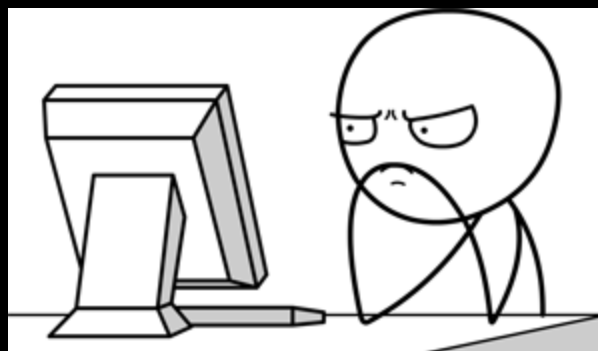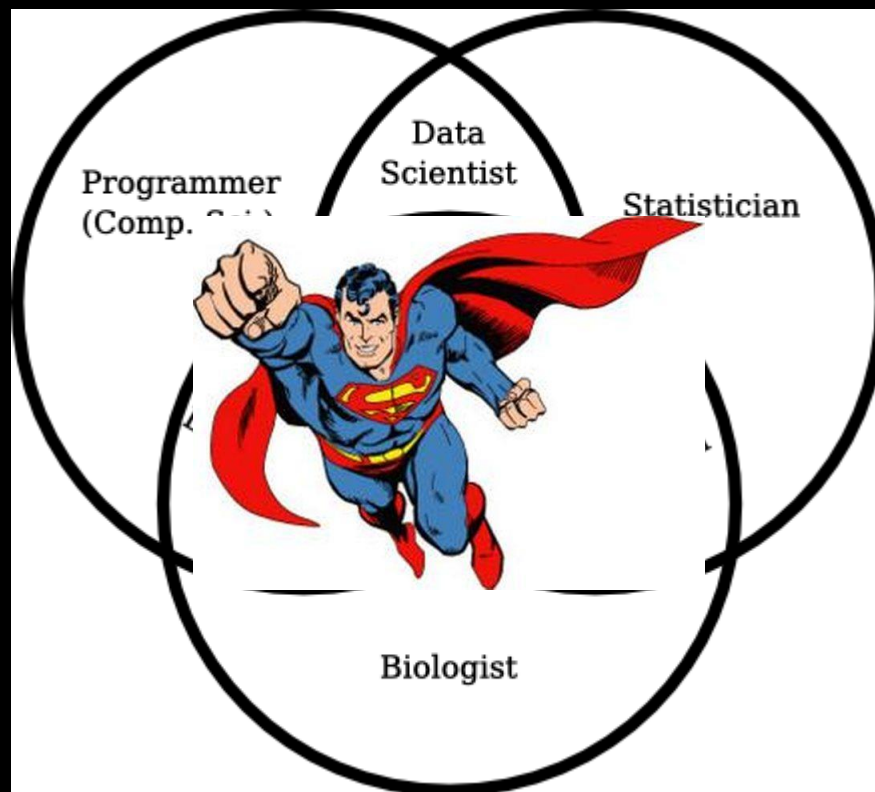MSIB32500 Advanced Bioinformatics Spring 2017

# RNAseq Data Analysis and Clinical Applications, Part I

Riyue Bao, Ph.D.

Research Assistant Professor (Bioinformatics)

Center for Research Informatics & Department of Pediatrics

The University of Chicago

# Outline

- Part I (05/27/2017)
  - Introduction to RNAseq technology and clinical applications
  - Hands on: From raw data to gene expression quantification

- Part II (06/03/2017)
  - Differential gene expression analysis and data visualization
  - Hands on: Identification of genes and pathways significantly changed under condition
  - *Homework assignment*

- Part III (06/08/2017)
  - 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
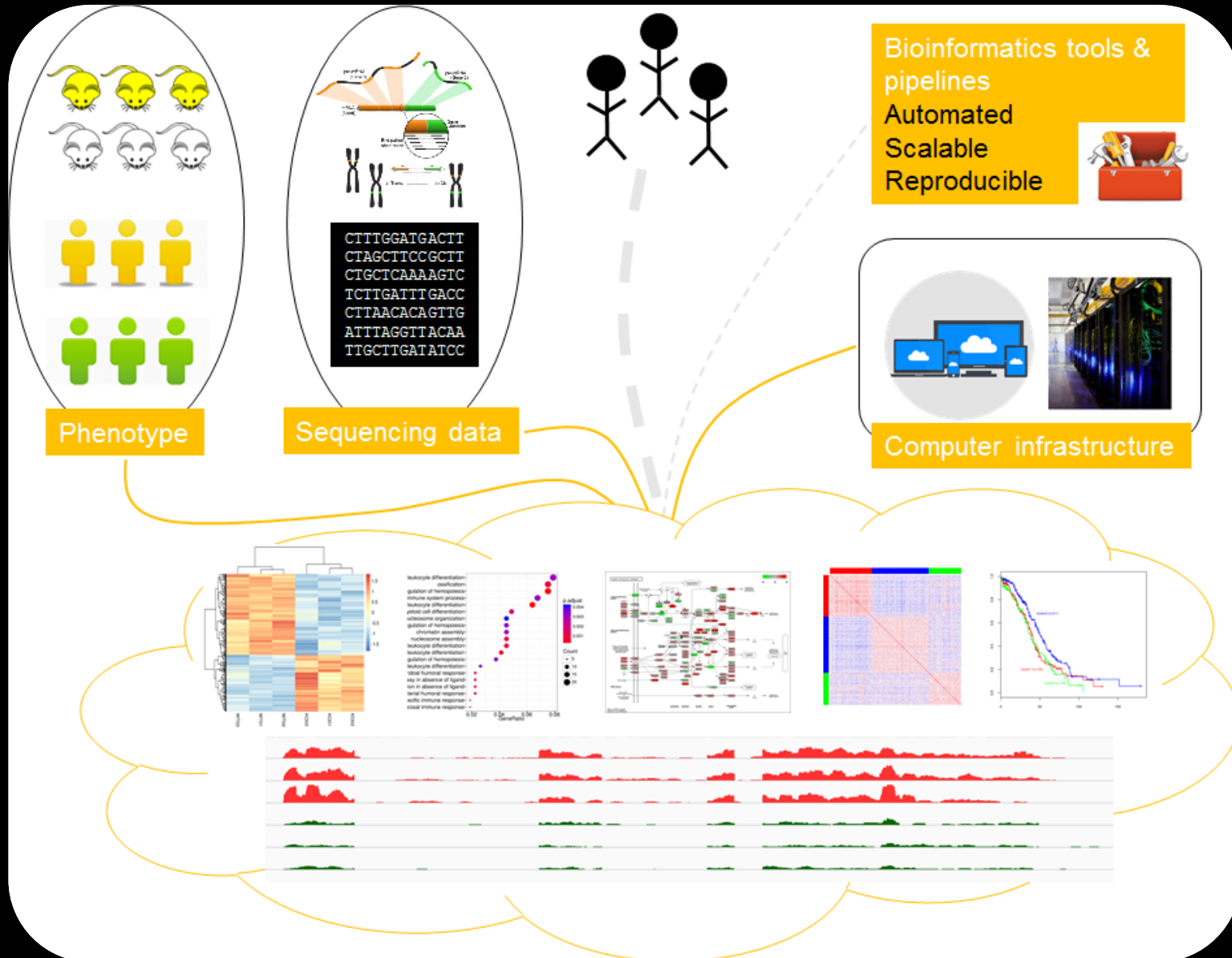
    - **lecture8.pdf**

    - **Handson8.Rmd**

- Tarbell high-performance computing (HPC) clusters (hands on practice)

# Objective

- Learn the good-practice RNAseq analysis pipeline
- Learn commonly used bioinformatics tools
- Practice the automated, scalable pipeline
- Explore the quality metrics and input/output of the RNAseq pipeline
- Visualize result files and quality plots

Phenotype

Sequencing data

CTTTGGATGACTT
CTAGCTTCCGCTT
CTGCTCAAAAGTC
TCTTGATTTGACC
CTTAACACAGTTG
ATTTAGGTTACAA
TTGCTTGATATCC

Bioinformatics tools & pipelines
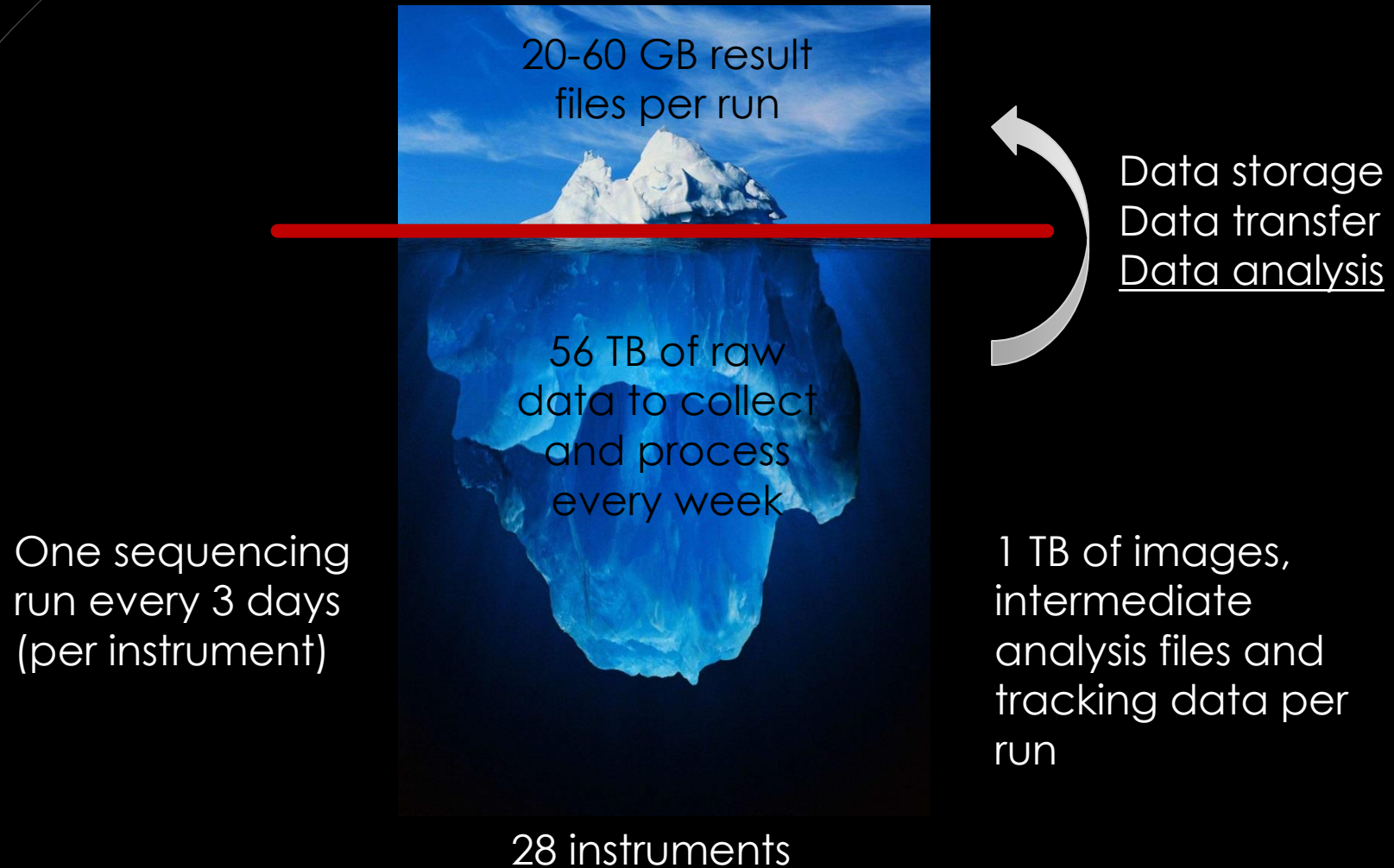Automated
Scalable
Reproducible

Computer infrastructure

# Biological and clinical questions

- I am interested in studying transcriptional landscape shift before and after drug treatment in cell lines

- I want to identify which pathways are affected after knocking down my favorite gene in mice

- I have expression data of clinical isolates collected at various time points, when patient's response changed. Why?

