# Clustering group exercise

- Each group will analyze a cluster of genes identified in the T cell expression table

- Analyze the table of top 100 genes **by variance** in 47 samples

- Cluster them using:
  - Group 1: UPGMA = 'linkage', 'average', 'RowPDistValue', 'euclidean',
  - Group 2: 'linkage', 'single', 'RowPDistValue', 'cityblock',
  - Group 3: 'linkage', 'average', 'RowPDistValue', 'correlation',
  - Group 4: UPGMA = 'linkage', 'single', 'RowPDistValue', 'euclidean',
  - Group 5: UPGMA = 'linkage', 'weighted', 'RowPDistValue', 'correlation',

- Use clustergram(…, 'Standardize','Row', 'linkage', ***as specified for your group,*** 'RowPDistValue' ***as specified for your group,*** 'RowLabels',gene_names1,'ColumnLabels', array_names)
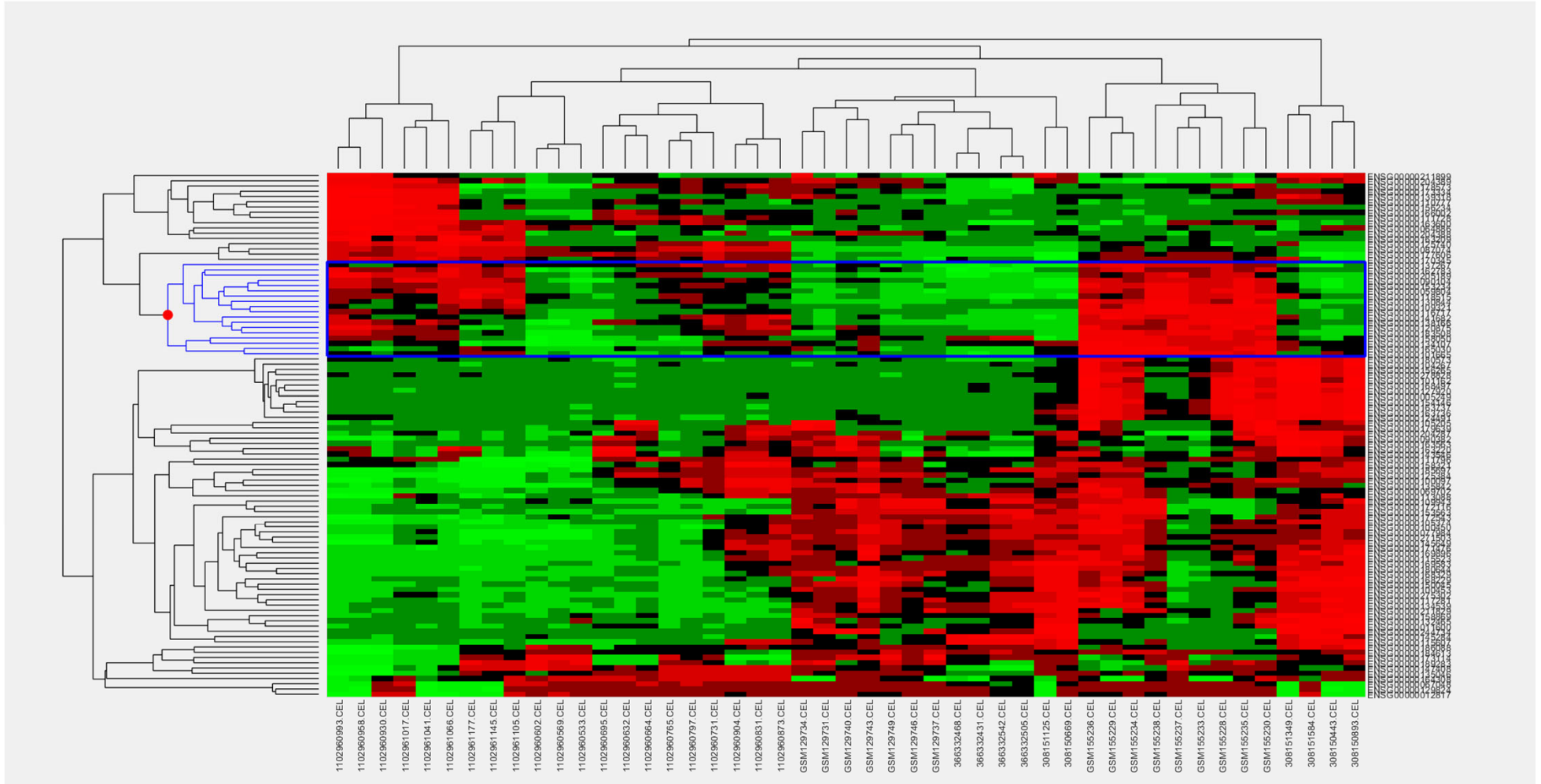
# Cluster analysis group exercise

- Which biological functions are overrepresented in different clusters?

