# Massive Data Analysis

# Introduction

**Department of Bioinformatics and Genomics, (BIG)
Centro de Investigación Príncipe Felipe (CIPF), and
Functional genomics node, (INB),
Valencia, Spain.**

**http://www.gepas.org.
http://www.babelomics.org**

INB

MINISTERIO DE EDUCACIÓN Y CIENCIA
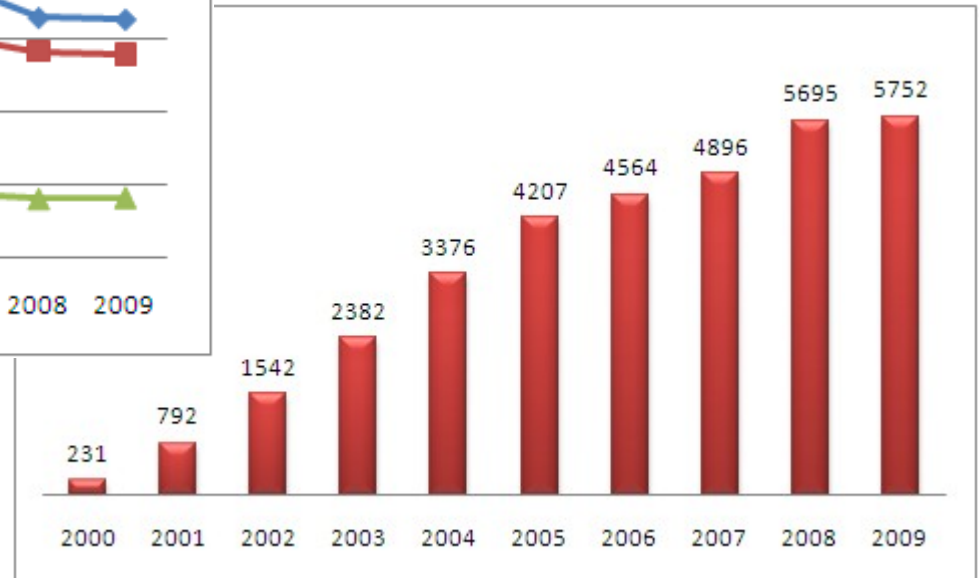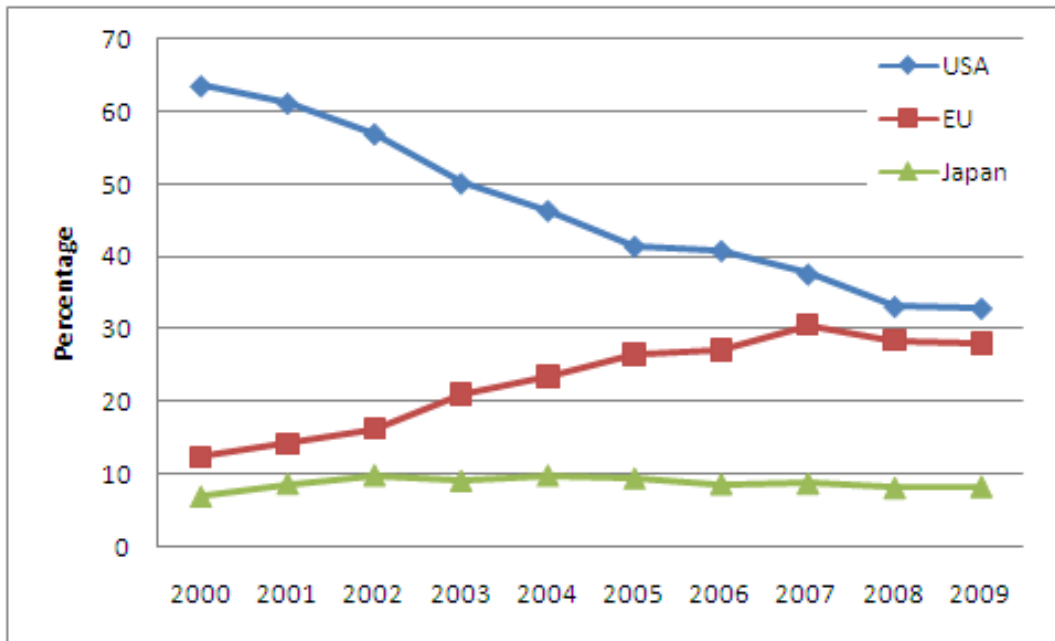
Rticc

ciberer

PRINCIPE FELIPE
CENTRO DE INVESTIGACION

INSTITUTO NACIONAL

# The Bioinformatics and Genomics Department at the Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain, and...
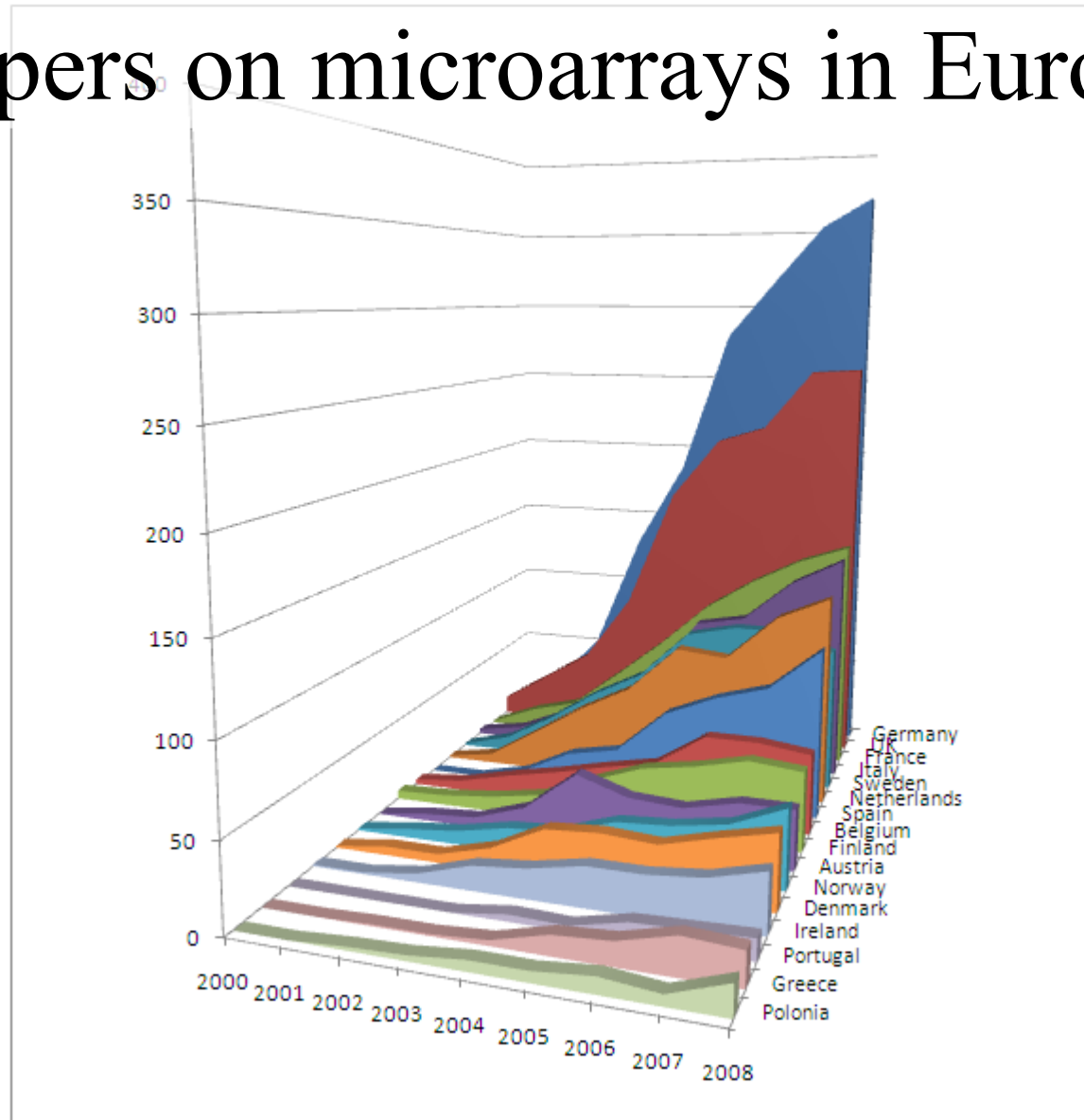
# Evolution of the percentages of published papers on microarrays



**Source Pubmed. Query:**
**date[Entrez Date] AND country[Affiliation]AND
microarray[Title/Abstract]**

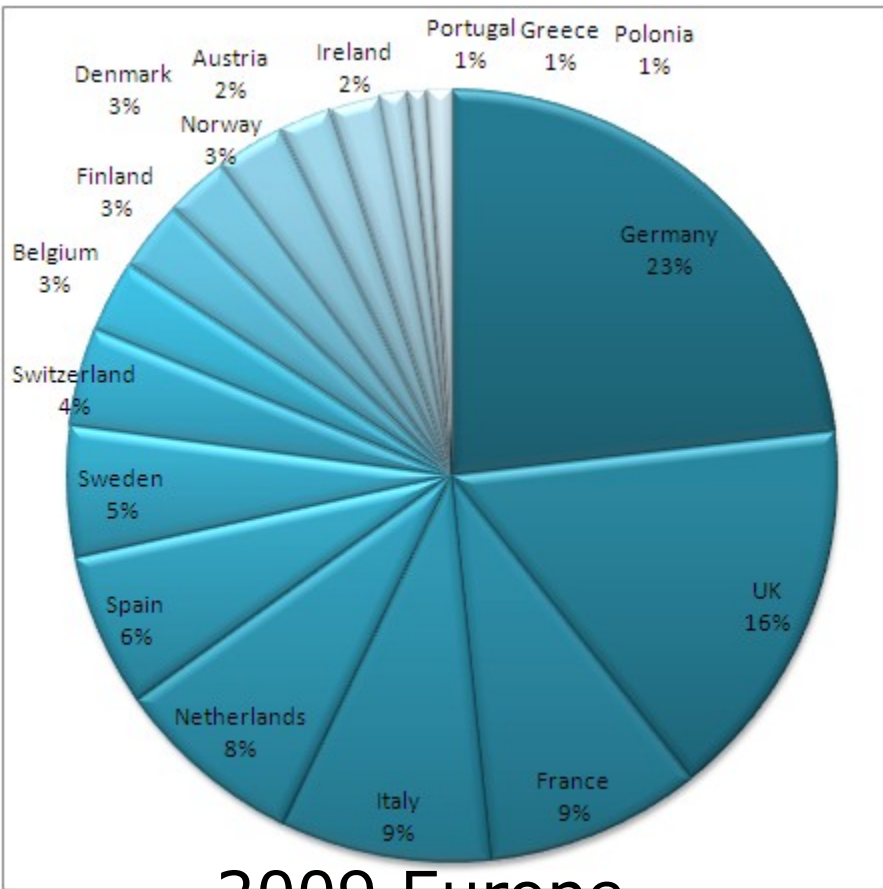# Evolution of the percentages of published papers on microarrays in Europe



**Source Pubmed. Query:**
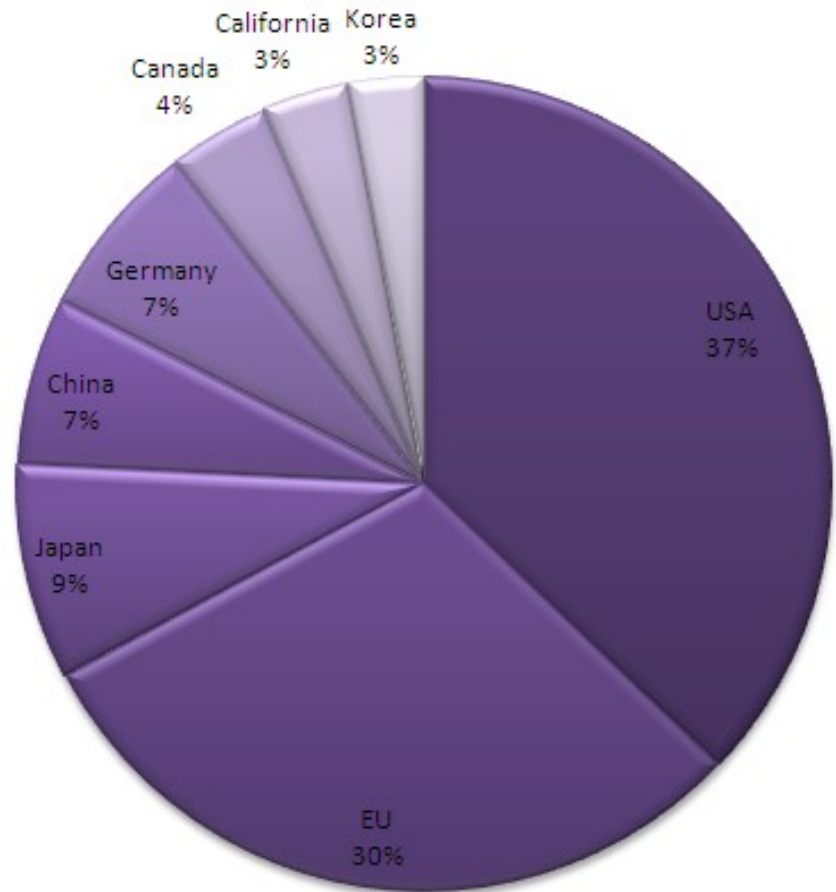**date[Entrez Date] AND country[Affiliation]AND microarray[Title/Abstract]**
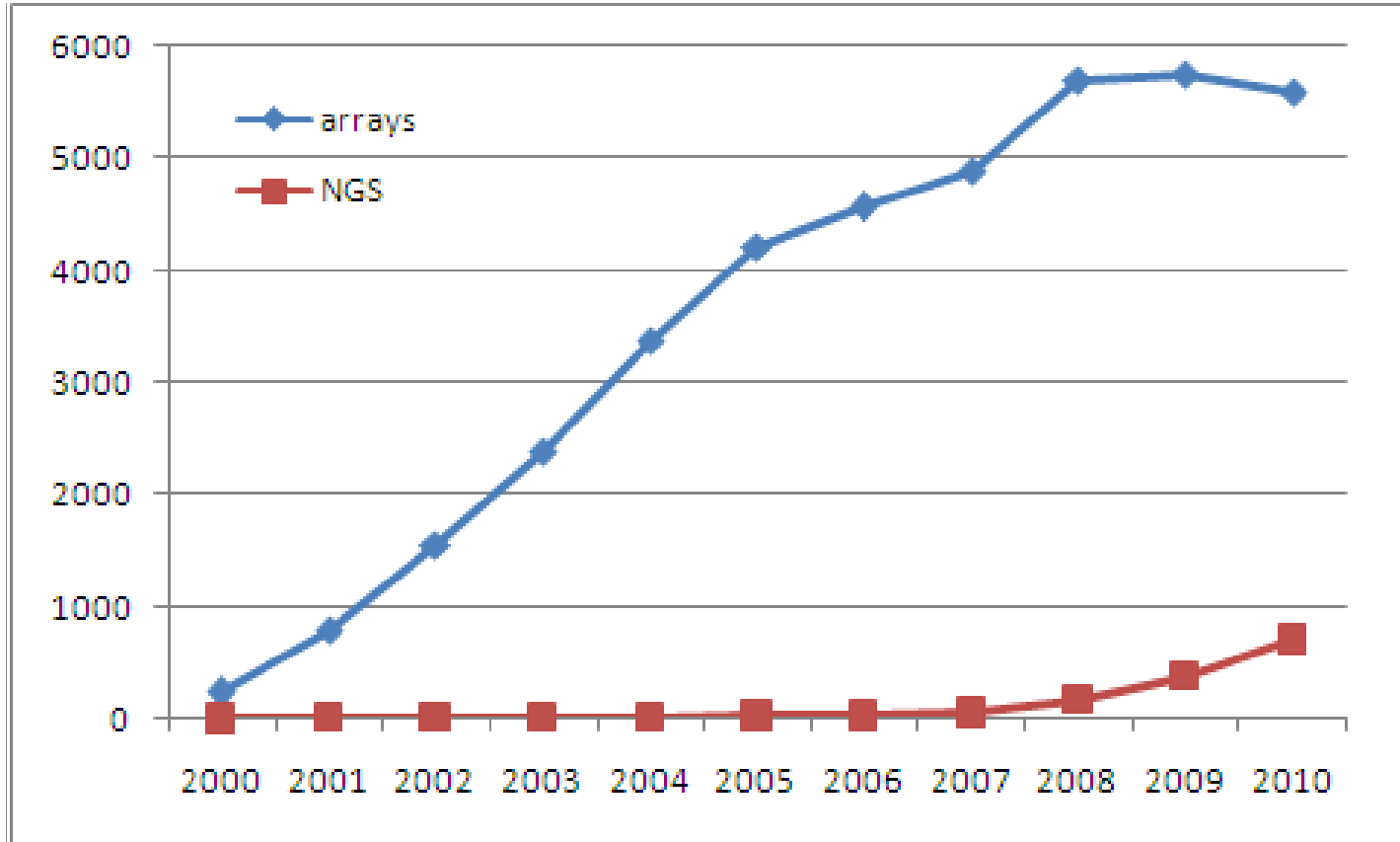
# Microarray publications



