# Bioinformatics for Pathway Enrichment Analysis Part 2

Veronique Voisin and Ruth Isserlin
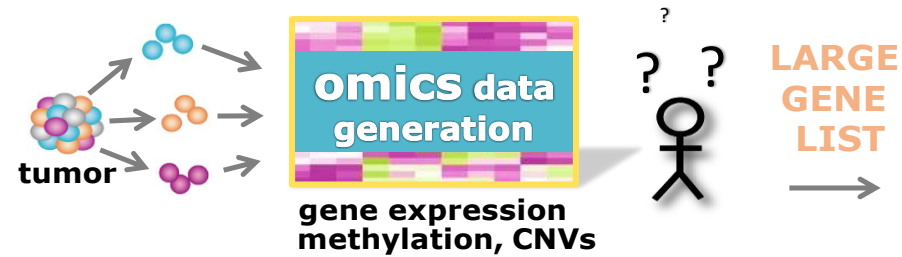
Computer Ontario Summer school 2025

# Course outline

- General concept of pathway enrichment analysis with a ranked list.

- How do we perform pathway enrichment analysis (steps)?

- Example :TCGA pancreatic cancer RNASeq expression

- Practical lab: pathway enrichment analysis using R and Cytoscape
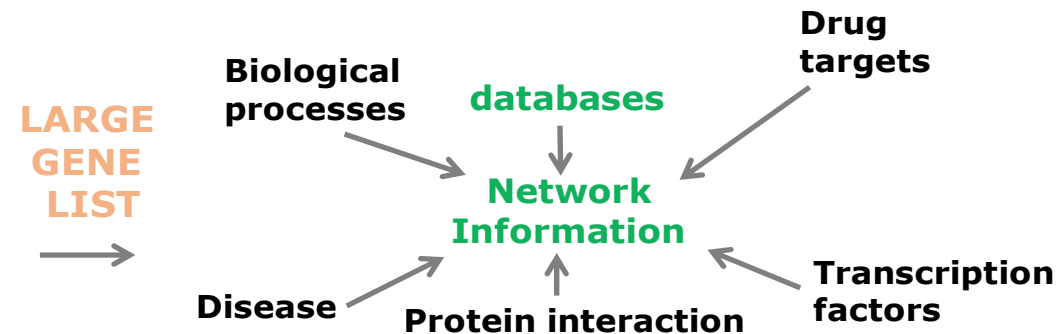
# Learning Objectives

- Be able to understand:

  - Be able to run GSEA(Gene Set Enrichment Analysis) on a <mark>ranked gene list</mark> and understand the main parameters and output results.
  - Understand the different ways we can visualize enrichment results
  - the advantages of these different pathway enrichment analysis and visualizations.
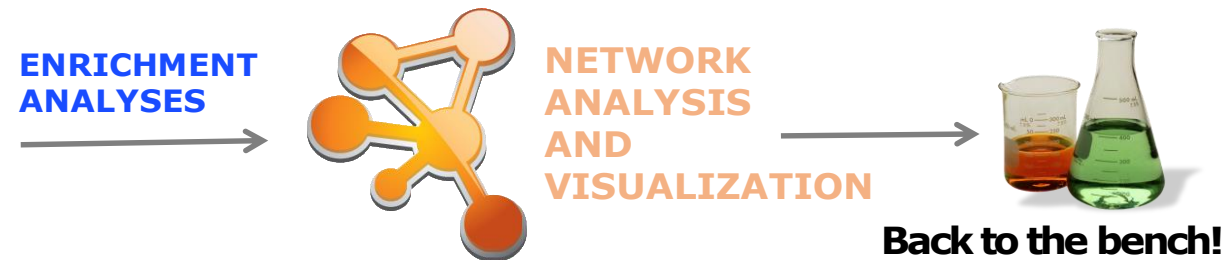
# General Workflow of pathway enrichment analysis

