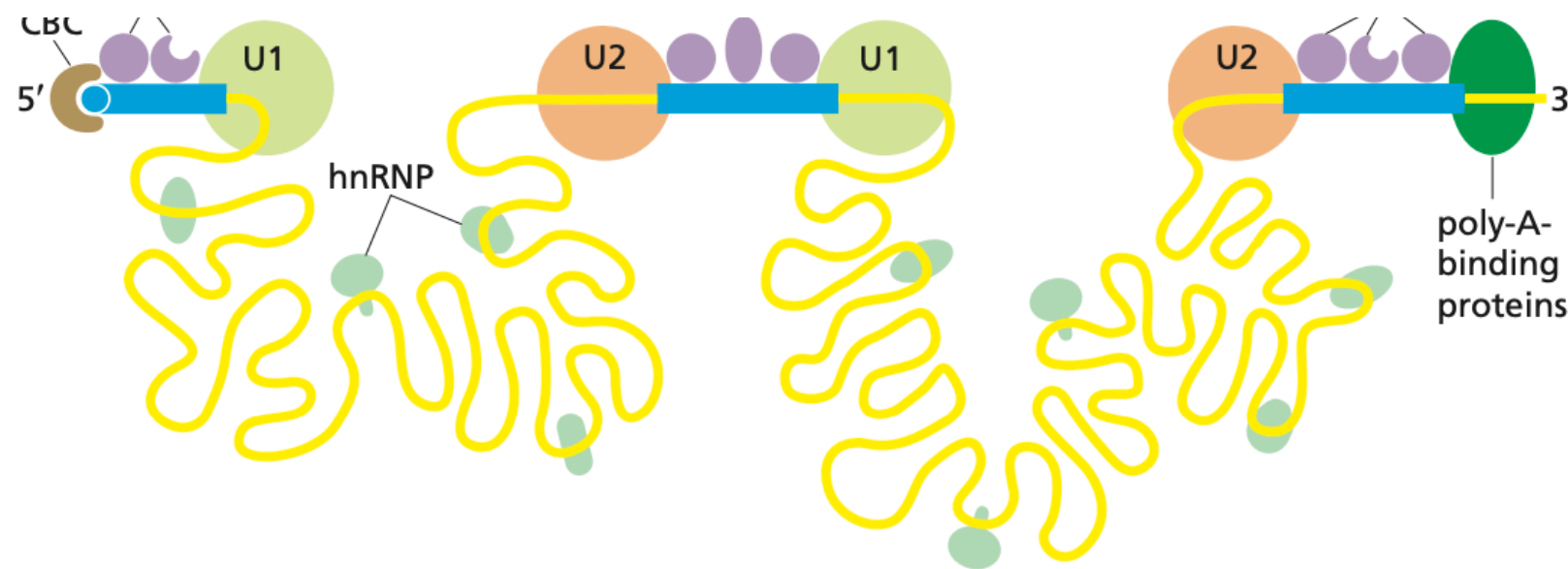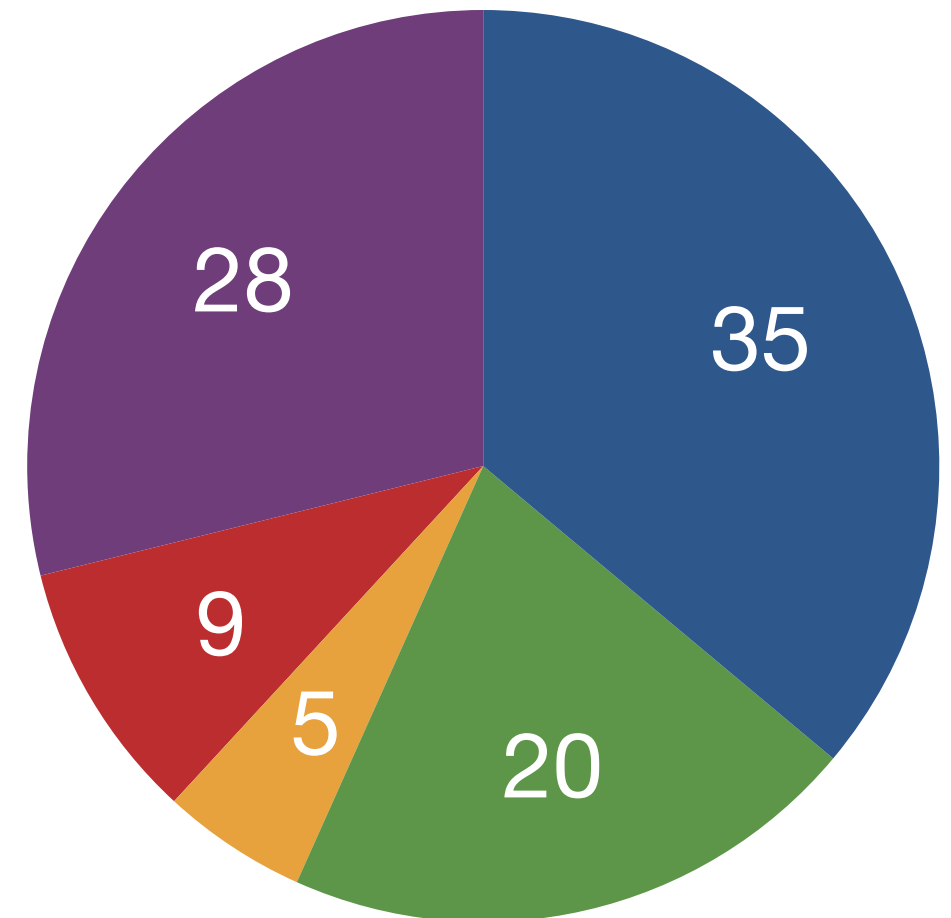# RNA sequencing

