

# A NGS data analysis course

## Introduction

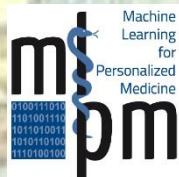
### Joaquín Dopazo

Computational Genomics Department,  
Centro de Investigación Príncipe Felipe (CIPF),  
Functional Genomics Node, (INB),  
Bioinformatics Group (CIBERER) and  
Medical Genome Project,  
Spain.

<http://bioinfo.cipf.es>  
<http://www.medicalgenomeproject.com>  
<http://www.babelomics.org>  
<http://www.hpc4g.org>  
 @xdopazo



University of Cambridge, 23-25 February 2015





PRINCIPE FELIPE  
CENTRO DE INVESTIGACION

@xdopazo  
@bioinfocipf

# 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), and...

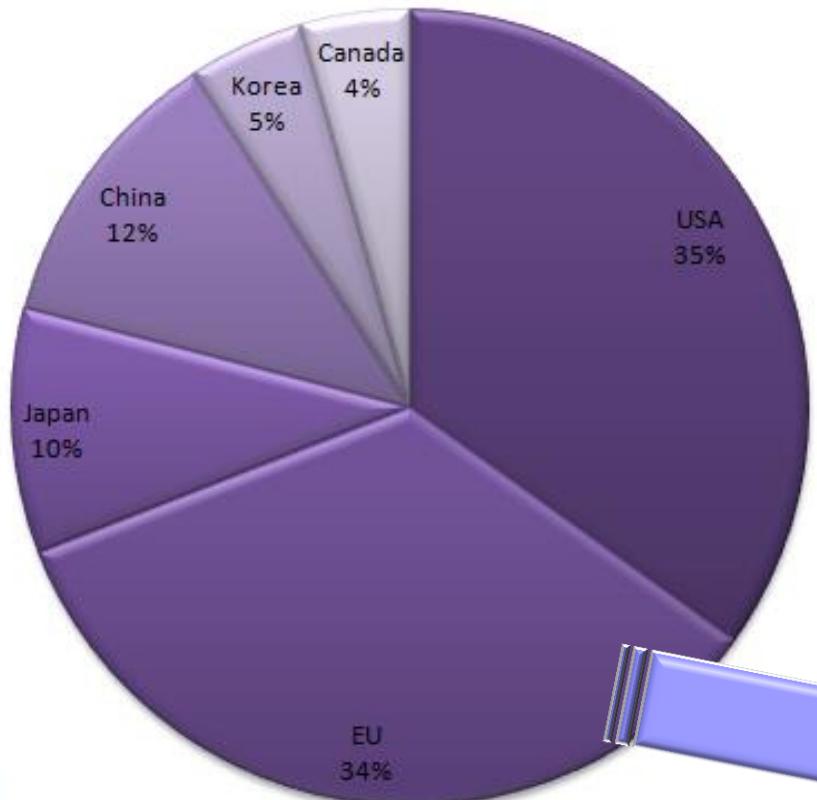


... Computational Biology Lab, HPCS  
University of Cambridge, UK...

...and EMBL-EBI



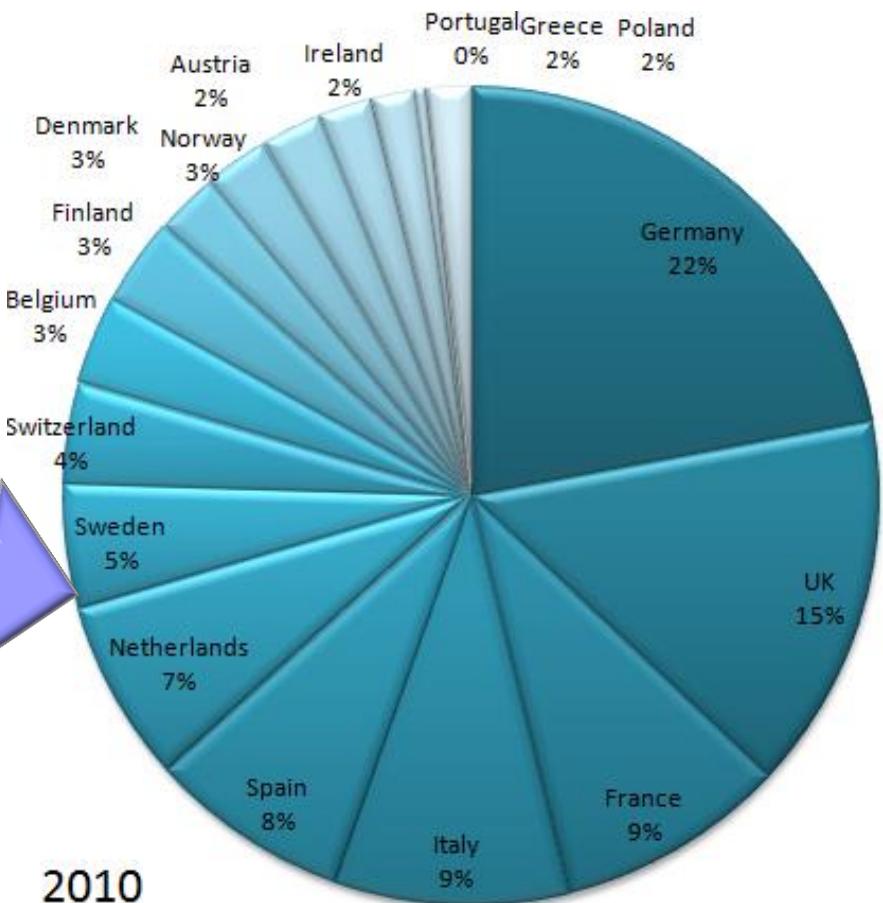
# Some bibliographic data: Microarray publications



2010 Worldwide

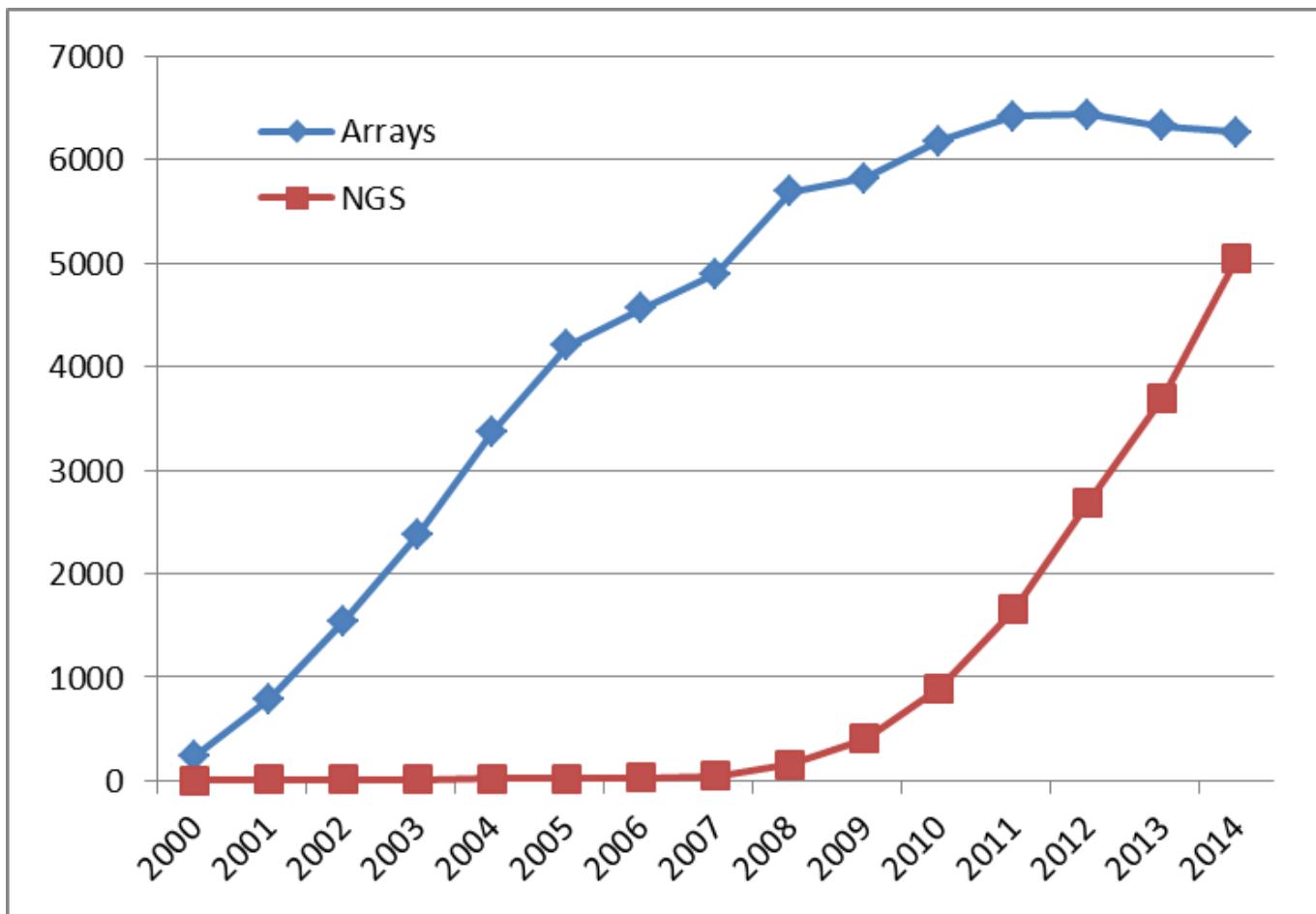
**Source Pubmed. Query:**  
**2009[Entrez Date] AND**  
**country[Affiliation]AND**  
**microarray[Title/Abstract]**

2010 Europe



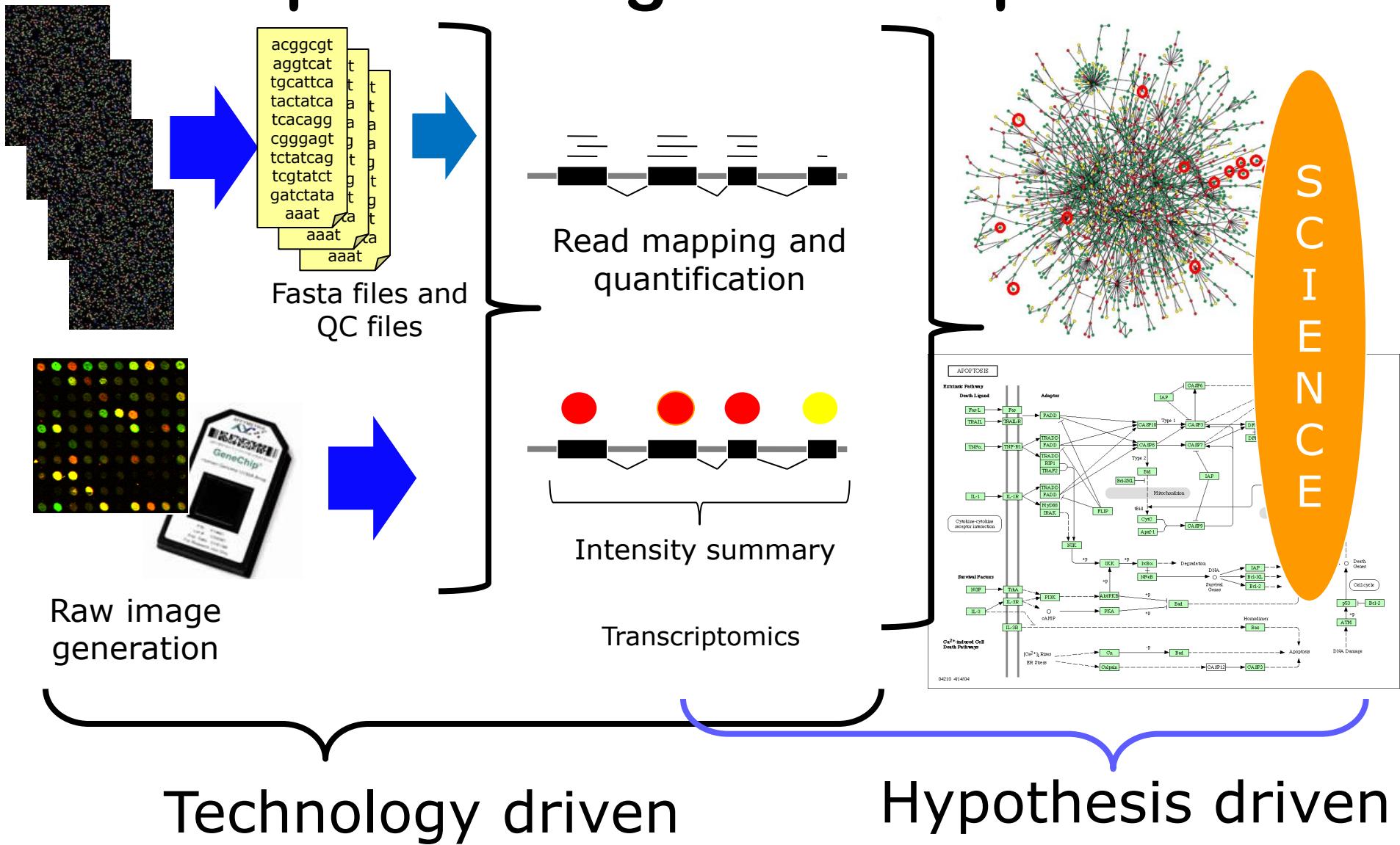
2010

# Evolution of the papers published in microarray and NGS technologies

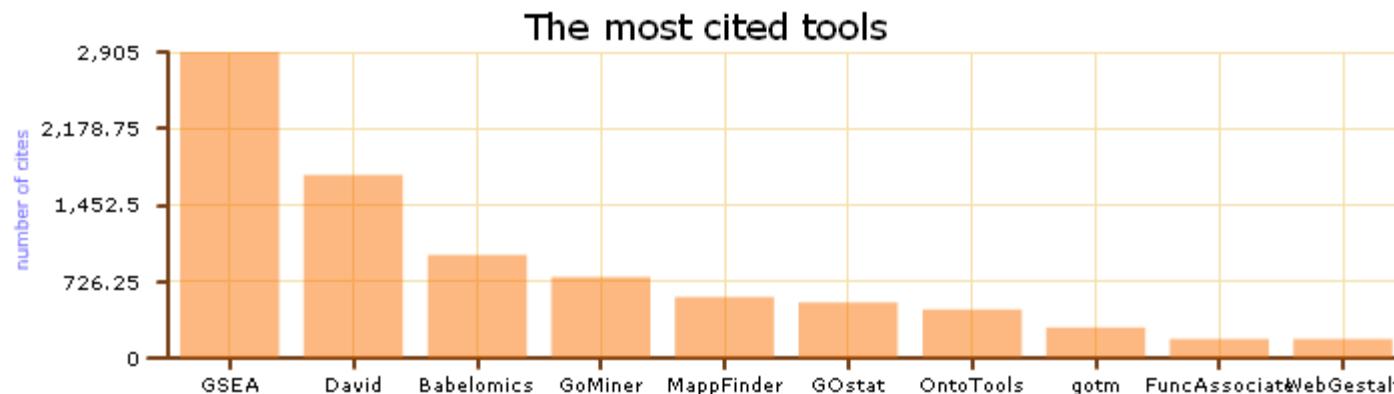
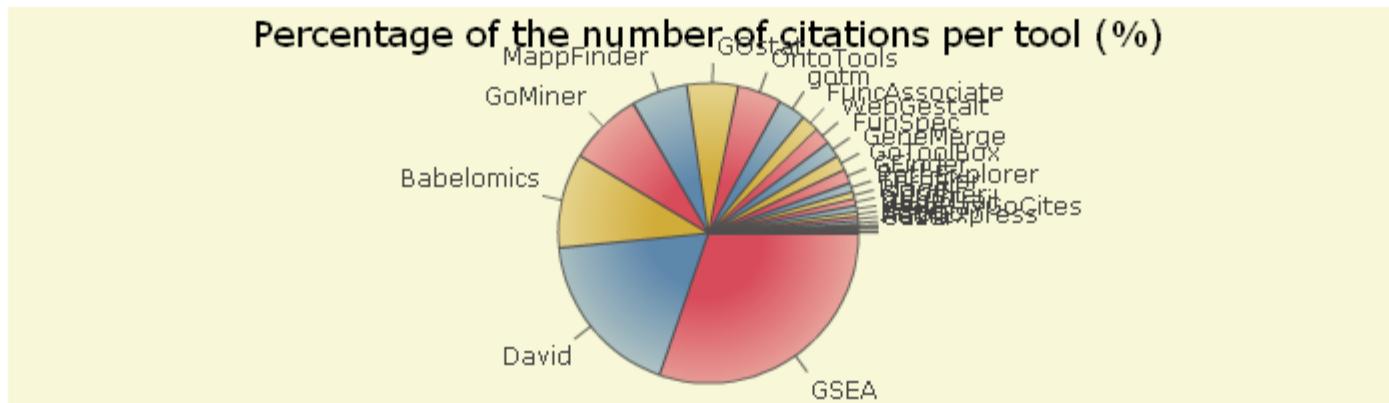


**Source Pubmed. Query:** "high-throughput sequencing"[Title/Abstract] OR "next generation sequencing"[Title/Abstract] OR "rna seq"[Title/Abstract]) AND year[Publication Date]

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



# Tools for genomic data analysis and functional profiling



# Some numbers

529 papers cite Babelomics (plus 798 FatiGO cites)

(source ISI Web of Knowledge, February 2015)

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

More than 1000 experiments per day.