- Pick a cluster:
  - Select a node on the tree of rows,
  - Right click
  - Choose "export group info" into the workspace
  - Name it gene_list

- Run the following two Matlab commands to display genes
  - g1=gene_list.RowNodeNames;
  - for m=1:length(g1); disp(g1{m}); end;

# Search for shared biological functions

- copy the list of displayed genes
- go to "Start Analysis" on https://david.ncifcrf.gov/tools.jsp
- Paste genes from gene list displayed by Matlab into the box in the left panel of the website
- select ENSEMBL_GENE_ID and "gene list" radio button
- Click "Functional Annotation Clustering"
- Select groups in "Annotation Summary Results" which have many genes from your list. Definitely select "PUBMED_ID" and interaction databases like "Biogrid"
- First look at "Functional Annotation Chart" rectangular button below to display all overrepresented terms. Sort by "Benjamini" correction for multiple hypotheses tetsting
- Select "Functional Annotation Clustering" rectangular button below to display annotation results for gene list broken into multiple groups (clusters) each with related biological functions
- Write down the # of genes in the cluster and the top functions in two most interesting clusters

# Using Group 1 options:
## 'linkage', 'average', 'RowPDistValue', 'euclidean',

**54 chart records**

| Sublist | Category | Term | RT | Genes | Count | % | P-Value | Benjamini |
|---|---|---|---|---|---|---|---|---|
| ☐ | GOTERM_CC_DIRECT | nucleus | RT | ▬▬▬▬▬▬▬ | 16 | 88.9 | 8.1E-7 | 3.7E-5 |
| ☐ | PIR_SUPERFAMILY | dual specificity protein phosphatase (MAP kinase phosphatase) | RT | ▬ | 3 | 16.7 | 4.0E-5 | 8.0E-5 |
| ☐ | GOTERM_MF_DIRECT | protein tyrosine/threonine phosphatase activity | RT | ▬ | 3 | 16.7 | 3.4E-5 | 1.3E-3 |
| ☐ | GOTERM_MF_DIRECT | MAP kinase tyrosine phosphatase activity | RT | ▬ | 3 | 16.7 | 3.4E-5 | 1.3E-3 |
| ☐ | GOTERM_MF_DIRECT | MAP kinase tyrosine/serine/threonine phosphatase activity | RT | ▬ | 3 | 16.7 | 5.9E-5 | 1.5E-3 |
| ☐ | INTERPRO | Mitogen-activated protein (MAP) kinase phosphatase | RT | ▬ | 3 | 16.7 | 3.3E-5 | 1.9E-3 |
| ☐ | SMART | RHOD | RT | ▬ | 3 | 16.7 | 2.5E-4 | 4.8E-3 |
| ☐ | INTERPRO | Rhodanese-like domain | RT | ▬ | 3 | 16.7 | 2.2E-4 | 6.2E-3 |
| ☐ | SMART | DSPc | RT | ▬ | 3 | 16.7 | 8.4E-4 | 8.0E-3 |
| ☐ | INTERPRO | Dual specificity phosphatase, catalytic domain | RT | ▬ | 3 | 16.7 | 6.0E-4 | 9.2E-3 |
| ☐ | INTERPRO | Dual specificity phosphatase, subgroup, catalytic domain | RT | ▬ | 3 | 16.7 | 6.6E-4 | 9.2E-3 |
| ☐ | GOTERM_BP_DIRECT | endoderm formation | RT | ▬ | 3 | 16.7 | 5.6E-5 | 1.1E-2 |
| ☐ | UP_KW_CELLULAR_COMPONENT | Nucleus | RT | ▬▬▬▬▬▬ | 13 | 72.2 | 1.5E-3 | 1.3E-2 |
| ☐ | SMART | PTPc_motif | RT | ▬ | 3 | 16.7 | 2.3E-3 | 1.5E-2 |
| ☐ | GOTERM_MF_DIRECT | phosphoprotein phosphatase activity | RT | ▬ | 3 | 16.7 | 8.0E-4 | 1.5E-2 |
| ☐ | INTERPRO | Protein-tyrosine phosphatase, catalytic | RT | ▬ | 3 | 16.7 | 1.4E-3 | 1.6E-2 |
| ☐ | UP_KW_PTM | Ubl conjugation | RT | ▬▬▬ | 7 | 38.9 | 4.5E-3 | 1.9E-2 |
| ☐ | UP_KW_PTM | Isopeptide bond | RT | ▬▬▬ | 6 | 33.3 | 5.4E-3 | 1.9E-2 |
| ☐ | INTERPRO | Protein-tyrosine phosphatase, active site | RT | ▬ | 3 | 16.7 | 2.1E-3 | 2.0E-2 |
| ☐ | INTERPRO | Protein-tyrosine/Dual specificity phosphatase | RT | ▬ | 3 | 16.7 | 2.8E-3 | 2.3E-2 |
| ☐ | UP_SEQ_FEATURE | DOMAIN:Rhodanese | RT | ▬ | 3 | 16.7 | 1.9E-4 | 2.4E-2 |
| ☐ | KEGG_PATHWAY | MAPK signaling pathway | RT | ▬▬ | 5 | 27.8 | 5.9E-4 | 2.8E-2 |
| ☐ | GOTERM_MF_DIRECT | myosin phosphatase activity | RT | ▬ | 3 | 16.7 | 2.4E-3 | 3.6E-2 |
| ☐ | GOTERM_MF_DIRECT | protein tyrosine phosphatase activity | RT | ▬ | 3 | 16.7 | 4.2E-3 | 5.3E-2 |
| ☐ | GOTERM_CC_DIRECT | nucleoplasm | RT | ▬▬▬▬▬ | 10 | 55.6 | 2.3E-3 | 5.4E-2 |
| ☐ | GOTERM_BP_DIRECT | negative regulation of MAPK cascade | RT | ▬ | 3 | 16.7 | 7.0E-4 | 6.8E-2 |

