# Bioinformatical analysis of omics expression data Part 5

Dr. Michael Turewicz[1,2]

[1]Institut für Klinische Biochemie und Pathobiochemie,
Deutsches Diabetes-Zentrum, Leibniz-Zentrum für Diabetesforschung an der Heinrich-Heine-Universität, Düsseldorf, Deutschland
[2]Deutsches Zentrum für Diabetesforschung (DZD), München-Neuherberg, Deutschland

# Course schedule

DDZ
Deutsches Diabetes-Zentrum

- Part 1 (25.10.23)
  - Introduction (omics, example data, programming)
  - Data preprocessing (data inspection, normalization, missing values)
  - Exercises: R programming tutorial (part 1)

- Part 2 (08.11.23)
  - Differential expression analysis (statistics, volcano plot)
  - Exercises: R programming tutorial (part 2)

- Part 3 (15.11.23)
  - Machine learning I: Clustering (clustering, PCA)
  - Exercises: Customized hierarchical clustering & PCA in R

- Part 4 (22.11.23)
  - Overrepresentation analysis (GO, Reactome)
  - Exercises: Own GO- & Reactome analysis in R & other tools

- **Part 5 (29.11.23)**
  - **Network analysis (STRING, Cytoscape)**
  - **Exercises: Own network analysis in R & STRING**

- Part 6 (06.12.23)
  - Machine learning II: Classification algorithms

# Recap of previous part

- **Gene Ontology (GO)**
  - Organism-specific hierarchy of curated biol. terms
  - → directed acyclic graph (DAG) of terms (= nodes)
  - → edges: "is a"- & "part of"-relationships
  - close to DAG-"root" general terms & terminal nodes most specific
  - Organized in 3 GO domains (separate DAGs): biological process, molecular function, cellular component
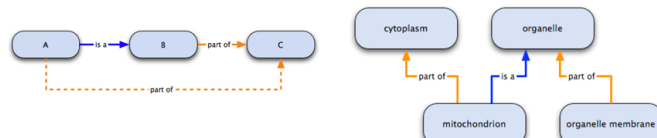
- **Overrepresentation analysis (ORA)**
  - **Basic idea:**
    1. Annotate input genes/proteins with the terms with which they are annotated in a biological database
    2. Return terms that are **statistically significantly (→ p-value) overrepresented** in input
  - **Statistical tests:** based on the urn model, e.g. Fisher's exact test, Kolmogorov-Smirnov test, (…)
  - Popular biological knowledge DBs: GO, Reactome, KEGG, WikiPathways, PhosphoSitePlus, (…)
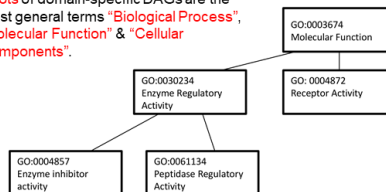
- **Programming (example: R)**
  - Own GO- & Reactome-based ORA + visualization in R & online tools

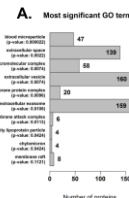- Based on "is a" or "part of" relationship

- Hierarchical relationship

Roots of domain-specific DAGs are the most general terms "Biological Process", "Molecular Function" & "Cellular Components".
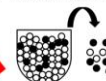
GO:0003674 Molecular Function
GO:0030234 Enzyme Regulatory Activity
GO:0004872 Receptor Activity
GO:0004857 Enzyme inhibitor activity
GO:0061134 Peptidase Regulatory Activity

More General
More Specific

**A.** Most significant GO terms
**B.**

Volcanoplot

- log10 p-value
log2 ratio of means

> Fisher's Exact Test (Hypergeometric Test)

The test implements the urn model.

What is the probability of getting 7 or more black balls?

$$P = 1 - \sum_{i=0}^{k-1} \frac{\binom{M}{i}\binom{N-M}{n-i}}{\binom{N}{n}}$$

- **N:** total number of proteins
- **M:** total number of proteins annotated with this term
- **n:** number of proteins in the set (all balls drawn)
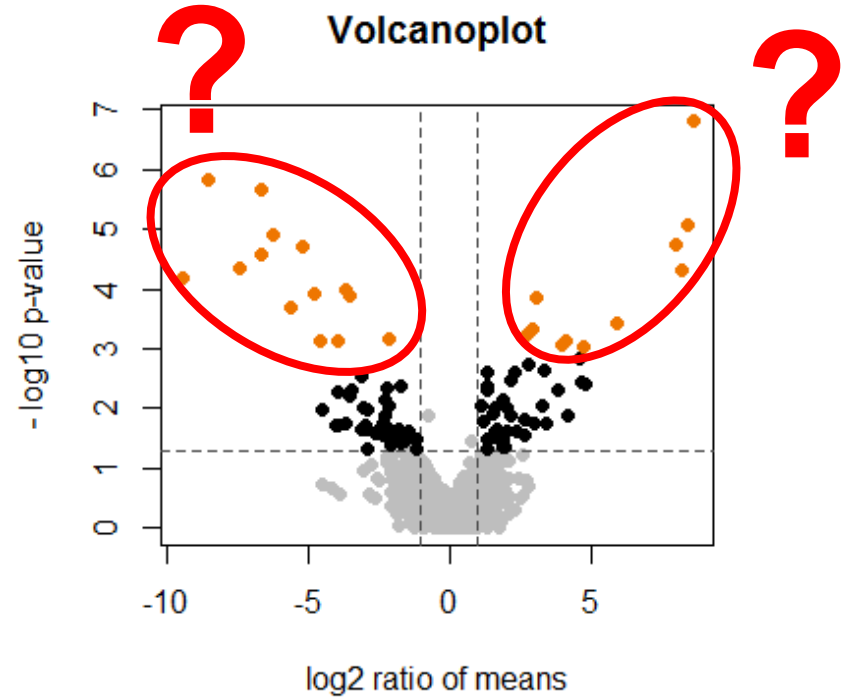- **k:** number of proteins in tested set annotated with this term (black balls drawn)

ca. 20,000 proteins

GO term 1 (many proteins)
GO term 3 (100)
tested protein set (25)
GO term 4 (20)
GO term 2 (average no. of proteins)