# Background

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

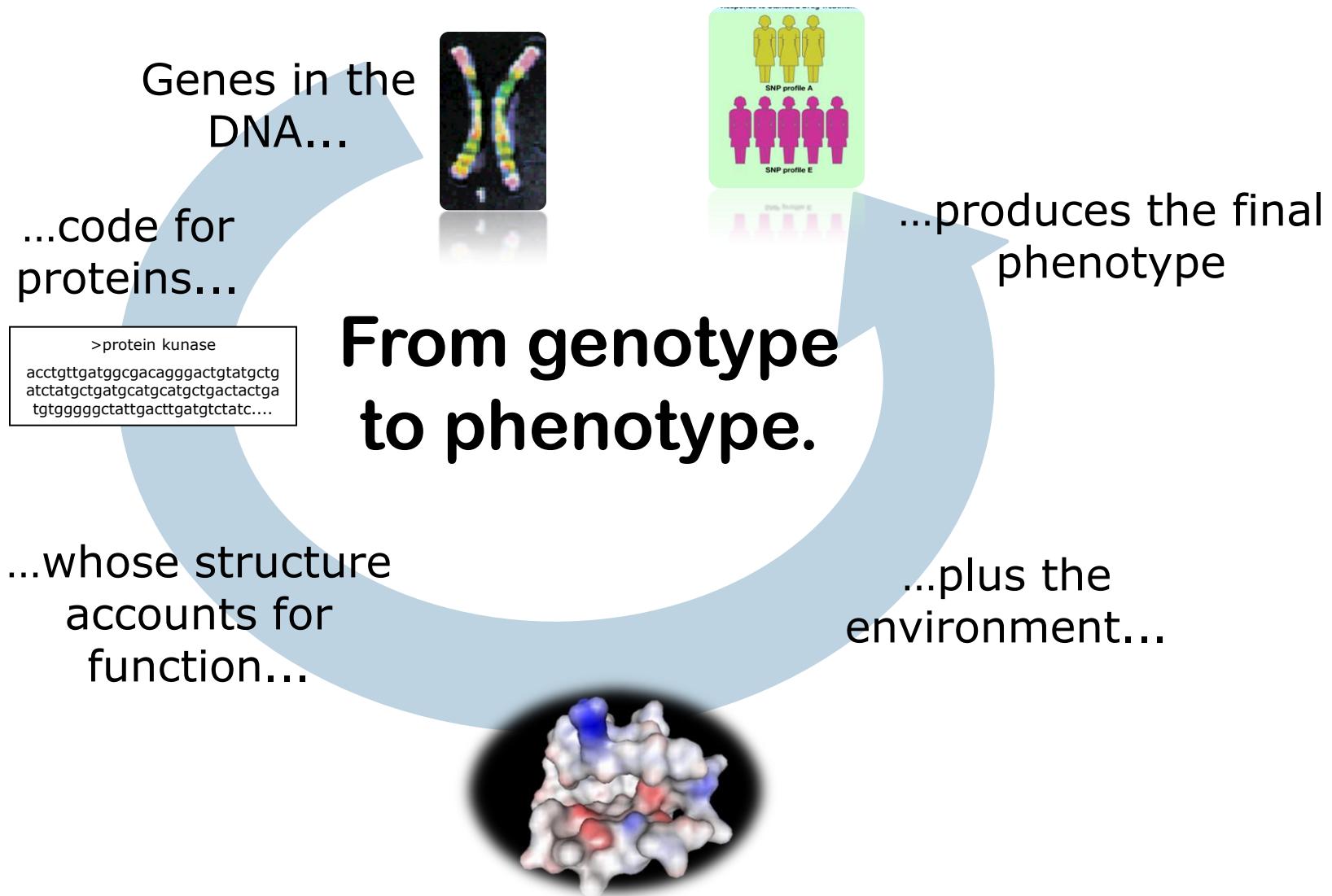
(*William Blake, 28 November 1757 – 12 August 1827, poet, painter, and printmaker*)



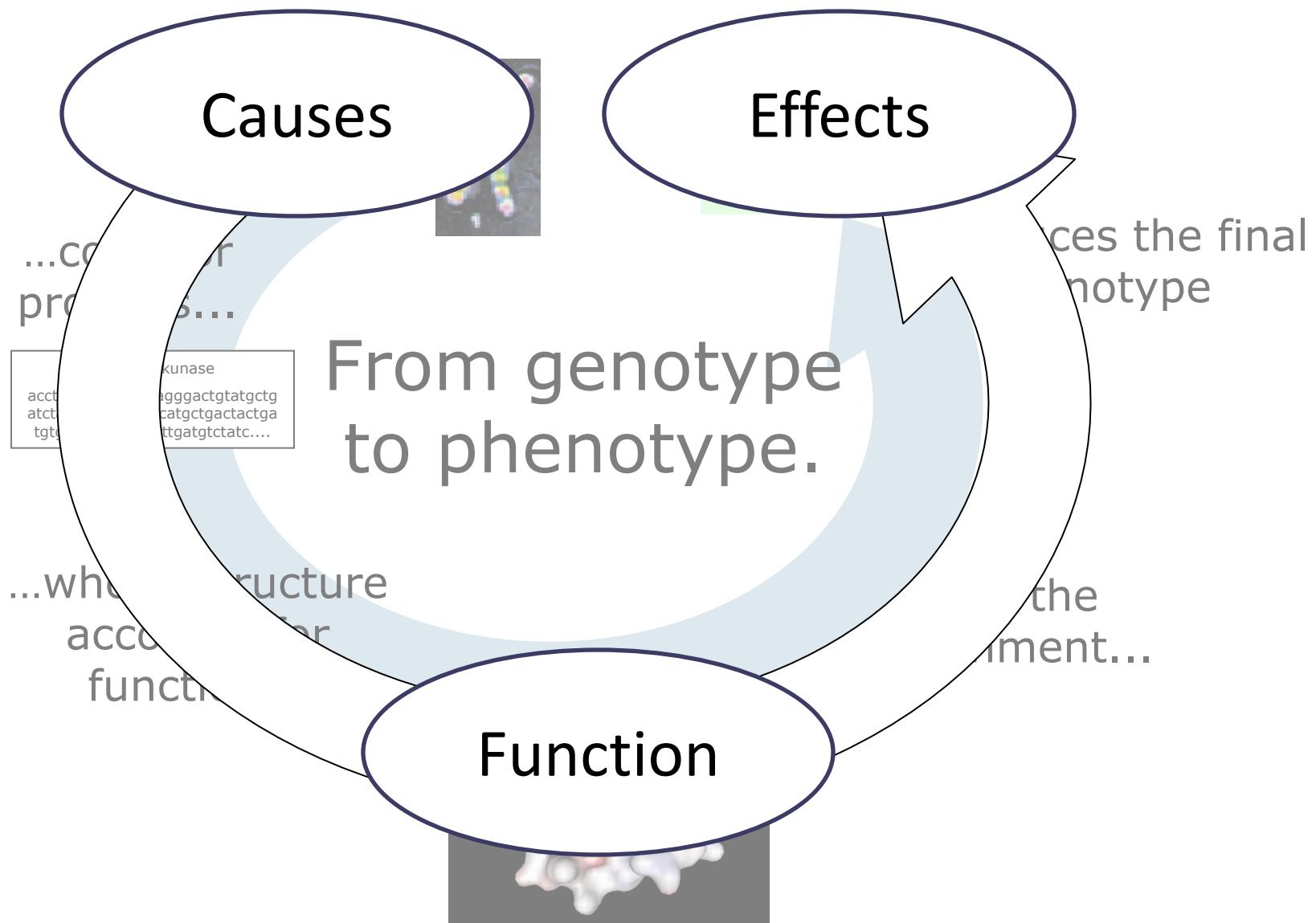
The introduction and popularisation of high-throughput techniques has drastically changed the way in which biological problems **can** be addressed and hypotheses **can** be tested.

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

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



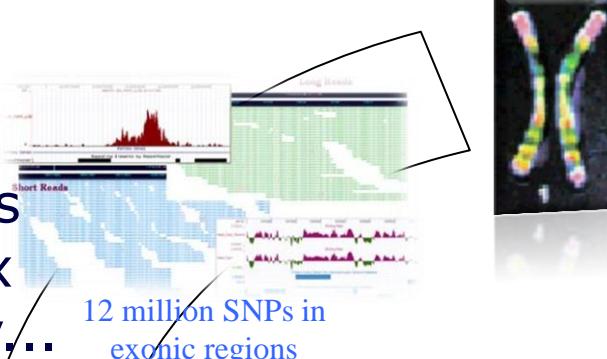
# Reductionistic approach to link causes (genome) to effects (phenotype) through actions (function)



Next Generation Sequencing

$10^9$ bp per round

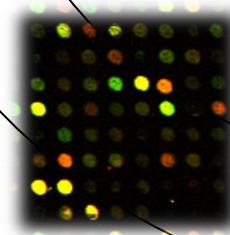
## Genes in the DNA...



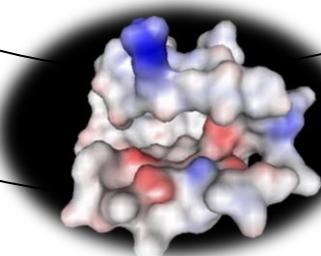
...whose final effect configures the phenotype...

## From genotype to phenotype.

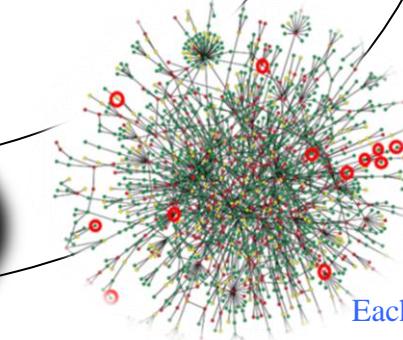
(in the post-genomics scenario)



...code for  
proteins...

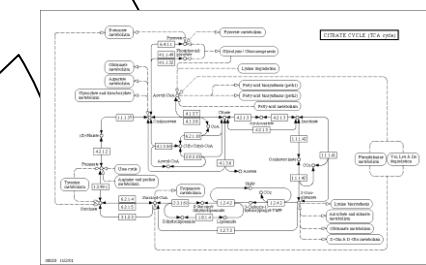


That undergo post-translational modifications, somatic recombination...  
100K-500K proteins



Each protein has an average of 8 interactions

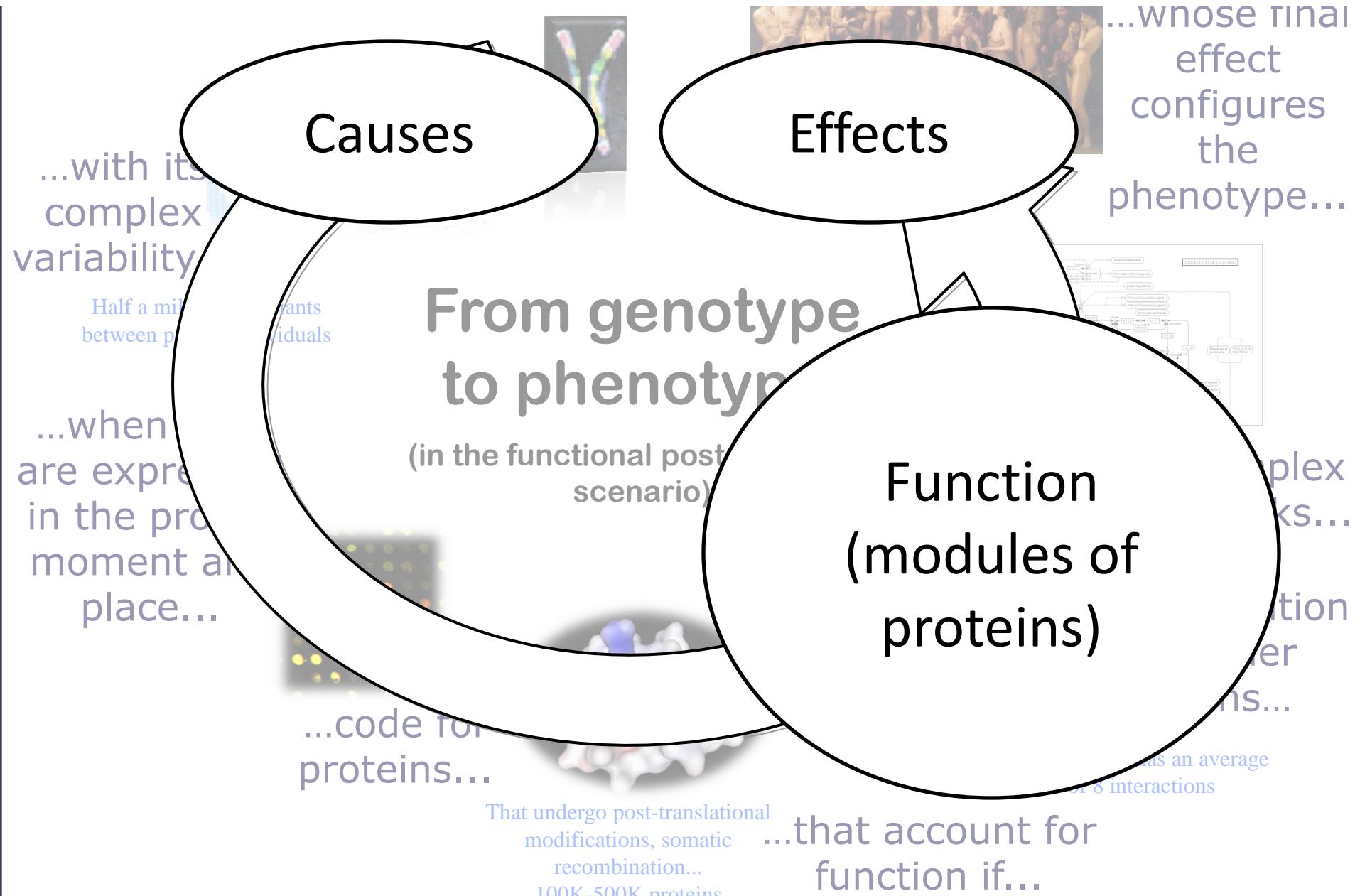
...that account for function if...



...conforming complex interaction networks...

...in cooperation with other proteins...

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



# High-throughput data for functional genomics

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

the DNA

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

## Genotyping

...with its complex variability...

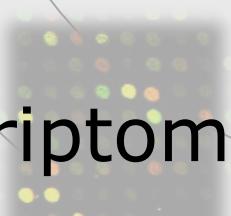
Half a million of variants between pairs of individuals

Genome wide

From genotype to phenotype.

