# MDA VII

## Joaquín Dopazo

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

**http://www.babelomics.org**
**http://bioinfo.cipf.es**

*Valencia, 21 March 2011*

# Who we are

## The Bioinformatics and Genomics Department at the Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain, and…
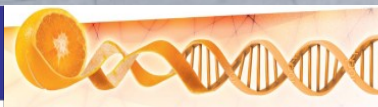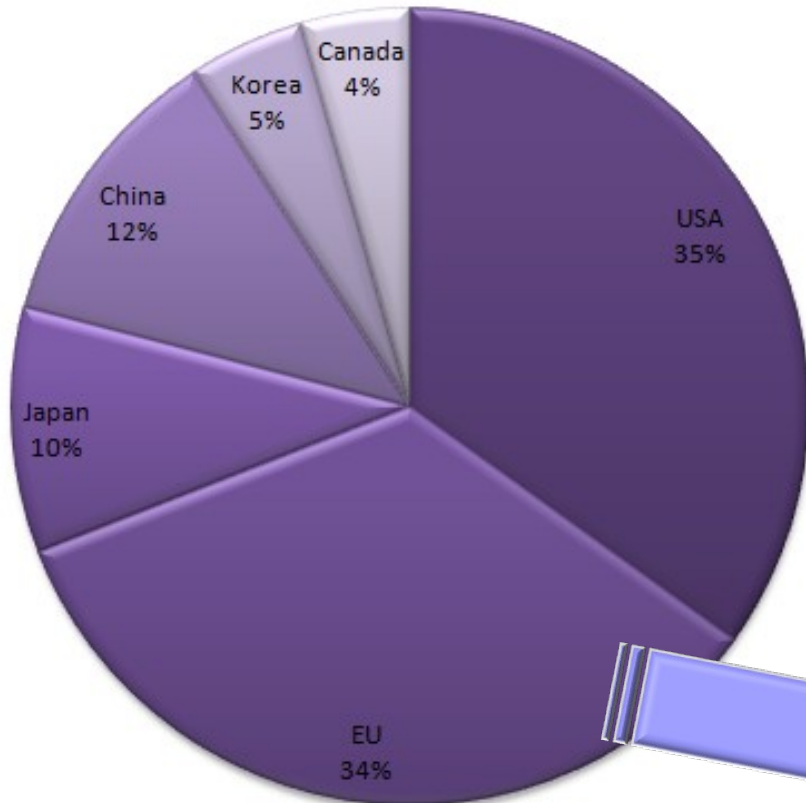
**…the INB, National Institute of Bioinformatics (Functional Genomics Node) and the CIBERER Network of Centers for Rare Diseases, and… …the Medical Genome Project (Sevilla)**
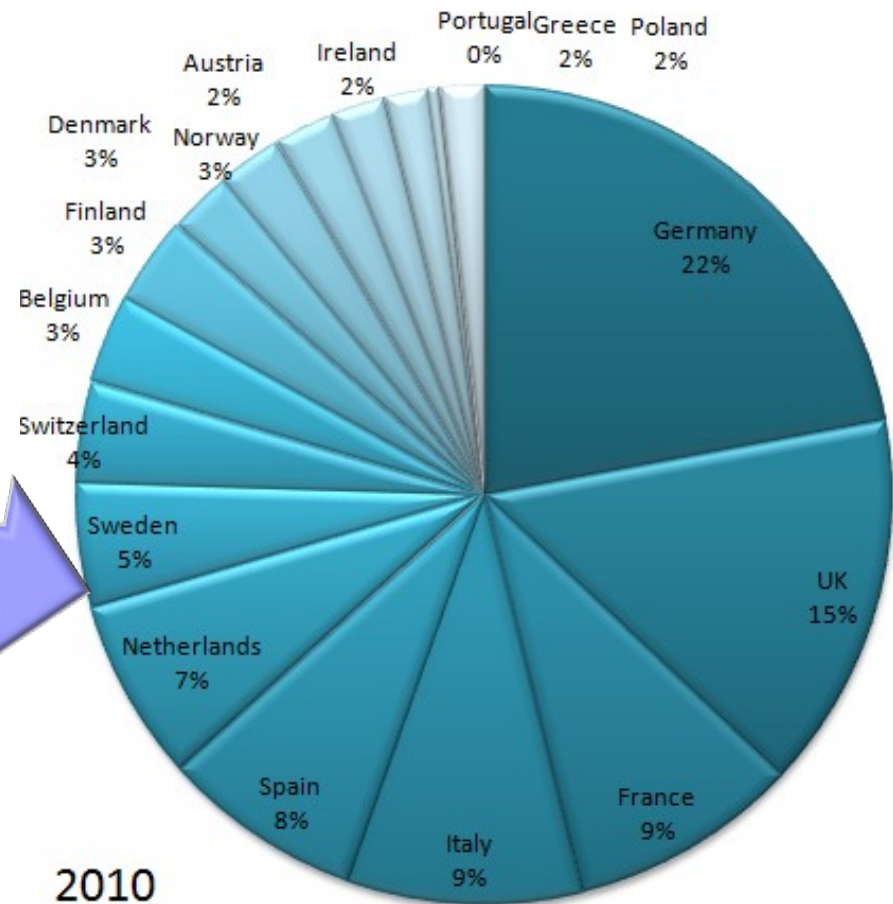
# Some bibliographic data Microarray publications



2010 Europe

2010 Worlwide

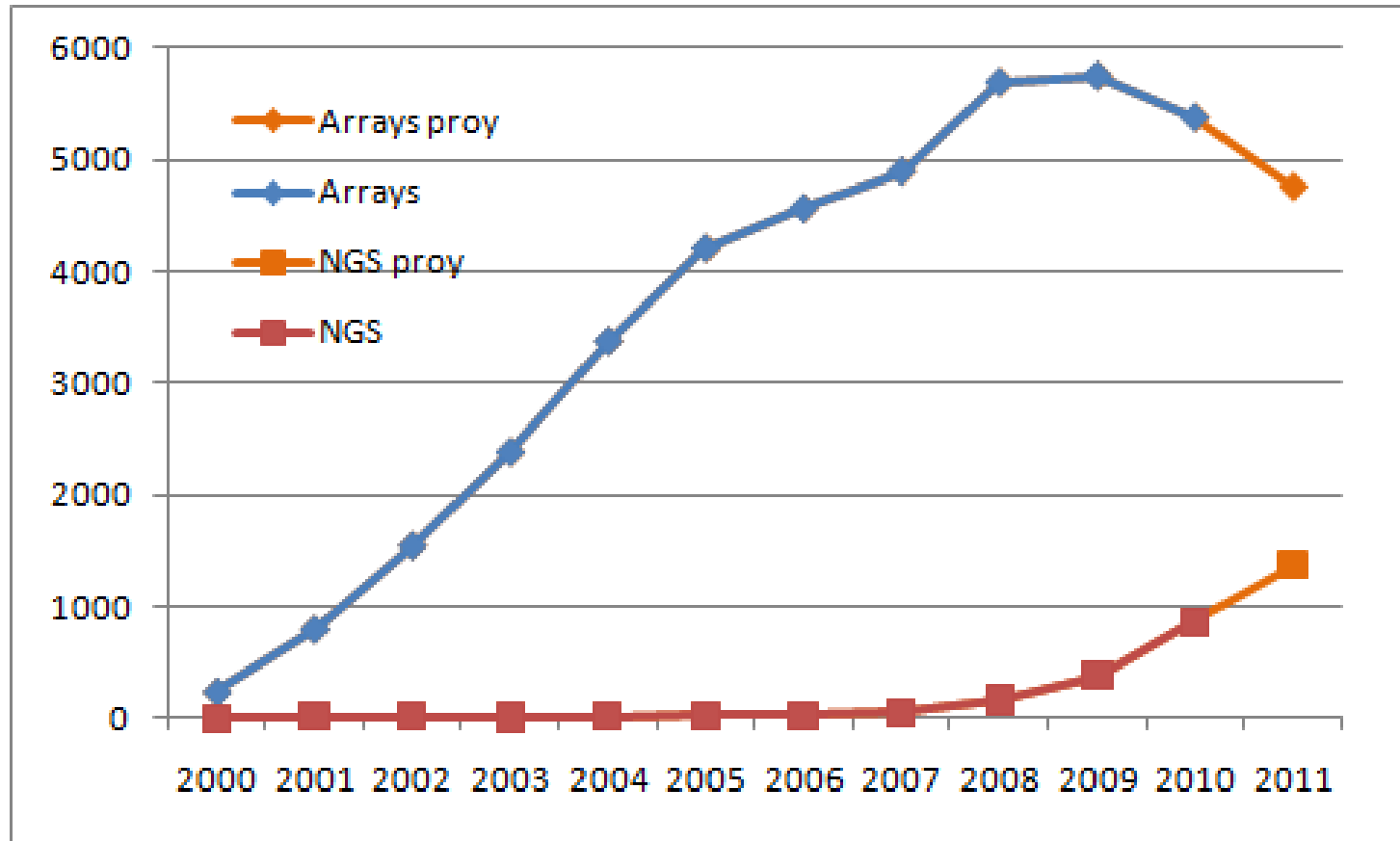**Source Pubmed. Query: 2009[Entrez Date] AND country[Affiliation]AND microarray[Title/Abstract]**