# List of candidates: what's next?
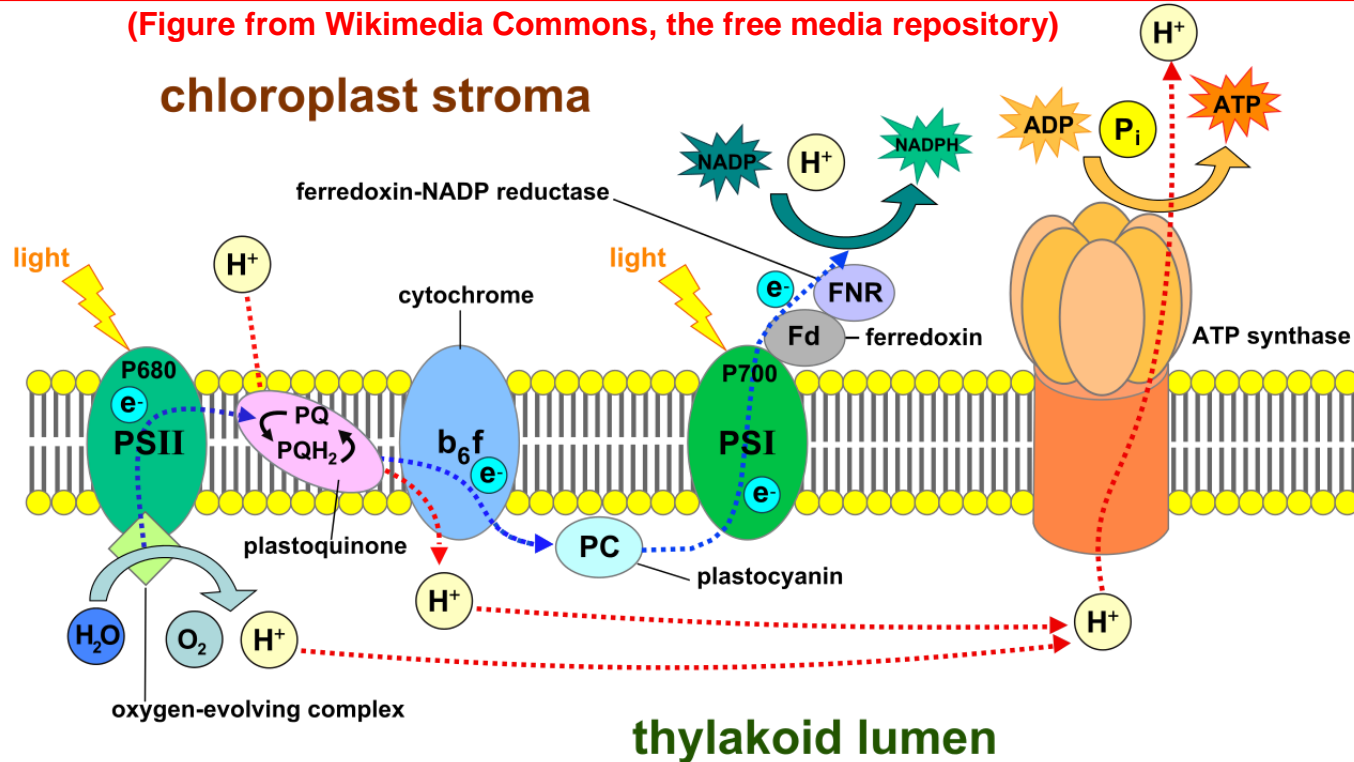
**DDZ**
Deutsches Diabetes-Zentrum

- We have learned to find a list of statistically significant differential candidates with p-values and fold changes

- How can we interpret these biologically?

- Are there biological connections that could explain a common occurrence?

- Can a common function be assumed?



3

# Motivation: Protein-protein interactions (PPIs)



**Example: Photosynthetic electron transport chain of the thylakoid membrane.**

(Figure from Wikimedia Commons, the free media repository)

# PPI networks: basic concepts

**PPI networks**

## Constructed from DBs

PPI networks can be constructed using knowledge from biological databases. Popular tools include:

- **STRING (free)** ✦STRING
- IntAct (free)
- Ingenuity Pathway Analysis (commercial)
- (…)

## Inferred from omics data

1) PPI prediction via amino acid sequences, esp. based on conserved sequences & well known PPIs in other species (**not discussed today**)

2) PPI prediction via the co-expression of genes/proteins. → **(weighted) gene/protein co-expression network analysis (WGCNA)**

# Networks from DBs: STRING

**DDZ**
Deutsches Diabetes-Zentrum

- **STRING (https://string-db.org/):**
  - online tool for PPI network analysis (& more).
  - widely used & free.

- STRING is a DB of **known & predicted PPIs**.

- Interactions are derived from:
  - Genomic context predictions
  - High-throughput lab experiments
  - (Conserved) co-expression
  - Text mining (PubMed, OMIM, …)
  - Knowledge in DBs (Reactome, DIP, BioGRID, MINT, Gene Ontology, KEGG, PDB, UniProt, …)

- Version 11.0:
  - 5090 organisms
  - > 24.6 millions proteins
  - > 3,000 millions interactions

- → **Search for PPIs for single or multiple proteins & visualization as a network.**

6

# STRING: search form



**Candidate list: gene/protein names or IDs (here: UniProt IDs)**

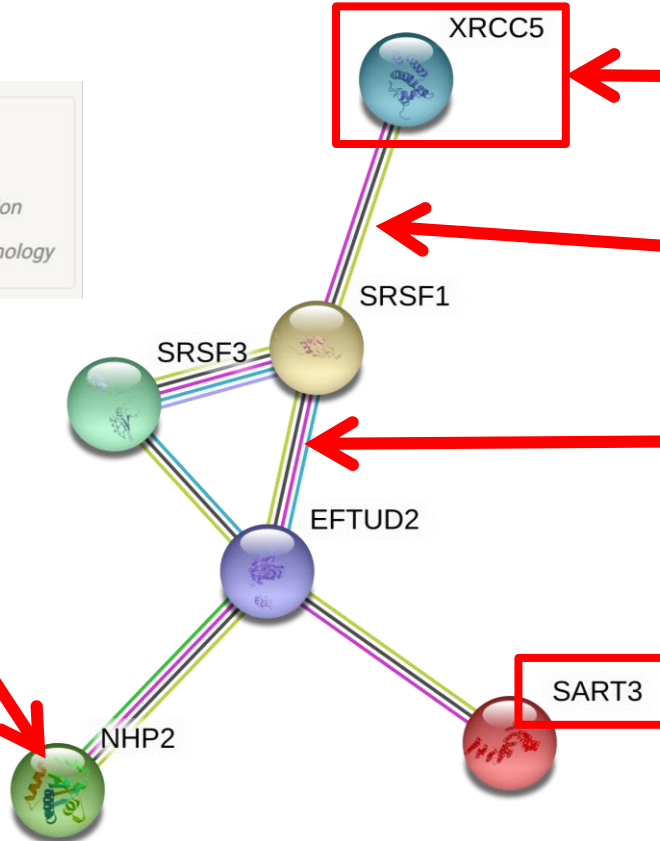**Various search modes available, e.g. "multiple proteins" to analyze list of candidates**

