

BABELOMICS

Functional Enrichment Analysis: FatiGO & Fatiscan

Martina Marbà

mmarba@cipf.es

Valencia, 25th March 2011

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





BABELOMICS

A systems biology web resource for the functional interpretation of genome-scale experiments.

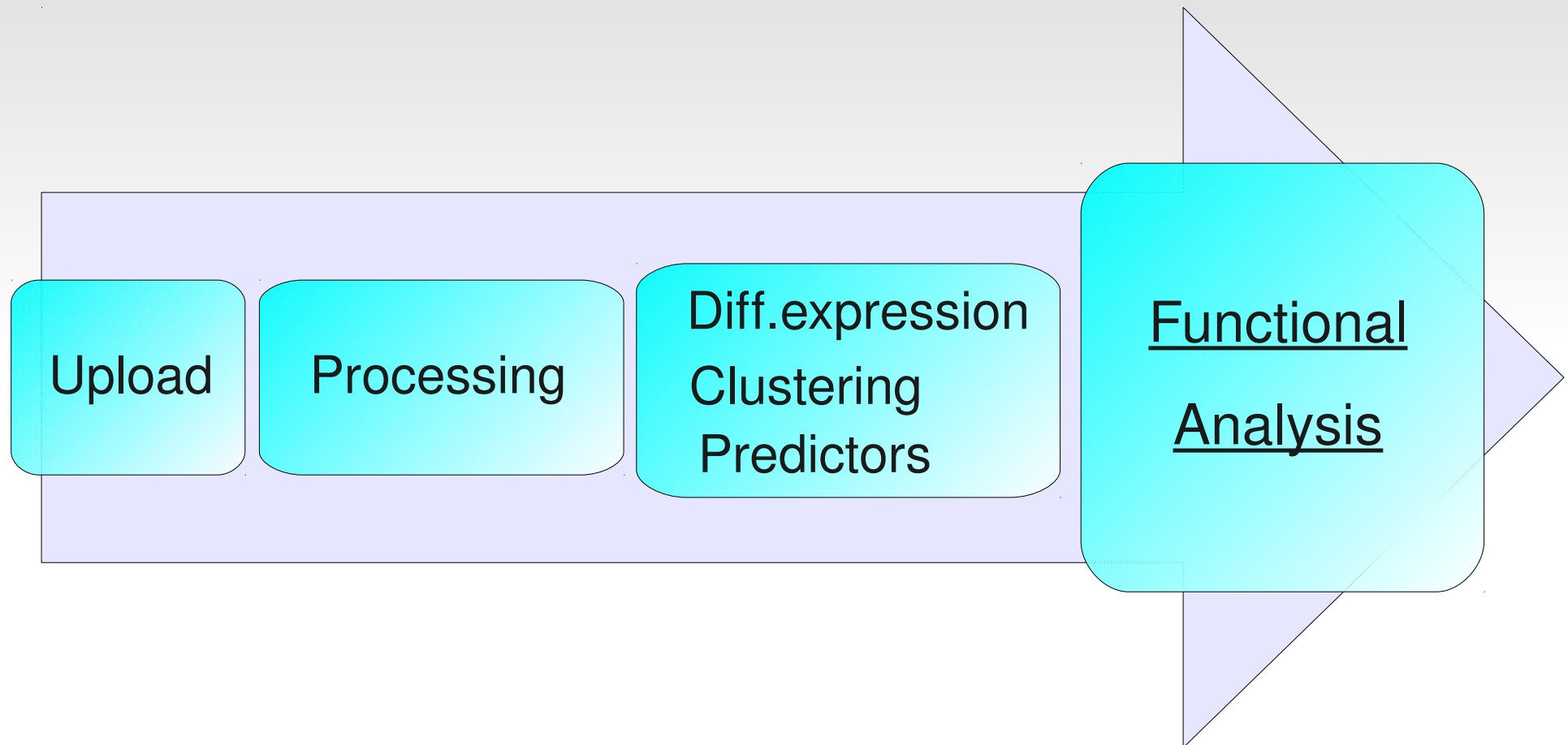
<http://www.babelomics.org>

Questions

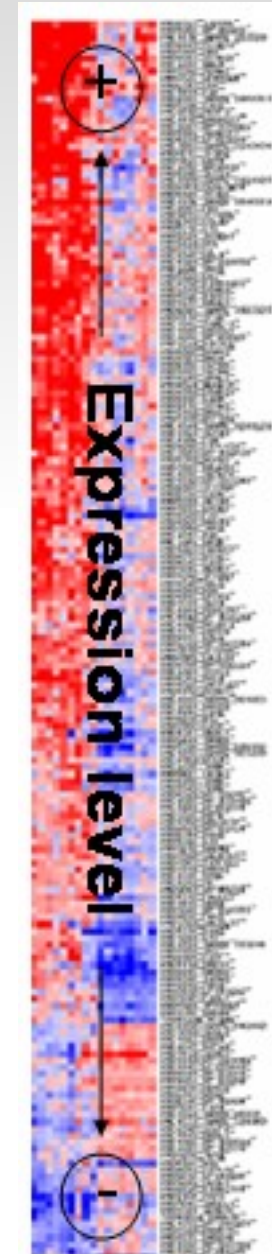
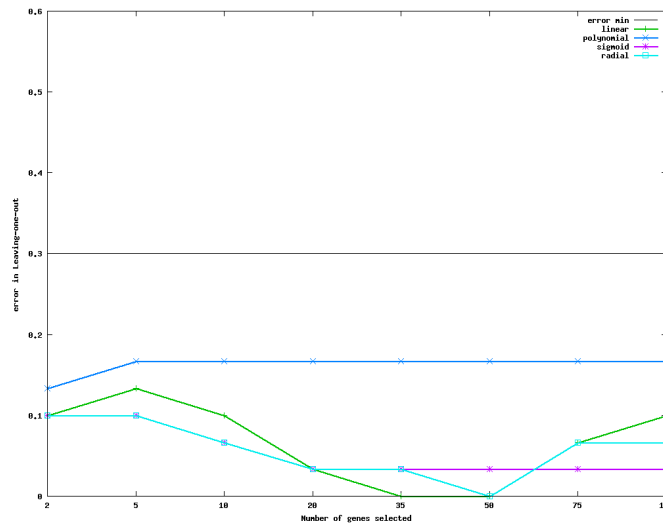
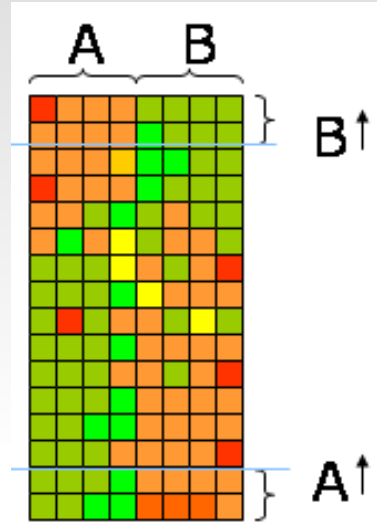
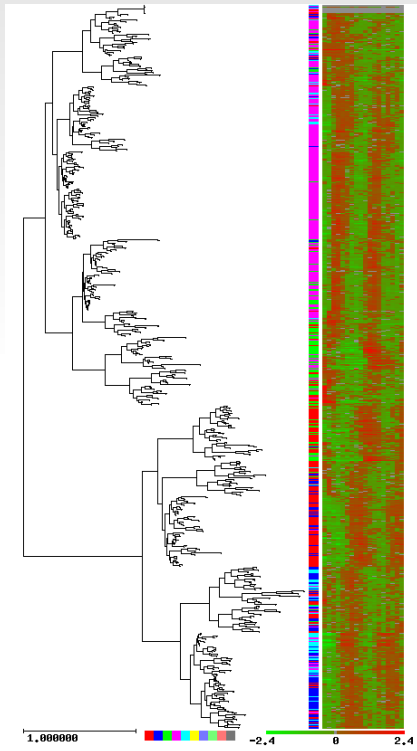
Questions that Functional enrichment analysis try to answer

- Is there any significant functional enrichment in my gene list?
- Are these genes involved in same pathways?
- What biological processes differentiate a healthy control from a diseased case?
- Do these genes share a specific microRNA regulation?
- Are they involved in the same disease?

