

Functional Profiling

Rubén Grillo Risco
Bioinformatics & Biostatistics Unit. CIPF



WODA

WEB-BASED OMICS DATA ANALYSIS



Unidad de
Bioinformática y
Bioestadística

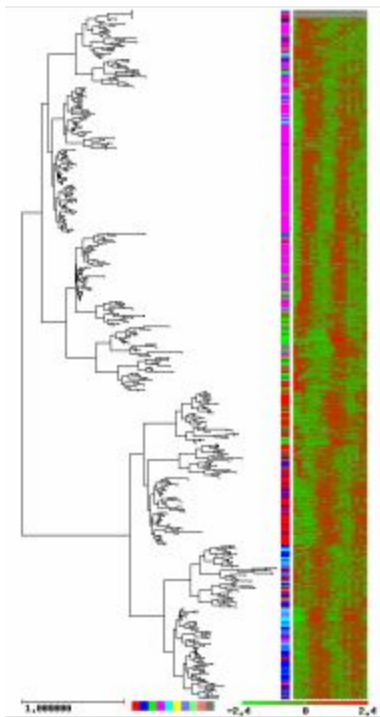


PRINCIPE FELIPE
CENTRO DE INVESTIGACION

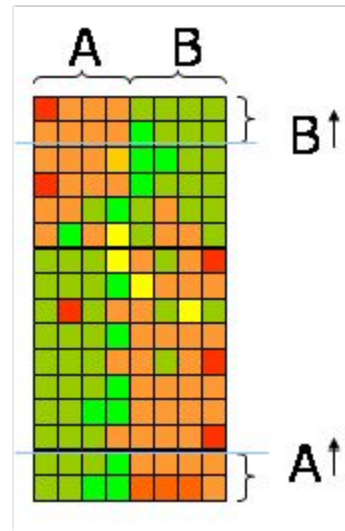
Outline

- Introduction
- Over-Representation Analysis (ORA)
- Gene Set Analysis (GSA)
- Network Analysis (NA)

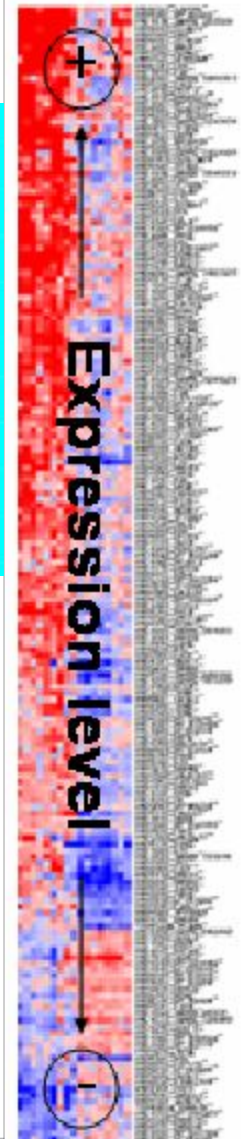
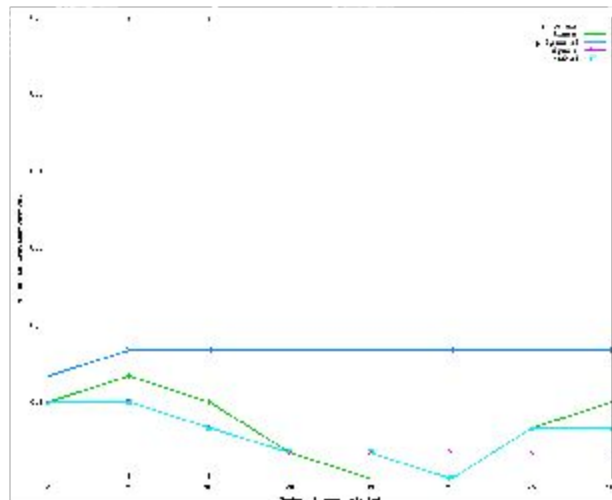
Genome-scale experiment output