**Organism**

# STRING: network (evidence view)



8

# STRING: basic settings

**Different network views:**

- **evidence view: edges show types of interaction evidence**
- **confidence view: edges represent confidence score**

meaning of network edges:

- evidence ( ⬤—⬤ line color indicates the type of interaction evidence )
- confidence ( ⬤—⬤ line thickness indicates the strength of data support )
- molecular action ( ⬤—⬤ line shape indicates the predicted mode of action )

**Interaction sources can be included/excluded into/from analysis**

active interaction sources:

- ☑ Textmining  ☑ Experiments  ☑ Databases  ☑ Co-expression
- ☑ Neighborhood  ☑ Gene Fusion  ☑ Co-occurrence

**Edges below this confidence score are not shown.**

minimum required interaction score:

medium confidence (0.400) ⇕

**Possible options:**
- **0.9 (highest)**
- **0.7 (high)**
- **0.4 (medium)**
- **0.15 (low)**
- **Custom value**

max number of interactors to show:

1st shell: - none / query proteins only - ⇕

2nd shell: - none - ⇕

**Optional network-extension by proteins not included in input list: in 2 layers a user-defined number of most confident interactors from STRING DB can be added.**
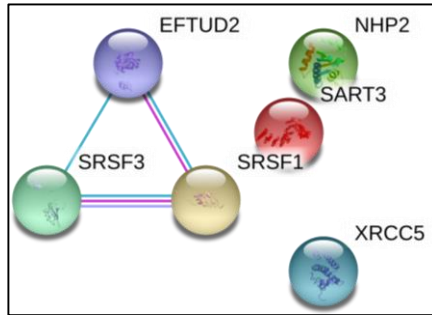
9

# STRING: network (confidence view)

# STRING: lax settings



Settings: "low confidence (0.150)" & all kinds of interactions

Low confidence (> 0.15)

# STRING: strict settings



Settings: "highest confidence (0.900)" & only known interactions (DBs & experiments)

12

# STRING: network statistics & ORA



13

# PPI networks: basic concepts

**DDZ**
Deutsches Diabetes-Zentrum

## PPI networks

### Constructed from DBs

PPI networks can be constructed using knowledge from biological databases. Popular tools include:

- **STRING (free)**
- IntAct (free)
- Ingenuity Pathway Analysis (commercial)
- (…)

### Inferred from omics data

1) PPI prediction via amino acid sequences, esp. based on conserved sequences & well known PPIs in other species (**not discussed today**)

2) PPI prediction via the co-expression of genes/proteins. → **(weighted) gene/protein co-expression network analysis (WGCNA)**

# Example dataset

- **Hepatocellular carcinoma (HCC) → liver cancer**

- **19 HCC vs. 19 controls (C) → healthy samples**

- **Obtained from Naboulsi et al., J. Proteome Res. 2016**

- **PRIDE: PXD002171**

- **Label-free quantification of 2,736 proteins via Progenesis QI software**

# Protein expression profiles



Q15029 and Q15020 vs. P13797

# Protein expression profiles



Q15029 and Q15020 vs. P13797

# Protein expression profiles



Q15029 and Q15020 vs. P13797

# Measure for co-expression: Spearman's rho



**Q15020 vs. Q15029**

# Measure for co-expression: Spearman's rho



r = 0.951
(both are involved in pre-mRNA splicing)