2009 Europe

2009 Worldwide

**Source Pubmed. Query: 2009[Entrez Date] AND country[Affiliation]AND microarray[Title/Abstract]**

# Trends in publications



**Source Pubmed. Query:** "high-throughput sequencing"[Title/Abstract] OR "next generation sequencing"[Title/Abstract] OR "rna seq"[Title/Abstract]) AND year[Publication Date]

# Some numbers

451 papers cite GEPAS (215 are SOTA cites)

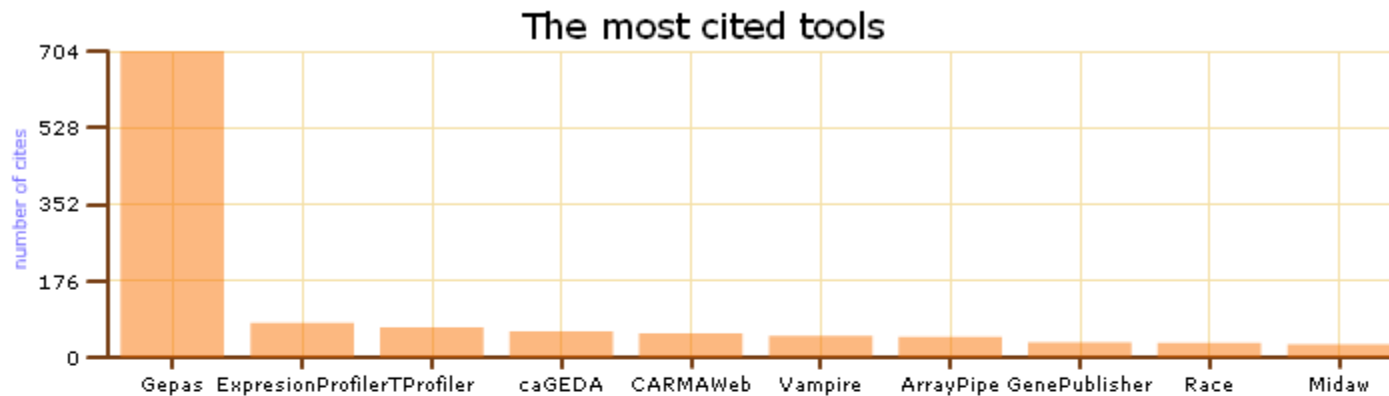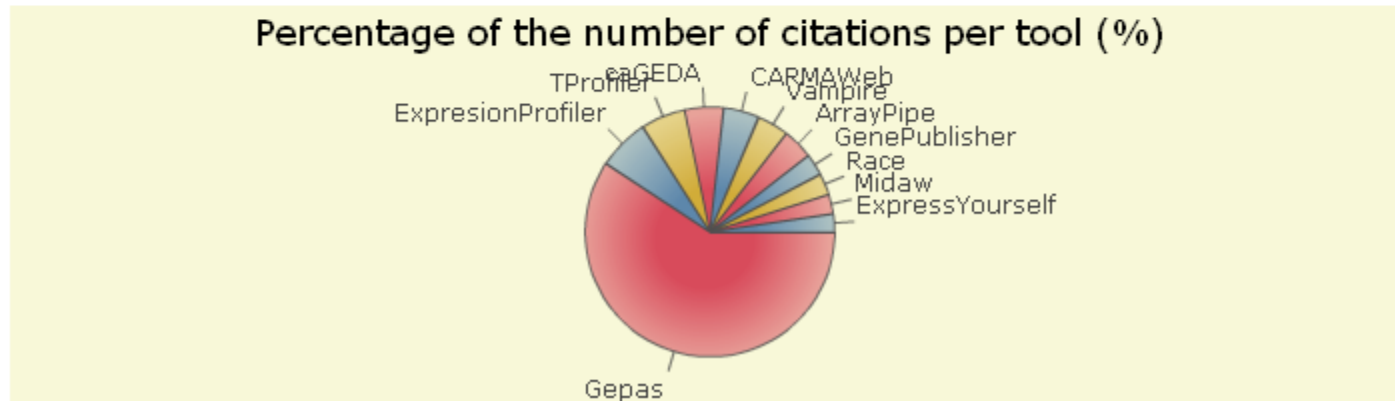632 papers cite Babelomics (442 are FatiGO cites)

*(source ISI Web of Knowledge, May 2010)*

More than 150,000 experiments analysed during the last year.
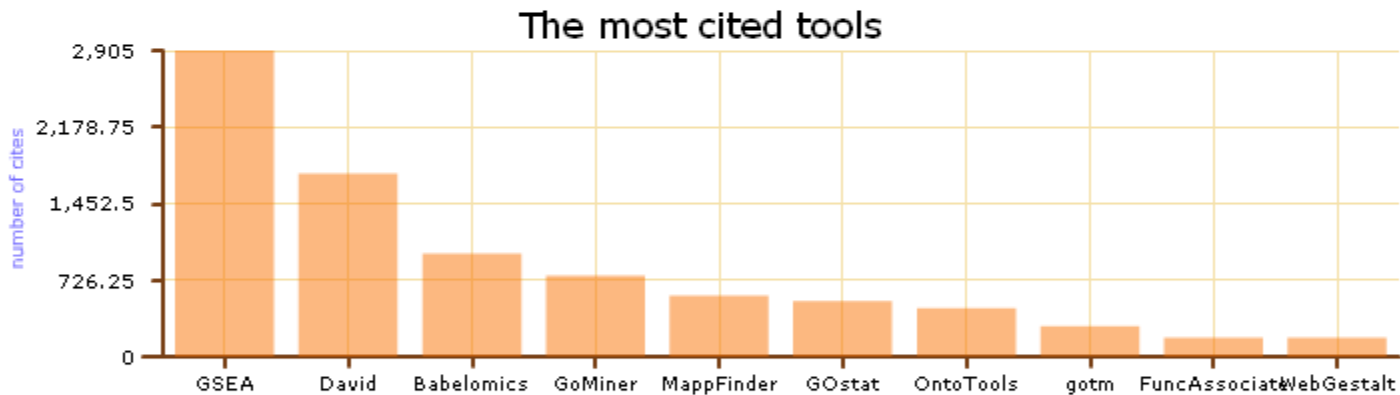
More than 1000 experiments per day.

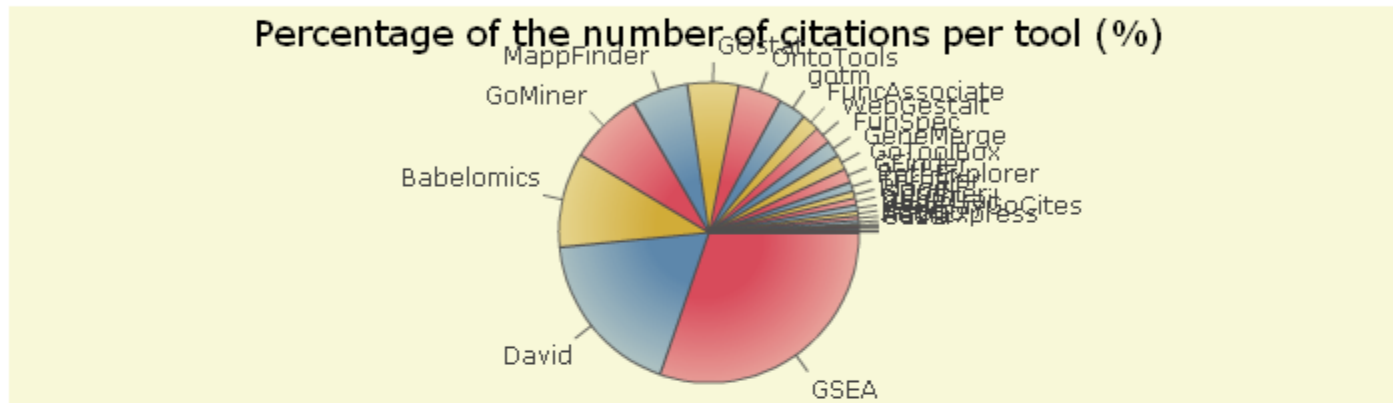# Tools for gene expression analysis



Percentage of the number of citations per tool (%)



The most cited tools

# Tools for functional profiling



Percentage of the number of citations per tool (%)



The most cited tools

# Structure of the course

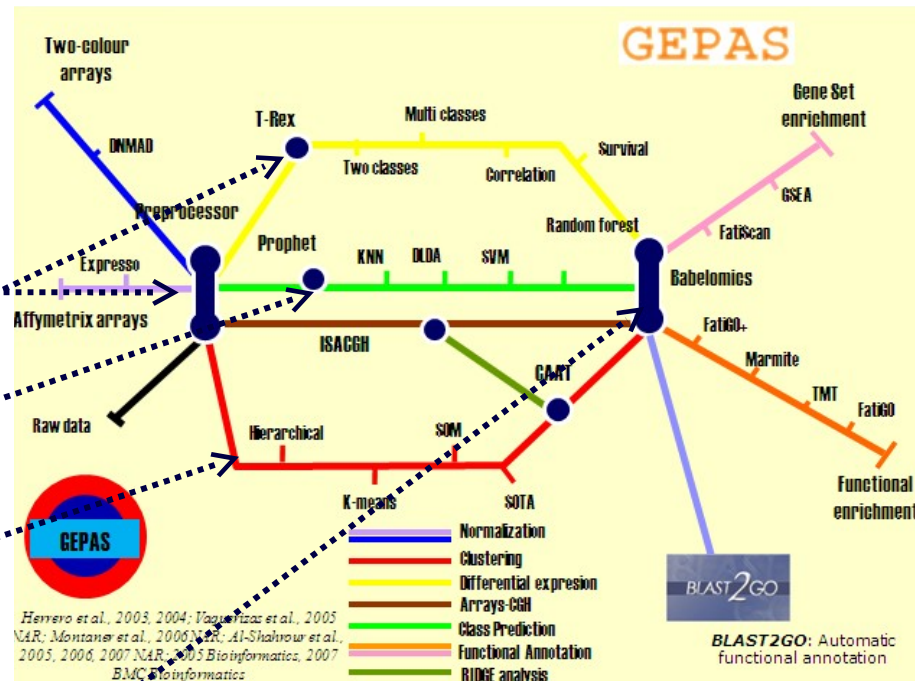| Theoretical | Hands-on **GEPAS** |
|---|---|