BEST1
BRCA2
FIT
BRCA1



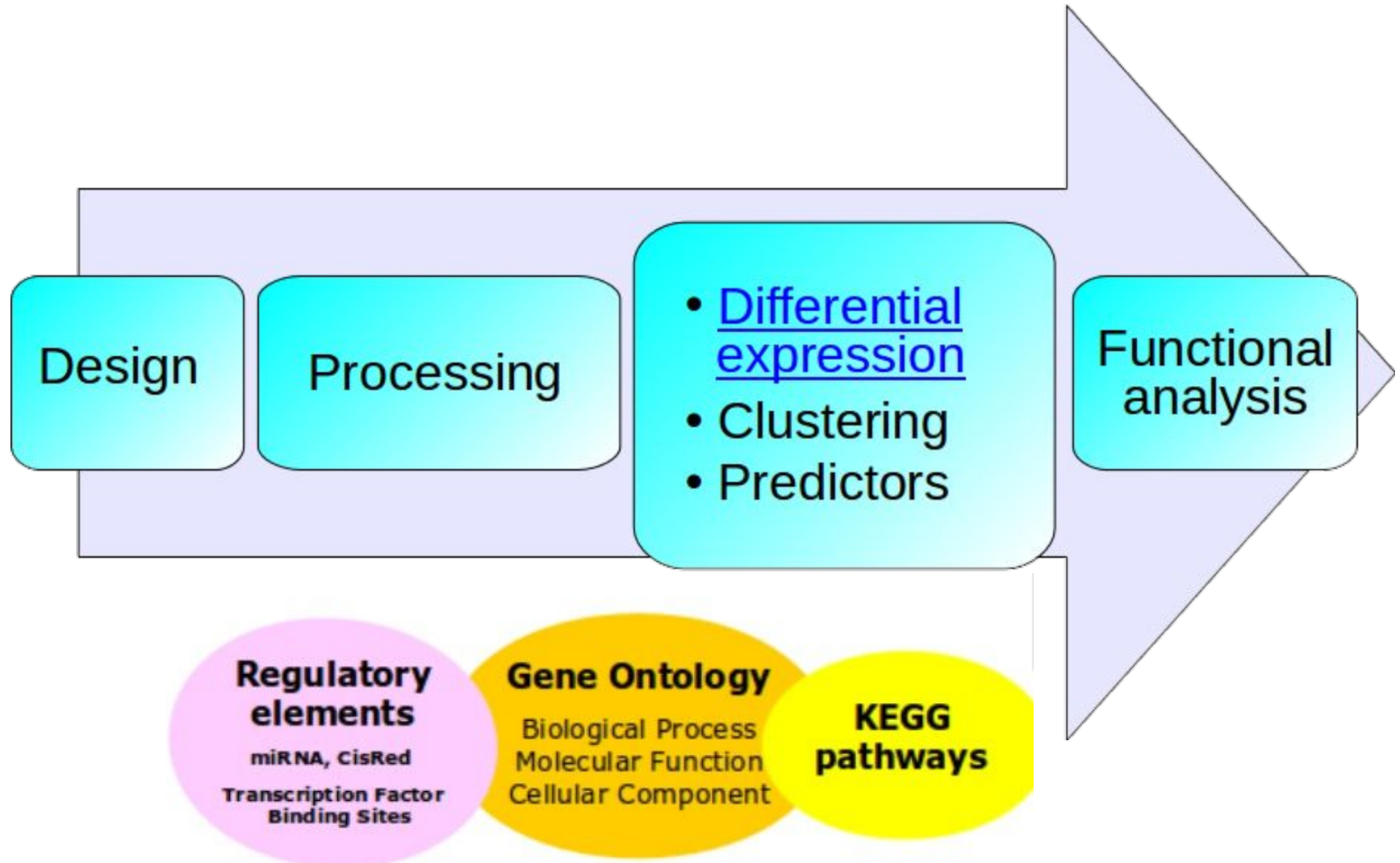
1007_s_at	12.4
1053_at	11.5
117_at	10.3
121_at	10.2
1255_g_at	9.9
1294_at	9.3
1316_at	8.2
1320_at	8.1
1405_i_at	7.7
1431_at	7.4



Questions we try to answer

- Is there any significant functional enrichment in my gene list / gene sets?
- Are these genes involved in common pathways?
- Do they share specific regulation?
- Are they involved in the same disease?

Omic data analysis pipeline



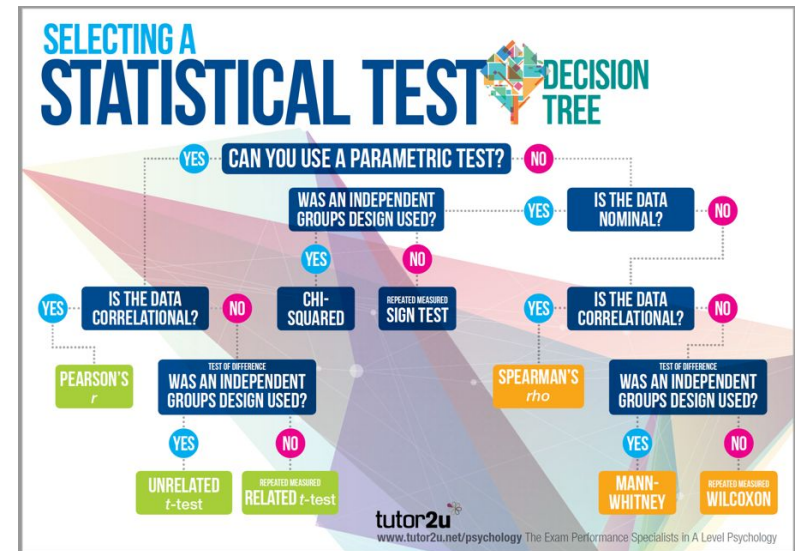
What do we need?

