MSIB32500 Advanced Bioinformatics Spring 2019

# RNAseq Data Analysis and Clinical Applications, Part II

Riyue Bao, Ph.D.

Research Assistant Professor (Bioinformatics)

Associate Director of Cancer Immunology (Bioinformatics)

Center for Research Informatics & Department of Pediatrics

The University of Chicago

# Outline

- Part I (05/25/2019)
  - Introduction to RNAseq technology and clinical applications
  - Hands on: From raw data to gene expression quantification

- Part II (06/01/2019)
  - Differential gene expression analysis and data visualization
  - Hands on: Identification of genes and pathways significantly changed under condition
  - ***Homework assignment (DUE 06/08/2019 11:59 PM; and optional tasks)***

- Part III (06/08/2019)
  - How to associate gene expression data with clinical outcome
  - Hands on: Use gene expression data to discover tumor subtypes and survival analysis

# Class materials

- GitHub
  - https://github.com/MScBiomedicalInformatics/MSIB32500
  - This lecture note contains the same contents as the notebook. In addition, the notebook also contains hands-on materials
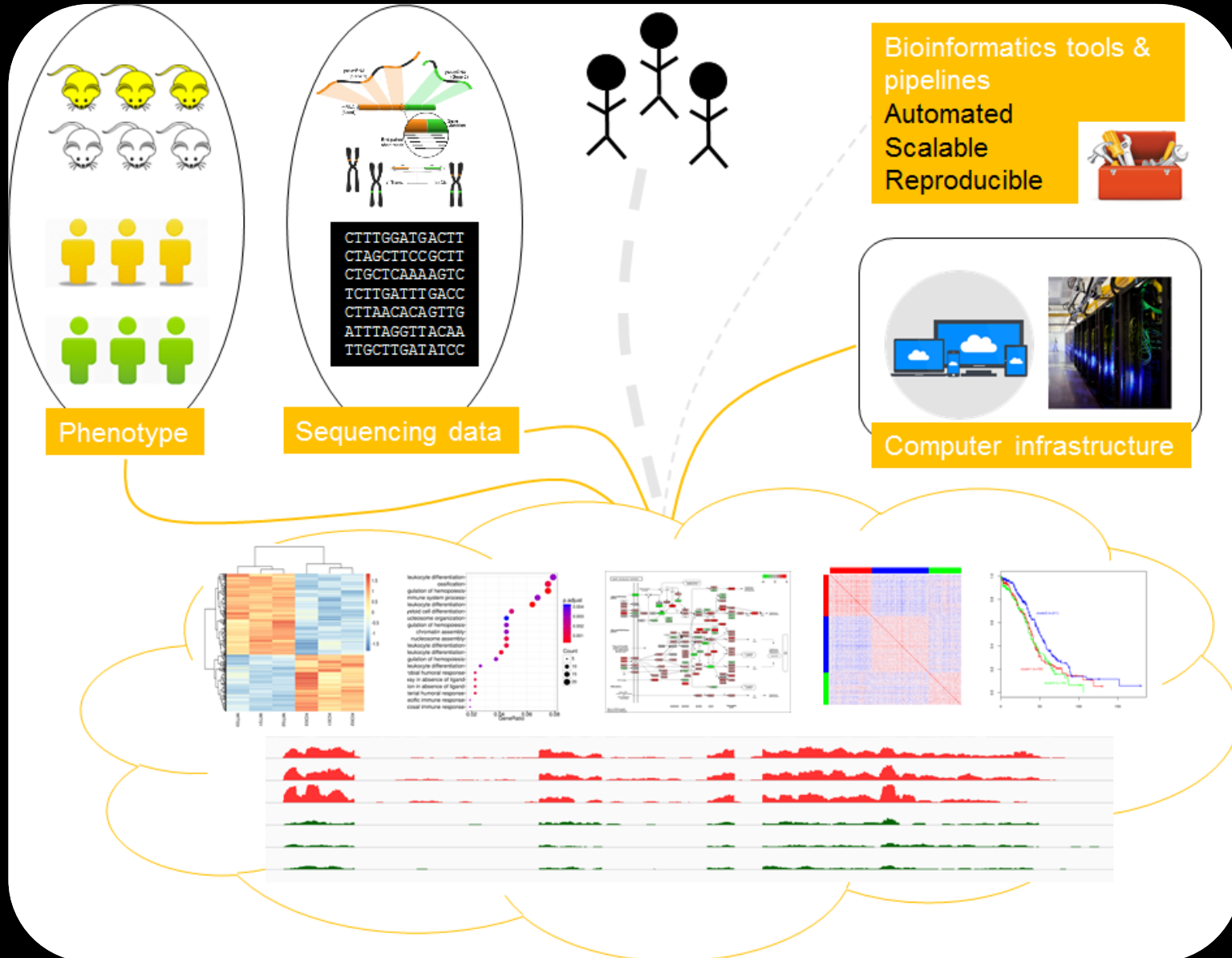  - **lecture9.pdf**
  - **Handson9.Rmd**

- Rstudio (or R console) on personal computers (hands on practice)

# Objective

- *(Recap from last class)*

- Detect genes differentially expressed between conditions

- Identify pathways / network enriched in genes of interest

- Generate high-quality figures for publication (PCA, heatmap, sample/gene cluster, GO/pathways, etc.)

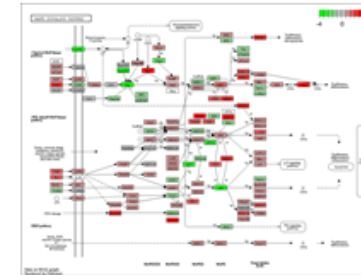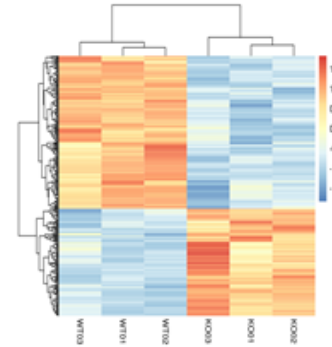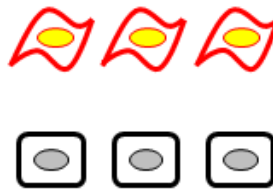- Become familiar with running commands in R / Rstudio

Phenotype

Sequencing data

CTTTGGATGACTT
CTAGCTTCCGCTT
CTGCTCAAAAGTC
TCTTGATTTGACC
CTTAACACAGTTG
ATTTAGGTTACAA
TTGCTTGATATCC

Bioinformatics tools & pipelines
Automated
Scalable
Reproducible

Computer infrastructure

# How to perform RNAseq analysis

The good-practice analysis protocol takes 8 major steps.

➡ **01-04**: From raw sequencing to transcript quantification

➡ **05-08**: DEG and pathway analysis (06/01, part II)



Raw sequencing data + sample group ➡ Differentially expressed genes and pathways

# How to perform RNAseq analysis

# IGV (Integrative Genome Viewer)

http://software.broadinstitute.org/software/igv/home



- Load existing genomes, or generate custom genomes

- Visualize standard file formats
  - BAM
  - BED
  - GTF
  - ... and more!

# Reference databases

- **Gene annotation database: GENCODE**
  - https://www.gencodegenes.org/
- **Ensembl database**
  - https://www.ensembl.org/index.html
- UCSC Genome Browser
  - https://genome.ucsc.edu/
- NCBI databases
  - https://www.ncbi.nlm.nih.gov/guide/genomes-maps/
- Genomic databases
  - GDC: https://portal.gdc.cancer.gov/
  - GTEx: https://gtexportal.org/home/
  - Single-cell RNAseq: https://portals.broadinstitute.org/single_cell

https://www.gencodegenes.org/
model organisms

# Release 28 (GRCh38.p12)

## GTF / GFF3 files

| Content | Regions | Description | Download |
|---|---|---|---|
| Comprehensive gene annotation | CHR | • It contains the comprehensive gene annotation on the reference chromosomes only<br>• This is the **main annotation file** for most users | GTF  GFF3 |
| Comprehensive gene annotation | ALL | • It contains the comprehensive gene annotation on the reference chromosomes, scaffolds, assembly patches and alternate loci (haplotypes)<br>• This is a **superset** of the main annotation file | GTF  GFF3 |
| Comprehensive gene annotation | PRI | • It contains the comprehensive gene annotation on the primary assembly (chromosomes and scaffolds) sequence regions<br>• This is a **superset** of the main annotation file | GTF  GFF3 |
| Basic gene annotation | CHR | • It contains the basic gene annotation on the reference chromosomes only<br>• This is a **subset** of the corresponding comprehensive annotation, including only those transcripts tagged as 'basic' in every gene | GTF  GFF3 |
| Basic gene annotation | ALL | • It contains the basic gene annotation on the reference chromosomes, scaffolds, assembly patches and alternate loci (haplotypes)<br>• This is a **subset** of the corresponding comprehensive annotation, including only those transcripts tagged as 'basic' in every gene | GTF  GFF3 |
| Long non-coding RNA gene annotation | CHR | • It contains the comprehensive gene annotation of lncRNA genes on the reference chromosomes<br>• This is a **subset** of the main annotation file | GTF  GFF3 |
| PolyA feature annotation | CHR | • It contains the polyA features (polyA_signal, polyA_site, pseudo_polyA) manually annotated by HAVANA on the reference chromosomes<br>• This dataset does **not** form part of the main annotation file | GTF  GFF3 |
| Consensus pseudogenes predicted by the Yale and UCSC pipelines | CHR | • 2-way consensus (retrotransposed) pseudogenes predicted by the Yale and UCSC pipelines, but not by HAVANA, on the reference chromosomes<br>• This dataset does **not** form part of the main annotation file | GTF  GFF3 |
| Predicted tRNA genes | CHR | • tRNA genes predicted by ENSEMBL on the reference chromosomes using tRNAscan-SE<br>• This dataset does **not** form part of the main annotation file | GTF  GFF3 |

