# Functional analysis of omics data
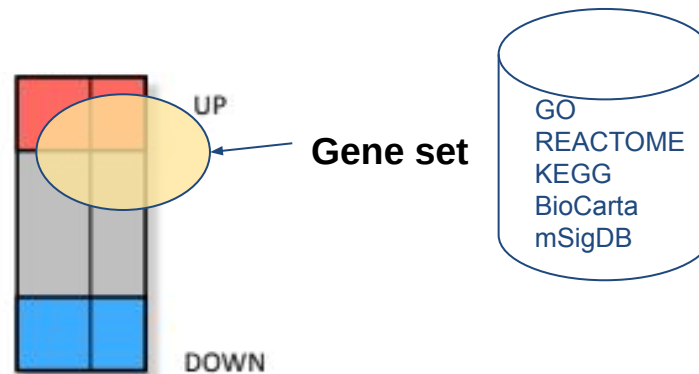
**february 2020**

**MAR**Genomics

# Summary

1. Summary of previous days

2. Example datasets: He et al, Am J Transl Res, 2018

    a. GEO

    b. TCGA

3. Hands on

    a. Differential Expression Analysis of Public Data

    b. DAVID

    c. GSEA Preranked

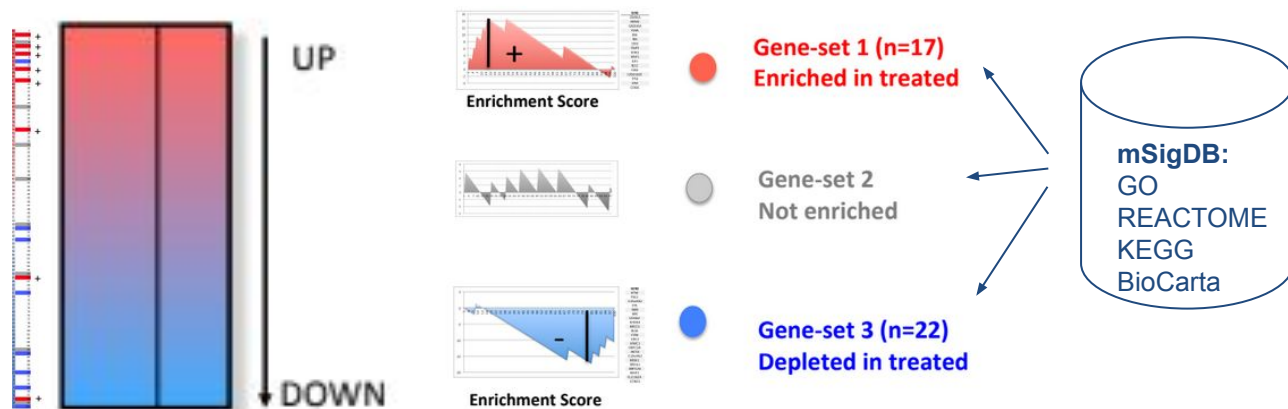    d. Cytoscape
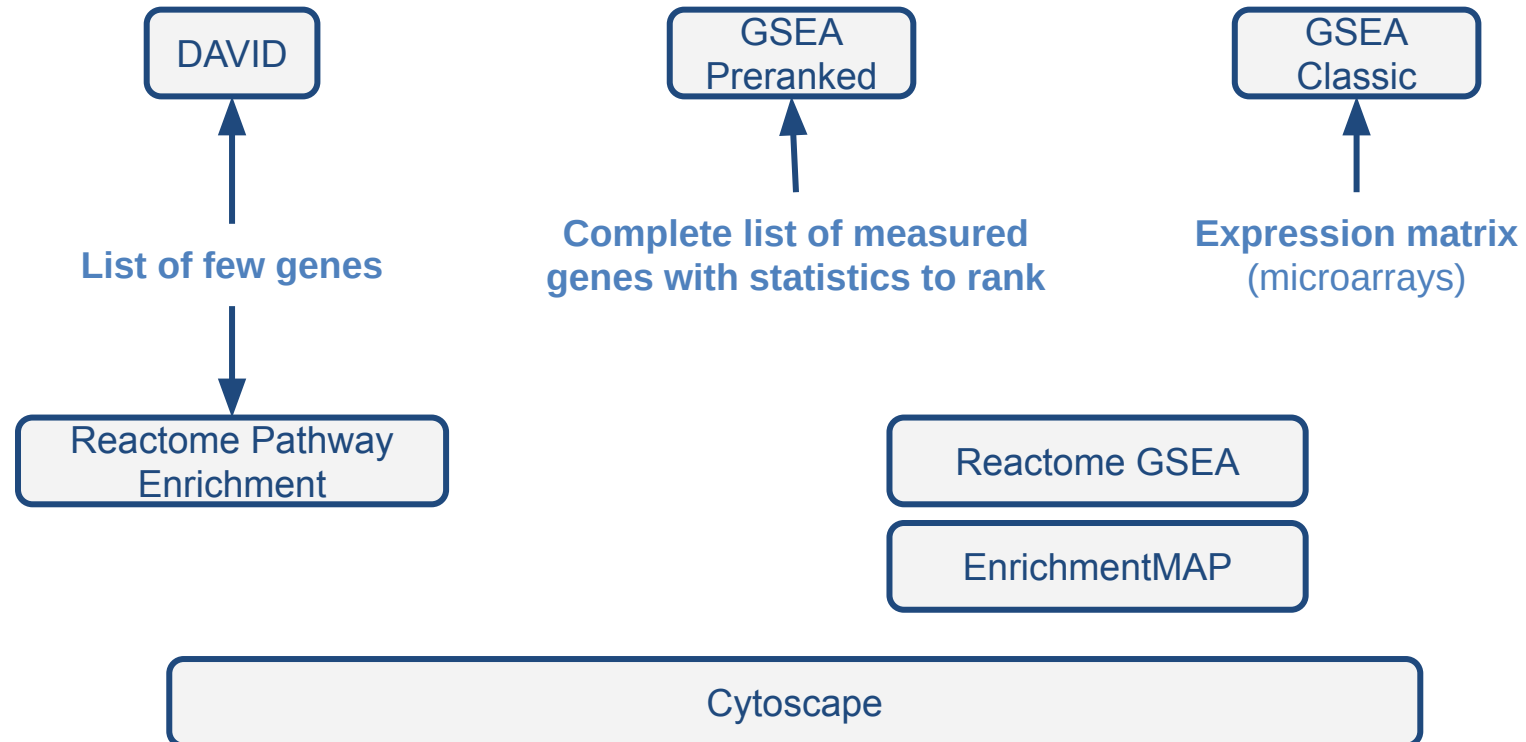
    e. Public Resources to "validate" candidates

# 1. Summary of enrichment strategies

- Gene list (e.g. expression change > 2-fold AND FDR < 0.05)



**Gene set**

GO
REACTOME
KEGG
BioCarta
mSigDB

- Ranked list (e.g. by -log10(p)*sign(logFC) )



**Enrichment Score**

Gene-set 1 (n=17)
Enriched in treated

Gene-set 2
Not enriched

Gene-set 3 (n=22)
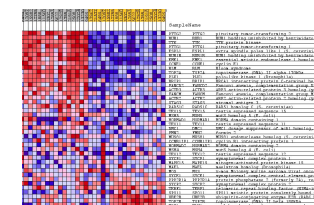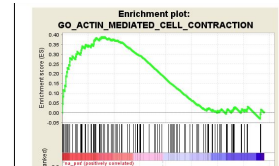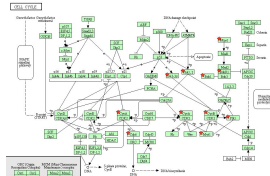Depleted in treated