2025-03-25

# Sequencing transcriptome



- Evaluate **expression** of genes/transcripts for:
  - All species of RNA
  - mRNA
  - small RNAs

- Evaluate expression levels of exons
  - Patterns of alternative splicing

- Evaluate transcriptional **alterations**
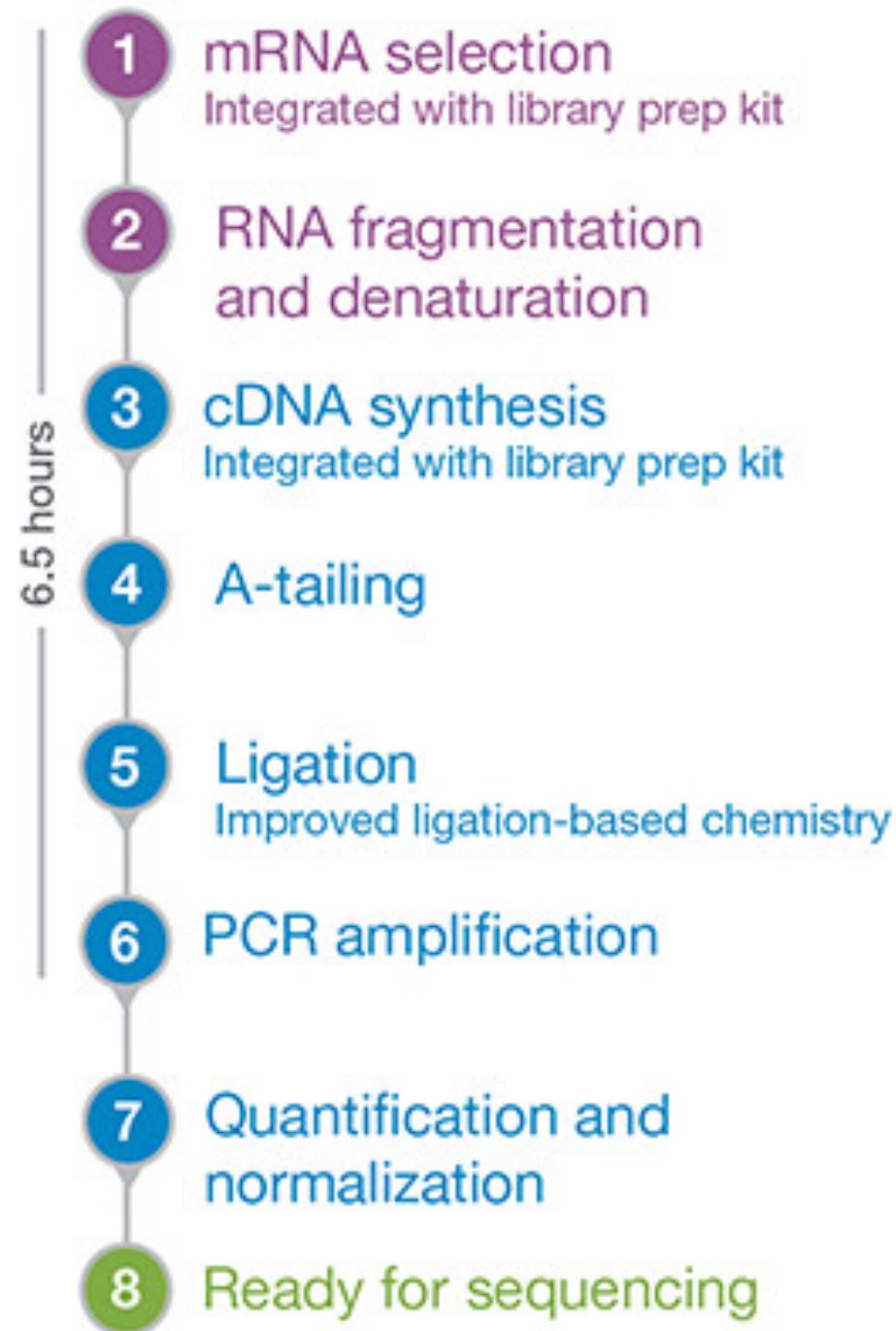
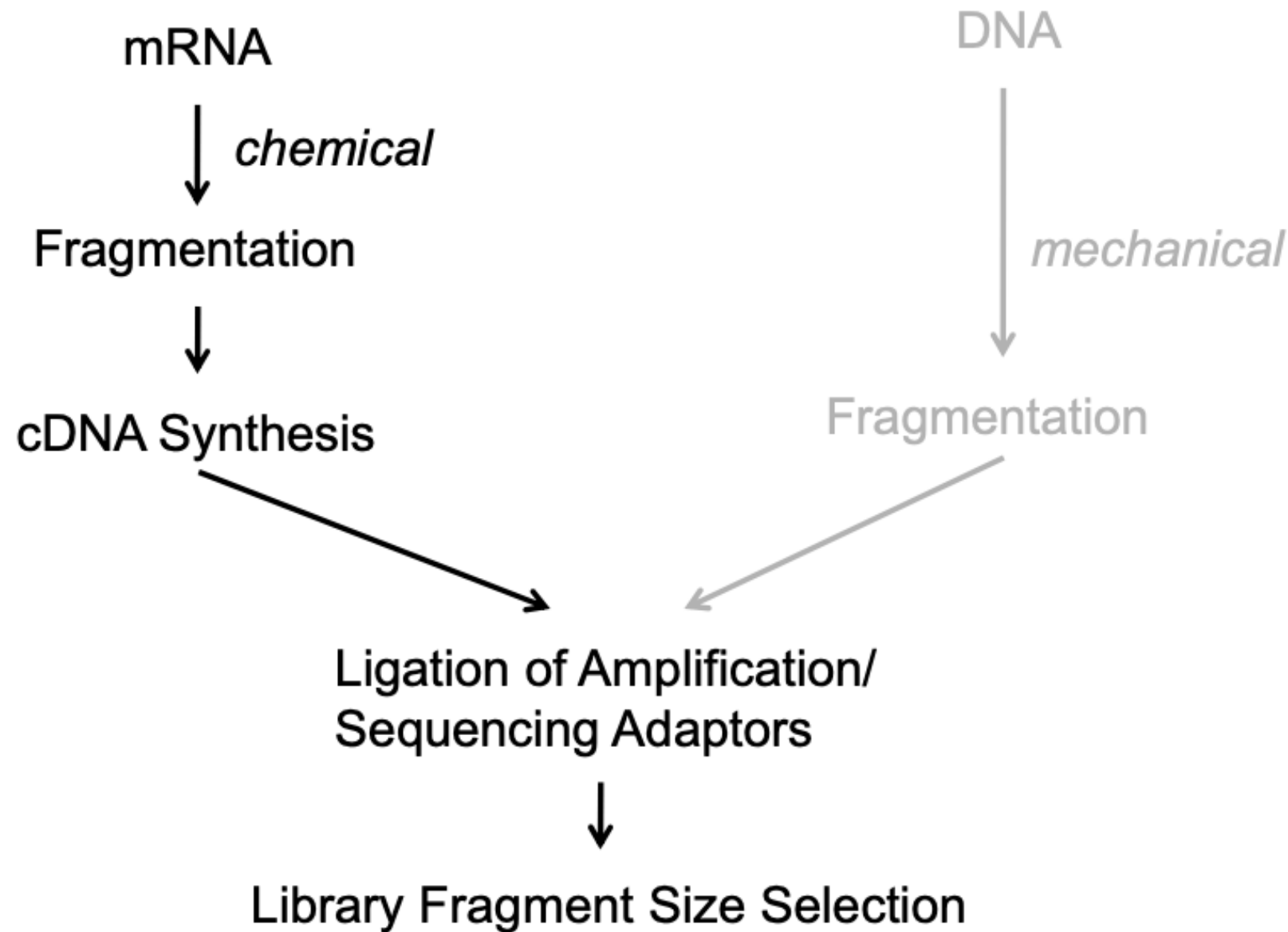- Annotate **regions** and **functional elements**

