

# CIBERER NGS course: from reads to candidate genes

## Introduction

Joaquín Dopazo

Computational Genomics Department,  
Centro de Investigación Príncipe Felipe (CIPF),  
Functional Genomics Node, (INB),  
Bioinformatics in Rare Diseases (BiER-CIBERER),  
Valencia, Spain.

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

 @xdopazo

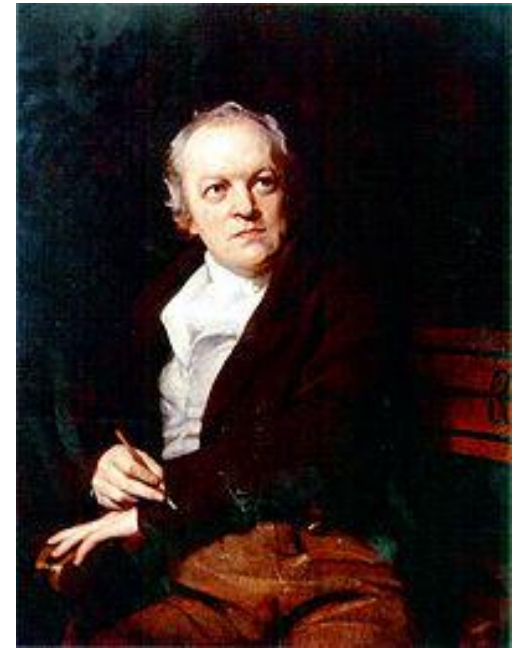
*CIPF, 28-30 September 2015*



# Background

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

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

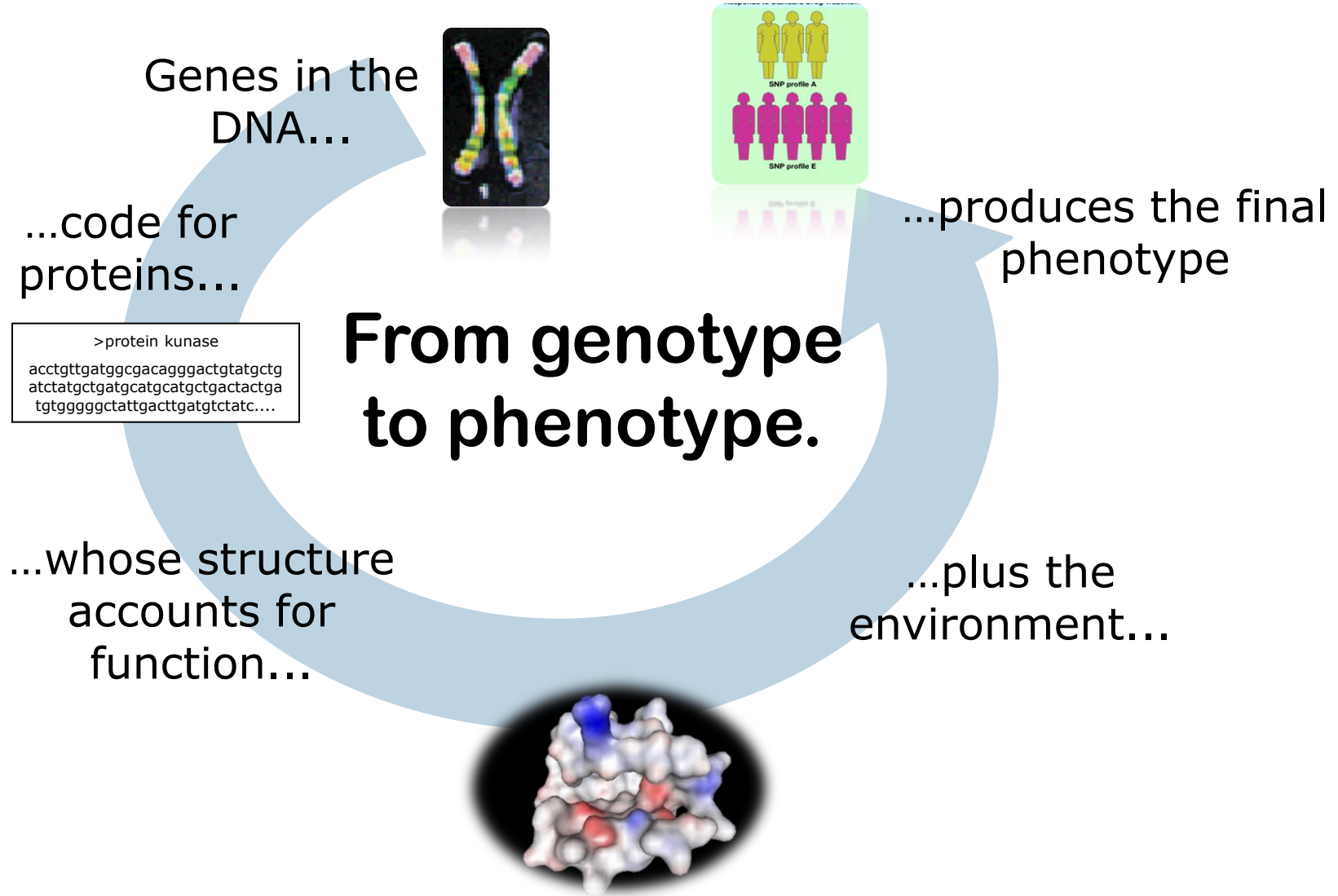


The introduction and popularisation of high-throughput techniques has drastically changed the way in which biological problems **can** be addressed and hypotheses **can** be tested.

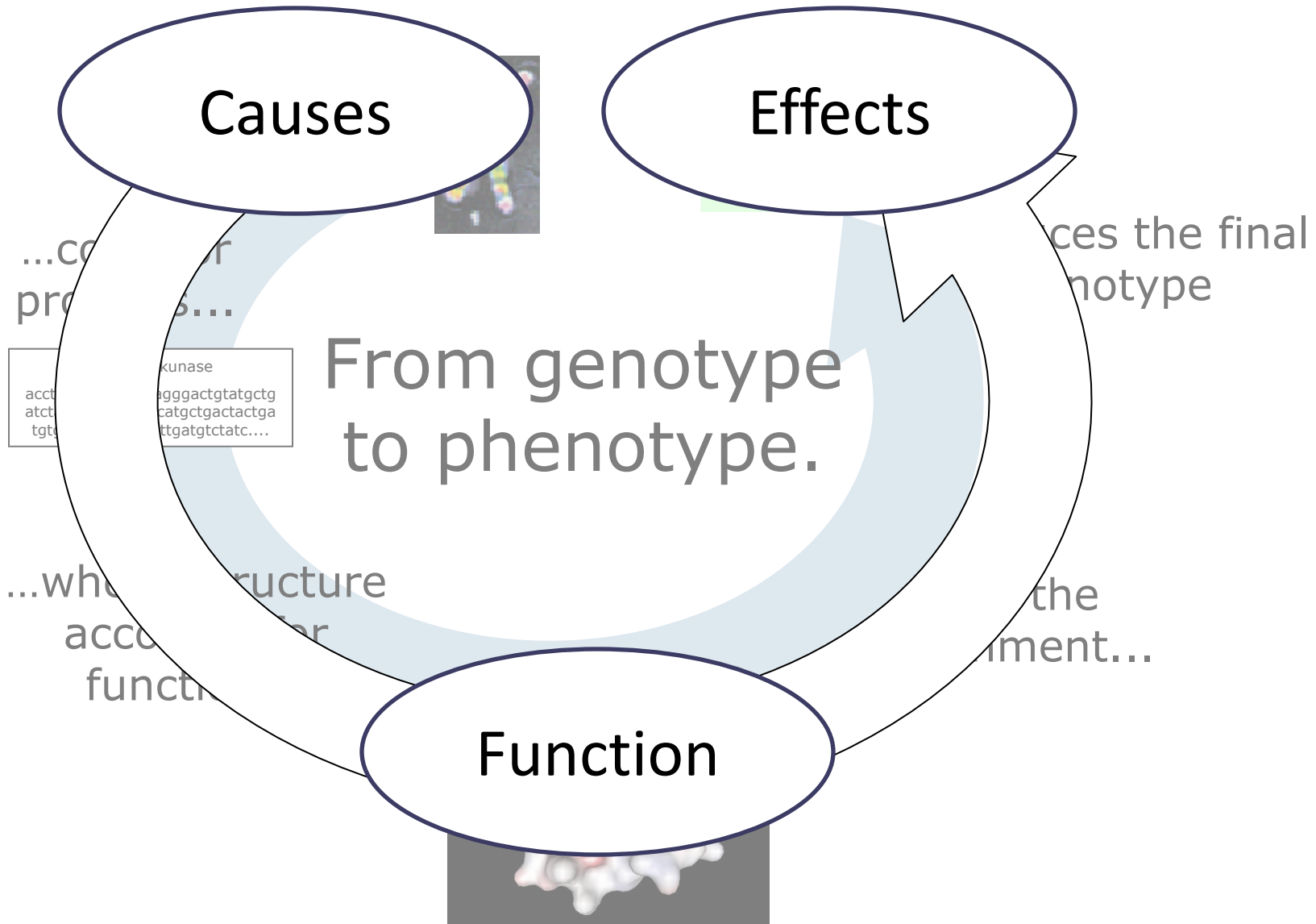
But not necessarily the way in which we really address or test them...

# Where do we come from?

## The pre-genomics paradigm

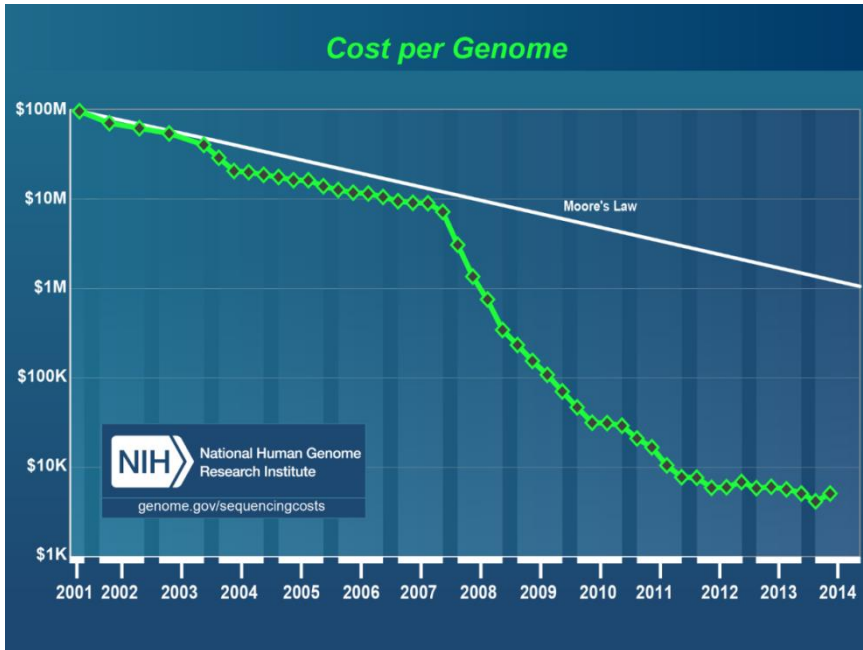


# Reduccionistic approach to link causes (genome) to effects (phenotype) through actions (function)

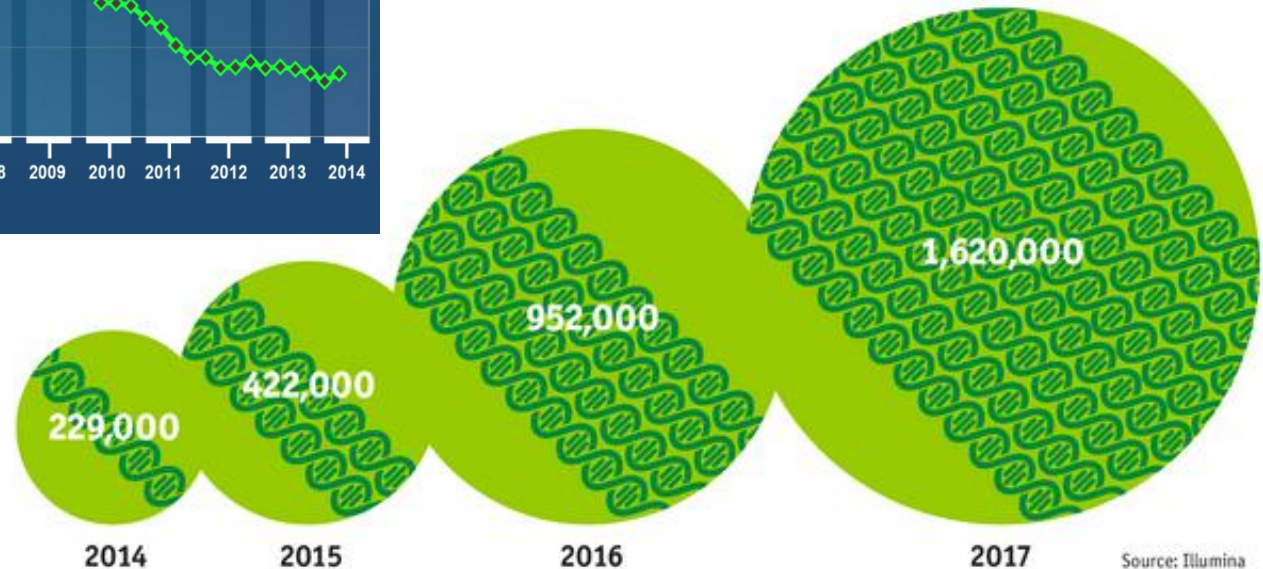


# The genome sequencing pace

<http://www.genome.gov/sequencingcosts/>



NGS is matching the cost of many conventional clinical tests



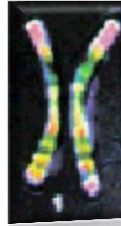
<http://www.economist.com/news/21631808-so-much-genetic-data-so-many-uses-genes-unzipped>

Next Generation Sequencing  
10<sup>9</sup>bp per round

Genes in  
the DNA...



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

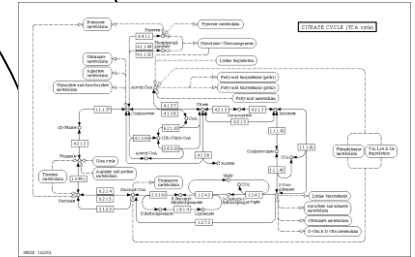


12 million SNPs in  
exonic regions

...with its  
complex  
variability...

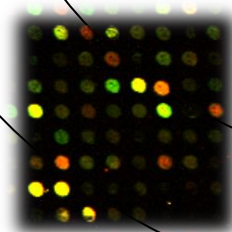
# From genotype to phenotype.

(in the post-genomics scenario)

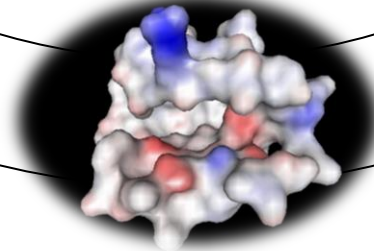


...conforming complex  
interaction networks...