**Introduction**

Normalization

Gene selection

Predictors

Clustering

Functional
interpretation

# Background
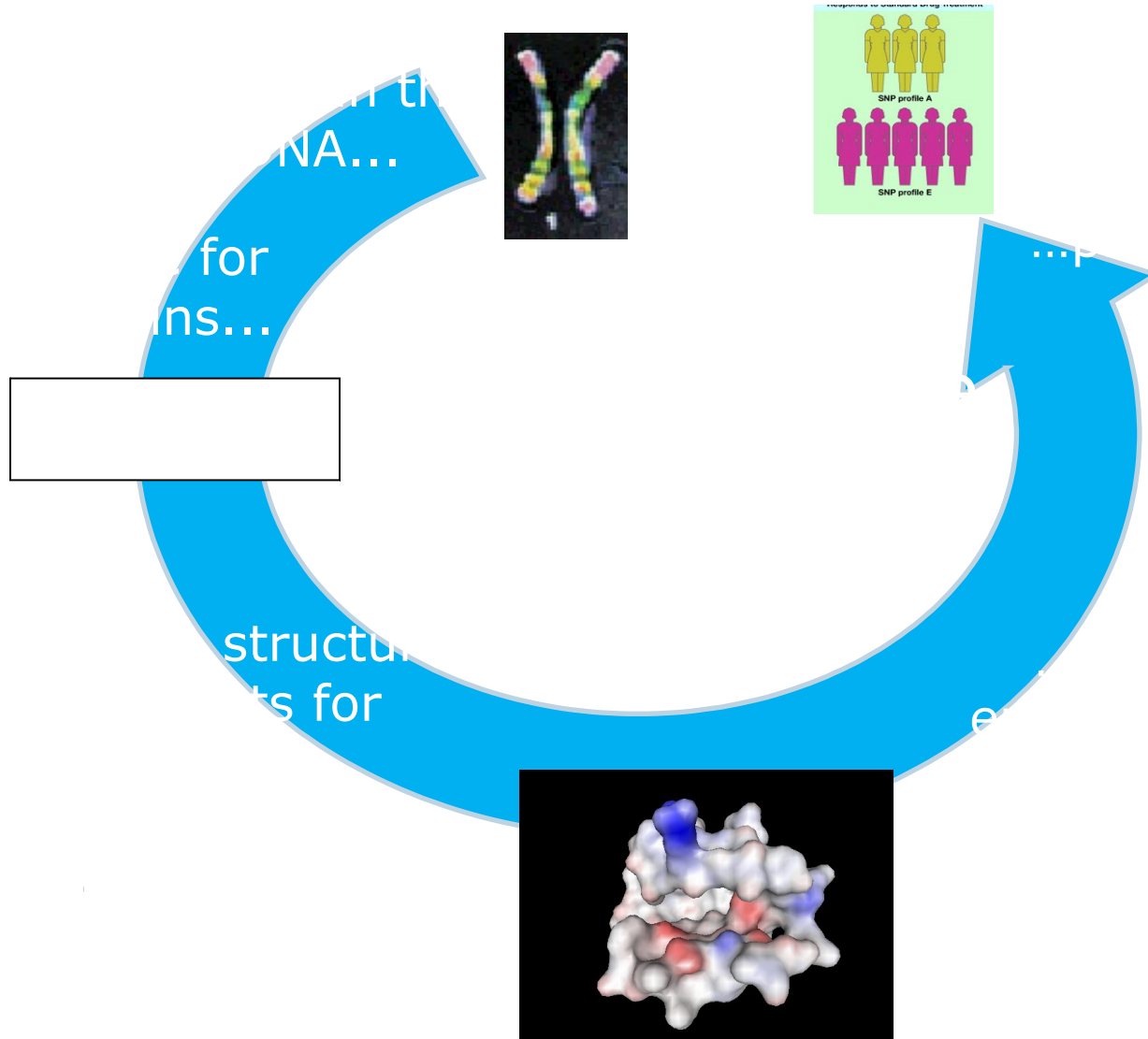


**The road of excess leads to the palace of wisdom**

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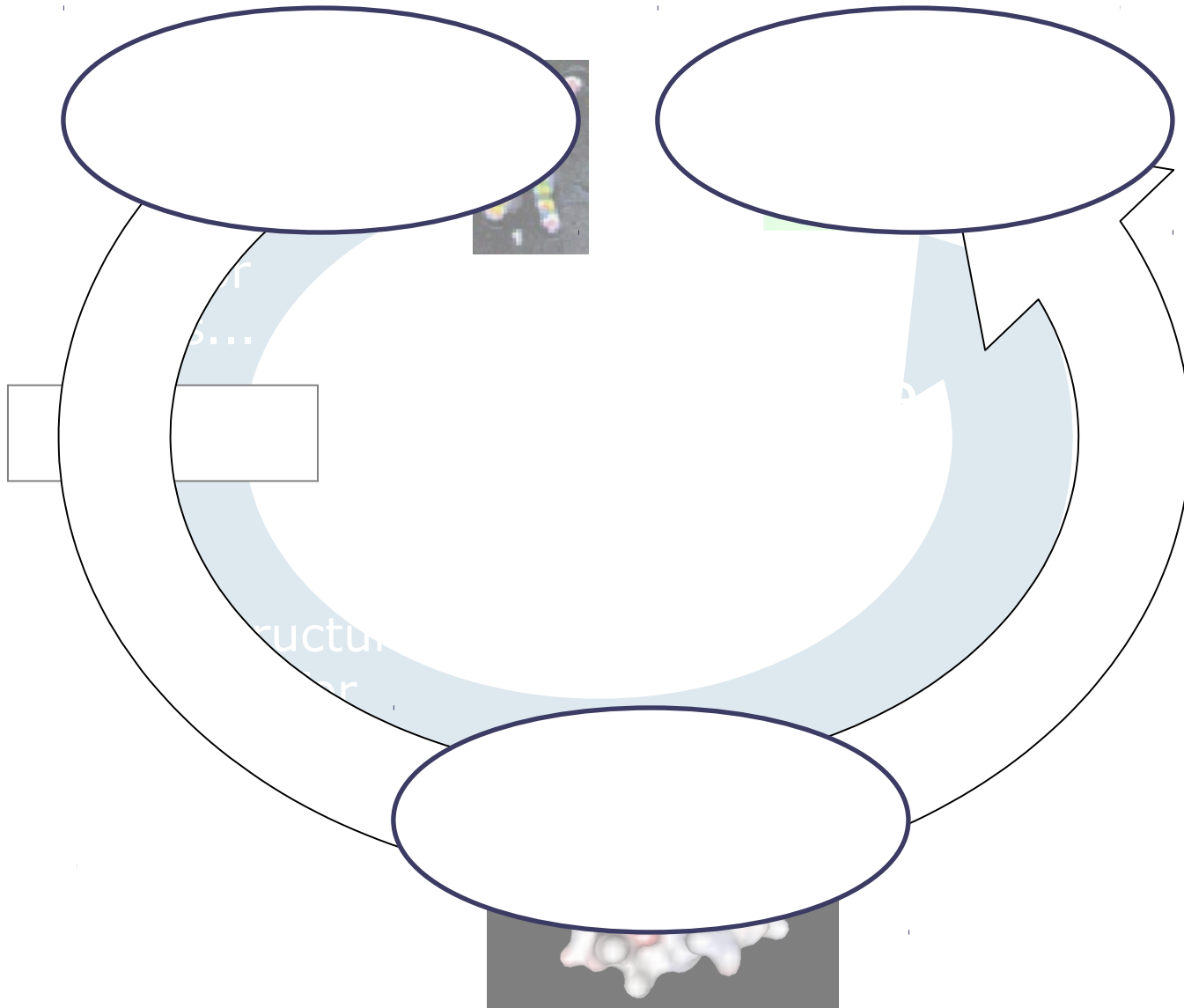
But not necessarily the way in which we really address or test them…

# Where do we come from?
# The pre-genomics paradigm

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

Next Generation Sequencing
$10^9$bp per round
($10^{10}$ expected soon)

Genes in the DNA.

>protein kunase
acctgttgatggcgacagggactgtatgctgatct
atgctgatgcatgcatgctgactactgatgtgggg
gctattgacttgatgtctatc....

…which can be different because of the variability.

15 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

From genotype to phenotype.

(in the functional post-genomics scenario)

…conforming complex interaction networks…

…code for proteins…

That undergo post-translational modifications, somatic recombination…
100K-500K proteins

…in cooperation with other proteins…

…that account for function if…

Each protein has an average of 8 interactions

...when ex proper place

A typical tissue is expressing among 5000 and 10000 genes

Causes

Effects

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

From genotype to phenotype

(in the functional genomics scen
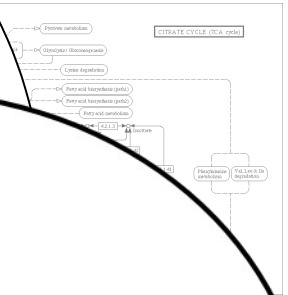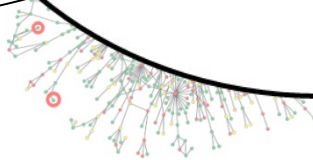
...code proteins.

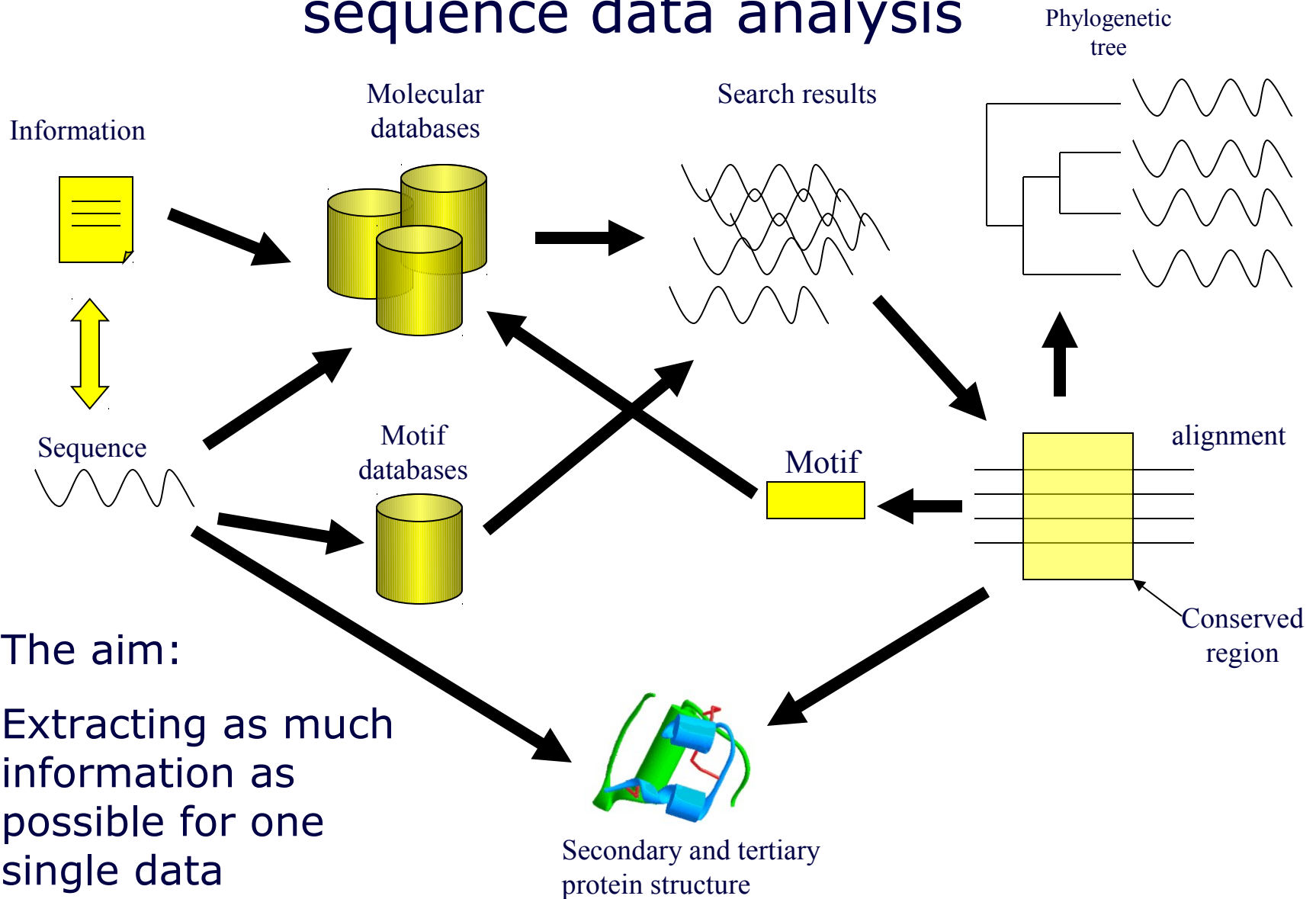That undergo post-translational modifications, somatic recombination... 100K-500K proteins

Function (modules of proteins)

...plex ks...

...whose structures account for function...

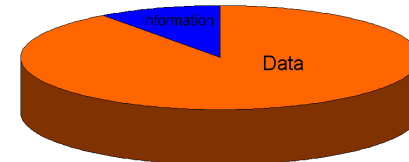Each protein has an average of 8 interactions

...ration other proteins...

# Bioinformatics tools for pre-genomic sequence data analysis

Information

Molecular databases

Search results

Phylogenetic tree

Sequence

Motif databases

Motif

alignment

Conserved region

The aim:

Extracting as much information as possible for one single data

Secondary and tertiary protein structure

# Post-genomic vision

## EMBL database growth (March 2009)

# Genome scale data and a note of caution on associations, correlations or patterns discovered:

Genome-wide technologies allows us to produce vast amounts of data.

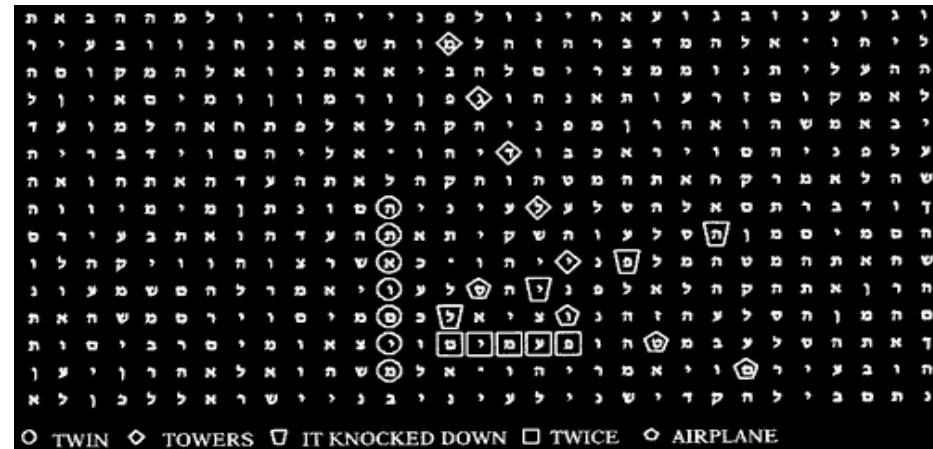But... dealing with many data (omic data) increase the occurrence of spurious associations due to chance

Hypothesis ➝ Experiment ➝ test

Is gene A involved in process B?

Experiment ➝ (sometimes) test ➝ Hypothesis

Is there any gene (or set of genes) involved in any process?

Sure, but... Is it real? (many hypotheses are rejected while this one is accepted *a posteriori*: numerology)
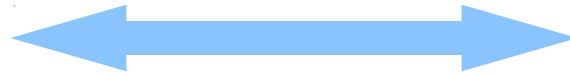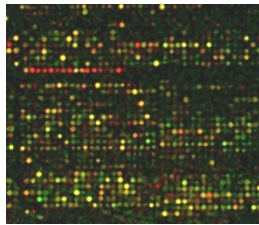
The test is dependent on the hypothesis and not *vice versa*



○ TWIN   ◇ TOWERS   ⬯ IT KNOCKED DOWN   □ TWICE   ◇ AIRPLANE

# Gene expression profiling. Historic perspective

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



• Classification of phenotypes / experiments. Can I distinguish among classes (either known or unknown), values of variables, etc. using molecular gene expression data? (sensitivity)

• Selection of differentially expressed genes among the phenotypes / experiments. Did I select the relevant genes, all the relevant genes and nothing but the relevant genes? (specificity)

• Biological roles the genes are carrying out in the cell. What general biological roles are really represented in the set of relevant genes? (interpretation)

# Microarrays arrive to an acceptable level of reproducibility



nature biotechnology

OCTOBER 2006
www.nature.com/nbt/journal/v24/n10s

Produced with support from

FDA  EPA United States Environmental Protection Agency

NATIONAL CANCER INSTITUTE  Agilent Technologies

The MicroArray Quality Control Consortium

Supplement to Nature Publishing Group

ARTICLES

nature biotechnology

## 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.

# FDA approves the first predictor based on microarrays

# DNA microarrays: the paradigm of a post-genomic technique



Competitive hybridization (two colors)

One color

# Primary analysis

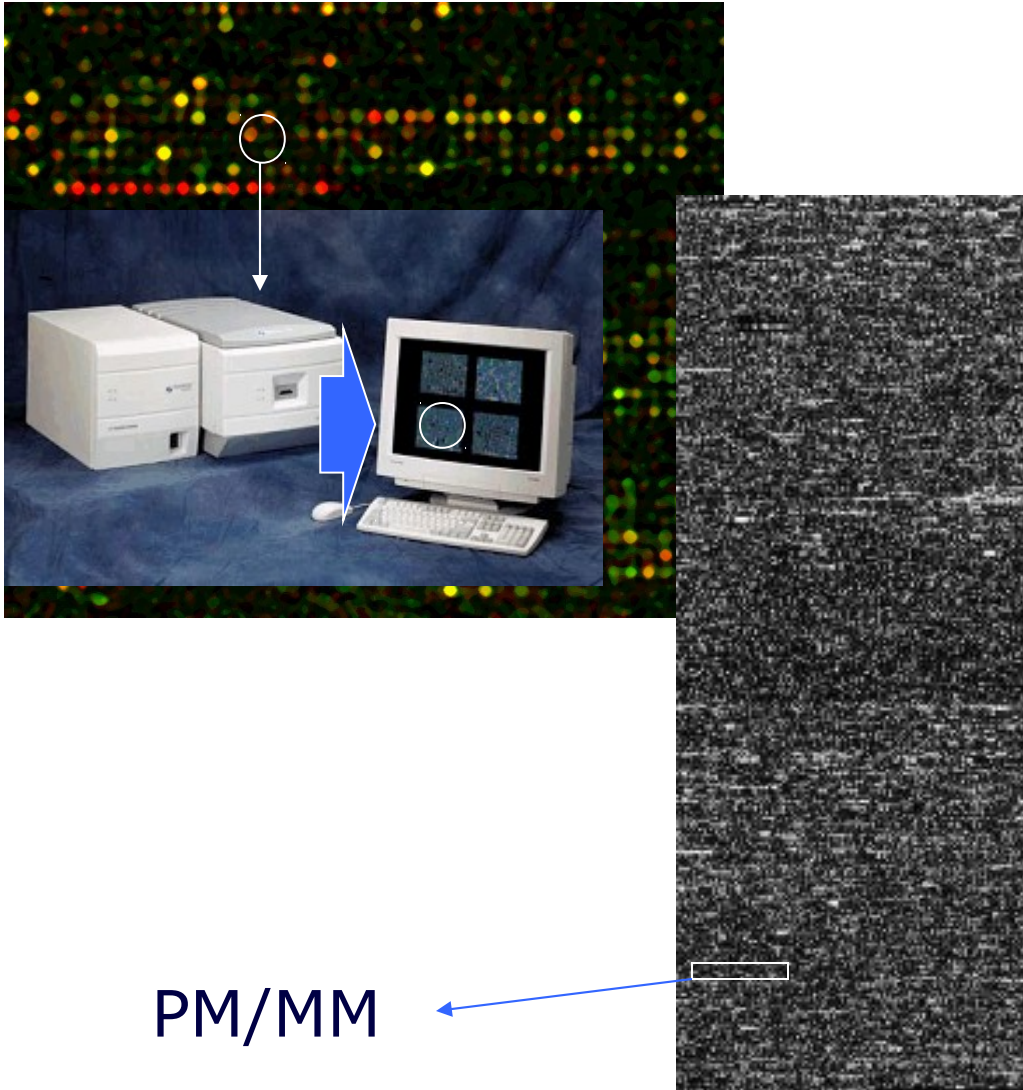•Transform images corresponding to hybridization intensities into numbers

•Convert all the numbers to a common scale that makes them comparable across experiments.

# Transforming images into numbers



**Two-color**
Test sample labeled red (**Cy5**)
Reference sample labeled green (**Cy3**)
Red : gene overexpressed in test sample
Green : gene underexpressed in test sample
**Yellow** - equally expressed
**red/green** - ratio of expression

**One color**
**Intensity** of a gene using the probes

**Affymetrix**
**Intensity** of a gene using the probes PM and in MM

Scanners generate a graphic file.