**mSigDB:**
GO
REACTOME
KEGG
BioCarta

# 1. Summary of enrichment strategies

- Gene list (e.g. expression change > 2-fold AND FDR < 0.05)

  - **Answers the question**: Are any gene sets surprisingly enriched (or depleted) in my gene list?
  - **Statistical test**: Fisher's Exact Test (aka Hypergeometric test)
  - **Tools**: DAVID
  - **Benefits:** simple, you only need a list of gene names
  - **Problems**: Possible loss of statistical power due to thresholding. Different results at different threshold settings

- Ranked list (e.g. by -log10(p)*sign(logFC) )

  - **Answers the question**: Are any gene sets ranked surprisingly high or low in my ranked list of genes?
  - **Statistical test**: GSEA
  - **Benefit**: use information of all genes measured. Increase of statistical power
  - **Problems**: more difficult to prepare files. You need the whole experiment

# 1. Summary of enrichment strategies



**DAVID**

**GSEA Preranked**

**GSEA Classic**

**List of few genes**

**Complete list of measured genes with statistics to rank**

**Expression matrix** (microarrays)

Reactome Pathway Enrichment

Reactome GSEA

EnrichmentMAP

Cytoscape

**NO TOPOLOGY**

**TOPOLOGY**

# 1. Summary of public sources (open data)

## Projects

ENCODE
TCGA
GTEx
Cancer Cell Line Encyclopedia
The Human Protein Atlas
FANTOM5
Researchers
Publications

...

## Web sources

GDAC  iCGC
cBioPortal
GEPIA   XENA
GEO  SRA
dbGaP

...

## Databases

GO
REACTOME
KEGG
BioCarta
mSigDB

# 2. Example datasets: [He (2018) Am J Transl Res](#)

## Analysis of differentially expressed genes, clinical value and biological pathways in prostate cancer

Zhaohui He,[1,*] Fucai Tang,[1,*] Zechao Lu,[2,*] Yucong Huang,[3] Hanqi Lei,[1] Zhibiao Li,[3] and Guohua Zeng[1]

► Author information  ► Article notes  ► Copyright and License information Disclaimer



Figure 1

**He_etal_2018.pdf**

www.imim.es

# 2. Example datasets: [He (2018) Am J Transl Res](#)

## Figure 2. GO Enrichment Analysis
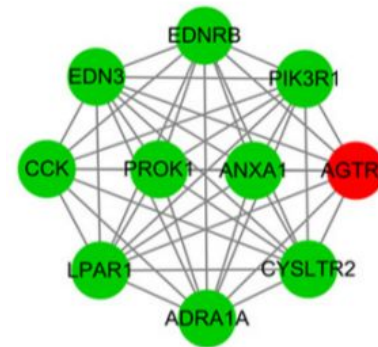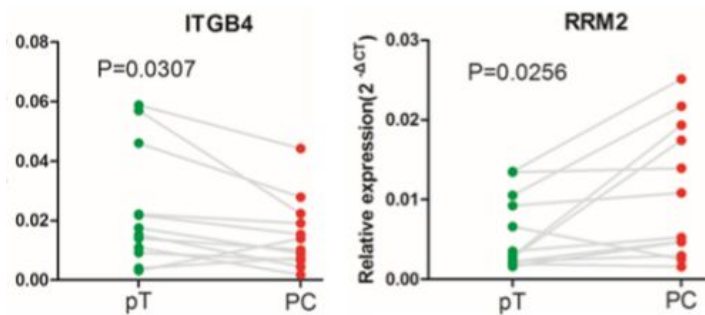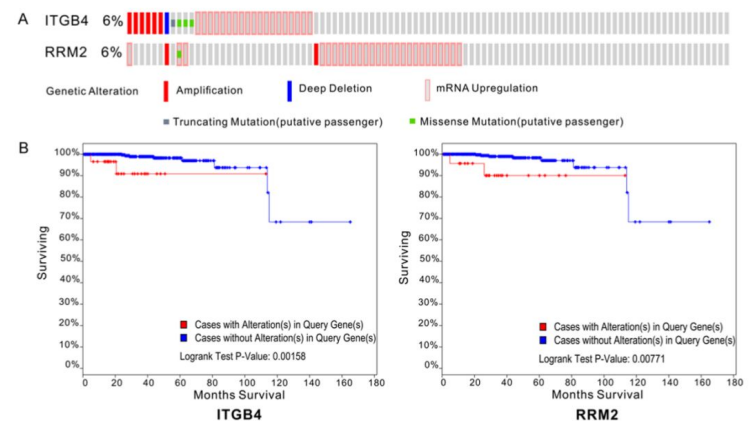


## Figure 3. PPI network of DEGs



## Figure 4. Validation of genes with RT-qPCR



**He_etal_2018.pdf**

## Figure 5. Validation of genes' prognostic value

# 2. Example datasets: [He (2018) Am J Transl Res](#)
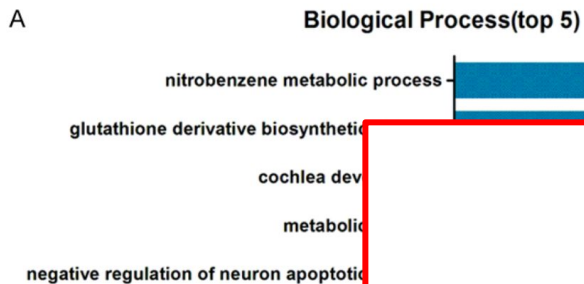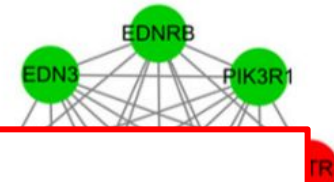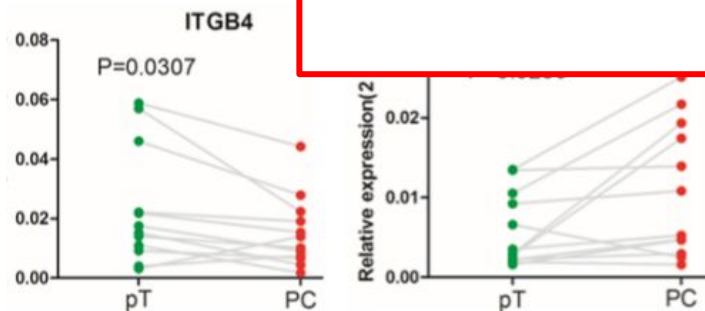
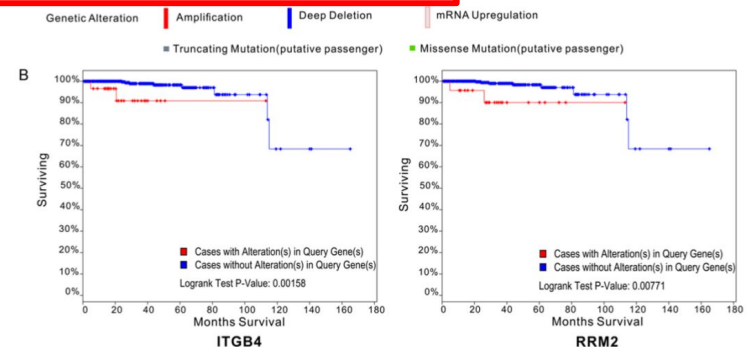**Figure 2. GO Enrichment Analysis**

**Figure 3. PPI network of DEGs**

**Biological Process(top 5)**

nitrobenzene metabolic process

glutathione derivative biosynthetic

cochlea dev

metabolic

negative regulation of neuron apoptotic

Are you able to reproduce this paper with the tools we have learnt during the course?