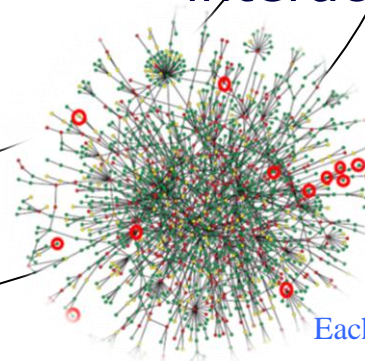
...when they  
are expressed  
in the proper  
moment and  
place...



...code for  
proteins...



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



...in cooperation  
with other  
proteins...

Each protein has an average  
of 8 interactions

...that account for  
function if...

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

Causes

Effects

From genotype to phenotype  
(in the functional post scenario)

Function  
(modules of proteins)

...with its complex variability

Half a million between p

ants individuals

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

...when are expressed in the pro moment at place...

...code for proteins...

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

...that account for function if...

plex KS...  
tion er  
ns...  
as an average of 8 interactions

# High-throughput data for functional genomics

Genotyping

Genome wide

Metabolomics

Transcriptomics

Proteomics

Almost-omics

From genotype to phenotype.

(in the functional post-genomics scenario)

Next Generation Sequencing  
10<sup>9</sup>bp per round

...with its complex variability...

Half a million of variants between pairs of individuals

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



...code for proteins...

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

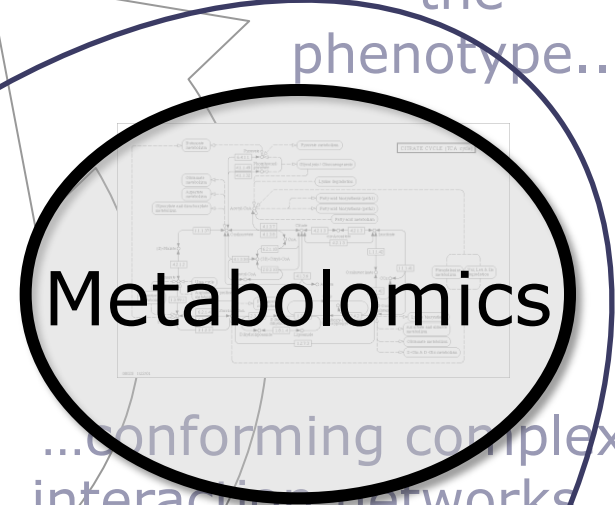
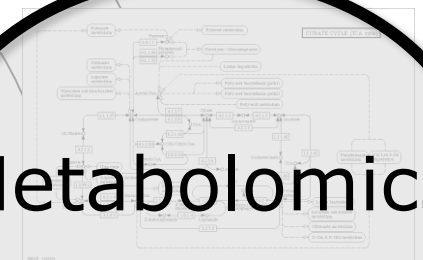
...that account function if...

...conforming complex interaction networks...

...in cooperation with other proteins...

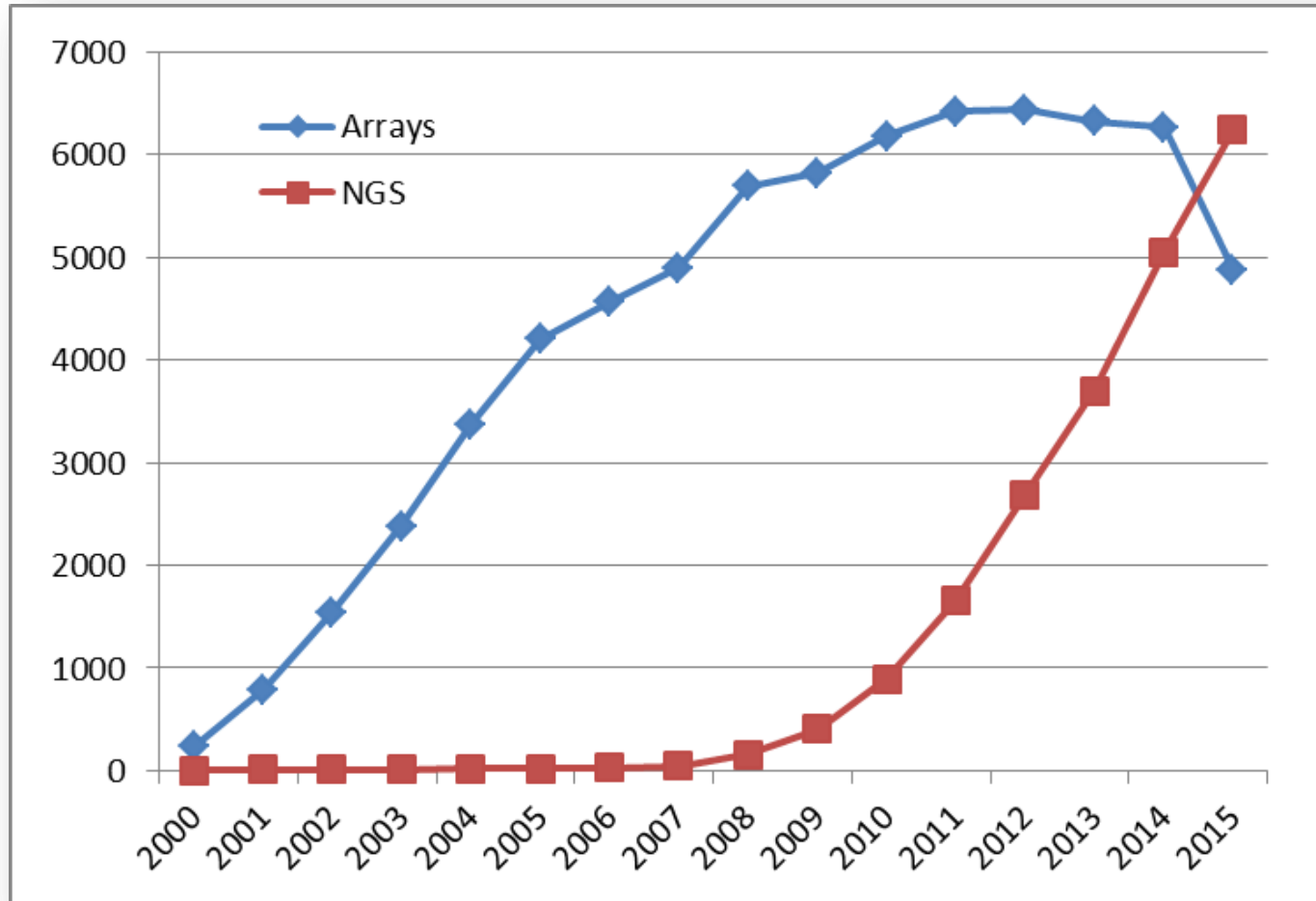
Each protein has an average of 8 interactions

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



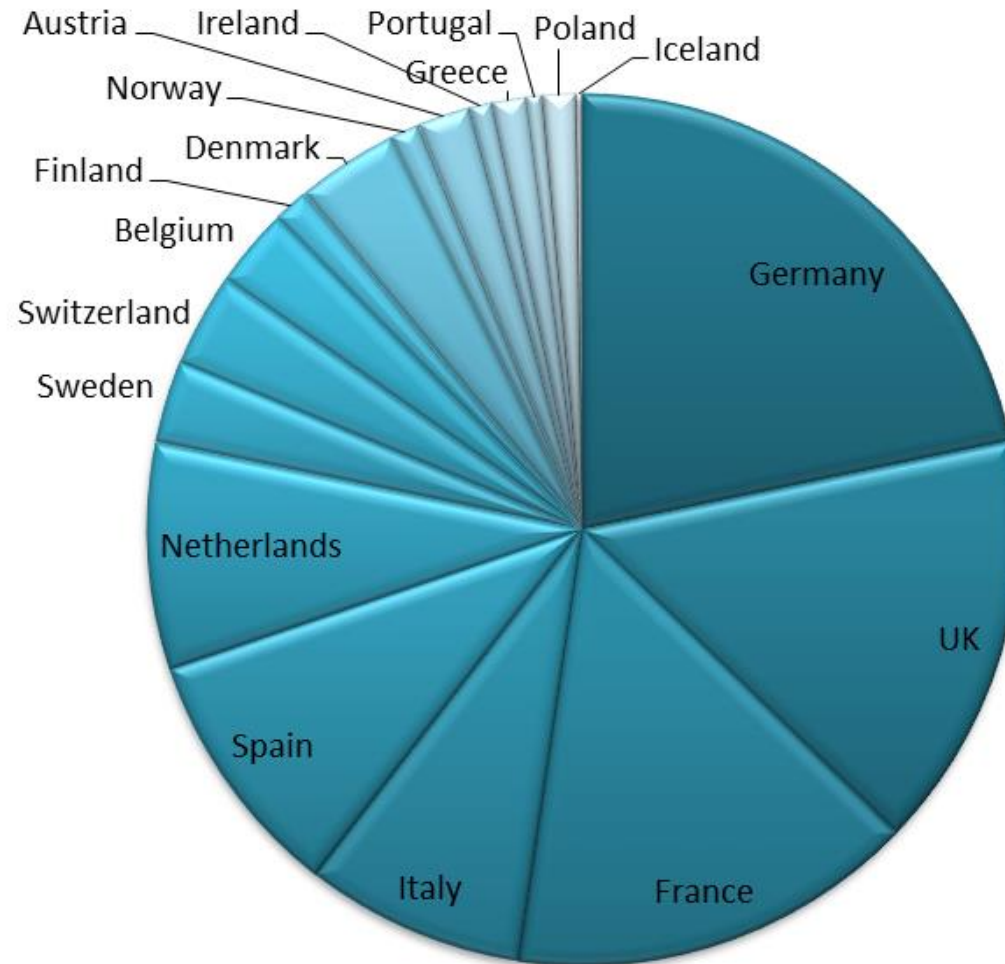


# Evolution of the papers published in microarray and NGS technologies



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

# Some bibliographic data: NGS publications

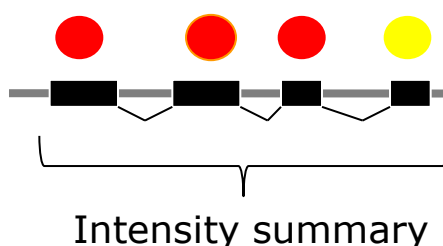
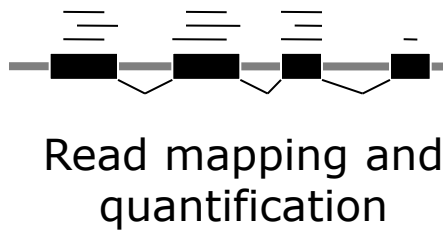
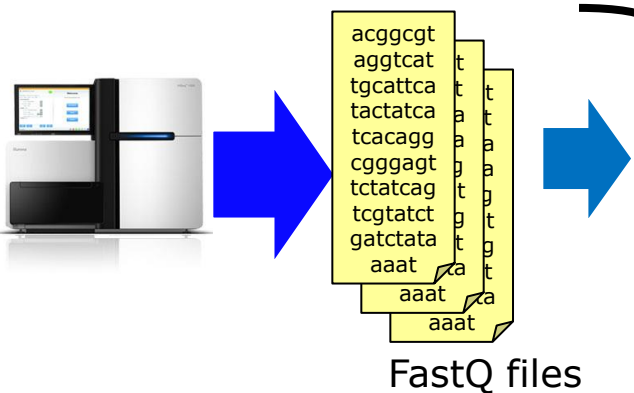


2015 Europe

**Source Pubmed. Query:**  
("high-throughput sequencing"[Title/Abstract] OR "next generation sequencing"[Title/Abstract] OR "rna seq"[Title/Abstract]) AND "2015"[Publication Date] AND country[Affiliation]

# Transcriptomics

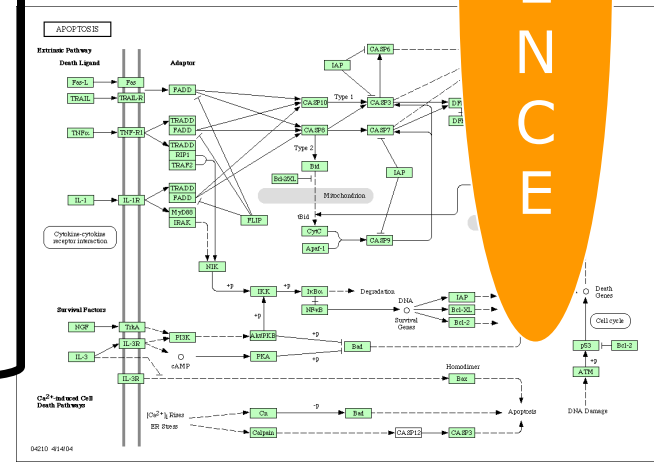
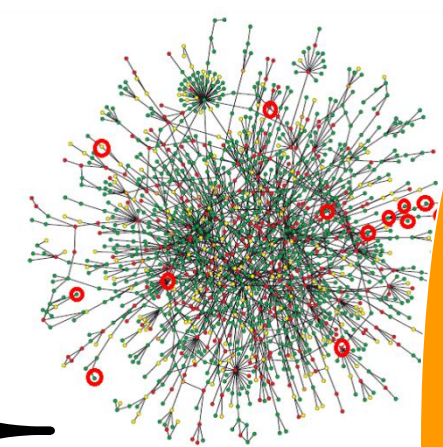
The double challenge:  
Data processing and interpretation



Raw image generation

Intensity summary

Transcriptomics



Technology driven

Hypothesis driven

# Before analysing your data you must know what is your question.

What is the aim? Class discovery? sample classification? gene selection? ...

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

Molecular classification of samples

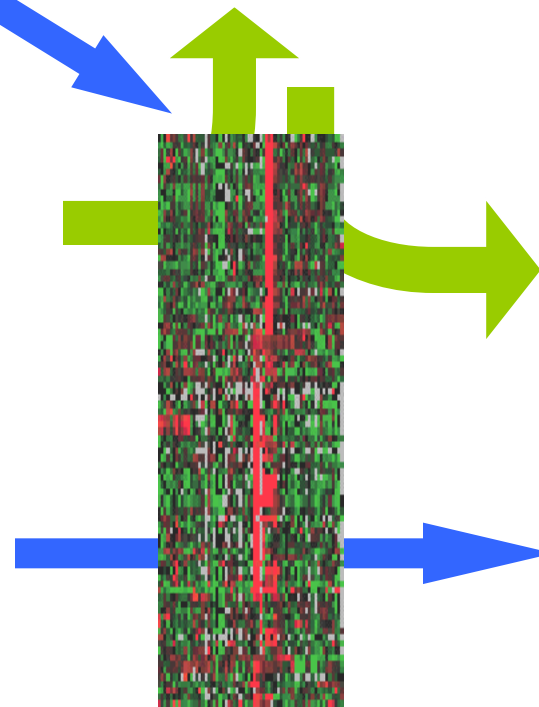
Co-expressing genes...

Different classes...



What genes are responsible for?