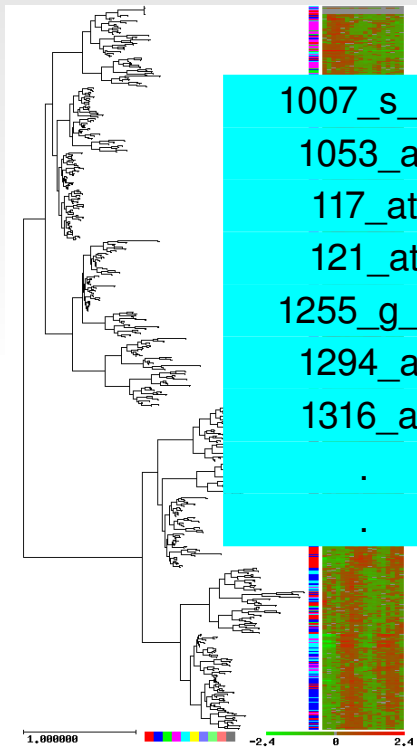
Data analysis workflow



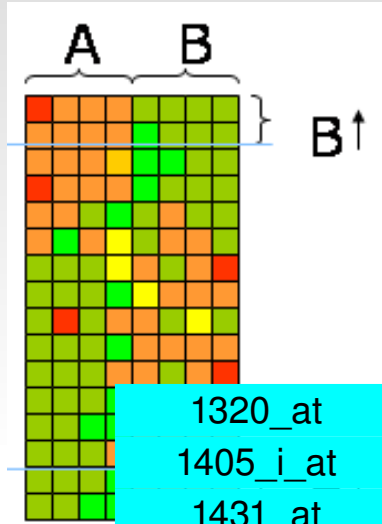
Genome-scale experiment output



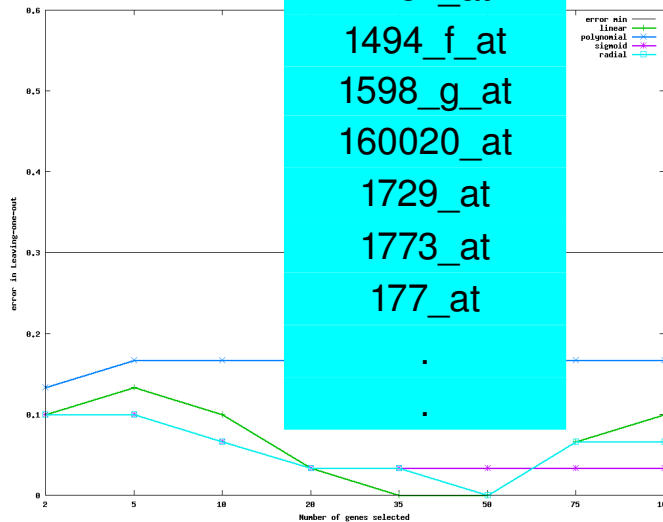
Genome-scale experiment output



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
 .
 .



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
.	.
.	.

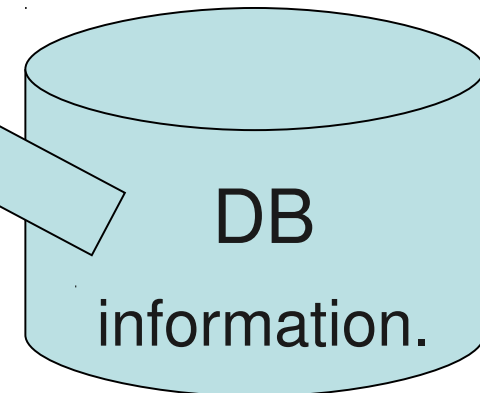
Functional interpretation

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

Recorded to:

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

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

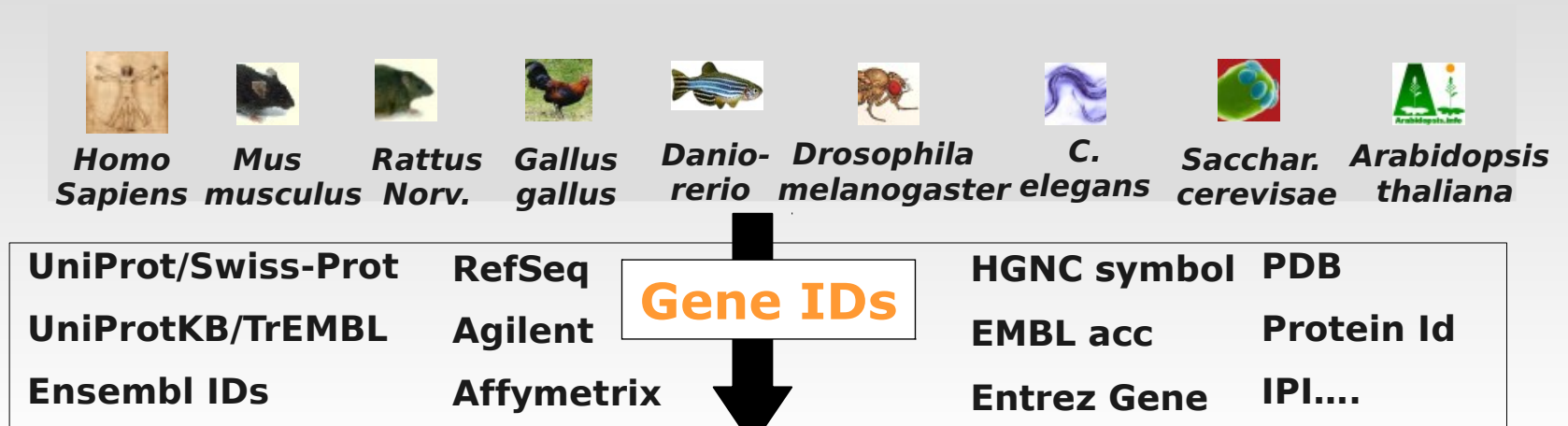


Already tested and
stored

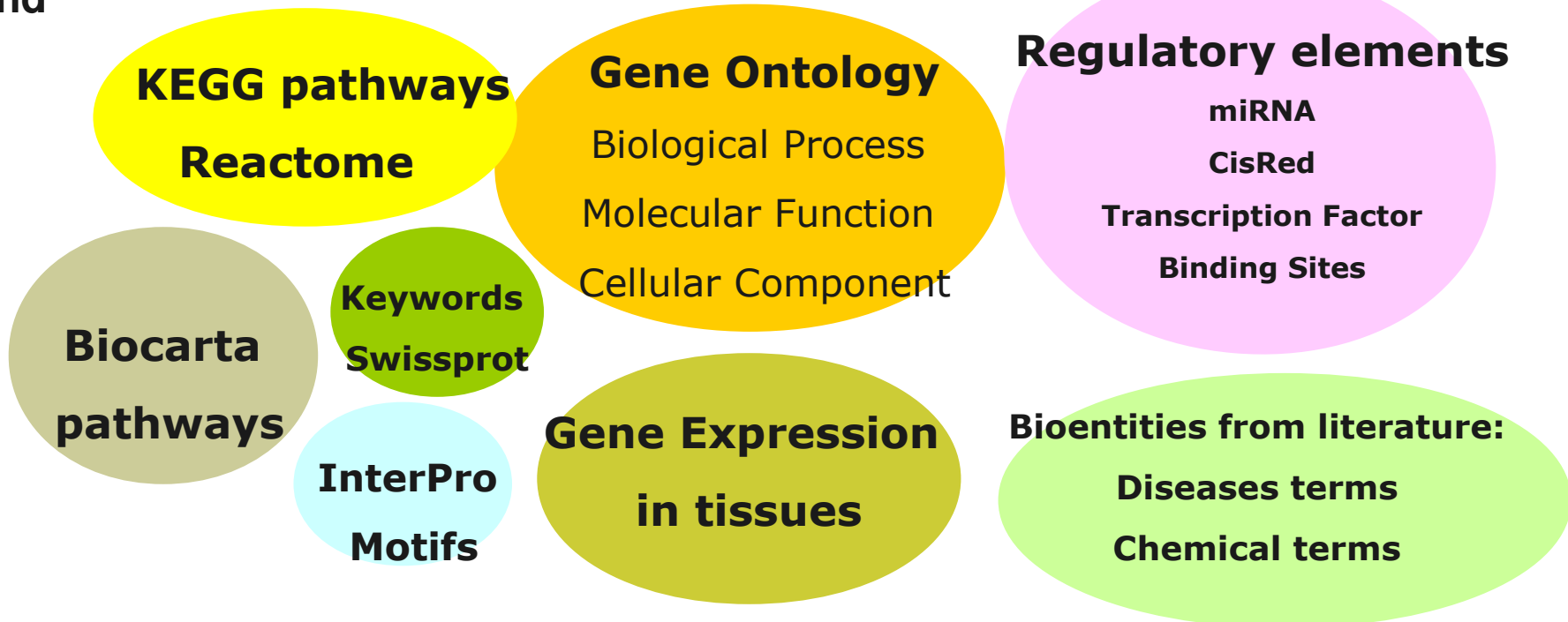
1320_at		12.4
1405_i_at		11.5
1431_at		10.3
1438_at		10.2
1487_at		9.9
1494_at		9.3
1598_at		8.2
1601_at		7.7
1701_at		7.4
1702_at		6.5
1703_at		6.2
1704_at		5.9
1705_at		5.8
1706_at		4.8
1707_at		4.7
1708_at		.
1709_at		.

Babelomics Databases

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

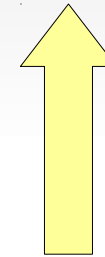
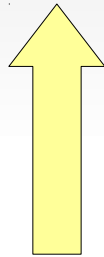


Functional databases



Functional interpretation FatiGO, Fatiscan

FatiGO and Fatiscan are web tools for:
statistical test, multiple test corrections, filtering ...



Lists of genes or ids, ie:
ranked differentially
expressed genes, a cluster
of genes, two particular
gene lists...

Integrated Biological DB
of Functional Annotation
(GO, GOSlimGOA, InterPro, KEGG,
Reactome, Biocarta, MiRNA targets,
Jaspar TFBS, ORegAnno)

FatiGO

- A web-based tool for the functional profiling of genome-scale experiments

The screenshot displays the Babelomics 4 web interface. At the top, the logo 'BABELOMICS 4' is shown with the subtitle 'gene expression and functional profiling analysis suite'. Below the logo is a navigation bar with tabs: 'Upload data', 'Processing data', 'Expression', 'Genomic', 'Functional analysis', and 'Utilities'. The 'Functional analysis' tab is circled in red. Below the navigation bar, a status bar indicates the user 'mmarba@cipf.es' is working on a project named 'Pre-processing Agilent' with 91.30 Mb of 1.00 Gb (8.92%) used and no active sessions. A green message box says 'Welcome to the new Babelomics 4, you can still use Babelomics 3 at: <http://babelomics3>'. The main content area is titled 'Functional analysis' and contains a list of tools under 'Single enrichment analysis'. The 'FatiGO' tool is circled in red, and a red arrow points from the 'Functional analysis' tab to it. Below 'FatiGO' is a description: 'Resource to show significant over-representation of GO terms.' Another tool, 'Marmite', is listed below with the description: 'Extracts blocks of related genes from an ordered list of genes by an associated value to the Marmite tool'.

BABELOMICS 4
gene expression and functional profiling analysis suite

Upload data Processing data Expression Genomic **Functional analysis** Utilities

mmarba@cipf.es working on project *Pre-processing Agilent* 91.30 Mb of 1.00 Gb (8.92%) no active

Welcome to the new Babelomics 4, you can still use Babelomics 3 at: <http://babelomics3>

Functional analysis

- Single enrichment analysis
 - **FatiGO**
Resource to show significant over-representation of GO terms.
 - Marmite
Extracts blocks of related genes from an ordered list of genes by an associated value to the Marmite tool

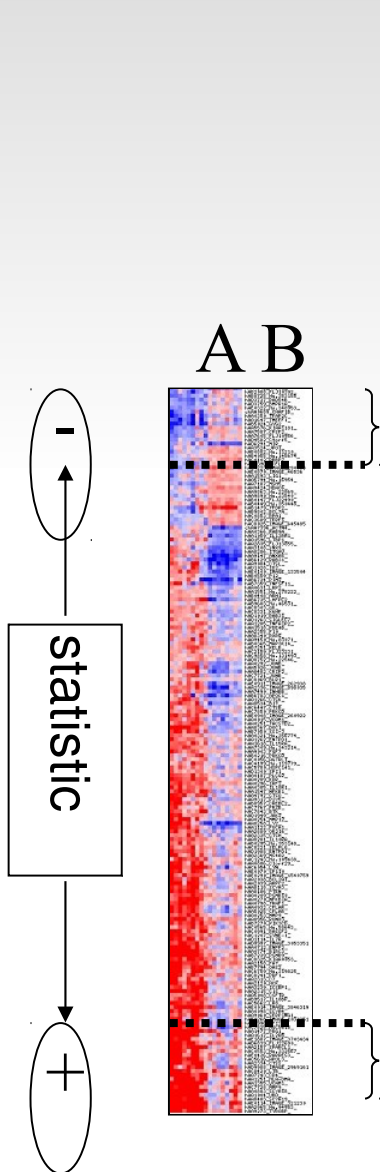
FatiGO features

- It allows us to compare functional annotation of:
 - **Two** list of genes
 - **One** list against the rest of genome
 - Lists of genes with user **submitted annotations**
- One statistical test for each Functional **Block** of annotation
 - Fisher's exact test
 - Multiple testing context (hundreds of annotation)
 - Filtering of annotation is convenient (the less tests the best correction)

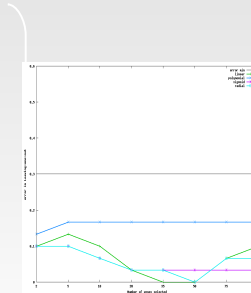
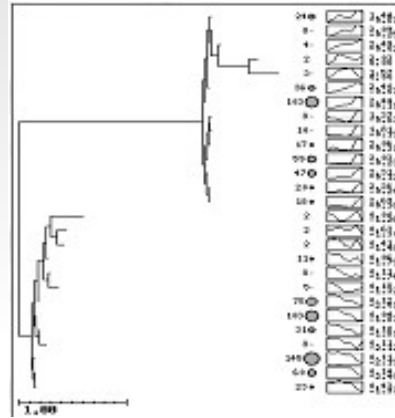
FatiGO features

- It allows us to compare functional annotation of:
 - **Two** list of genes
 - **One** list against the rest of genome
 - Lists of genes with user **submitted annotations**
- One statistical test for each Functional **Block** of annotation
 - Fisher's exact test
 - Multiple testing context (hundreds of annotation)
 - Filtering of annotation is convenient (the less tests the best correction)

Simple enrichment analysis

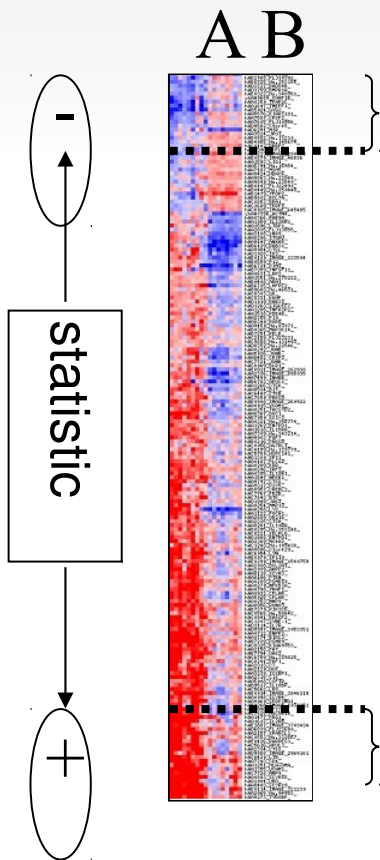


1007_s_at
 1053_at
 117_at
 121_at
 1255_g_at
 1294_at
 1316_at
 .
 .

