

Variant prioritization

Joaquín Dopazo

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

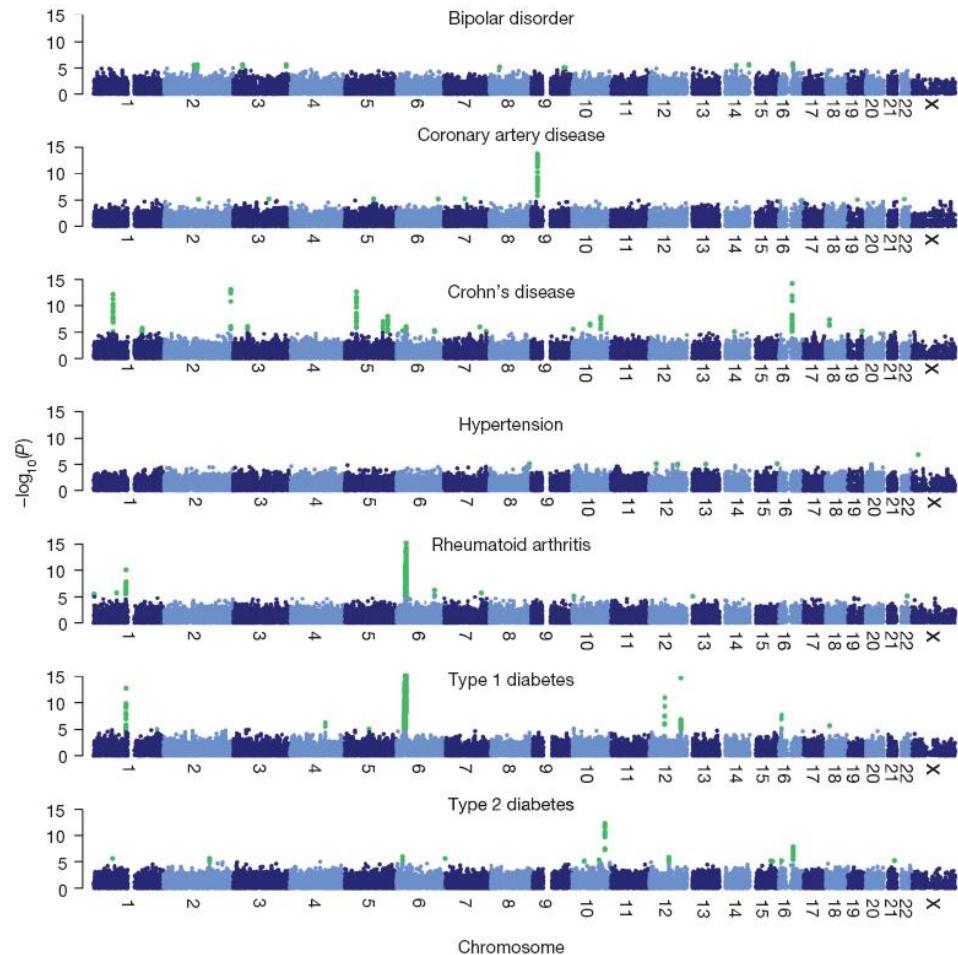
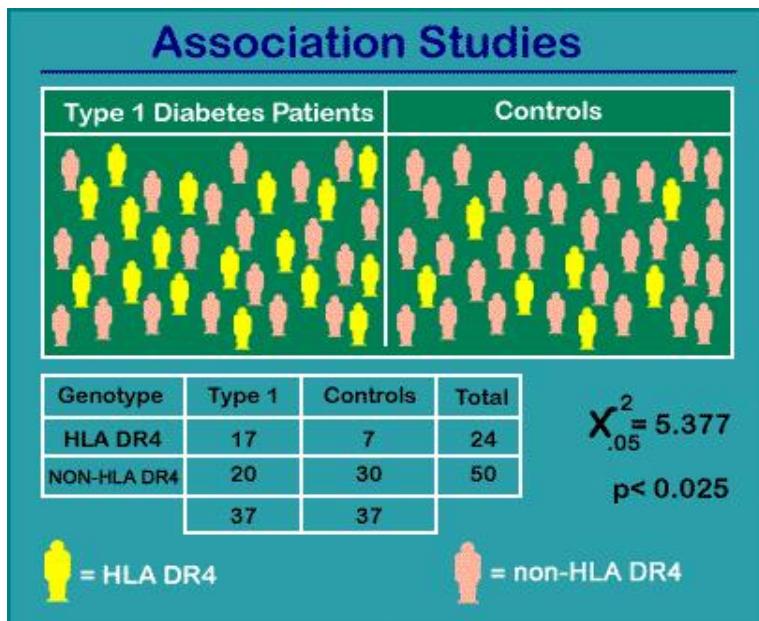
<http://bioinfo.cipf.es>
<http://www.medicalgenomeproject.com>
<http://www.babelomics.org>
<http://www.hpc4g.org>
 @xdopazo

University of Cambridge, 23-25 February 2015

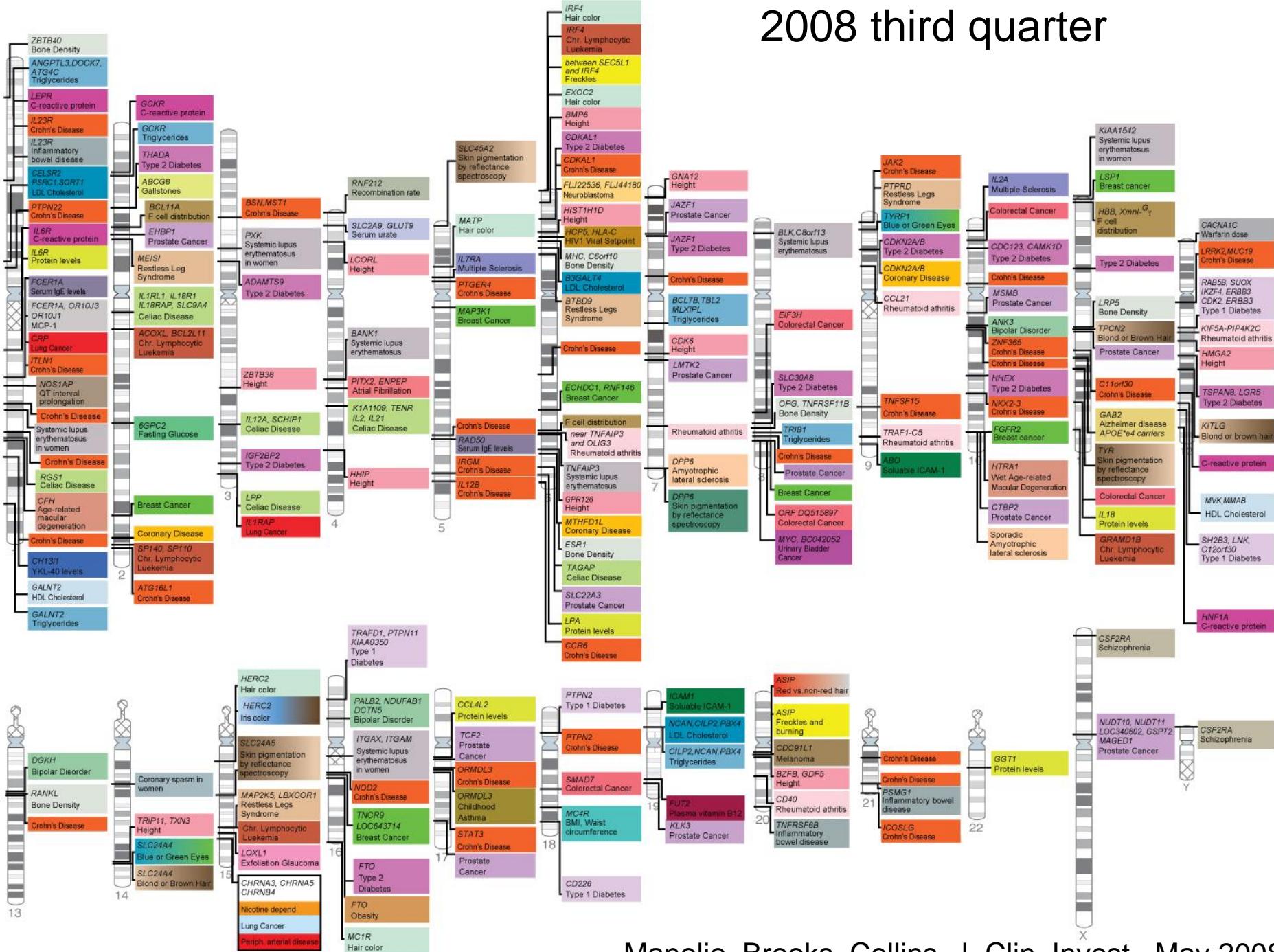


The dawn of genomic data production

Candidate gene studies using GWAS



2008 third quarter

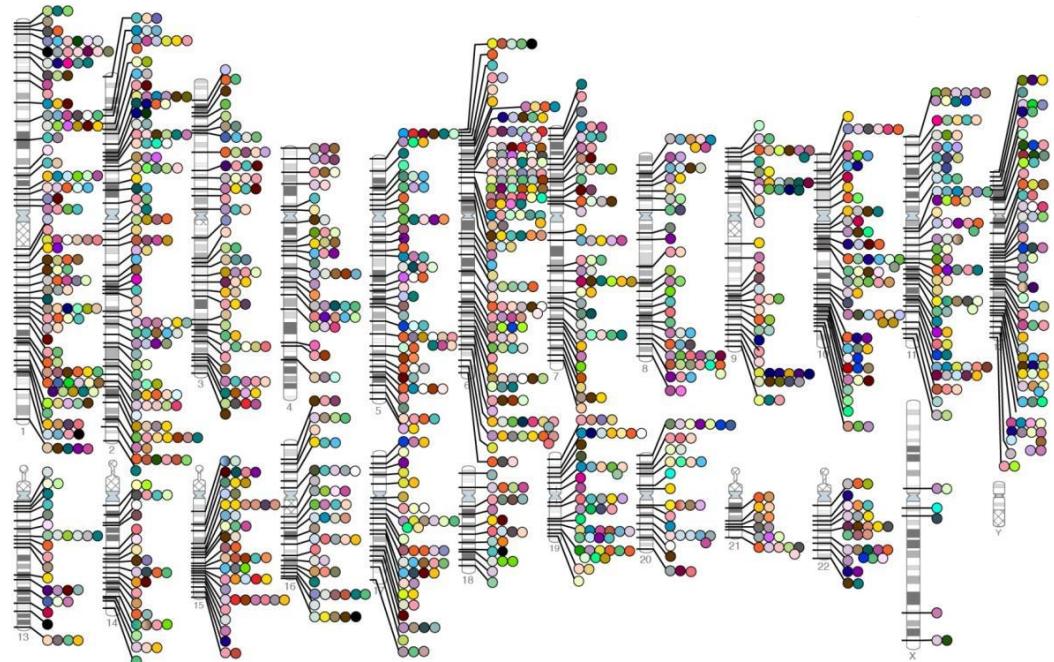


Published Genome-Wide Associations

By the time of the completion of the human genome sequence, in **2005**, just a **few** genetic variants were known to be significantly associated to diseases.

When the first exhaustive catalogue of GWAS was compiled, in 2008, only three years later, more than **500** single nucleotide polymorphisms (SNPs) were associated to traits.

Today, the catalog has collected more than 1,900 papers reporting **14,012** SNPs significantly associated to more than **1,500** traits.



NHGRI GWA Catalog
www.genome.gov/GWASStudies

Lessons learned from GWAS

- **Many loci/variants** contribute to complex-trait variation
- There is evidence for **pleiotropy**, i.e., that the same **loci/variants** are associated with multiple traits.
- Much of the **heritability** of the trait **cannot be explained** by the **individual** loci/variants found associated to the trait.

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 46/8 October 2009 doi:10.1038/nature08494 nature REVIEWS

Finding the missing heritability of complex diseases

Teri A. Manolio¹, Francis S. Collins², Nancy J. Cox³, David B. Goldstein⁴, Lucia A. Hindorff⁵, David J. Hunter⁶, Mark I. McCarthy⁷, Erin M. Ramos⁵, Lon R. Cardon⁸, Aravinda Chakravarti⁹, Judy H. Cho¹⁰, Alan E. Guttmacher¹, Augustine Kong¹¹, Leonid Kruglyak¹², Elaine Mardis¹³, Charles N. Rotimi¹⁴, Montgomery Slatkin¹⁵, David Valle⁹, Alice S. Whittemore¹⁶, Michael Boehnke¹⁷, Andrew G. Clark¹⁸, Evan E. Eichler¹⁹, Greg Gibson²⁰, Jonathan L. Haines²¹, Trudy F. C. Mackay²², Steven A. McCarroll²³ & Peter M. Visscher²⁴

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 this problem?

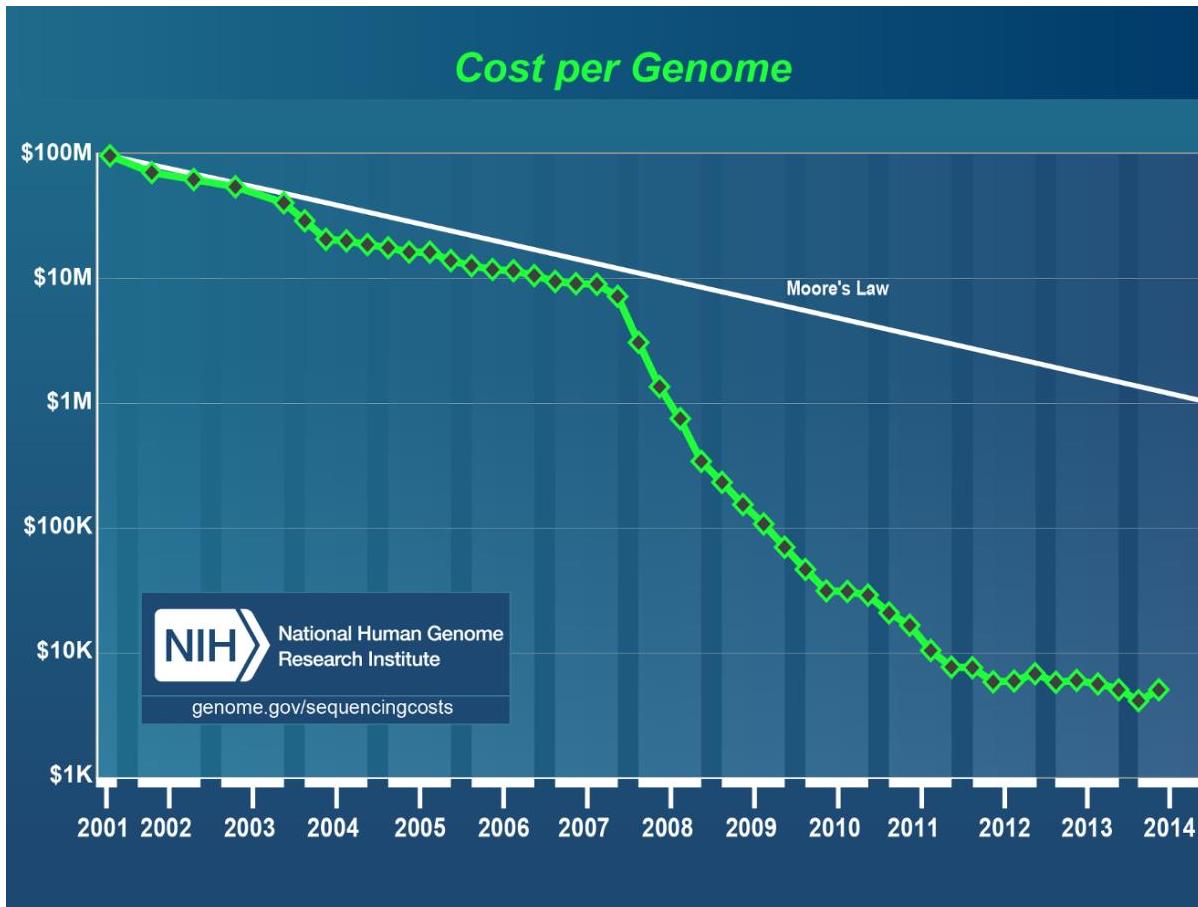
Rare Variants, rare CNVs, epigenetics or.. epistatic effects?

