | 82335 |
| Long non-coding RNA genes | 15779 | - full length protein-coding: | 56541 |
| Small non-coding RNA genes | 7569 | - partial length protein-coding: | 25794 |
| Pseudogenes | 14723 | Nonsense mediated decay transcripts | 14889 |
| - processed pseudogenes: | 10693 | Long non-coding RNA loci transcripts | 28468 |
| - unprocessed pseudogenes: | 3519 | | |
| - unitary pseudogenes: | 218 | | |
| - polymorphic pseudogenes: | 38 | | |
| - pseudogenes: | 18 | | |

https://www.gencodegenes.org/stats/current.html

7                                                                Bio & Quality
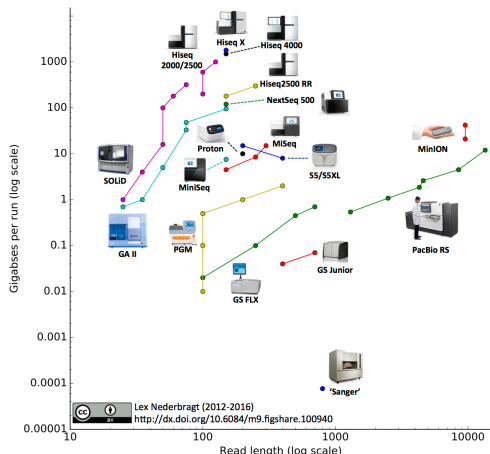
---

## What is « new » with RNA-Seq ?

- No prior knowledge of sequence needed
- Specificity of what is measured
- Increased dynamic range of measure, more sensitive detection
- Direct quantification
- Good reproducibility
- Different levels : genes, transcripts, allele specificity, structure variations
- New feature discovery: transcripts, isoforms, ncRNA, structures (fusion...)
- Possible detection of SNPs, ...

8                                                                Bio & Quality

---

## Sequencing platforms

https://flxlexblog.wordpress.com/2016/07/08/developments-in-high-throughput-sequencing-july-2016-edition/#more-790
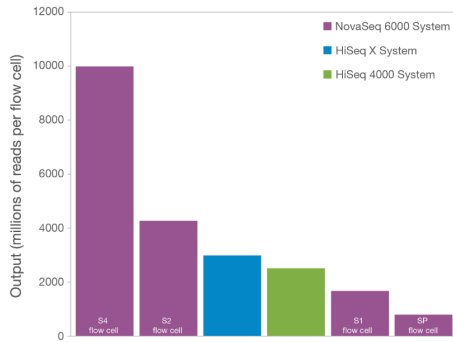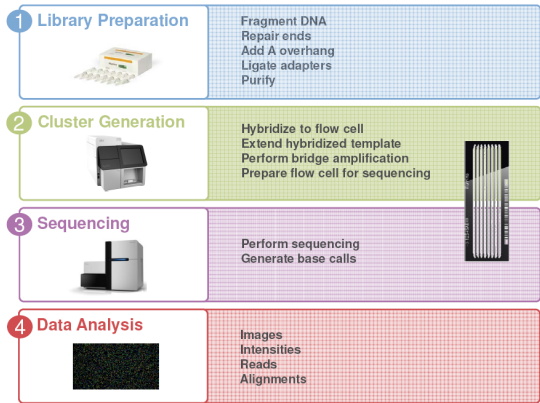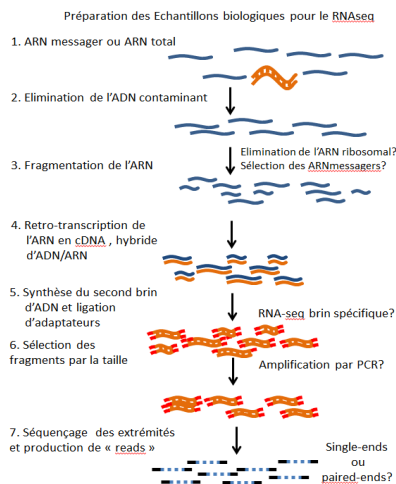
## Illumina Sequencing platforms