Gene list being analyzed

Clustering options and stringency

score for the group based on the EASE scores of each term members. The higher, the more enriched.

ALL genes involved in this annotation cluster

Every term in the annotation cluster

Genes involved in individual term

Related Term Search

A group of terms having similar biological meaning due to sharing similar gene members

**Functional Annotation Clustering**
**Current Gene List: demolist1**
**171 DAVID IDs**

⊞ **Options**   **Classification Stringency** High ▾

[ Rerun using options ]   [ Create Sublist ]            Download File

| | | | | | | |
|---|---|---|---|---|---|---|
| **Annotation Cluster 1** | | Enrichment Score: 3.69 | | | | |
| ☐ | SP_PIR_KEYWORDS | chromoprotein | RT | ▬ | 7 | 1.1E-5 |
| ☐ | SP_PIR_KEYWORDS | metalloprotein | RT | ▬ | 8 | 4.7E-5 |
| ☐ | SP_PIR_KEYWORDS | iron | RT | ▬ | 9 | 2.1E-4 |
| ☐ | GOTERM_MF_ALL | iron ion binding | RT | ▬ | 10 | 2.5E-4 |
| ☐ | SP_PIR_KEYWORDS | heme | RT | ▬ | 7 | 3.5E-4 |
| ☐ | GOTERM_MF_ALL | tetrapyrrole binding | RT | ▬ | 6 | 1.3E-3 |
| ☐ | GOTERM_MF_ALL | heme binding | RT | ▬ | 6 | 1.3E-3 |
| **Annotation Cluster 2** | | Enrichment Score: 3.52 | | | | |
| ☐ | SP_PIR_KEYWORDS | antibiotic | RT | ▬ | 5 | 2.2E-4 |
| ☐ | SP_PIR_KEYWORDS | antimicrobial | RT | ▬ | 5 | 2.4E-4 |
| ☐ | GOTERM_BP_ALL | defense response to bacteria | RT | ▬ | 6 | 5.4E-4 |
| **Annotation Cluster 3** | | Enrichment Score: 2.66 | | | | |
| ☐ | UP_SEQ_FEATURE | domain:Ig-like C2-type 1 | RT | ▬ | 8 | 5.4E-4 |
| ☐ | UP_SEQ_FEATURE | domain:Ig-like C2-type 2 | RT | ▬ | 8 | 5.4E-4 |
| ☐ | INTERPRO_NAME | Immunoglobulin | RT | ▬ | 6 | 3.6E-2 |
| **Annotation Cluster 4** | | Enrichment Score: 2.63 | | | | |

EASE Score, the modified Fisher Exact P-Value. They are identical to that in the Chart Report. The smaller, the more enriched.