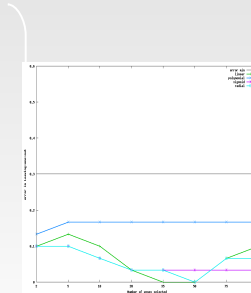
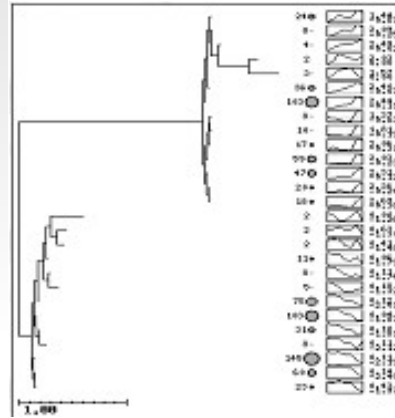


1320_at
 1405_i_at
 1431_at
 1438_at
 1487_at
 1494_f_at
 1598_g_at
 160020_at
 1729_at
 1773_at
 177_at
 .
 .

Simple enrichment analysis



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



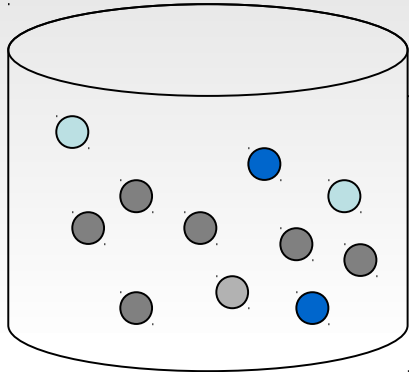
GO

$$4/7 \sim 2/11$$

- 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
- .
- .

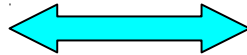
FatiGO test

One Gene List (A)

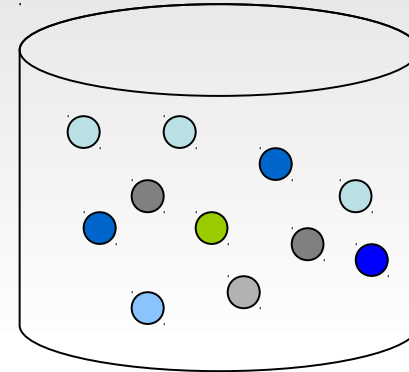


Biosynthesis 60% ●

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



The other list (B)



Biosynthesis 20% ●

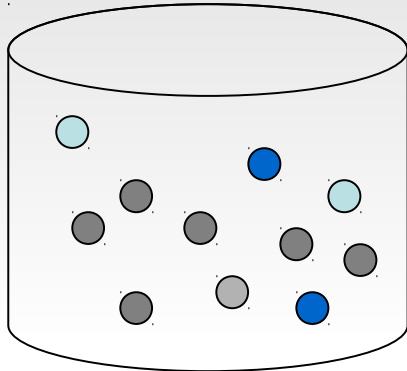
Genes in group A have significantly to do with **biosynthesis**.

	A	B
Biosynthesis	60	20
No biosynthesis	40	80

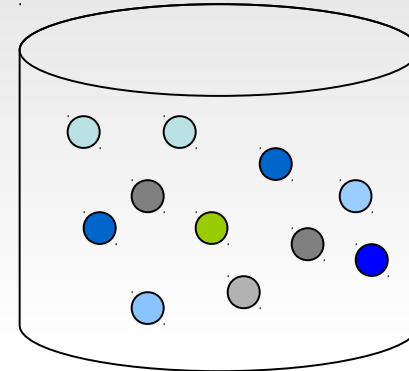
p.val = 5.318e-09

FatiGO test

One Gene List (A)



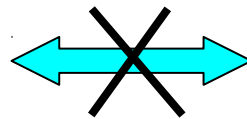
The other list (B)



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

Sporulation 20% ●

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



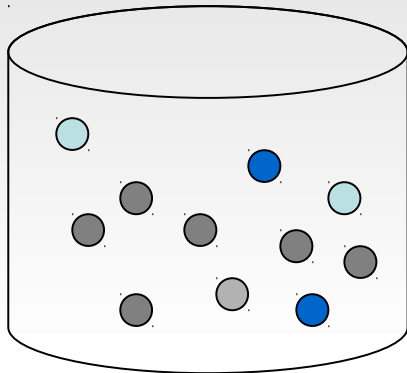
Sporulation 30% ●

	A	B
Sporulation	20	30
No sporulation	80	70

p.val = 0.964

FatiGO test

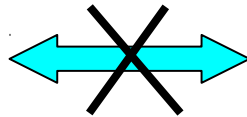
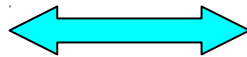
One Gene List (A)



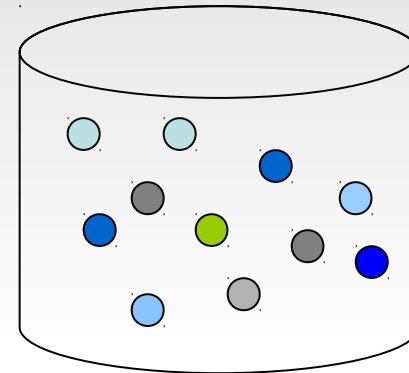
Biosynthesis 60% ●

Sporulation 20% ●

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



The other list (B)



Biosynthesis 20% ●

Sporulation 30% ●

	A	B
Biosynthesis	60	20
No biosynthesis	40	80

	A	B
Sporulation	20	30
No sporulation	80	70

We do this for each term (GO, miRNA, Interpro , ...)!!!

FatiGO form

Define your comparison

- Id list vs Id list
- Id List vs Rest of genome
- Id List vs Rest of ids contained in your annotations (complementary list)

Select your data

List 1 :

no data selected.
Or go to Upload Data form: [Upload \[idlist\]](#)

List 2 :

no data selected.
Or go to Upload Data form: [Upload \[idlist\]](#)

Options

Fisher exact test

Remove duplicates?

Databases

Organism

GO biological process [\[options\]](#)

FatiGO form

Do you want to compare 2 conditions or one vs the rest of genome ?

Define your comparison

Id list vs Id list
 Id List vs Rest of genome
 Id List vs Rest of ids contained in your annotations (c

Select your data

List 1 : no data selected.
Or go to Upload Data form: [Upload \[idlist\]](#)

List 2 : no data selected.
Or go to Upload Data form: [Upload \[idlist\]](#)

Options

Fisher exact test
Remove duplicates?

Databases

Organism
 GO biological process [\[options\]](#)

eg. Compare 2 tissues or responder genes vs. non-responders

eg: genes that respond to one treatment against the genome

FatiGO form

Define your comparison

- Id list vs Id list
- Id List vs Rest of genome
- Id List vs Rest of ids contained in your annotations (c

Upload first at Data Upload

Select your data

List 1 : no data selected.
Or go to Upload Data form: [Upload \[idlist\]](#)

List 2 : no data selected.
Or go to Upload Data form: [Upload \[idlist\]](#)

Data
selection

“txt” file with
gene lists:

```
gene1  
gene2  
gene3  
...
```

Options

Fisher exact test

Remove duplicates?