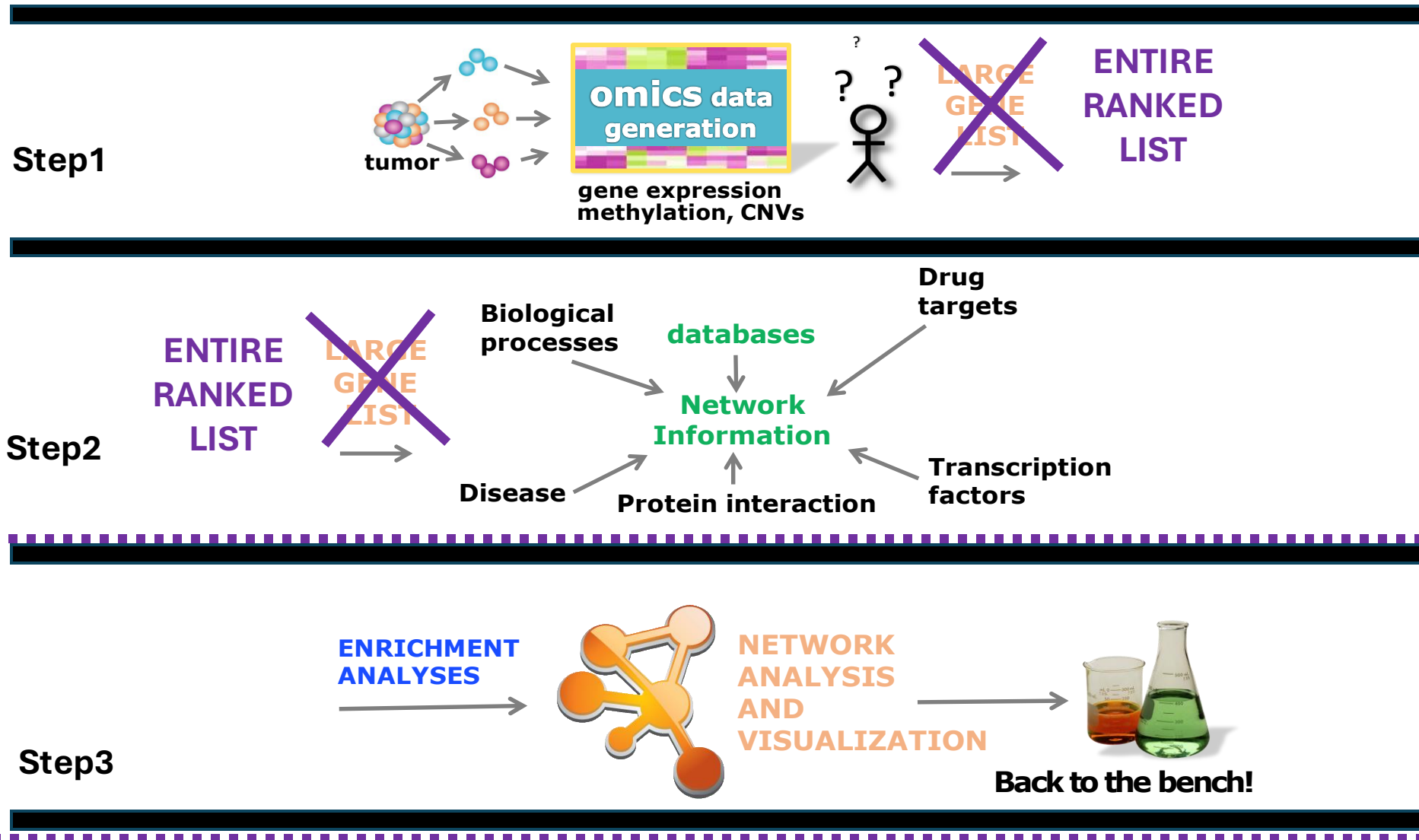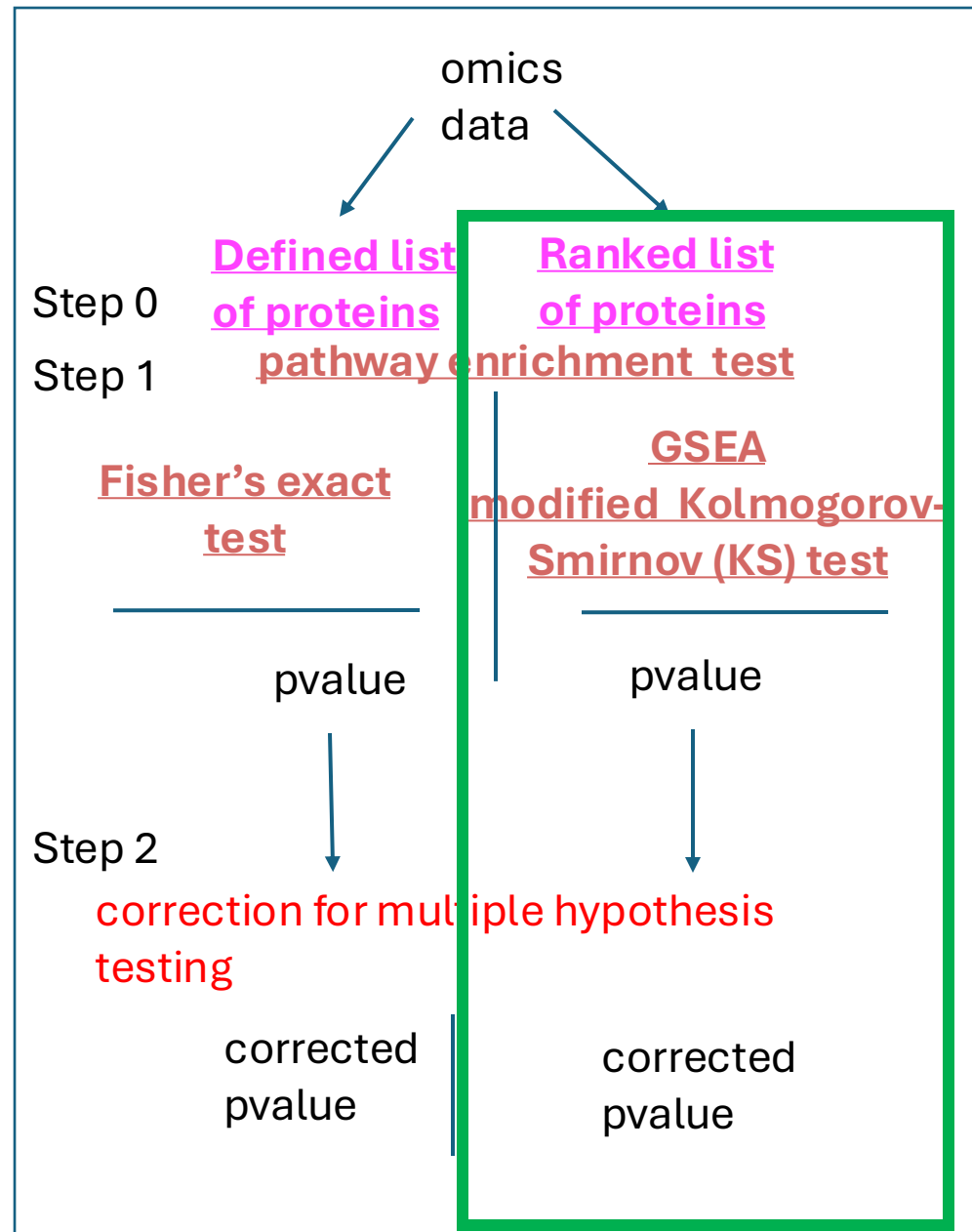**Step1**