**Figure 4. Validation** ...s' prognostic value

ITGB4
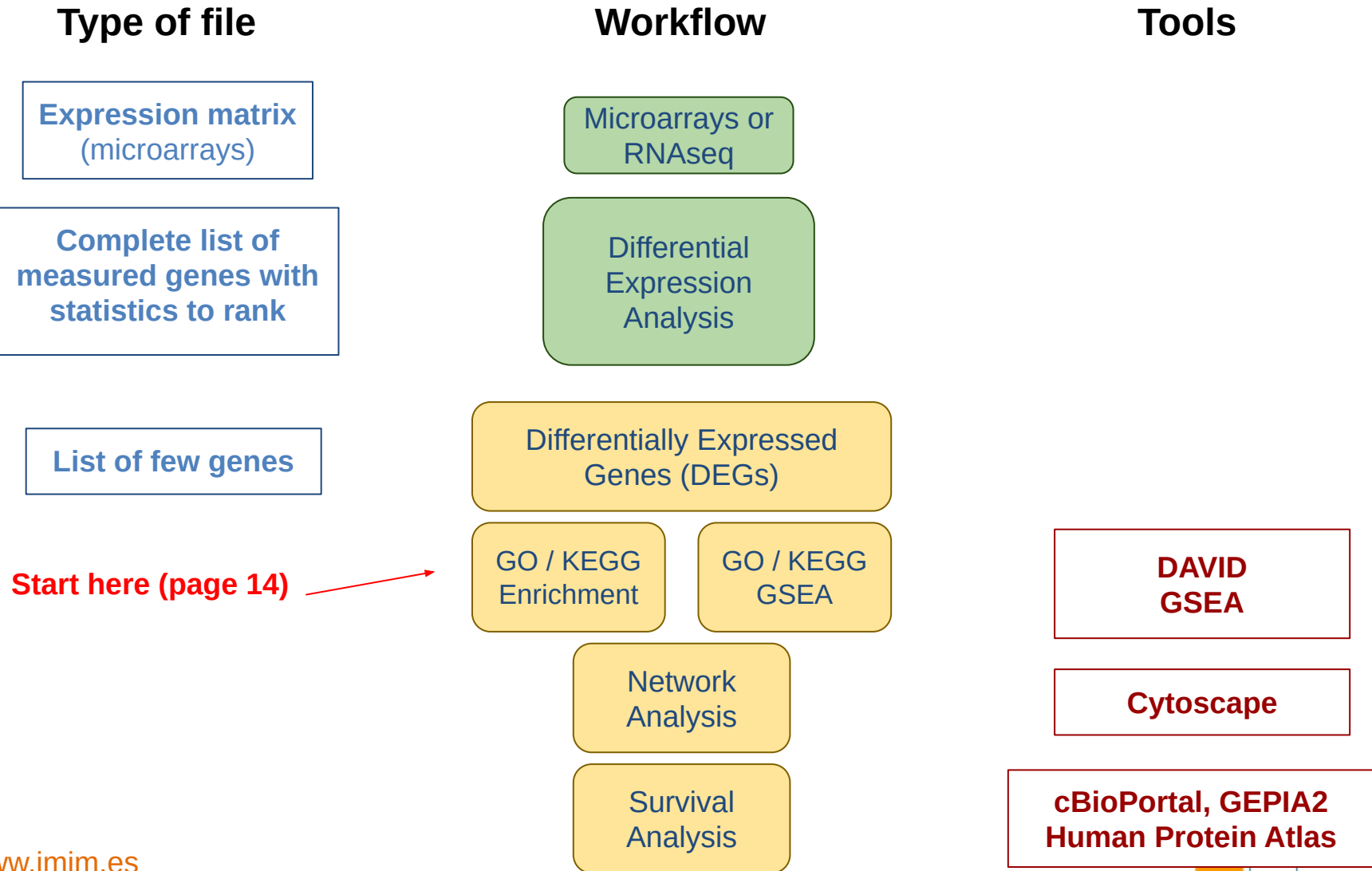P=0.0307

**He_etal_2018.pdf**

# 3. Hands on: [He et al (2018) Am J Transl Res](#)

## Type of file

**Expression matrix** (microarrays)

**Complete list of measured genes with statistics to rank**

**List of few genes**

## Workflow

GEO (GSE46602)

TCGA (PRAD)

Differential Expression Analysis

Differential Expression Analysis

Consensus Differentially Expressed Genes (DEGs)

GO / KEGG Enrichment

GO / KEGG GSEA

Network Analysis

Survival Analysis

## Tools

**GEO GEPIA2**

**GEO2R GEPIA2 (R)**

**Draw Venn Venny**

**DAVID GSEA**

**Cytoscape**

**cBioPortal, GEPIA2 Human Protein Atlas**

# 3. Hands on: Your own data

**Type of file**

Expression matrix (microarrays)

Complete list of measured genes with statistics to rank

List of few genes

**Start here (page 14)** →

**Workflow**

Microarrays or RNAseq

Differential Expression Analysis

Differentially Expressed Genes (DEGs)

GO / KEGG Enrichment

GO / KEGG GSEA

Network Analysis

Survival Analysis

**Tools**

DAVID
GSEA

Cytoscape

cBioPortal, GEPIA2
Human Protein Atlas

# 3. Hands on: Differential Expression

Imagine you are the authors of the paper, so you are interested in studying key genes and functions involved in prostate cancer. Take advantage of published datasets by using the tools we have learnt during this course.

- Perform differential expression analysis for GSE46602 using GEO2R
    - control: 14 samples
    - tumor: 36 samples
    - Results saved in file: GSE46602_GEO2R_all_results.txt

# 3. Hands on: Differential Expression

Imagine you are the authors of the paper, so you are interested in studying key genes and functions involved in prostate cancer. Take advantage of published datasets by using the tools we have learnt during this course.

- Use **GEPIA2** to perform differential expression analysis of RNAseq data from TCGA
  - Results saved in file: PRAD_GEPIA2_table_degenes.txt
  - Add column names: Gene, ENSEMBL, median (tumor), median (normal), log2FC, adj.p.val
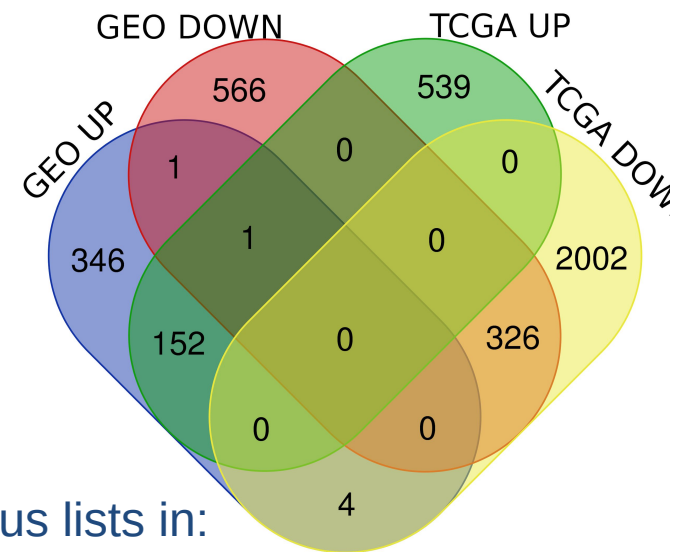
# 3. Hands on: Intersection of lists

Now, you can use the lists of genes to perform downstream functional analyses.