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 ⁷²	5	50%
Crohn's disease ²¹	32	20%
Systemic lupus erythematosus ⁷³	6	15%
Type 2 diabetes ⁷⁴	18	6%
HDL cholesterol ⁷⁵	7	5.2%
Height ¹⁵	40	5%
Early onset myocardial infarction ⁷⁶	9	2.8%
Fasting glucose ⁷⁷	4	1.5%

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

If rare variants eluded detection because were under represented among the SNPs, genomic sequencing would reveal them.



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^{1,10}, Kati J Buckingham^{2,10}, Choli Lee¹, Abigail W Bigham², Holly K Tabor^{2,3}, Karin M Dent⁴, Chad D Huff⁵, Paul T Shannon⁶, Ethylin Wang Jabs^{7,8}, Deborah A Nickerson¹, Jay Shendure¹ & Michael J Bamshad^{1,2,9}

We demonstrate the first successful application of exome sequencing to discover the gene for a rare mendelian disorder of unknown cause, Miller syndrome (OMIM 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*†, Sarah B. Ng‡, Abigail W. Bigham *§, Holly K. Tabor*||, Mary J. Emond¶, Deborah A. Nickerson† and Jay Shendure†

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.

OPEN ACCESS Freely available online

PLOS GENETICS

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

Evgeny A. Glazov^{1,*}, Andreas Zankl^{2,3}, Marina Donskoi¹, Tony J. Kenna¹, Gethin P. Thomas¹, Graeme R. Clark¹, Emma L. Duncan^{1,3}, Matthew A. Brown^{1*}

¹ University of Queensland Diamantina Institute, Princess Alexandra Hospital, Woolloongabba, Australia, ² Centre for Clinical Research, The University of Queensland, St. Lucia, Australia, ³ Murdoch Childrens Research Institute, Melbourne, Victoria, Australia

European Journal of Human Genetics (2011) 19, 115–117
© 2011 Macmillan Publishers Limited All rights reserved 1088-4813/11
www.nature.com/ejhg



small pedigrees
skeletal dysplasia,
dysplasia, brachydactyly.
The two
a rare form of
sequencing.
encodes a core
the *RMRP* RNA
and activity of
by which *POP1*

? Mutations As the
sense, which permits

SHORT REPORT

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

Kimia Kahrizi¹, Cougar Hao Hu², Masoud Garshabi², Seyedeh Sedigheh Abedini¹, Shirin Ghadami¹, Roxana Kariminejad³, Reinhard Ullmann⁴, Wei Chen², H-Hilger Ropers², Andreas W Kuss², Hossein Najmabadi¹ and Andreas Tschach^{*2}

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.203del
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

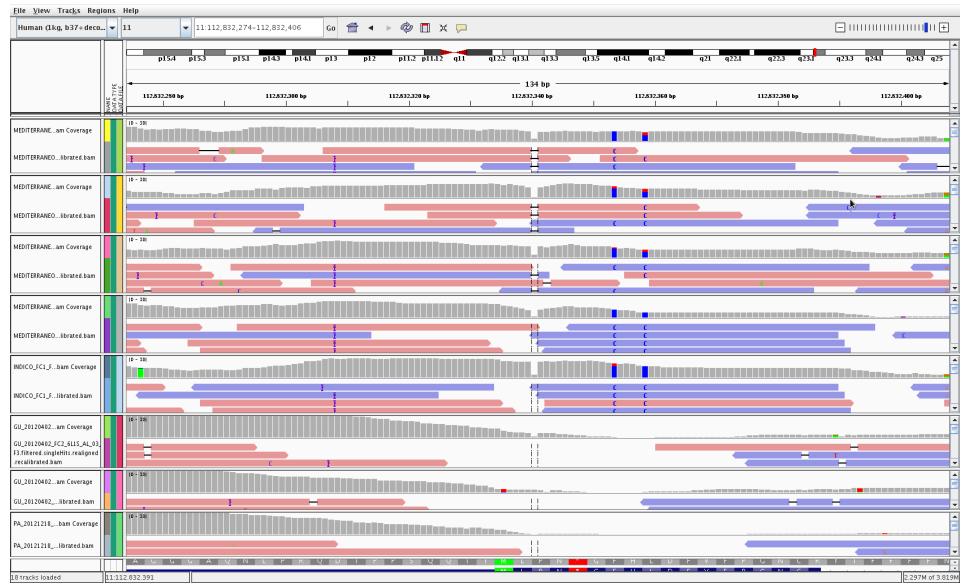
Whole-exome sequencing identifies novel compound heterozygous mutations in *USH2A* in Spanish patients with autosomal recessive retinitis pigmentosa

Cristina Méndez-Vidal,^{1,2} María González-del Pozo,^{1,2} Alicia Vela-Boza,³ Javier Santoyo-López,³ Francisco J. López-Domínguez,³ Carmen Vázquez-Marouschek,⁴ Joaquín Dopazo,^{3,5,6} Salud Borrego,^{1,2} Guillermo António,^{1,2,3}

¹Department of Genetics, Reproduction and Fetal Medicine, Institute of Biomedicine of Seville, University Hospital Virgen del Rocío/CSIC/University of Seville, Seville, Spain; ²Centro de Investigación Biomédica en Red de Enfermedades Raras (CIBERER), Seville, Spain; ³Medical Genome Project, Genomics and Bioinformatics Platform of Andalucía (GBPA), Seville, Spain;

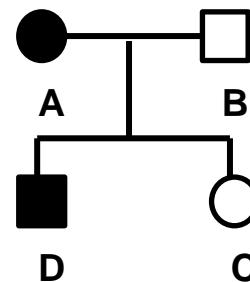
⁴Department of Ophthalmology, University Hospital Virgen del Rocío, Seville, Spain; ⁵Department of Bioinformatics, Centro de Investigación Príncipe Felipe, Valencia, Spain; ⁶Functional Genomics Node (INB), Centro de Investigación Príncipe Felipe, Valencia, Spain

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

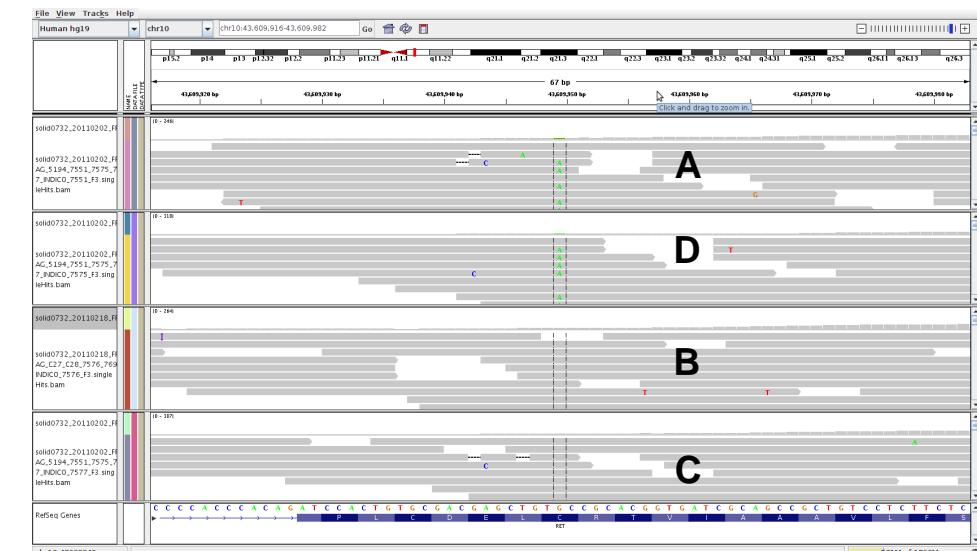


Cases

Controls



Segregation
within a
pedigree



Variant/gene prioritization by heuristic 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

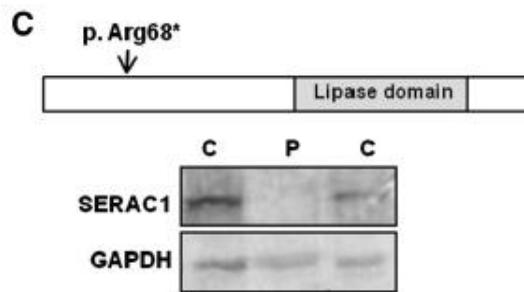
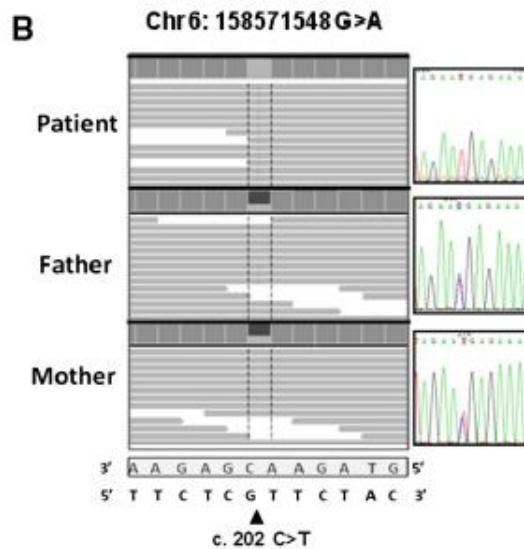
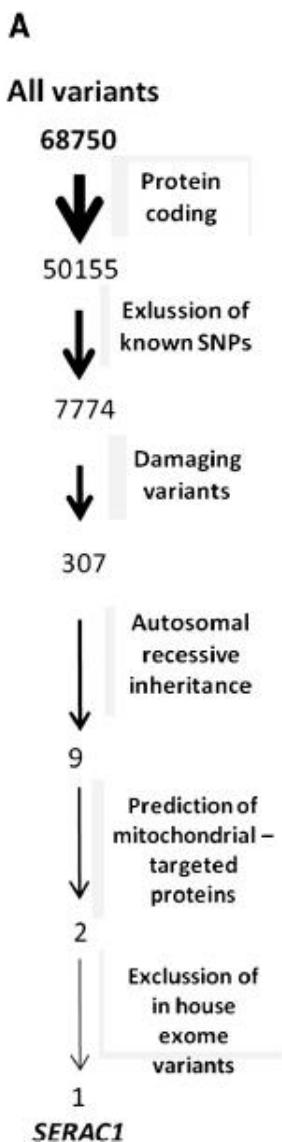
Control of sequencing errors (missing values)

Testing strategies

Heuristic 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 ^{a,b}, María Teresa García-Silva ^c, Xènia Ferrer-Cortès ^a, Aleix Navarro-Sastre ^{a,b}, Judith García-Villoria ^{a,b}, María Josep Coll ^{a,b}, Enrique Vidal ^d, Jorge Jiménez-Almazán ^d, Joaquín Dopazo ^{d,e,f}, Paz Briones ^{a,b,g}, Orly Elpeleg ^h, Antonia Ribes ^{a,b,*}

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

^b CIBER de Enfermedades Raras (CIBERER), Barcelona, Spain

^c Unidad de Enfermedades Mitochondriales- Enfermedades Metabólicas Hereditarias, Servicio de Pediatría, Hospital 12 de Octubre, Madrid, Spain

^d BIER, CIBERER, Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain

^e Computational Medicinal Institute, Centro de Investigación Príncipe Felipe (CIPF), Valencia, Spain

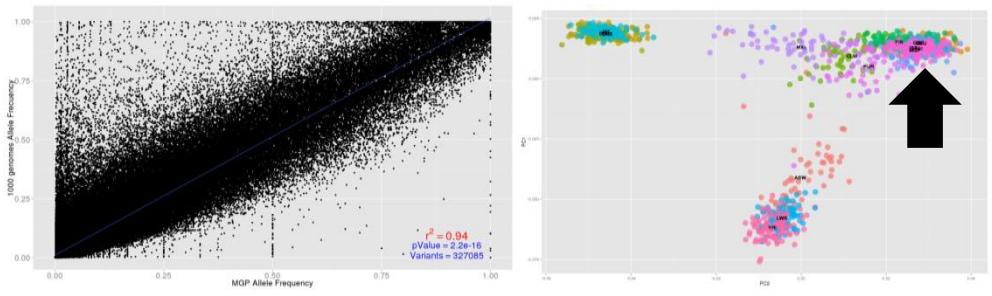
^f Functional Genomics Node, (INB) at CIPF, Valencia, Spain

^g Consejo Superior de Investigaciones Científicas (CSIC), Barcelona, Spain

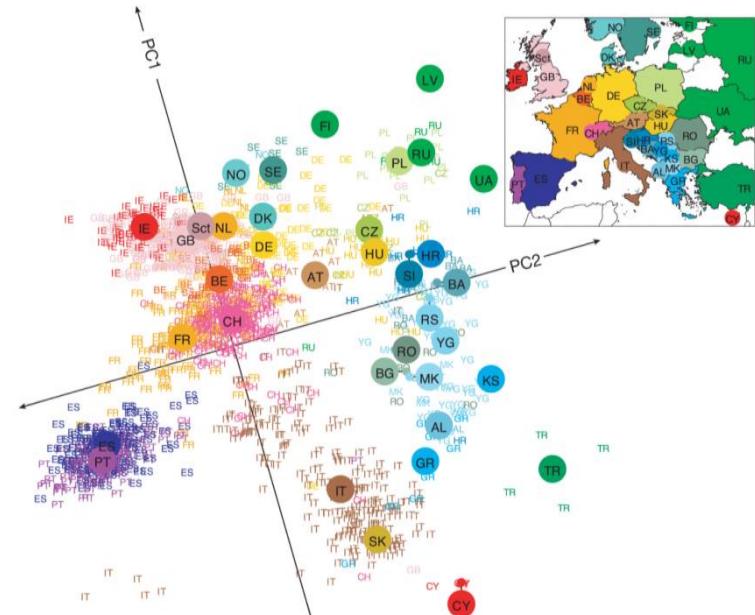
^h Monique and Jacques Roboh Department of Genetic Research, Hadassah, Hebrew University Medical Center, Jerusalem, Israel

Use known variants and their population frequencies to filter out false candidates