# Functional Annotation Clustering

**Current Gene List: List_3**
**Current Background: Homo sapiens**
**18 DAVID IDs**

⊞ **Options**   **Classification Stringency** [Medium ∨]

[ Rerun using options ]   [ Create Sublist ]

## 25 Cluster(s)

💾 **Download File**

| Annotation Cluster 1 | | Enrichment Score: 5.2 | **G** | | Count | P_Value | Benjamini |
|---|---|---|---|---|---|---|---|
| ☐ | DISGENET | Juvenile arthritis | RT | | 7 | 1.5E-8 | 4.7E-7 |
| ☐ | DISGENET | Juvenile psoriatic arthritis | RT | | 7 | 1.5E-8 | 4.7E-7 |
| ☐ | DISGENET | Polyarthritis, Juvenile, Rheumatoid Factor Negative | RT | | 7 | 1.5E-8 | 4.7E-7 |
| ☐ | DISGENET | Polyarthritis, Juvenile, Rheumatoid Factor Positive | RT | | 7 | 1.5E-8 | 4.7E-7 |
| ☐ | DISGENET | Juvenile-Onset Still Disease | RT | | 7 | 1.8E-8 | 4.7E-7 |
| ☐ | KEGG_PATHWAY | MAPK signaling pathway | RT | | 5 | 5.9E-4 | 2.8E-2 |
| ☐ | BIOGRID_INTERACTION | mitogen-activated protein kinase 1(MAPK1) | RT | | 4 | 3.8E-3 | 1.0E0 |
| ☐ | WIKIPATHWAYS | MAPK signaling pathway | RT | | 3 | 5.8E-2 | 6.9E-1 |
| ☐ | GAD_DISEASE_CLASS | UNKNOWN | RT | | 5 | 1.5E-1 | 9.9E-1 |

| Annotation Cluster 2 | | Enrichment Score: 2.83 | **G** | | Count | P_Value | Benjamini |
|---|---|---|---|---|---|---|---|
| ☐ | INTERPRO | Mitogen-activated protein (MAP) kinase phosphatase | RT | | 3 | 3.3E-5 | 1.9E-3 |
| ☐ | GOTERM_MF_DIRECT | protein tyrosine/threonine phosphatase activity | RT | | 3 | 3.4E-5 | 1.3E-3 |
| ☐ | GOTERM_MF_DIRECT | MAP kinase tyrosine phosphatase activity | RT | | 3 | 3.4E-5 | 1.3E-3 |
| ☐ | PIR_SUPERFAMILY | dual specificity protein phosphatase (MAP kinase phosphatase) | RT | | 3 | 4.0E-5 | 8.0E-5 |
| ☐ | GOTERM_BP_DIRECT | endoderm formation | RT | | 3 | 5.6E-5 | 1.1E-2 |
| ☐ | GOTERM_MF_DIRECT | MAP kinase tyrosine/serine/threonine phosphatase activity | RT | | 3 | 5.9E-5 | 1.5E-3 |
| ☐ | PUBMED_ID | 27880917 | RT | | 4 | 1.7E-4 | 2.5E-2 |
| ☐ | UP_SEQ_FEATURE | DOMAIN:Rhodanese | RT | | 3 | 1.9E-4 | 2.4E-2 |
| ☐ | INTERPRO | Rhodanese-like domain | RT | | 3 | 2.2E-4 | 6.2E-3 |
| ☐ | SMART | RHOD | RT | | 3 | 2.5E-4 | 4.8E-3 |

