

Babelomics

Experimental Data Analysis

Jun 2014 in Cambridge

David Montaner

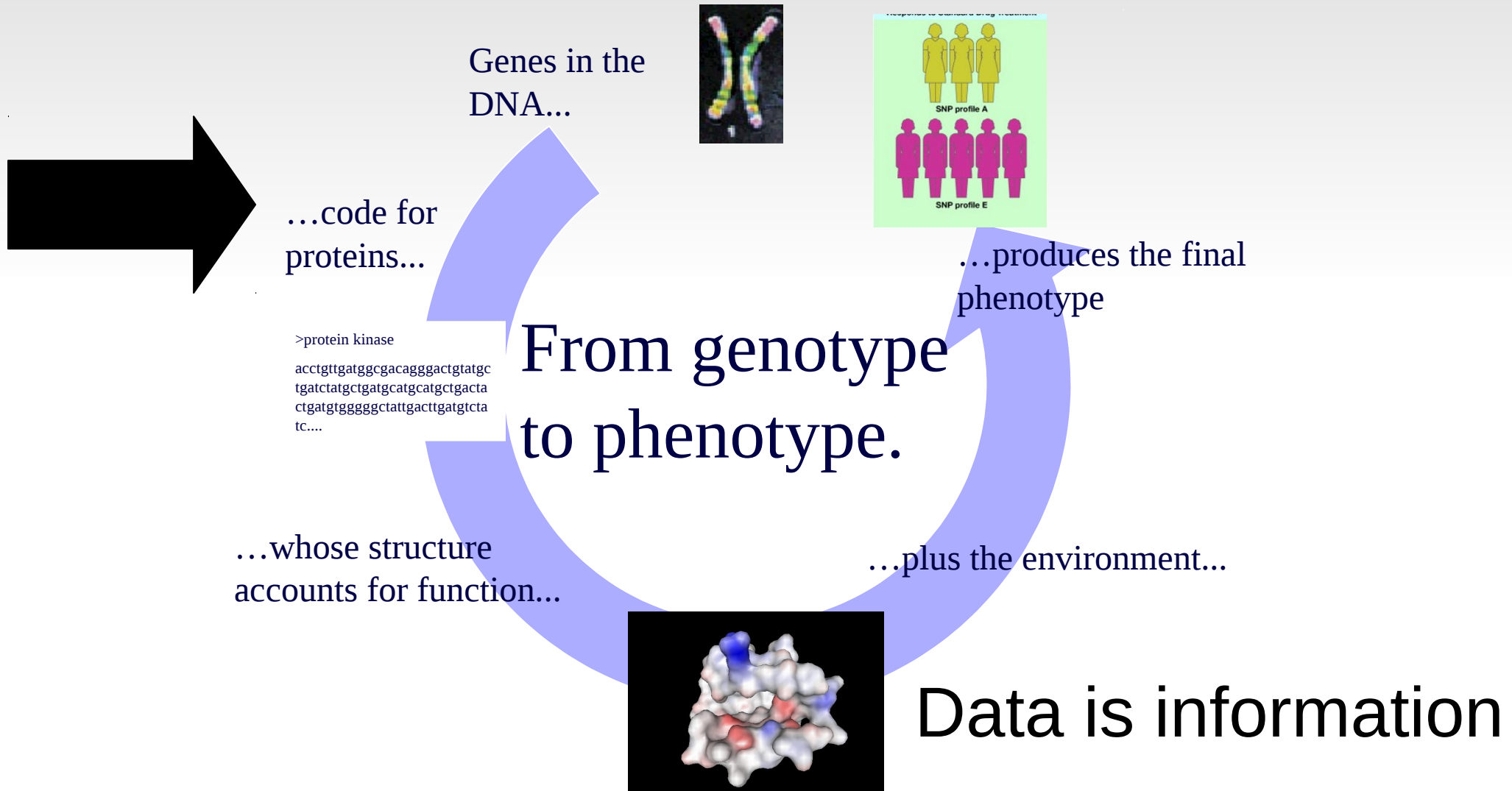
dmontaner@cipf.es

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

Bioinformatics and Genomics Department
Centro de Investigacion Principe Felipe (CIPF)
(Valencia, Spain)

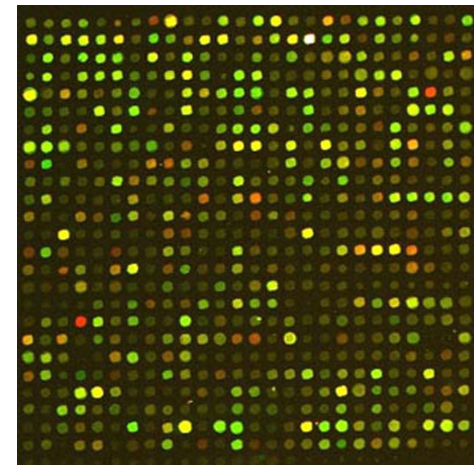
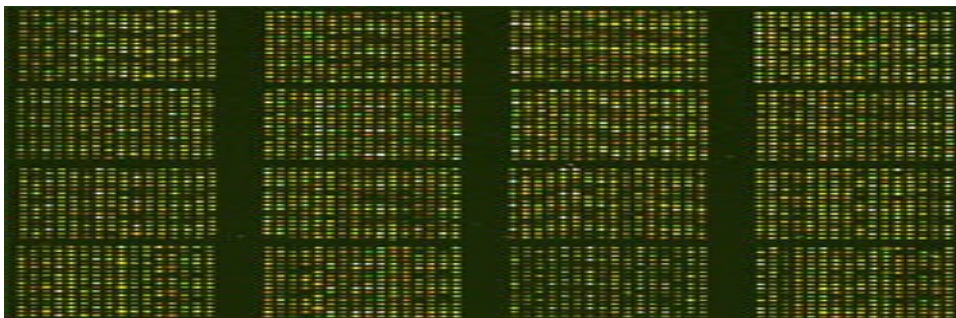


Genetic Research



High Throughput Technologies

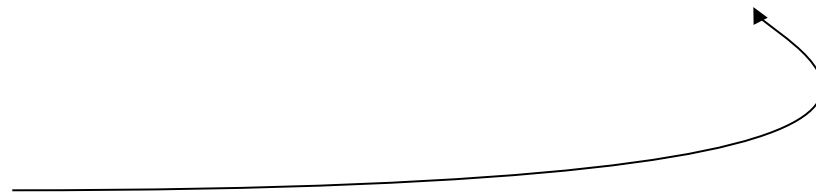
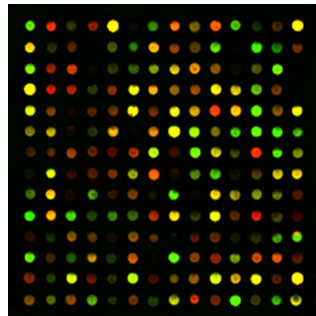
- 1988 arrayed DNAs were used
- 1991 oligonucleotides are synthesized on a glass slide through photolithography
- 1995 DNA Microarrays
- 1997 Genome wide Yeast Microarray
- 2005 First next-generation sequencing system



High Throughput Technologies

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

William Blake
Proverbs of Hell
(1790–1793)



Next Generation Sequencing
SOLID 6Gbp per round

Genes in the
DNA...

```
>protein kinase  
acctgttgatggcgacagggactgtatgctgac  
tatgctgatgcatgcatgctgactactgatgtggg  
ggctattgactgatgtctatc....
```



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

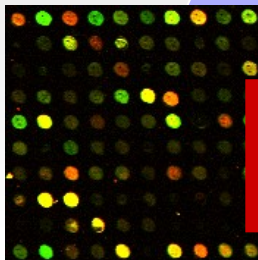
10 million
SNPs



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

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

A typical tissue is
expressing among
5000 and 10000
genes



Data \neq information

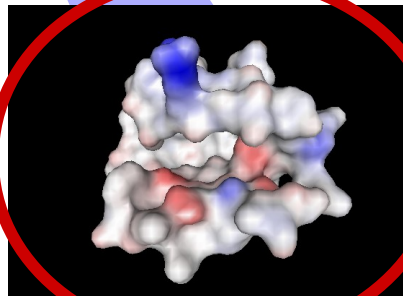
(in the functional post-genomic
scenario)

...code for
proteins...

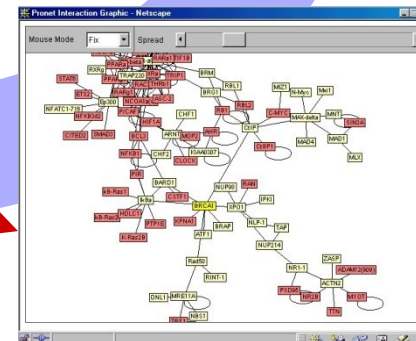
That undergo
post-translational
modifications, somatic
recombination...

100K-500K proteins

...whose structures account
for function...



...conforming complex
interaction networks...



...in cooperation
with other
proteins...

Each protein has an average
of 8 interactions


Bioinformatics

- 2003 GEPAS
- Bioinformatics department open at CIPF
- 2006 BABELOMICS
- 2008 Blast2GO
- 2010 GEPAS and some characteristics of Blast2GO are included into Babelomics.

Microarrays arrive to an acceptable level of reproducibility



ARTICLES



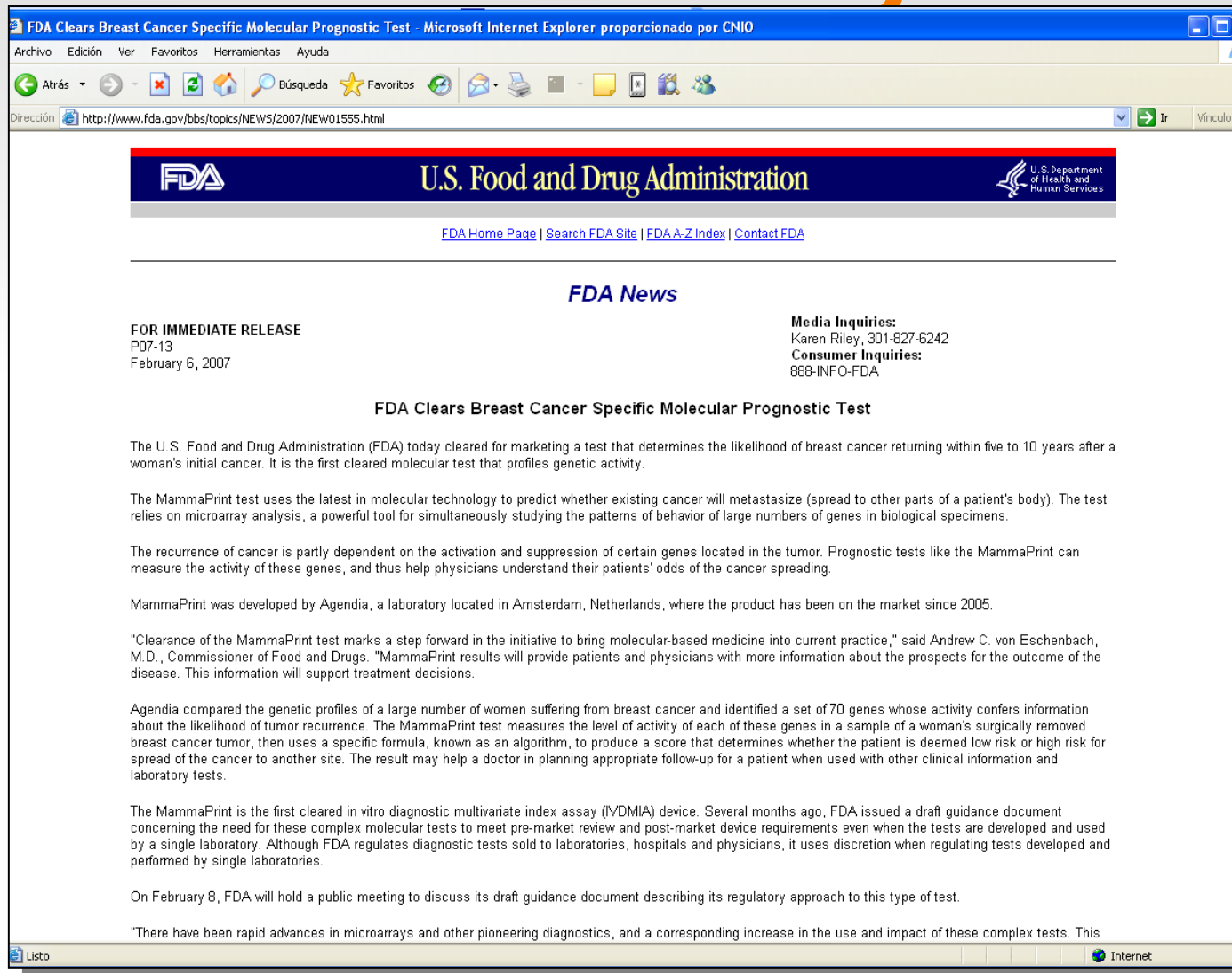
The MicroArray Quality Control (MAQC) project shows inter- and intraplatform reproducibility of gene expression measurements

MAQC Consortium*

Over the last decade, the introduction of microarray technology has had a profound impact on gene expression research. The publication of studies with dissimilar or altogether contradictory results, obtained using different microarray platforms to analyze identical RNA samples, has raised concerns about the reliability of this technology. The MicroArray Quality Control (MAQC) project was initiated to address these concerns, as well as other performance and data analysis issues. Expression data on four titration pools from two distinct reference RNA samples were generated at multiple test sites using a variety of microarray-based and alternative technology platforms. Here we describe the experimental design and probe mapping efforts behind the MAQC project. We show intraplatform consistency across test sites as well as a high level of interplatform concordance in terms of genes identified as differentially expressed. This study provides a resource that represents an important first step toward establishing a framework for the use of microarrays in clinical and regulatory settings.