- Typically dbSNP, 1000 genomes and the 6515 exomes from the ESP are used as sources of population frequencies.
- We sequenced **300 healthy controls** (rigorously phenotyped) to add an extra filtering step to the analysis pipeline



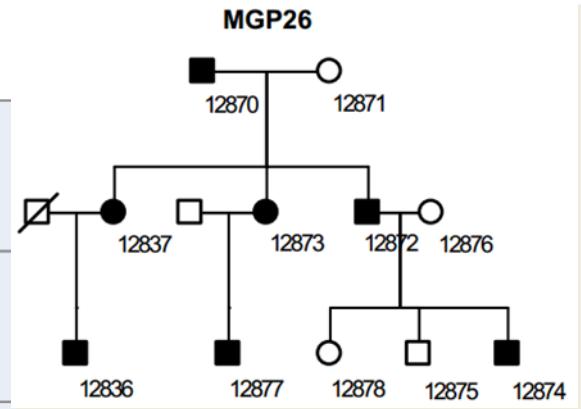
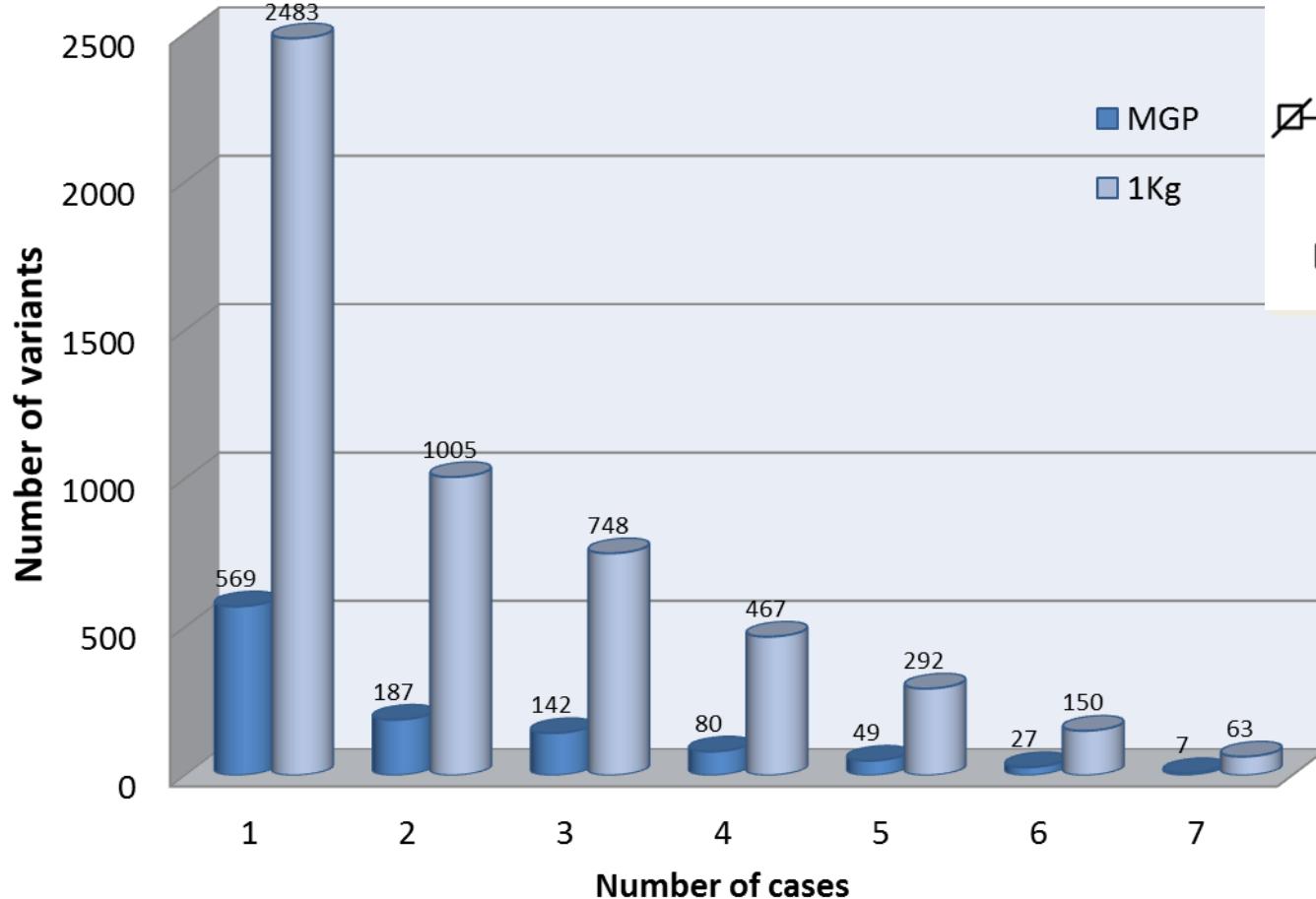
How important do you think is local information to detect disease genes?



Novembre et al., 2008. Genes mirror geography within Europe. Nature

Filtering with or without local variants

Number of genes as a function of individuals in the study of a dominant disease
Retinitis Pigmentosa autosomal dominant



The use of local variants makes an enormous difference

The CIBERER Exome Server (CES): the first repository of variability of the Spanish population

Used for more than one year and a half within CIBERER to discover new disease variants and genes.

Other similar initiatives arise. E.g.
the GoNL <http://www.nlgenome.nl/>

EJHG Open

European Journal of Human Genetics (2014) 22, 221–227
© 2014 Macmillan Publishers Limited. All rights reserved 1018-4813/14
www.nature.com/ejhg

ARTICLE

The Genome of the Netherlands: design, and project goals

Dorret I Boomsma^{*1,22}, Cisca Wijmenga^{2,22}, Eline P Slagboom^{3,22}, Morris A Swertz^{2,22}, Lennart C Karssemen⁴, Abdel Abdellaoui¹, Kai Ye⁵, Victor Guryev^{5,6}, Martijn Vermaat^{7,8,9}, Freerk van Dijk², Laurent C Franciolli¹⁰, Jouke Jan Hottenga¹, Jeroen FJ Laros^{7,8,9}, Qibin Li¹¹, Yingrui Li¹¹, Hongzhi Cao¹¹, Ruoyan Chen¹¹, Yuanping Du¹¹, Ning Li¹², Sujie Cao¹², Jessica van Setten¹⁰, Androniki Menelaou¹⁰, Sara L. Pulit¹⁰, Jayne Y Hehir-Kwa¹⁵, Marian Beekman¹⁶, Clara C Elbers¹⁰, Heorhiy Byelas², Anton JM de Craen¹⁶, Patrick Deelen², Martijn Dijkstra², Johan T den Dunnen^{8,9}, Peter de Knijff^{8,9}, Jeanine Houwing-Duistermaat¹⁷, Vyacheslav Koval¹⁸, Karol Estrada¹⁸, Albert Hofman⁴, Alexandros Kanterakis², David van Eckvort⁷, Hailiang Mai⁷, Mathijs Kattenberg¹, Elisabeth M van Leeuwen⁴, Pieter BT Neerinckx², Ben Oostra¹⁹, Fernando Rivadeneira¹⁸, Elka HD Suchiman³, Andre G Uitterlinden¹⁸, Gonnieke Willenssen¹, Bruce H Wolffenbuttel²⁰, Jun Wang^{11,13,14,22}, Paul IW de Bakker^{10,22}, Gert-Jan van Ommen^{21,22} and Cornelia M van Duijn^{4,22}

Exome Server															
Filters			Variant Info												
Variant	Alleles	Gene	BIER				1000G				EVs				
			0/0	0/1	1/1	./-	MAF	0/0	0/1	1/1	./-	MAF	0/0	0/1	1/1
10:43572699	C>A	RET	74	.	1	.	0.013
10:43572721	G>A	RET	74	1	.	.	0.007
10:43596968	A>G	RET	4	20	51	.	0.187	107	405	580	.	0.283	296	1893	4317
10:43596003	G>A	RET	74	1	.	.	0.007
10:43596179	G>A	RET	52	18	5	.	0.187	755	289	48	.	0.176	4133	1997	356
10:43596182	G>A	RET	73	2	.	.	0.013	1089	3	.	.	0.001	6456	26	1
10:43597827	C>A	RET	72	3	.	.	0.020	1074	18	.	.	0.008	6391	113	2
10:43600372	G>A	RET	74	.	1	.	0.013
10:43600517	G>A	RET	74	1	.	.	0.007
10:43600521	G>C	RET	74	1	.	.	0.007
10:43600689	A>G	RET	39	27	9	.	0.300	421	499	172	.	0.386	2560	3023	911
10:43601985	C>A	RET	74	1	.	.	0.007
10:43602007	G>A	RET	74	1	.	.	0.007
10:43606650	C>T	RET	74	1	.	.	0.007	1037	53	2	.	0.026	6085	396	25
10:43606687	A>G	RET	8	42	25	.	0.387	73	347	672	.	0.226	465	2343	3698
10:43607516	C>A	RET	74	1	.	.	0.007
10:43607590	C>T	RET	74	.	1	.	0.013
10:43607600	C>A	RET	74	.	1	.	0.013
10:43608966	G>A	RET	74	1	.	.	0.013
10:43609008	C>A	RET	74	1	.	.	0.007
10:43609163	C>A	RET	74	1	.	.	0.007
10:43609889	C>A	RET	74	1	.	.	0.013
10:43610119	G>A	RET	48	25	2	.	0.193	790	256	36	.	0.155	4636	1695	175
10:43612225	G>C	RET	74	1	.	.	0.007	1083	9	.	.	0.004	6446	56	1
10:43613943	G>T	RET	3	23	49	.	0.193	110	385	597	.	0.277	259	2052	4195

<http://ciberer.es/bier/exome-server/>

Information provided

Genomic coordinates,
variation, and gene.

SNPid
if any

Genotypes in the
different reference
populations

Exome Server

Summary Variants Genome Viewer

Filters

Reload Clear Search

Region/Gene

Region Gene

Enter genes (comma separated)
RET

Controls +

Variant Info

Variant	Alleles	SNP Id	Gene	BIER				1000G				EVS				Poly			
				Genotypes				MAF	Genotypes				MAF	Genotypes				MAF	
				0/0	0/1	1/1	./.		0/0	0/1	1/1	./.		0/0	0/1		1/1		./.
10:43572699	C>A .	RET	74	.	1	.	0.013				
10:43572721	G>A .	RET	74	1	.	.	0.007				
10:43595968	A>G rs1800858	RET	4	20	51	.	0.187	107	405	580	.	0.283	296	1893	4317	.	0.191		
10:43596003	G>A .	RET	74	1	.	.	0.007		
10:43596179	G>A rs2435351	RET	52	18	5	.	0.187	755	289	48	.	0.176	4133	1997	356	.	0.209		
10:43596182	G>A rs200468424	RET	73	2	.	.	0.013	1089	3	.	.	0.001	6456	26	1	.	0.002		
10:43597827	C>A rs1800859	RET	72	3	.	.	0.020	1074	18	.	.	0.008	6391	113	2	.	0.009		
10:43600372	G>A .	RET	74	.	1	.	0.013		
10:43600517	G>A .	RET	74	1	.	.	0.007		
10:43600521	G>C .	RET	74	1	.	.	0.007		
10:43600689	A>G rs2435352	RET	39	27	9	.	0.300	421	499	172	.	0.386	2560	3023	911	.	0.373		
10:43601985	C>A .	RET	74	1	.	.	0.007		
10:43602007	G>A .	RET	74	1	.	.	0.007		

Page 1 of 2 | ► | 🔍

Variants 1 - 25 of 36

Columns

Occurrence of pathological variants in “normal” population

Exome Server

Summary Variants Genome Viewer

Filters
Reload Clear Search
Region/Genome
 Region Gene
Enter genes (comma separated)
BBS2

Controls +

Variant	Alleles	Gene	BIER				1000G				EVS				Polyphen	SIFT	Phenotype				
			Genotypes		MAF	Genotypes		MAF	Genotypes		MAF	Genotypes		MAF							
			0/0	0/1		1/1	./.		0/0	0/1		1/1	./.								
16:56501806	C>T	OGFOD1,BBS2	74	1	.	.	.	0.007	736	316	40	.	0.181	4578	1758	165	.	0.160	BODY MASS INDEX,Height,Two-hour glu...		
16:56504724	G>C	OGFOD1,BBS2	40	28	7	.	.	0.280	830	237	25	.	0.131	4202	2024	275	.	0.198	BODY MASS INDEX,Height,Two-hour glu...		
16:56508721	T>C	OGFOD1,BBS2	74	1	.	.	.	0.007		
16:56508883	C>T	OGFOD1,BBS2	74	1	.	.	.	0.007		
16:56509441	T>G	BBS2,OGFOD1	73	2	.	.	.	0.013		
16:56510072	A>C	BBS2,OGFOD1	73	2	.	.	.	0.013	1085	7	.	.	0.003	6375	125	1	.	0.010			
16:56533804	T>G	BBS2	74	1	.	.	.	0.007	1076	15	1	.	0.008	6339	161	1	.	0.012			
16:56535193	C>T	BBS2	72	3	.	.	.	0.020			
16:56535207	AG...>...	BBS2	74	1	.	.	.	0.007			
16:56543827	A>G	BBS2	74	1	.	.	.	0.007	1084	8	.	.	0.004	6402	98	1	.	0.007			
16:56545175	T>C	BBS2	19	23	1	.	.	0.100	601	101	67	.	0.265	1193	2801	241	.	0.195	Fasting proinsulin/secretive cotts,Tot...		
16:56548501	C>T	BBS2	.	.	75	.	.	0.000	.	9	1083	.	0.004	1	73	6427	.	0.006	0.001	1	BARDET-BIEDL SYNDROME 2,Bardet-Bie...
16:56553814	A>G	BBS2	74	1	.	.	.	0.007	910	163	19	.	0.092	5706	741	50	.	0.065			
16:56553816	A>C	BBS2	74	1	.	.	.	0.007	910	163	19	.	0.092	5744	704	49	.	0.062			