Annotation file + Statistical tests

gene → function

FunctionA: Gene1, Gene2, Gene3...

FunctionB: Gene4, Gene2, Gene5...



Functional databases



Homo sapiens



Mus musculus



Rattus norvegicus



Gallus gallus



Danio rerio



Drosophila melanogaster



C. elegans



Saccharomyces cerevisiae



Arabidopsis thaliana

UniProt/Swiss-Prot

UniProtKB/TrEMBL

Ensembl IDs

EntrezGene

Affymetrix

Agilent



Genes IDs

HGNC symbol

EMBL acc

RefSeq

PDB

Protein Id

IPI....

Biological databases

KEGG pathways

Biocarta pathways

Keywords Swissprot

Gene Ontology

Biological Process
Molecular Function Cellular Component

Gene Expression in tissues

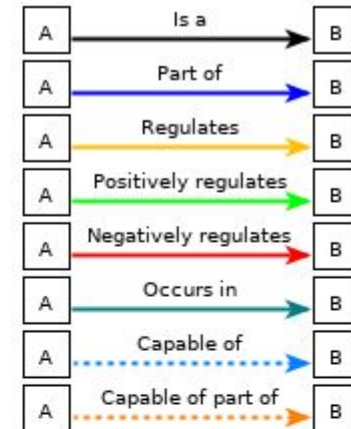
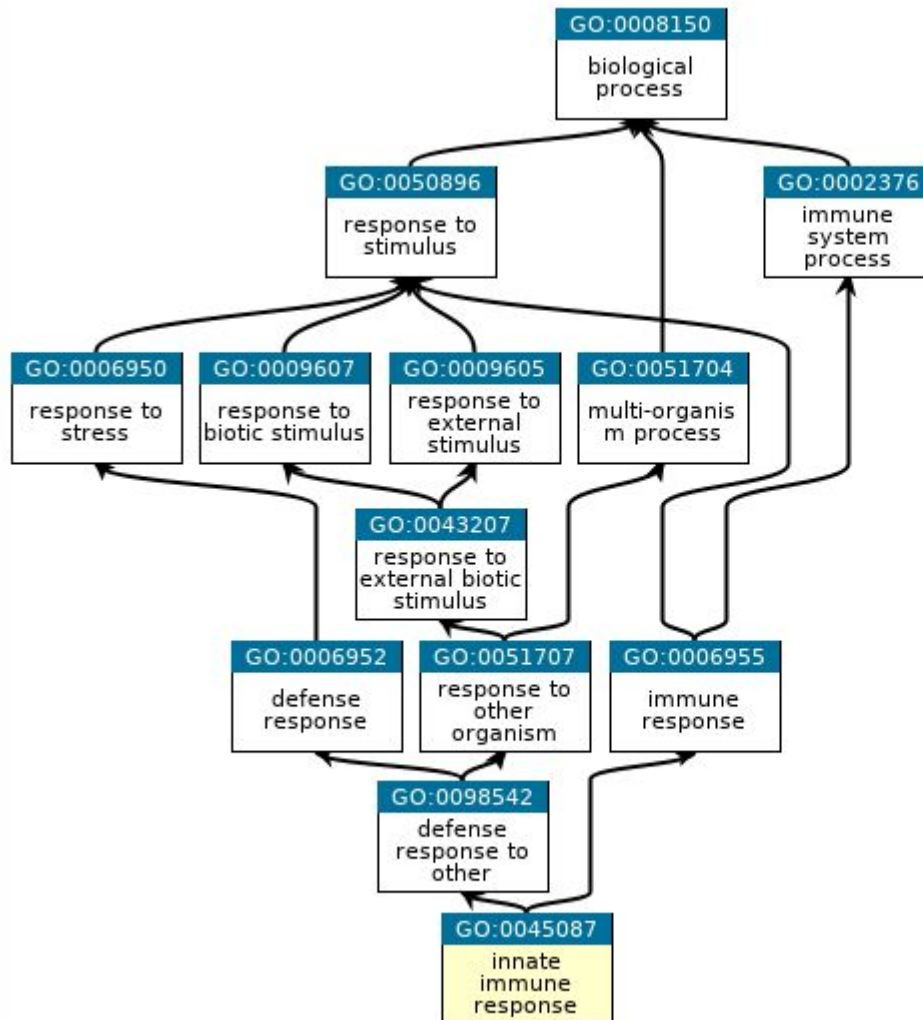
Regulatory elements