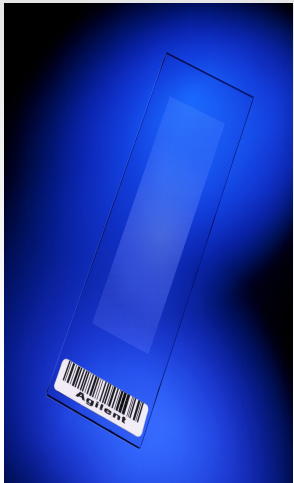
Shing Group <http://www.nature.com/naturebiotechnology>

FDA approves the first predictor based on microarrays



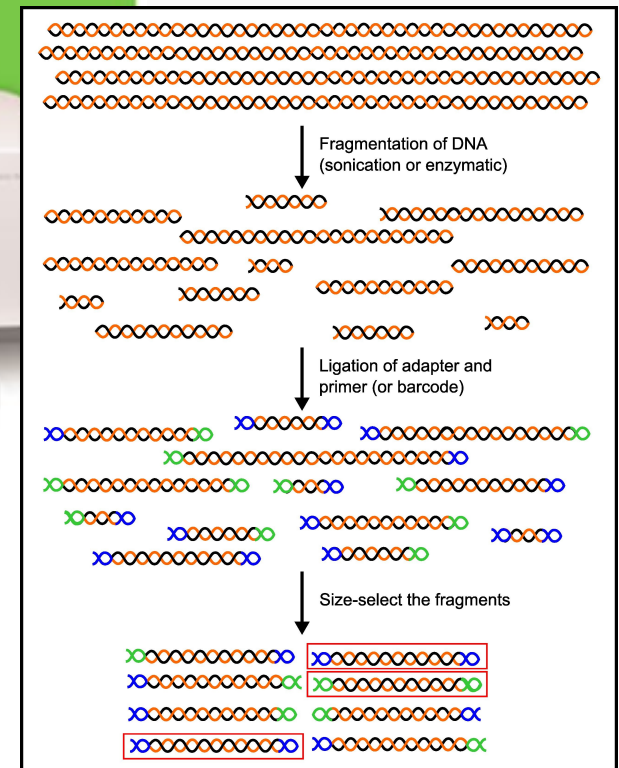
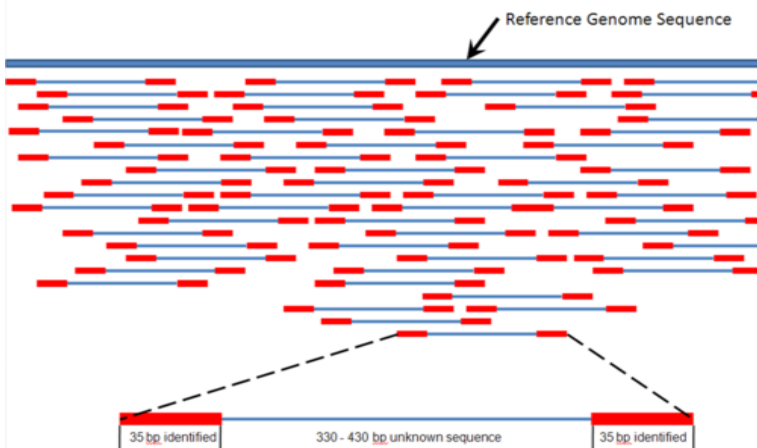
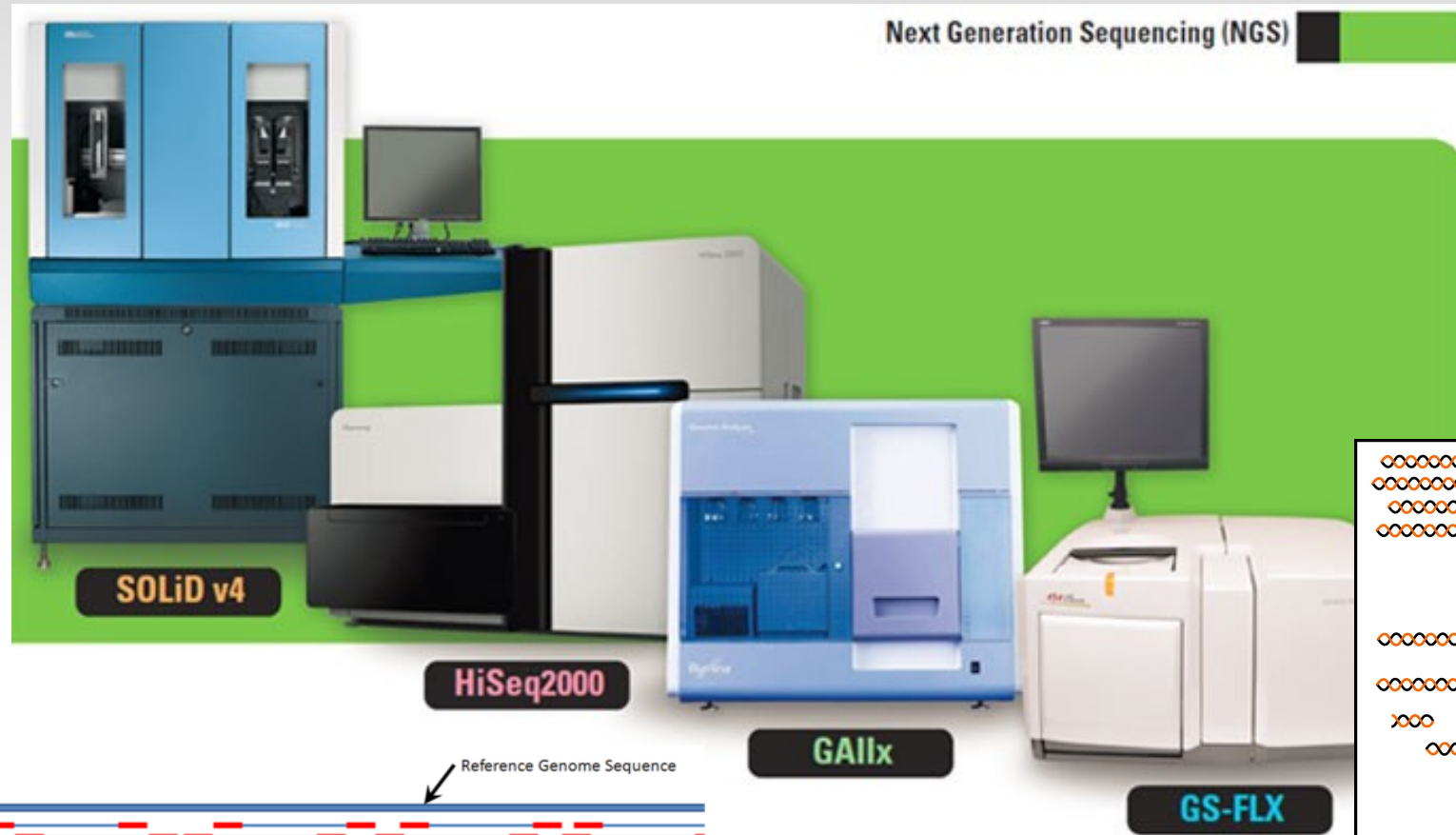
February 2007: MamaPrint

DNA Microarrays: the paradigm of a post-genomic technique

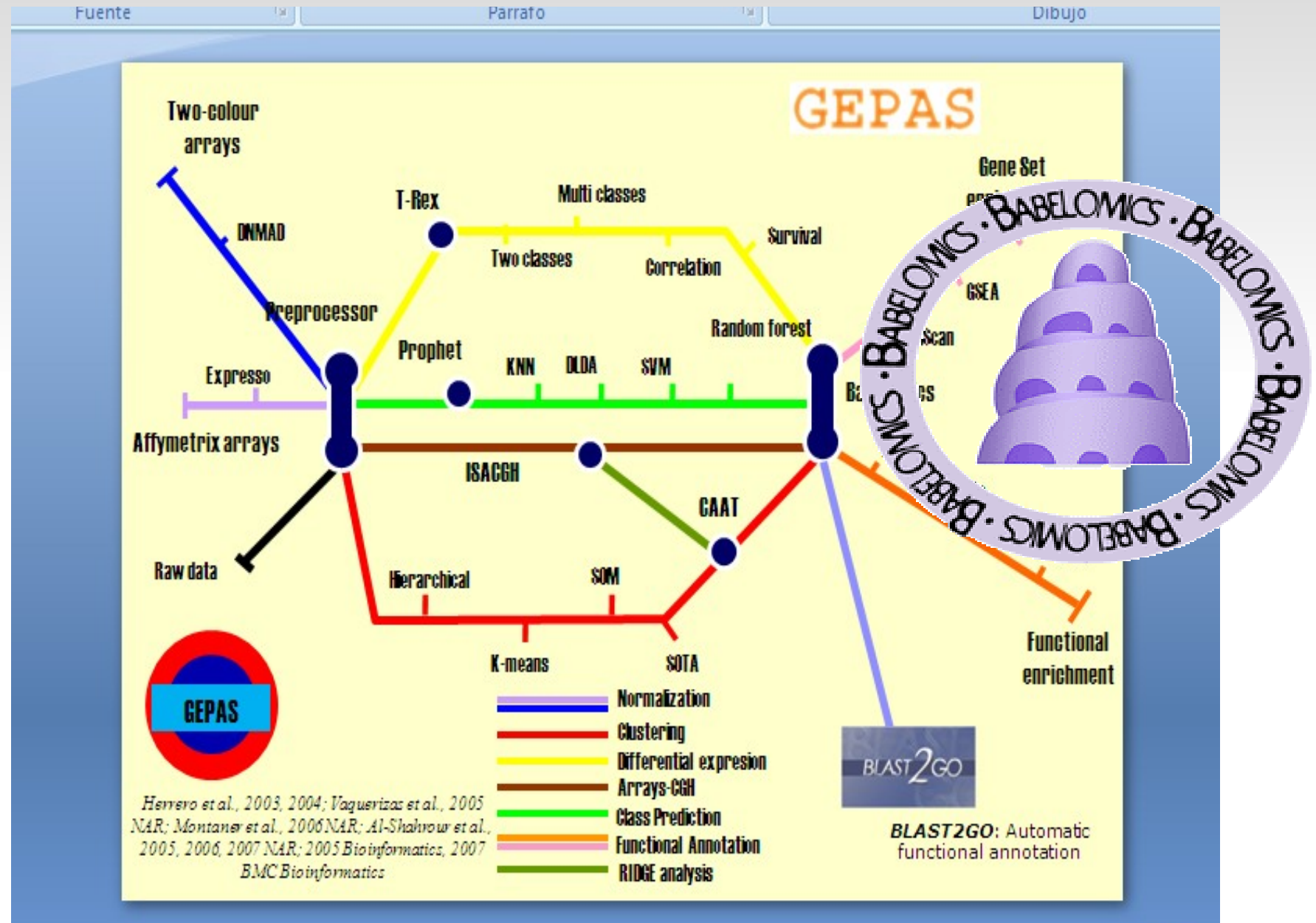


gene1	10.23	9.98	10.41	10.55	10.65	9.69
gene2	10.51	9.74	10.65	10.63	10.43	10.35
gene3	9.89	10.02	9.89	11.03	10.21	10.77
gene4	10.25	10.83	8.94	10.16	10.49	10.46
gene...

NGS



Functional interpretation



Functional profiling of genome-scale experiments in the post-genomic era

My data...

How are structured?

What are these groups?

What is this gen?

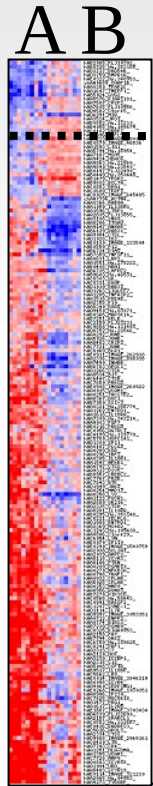
Microsoft Excel

Archivo Edición Herramientas Datos Ventana Ayuda

Expresión Analysis: Pivot Tab

Totaldata.xls [Solo lectura]

