# RNA-seq data analysis workshop for biologists, part I

## from raw data to read counts

**Genome**

Database of all the DNA
of the organism

new transcripts
gene isoforms

**Transcriptome**

Database of all known
transcripts for the organism

more accurate quantification

1. **Genome or transcriptome?**
2. Where can I find the reference?
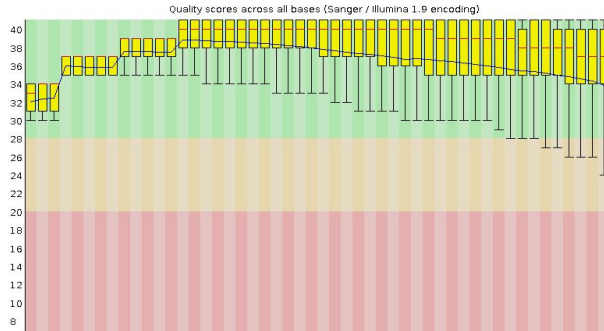
# What are the steps of an RNA-Seq analysis?



1. **Quality control**
2. Alignment or mapping
3. Count reads
4. Differential expression

# What are the steps of an RNA-Seq analysis?

RNA-Seq reads



Align reads to
genome

Genome

B. Haas, M. Zody. Advancing RNA-Seq analysis. Nature biotechnology. 28:5. 2010.
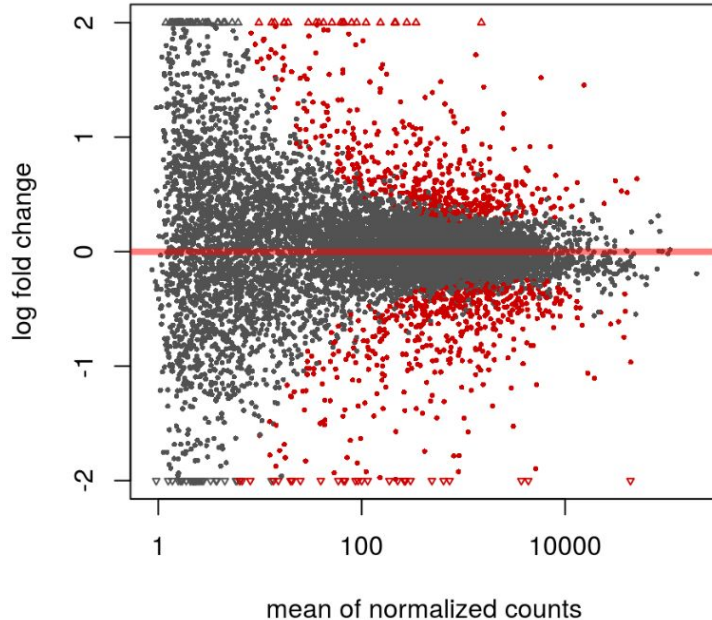
1. Quality control
2. **Alignment or mapping**
3. Count reads
4. Differential expression

What are the steps of an RNA-Seq analysis?

| Gene | Sample1 | Sample2 | Sample3 |
|---|---|---|---|
| ENSDART00000151582 | 462 | 4 | 454 |
| ENSDART00000146024 | 31 | 5408 | 41 |
| ENSDART00000052082 | 353 | 42 | 4 |
| ENSDART00000183148 | 6 | 702 | 56 |
| ENSDART00000077539 | 1246 | 42 | 12 |
| ENSDART00000178294 | 8 | 116 | 600 |
| ENSDART00000190290 | 185 | 468 | 691 |
| ENSDART00000129730 | 374 | 733 | 348 |
| ENSDART00000030215 | 825 | 25 | 520 |

1. Quality control
2. Alignment or mapping
3. **Reads counting**
4. Differential expression

# What are the steps of an RNA-Seq analysis?

**Will be covered in second part of the workshops**

1. Quality control
2. Alignment or mapping
3. Count reads
4. **Differential expression**

Tools:

**FastQC**

**MultiQC**

What statistics we are interested in?

Basic Statistics

Per base sequence quality

Per tile sequence quality

Per sequence quality scores

Per base sequence content

Per sequence GC content

Per base N content

Sequence Length Distribution

Sequence Duplication Levels
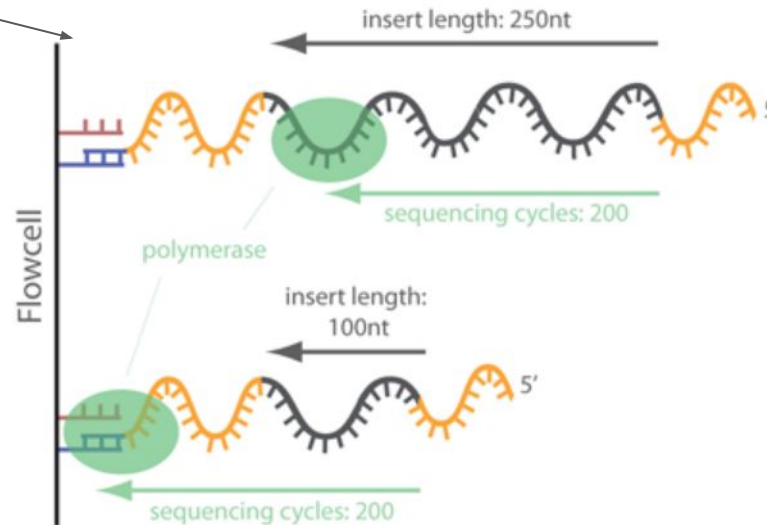
Overrepresented sequences