| Annotation Cluster 3 | | Enrichment Score: 2.43 | G | | Count | P_Value | Benjamini |
|---|---|---|---|---|---|---|---|
| | DISGENET | Arsenic Poisoning, Inorganic | RT | | 3 | 3.5E-3 | 4.6E-2 |
| | DISGENET | Nervous System, Organic Arsenic Poisoning | RT | | 3 | 3.5E-3 | 4.6E-2 |
| | DISGENET | Arsenic Poisoning | RT | | 3 | 3.5E-3 | 4.6E-2 |
| | DISGENET | Arsenic Encephalopathy | RT | | 3 | 3.5E-3 | 4.6E-2 |
| | DISGENET | Arsenic Induced Polyneuropathy | RT | | 3 | 3.5E-3 | 4.6E-2 |
| | DISGENET | Dermatologic disorders | RT | | 3 | 5.1E-3 | 5.6E-2 |
| **Annotation Cluster 4** | | **Enrichment Score: 2.26** | G | | **Count** | **P_Value** | **Benjamini** |
| | PUBMED_ID | 19322201 | RT | | 7 | 1.3E-8 | 5.9E-6 |
| | BIOGRID_INTERACTION | ELAV like RNA binding protein 1(ELAVL1) | RT | | 7 | 4.4E-3 | 1.0E0 |
| | UCSC_TFBS | CEBPA | RT | | 7 | 1.8E-1 | 1.0E0 |
| | UCSC_TFBS | CDPCR3HD | RT | | 7 | 6.5E-1 | 1.0E0 |
| | UCSC_TFBS | FOXD3 | RT | | 5 | 7.4E-1 | 1.0E0 |
| **Annotation Cluster 5** | | **Enrichment Score: 2.14** | G | | **Count** | **P_Value** | **Benjamini** |
| | GOTERM_BP_DIRECT | negative regulation of transcription from RNA polymerase II promoter | RT | | 6 | 1.4E-3 | 9.1E-2 |
| | BIOGRID_INTERACTION | retinoid X receptor alpha(RXRA) | RT | | 3 | 6.1E-3 | 1.0E0 |
| | GOTERM_MF_DIRECT | protein heterodimerization activity | RT | | 3 | 4.5E-2 | 3.7E-1 |
| **Annotation Cluster 6** | | **Enrichment Score: 1.95** | G | | **Count** | **P_Value** | **Benjamini** |
| | REACTOME_PATHWAY | Generic Transcription Pathway | RT | | 7 | 2.8E-3 | 1.7E-1 |
| | REACTOME_PATHWAY | RNA Polymerase II Transcription | RT | | 7 | 4.6E-3 | 1.7E-1 |
| | REACTOME_PATHWAY | Gene expression (Transcription) | RT | | 7 | 8.2E-3 | 2.0E-1 |
| | GAD_DISEASE_CLASS | UNKNOWN | RT | | 5 | 1.5E-1 | 9.9E-1 |
| **Annotation Cluster 7** | | **Enrichment Score: 1.76** | G | | **Count** | **P_Value** | **Benjamini** |
| | PUBMED_ID | 18029348 | RT | | 6 | 1.8E-5 | 3.4E-3 |
| | UP_KW_PTM | Isopeptide bond | RT | | 6 | 5.4E-3 | 1.9E-2 |
| | PUBMED_ID | 15342556 | RT | | 3 | 7.9E-3 | 4.8E-1 |
| | PUBMED_ID | 26496610 | RT | | 3 | 1.0E-1 | 1.0E0 |
| | GOTERM_MF_DIRECT | metal ion binding | RT | | 4 | 4.5E-1 | 1.0E0 |
| | UCSC_TFBS | TAL1ALPHAE47 | RT | | 3 | 7.9E-1 | 1.0E0 |