tumor → omics data generation → LARGE GENE LIST

gene expression methylation, CNVs

**Step2**

LARGE GENE LIST →

Biological processes → Network Information

databases → Network Information

Drug targets → Network Information

Disease → Network Information

Protein interaction → Network Information

Transcription factors → Network Information

**Step3**

ENRICHMENT ANALYSES → NETWORK ANALYSIS AND VISUALIZATION → Back to the bench!

# General Workflow of pathway enrichment analysis



**Step1**

**Step2**

**Step3**

# Two types of gene lists : defined or ranked list.

- Fisher's Exact Test, aka Hypergeometric Test
- GSEA for ranked lists.
- Multiple test corrections:
  - Bonferroni correction
  - False Discovery Rate computation using Benjamini-Hochberg procedure

omics data

Step 0
Step 1

**Defined list of proteins**

**Ranked list of proteins**

**pathway enrichment test**

**Fisher's exact test**

**GSEA modified Kolmogorov-Smirnov (KS) test**

pvalue

pvalue

Step 2

correction for multiple hypothesis testing

corrected pvalue

corrected pvalue

# Why test enrichment with ranked gene lists?

- Possible problems with thresholded gene list test include:
  - No "natural" value for the threshold.  Avoids arbitrary cut-offs.
  - Different results at different threshold settings
  - Useful for Noisy or Weak Signals
    - Particularly valuable when no genes are strongly differentially expressed, but **pathways** still show directional change.
  - Possible loss of statistical power due to the thresholding
    - No resolution between significant signals with different strengths
    - Weak signals are ignored.