Databases

Organism

GO biological process [\[options\]](#)

FatiGO form

Define your comparison

- Id list vs Id list
- Id List vs Rest of genome
- Id List vs Rest of ids contained in your annotations (c

Select your data

List 1 :

no data selected.
Or go to Upload Data form: [Upload \[idlist\]](#)

List 2 :

no data selected.
Or go to Upload Data form: [Upload \[idlist\]](#)

Options

Fisher exact test

Remove duplicates?

Databases

Organism

GO biological process [\[options\]](#)

HINT:

- *Two tailed* for 2 lists
- *One tailed* for 1list vs rest of genome (or your annotations)

Algorithm options

Removing duplicates:

- To choose one or other option depends on where gene lists come from.

FatiGO form

Define your comparison

- Id list vs Id list
- Id List vs Rest of genome
- Id List vs Rest of ids contained in your annotations (c

Select your data

List 1 :

no data selected.

Or go to Upload Data form: [Upload \[idlist\]](#)

List 2 :

no data selected.

Or go to Upload Data form: [Upload \[idlist\]](#)

Options

Fisher exact test

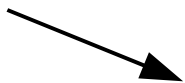
Remove duplicates?

Databases

Organism

GO biological process [\[options\]](#)

Which type of functional information?



?

FatiGO form

Which type of functional information?

Databases

Organism ▼

- 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 [\[options\]](#)

Your annotations no data selected.
Or go to Upload Data form: [Upload \[annotation\]](#)

Use one or more of the given databases

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

FatiGO form

A.

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

FatiGO form

B.

Which type of functional information?

Databases

Organism

- 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 [\[options\]](#)
- Your annotations no data selected.
Or go to Upload Data form: [Upload \[annotation\]](#)

Job

job name:

job description:

Your annotations: useful when you work with your own annotations OR with an organism that is not in Babelomics

Upload first at Data Upload

Example (your annotations):

38969_at	GO:0003677
37639_at	GO:0006306
37149_s_at	GO:0004674
37149_s_at	GO:0005525
37639_at	GO:0006306
37149_s_at	GO:0004674
...	...

FatiGO form

Databases

Organism

- 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 [\[options\]](#)
- Jaspas TFBS [\[options\]](#)
- ORegAnno [\[options\]](#)

Your annotations no data selected.
Or go to Upload Data form: [Upload \[annoti](#)

What's
your
job name?

Job

job name:

job description:

Set up a job name
and optionally,
give a description.

FatiGO results

Summary results:

■ Id annotations per DB :

<i>DB</i>	<i>List1</i>	<i>Genome</i>
GO biological process (levels from 3 to 9)	350 of 500 (70%) 11.26 annotations/id	11716 of 23198 (2343.2%) 5.08 annotations/id
GO molecular function (levels from 3 to 9)	344 of 500 (68.8%) 3.05 annotations/id	11370 of 23198 (2274%) 1.92 annotations/id

Tables significant terms:

▼ Significant Results

■ Number of significant terms per DB :

<i>DB</i>	<i>Number of significant terms</i>
GO biological process (levels from 3 to 9)	142
GO molecular function (levels from 3 to 9)	30

FatiGO results

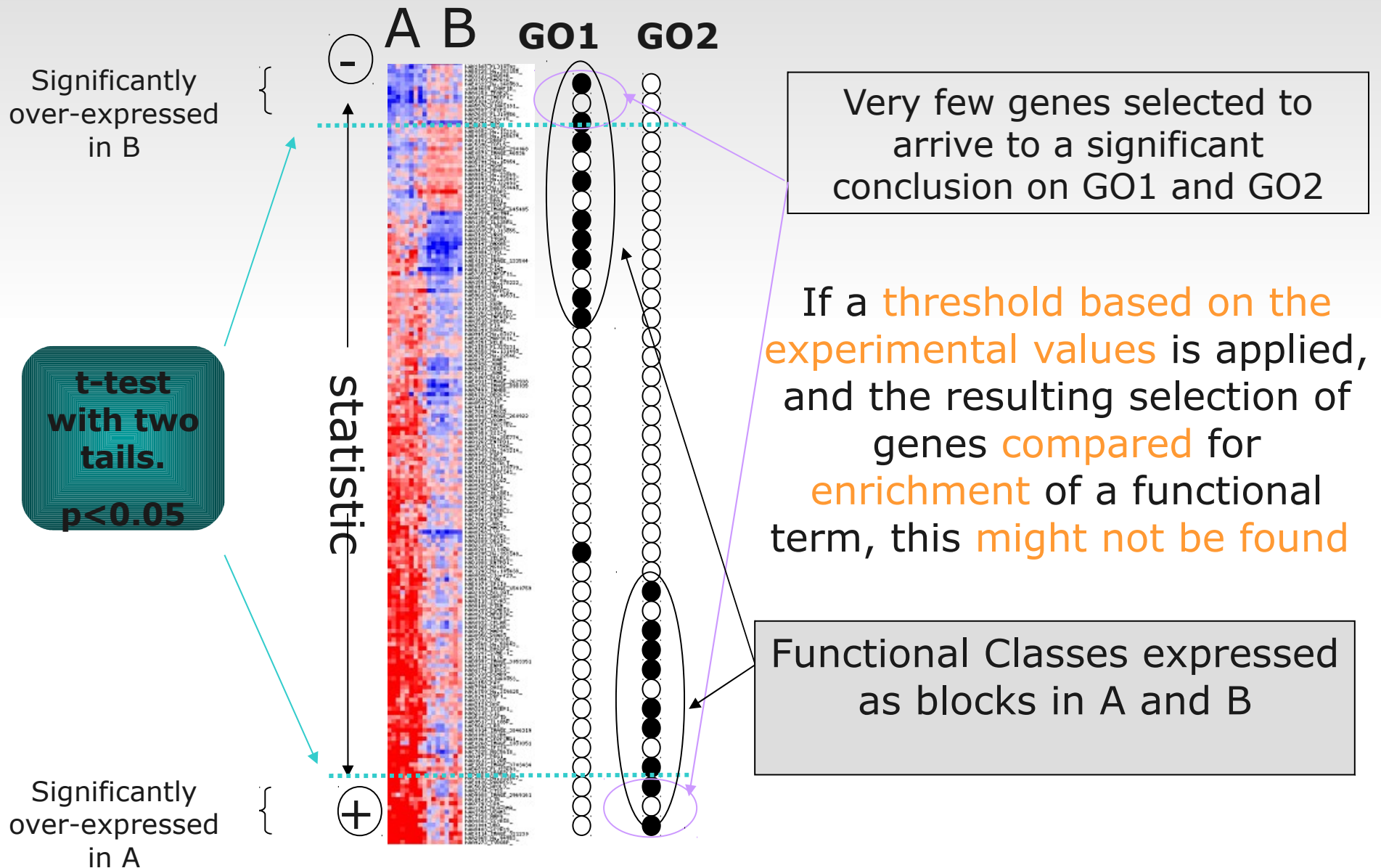
Significant results:

Term	Term size	Term size (in genome)	Term annotation % per list	Annotated ids	Odds ratio (log _e)	pvalue	Adjusted pvalue ▲
negative regulation of apoptosis (GO:0043066)	412	403	list 1: 7.2% list 2: 1.62%	list 1: 205225_at,20979... list 2: ENSG00000001084,ENSG...	1.5495	7.006e-13	7.65e-10
negative regulation of programmed cell death (GO:0043069)	418	409	list 1: 7.2% list 2: 1.65%	list 1: 205225_at,20979... list 2: ENSG00000001084,ENSG...	1.5334	1.074e-12	7.65e-10
cellular amino acid derivative metabolic process (GO:0006575)	182	173	list 1: 4.8% list 2: 0.68%	list 1: 209604_s_at,209... list 2: ENSG00000001084,ENSG...	1.995	9.24e-13	7.65e-10
cellular amino acid and derivative metabolic process (GO:0006519)	447	447	list 1: 7.4% list 2: 1.77%	list 1: 209604_s_at,209... list 2: ENSG00000001084,ENSG...	1.491	1.7e-12	9.082e-10

↑
Enriched class

↑
Annotated genes per GO from each list

FatiGO approach may not be very powerful



FatiScan

Upload data Processing data Expression Genomic **Functional analysis**

Functional analysis

- Single enrichment analysis
 - **FatiGO**
Provides significant
 - **Marmite**
Single enrichment a
- Set enrichment analysis
 - **Gene set analysis**
Finds gene-sets with gene-set analysis
 - **MarmiteScan**
Implements gene-se
 - **GeSBAP**
Gene set analysis (a with SNPs or CNVs)

Gene set analysis

► [Online examples \(test the form with example data\)](#)

Select your ranked list

no data selected.