**Figure 2: The NovaSeq 6000 System offers the broadest output range—**The NovaSeq 6000 System generates from 80 Gb and 800 M reads to 3 Tb and 10 B reads of data in single flow cell mode. In dual flow cell mode, output can be up to 6 Tb and 20 B reads. The tunable output makes the NovaSeq 6000 System accessible for a wide range of applications.

https://www.illumina.com/content/dam/illumina-marketing/documents/products/datasheets/novaseq-6000-system-specification-sheet-770-2016-025.pdf
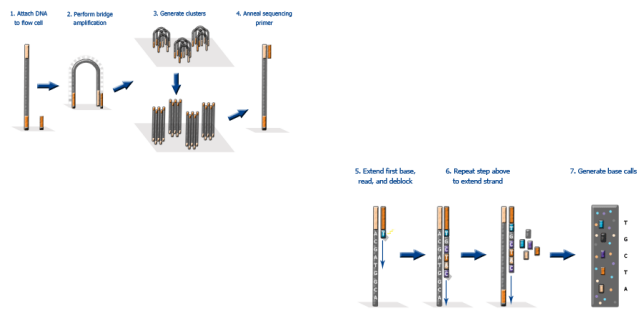
## Illumina RNA-Seq protocol



## RNA-Seq library preparation

## Clusters generation / Sequencing



https://www.eurofinsgenomics.co.in/en/eurofins-genomics/product-faqs/next-generation-sequencing/general-techni
cal-questions/what-is-the-principal-of-the-illumina-sequencing-technology.aspx

---

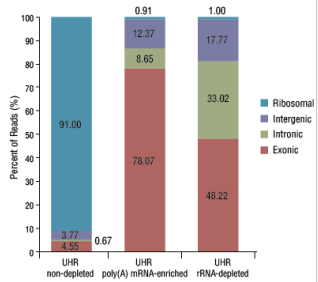## How to define experimental protocol ?

- Ribo-depletion or polyA-selection ?
- Single-end or paired-end ?
- How long should my reads be ?
- How many replicates ?
  - Technical or/and biological replicates ?
- How many reads for each sample?
- How many conditions for a full transcriptome ?

---

## Déplétion / Enrichissement ?

- Similar results

*Comparison of RNA-Seq by poly (A) capture, ribosomal RNA depletion, and DNA microarray for expression profiling, BMC Genomics , 2014*
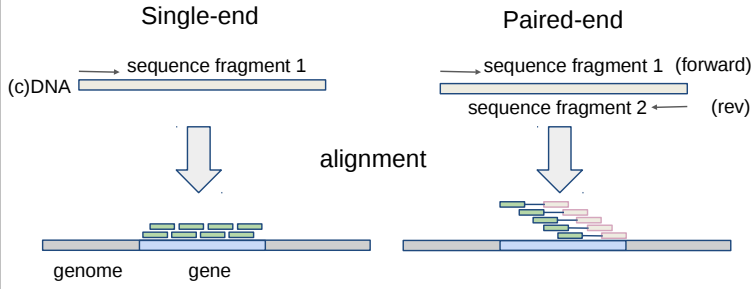
- RNA depletion:
  - For bacterial
  - ARN more varied
  - CircRNA
  - Some ncRNA

- polyA enrichment:
  - More reads into exons
  - Less biological material
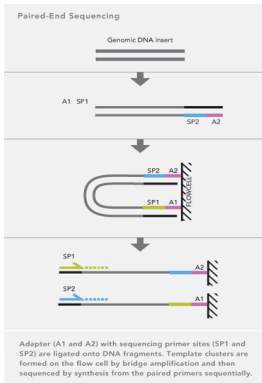  - No transcript without PolyA or partially degraded
  - No circRNA biais



https://content.neb.com/products/e6310-nebnext-rrna-depletion-kit-human-mouse-rat

## Paired-end VS single-end

Single-end

Paired-end

sequence fragment 1

(c)DNA

sequence fragment 1  (forward)

sequence fragment 2  ⟵  (rev)

alignment

genome      gene

- The cDNA size hive the insert size (ex. 200-500 pb).
- The fragment are usually forward-reverse.