# Measure for co-expression: Spearman's rho
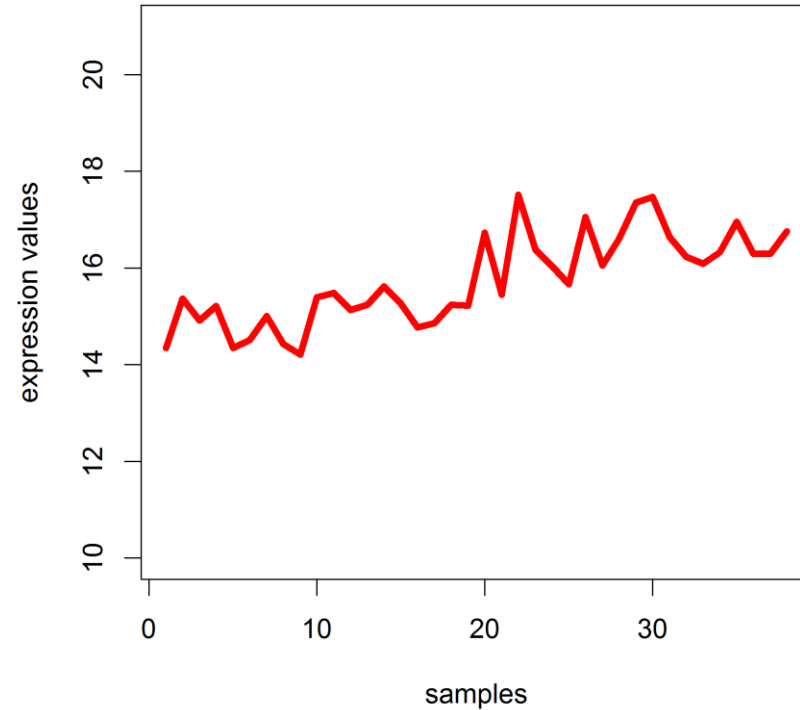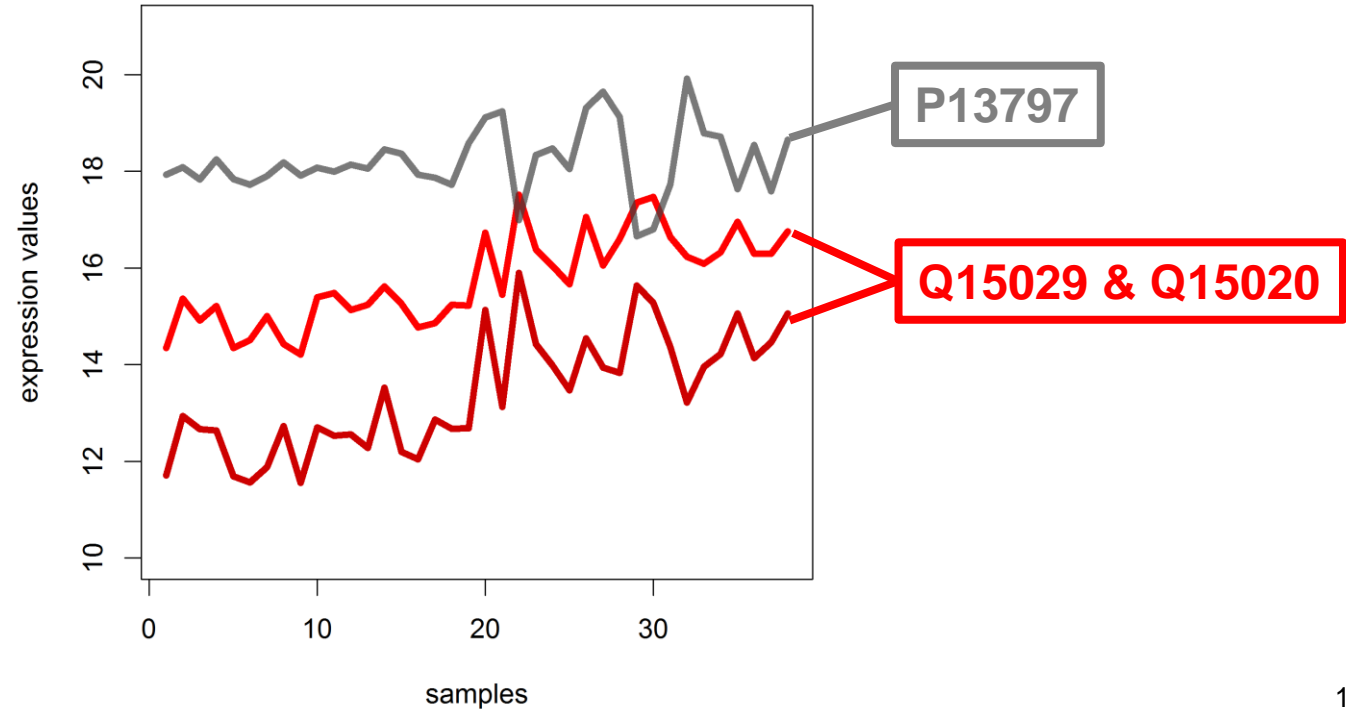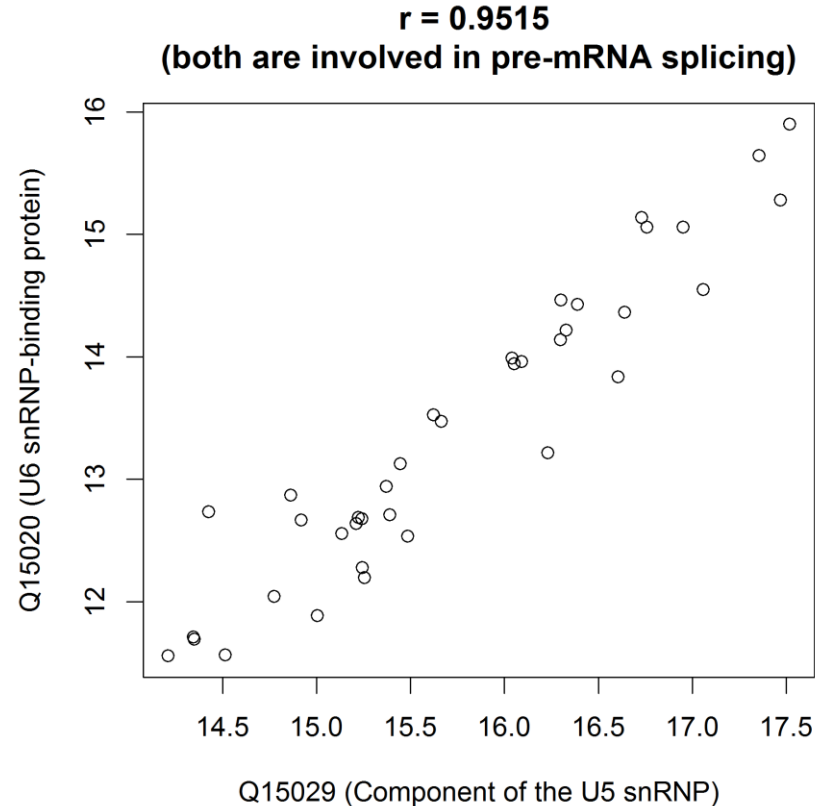