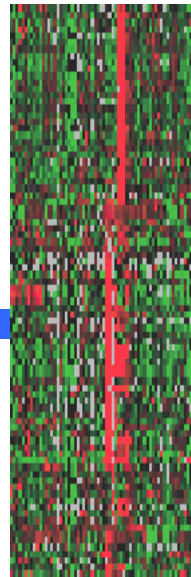
What do they have in common?



# 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...



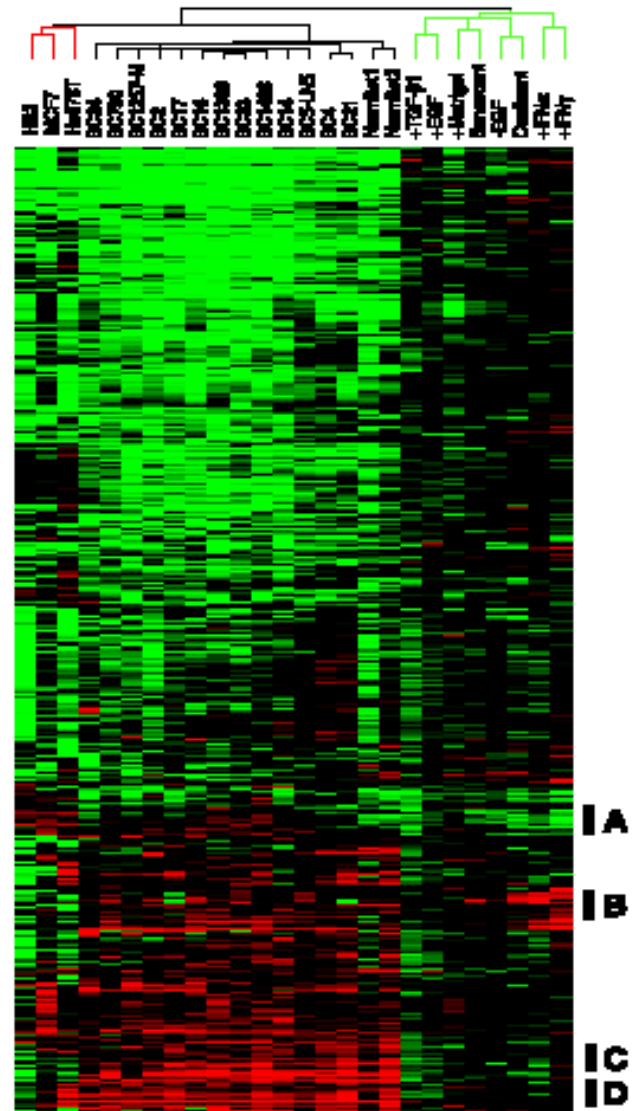
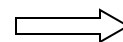
- What genes co-express?
- How many different expression patterns do we have?
- What do they have in common?
- Etc.

# 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 (1999)

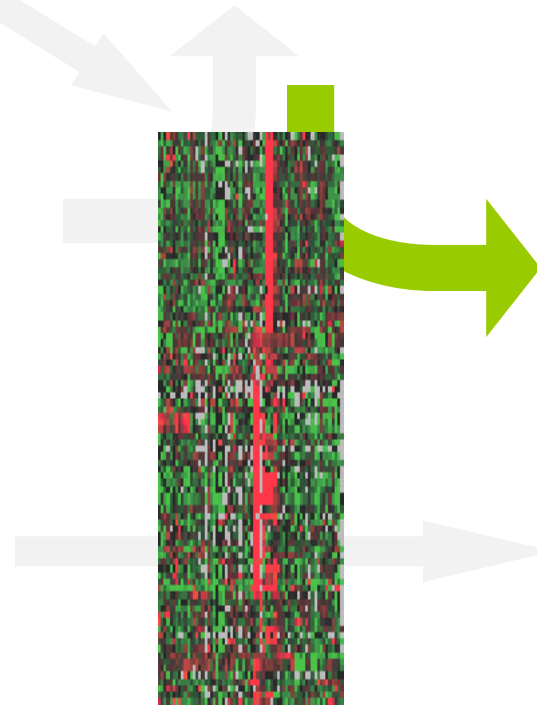
# Supervised problems.

## Differential gene expression

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?

# Differential gene expression

The simplest way: univariate gene-by-gene.  
Other multivariate approaches can be used

- **Two classes**

- T-test

- Limma

- Fold-change

- **Multiclass**

- Anova

- Limma

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

- Pearson

- Spearman

- Regression

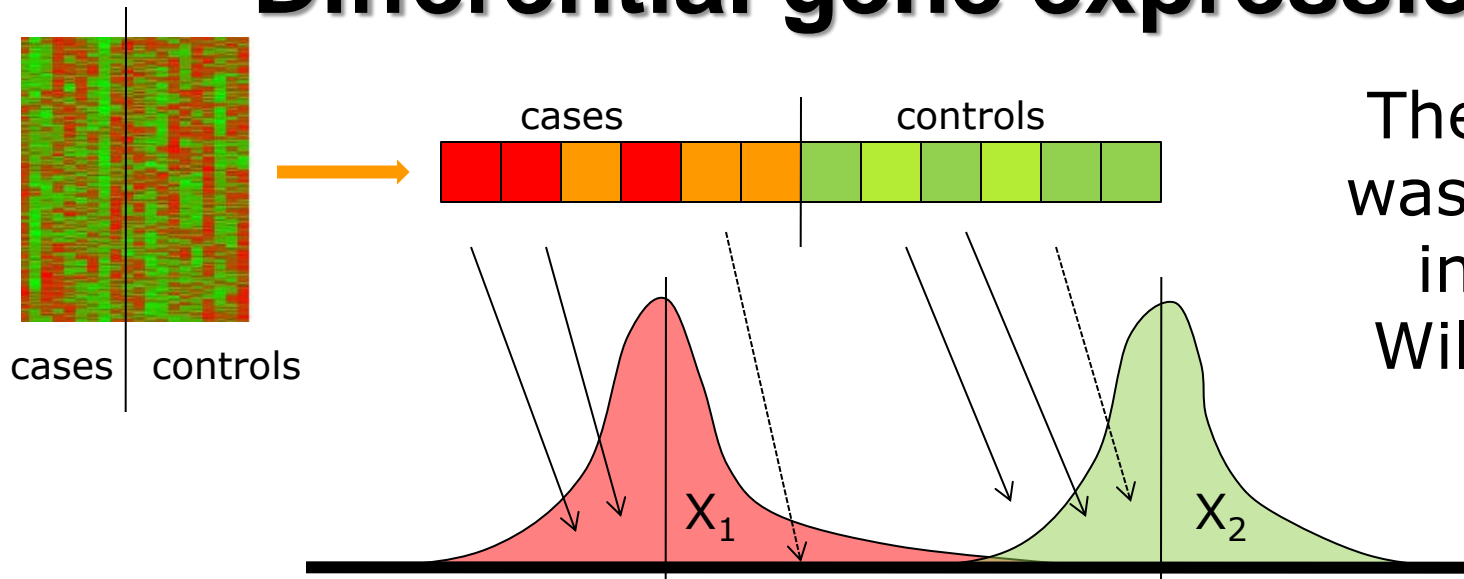
- **Survival**

- Cox model

- **Time Course**



# Differential gene expression



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

Significantly different



Non significantly different

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{X_1 X_2} \cdot \sqrt{\frac{2}{n}}} \quad \text{being} \quad S_{X_1 X_2} = \sqrt{\frac{S_{X_1}^2 + S_{X_2}^2}{2}}$$

# Supervised problems.

## sample classification

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

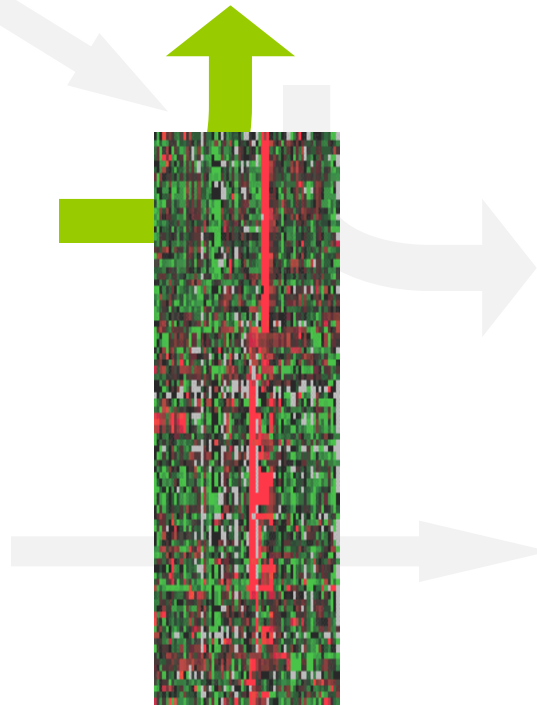
Different classes...

Molecular classification of samples

What genes are responsible for?

Co-expressing genes...

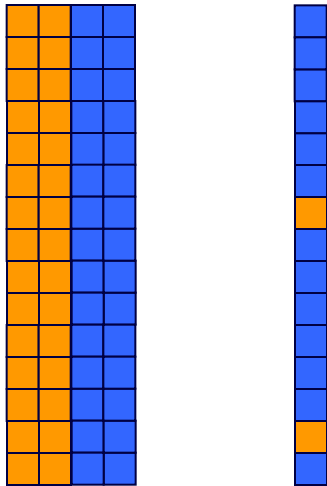
What do they have in common?



# Predictors and molecular signatures

What is a predictor?

A B X

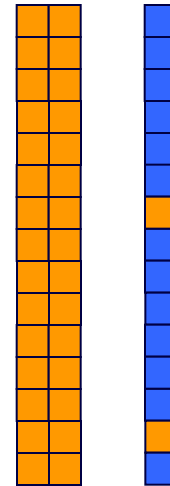


Is X, A or B?

$$\text{Diff (B, X)} = 2$$

$$\text{Diff (A, X)} = 13$$

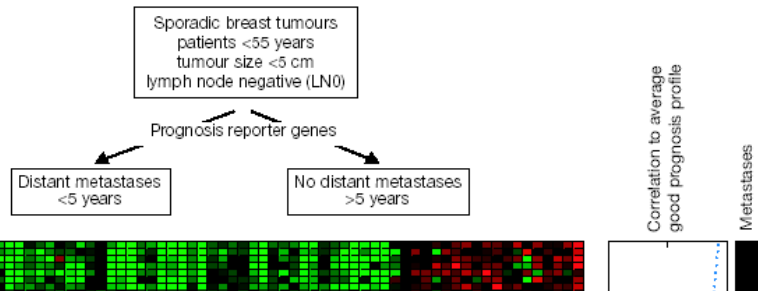
Intuitive notion:



Most probably X belongs to class B

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

# Predictor of clinical outcome in breast cancer



Genes are arranged to their correlation with the prognostic groups

agendia > Home - Microsoft Internet Explorer

Archivo Edición Ver Favoritos Herramientas Ayuda

Búsqueda Favoritos Multimedia

Dirección http://www.agendia.com/

Links Ensembl Genome Browser NCBI Home Page Google Scholar Bioinformatics - Manuscript Central [TM] MailSite Express Norton Internet Security

**Agendia.**

The future is now IS diagnostics by genomics

Now

Enter now . Unveiling Tumor behavior . Descifrar

**Tumor profiling to improve cancer treatment**

Agendia is a world leader in gene expression analysis-based diagnostics. By focusing on the genetic properties of a tumor, Agendia is able to unlock essential information on the risk of cancerous spread, cancer recurrence, the response to certain drugs or the primary site of a tumor. This information assists oncologists and physician to design tailor-made treatment plans that enhance the chances of success for cancer patients.

In addition, Agendia's expertise in gene expression profiling offers opportunities for pharmaceutical companies to improve their drug development tracks.

Internet

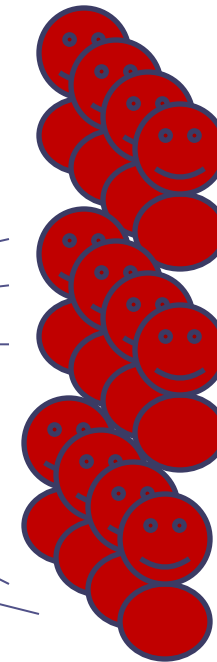
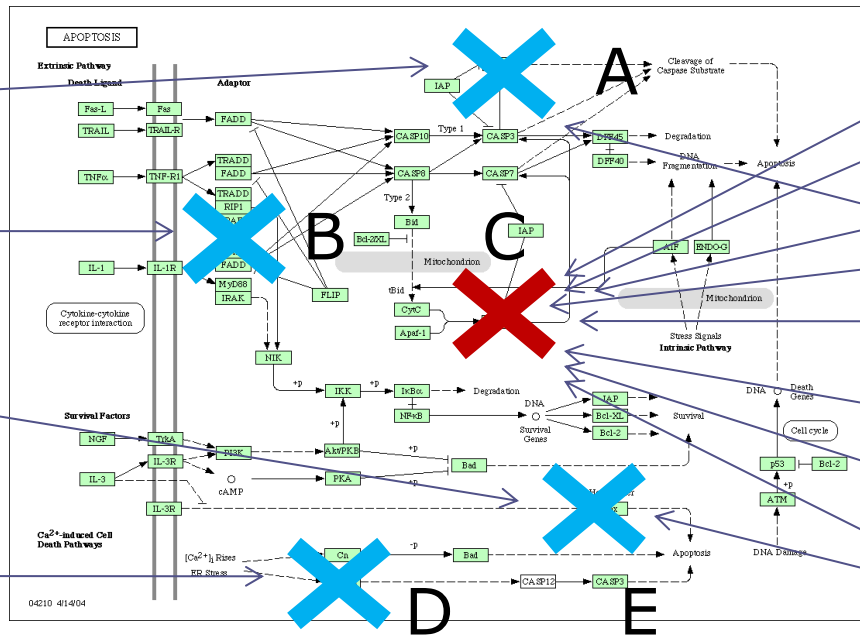
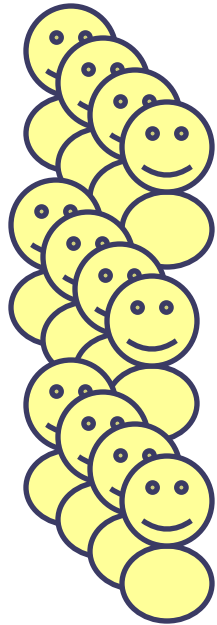
The screenshot shows a Microsoft Internet Explorer browser window displaying the Agendia website. The browser's address bar shows 'http://www.agendia.com/'. The website features the Agendia logo and a large graphic with the text 'The future is now IS diagnostics by genomics' and 'Now'. Below this, there is a section titled 'Tumor profiling to improve cancer treatment' with a brief description of Agendia's services. A blue arrow points from the text on the right towards the website content.