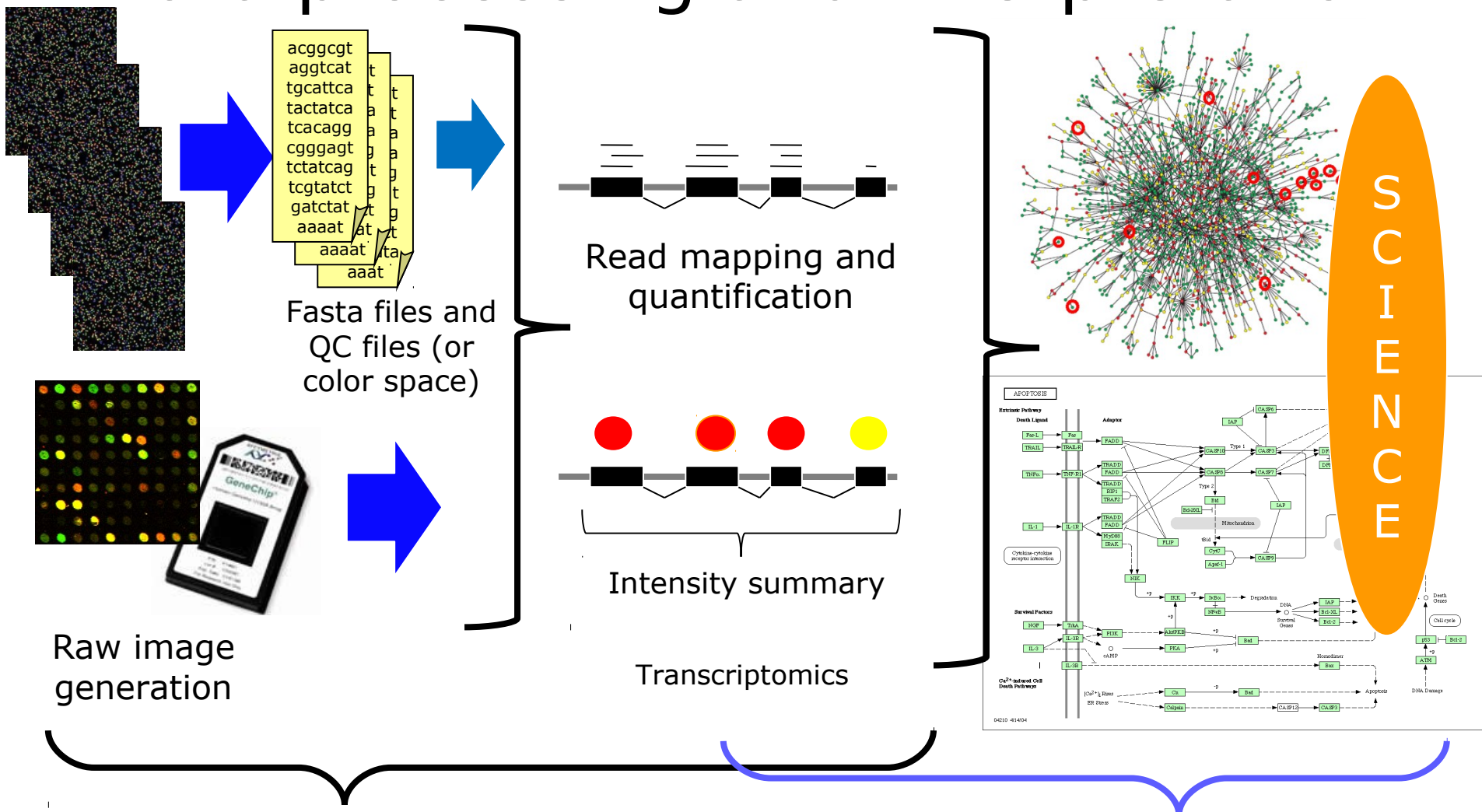# Evolution of the papers published in microarray and next gen technologies



**Source Pubmed. Query:** "high-throughput sequencing"[Title/Abstract] OR "next generation sequencing"[Title/Abstract] OR "rna seq"[Title/Abstract]) AND year[Publication Date]
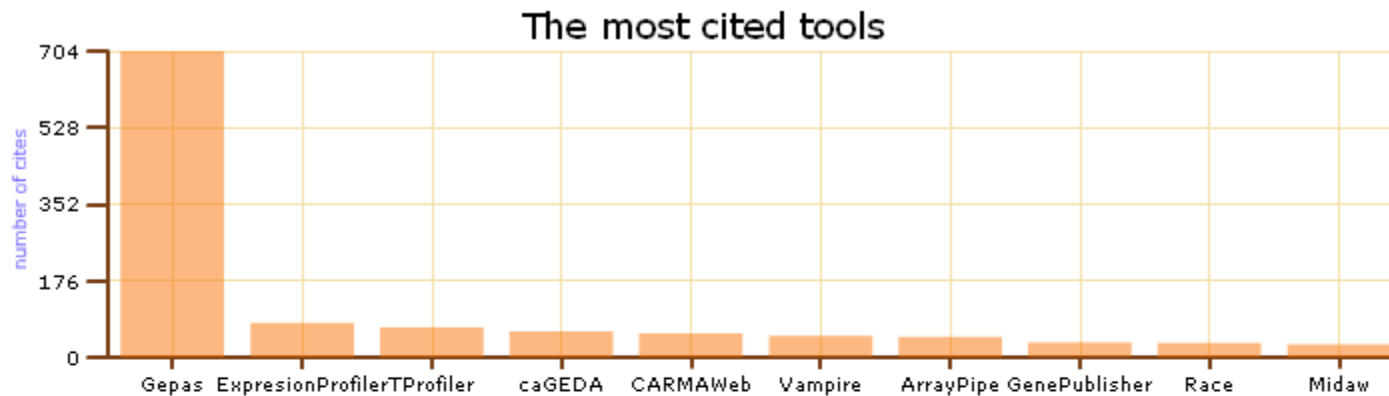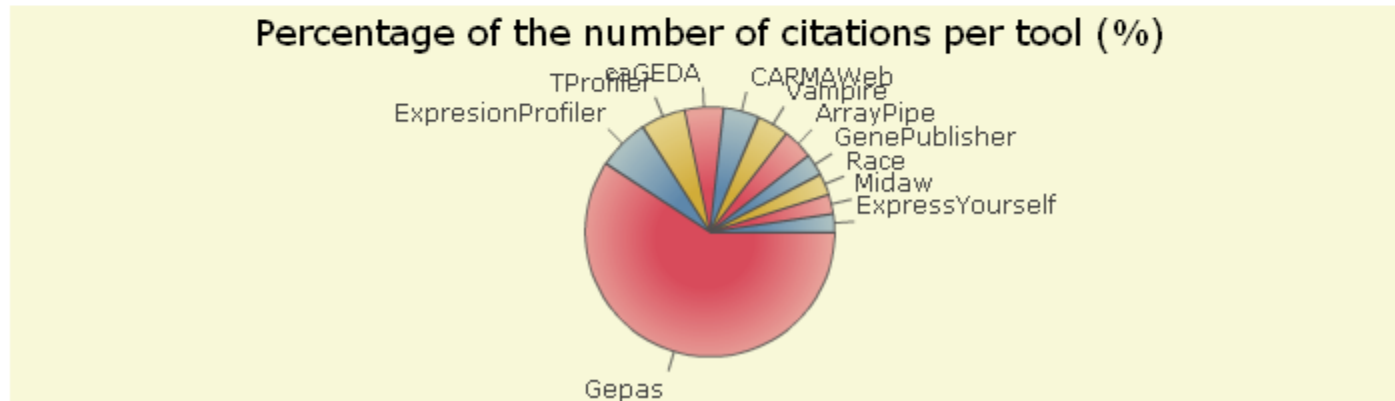**Projections 2011** based on January and February

# Genomic data, the double challenge:
## Data processing and interpretation



Fasta files and QC files (or color space)

Raw image generation

Read mapping and quantification

Intensity summary

Transcriptomics

S C I E N C E

Technology driven

Hypothesis driven

# 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

# 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.
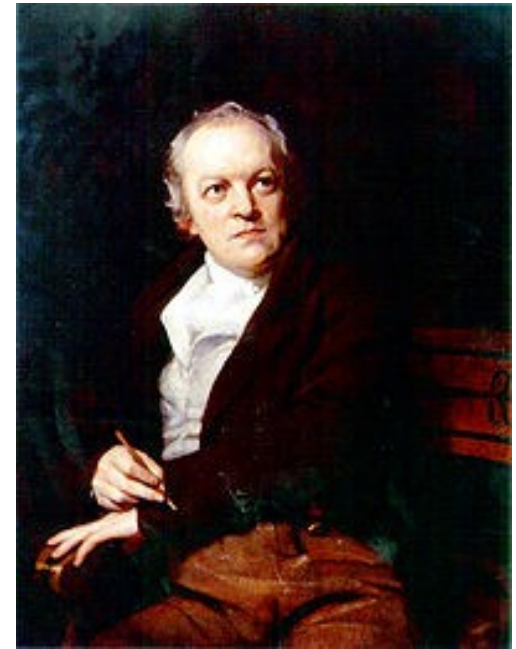
More than 1000 experiments per day.