Page 1 of 1 | << << >> >> | Columns | Variants 1 - 14 of 14

Reference genome is mutated

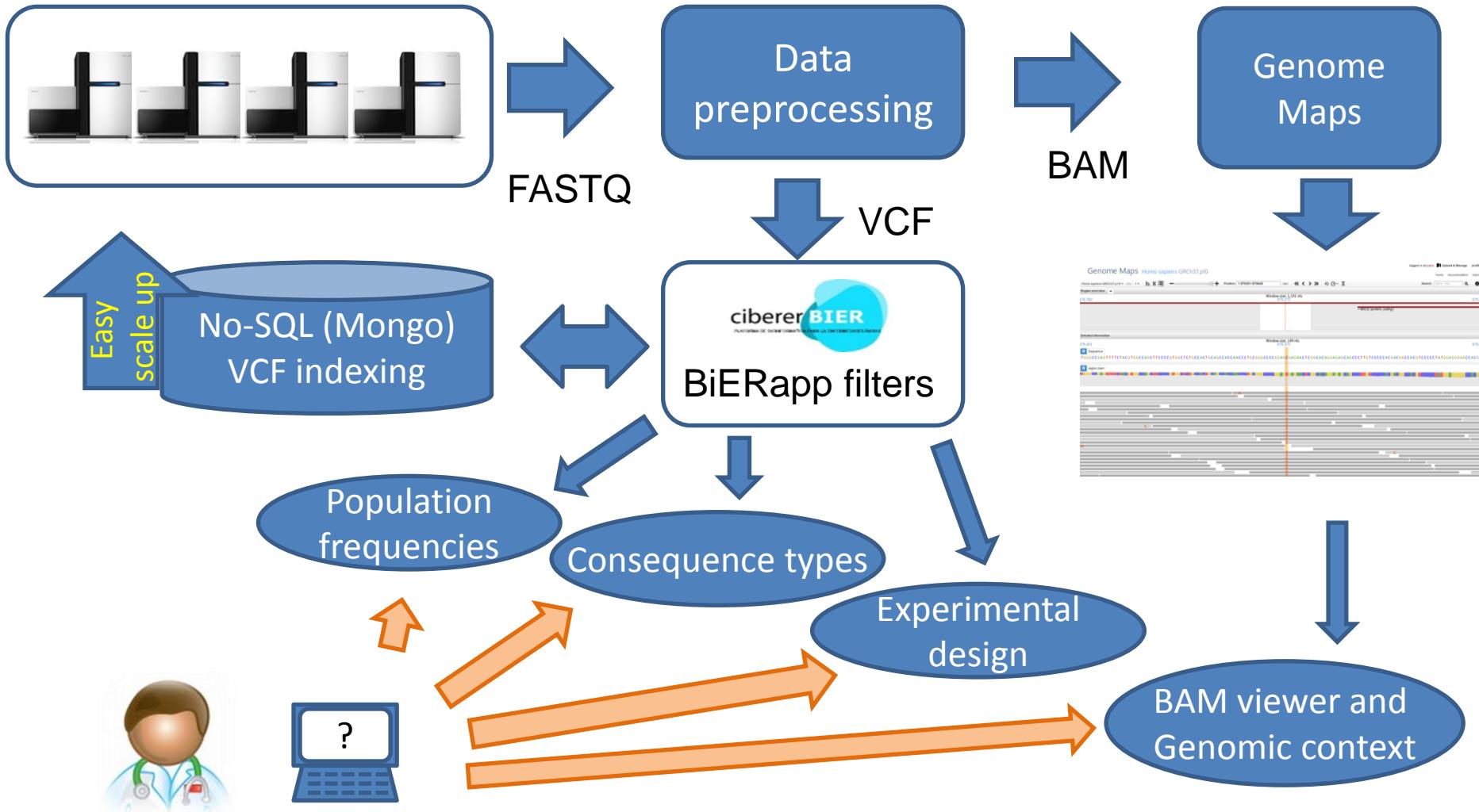
Nine carriers in 1000 genomes

One affected and 73 carriers in EVS

An example of end user's tool

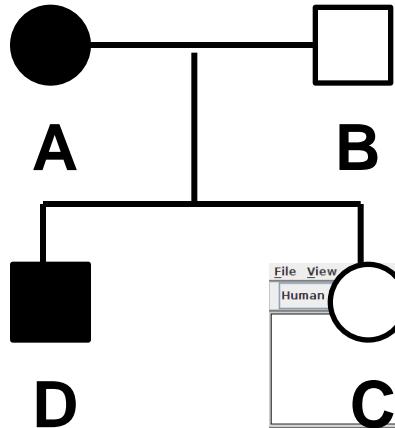
BiERapp: interactive web-based tool for easy candidate prioritization by heuristic filtering

SEQUENCING CENTER

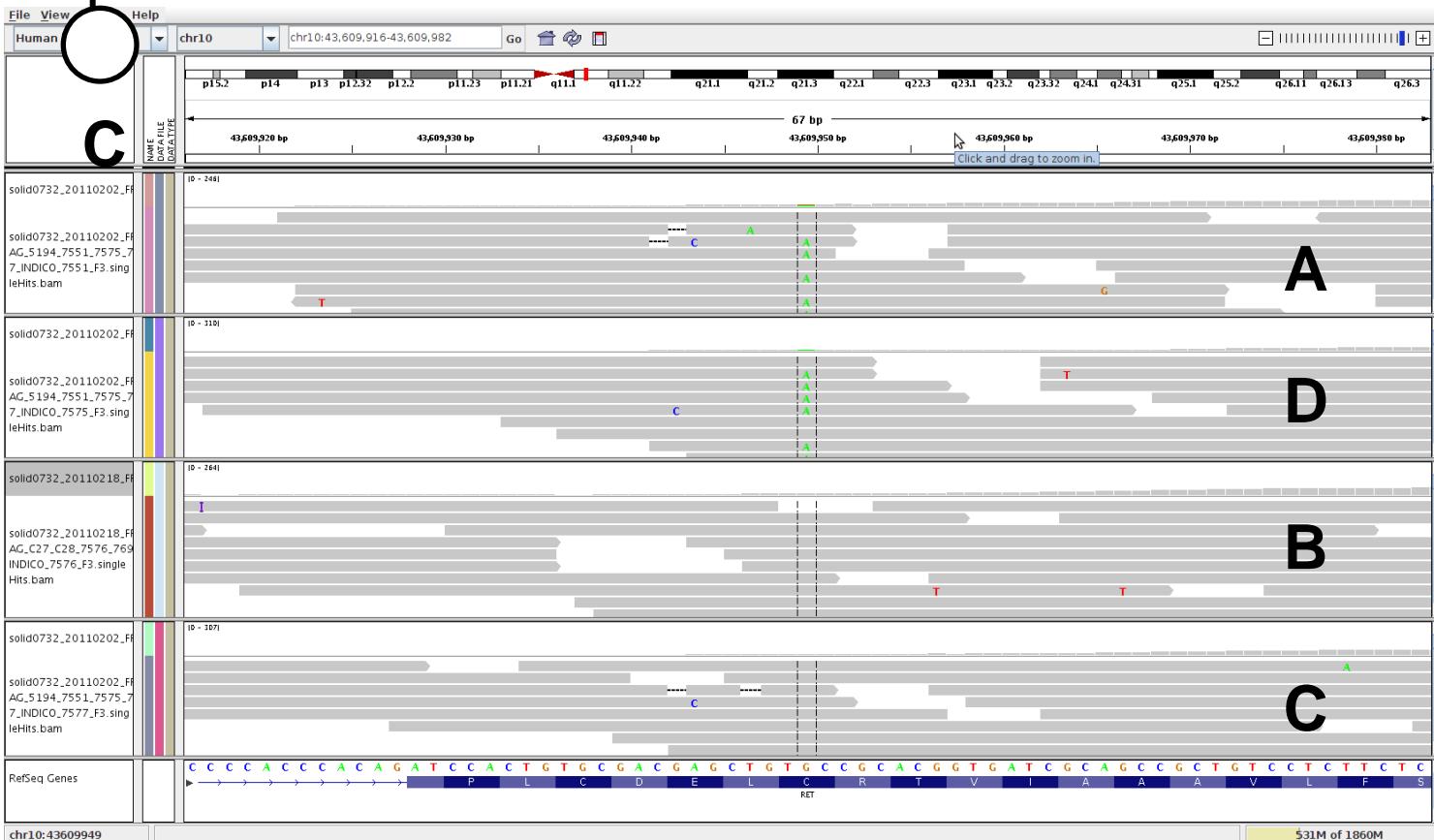


How efficient is exome/genome sequencing?

Low rate of false negatives. An example with MTC



Dominant:
Heterozygotic in A and D
Homozygotic reference allele in B and C
Homozygotic reference allele in controls



The
codon
634
mutation

Heuristic filtering approach. Exome sequencing produces many false positives

Table 1 | Mean number of coding variants in two populations

Variant type	Mean number of variants (\pm sd) in African Americans	Mean number of variants (\pm sd) in European Americans
Novel variants		
Missense	303 (\pm 32)	192 (\pm 21)
Nonsense	5 (\pm 2)	5 (\pm 2)
Synonymous	209 (\pm 26)	109 (\pm 16)
Splice	2 (\pm 1)	2 (\pm 1)
Total	520 (\pm 53)	307 (\pm 33)
Non-novel variants		
Missense	10,828 (\pm 342)	9,319 (\pm 233)
Nonsense	98 (\pm 8)	89 (\pm 6)
Synonymous	12,567 (\pm 416)	10,536 (\pm 280)
Splice	36 (\pm 4)	32 (\pm 3)
Total	23,529 (\pm 751)	19,976 (\pm 505)
Total variants		
Missense	11,131 (\pm 364)	9,511 (\pm 244)
Nonsense	103 (\pm 8)	93 (\pm 6)
Synonymous	12,776 (\pm 434)	10,645 (\pm 286)
Splice	38 (\pm 5)	34 (\pm 4)
Total	24,049 (\pm 791)	20,283 (\pm 523)

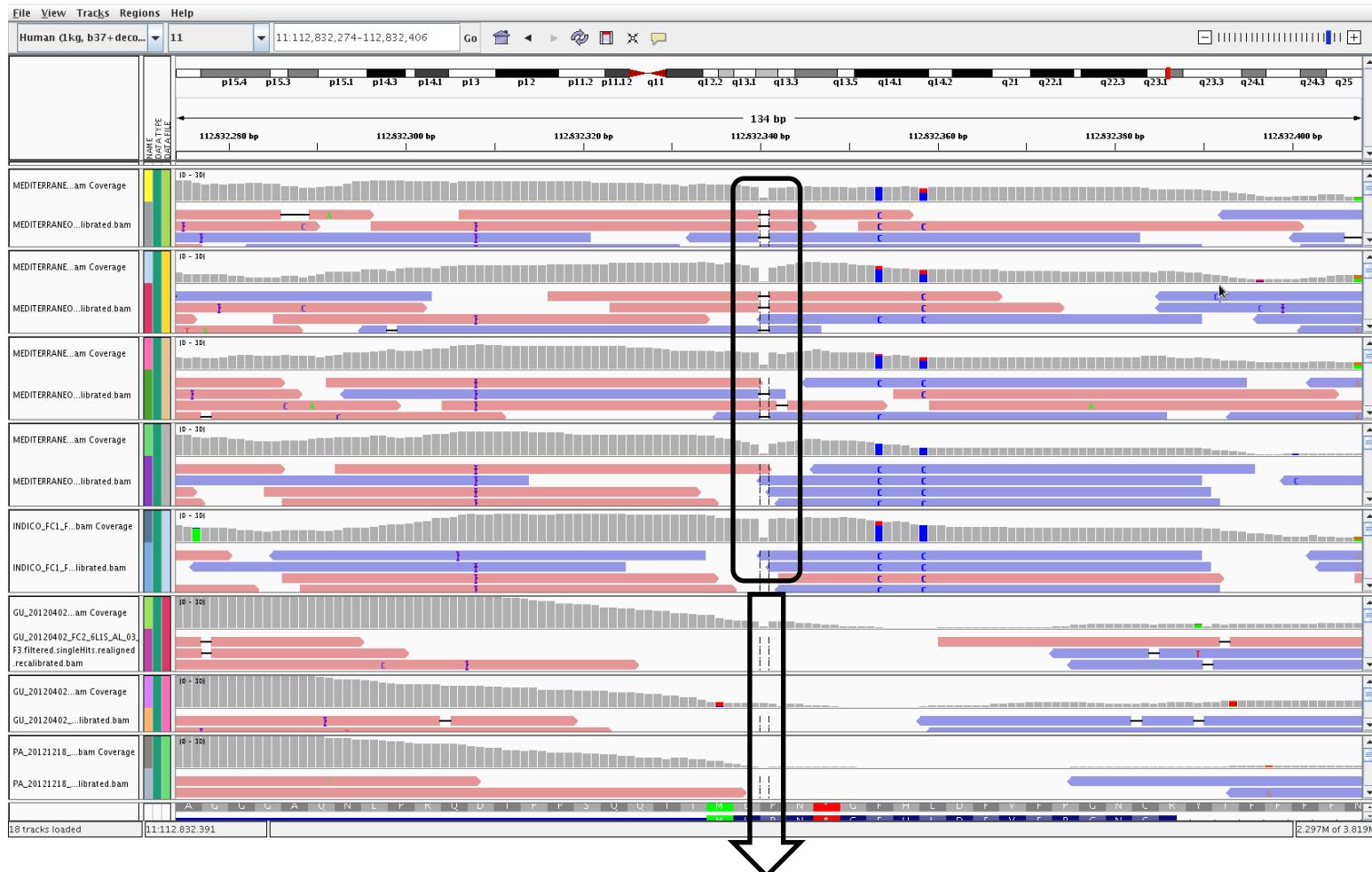
The table lists the mean number (\pm standard deviation (sd)) of novel and non-novel coding single nucleotide variants from 100 sampled African Americans and 100 European Americans. Non-novel variants refer to those found in dbSNP131 or in 200 other control exomes. Capture was performed using the Nimblegen V2 target. The analysis pipeline consisted of: alignment using the Burrows-Wheeler alignment tool; recalibration; realignment around insertion-deletions and merging with the Genome Analysis Toolkit (GATK)⁹¹; and removal of duplicates with PICARD. Variants were called using the following parameters: quality score > 50, allele balance ratio < 0.75; homopolymer run > 3; and quality by depth < 8. Variants were called from a RefSeq37.2 target (35,804,408 bp).

Average values obtained per exome (>800)

After filtering by:	SNVs
Conventional filter QC, coverage...	60,000
Mapping and haplotype coherence, missing sites...	30,000
Nonsynonymous (nonsense and missense)	5,000
Unknown (not present in controls)	150-300
Segregate with the families	< 100

We can detect the disease mutation(s)... along with many other unrelated variants

Some false positives are errors that can easily be avoided. E.g. missing positions



The promising variant (a frameshift present in all patients but not detected in controls) was not real. It was not properly covered by reads in controls.

**And there are many real variants
with potential phenotypic effect**

Findings:

20.000 total variants

1000 new variants

300-500 LOF variants (>50 homozygous)

100 known variants associated to disease

A report must contain:

- 1) Diagnostic variants
 - 2) Therapy-related variants
 - 3) Susceptibility variants
 - 4) Incidental findings with risk for the patient

My first exome...

