# RNA-Seq data analysis

Bioinformatics Group

2 May 2017

مدينة الملك عبدالعزيز
للعلوم والتقنية KACST

# Materials and Software

- Bowtie2 (http://bowtie-bio.sourceforge.net/bowtie2/index.shtml).
- Samtools (http://samtools.sourceforge.net).
- Tophat2 (http://ccb.jhu.edu/software/tophat/index.shtml).
- STAR (https://github.com/alexdobin/STAR).
- Cufflinks (http://cole-trapnell-lab.github.io/cufflinks/).
- Rstudio (https://www.rstudio.com).
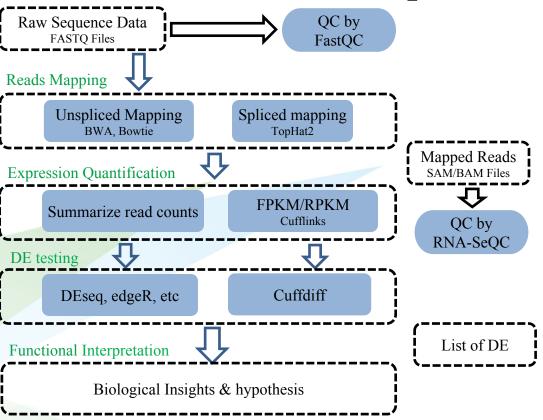- CummeRbund package.
- DESeq2 package.

# *Benefits & Challenges*

**<u>Benefits:</u>**

- Independence on prior knowledge

- High resolution, sensitivity and large dynamic range

**<u>Challenge:</u>**

- Interpretation is not straightforward

- Procedures continue to evolve

# From reads to differential expression

# FASTQ file

**Line1:** Sequence identifier
**Line2:** Raw sequence
**Line3:** meaningless
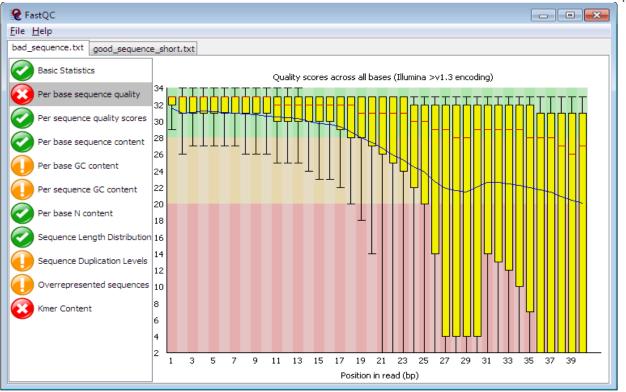**Line4:** quality values for the sequence

# Sequencing QC

## Information we need to check

- Basic information( total reads, sequence length, etc.)
- Per base sequence quality
- Overrepresented sequences
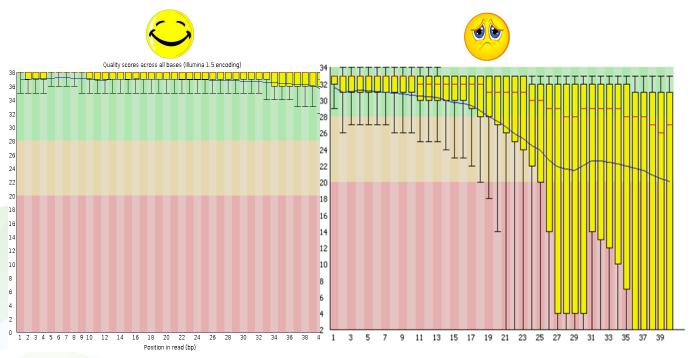- GC content
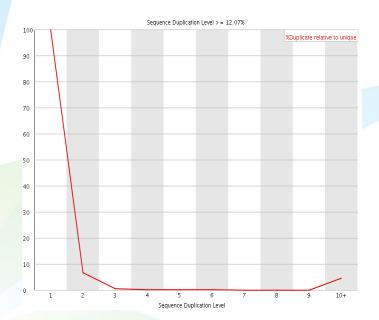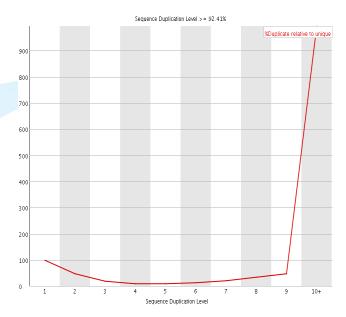- Duplication level
- Etc.