But, which list of genes are you going to work with? In the article they use adj.P.Val<0.05 AND |logFC|>1

- You can check the overlap using [Draw Venn](), [Venny]() or [VennDiagrams]()
- With the data I downloaded, I get this venn:
  - 326 consensus DOWN genes
  - 152 consensus UP genes
  - Total: 478 DEG

  **The numbers are not exactly the same as in the publication (168 UP, 316 DOWN). Any ideas why?**



- You can find the UP,DOWN and ALL consensus lists in:
  - Consensus_UP.txt
  - Consensus_DOWN.txt
  - Consensus_ALL.txt

# 3. Hands on: Enrichment analysis

Using the list of genes, you can perform a simple gene set enrichment using DAVID. Do you find similar gene sets enriched as in the paper?

# 3. Hands on: Enrichment analysis

Using the list of genes, you can perform a simple gene set enrichment using DAVID. Do you find similar gene sets enriched as in the paper?

| Sublist | Category | Term | RT | Genes | Count | % | P-Value | Benjamini |
|---------|----------|------|-----|-------|-------|---|---------|-----------|
| ☐ | GOTERM_BP_DIRECT | hemidesmosome assembly | RT | | 7 | 1.5 | 1.8E-7 | 4.3E-4 |
| ☐ | GOTERM_BP_DIRECT | cell differentiation | RT | | 28 | 5.9 | 3.3E-5 | 3.9E-2 |
| ☐ | GOTERM_BP_DIRECT | angiogenesis | RT | | 18 | 3.8 | 4.0E-5 | 3.2E-2 |
| ☐ | GOTERM_BP_DIRECT | extracellular matrix organization | RT | | 15 | 3.1 | 3.7E-4 | 2.0E-1 |
| ☐ | GOTERM_BP_DIRECT | negative regulation of epithelial cell proliferation | RT | | 8 | 1.7 | 4.3E-4 | 1.9E-1 |
| ☐ | GOTERM_BP_DIRECT | negative regulation of protein phosphorylation | RT | | 8 | 1.7 | 7.3E-4 | 2.5E-1 |
| ☐ | GOTERM_BP_DIRECT | positive regulation of MAPK cascade | RT | | 9 | 1.9 | 8.5E-4 | 2.6E-1 |
| ☐ | GOTERM_BP_DIRECT | brown fat cell differentiation | RT | | 6 | 1.3 | 1.0E-3 | 2.7E-1 |
| ☐ | GOTERM_BP_DIRECT | response to hypoxia | RT | | 13 | 2.7 | 1.2E-3 | 2.7E-1 |
| ☐ | GOTERM_BP_DIRECT | positive reg | | | | | | |
| ☐ | GOTERM_BP_DIRECT | kidney deve | | | | | | |
| ☐ | GOTERM_BP_DIRECT | activation o | | | | | | |
| ☐ | GOTERM_BP_DIRECT | digestive tra | | | | | | |

Enrichment results of GO BP are not very similar to Figure 2….what can be the reason?
- Different list of genes
- Database versions

Also, I can be advisable to do it separately for up and down-regulated genes

# 3. Hands on: Enrichment analysis

Using the list of genes, you can perform a simple gene set enrichment using DAVID. Do you find similar gene sets enriched as in the paper?

| Sublist | Category | Term | RT | Genes | Count | % | P-Value | Benjamini |
|---------|----------|------|-----|-------|-------|-----|---------|-----------|
| ☐ | KEGG_PATHWAY | Focal adhesion | RT | ▪ | 17 | 3.6 | 1.5E-4 | 3.4E-2 |
| ☐ | KEGG_PATHWAY | PI3K-Akt signaling pathway | RT | ▬ | 21 | 4.4 | 1.1E-3 | 1.2E-1 |
| ☐ | KEGG_PATHWAY | Chemical carcinogenesis | RT | ▪ | 9 | 1.9 | 1.4E-3 | 1.0E-1 |
| ☐ | KEGG_PATHWAY | Amoebiasis | RT | ▪ | 10 | 2.1 | 2.4E-3 | 1.3E-1 |
| ☐ | KEGG_PATHWAY | Drug metabolism - cytochrome P450 | RT | ▪ | 8 | 1.7 | 2.4E-3 | 1.1E-1 |
| ☐ | KEGG_PATHWAY | Glutathione metabolism | RT | ▪ | 7 | 1.5 | 2.5E-3 | 9.1E-2 |
| ☐ | KEGG_PATHWAY | Protein digestion and absorption | RT | ▪ | 9 | 1.9 | 2.7E-3 | 8.4E-2 |
| ☐ | KEGG_PATHWAY | Metabolism of xenobiotics by cytochrome P450 | RT | ▪ | 8 | 1.7 | 3.9E-3 | 1.1E-1 |
| ☐ | KEGG_PATHWAY | Pathways in cancer | RT | ▬ | 21 | 4.4 | 5.1E-3 | 1.2E-1 |
| ☐ | KEGG_PATHWAY | ECM-receptor interaction | RT | ▪ | 8 | 1.7 | 9.4E-3 | 2.0E-1 |
| ☐ | KEGG_PATHWAY | AMPK signaling pathway | RT | ▪ | 9 | 1.9 | 1.9E-2 | 3.3E-1 |

KEGG seems to give more similar results

# 3. Hands on: GSEA

Use the complete ranked list of genes (TCGA data) to perform GSEA analysis.

- GSEA classic
  - You need to generate two files (.gct and .cls), which can be a bit tricky. Also, it is only available for microarray data

- **GSEA Preranked**
  - You need to generate one file (.rnk). Easy to generate with a simple formula (-log10*sign(FC)) and suitable for microarrays and RNAseq data

# 3. Hands on: GSEA

We will perform a GSEA Preranked using PRAD GEPIA2 data:

- Prepare rnk input file:
  - take column Gene and Rank from file PRAD_GEPIA2_table_degenes.txt and save it in a new file. Call it PRAD.rnk.
  - Or use the provided PRAD.rnk file
- Load data → Browse files → PRAD.rnk
- Run GSEAPreranked
- Select c5.bp.v7.0.symbols.gmt gene set database (or any other you want)
- Select "No collapse"
- Run and wait (output in folder my_analysis.GseaPreranked.1582108828269)

# TOP 10 GO BP terms: DAVID vs GSEA

## Can you make these plots in excel?

**TOP 10 GO terms DAVID**



Observe the different top results between the two methods.

Are they the same as in Figure 3?

Reproducibility is an issue!!!