(in the functional post-genomics scenario)

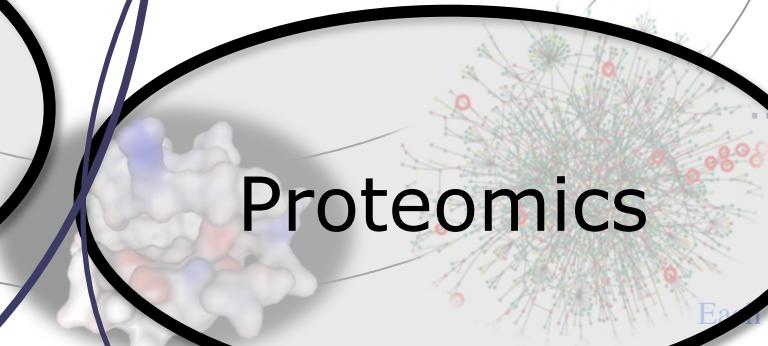
## Transcriptomics



...code for proteins...

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

## Proteomics



...that account for function if...

Almost-omics

Each protein has an average of 8 interactions

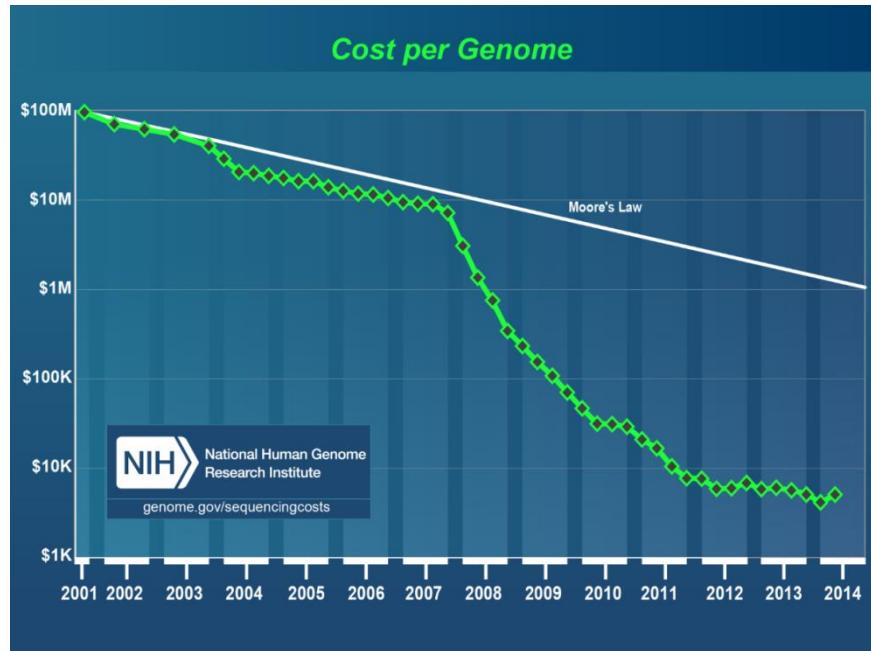
## Metabolomics

...conforming complex interaction networks...

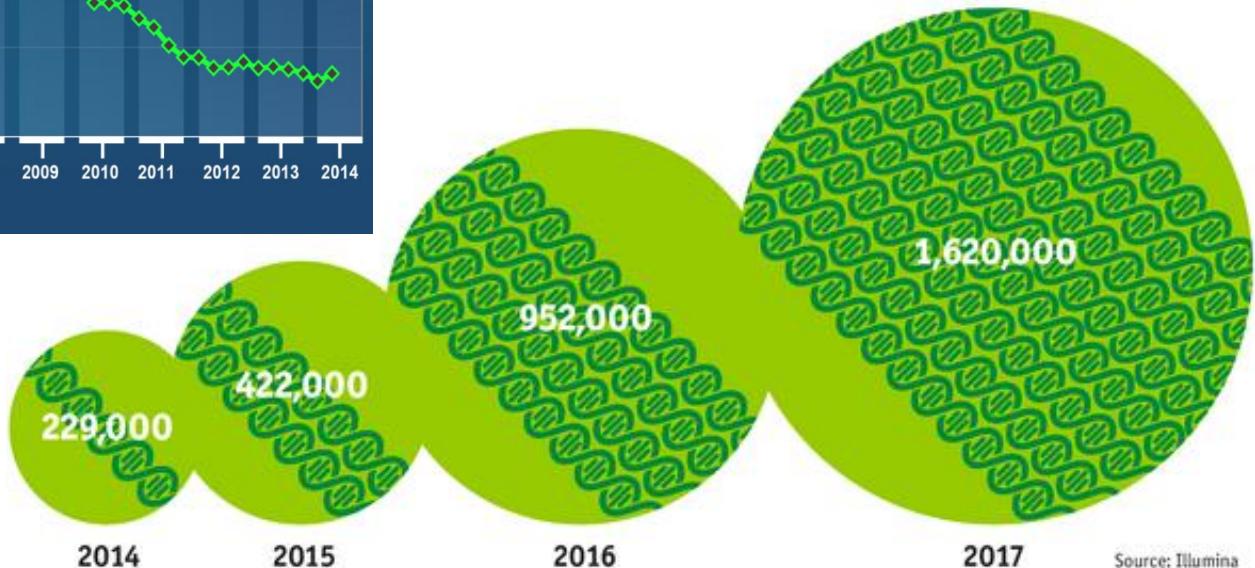
...in cooperation with other proteins...

# The genome sequencing pace

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



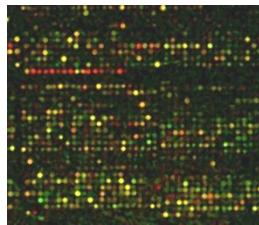
NGS is matching the cost of many conventional clinical tests



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

# 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**)

# Studies must be hypothesis driven.

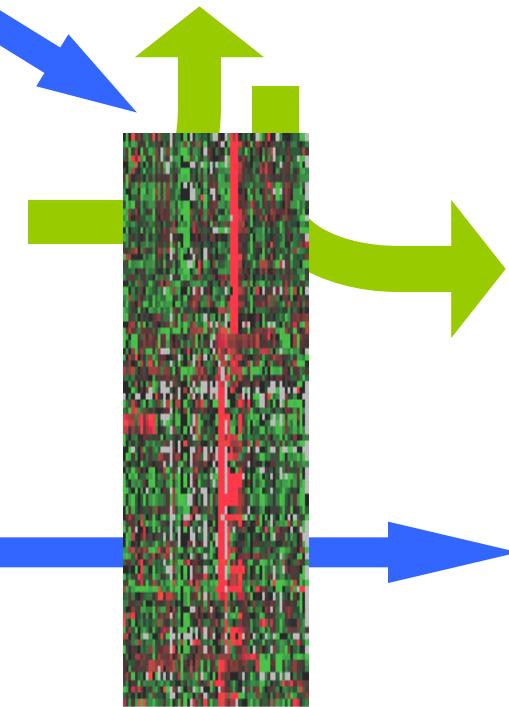
What is the aim? Class discovery? sample classification? gene selection? ...

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

Molecular classification of samples

Co-expressing genes...

Different classes...



Unsupervised



Supervised

What genes are responsible for?

What do they have in common?

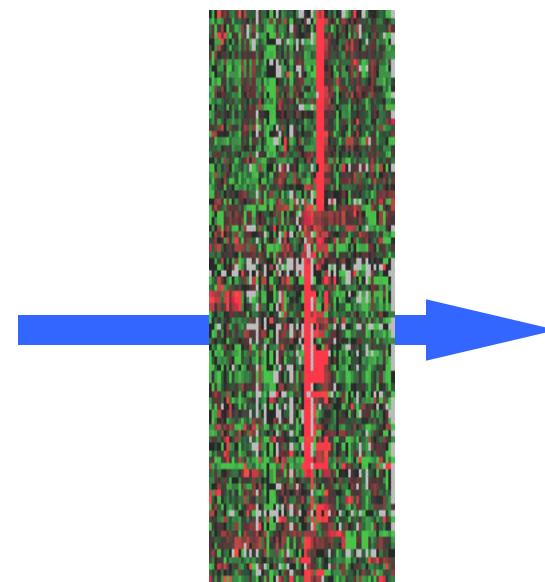
# Unsupervised problem: class discovery

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...



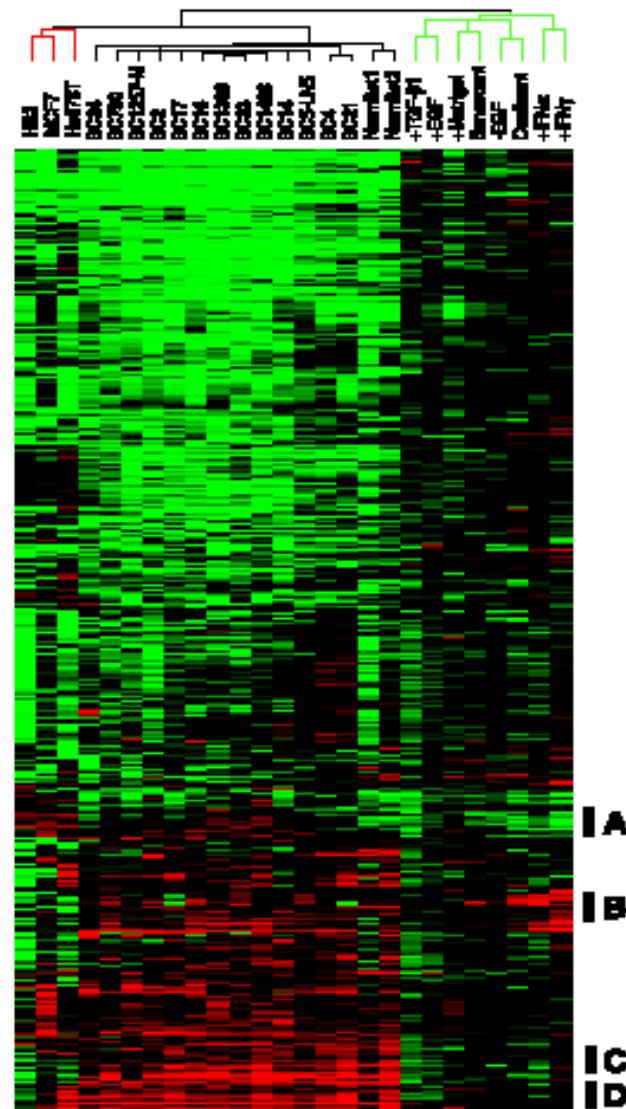
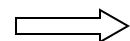
- What genes co-express?
- How many different expression patterns do we have?
- What do they have in common?
- Etc.

# 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?

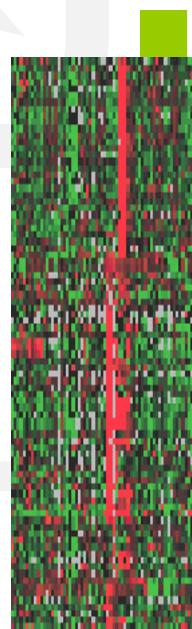
**Different classes...**

Molecular  
classification of  
samples

**What genes are  
responsible for?**

Co-expressing genes...

What do they  
have in  
common?



# Differential gene expression

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

- **Two classes**

- T-test
- Limma
- Fold-change

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

- Pearson
- Spearman
- Regression

- **Multiclass**

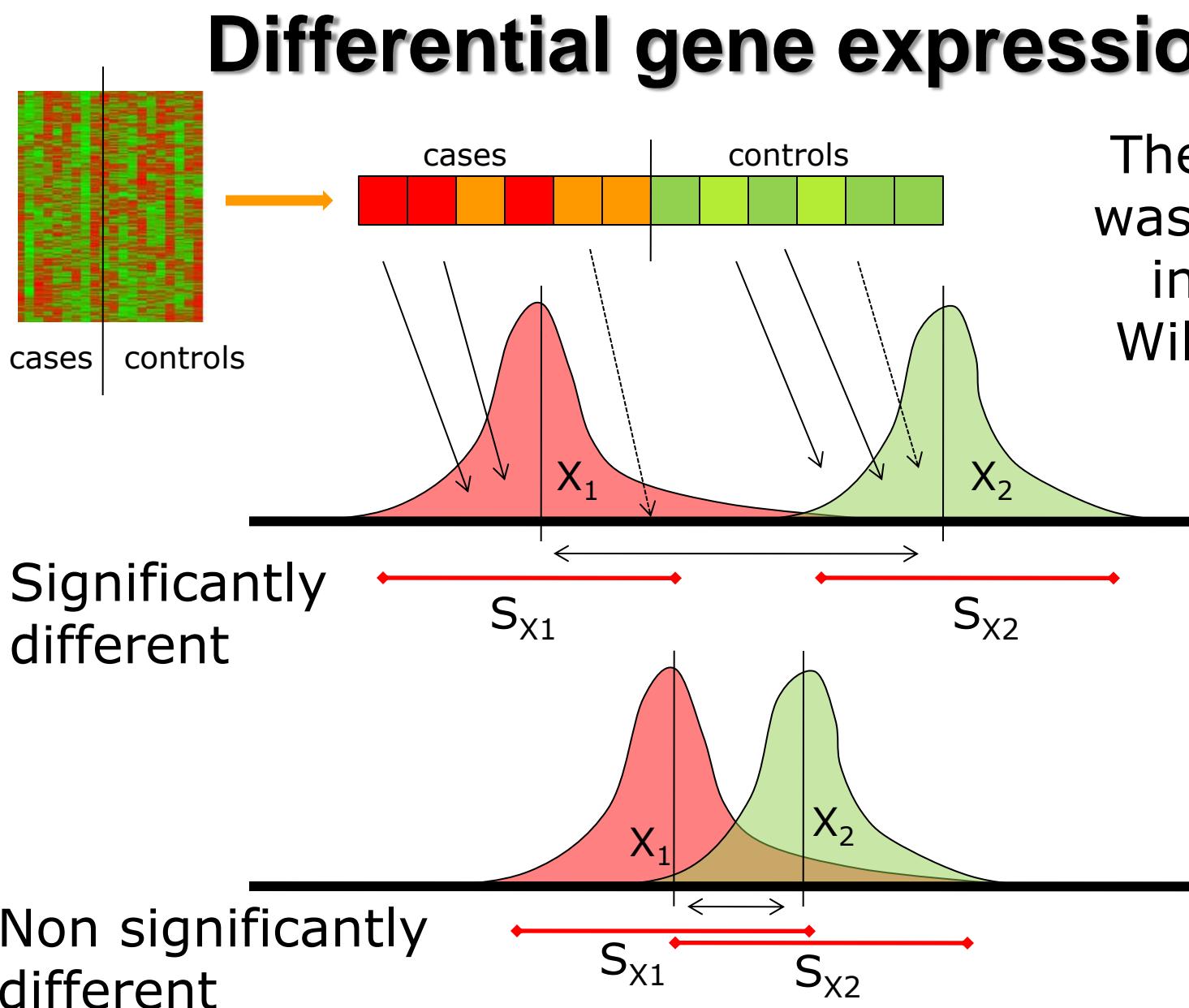
- Anova
- Limma

- **Survival**

- Cox model

- **Time Course**

# Differential gene expression



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

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

# **Supervised problems.**

## **sample classification**

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

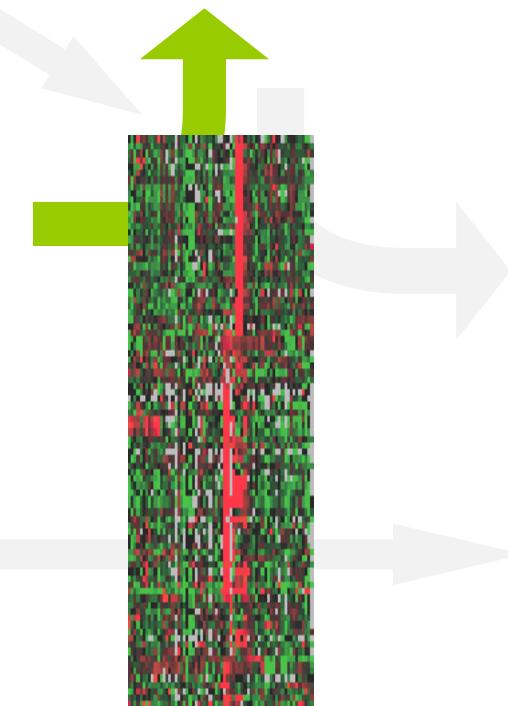
Different classes...