# Version 28 (November 2017 freeze, GRCh38) - Ensembl 92
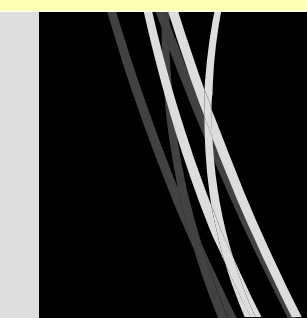
## General stats

| | | | |
|---|---|---|---|
| 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 | | |
| Immunoglobulin/T-cell receptor gene segments | | Total No of distinct translations | 61132 |
| - protein coding segments: | 408 | Genes that have more than one distinct translations | 13641 |
| - pseudogenes: | 237 | | |

## Fasta files

| Content | Regions | Description | Download |
|---|---|---|---|
| Transcript sequences | CHR | • Nucleotide sequences of all transcripts on the reference chromosomes | Fasta |
| Protein-coding transcript sequences | CHR | • Nucleotide sequences of coding transcripts on the reference chromosomes<br>• Transcript biotypes: protein_coding, nonsense_mediated_decay, non_stop_decay, IG_*_gene, TR_*_gene, polymorphic_pseudogene | Fasta |
| Protein-coding transcript translation sequences | CHR | • Amino acid sequences of coding transcript translations on the reference chromosomes<br>• Transcript biotypes: protein_coding, nonsense_mediated_decay, non_stop_decay, IG_*_gene, TR_*_gene, polymorphic_pseudogene | Fasta |
| Long non-coding RNA transcript sequences | CHR | • Nucleotide sequences of long non-coding RNA transcripts on the reference chromosomes | Fasta |
| Genome sequence (GRCh38.p12) | ALL | • Nucleotide sequence of the GRCh38.p12 genome assembly version on all regions, including reference chromosomes, scaffolds, assembly patches and haplotypes<br>• The sequence region names are the same as in the GTF/GFF3 files | Fasta |
| Genome sequence, primary assembly (GRCh38) | PRI | • Nucleotide sequence of the GRCh38 primary genome assembly (chromosomes and scaffolds)<br>• The sequence region names are the same as in the GTF/GFF3 files | Fasta |

- Genome-based alignment: STAR
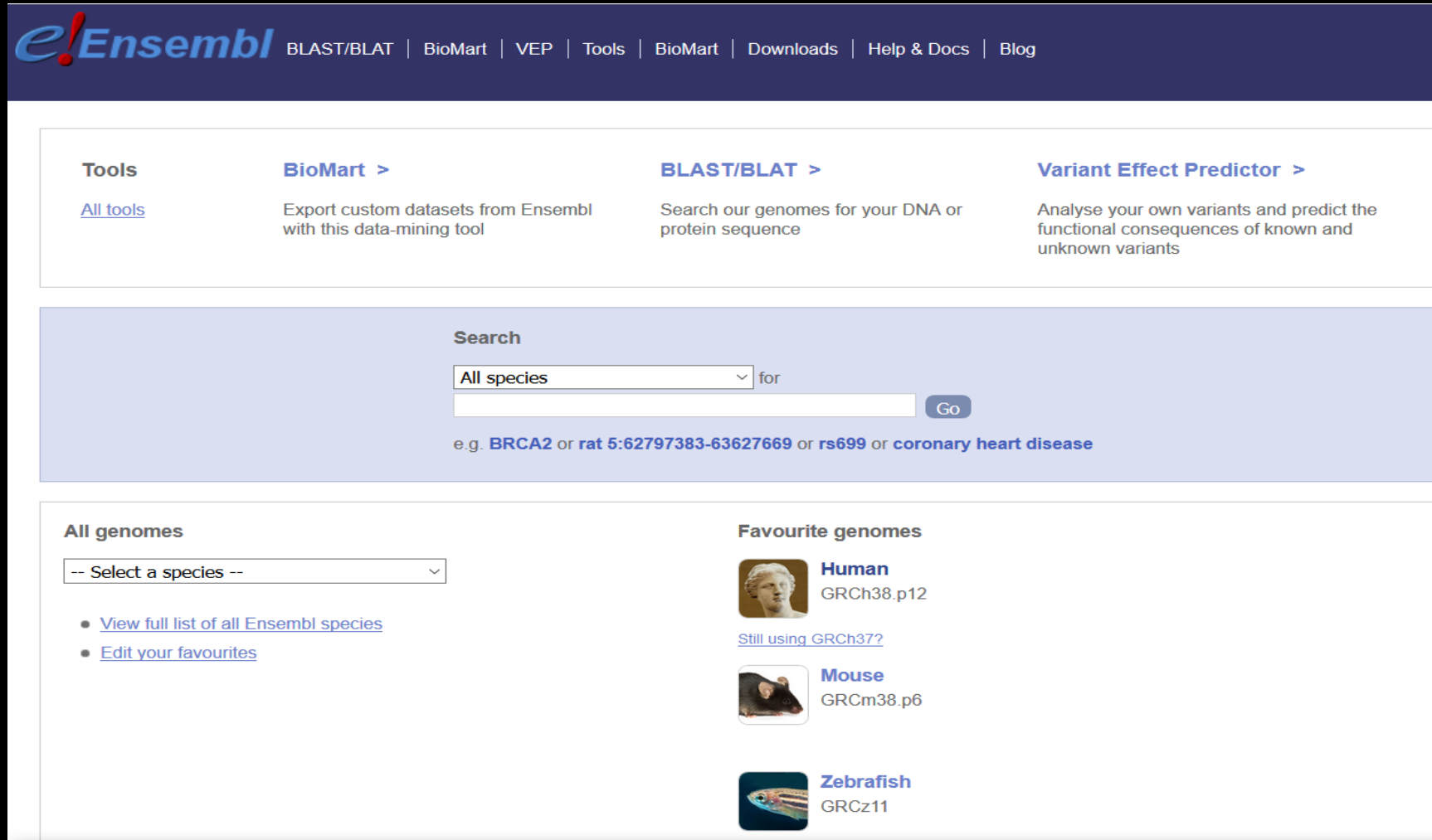- Transcriptome-based (pseudo)alignment: Kallisto, Salmon

## Metadata files

| Content | Regions | Description | Download |
|---------|---------|-------------|----------|
| Annotation remarks | ALL | • Remarks made during the manual annotation of the transcript | Metadata |
| Entrez gene ids | ALL | • Entrez gene ids associated to GENCODE transcripts (from Ensembl xref pipeline) | Metadata |
| Exon annotation evidence | ALL | • Piece of evidence used in the annotation of an exon (usually peptides, mRNAs, ESTs) | Metadata |
| Gene source | ALL | • Source of the gene annotation (Ensembl, Havana, Ensembl-Havana merged model or imported in the case of small RNA and mitochondrial genes) | Metadata |
| Gene symbol | ALL | • HGNC approved gene symbol (from Ensembl xref pipeline) | Metadata |
| PDB id | ALL | • PDB entries associated to the transcript (from Ensembl xref pipeline) | Metadata |
| PolyA features | ALL | • Manually annotated polyA features overlapping the transcript 3'-end | Metadata |
| PubMed id | ALL | • Pubmed ids of publications associated to the transcript (from HGNC website) | Metadata |
| RefSeq | ALL | • RefSeq RNA and/or protein associated to the transcript (from Ensembl xref pipeline) | Metadata |
| Selenocysteine | ALL | • Amino acid position of a selenocysteine residue in the transcript | Metadata |
| SwissProt | ALL | • UniProtKB/SwissProt entry associated to the transcript (from Ensembl xref pipeline) | Metadata |
| Transcript source | ALL | • Source of the transcript annotation | Metadata |
| Transcript annotation evidence | ALL | • Piece of evidence used in the annotation of the transcript | Metadata |
| TrEMBL | ALL | • UniProtKB/TrEMBL entry associated to the transcript (from Ensembl xref pipeline) | Metadata |

ID conversion between different annotation databases (e.g. NCBI/RefSeq, Ensembl)