- I have a cohort of patients and want to discover which gene signature predicts patient's response to treatment

- I want to detect gene fusions, expressed mutations, and disrupted isoforms in tumors that may be related to disease

... and more!

# The Sequencing Iceberg

20-60 GB result files per run

56 TB of raw data to collect and process every week

Data storage
Data transfer
Data analysis

One sequencing run every 3 days (per instrument)

1 TB of images, intermediate analysis files and tracking data per run

28 instruments

# Current Sequencing Technologies

**Illumina**
HiSeq 2000/2500/X *TEN*



**Illumina**
MiSeq



**Ion Torrent** (Life Technologies)
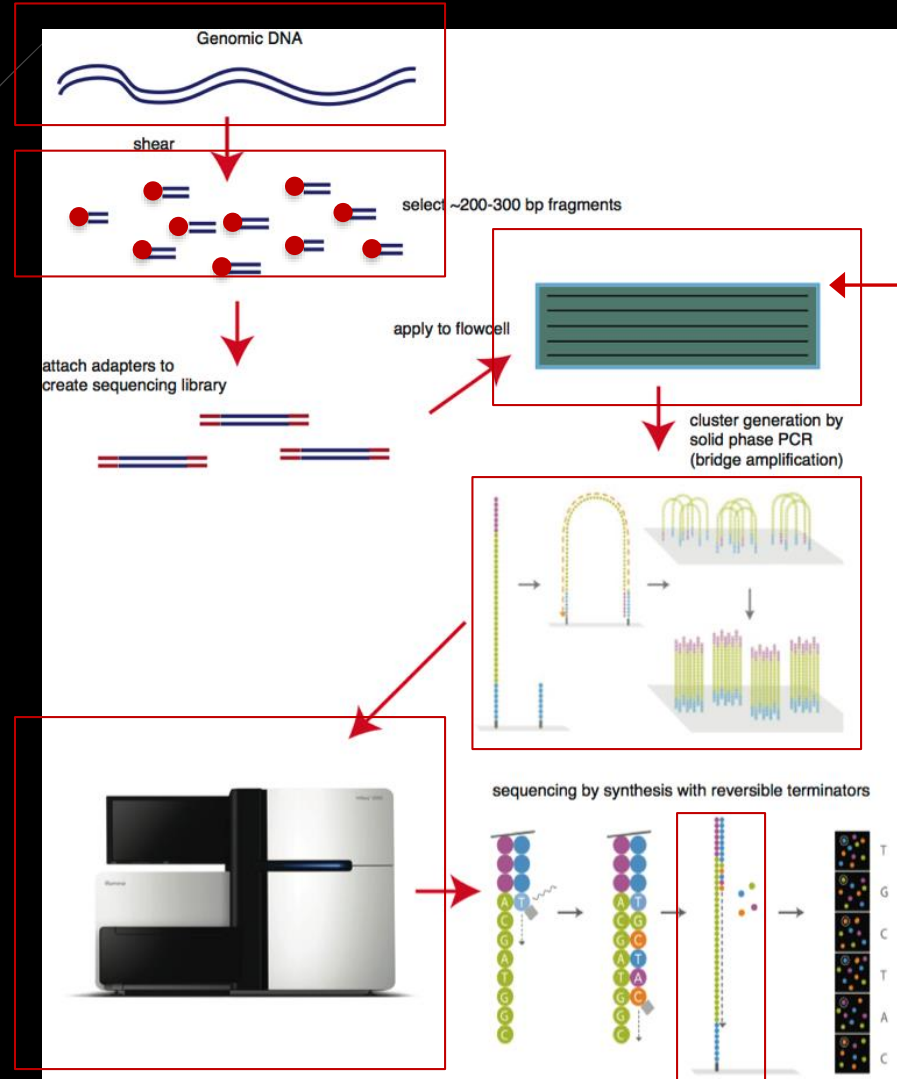Ion Proton
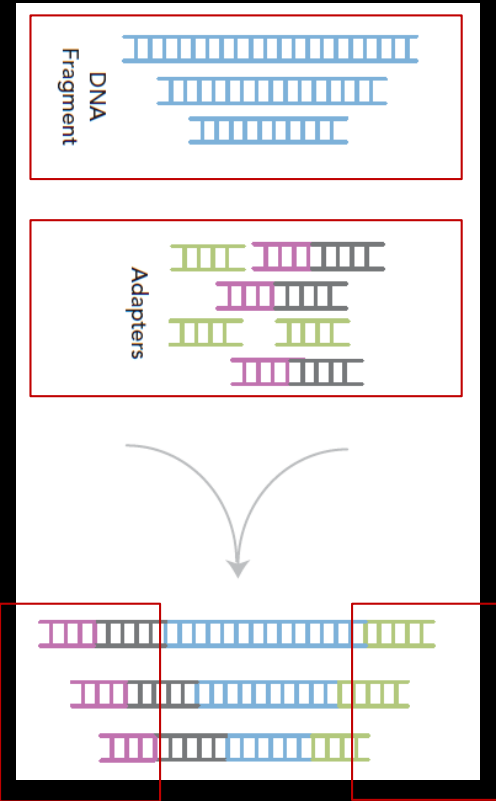


**PacBio**
RS

# Illumina



- Sample
- Library
- DNA fragment
- Barcode
- Run
- Flow cell
- Lane
- Cluster
- Read
- Adapter/Primer

Though we are talking about Illumina here, many of those terms can be applied to other sequencing technologies.
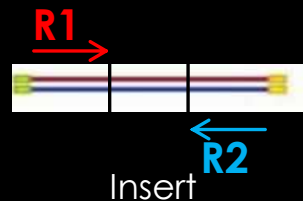
# Illumina

DNA fragment

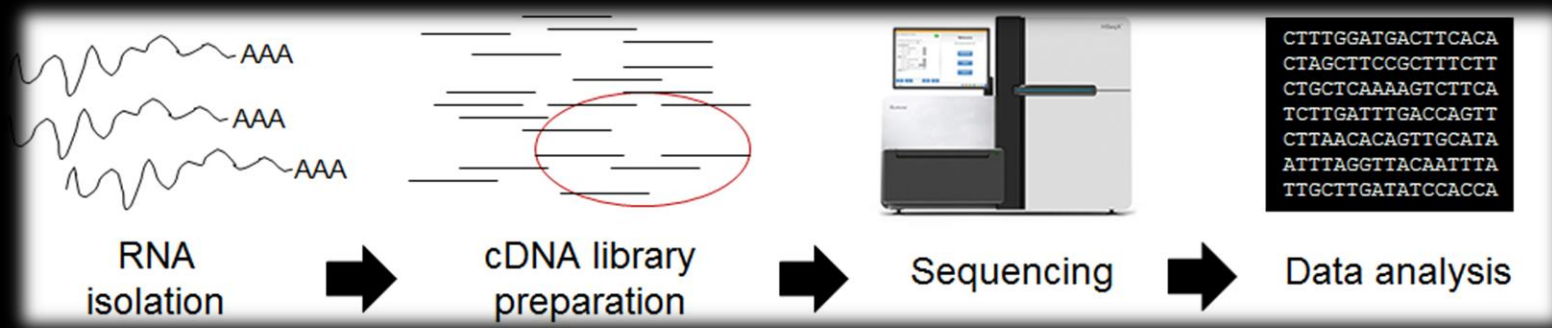Adapter/Primer

Single-End

Paired-End

R1

R2

Insert

- Sample
- Library
- DNA fragment
- Barcode
- Run
- Flow cell
- Lane
- Cluster
- Read
- Adapter/Primer

Fragment size = R1+R2+insert size

# What is RNAseq?

**High-throughput sequencing of RNA: Profile, identify or assemble transcripts**

➡ Detect gene expression changes between conditions

➡ Identify novel splice sites / exons, mutations, fusion genes, etc.