  ** reduces bias, increases reproducibility, and improves sensitivity for detecting real biological patterns

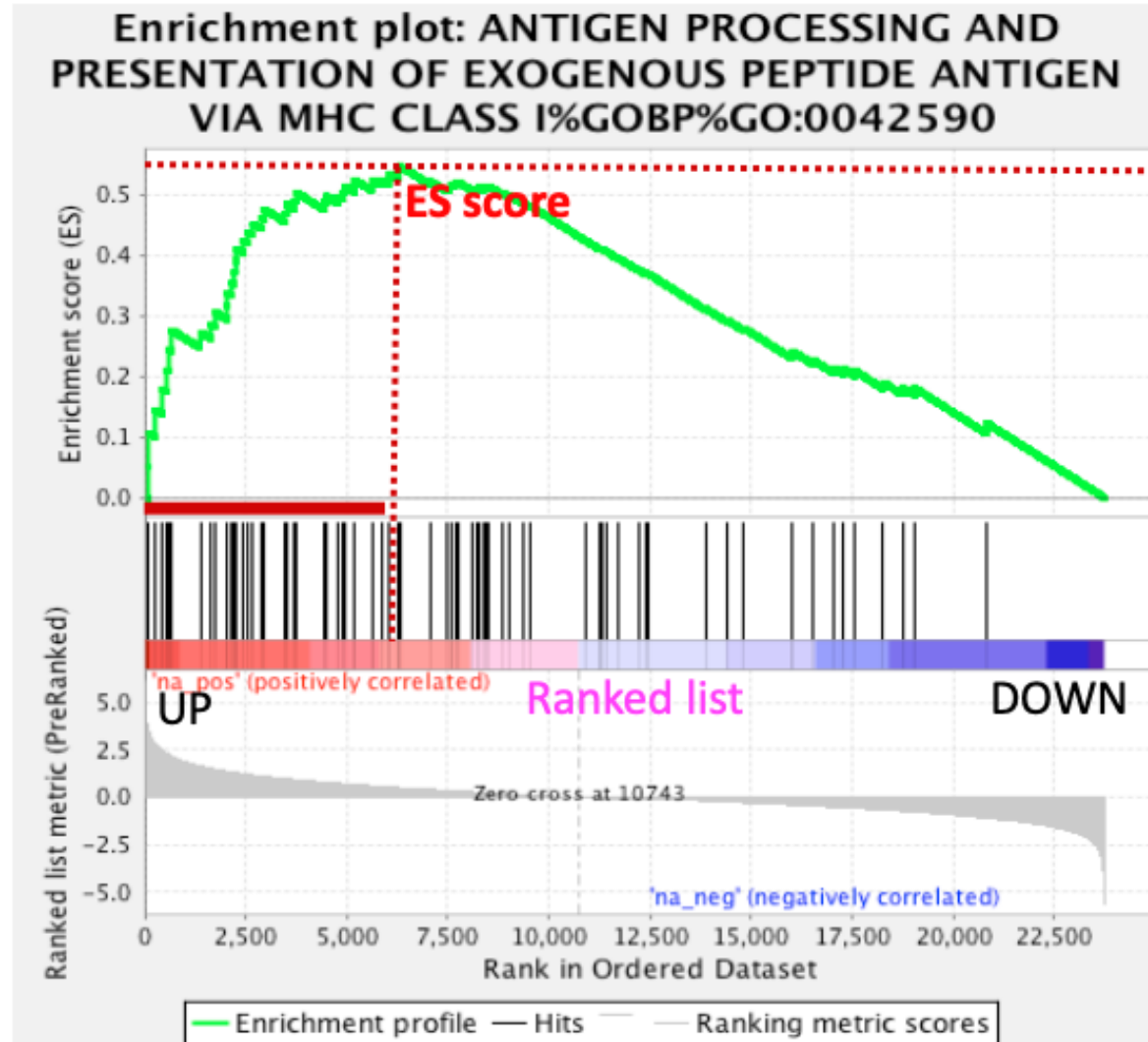# Gene Set Enrichment Analysis - GSEA

- Original paper published in Nature Genetics 2003 (Mootha el al) and modified and republished in PNAS 2005 (Subramanian et al)

- In the original paper, Mootha et al (2003) studied diabetes and identified that their gene list was significantly enriched in a pathway called "oxidative phosphorylation".

- The particularity of this finding was that the individual genes in this pathway were only down-regulated by a small amount but the addition of all these subtle signals had a great impact on the pathway.

- They validated their finding experimentally.