| RNA Transcription | RNA-Protein Interactions | RNA Modifications | RNA Structure | Low-Level RNA Detection |
|---|---|---|---|---|
| RNA-Seq | Ribo-Seq | MeRIP-Seq | SHAPE-Seq | scRNA-Seq |
| CaptureSeq | RIP-Seq | miCLIP-m6A | icSHAPE | SUPeR-Seq |
| RASL-Seq | CLIP-Seq | PSI-Seq | CIRS-Seq | UMI |
| ClickSeq | Pol II CLIP | Pseudo-Seq | SHAPE-MaP | Digital RNA Sequencing |
| 3Seq | miR-CLIP | ICE | DMS-Seq | MARS-Seq |
| cP-RNA-Seq | eCLIP | | SPARE | Quartz-Seq |
| 3P-Seq | irCLIP | | PARS-Seq | DP-Seq |
| 2P-Seq | PAR-CLIP | | Cap-Seq | Smart-Seq |
| 3'-Seq | iCLIP | | CIP-TAP | FRISCR |
| TIF-Seq | BrdU-CLIP | | | CEL-Seq |
| PEAT | AGO-CLIP | | | STRT-Seq |
| SMORE-Seq | PIP-Seq | | | TCR Chain Pairing |
| TL-Seq | hiCLIP | | | TCR-LA-MC PCR |
| TATL-Seq | RBNS | | | CirSeq |
| RARseq | TRIBE | | | TIVA |
| TAIL-Seq | HiTS-RAP | | | PAIR |
| PAL-Seq | TRAP-Seq | | | CLaP |
| FRT-S wcell | DLAF | | | CytoSeq |
| ChIRP | miTRAP | | | Drop-Seq: |
| CHART | CLASH | | | Hi-SCL |
| RAP | | | | InDrop |
| GRO-seq | | | | snRNA-Seq |
| Bru-Seq | | | | Nuc-Seq |
| BruChase-Seq | | | | Div-Seq |
| 5'-GRO-Seq | | | | SCRB-Seq |
| BruDRB-Seq | | | | G&T-Seq |
| 4sUDRB-Seq | | | | scM&T-Seq |
| PRO-Seq | | | | scTrio-seq |
| PRO-Cap | | | | |
| CAGE | | | | |
| 3'NT Method | | | | |
| NET-Seq | | | | |
| mNET-Seq | | | | |
| PARE-Seq | | | | |
| GMUCT | | | | |

# illumina®

## Protocolli di sequenziamento



- 35 — RNA transcription
- 20 — RNA-Protein interactions
- 5 — RNA modifications
- 9 — RNA structure
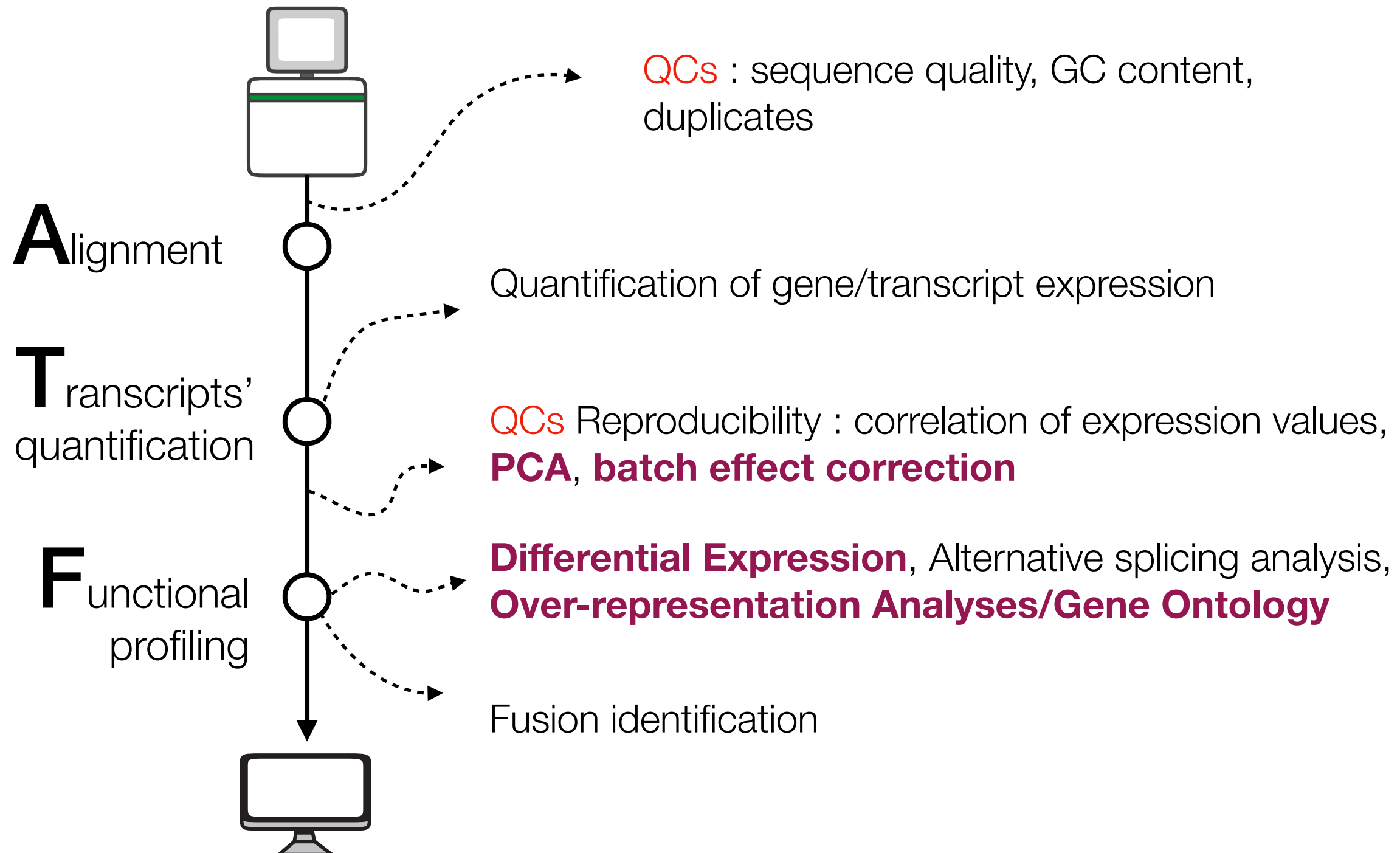- 28 — Low-level RNA detection

