

Hacia un cambio de paradigma en el análisis de datos genómicos

Joaquín Dopazo

**Department of Bioinformatics and Genomics,
Centro de Investigación Príncipe Felipe (CIPF),
Functional Genomics Node, (INB), and
Bioinformatics Group (CIBERER)
Valencia, Spain.**

<http://www.gepas.org>
<http://www.babelomics.org>
<http://bioinfo.cipf.es>



INSTITUTO NACIONAL
DE BIOINFORMÁTICA



Background

**The road of excess leads to
the palace of wisdom**

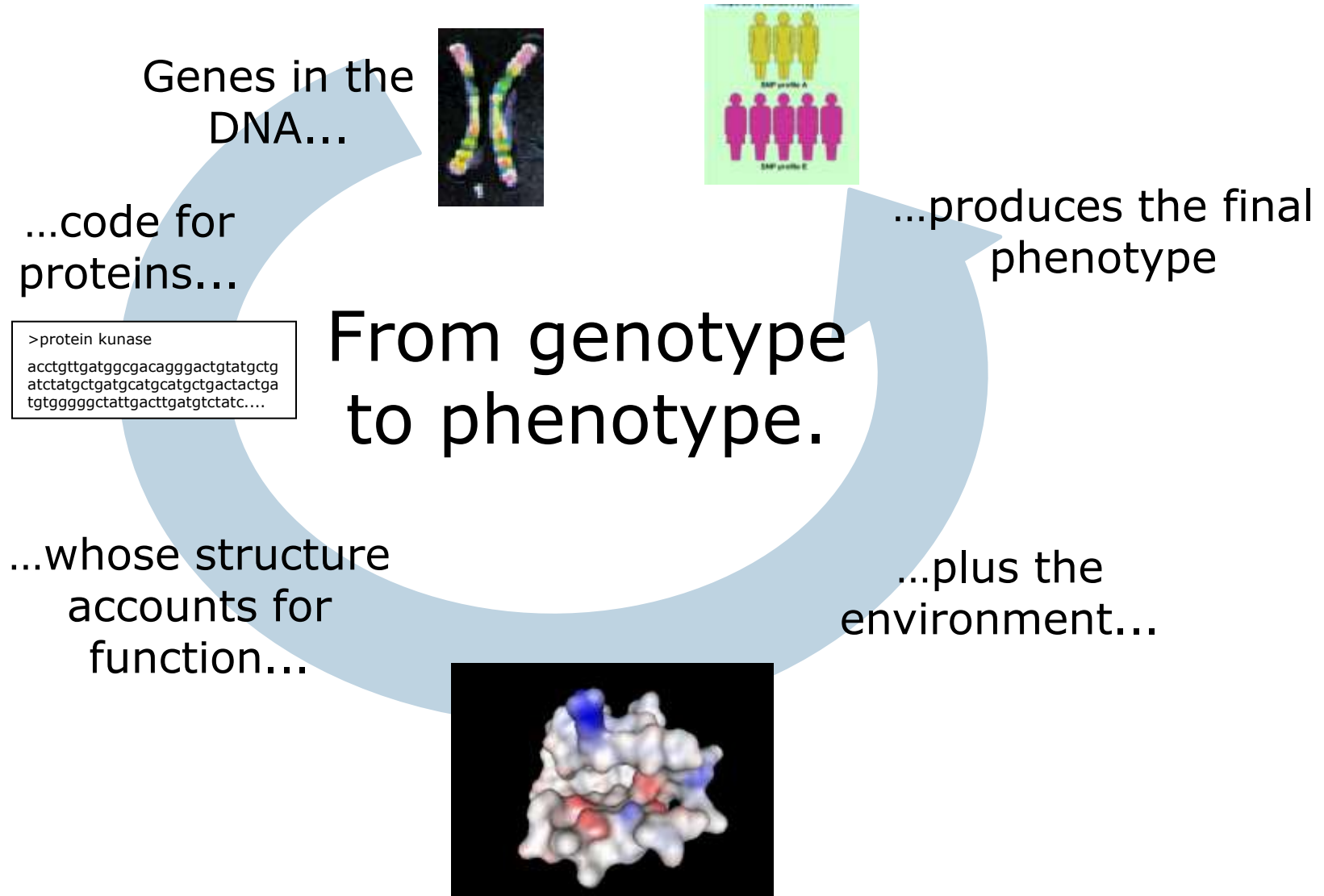
*(William Blake, 28 November 1757 – 12
August 1827, poet, painter, and printmaker)*



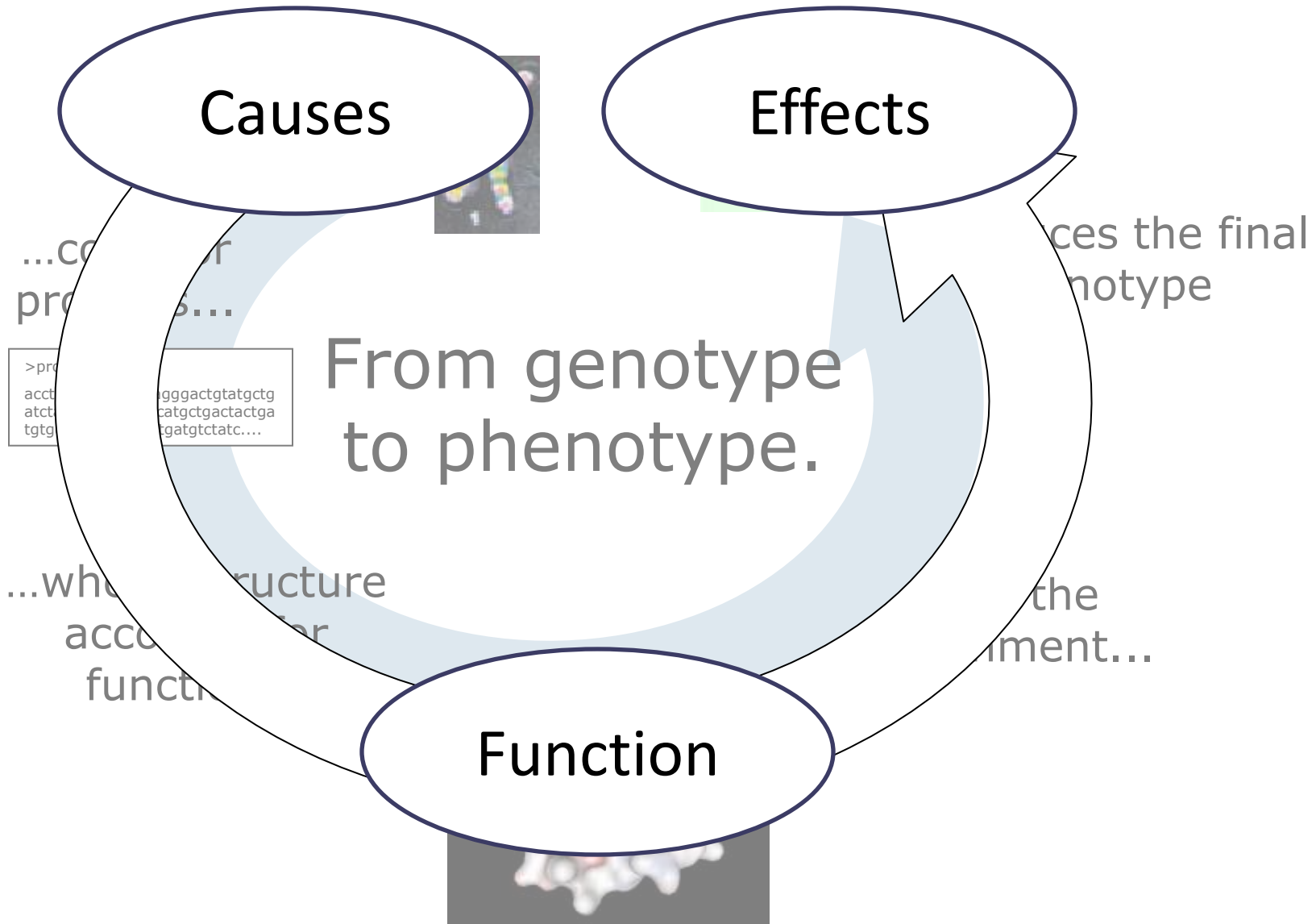
The introduction and popularisation of high-throughput techniques has drastically changed the way in which biological problems **can** be addressed and hypotheses **can** be tested.

But not necessarily the way in which we really address or test them...

Where do we come from? The pre-genomics paradigm



Reduccionistic approach to link causes (genes) to effects (phenotype) through actions (function)



Next Generation Sequencing
10⁹bp per round

Genes in the DNA.

...which can be different because of the variability.

15 million SNPs

>protein kinase

```
acctgttgatggcgacagggactgtatgctgatct
atgctgatgcatgcatgctgactactgatgtgggg
gctattgacttgatgtctatc....
```

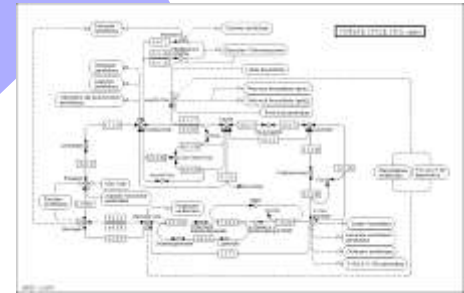


...whose final effect configures the phenotype...

...when expressed in the proper moment and place...

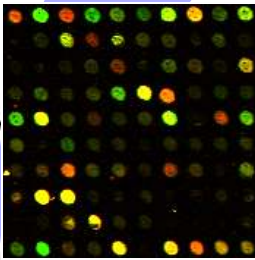
From genotype to phenotype.

(in the functional post-genomics scenario)



...conforming complex interaction networks...

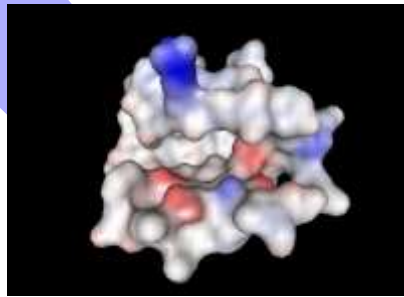
A typical tissue is expressing among 5000 and 10000 genes



...code for proteins...

That undergo post-translational modifications, somatic recombination...

100K-500K proteins



...that account for function if...

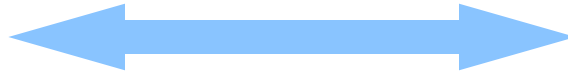
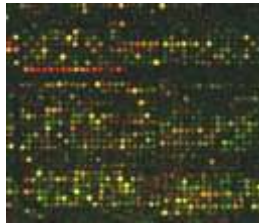


...in cooperation with other proteins...