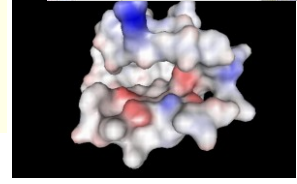
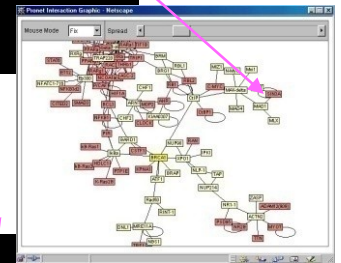
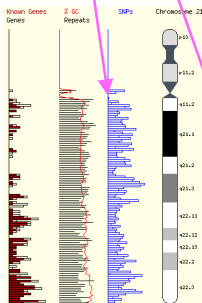
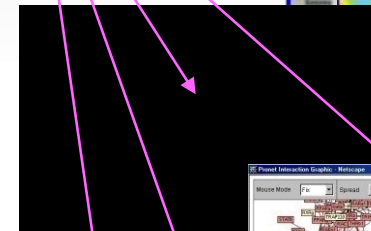
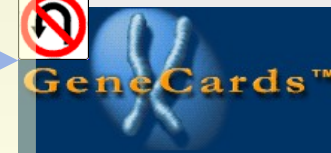
	E	F	G	H	I	J	K	L	M	N
65	578.6 *		1.4		0.26	M12481	Mouse cytoplasmic beta-actin mRNA (5'...			
66	534.9 *	~1.6			0.22	M12481	Mouse cytoplasmic beta-actin mRNA (5'...			
67	403.6 *	~1.5			0.15	X61388	SGD: YEL002C Yeast S. cerevisiae WBP1 Oligosacchar...			
68	535.2 *	~1.6			0.22	U18530	SGD: YEL018W Yeast S. cerevisiae Protein of unknown fu...			
69	567.7 *	~1.6			-0.27	M23316	SGD: YEL024C Yeast S. cerevisiae RIP1 Rieske iron-sulfur...			
70	-114.5 *	~1.1			-0.03	K02207	SGD: YEL021W Yeast S. cerevisiae URA3 gene coding for...			
71	-125.4 *		-1		-0.01	Cluster Incl	M16465 Calpactin I light chain /cds=(68,361) /gb=M16...			
72	-1091.6		-1.2		-0.14	Cluster Incl	Z67748 M. musculus spermidine synthase gene /cds=(...			
73	-757.2		-1.3		-0.17	Cluster Incl	X12973 Murine MLC1FAMLC3F gene for myosin alkali I...			
74	9626.6		1.3		0.83	Cluster Incl	A849035 UHM-AH1-agw-a06-0-U1.s1 Mus musculus c...			
75	-847.4		-1.3		-0.21	Cluster Incl	AW123542 UHM-BH2.1-agf-F01-0-U1.s1 Mus musculus c...			
76	2583.1		1.1		0.09	Cluster Incl	AF055983 Mus musculus proteasome alpha7/CS8 subu...			
77	192.5 *	~1.2			0.05	Cluster Incl	AB006361 Mus musculus mRNA for prostaglandin D s...			
78	2980.2 *	~4.4			1.63	Cluster Incl	AB006361 Mus musculus mRNA for prostaglandin D s...			
79	-20.1		-1		0	Cluster Incl	AB011081 Mus musculus mRNA for huntingtin interact...			
80	1380.9 *	~2.6			1.81	Cluster Incl	AB011081 Mus musculus mRNA for huntingtin interact...			
81	753.2 *		1.2		0.1	Cluster Incl	U97170 Mus musculus protein kinase inhibitor gamma...			
82	-2774.7		-1.9		-1.43	Cluster Incl	M36120 Keratin complex 1, acidic, gene 19 /cds=(0,12...			
83	3614.4 *	~5.1			1.98	Cluster Incl	U19604 DNA ligase 1, ATP-dependent /cds=(304,3054)			
84	0 *	~0.0			0	Cluster Incl	A851492 UHM-BH0-ajp-04-0-U1.s2 Mus musculus cD...			
85	3310.9		1.2		0.24	Cluster Incl	AB025408 Mus musculus mRNA for sid478g, complet...			
86	-1291		-1.5		-0.42	Cluster Incl	AF059735 Mus musculus C-terminal binding protein 2...			
87	-263.3 *	~1.3			-0.09	Cluster Incl	AF063454 Mus musculus tetraspan TM4SF (Tspan-6)			
88	77.5 *		1.1		0.01	Cluster Incl	D45650 Hydroxysteroid 17-beta-dehydrogenase 1 /cds...			
89	2047.2 *	~3.3			1.1	Cluster Incl	AF039299 Mus musculus 17-beta-hydroxysteroid dehy...			
90	809.9 *	~1.9			0.38	Cluster Incl	M84457 Vascular cell adhesion molecule 1 /cds=(57,2...			
91	-124.3 *	~1.1			-0.03	Cluster Incl	U12884 Mus musculus C57BL/6 vascular cell adhesio...			
92	-675.5 *	~1.8			-0.37	Cluster Incl	U12884 Mus musculus C57BL/6 vascular cell adhesio...			
93	1465.4 *	~2.7			0.76	Cluster Incl	AJ238636 Mus musculus mRNA for nucleoside diphos...			
94	838.2		1.1		0.1	Cluster Incl	U70475 Nuclear, factor, erythroid derived 2, like 2 /cds...			
95	4969.4 *	~6.7			8.84	Cluster Incl	AF045673 Mus musculus FLH-LRR associated protein...			
96	148.3 *	~1.2			0.04	Cluster Incl	A891475 u59a06.x1 Mus musculus cDNA, 3' end /cd...			



Cell cycle...

DBs Information

I M19380 Calmodulin 3 /cds=(109,558) /gb=M19380 /gi=469419
 I A842326 UHM-AH1-afz-b-11-0-U1.s1 Mus musculus cDNA, 3'
 I AJ242663 Mus musculus mRNA for cathepsin Z precursor (cts
 I U12620 Dipeptidyl/peptidase 4 /cds=(117,2399) /gb=U12620 /g
 I M13444 Mouse alpha-tubulin isotype M-alpha-4 mRNA, compl
 I U11027 Mus musculus C57BL/6J Sec61 protein complex gam
 I J05926 Phosphofructokinase, liver, B-type /cds=(42,2384) /gb=
 I Z67748 M. musculus mRNA for phosphatase 2A catalytic subu
 I U80932 Serine/threonine kinase 6 /cds=(48,1235) /gb=U80932
 I U47024 Maternal embryonic message 3 /cds=(137,2401) /gb=I
 I AF075136 Mus musculus Sin3-associated protein (sap30) mR
 I M25944 Mouse carbonic anhydrase II (CAII) mRNA, 3' end /cd
 I X74671 Neurofibromatosis 2 /cds=(576,2366) /gb=X74671 /gi=
 I M12848 Mouse myb proto-oncogene mRNA encoding 71 kd m
 I AW125458 UHM-BH2.2-agm-a-07-0-U1.s1 Mus musculus cDNA
 I U84903 Ribosomal protein L23 /cds=(61,501) /gb=U84903 /gi=
 I U35141 Mus musculus retinoblastoma-binding protein (mRbAp
 I U19621 Mus musculus vesicle transport protein (munc-18c) m
 I M15268 Aminolevulinic acid synthase 2, erythroid /cds=(0,178
 I M25149 Tissue specific transplantation antigen F91A /cds=(0,
 I X65449 Calcyclin /cds=(159,426) /gb=X65449 /gi=50271 /uqph

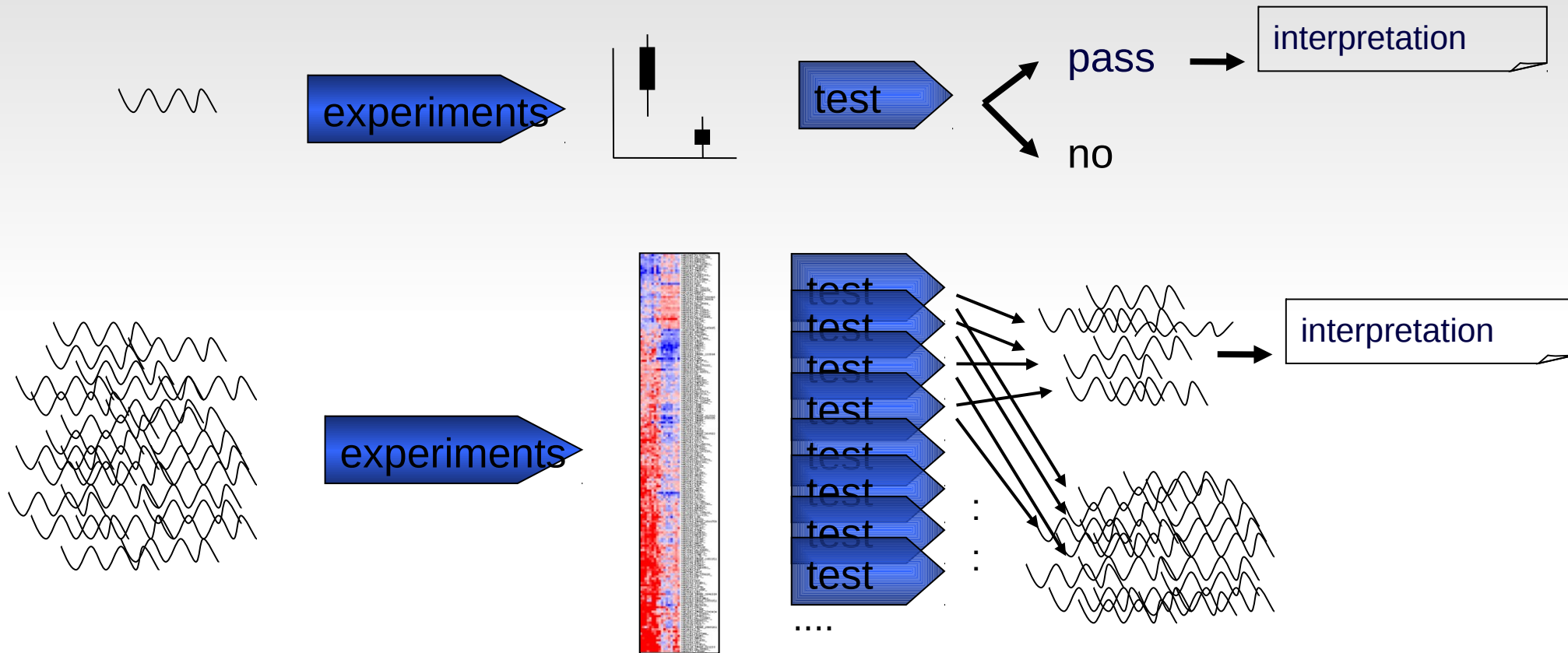


Analysis

Functional profiling

Links

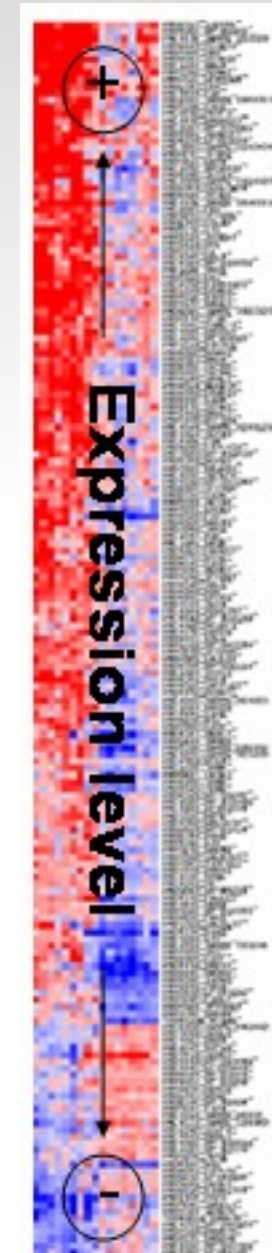
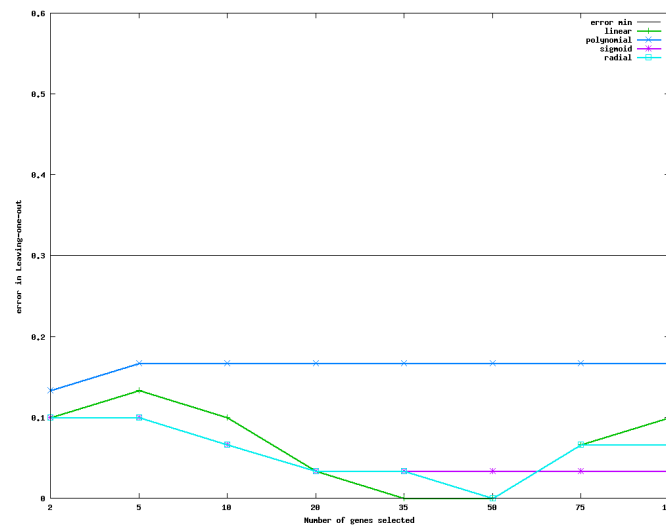
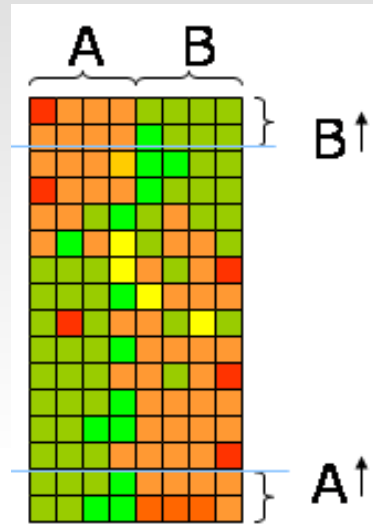
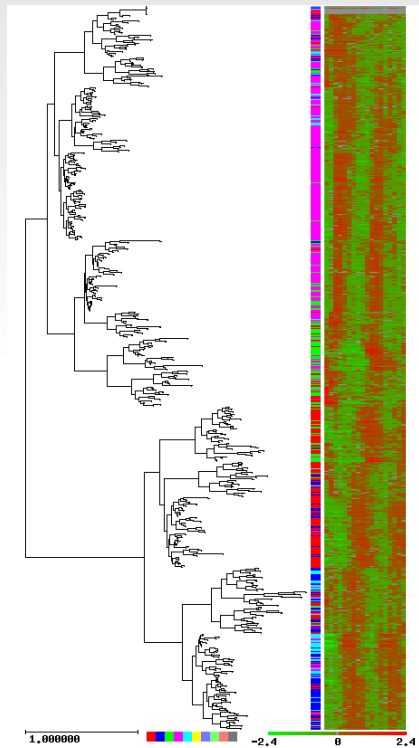
Functional enrichment approach reproduces pre-genomics paradigms



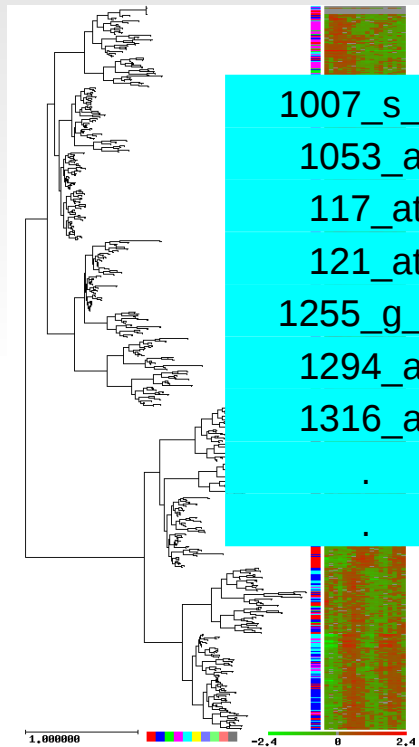
Context and cooperation between genes is not ignored