16

Bio & Quality

---

## Paired-end sequencing

- Modification of the standard single-read DNA library preparation facilitates reading both ends of each fragment

- Improvement of mapping

- Help to detect structural variations in the genome like insertions or deletions, copy number variations, and genome rearrangements
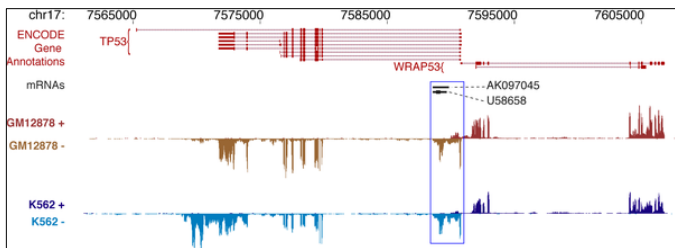
Paired-End Sequencing

Genomic DNA Insert

A1  SP1

SP2  A2

SP2  A2

SP1  A1

SP1

A2

SP2

A1

Adapter (A1 and A2) with sequencing primer sites (SP1 and SP2) are ligated onto DNA fragments. Template clusters are formed on the flow cell by bridge amplification and then sequenced by synthesis from the paired primers sequentially.

---

## Strand specific RNA-Seq protocol

Nat Methods. 2010 Sep;7(9):709-15. Epub 2010 Aug 15.

**Comprehensive comparative analysis of strand-specific RNA sequencing methods.**

Levin JZ, Yassour M, Adiconis X, Nusbaum C, Thompson DA, Friedman N, Gnirke A, Regev A.

Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts, USA.
jlevin@broadinstitute.org

**Abstract**

chr17:   7565000      7575000      7585000      7595000      7605000

ENCODE Gene Annotations   TP53

WRAP53

mRNAs         AK097045
              U58658

GM12878 +
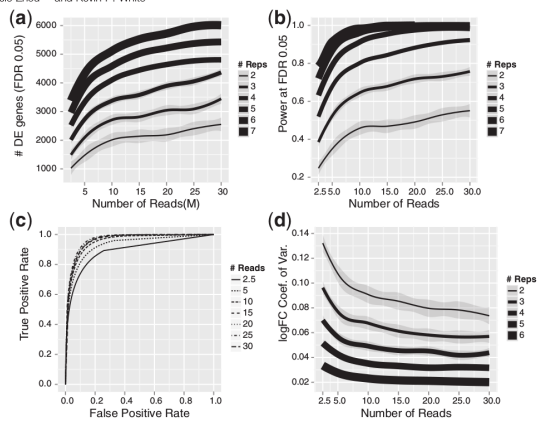GM12878 -

K562 +
K562 -

18

Bio & Quality

## Experimental protocol:
### Depth VS Replicates

- Encode (2016):
  - Experiments should be performed with **two or more biological replicates**, unless there is a compelling reason why this is impractical or wasteful
  - Replicate concordance: the gene level quantification should have a Spearman correlation of >0.9 between isogenic (same donor) replicates and >0.8 between anisogenic (different donor) replicates.
- Between **30M and 100M reads** per sample depending on the study.

https://www.encodeproject.org/documents/cede0cbe-d324-4ce7-ace4-f0c3eddf5972/@@download/attachment/ENC_ODE%20Best%20Practices%20for%20RNA_v2.pdf

## Experimental protocol:
### Depth VS Replicates

*Gene expression*     Advance Access publication December 6, 2013

**RNA-seq differential expression studies: more sequence or more replication?**

Yuwen Liu[1,2], Jie Zhou[1,3] and Kevin P. White[1,2,3,*]

## Retrieve public data

**Why ?**

- Because there is a lot of public data that would be sufficient for your analysis