Each protein has an average of 8 interactions

Functional genomics

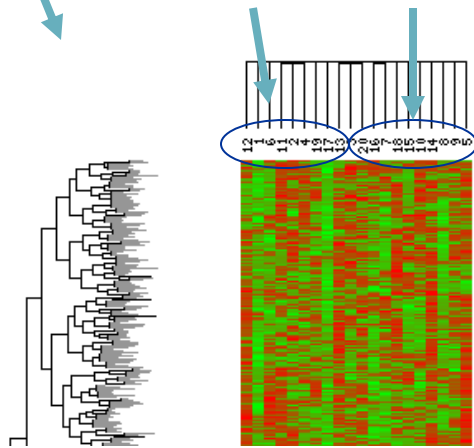
Differences at phenotype level are the visible cause of differences at molecular level which, in many cases, can be detected by measuring genomic mutations or the levels of gene expression, etc. The same holds for different experiments, treatments, strains, etc.



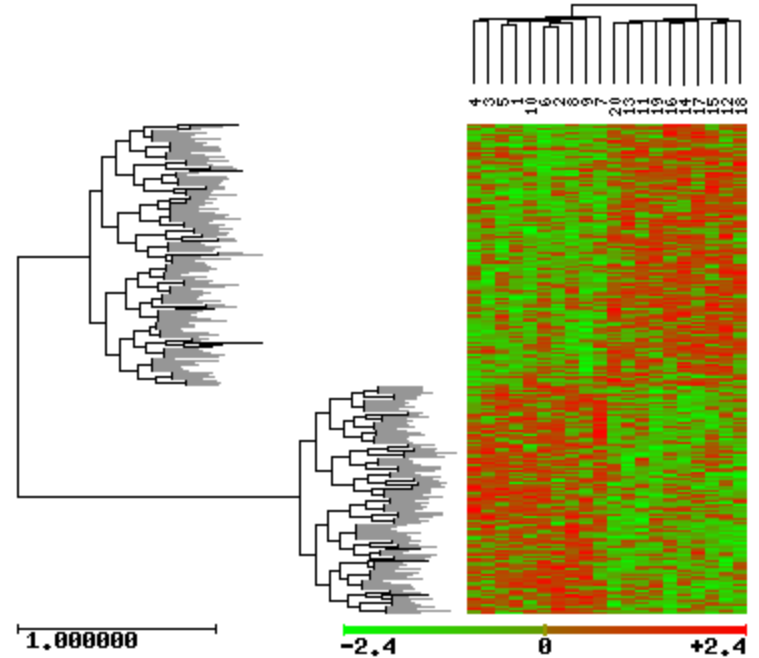
Lets have a closer look to the way we work in functional genomics with high-throughput data.

A closer look to a simple problem. Finding signatures , which implies gene selection for class discrimination

~25,000 genes
Case(10)/control(10)



thebest - [04/10/2003 18:57:43 GMT]



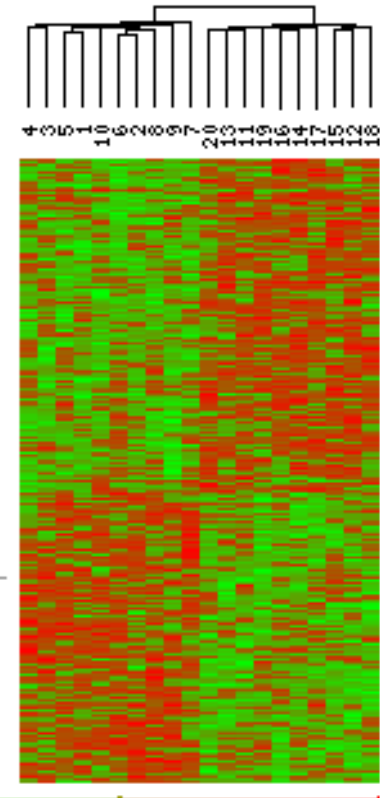
Genes differentially expressed
among classes (t-test), with a
p-value < 0.05

Signature

Sorry... the data were a collection of random numbers labelled for two classes

This is a multiple-testing statistic contrast.

GMT1



Row	ID	unadj_p	adj_p	FDR_indep	FDR_dep	obs_stat
1630	1630	0.00019998	0.152685	0.49995	1	5.47044
4148	4148	0.00019998	0.746225	0.49995	1	4.49902
3348	3348	0.0009999	0.983002	0.861025	1	4.01726
3449	3449	0.00149985	0.986401	0.861025	1	3.99374
3443	3443	0.00129987	0.9959	0.861025	1	3.86046
4703	4703	0.00169983	0.9996	0.861025	1	3.7251
2731	2731	0.00169983	0.9996	0.861025	1	3.66628
14	14					62427
39	39					60596
1062	1062					58109
3738	3738					52935
1840	1840					43721
1007	1007					41937
1542	1542					41428
1360	1360					.4025
844	844					40212
4631	4631					37412
11	11	0.00539946	1	0.8888	1	3.36813
4102	4102	0.00219978	1	0.861025	1	3.35909
4195	4195	0.0019999	1	0.861025	1	3.35235
4716	4716	0.00439956	1	0.8888	1	3.28286
4420	4420	0.00619933	1	0.8888	1	3.2427
4198	4198	0.00539946	1	0.8888	1	3.23225
3793	3793	0.00279972	1	0.861025	1	3.22175
3916	3916	0.00279972	1	0.861025	1	3.19595
372	372	0.00359958	1	0.8888	1	3.19547
3488	3488	0.0069993	1	0.8888	1	3.12957
3992	3992	0.00819945	1	0.8888	1	3.0987
1248	1248	0.00779922	1	0.8888	1	3.09834

You were not interested *a priori* in the first (whatever), best discriminant, gene.

Adjusted p-values must be used!

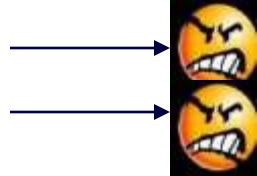
Some random sequences of values might look like if they were differentially distributed among two classes

The rationale for multiple testing



= 10 heads. $P=0.5^{10}=0.00098$

Take one coin, flip it 10 times. Got 10 heads? Use it for betting



→ 10 heads !!!

⋮



1000 coins

$$P = 1 - (1 - 0.5^{10})^{1000} = 0.62$$

It is not the same getting 10 heads with **my** coin than getting 10 heads in **one among** 1000 coins

Only a perturbing comment: do genes behave like coins?

The curse of dimensionality

The more we see the less we can be confirm

As we historically have moved from gene-targeted studies to GWAS or microarray-based gene expression studies we have gained a 1000x or even more in resolution but we have lost a lot in testing power.

In other words: our gene signatures, associated SNPs, etc., constitute an under-representation of the biologically relevant genes

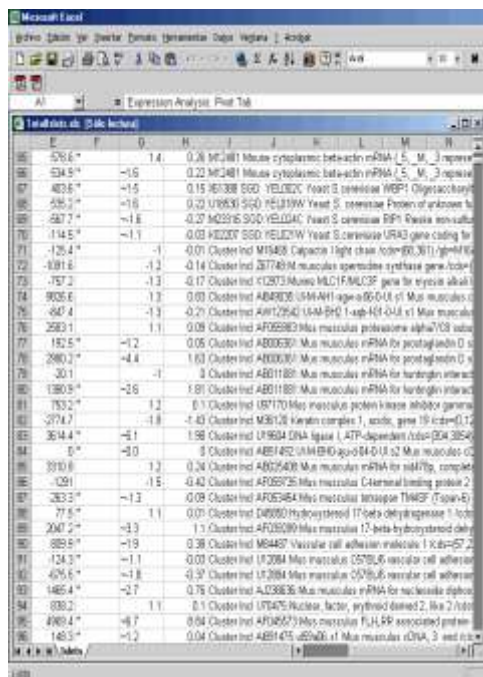
So far so good... we know the genes but, what are they doing in the cell?

The data...

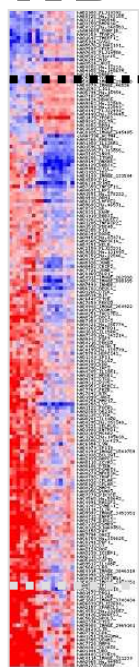
...how are structured?

What are these groups?

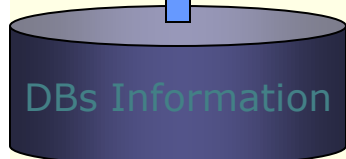
What is this gen?



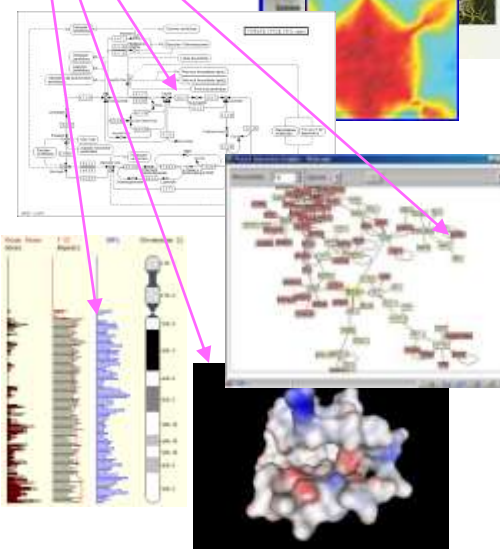
A B



Cell cycle...



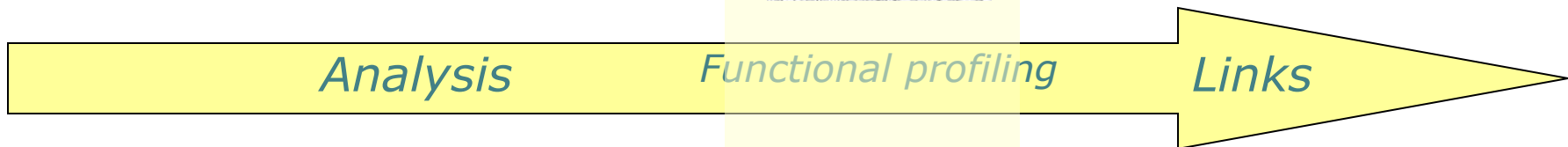
U12084 Mus musculus C57BL/6 vascular cell adhesion...
 AF02530 Mus musculus 17-beta dehydrogenase 2...
 M64437 Vascular cell adhesion molecule-1 (VcAM-1)
 U12084 Mus musculus C57BL/6 vascular cell adhesion...
 AJ33636 Mus musculus mRNA for nucleolar dipole...
 U7475 Nucleolar factor, synthetase domain 2, like 2 (Nof...
 AF04603 Mus musculus FHL-RR associated protein...
 AF1475 45S ribosomal L1 Mus musculus rDNA, 3' end (rD...