If you want, try to do this comparison with KEGG

**TOP 10 GO terms GSEA UP**

# 3. Hands on: Network analysis

Upload the list of 478 DEGs in Cytoscape and populate it with a database (eg.STRING):
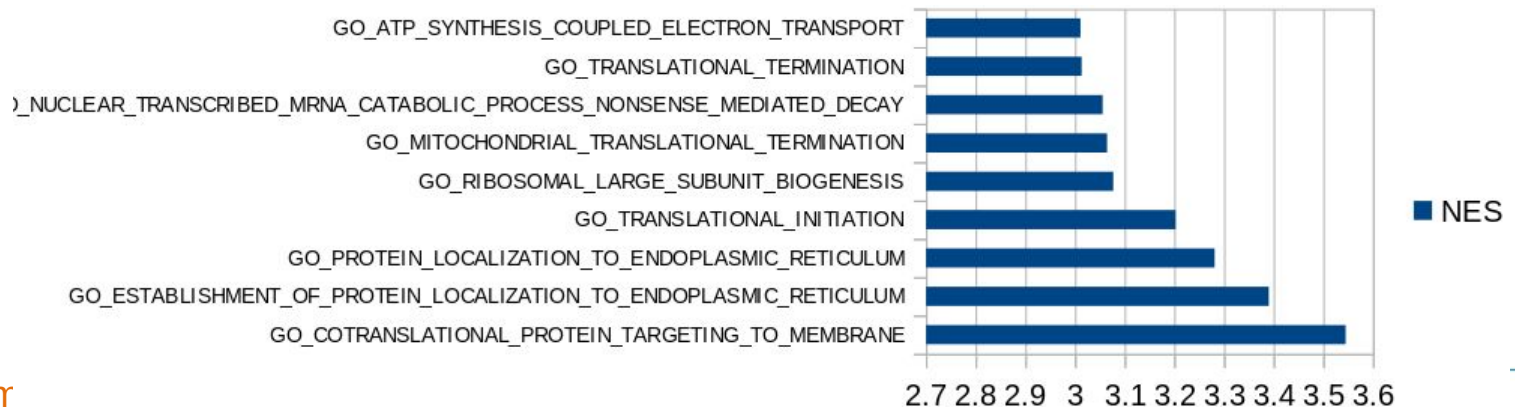
1. Install stringApp and MCODE app
2. File → New Network → New Network from Public Databases → Data Source=STRING:protein query
3. Copy/Paste the list of DEG (Up and Down) AND confidence score > 0.7
4. Tools → NetworkAnalyzer → Network Analysis → Analyse network
5. Select → Node:Degree "is" between 1 and 1000 → New Network from Selection (how many genes have you filtered? )
6. Import Table from File: PRAD_GEPIA2_table_degenes.txt
7. Style
   a. Default
   b. Fill Color using log2FC column
8. Save this network as "Network_STRING_0.7_1degree_GEPIA.cys"

# 3. Hands on: Network analysis

Upload the list of 478 DEGs in Cytoscape and populate it with a database (eg.STRING)

1. Install STR
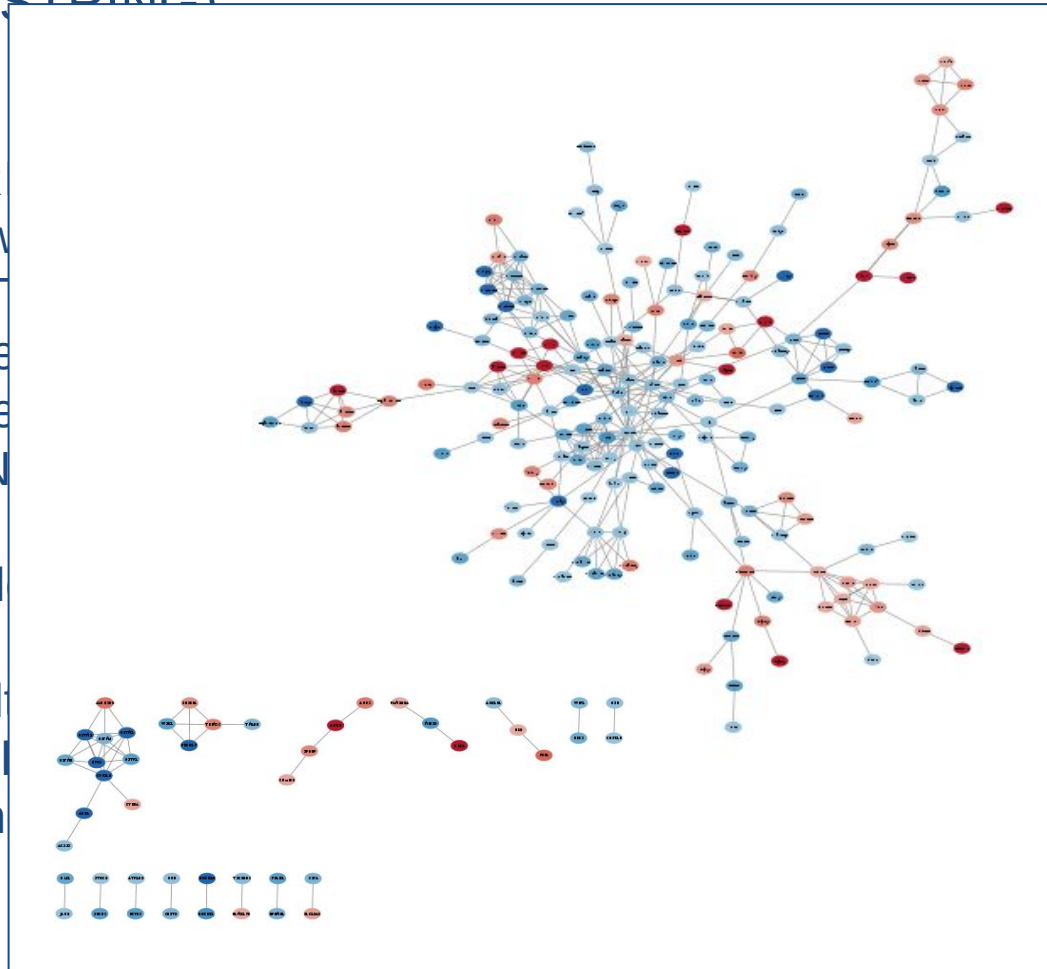2. File → New                                                     → Data
   Source=ST
3. Copy/Paste                                                     core > 0.7
4. Tools → Ne                                                     twork
5. Select → N                                                     vork from
   Selection
6. Import Tabl
7. Style
   a. Defaul
   b. Fill Col
8. Save this n                                                   A.cys"

# 3. Hands on: Network analysis

Integrate Mutations and CNA:

- Search in cBioPortal the frequency of mutations and copy number in PRAD and integrate the information to the network
  - Select study: Prostate Adenocarcinoma (TCGA, Firehose Legacy) → Explore selected studies
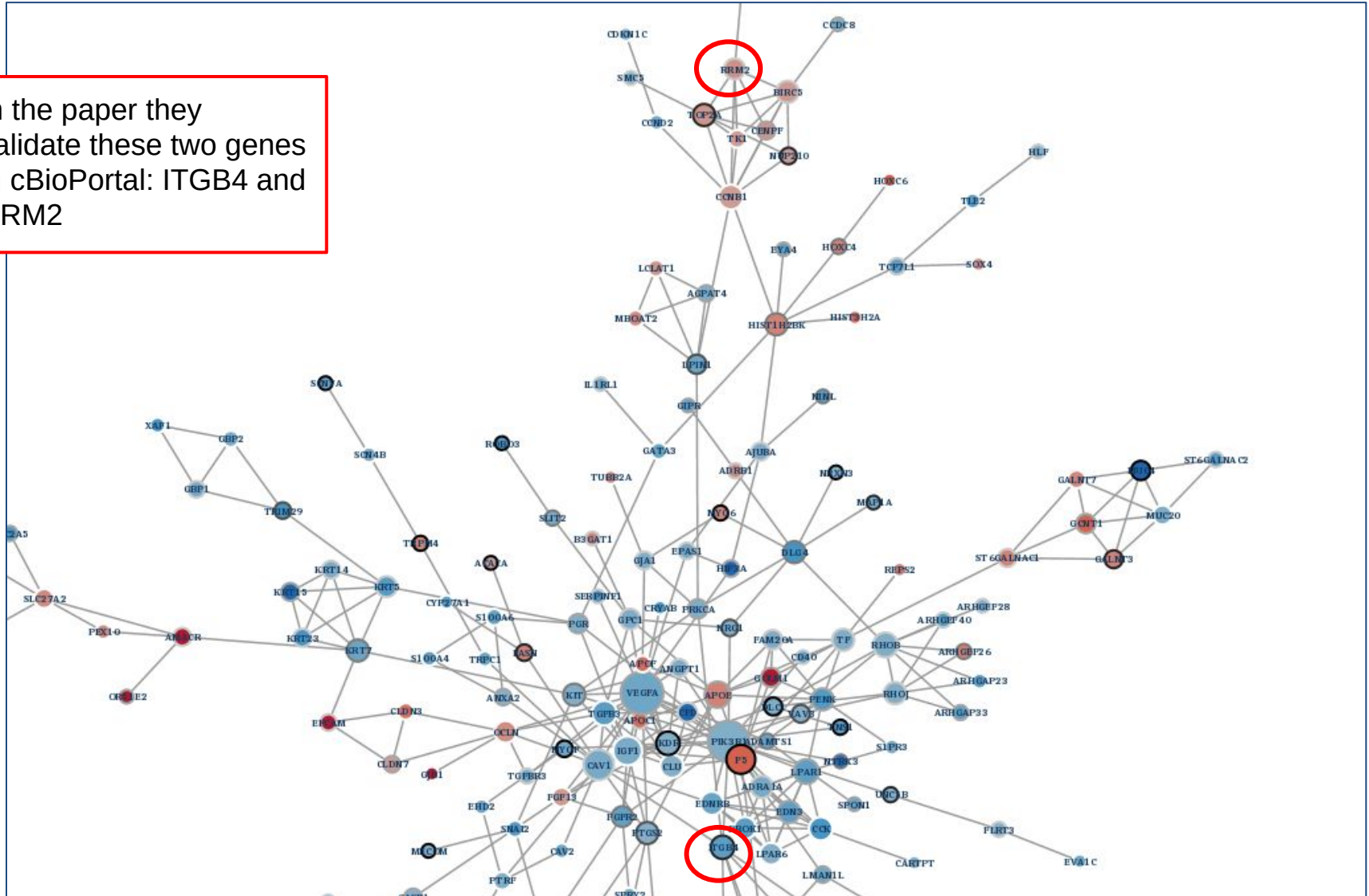
# 3. Hands on: Network analysis