# Structure of the course

## Theoretical

| Monday | Tuesday | Wednesday | Thursday | Friday |
|---|---|---|---|---|
| Course Reception Course Overview | Transcript Assembly | Statistical Reminder | Differential expression for microarrays | Gene-Set Methodologies |
| Coffee break | Coffee break | Coffee break | Coffee break | Coffee break |
| Introduction to Linux Introduction to NGS Technologies | Extracting RNAseq Counts in NGS Studies | Microarray Data Normalization | Predictors | Biological Networks |
| LUNCH | LUNCH | LUNCH | LUNCH | LUNCH |
| NGS Data Preprocessing | Differential Expression in RNAseq Studies | Genomic SNP data analysis | Clustering Methods | Closing |
| Coffee break | Coffee break | Coffee break | Coffee break | |
| NGS Read Mapping | Functional annotation | Genome Wide Association Studies | Biological Databases | |

## Theoretical and Hands-on:

# Background
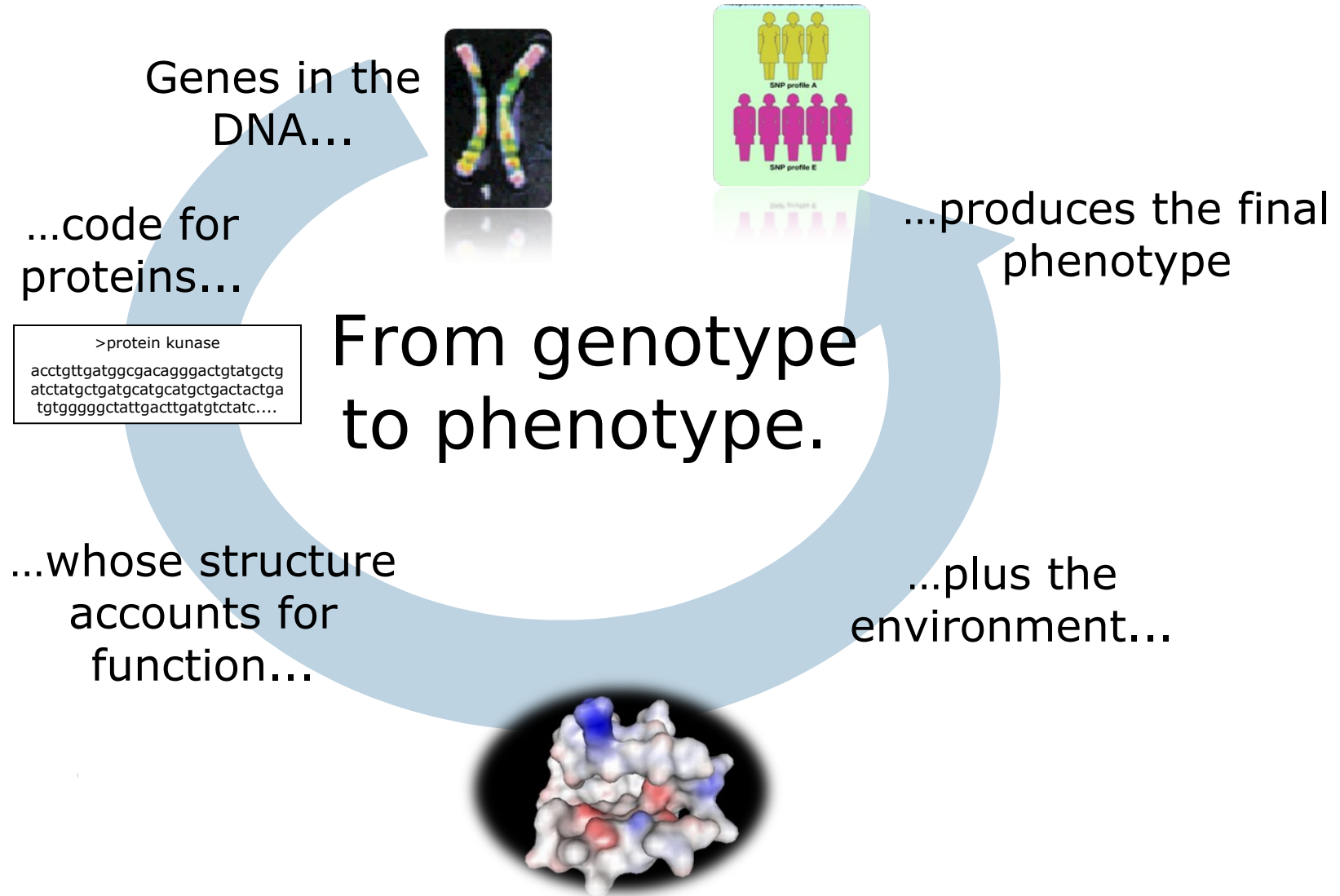
## The road of excess leads to the palace of wisdom

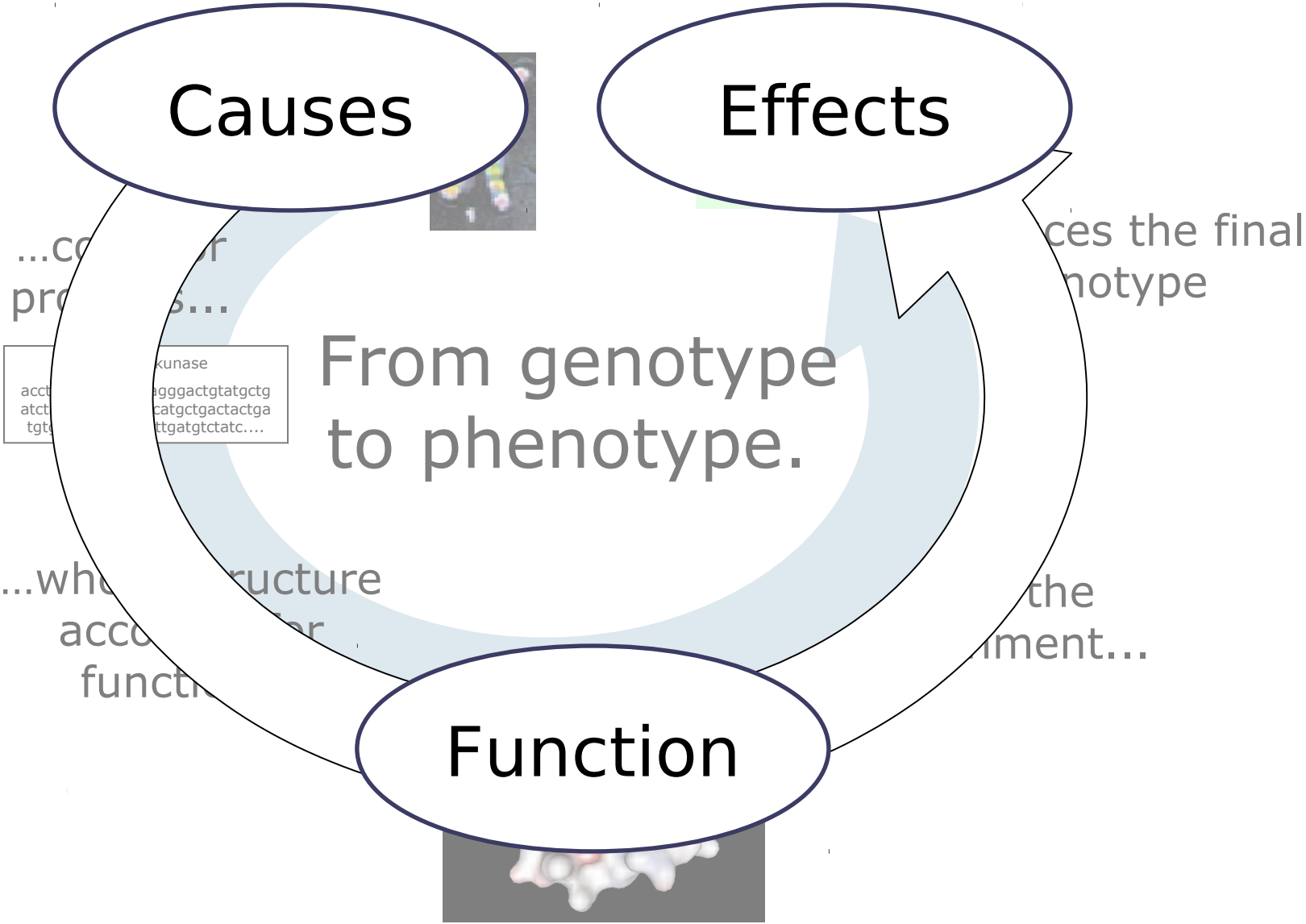(*William Blake, 28 November 1757 – 12 August 1827, poet, painter, and printmaker*)

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



Genes in the DNA...