C.G.B.

UNIVERSITÀ DEGLI STUDI DI MILANO

# Key steps in sequencing

# Deciphering gene expression

RNA-seq data analysis workflow:



**A**lignment

**T**ranscripts' quantification

**F**unctional profiling

QCs : sequence quality, GC content, duplicates

Quantification of gene/transcript expression

QCs Reproducibility : correlation of expression values, **PCA**, **batch effect correction**

**Differential Expression**, Alternative splicing analysis, **Over-representation Analyses/Gene Ontology**

Fusion identification

C.G.B.

UNIVERSITÀ DEGLI STUDI DI MILANO

# Deciphering gene expression



Raw reads
.fastq

Mapping
STAR

Aligned reads
.sam/.bam

Counting
featureCounts

Read count table
.txt

Descriptive plots

Normalizing
DESeq2

Normalized read count table
.Robj

DE test & multiple testing correction
DESeq2

List of fold changes & statistical values
.Robj, .txt

Filtering
Customized scripts

Downstream analyses on DE genes

C.G.B.

UNIVERSITÀ
DEGLI STUDI
DI MILANO

# Gene expression level distribution



https://academic.oup.com/nar/article/48/4/1730/5691219

# Gene expression level distribution

Samples

$a_1$ $a_2$ $a_3$



Genes

$g_1$
$g_2$
$g_3$
...
$g_z$

High
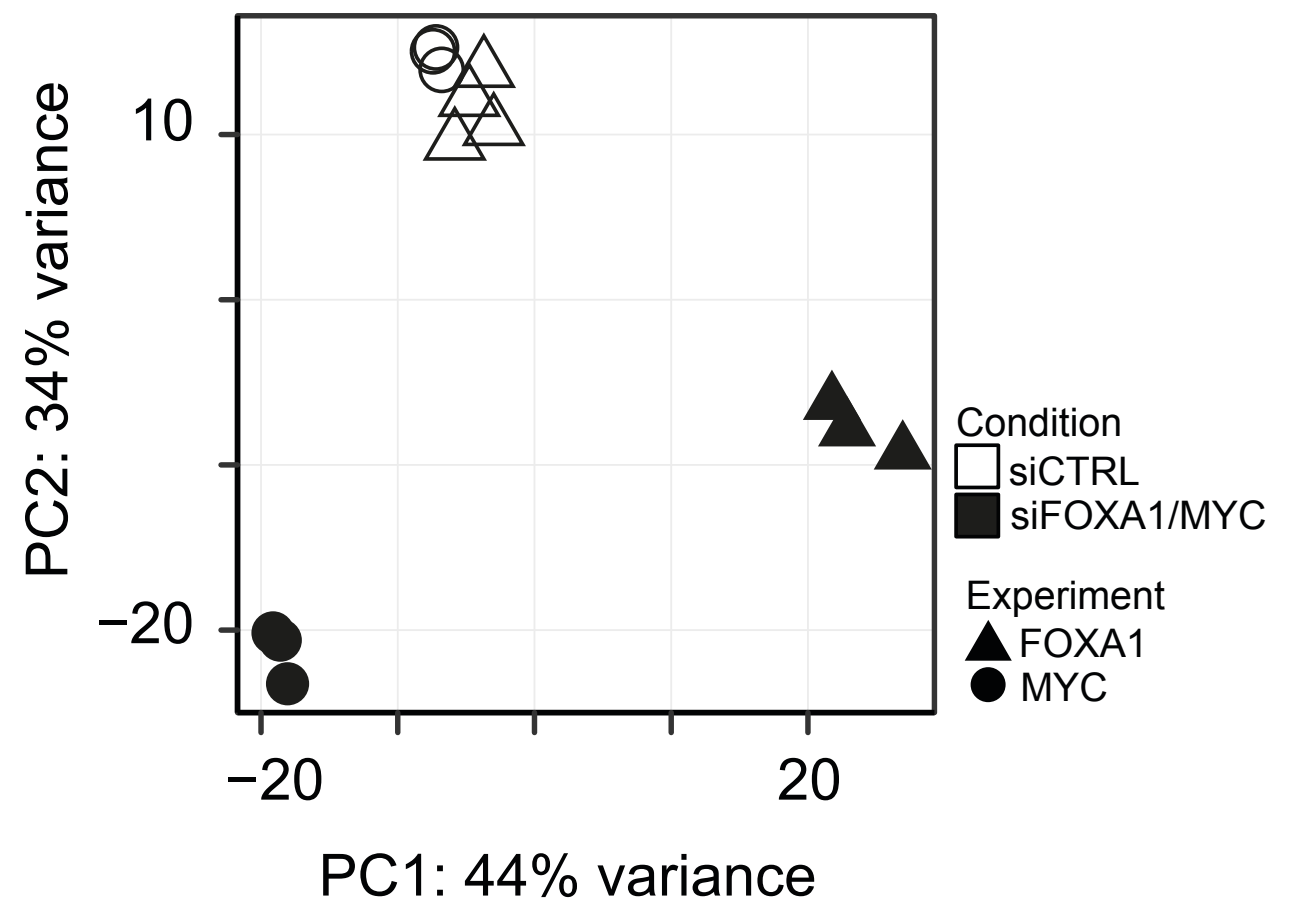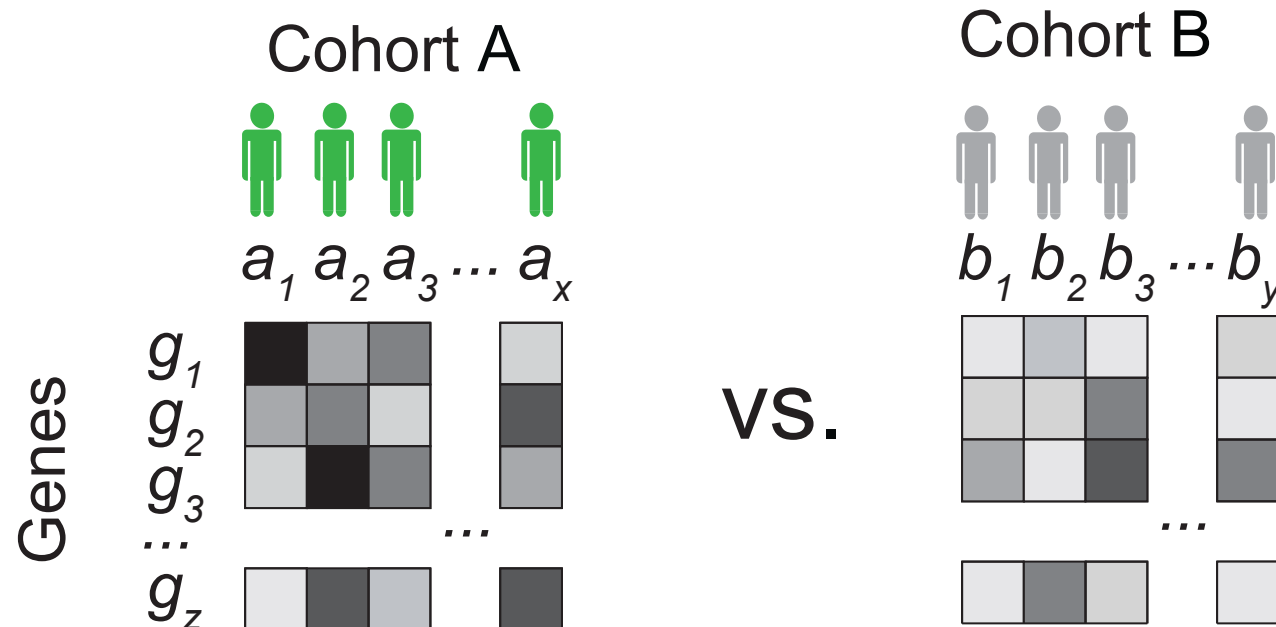
Expression

Low

**Reproducibility**

Principal Component Analysis (PCA)

PCA is a statistical technique for dimensionality reduction. We use PCA when a dataset presents a high number of features (genes in this case). It is like compressing information about ~20,000 in two dimensions or some more if we need it.
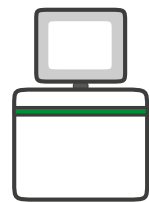


PC2: 34% variance

10

−20

Condition
☐ siCTRL
■ siFOXA1/MYC

Experiment
▲ FOXA1
● MYC

−20                    20

PC1: 44% variance

C.G.B.