# GSEA score calculation

## modified Kolmogorov Smirnov test (KS test)



**Ranked gene list**

UP

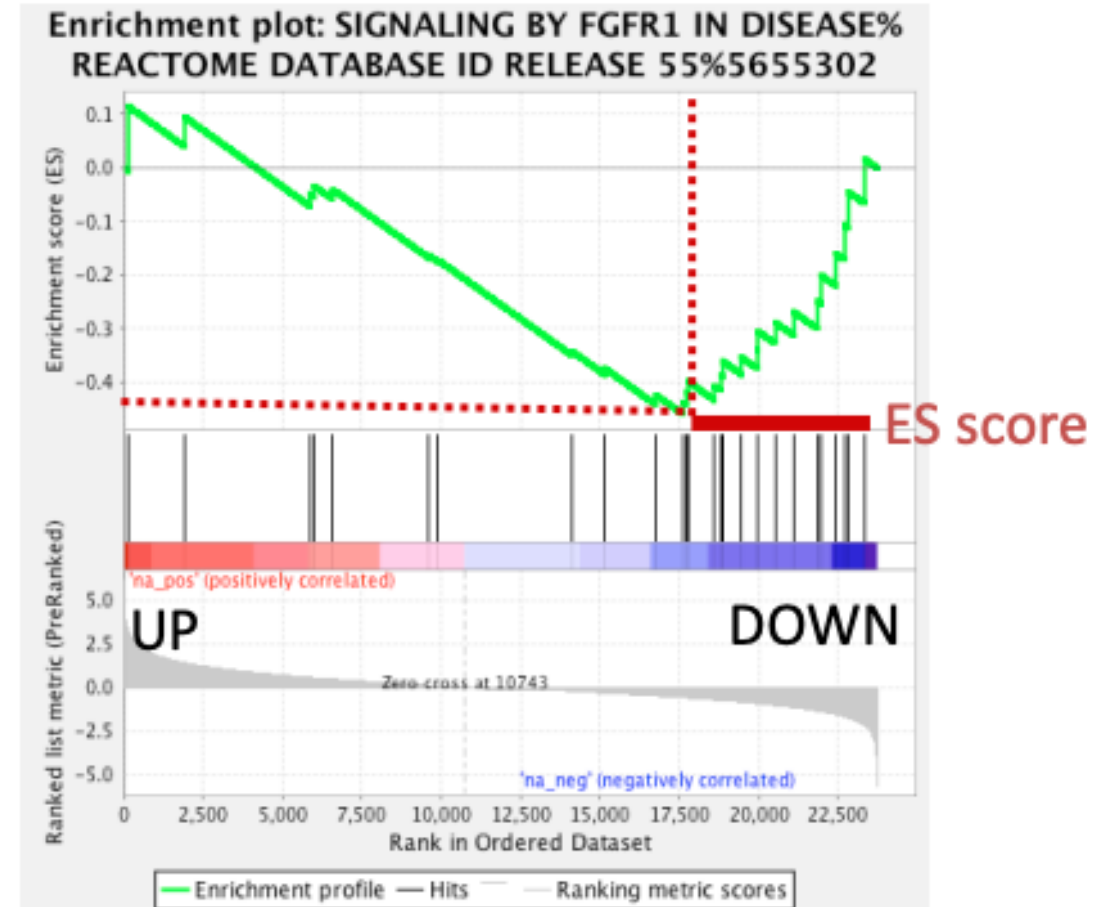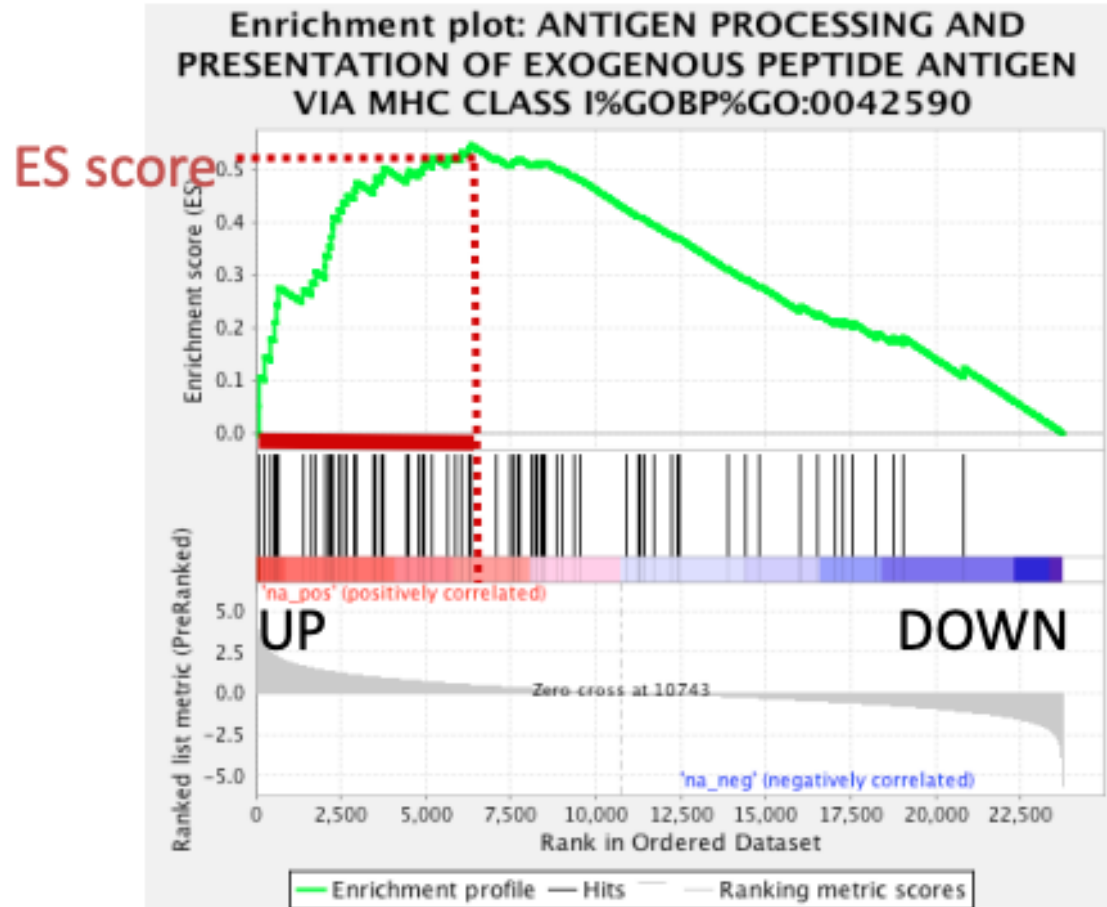| | |
|---|---|
| BGN | 32.76 |
| ANTXR1 | 30.36 |
| FZD1 | 29.36 |
| COL16A1 | 28.88 |
| KLF3 | 1.08 |
| RASEF | 0.05 |
| ... | ... |
| ... | ... |
| ISOC1 | 0.05 |
| ANO1 | 0.04 |
| CBWD3 | -1.09 |
| GBP4 | -15.6 |
| TAP1 | -19 |
| PSMB9 | -19.7 |