# https://www.ensembl.org/index.html all organisms

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

### General stats

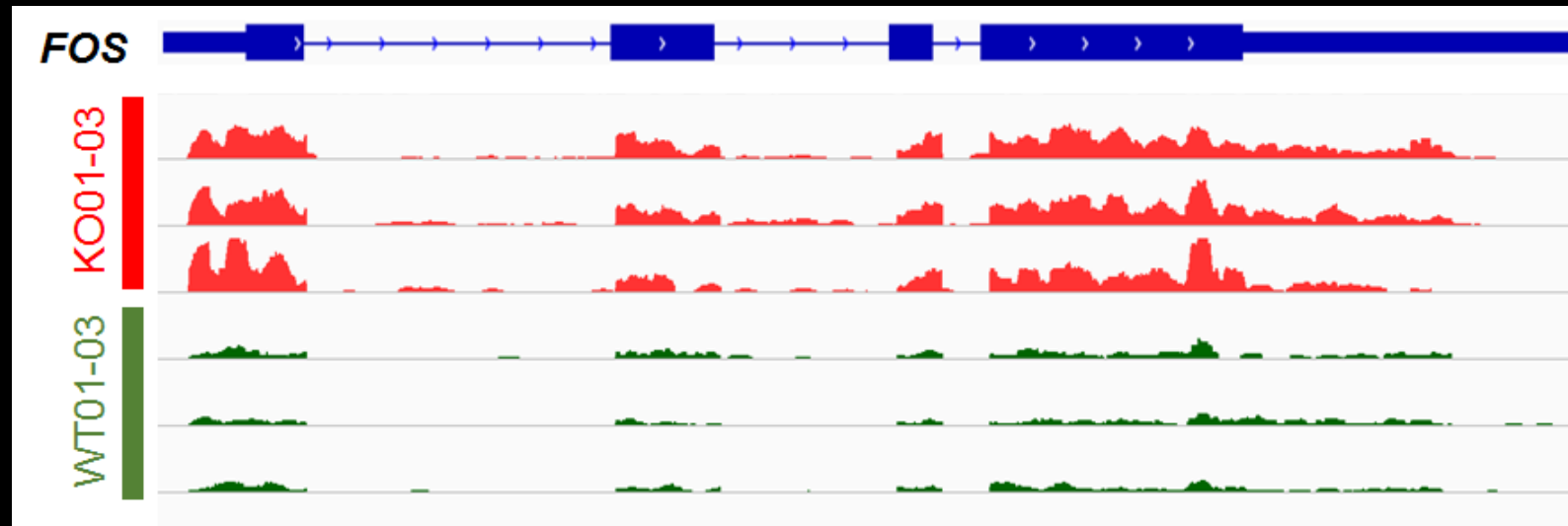| | | | |
|---|---|---|---|
| Total No of Genes | 58381 | Total No of Transcripts | 203835 |
| Protein-coding genes | 19901 | Protein-coding transcripts | 82335 |
| Long non-coding RNA genes | 15779 | - full length protein-coding: | 56541 |
| Small non-coding RNA genes | 7569 | - partial length protein-coding: | 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 | | |
| Immunoglobulin/T-cell receptor gene segments | | Total No of distinct translations | 61132 |
| - protein coding segments: | 408 | Genes that have more than one distinct translations | 13641 |
| - pseudogenes: | 237 | | |

Gencode and Ensembl are generally in sync

# How to perform RNAseq analysis

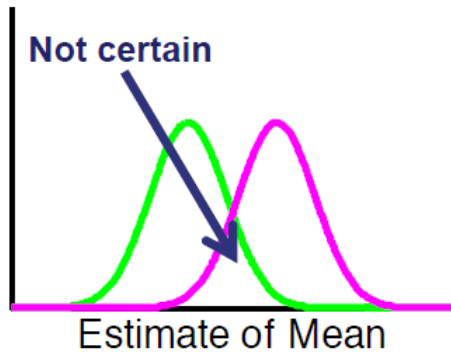# 05-08: Identify differentially expressed genes and pathways: DESeq2, clusterProfiler

- After steps 01-04, we have generated read alignment and counts for every annotated gene on the genome

- The next step is to utilize the read counts data to detect DEGs

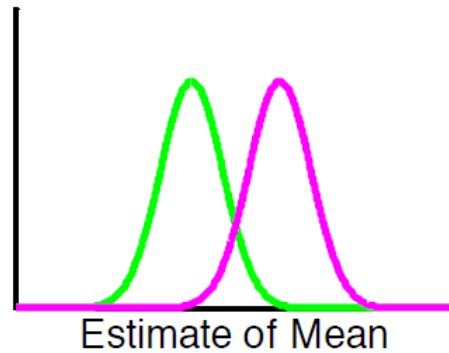- For example, if we visualize *FOS* gene across 6 samples in genome browser



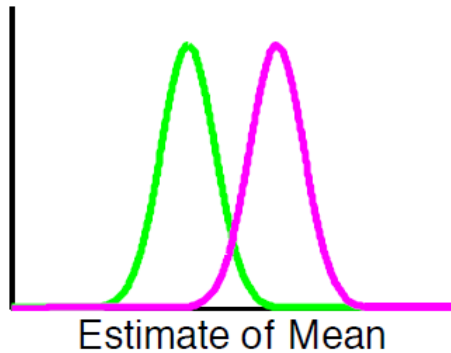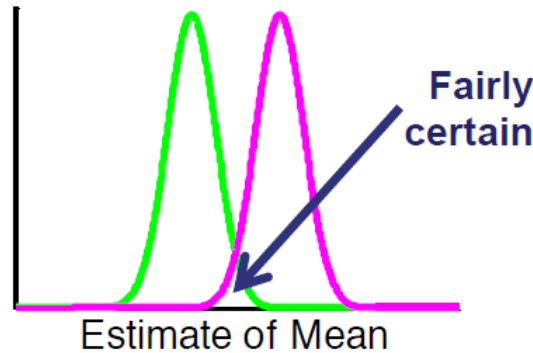*FOS = Fos proto-oncogene, AP-1 transcription factor subunit*