Software analyzes the file: GenePix Pro (by Axon Instruments, Inc.) or Imagene (By Biodiscovery, Inc.)
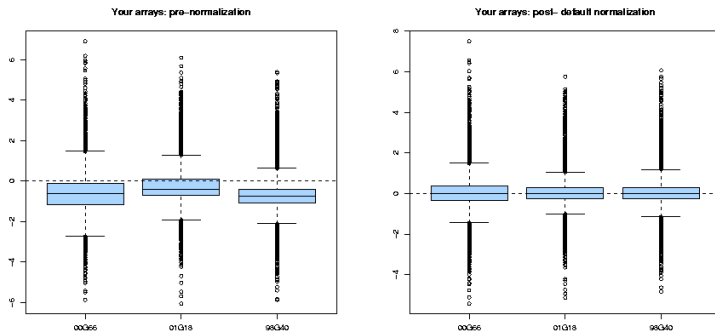There are free systems too: TIGR Spotfinder, ScanAlyze, etc

PM/MM

# Normalisation



A

B

C

Before (left) and after (right) normalisation. A) BoxPlots, B) BoxPlots of subarrays and C) MA plots (ratio versus intensity)

(a) After normalization by average (b) after print-tip lowess normalization (c) after normalisation taking into account spatial effects

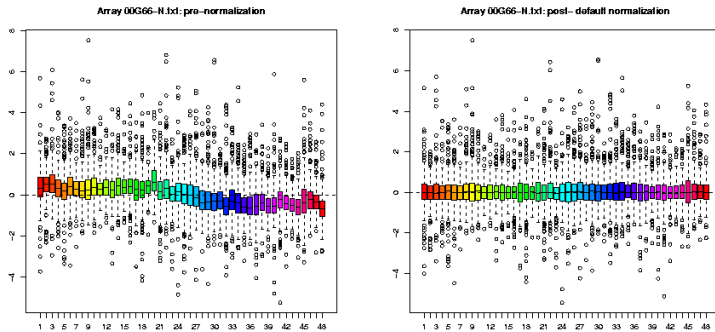There are many sources of error that can affect and seriously bias the interpretation of the results. Differences in the efficiency of labelling, the hybridisation, local effects, etc.

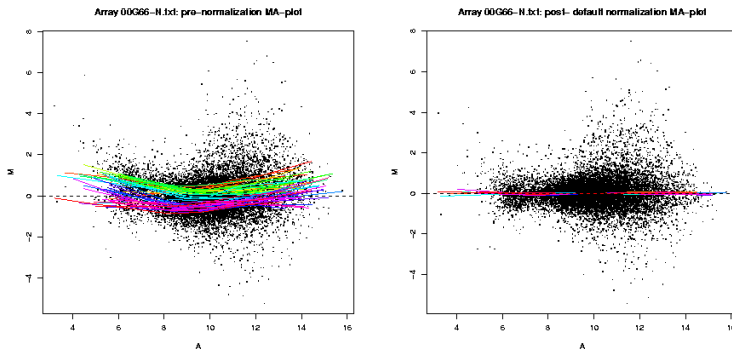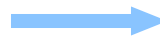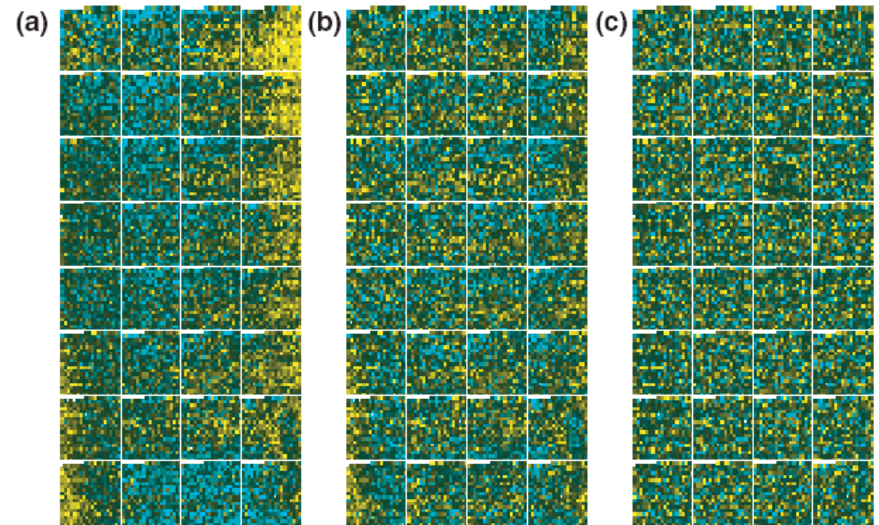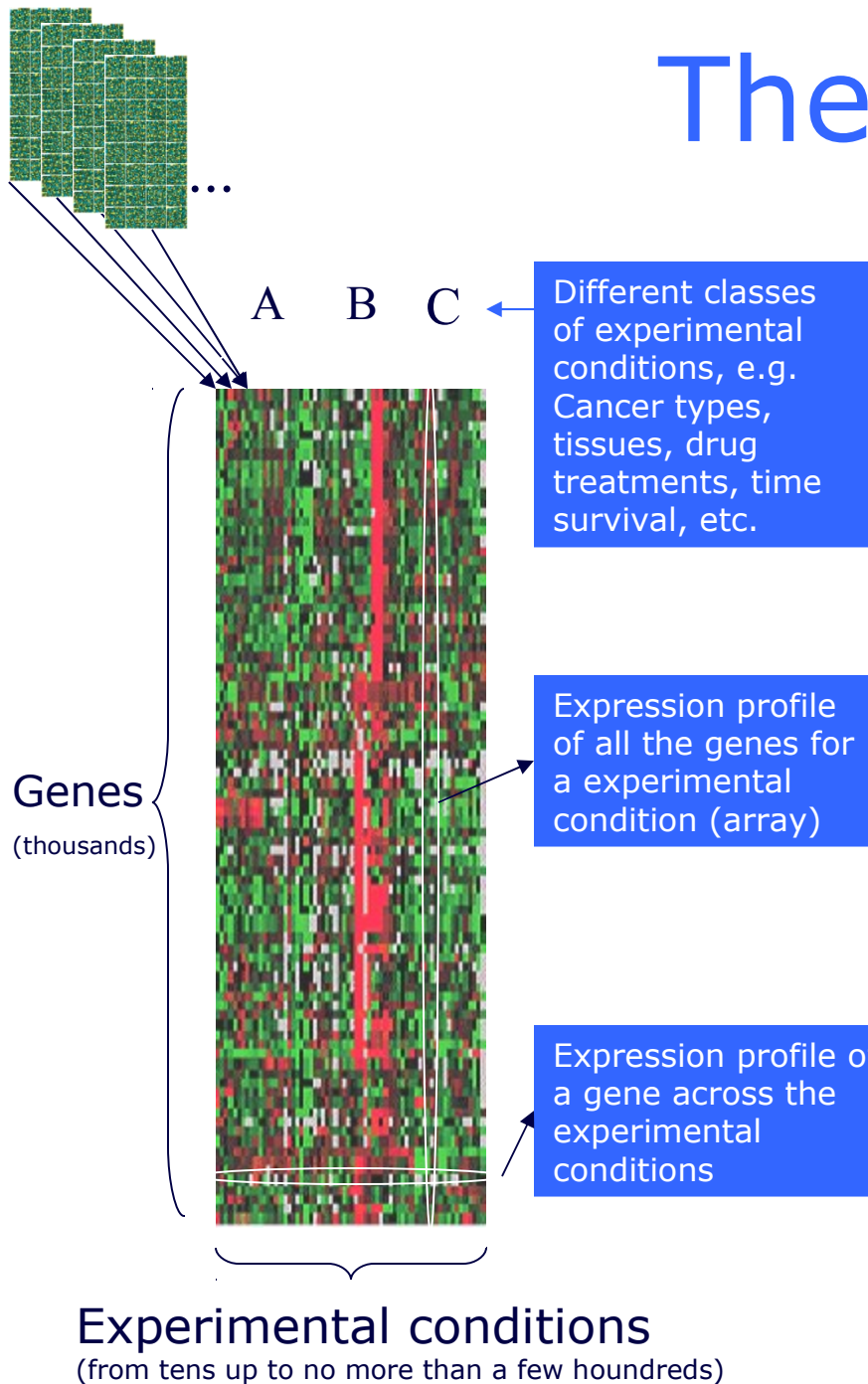Normalisation is a necessary step before proceeding with the analysis

# Secondary analysis

Once the measurements are in a common, comparable scale the results can be studied.
Diferent studies can be made that include class discovery, classification, gene selection, etc.

# The data



A  B  C

Different classes of experimental conditions, e.g. Cancer types, tissues, drug treatments, time survival, etc.

Genes
(thousands)

Expression profile of all the genes for a experimental condition (array)

Expression profile of a gene across the experimental conditions

Experimental conditions
(from tens up to no more than a few houndreds)

**Characteristics of the data:**

• We NEVER deal with individual arrays, we deal with collections of arrays obtained for a given experimental design

• Most of the genes are not informative with respect to the trait we are studying (account for unrelated physiological conditions, etc.)

• Number of variables (genes) is several orders of magnitude larger than the number of experiments

• Low signal to noise ratio
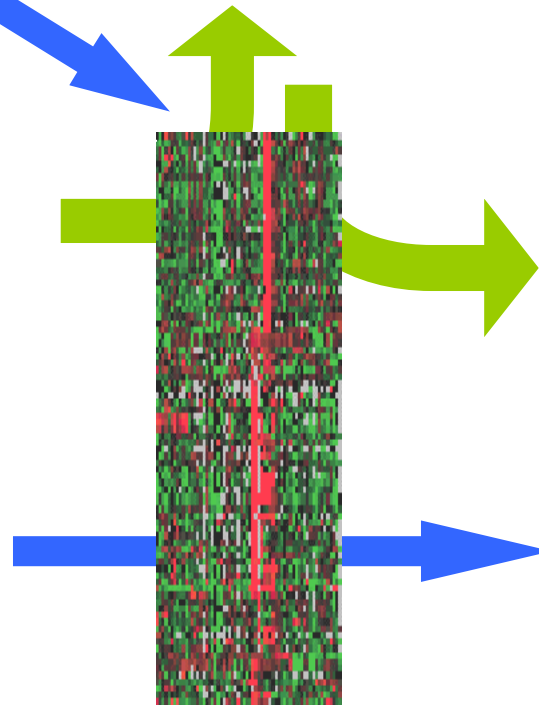
# Studies must be hypothesis driven.

What is our aim? Class discovery? sample classification?  gene selection? …

Can we find groups of experiments with similar gene expression profiles?

Different classes…

Unsupervised

Supervised

Molecular classification of samples

What genes are responsible for?

Co-expressing genes…
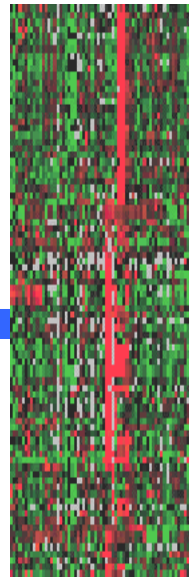
What do they have in common?

# Unsupervised problem: class discovery

Our interest is in discovering clusters of items (genes or experiments) which we do not know beforehand

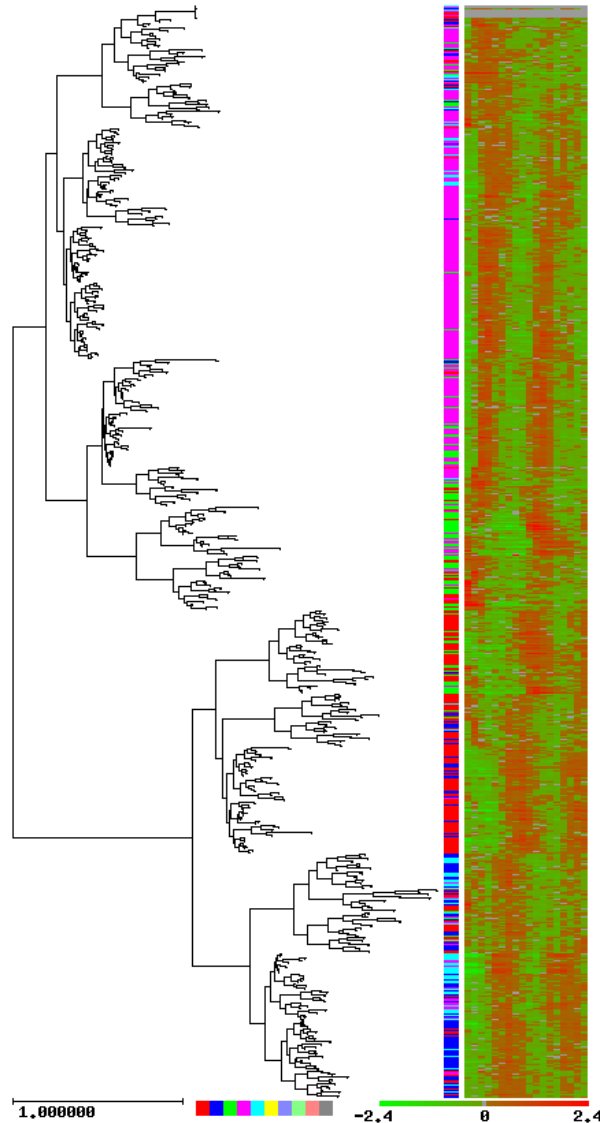Can we find groups of experiments with similar gene expression profiles?



Co-expressing genes...

- What genes co-express?

- How many different expression patterns do we have?

- What do they have in common?

- Etc.

# Unsupervised clustering methods: Method + distance: produce groups of items based on its <u>global</u> similarity

**Non hierarchical**     **hierarchical**

K-means, PCA     UPGMA

SOM     SOTA

Different levels of information

# An unsupervised problem: clustering of genes.



- Gene clusters are previously unknown

- Distance function

- Cluster gene expression patterns based uniquely on their similarities.

- Results are subjected to further interpretation (if possible)

# Clustering of experiments: The rationale
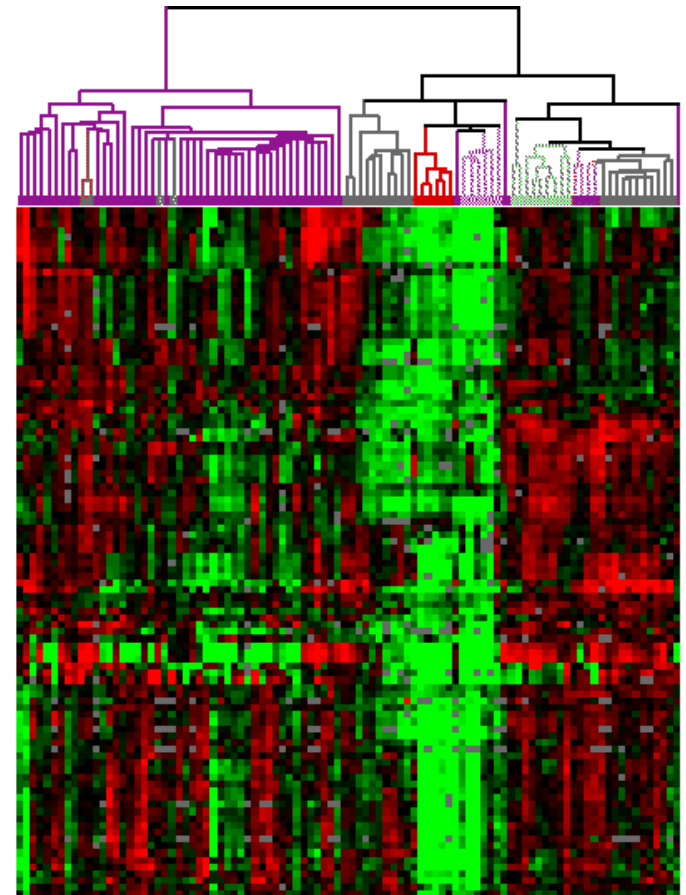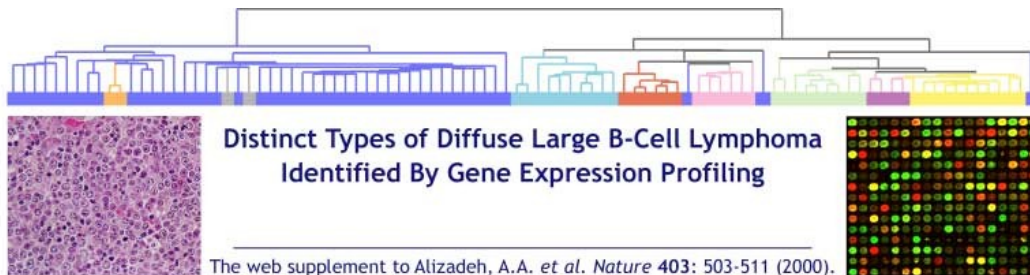
If enough genes have their expression levels altered in the different experiments, we might be able of finding these classes by comparing gene expression profiles.