...code for proteins...

>protein kunase

acctgttgatggcgacagggactgtatgctg
atctatgctgatgcatgcatgctgactactga
tgtgggggctattgacttgatgtctatc....

From genotype to phenotype.

...produces the final phenotype

...whose structure accounts for function...

...plus the environment...

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

Next Generation Sequencing
$10^9$bp per round

Genes in
the DNA...

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

...with its
complex
variability...

12 million SNPs in
exonic regions

From genotype
to phenotype.

(in the post-genomics scenario)

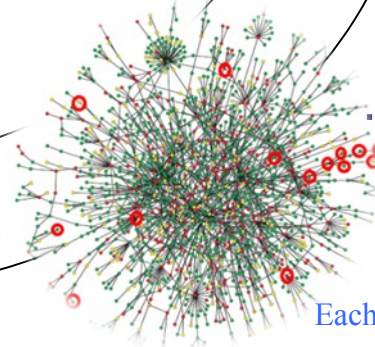...conforming complex
interaction networks...

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

...in cooperation
with other
proteins...

...code for
proteins...

Each protein has an average
of 8 interactions
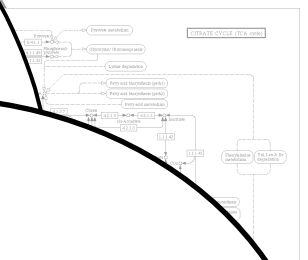
That undergo post-translational
modifications, somatic
recombination...
100K-500K proteins

...that account for
function if...

# Holistic approach. Causes and effects remain essentially the same. The concept of function has changed

Causes

Effects

…whose final effect configures the phenotype…

…with its complex variability

Half a mi... ...ants between p... ...iduals

From genotype to phenotype

…when ar ir n

All *science* is either physics or *stamp collecting*

*Ernst Rutherford*

(in the functional ...omics scen...

Function (modules of proteins)

...plex ...ks...

...tion ...er ...ns...

…code for proteins…

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

...s an average ...f 8 interactions

…that account for function if…

# Technologies for transcriptomics and genotyping and the corresponding bioinformatics support
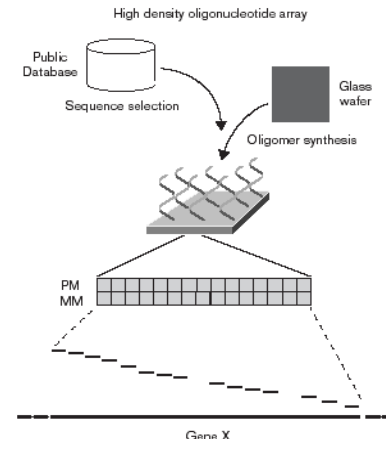
Microarray

User-friendly Babelomics

R and scripting
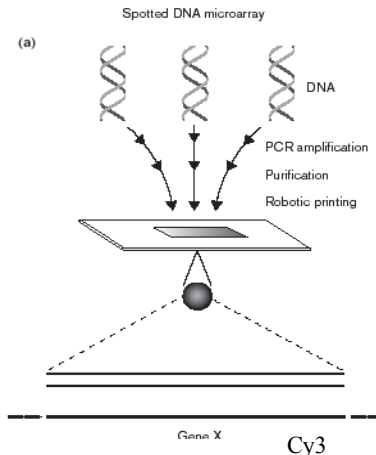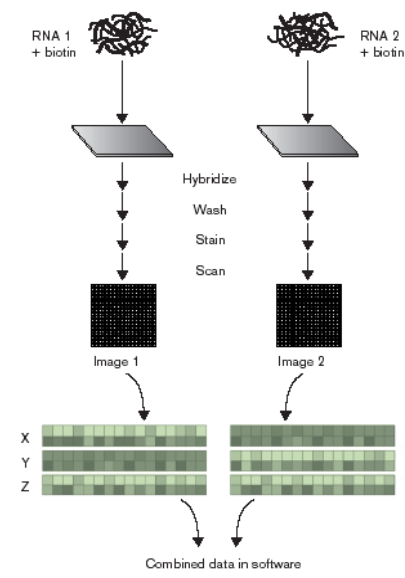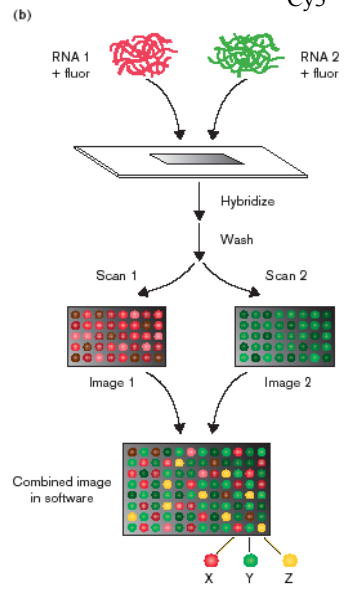
NGS

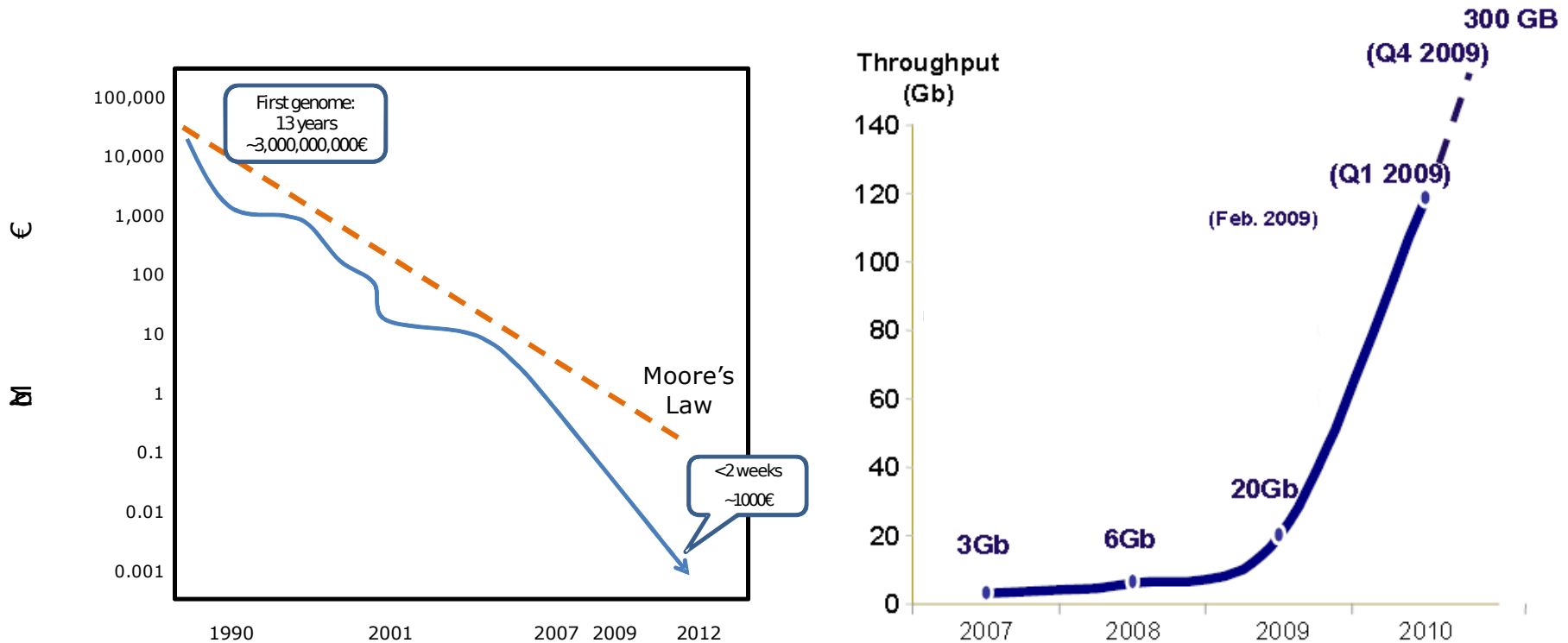# DNA expression microarrays. Strategies of hybridization