Or go to Upload Data form: [Upload \[idlist:ranked\]](#)

Options

Logistic model

Fatiscan

Fisher exact test

Remove duplicates?

Databases

Organism

GO biological process [\[options\]](#)

GO molecular function [\[options\]](#)

GO cellular component [\[options\]](#)

Fatiscan features

- Interpret a **ranked list of genes**.
- There is not need for choosing a cut-off. All information is included.
- One statistical test for each Functional **Block** of annotation
 - Fisher's exact test
 - Multiple testing context (hundreds of annotation)
 - Filtering of annotation is convenient (the less tests the best correction)

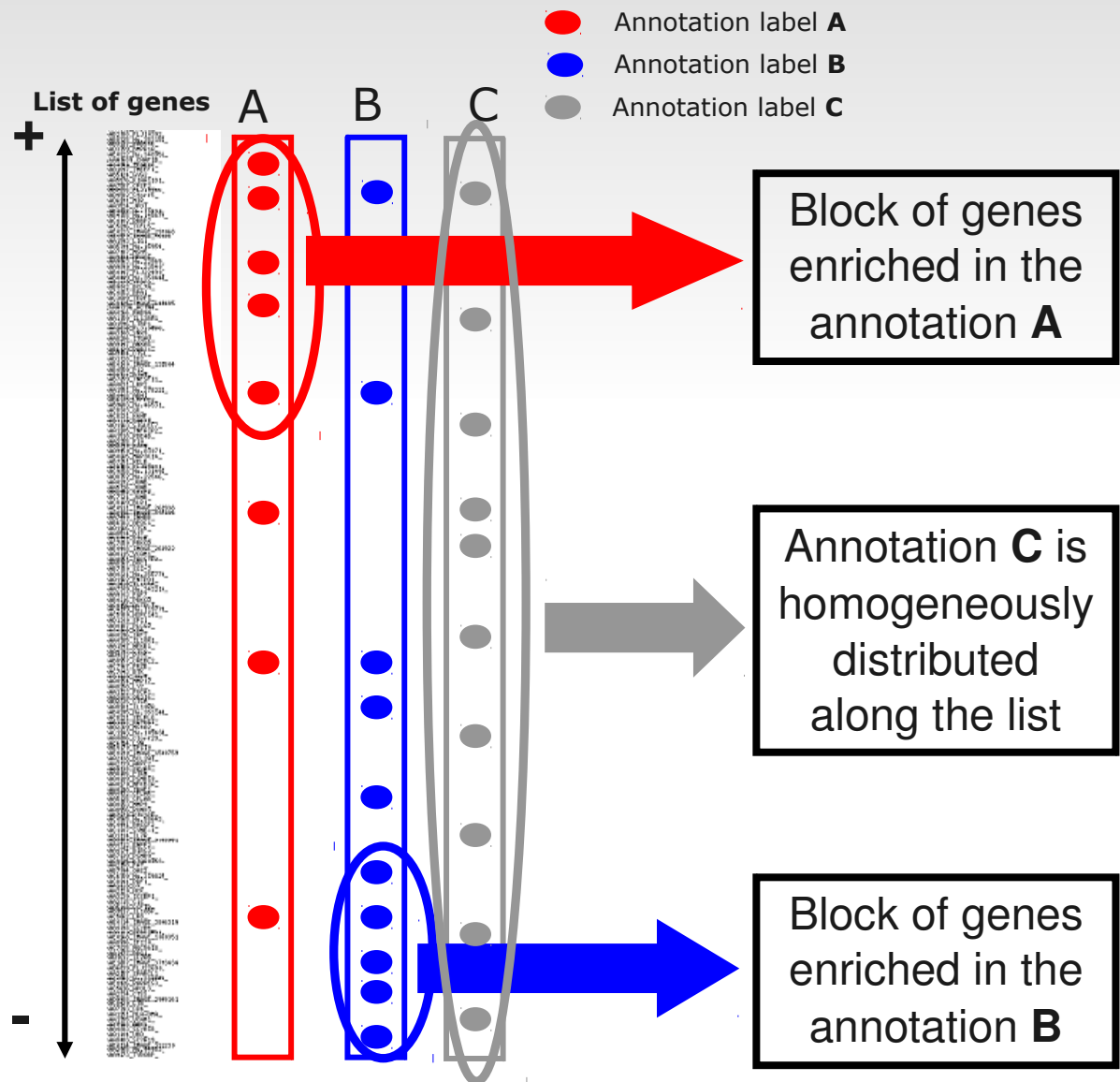
Fatiscan features

- Interpret a ranked list of genes.
- There is not need for choosing a cut-off. All information is included.
- One statistical test for each Functional **Block** of annotation
 - Fisher's exact test
 - Multiple testing context (hundreds of annotation)
 - Filtering of annotation is convenient (the less tests the best correction)

FatiScan

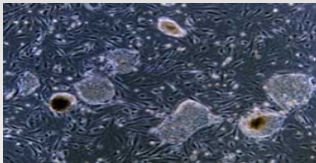
Testing along an ordered list

- 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.

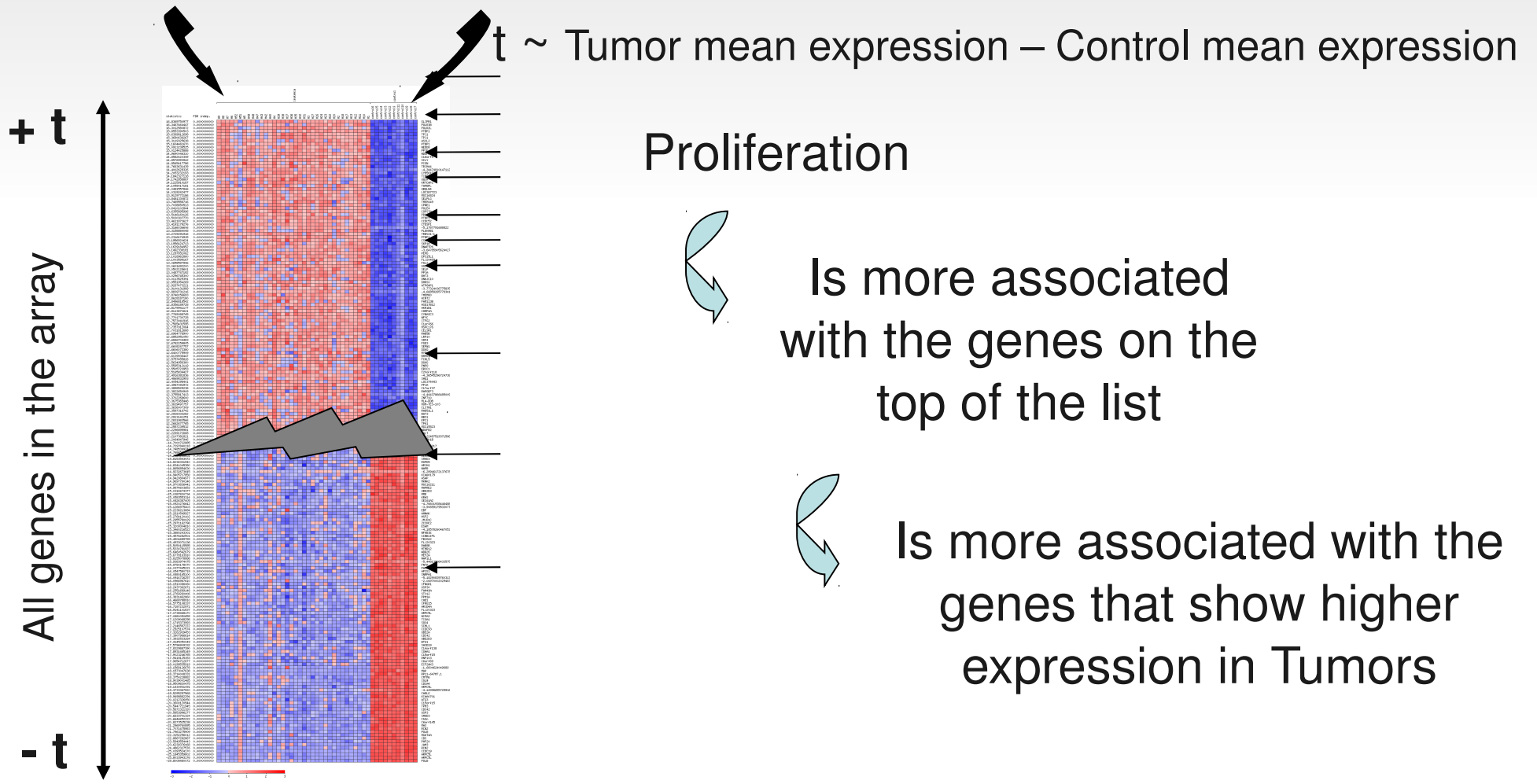
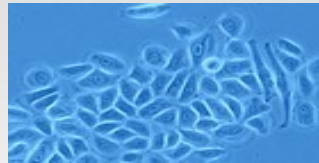