**Distinctive gene expression patterns in human mammary epithelial cells and breast cancers**

Overview of the combined *in vitro* and breast tissue specimen cluster diagram. A scaled-down representation of the 1,247-gene cluster diagram The black bars show the positions of the clusters discussed in the text: (*A*) proliferation-associated, (*B*) IFNregulated, (*C*) B lymphocytes, and (*D*) stromal cells.

*Perou et al., PNAS 96 (1999)*

# Clustering of experiments: The problems

Any gene (regardless its relevance for the classification) has the same weight in the comparison.

If relevant genes are not in overwhelming majority we will find:
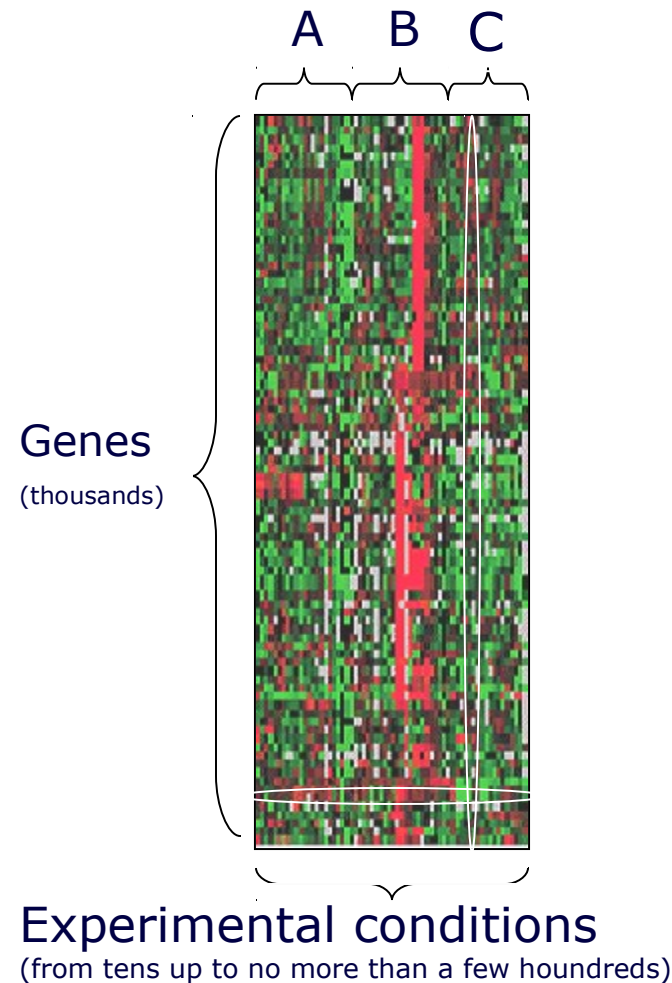
Noise

and/or

irrelevant trends



**Distinct Types of Diffuse Large B-Cell Lymphoma Identified By Gene Expression Profiling**

The web supplement to Alizadeh, A.A. *et al. Nature* **403**: 503-511 (2000).

# Supervised problems: Class prediction and gene selection, based on gene expression profiles

Information on classes (<u>defined on criteria external to the gene expression measurements</u>) is used.

A   B   C

Genes
(thousands)

Experimental conditions
(from tens up to no more than a few houndreds)

Problems:

How can classes A, B, C... be distinguished based on the corresponding profiles of gene expression?

How a continuous phenotypic trait (resistance to drugs, survival, etc.) can be predicted?

Class prediction

And

Which genes among the thousands analysed are relevant for the classification?

Gene selection

# Studies must be hypothesis driven.

## gene selection

Can we find groups
of experiments with
similar gene
expression profiles?

**Different classes...**

Molecular
classification of
samples

**What genes are
responsible for?**

Co-expressing genes...

What do they
have in
common?

# Gene selection.

The simplest way: univariant gene-by-gene.
Other multivariant approaches can be used

- **One class**
  Limma
- **Two classes**
  T-test
  Limma
  Fold-change

- **Multiclass**
  Anova
  Limma

- **Continuous variable (e.g. level of a metabolite)**
  Pearson
  Spearmam
  Regression

- **Survival**
  Cox model

- **Time Course**

# Gene selection



The t-statistic was introduced in 1908 by William Sealy Gosset

cases | controls

cases    controls

$X_1$    $X_2$

**Significantly different**

$S_{X1}$    $S_{X2}$

$X_1$    $X_2$

**Non significantly different**

$S_{X1}$    $S_{X2}$

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{X_1 X_2} \cdot \sqrt{\frac{2}{n}}}$$

being

$$S_{X_1 X_2} = \sqrt{\frac{S_{X_1}^2 + S_{X_2}^2}{2}}.$$

# A simple problem: gene selection for class discrimination



~15,000 genes

Case(10)/control(10)

thebest – [04/10/2003 18:57:43  GMT]

1.000000

-2.4    0    +2.4

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

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



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

So... Why do we find good p-values?

| unadj.p | adj_p | FDR_indep | FDR_dep | obs_stat |
|---|---|---|---|---|
| 0.00019998 | 0.152685 | 0.49995 | 1 | 5.47044 |
| 0.00019998 | 0.746225 | 0.49995 | 1 | 4.49902 |
| 0.0009999 | 0.983002 | 0.861025 | 1 | 4.01726 |
| 0.00149985 | 0.986401 | 0.861025 | 1 | 3.99374 |
| 0.00129987 | 0.9959 | 0.861025 | 1 | 3.86046 |
| 0.00169983 | 0.9996 | 0.861025 | 1 | 3.7251 |

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

**Adjusted p-values must be used!**

| 1840 | 1840 | | | |
| 1007 | 1007 | | | |
| 1542 | 1542 | | | |
| 1360 | 1360 | | | |
| 844 | 844 | | | |
| 4631 | 4631 | | | |
| 11 | 11 | 0.00539946 | 1 | 0.8888 | 1 | 3.36813 |
| 4102 | 4102 | 0.00219978 | 1 | 0.861025 | 1 | 3.35909 |
| 285 | 285 | 0.0029997 | 1 | 0.861025 | 1 | 3.35235 |
| 4716 | 4716 | 0.00439956 | 1 | 0.8888 | 1 | 3.28286 |
| 4430 | 4430 | 0.00669933 | 1 | 0.8888 | 1 | 3.2427 |
| 4398 | 4398 | 0.00559944 | 1 | 0.8888 | 1 | 3.23225 |
| 3793 | 3793 | 0.00279972 | 1 | 0.861025 | 1 | 3.22175 |
| 3462 | 3462 | 0.0042957 | 1 | 0.8888 | 1 | 3.19595 |
| 972 | 972 | 0.0039996 | 1 | 0.8888 | 1 | 3.19547 |
| 3488 | 3488 | 0.0069993 | 1 | 0.8888 | 1 | 3.12957 |
| 3992 | 3992 | 0.00849915 | 1 | 0.8888 | 1 | 3.0987 |
| 1248 | 1248 | 0.00779922 | 1 | 0.8888 | 1 | 3.09834 |

# On the problem of multiple testing

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

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

---

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

10 heads !!!

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

1000 coins

Will you still use this coin for betting?

# Studies must be hypothesis driven.

## sample classification

Can we find groups of experiments with similar gene expression profiles?

Different classes…

Molecular classification of samples



What genes are responsible for?

Co-expressing genes…

What do they have in common?

# Context: personalized medicine and what is the future

Big challenge for the pharma industry in the 21ˢᵗ century

Driven by academy and regulatory authorities

Relies or pharmacogenomic tests that properly stratifies patients

In the years coming, new tests based on different "omics" methodologies will open new avenues for new personalized drugs and treatments

# FDA approves the first predictor based on microarrays

# The MicroArray Quality Control (MAQC) Project: An FDA-Led Effort Toward Personalized Medicine

**MAQC Website: http://edkb.fda.gov/MAQC/**
*MAQC-II Objective:*
*Reaching consensus on the "best practices" (Data Analysis Protocol, DAP) in developing and validating microarray-based predictive models (classifiers) for clinical and preclinical applications.*
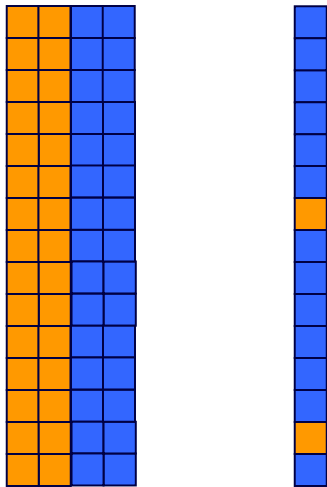
A international consortium of 36 data analysis teams submitted prediction results from 18,202 models for 6 datasets to the MAQC-II

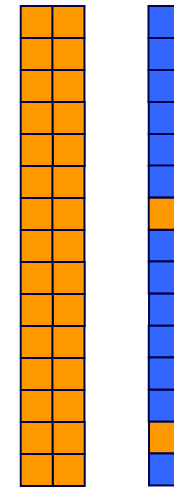# Of predictors and molecular signatures

## What is a predictor?

A  B     X

Intuitive notion:



Is X, A or B?
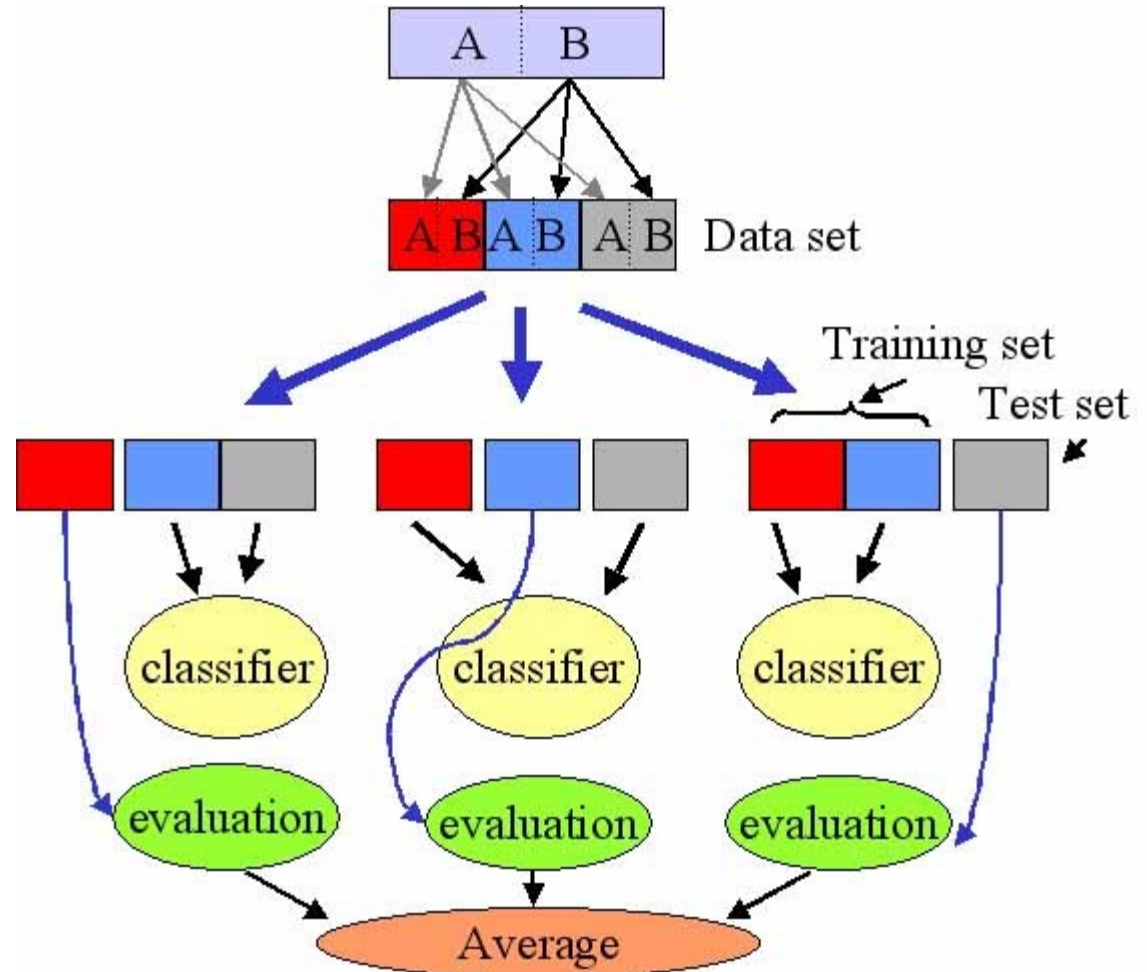
Diff (B, X) = 2     Diff (A, X) = 13

Most probably X belongs to class B

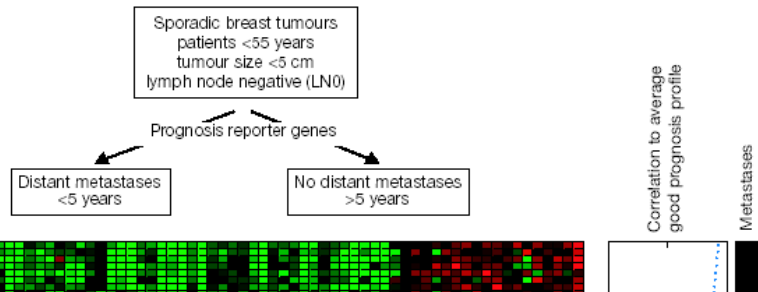Algorithms: DLDA, KNN, SVM, random forests, PAM, etc.

# Cross-validation

The efficiency of a classifier can be estimated through a process of cross-validation.

Typical are three-fold, ten-fold and leave-one-out (LOO), in case of few samples for the training

# Predictor of clinical outcome in breast cancer



Genes are arranged to their correlation eith the pronostic groups

Pronostic classifier with optimal accuracy

*van't Veer et al., Nature, 2002*

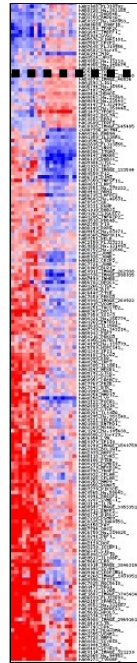# Functional profiling of genome-scale experiments in the post-genomic era



My data...

How are structured?

What are these groups?

What is this gen?

A B

?

Cell cycle...

DBs Information

GeneCards™

*Analysis*

*Functional profiling*

*Links*

# **Gene Ontology** CONSORTIUM

**http://www.geneontology.org**

- The objective of GO is to provide controlled vocabularies for the description of the molecular function, biological process and cellular component of gene products.

- These terms are to be used as attributes of gene products by collaborating databases, facilitating uniform queries across them.

- The controlled vocabularies of terms are structured

# Two-steps functional interpretation

1. Genes are selected based on their experimental values and...

2. Enrichment in functional terms is tested (FatiGO, GoMiner, etc.)

# Testing two GO terms
## (remember, we have to test thousands)

Group A

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

Group B

Biosynthesis   Other

| 6 | 4 | A |
|---|---|---|
| 2 | 8 | B |

The popular Fisher's test

Biosynthesis 60% ●  ⟷  Biosynthesis 20% ●

Sporulation   20% ●  ⤬  Sporulation   20% ●

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

# GO terms found in sets of 50 genes

| GO | Definition | p-value | Adjusted p-value |
|---|---|---|---|
| GO:0006790 | sulfur metabolism | 0.0595683 | 1 |
| GO:0042592 | homeostasis | 0.0157944 | 0.300094 |
| GO:0016265 | death | 0.116317 | 1 |
| GO:0050874 | organismal physiological process | 0.151987 | 1 |
| GO:0008152 | metabolism | 0.129865 | 1 |
| GO:0019058 | viral infectious cycle | 0.016503 | 0.181353 |
| GO:0019059 | initiation of viral infection | 0.0123062 | 0.459417 |
| GO:0009056 | catabolism | 0.0276032 | 1 |
| GO:0006766 | vitamin metabolism | 0.00875837 | 0.604328 |
| GO:0007155 | cell adhesion | 0.122953 | 1 |

Each row corresponds to a random selection of 50 genes from the *E. coli* genome, compared with respect to the rest of the genome.