- The authors often use only part of the data to answer their own problems

- Perhaps you don't need to sequence your own data

## Slide 22

**Retrieve public data**

ENA https://www.ebi.ac.uk/ena



22    Bio & Quality

## Slide 23

**Retrieve public data**

SRA https://www.ncbi.nlm.nih.gov/sra



23    Bio & Quality

## Slide 24

**Retrieve public data**



**Accession : SRX/ERX/DRX**

**SRPxxxxxx : Project**
**SRXxxxxxx : Experiement**
**SRRxxxxxx : Run**

**GSMxxxxxx : GEO id**

http://bioinfo.genotoul.fr/index.php/faq/bioinfo_tips_faq/

24    Bio & Quality

## Retrieve public data



## Retrieve public data



## Retrieve public data

- On genologin, use sratoolkit to :
  - download raw file
  - and convert format.

```
mkdir ~/work/ncbi
ln -s ~/work/ncbi ~/ncbi
module load bioinfo/sratoolkit.2.8.2-1
prefetch <sra_accession> --max-size
(20G by default)
Files are created into:
~/work/ncbi/public/sra/
Convertion
fastq-dump --gzip  sra_file.sra
```

## **Summary -** Sequence quality

- Known RNAseq biais

- How to check the quality ?

- How to clean the data ?

---

## **RNAseq specific bias**

• Influence of the library preparation

• Random hexamer priming

• Positional bias and sequence specificity bias.
  *Robert et al. Genome Biology, 2011,12:R22*

• Transcript length bias

• « Mappability »

29                                                          Bio & Quality

---

## **Hexamer random priming bias**

Préparation des Echantillons biologiques pour le RNAseq

1. ARN messager ou ARN total

2. Elimination de l'ADN contaminant

3. Fragmentation de l'ARN
   Elimination de l'ARN ribosomal?
   Sélection des ARNmessagers?

Random priming
→ not so random

4. Retro-transcription de l'ARN en cDNA , hybride d'ADN/ARN

5. Synthèse du second brin d'ADN et ligation d'adaptateurs
   RNA-seq brin spécifique?

6. Sélection des fragments par la taille
   Amplification par PCR?

7. Séquençage des extrémités et production de « reads »
   Single-ends ou paired-ends?

30                                                          Bio & Quality

10

# Hexamer random priming bias

**Biases in Illumina transcriptome sequencing caused by random hexamer priming**
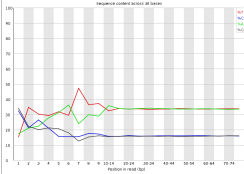
Kasper D. Hansen[1,*], Steven E. Brenner[2] and Sandrine Dudoit[1,3]

**ABSTRACT**

Generation of cDNA using random hexamer priming induces biases in the nucleotide composition at the beginning of transcriptome sequencing reads from the Illumina Genome Analyzer. The bias is independent of organism and laboratory and impacts the uniformity of the reads along the transcriptome. We provide a read count reweighting scheme, based on the nucleotide frequencies of the reads, that mitigates the impact of the bias.

−A strong distinctive pattern in the nucleotide frequencies of the first 13 positions at the 5'-end :
- sequence specificity of the polymerase
- due to the end repair performed

−Reads beginning with a hexamer over-represented in the hexamer distribution at the beginning relative to the end are down-weighted

---

# Transcript length bias

**Transcript length bias in RNA-seq data confounds systems biology.**
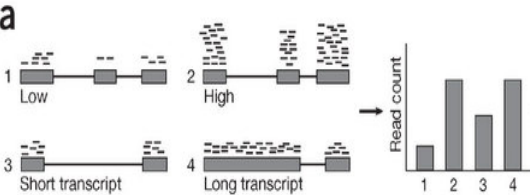
Oshlack A, Wakefield MJ.

**Abstract**
**Background:** Several recent transcriptome analysis (RNA genome transcriptional profil genomic sequences. As yet, a still in the stages of exploring
**Results:** We investigated the published data sets. For stand call differentially expressed g transcript.
**Conclusion:** Transcript leng current protocols for RNA-s expressed genes, and in parti other multi-gene systems biol
**Reviewers:** This article was Cloonan (nominated by Mark

− *the differential expression of longer transcripts is more likely to be identified than that of shorter transcripts*

**Length bias correction for RNA-seq data in gene set analyses**
Liyan Gao[1,†], Zhide Fang[2,†], Kui Zhang[1], Degui Zhi[1] and Xiangqin Cui[1,*]