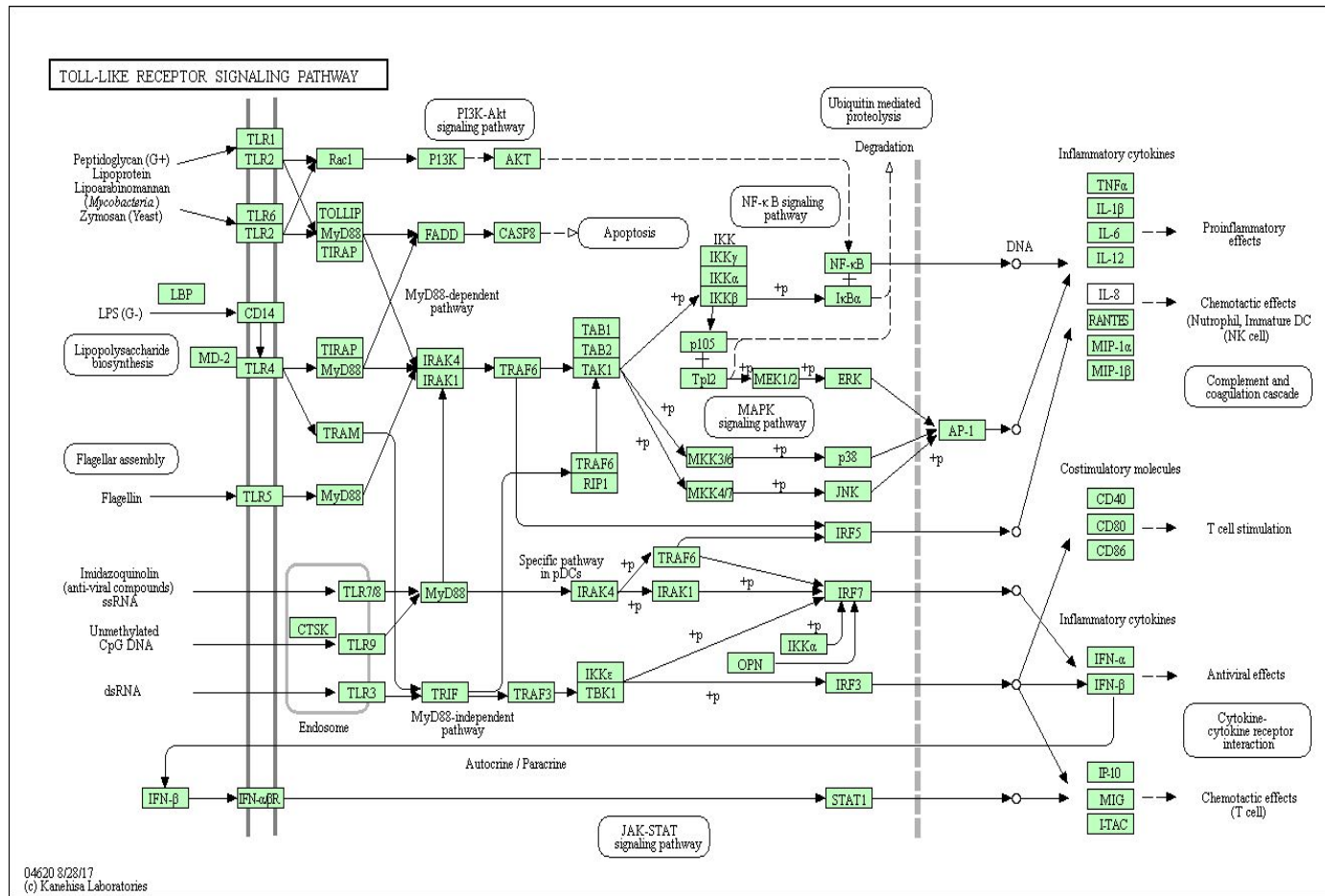
MiRNA, CisRed
Transcription Factor Binding Sites

Bioentities from literature:

**Diseases terms
Chemical terms**

Gene Ontology





Reactome

reactome 3.6 Pathways for: Homo sapiens

Analysis: Tour: Layout:

Event Hierarchy:

- Cell Cycle
- Cell-Cell communication
- Cellular responses to external stimuli
- Chromatin organization
- Circadian Clock
- Developmental Biology
- Digestion and absorption
- Disease
- DNA Repair
- DNA Replication
- Extracellular matrix organization
- Gene expression (Transcription)
- Hemostasis
 - Platelet homeostasis
 - Prostacyclin signalling through prostacyclin receptor**
 - Nitric oxide stimulates guanylate cyclase
 - Binding of ATP to P2X receptors
 - Platelet calcium homeostasis
 - Platelet sensitization by platelet calcium homeostasis
 - PAFAH2 hydrolyses PAF to lyso-PAF
 - Platelet Adhesion to exposed collagen
 - Platelet activation, signalling and aggregation
 - Formation of Fibrin Clot (Clotting Cascade)
 - Dissolution of Fibrin Clot
 - Cell surface interactions at the vascular wall
 - Factors involved in megakaryocyte development
 - Immune System
 - Metabolism
 - Metabolism of proteins
 - Metabolism of RNA
 - Mitophagy
 - Muscle contraction
 - Neuronal System

Search for a term, e.g. pten ...

Description: Prostacyclin signalling through prostacyclin receptor Id: R-HSA-392851.2 Species: Homo sapiens

Summation

Prostacyclin (PGI₂) is continuously produced by healthy vascular endothelial cells. It inhibits platelet activation through interaction with the G_s-coupled receptor PTGIR, leading to increased cAMP, a consequent increase in cAMP-dependent protein kinase activity which prevents increases of cytoplasmic [Ca²⁺]_i necessary for activation (Wouffe et al. 2001). PGI₂ is also an effective vasodilator. These effects oppose the effects of thromboxane (TXA₂), another eicosanoid, creating a balance of blood circulation and platelet activation.

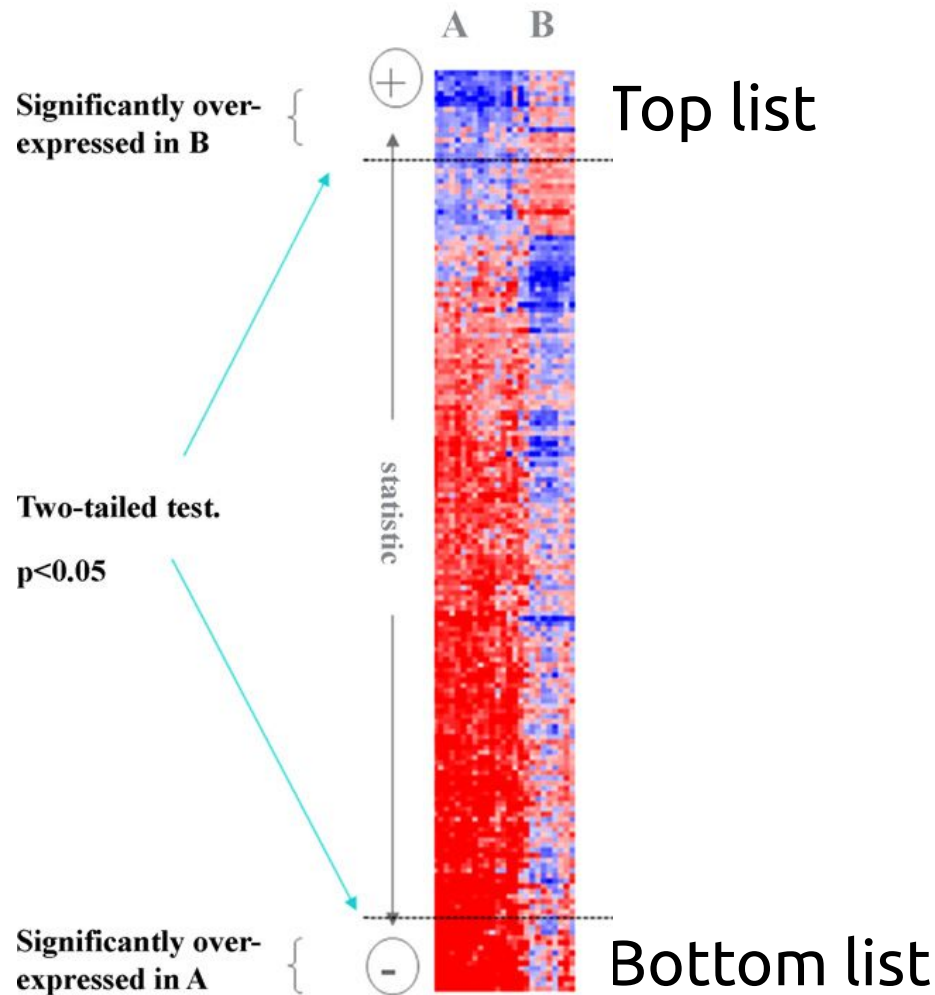
Background literature references...