Cy5

Cy3

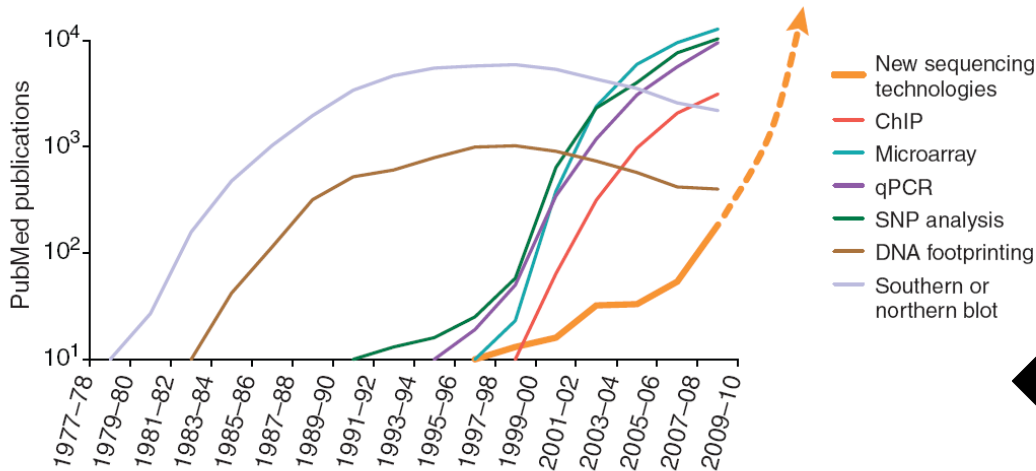Competitive hybridization (two colors)

One color

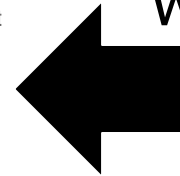# Next generation sequencing technologies are here



The cost goes down, while the amount of data to manage and its complexity raise exponentially.
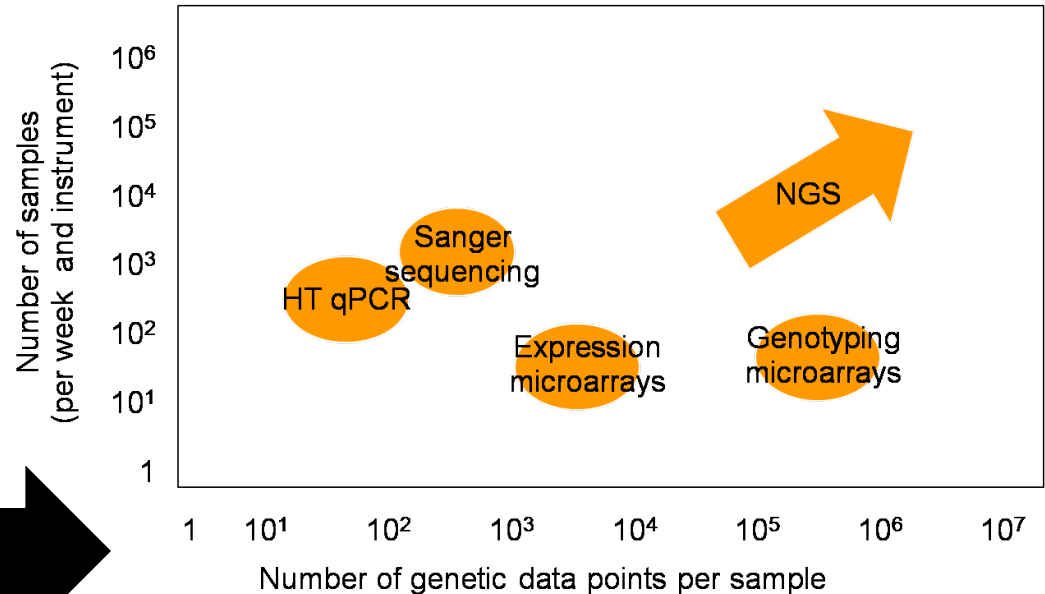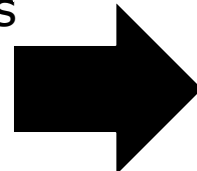
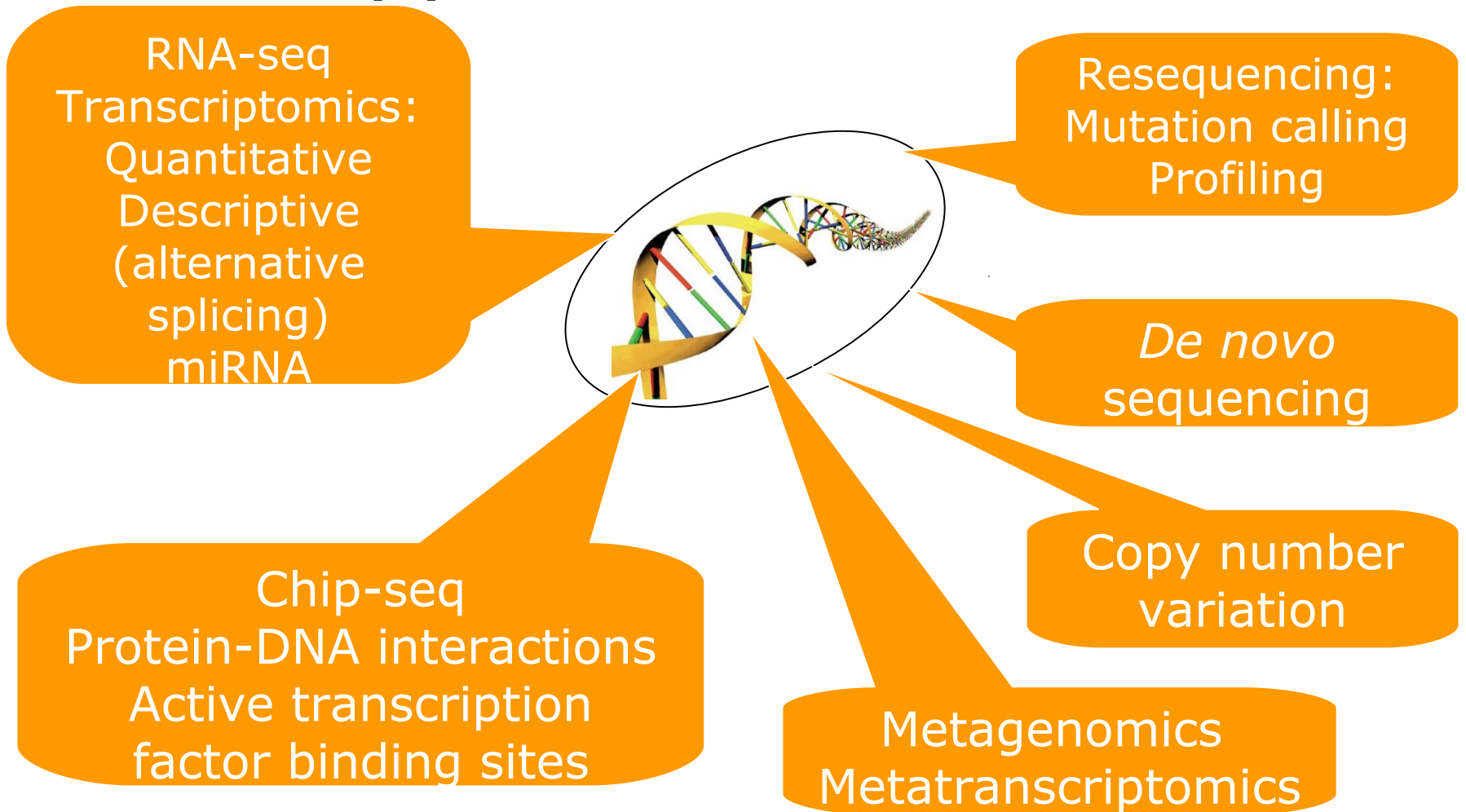# Next generation sequencing technologies are here



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
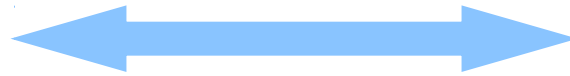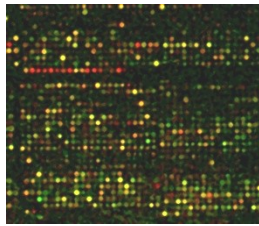
# Some of the most common applications of NGS

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

# 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 we 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 we 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)

# Primary analysis

• Transform images corresponding to hybridization intensities (microarrays) or to read counts (NGS) into numbers

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

# Secondary analysis

Once the measurements are in a common, comparable scale the results can be studied.