← Pronostic classifier with optimal accuracy

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

# Genotyping/Resequencing: Finding 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	0	0	0	1	0	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	0	0	0	0
Gene C	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1
Gene D	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	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	0	0	0	0	0	1	0	0	0	0



# Primary data analysis tools



**Fastq file, up to hundreds of GB per run**

**QC and preprocessing**

QC stats, filtering and preprocessing options

**HPG Aligner, short read aligner**

Double mapping strategy:  
Burrows-Wheeler Transform (*GPU Nvidia CUDA*) + Smith-Waterman (*CPU OpenMP+SSE/AVX*)

**SAM/BAM file**

**QC and preprocessing**

QC stats, filtering and preprocessing options

**Variant calling analysis**

GATK and SAM mPileup HPC Implementation.  
Statistics genomic tests

**VCF file**

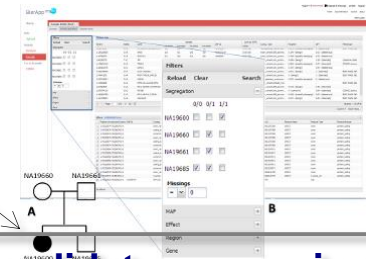
**QC and preprocessing**

QC stats, filtering and preprocessing options



**Genome Maps. Variant VCF viewer**

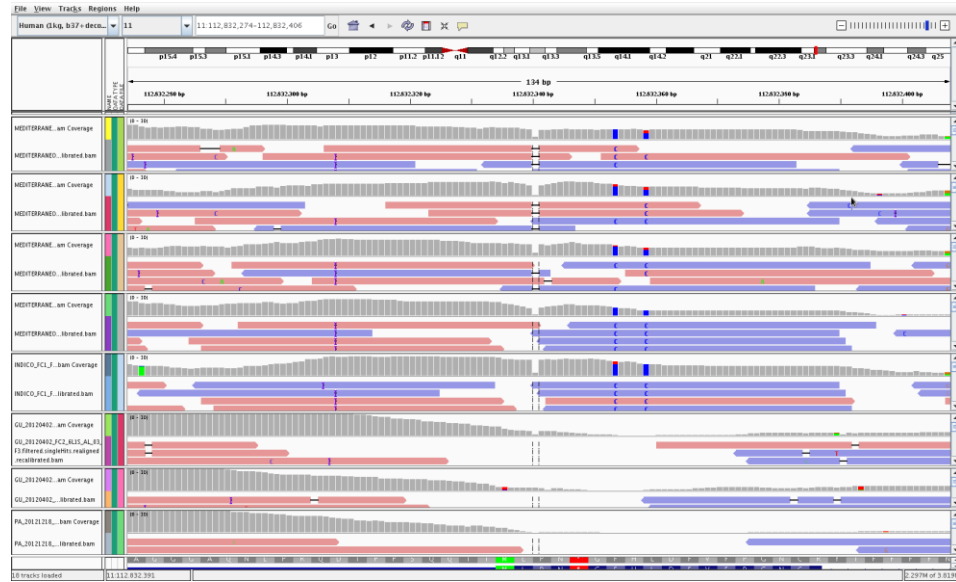
HTML5+SVG Web based viewer



**BiERApp candidate gene prioritization**

Consequence type, pathogenicity, Population frequencies, etc

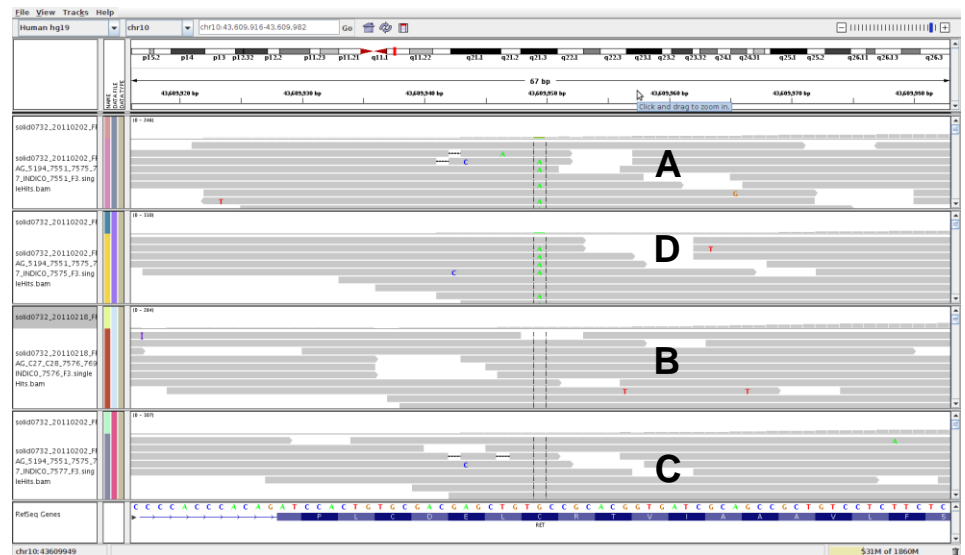
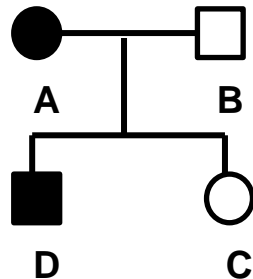
# The principle: comparison of patients to reference controls or segregation within families



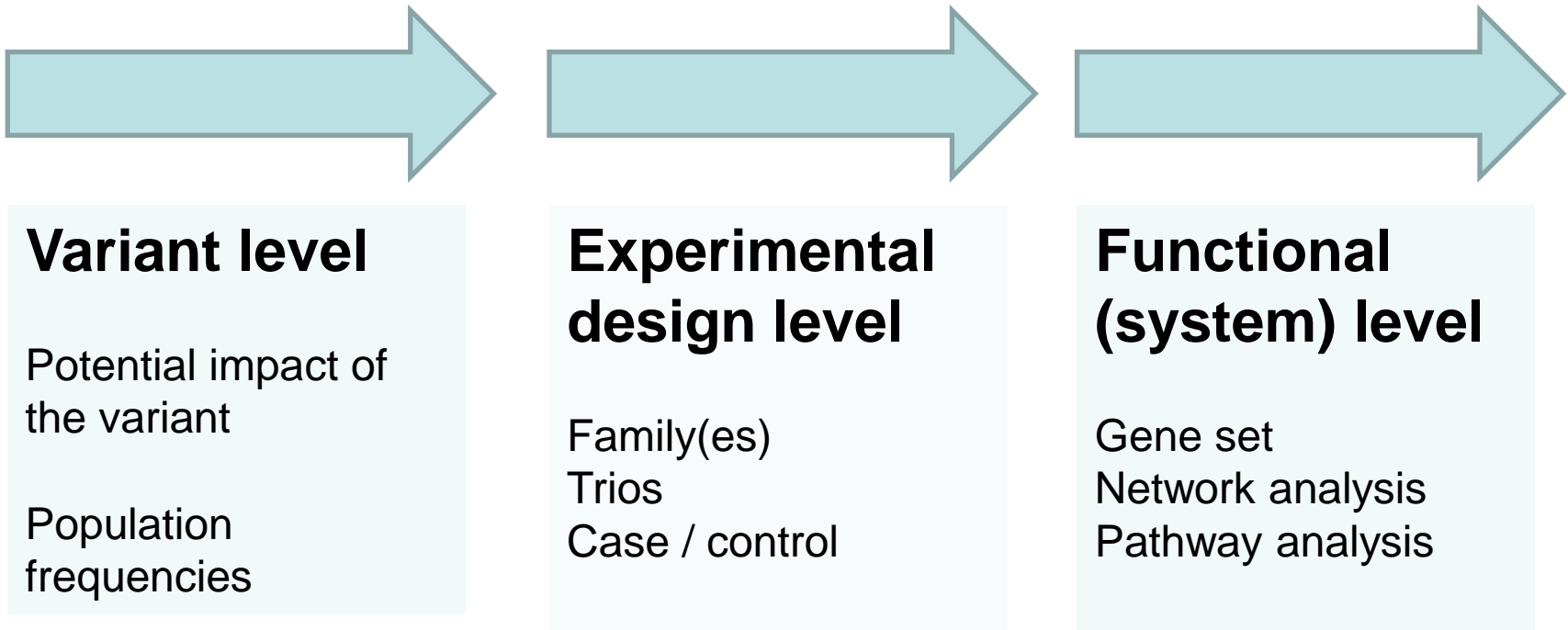
Cases

Controls

Segregation within a pedigree



# Variant/gene prioritization by successive filtering



Control of sequencing errors (missing values)

Testing strategies



# Pipeline of data analysis



Initial QC	Mapping + QC	Variant calling + QC	Variant and gene prioritization + QC
Sequence cleansing Base quality Remove adapters Remove duplicates  FASTQ file	Mapping (HPG) Remove multiple mapping reads Remove low quality mapping reads Realigning Base quality recalibrating  BAM file	Calling and labeling of missing values Calling SNVs and indels (GATK) using 6 statistics based on QC, strand bias, consistence (poor QC callings are converted to missing values as well)  Create multiple VCF with missing, SNVs and indels  VCF file	Counts of sites with variants Variant annotation (function, putative effect, conservation, etc.) Inheritance analysis (including compound heterozygotes in recessive inheritance) Filtering by frequency with external controls ( <b>Spanish controls</b> , dbSNP, 1000g, 5500g) and annotation Multi-family intersection of genes and variants Network-based prioritization Report

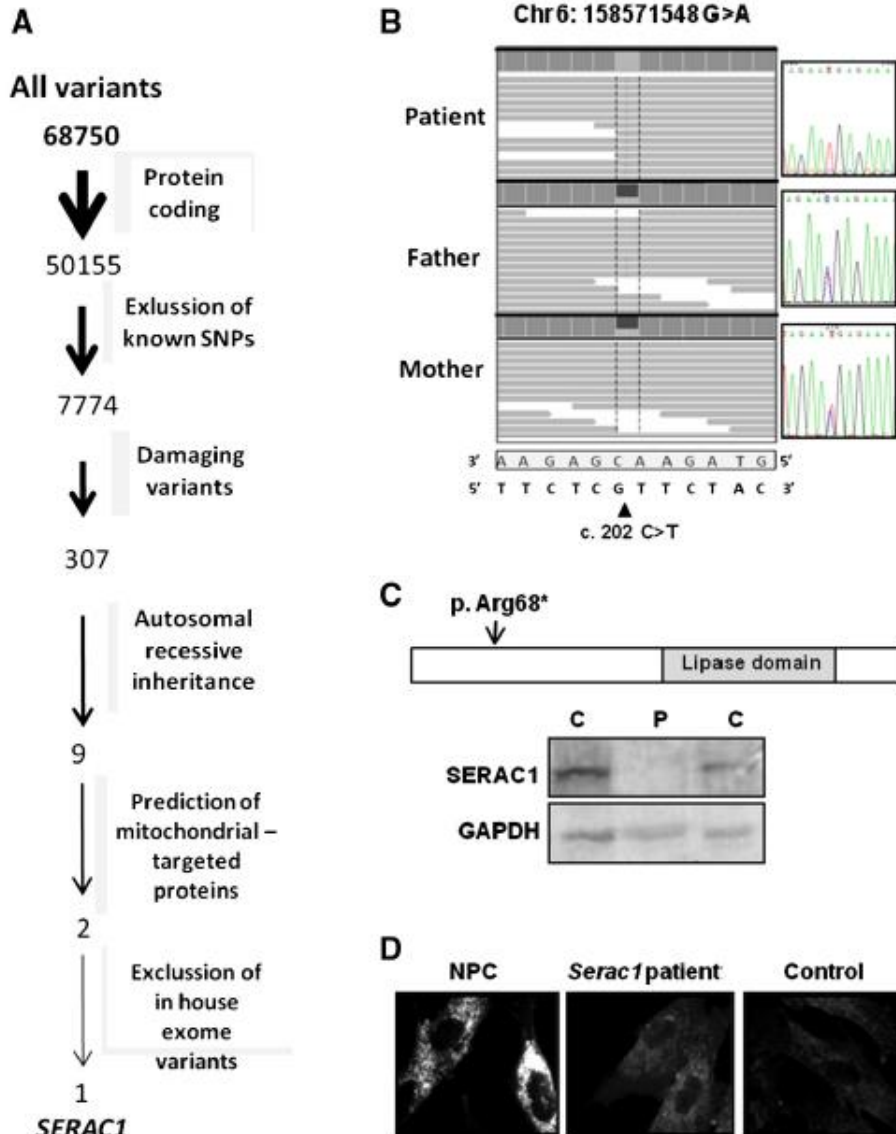
Primary analysis

Gene prioritization

# Successive Filtering approach

## An example with 3-Methylglutaconic aciduria syndrome

F. Tort et al. / Molecular Genetics and Metabolism xxx (2013) xxx–xxx



3-Methylglutaconic aciduria (3-MGA-uria) is a heterogeneous group of syndromes characterized by an increased excretion of 3-methylglutaconic and 3-methylglutaric acids. WES with a consecutive filter approach is enough to detect the new mutation in this case.

Contents lists available at SciVerse ScienceDirect

Molecular Genetics and Metabolism

journal homepage: [www.elsevier.com/locate/ymgme](http://www.elsevier.com/locate/ymgme)

Exome sequencing identifies a new mutation in *SERAC1* in a patient with 3-methylglutaconic aciduria

Frederic Tort<sup>a,b</sup>, María Teresa García-Silva<sup>c</sup>, Xènia Ferrer-Cortès<sup>a</sup>, Aleix Navarro-Sastre<sup>a,b</sup>, Judith García-Villoria<sup>a,b</sup>, Maria Josep Coll<sup>a,b</sup>, Enrique Vidal<sup>d</sup>, Jorge Jiménez-Almazán<sup>d</sup>, Joaquín Dopazo<sup>d,e,f</sup>, Paz Briones<sup>a,b,g</sup>, Orly Elpeleg<sup>h</sup>, Antonia Ribes<sup>a,b,\*</sup>

<sup>a</sup> Secció d'Errors Congènits del Metabolisme, Servei de Bioquímica i Genètica Molecular, Hospital Clínic, IDIBAPS, 08028, Barcelona, Spain

<sup>b</sup> CIBER de Enfermedades Raras (CIBERER), Barcelona, Spain

<sup>c</sup> Unidad de Enfermedades Mitocondriales- Enfermedades Metabólicas Hereditarias, Servicio de Pediatría, Hospital 12 de Octubre, Madrid, Spain

<sup>d</sup> IBER, CIBERER, Centro de Investigación Príncipe Felipe, Valencia, Spain

<sup>e</sup> Computational Medicine Institute, Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain

<sup>f</sup> Functional Genomics Node, (INB) at CIPF, Valencia, Spain

<sup>g</sup> Consejo Superior de Investigaciones Científicas (CSIC), Barcelona, Spain

<sup>h</sup> Maniague and Jacques Robson Department of Genetic Research, Hadassah, Hebrew University Medical Center, Jerusalem, Israel

\* Corresponding author. E-mail: [aribes@cipf.es](mailto:aribes@cipf.es)

# Exome sequencing has been systematically used to identify Mendelian disease genes

## ARTICLES

nature  
genetics

### Exome sequencing identifies the cause of a mendelian disorder

Sarah B Ng<sup>1,10</sup>, Kati J Buckingham<sup>2,10</sup>, Choli Lee<sup>1</sup>, Abigail W Bigham<sup>2</sup>, Holly K Tabor<sup>2,3</sup>, Karin M Dent<sup>4</sup>, Chad D Huff<sup>5</sup>, Paul T Shannon<sup>6</sup>, Ethilyn Wang Jabs<sup>7,8</sup>, Deborah A Nickerson<sup>1</sup>, Jay Shendure<sup>1</sup> & Michael J Bamshad<sup>1,2,9</sup>

We demonstrate the first successful application of exome sequencing to discover the gene for a rare mendelian disorder of unknown cause, Miller syndrome (MIM%263750). For four affected individuals in three independent kindreds, we captured and sequenced coding regions to a mean coverage of 40x and sufficient depth to call variants at ~97% of each targeted exome. Filtering against public SNP databases and eight HapMap exomes for genes with two previously unknown variants in each of the four individuals identified a single candidate gene, *DHODH*, which encodes a key enzyme in the pyrimidine *de novo* biosynthesis pathway. Sanger sequencing confirmed the presence of *DHODH* mutations in three additional families with Miller syndrome. Exome sequencing of a small number of unrelated affected individuals is a powerful, efficient strategy for identifying the genes un-

## REVIEWS

### TRANSLATIONAL GENETICS

### Exome sequencing as a tool for Mendelian disease gene discovery

Michael J. Bamshad<sup>\*\*†</sup>, Sarah B. Ng<sup>†</sup>, Abigail W. Bigham<sup>\*\*§</sup>, Holly K. Tabor<sup>\*\*||</sup>, Mary J. Emond<sup>†</sup>, Deborah A. Nickerson<sup>†</sup> and Jay Shendure<sup>†</sup>

Abstract | Exome sequencing — the targeted sequencing of the subset of the human genome that is protein coding — is a powerful and cost-effective new tool for dissecting the genetic basis of diseases and traits that have proved to be intractable to conventional gene-discovery strategies. Over the past 2 years, experimental and analytical approaches relating to exome sequencing have established a rich framework for discovering the genes underlying unsolved Mendelian disorders. Additionally, exome sequencing is being adapted to explore the extent to which rare alleles explain the heritability of complex diseases and health-related traits. These advances also set the stage for applying exome and whole-genome sequencing to facilitate clinical diagnosis and personalized disease-risk profiling.