UNIVERSITÀ
DEGLI STUDI
DI MILANO

# Differential expression



Cohort A

$a_1 \; a_2 \; a_3 \; ... \; a_x$

Genes $g_1$ $g_2$ $g_3$ ... $g_z$

vs.

Cohort B

$b_1 \; b_2 \; b_3 \; ... \; b_y$

Two are the main goals of a differential expression (DE) analysis:

1. Estimate the **entity of variation** between the two conditions, i.e. calculate Fold Change (FC)

2. Estimate the **significance of the difference**, i.e. p-value, and correct it for multiple testing (p-adjusted).
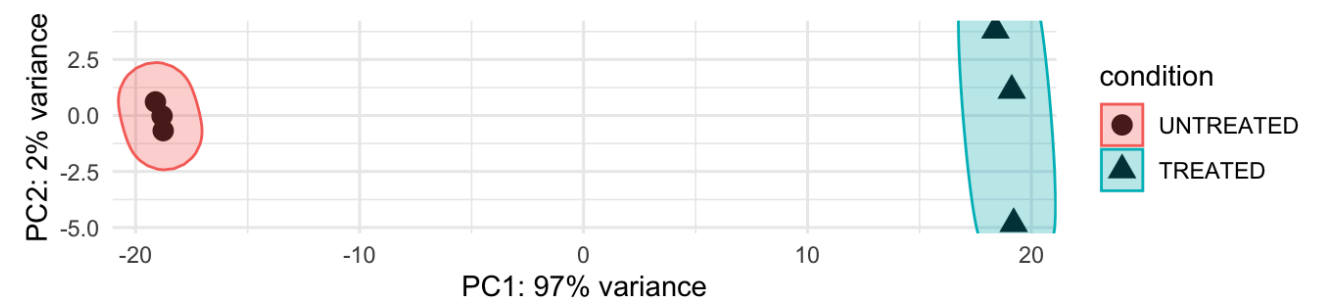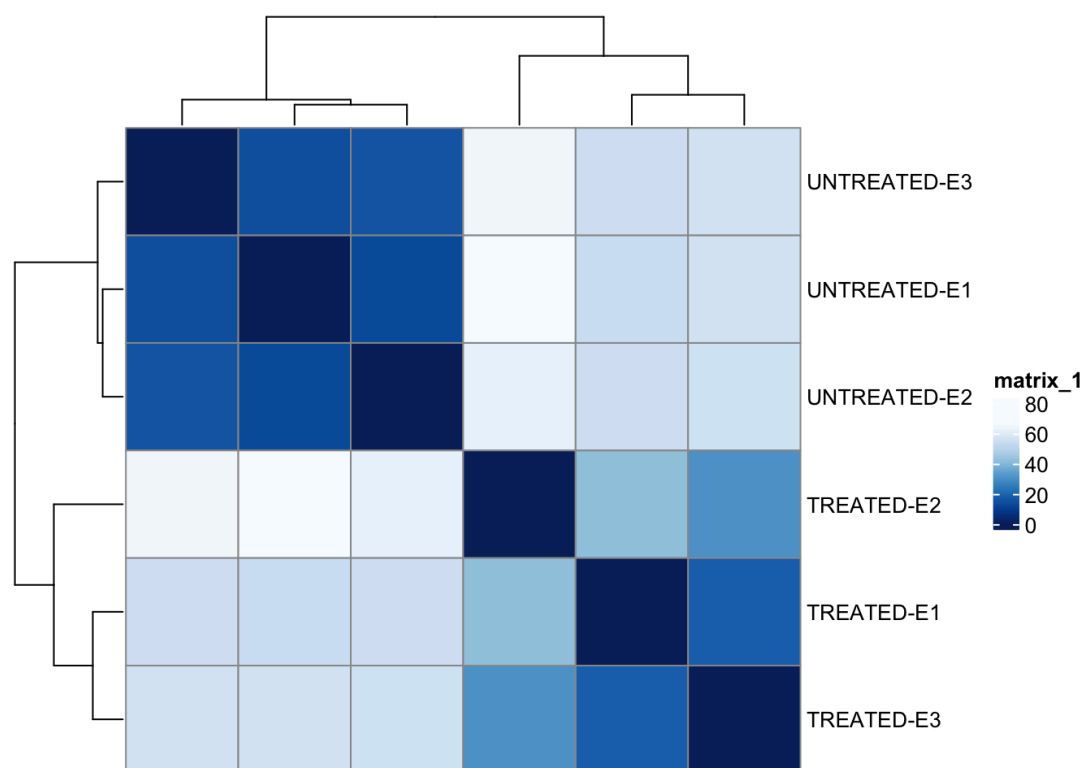
C.G.B.