➡ Broad detection range, high sensitivity, low requirement of RNA amount

➡ Available for all species (reference genome is optional): reference genome-guided alignment or *de novo* assembly
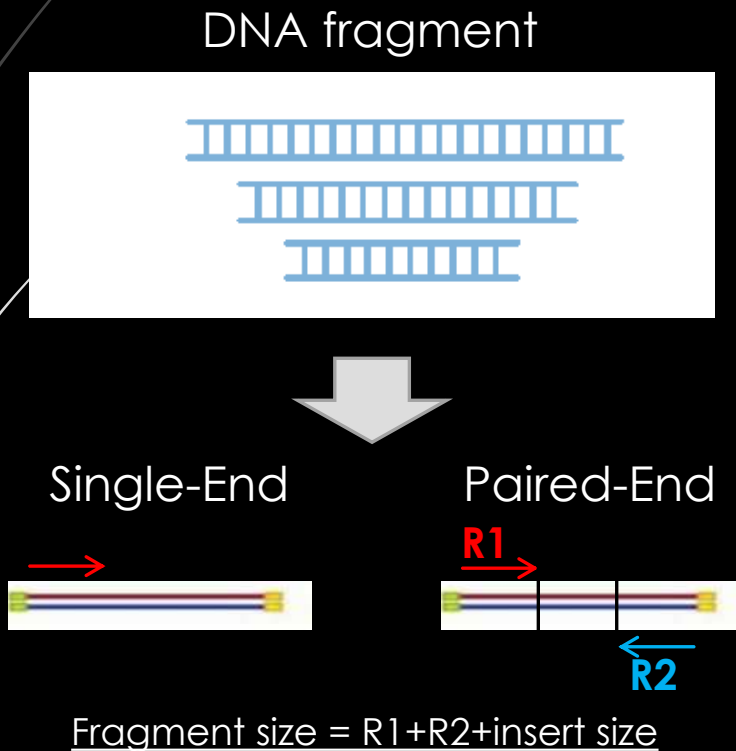


*Biological sample or clinical specimen*
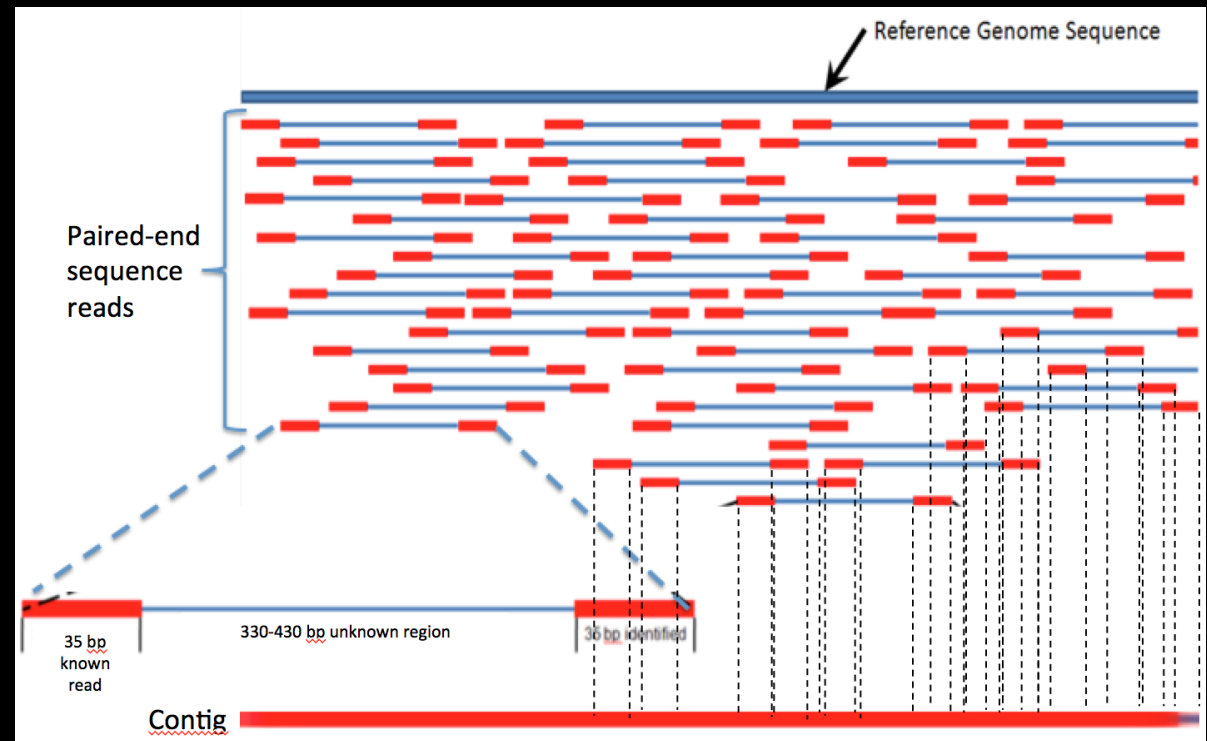
*Millions of reads!*

# Which factors to consider if I want to initiate an RNAseq experiment?

- Sample source (fresh cells, frozen, FFPE)
- Sample quality (RNA quality, DNA contamination, tumor purity)
- RNA concentration (e.g. 100ng total RNA)
- Ribosome RNA depletion (accounts for 80% of total RNA)
- **Library type (single-end or paired-end reads)**
- **Library strandness (unstranded or stranded/directional)**
- **Sequencing depth (20 million, 50 million, or >100 million reads)**

# Library type: Single-end vs paired-end reads

DNA fragment

PE reads: Higher alignment rate and accuracy

Single-End

Paired-End

R1

R2

Fragment size = R1+R2+insert size

Reference Genome Sequence

Paired-end sequence reads

35 bp known read

330-430 bp unknown region

36 bp identified

Contig

*2 x 50bp PE reads >> 1 x 100 bp SE reads!*

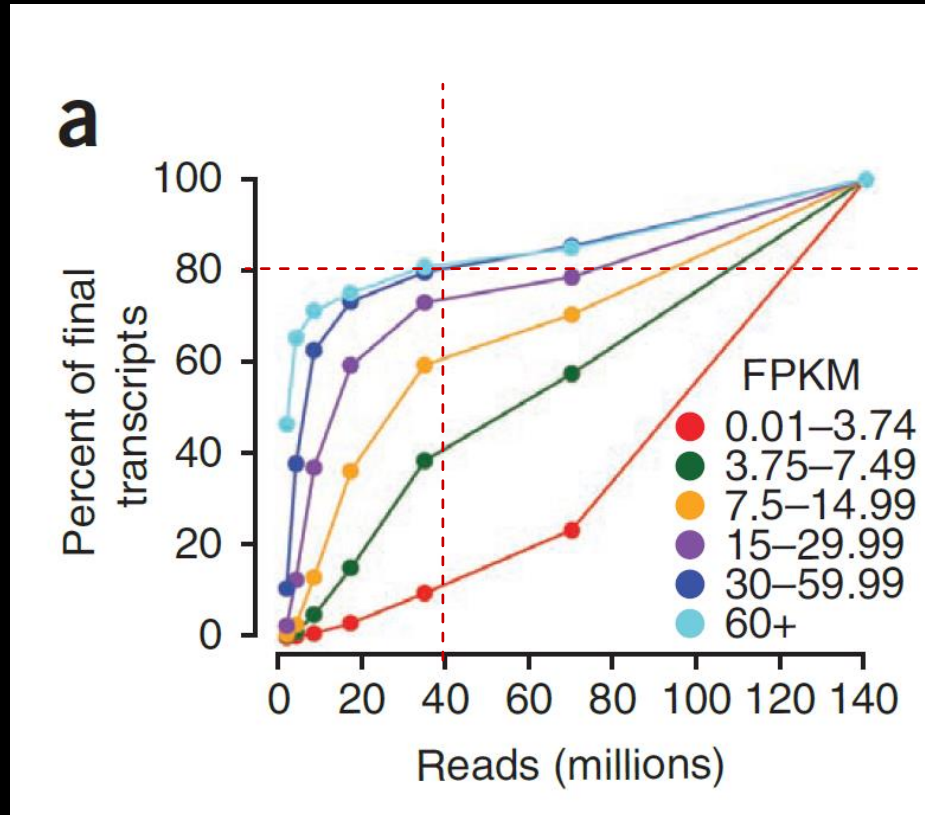# Library strandness: stranded-protocol is always recommended!



- Non-stranded/Non-directional
- Forward, also known as second-stranded (coding strand)
- Reverse, also known as first-stranded (template strand)

# Sequencing depth: How many reads do I need?

Coverage affects downstream gene discovery, expression estimation, and power of statistical analysis
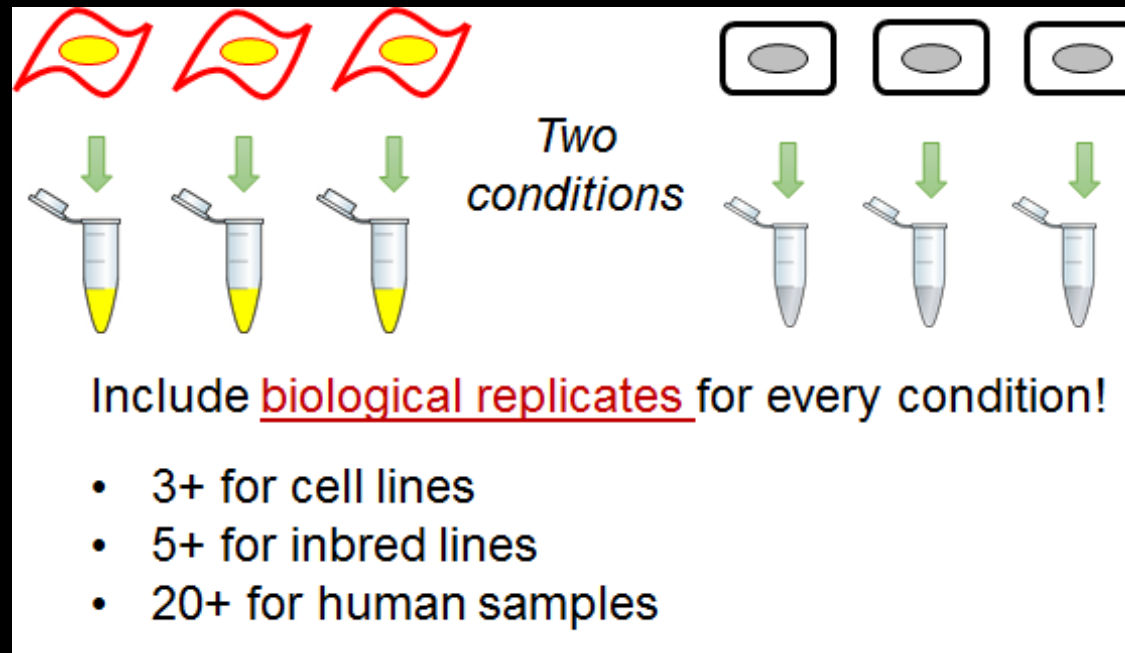
*The more, the better!*



Trapnell et al. Nature Protocol 2012

**ENCODE saturation analysis**

- 214 million 2x100bp PE reads
- H1 human embryonic stem cells
- 80% of the genes with FPKM≥10 are detected by ~36 million mapped reads per sample
- Genes with FPKM<10: ~80 million mapped reads per sample
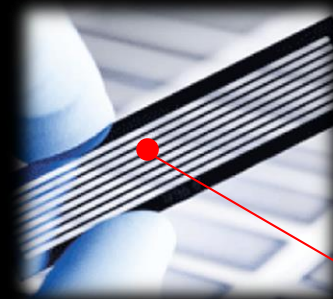
Sims et al., 2014 Nature Reviews

# Experimental design: Biological replicates

Include biological replicates for **increased discovery power** and reduced false positives/negatives!



Two conditions

Include biological replicates for every condition!

- 3+ for cell lines
- 5+ for inbred lines
- 20+ for human samples

# Experimental design: Multiplexing and Randomization

➡ **Multiplexing**: simultaneously measures multiple libraries in one sequencing lane. Unique barcodes are added to label DNA molecules from each library

➡ **Randomization**: Avoid loading samples from the same biological group in the same sequencing lane. Minimizes technical bias and lane-specific effects.