# DEG detection



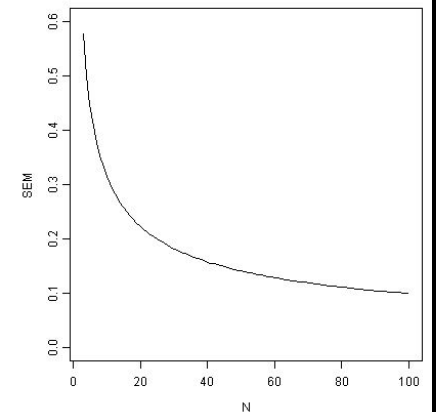$$SE = \frac{\sigma}{\sqrt{n}}$$

← Standard deviation

← Number of samples

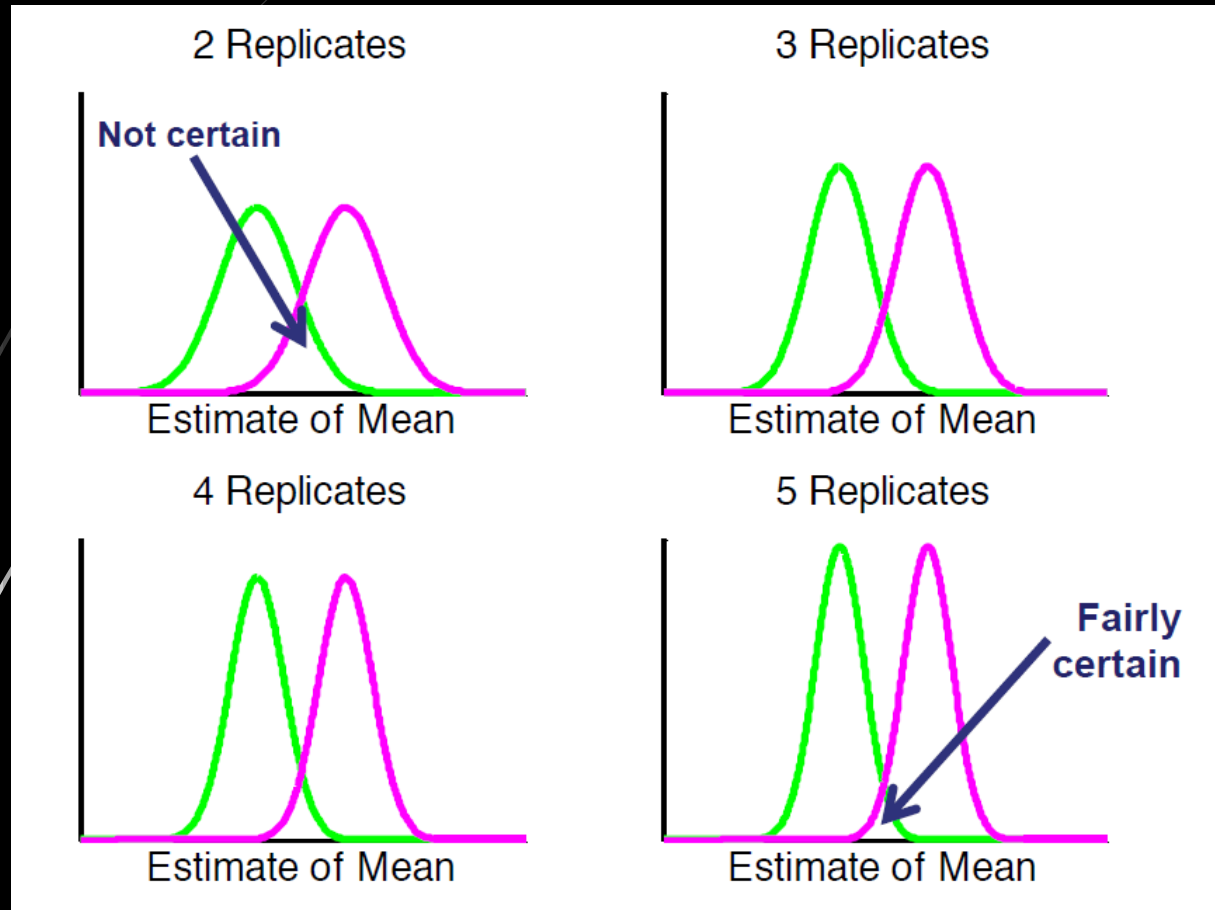### Standard Error of the Mean

$$SE_{M_x} = \frac{\sigma}{\sqrt{N}}$$

- This equation implies that sampling error decreases as sample size increases.
- This is important because it suggests that if we want to make sampling error as small as possible, we need to use as large of a sample size as we can manage.
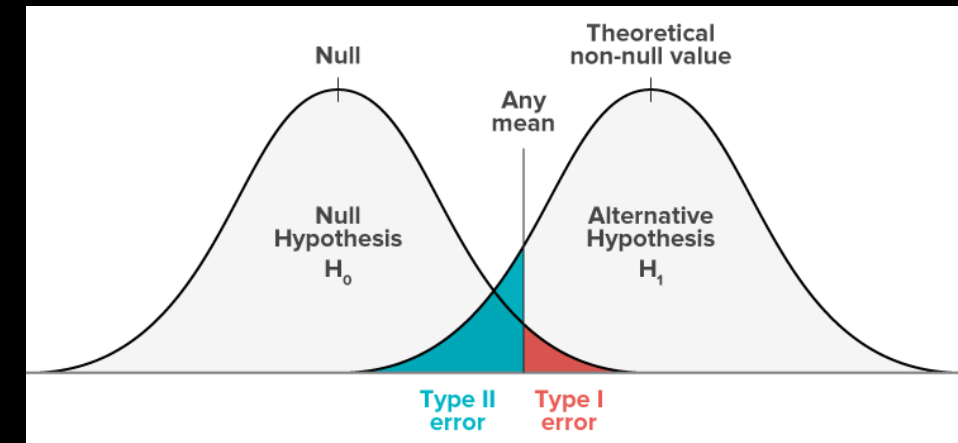
More biological replicates per group lead to higher discovery power, sensitivity and specificity.

# DEG detection

More biological replicates per group lead to higher discovery power, sensitivity and specificity.

# DEG detection



Measurement Uncertainty In Different Types of Replicates

- Biological Replicates
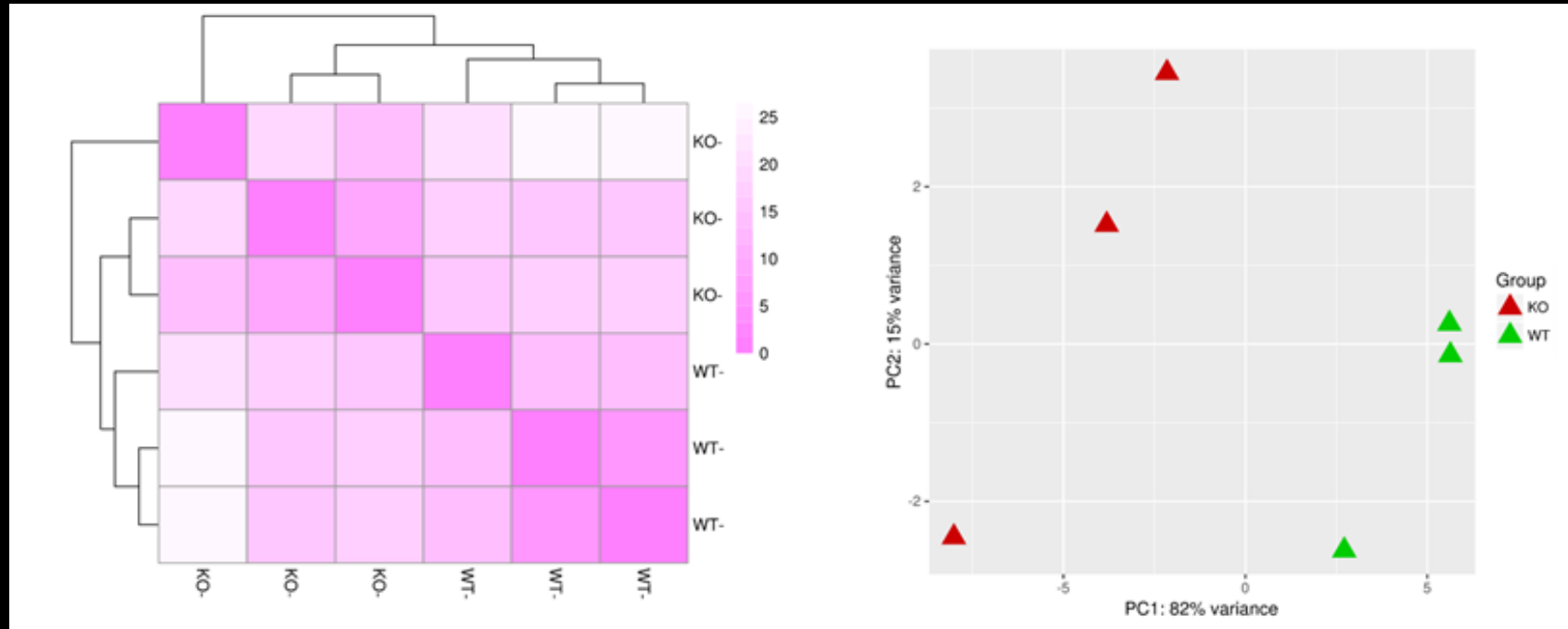- Technical Replicates
- Poisson Only

Replicate 2 / Replicate 1

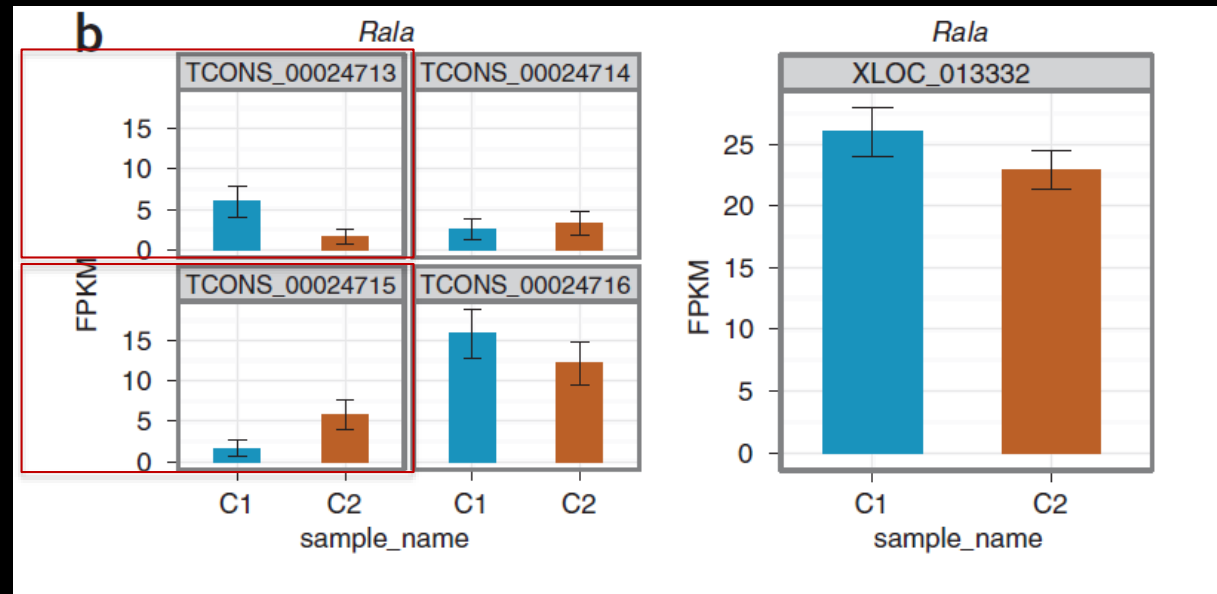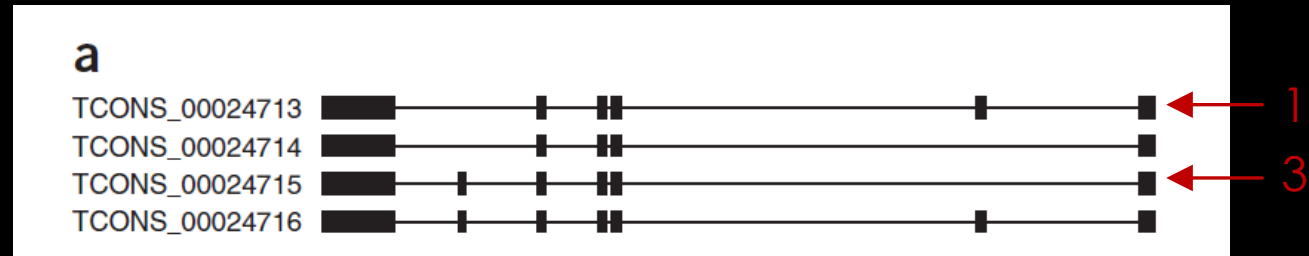High correlation is expected between biological replicates.

If one sample is an outlier, it can be identified if multiple replicates are included in an experiment.