| Annotation Cluster 3 | | Enrichment Score: 2.43 | G | | Count | P_Value | Benjamini |
|---|---|---|---|---|---|---|---|
| DISGENET | Arsenic Poisoning, Inorganic | RT | | | 3 | 3.5E-3 | 4.6E-2 |
| DISGENET | Nervous System, Organic Arsenic Poisoning | RT | | | 3 | 3.5E-3 | 4.6E-2 |
| DISGENET | Arsenic Poisoning | RT | | | 3 | 3.5E-3 | 4.6E-2 |
| DISGENET | Arsenic Encephalopathy | RT | | | 3 | 3.5E-3 | 4.6E-2 |
| DISGENET | Arsenic Induced Polyneuropathy | RT | | | 3 | 3.5E-3 | 4.6E-2 |
| DISGENET | Dermatologic disorders | RT | | | 3 | 5.1E-3 | 5.6E-2 |
| **Annotation Cluster 4** | | **Enrichment Score: 2.26** | G | | Count | P_Value | Benjamini |
| PUBMED_ID | 19322201 | RT | | | 7 | 1.3E-8 | 5.9E-6 |
| BIOGRID_INTERACTION | ELAV like RNA binding protein 1(ELAVL1) | RT | | | 7 | 4.4E-3 | 1.0E0 |
| UCSC_TFBS | CEBPA | RT | | | 7 | 1.8E-1 | 1.0E0 |
| UCSC_TFBS | CDPCR3HD | RT | | | 7 | 6.5E-1 | 1.0E0 |
| UCSC_TFBS | FOXD3 | RT | | | 5 | 7.4E-1 | 1.0E0 |
| **Annotation Cluster 5** | | **Enrichment Score: 2.14** | G | | Count | P_Value | Benjamini |
| GOTERM_BP_DIRECT | negative regulation of transcription from RNA polymerase II promoter | RT | | | 6 | 1.4E-3 | 9.1E-2 |
| BIOGRID_INTERACTION | retinoid X receptor alpha(RXRA) | RT | | | 3 | 6.1E-3 | 1.0E0 |
| GOTERM_MF_DIRECT | protein heterodimerization activity | RT | | | 3 | 4.5E-2 | 3.7E-1 |
| **Annotation Cluster 6** | | **Enrichment Score: 1.95** | G | | Count | P_Value | Benjamini |
| REACTOME_PATHWAY | Generic Transcription Pathway | RT | | | 7 | 2.8E-3 | 1.7E-1 |
| REACTOME_PATHWAY | RNA Polymerase II Transcription | RT | | | 7 | 4.6E-3 | 1.7E-1 |
| REACTOME_PATHWAY | Gene expression (Transcription) | RT | | | 7 | 8.2E-3 | 2.0E-1 |
| GAD_DISEASE_CLASS | UNKNOWN | RT | | | 5 | 1.5E-1 | 9.9E-1 |
| **Annotation Cluster 7** | | **Enrichment Score: 1.76** | G | | Count | P_Value | Benjamini |
| PUBMED_ID | 18029348 | RT | | | 6 | 1.8E-5 | 3.4E-3 |
| UP_KW_PTM | Isopeptide bond | RT | | | 6 | 5.4E-3 | 1.9E-2 |
| PUBMED_ID | 15342556 | RT | | | 3 | 7.9E-3 | 4.8E-1 |
| PUBMED_ID | 26496610 | RT | | | 3 | 1.0E-1 | 1.0E0 |
| GOTERM_MF_DIRECT | metal ion binding | RT | | | 4 | 4.5E-1 | 1.0E0 |
| UCSC_TFBS | TAL1ALPHAE47 | RT | | | 3 | 7.9E-1 | 1.0E0 |

QUESTIONS
FOUND IN GOOGLE AUTOCOMPLETE

Credit: XKCD comics

# Reminder from the first lecture

# Sea urchin embryonic development (from endomesoderm up to 30 hours) by Davidson's lab

May 29, 2007



Ubiq=ubiquitous; Mat = maternal; activ = activator; rep = repressor;
unkn = unknown; Nucl. = nuclearization; χ = β-catenin source;
nβ-TCF = nuclearized b-β-catenin-Tcf1; ES = early signal;
ECNS = early cytoplasmic nuclearization system; Zyg. N. = zygotic Notch.

Copyright © 2001–2007 Hamid Bolouri and Eric Davidson

# Protein-Protein binding
# IntAct Database (Dec 2015)
# Interactions: 577,297   Proteins: 89,716



Baker's yeast *S. cerevisiae* (only nuclear proteins shown)
From S. Maslov, K. Sneppen, Science 2002

Worm *C. elegans*
From S. Lee et al , Science 2004

# Metabolic pathway chart by ExPASy: 5702 reactions as of December 2015

# Basic concepts of network analysis

# Degree of a node – its # of neighbors



Degree
$K_2=2$

2

1

Degree
$K_1=4$

# Directed networks have
# in- and out- degrees

# How to find "important" nodes?

- By their degree
- Hubs = important
- Example: Google's PageRank

# How Google PageRank algorithm works?

- Google was solving the following problem in mid-1990s: too many websites match a typical search query: need to rank websites.