The unit of interest in the study is shifted from gene to function

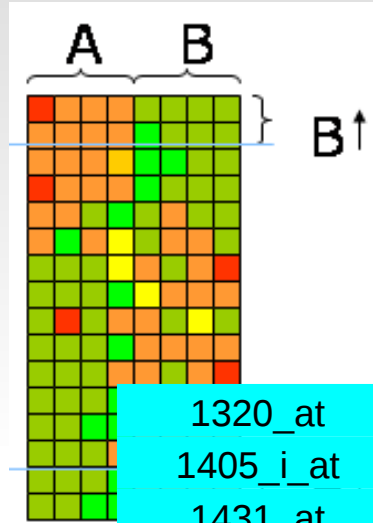
Genome-scale experiment output



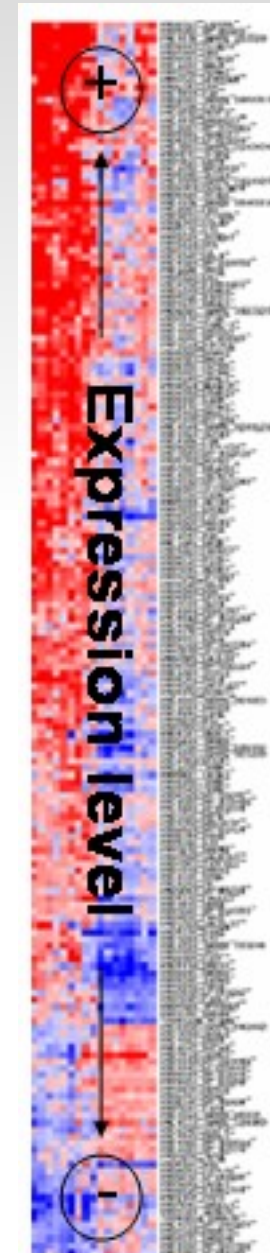
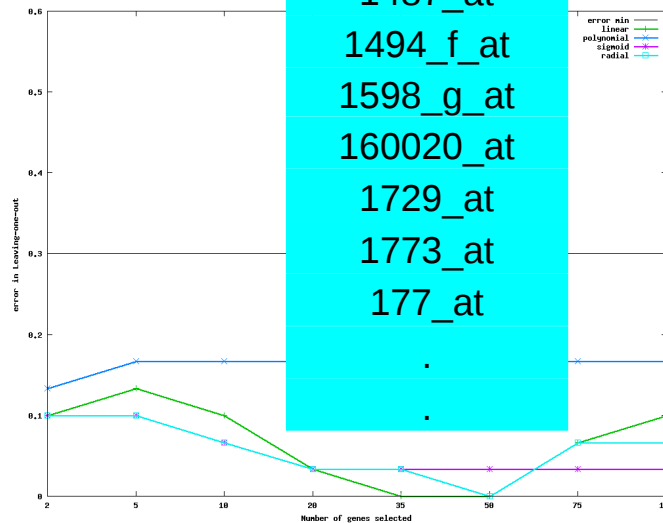
Genome-scale experiment output



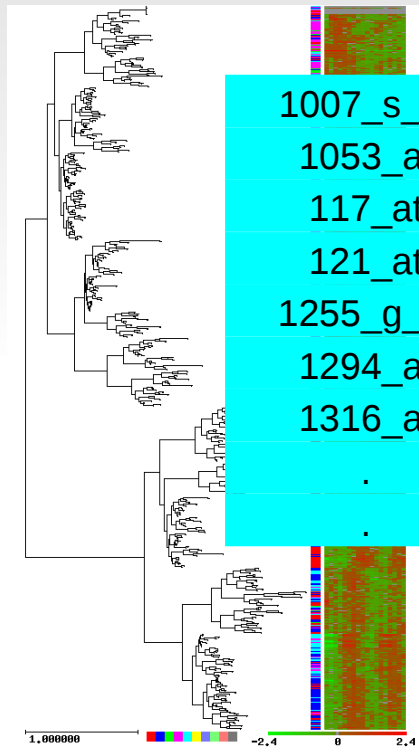
1007_s_at
1053_at
117_at
121_at
1255_g_at
1294_at
1316_at



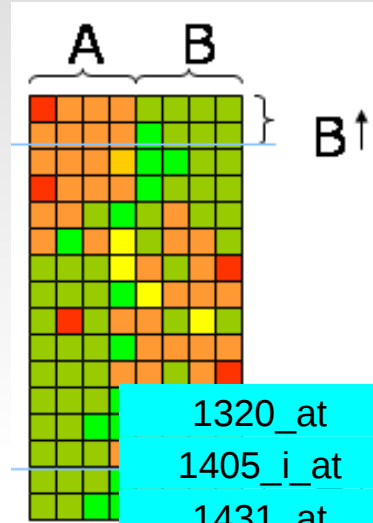
1320_at
1405_i_at
1431_at
1438_at
1487_at
1494_f_at
1598_g_at
160020_at
1729_at
1773_at
177_at



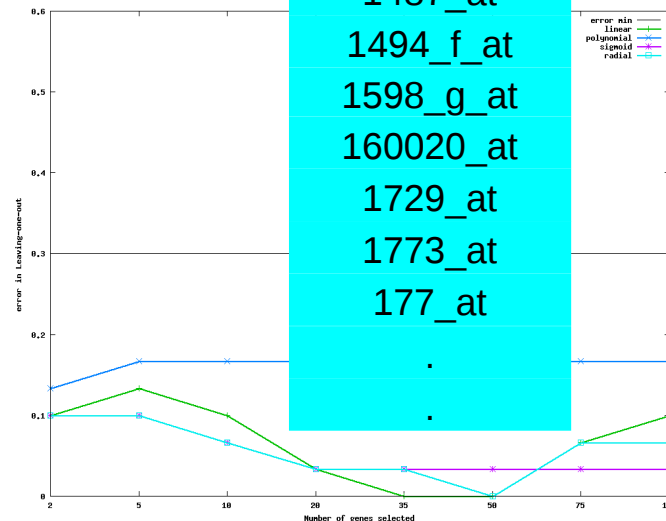
Genome-scale experiment output



1007_s_at
1053_at
117_at
121_at
1255_g_at
1294_at
1316_at



1320_at
1405_i_at
1431_at
1438_at
1487_at
1494_f_at
1598_g_at
160020_at
1729_at
1773_at
177_at



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
1438_at	6.5
1487_at	6.2
1494_f_at	5.9
1598_g_at	5.8
160020_at	4.8
1729_at	4.7

Functional interpretation

1007_s_at
1053_at
117_at
121_at
1255_g_at
1294_at
1316_at
.
.

1320_at
1405_i_at
1431_at
1438_at
1487_at
1494_f_at
1598_g_at
160020_at
1729_at
1773_at
177_at
.
.

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
1438_at	6.5
1487_at	6.2
1494_f_at	5.9
1598_g_at	5.8
160020_at	4.8
1729_at	4.7
.	.
.	.

Functional interpretation

Experimental results
observed in the lab
(not always a wet-lab)

Recorded to:

- Test a hypothesis.
- Get a first insight of a biological process.

1007_s_at

1053_at

117_at

121_at

1255_g_at

1294_at

1316_at

12.4

11.5

10.3

10.2

9.9

9.3

8.2

8.1

7.7

7.4

6.5

6.2

5.9

5.8

4.8

4.7

.

.

1

1

14

14

1494

1598

160020

1729

1773

177

.

.

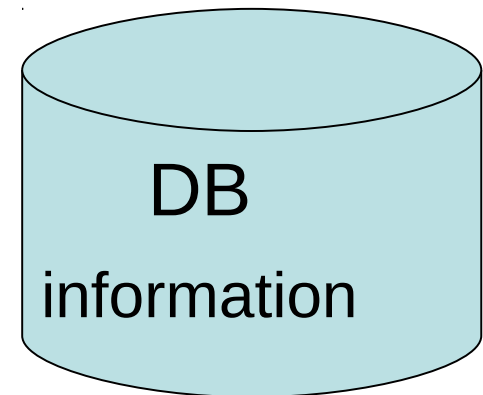
Functional interpretation

Experimental results
observed in the lab
(not always a wet-lab)

Recorded to:

- Test a hypothesis.
- Get a first insight of
a biological process.

To “interpret”
experimental results is to
use **current knowledge**
to rearrange them in a
meaningful way.



Already tested and
stored

Functional interpretation

1007_s_at

1053_at

117_at

121_at

1255_g

1294

1316

Experimental results
observed in the lab
(not always a wet-lab)

Recorded to:

- Test a hypothesis.
- Get a first insight of
a biological process.

12.4

11.5

10.3

10.2

9.9

9.3

8.2

7.7

7.4

6.5

6.2

5.9

5.8

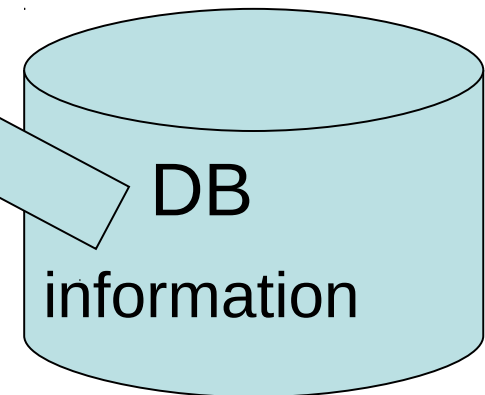
4.8

4.7

.

.

To “interpret”
experimental results is to
use **current knowledge**
to rearrange them in a
meaningful way.



Already tested and
stored

Babelomics Databases

Some of the biological databases contains **Functional Information** of the genes and sequences



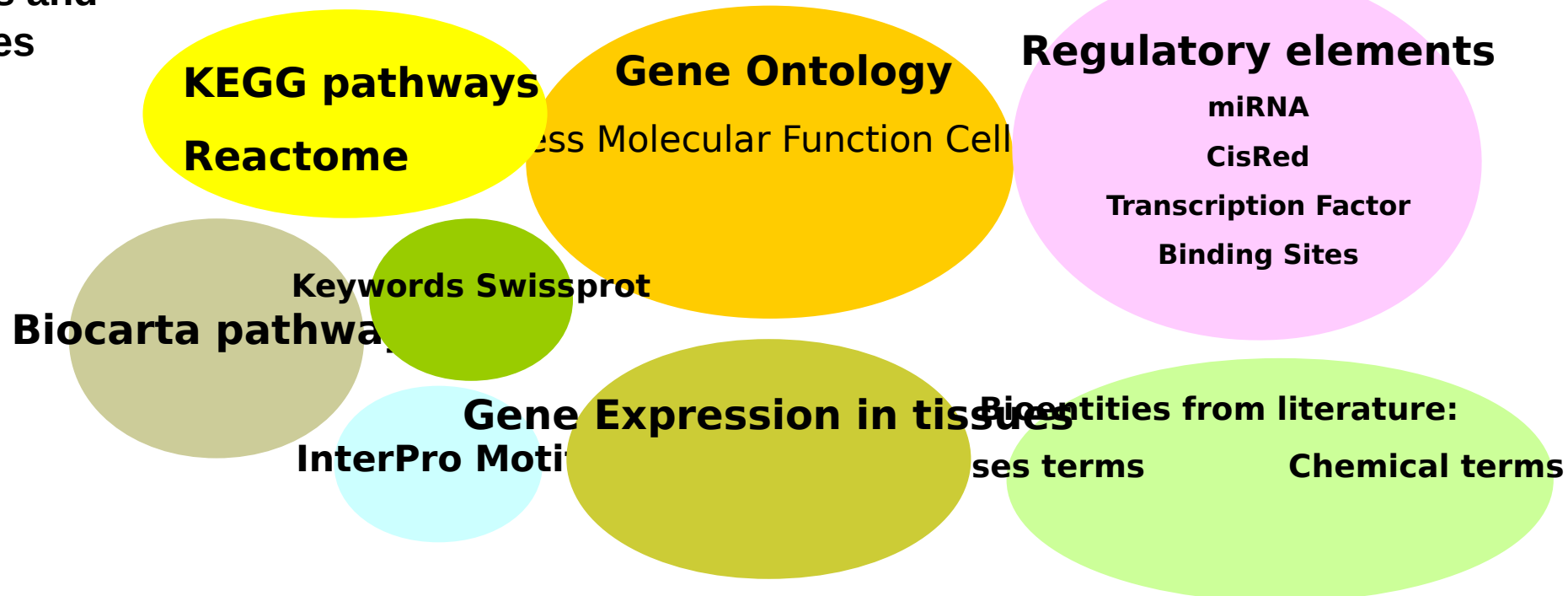
UniProt/Swiss-Prot
UniProtKB/TrEMBL
Ensembl IDs

EntrezGene
Affymetrix
Agilent

HGNC symbol
EMBL acc
RefSeq
PDB
Protein Id
IPI....

Gene IDs

Functional databases



Babelomics Databases

Which type of functional information?

The screenshot shows the Babelomics Databases interface. At the top, there is a 'Databases' section with a dropdown menu for 'Organism' set to 'Human (homo sapiens)'. Below this is a list of databases with checkboxes and links to 'options' for each:

- ☐ GO biological process [\[options\]](#)
- ☐ GO molecular function [\[options\]](#)
- ☐ GO cellular component [\[options\]](#)
- ☐ GOSlim GOA [\[options\]](#)
- ☐ Interpro [\[options\]](#)
- ☐ KEGG pathways [\[options\]](#)
- ☐ Reactome [\[options\]](#)
- ☐ Biocarta [\[options\]](#)
- ☐ miRNA targets [\[options\]](#)
- ☐ Jaspar TFBS [\[options\]](#)
- ☐ ORegAnno [\[options\]](#)

Below the list, there is a checkbox for 'Your annotations'. To its right is a 'browse server' button and the text 'no data selected.'. Below this is a link 'Or go to Upload Data form: [Upload \[annotation\]](#)'. A large curly bracket on the right side of the database list points to a text box. Another curly bracket on the right side of the 'Your annotations' section points to another text box.

The 'Job' section at the bottom has two text input fields: 'job name:' and 'job description:'. At the very bottom is a 'Run' button.

Use one or more of the given databases

If it is not in the databases, use your annotations option.