# How to identify an outlier?

- PCA plot (visualization)
- Unsupervised sample clustering based on all genes or *top variable genes (e.g. 1500)*

# Transcript vs gene level quantification



Isoform 1 *

Isoform 3 *

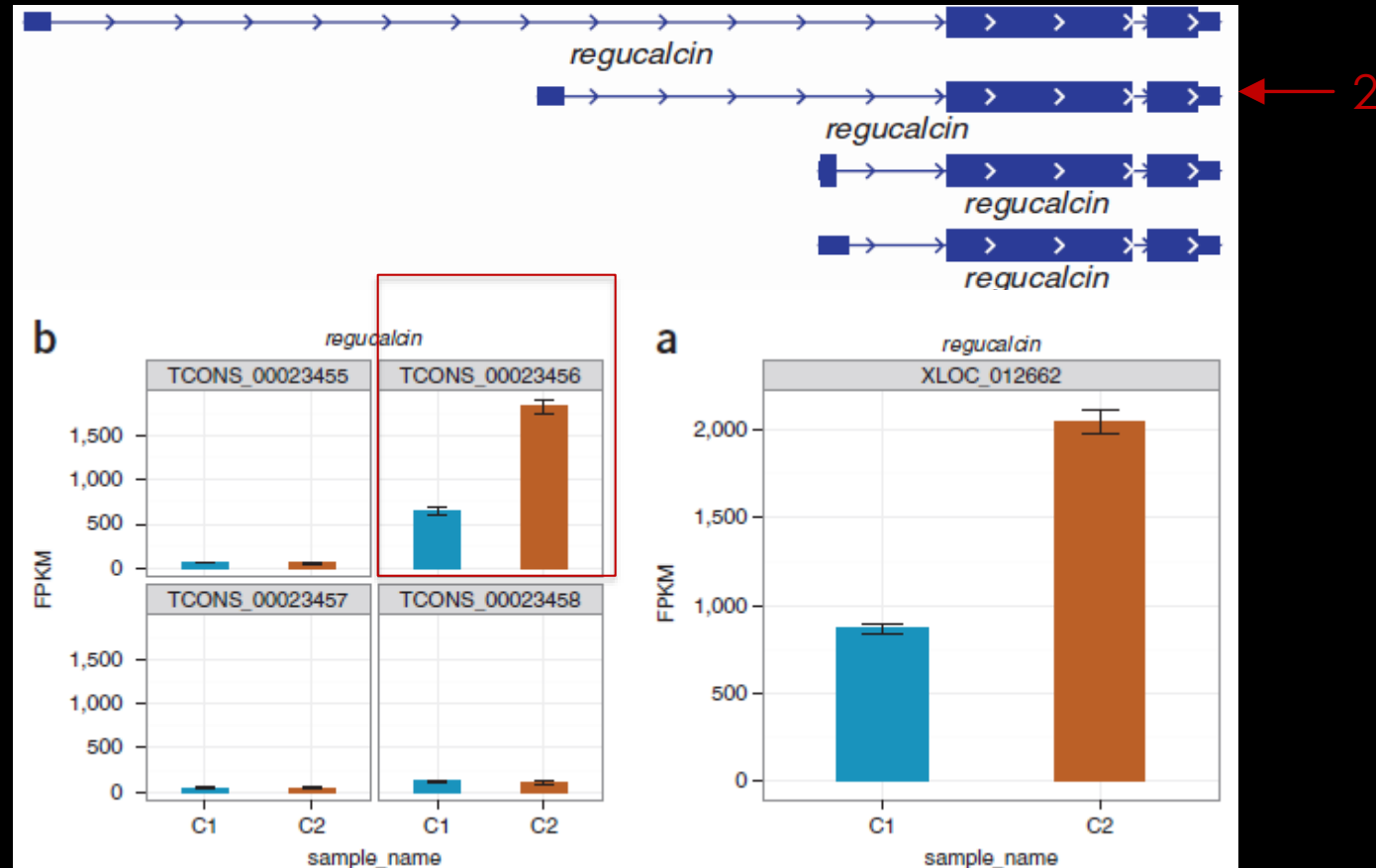Isoforms 1 – 4                                        Gene

**Difference in gene-level expression is not significant due to variability of isoforms**
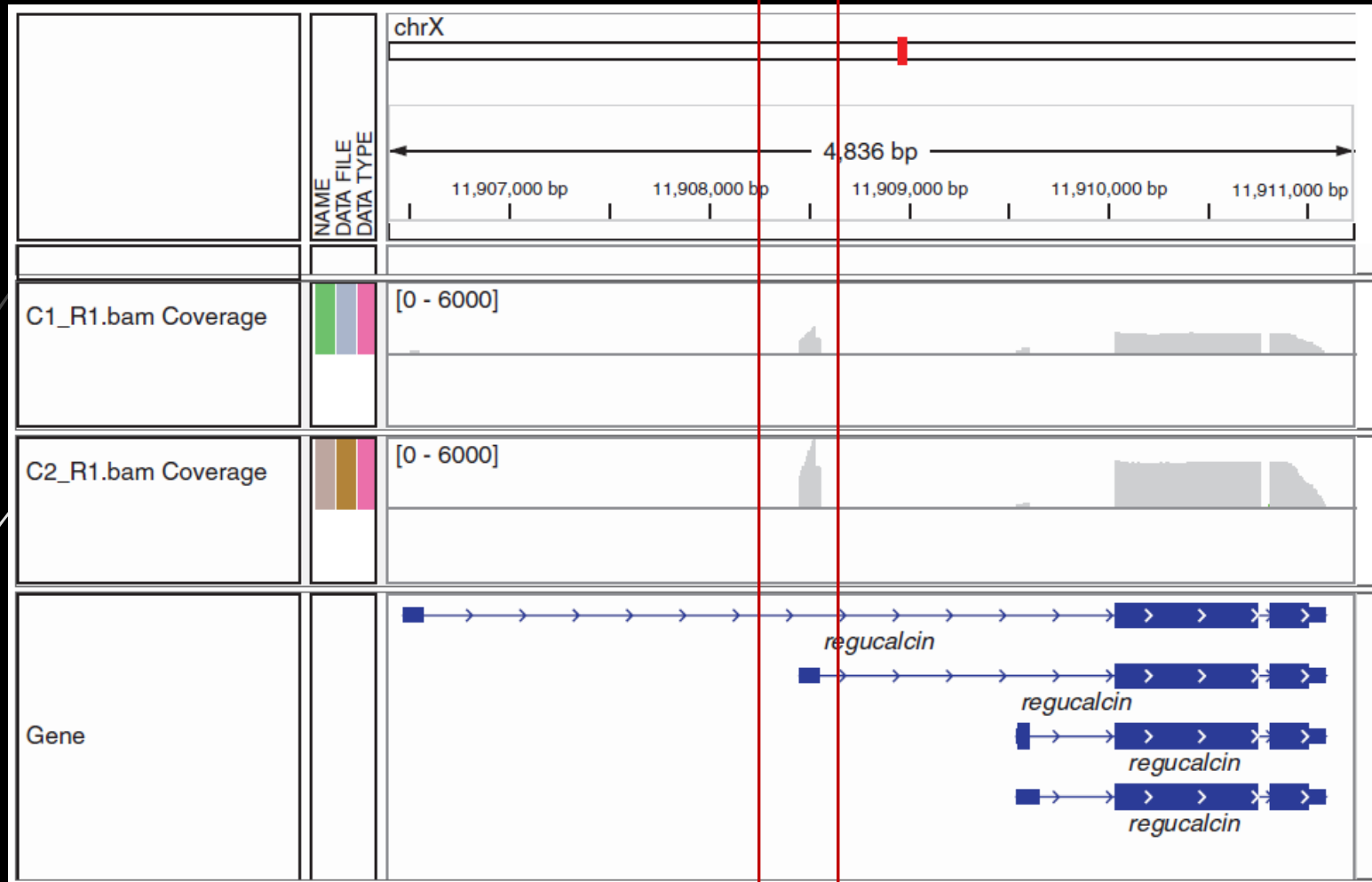
# Transcript vs gene level quantification
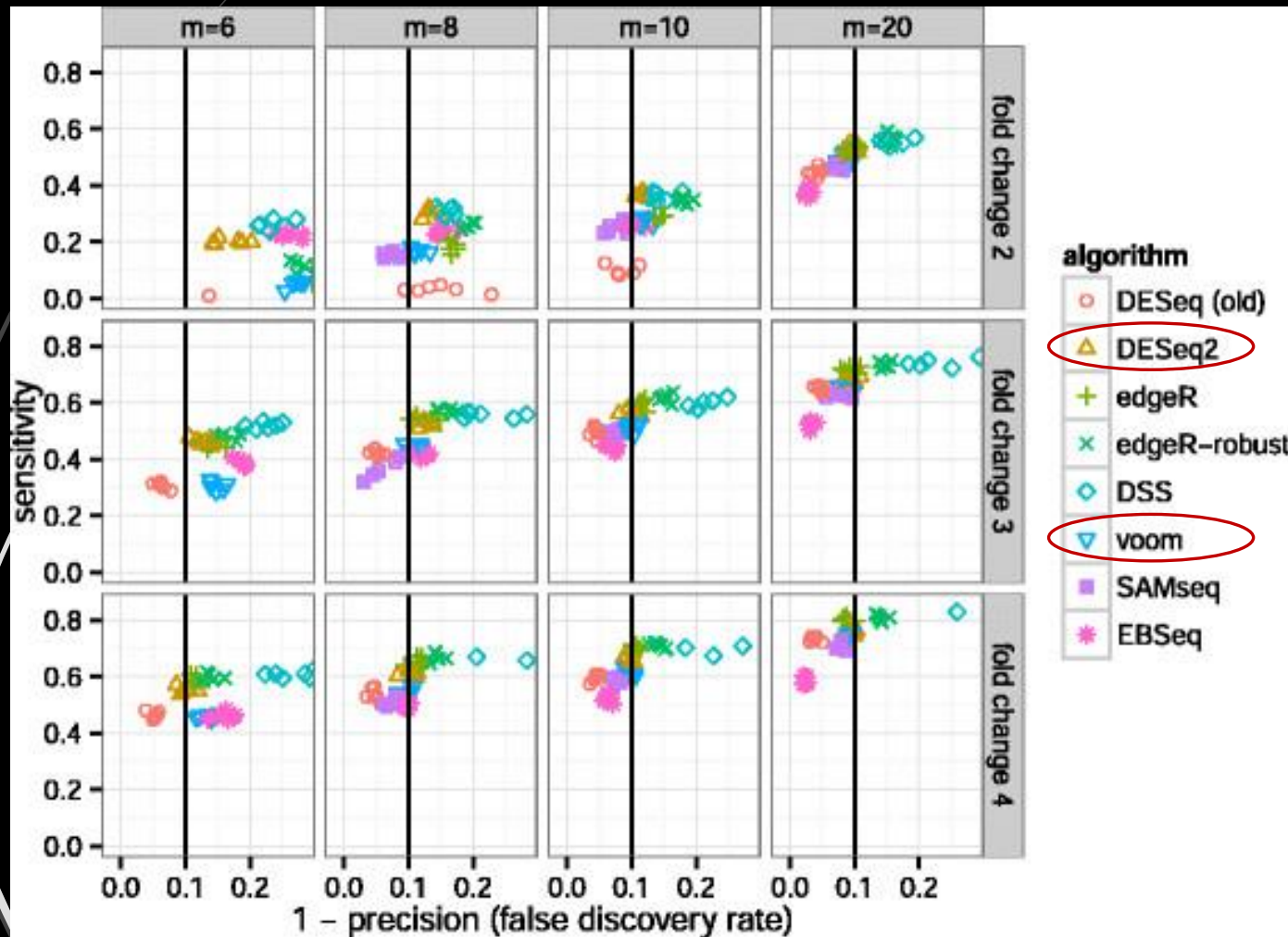