- Other popular search engines (e.g. Altavista) count the # of times a query word appears in website's text. Websites respond by putting lots of invisible words

- One could rank the importance of webpages by number of hyperlinks pointing to it (in-degree $K_{in}$) but:
  - Too democratic: It doesn't take into account the importance of webpages sending hyperlinks
  - it's easy to trick and artificially boost the rank

- Google's solution: simulate the behavior of many "random surfers" and then count the number of times they visited each webpage = it's PageRank
  - Popular pages send more surfers your way → the PageRank weight is proportional to $K_{in}$ but weighted by popularity

# PageRank algorithm is Google's $2.8T idea

- PageRank assigns to every webpage an importance score $G_i$

- The meaning of $G_i$ – how often random surfers visit this website

- To determine solves a self-consistent Eq.:

$G_i \sim \sum_j T_{ij} G_j.$ Here
$T_{ij} = A_{ij}/K_{out}(j)$ is the normalized adjacency matrix

- It finds the principal eigenvector (the one with the largest eigenvalue).

# Problem with PageRank algorithm and how Google solved it

- Problem: surfers can be trapped in infinite loops with one or more entrances and no exits

- Model with random jumps mimicking surfers getting bored when following a chain of links

$$G_i \sim (1-\alpha)\sum_j T_{ij} G_j + \alpha \sum_j G_j$$

- $\alpha=0.15$ meaning that an average web surfer (circa 1995) on average jumped around $1/\alpha \approx 6$ webpages before going somewhere else

# How to find "important" nodes?

- By their connectivity
- Connectors = important
- Betweenness-centrality

# Betweenness centrality: definition

- Take a node i
- There are (N-1)*(N-2)/2 pairs of other nodes
- For each pair find the shortest path on the network
- If more than one shortest path, sample them equally
- Betweenness-centrality C(i) ~ the number of shortest paths going through node i

# How is it connected to expression data analysis?

# T-cell expression data

- The matrix contains 47 expression samples from Lukk et al, Nature Biotechnology 2011
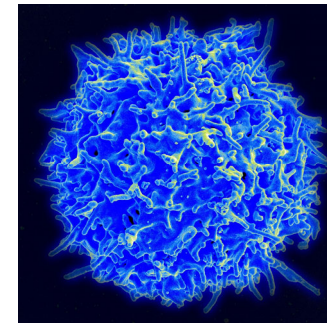
- All samples are normal T-cells from different individuals

- Only the top 3000 genes with the largest variability were used

- The value is log2 of gene's expression level in a given sample as measured by microarray technology

A global map of human gene expression

Margus Lukk, Misha Kapushesky, Janne Nikkilä, Helen Parkinson, Angela Goncalves, Wolfgang Huber, Esko Ukkonen & Alvis Brazma

Affiliations | Corresponding author

Although there is only one human genome sequence, different genes are expressed in many different cell types and tissues, as well as in different developmental stages or diseases. The structure of this 'expression space' is still largely unknown, as most transcriptomics experiments focus on sampling small regions. We have constructed a global gene expression map by integrating microarray data from 5,372 human samples representing 369 different cell and tissue types, disease states and cell lines. These have been compiled in an online resource (http://www.ebi.ac.uk/gxa/array/U133A) that allows the user to search for a gene of interest and

# Correlated pairs

**plausible biological connection based**

**on short description**

 **g1=1994; g2=188; <span style="color:red">group 1</span>**

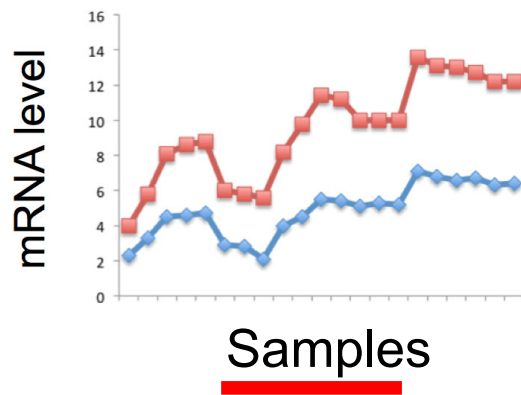 **g1=2872; g2=1269; <span style="color:red">group 2</span>**

 **g1=1321; g2=10; <span style="color:red">group 3</span>**

 **g1= 886; g2=819; <span style="color:red">group 4</span>**

 **g1=2138; g2=1364; <span style="color:red">group 5</span>**
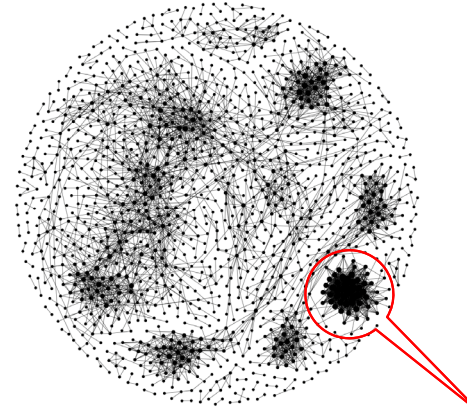
**no obvious biological common function**