Integrate Mutations and CNA:

- Search in cBioPortal the frequency of mutations and copy number in PRAD and integrate the information to the network

    – Select study: Prostate Adenocarcinoma (TCGA, Firehose Legacy) → Explore selected studies
    – Download mutations and CN data
    – Import Table from File: PRAD_Mutated_Genes.txt
    – Import Table from File: PRAD_CNA_Genes.txt
    – If you don't manage, open the network provided in file: "Network_STRING_0.7_1degree_GEPIA_cBioPortal.cys"

- Try to style the network using:
    – Style: Ripple
    – Size → Degree
    – Border Paint: # Mut
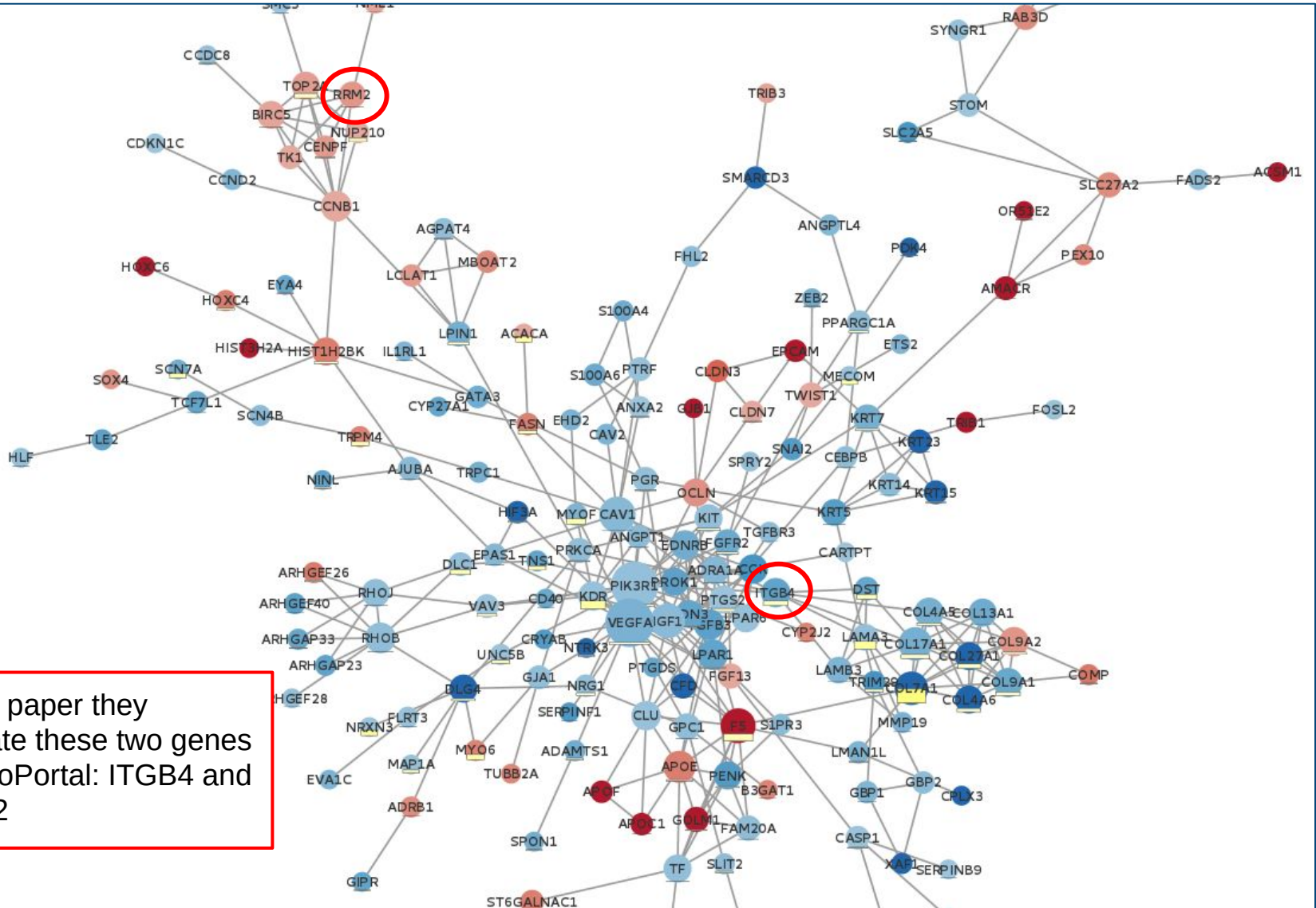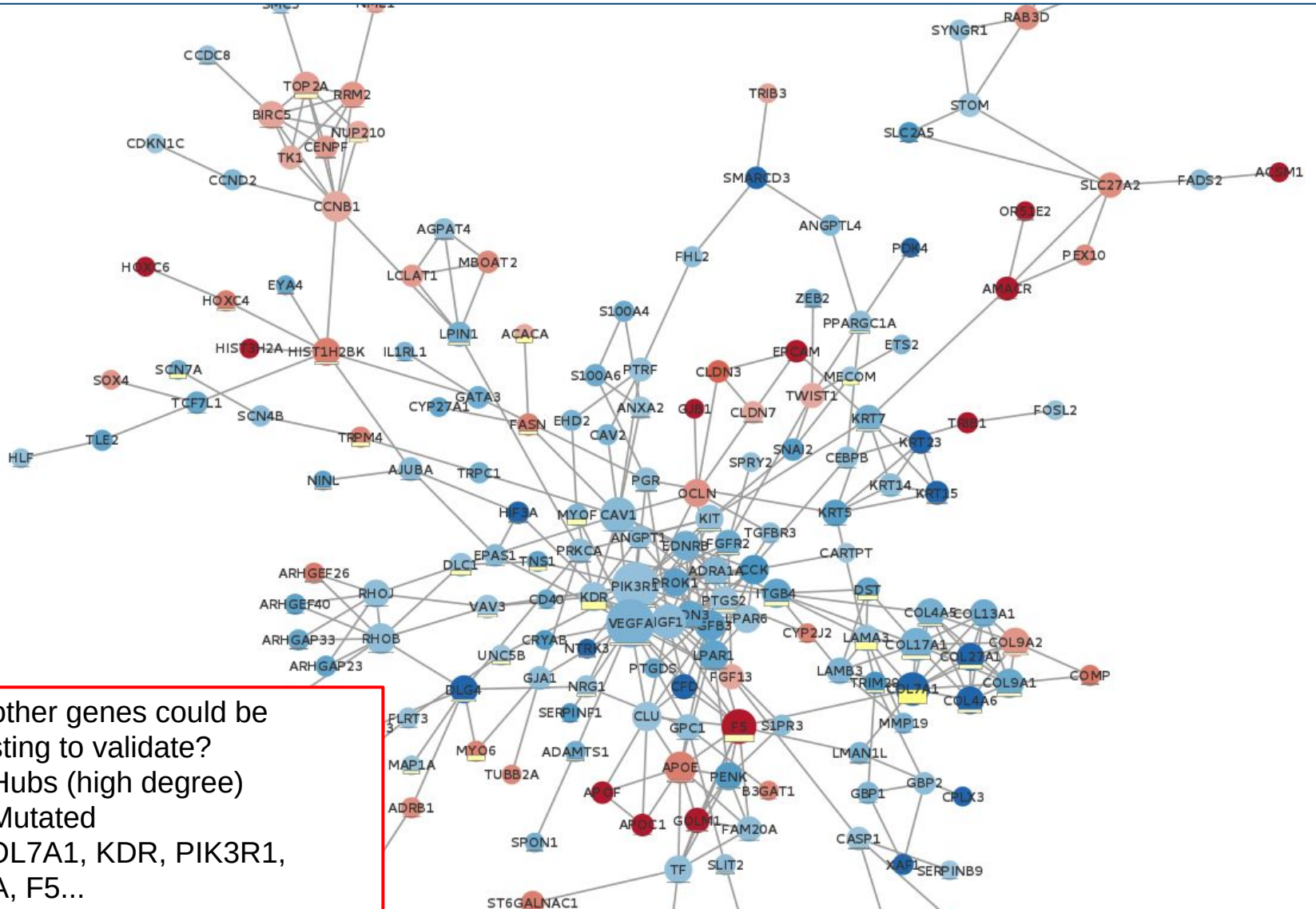    – Or Image/Chart → Bar → # Mut