FatiScan Example - two classes

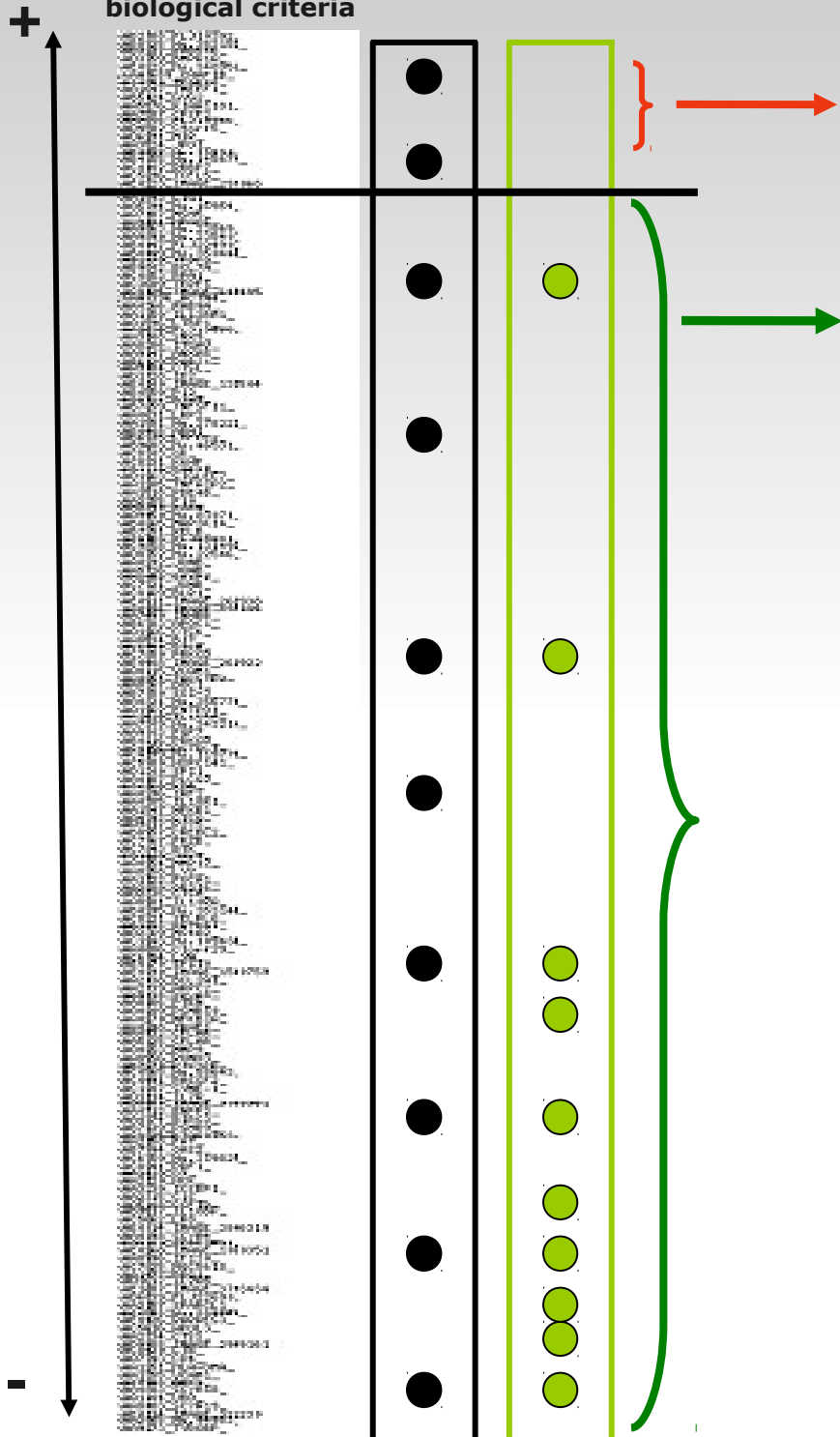
Tumor



Control

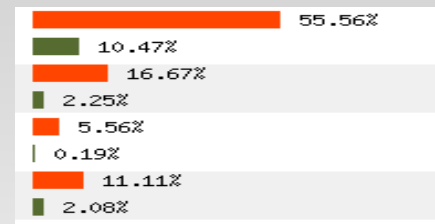


List of genes ranked by biological criteria



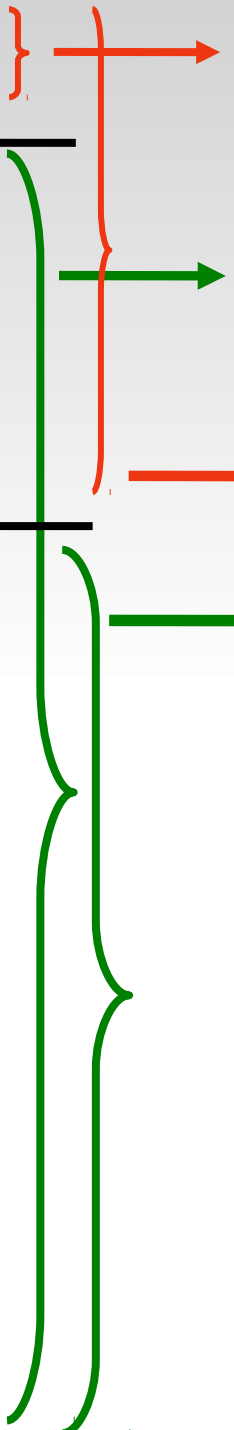
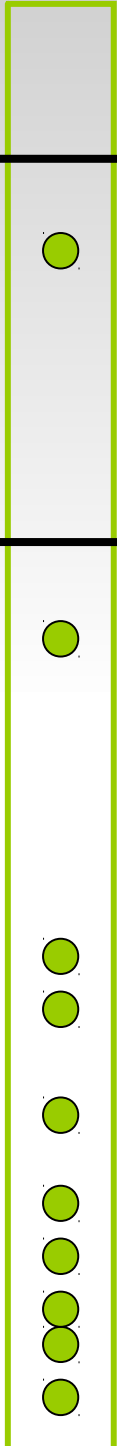
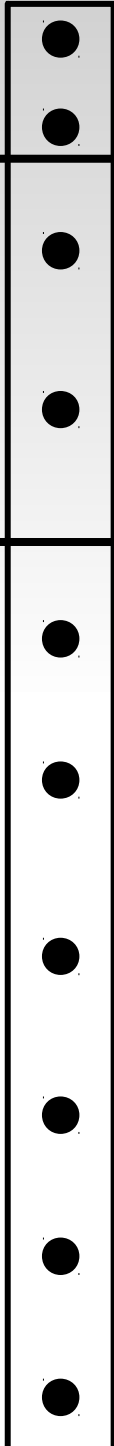
Fisher's test