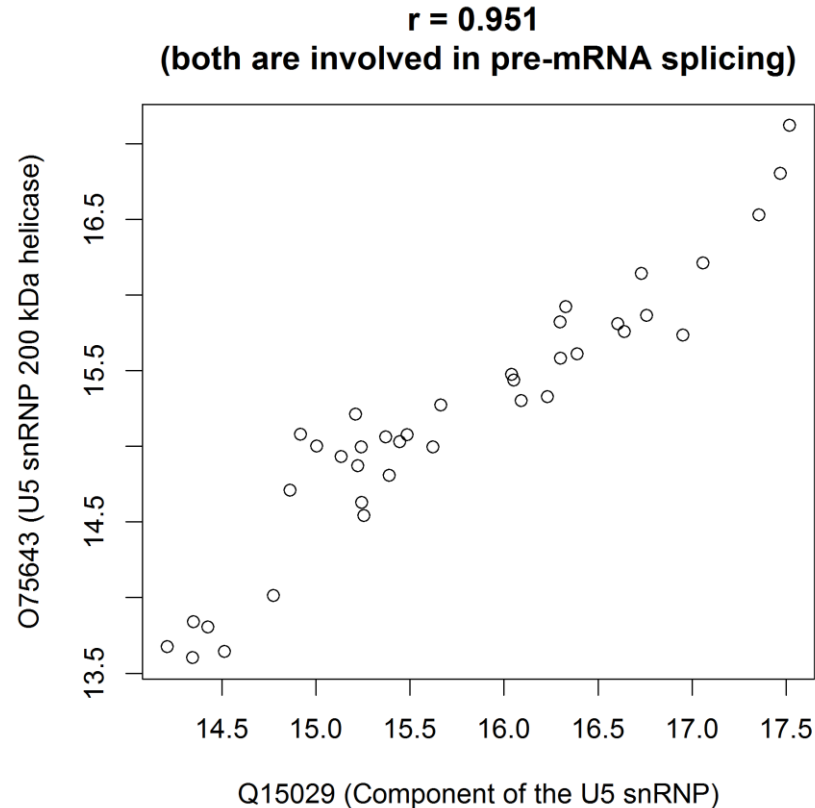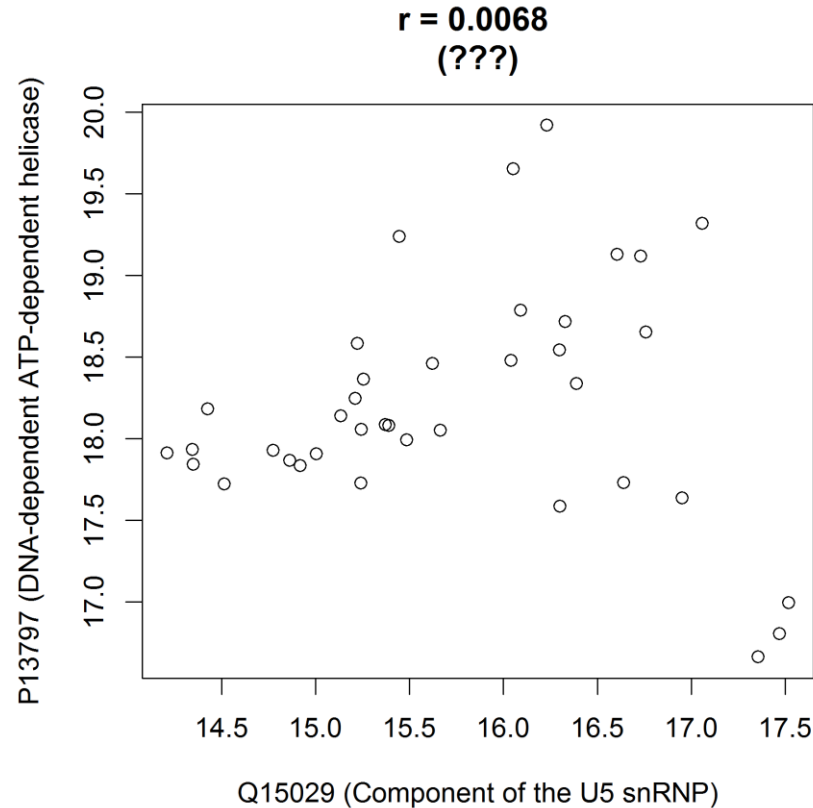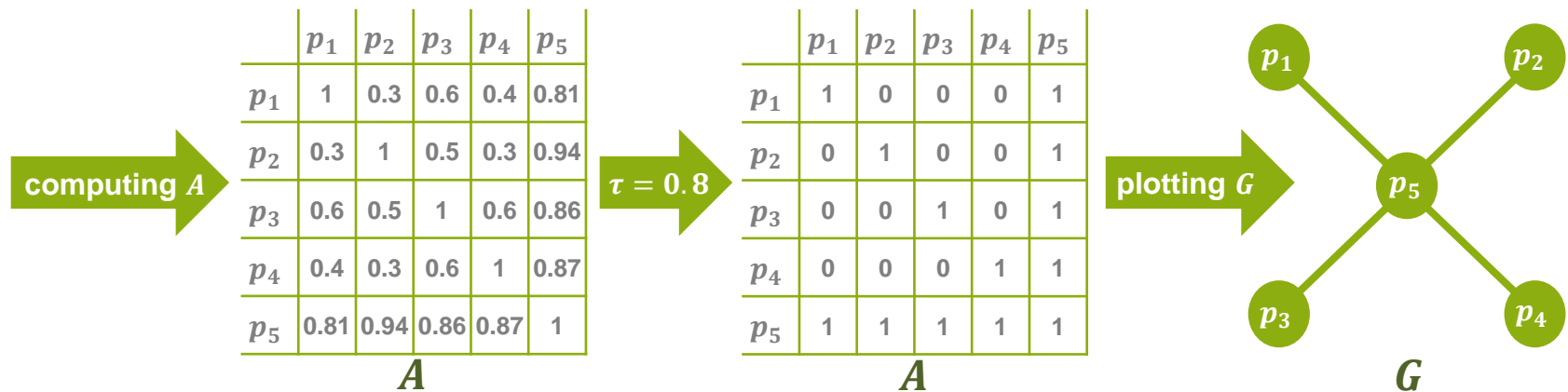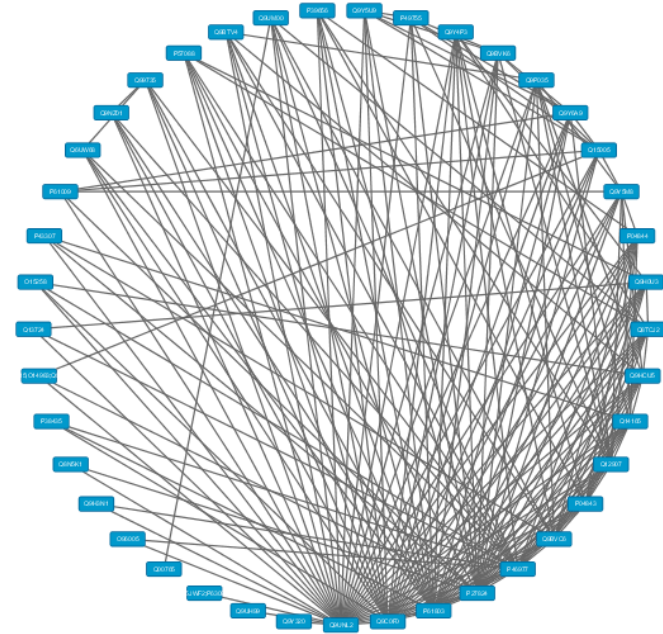