List of variants

A high level of deleterious variability exists in the human genome

- Variants predicted to severely affect the function of human protein coding genes known as loss-of-function (LOF) variants were thought:
 - To have a potential deleterious effect
 - To be associated to severe Mendelian disease
- However, an unexpectedly large number of LOF variants have been found in the genomes of apparently healthy individuals: 281-515 missense substitutions per individual, 40-85 of them in homozygous state and predicted to be highly damaging.
- A similar proportion was observed in miRNAs and possibly affect to any functional element in the genome

ARTICLE

Deleterious- and Disease-Allele Prevalence in Healthy Individuals: Insights from Current Predictions, Mutation Databases, and Population-Scale Resequencing

Yali Xue,¹ Yuan Chen,¹ Qasim Ayub,¹ Ni Huang,¹ Edward V. Ball,² Matthew Mort,² Andrew D. Phillips,² Katy Shaw,² Peter D. Stenson,² David N. Cooper,² Chris Tyler-Smith,^{1,*} and the 1000 Genomes Project Consortium

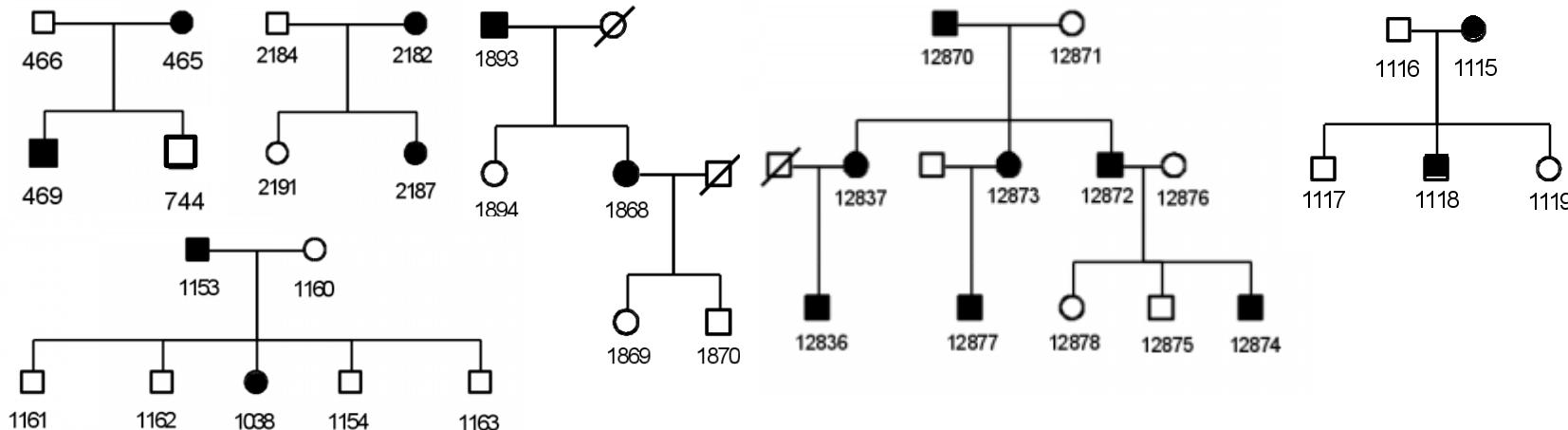
We have assessed the numbers of potentially deleterious variants in the genomes of apparently healthy humans by using (1) low-coverage whole-genome sequence data from 179 individuals in the 1000 Genomes Pilot Project and (2) current predictions and databases of deleterious variants. Each individual carried 281–515 missense substitutions, 40–85 of which were homozygous, predicted to be highly damaging. They also carried 40–110 variants classified by the Human Gene Mutation Database (HGMD) as disease-causing mutations (DMs), 3–24 variants in the homozygous state, and many polymorphisms putatively associated with disease. Whereas many of these DMs are likely to represent disease-allele-annotation errors, between 0 and 8 DMs (0–1 homozygous) per individual are predicted to be highly damaging, and some of them provide information of medical relevance. These analyses emphasize the need for improved annotation of disease alleles both in mutation databases and in the primary literature; some HGMD mutation data have been recategorized on the basis of the present findings, an iterative process that is both necessary and ongoing. Our estimates of deleterious-allele numbers are likely to be subject to both overcounting and undercounting. However, our current best mean estimates of ~400 damaging variants and ~2 bona fide disease mutations per individual are likely to increase rather than decrease as sequencing studies ascertain rare variants more effectively and as additional disease alleles are discovered.

The screenshot shows the homepage of the **Genome Medicine** journal. At the top right is a search bar with the placeholder "Search Genome Medicine for". Below the search bar are navigation links: Home, Articles (which is highlighted in green), Authors, Reviewers, About this journal, My Genome Medicine, and Subscriptions. In the center, there is a featured article titled "A map of human microRNA variation uncovers unexpectedly high levels of variability". The authors listed are Jose Carbonell, Eva Alloza, Pablo Arce, Salud Borrego, Javier Santoyo, Macarena Ruiz-Ferrer, Ignacio Medina, Jorge Jimenez-Almazan, Cristina Mendez-Vidal, Maria Gonzalez-del Pozo, Alicia Vela, Shomi S Bhattacharya, Guillermo Antinolo and Joaquin Dopazo. A red button labeled "Highly accessed" and a blue button labeled "Open Access" are visible above the article title.

Such apparently deleterious mutation must be first detected and then distinguished from real pathological mutations

Moreover, even Mendelian genes can be elusive.

Intuitive belief: multiple family information should help



	Families					
	1	2	3	4	5	6
Variants	3403	82	4	0	0	0
Genes	2560	331	35	8	1	0

Observation: this is not always true, not even in cases of Mendelian diseases

Is the single-gene approach realistic? Can we easily detect disease-related variants?

There are several problems:

- a) Interrogating 60Mb sites (3000 Mb in genomes) produces too many variants. A large number of these segregating with our experimental design
- b) There is a non-negligible amount of **apparently deleterious** variants that (apparently) has no pathologic effect
- c) In many cases we are not targeting rare but **common** variants (which occur in normal population)
- d) In many cases only one variant does not explain the disease but rather a **combination** of them (epistasis)
- e) Consequently, the few individual variants found associated to the disease usually account for a **small portion** of the trait **heritability**

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?**

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 ⁷²	5	50%
Crohn's disease ²¹	32	20%
Systemic lupus erythematosus ⁷³	6	15%
Type 2 diabetes ⁷⁴	18	6%
HDL cholesterol ⁷⁵	7	5.2%
Height ¹⁵	40	5%
Early onset myocardial infarction ⁷⁶ *	9	1.9%
Fasting glucose ⁷⁷	4	1.5%

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

Human
genetics

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.

2010 Nature America, Inc. All rights reserved.

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

Jian Yang¹, Beben Benyamin¹, Brian P McEvoy¹, Scott Gordon¹, Anjali K Henders¹, Dale R Nyholt¹, Pamela A Madden², Andrew C Heath², Nicholas G Martin¹, Grant W Montgomery¹, Michael E Goddard³ & Peter M Visscher¹

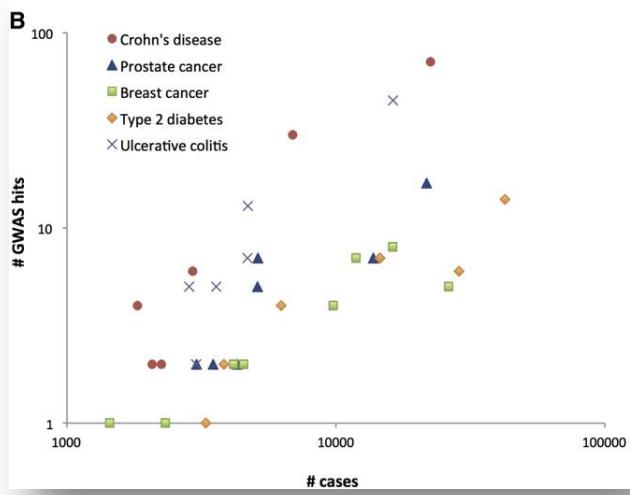
SNPs discovered by genome-wide association studies (GWASs) account for only a small fraction of the genetic variation of complex traits in human populations. Where is the remaining heritability? We estimated the proportion of variance for human height explained by 294,831 SNPs genotyped in 3,925 unrelated 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.

At the end, most of the heritability was there... 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 it with the

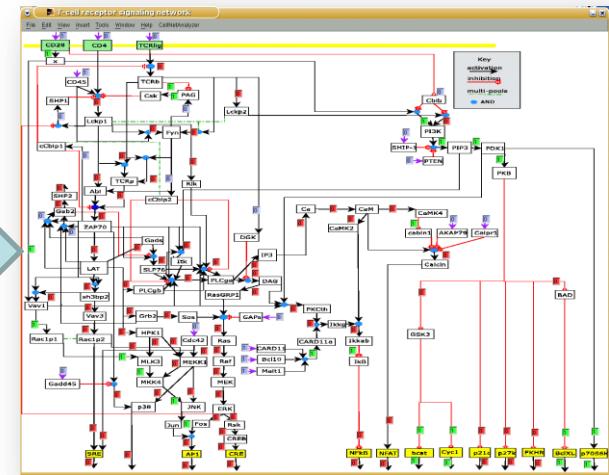
Height in humans is a classical quantitative trait, easy to measure and studied for well over a century as a model for investigating the genetic basis of complex traits. The heritability of height has been estimated to be ~0.8 (refs. 9,11–13). Rare mutations that cause extreme short or tall stature have been found^{14,15}, 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^{16–19}.

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

At the crossroad: how detection power of genomic technologies can be increased?



There are two (non mutually exclusive) ways



Scaling up: by increasing sample size.

It is known that larger size allows detecting more individual gene (biomarker) associations.

Limitations: Budget, patients availability and the own nature of the disease.

Changing the perspective: systems approach to understand variation

Interactions, multigenicity can be better detected and the role of variants understood in the context of disease mechanism.

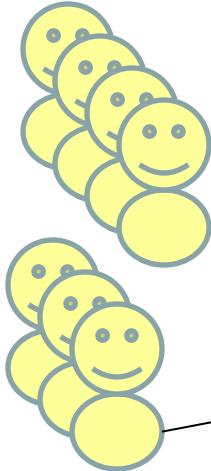
Limitations: Available information

Modular nature of human genetic diseases

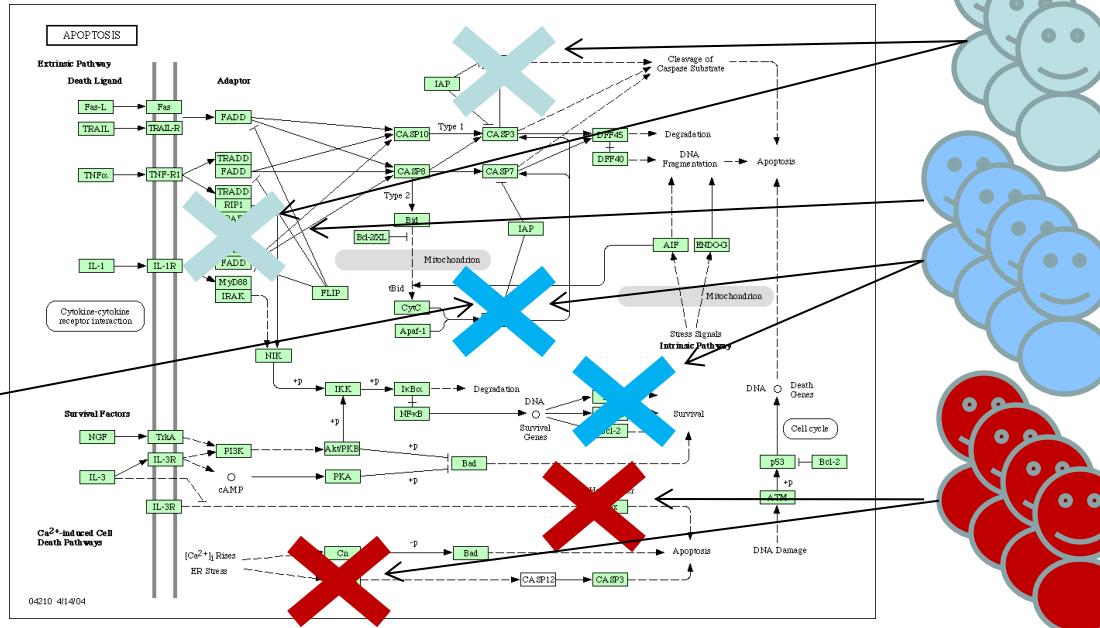
- With the development of **systems biology**, studies have shown that phenotypically similar diseases are often caused by **functionally related genes**, being referred to as the **modular nature of human genetic diseases** (Oti and Brunner, 2007; Oti et al, 2008).
- This modularity suggests that **causative genes** for the same or phenotypically similar diseases may generally reside in the same **biological module**, either a **protein complex** (Lage et al, 2007), a **sub-network** of protein interactions (Lim et al, 2006) , or a **pathway** (Wood et al, 2007)

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

Controls



Cases



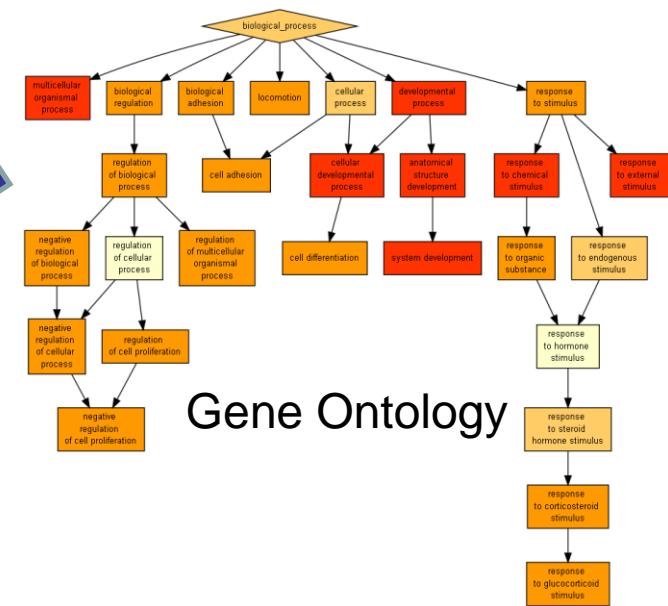
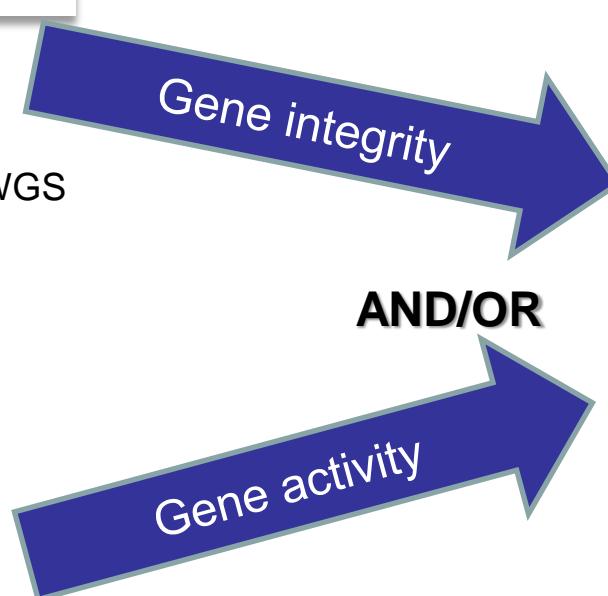
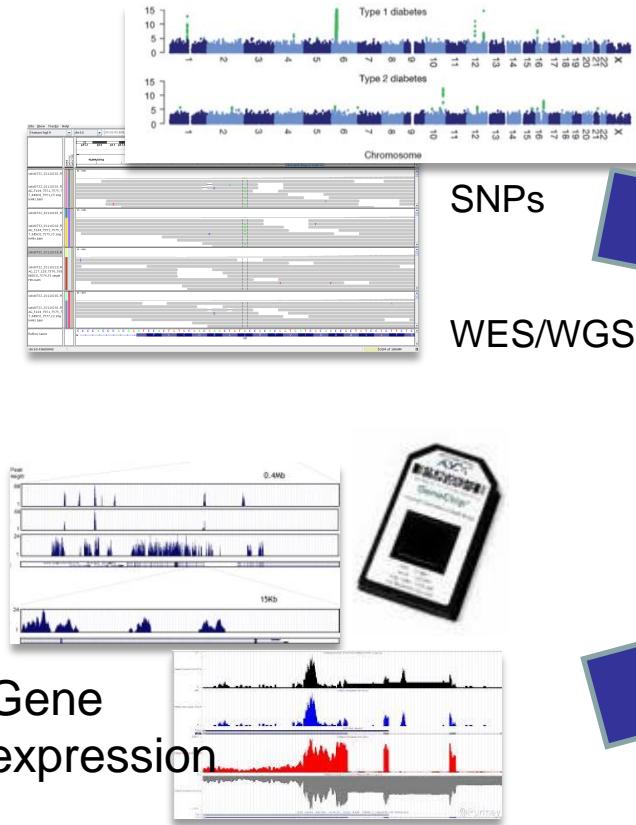
Affected **cases** in complex diseases will be a **heterogeneous** population with different mutations (or combinations).