Molecular  
classification of  
samples

What genes are  
responsible for?

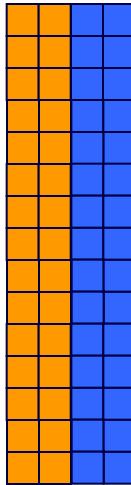
Co-expressing genes...

What do they  
have in  
common?



# Predictors and molecular signatures

A   B   X



Is X,  
A  
or B?

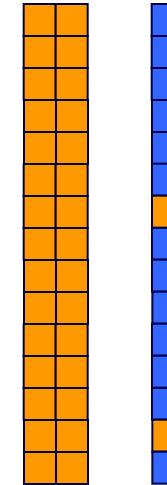


Diff (B, X) = 2



What is a predictor?

Intuitive notion:



Diff (A, X) = 13

Most probably X belongs to class B

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

# Predictor of clinical outcome in breast cancer



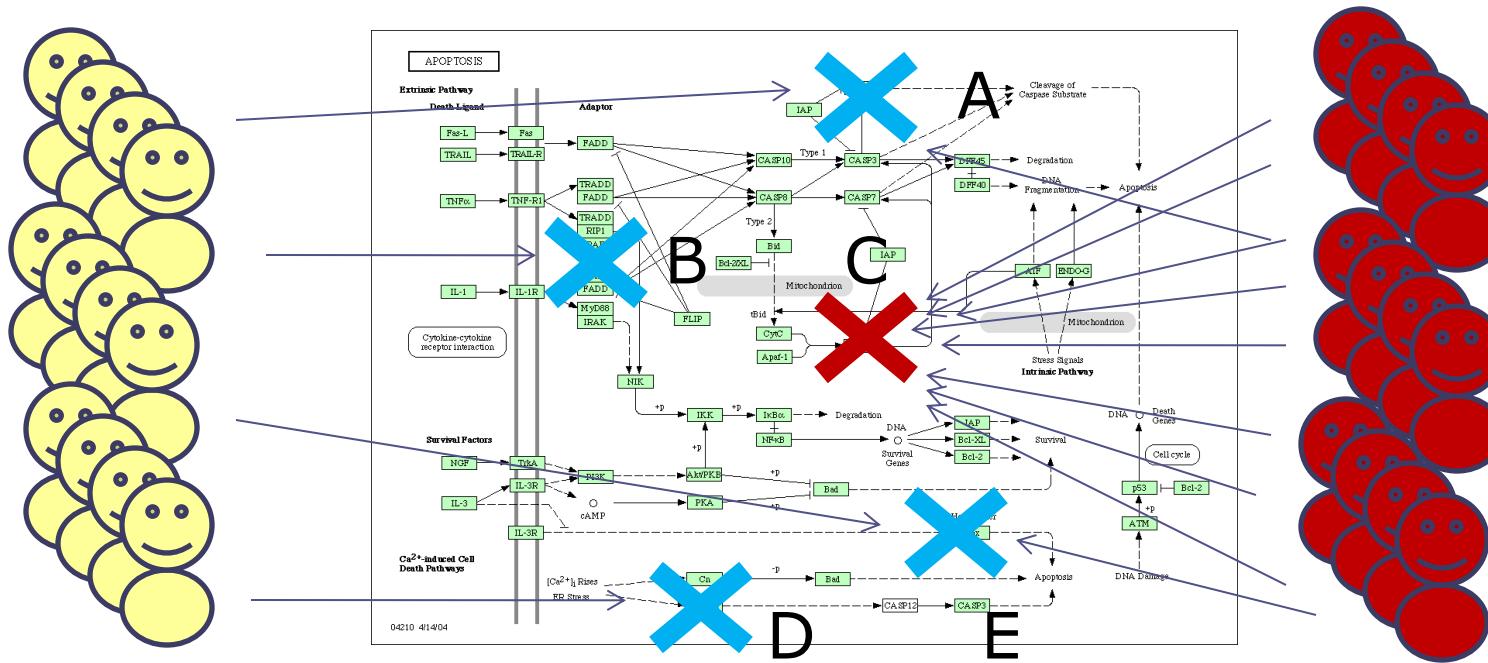
Genes are arranged to their correlation with the prognostic groups

← Pronostic classifier with optimal accuracy

van't Veer et al.,  
Nature, 2002

# Genotyping: finding mutations associated to diseases

The simplest case: monogenic disease



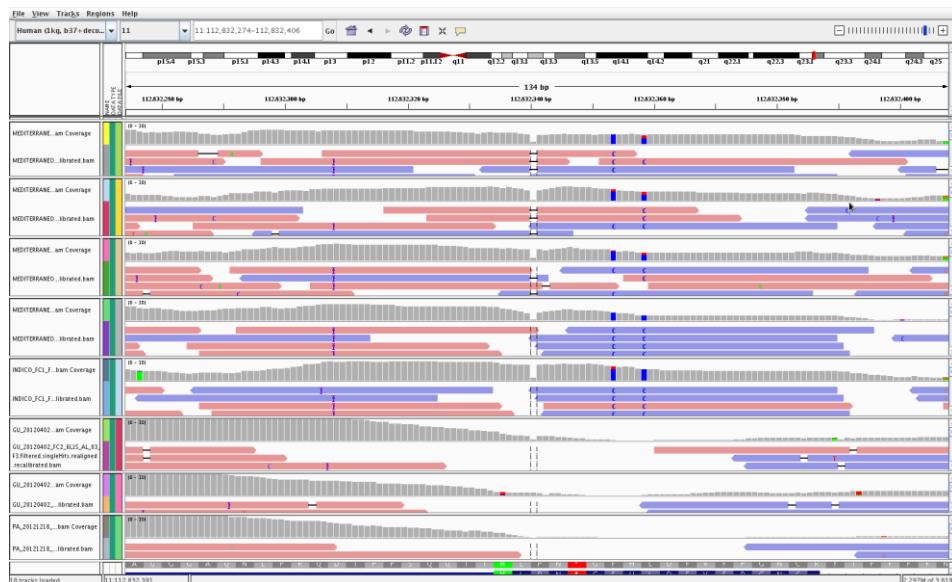
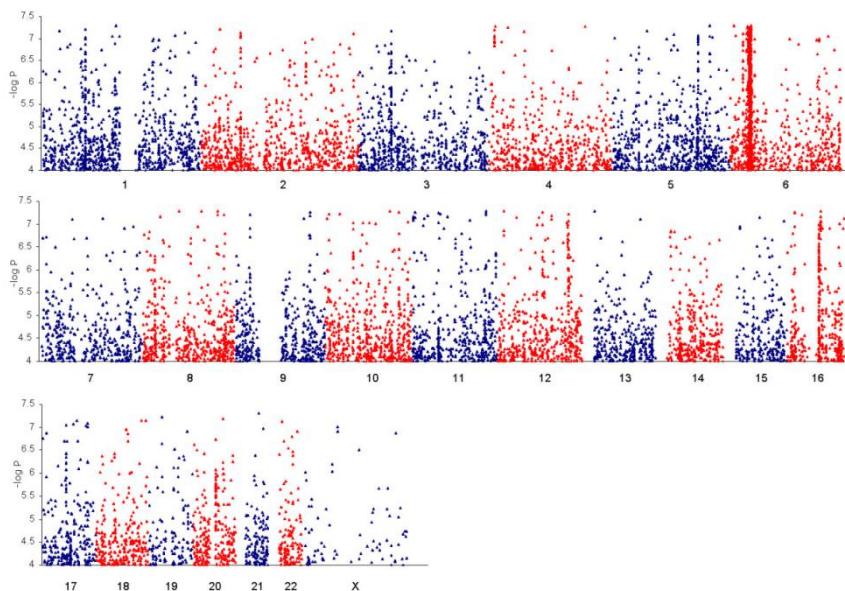
Controls

Cases

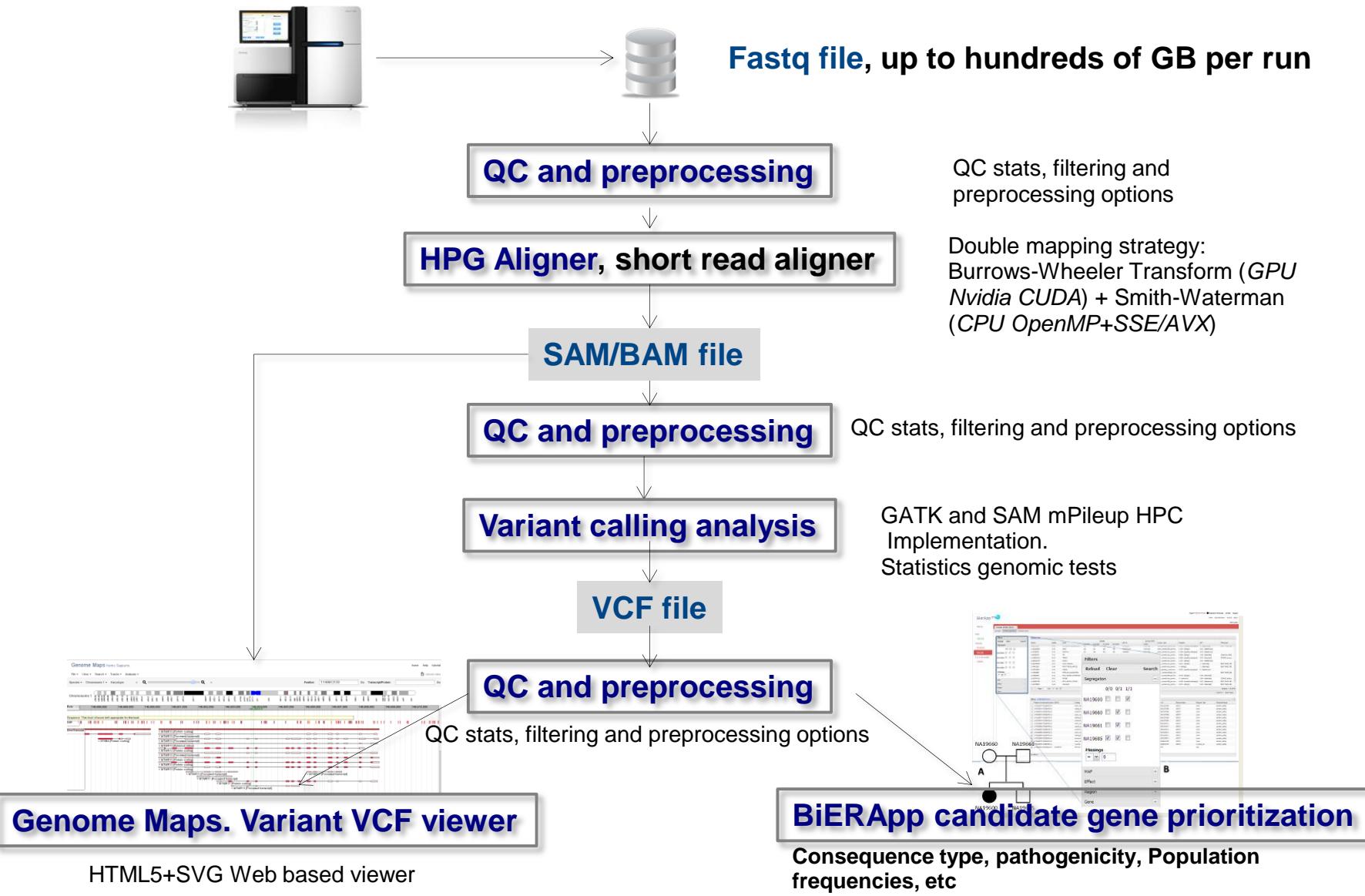
Gene A	1 0 0 0 0 0 0 0 0 0 0 0	0 0 0 0 1 0 0 0 0 0 0 0
Gene B	0 0 0 1 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0 0 0 0 0
Gene C	0 0 0 0 0 0 0 0 0 0 0 0	1 1 1 1 1 1 1 1 1 1 1 1
Gene D	0 0 0 0 0 0 0 0 0 0 1 0	0 0 0 0 0 0 0 0 0 0 0 0
Gene E	0 0 0 0 0 1 0 0 0 0 0 0	0 0 0 0 0 0 1 0 0 0 0 0



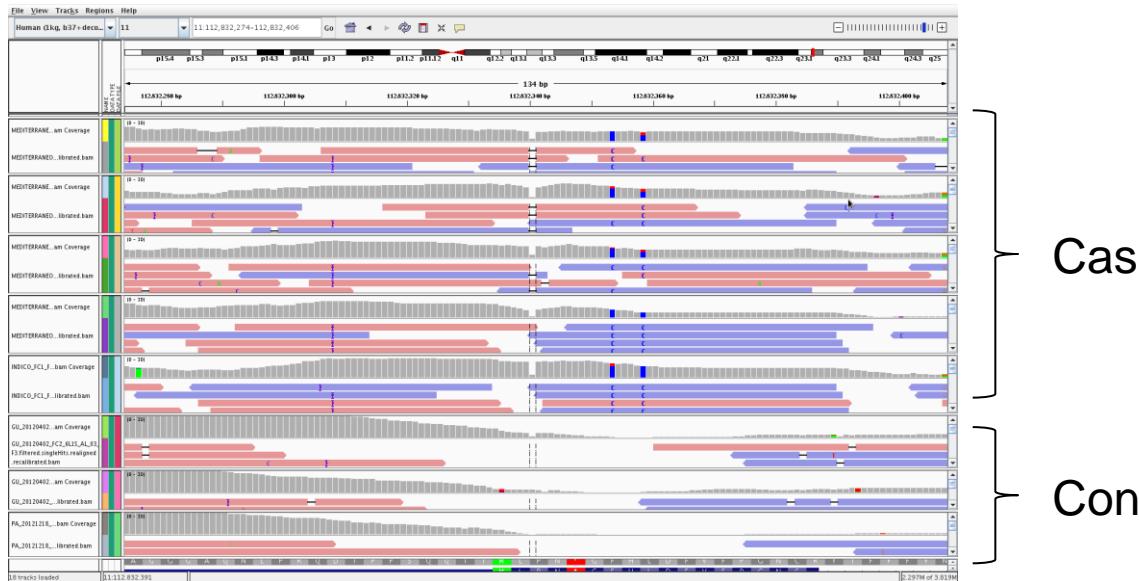
# Genotyping and nowadays NGS



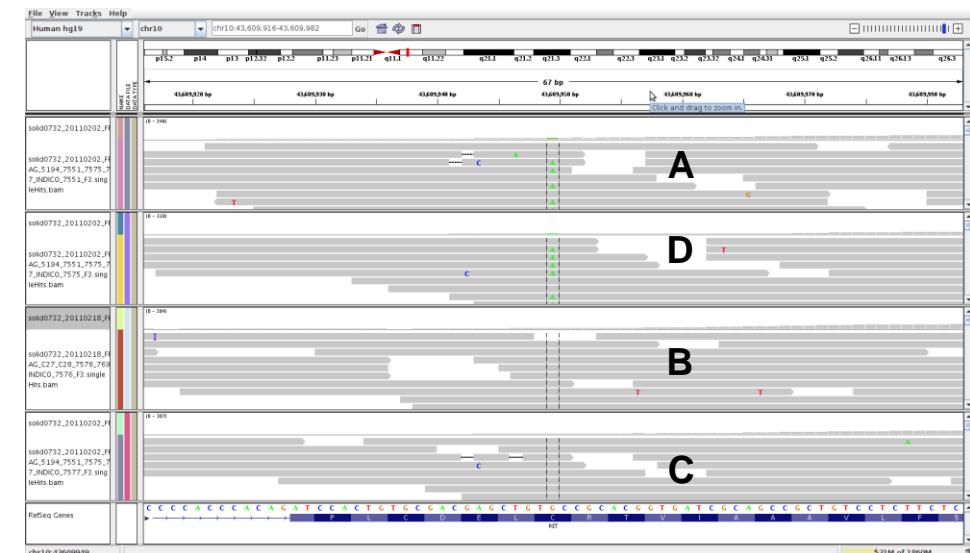
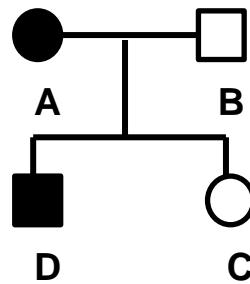
# Primary data analysis tools