P13797 vs. Q15020

# Protein co-expression network inference

- Let $P = \{p_1, \dots, p_5\}$ be a set of protein expression profiles (from quant. proteomics)
- Consider pairs of protein expression profiles $p_i, p_j \epsilon P$.
- Correlation-based measure for co-expression: $a_{ij} = \left| cor(p_i, p_j) \right|^\beta, \beta \epsilon \{1,2,3, \dots\}$
- Adjacency matrix: $A = [a_{ij}] = \begin{pmatrix} a_{11} & \cdots & a_{15} \\ \vdots & \ddots & \vdots \\ a_{51} & \cdots & a_{55} \end{pmatrix}$
- Unweighted network: $G = (P, A)$ after setting $a_{ij} = \begin{cases} 1 \ \forall a_{ij} \geq \tau \\ 0 \ \forall a_{ij} < \tau \end{cases}$

1 = "edge"
0 = "no edge"



**computing $A$**

|       | $p_1$ | $p_2$ | $p_3$ | $p_4$ | $p_5$ |
|-------|-------|-------|-------|-------|-------|
| $p_1$ | 1     | 0.3   | 0.6   | 0.4   | 0.81  |
| $p_2$ | 0.3   | 1     | 0.5   | 0.3   | 0.94  |
| $p_3$ | 0.6   | 0.5   | 1     | 0.6   | 0.86  |
| $p_4$ | 0.4   | 0.3   | 0.6   | 1     | 0.87  |
| $p_5$ | 0.81  | 0.94  | 0.86  | 0.87  | 1     |

$A$

**$\tau = 0.8$**

|       | $p_1$ | $p_2$ | $p_3$ | $p_4$ | $p_5$ |
|-------|-------|-------|-------|-------|-------|
| $p_1$ | 1     | 0     | 0     | 0     | 1     |
| $p_2$ | 0     | 1     | 0     | 0     | 1     |
| $p_3$ | 0     | 0     | 1     | 0     | 1     |
| $p_4$ | 0     | 0     | 0     | 1     | 1     |
| $p_5$ | 1     | 1     | 1     | 1     | 1     |

$A$

**plotting $G$**

$G$

22

# Protein co-expression networks

- **Module:** Cluster of interconnected (i.e., co-expressed) proteins.