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.

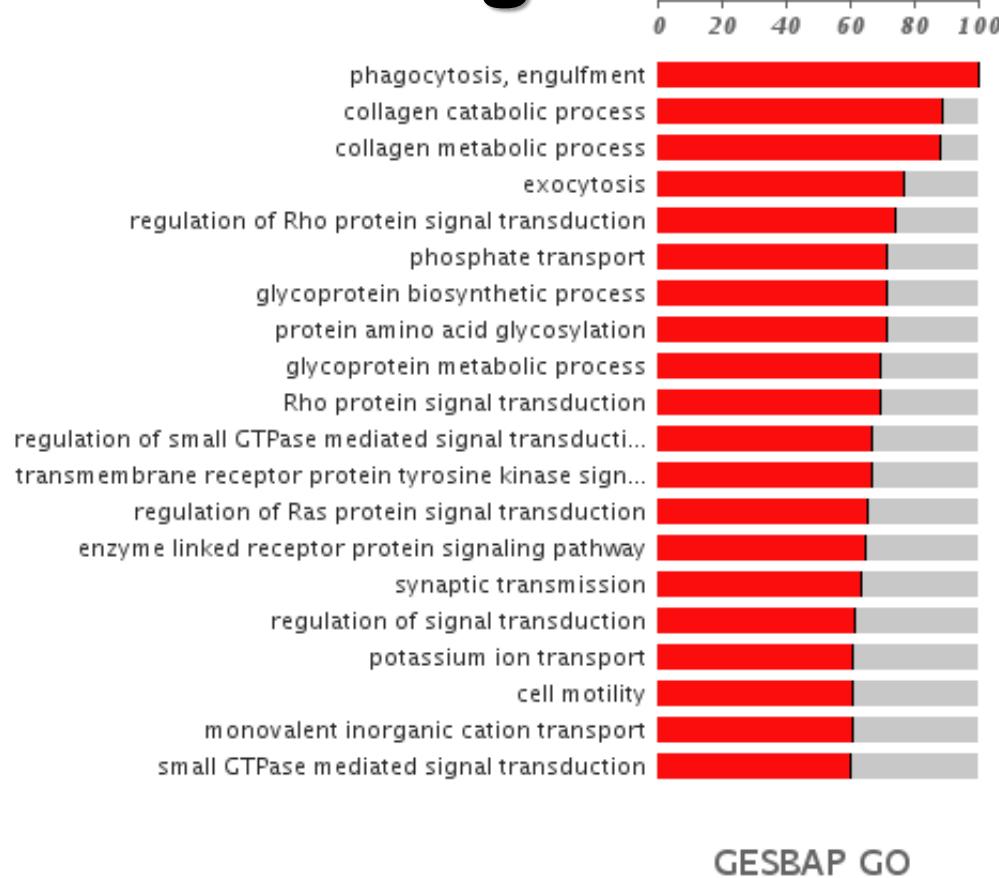
The CGEMS initiative. (Hunter et al. Nat Genet 2007)

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

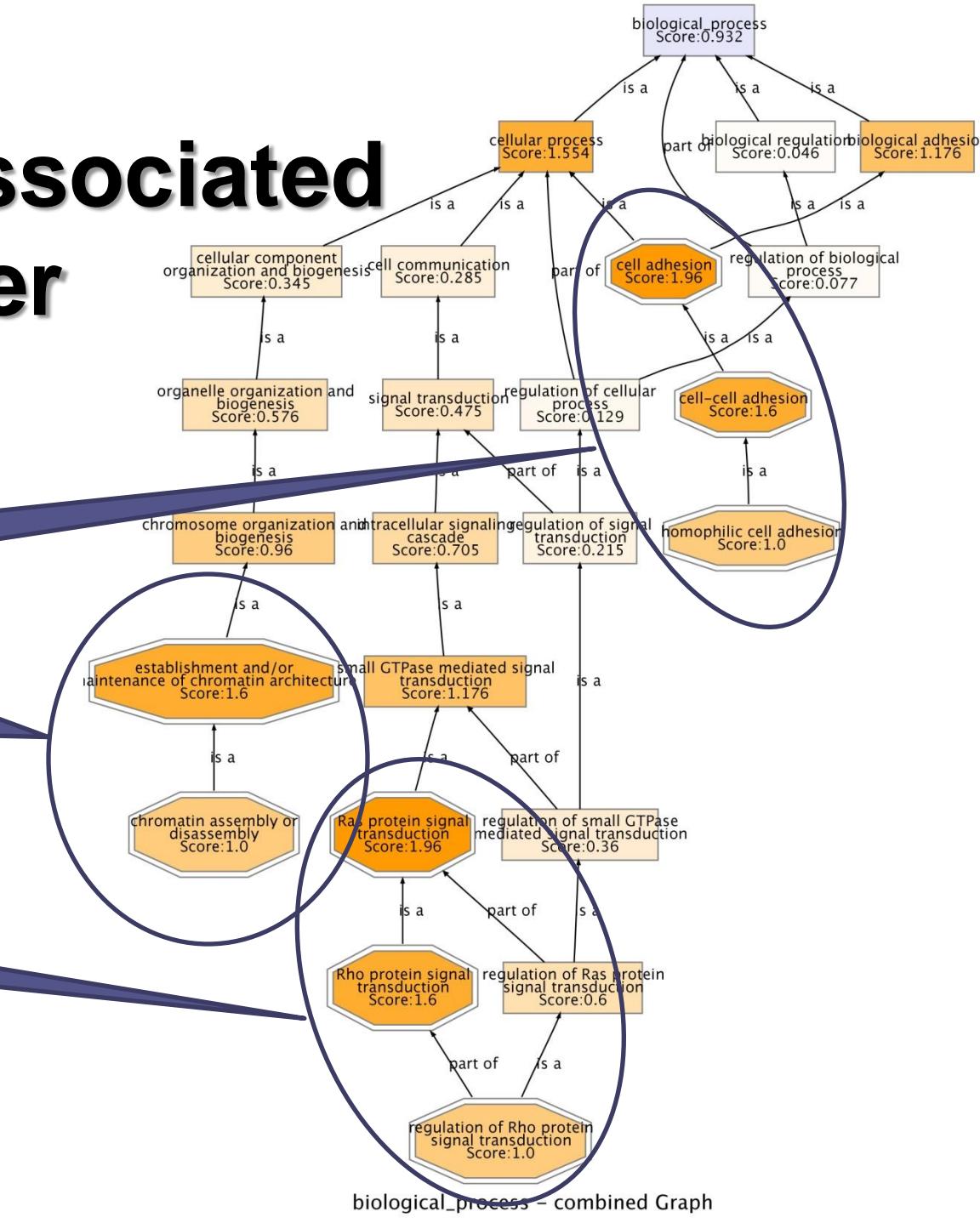
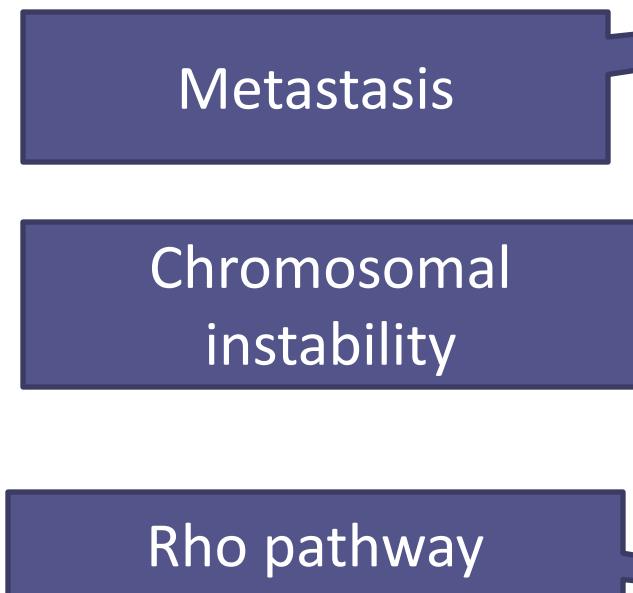
CGEMS initiative.
(Hunter et al. Nat Genet 2007)

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=4.45x10⁻⁰³) in which FGFR2 is included.

GO processes significantly associated to breast cancer



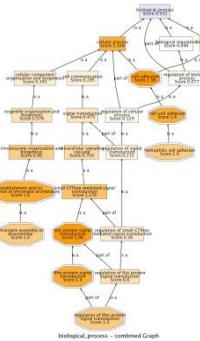
From gene-based to function-based perspective

SNPs,
Gene expression

Gene₁
Gene₂
Gene₃
Gene₄
:
:
:
Gene₂₂₀₀₀



Gene
Ontology

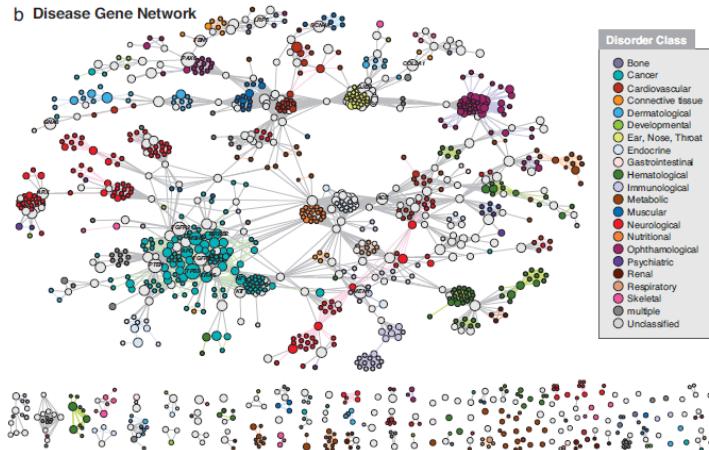


	SNPs, gene exp.	GO
Detection power	Low (only very prevalent genes)	high
Annotations available	many	many
Use	Biomarker	Illustrative, give hints

Can the interactome help to find disease mutations?

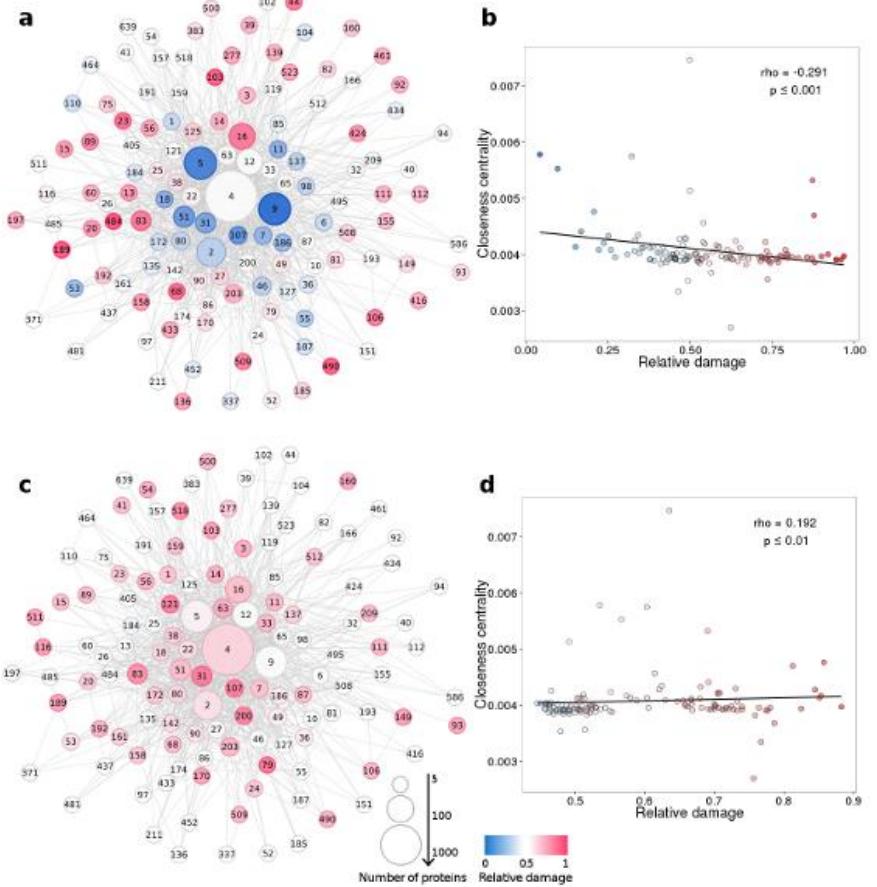
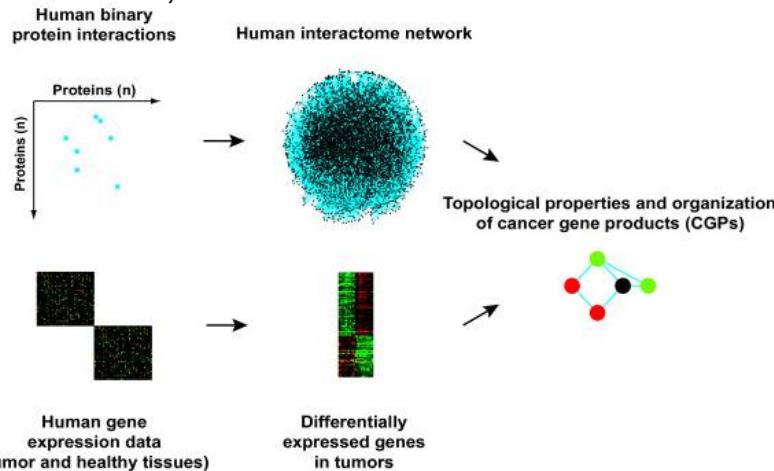
Disease genes are close in the interactome