*Sequencing lanes*

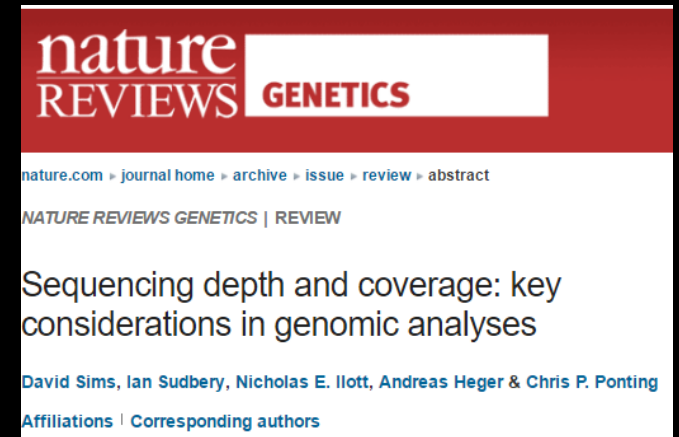# Experimental design: Multiplexing and Randomization

# Challenges and limitations

- Relatively poor RNA quality for tumor FFPE samples

- Contamination from adjacent normal tissue

- Still more expensive than targeted-panel sequencing such as NanoString

- 40 million mapped reads are usually sufficient for gene profiling, but > 80 million are required to detect bottom 1% lowly expressed genes

# Whom to contact if I want to initiate an NGS (e.g. RNAseq) experiment?

➡ The University of Chicago Genomics Facility

Pieter W. Faber, Ph.D.
Technical Director
pfaber@bsd.uchicago.edu
773.834.8420


➡ The Center for Research Informatics Bioinformatics Core

Jorge Andrade, Ph.D.
Director of Bioinformatics
jandrade@bsd.uchicago.edu

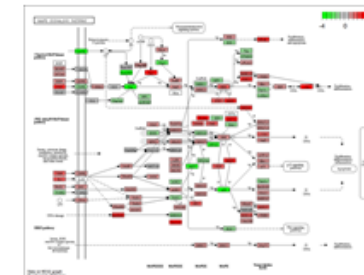# How to perform RNAseq analysis

The good-practice analysis protocol takes 8 major steps.

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

➡ **05-08**: DEG and pathway analysis (06/03, part II)



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

# How to perform RNAseq analysis

# 01 Quality assessment of raw sequencing reads: FastQC

➡ Raw sequencing reads are stored in **FastQ** format (e.g. *KO01.fastq.gz*), where each read is presented by 4 lines

| | |
|---|---|
| Sequence header -- | @SRR1205282.43583628 |
| Sequence string -- | GACTATCTTGGCCAACATGGTGAAACCCCGTCTCTACTAAAAATACAAA |
| Quality header -- | + |
| Quality string -- | CCCFFFFFHHHHHIIIIIIIHHIIIIIIIIIIIIIIIIIIIIIIIIIIII |

➡ Check if the reads are of high quality

➡ Check if any preprocessing step is required (e.g. base trimming, adapter clipping, read filtering)

# 01 Quality assessment of raw sequencing reads: FastQC

- Method
  - FastQC version 0.11.5
  - Scan raw sequencing reads and produce QC reports for evaluation
- Our data
  - Read quality pretty good (baseQ >= 30 in all base positions)
  - Preprocessing is optional

```
fastqc --extract -o $out.dir -t 2 --nogroup $r1.fastq.gz
```

# 01 Quality assessment of raw sequencing reads: FastQC

# 02 Preprocessing: Trimmomatic

➡ Preprocess reads to improve mapping rate and accuracy

➡ Trim low-quality bases, clip adapters, etc.

➡ Avoid over-trimming in RNAseq!

➡ Clean up reads for improved alignment rate and accuracy (for the next step)

# 02 Preprocessing: Trimmomatic

- Method
  - Trimmomatic version 0.36
  - Clip adapters
  - Trim leading/trailing low quality or N bases
  - Trim reads to a specific length
  - Filter out reads of low average quality / of specific length
  - Convert base quality scores between Phred33 and Phred64 FastQ format
- Our data (KO01 as an example)
  - 74,126,025 reads total. Survived: 73,636,793 (99.34%) Dropped: 489,232 (0.66%)

```
java -Xmx4G -jar trimmomatic-0.36.jar SE -threads 4 -phred33
$r1.fastq.gz $r1.trim.fastq.gz ILLUMINACLIP:TruSeq3-
SE.fa:2:30:10 LEADING:5 TRAILING:5 MINLEN:36 SLIDINGWINDOW:4:15
```

➡ Read mapping identifies the location in the genome where a sequencing read comes from

➡ **Splice-aware aligner** (e.g. *STAR*)

Each horizontal bar represents one read. Red/blue indicates reads aligned to plus/minus strand on the genome, respectively.



780 reads map to *FOS* gene

# 03 Map reads to reference genome (GRCh38): STAR

- A read may map to one, or multiple genomic locations (Mapping quality: MAPQ)

- Accurate mapping result is the key for downstream differential gene expression identification

- Read alignment results are stored in **SAM or BAM (Binary SAM)** format

# 03 Map reads to reference genome (GRCh38): STAR

- Method

  - STAR version 2.5.3a

  - splice-aware aligner

  - Ultrafast (~15 minutes for 50m reads), require lots of memory (~36GB for human genome)

  - Flexible options to allow canonical/non-canonical junctions, with/wo known gene annotations, etc.

- Different aligners may generate very different results! (Engström et al. *Nature methods* 2013. Systematic evaluation of spliced alignment programs for RNA-seq data)

```
STAR --runMode alignReads --genomeLoad NoSharedMemory --
outFileNamePrefix $out.prefix --readFilesCommand zcat --
genomeDir $refgenome.dir --readFilesIn $r1.trim.fastq.gz --
runThreadN 2 --outSAMstrandField intronMotif --
outFilterIntronMotifs RemoveNoncanonicalUnannotated --
outSAMtype BAM SortedByCoordinate
```

# 04 Quantify transcript abundance: featureCounts

- Estimate number of reads mapped to gene features (e.g. gene, exon, etc.)

- Method

  - featureCounts version 2.5.2b

  - Ultrafast (~10 minutes for 50m reads), require low amount of memory (~4GB for human genome)

  - Flexible options to count the reads based on specific mapping criteria or study purposes

    - Gene-level or exon-level (for isoforms); Uniquely mapped reads; Primary alignment; Properly paired reads (if reads are paired-end); Mapping quality thresholds

- Choose the option accordingly based on your experimental design!

```
featureCounts -s 0 -F GTF -t exon -g gene_id -Q 255 -J --primary -a
$refgeneanno.gtf -T 2 -G $refgenome.fa -o
$sample.star.featurecounts.raw_counts.txt $sample.star.merged.bam
```

```
//========================== featureCounts setting ==========================\\   //========================== Running ==========================\\
||                                                                            ||   ||                                                                            ||
||             Input files : 1 BAM file                                       ||   || Load annotation file gencode.v24.primary_assembly.annotation.chr11.gtf     ||
||                           S results/rnaseq/DLBC_samples_grch38/KO01/al ...  ||   ||    Features : 1177678                                                       ||
||                                                                            ||   ||    Meta-features : 60725                                                    ||
||             Output file : results/rnaseq/DLBC_samples_grch38/KO01/read ...  ||   ||    Chromosomes/contigs : 59                                                 ||
||             Annotations : gencode.v24.primary_assembly.annotation.chr11...  ||   ||                                                                            ||
||       Junction Counting : <output_file>.jcounts                            ||   || Loading FASTA contigs : GRCh38.primary_assembly.genome.chr11.fa            ||
||                                                                            ||   ||    194 contigs were loaded                                                  ||
||                 Threads : 12                                               ||   ||                                                                            ||
||                   Level : meta-feature level                               ||   || Process BAM file results/rnaseq/DLBC_samples_grch38/KO01/alignment/KO0 ...  ||
||              Paired-end : yes                                              ||   ||    Single-end reads are included.                                          ||
||         Strand specific : no                                              ||   ||    Assign reads to features...                                             ||
||      Multimapping reads : primary only                                     ||   ||    Total reads : 93186698                                                   ||
||  Multi-overlapping reads : not counted                                     ||   ||    Successfully assigned reads : 55910832 (60.0%)                           ||
||        Read orientations : fr                                              ||   ||    Running time : 10.23 minutes                                            ||
||                                                                            ||   ||                                                                            ||
||          Chimeric reads : not counted                                     ||   ||             Found 230830 junctions in all the input files.                 ||
||         Both ends mapped : not required                                    ||   ||                                                                            ||
||                                                                            ||   ||                  Read assignment finished.                                 ||
||                                                                            ||   ||                                                                            ||
\\==================== http://subread.sourceforge.net/ ====================//   ||                                                                            ||
                                                                                 \\==================== http://subread.sourceforge.net/ ====================//
```