- Modules can represent **pathways** & **hub proteins** (most interconnected proteins) can represent their modulators.

- **Topological Overlap:** Similarity measure for the "interconnectedness" between two proteins based on the number of shared neighbors.

- **TOM (Topological Overlap Matrix):** Matrix $\Omega = [\omega_{ij}]$ containing the topological overlap between all proteins. Used for the detection of modules.
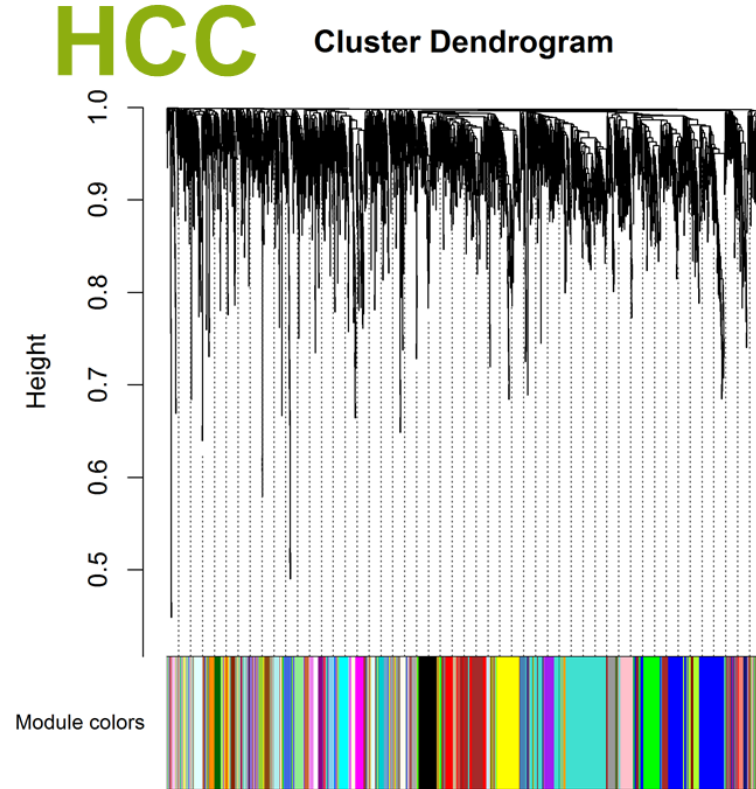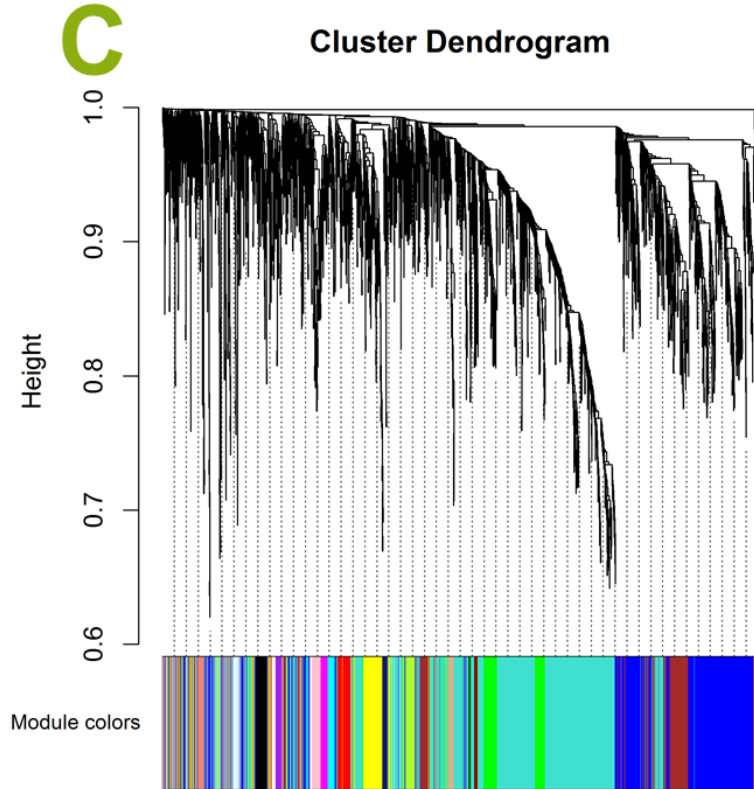


**Distance metric for clustering.**

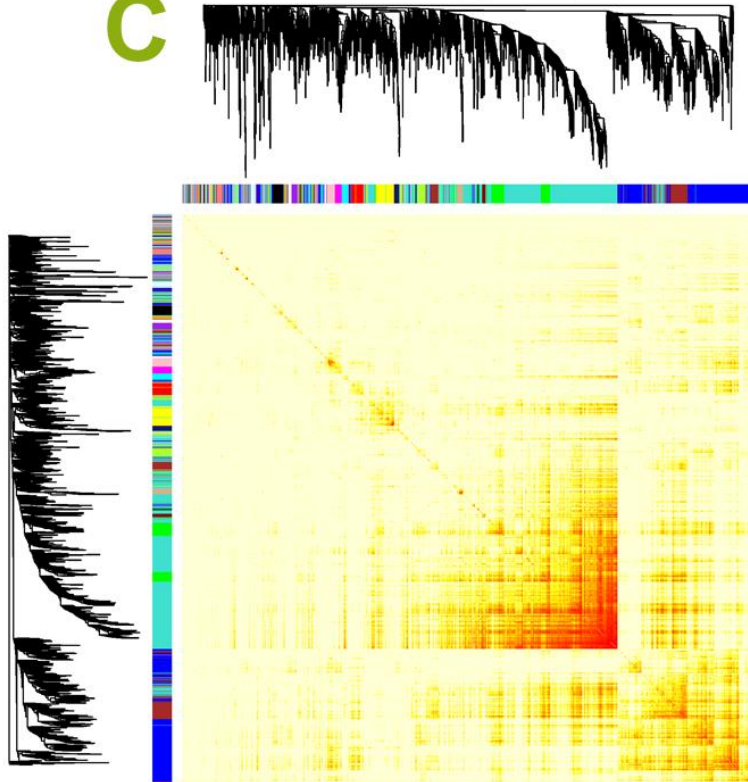$$\omega_{ij} = \frac{\sum_u a_{iu} a_{uj} + a_{ij}}{min\{k_i, k_j\} + 1 - a_{ij}}$$