Goh 2007 PNAS



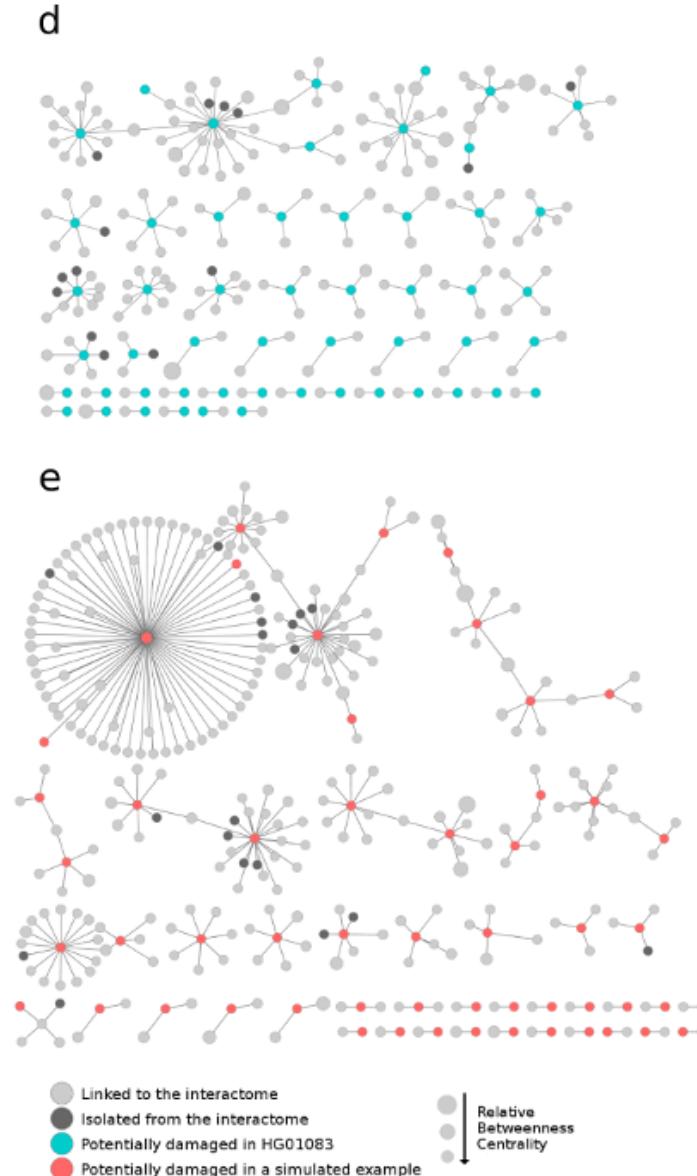
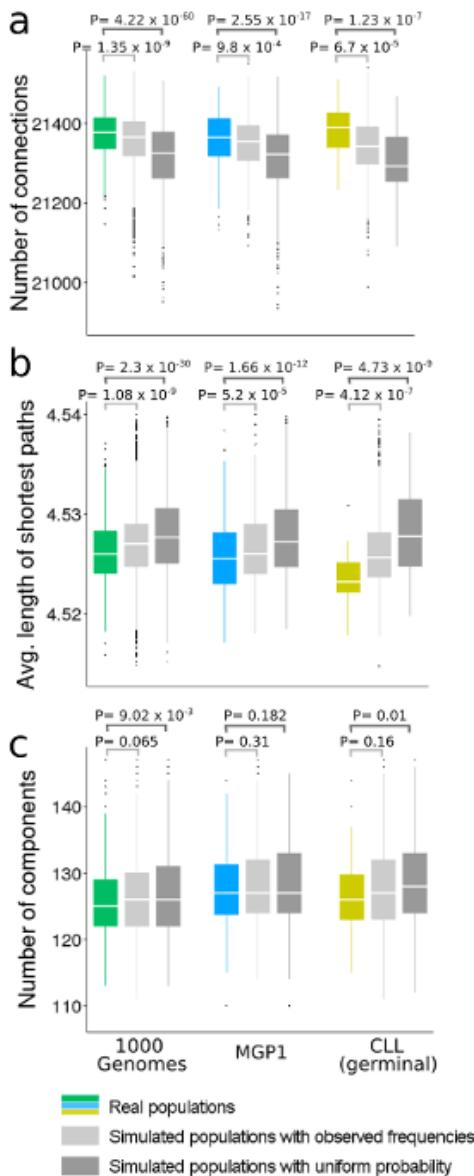
Cancer genes are central.

Hernandez, 2007 BMC Genomics



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

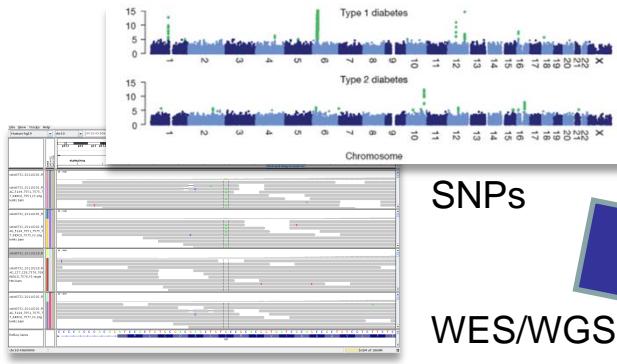
The role of interactome in buffering the deleteriousness of LoF mutations



Comparison of the interactome damage between real and random individuals after removing the nodes corresponding to proteins containing deleterious variants in both alleles (homozygote). Two different scenarios are simulated: Simulated populations with **uniform probability**, where proteins are randomly removed, and Simulated populations with **observed frequencies**, where proteins are removed with a probability proportional to the frequency of variation in the 1000 genomes population

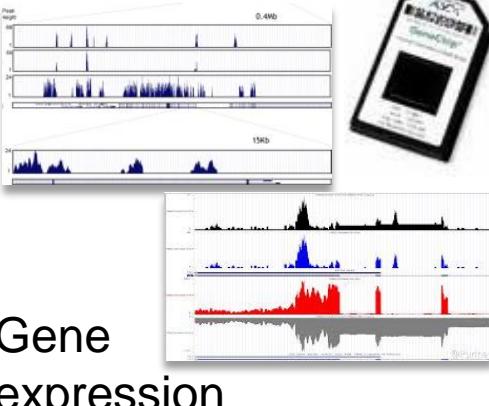
Garcia-Alonso 2014 Mol Syst Biol

From gene-based to function-based perspective

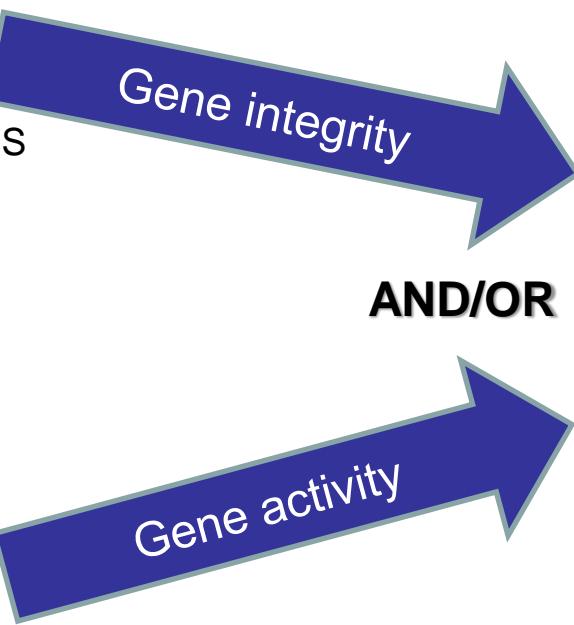


SNPs

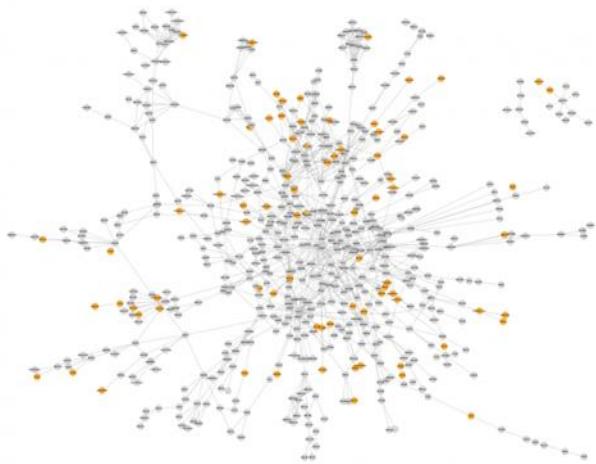
WES/WGS



Gene expression



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^{1,2}, Marta Bleda^{2,3}, Rocío Núñez-Torres^{1,2}, Ignacio Medina^{3,4}, Berta Luzón-Toro^{1,2}, Luz García-Alonso³, Ana Torroglosa^{1,2}, Martina Marbà^{3,4}, Ma Valle Enguix-Riego^{1,2}, David Montaner³, Guillermo Antíñolo^{1,2}, Joaquín Dopazo^{2,3,4*} and Salud Borrego^{1,2*}

* Corresponding authors: Joaquín Dopazo idopazo@cipf.es - Salud Borrego 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-Alonso¹, Roberto Alonso¹, Enrique Vidal¹, Alicia Amadoz¹, Alejandro de María¹, Pablo Minguez², Ignacio Medina^{1,3} and Joaquín Dopazo^{1,3,4,*}

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

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

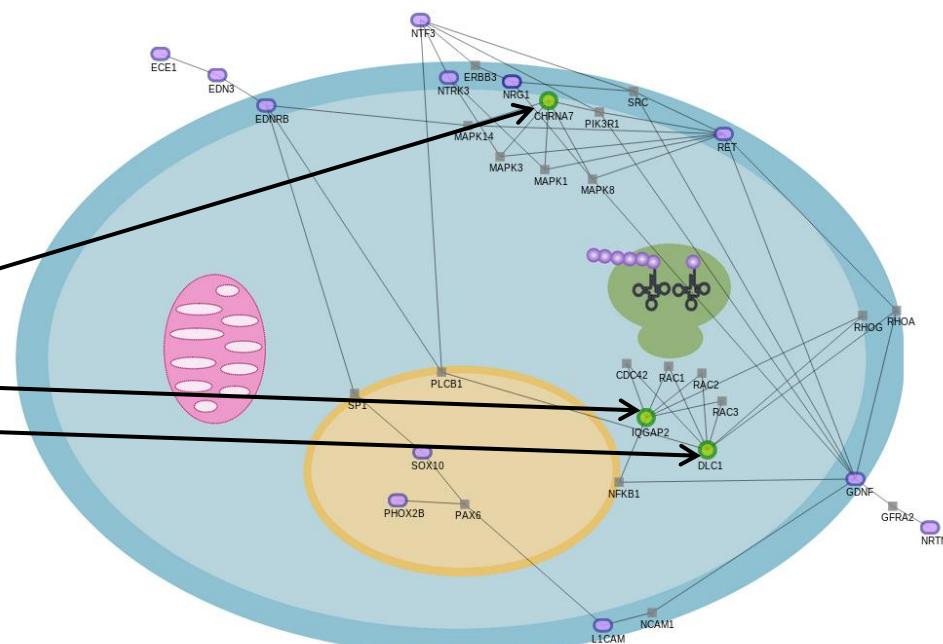
SNOW, a web-based tool for the statistical analysis of protein-protein interaction networks

Pablo Minguez¹, Stefan Götz^{1,2}, David Montaner¹, Fatima Al-Shahrour¹ and Joaquin Dopazo^{1,2,3,*}

¹Department of Bioinformatics and Genomics, Centro de Investigación Príncipe Felipe (CIPF),

²CIBER de Enfermedades Raras (CIBERER) and ³Functional Genomics Node (INB) at CIPF, Valencia, Spain

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



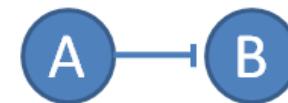
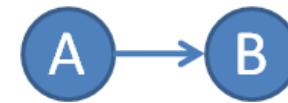
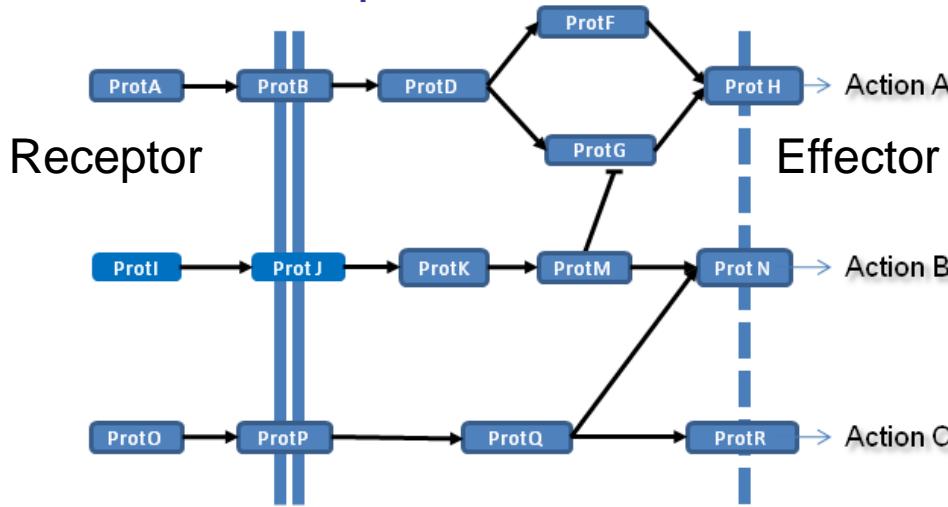
From gene-based to function-based perspective