Different studies can be made that include class discovery, classification, gene selection, variant calling, etc.

# Studies must be hypothesis driven.

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

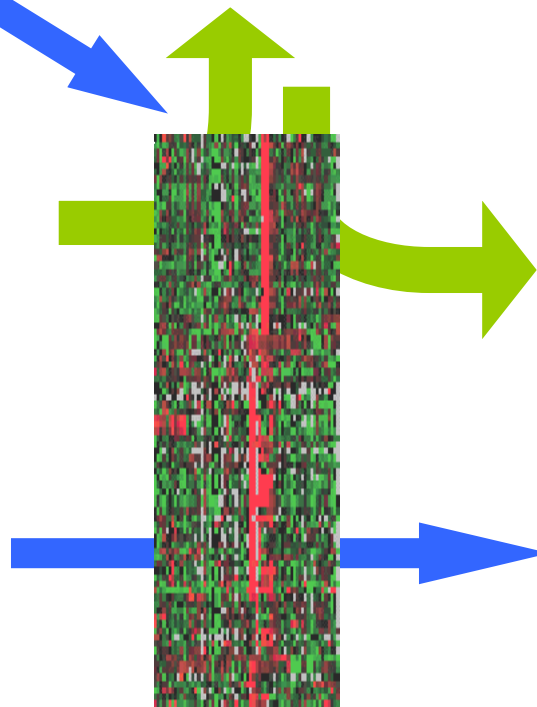Can we find groups of experiments with similar gene expression profiles?

Different classes…

| | |
|---|---|
| ▬ (blue) | Unsupervised |
| ▬ (green) | 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?
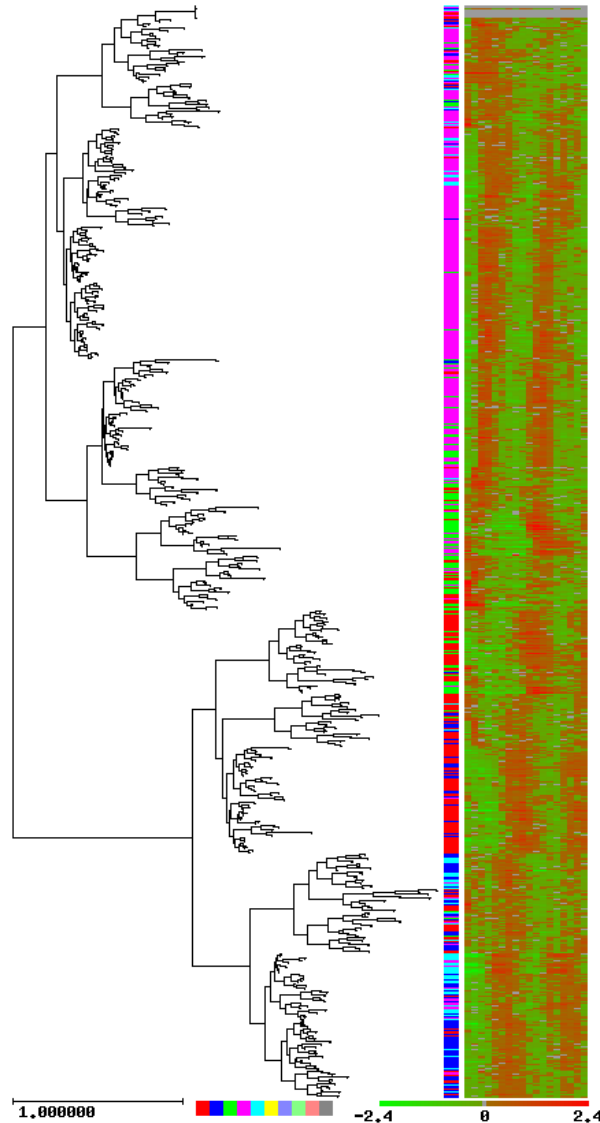
Co-expressing genes...

- 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?

And

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

Class prediction

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

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

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

---



10 heads !!!

1000 coins

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

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

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?

# 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*

# Genotyping

# Genotyping to find mutations associated to diseases

## The simplest case: monogenic disease



Controls

Cases

| | | |
|---|---|---|
| Gene A | 1 0 0 0 0 0 0 0 0 0 0 0 | 0 0 0 0 1 0 0 0 0 0 0 0 |
| Gene B | 0 0 0 1 0 0 0 0 0 0 0 0 | 0 0 0 0 0 0 0 0 0 0 0 0 |
| Gene C | 0 0 0 0 0 0 0 0 0 0 0 0 | 1 1 1 1 1 1 1 1 1 1 1 1 |
| Gene D | 0 0 0 0 0 0 0 0 0 0 1 0 | 0 0 0 0 0 0 0 0 0 0 0 0 |
| Gene E | 0 0 0 0 0 1 0 0 0 0 0 0 | 0 0 0 0 0 0 1 0 0 0 0 |

# The real life in GWAS

Our analysis of Hirschsprung´s disease

54 trios of short-segment Hirschsprung´s disease Affy 6.0 (1million SNPs)

Conventional TDT test reports only 4 significant SNPs mapping only on one gene: RET, already knowk to be associated to the disease

This is not a matter of sample size: an example of GWAS in Breast Cancer.

The CGEMS initiative. (Hunter et al. Nat Genet 2007)

1145 cases 1142 controls. Affy 500K

Conventional association test reports only significant 4 SNPs mapping only on one gene: FGFR2

Conclusions: conventional tests are not providing much resolution. What is the problem with them? Are there solutions?

# Clear individual gene associations are difficult to find in multifactorial diseases

Controls

Cases



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

# 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...

GeneCards™

DBs Information

*Analysis*

*Functional profiling*

*Links*

# 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

| Cell cycle | 60% | ● | ⟷ | Cell cycle | 20% | ● |
| Apoptosis | 20% | ● | ⟷ | Apoptosis | 20% | ● |

Genes in group A have significantly to do with cell cycle, but not with apoptosis.

# 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*

# 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

0    20    40    60

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%

0.1702  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)

This testing strategy is very strict in controlling:

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

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

But, we forget about

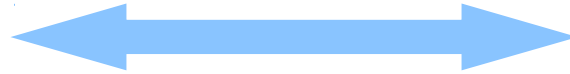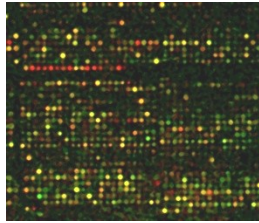Type III error : get the right answer having asked the wrong question!

The testing strategy we are conducting is implicitly answering a question different to the one we want to ask.
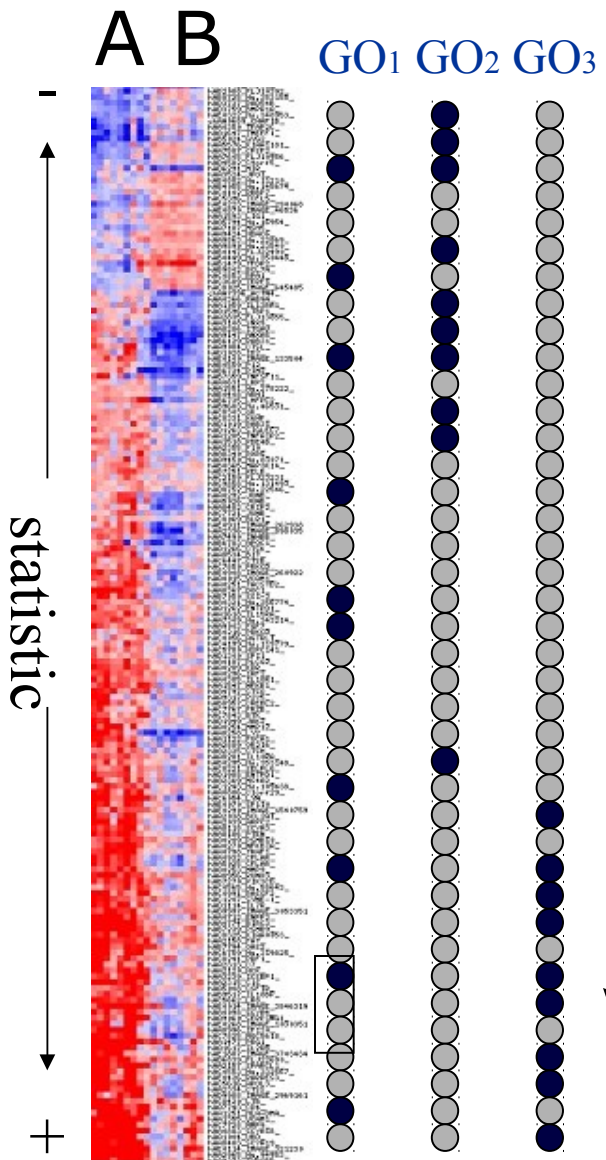
# Functional genomics.
## Historic perspective and future