Babelomics Databases

Databases

Organism: Human (homo sapiens)

☒ GO biological process [options]

☐ GO molecular function [options]

☐ GO cellular component [options]

First select an organism

OPTIONS:

Test all the GO or only annotated terms

Discard functions with too few or too many genes?

If you have an hypothesis, better test this first!!!!!!

GO biological process options

GO parameters

Select annotation through ontology levels

☐ Propagate annotation to upper levels

☐ Direct annotation

GO level must be among levels and

Filter terms by number of annotated ids in DB

Minimum (typically 5-20)

Maximum (typically 500-Inf)

Number of annotated ids is computed from

☒ Genome

☐ Your input ids

Filter terms by keywords

Keywords (e.g. metabolism cancer)

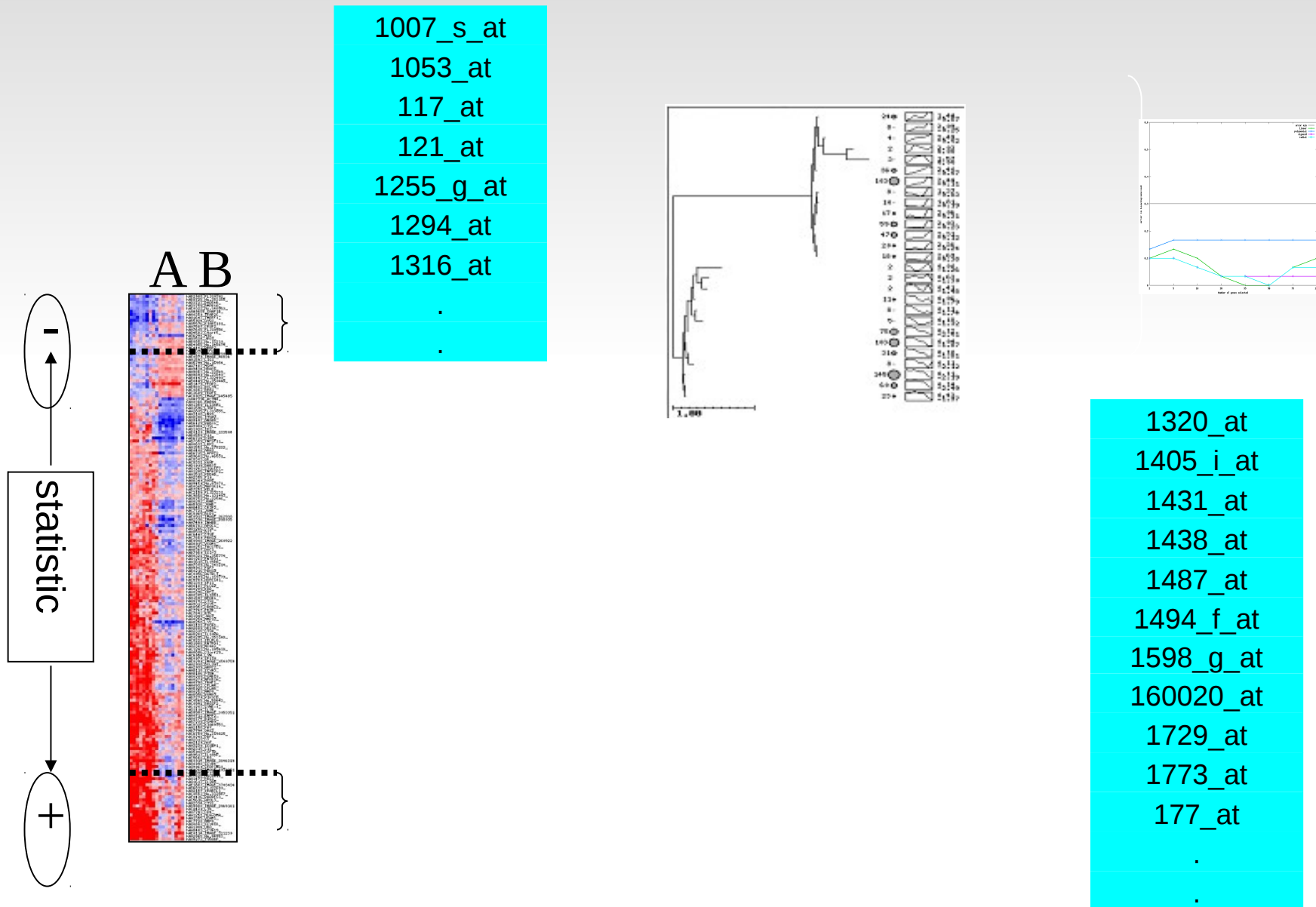
Your search must match

☒ all keywords

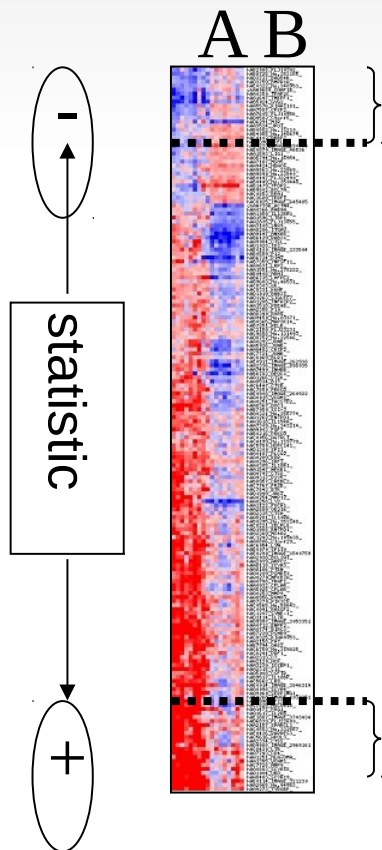
☐ any keyword

☒ Add children of selected terms

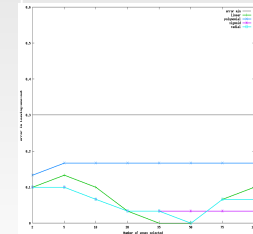
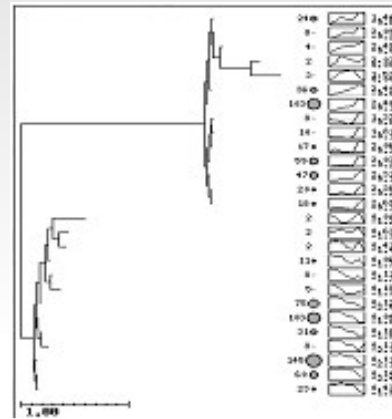
Simple enrichment analysis



Simple enrichment analysis



1007_s_at
1053_at
117_at
121_at
1255_g_at
1294_at
1316_at
.
.



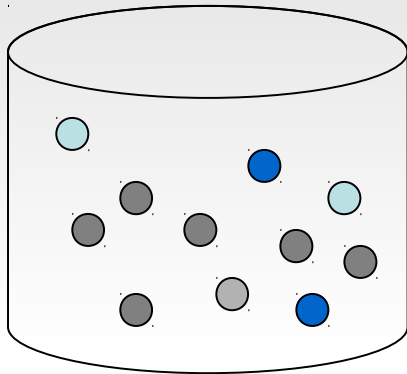
GO

4/7 ~ 2/11

1320_at
1405_i_at
1431_at
1438_at
1487_at
1494_f_at
1598_g_at
160020_at
1729_at
1773_at
177_at
.
.

FatiGO test

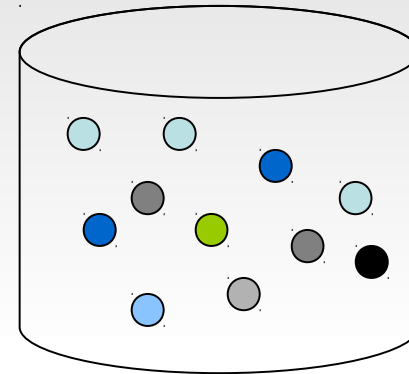
One Gene List (A)



Biosynthesis 60% ●

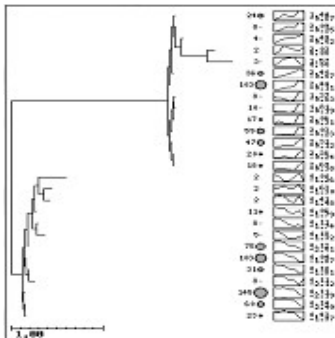
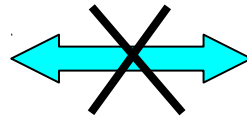
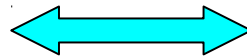
Sporulation 20% ●

The other list (B)