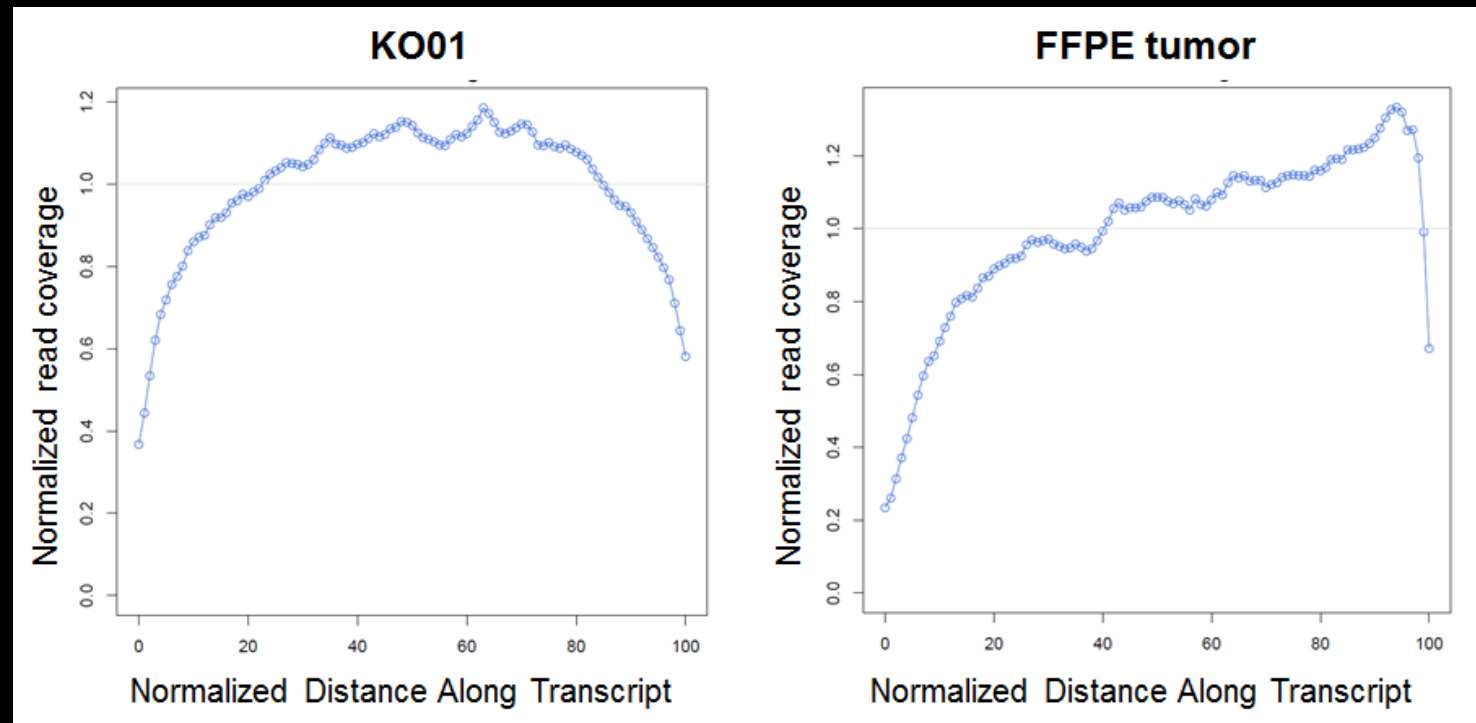| Geneid | Chr | Start | End | Strand | Length | KO01 |
|--------|-----|-------|-----|--------|--------|------|
| ENSG00000019485.12 | chr11;chr11;chr11;chr11;chr11;chr11;chr11;chr11;chr11;chr11;chr11;chr11;chr11 | 45095806;45146675;45147343;45181761;45182246;45182861;45204711;45208874;45212469;45219570;45224217;45225995 | 45095901;45146877;45147581;45181885;45182349;45183123;45204778;45209295;45214713;45219757;45225087;45235124 | +;+;+;+;+;+;+;+;+;+;+;+ | 13954 | 710 |
| ENSG00000170345.9 | chr14;chr14;chr14 | 75278774;75279237;75279643 | 75279128;75279531;75282230 | +;+;+ | 3238 | 780 |
| ENSG00000119660.4 | chr14 | 75292131 | 75292495 | - | 365 | 0 |
| ENSG00000259687.1 | chr14;chr14;chr14 | 75294404;75294677;75296120 | 75294541;75294737;75296638 | +;+;+ | 718 | 0 |
| ENSG00000259319.1 | chr14 | 75423683 | 75427741 | - | 4059 | 25 |

# RNAseq quality metrics & coverage: Picardtools, bedtools, RSeQC

➡ Evaluate the quality of reads and alignments

➡ Identify potential problems regarding (RNA) sample quality

  ➡ Is the RNA highly degraded?

  ➡ Is there high-level genomic DNA contamination?

  ➡ Was ribosome RNA successfully depleted?

  ➡ How do reads distribute on the genome? (exons, introns, intergenic, etc.)

  ➡ Is the strandness of read alignment consistent with library type? non-stranded or forward/reverse strand-specific libs

  ➡ Does the target gene that was knocked down in KO samples have reduced expression, compared to WT?

# a) RNA degradation

➡ RNA quality: Fresh samples (e.g. cell line) > frozen samples (e.g. mouse tissues) >> FFPE samples (e.g. human tumors)

➡ During lib prep, RNA quality is inferred by The RNA integrity number (RIN) evaluated using the 28S to 18S rRNA ratio (e.g. RIN > 4.5)

➡ However, studies have shown that RIN can be quite inaccurate for FFPE samples

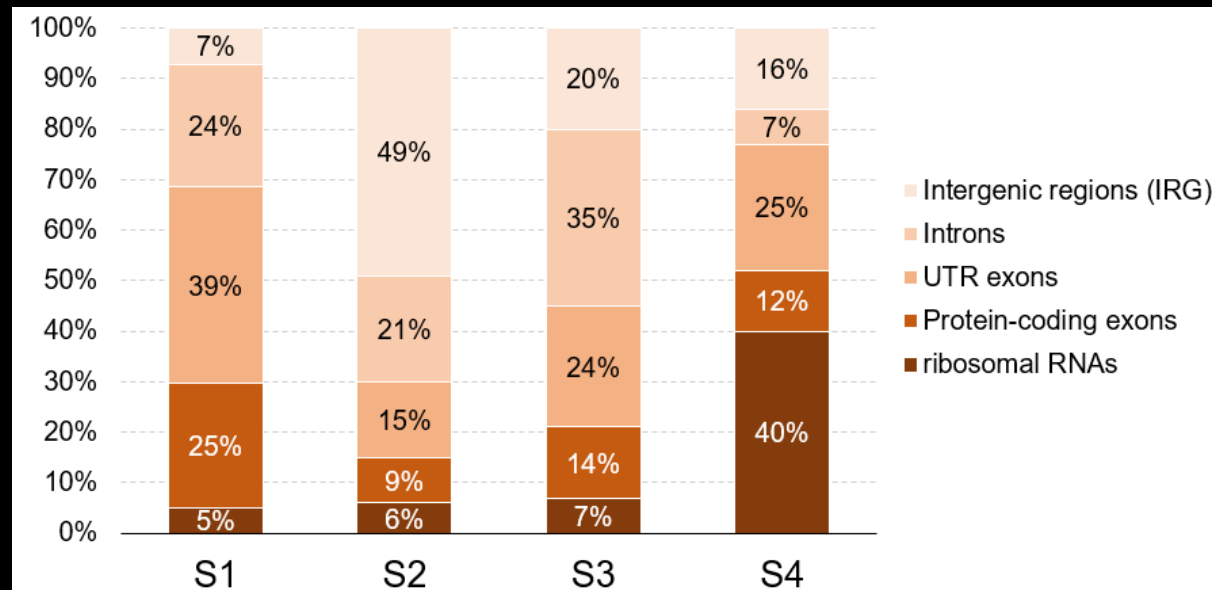➡ *Gene body coverage plot* (Picardtools). FFPE samples often has 5' degradation

# b) Genomic DNA contamination

- In lib prep, DNase digestion removes genomic DNA

- While fresh samples often has good-quality RNAs, highly-degraded samples (e.g. FFPE tumors) often has a higher degree of genomic DNA contamination

- Sometimes DNA contamination could occupy 70% of sequencing reads, greatly reducing the discovery power of DEG analysis

- Good assessment to identify and estimate genomic DNA contamination includes

  - Read distribution in genomic features: high fraction of intergenic reads indicates DNA contamination

  - Visualize intergenic region in genome browser (e.g. IGV)

# b) Genomic DNA contamination

**Q1: Which sample (S1-4) has the most severe genomic DNA contamination?**

*Hint: higher percentage of intergenic reads indicates more severe DNA contamination in RNA samples*
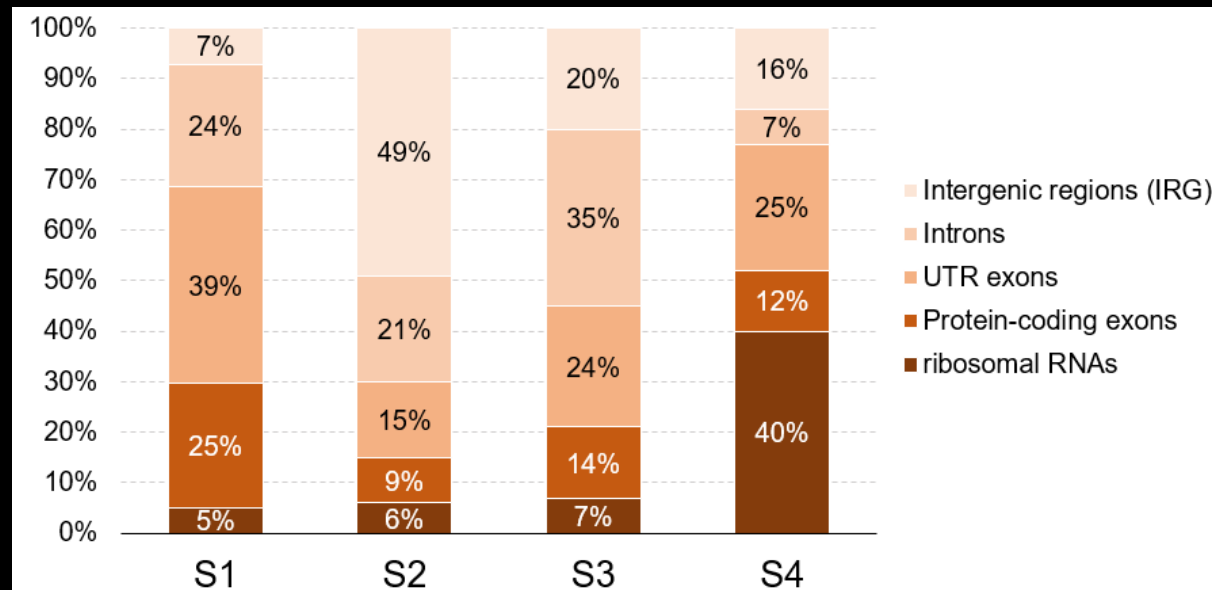
# c) Ribosome RNA fraction

- rRNA accounts for > 80% of the transcriptome

- Mmost RNAseq lib prep protocol includes an ribosomal RNA depletion step

- However, if RNA quality is relatively poor (e.g. FFPE tumors), rRNA depletion often is efficient

- Accessing if the depletion step is successful through read distribution in genomic features

- How many reads map to ribosome RNA regions on the genome?