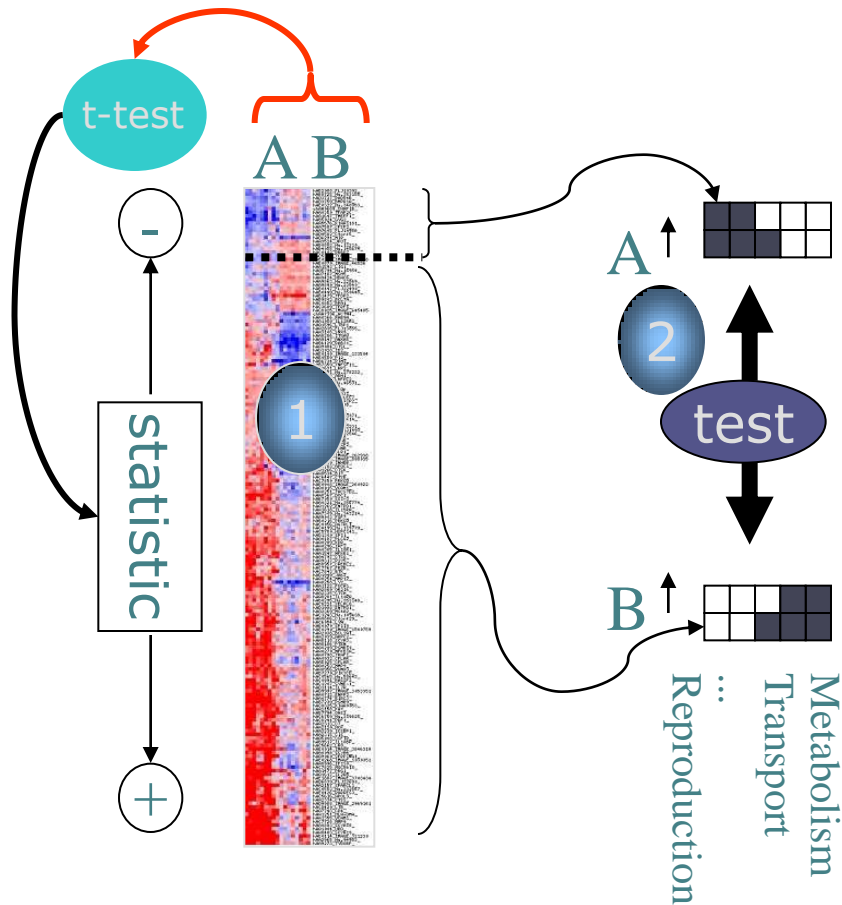
Analysis

Functional profiling

Links



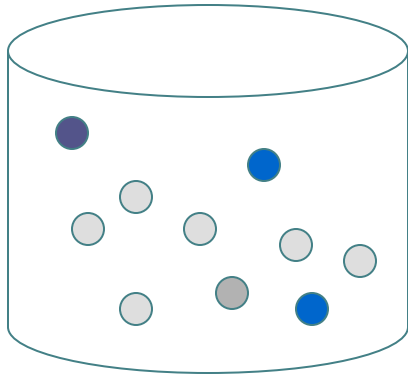
Testing for functional enrichment



- 1 Genes are selected based on their experimental values and...
- 2 Enrichment in functional terms is tested (FatiGO, GoMiner, etc.)

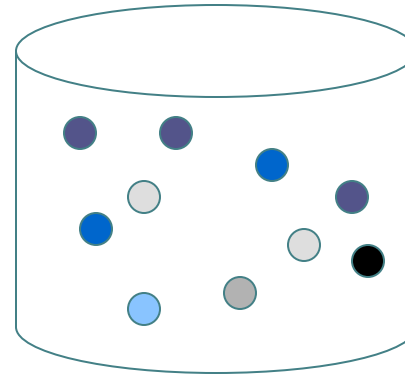
Testing for functional enrichment

Selected (A)



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

Background (B)



		Biosynthesis	
		Other	
6	4		A
2	8		B

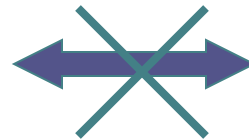
The popular Fisher's test

Biosynthesis 60% ○



Biosynthesis 20% ○

Metabolism 20% ●



Metabolism 20% ●

Genes in group A have significantly associated to biosynthesis, but not to metabolism.

- Aedes aegypti*
home page | site map
- Anopheles gambiae*
home page | site map
- Bos taurus*
home page | site map
- Caenorhabditis elegans*
home page | site map
- Canis familiaris*
home page | site map
- Cavia porcellus*
home page | site map
- Ciona intestinalis*
home page | site map
- Ciona savignyi*
home page | site map
- Danio rerio*
home page | site map
- Dasyatis novemcinctus*
home page | site map
- Drosophila melanogaster*
home page | site map
- Microcebus murinus*
home page | site map
- Monodelphis domestica*
home page | site map
- Mus musculus*
home page | site map
- Myotis lucifugus*
home page | site map
- Ochotona princeps*
home page | site map
- Ornithorhynchus anatinus*
home page | site map
- Oryzotilus cuniculus*
home page | site map
- Oryzias latipes*
home page | site map
- Otolemur garnettii*
home page | site map
- Pan troglodytes*
home page | site map
- Pongo pygmaeus*
home page | site map
- Echinops teliaeri*
home page | site map
- Equus caballus*
home page | site map
- Erinaceus europaeus*
home page | site map
- Felis catus*
home page | site map
- Gallus gallus*
home page | site map
- Gasterosteus aculeatus*
home page | site map
- Homo sapiens*
home page | site map
- Leopoldia africana*
home page | site map
- Macaca mulatta*
home page | site map

Genome Annotation

Structural Annotation

Functional Annotation

Biological Databases

Gene Annotation

Gene Set Annotation

Protein-Protein interactions

Protein Structure

KEGG pathways

Keywords Swissprot

Biocarta pathways

Motifs Domains

Bioentities from literature

Gene Ontology
Biological Process
Molecular Function
Cellular Component

Gene Expression Modules

Reactome

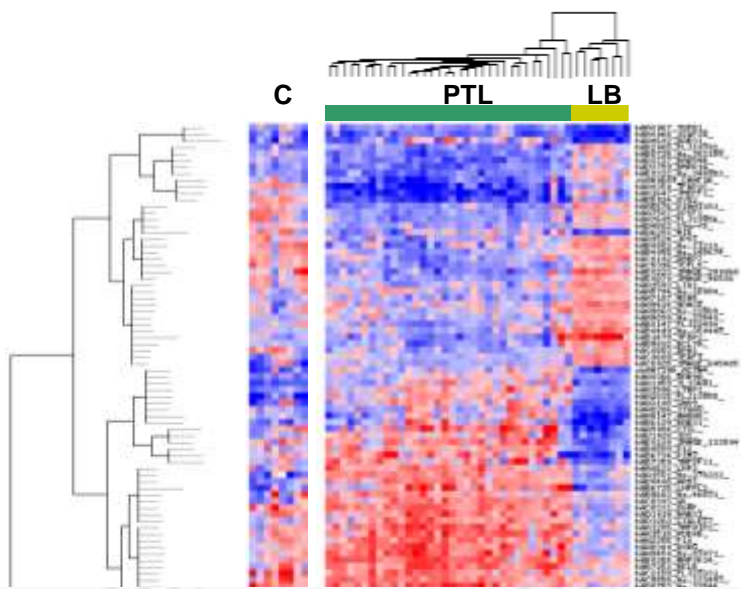
Regulatory elements
miRNA
CisRed
Transcription Factor Binding Sites

mSigDB



Understanding why genes differ in their expression between two different conditions

Lymphomas from mature lymphocytes



Gene Ontology Term

response to external stimulus

response to stress

signal transduction

cell motility

resistance to pathogenic bacteria

viral replication

cell death

regulation of gene expression, epigenetic



p-values(*)

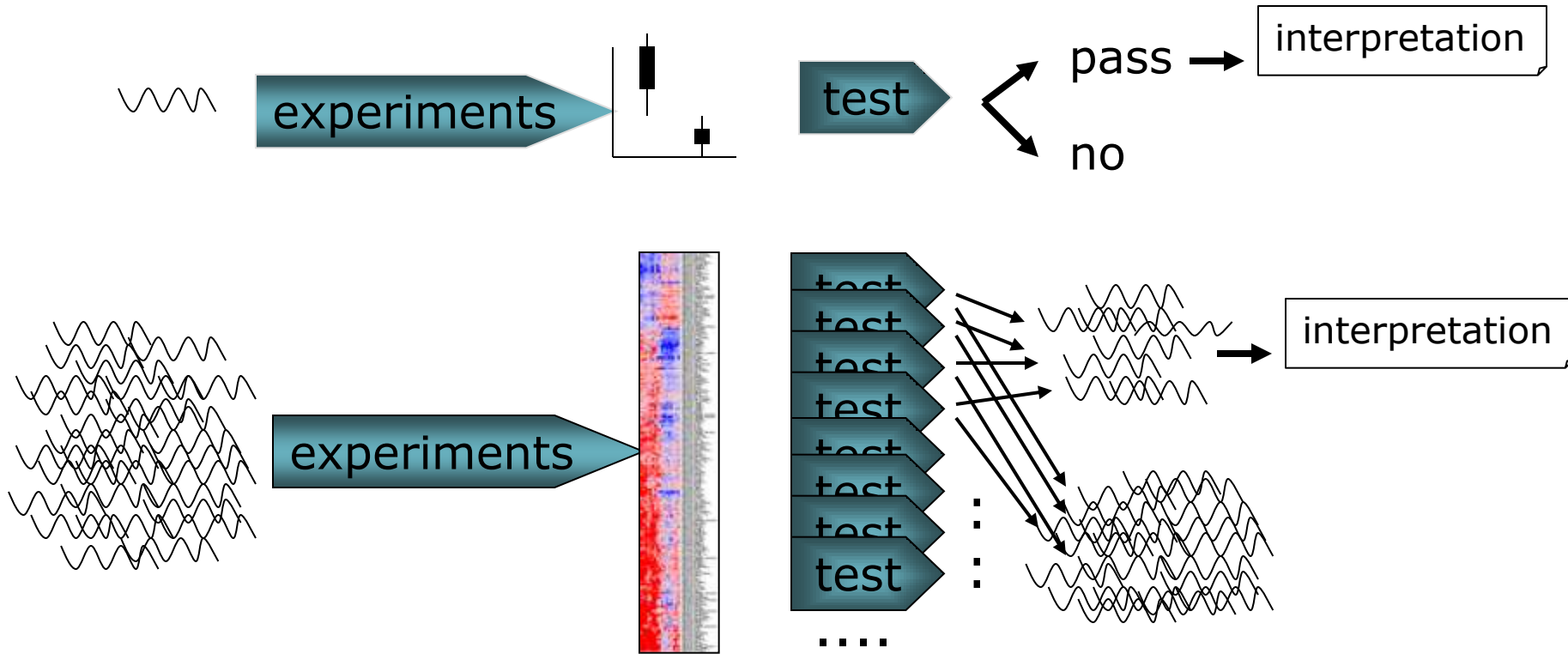
Term	0	0.0001	0	0
response to external stimulus	0	0.0001	0	0
response to stress	0	0.0002	0.0004	0.0019
signal transduction	0.0065	0.1172	0.1492	0.7191
cell motility	0.0216	0.3464	0.3734	1
resistance to pathogenic bacteria	0.0642	0.8945	0.8857	1
viral replication	0.1529	0.9887	1	1
cell death	0.1702	0.9912	1	1
regulation of gene expression, epigenetic	0.1806	0.9940	1	1

lymphocyte

enriched,
genes in the

enriched among
genes associated to
lymphomas in mature

Functional enrichment tests reproduce pre-genomics paradigms



Context and cooperation between genes is ignored

So, what is wrong with what we are doing?