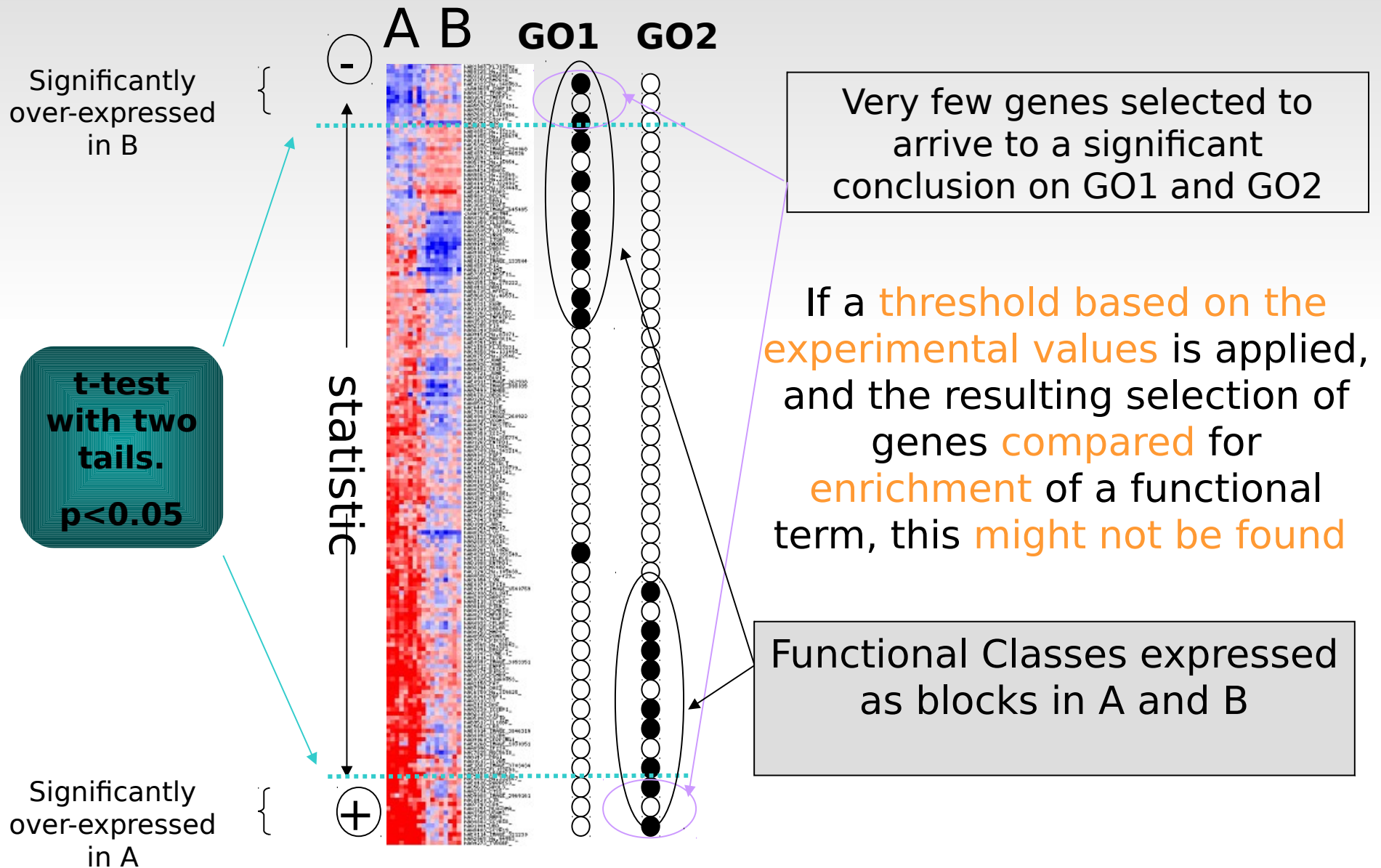
Biosynthesis 20% ●

Sporulation 20% ●



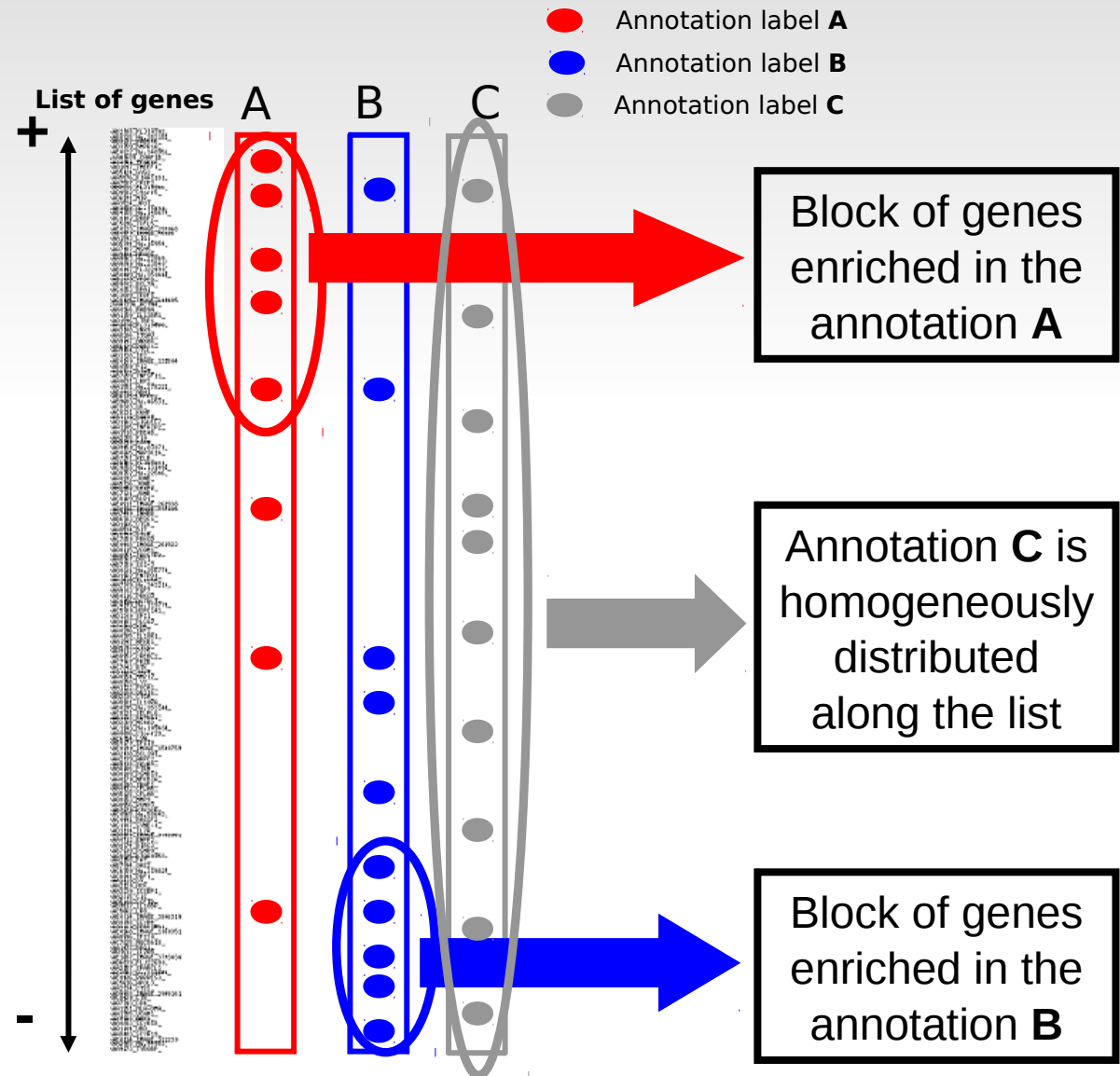
	A	B
Biosynthesis	6	2
No biosynthesis	4	8

FatiGO approach may not be very powerful

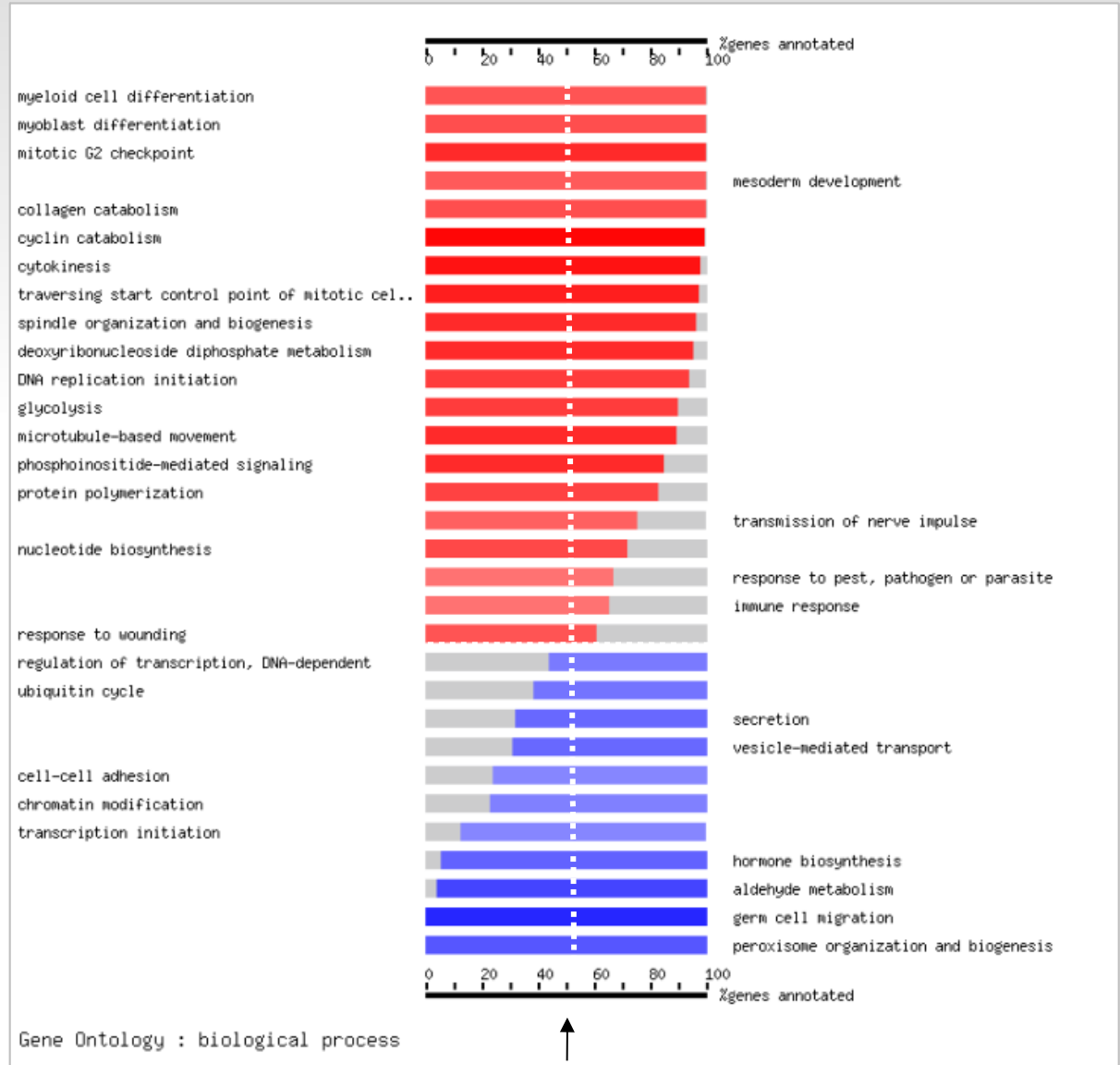
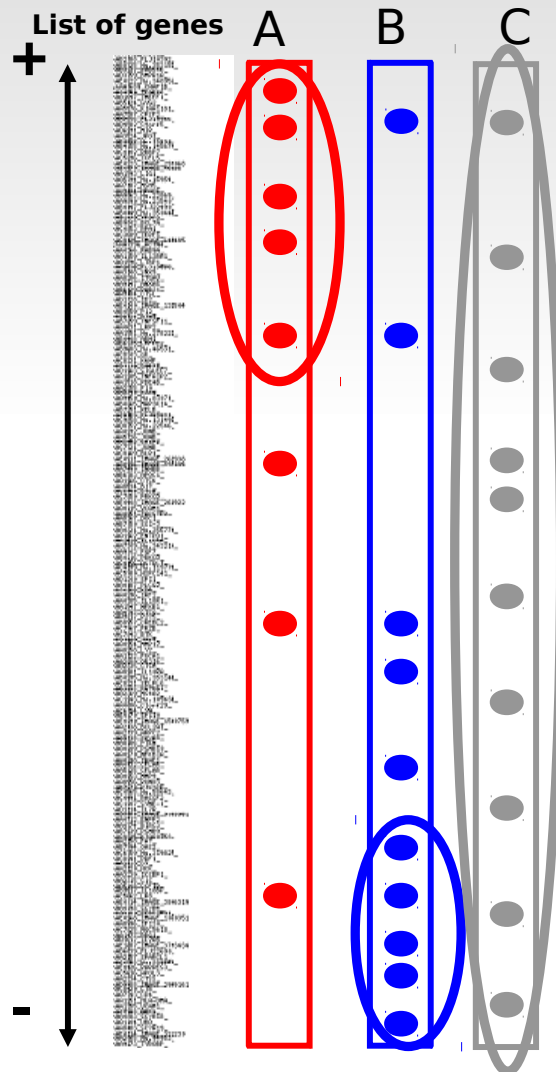


FatiScan, testing along an ordered list

- Index ranking genes according to some biological aspect under study.
- Database that stores gene class membership information.
- **FatiScan** searches over the whole ordered list, trying to find runs of functionally related genes.



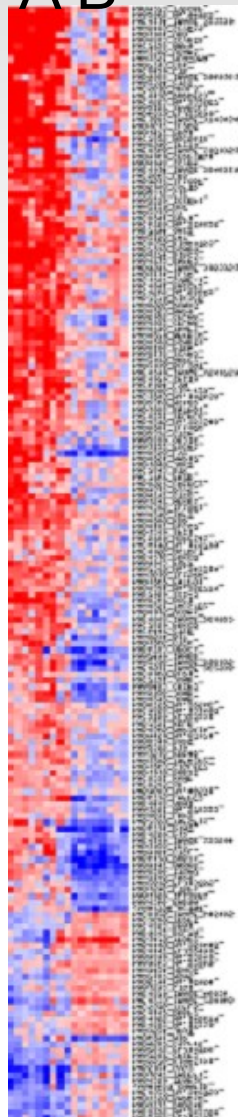
FatiScan results



Babelomics

Fatiscan results

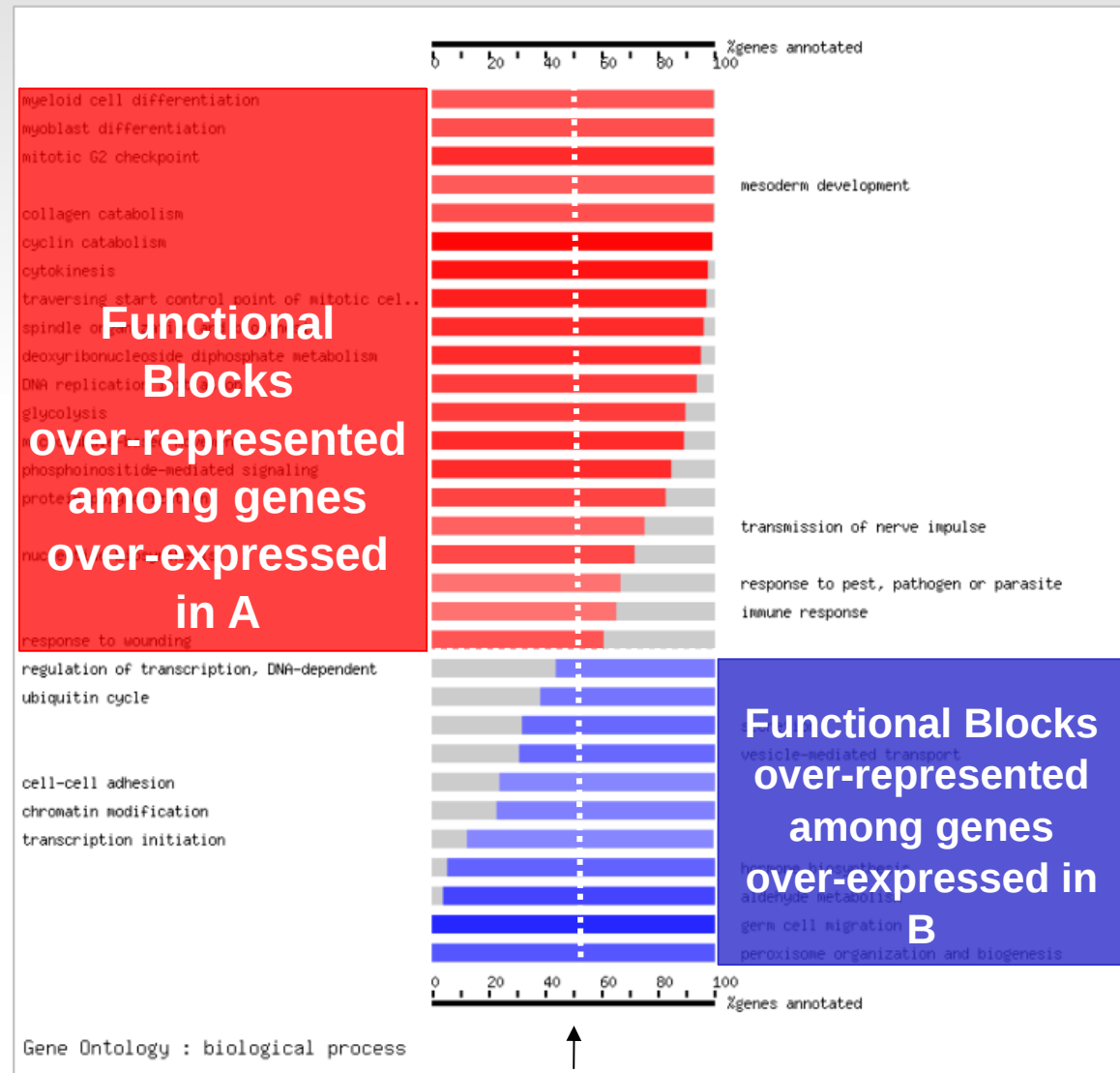
A B



Gene ranking index

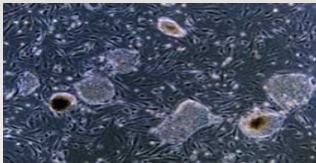
↑ (+)

↓ (-)