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
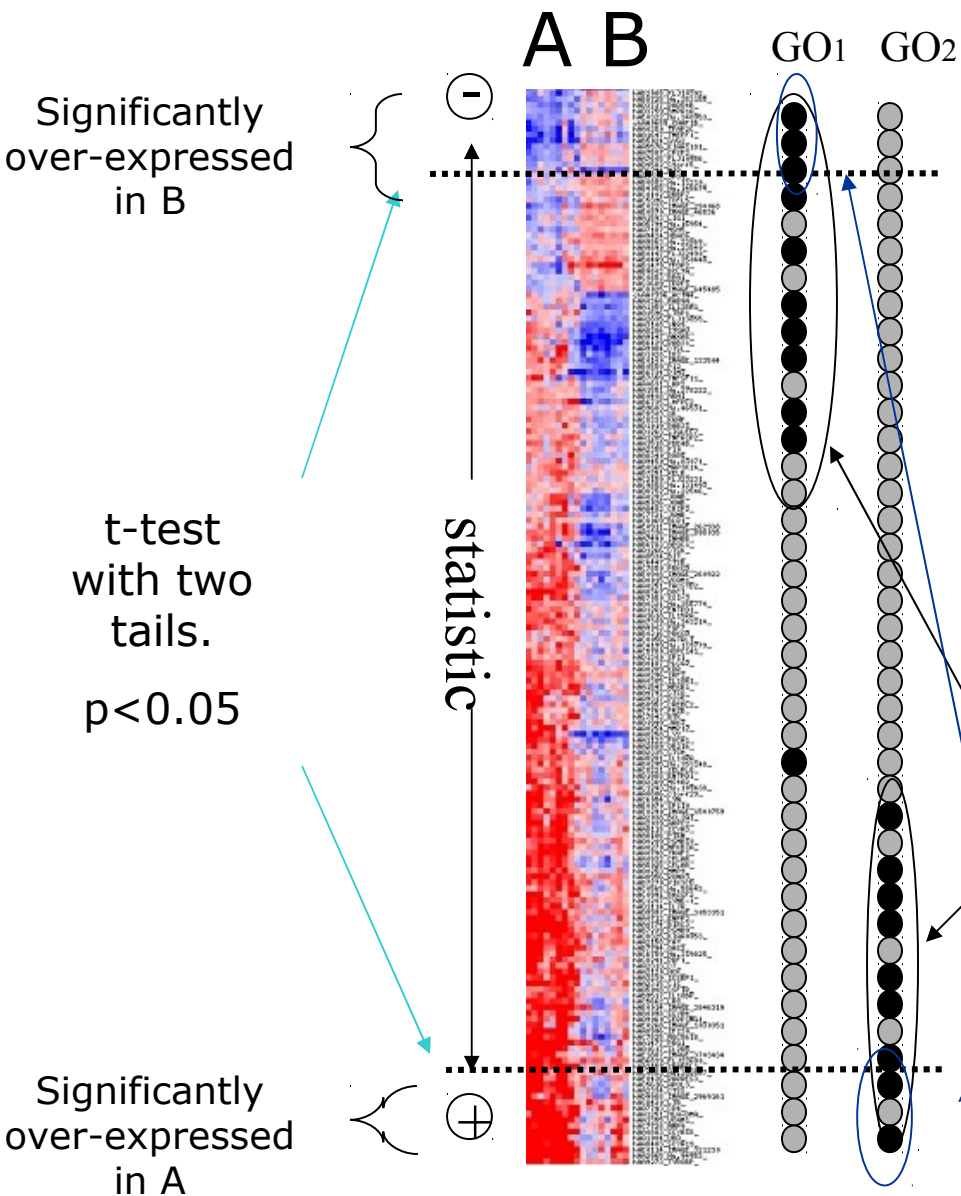


**Ranking**: A list of genes is ranked by their differential expression between two experimental conditions **A** and **B** (using fold change, a t-test, etc.)

**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

statistic

Significantly over-expressed in A

A B    GO₁   GO₂

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.
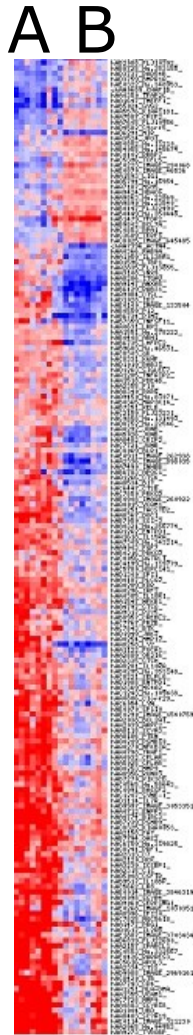
Classes expressed as blocks in A and B

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

# GSA 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)*

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

| Healthy vs diabetic | Functional class | GO | KEGG |
|---|---|---|---|
| Up-regulated | Oxidative phosphorylation | X | X |
| | ATP synthesis | | X |
| | Ribosome | | X |
| | Mitochondrion | X | |
| | Nucleotide biosynthesis | X | |
| | NADH dehidrogenase (ubiquinone activity) | X | |
| | Nuclease activity | X | |
| Down-regulated | Insulin signalling pathway | | X |

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

# Protein-protein interaction networks

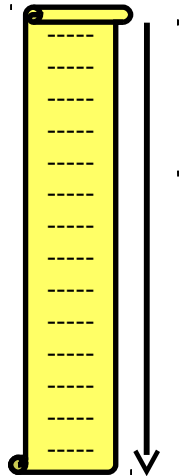## Evaluation of the cooperative behaviour of a list of genes