GO terms in blue (p-value < 0.05 in individual test) have assymetrical distributions by chance (see adjusted p-values).

# How to test significant differences in the distribution of biological tems between groups of genes?
## FatiGO: GO-driven data analysis
Provides a statistical framework able to deal with multiple-testing hipothesis



*Al-Shahrour et al., 2004 Bioinformatics (3rd most cited paper in computing sciences. Source: ISI Web of knowledge.)*

*Al-Shahrour et al., 2005 Bioinformatics. Al-Shahrour et al., 2005 NAR*

*Al-Shahrour et al., 2006 NAR. Al-Shahrour et al., 2007 BMC Bioinformatics*

*Al-Shahrour et al., 2007 NAR*

# Compilation of tools for functional interpretation of sets of genes

| Tool | Statistical model | Correction for multiple experiments | Functional labels | Site (web-based applications) | Reference |
|------|-------------------|-------------------------------------|-------------------|-------------------------------|-----------|
| Babelomics | Fisher's exact test, t-test, Kolmogorov-Smirnov | **FDR, q-value** | **GO, KEGG, protein domains, swissprot keywords, Transfac motifs, CisRed motifs, chromosomal location, tissues, bioentities (text-mining)** | http://www.babelomics.org | **(Al-Shahrour et al., 2006; Al-Shahrour et al., 2005)** |
| BayGO | hypergeometric | bayesian | **GO** | | **(Vencio et al., 2006)** |
| DAVID / EASEonline | Fisher's exact test | Bonferroni | **GO, pathways, diseases, protein domains, interactions** | http://david.abcc.ncifcrf.gov/ | **(Dennis et al., 2003; Hosack et al., 2003)** |
| FatiGO+ | Fisher's exact test | step-down minP, FDR | **GO, KEGG, protein domains, swissprot keywords, Transfac motifs, CisRed motifs, chromosomal location, tissues** | http://www.fatigo.org | **(Al-Shahrour et al., 2004)** |
| FuncSpec | hypergeometric | **Bonferroni** | **GO, phenotypes, protein interactions, etc. (only for yeast)** | http://funspec.med.utoronto.ca/ | **(Robinson et al., 2002)** |
| GeneMerge | hypergeometric | Bonferroni | **GO, KEGG, chromosomal location, other.** | http://genemerge.bioteam.net/ | **(Castillo-Davis & Hartl, 2003)** |
| GO:TermFinder | hypergeometric | Bonferroni | **GO** | | **(Boyle et al., 2004)** |
| GoMiner | Fisher's exact test | FDR | **GO** | | **(Zeeberg et al., 2003; Zeeberg et al., 2005)** |
| GOstat | **X2** Fisher's exact test | FDR, Holm | **GO** | http://gostat.wehi.edu.au/ | **(Beissbarth & Speed, 2004)** |
| GoSurfer | **X2** | q-value | **GO** | | **(Zhong et al., 2004)** |
| GOToolBox | hypergeometric, binomial, Fisher's exact test | Bonferroni | **GO** | http://crfb.univ-mrs.fr/GOToolBox/index.php | **(Martin et al., 2004)** |
| Ontology Traverser | hypergeometric | FDR | **GO** | http://franklin.imgen.bcm.tmc.edu/rho-old/services/OntologyTraverser/ | **(Young et al., 2005)** |
| Onto-Tools | X2, binomial, hypergeometric Fisher's exact test | Sidak, Holm, Bonferroni, FDR | **GO, KEGG** | http://vortex.cs.wayne.edu/projects.htm | **(Draghici et al., 2003; Khatri et al., 2005)** |
| FuncAssociate | Fisher's exact test | -- | **GO** | http://llama.med.harvard.edu/cgi/func/funcassociate | **(Berriz et al., 2003)** |
| GOTM | hypergeometric | -- | **GO** | http://bioinfo.vanderbilt.edu/gotm/ | **(Zhang et al., 2004)** |
| CLENCH | Hypergeometric, X2, binomial | -- | **GO (only for *A. thaliana*)** | -- | **(Shah & Fedoroff, 2004)** |

# Understanding why genes differ in their expression between two different conditions

Limphomas from mature lymphocytes (LB) and precursor T-lymphocyte (PTL).

Genes differentially expressed, selected among the ~7000 genes in the CNIO oncochip

Genes differentially expressed among both groups were mainly related to immune response (activated in mature lymphocytes)

*Martinez et al.,* Clinical Cancer Research. **10**: 4971-4982.

# Biological processes shown by the genes differentially expressed among PTL-LB

| | Cluster Query | Cluster Reference |
|---|---|---|
| Total number of initial genes: | 162 | 4764 |
| Total number of genes no repeated: | 129 | 4731 |
| Total number of Cluster IDs retired - their currents Cluster IDs | 7 - 23 | 449 - 1627 |
| Total number of genes no repeated with current Cluster IDs: | 145 | 5909 |
| Total number of genes no repeated with GO at level 3 and biological_process: | 88 | 2610 |
| Total number of genes no repeated with GO but NOT at level 3 and ontology | | |
| Total number of genes no repeated without GO annotated: | | |

Gene Ontology Term

response to external stimulus — 36.36%
11.65%

response to stress — 21.59%
6.86%

signal transduction — 39.77%
26.05%

cell motility — 9.09%
3.79%

resistance to pathogenic bacteria — 1.14%
0.04%

viral replication — 1.14%
0.15%

cell death — 9.09%
5.75%

regulation of gene expression, epigenetic — 1.14%
0.19%

0.1702  0.9912  1  1

0.1806  0.9940  1  1

## Obvious? NO

1) You now know that there are no other co-variables (e.g. age, sex, etc)

2) If you do not have previously a strong biological hypothesis, now you have an explanation

# Weaknesses of the two-steps, functional enrichment approach

Low sensitivity of conventional gene selection methods

A B

**A**

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

B

17 with normal tolerance to glucose (**NTG**)

*(Mootha et al., 2003)*

Instability of molecular signatures. Variable selection with microarray data can lead to many solutions that are equally good from the point of view of prediction rates, but that share few common genes (Ein-Dor 2006 PNAS)

Platform comparison. There are still some concerns with the cross-platform coherence of results. Paradoxically, despite the fact that gene-by-gene results are not always the same, the biological themes emerging from the different platforms are increasingly consistent (Bammler 2005 Nat Methods)

# Functional enrichment approach reproduces pre-genomics paradigms



Context and cooperation between genes is ignored

# So, what is wrong with what we are doing?

We seek for the functions activated/deactivated in our experiment

To find them we firstly seek for genes activated/deactivated one at a time (independently)

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

Therefore… is all wrong with this. The test we conduct is implicitly answering a question different to the one we want to ask.

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

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

# The true proxies of function

Are we asking the proper questions?

Why do we think in terms of genes?

What are the real bricks that account for the cellular behaviour and for the phenotype or the response to stimulus represented in our experiment?  The genes or other higher level units?



Function

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

Experiment



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

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**)

Therefore, why not to observe the elephant?

# Functional genomics.
## Historic perspective and future

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



- Classification of phenotypes / experiments. Sensitivity

- Selection of differentially expressed genes Specificity

- Biological roles the genes are carrying out in the cell. Interpretation

- Reformulating the questions. Are we asking the proper questions? What are the real bricks that account for the cellular behaviour and for the phenotype or the response to environmental stimuli?  The genes or other higher level units?

# Cooperative activity of genes can be detected and related to a macroscopic observation



A B

GO$_1$ GO$_2$ GO$_3$

statistic

−

+

**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.)
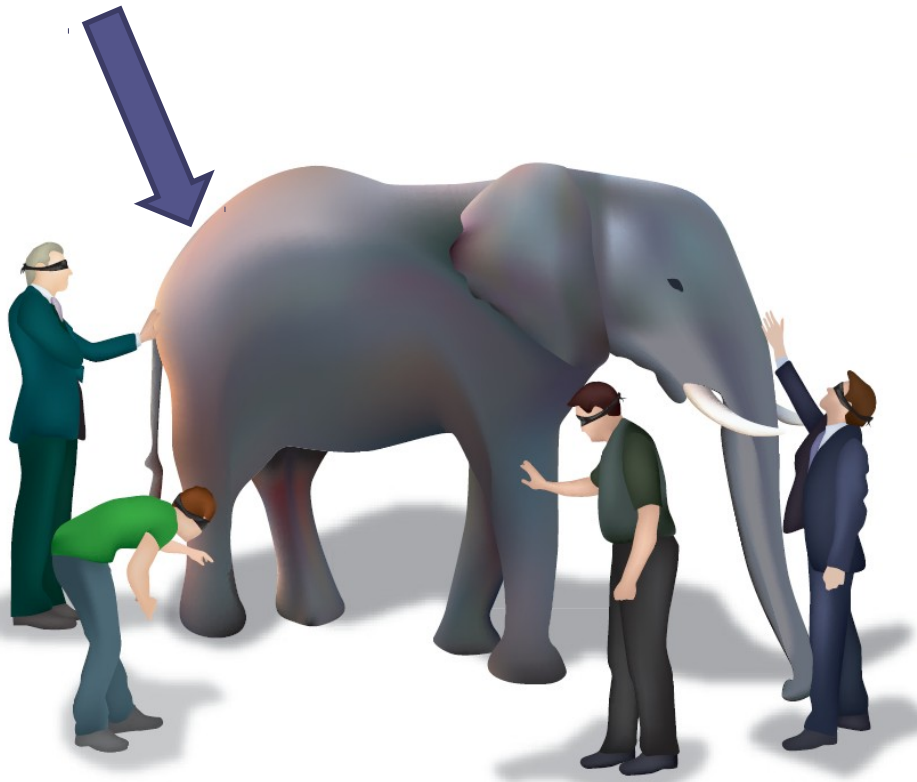
**Distribution of GO**: Rows GO**1**, GO**2** and GO**3** represent the position of the genes belonging to three different GO terms across the ranking.

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



Significantly over-expressed in B

t-test with two tails.

p<0.05

Significantly over-expressed in A

A B

GO₁ GO₂

statistic

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

Classes expressed as blocks in A and B

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

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

A B    GO₁  GO₂



Significantly over-expressed in B

t-test with two tails.

p<0.05

statistic

Significantly over-expressed in A

The main problem is that the two-steps approach cannot distinguish between these two different cases.

We put both sides of the partition into two bags and destroy the structure of the data.

|      | up | down |
|------|----|------|
| GO   | 3  | 9    |
| no GO| 0  | 25   |

Same contingency table for GO₁ and GO₂ !!

# Gene-set enrichment methods



GSEA

A B

FatiScan

Measure ES for each gene set

Permute class labels (1,000 times)

statistic

Gene set

background

**Independent of the experimental design**

# FatiScan, a segmentation test, provides an easy approach to directly testing functional terms



E.g., term $GO_2$, partition $p_1$

| | up | down |
|---|---|---|
| GO | 4 | 6 |
| no GO | 2 | 30 |

GOs can be directly tested by a segmentation test. A series of partitions of the list are performed (**p1, p2, p3...**) and the GO terms for each functional class in the upper part are compared to the corresponding ones in the lower part by a Fisher test. Asymmetrical distributions of terms towards the extremes of the list will produce significant values of the test.

Finally, p-values are adjusted by FDR

*Al-Shahrour et al., 2005 Bioinformatics*

# Obtaining significant results



For each GO term (T), different partitions (P) are tested.

TxP p-values of tests to be adjusted for multiple testing.

Empirical results suggest that 20 to 50 partitions optimally find significant asymmetrical distributions of terms

*Al-Shahrour et al., 2005 Bioinformatics*

# Case study: functional differences in a class comparison experiment

A B

**A**

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

B

17 with normal tolerance to glucose (**NTG**)

*(Mootha et al., 2003)*

FatiScan

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

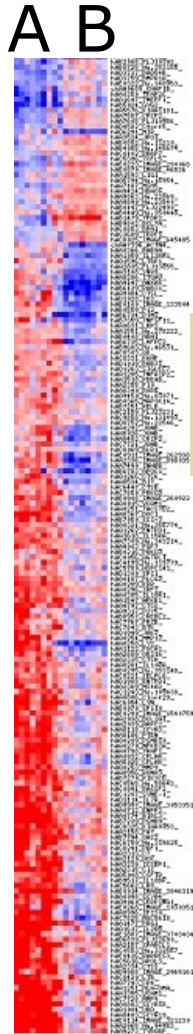| Healthy vs diabetic | Functional class | Repository | | |
| --- | --- | --- | --- | --- |
| | | GO | KEGG | Swissprot keyword |
| Up-regulated | Oxidative phosphorylation | X | X | |
| | ATP synthesis | | X | |
| | Ribosome | | X | |
| | Ubiquinone | | | X |
| | Ribosomal protein | | | X |
| | Ribonucleoprotein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

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

# Beyond discrete variables: Survival data

**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

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
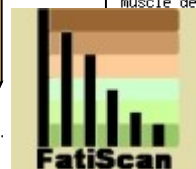
**- Survival**

| 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**



1 Over-represented terms associated with top values    2 Under-represented terms associated wit
4 Under-represented terms associated with top values    3 Over-represented terms associated with

%genes annotated

antigen processing
antigen presentation
M phase of mitotic cell cycle
regulation of cell cycle
cellular localization
macromolecule metabolism
primary metabolism
nervous system development
muscle development
regulation of organismal physiologic
muscle contraction
epidermis development