# FastQC

# Per base sequence quality

# Duplication level
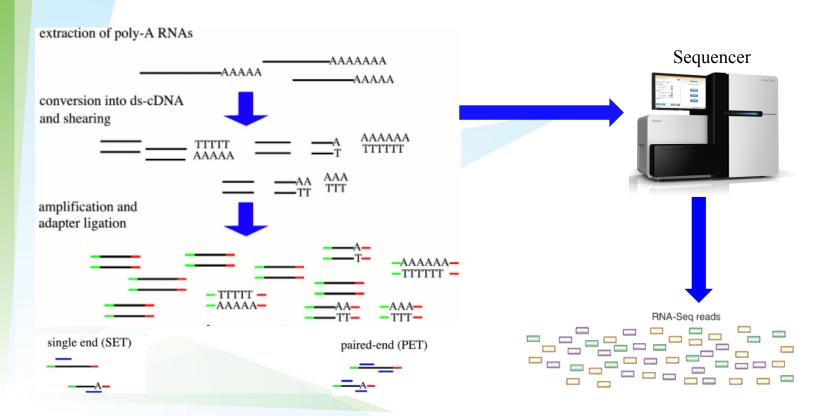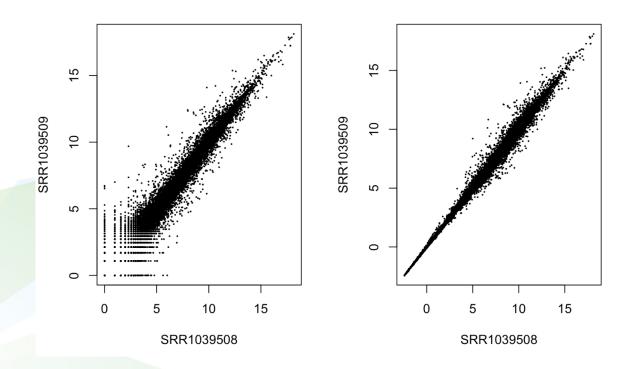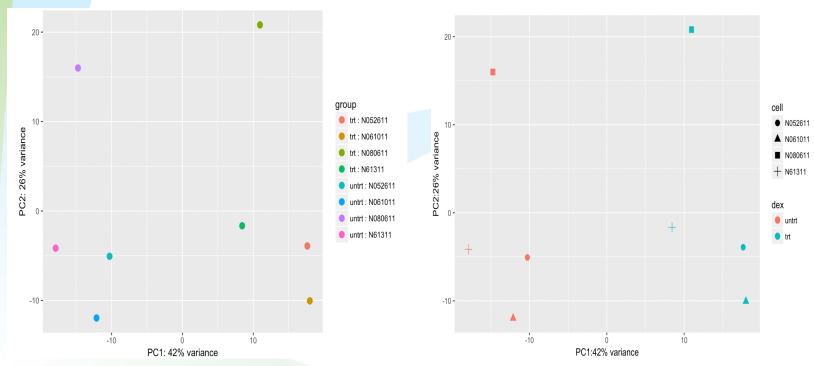
# Overview of RNA-Seq
## Transcriptome profiling using NGS

# Before and After Normalization
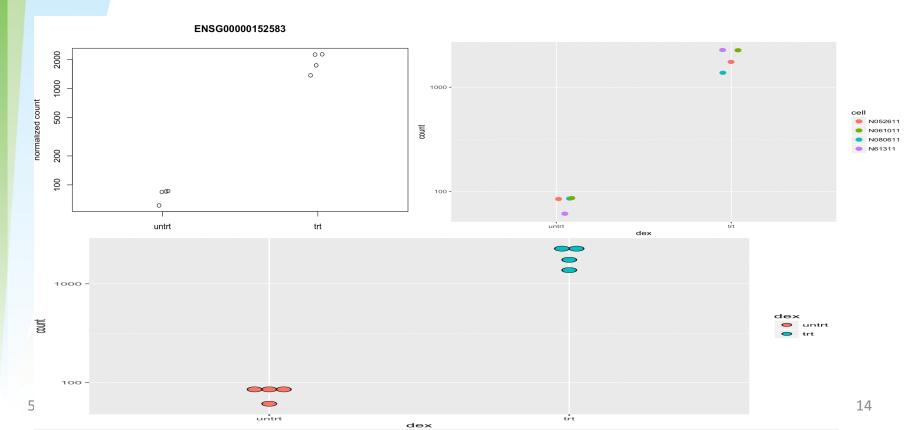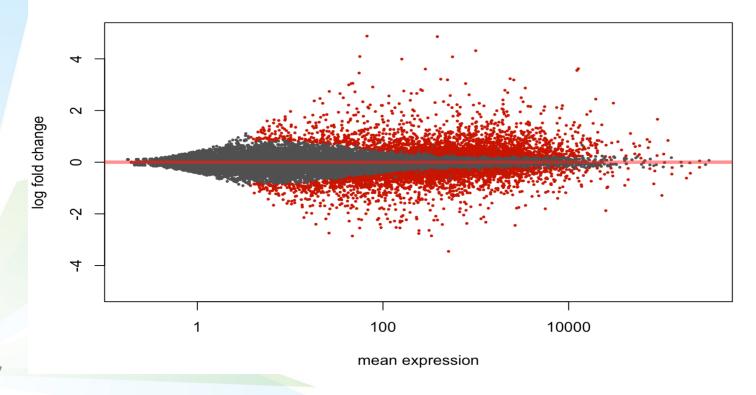
# Principal component analysis (PCA) Plots

# Principal component analysis (PCA) Plots

# MP-Plot