### Whole-Exome Re-Sequencing in a Family Quartet Identifies *POP1* Mutations As the Cause of a Novel Skeletal Dysplasia

Evgeny A. Glazov<sup>1\*§</sup>, Andreas Zanki<sup>2,3</sup>, Marina Donskoi<sup>1</sup>, Tony J. Kenna<sup>1</sup>, Gethin P. Thomas<sup>1</sup>, Graeme R. Clark<sup>†</sup>, Emma L. Duncan<sup>1,2</sup>, Matthew A. Brown<sup>1\*</sup>

1 University of Queensland Diamantina Institute, Princess Alexandra Hospital, Woolloongabba, Australia, 2 Centre for Clinical Research, The University of Queensland, Herston, Australia, 3 School of Medicine, Faculty of Health Sciences, The University of Queensland, Herston, Australia

#### Abstract

Recent advances in DNA sequencing have enabled mapping of genes for monogenic traits in families with small pedigrees and even in unrelated cases. We report the identification of disease-causing mutations in a rare, severe, skeletal dysplasia.

European Journal of Human Genetics (2011) 19, 115–117  
© 2011 Macmillan Publishers Limited All rights reserved 1018-4813/11  
www.nature.com/ejhg



The two forms of quencing, as a core MRP RNA activity of rich *POP1*

### SHORT REPORT

### Next generation sequencing in a family with autosomal recessive Kahrizi syndrome (OMIM 612713) reveals a homozygous frameshift mutation in *SRD5A3*

Kimia Kahrizi<sup>1</sup>, Cougar Hao Hu<sup>2</sup>, Masoud Garshashi<sup>2</sup>, Seydeh Sedigheh Abedini<sup>1</sup>, Shirin Ghadami<sup>1</sup>, Roxana Kariminejad<sup>1</sup>, Reinhard Ullmann<sup>2</sup>, Wei Chen<sup>2</sup>, H-Hilger Ropers<sup>2</sup>, Andreas W Kuss<sup>2</sup>, Hossein Najmabadi<sup>1</sup> and Andreas Tzschach<sup>1,2</sup>

As part of a large-scale, systematic effort to unravel the molecular causes of autosomal recessive mental retardation, we have previously described a novel syndrome consisting of mental retardation, coloboma, cataract and kerchia (Kahrizi syndrome).

OMIM 612713  
array-based  
(c.203del)  
interval  
essential  
families  
and eye  
potential  
European

Keyword:  
consanguinity



Molecular Vision 2013; 19:2187-2195 <<http://www.molvis.org/molvis/v19/2187>>  
Received 21 May 2013 | Accepted 5 November 2013 | Published 7 November 2013

© 2013 Molecular Vision

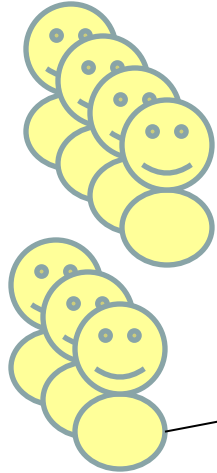
### Whole-exome sequencing identifies novel compound heterozygous mutations in *USH2A* in Spanish patients with autosomal recessive retinitis pigmentosa

Cristina Méndez-Vidal<sup>1,2</sup>, María González-del Pozo<sup>1,2</sup>, Alicia Vela-Boza<sup>3</sup>, Javier Santoyo-López<sup>2</sup>, Francisco J. López-Domingo<sup>3</sup>, Carmen Vázquez-Marouschek<sup>4</sup>, Joaquín Dopazo<sup>3,5,6</sup>, Salud Borrego<sup>1,2</sup>, Guillermo Antiñolo<sup>1,2,3</sup>

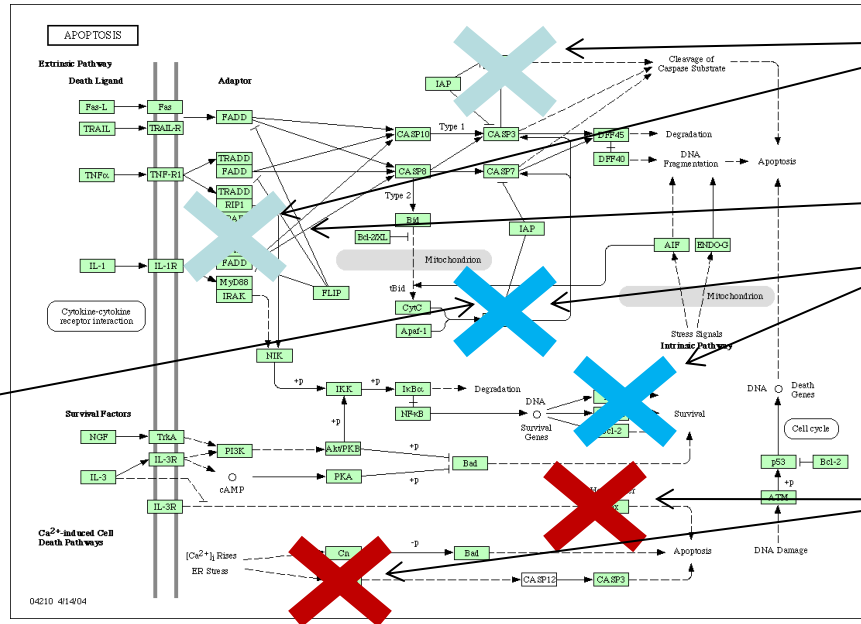
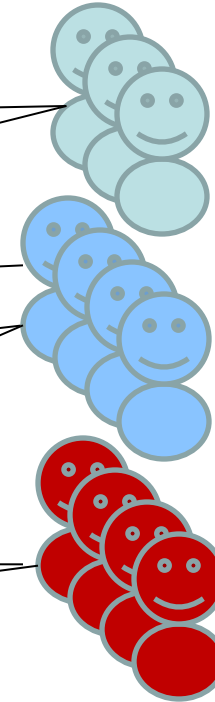
<sup>1</sup>Department of Genetics, Reproduction and Fetal Medicine, Institute of Biomedicine of Seville, University Hospital Virgen del Rocío/CSIC/University of Seville, Seville, Spain; <sup>2</sup>Centro de Investigación Biomédica en Red de Enfermedades Raras (CIBERER), Seville, Spain; <sup>3</sup>Medical Genome Project, Genomics and Bioinformatics Platform of Andalusia (GBPA), Seville, Spain; <sup>4</sup>Department of Ophthalmology, University Hospital Virgen del Rocío, Seville, Spain; <sup>5</sup>Department of Bioinformatics, Centro de Investigación Príncipe Felipe, Valencia, Spain; <sup>6</sup>Functional Genomics Node (INB), Centro de Investigación Príncipe Felipe, Valencia, Spain

# An approach inspired on systems biology can help in detecting causal genes

Controls



Cases



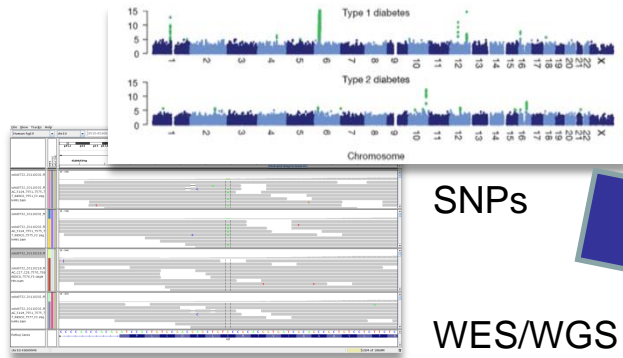
Affected **cases** in complex diseases will be a **heterogeneous** population with different mutations (or combinations).

Many cases and controls are needed to obtain significant associations.

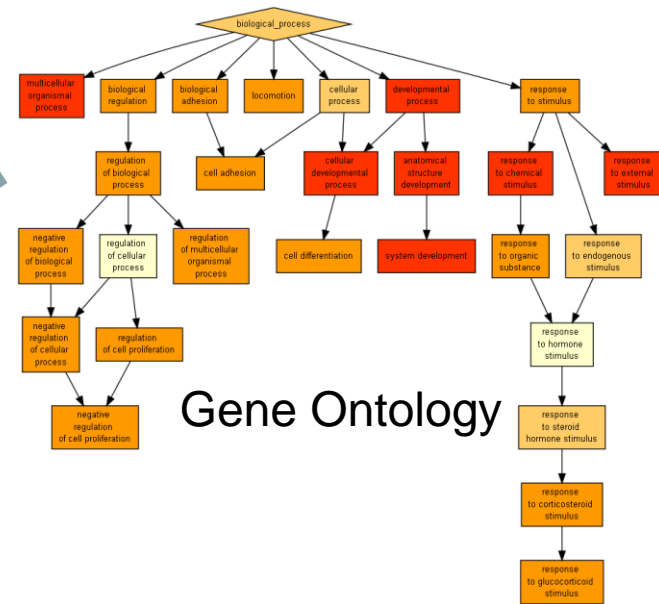
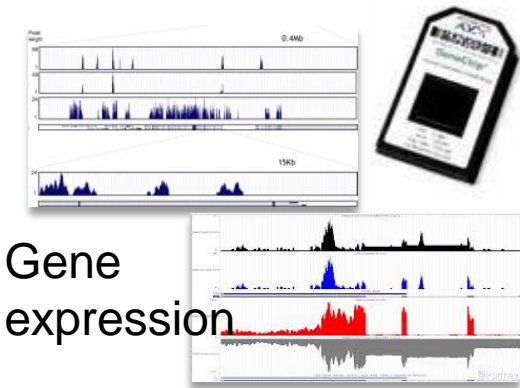
The only **common element** is the (know or unknown) **pathway affected**.

**Disease understood as the failure of a functional module**

# From gene-based to function-based perspective



AND/OR



**Gene Ontology** are **labels** to genes that describe, by means of a controlled vocabulary (ontology), the **functional role(s)** played by the genes in the cell. A set of genes **sharing** a **GO** annotation can be considered a **functional module**.

# An example of GWAS

GWAS in Breast Cancer.

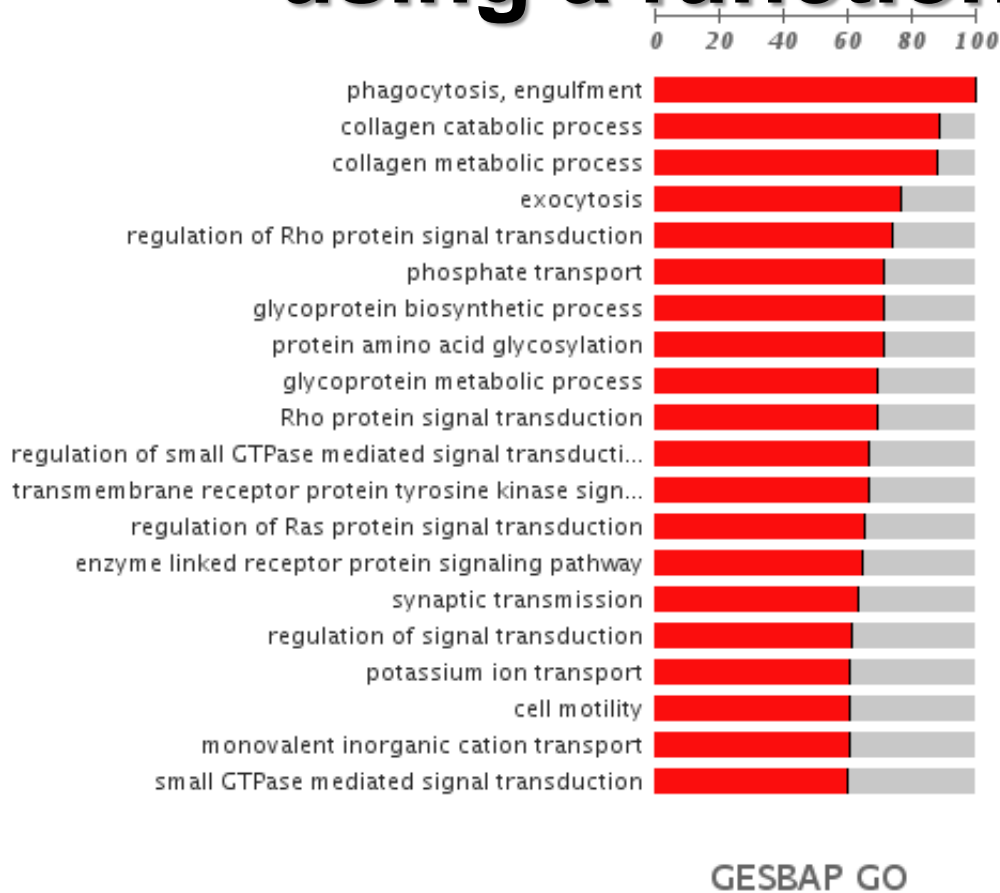
The CGEMS initiative. (Hunter et al. Nat Genet 2007)

1145 cases 1142 controls. Affy 500K

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

Conclusions: **conventional SNP-based** or **gene-based tests** are not providing much resolution.