# c) Ribosome RNA fraction

**Q2: Which sample (S1-4) has the least efficient depletion of ribosome RNAs?**

*Hint: rRNAs account for > 80% of the whole transcriptome. If not removed, the majority of the sequencing reads will be derived from rRNA*

# c) Library standness: fraction of correctly oriented reads



RSeQC: infer_experiment.py

For pair-end RNA-seq, there are two different ways to strand reads (such as Illumina ScriptSeq protocol):

1. 1++,1–,2+-,2-+

- read1 mapped to '+' strand indicates parental gene on '+' strand
- read1 mapped to '-' strand indicates parental gene on '-' strand
- read2 mapped to '+' strand indicates parental gene on '-' strand
- read2 mapped to '-' strand indicates parental gene on '+' strand

2. 1+-,1-+,2++,2–

- read1 mapped to '+' strand indicates parental gene on '-' strand
- read1 mapped to '-' strand indicates parental gene on '+' strand
- read2 mapped to '+' strand indicates parental gene on '+' strand
- read2 mapped to '-' strand indicates parental gene on '-' strand

For single-end RNA-seq, there are also two different ways to strand reads:

1. ++,–

- read mapped to '+' strand indicates parental gene on '+' strand
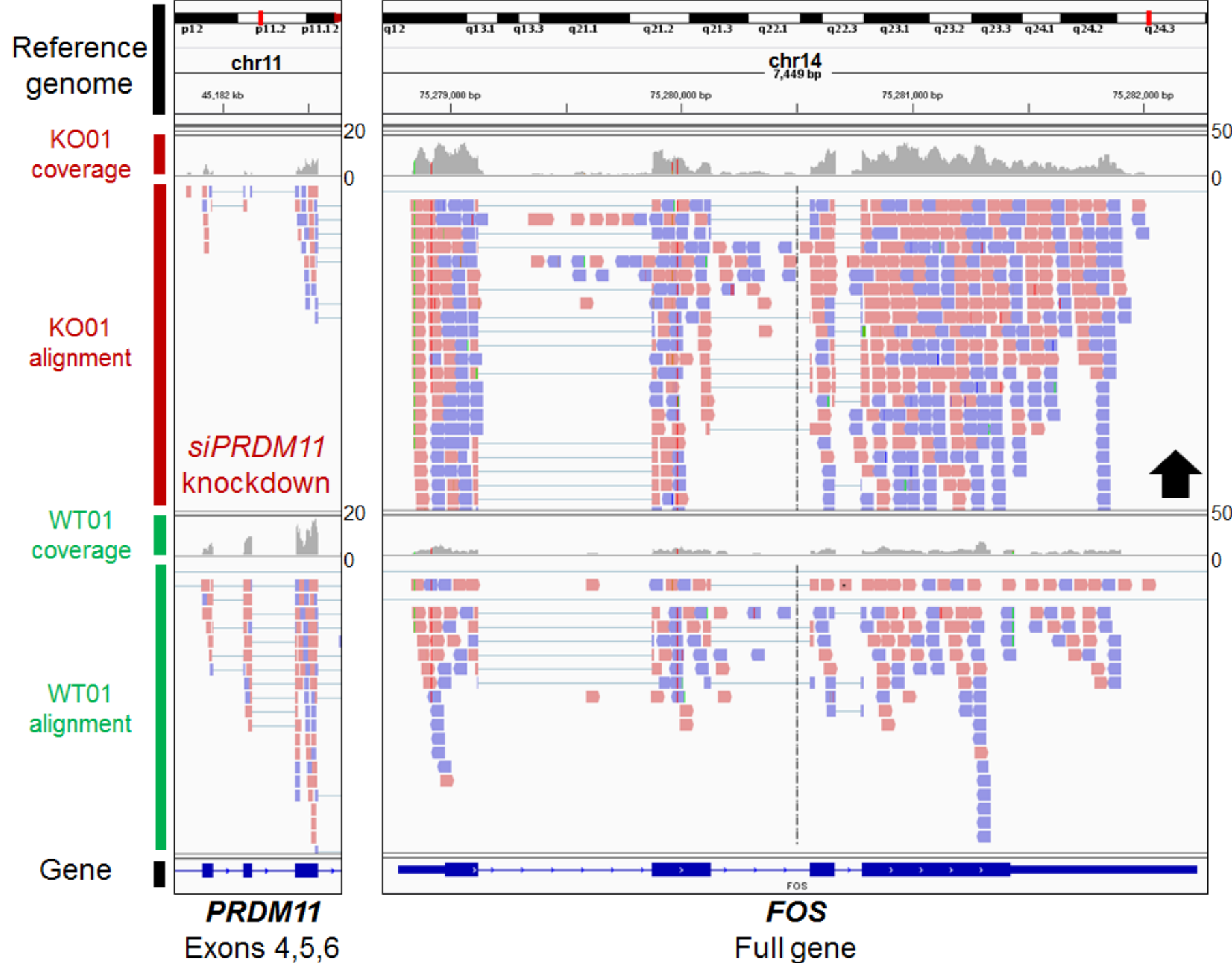- read mapped to '-' strand indicates parental gene on '-' strand

2. +-,-+

- read mapped to '+' strand indicates parental gene on '-' strand
- read mapped to '-' strand indicates parental gene on '+' strand

# e) Confirmation of reduced/increased expression for knockdown/overexpressed genes in KO/OE samples
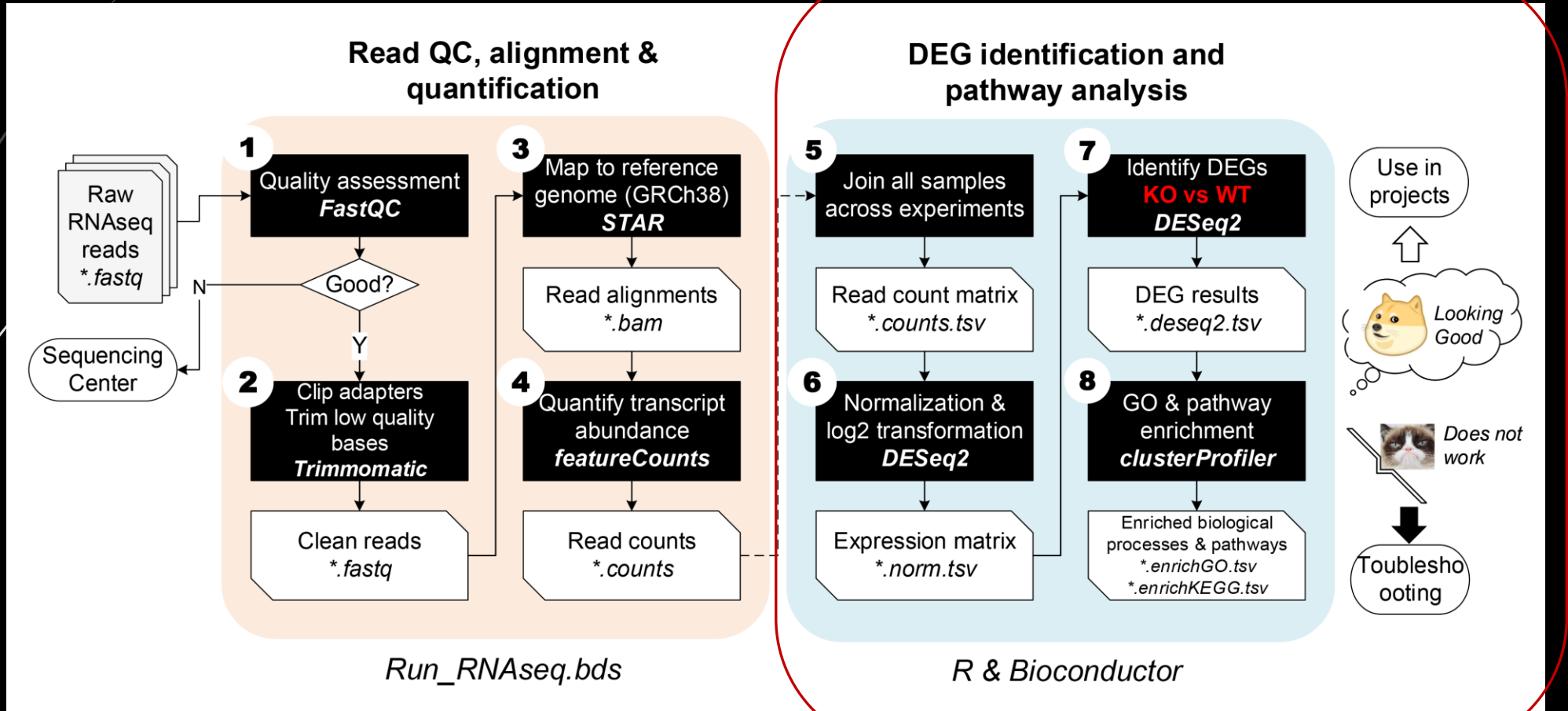
- Is there a gene that is expected/known to be repressed or overexpressed?

- Does the RNAseq result reflect expression change of this gene?

- In our data, *PRDM11* expression is expected to be down since this is the gene knocked down in U2932 cells

- NOT all exons are affected!

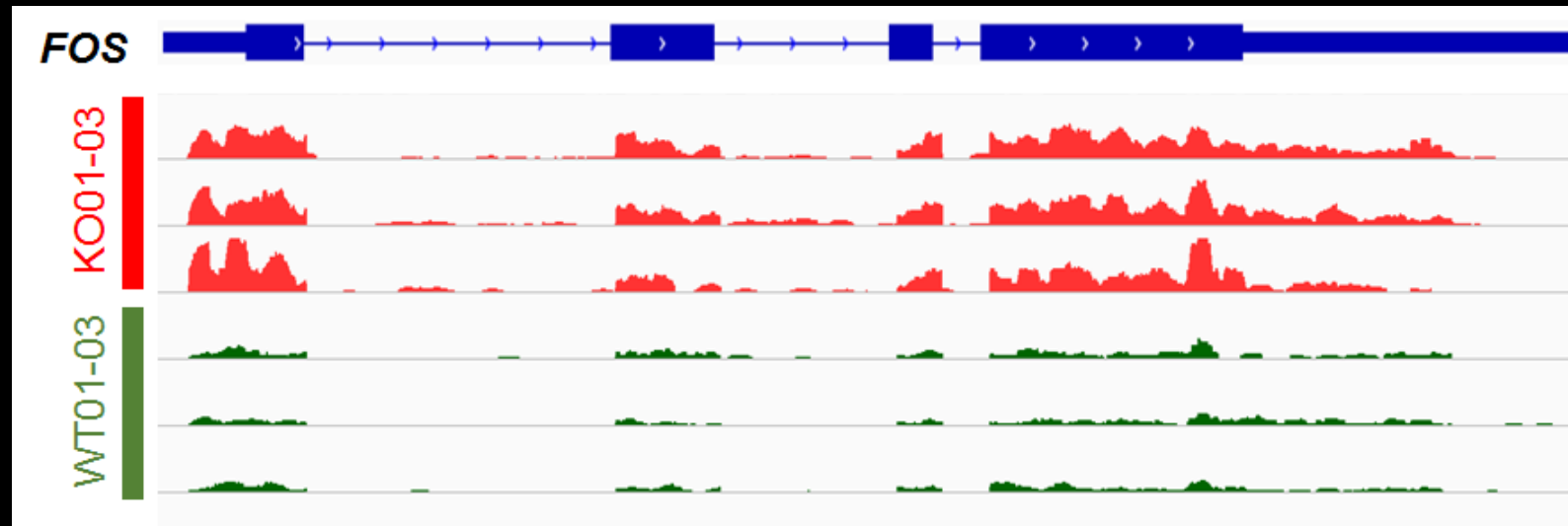- Only those targeted by siRNA (small RNA interference) will be affected