Adapter Content

1. **Quality visualization**
2. Reads filtering / trimming
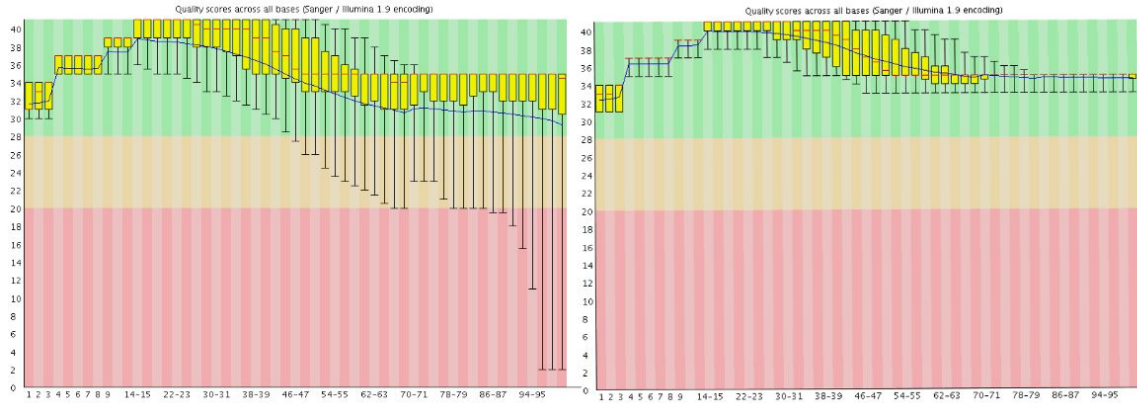
Quality control

Why we filter or trimm our precious data:
- adapters
- low read quality
- rRNA
- mtDNA

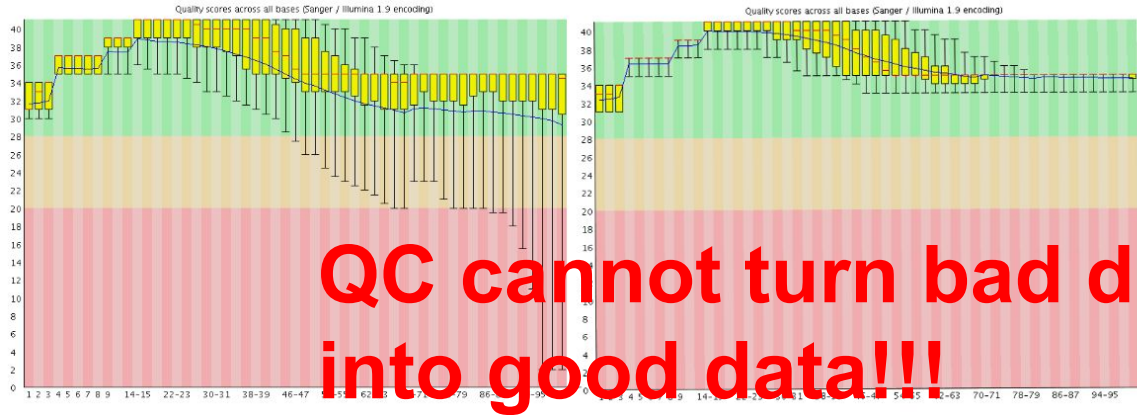1. Quality visualization
2. **Reads filtering / trimming**



http://www.ecseq.com/support/ngs/trimming-adapter-sequences-is-it-necessary

# Quality control

## Quality trimming



1. Quality visualization
2. **Reads filtering / trimming**

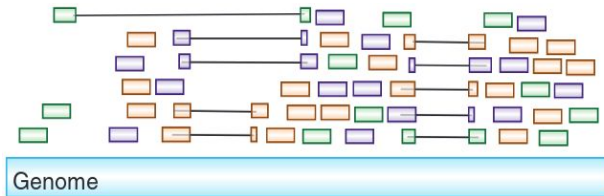## Adapter trimming
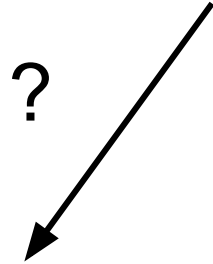
# Quality control

## Quality trimming



1. Quality visualization
2. **Reads filtering / trimming**

**QC cannot turn bad data into good data!!!**

## Adapter trimming

Alignments or mapping

RNA-Seq reads

**1.** **How does it work?**
2. What are the options?

?

In RNA-seq we are interested in quantification

Genome

**Alignment methods**

- STAR
- HiSat2
- BWA
- BBMap
- Subjunc

**Mapping methods**

- Salmon
- Kalisto

1. How does it work?
2. **What are the options?**

# Reference



ENSEMBL



UCSC Genome Browser

1. Genome or transcriptome?
2. **Where can I find the reference?**