# 3. Hands on: Network analysis

In the paper they validate these two genes in cBioPortal: ITGB4 and RRM2
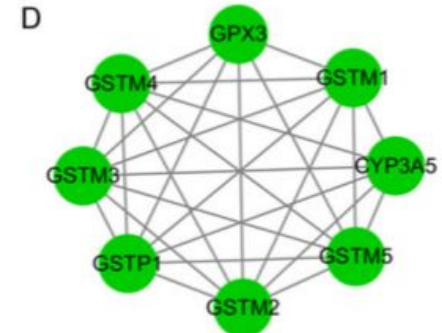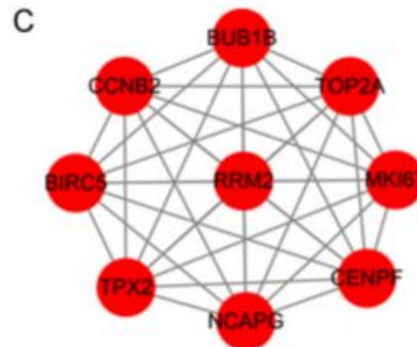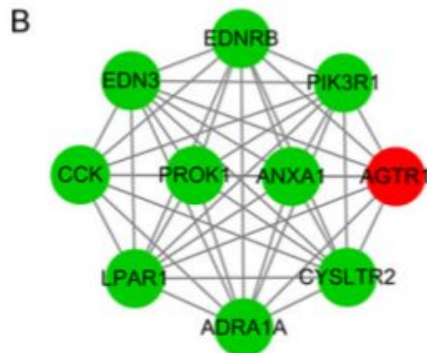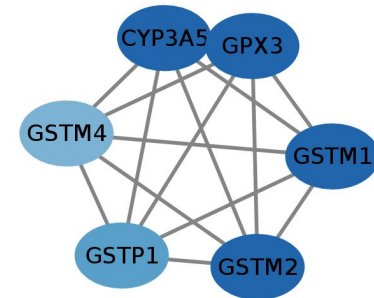
# 3. Hands on: Network analysis



In the paper they validate these two genes in cBioPortal: ITGB4 and RRM2

# 3. Hands on: Network analysis



What other genes could be
Interesting to validate?
- Hubs (high degree)
- Mutated
Eg: COL7A1, KDR, PIK3R1,
TOP2A, F5...

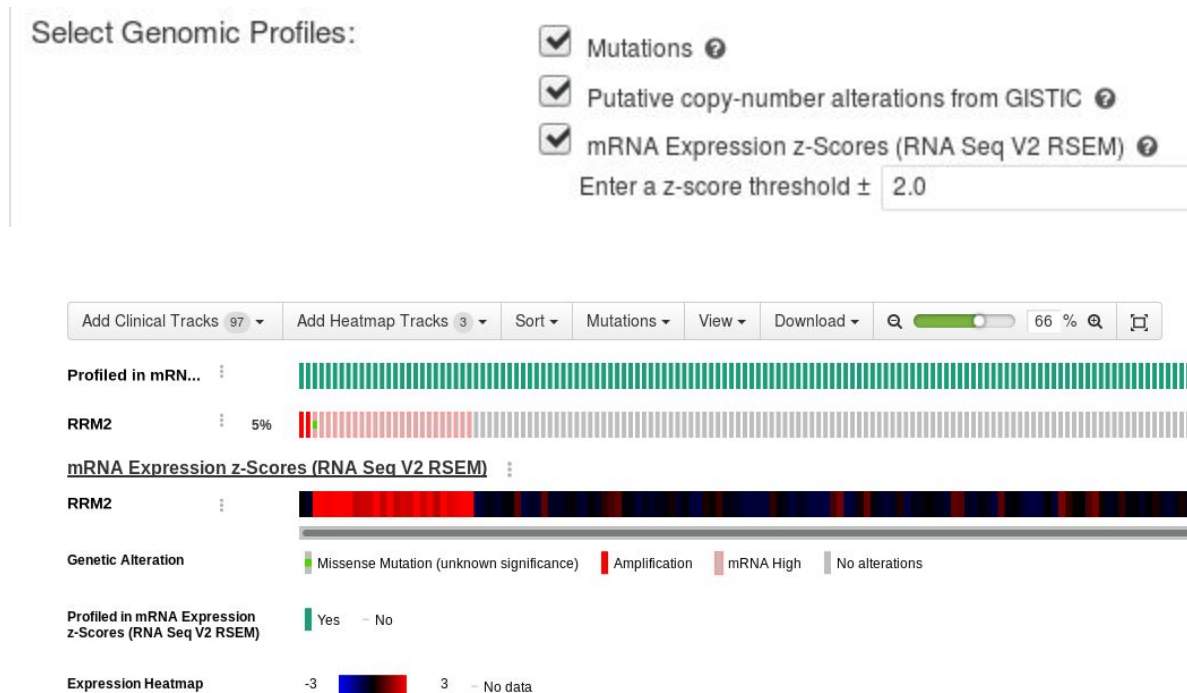# 3. Hands on: Network analysis

- MCODE → + → Analyze current network (default parameters)
- Do you find the modules in Figure 3?