---

# Bias "mappability"

- Quality of the reference genome influence results
  - assembly
  - finishing

- Sequence composition
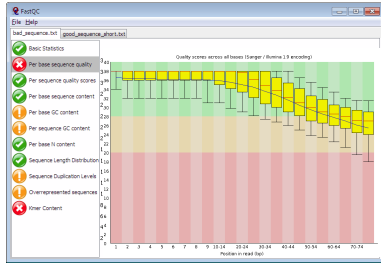- Repeated sequences
- Annotation quality

## Verifying RNA-Seq quality

**FastQC** :
*http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/*



*Has been developed for genomic data*

---

## fastq format

- Standard for storing outputs of HTS
- A text-based format for storing a read and its corresponding quality scores
- 1 read  <-> 4 lines

```
@HWI-ST218:596:C90JYANXX:8:1101:1293:2188 1:N:0:ATTCAGAATAATCTTA
NCTAAGTGTTAGGGGGTTTCCGCCCCTTAGTGCTGCAGCTAACGCATTAAGCACTCCGCCTGGGGAGTACGGTCGCAAGACTGAAAA
+
#<3?BFGGGGGGEGGGGGGGEGGGGGG@F1FGGGGGGDDGG1FB</9FE=EGGGGGGGG>GGGGBGGGGG<<C/BDGGGGGGC=GGG
```

1. Begins with '@' character and is followed by a sequence identifier
2. The raw sequence
3. Begins with a  '+' character and  is optionally followed by the same sequence identifier
4. Encodes the quality values for th read , contains the same number of symbols as letters in the read

---

## fastq format

- Sequence identifier

**@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG**

1. Begins with '@' character and is followed by a sequence identifier

| | |
|---|---|
| **EAS139** | the unique instrument name |
| **136** | the run id |
| **FC706VJ** | the flowcell id |
| **2** | flowcell lane |
| **2104** | tile number within the flowcell lane |
| **15343** | 'x'-coordinate of the cluster within the tile |
| **197393** | 'y'-coordinate of the cluster within the tile |
| **1** | the member of a pair, 1 or 2 *(paired-end or mate-pair reads only)* |
| **Y** | Y if the read is filtered, N otherwise |
| **18** | 0 when none of the control bits are on, otherwise it is an even number |
| **ATCACG** | index sequence |

## fastq format

- Base quality (Sanger standard)

```
@HWI-ST218:596:C90JYANXX:8:1101:1293:2188 1:N:0:ATTCAGAATAATCTTA
NCTAAGTGTTAGGGGGTTTCCGCCCCTTAGTGCTGCAGCTAACGCATTAAGCACTCCGCCTGGGGAGTACGGTCGCAAGACTGAAAA
+
#<3?BFGGGGGGEGGGGGGGEGGGGGG@F1FGGGGGGDDGG1FB</9FE=EGGGGGGGGG>GGGGBGGGGG<<C/BDGGGGGGC=GGG
```

ASCII-encoded version of the PHRED quality given by $Q_{\mathrm{PHRED}} = -10 \times \log_{10}(P_e)$

SANGER=PHRED+33 : H=ASCII(40+33)  $Q = -10 \log_{10} P \Leftrightarrow P = 10^{\frac{-Q}{10}}$

| Score de qualité phred | Probabilité d'une identification incorrecte | Précision de l'identification d'une base |
|---|---|---|
| 10 | 1 pour 10 | 90 % |
| 20 | 1 pour 100 | 99 % |
| 30 | 1 pour 1000 | 99.9 % |
| 40 | 1 pour 10000 | 99.99 % |
| 50 | 1 pour 100000 | 99.999 % |

---

## fastq format



---

## fastqQC Report

### Summary

- Basic Statistics
- Per base sequence quality
- Per sequence quality scores
- Per base sequence content
- Per base GC content
- Per sequence GC content
- Per base N content
- Sequence Length Distribution
- Sequence Duplication Levels
- Overrepresented sequences
- Kmer Content

The analysis in FastQC is performed by a series of analysis modules.