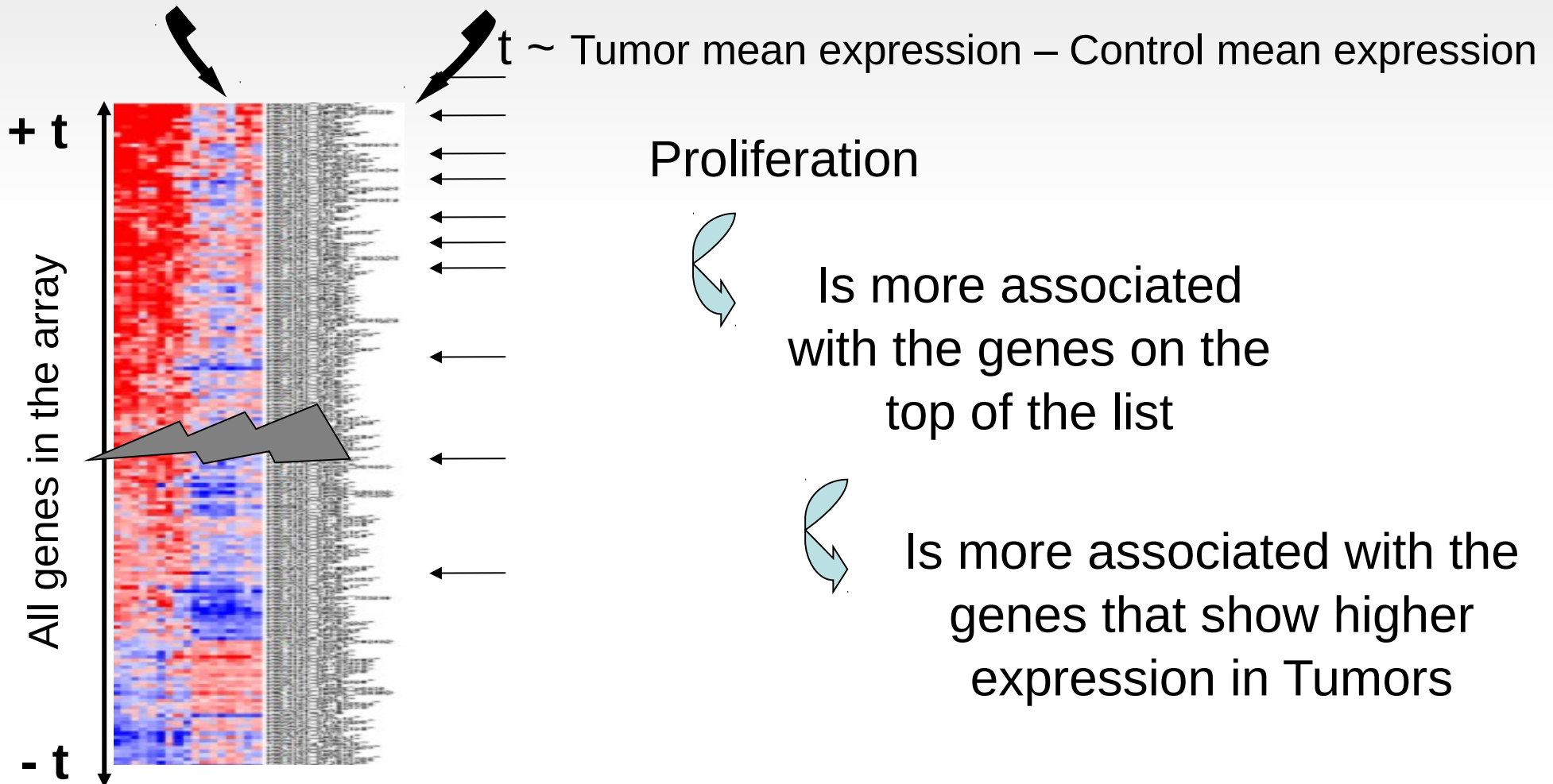
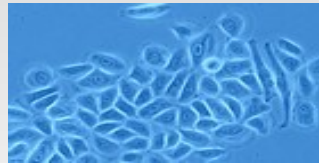


FatiScan Example - two classes

Tumor



Control



FatiScan Example - Survival Analysis

- Cromer et al. **Identification of genes associated with tumorigenesis and metastatic potential of hypopharyngeal cancer by microarray analysis.** *Oncogene* 2004, **23**(14) : 2484-2498.
- 34 hypopharyngeal cancer samples taken from patients undergoing surgery.
- Analyzed using Affymetrix HG-U95A microarrays (~12650 distinct transcription features).
- Disease free survival time after intervention was recorded

Cox proportional hazards model

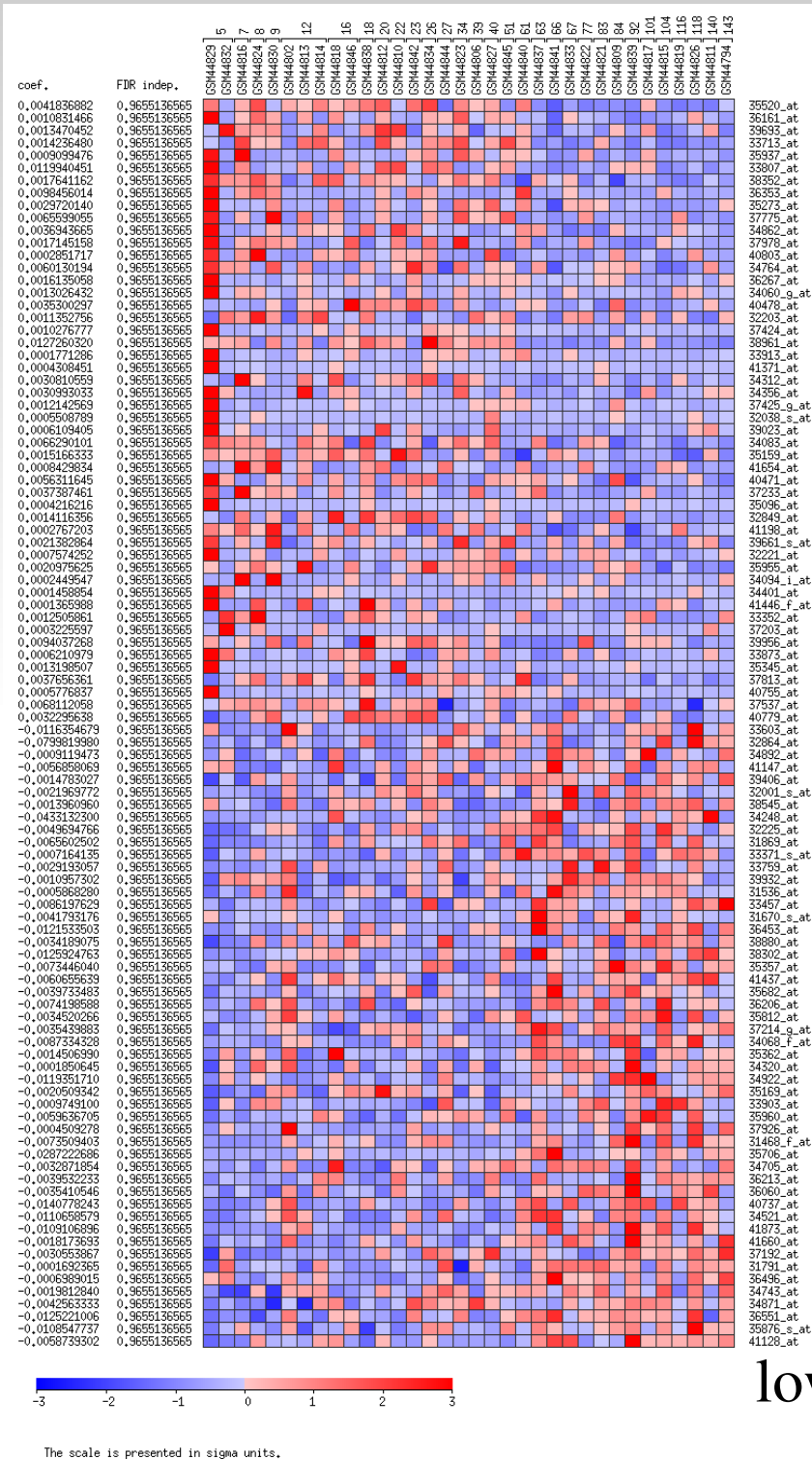
$$h(t) = h_0(t) * \exp(\beta * \text{gene expression})$$

Hazard
increased
with
expression

+ β

Hazard
decreased
with
expression

- β

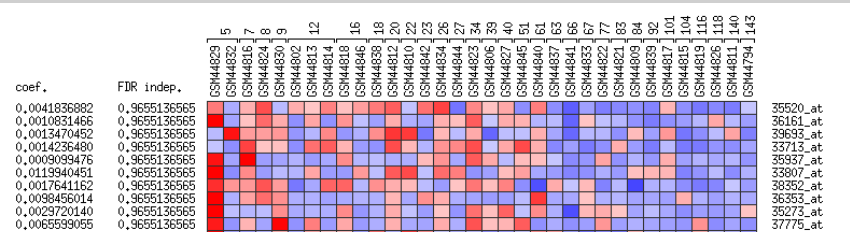


Gene Ontology:
biological process

lowest p-value = 0.96

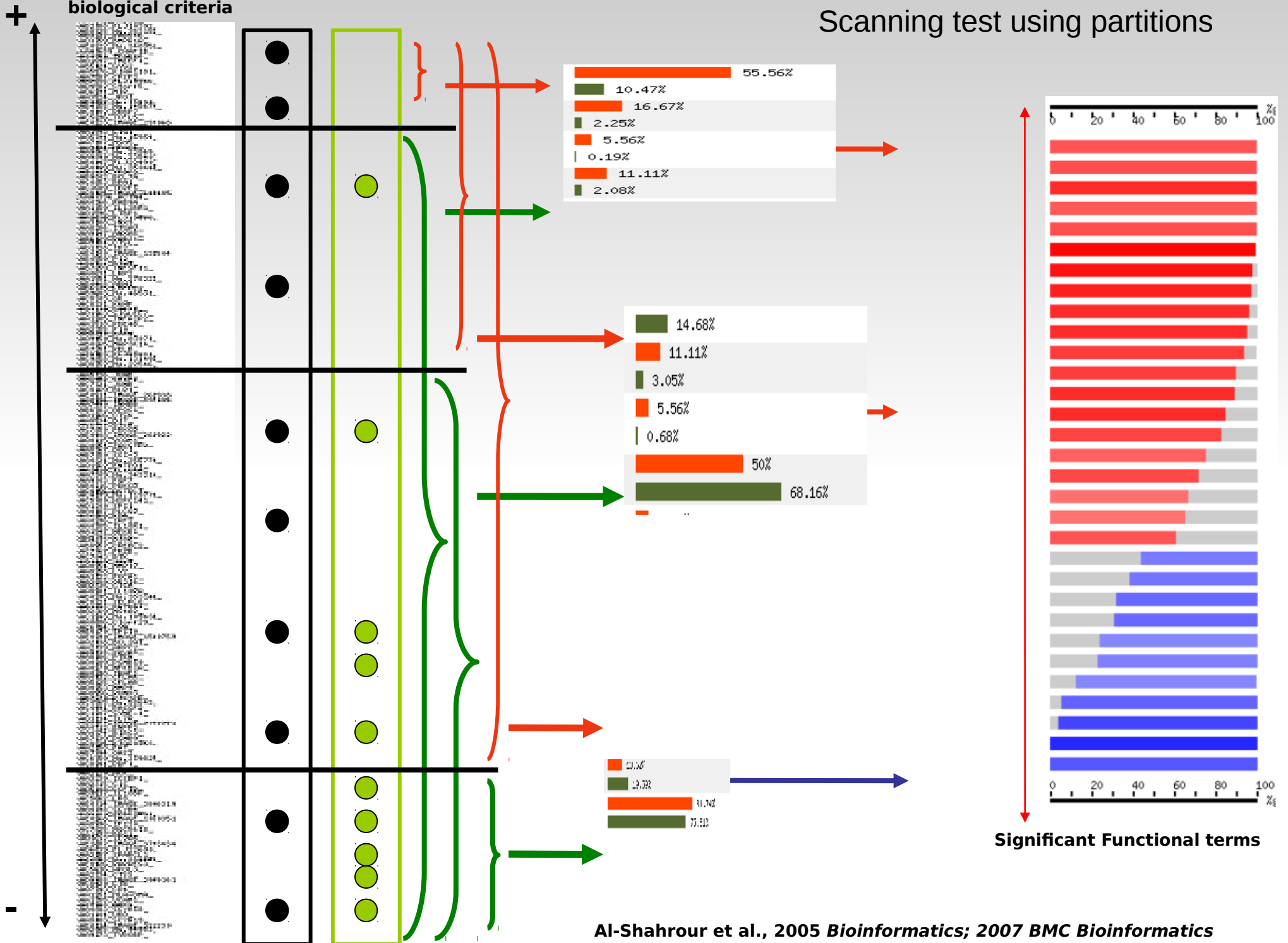
Hazard
increased
with
expression

$+\beta$



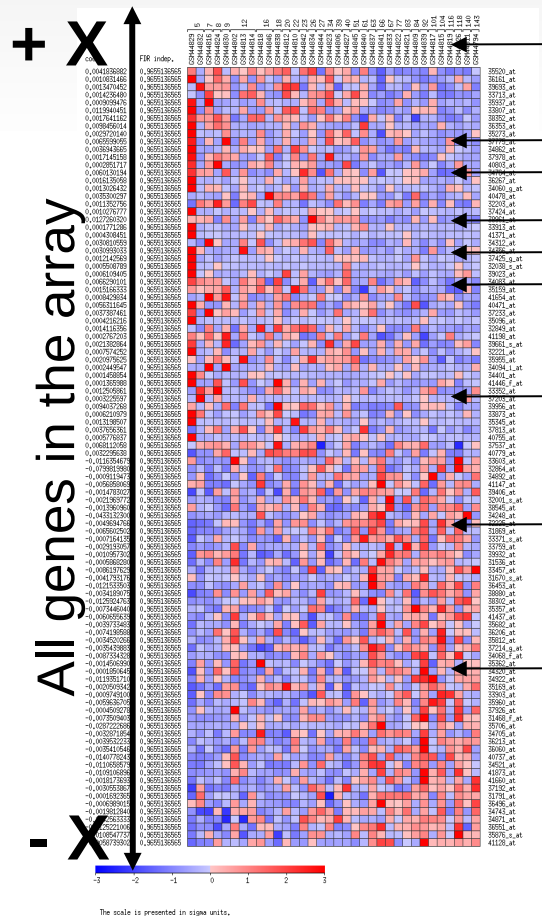
List of genes ranked by biological criteria