DOWN

1. Walk down the ranked list of all genes.
2. **Increase a running-sum** when a gene is in the gene set.
3. **Decrease the running-sum** when a gene is not.
4. The **maximum deviation from zero** of this running-sum is the **enrichment score (ES)**
5. Any gene contributing the ES before the maximum deviation is considered part of the **leading edge.**
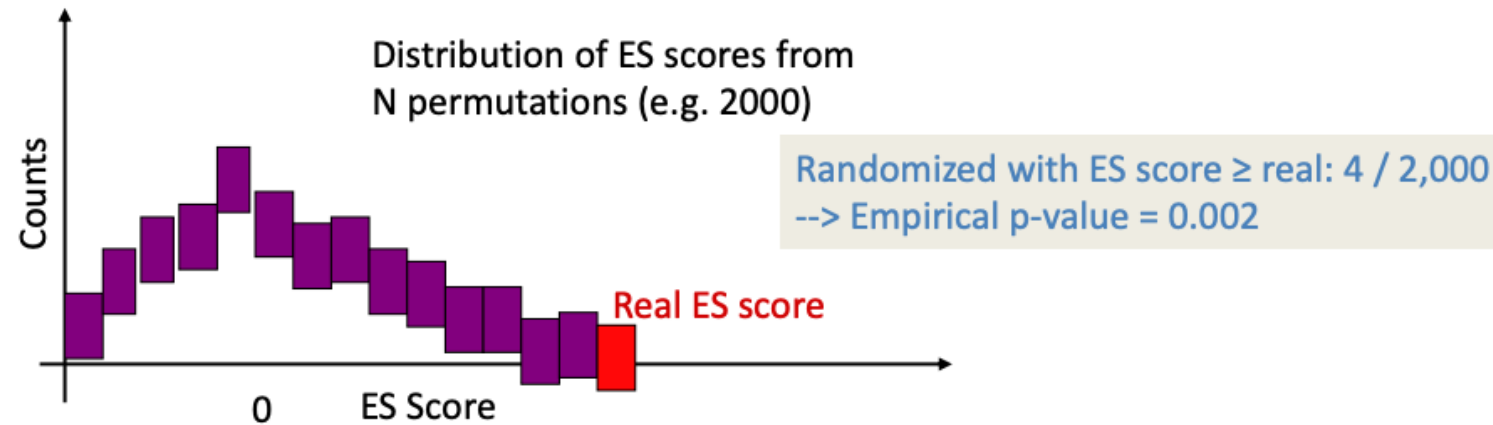
# **Weighted -** Kolmogorov Smirnov test (KS test)



Weight
p=1 (default)

Options:
p=0 no weighting (all genes contribute equally, like Kolmogorov–Smirnov test)
p=1 linear weighting (default)
p>1 more emphasis on top-ranked genes

Found/not found in gene-set (pathway)

running sum

enirchment score

0.4
0.3
0.2
0.1
0

UP    32.76  30.36  25.62  6.04  5.04    4.03    3.01  2.96    2.88  2.77  2.66    **Ranks**

# Positive and Negative enrichment scores

# Going from ES, NES, p-value and FDR



Distribution of ES scores from
N permutations (e.g. 2000)

Randomized with ES score ≥ real: 4 / 2,000
--> Empirical p-value = 0.002

Real ES score

Counts

0    ES Score

- Randomize the ranked list and compute ES values for each pathway
- Repeat n times (by default 1000, but consider using higher)
- P-value = empirical p-value → (number of randomized ES scores > real value) / num of permutations
- NES = Normalized ES → $$NES = \frac{ES}{\text{mean}(|ES_{null}|)}$$

- FDR = false discovery Rate is a multiple testing correction that estimates **how likely it is that a gene set with a given NES is a false positive.**

$$FDR(NES_{obs}) = \frac{\text{Number of null NES} \geq \text{NES}_{obs}}{\text{Number of observed NES} \geq \text{NES}_{obs}}$$

# Which files do we need to run GSEA?

- A ranked list of genes called the rank file
  - this is a text file (tab separated) that should be renamed to end with the extension .rnk
  - This file has 2 columns :
    - gene identifier
    - ranking values

- A file called a .gmt file that contains the pathway data base (the gene-sets)
  - this is a text file (tab separated) that should end with the extension .gmt
  - the first column contains gene-set names and the additional columns contains the gene names included in each gene-set

# What is a ranked gene list?

## Two-class design: ranked gene list



**Expression Matrix**

Class-1   Class-2

**Genes Ranked by Differential Statistic**

UP

DOWN

E.g.:
Fold change
Log (ratio)
t values from t-test

**Ranking score = sign(logFC)*-log10(pvalue)**

|        | LogFC | Pvalue   | score  |
|--------|-------|----------|--------|
| BGN    | +1    | 1.73E-33 | 32.76  |
| ANTXR1 | +1    | 4.39E-31 | 30.36  |
| FZD1   | +1    | 4.41E-30 | 29.36  |
| COL16A1| +1    | 1.33E-29 | 28.88  |
| KLF3   | +1    | 8.32E-02 | 1.08   |
| RASEF  | +1    | 9.01E-01 | 0.05   |
| ISOC1  | +1    | 9.01E-01 | 0.05   |
| ANO1   | +1    | 9.01E-01 | 0.04   |
| CBWD3  | -1    | 8.18E-02 | -1.09  |
| GBP4   | -1    | 2.45E-16 | -15.61 |
| TAP1   | -1    | 1.04E-19 | -18.98 |
| PSMB9  | -1    | 1.84E-20 | -19.73 |