UNIVERSITÀ DEGLI STUDI DI MILANO

# DESeq2

# Normalization

Normalising data is fundamental. If we skip this step we introduce biases in our analysis.

| Normalization method | Description | Accounted factors | Recommendations for use |
|---|---|---|---|
| **CPM** (counts per million) | counts scaled by total number of reads | sequencing depth | gene count comparisons between replicates of the same samplegroup; **NOT for within sample comparisons or DE analysis** |
| **TPM** (transcripts per kilobase million) | counts per length of transcript (kb) per million reads mapped | sequencing depth and gene length | gene count comparisons within a sample or between samples of the same sample group; **NOT for DE analysis** |
| **RPKM/FPKM** (reads/fragments per kilobase of exon per million reads/fragments mapped) | similar to TPM | sequencing depth and gene length | gene count comparisons between genes within a sample; **NOT for between sample comparisons or DE analysis** |
| DESeq2's **median of ratios** [1] | counts divided by sample-specific size factors determined by median ratio of gene counts relative to geometric mean per gene | sequencing depth and RNA composition | gene count comparisons between samples and for **DE analysis; NOT for within sample comparisons** |
| EdgeR's **trimmed mean of M values (TMM)** [2] | uses a weighted trimmed mean of the log expression ratios between samples | sequencing depth, RNA composition, and gene length | gene count comparisons between and within samples and for **DE analysis** |

https://docs.gdc.cancer.gov/Data/Bioinformatics_Pipelines/Expression_mRNA_Pipeline/#mrna-expression-transformation

https://hbctraining.github.io/DGE_workshop/lessons/02_DGE_count_normalization.html

C.G.B.

UNIVERSITÀ DEGLI STUDI DI MILANO

# Functional annotation

Once identified differentially expressed genes, we can ask if they belong to some particular groups of genes, i.e. if they have common functionalities.

We can perform a gene ontology/over-representation analysis/gene set enrichment analysis



| | |
|---|---|
| **Molecular Function** | Molecular-level activities performed by gene products. Molecular function terms describe activities that occur at the molecular level, such as "catalysis" or "transport". GO molecular function terms represent activities rather than the entities (molecules or complexes) that perform the actions, and do not specify where, when, or in what context the action takes place. Molecular functions generally correspond to activities that can be performed by individual gene products (*i.e.* a protein or RNA), but some activities are performed by molecular complexes composed of multiple gene products. Examples of broad functional terms are *catalytic activity* and *transporter activity*; examples of narrower functional terms are *adenylate cyclase activity* or *Toll-like receptor binding*. To avoid confusion between gene product names and their molecular functions, GO molecular functions are often appended with the word "activity" (a *protein kinase* would have the GO molecular function *protein kinase activity*). |
| **Cellular Component** | The locations relative to cellular structures in which a gene product performs a function, either cellular compartments (*e.g.*, *mitochondrion*), or stable macromolecular complexes of which they are parts (*e.g.*, the *ribosome*). Unlike the other aspects of GO, cellular component classes refer not to processes but rather a cellular anatomy. |
| **Biological Process** | The larger processes, or 'biological programs' accomplished by multiple molecular activities. Examples of broad biological process terms are *DNA repair* or *signal transduction*. Examples of more specific terms are *pyrimidine nucleobase biosynthetic process* or *glucose transmembrane transport*. Note that a biological process is not equivalent to a pathway. At present, the GO does not try to represent the dynamics or dependencies that would be required to fully describe a pathway. |

http://geneontology.org/docs/ontology-documentation/

# Functional annotation

Once identified differentially expressed genes, we can ask if they belong to some particular groups of genes, i.e. if they have common functionalities.

We can perform a gene ontology/over-representation analysis/gene set enrichment analysis

## Molecular Signatures Database

**Human Collections**

| | |
|---|---|
| **H** — **hallmark gene sets** are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes. | **C5** — **ontology gene sets** consist of genes annotated by the same ontology term. |
| **C1** — **positional gene sets** corresponding to human chromosome cytogenetic bands. | **C6** — **oncogenic signature gene sets** defined directly from microarray gene expression data from cancer gene perturbations. |
| **C2** — **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts. | **C7** — **immunologic signature gene sets** represent cell states and perturbations within the immune system. |
| **C3** — **regulatory target gene sets** based on gene target predictions for microRNA seed sequences and predicted transcription factor binding sites. | **C8** — **cell type signature gene sets** curated from cluster markers identified in single-cell sequencing studies of human tissue. |
| **C4** — **computational gene sets** defined by mining large collections of cancer-oriented microarray data. | MSigDB Molecular Signatures Database |

https://www.gsea-msigdb.org/gsea/msigdb/