External identifiers

Outline

- Introduction
- Over-Representation Analysis (ORA)
- Gene Set Analysis (GSA)
- Network Analysis (NA)

Over-Representation Analysis



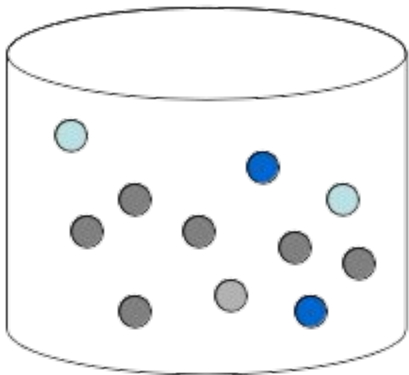
1. List is created using a certain threshold or criteria.
2. For each pathway, input genes that are part of the pathway are counted.
3. This process is repeated for an appropriate background list of genes
4. Every pathway is tested for over- or under-representation in the list of input genes.

Over-Representation Analysis

FatiGO test

Finding significant associations of **Gene Ontology**

One Gene List (A)



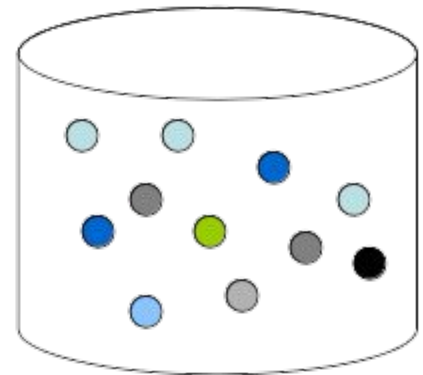
Biosynthesis 60% ●

Sporulation 20% ●

Are this two
groups of genes
carrying out
different
biological roles?



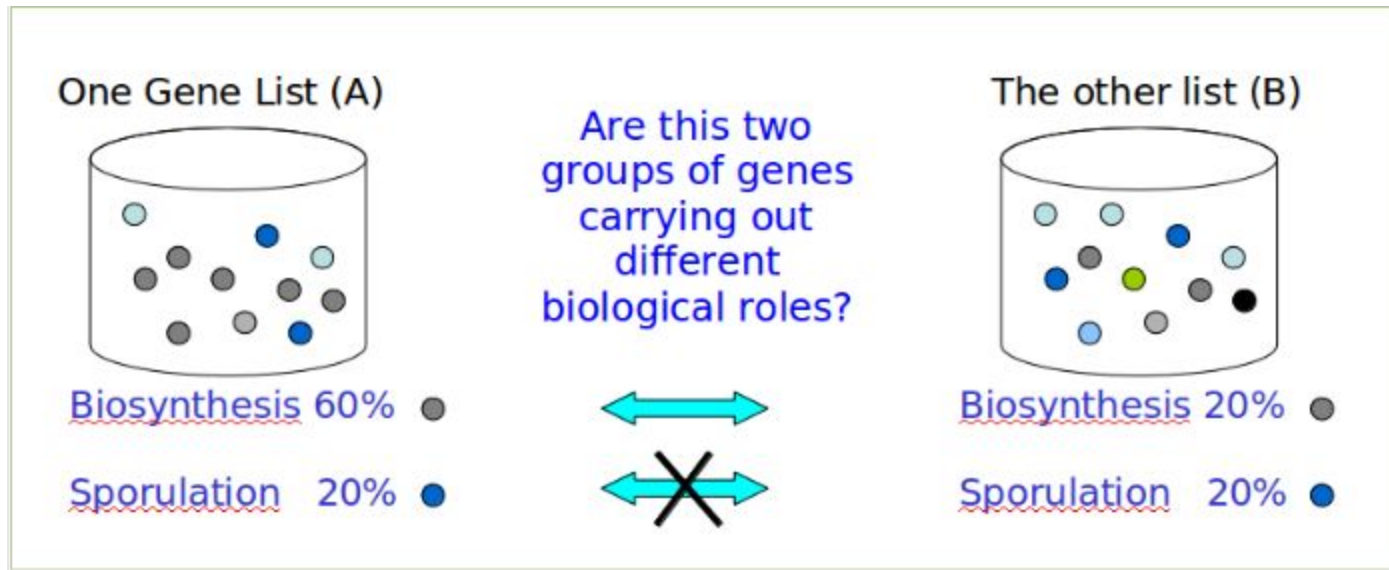
The other list (B)



Biosynthesis 20% ●

Sporulation 20% ●

Over-Representation Analysis



Genes in group A have significantly to do with biosynthesis, but not with sporulation.