UP

DN

# What does a .gmt file look like?

- Already introduced in part 1 but just for a refresher
  - It is a tab delimited text file with extension .gmt that contains pathways and their associated genes

| name | description | genes | | | | | |
|------|-------------|-------|--|--|--|--|--|
| HALLMARK_P53_PATHWAY%MSIGDBHALLMARK%HALLMARK_P53_PATHWAY | HALLMARK_P53_PATHWAY | RB1 | CYFIP2 | APP | STEAP3 | CCNK | HEL-S-130P |
| HALLMARK_MITOTIC_SPINDLE%MSIGDBHALLMARK%HALLMARK_MITOTIC_SPINDLE | HALLMARK_MITOTIC_SPINDLE | HCTP4 | CCNB2 | RHOT2 | BCL2L11 | KIF3B | STK38L |
| HALLMARK_PI3K_AKT_MTOR_SIGNALING%MSIGDBHALLMARK%HALLMARK_PI3K_AKT_MTOR_SIGNA | HALLMARK_PI3K_AKT_MTOR_SIGNALING | PTEN | PRKAG1 | SLA | ECSIT | HEL-S-125m | RPTOR |
| HALLMARK_EPITHELIAL_MESENCHYMAL_TRANSITION%MSIGDBHALLMARK%HALLMARK_EPITHELIAL_ | HALLMARK_EPITHELIAL_MESENCHYMAL_TRANSITION | NTM | HTRA1 | TNFRSF11B | THY1 | NID2 | FBLN2 |
| HALLMARK_GLYCOLYSIS%MSIGDBHALLMARK%HALLMARK_GLYCOLYSIS | HALLMARK_GLYCOLYSIS | PSMC4 | CTH | AGRN | POLR3K | PFKP | QSCN6 |
| HALLMARK_WNT_BETA_CATENIN_SIGNALING%MSIGDBHALLMARK%HALLMARK_WNT_BETA_CATEN | HALLMARK_WNT_BETA_CATENIN_SIGNALING | HDAC5 | HDAC2 | MAML1 | HDAC11 | LEF1 | PSEN2 |
| HALLMARK_SPERMATOGENESIS%MSIGDBHALLMARK%HALLMARK_SPERMATOGENESIS | HALLMARK_SPERMATOGENESIS | CCNB2 | CDC2 | CLPB | CHRM4 | GFI1 | ACRBP |
| HALLMARK_APOPTOSIS%MSIGDBHALLMARK%HALLMARK_APOPTOSIS | HALLMARK_APOPTOSIS | CTH | RHOT2 | BCL2L11 | LEF1 | PSEN2 | GADD45B |
| HALLMARK_G2M_CHECKPOINT%MSIGDBHALLMARK%HALLMARK_G2M_CHECKPOINT | HALLMARK_G2M_CHECKPOINT | HCTP4 | HMGA1 | CCNB2 | FANCC | FOXN3 | CDC6 |
| HALLMARK_COMPLEMENT%MSIGDBHALLMARK%HALLMARK_COMPLEMENT | HALLMARK_COMPLEMENT | TFPI2 | IL6 | DOCK4 | HEL-S-130P | TNFAIP3 | CXCL1 |
| HALLMARK_E2F_TARGETS%MSIGDBHALLMARK%HALLMARK_E2F_TARGETS | HALLMARK_E2F_TARGETS | HMGA1 | POLD2 | CCNB2 | EXOSC8 | PPP1R8 | USP1 |
| HALLMARK_FATTY_ACID_METABOLISM%MSIGDBHALLMARK%HALLMARK_FATTY_ACID_METABOLISM | HALLMARK_FATTY_ACID_METABOLISM | SETD8 | ACSM3 | ECI2 | CPOX | AQP7 | UBE2L6 |
| HALLMARK_TNFA_SIGNALING_VIA_NFKB%MSIGDBHALLMARK%HALLMARK_TNFA_SIGNALING_VIA_ | HALLMARK_TNFA_SIGNALING_VIA_NFKB | GADD45B | MYC | INHBA | IL6 | ID2 | CXCL6 |
| HALLMARK_HEDGEHOG_SIGNALING%MSIGDBHALLMARK%HALLMARK_HEDGEHOG_SIGNALING | HALLMARK_HEDGEHOG_SIGNALING | SLIT1 | PML | NRCAM | NKX6-1 | THY1 | L1CAM |
| HALLMARK_CHOLESTEROL_HOMEOSTASIS%MSIGDBHALLMARK%HALLMARK_CHOLESTEROL_HOMEO | HALLMARK_CHOLESTEROL_HOMEOSTASIS | SQLE | ERRFI1 | GLDC | HEL-S-7 | TMEM97 | CLU |
| HALLMARK_INFLAMMATORY_RESPONSE%MSIGDBHALLMARK%HALLMARK_INFLAMMATORY_RESPO | HALLMARK_INFLAMMATORY_RESPONSE | MYC | INHBA | IL6 | CXCL6 | CXCL8 | ITGB3 |
| HALLMARK_TGF_BETA_SIGNALING%MSIGDBHALLMARK%HALLMARK_TGF_BETA_SIGNALING | HALLMARK_TGF_BETA_SIGNALING | TMEPAI | SLC20A1 | KLF10 | SPTBN1 | TGIF1 | MAP3K7 |
| HALLMARK_INTERFERON_GAMMA_RESPONSE%MSIGDBHALLMARK%HALLMARK_INTERFERON_GAM | HALLMARK_INTERFERON_GAMMA_RESPONSE | HLA-DMA | SLC25A28 | PSMB2 | STAT4 | FCGR1A | PFKP |
| HALLMARK_UV_RESPONSE_DN%MSIGDBHALLMARK%HALLMARK_UV_RESPONSE_DN | HALLMARK_UV_RESPONSE_DN | FBLN5 | BDNF | MYC | IRS1 | COL1A1 | PRDM2 |
| HALLMARK_IL6_JAK_STAT3_SIGNALING%MSIGDBHALLMARK%HALLMARK_IL6_JAK_STAT3_SIGNALIN | HALLMARK_IL6_JAK_STAT3_SIGNALING | CD36 | TNF | PTPN1 | STAT1 | STAT3 | CSF2 |
| HALLMARK_DNA_REPAIR%MSIGDBHALLMARK%HALLMARK_DNA_REPAIR | HALLMARK_DNA_REPAIR | POLA2 | DUT | SAC3D1 | SSRP1 | POLE4 | NUDT21 |
| HALLMARK_ADIPOGENESIS%MSIGDBHALLMARK%HALLMARK_ADIPOGENESIS | HALLMARK_ADIPOGENESIS | EPHX2 | CHCHD10 | RMDN3 | COQ9 | DNAJC15 | COQ5 |
| HALLMARK_INTERFERON_ALPHA_RESPONSE%MSIGDBHALLMARK%HALLMARK_INTERFERON_ALPHA | HALLMARK_INTERFERON_ALPHA_RESPONSE | SLC25A28 | IFITM3 | GBP4 | CD74 | MX1 | CSF1 |