Quick evaluation of whether the results of the module seem :
- entirely normal (green tick),
- slightly abnormal (orange triangle)
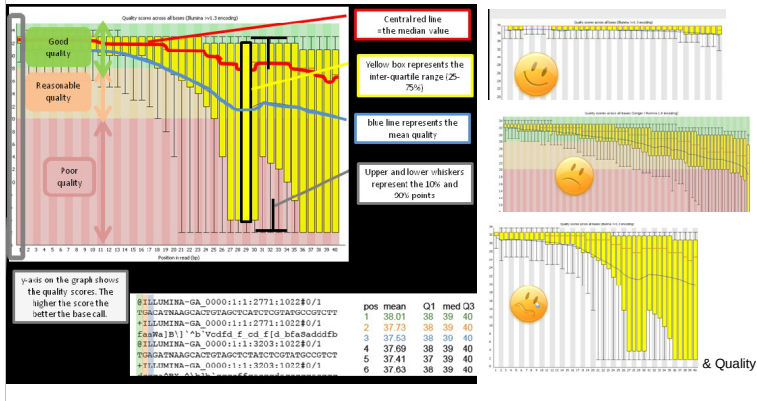- or very unusual (red cross).

These evaluations must be taken in the context of what you expect from your library. A 'normal' sample as far as FastQC is concerned is random and diverse.

## fastqQC Report

**Statistics per Base Sequence Quality**

This view shows an overview of the range of quality values across all bases at each position in the FastQ file.
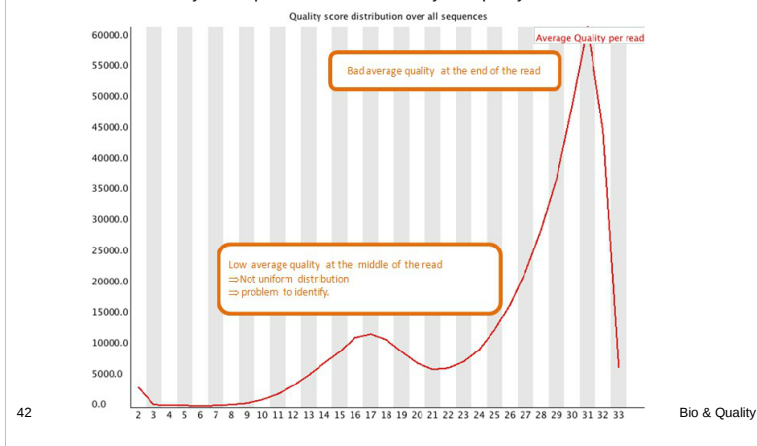
Common to see base calls falling into the orange area towards the end of a read.



## fastqQC Report

**Statistics per Sequence Quality Score**

See if a subset of your sequences have universally low quality values.
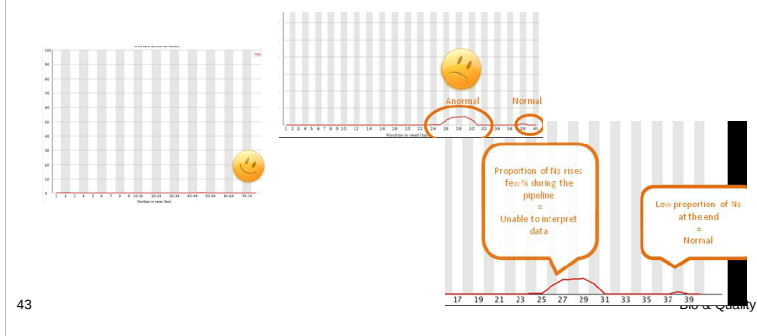


Bio & Quality

## fastqQC Report

**Statistics per Base N Content**

This module plots out the percentage of base calls at each position for which an N was called.

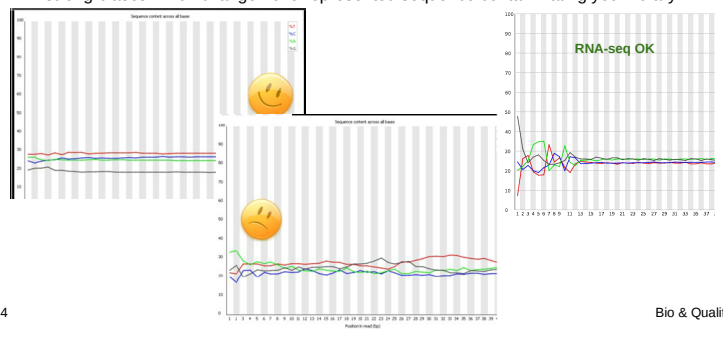Usual to see a very low proportion of Ns appearing nearer the end of a sequence.