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



Segregation  
within a  
pedigree



# **Variant/gene prioritization by successive filtering**



## **Variant level**

Potential impact of  
the variant

Population  
frequencies

## **Experimental design level**

Family(es)  
Trios  
Case / control

## **Functional (system) level**

Gene set  
Network analysis  
Pathway analysis

Control of sequencing errors (missing values)

Testing strategies



# Pipeline of data analysis

## Initial QC

Sequence cleansing  
Base quality  
Remove adapters  
Remove duplicates

FASTQ file

## Mapping + QC

Mapping (HPG)  
Remove multiple mapping reads  
Remove low quality mapping reads  
Realigning  
Base quality recalibrating

BAM file

## Variant calling + QC

Calling and labeling of missing values  
Calling SNVs and indels (GATK) using 6 statistics based on QC, strand bias, consistency (poor QC callings are converted to missing values as well)  
Create multiple VCF with missing, SNVs and indels

VCF file

## Variant and gene prioritization + QC

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 (**Spanish controls**, dbSNP, 1000g, 5500g) and annotation  
Multi-family intersection of genes and variants  
Network-based prioritization  
**Report**

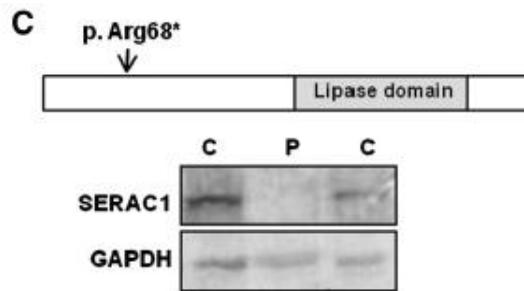
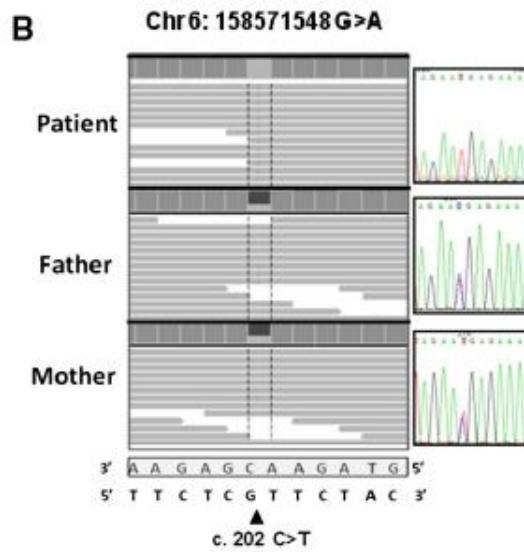
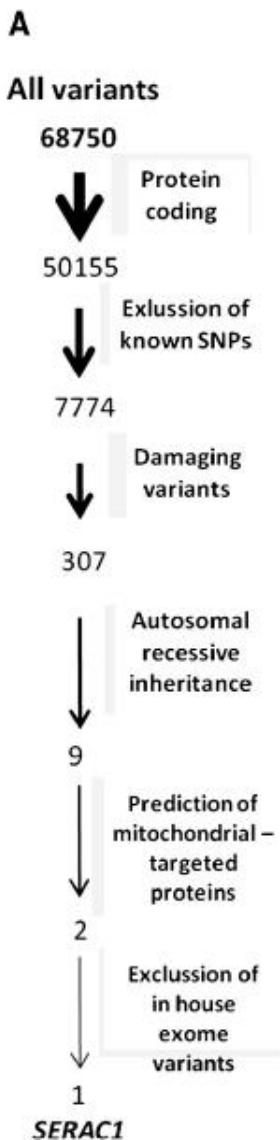
Primary analysis

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-MGAuria) 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.



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>, María 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 Mitochondriales- Enfermedades Metabólicas Hereditarias, Servicio de Pediatría, Hospital 12 de Octubre, Madrid, Spain

<sup>d</sup> BIER, CIBERER, Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain

<sup>e</sup> Computational Medicinal 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> Monique and Jacques Roboh Department of Genetic Research, Hadassah, Hebrew University Medical Center, Jerusalem, Israel

# 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>, Ethylin 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 exon. 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

## 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.

IN THIS ISSUE | May 2011 | Volume 19 | Issue 5

PLOS GENETICS

### 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 Zankl<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>1</sup>, Emma L. Duncan<sup>1,3</sup>, Matthew A. Brown<sup>1\*</sup>

<sup>1</sup> University of Queensland Diamantina Institute, Princess Alexandra Hospital, Woolloongabba, Australia, <sup>2</sup> Centre for Clinical Research, The University of Queensland, Herston, Australia, <sup>3</sup> 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 108-4813/11  
www.nature.com/ejhg



. The two  
re form of  
sequencing  
as a core  
lMPR RNA  
activity of  
which *POP1*

lations As the  
which permits

establishment  
the Rebecca  
by a National  
The funders

#### 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 Garshabi<sup>2</sup>, Seyedeh Sedigheh Abedini<sup>1</sup>, Shirin Ghadami<sup>1</sup>, Roxane Kariminejad<sup>3</sup>, Reinhard Ullmann<sup>4</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 Tschach<sup>\*2,5</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 kyphosis (Kahrizi syndrome)

OMIM 612713

array-based

(c.203 kb)

interval.

essential

families

and eye

potential

European

Keywords:

consanguinity

MV 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>3</sup> Francisco J. López-Domínguez,<sup>3</sup> Carmen Vázquez-Marouschek,<sup>4</sup> Joaquín Dopazo,<sup>3,5,6</sup> Salud Borrego,<sup>1,2</sup> Guillermo António,<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 Andalucía (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

# Where did the heritability go?

The missing heritability problem: individual genes cannot explain the heritability of traits

NEWS FEATURE PERSONAL GENOMES NATURE/Vol 456/November 2008



The case of the missing heritability

When scientists opened up the human genome, they expected to find the genetic components of common traits and diseases. But they were nowhere to be seen. Brendan Maher shines a light on six places where the missing loot could be stashed away.

Vol 461/8 October 2009 doi:10.1038/nature08494 nature REVIEWS

## Finding the missing heritability of complex diseases

Teri A. Manolio<sup>1</sup>, Francis S. Collins<sup>2</sup>, Nancy J. Cox<sup>3</sup>, David B. Goldstein<sup>4</sup>, Lucia A. Hindorff<sup>5</sup>, David J. Hunter<sup>6</sup>, Mark I. McCarthy<sup>7</sup>, Erin M. Ramos<sup>5</sup>, Lon R. Cardon<sup>8</sup>, Aravinda Chakravarti<sup>9</sup>, Judy H. Cho<sup>10</sup>, Alan E. Guttmacher<sup>1</sup>, Augustine Kong<sup>11</sup>, Leonid Kruglyak<sup>12</sup>, Elaine Mardis<sup>13</sup>, Charles N. Rotimi<sup>14</sup>, Montgomery Slatkin<sup>15</sup>, David Valle<sup>9</sup>, Alice S. Whittemore<sup>16</sup>, Michael Boehnke<sup>17</sup>, Andrew G. Clark<sup>18</sup>, Evan E. Eichler<sup>19</sup>, Greg Gibson<sup>20</sup>, Jonathan L. Haines<sup>21</sup>, Trudy F. C. Mackay<sup>22</sup>, Steven A. McCarroll<sup>23</sup> & Peter M. Visscher<sup>24</sup>

Genome-wide association studies have identified hundreds of genetic variants associated with complex human diseases and traits, and have provided valuable insights into their genetic architecture. Most variants identified so far confer relatively small increments in risk, and explain only a small proportion of familial clustering, leading many to question how the remaining, 'missing' heritability can be explained. Here we examine potential sources of missing heritability and propose research strategies, including and extending beyond current genome-wide association approaches, to illuminate the genetics of complex diseases and enhance its potential to enable effective disease prevention or treatment.

How to explain all this?  
Rare Variants, rare CNVs, epigenetics or.. epistatic effects?

Table 1 | Estimates of heritability and number of loci for several complex traits

Disease	Number of loci	Proportion of heritability explained
Age-related macular degeneration <sup>72</sup>	5	50%
Crohn's disease <sup>21</sup>	32	20%
Systemic lupus erythematosus <sup>73</sup>	6	15%
Type 2 diabetes <sup>74</sup>	18	6%
HDL cholesterol <sup>75</sup>	7	5.2%
Height <sup>15</sup>	40	5%
Early onset myocardial infarction <sup>76</sup>	9	2.8%
Fasting glucose <sup>77</sup>	4	1.5%

\* Residual is after adjustment for age, gender, diabetes.

# Is the heritability missing or are we looking at the wrong place?

How to explain missing heritability?  
Rare Variants, rare CNVs, epigenetics or.. **epistatic effects?**

NEWS FEATURE PERSONAL GENOMES

NATURE Vol 456 November 2008

The case of the missing heritability

When scientists opened up the human genome, they expected to find the genetic components of common traits and diseases. But they were nowhere to be seen. Brendan Maher shines a light on six places where the missing loot could be stashed away.

Table 1 | Estimates of heritability and number of loci for several complex traits

Disease	Number of loci	Proportion of heritability explained
Age-related macular degeneration <sup>72</sup>	5	50%
Crohn's disease <sup>21</sup>	32	20%
Systemic lupus erythematosus <sup>73</sup>	6	15%
Type 2 diabetes <sup>74</sup>	18	6%
HDL cholesterol <sup>75</sup>	7	5.2%
Height <sup>15</sup>	40	5%
Early onset myocardial infarction <sup>76</sup> *	9	1.8%
Fasting glucose <sup>77</sup>	4	1.5%

\* Residual is after adjustment for age, gender, diabetes.

genetics

Common SNPs explain a large proportion of the heritability for human height

Jian Yang<sup>1</sup>, Beben Benyamin<sup>1</sup>, Brian P McEvoy<sup>1</sup>, Scott Gordon<sup>1</sup>, Anjali K Henders<sup>1</sup>, Dale R Nyholt<sup>1</sup>, Pamela A Madden<sup>2</sup>, Andrew C Heath<sup>2</sup>, Nicholas G Martin<sup>1</sup>, Grant W Montgomery<sup>1</sup>, Michael E Goddard<sup>3</sup> & Peter M Visscher<sup>1</sup>

SNPs discovered by genome-wide association studies (GWASs) account for only a small fraction of the genetic variation of complex traits in human populations. Where is the remaining heritability? We estimated the proportion of variance for human height explained by 294,831 SNPs genotyped in 3,925 unrelated individuals using a linear model analysis, and validated the estimation method with simulations based on the observed genotype data. We show that 45% of variance can be explained by considering all SNPs simultaneously. Thus, most of the heritability is not missing but has not previously been detected because the individual effects are too small to pass stringent significance tests. We provide evidence that the remaining heritability is due to incomplete linkage disequilibrium between causal variants and genotyped SNPs, exacerbated by causal variants having lower minor allele frequency than the SNPs explored to date.

2010 Nature America, Inc. All rights reserved.

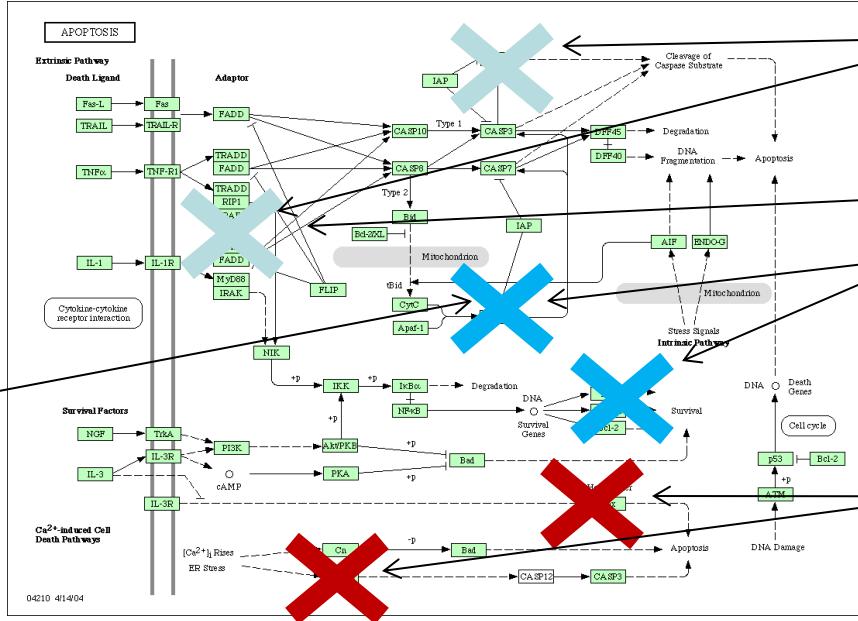
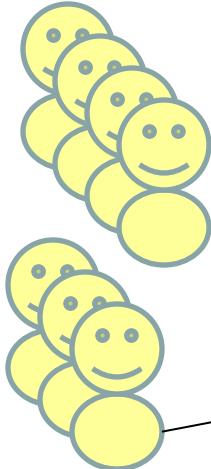
At the end, most of the heritability was there... of variation that their effects do not reach stringent significance thresholds and/or the causal variants are not in complete linkage disequilibrium (LD) with the SNPs that have been genotyped. Lack of complete LD might, for example, occur if causal variants have a minor allele frequency (MAF) than genotyped SNPs. Here we test these two hypotheses by estimating the contribution of each to the heritability of human height and comparing them.

Height in humans is a classical quantitative trait, easy to measure and studied for well over a century as a model for investigating the genetic basis of complex traits. The heritability of height has been estimated to be ~0.8 (refs. 9,11–13). Rare mutations that cause extreme short or tall stature have been found<sup>14,15</sup>, but these do not explain much of the variation in the general population. Recent GWASs on tens of thousands of individuals have detected ~50 variants that are associated with height in the population, but these in total account for only ~5% of phenotypic variance<sup>16–19</sup>.

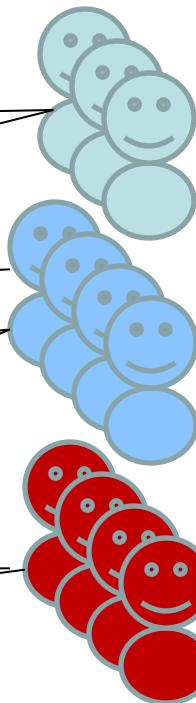
Data from a GWAS that are collected to detect statistical associations between SNPs and complex traits are usually analyzed by testing each

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

## Controls



## Cases



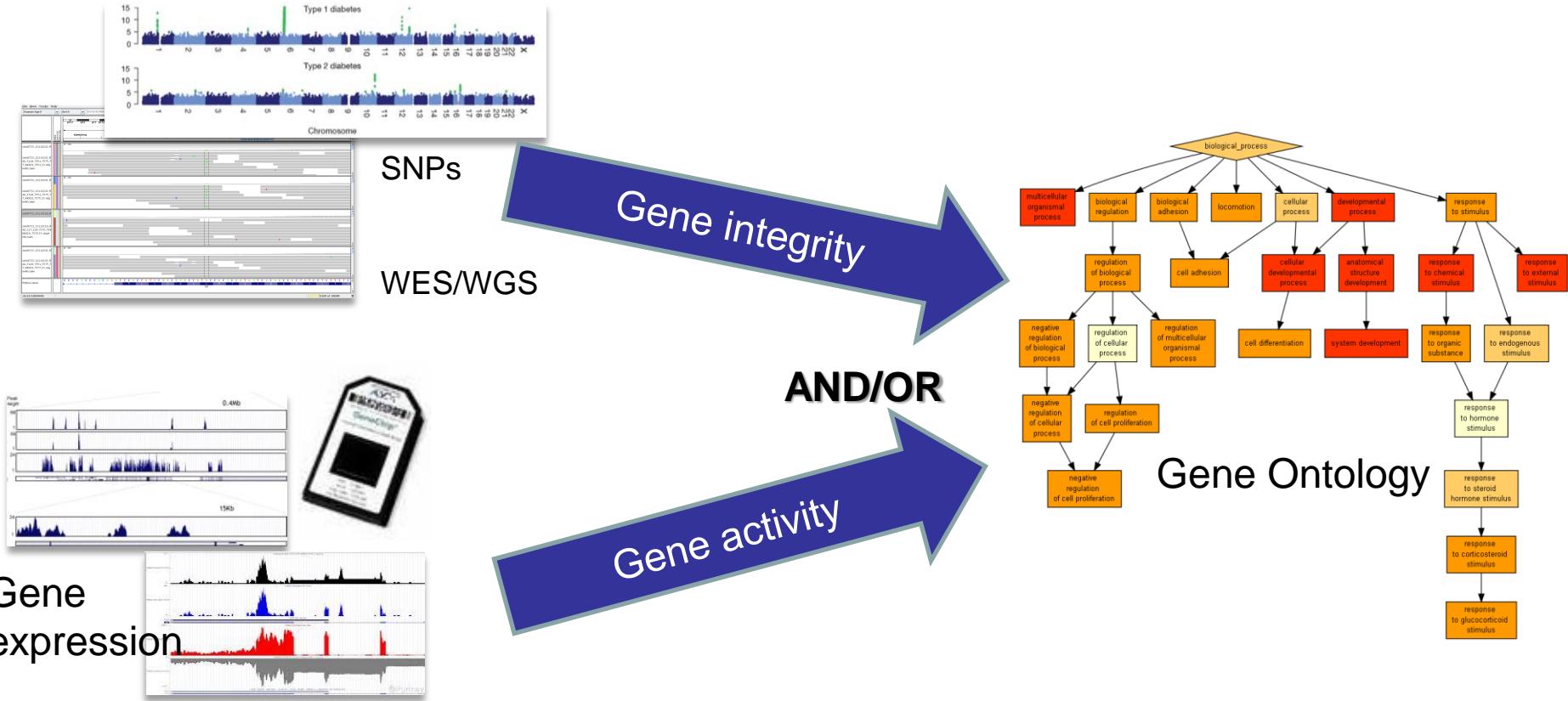
Affected **cases** in complex diseases will be a **heterogeneous** population with different mutations (or combinations).