%genes annotated

Gene Ontology : biological process

# Comparison of gene set methods at a glance

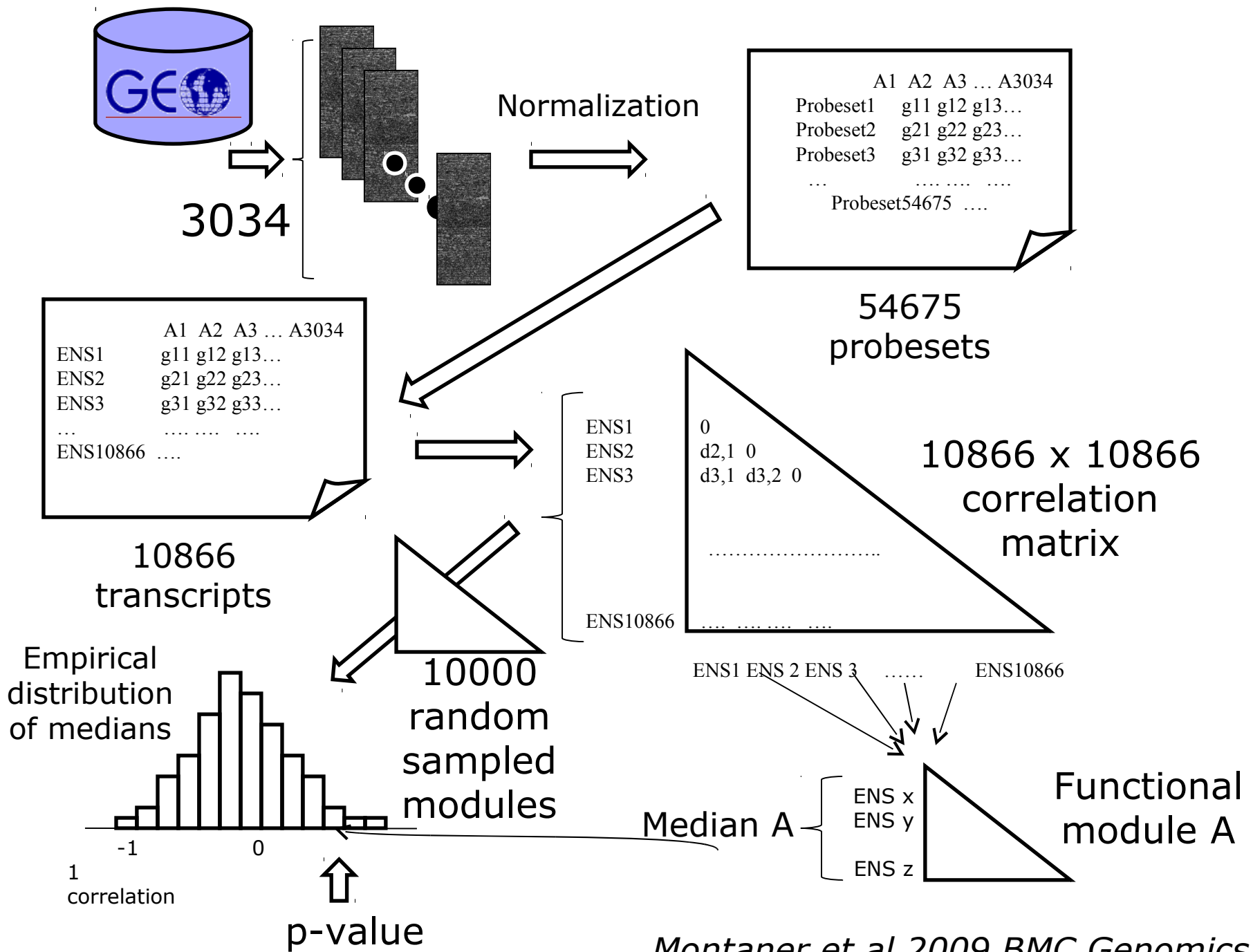| Healthy vs diabetic | Functional class | Repository | | | | Method | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | GO | KEGG | Swissprot keyword | Defined in GSEA | FatiScan | GSEA | PAGE | Tian et al. |
| Up-regulated | Oxidative phosphorylation | + | + | | + | yes | yes | yes | yes |
| | ATP synthesis | | + | | | yes | - | - | - |
| | Ribosome | | + | | | yes | - | - | - |
| | Ubiquinone | | | + | | yes | - | - | - |
| | Ribosomal protein | | | + | | yes | - | - | - |
| | Ribonucleoprotein | | | + | | yes | - | - | - |
| | Mitochondrion | + | | + | + | yes | yes | yes | yes |
| | Transit peptide | | | + | | yes | - | - | - |
| | Nucleotide biosynthesis | + | | | + | yes | yes | yes | yes |
| | NADH dehidrogenase (ubiquinone) activity | + | | | | yes | - | - | - |
| | Nuclease activity | + | | | | yes | - | - | - |
| Dow-regulated | Insulin signalling pathway | | + | | | yes | - | - | - |

Terms from distinc repositories, reported by different methods in the diabetes dataset (Mootha et al., 2003)

# Still one more problem…
# are functional modules defining
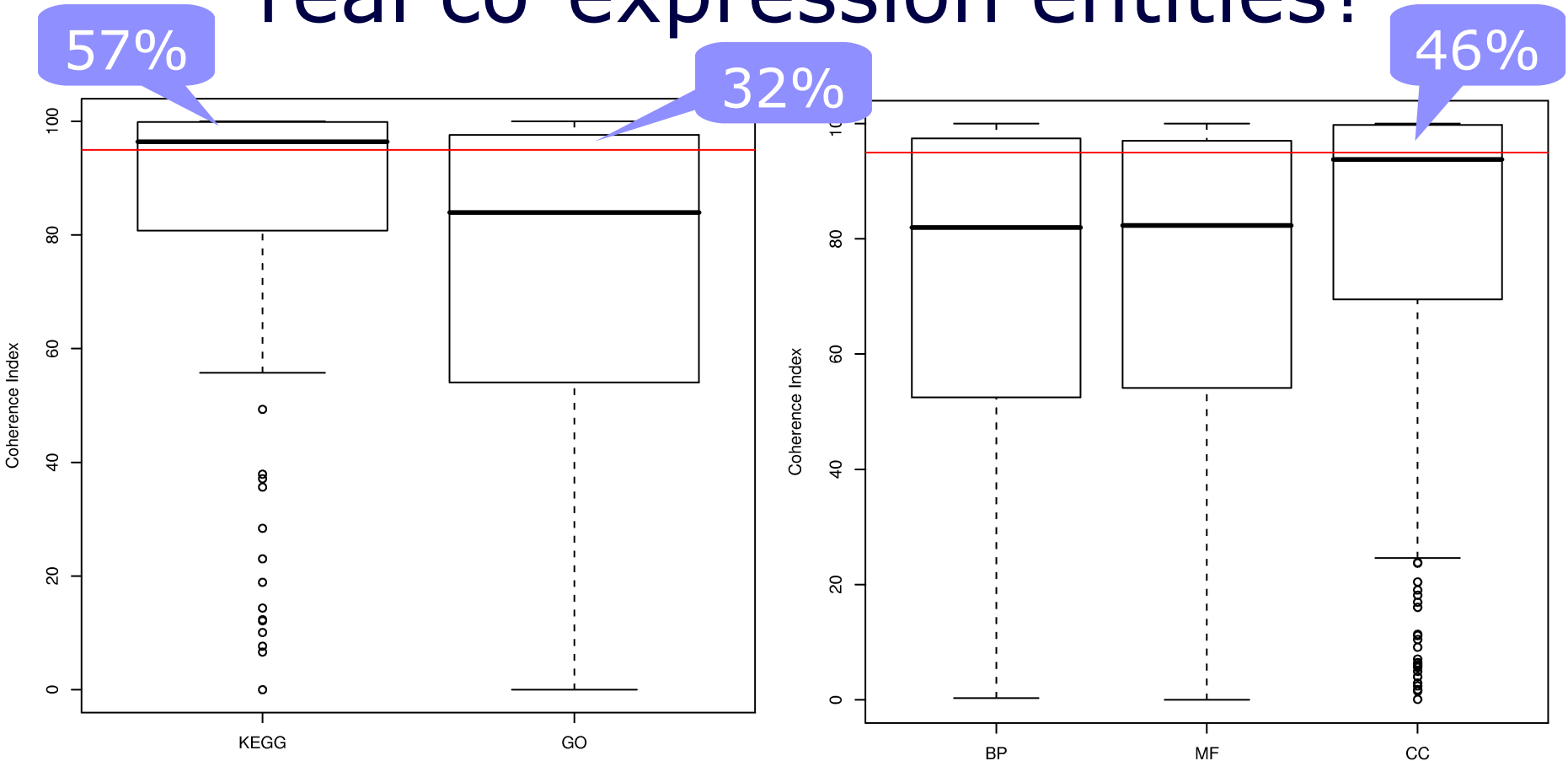# real co-expression classes?

Not a naïve and trivial question.

Functional enrichment methods and gene set analysis methods rely on the assumption that the modules tested do **coexpress**

There are tens of thousands GO terms and hundreds of KEGG pathways

*Montaner et al 2009 BMC Genomics*

# But are functional modules defining real co-expression entities?



Coherence index: (1-p-value)*100.
CI > 95% means internal co-expression significantly higher than random co-expression

*Montaner et al 2009 BMC Genomics*

# Weighting gene module membership by co-expression

| KEGG pathway | Unweighted test | | | Weighted test | | |
|---|---|---|---|---|---|---|
| | statistic | p-value | adjusted p-value | statistic | p-value | Adjusted p-value |
| Caprolactam degradation | 2.741 | 0.059 | 0.289 | 3.124 | 0.003 | 0.034 |
| Cell cycle | 2.588 | 0 | 0 | 2.711 | 0 | 0 |
| Maturity onset diabetes of the young | 2.517 | 0.075 | 0.289 | 2.734 | 0.008 | 0.034 |
| RNA polymerase | 2.497 | 0.077 | 0.289 | 2.657 | 0.009 | 0.034 |
| One carbon pool by folate | 2.497 | 0.077 | 0.289 | 2.766 | 0.007 | 0.034 |
| Urea cycle and metabolism of amino groups | 2.497 | 0.077 | 0.289 | 2.674 | 0.009 | 0.034 |
| Heparan sulfate biosynthesis | 2.478 | 0.078 | 0.289 | 2.818 | 0.006 | 0.034 |
| Alanine and aspartate metabolism | 2.386 | 0.087 | 0.289 | 2.497 | 0.012 | 0.04 |
| Amyotrophic lateral sclerosis (ALS) | 2.386 | 0.087 | 0.289 | 2.91 | 0.005 | 0.034 |
| beta-Alanine metabolism | 2.318 | 0.094 | 0.289 | 2.668 | 0.009 | 0.034 |
| Basal transcription factors | 2.125 | 0.116 | 0.298 | 2.431 | 0.014 | 0.04 |
| Benzoate degradation via CoA ligation | 2.072 | 0.123 | 0.298 | 2.468 | 0.013 | 0.04 |
| Limonene and pinene degradation | 1.986 | 0.135 | 0.298 | 2.306 | 0.018 | 0.048 |

Very simple weight schema:
W=2 if correlation is positive
W=0.5 if negative
W=1 if not in the class



*Montaner et al 2009 BMC Genomics*

# Future directions



Testing hierarchies is better
Functions and pathways are correlated.
Testing models will increase our sensitivity
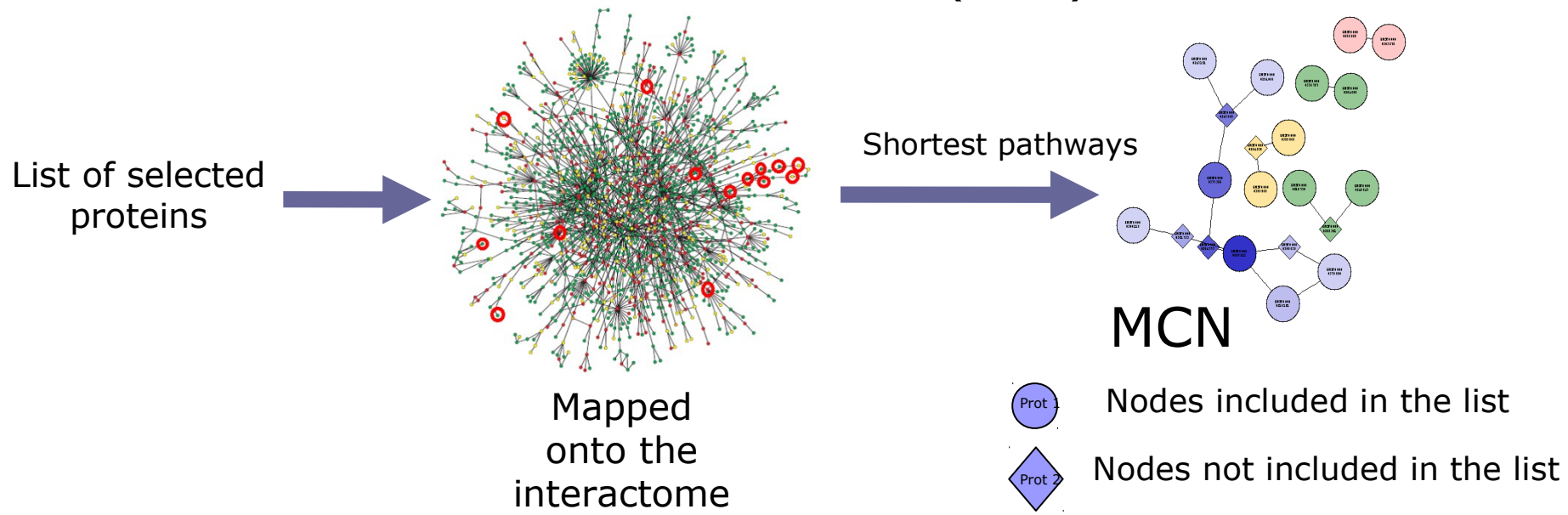
Pathways are not categorical variables

In general (systems) biology is behind. Our questions must be inspired directly by biology
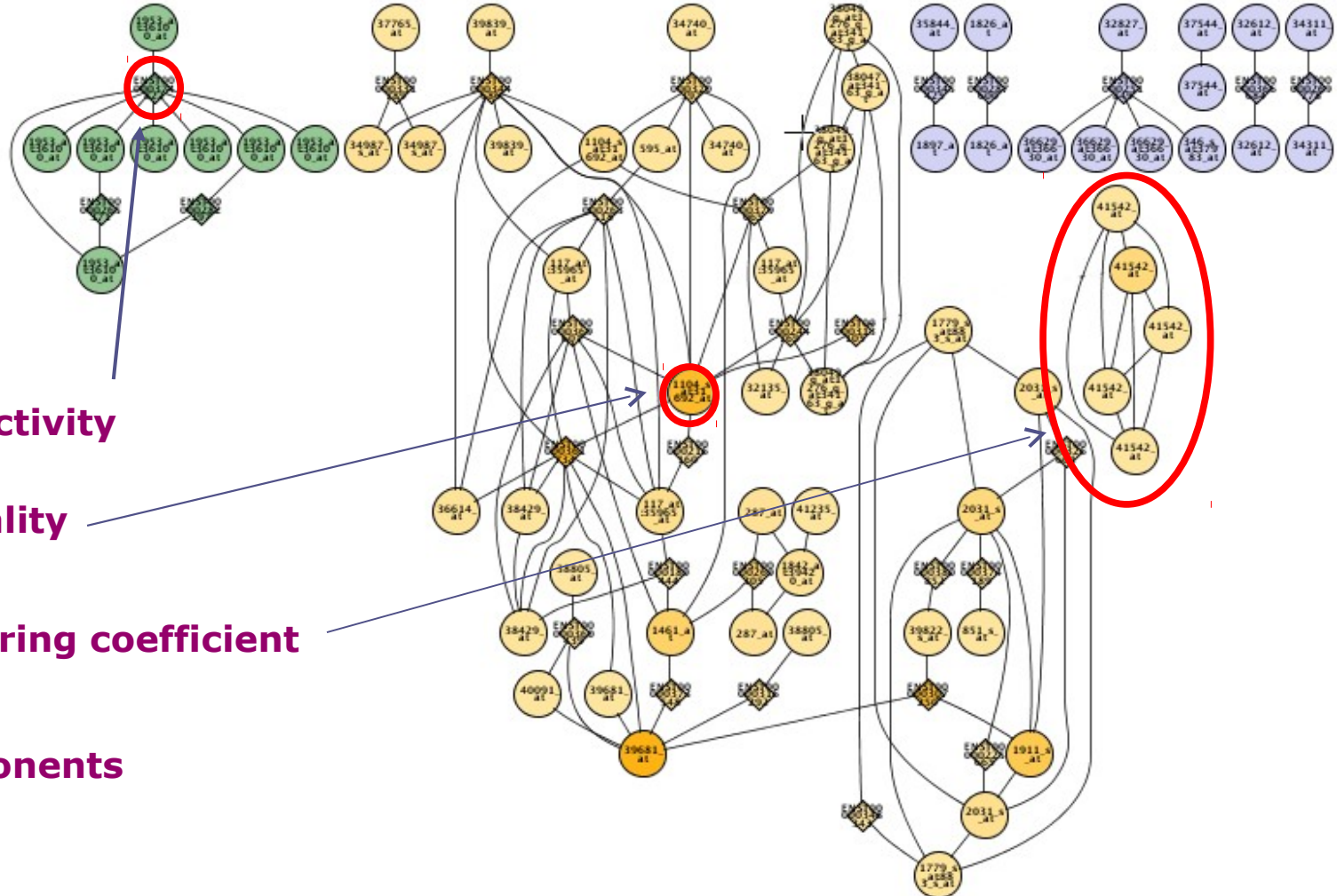
# Protein-protein interaction networks

## 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)



List of selected proteins

Mapped onto the interactome

Shortest pathways

MCN

Nodes included in the list

Nodes not included in the list

# Network parameters



1. **Connectivity**

2. **Centrality**

3. **Clustering coefficient**

4. **Components**

# Evaluation of the Minimum Connection Network (MCN)

**Parameters to evaluate:** connectivity, centrality , clustering coeficient, components

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



**Network parameters**

| Clustering Coeff: | Betweenness: | Connections: |
| List1 > Random pval=1e-04 | List1 > Random pval=2e-04 | List1 > Random pval=0 |

Number of components [95% confidence interval]:   List1: 38 [46-79]
Number of components with more than 1 node:   List1: 8
Number of Bicomponents:   List1: 41
Articulation points:   List1: 56

# Significant connections

# Babelomics



Since May 1st, Babelomics 4.0

# Some numbers

451 papers cite GEPAS (215 are SOTA cites)

632 papers cite Babelomics (442 are  FatiGO cites)

*(source ISI Web of Knowledge, May 2010)*

More than 150,000 experiments analysed during the last year.