**Our data**

- *PRDM11* knockdown U2932 cells in triplicates (KO01-03 vs WT01-03)
- NOT the full PRDM11 gene is knockdown!
- siRNAs target exons 4,5,6 of PRDM11, thus only those three exons show a reduction of expression in the KO samples; other exons are not affected
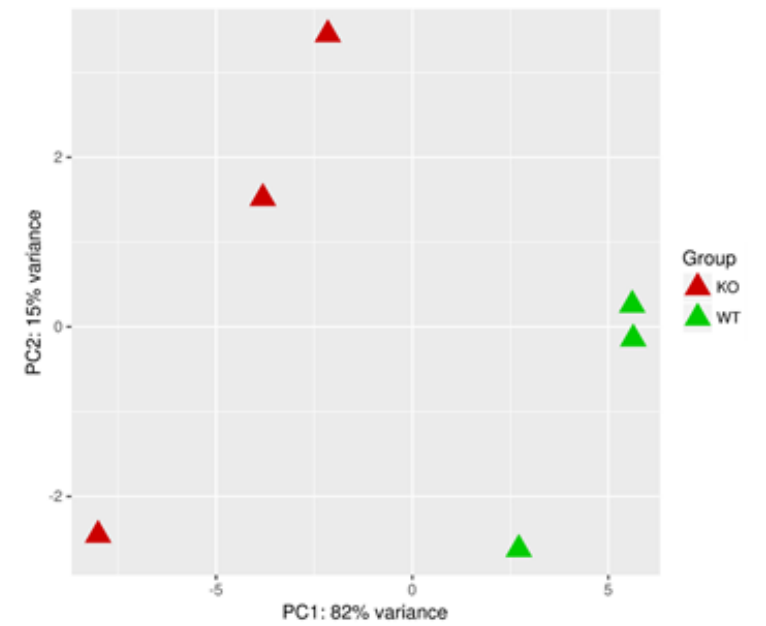- PRDM11 knockdown leads to upregulation of FOS expression

Figure labels:

Reference genome — chr11 — chr14 7,449 bp

KO01 coverage

KO01 alignment — *siPRDM11* knockdown

WT01 coverage

WT01 alignment

Gene

*PRDM11* Exons 4,5,6
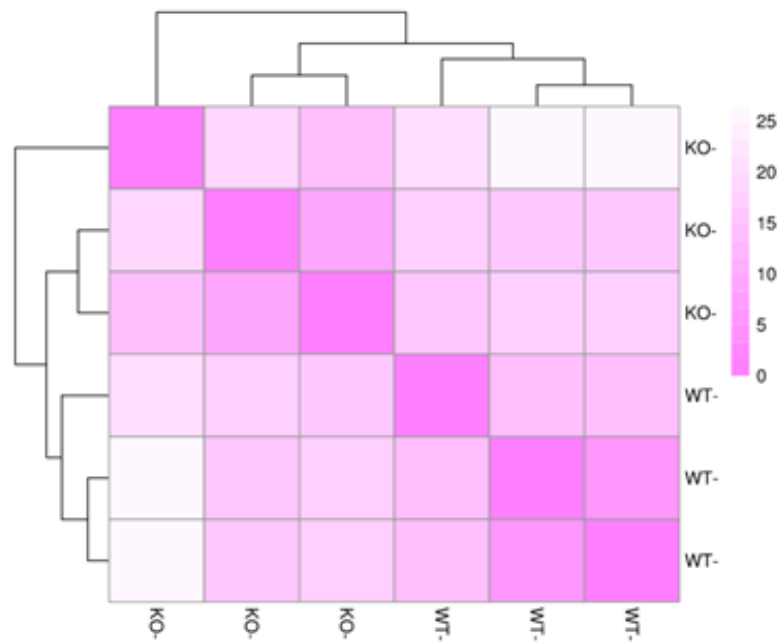
*FOS* Full gene

# How to perform RNAseq analysis

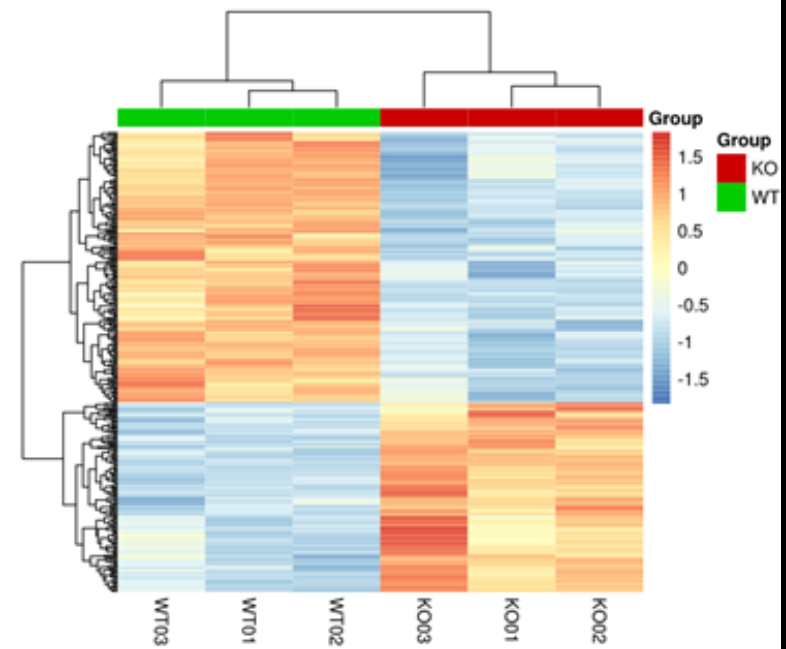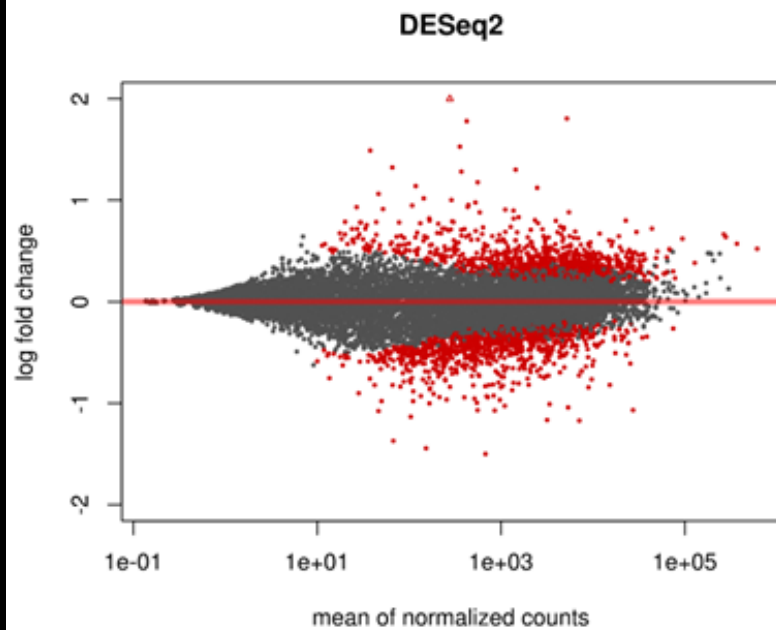# 05-08: Identify differentially expressed genes and pathways: DESeq2, clusterProfiler

- 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

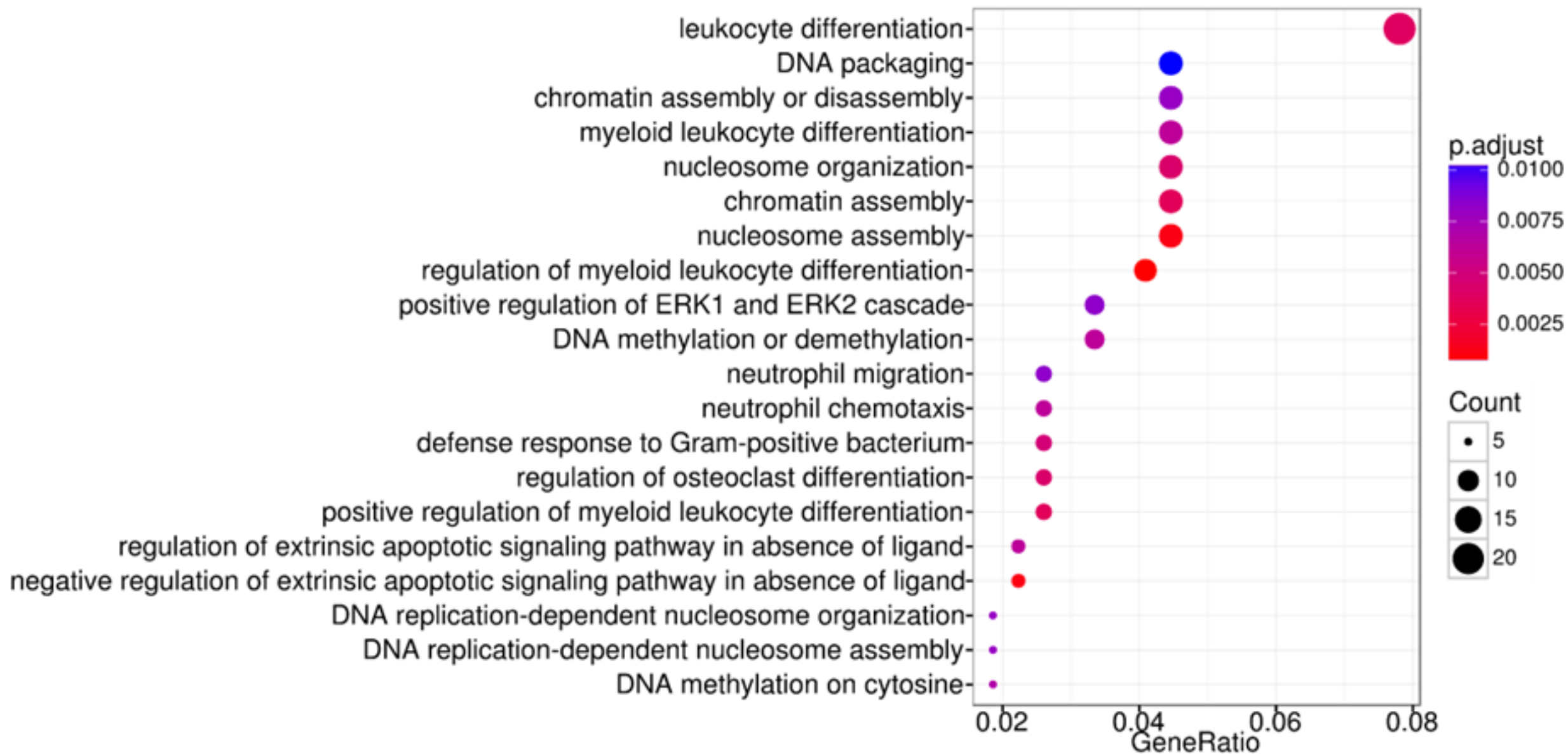- For example, if we visualize *FOS* gene across 6 samples in genome browser