# The same GWAS data re-analyzed using a function-based test



Breast Cancer

CGEMS initiative.

(Hunter et al. Nat

Genet 2007)

1145 cases 1142

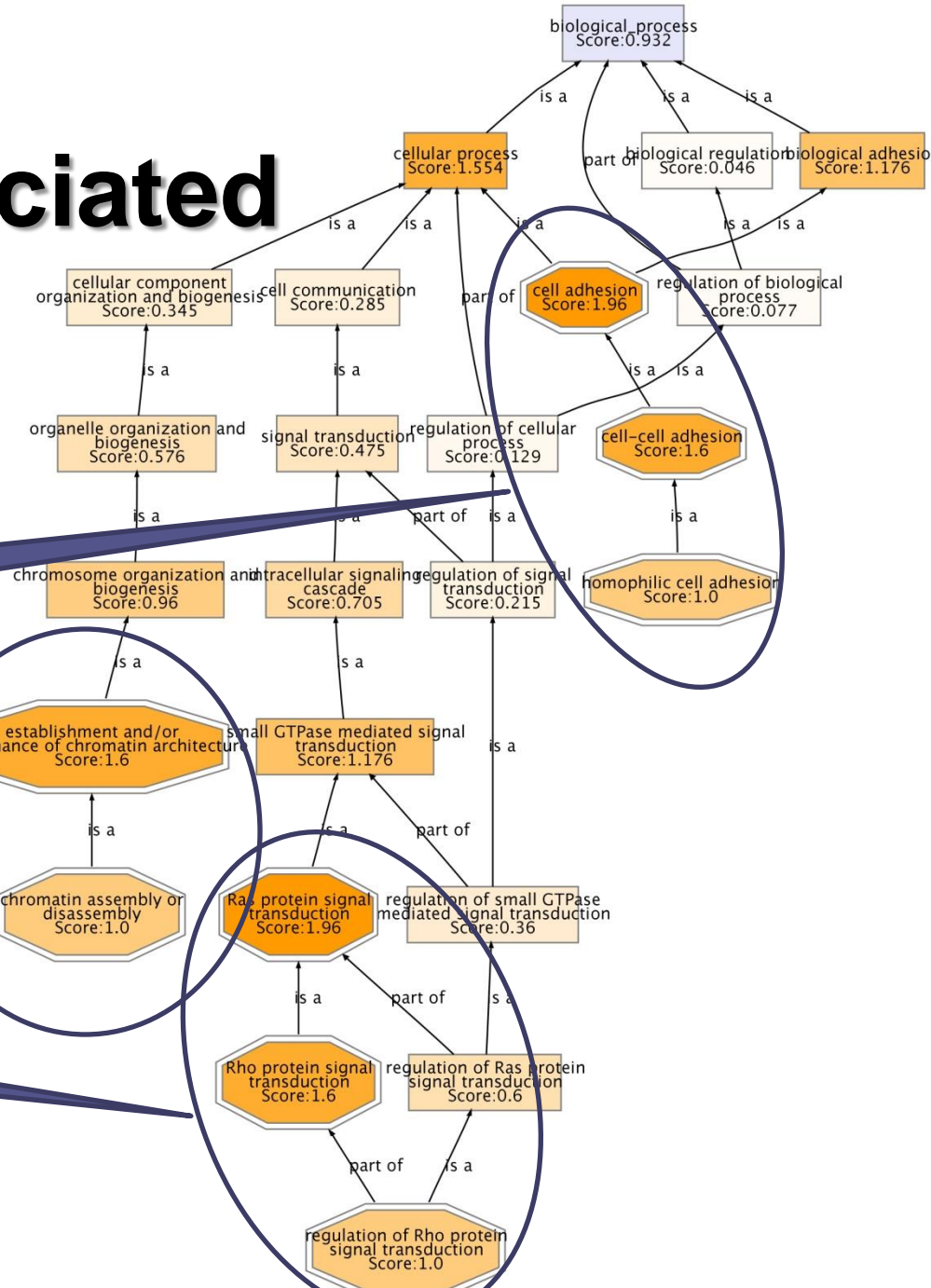
controls. Affy 500K

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

FGFR2

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

# GO processes significantly associated to breast cancer



Metastasis

Chromosomal instability

Rho pathway



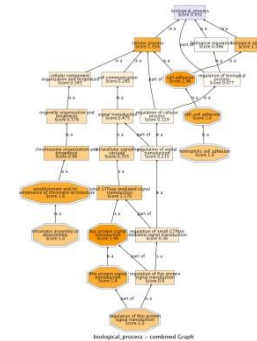
# From gene-based to function-based perspective

SNPs,  
Gene expression

Gene<sub>1</sub>  
Gene<sub>2</sub>  
Gene<sub>3</sub>  
Gene<sub>4</sub>  
:  
:  
:  
:  
Gene<sub>22000</sub>

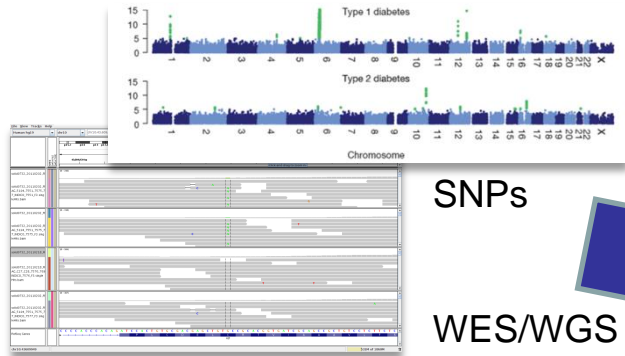


Gene  
Ontology



	SNPs, gene exp.	GO
Detection power	Low (only very prevalent genes)	high
Annotations available	many	many
Use	Biomarker	Illustrative, give hints

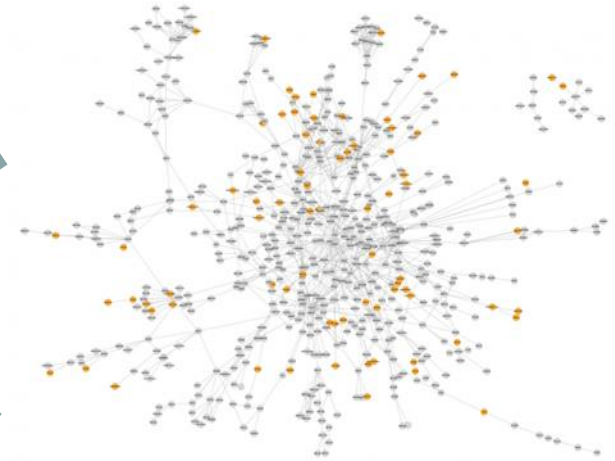
# From gene-based to function-based perspective



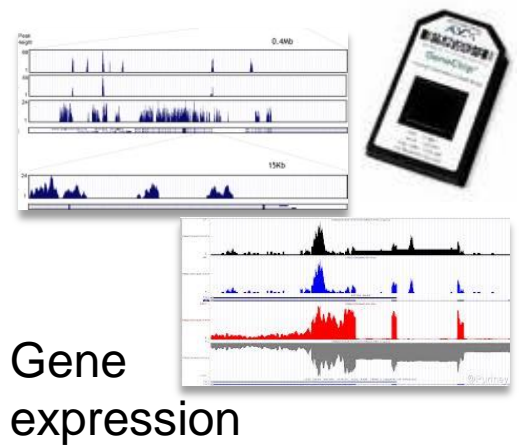
AND/OR



Using protein interaction networks as an scaffold to interpret the genomic data in a functionally-derived context



What part of the interactome is active and/or is damaged



# Network analysis helps to find disease genes in complex diseases

Research

Open Access

## Four new loci associations discovered by pathway-based and network analyses of the genome-wide variability profile of Hirschsprung's disease

Raquel Ma Fernández<sup>1,2</sup>, Marta Bleda<sup>2,3</sup>, Rocío Núñez-Torres<sup>1,2</sup>, Ignacio Medina<sup>3,4</sup>, Berta Luzón-Toro<sup>1,2</sup>, Luz García-Alonso<sup>3</sup>, Ana Torroglosa<sup>1,2</sup>, Martina Marbà<sup>3,4</sup>, Ma Valle Enguix-Riego<sup>1,2</sup>, David Montaner<sup>3</sup>, Guillermo Antiñolo<sup>1,2</sup>, Joaquín Dopazo<sup>2,3,4\*</sup> and Salud Borrego<sup>1,2\*</sup>

\* Corresponding authors: Joaquín Dopazo [jdopazo@cipf.es](mailto:jdopazo@cipf.es) - Salud Borrego [salud.borrego.sspa@juntadeandalucia.es](mailto:salud.borrego.sspa@juntadeandalucia.es)

► Author Affiliations

For all author emails, please [log on](#).

Orphanet Journal of Rare Diseases 2012, 7:103 doi:10.1186/1750-1172-7-103

Published: 28 December 2012

Published online 27 July 2012

Nucleic Acids Research, 2012, Vol. 40, No. 20 e158  
doi:10.1093/nar/gks699

## Discovering the hidden sub-network component in a ranked list of genes or proteins derived from genomic experiments

Luz García-Alonso<sup>1</sup>, Roberto Alonso<sup>1</sup>, Enrique Vidal<sup>1</sup>, Alicia Amadoz<sup>1</sup>, Alejandro de María<sup>1</sup>, Pablo Minguéz<sup>2</sup>, Ignacio Medina<sup>1,3</sup> and Joaquín Dopazo<sup>1,3,4,\*</sup>

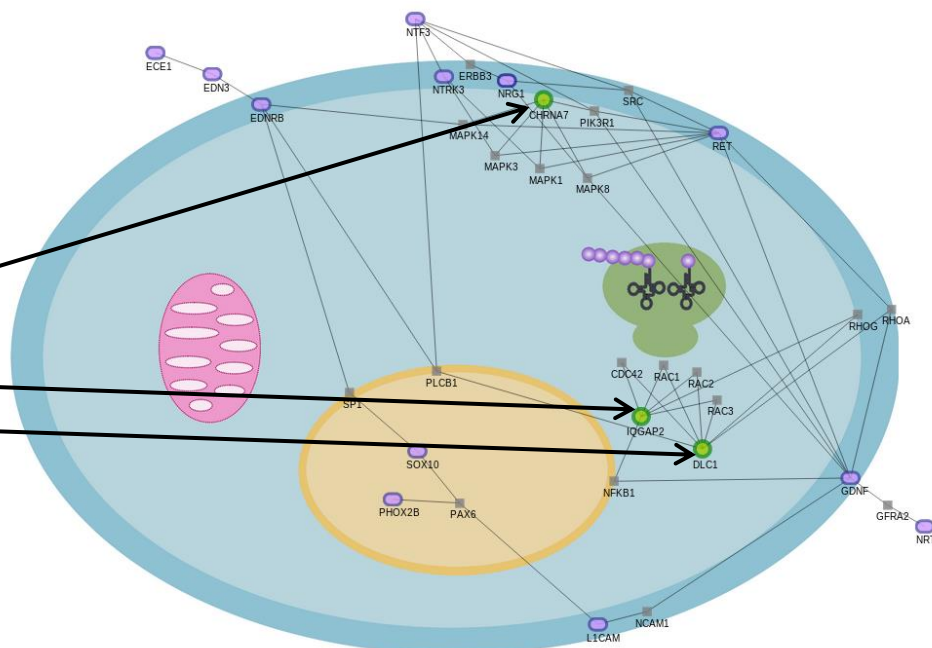
<sup>1</sup>Department of Bioinformatics, Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain, <sup>2</sup>European Molecular Biology Laboratory, Meyerhofstrasse 1, 69117 Heidelberg, Germany, <sup>3</sup>Functional Genomics Node (INB) at CIPF, Valencia and <sup>4</sup>CIBER de Enfermedades Raras (CIBERER), Valencia, Spain

Received March 14, 2012; Revised June 1, 2012; Accepted June 26, 2012

*CHRNA7* (rs2175886 p = 0.000607)  
*IQGAP2* (rs950643 p = 0.0003585)  
*DLC1* (rs1454947 p = 0.007526)  
*RASGEF1A*\* (rs1254964 p = 3.856x10<sup>-05</sup>)

\*no interactions known (yet)

SNPs validated in independent cohorts



Nucleic Acids Research Advance Access published May 19, 2009

Nucleic Acids Research, 2009, 1-6  
doi:10.1093/nar/gkp402

## SNOW, a web-based tool for the statistical analysis of protein-protein interaction networks

Pablo Minguéz<sup>1</sup>, Stefan Götz<sup>1,2</sup>, David Montaner<sup>1</sup>, Fatima Al-Shahrour<sup>1</sup> and Joaquin Dopazo<sup>1,2,3,\*</sup>

<sup>1</sup>Department of Bioinformatics and Genomics, Centro de Investigación Príncipe Felipe (CIPF), <sup>2</sup>CIBER de Enfermedades Raras (CIBERER) and <sup>3</sup>Functional Genomics Node (INB) at CIPF, Valencia, Spain

Received January 21, 2009; Revised April 22, 2009; Accepted May 2, 2009

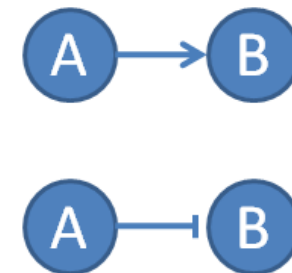
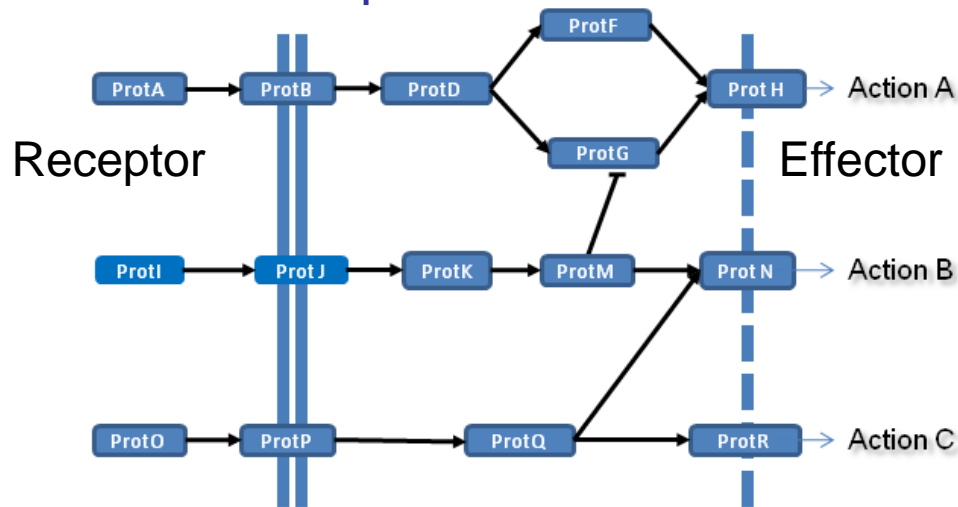
# From gene-based to function-based perspective