# EXAMPLE WITH A RANKED GENE LIST

**Dataset Pancreatic Ductal Adenocarc (TCGA)**

Tumor subtype

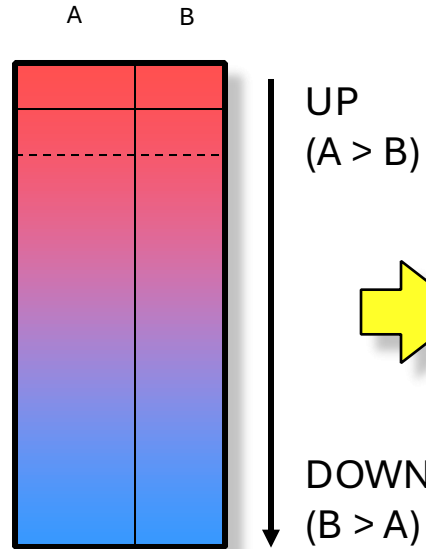**Basal** vs **Classical**

Differential expression (edgeR)

Rank file

Gene-Set Enrichment Analysis (fGSEA)

Classical
Basal-like

Moffitt, R., Marayati, R., Flate, E. *et al.* Virtual microdissection identifies distinct tumor- and stroma-specific subtypes of pancreatic ductal adenocarcinoma. *Nat Genet* **47**, 1168–1178 (2015)
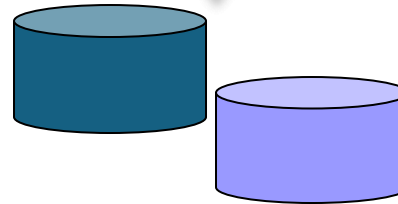
# Pathway Enrichment Analysis



Merico D, Isserlin R, Stueker O, Emili A, Bader GD
Enrichment map: a network-based method for gene-set enrichment visualization and interpretation PLoS One. 2010 Nov 15;5(11):e13984

# Network Basics

**Node (molecule/entity)**

- Gene
- Protein
- Transcript
- Drug
- MicroRNA
- ...

**Edge (interaction/relationship)**

- Genetic interaction
- Physical protein interaction
- Co-expression
- Metabolic reaction
- DNA-binding
- ...

# Enrichment Map Basics

**Node (molecule/entity)**

- Pathway or geneset

- Size is correlated to number of genes in set

- Color indicates class in below example (for example Up/Down, classA/classB)

**Edge (interaction/relationship)**

- Degree of overlap between two genesets

- The more genes two pathways have in common the thicker/stronger the connection

# Enrichment Map



Enriched in phenotype
A [red-white-blue gradient] B

$$\frac{|A \cap B|}{\min(|A|,|B|)}$$

# Typical Output



### Network Visualization

Each row is a gene-set (pathway).

It displays:
- a score associated with the magnitude of overlap between gene-set and gene list.
- a pvalue that estimates the significance of the enrichment (by chance or not).
- a corrected pvalue (FDR) that corrects for multiple hypothesis testing.

# During the practical lab:

- We will run fGSEA from R
- Examine and Visualize the results in R
- Create output files that we can then use in Cytoscape

This module is taken from
CBW Pathway and Network Analysis Workshop 2024

# 2025 Canadian Bioinformatics Workshop
## Please visit the website to see workshops presently offered



https://bioinformatics.ca/workshops/current-workshops/

# Time for practical lab

- https://baderlab.github.io/ComputeOntario_Pathways_2025/gsea-lab.html