Our aim:

We seek for the functions activated/deactivated in our experiment.

What we do:

We firstly seek for genes activated/deactivated one at a time (independently)

In a second step we look among them for enrichment in functions (cooperative activities) using a second test that consider functions independent.

So, what is wrong with what we are doing? (II)

This testing strategy is very strict in controlling:

Type I error (α): reject the null hypothesis when the null hypothesis is true, (false positive)

Type II error (β): fail to reject the null hypothesis when the null hypothesis is false (false negative)

But, we forget about

Type III error : get the right answer having asked the wrong question!

The testing strategy we are conducting is implicitly answering a question different to the one we want to ask.

What is the entity that accounts for functionality at the cell level?

Experiment

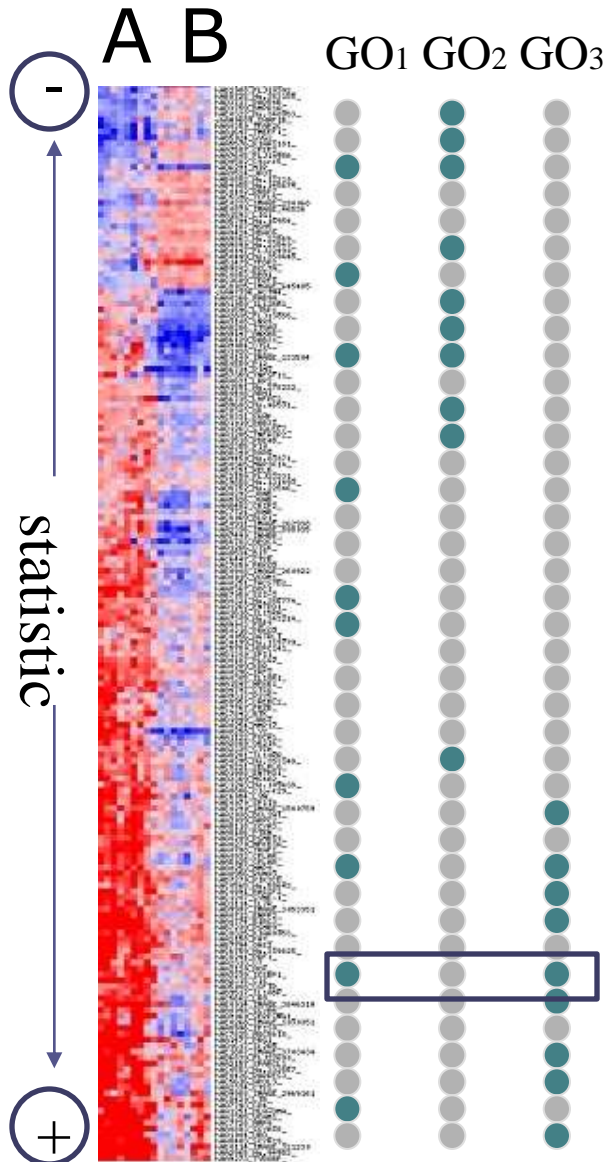


Blindfolded men (**dots in the array**) are the reporters of the individual parts (**genes**), but the reaction (**function altered**) is carried out by the elephant (**functional module, e.g. pathway**)

The wise but blindfolded men could not agree on a description of the elephant's phenotype

Therefore, why not to observe the elephant?

Cooperative activity of genes (**modules**) can be detected and related to a macroscopic observation



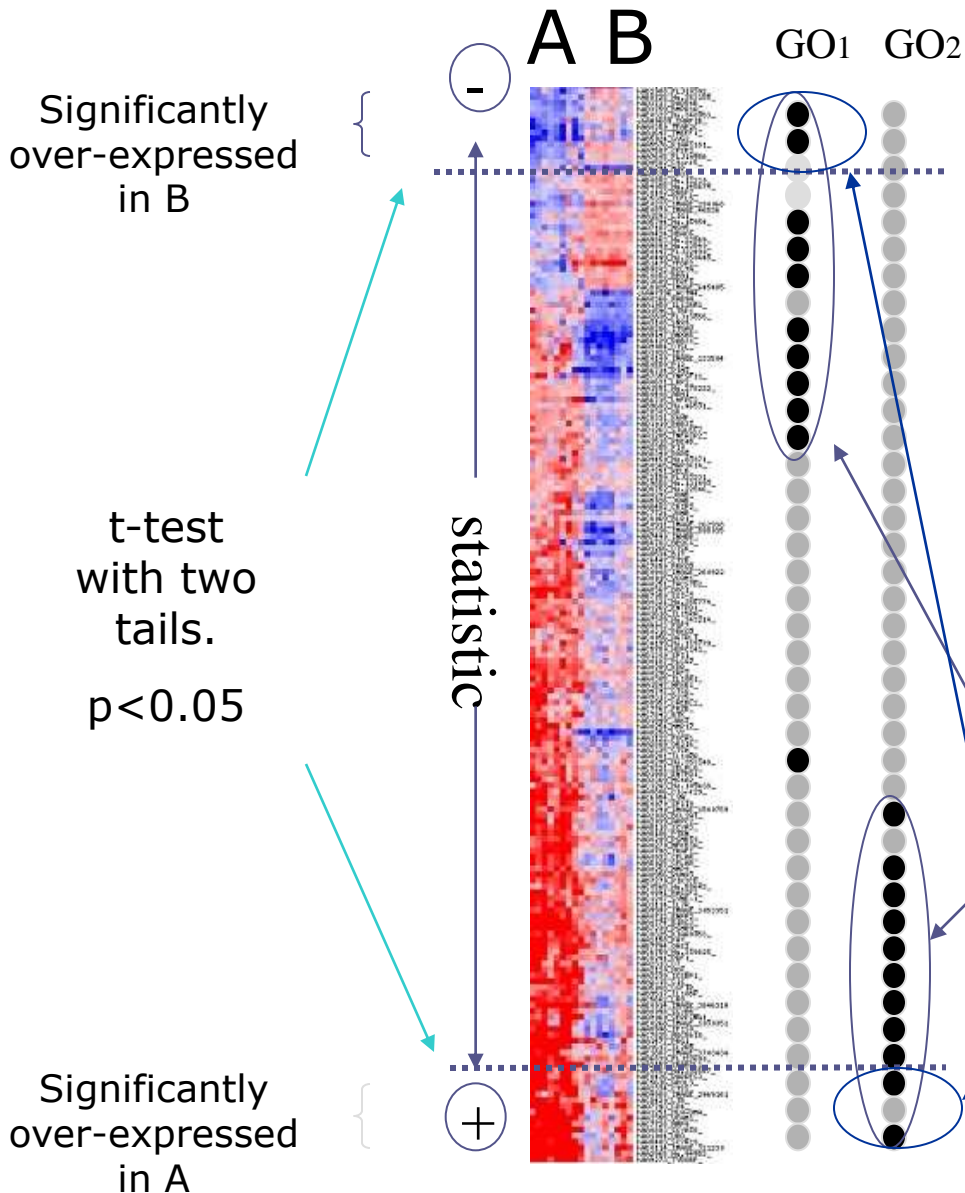
Ranking: A list of genes is ranked by their differential expression between two experimental conditions **A** and **B** (using fold change, a t-test, etc.)

Distribution of GO: Rows GO1, GO2 and GO3 represent the position of the genes belonging to three different GO terms (**modules**) across the ranking.

The first GO term is completely uncorrelated with the arrangement, while GOs **2** and **3** are clearly associated to high expression in the experimental conditions **B** and **A**, respectively.

Note that genes can be multi-functional

A previous step of gene selection causes loss of information and makes the test insensitive



If a threshold based on the experimental values is applied, and the resulting selection of genes compared for over-abundance of a functional module, this might not be found.