More than 1000 experiments per day.

# Tools for gene expression analysis

# Tools for functional profiling



Percentage of the number of citations per tool (%)



The most cited tools

# Other tools (non-commertial)

To cover more specific analysis requirements

Bioconductor:
http://www.bioconductor.org

BRB tools:
http://linus.nci.nih.gov/BRB-ArrayTools.html

TM4 (MeV):
http://www.tm4.org/mev.html

Easier use

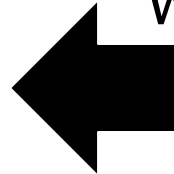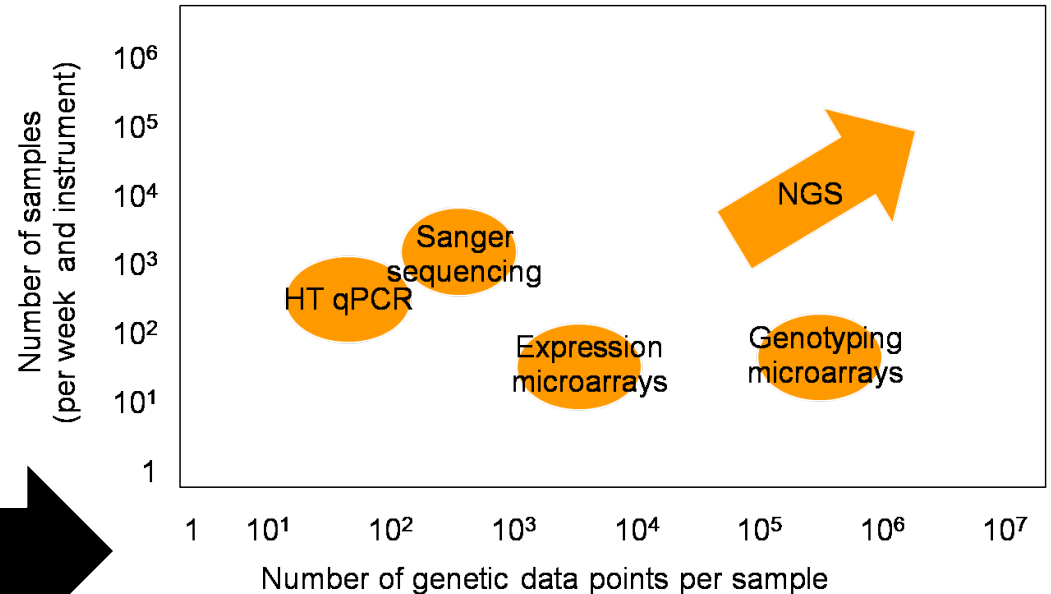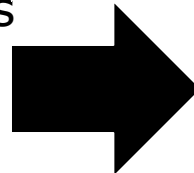Power of tool

# What is next?

# Next generation technology is here



Observed and expected trend of publications in which NGS is being used.

Relative throughput of the different technologies. NGS emerges with a potential of data production that will, eventually wipe out conventional HT technologies in the years coming

# Some of the most common applications of NGS

RNA-seq
Transcriptomics:
Quantitative
Descriptive
(alternative
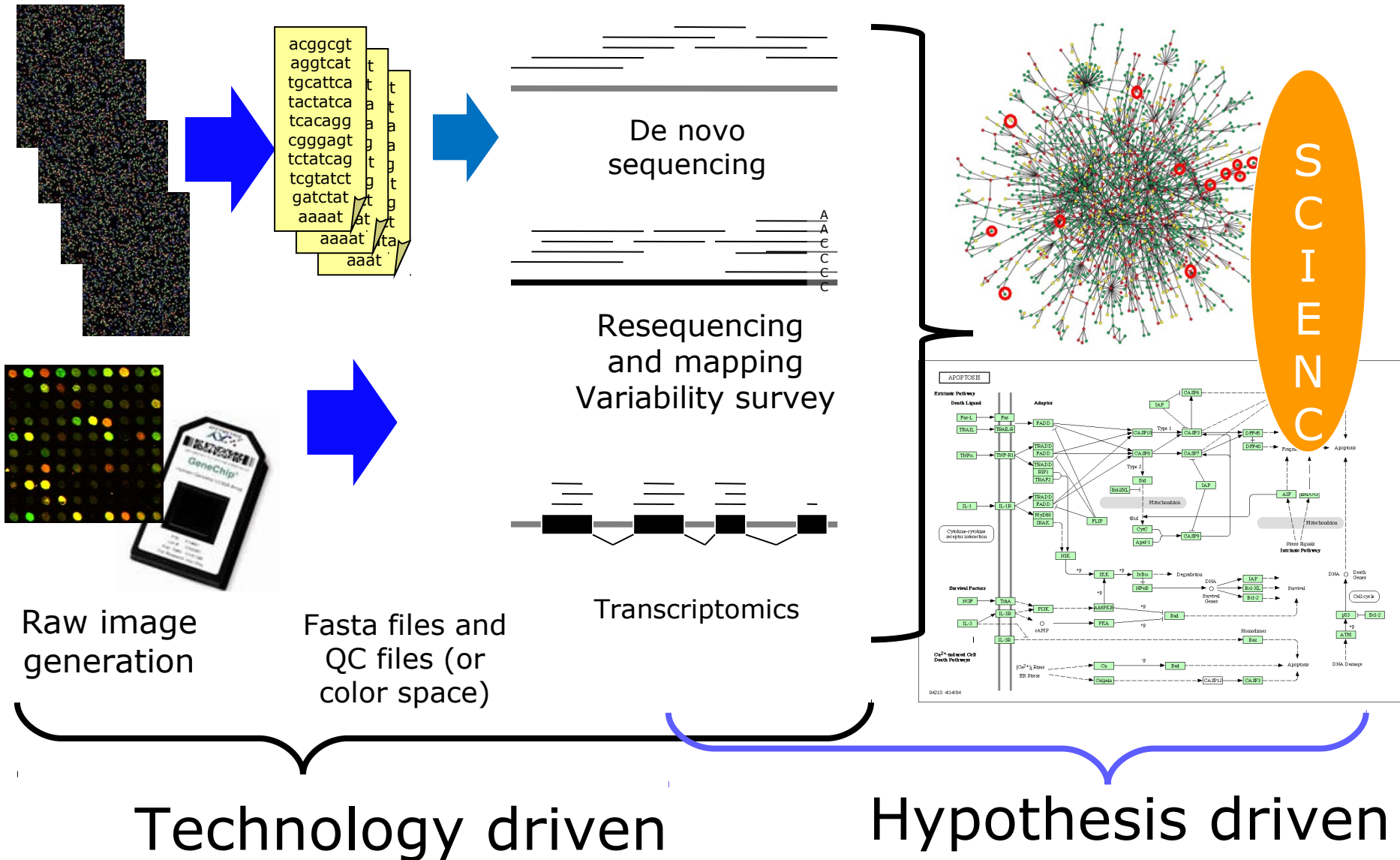splicing)
miRNA

Resequencing:
Mutation calling
Profiling

*De novo*
sequencing

Chip-seq
Protein-DNA interactions
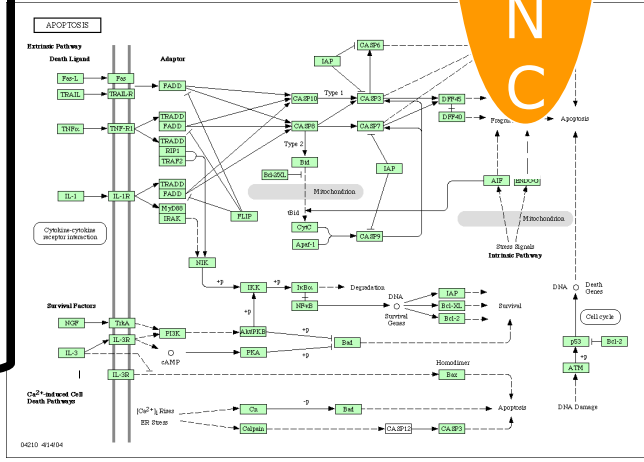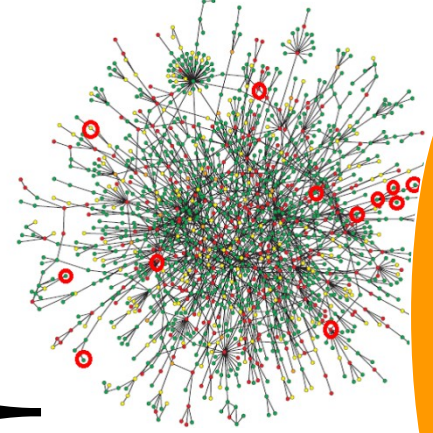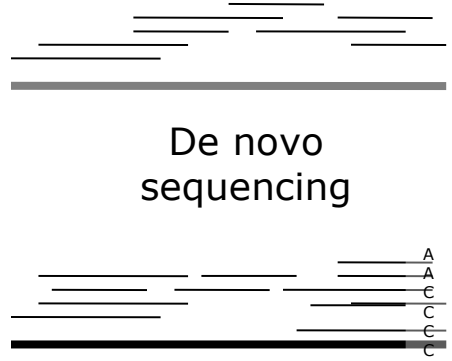Active transcription
factor binding sites

Copy number
variation

Metagenomics
Metatranscriptomics

# Pipeline general of analysis



acggcgt
aggtcat
tgcattca
tactatca
tcacagg
cgggagt
tctatcag
tcgtatct
gatctat
aaaat
aaaat
aaat

De novo
sequencing

A
A
C
C
C
C

Resequencing
and mapping
Variability survey

Transcriptomics

S
C
I
E
N
C

APOPTOSIS

Raw image
generation

Fasta files and
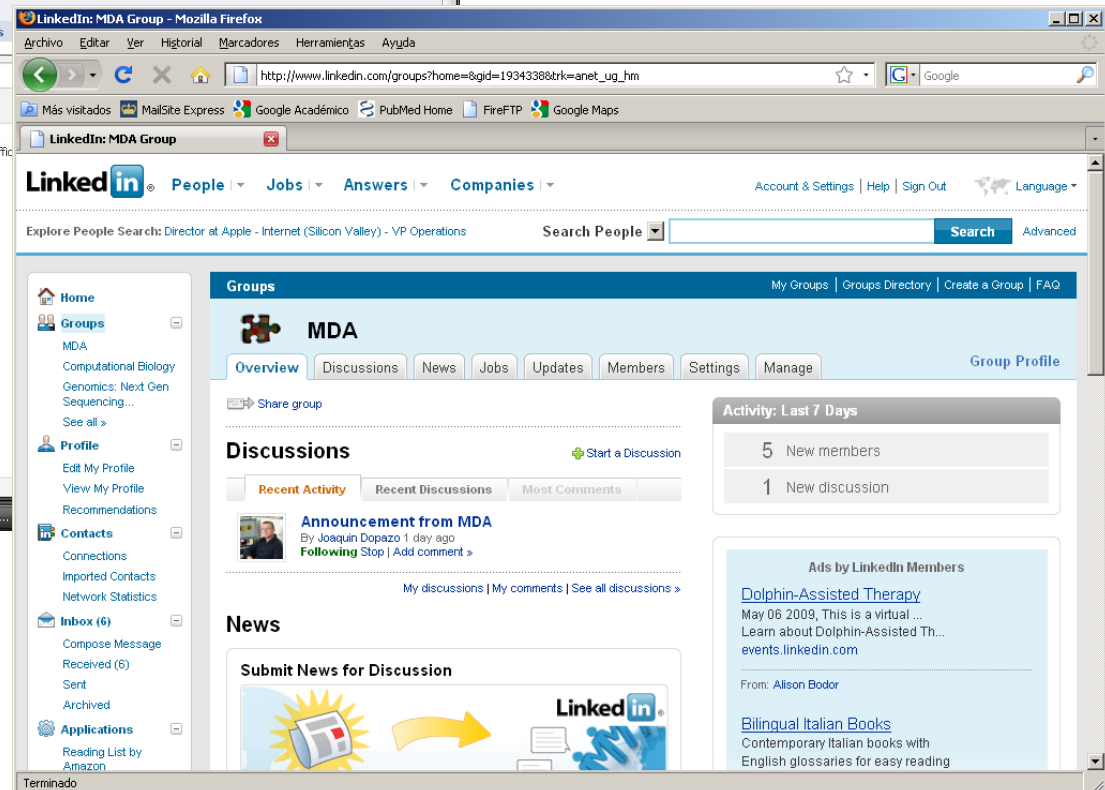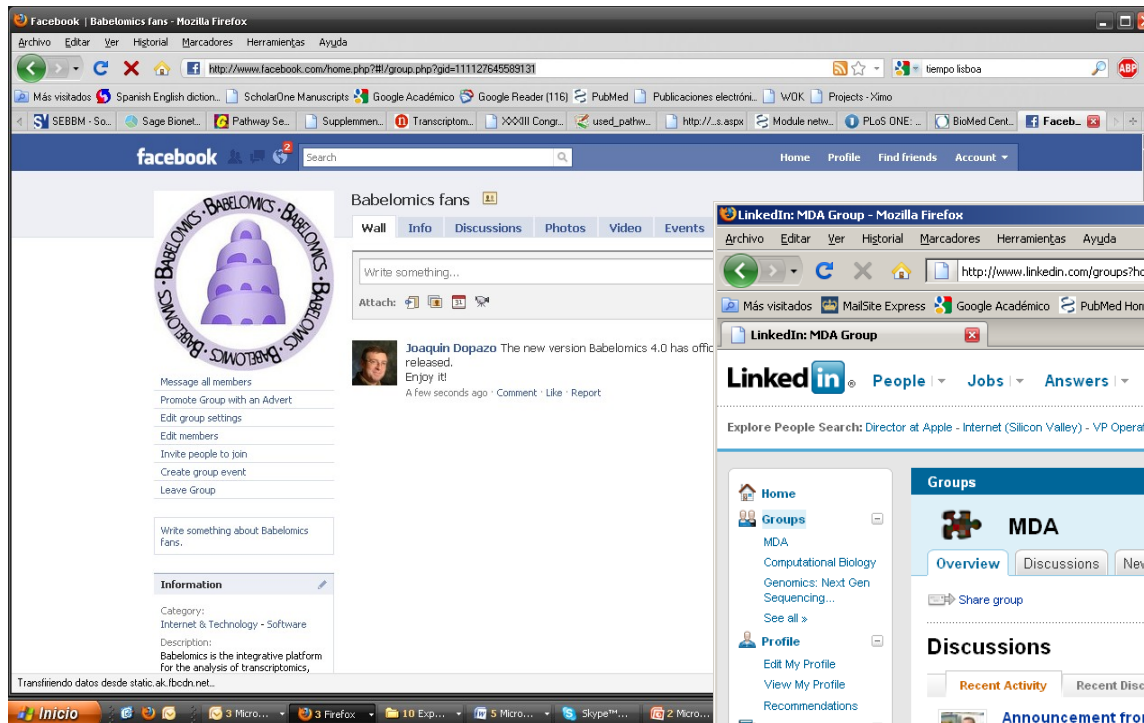QC files (or
color space)

Technology driven

Hypothesis driven

# SOCIAL:
## MDA group in Linked-in
## Babelomics group in Facebook

# 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 Bau
Emidio Capriotti
Jose Carbonell
Ana Conesa
Adriana Cucchi
Hernán Dopazo
Pablo Escobar
Francisco García
Stefan Goetz
Martina Marbà
Marc Martí
Ignacio Medina
Pablo Minguez
David Montaner
Marina Naval
Luis Pulido
Javier Santoyo
Patricia Sebastian
François Serra
Sonia Tarazona
Joaquín Tárraga

## ...the INB, National Institute of Bioinformatics (Functional Genomics Node) and the CIBERER Network of Centers for Rare Diseases