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



**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.

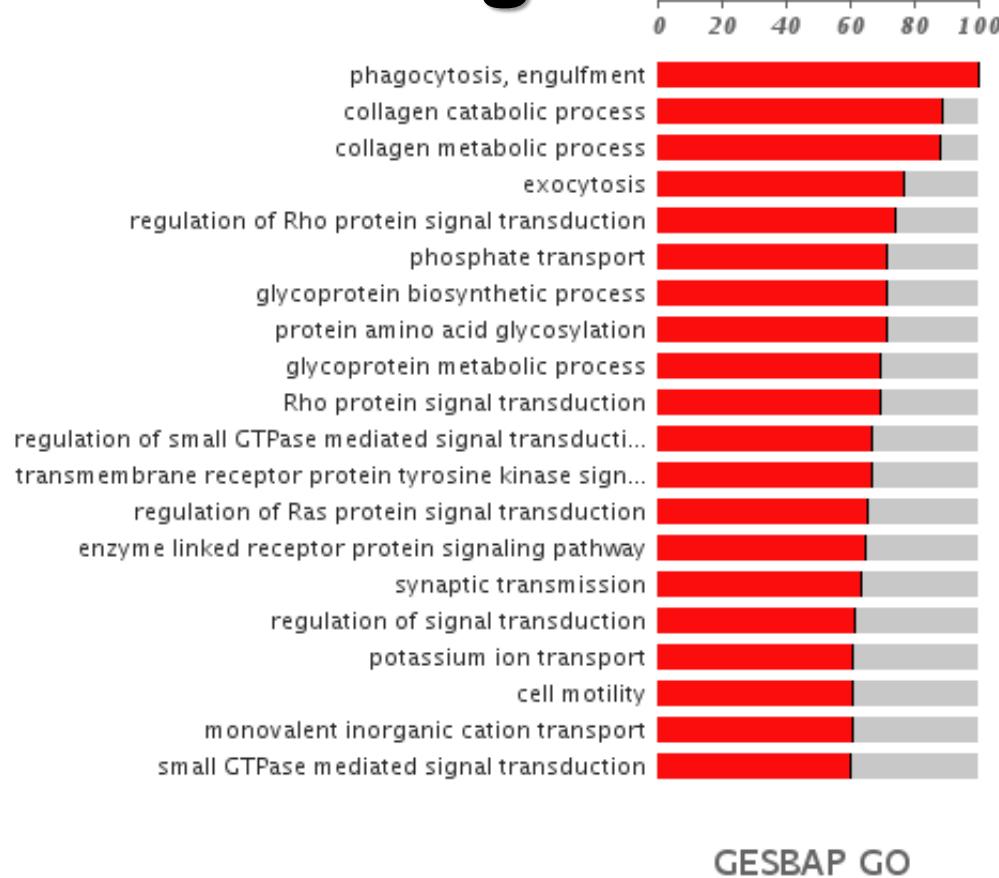
The CGEMS initiative. (Hunter et al. Nat Genet 2007)

1145 cases 1142 controls. Affy 500K

Conventional association test reports only 4 SNPs  
significantly mapping 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

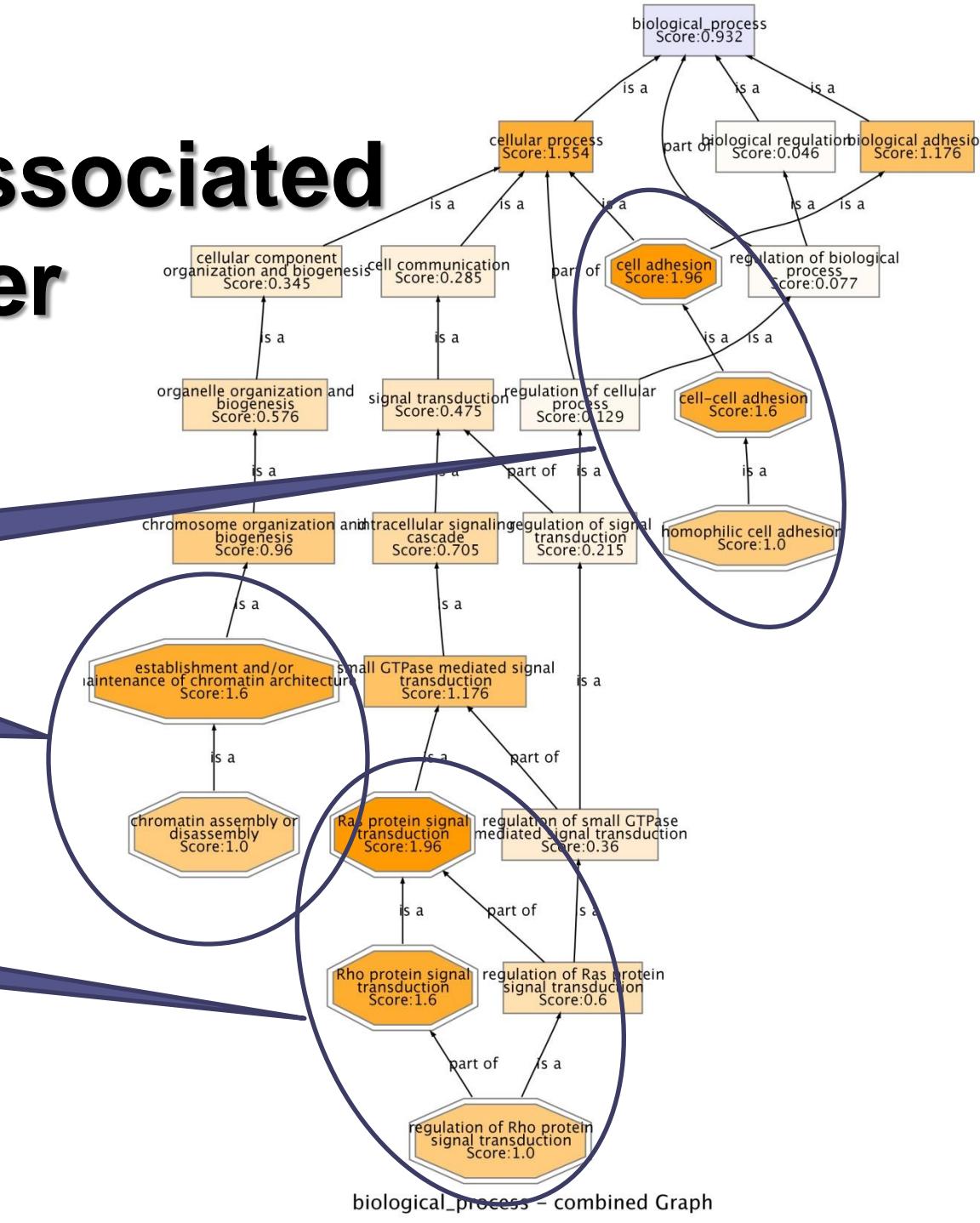
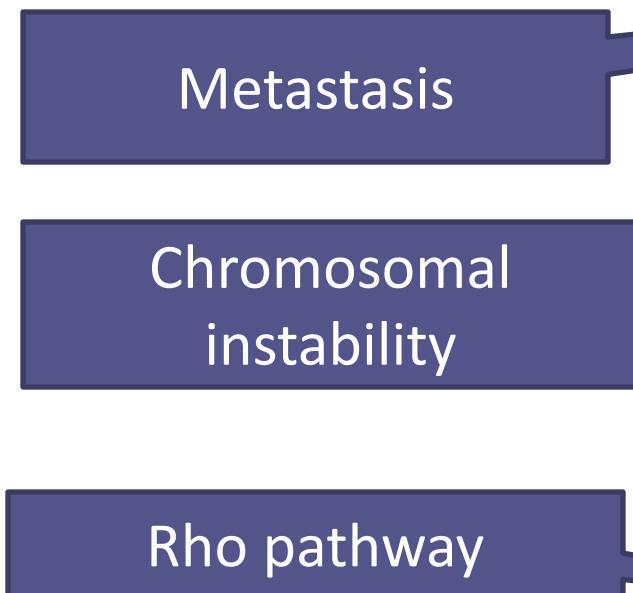
CGEMS initiative.  
(Hunter et al. Nat Genet 2007)

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=4.45x10<sup>-03</sup>) in which FGFR2 is included.

# GO processes significantly associated to breast cancer



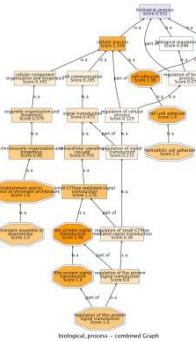
# From gene-based to function-based perspective

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

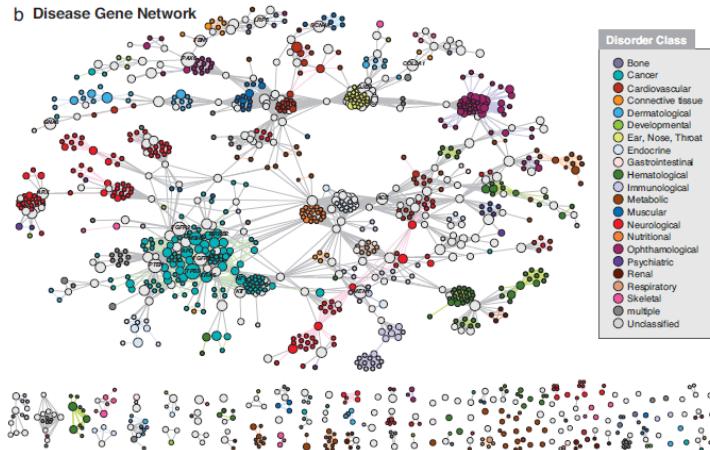


	SNPs, gene exp.	GO
<b>Detection power</b>	Low (only very prevalent genes)	high
<b>Annotations available</b>	many	many
<b>Use</b>	Biomarker	Illustrative, give hints

# Can the interactome help to find disease mutations?

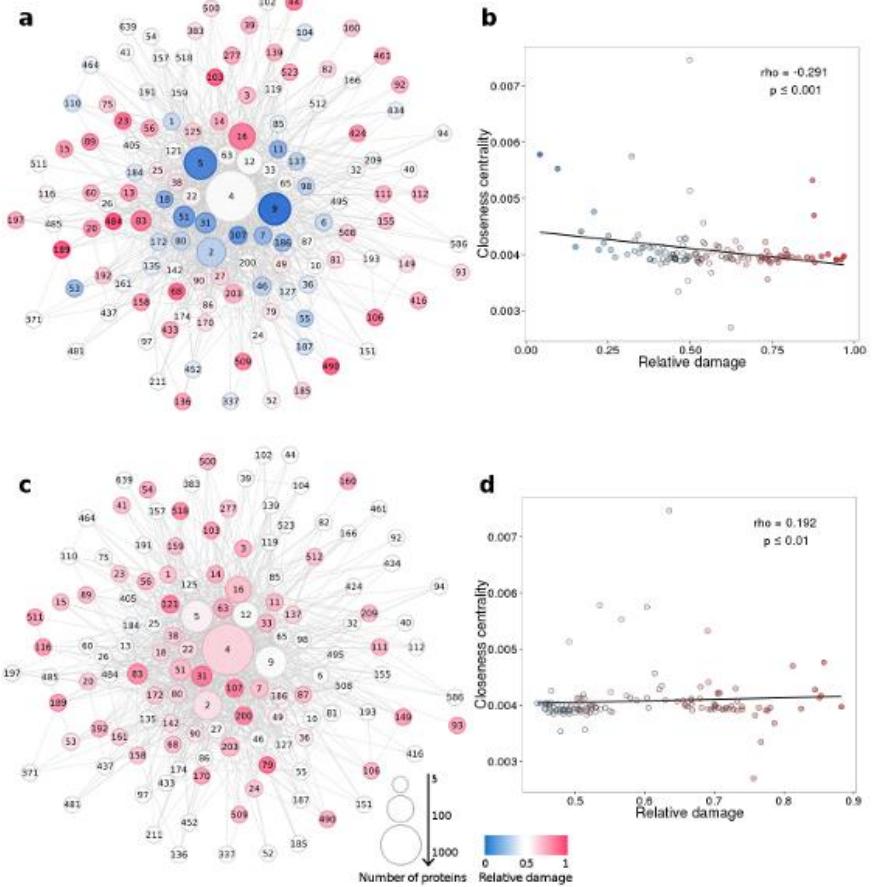
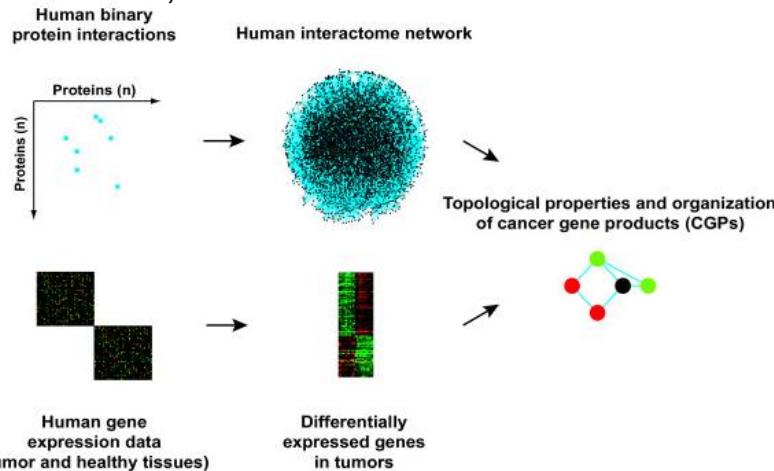
Disease genes are close in the interactome

Goh 2007 PNAS



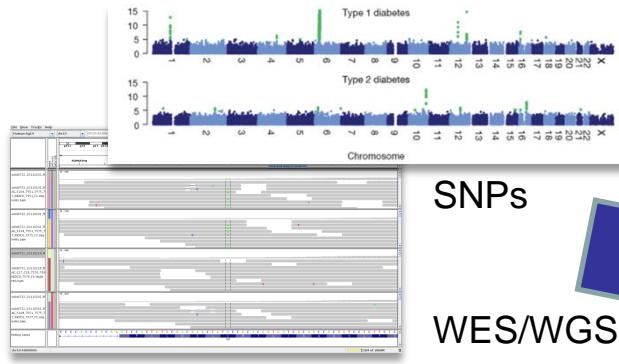
Cancer genes are central.