	A	B
Biosynthesis	6	2
No biosynthesis	4	8

**We do this for each term (GO, miRNA, Interpro, ...)
Thousand of terms, so Multiple Test Correction is needed!!!**

Web tools for ORA

- Babelomics (FatiGO): <http://babelomics.bioinfo.cipf.es/>
- Panther: <http://www.pantherdb.org/>
- DAVID: <https://david.ncifcrf.gov/>
- Reactome: <https://reactome.org/>
- g:Profiler: <https://biit.cs.ut.ee/gprofiler/gost>
- FuncAssociate: <http://llama.mshri.on.ca/funcassociate/>
- WebGestalt: <http://www.webgestalt.org/>

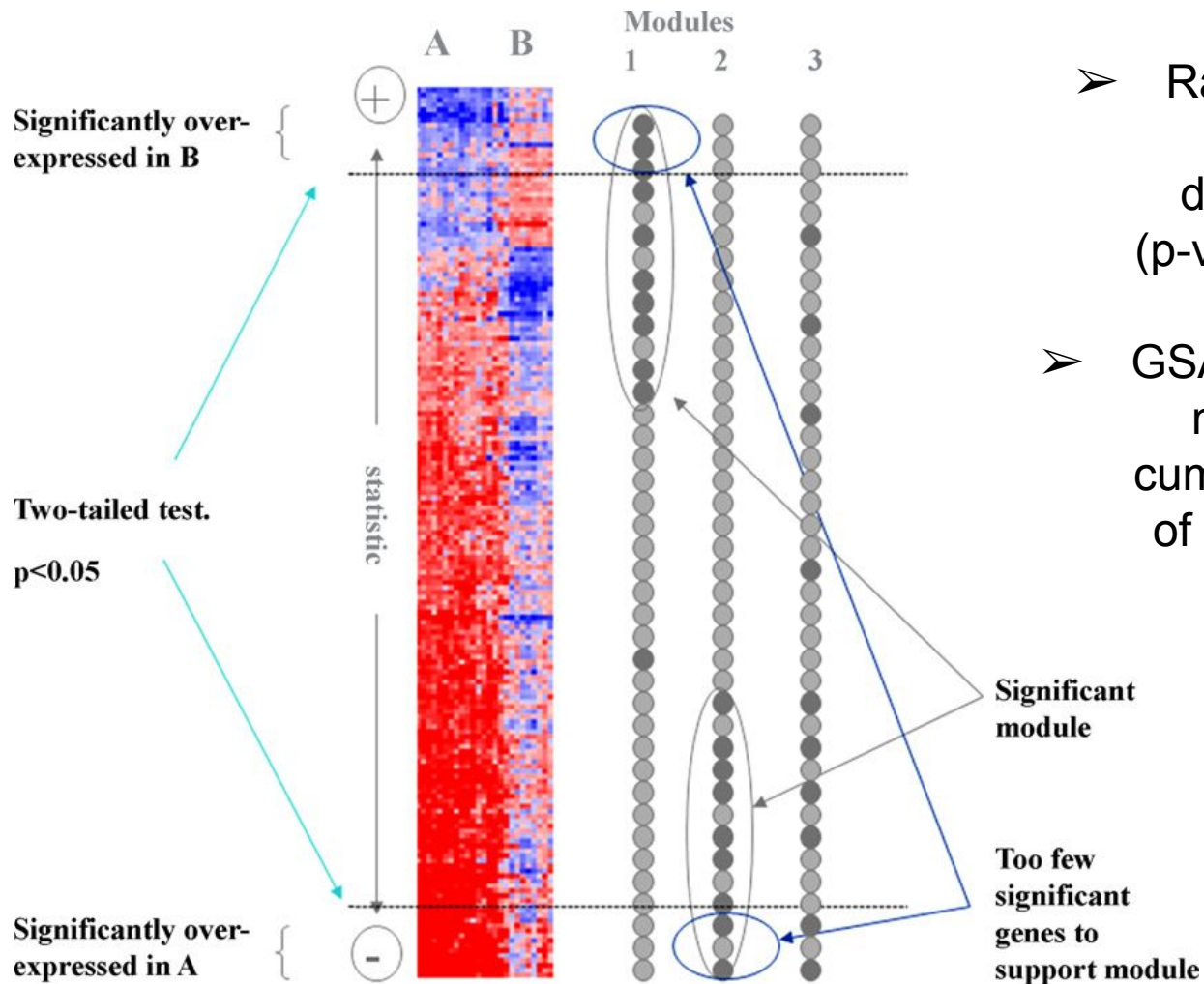
Outline

- Introduction
- Over-Representation Analysis (ORA)
- **Gene Set Analysis (GSA)**
- Network Analysis (NA)