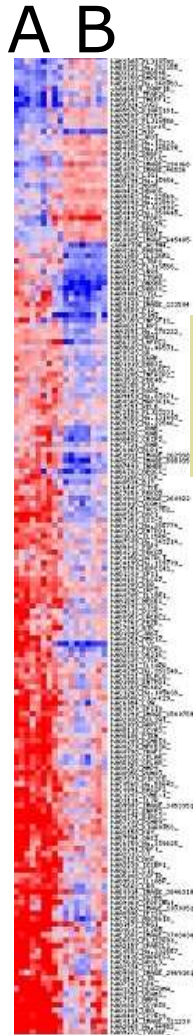
Modules expressed as blocks in A and B

Very few genes selected to arrive to a significant conclusion on GOs 1 and 2

Case study: functional differences in a class comparison experiment

A

8 with impaired tolerance (IGT) + 18 with type 2 diabetes mellitus (DM2)



No one single gene shows **significant** differential expression upon the application of a t-test



Healthy vs diabetic	Functional class	Repository		
		GO	KEGG	Swissprot keyword
Up-regulated	Oxidative phosphorylation	X	X	
	ATP synthesis		X	
	Ribosome		X	
	Ubiquinone			X
	Ribosomal protein			X
	Ribonucleoprotein			X
	Mitochondrion	X		X
	Transit peptide			X
	Nucleotide biosynthesis	X		
	NADH dehydrogenase (ubiquinone) activity	X		
Nuclease activity	X			
Dow-regulated	Insulin signalling pathway		X	

B

17 with normal tolerance to glucose (NTG)

Nevertheless, many pathways, and functional blocks are **significantly** activated/deactivated

Beyond discrete variables: Survival data

Since FatiScan depends only on a list of ordered genes, and not on the original experimental values, it can be applied to different experimental designs

Microarrays
34 samples from
tumours of
hypopharyngeal
cancer (GEO
GDS1070)



Gene
selection

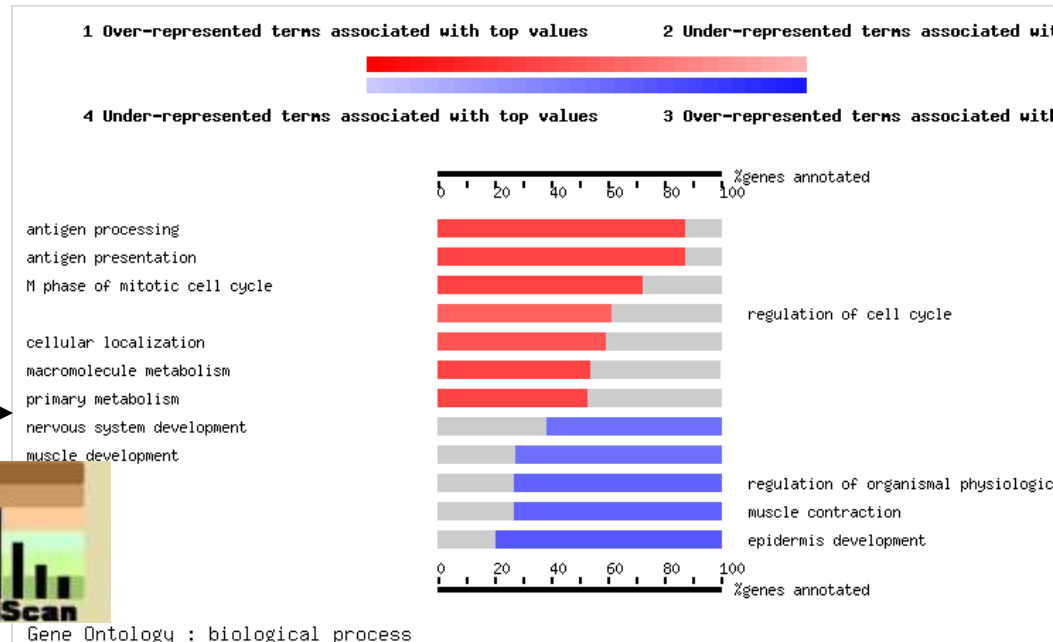
Cox Proportional-Hazards model to study how the expression of each gene across patients is related to their survival

- Survival

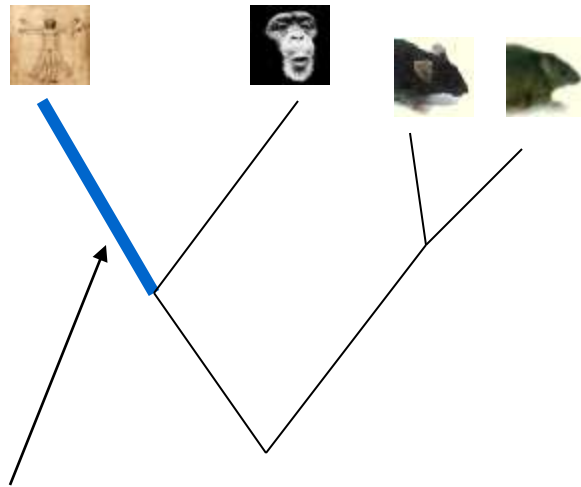
Gen	risk
Gen1	5.8
Gen2	5.6
Gen3	5.4
Gen4	5.2
Gen5	5.2
Gen6	5.0
.....
.....
.....
Gen1000	-6.0
Gen1001	-6.3



+ Survival



Beyond arrays: evolutionary systems biology



We are interested in the human lineage

Mutations occur on single genes but natural selection acts on phenotypes by operating on whole sub-cellular systems (represented by GO).

Comparison of the relative rates of synonymous (**Ks**) and non-synonymous (**Ka**) substitutions. The ratio of these values, the ($\omega = \mathbf{Ka/Ks}$) is a widely accepted measure of the selective pressure

20,469 known Ensembl human protein-coding genes from the Ensembl v.30.35h were used

Gene-set analysis of GO terms positively selected in humans

GO term

p-value

sensory perception of smell (GO:0007608)

1.3×10^{-5}

sensory perception of chemical stimulus (GO:0007606)

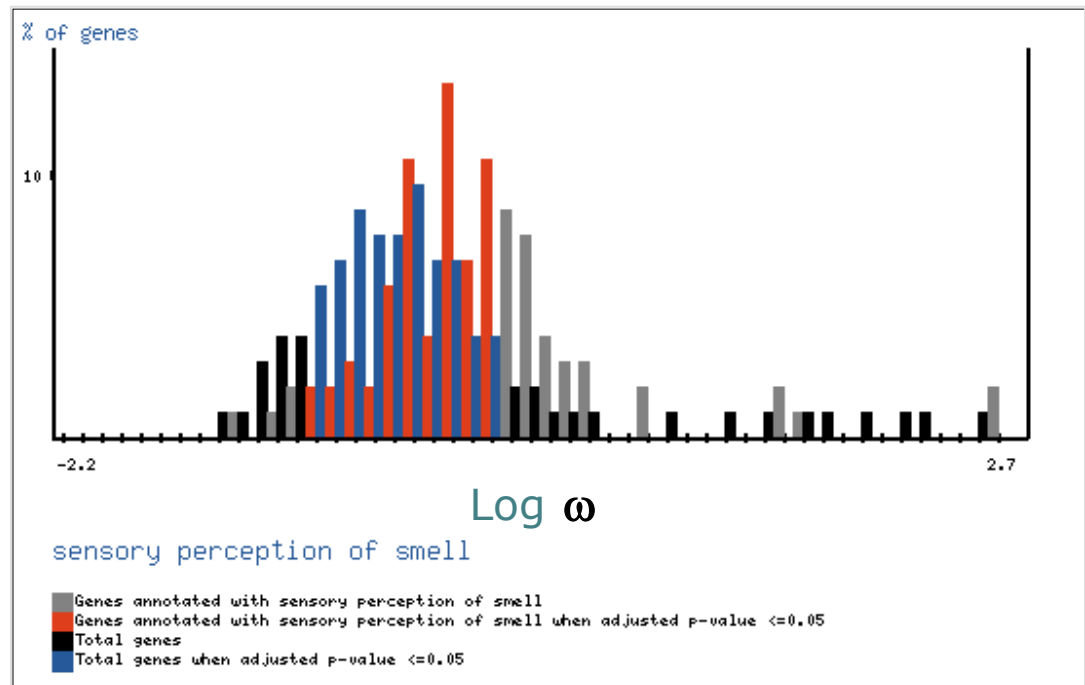
0.0014

G-protein coupled receptor protein signalling pathway (GO:0007186)

0.0095

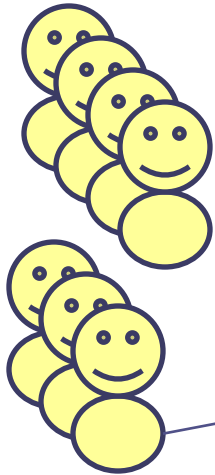
Gene-set enrichment is applied to the list of human genes ordered according to ω values

If genes positively selected are firstly tested (one at a time) and then analysed for significant enrichment of GO (functional enrichment), no results are found

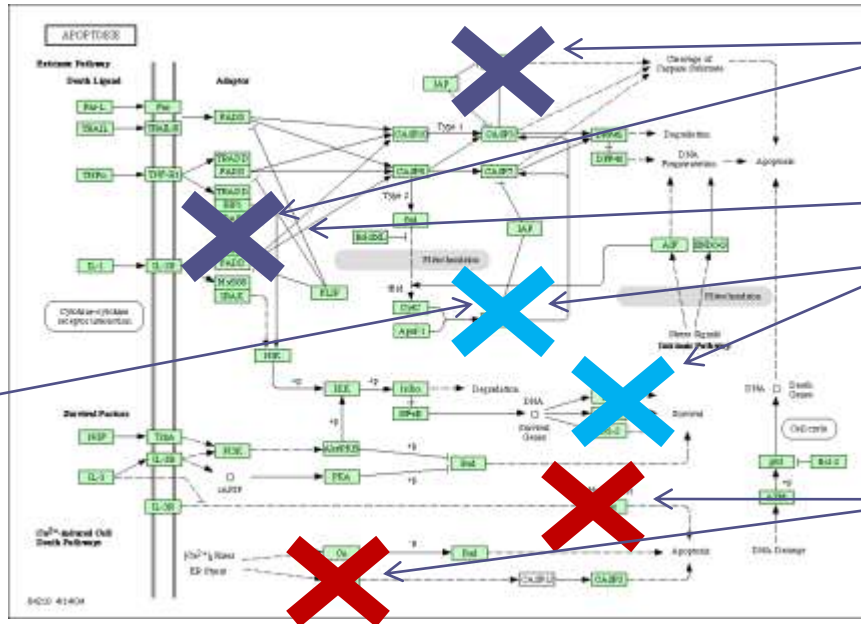
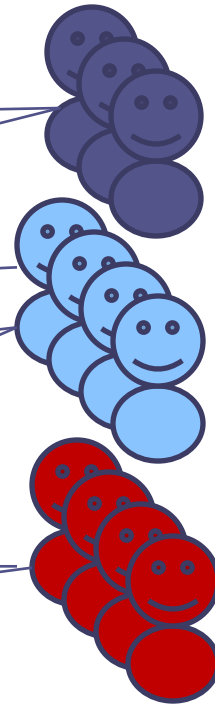


Expanding the concept of gene-set analysis to GWAS

Controls



Cases



The cases of the multifactorial disease will have different mutations (or combinations). Many cases have to be used to obtain significant associations to many markers. The only common element is the pathway (unknown at this moment) affected.

Gene-set analysis of GWAS

SNPs are mapped to
genes in LD.

Genes are arranged by
the highest association
value among the
corresponding SNPs

Gene-set analysis (or
**Pathway-Based
analysis, PBA**) is
conducted on the gene
ranked list.