	<b>SNPs, gene expression, etc.</b>	<b>GO</b>	<b>Protein interaction networks</b>
<b>Detection power</b>	Low (only very prevalent genes)	High	High
<b>Information coverage</b>	Almost all	Almost all	Less (~9000 genes in human)
<b>Use</b>	Biomarker	Illustrative, give hints	Biomarker*

*\*Need of extra information (e.g. GO) to provide functional insights in the findings*

# From gene-based to mechanism-based perspective

**Transforming gene expression** values into another value that **accounts for a function**. Easiest example of modeling function: **signaling pathways**. Function: transmission of a signal from a receptor to an effector



**Activations  
and  
repressions  
occur**

	ProtH	ProtN	ProtR
ProtA	1	0	0
ProtI	1	1	0
ProtQ	0	1	1
function	Action A	Action B	Action C

# Modeling pathways

Sebastian-Leon et al. *BMC Systems Biology* 2014, **8**:121  
<http://www.biomedcentral.com/1752-0509/8/121>



METHODOLOGY ARTICLE

Open Access

## Understanding disease mechanisms with models of signaling pathway activities

Patricia Sebastián-Leon<sup>1</sup>, Enrique Vidal<sup>1,2,3</sup>, Pablo Minguez<sup>1,4</sup>, Ana Conesa<sup>1</sup>, Sonia Tarazona<sup>1</sup>, Alicia Amadoz<sup>1</sup>, Carmen Armero<sup>5</sup>, Francisco Salavert<sup>1,2</sup>, Antonio Vidal-Puig<sup>6</sup>, David Montaner<sup>1</sup> and Joaquín Dopazo<sup>1,2,7\*</sup>

Published online 8 June 2013

*Nucleic Acids Research*, 2013, Vol. 41, Web Server issue W213–W217  
 doi:10.1093/nar/gkt451

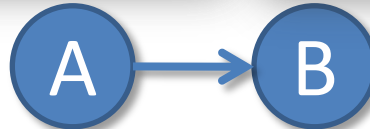
## Inferring the functional effect of gene expression changes in signaling pathways

Patricia Sebastián-León<sup>1</sup>, José Carbonell<sup>1</sup>, Francisco Salavert<sup>1,2</sup>, Rubén Sanchez<sup>3</sup>, Ignacio Medina<sup>1</sup> and Joaquín Dopazo<sup>1,2,4,\*</sup>

<sup>1</sup>Department of Computational Genomics, Centro de Investigación Príncipe Felipe (CIPF), Valencia 46012, Spain, <sup>2</sup>CIBER de Enfermedades Raras (CIBERER), Valencia 46012, Spain, <sup>3</sup>Genometra S.L., Valencia, Spain and <sup>4</sup>Functional Genomics Node (INB) at CIPF, Valencia 46012, Spain

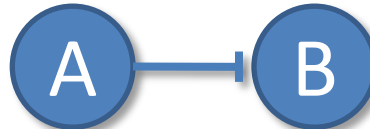
Received March 3, 2013; Revised April 16, 2013; Accepted May 2, 2013

Activation



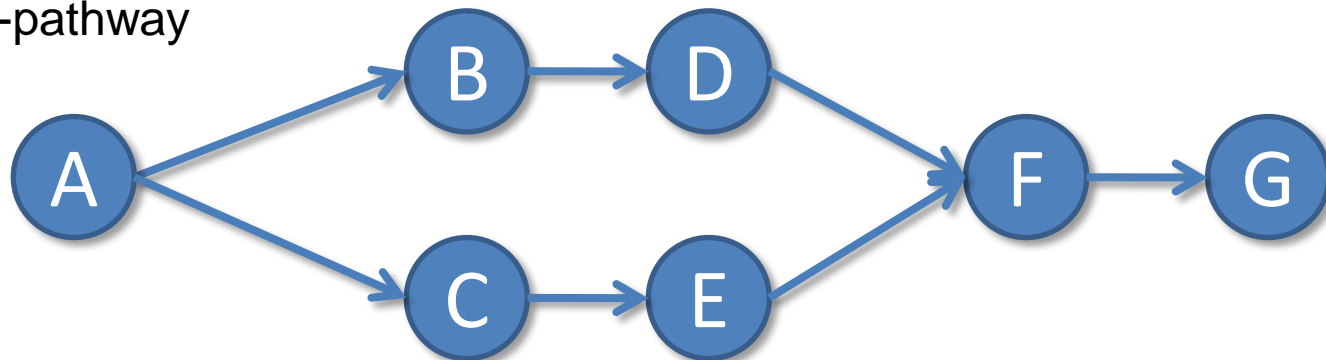
Prob. =  $P(A \text{ activated})P(B \text{ activated})$

Inhibition



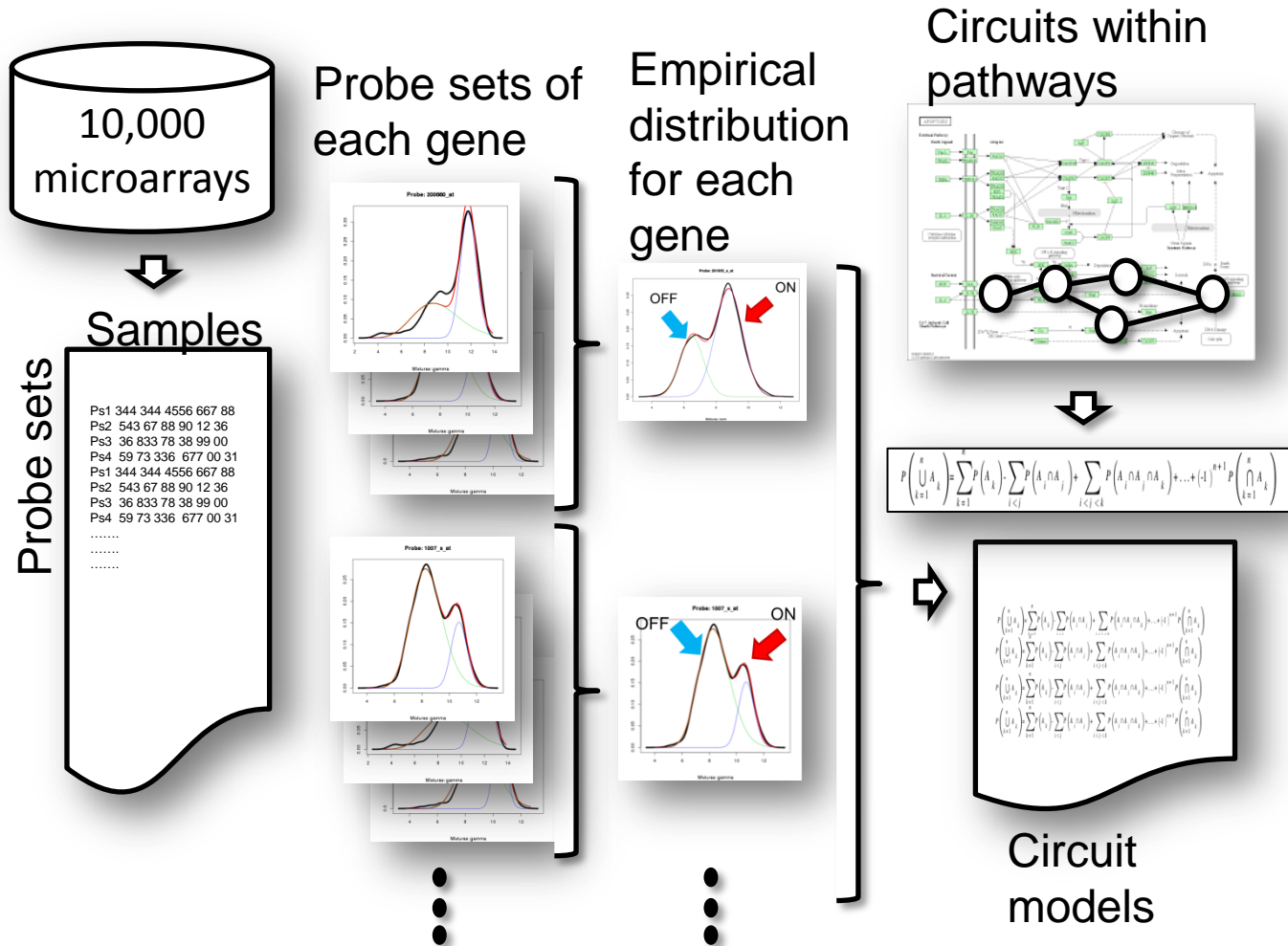
Prob. =  $[1 - P(A \text{ activated})]P(B \text{ activated})$

Sub-pathway



$$P(A \rightarrow G \text{ activated}) = P(A)P(B)P(D)P(F)P(G) + P(A)P(C)P(E)P(F)P(G) - P(A)P(F)P(G)P(B)P(C)P(D)P(E)$$

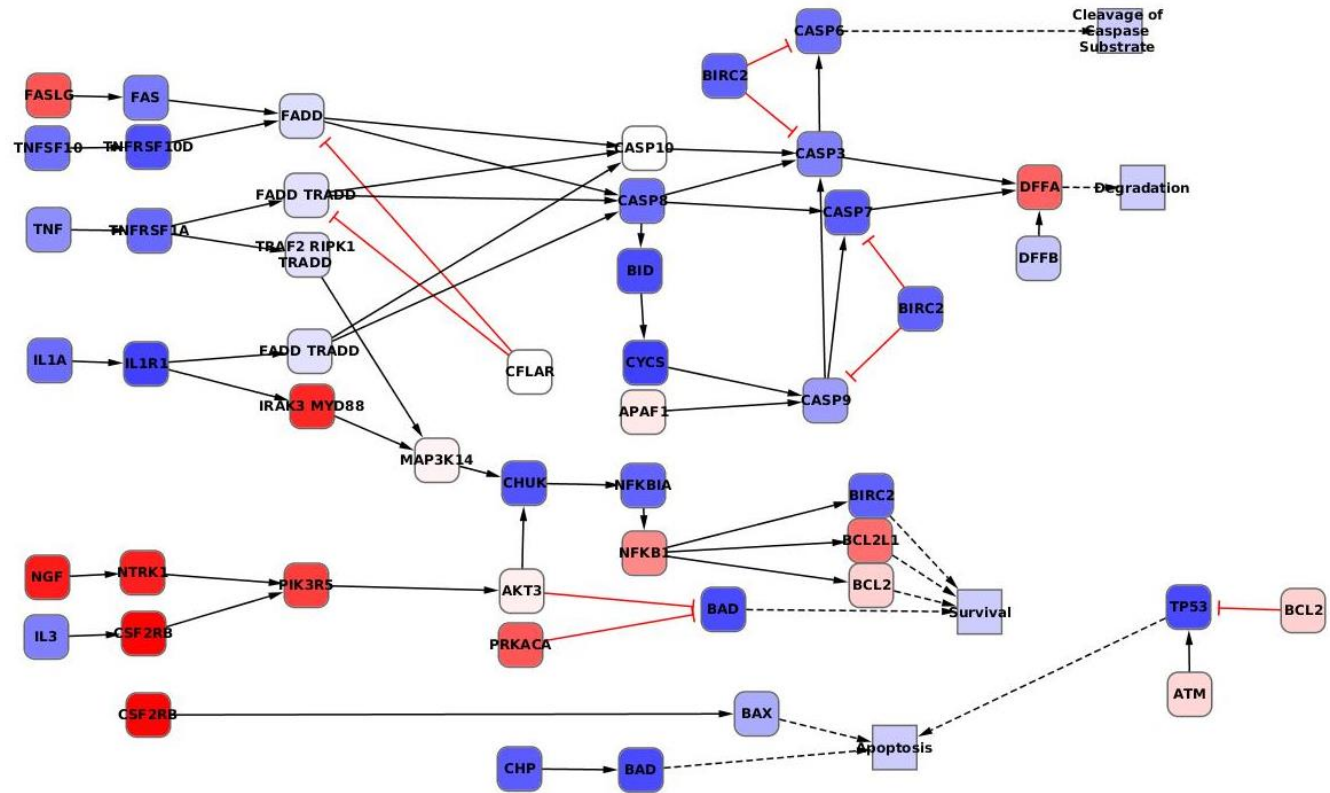
# Modeling activation probabilities of circuits



# The effects of changes in gene activity are not obvious

What would you predict about the consequences of gene activity changes in the apoptosis pathway in a case control experiment of colorectal cancer?

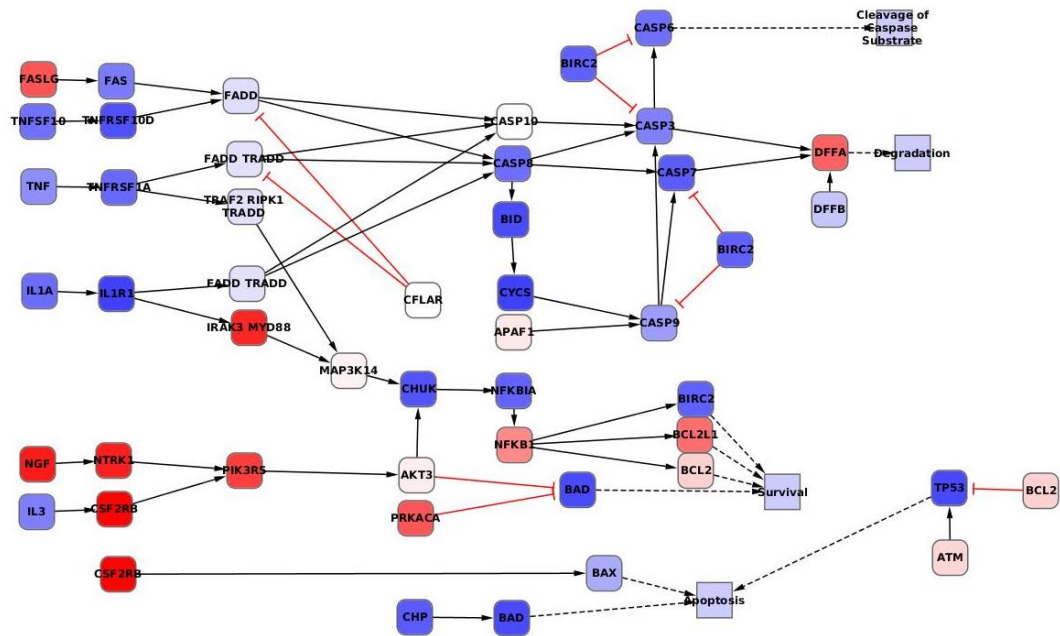
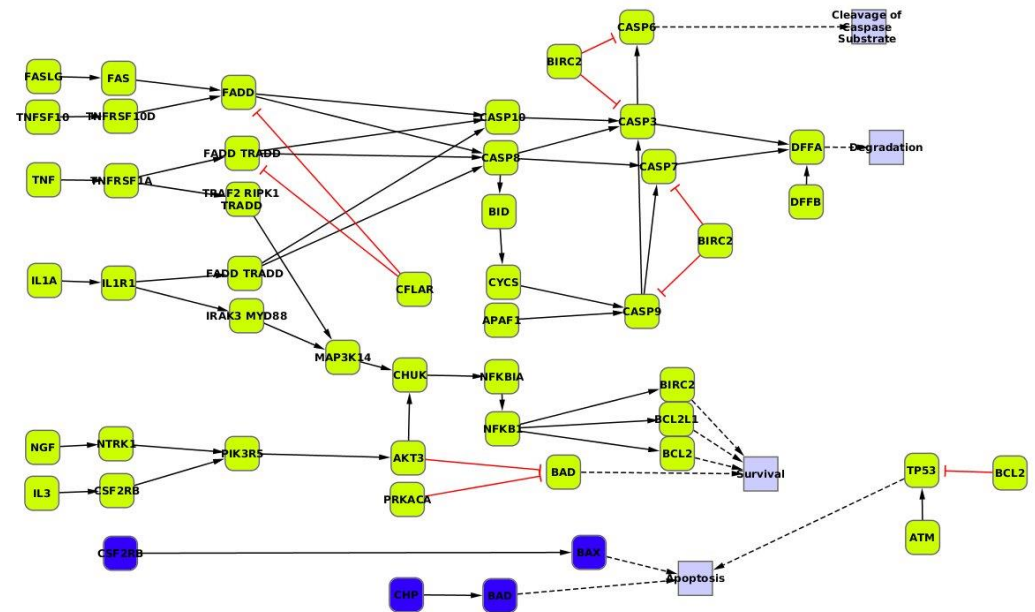
The figure shows the gene up-regulations (red) and down-regulations (blue)