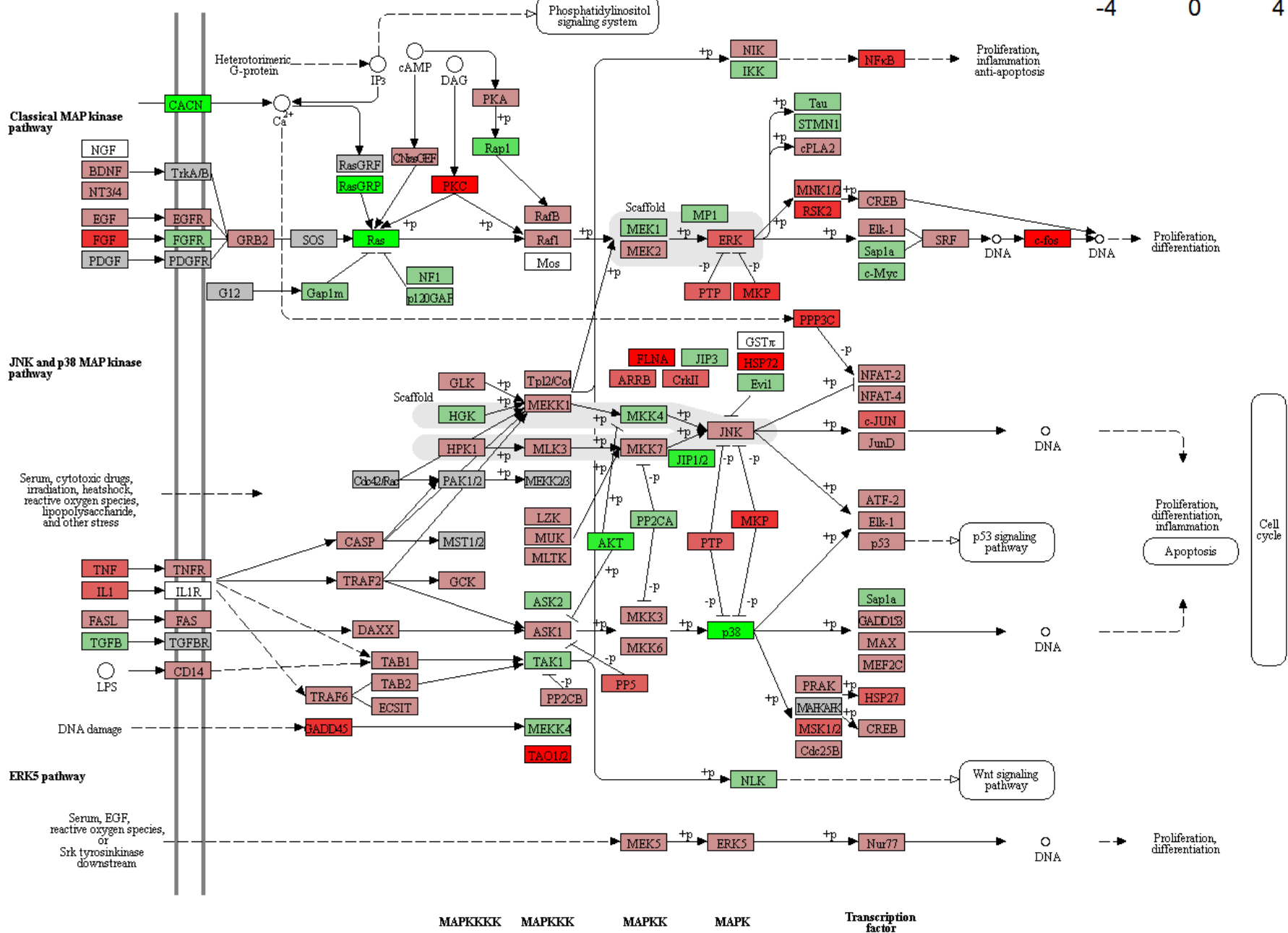*FOS = Fos proto-oncogene, AP-1 transcription factor subunit*

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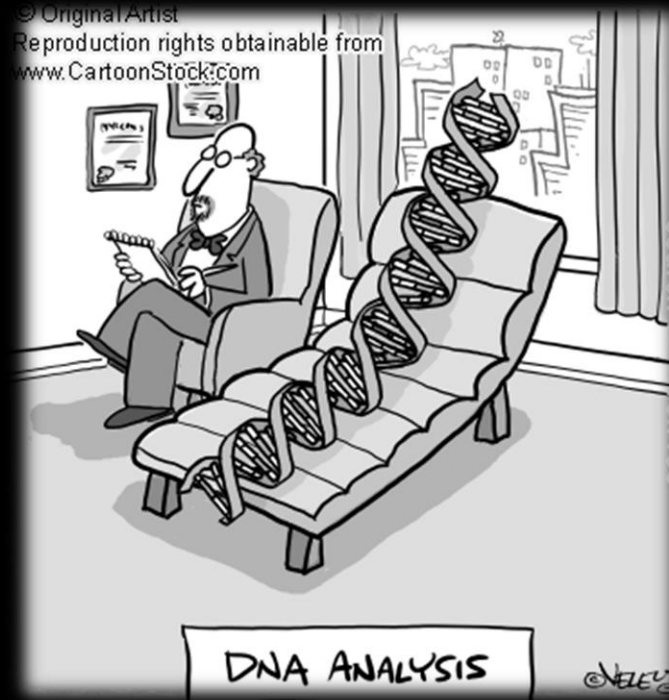[1] "Heatmap = 296 genes on the row, 6 samples on the column"

DESeq2
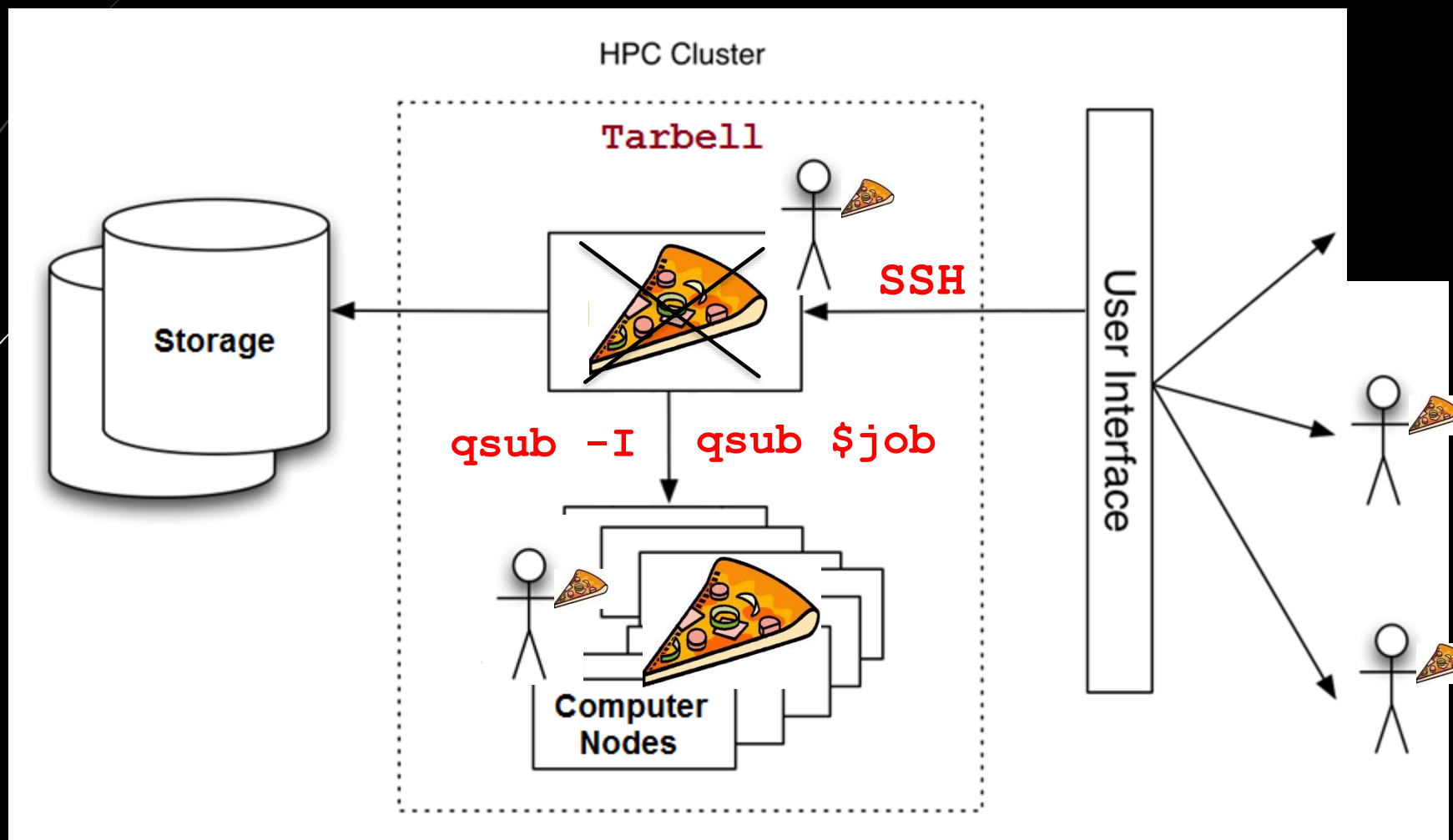
# Thank you!



Questions

# Hands-on practice START

➤ **Open your handson.Rmd on the Github**

➤ *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*

*Job running is prohibited on the head node!*

# Ten essential Linux commands

**cd**        Change directory

**ls**        List contents

**cp**        Copy

**mv**        Move/rename

**rm**        Delete

**pwd**       Print the current path

**head**      Show the first few lines of a file

**more**      View a file by page

**mkdir**     Create a new directory

**man**       Show help menu for a command

# Five useful HPC/Shell commands

**qsub**     Submit a job

**qstat**    List submitted jobs from the user

**qdel**     Delete a job

**showq**    Show all jobs on the cluster

**ssh**      Log into a server using Secure Shell (SSH)