Bio & Quality

## fastqQC Report

**Statistics  Per Base Sequence Content**

Per Base Sequence Content plots out the proportion of each base position in a file for which each of the four normal DNA bases has been called.

In a random library : little/no difference between the different bases of a sequence run, so the lines in this plot should run parallel with each other.

If strong biases which change : overrepresented sequence contaminating your library.
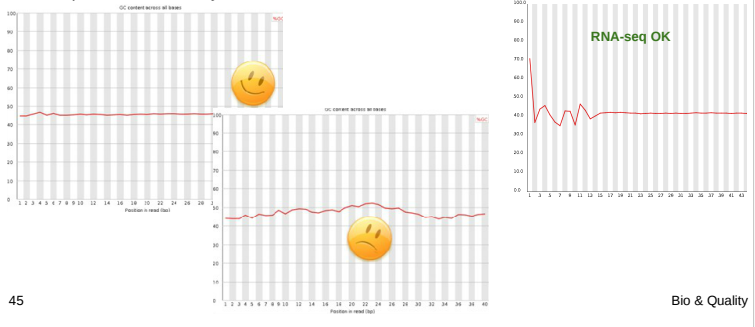


**RNA-seq OK**

44

## fastqQC Report

**Statistics per Base GC Distribution**

Per Base GC Content plots out the GC content of each base position in a file.

Random library : little/no difference between the different bases of a sequence run
=> plot horizontally.
The overall GC content should reflect the GC content of the underlying genome.

GC bias: changes in different bases, overrepresented sequence contaminating your library.
=> plot not horizontally.
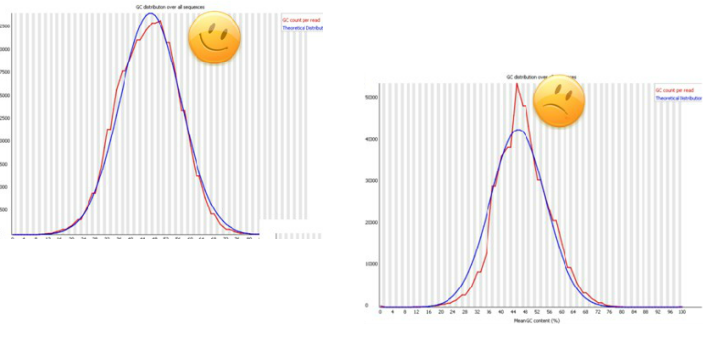


**RNA-seq OK**

45

## fastqQC Report

**Statistics per Sequence GC Content**

This module measures the GC content across the whole length of each sequence in a file and compares it to a modeled normal distribution of GC content.
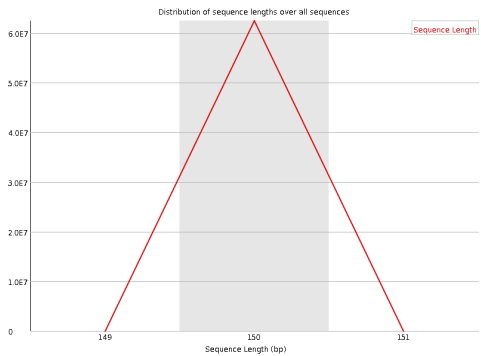


46

15

Bio & Quality

## fastqQC Report

**Statistics per Sequence Length Distribution**
Some sequence fragments contain reads of wildly varying lengths.

Even within uniform length libraries some pipelines will trim sequences to remove poor quality base calls from the end.
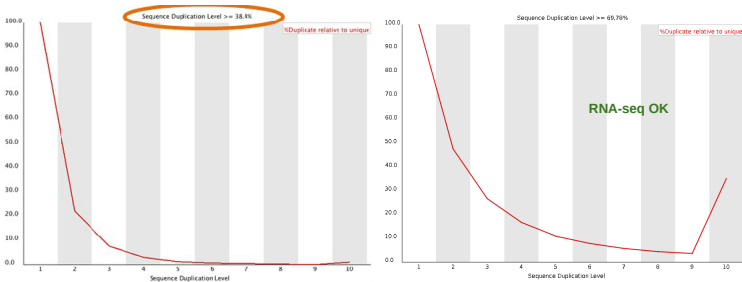