Scanning test using partitions



Logistic test

- Not using partitions
- But logistic regression model



$$\ln \left(\frac{P(g \in F)}{P(g \notin F)} \right) = K + \alpha X$$

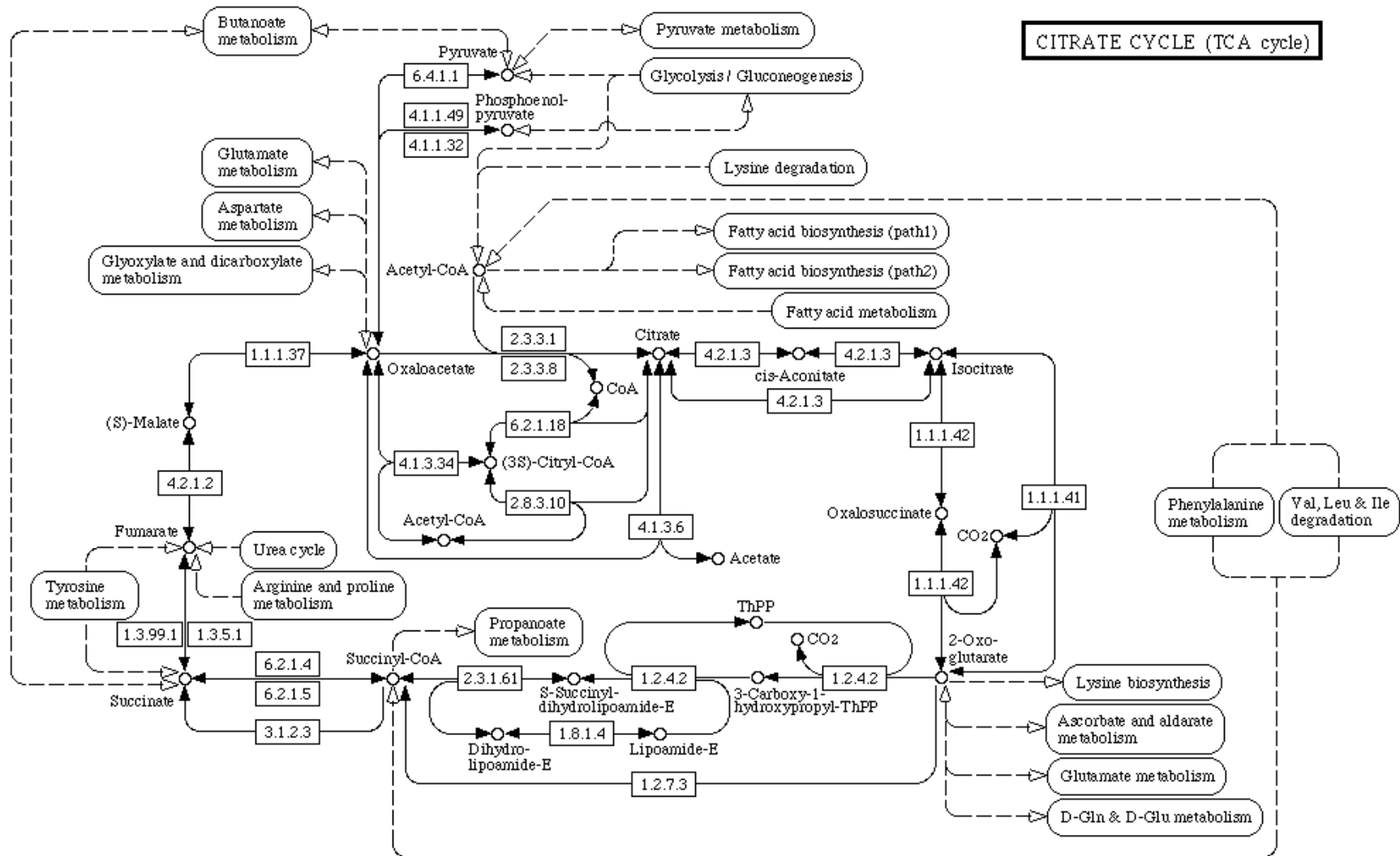
alpha > 0 : increasing X increases the probability of the gen to be annotated

alpha < 0 : decreasing X increases the probability of the gen to be annotated

Remarks

- The unit of information over which we test is shifted from genes to functional blocks (multiple testing again)
- We do one statistical test for each block
- All genes in the block are treated equally

Network modeling



Remarks

- The unit of information over which we test is shifted from genes to functional blocks (multiple testing again)
- We do one statistical test for each block
- All genes in the block are treated equally
- Annotation information is 0, 1

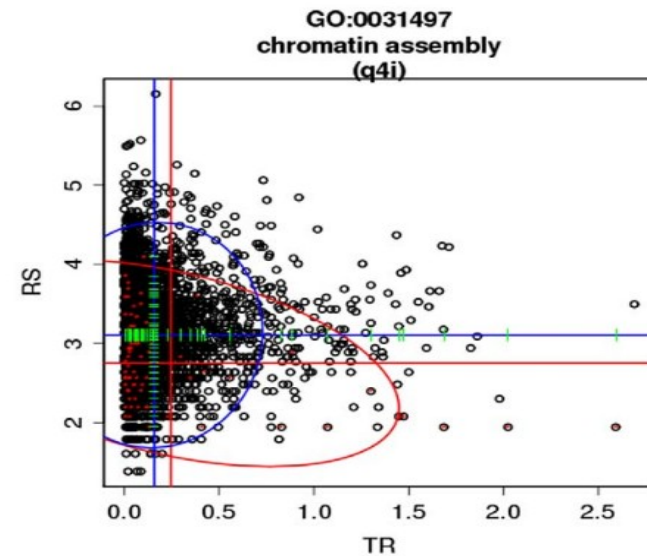
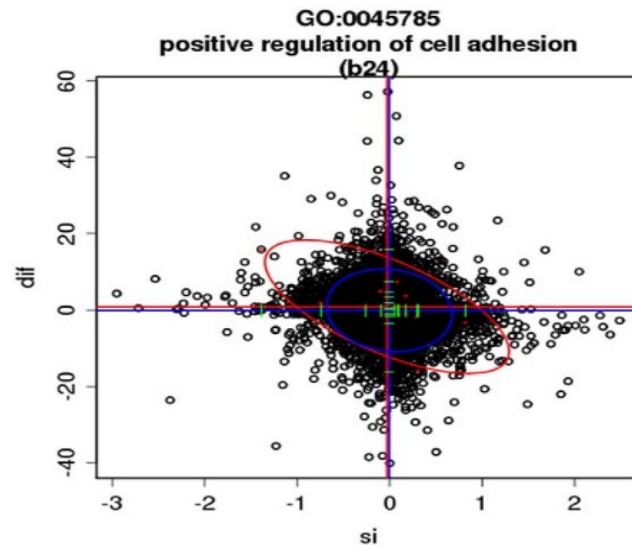
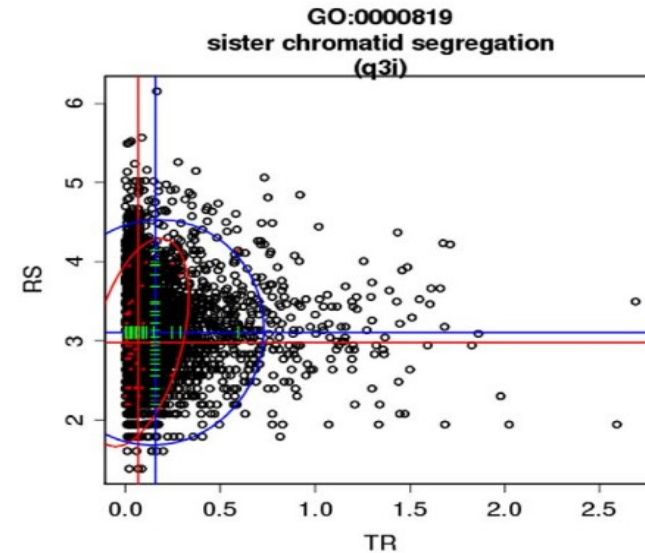
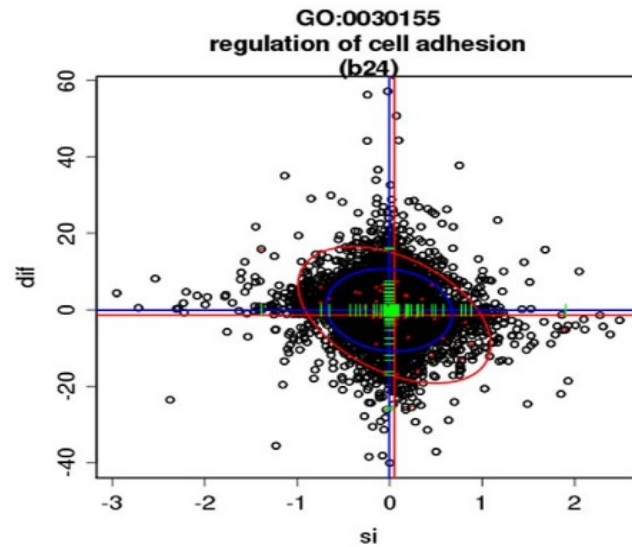
Remarks

- The unit of information over which we test is shifted from genes to functional blocks (multiple testing again)
- We do one statistical test for each block
- All genes in the block are treated equally
- Annotation information is 0, 1
- Genes independently may not show a strong pattern of association but the block coordinately does.

Remarks

- The unit of information over which we test is shifted from genes to functional blocks (multiple testing again)
- We do one statistical test for each block
- All genes in the block are treated equally
- Annotation information is 0, 1
- Genes independently may not show a strong pattern of association but the block coordinately does.
- Only ranking genes according to a unique condition

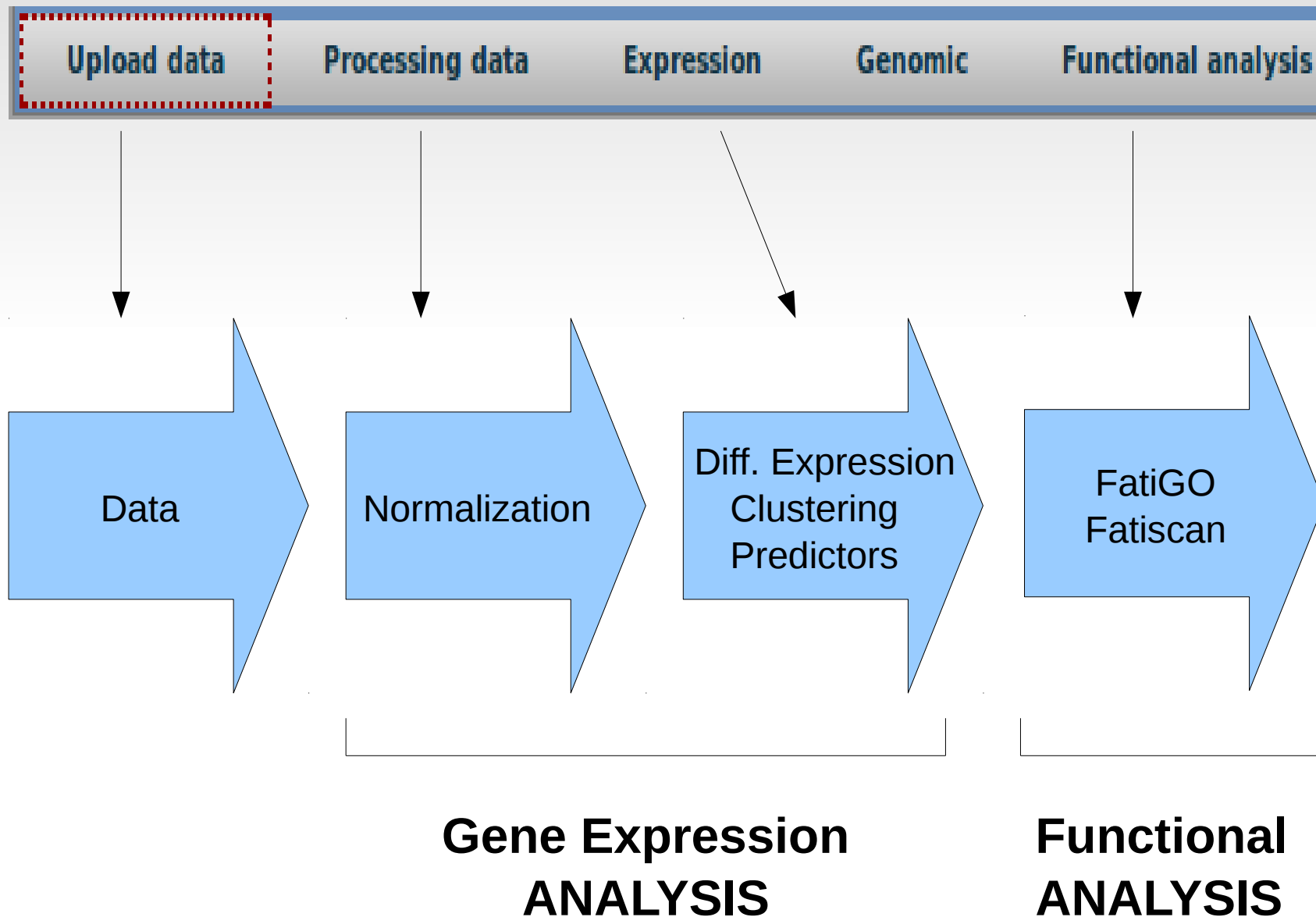
Multidimensional Analysis



Gene Set Methods - 2 general Approaches

- Competitive Hypothesis – Babelomics
 - Each functional block is compared to the remaining genes of the genome (or the second list).
 - Independent of the test used to derive the ranking.
- Self Contained Hypothesis – Goeman 2004
 - Checks that the block it self is differentially expressed, correlated with phenotype, associated to survival...
 - Has to be developed with the test creating the ranking of the genes.

Modules / Analysis progress



Thank You

www.babelomics.org