Web tools for GSA

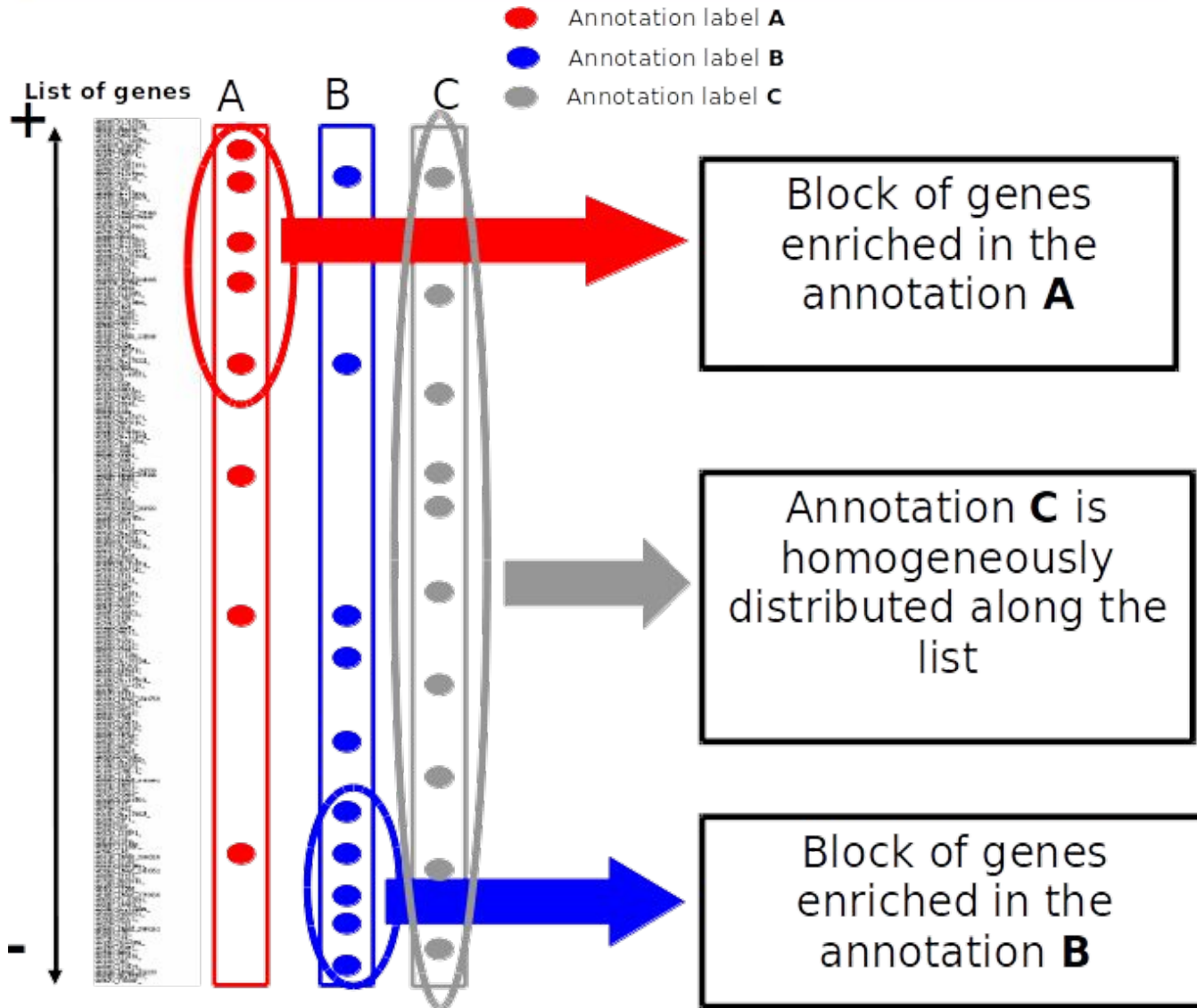
- GSEA (Gene Set Enrichment Analysis): <http://software.broadinstitute.org/gsea/index.jsp>
- Babelomics (logistic model): <http://babelomics.bioinfo.cipf.es/>
- Panther: <http://www.pantherdb.org/>
- FuncAssociate: <http://llama.mshri.on.ca/funcassociate/>
- WebGestalt: <http://www.webgestalt.org/>

Gene Set Analysis



- Ranked list of all genes according to their differential expression (p-value, t-statistic, logFC)
- GSA directly tests for gene modules significantly cumulated in the extremes of a ranked list of genes.

Gene Set Analysis



Any question?



Activities

1. Over-representation and GSA exercises:
<http://bioinfo.cipf.es/WODA19/doku.php/bbdd>
2. Protein-protein interaction exercises:
http://bioinfo.cipf.es/WODA19/doku.php/ex_ppi