23

# Comparing group-specific networks
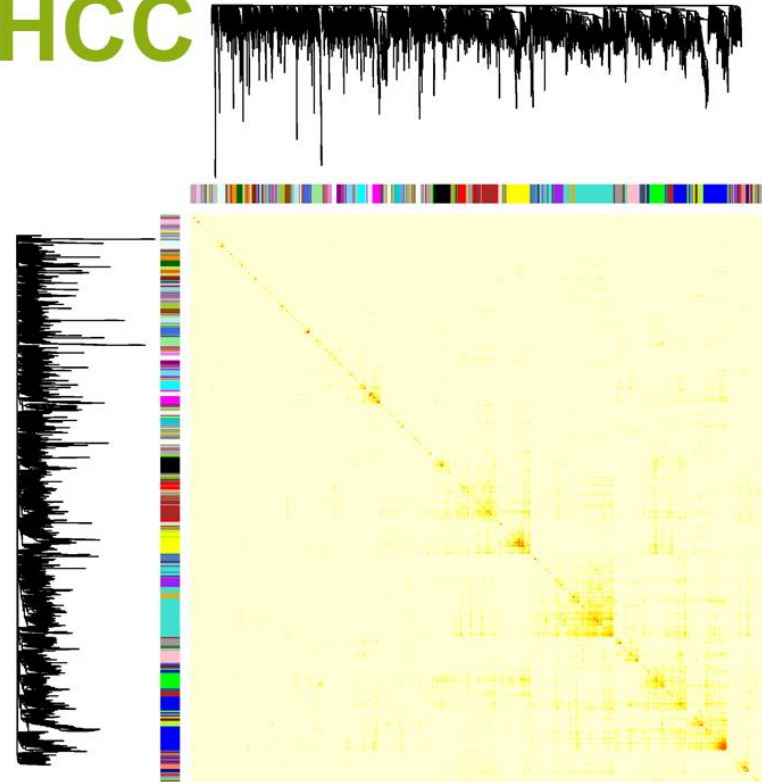
# Comparing group-specific networks

# Comparing group-specific networks
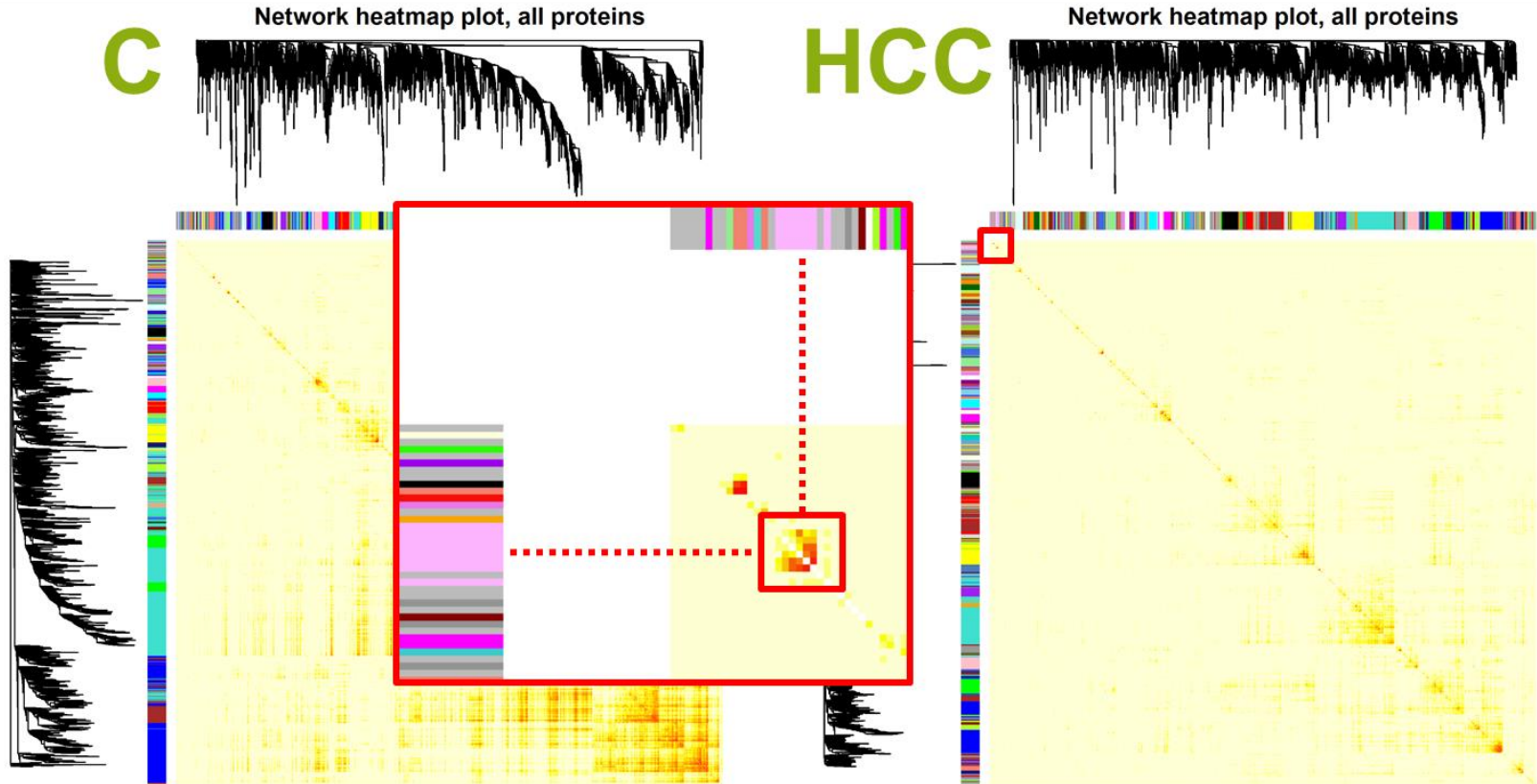


Network heatmap plot, all proteins

Network heatmap plot, all proteins

C

HCC

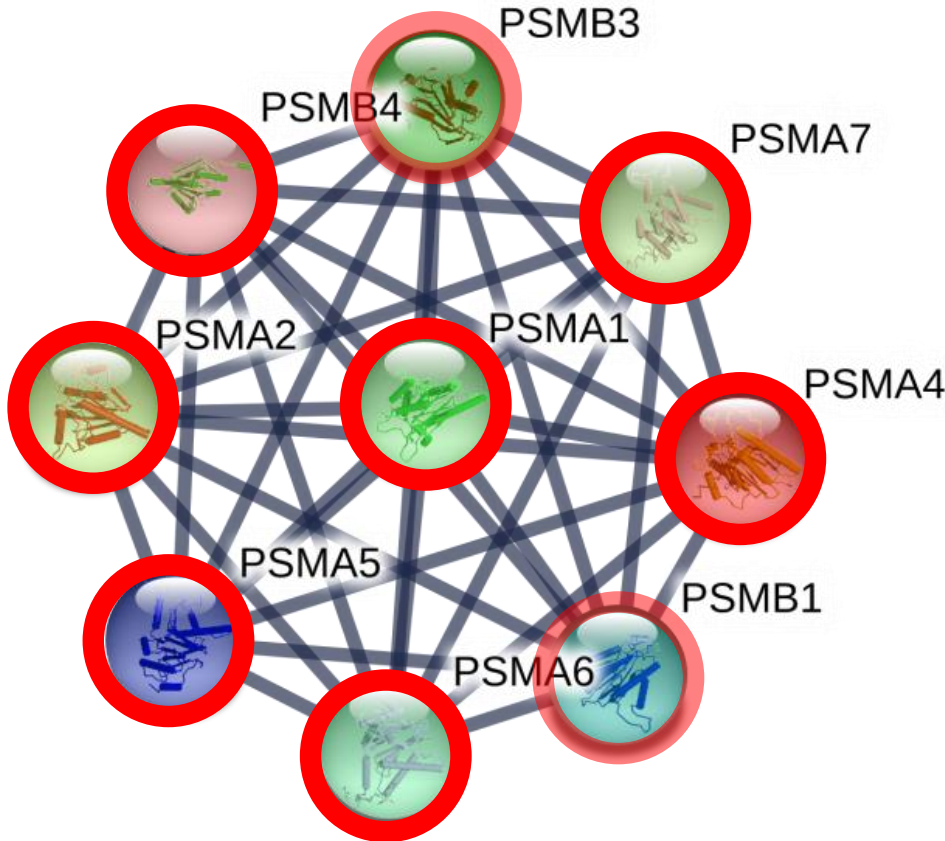# (Sub-)network visualization



Hub protein

# Network interpretation (STRING)



**Enrichment results (adjusted p-values):**

- GO (cellular component): "**proteasome core complex, alpha subunit complex**" (p = 3.66e-14)
- KEGG pathway: "**proteasome**" (p = 1.54e-12)
- PFAM domains: "**proteasome subunit**" (p = 1.3e-17)
- INTERPRO domains: "**Proteasome, subunit alpha/beta**" (p = 2.77e-17)

**Components of the ubiquitin-proteasome pathway are known targets for cancer therapy (proteasome inhibitors) – also in discussion for HCC…**
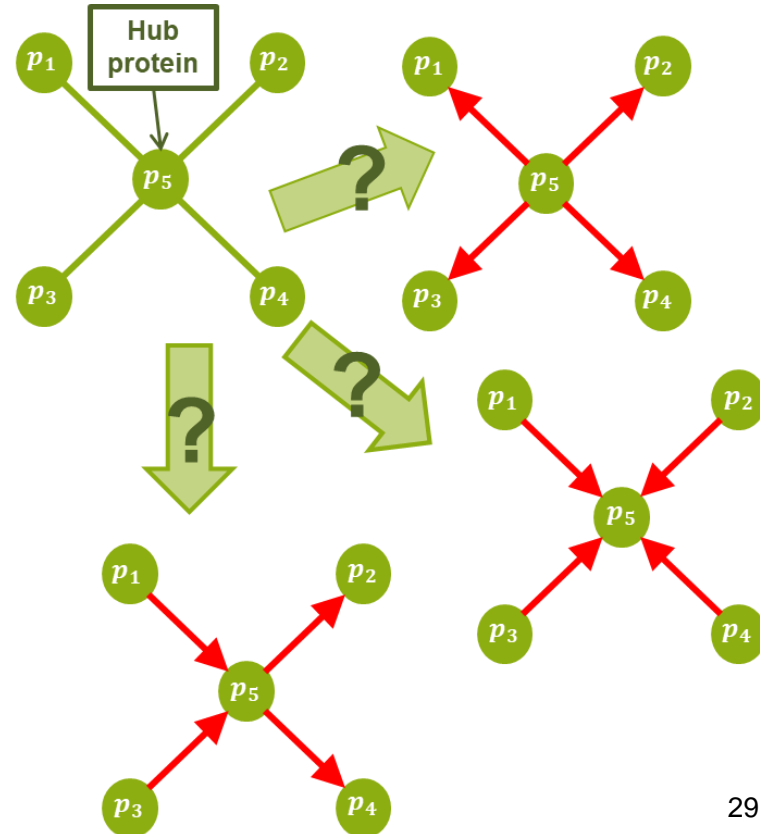
○ **= proteasome subunit**

# Network interpretation

**DDZ** Deutsches Diabetes-Zentrum

**Potential reasons for co-expression:**
- Direct protein-protein interaction
- Target protein ←→ protease
- Target protein ←→ kinase/phosphorylase
- Receptor protein ←→ effector protein
- Signaling complexes
- Scaffold protein complexes
- Target protein ←→ transcription factor
- (…)

**Elucidation by:**
- Best way: experiments!
- Protein annotation (e.g., GO-/Reactome-ORA)
- PPI annotation
- Protein module identification & characterization
- Identification & characterization of hub proteins
- Differential network analysis
- (…)



29

# Hands on part!

# Exercises

- ## **Exercise 5**
  - https://drive.google.com/drive/folders/1vmewprs0gkpakU8idbgtexDIwmGVUJz3?usp=sharing
  - Use our example dataset from GitHub for the following exercises
  - **Exercise 5.1:** Perform an own STRING network analysis in R using differential candidates and confidence view (confidence threshold = 0.7). Visualize the network without not connected nodes.
  - **Exercise 5.2:** Perform an own WCGNA in R and find an biologically interesting module. You can use STRING analysis (via searching the UniProt IDs in the STRING web application) to quickly check whether a module is interesting (i.e. highly confident STRING interactions & interesting ORA results in STRING).
  - Please send me your solutions as an ".R"-file

# Thank you!