# Apoptosis inhibition is not obvious from gene expression

Two of the three possible sub-pathways leading to apoptosis are inhibited in colorectal cancer. Upper panel shows the inhibited sub-pathways in blue. Lower panel shows the actual gene up-regulations (red) and down-regulations (blue) that justify this change in the activity of the sub-pathways



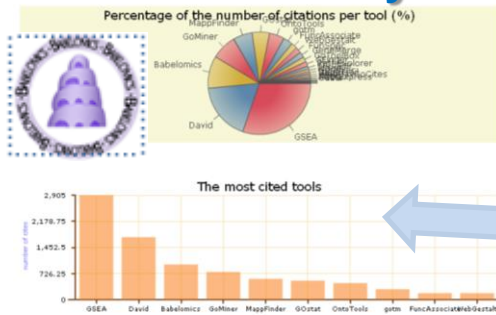
# From gene-based to function-based perspective

	SNPs, gene expression, etc.	GO	Protein interaction networks	Models of cellular functions
Detection power	Low (only very prevalent genes)	High	High	Very high
Information coverage	Almost all	Almost all	Low (~9000 genes in human)	Low (~6700 genes in human)*
Use	Biomarker	Illustrative, give hints	Biomarker	Biomarker that explain disease mechanism

\*Only ~800 genes in human signaling pathways

# Software development

## Functional analysis



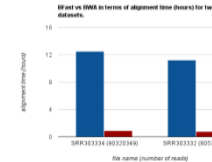
Babelomics is the third most cited tool for functional analysis. Includes more than 30 tools for advanced, systems-biology based data analysis



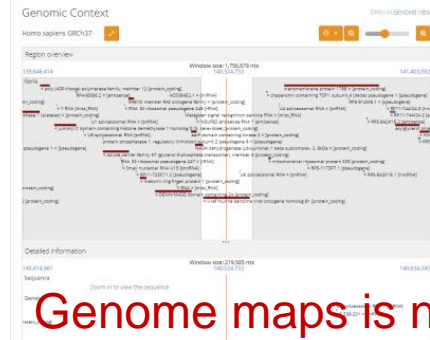
See interactive map of for the last 24h use <http://bioinfo.cipf.es/toolsusage>

## Mapping

HPC on CPU, SSE4,  
GPUs on NGS data  
processing  
Speedups up to 40X



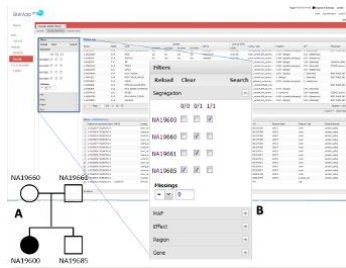
## Visualization



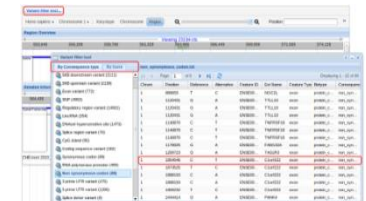
Ultrafast  
genome  
viewer with  
google  
technology

Genome maps is now part  
of the ICGC data portal

## Variant prioritization

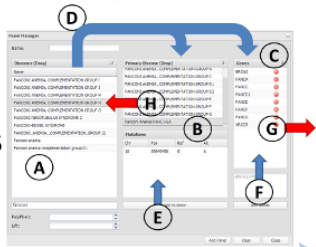


## Variant annotation

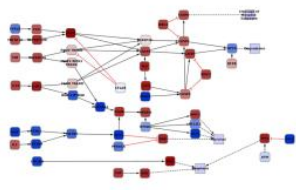


## Diagnostic

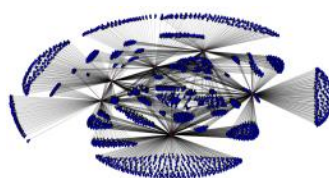
NGS  
panels



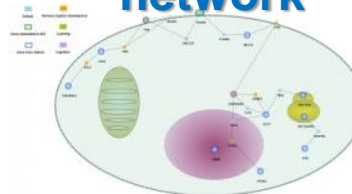
## Signaling network



## Regulatory network



## Interaction network



## CellBase

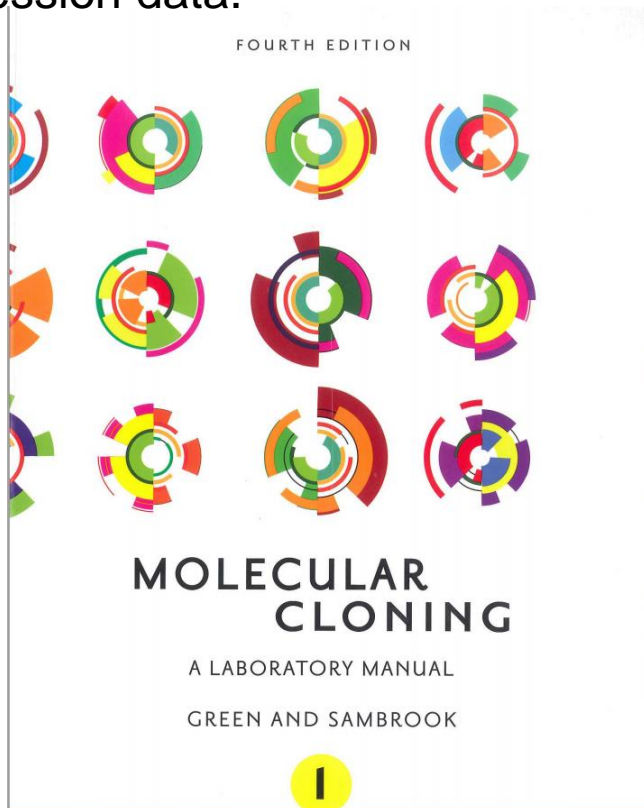


Knowledge  
database

More than 150.000 experiments were analyzed in our tools during the last year

# Babelomics in the Maniatis

The Babelomics suite of programs becomes a classic. Now is cited as a method in the last edition of **Molecular Cloning**, the popular **Maniatis**. The protocol 4 of chapter 8, Expression Profiling by Microarray and RNA-seq, contains a description on how to use Babelomics to analyze expression data.



578 / Chapter 8

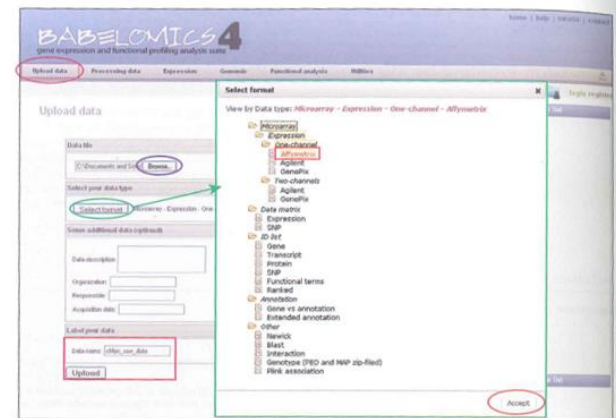


FIGURE 1. Babelomics data uploading form. Click "Browse" to upload the data file named cMyc.zip. Select "Affymetrix" as the format, click "Accept" in the pop-up "Select format" panel, and assign the name as "cMyc\_raw\_data." Click "Upload" to submit the file.

- iii. Assign "cMyc\_raw\_data" as the data name.
  - iv. Click "Upload" to submit the files and wait for the validation to complete. All submitted data are listed in the "Data list" panel.
3. When the data submission is finished, its status in the "Data list" panel changes to "valid."
- i. To check the expression intensity of the raw data before normalization, click the "Microarray raw-data plot" link in the "Utilities" tab.
  - ii. In the page followed by the link, click "browse server," select "Uploaded data" → "cMyc\_raw\_Data," and click "Accept."
  - iii. Set the job name as "CMyc\_original\_boxplot" and click "Run."
  - iv. After the job is finished, click it in the "Job list" panel and the "Box-plots" link to view the box plots as shown in Figure 2.

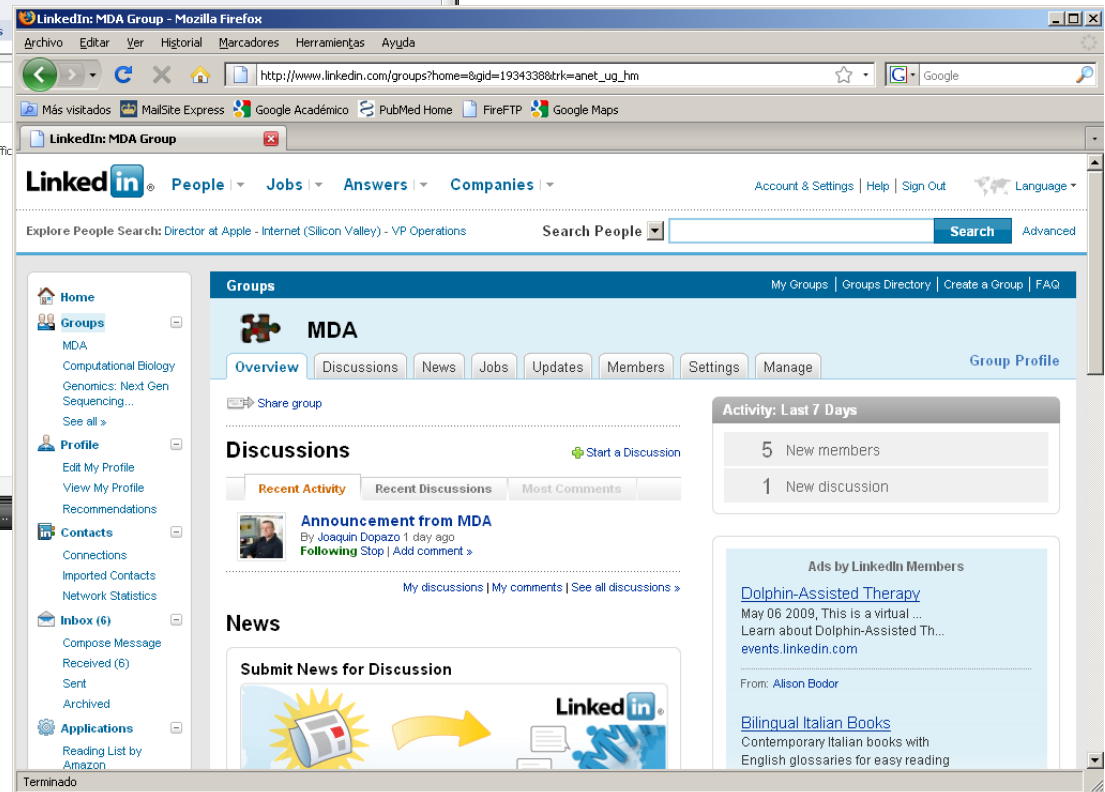
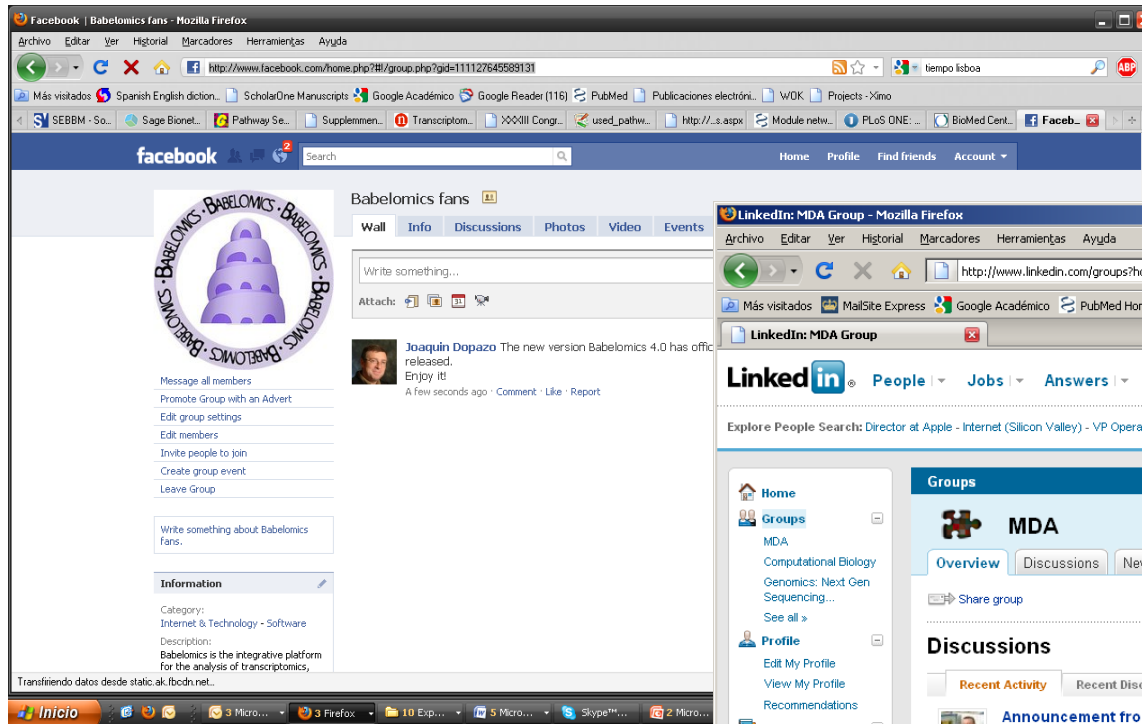
*Each box plot displays summary statistics of a sample, with the box containing the middle 50% of the data, the upper (lower) edge of the box indicating 75th (25th) percentile of the data, and the vertical lines (whiskers) indicating maximum and minimum values. We can see that the eight data sets in our example have systematically different distributions of intensities.*

High impact developments

# SOCIAL:

## MDA group in Linked-in

## Babelomics group in Facebook and twitter



@babelomics

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



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

Follow us on twitter



@xdopazo

@bioinfocipf