47     Bio & Quality

## fastqQC Report

**Statistics per Duplicate Sequences**

High level of duplication indicate an enrichment biais.
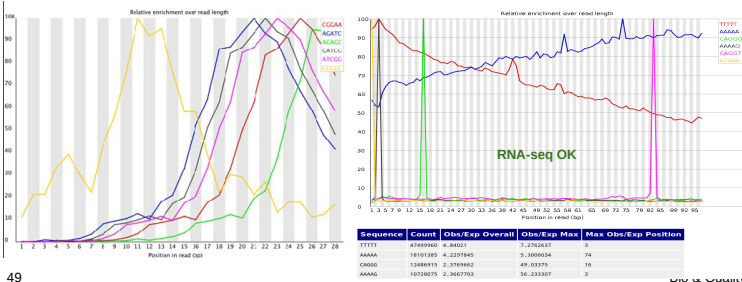


48     Bio & Quality

## fastqQC Report

**Overrepresented Kmers**
- A kmer is a subsequence of length k
- Should spot overrepresented sequences, give a good impression of any contamination.
- Kmers showing a rise towards the end of the library indicate progressive contamination with adapters.
- Check for adaptor sequence or poly-A sequence



49     Bio & Quality

16

Bio & Quality

## Take home message on quality analysis

Elements to be checked :
- Random priming effect
- K-mer (polyA, polyT)
- Adaptor presence

Alignment on reference for the second quality check and filtering.

A good run?:
- Expected number of reads produced,
- Length of the reads expected (150pb),
- Random selection of the nucleotides and the GC%,
- Good alignment: very few unmapped reads, pairs mapped on opposite strands.

50                                                          Bio & Quality

## Cleaning analysis

- Cleaning :
  - Low quality bases
  - Adaptors

- Software :
  - Trim_galore
  - Cutadapt
  - Trimmomatic
  - Sickle
  - PRINSEQ
  - ...

51                                                          Bio & Quality
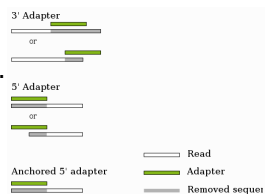
## Cutadapt

- Searches & removes adapter & tag in all reads.
- Trim quality
- Filter too short or untrimmed reads (in a separate output file).

```
module load bioinfo/cutadapt-1.8.3-python-2.7.2
cutadapt -a ADAPTER_FWD -A ADAPTER_REV -o out1.fatsq -p
out2.fastq  reads1.fastq reads2.fastq
```

Ex.: cutadapt -a AACCGGTT -o output.fastq input.fastq
(3' adapter, single read)
Input file : fasta, fastq or compressed (gz, bz2, xz).

Source : http://cutadapt.readthedocs.io/en/stable/guide.html

53                                                          Bio & Quality

17

## trim_galore

- Detect automatically adaptor
- Trim adaptor
- Trim low quality bases
- Trim N bases
- Remove read with length lower than 20b

```
module load bioinfo/cutadapt-1.14-python-2.7.2
module load bioinfo/FastQC_v0.11.7
module load bioinfo/TrimGalore-0.4.5
mkdir DIR
trim_galore   --fastqc
              --stringency 3
              --length 25
              --trim-n
              -o DIR
              --paired <read1> <read2>
```

## Hands-on: quality control

**Data for the exercises:**

- from Mohammed Zouine (ENSAT)
- tomato wild type and mutant type (without seeds) with the transcription factor Sl-ARF8 (auxine response factor 8) overexpressed
- clonal lineage
- paired, 100 pb non stranded
- triplicated
- in the publication process
- subsampled on chromosome 6 for faster analysis

*Use FastQC and trim_galore*

   *Exercise 1 :  quality control of used datasets*
                *cleaning used datasets*