Isoform 2 *

Isoforms 1 – 4          Gene

**Difference in gene-level expression is significant, which is largely due to a great increase in the expression of isoform 2**

**Difference in gene-level expression is significant, which is largely due to a great increase in the expression of isoform 2**

# Comparison of different DEG identification methods



Sensitivity and precision of algorithms across combinations of sample size and effect size.
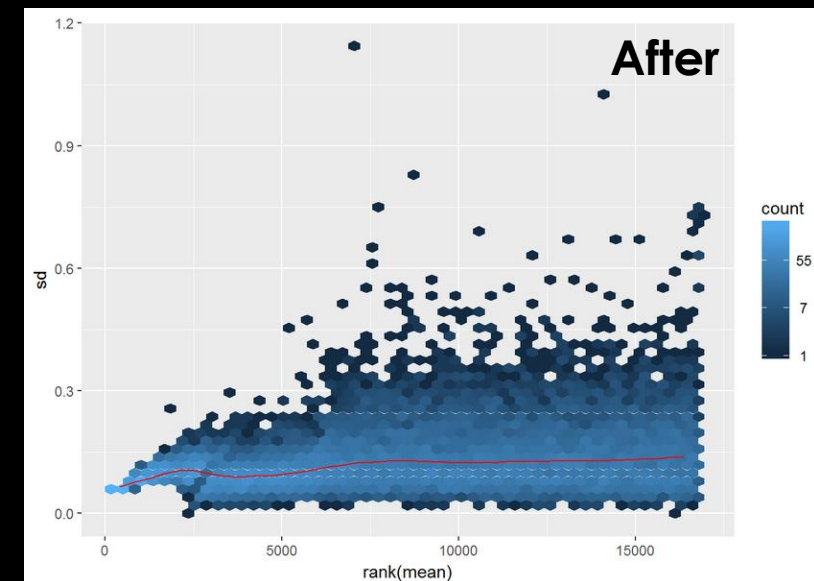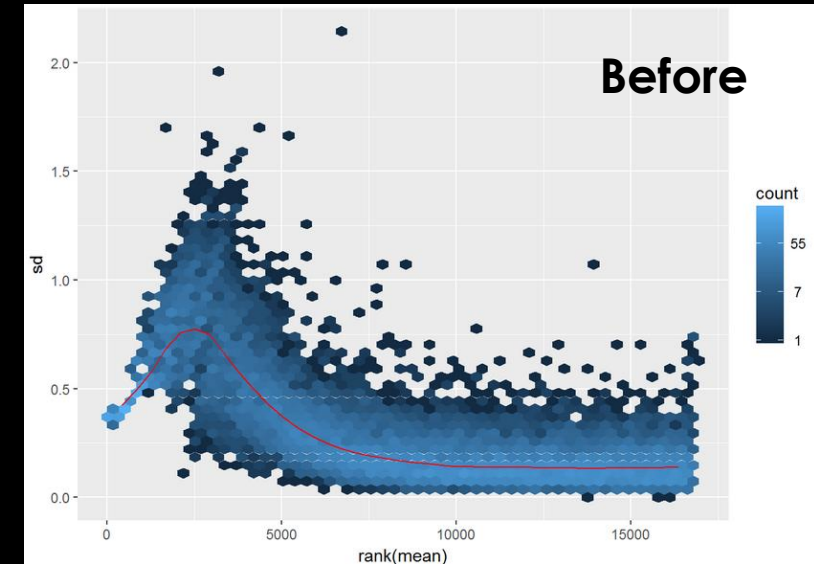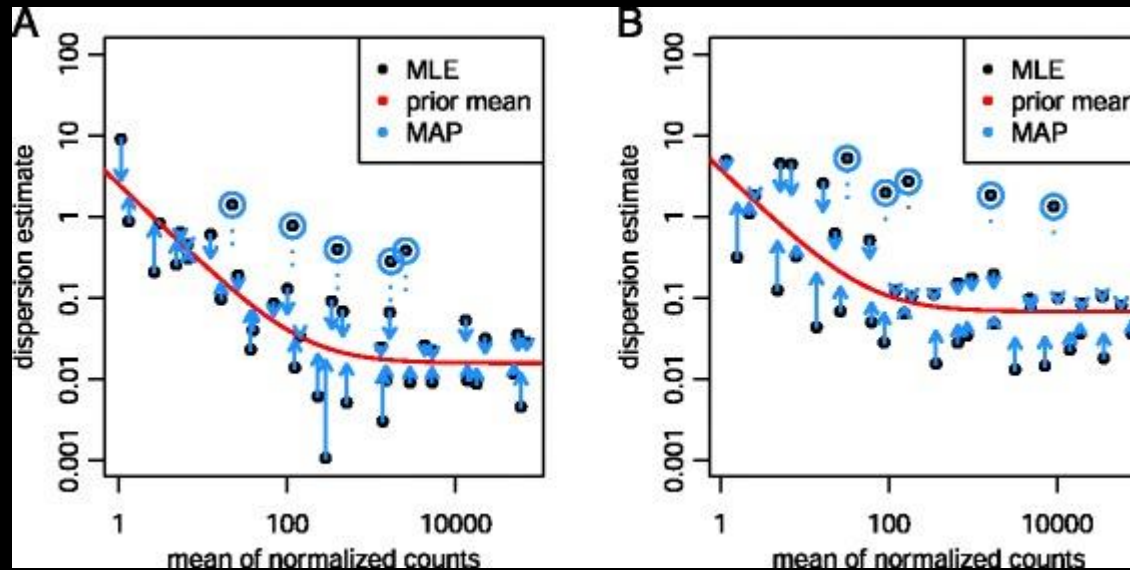
DESeq2 and edgeR often had the highest sensitivity of those algorithms that controlled the FDR, i.e., those algorithms which fall on or to the left of the vertical black line.