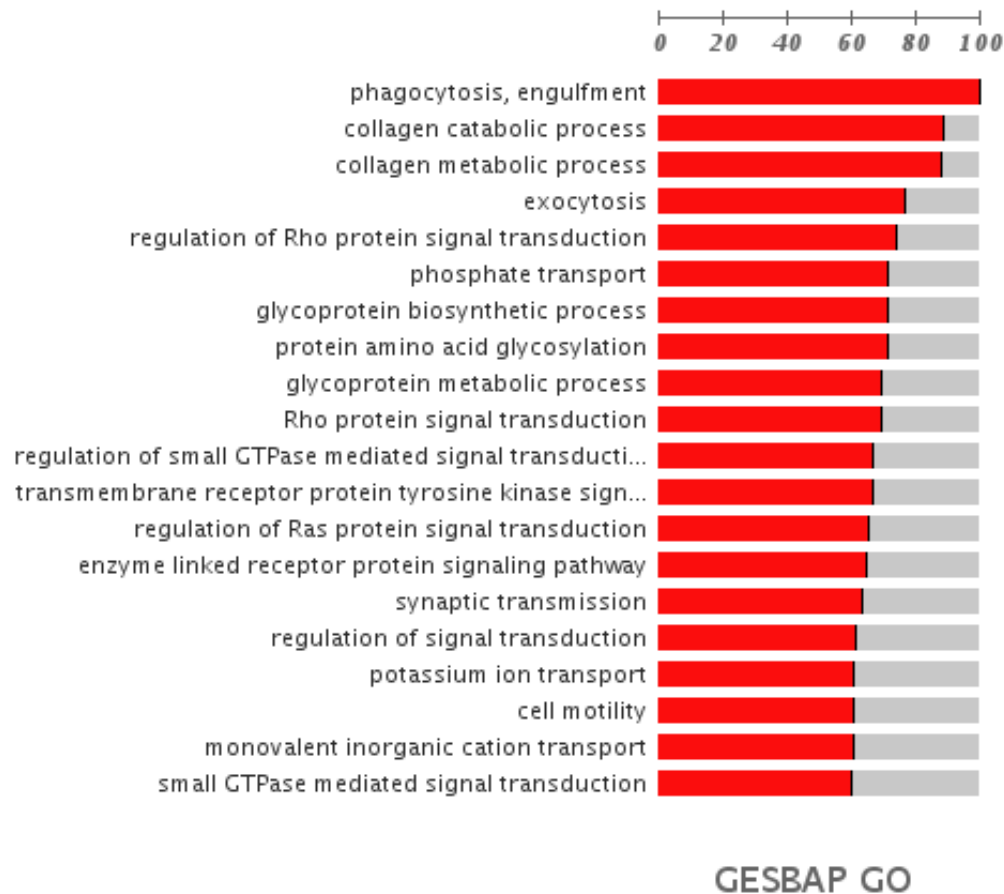
GESBAP can do that

<http://bioinfo.cipf.es/gesbap>



The screenshot displays the GESBAP web application interface. At the top, there is a navigation bar with 'Start', 'View Jobs', and 'Documentation' links. The main content area is titled 'Test your data' and includes several sections: 'Select your organism' (set to 'human'), 'Select your data' (with options for 'Load a gene data example', 'Select data type' (SNP, Gene, Genotype), and 'Association file' upload), 'Databases' (with checkboxes for GO biological process, KEGG pathways, and BioCarta), 'GO biological process options' (including 'GO parameters' with 'Minimum level' and 'Maximum level' dropdowns, and 'Filter terms by number of associated genes in BB' with 'Minimum (typical 5-20)' and 'Maximum (typical 100-inf)' dropdowns), and 'Filter terms by keyword (e.g. metastasis cancer)' with a 'Keywords' field and a dropdown menu. At the bottom, there is a 'Job name' field with the value 'affet' and a 'Run' button. The footer contains the text 'GESBAP 1.0.0 - bioInfo 2009. CIPF - Last update Wed Dec 02 15:26:03 CET 2009' and 'Terminado'.

An example of GSA in GWAS



Breast Cancer

CGEMS initiative.

(Hunter et al. Nat

Genet 2007)

1145 cases 1142

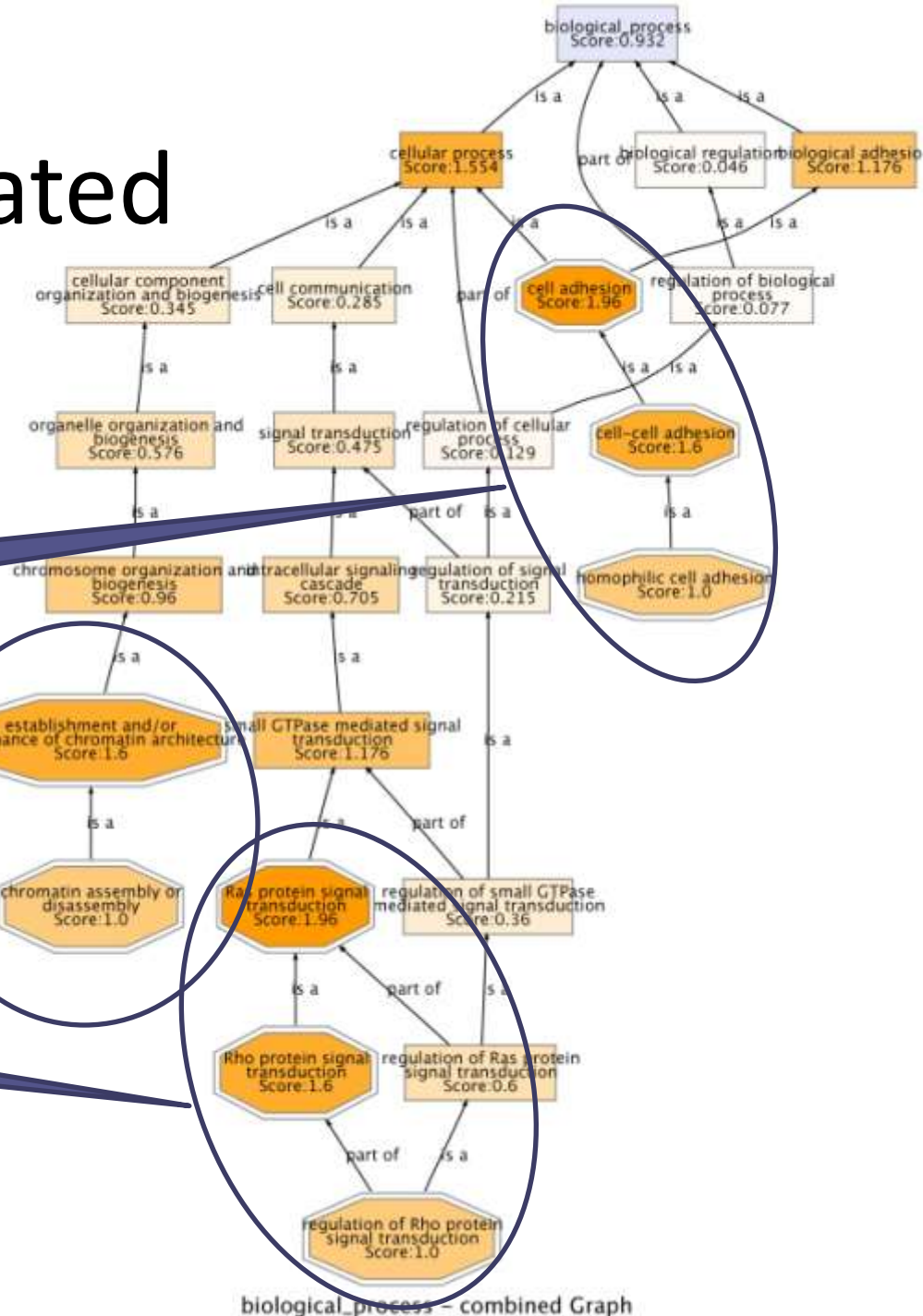
controls. Affy 500K

Only 4 SNPs were significantly associated, mapping only in one gene:

FGFR2

PBA reveals 19 GO categories including *regulation of signal transduction* (FDR-adjusted p-value= 4.45×10^{-03}) in which FGFR2 is included.