# 3. Hands on: Network analysis
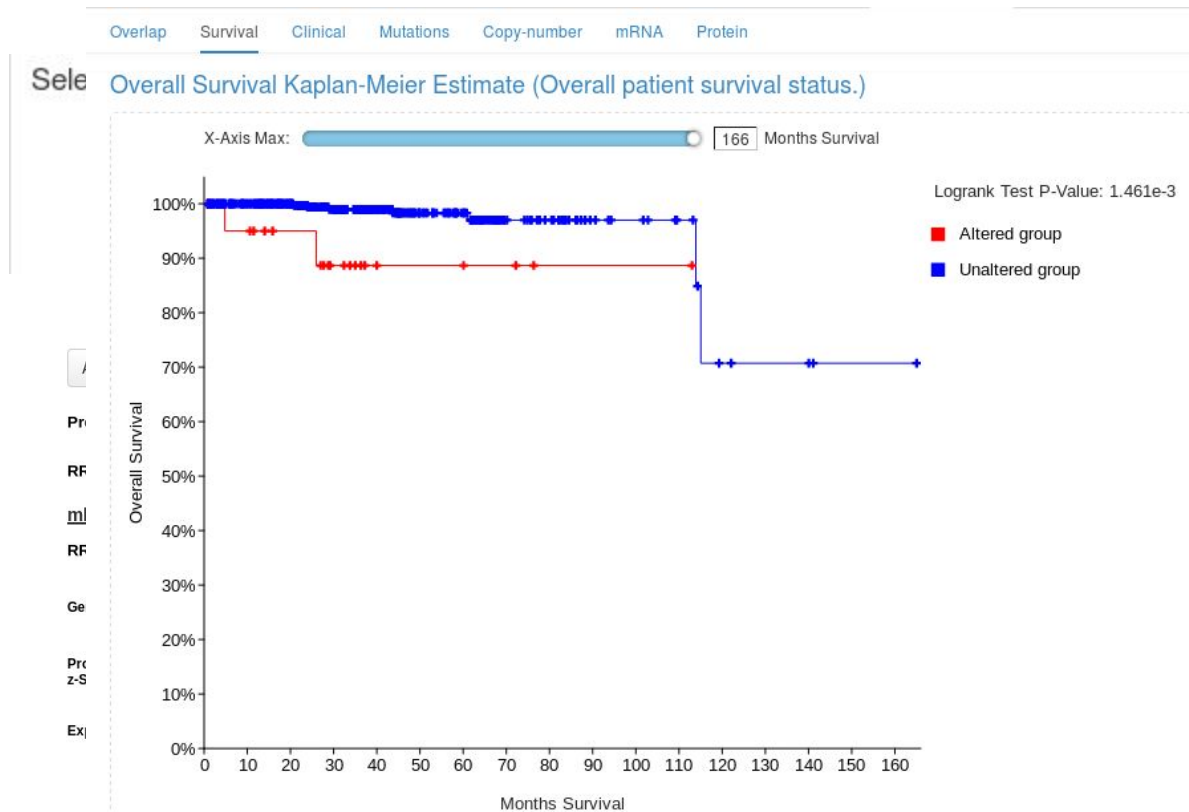
Prognostic value of identified genes:

- Search in cBioPortal

  – Select study: Prostate Adenocarcinoma (TCGA, Firehose Legacy) → Query Genes → **RRM2**

# 3. Hands on: Network analysis
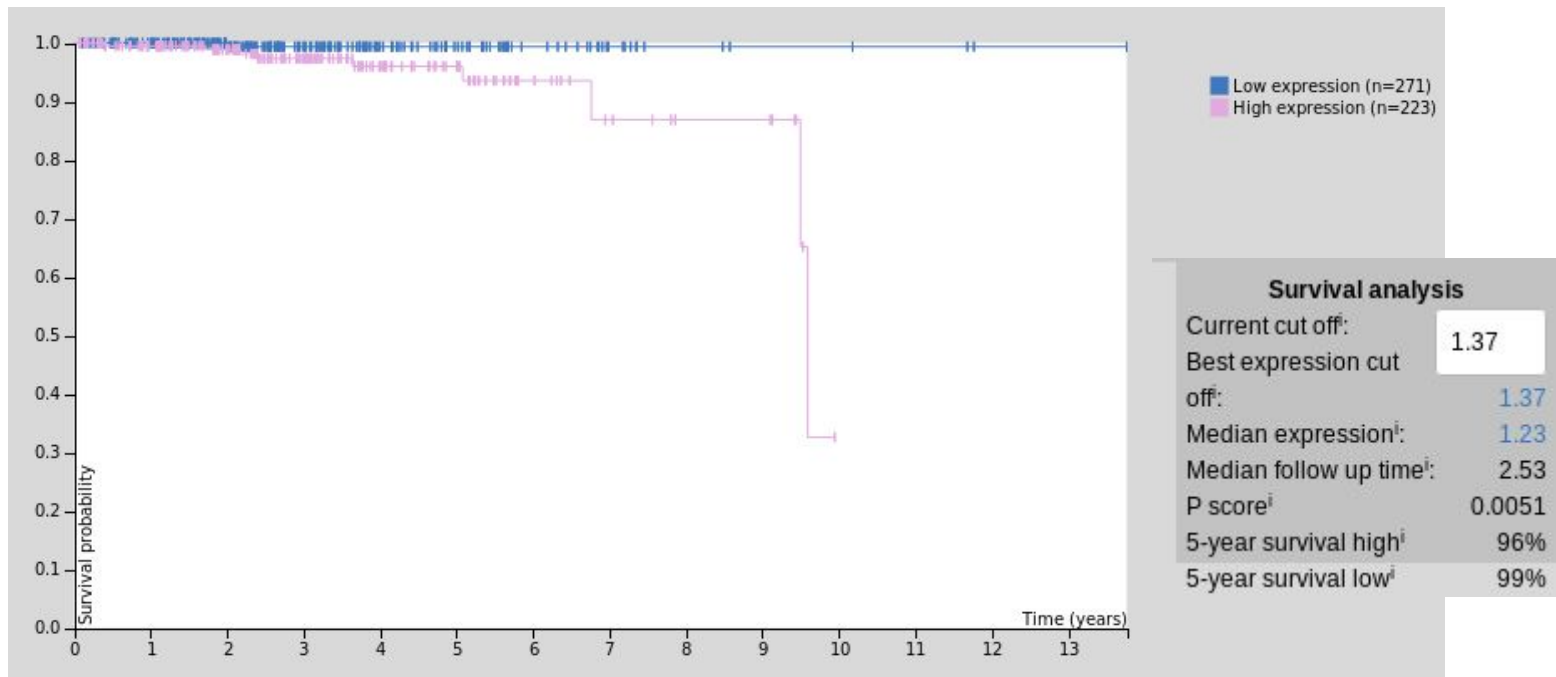
Prognostic value of identified genes:

- Search in cBioPortal

  – Go to Survival Tab

# 3. Hands on: Network analysis

Prognostic value of identified genes:

- Search in [The Human Protein Atlas](#) → RMM2 → Pathology → click on Prostate Cancer boxplot



**What is the difference between the two Kaplan-Meier curves?**

# Summary

After the course you should be able to…

- select the appropriate enrichment test for your data
- be aware of the different gene sets databases
- perform an enrichment test with a list of genes using DAVID
- perform GSEA
- do basic network manipulation with Cytoscape
- use some apps to extend the functionality
- understand that reproducibility is a huge issue
- publish a bioinformatics paper such as He et al!