*m: total sample size; split into two even-sized groups for comparison*

# DESeq2

- Count matrix data

- Assume data follow negative binomial distribution (sometimes also called a gamma-Poisson distribution) with mean (μ) and dispersion (a) parameters

- Within-group variability, i.e., the variability between replicates, is modeled by the dispersion parameter alpha, which describes the *variance* of counts

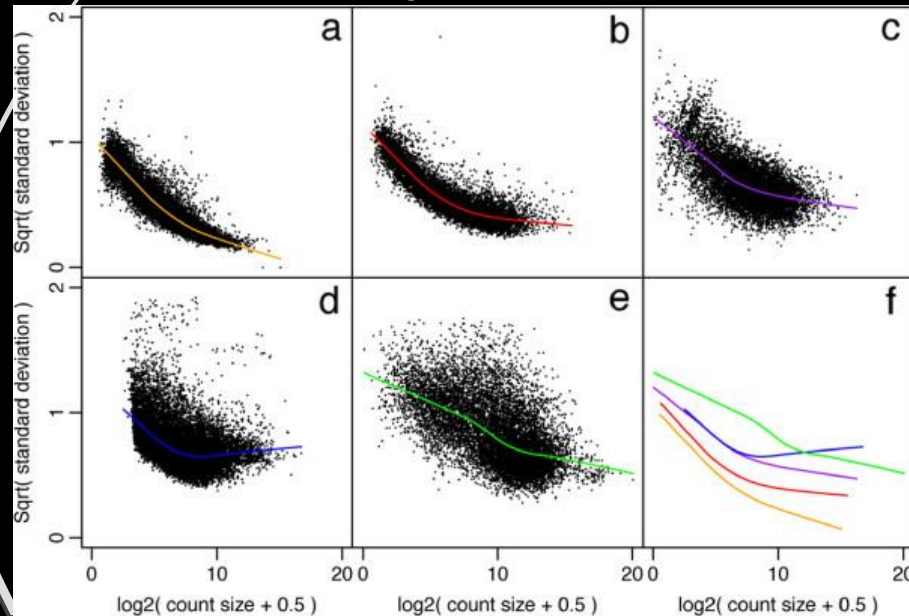- Empirical Bayes shrinkage for dispersion estimation

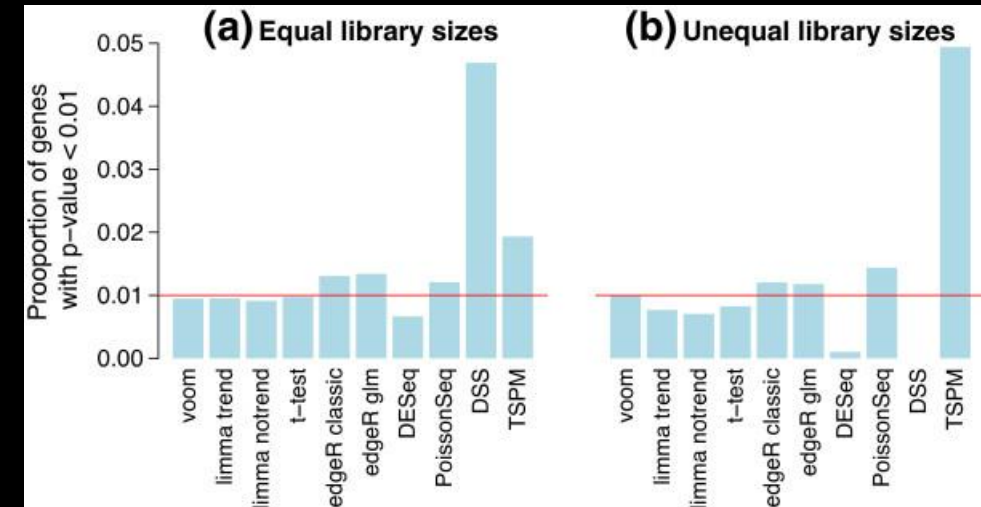MAP, maximum a posteriori; MLE, maximum-likelihood estimate

# Limma voom (*weighted* algorithm)

- To model the mean-variance relationship than to specify the exact probabilistic distribution of the counts (e.g. NB or Poisson)

- Provide accurate Type I (alpha) and Type II error (beta) control compared to other methods, especially when sample size is small
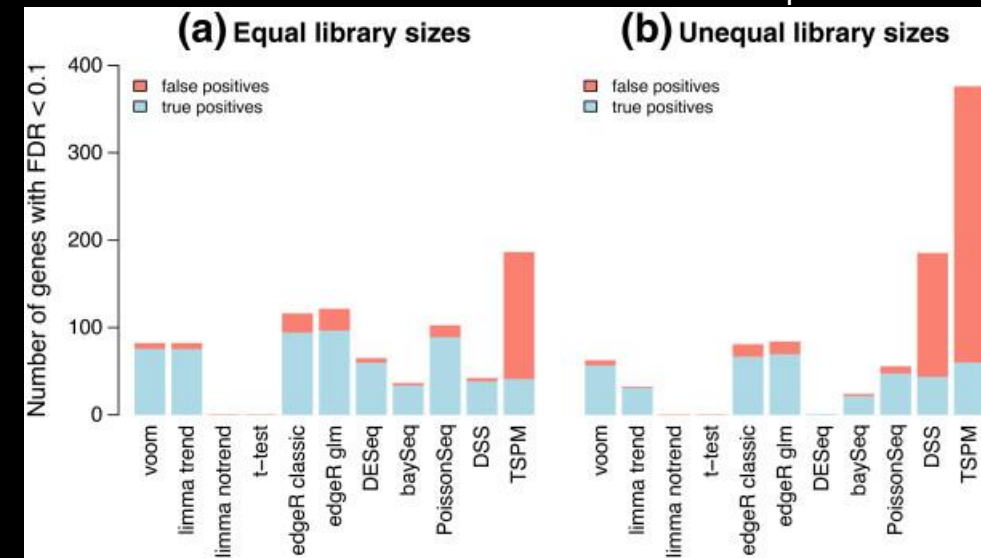
- *Voom with sample quality weights*

Law et al., Genome Biology 2014

Type I error rates in the absence of true differential expression



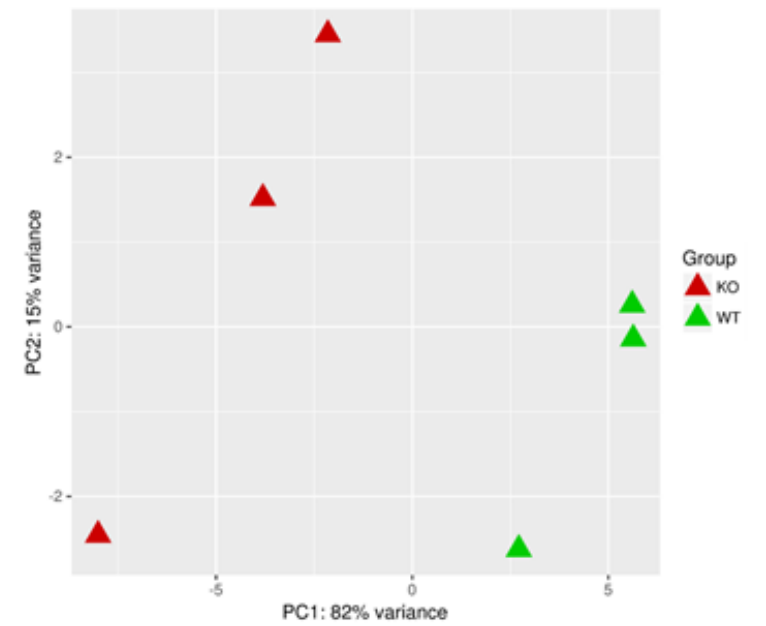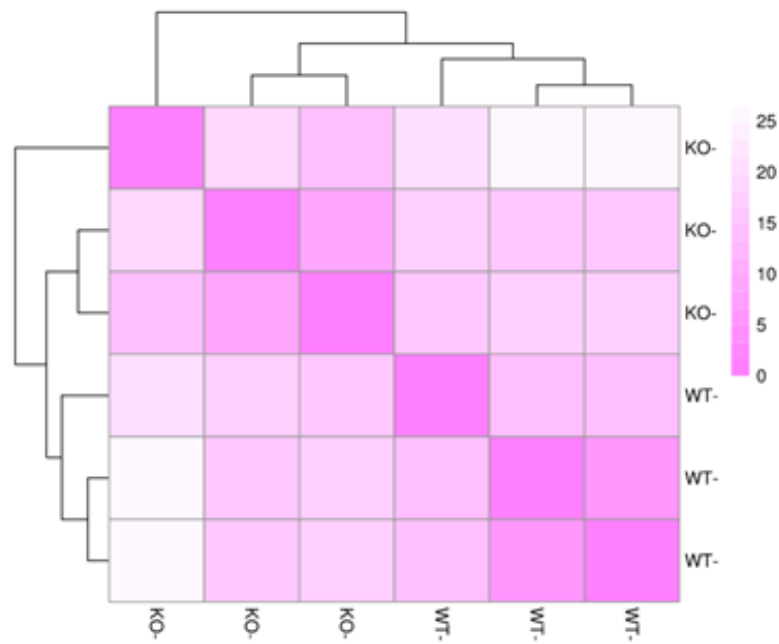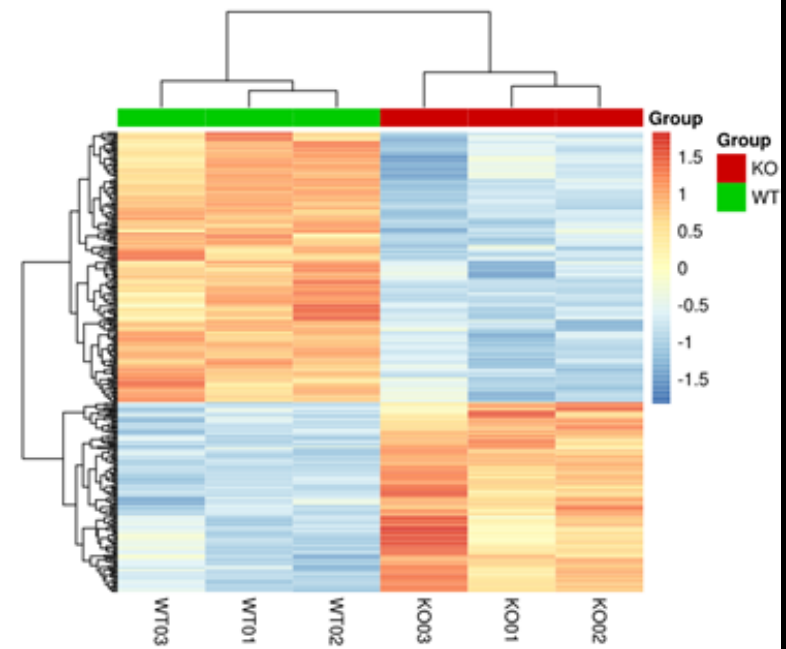Power to detect true differential expression