GO processes significantly associated to breast cancer



Metastasis

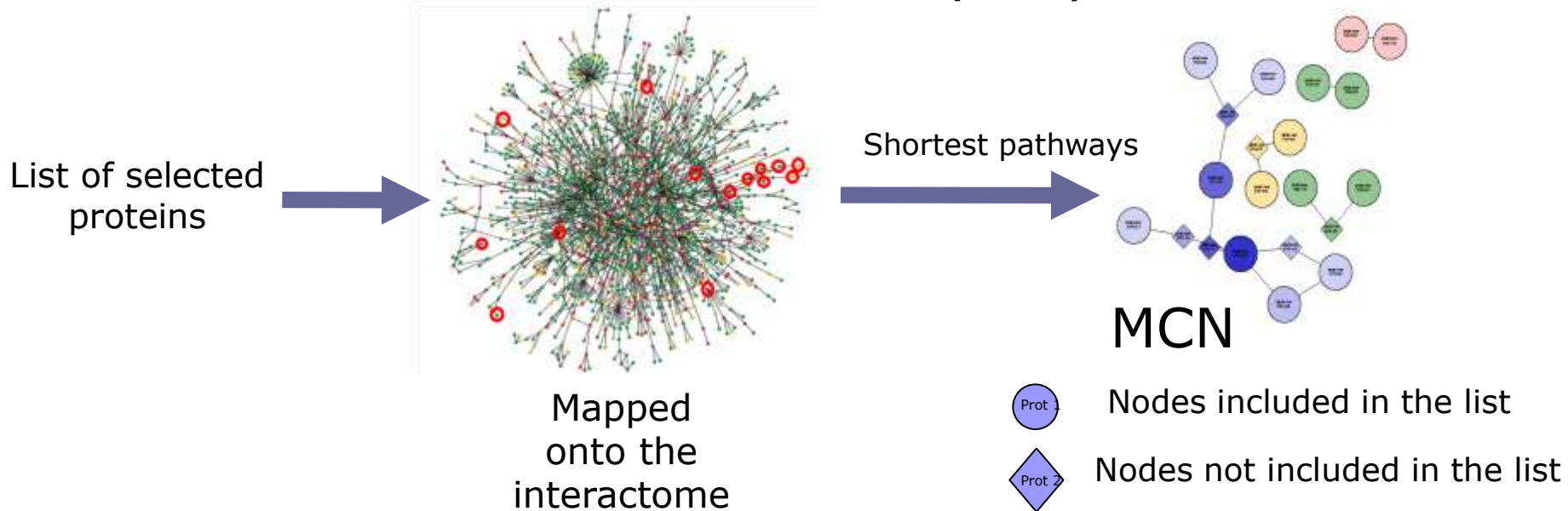
Chromosomal instability

Rho pathway

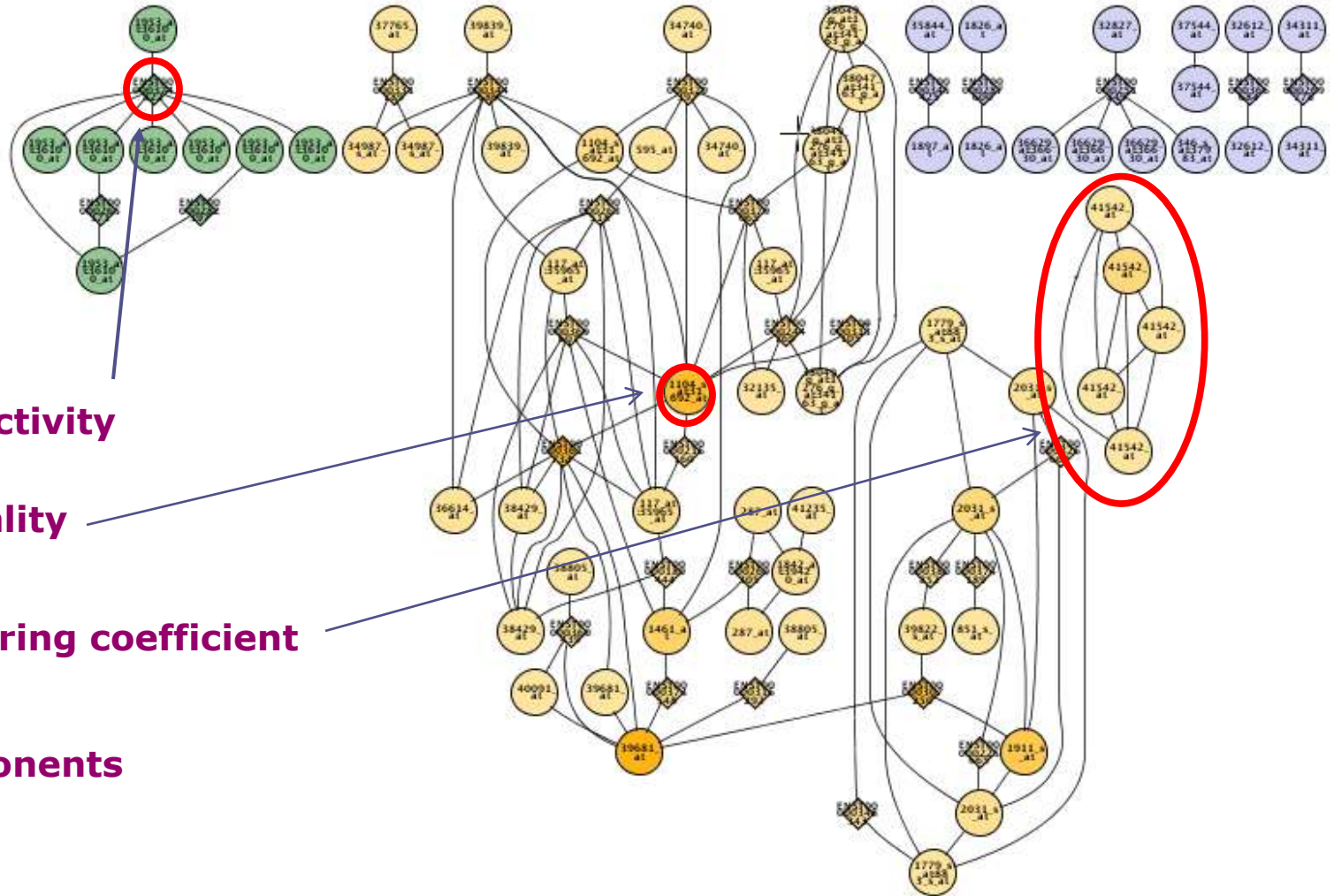
Using other gene modules: Protein-protein interaction networks

Evaluation of the cooperative behaviour of a list of genes

Shortest pathways between all pairs of nodes in the list.
The minimum connection network (MCN)



Network parameters



1 Connectivity

2 Centrality

3 Clustering coefficient

4 Components

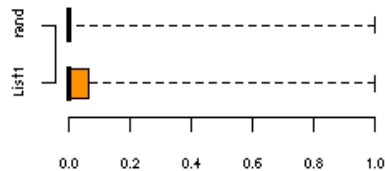
Evaluation of the Minimum Connection Network (MCN)

Parameters to evaluate: connectivity, centrality, clustering coefficient, components

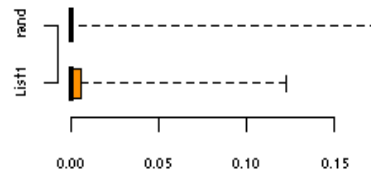
Distribution of the parameters' values versus distribution in random MCNs (compared through Kolmogorov-Smirnov tests)

Network parameters

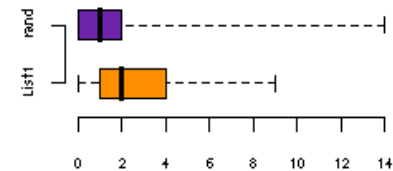
Clustering Coeff:
List1 > Random pval=**1e-04**



Betweenness:
List1 > Random pval=**2e-04**



Connections:
List1 > Random pval=**0**

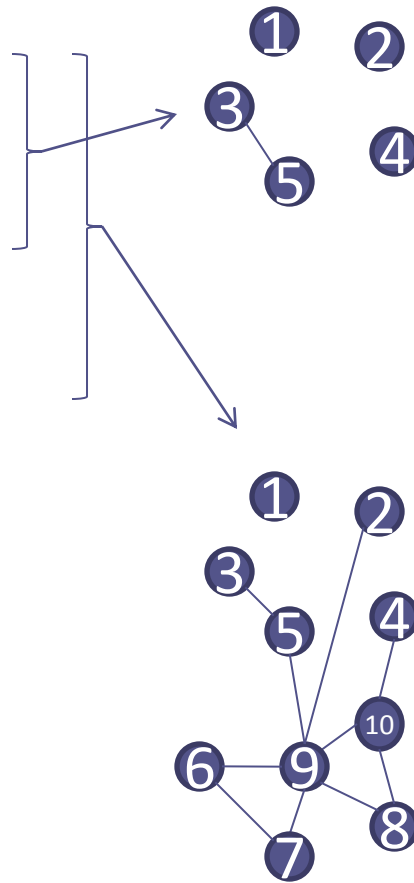


Number of components [95% confidence interval]:
Number of components with more than 1 node:
Number of Bicomponents:
Articulation points:

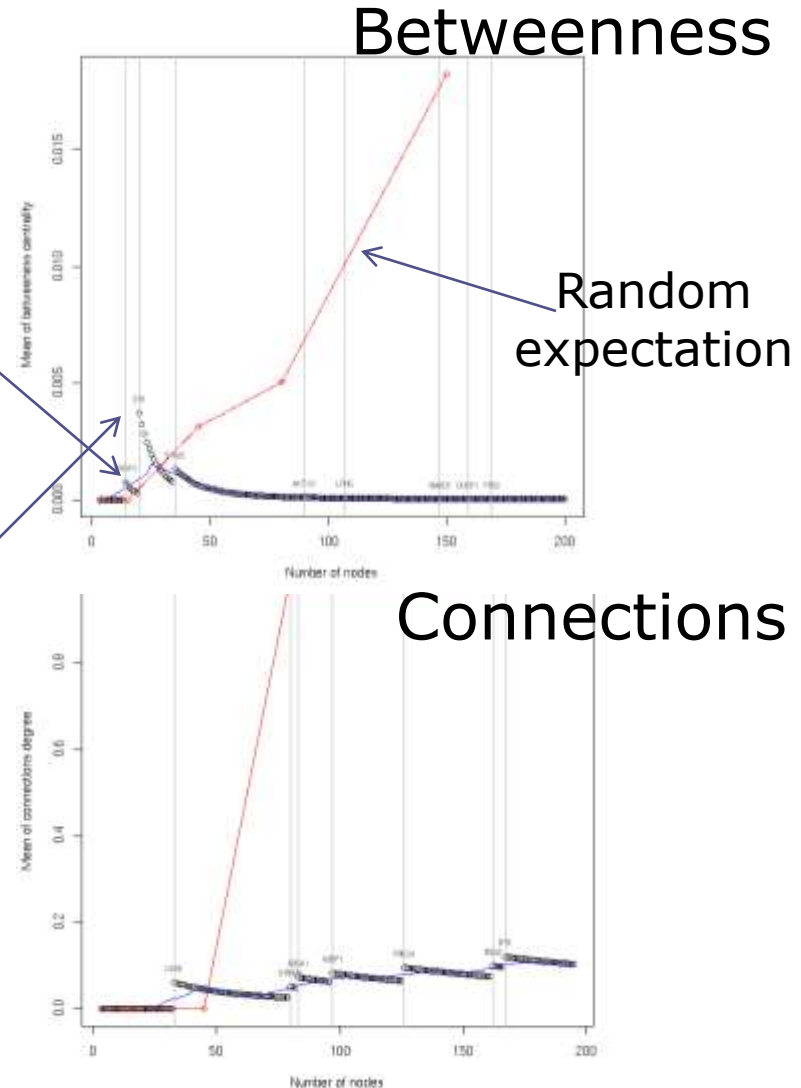
List1: 38 [46-79]
List1: 8
List1: 41
List1: 56

Study of relevant network parameters along the list of genes ranked by the most associated SNP