**g1=1+floor(rand.*3000); g2=1+floor(rand.*3000);**
**disp([g1, g2])**

To analyze
correlations in expression
for all pairs of genes:
Co-expression networks

# How to construct a co-expression network?

A co-expression network



Functional modules

- Start with a matrix of log2 of expression levels of
  N genes in K samples (conditions): for our T-cell data N=3000, K=47

- For each of N(N-1)/2 pairs of genes i and j calculate
  the correlation coefficient $\rho_{ij}=\sigma_{ij}/\sigma_i\sigma_j$ of gene levels across K samples

- Put a threshold, e.g. $\rho_{ij}>0.85$, or otherwise select
  the most correlated pairs of genes (~4500 in our case).
  Now you have a weighted network.

- Identify densely interconnected functional modules in
  this network.

- Modules can be used to infer unknown functions of genes via
  "Guilt by Association" principle.

# How to install Gephi software for network analysis?

- Install Gephi from: https://gephi.org/users/download/
- One of the common problems with installation is the version of Java on your computer. One possible solution is here: https://github.com/gephi/gephi/issues/1787.

Sometimes after installation Gephi may complain that it cannot find java version 1.8 or higher. In this case you need to go to  C:\Program Files\Gephi-0.9.2\etc
Open file gephi.conf using notepad.exe (MS Word does not work!).
Add a line jdkhome="C:\Program Files (x86)\Java\jre1.8.0_231"
(the numbers in …jre1.8.0_231  may be changed to reflect the actual directory where Java is installed on your computer).
If JDK is not installed on your computer, you need to install itfirst from https://www.java.com/en/download/win10.jsp"

# Co-expression network analysis exercise

- Start Gephi and open
  <span style="color:red">coexpression_network_random_start.gephi</span>

- Run "Layout" → Fruchterman Reingold → Speed 10.0

- <u>Run "Average degree", "Network diameter", "Modularity"</u> in the Statistics tab in the right panel.

- <u>Color nodes by "modularity class"</u>:
  Appearance → Nodes → Partition → Palette Icon → Modularity class

- <u>Size nodes first by "degree".</u>
  Appearance → Nodes → Ranking → Multiple Circles Icon → Degree
  - If the nodes are too small, select "Min size": 10 and "Max size":80
  - Nodes in large tightly connected clusters have large degree

- <u>Then size nodes by "betweenness-centrality"</u>
  Appearance → Nodes → Ranking → Multiple Circles Icon → Betweenness-centrality
  - Large circles are "coordinator" genes connecting different co-expressed clusters to each other. Potentially biologically interesting

# Disease-disease similarity network

- Based on the table summarizing **<u>all current medical knowledge</u>** of genes implicated in diseases:
  - Rows: 516 common human diseases
  - Columns: 25,000 human genes
  - Matrix element $D_{i\alpha}$ =1 if the gene $\alpha$ is known to be involved in the disease i . 0 – otherwise
- Constructed disease-disease similarity network:
  - Weight of the edge - # of shared genes between two diseases
  - Easy to construct: the adjacency matrix A of the network is simply $A = D \cdot D^+$

# Disease network analysis exercise

- Start Gephi and open disease_disease_random_start.gephi

- Run "Layout" → Fruchterman Reingold → Speed 10.0
  Observe how clusters emerge.

- Run "Average degree", "Network diameter", "Modularity" analysis tools in the right panel.

- Color nodes with medical term:  "disorder class"
  Appearance → Nodes → Partition → Palette Icon → Disorder class

- Then color nodes by "modularity class". See how well it agrees with the previous color.
  Appearance → Nodes → Partition → Palette Icon → Modularity class

- Size nodes first by "degree".
  Appearance → Nodes → Ranking → Multiple Circles Icon → Degree
  - Which disease has the largest degree?

-  Size nodes by "betweenness centrality"
  Appearance → Nodes → Ranking → Multiple Circles Icon → Degree
  - Which diseases have the largest betweenness-centrality?

  These "connector" diseases linking different diseases clusters to each other. They highlight potentially interesting connections between diseases

QUESTIONS
FOUND IN GOOGLE AUTOCOMPLETE