Scanning test using partitions



List of genes ranked by biological criteria

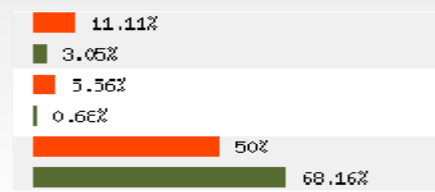
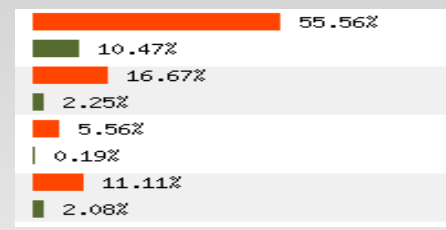
+
-



Fisher's test

Fisher's test

Scanning test using partitions

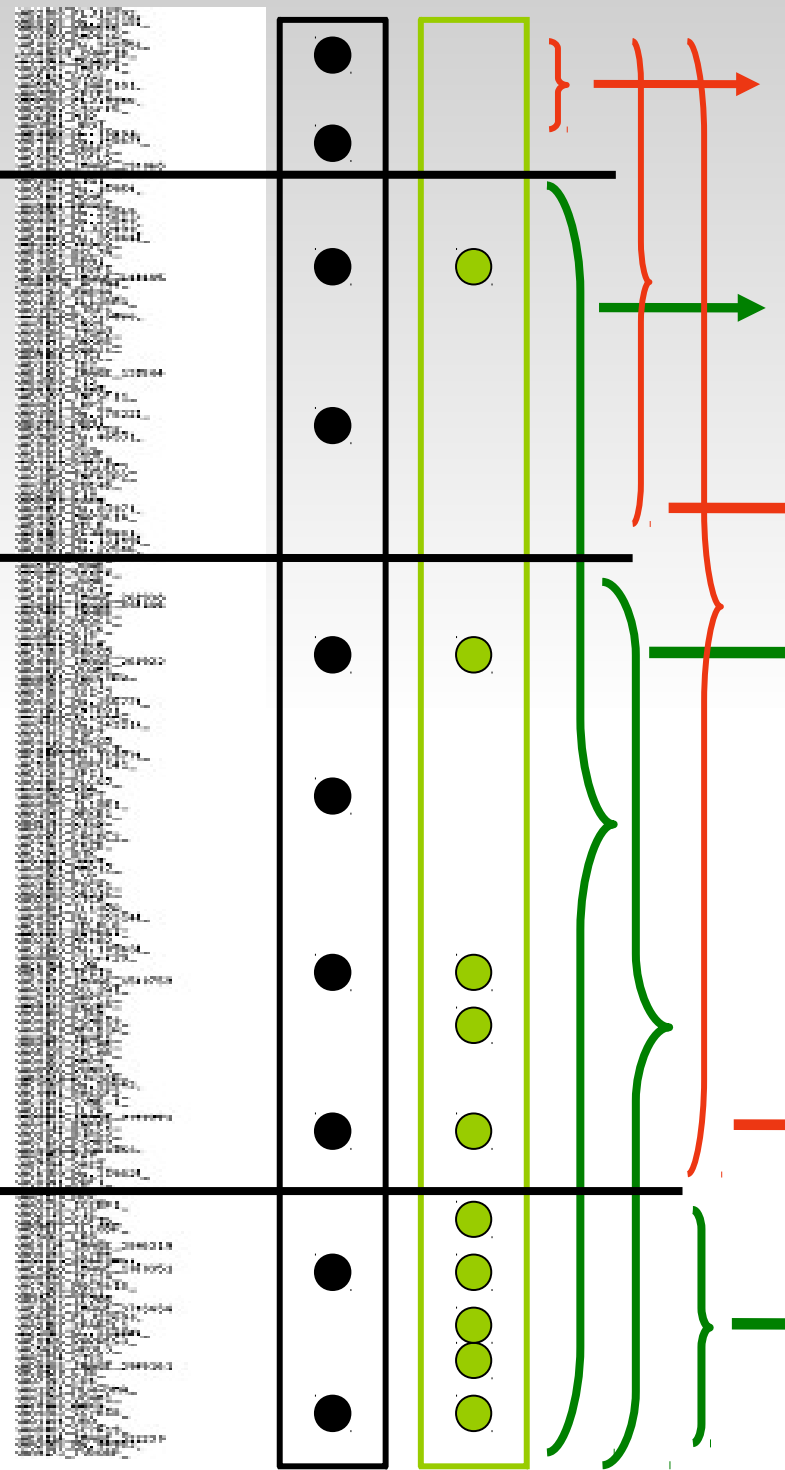


List of genes ranked by biological criteria

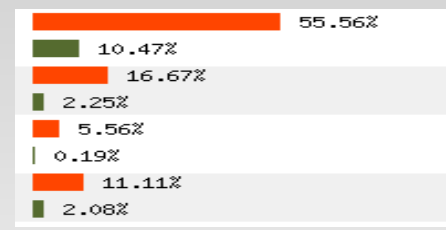
Scanning test using partitions

+

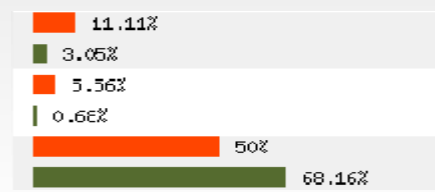
-



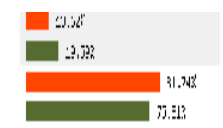
Fisher's test



Fisher's test

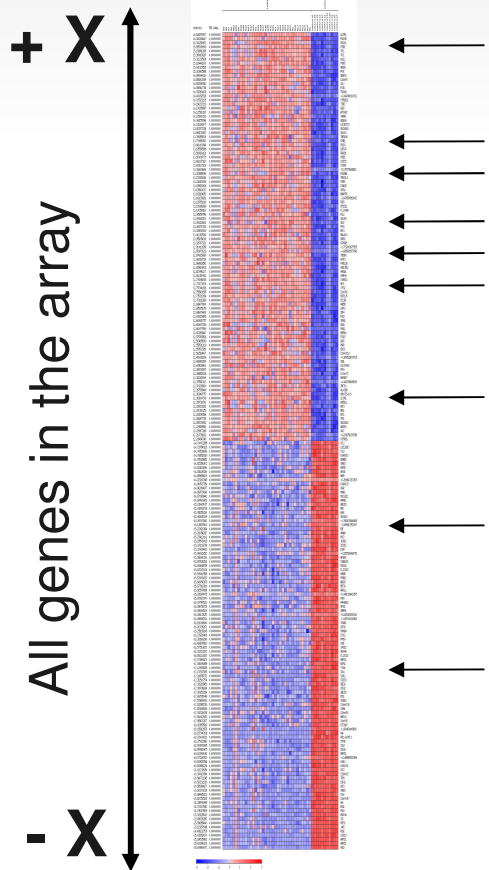


Fisher's test



New 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

Babelomics Tools



MARMITE: Finds differential distributions of bioentities extracted from **PubMed** between two groups of genes.



MarmiteScan: Use chemical and disease-related information to detect related blocks of genes in a gene list with associated values.

Annotation is not 0, 1 may take any value.

Number of articles in which a gene is associated to a:

- Chemical product
- Disease associated
- Drug
- Gene
- Symptom

Take Home Message

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