Gene	p-value
Gene ₁	p ₁
Gene ₂	p ₂
Gene ₃	p ₃
⋮	⋮
⋮	⋮
Gene _n	p _n

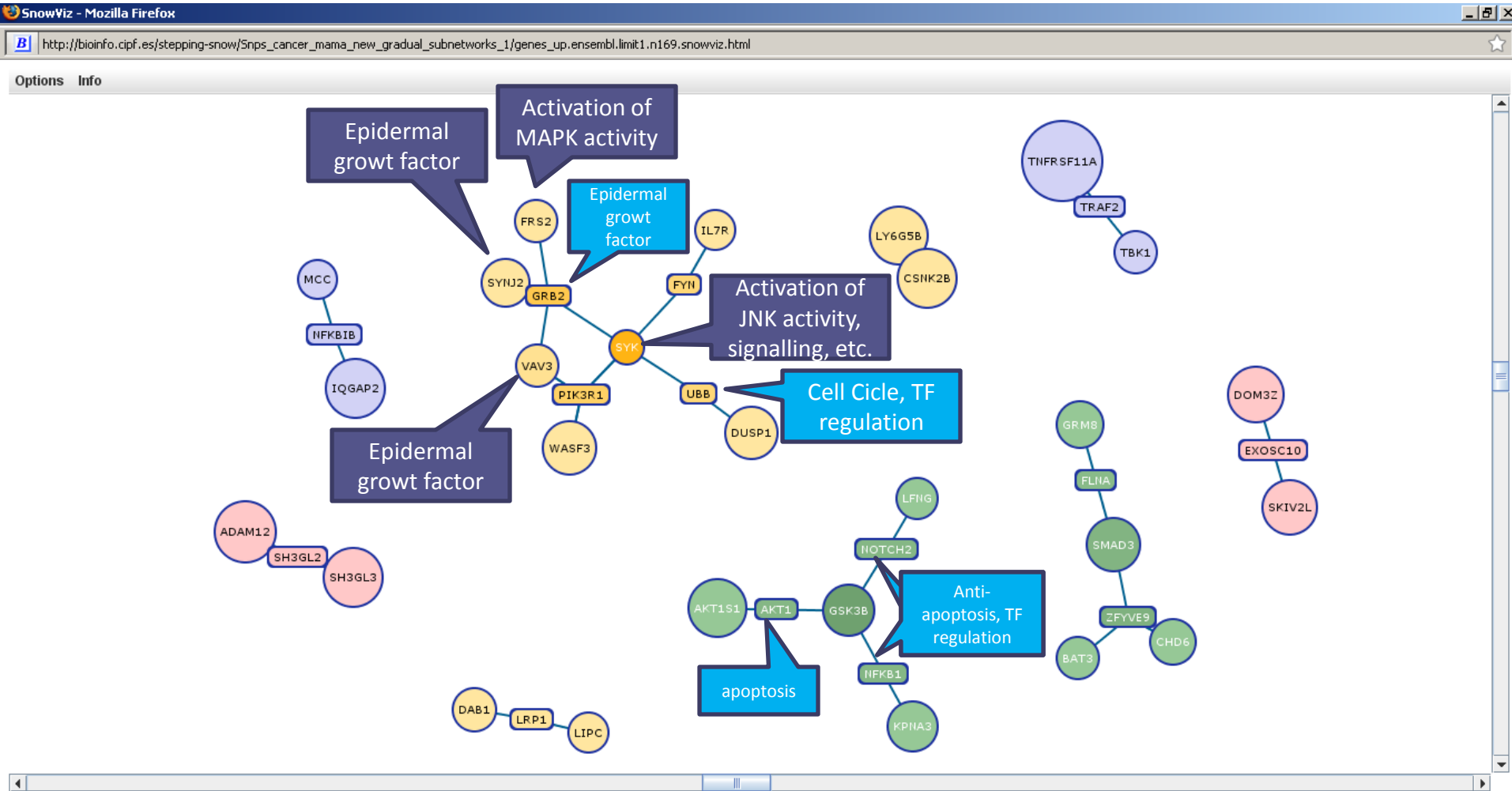


Threshold-free approach

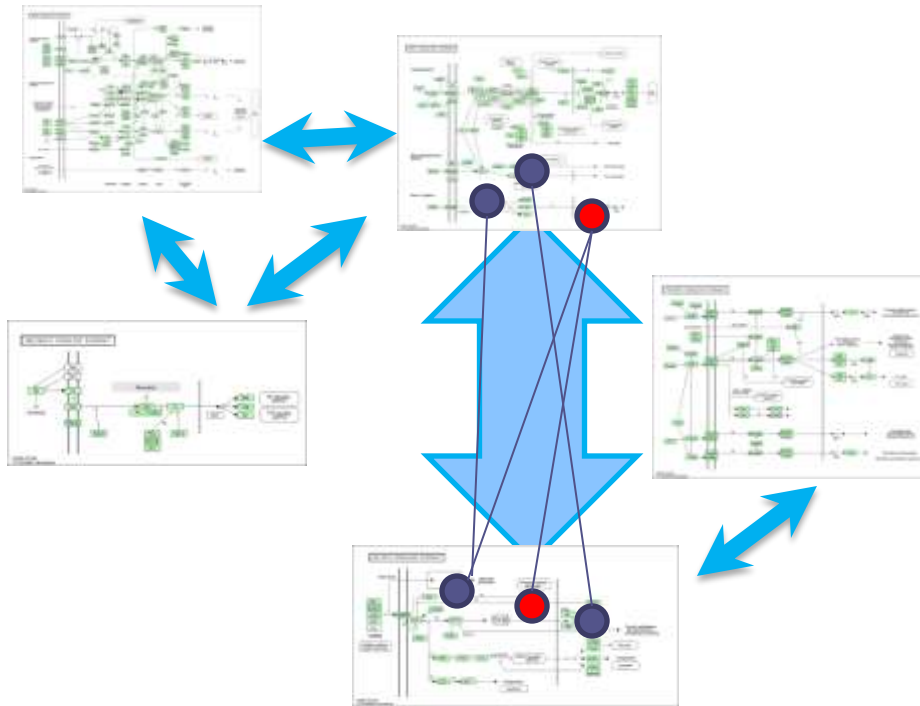


Significant connections

Breast Cancer
CGEMS initiative.
(Hunter et al. Nat
Genet 2007)



Towards a higher level of organization: relationships between modules (supermodules)



Proteins might be up or down in different experiments affecting the connectivity of the pathways.

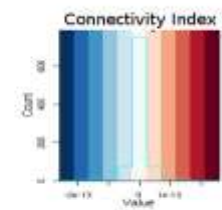
Rationale: inter-pathway PPIs implies physical proximity and suggests functional link

Towards a higher level of organization: relationships between modules (supermodules)

Changes in interactions between pathways in cancer Cellular Processes (gains in cancer)

Prostate

Mammary Gland



Redish colors mean physical connections lost in cancer
Bluish colors mean physical connections gained in cancer

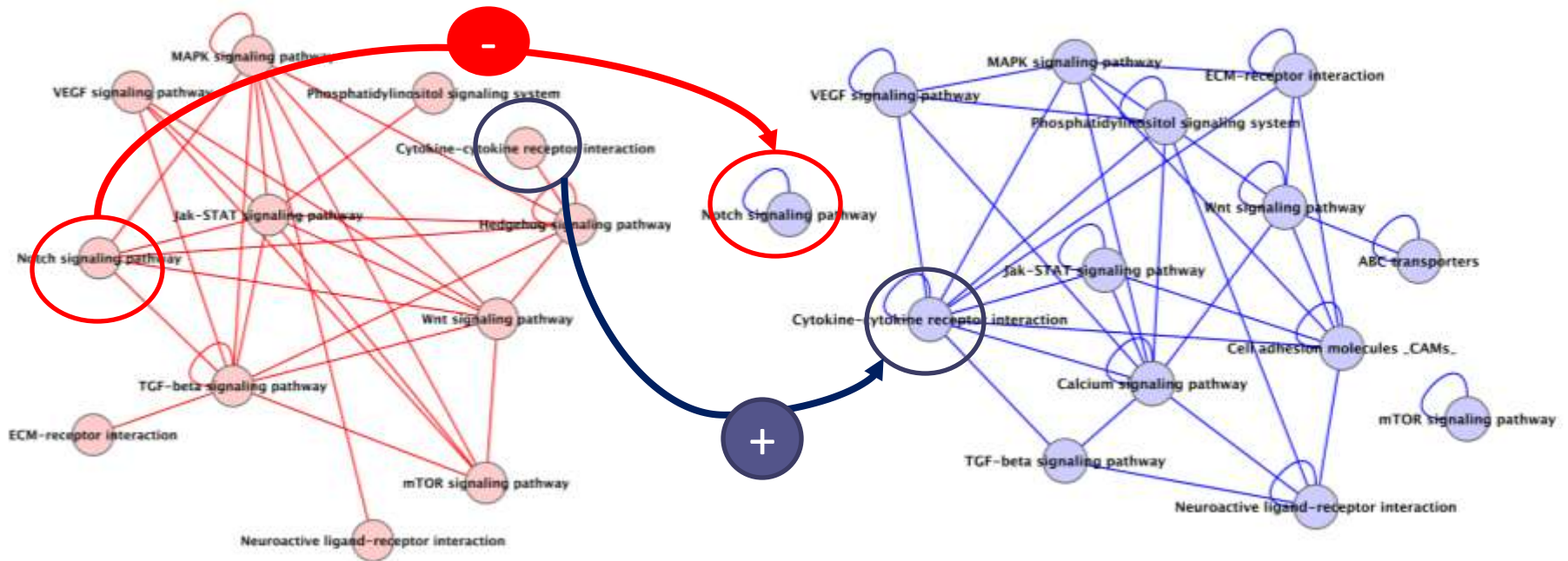
- 1 - auto-connections in Cell cycle
- 2 - Cell cycle - Tight junction
- 3 - Gap junction - Insulin signaling pathway
- 4 - Gap junction - Fc epsilon RI signaling pathway

- 5 - Toll like receptor signaling pathway auto-connections
- 6 - Toll like receptor signaling pathway - B cell receptor signaling
- 7 - Insulin signaling pathway - Melanogenesis

Relationships between pathways

Normal aoesophagus

Cancer



Diseased state causes a rewiring both within and among pathways. Based on the presence/absence of transcripts, mapped on the interactome and then in the pathways context

The babelomics suite for functional profiling of genomic experiments



Over 3000 registered users. More than 1000 experiments analysed daily

<http://www.babelomics.org>

Al-Shahrour et al., 2005, 2006, 2007, 2008 NAR; 2004, 2005 Bioinformatics, 2007 BMC Bioinformatics;

Biological information from:

- GO
- Interpro motifs
- KEGG pathways
- Biocarta pathways
- Swissprot keywords
- TFBSs (Transfac)
- Regulatory motifs (CisRED)
- miRNAs
- Protein interactions
- Tissues
- Text-mining
- Chromosomal location

For

Human, mouse, rat, chicken, cow, fly, worm, yeast, *A. thaliana* and bacteria

Tests for

- functional enrichment
- gene set enrichment
- network enrichment

The Bioinformatics and Genomics Department at the Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain, and...

Joaquín Dopazo
Eva Alloza
Leonardo Arbiza
Fátima Al-Shahrour
Davide Baù
Emidio Capriotti
Jose Carbonell
Ana Conesa
Hernán Dopazo
Pablo Escobar
Francisco García
Stefan Goetz
Martina Marbà
Marc Martí
Ignacio Medina
Pablo Minguez
David Montaner
Marina Naval
Javier Santoyo
Patricia Sebastian
François Serra
Sonia Tarazona
Joaquín Tárraga
Adriana Cucchi



...the INB, National Institute of Bioinformatics (Functional Genomics Node) and the CIBER-ER Network of Centers for Rare Diseases