Hernandez, 2007 BMC Genomics



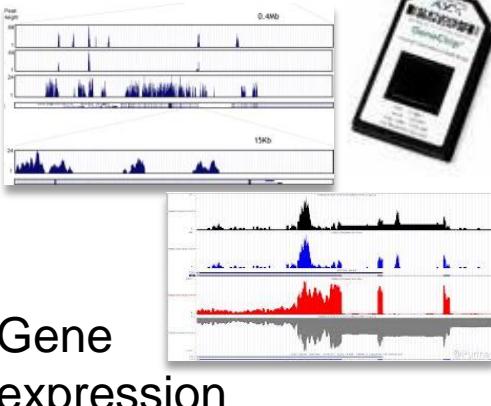
Deleterious mutations in 1000g (up) and somatic CLL deleterious mutations (down)  
Garcia-Alonso 2014 Mol Syst Biol

# From gene-based to function-based perspective



SNPs

WES/WGS



Gene expression

*Gene integrity*

AND/OR

*Gene activity*

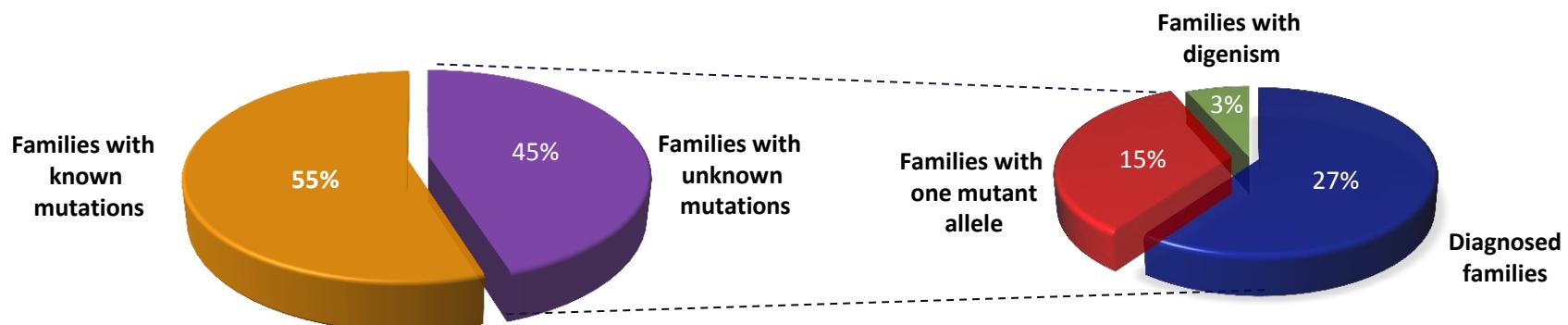
Using protein interaction networks as an scaffold to interpret the genomic data in a functionally-derived context



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

# Example with Inherited Retinal Dystrophies

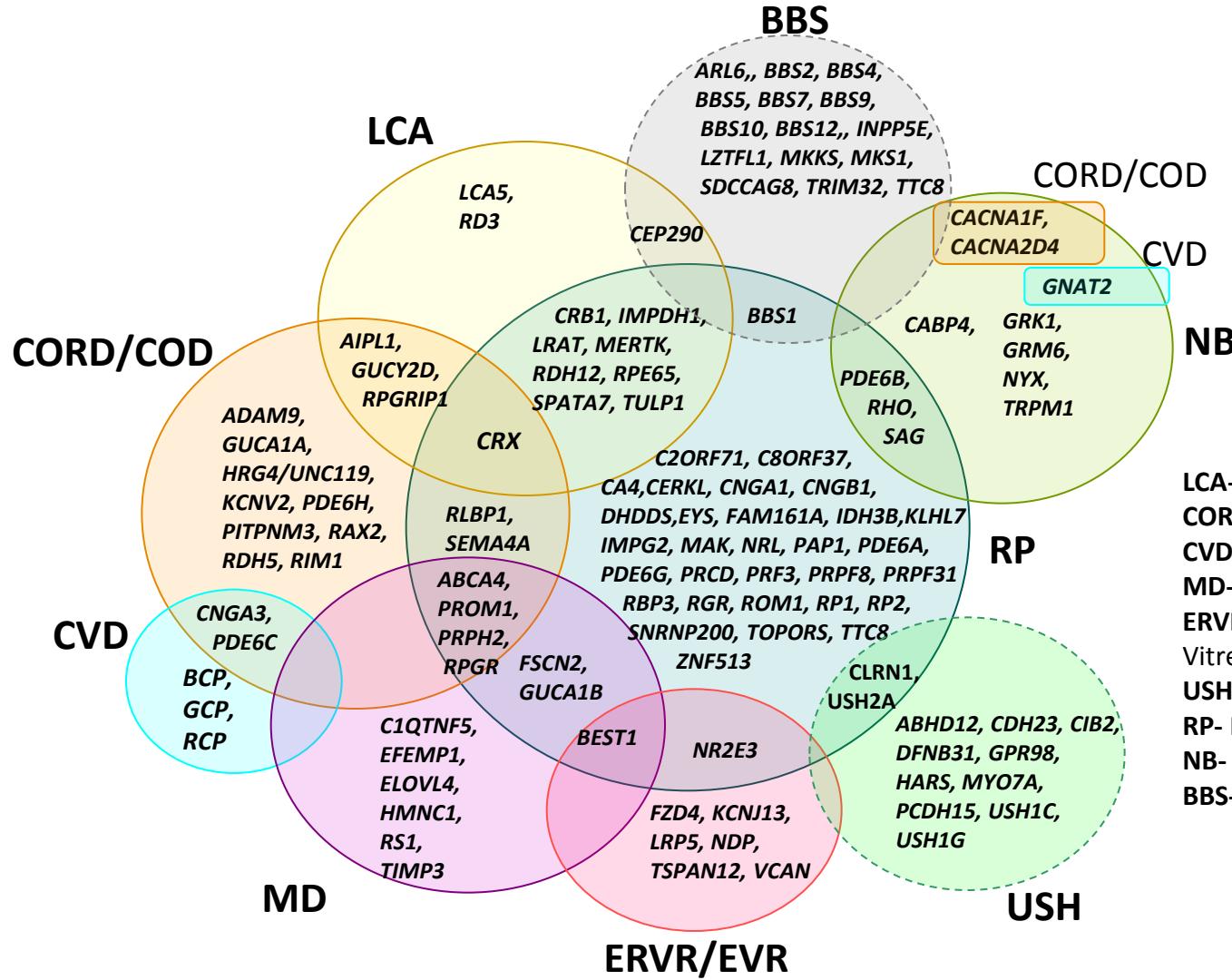
- Prevalence 1 in 3000
- Clinically and genetically very heterogeneous
- 190 GENES account for approx. 50% of IRDs.



Novel variants and genes remain to be discovered

# Network analysis.

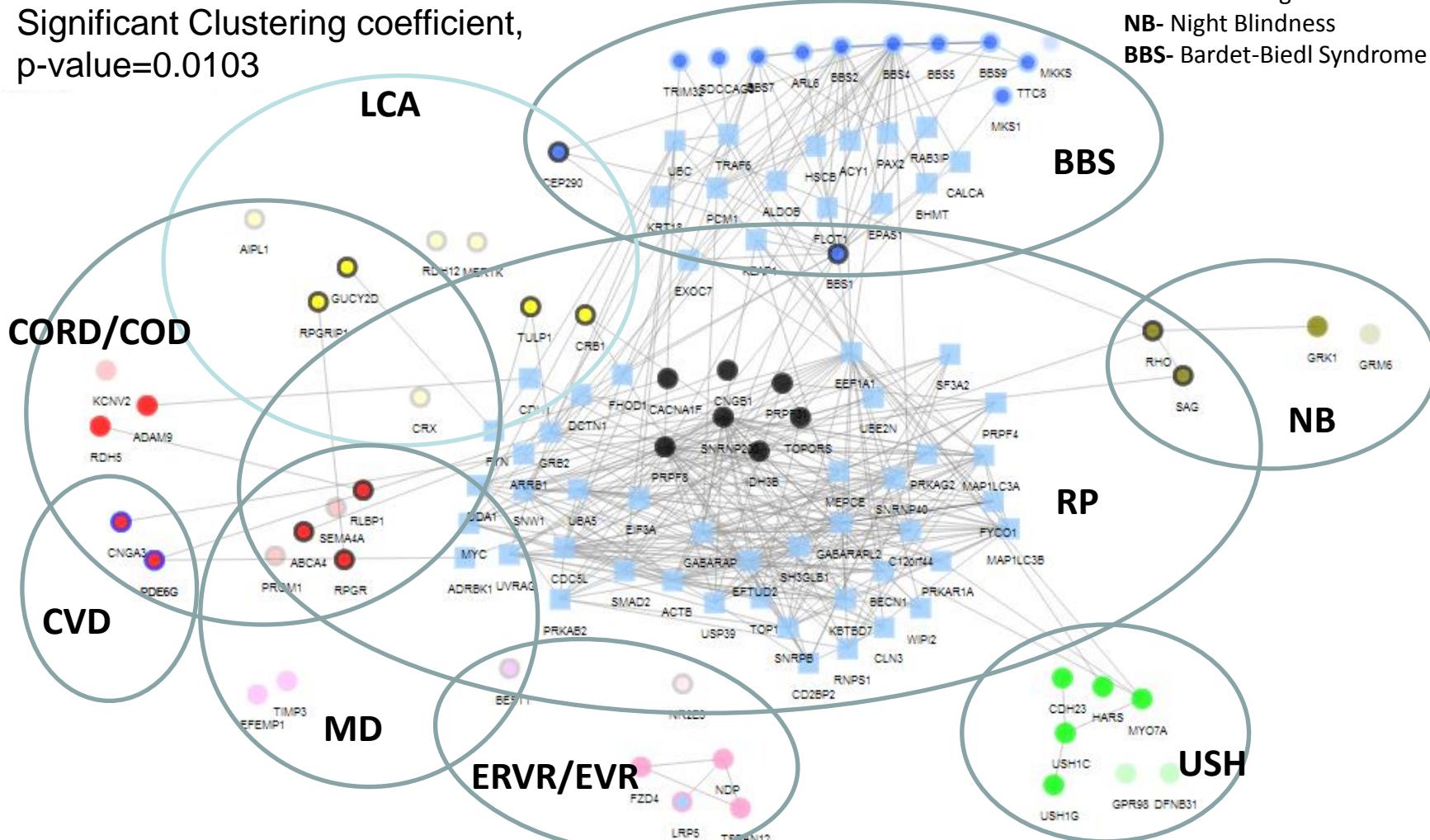
Is genetic overlapping among IRDs related to protein interaction?



**LCA**-Leber Congenital Amaurosis  
**CORD/COD**- Cone and cone-rod dystro.  
**CVD**- Colour Vision Defects  
**MD**- Macular Degeneration  
**ERVR/EVR**- Erosive and Exudative Vitreoretinopathies  
**USH**- Usher Syndrome  
**RP**- Retinitis Pigmentosa  
**NB**- Night Blindness  
**BBS**- Bardet-Biedl Syndrome

# Connectivity among IRDs explain genetic overlap

Significant Clustering coefficient,  
p-value=0.0103



LCA-Leber Congenital Amaurosis  
 CORD/COD- Cone and cone-rod dystro.  
 CVD- Colour Vision Defects  
 MD- Macular Degeneration  
 ERVR/EVR- Erosive and Exudative  
 Vitreoretinopathies  
 USH- Usher Syndrome  
 RP- Retinitis Pigmentosa  
 NB- Night Blindness  
 BBS- Bardet-Biedl Syndrome

# 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 [idopazo@cipf.es](mailto:idopazo@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 Minguez<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, 37–6  
doi:10.1093/nar/gkp402

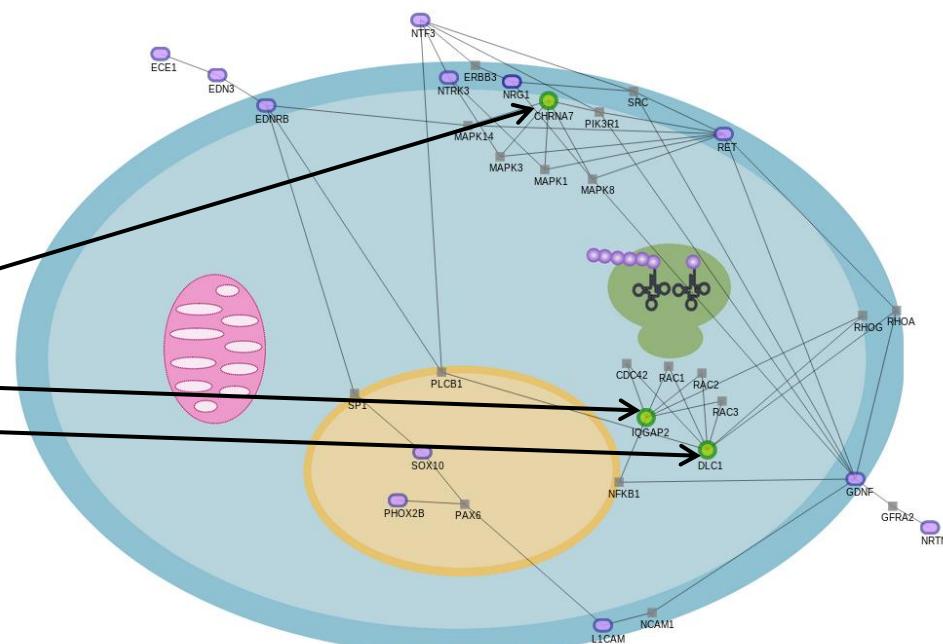
## SNOW, a web-based tool for the statistical analysis of protein-protein interaction networks

Pablo Minguez<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



# Software that detects the most significant sub-network within a list of genes

Nucleic Acids Research Advance Access published May 19, 2009  
Nucleic Acids Research, 2009, 37–6  
doi:10.1093/nar/gkp402

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

Pablo Minguez<sup>1</sup>, Stefan Götz<sup>1,2</sup>, David Montane<sup>1</sup>, Fatima Al-Shahrour<sup>1</sup> and Joaquín 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