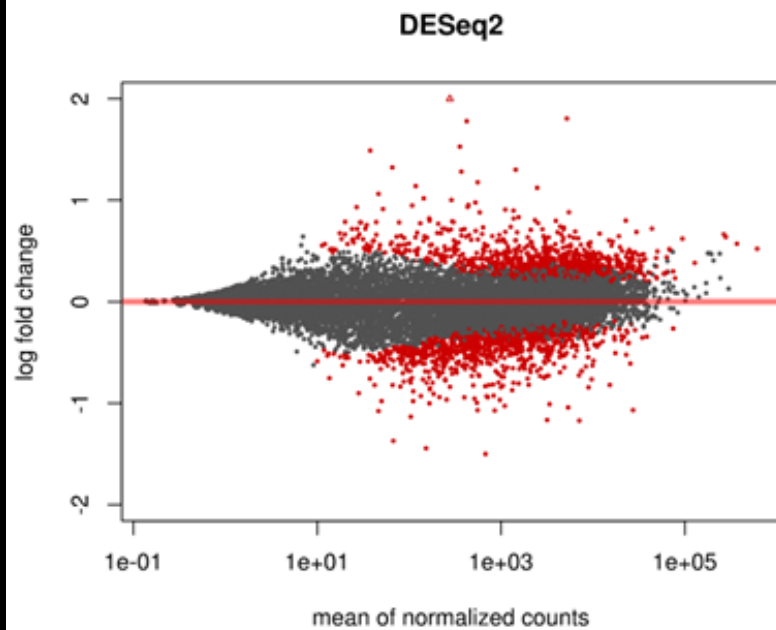# More databases!

- Gene annotation database: GENCODE

  - https://www.gencodegenes.org/

- Gene Ontology (GO) database: Gene Ontology Consortium

  - http://www.geneontology.org/

- Pathway database: KEGG

  - http://www.genome.jp/kegg/

- Predefined gene sets: MSigDB

  - http://software.broadinstitute.org/gsea/msigdb/

[1] "Genes significant = 296 (fc, 1.5, fdr 0.05)"
[1] "Heatmap = 296 genes on the row, 6 samples on the column"

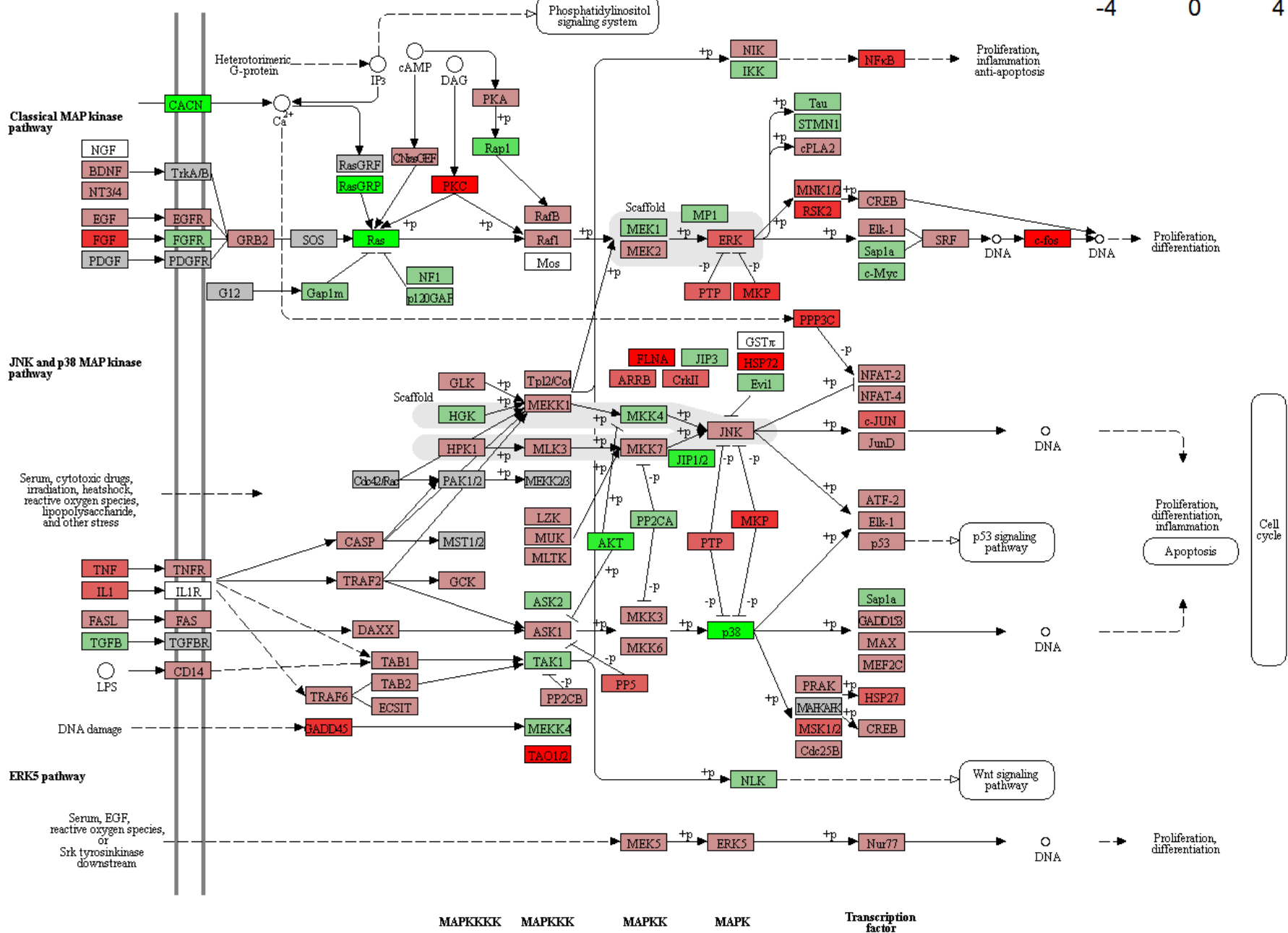MAPK SIGNALING PATHWAY

Data on KEGG graph
Rendered by Pathview
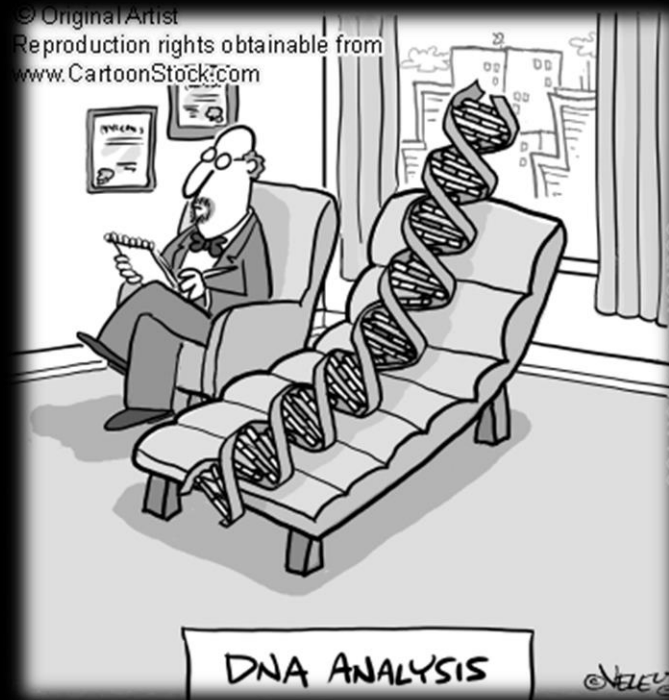
# Thank you!



Questions

# Hands-on practice START

- Open your handson.Rmd on the Github or download to local computer
- https://github.com/MScBiomedicalInformatics/MSIB32500/blob/master/lectures/handson9.html
- *Dataset: two groups (PRDM11 KO vs WT, human U2932 cells), 6 samples*
- *Single-end reads, unstranded libraries*

| Sample | Group | Sequencing File | Sequencing Data |
|--------|-------|-----------------|-----------------|
| KO01 | KO | KO01.fastq.gz | 74,126,025 reads |
| KO02 | KO | KO02.fastq.gz | 64,695,948 reads |
| KO03 | KO | KO03.fastq.gz | 52,972,573 reads |
| WT01 | WT | WT01.fastq.gz | 55,005,729 reads |
| WT01 | WT | WT02.fastq.gz | 61,079,377 reads |
| WT01 | WT | WT03.fastq.gz | 66,517,156 reads |

Fog. et al. 2015. Loss of *PRDM11* promotes MYC-driven lymphomagenesis. Blood 125(8):1272-81

*PRDM11 = PR/SET domain 11*