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

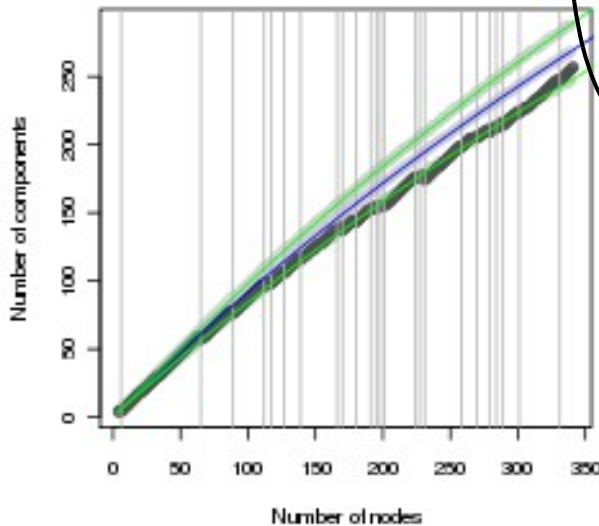Nodes included in the list

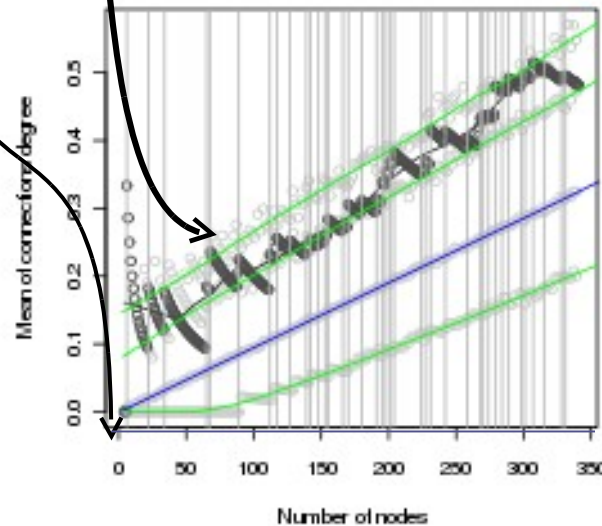Nodes not included in the list

# Gene-set-like network analysis

The list is traversed from higher to lower parameter values and the network properties are compared to their random expectations
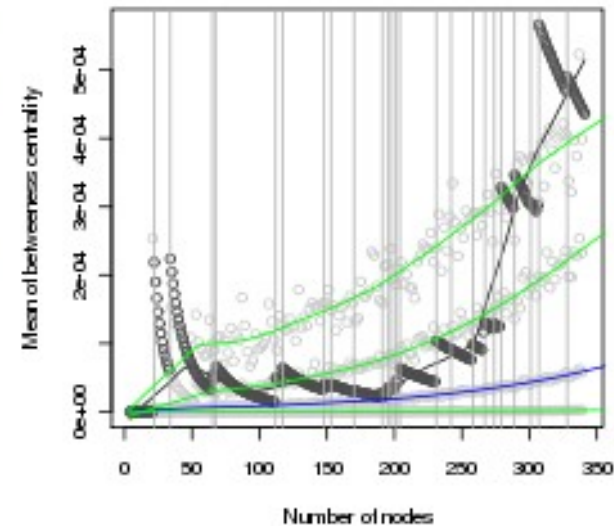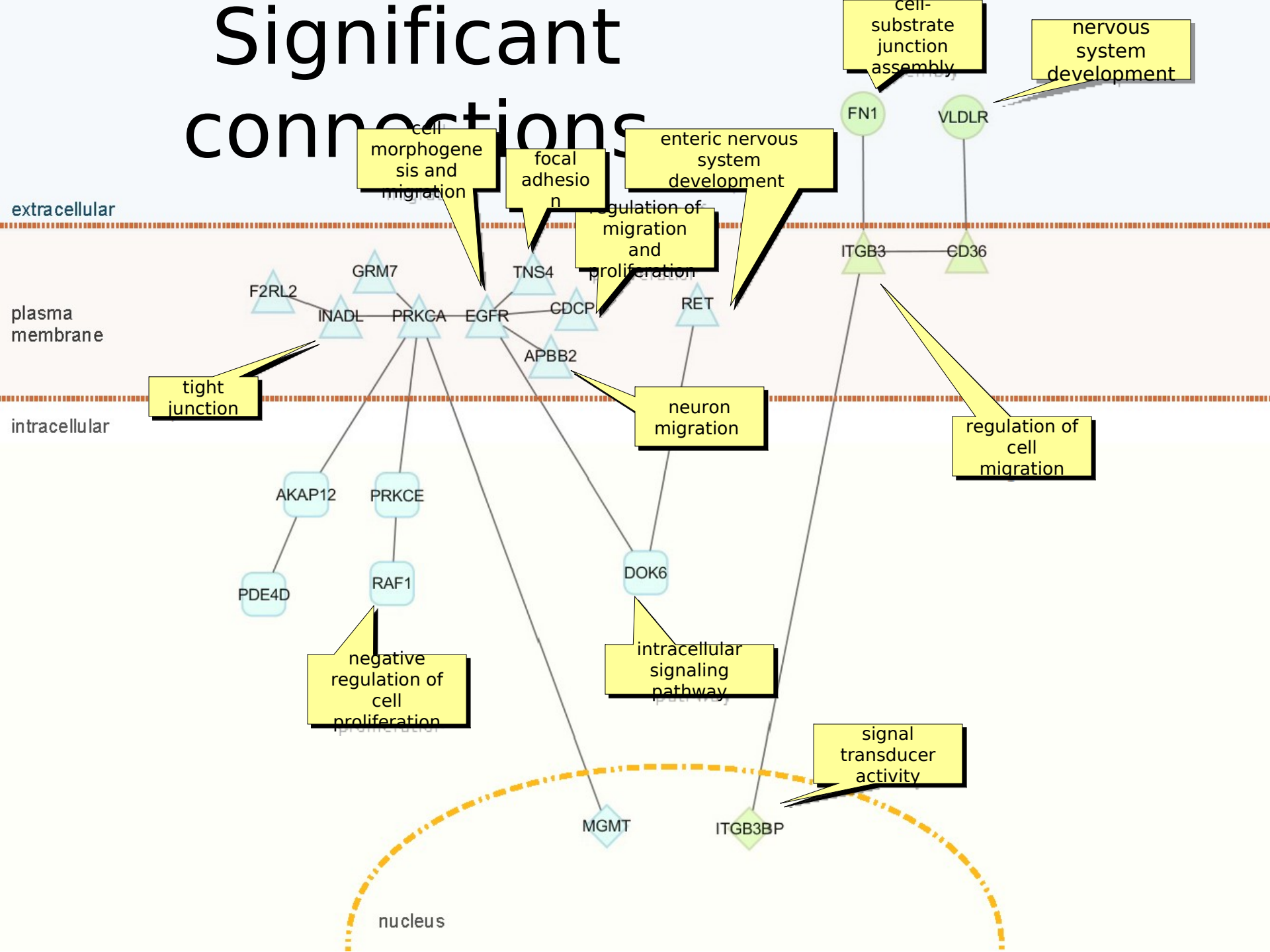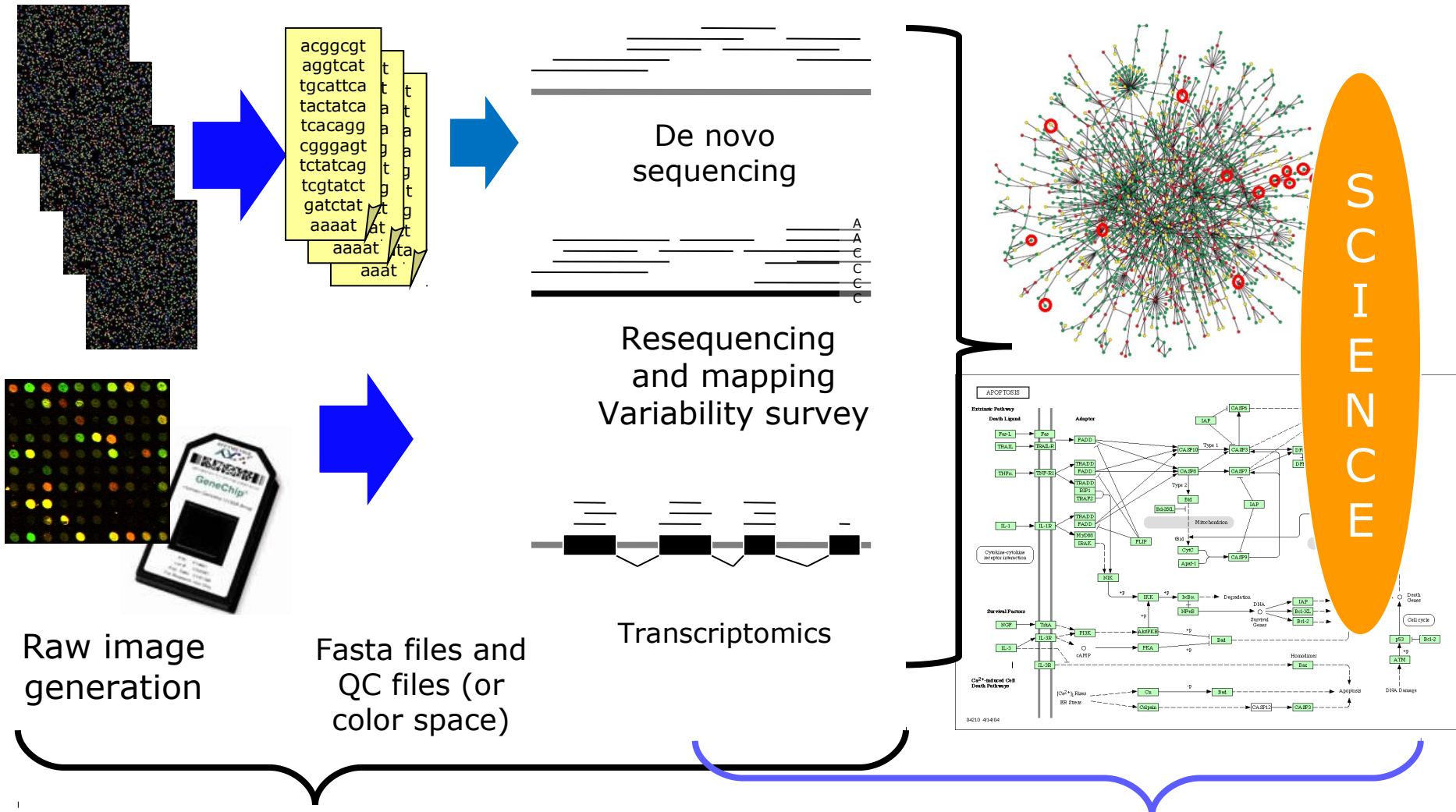
# Significant connections

# What is next?

## Functional classes have internal structure. Exploiting function and internal structure by modeling pathways

| Method | Gene-based selection | Function-based selection | Function | Relationships among components |
|---|---|---|---|---|
| Functional enrichment | X | | X | |
| Gene-set analysis | | X | X | |
| Network enrichment | X | | | X |
| Network enrichment analysis | | X | | X |
| Pathway modeling | | X | X | X |

# Example:

Dysregulated gene expression networks in human acute myelogenous leukemia stem cells

# Pipeline general of analysis



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

De novo
sequencing

Resequencing
and mapping
Variability survey

Transcriptomics

SCIENCE

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