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

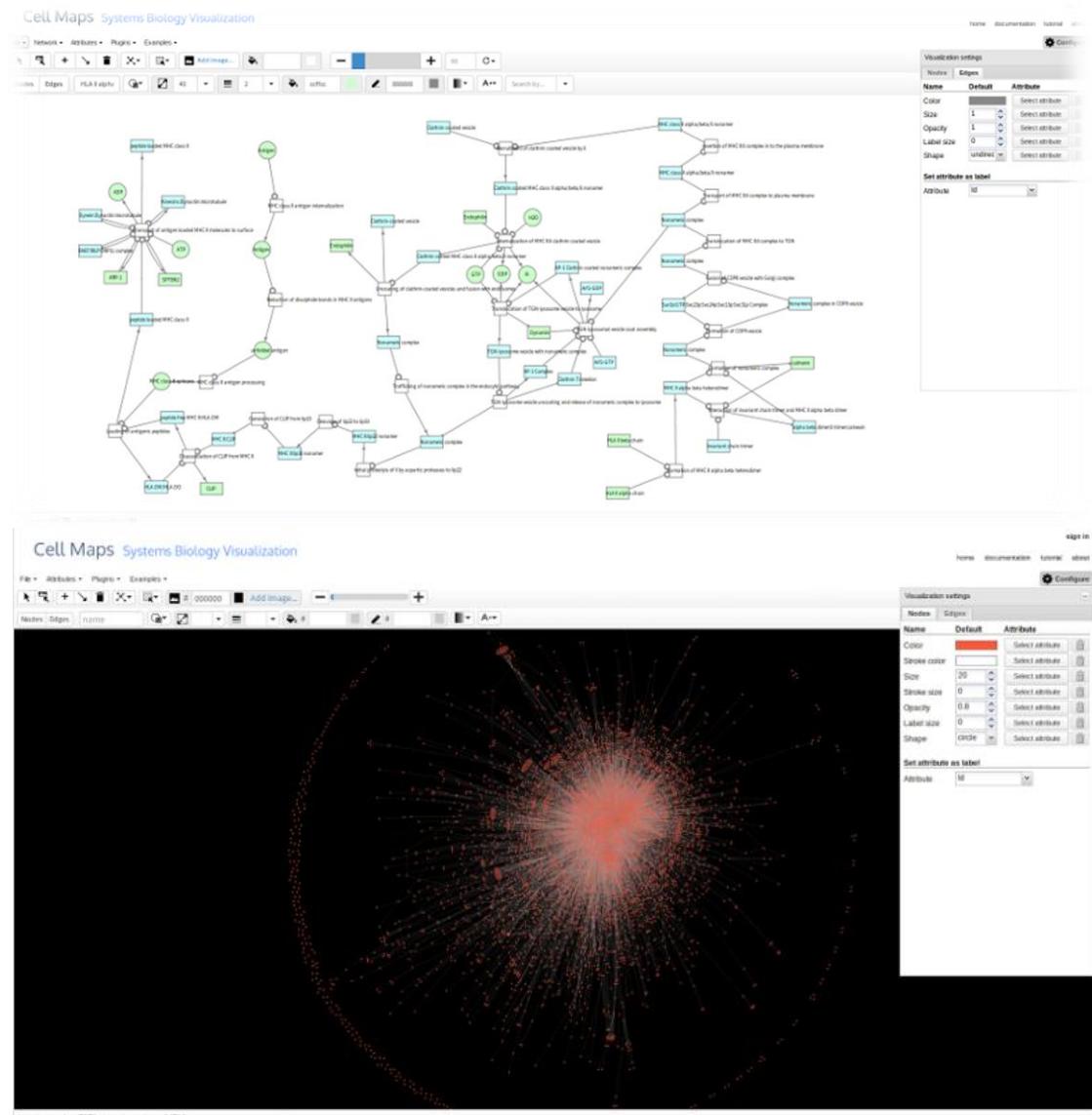
Luz García-Alonso<sup>1</sup>, Roberto Alonso<sup>1</sup>, Enrique Vidal<sup>1</sup>, Alicia Amadoz<sup>1</sup>, Alejandro de María<sup>1</sup>, Pablo Minguez<sup>2</sup>, Ignacio Medina<sup>1,3</sup> and Joaquín Dopazo<sup>1,3,4,\*</sup>

<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

Network analysis and the interactive, web-based network viewer CellMaps

<http://cellmaps.babelomics.org/>



# From gene-based to function-based perspective

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

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

# 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 Sebastian-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 Amadorz<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

## Inferred 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

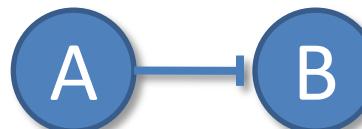
Received March 3, 2013; Revised April 18, 2013; Accepted May 2, 2013

Activation



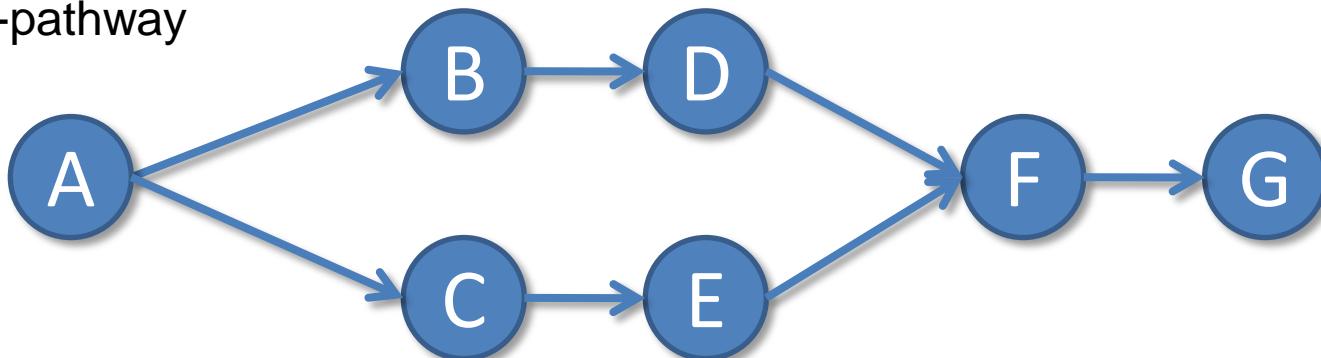
$$\text{Prob.} = P(\text{A activated})P(\text{B activated})$$

Inhibition



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

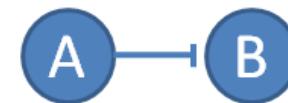
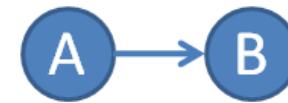
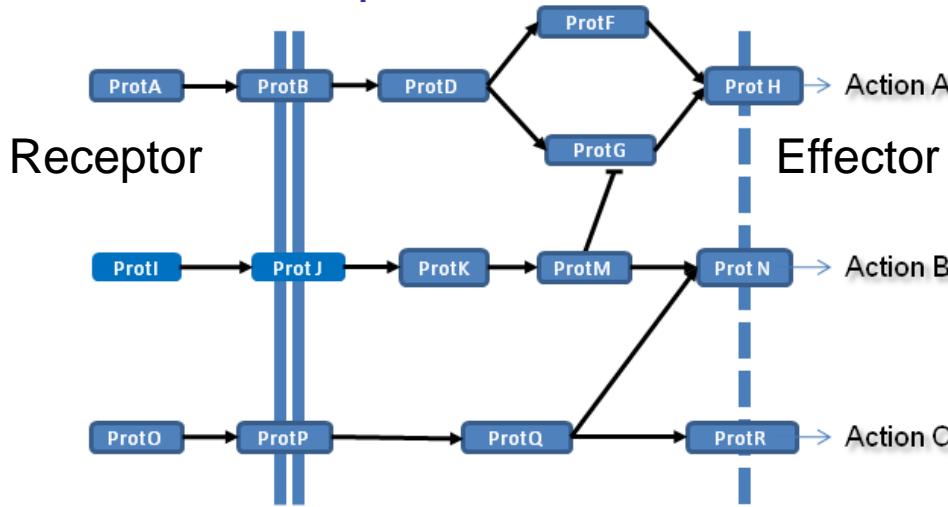
Sub-pathway



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

# 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

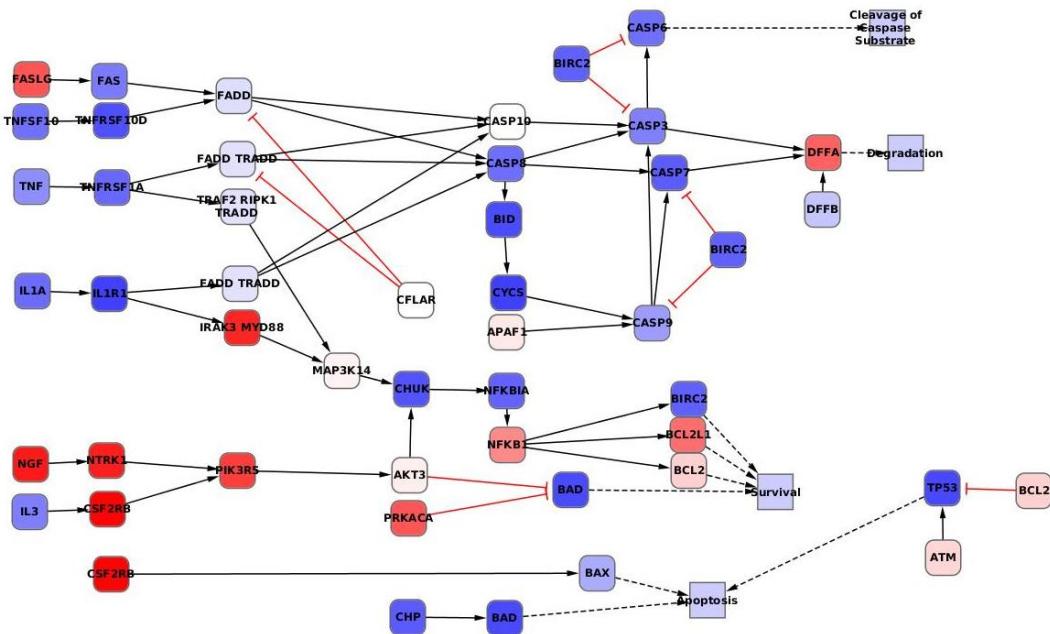
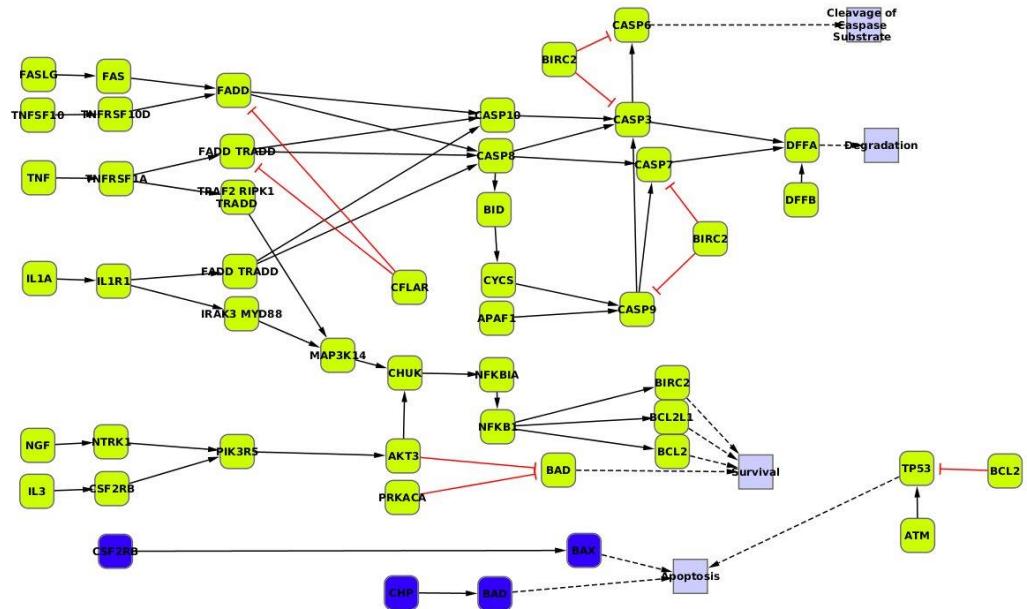


**Activations  
and  
repressions  
occur**

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

# 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



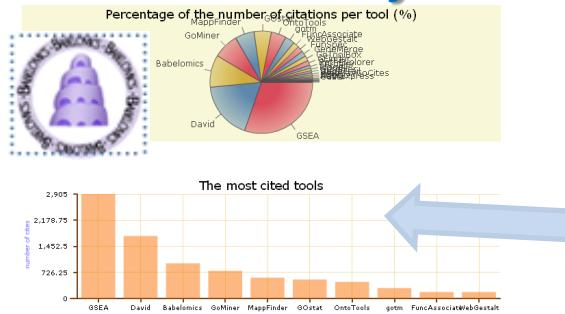
# From gene-based to function-based perspective

	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

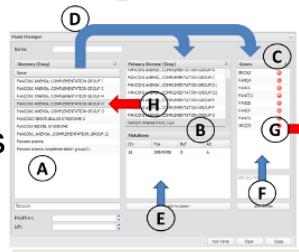
## Functional analysis



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

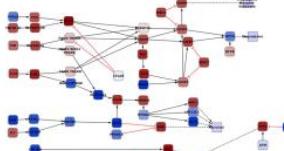


## Diagnostic

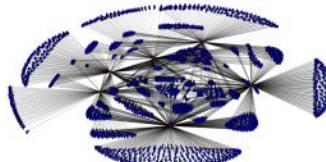


NGS  
panels

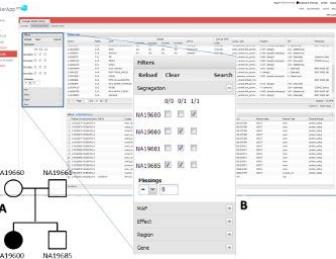
## Signaling network



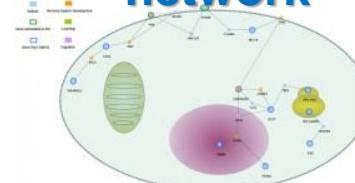
## Regulatory network



## Variant prioritization



## Interaction network



## Variant annotation



## Mapping

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

## Visualization



Genome maps is now part of the ICGC data portal

## CellBase

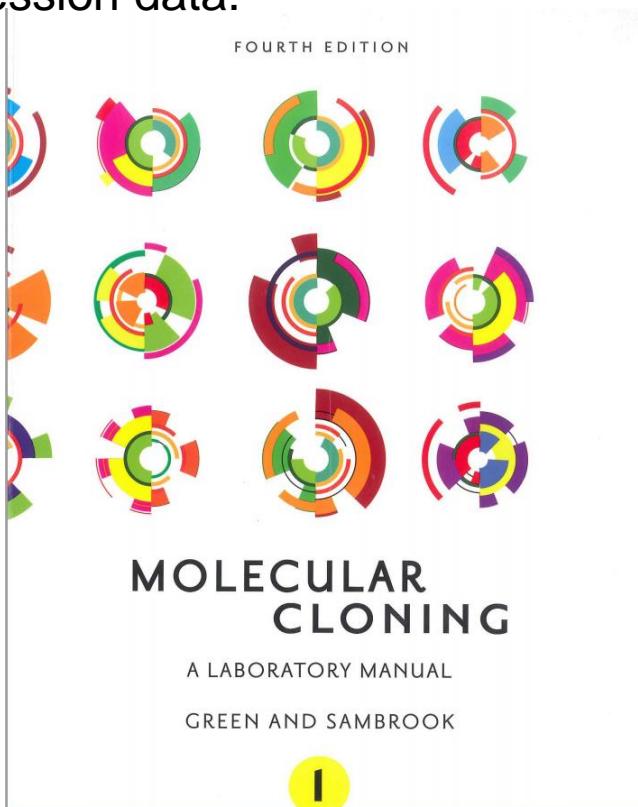


Ultrafast  
genome  
viewer with  
google  
technology

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.



High impact developments

A screenshot of the Babelomics 4 web interface. The top navigation bar includes links for "Upload data", "Processing data", "Expression", "Genomic", "Functional analysis", and "Utilities". The main panel is titled "Upload data" and shows a file selection dialog with "cMyc.zip" selected. A green arrow points from this dialog to a "Select format" dropdown menu on the right. The "Select format" menu is expanded, showing various options like "Microarray", "Expression", "One-channel", etc., with "One-channel" highlighted. A red box highlights the "Data type" dropdown set to "One-channel". A red box also highlights the "Data name" input field containing "cMyc\_raw\_data". At the bottom right of the upload panel is a red-bordered "Upload" button. To the right of the upload panel is a "Data list" panel showing a single item: "cMyc\_raw\_data". At the bottom right of the page is a green "Accept" button.

FIGURE 1. Babelomics data uploading form. Click "Browse" to upload the data file named cMyc.zip. Select "Affy-metrix" 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.

# SOCIAL:

## MDA group in LinkedIn

## Babelomics group in Facebook and twitter

Facebook | Babelomics fans - Mozilla Firefox

LinkedIn: MDA Group - Mozilla Firefox

Babelomics fans

Wall Info Discussions Photos Video Events

Write something...

Joquín Dopazo The new version Babelomics 4.0 has officially released.  
Enjoy it!  
A few seconds ago Comment Like Report

Message all members Promote Group with an Advert Edit group settings Edit members Invite people to join Create group event Leave Group

Information Category: Internet & Technology - Software Description: Babelomics is the integrative platform for the analysis of transcriptomics, ... Transfiriendo datos desde static.fbcdn.net...

Inicio 3 Micro... 3 Firefox 10 Exp... 5 Micro... Skype™ 2 Micro...

@babelomics



PRINCIPE FELIPE  
CENTRO DE INVESTIGACION

@xdopazo  
@bioinfocipf

# 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), and...

... Ignacio Medina  
Head of Computational Biology  
Lab, HPCS  
University of Cambridge, UK



EMBL-EBI Scientific Collaborator



iSCB  
INTERNATIONAL SOCIETY FOR COMPUTATIONAL BIOLOGY  
AFFILIATED REGIONAL INSTITUTE