	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

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



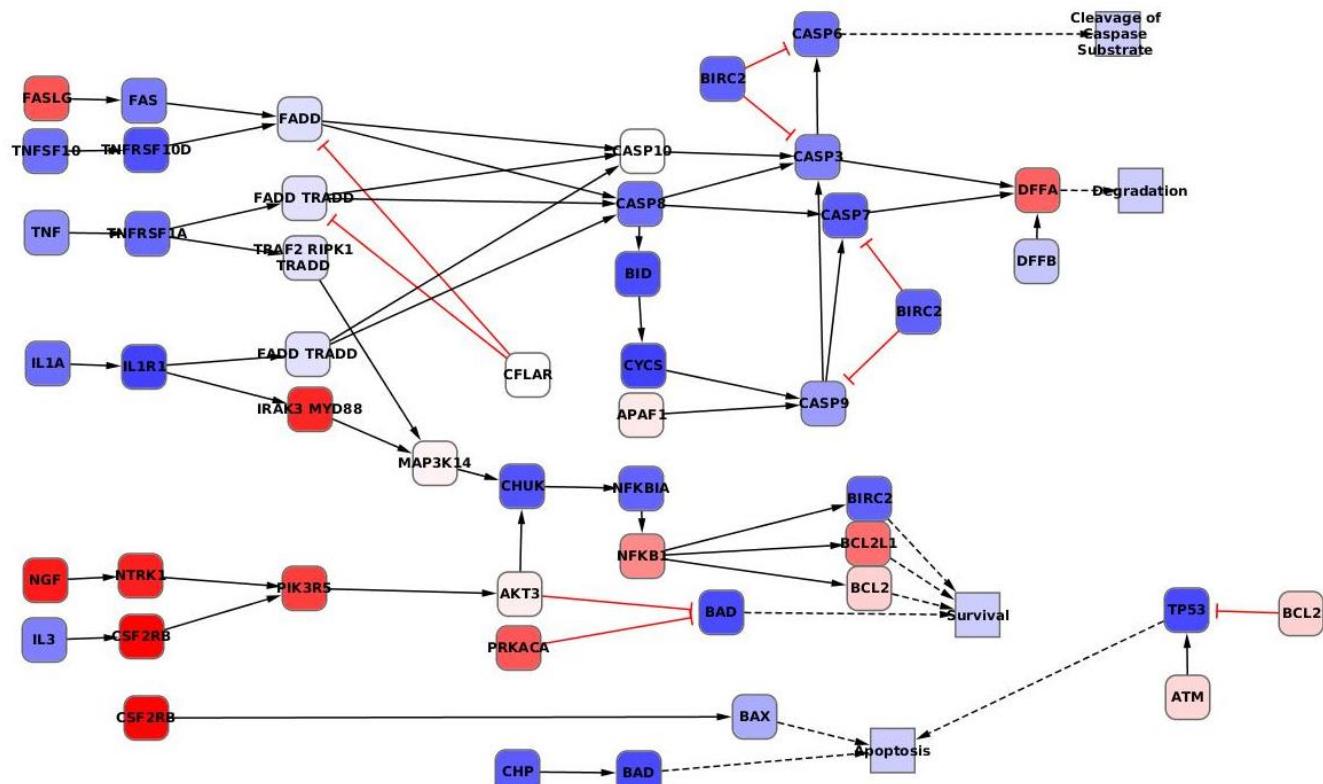
**Activations
and
repressions
occur**

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

The effects of changes in gene activity are not obvious

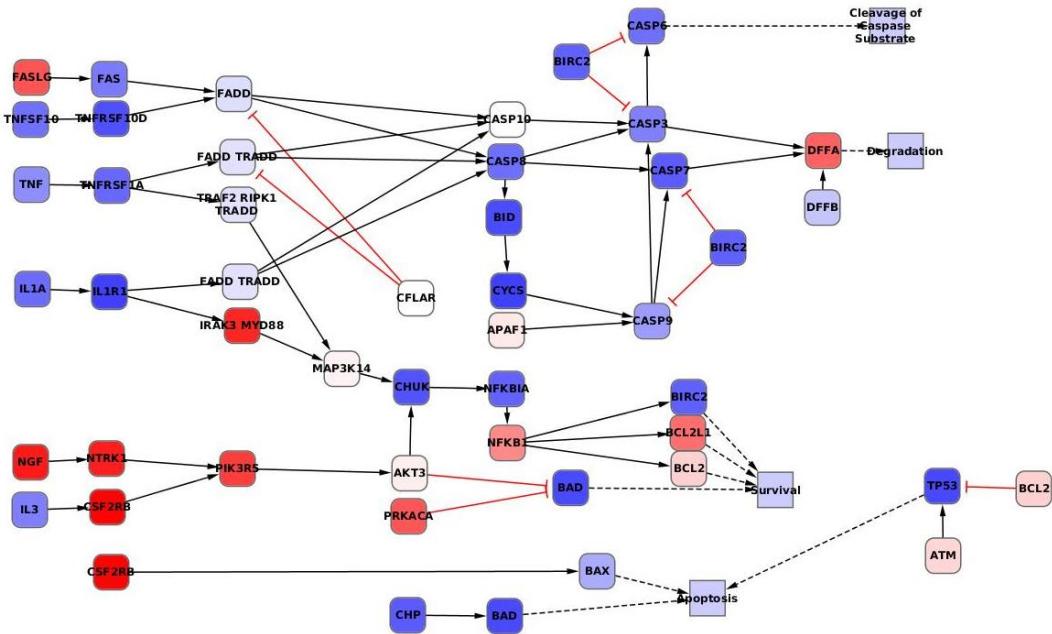
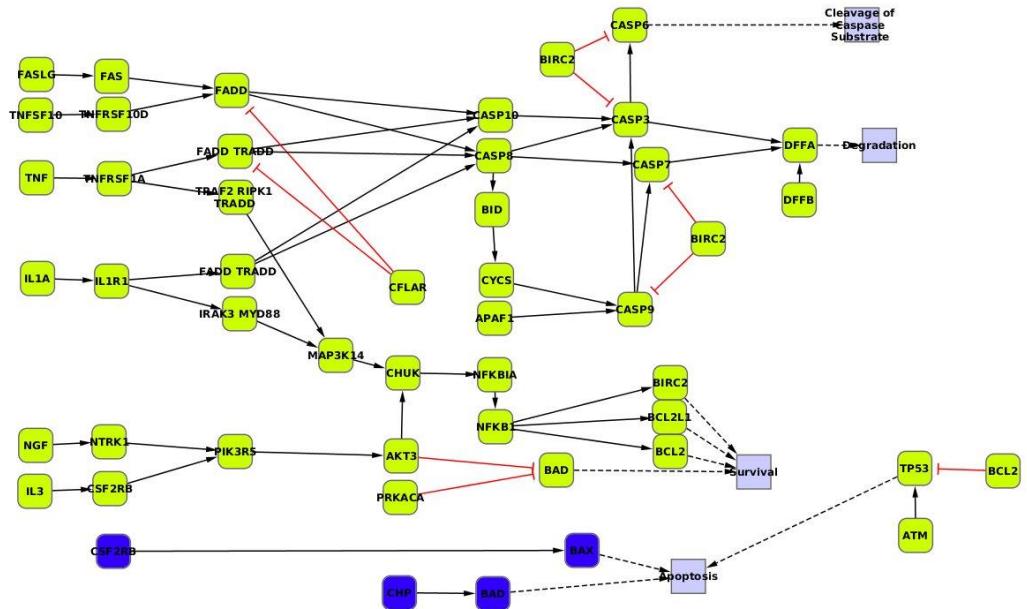
What would you predict about the consequences of gene activity changes in the apoptosis pathway in a case control experiment of colorectal cancer?

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

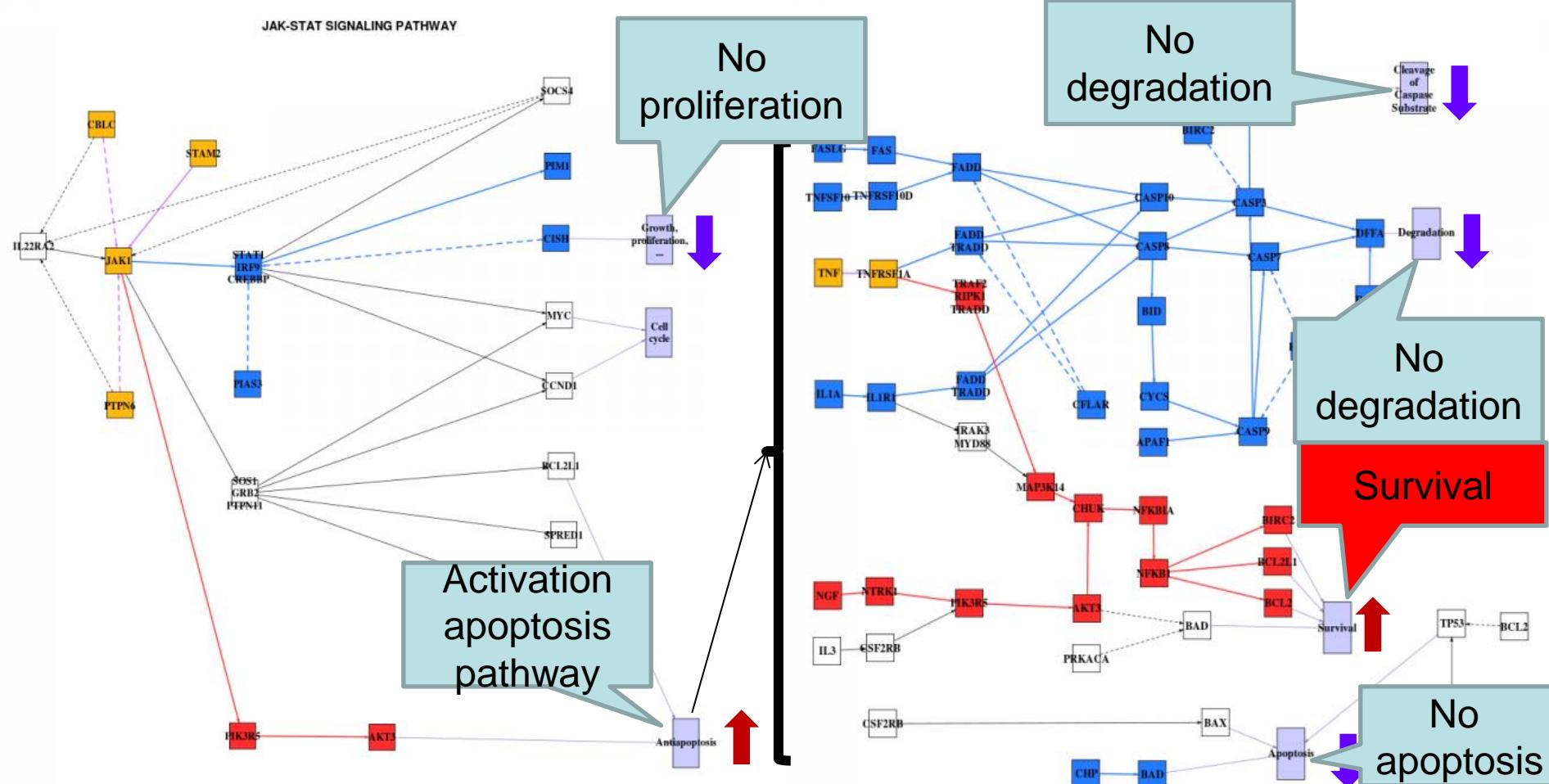


Apoptosis inhibition is not obvious from gene expression

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



Different pathways cross-talk to deregulate programmed death in Fanconi anemia



FA is a rare chromosome instability syndrome characterized by aplastic anemia and cancer and leukemia susceptibility. It has been proposed that disruption of the apoptotic control, a hallmark of FA, accounts for part of the phenotype of the disease.

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

Future prospects

Hospital Universitario La Paz

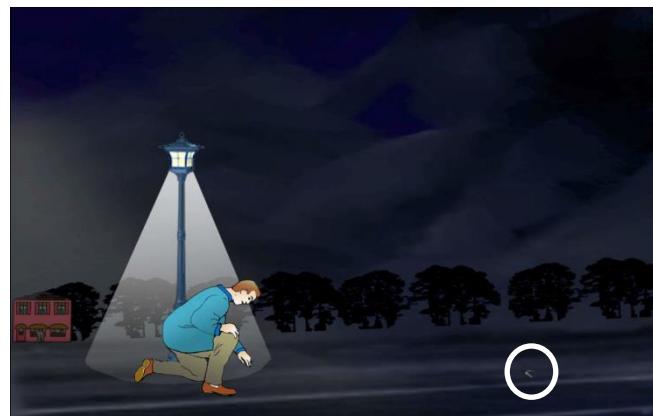
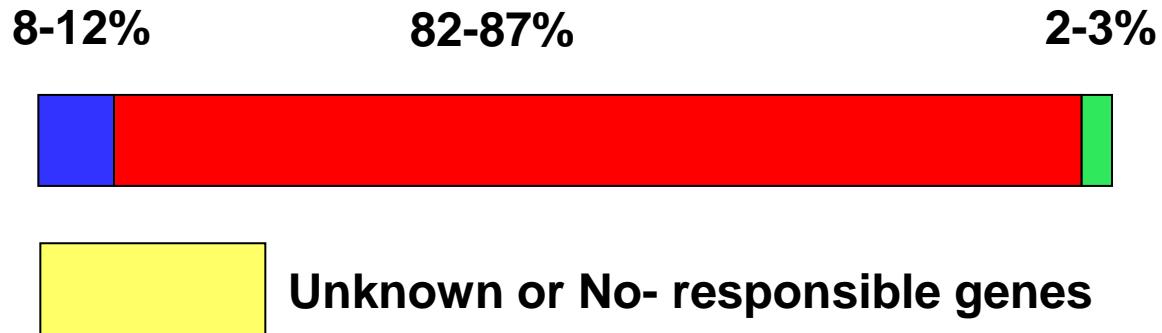
Known causes of Human Genetic Diseases



Pablo
Lapunzina,
Personal
communication

All genetic/genomic
or epigenetic
diseases with known
cause:
~ 5000 disorders

5 Kb- ? Mb 1 bp- 200 bp No dosage changes
GENOMICS **GENETICS** **EPIGENETICS**



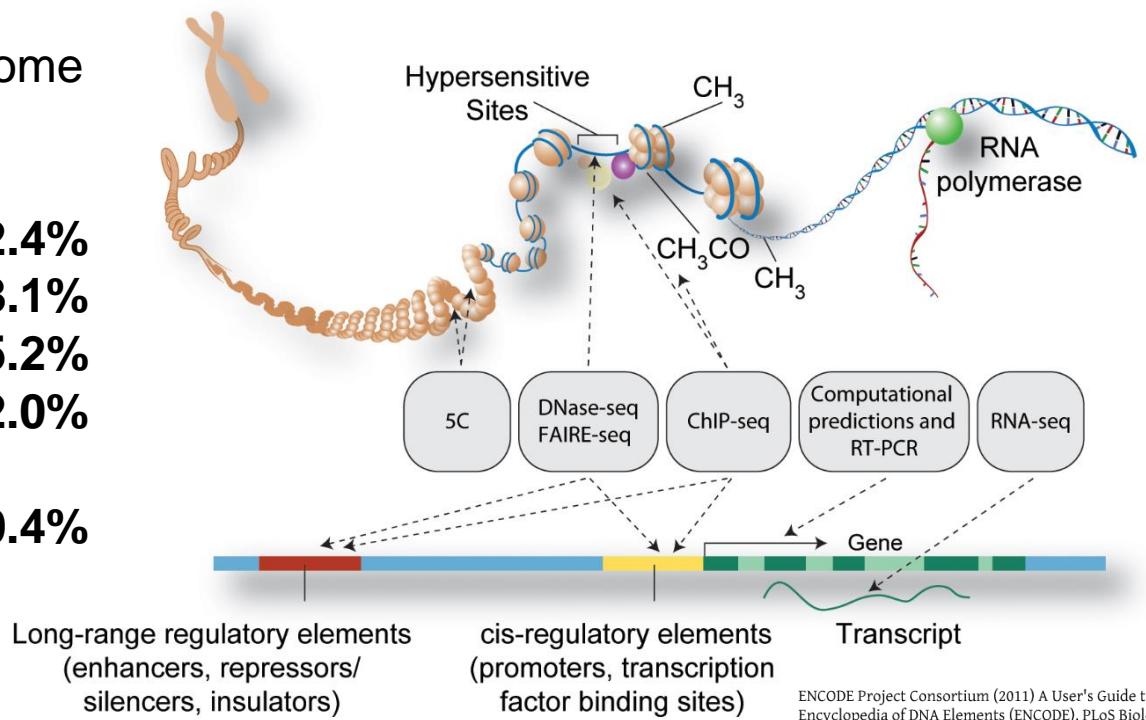
Fact: exons represent a comparatively small part of the complete genome
Other fact: there is still a lot of missing heritability

The ENCODE project suggests a functional role for a large fraction of the genome

Which percentage of the genome is occupied by:

Coding genes:	2.4%
TFBSs	8.1%
Open chromatin regions	15.2%
Different RNA types	62.0%

Total annotated elements: **80.4%**



ENCODE Project Consortium (2011) A User's Guide to the Encyclopedia of DNA Elements (ENCODE). PLoS Biology 9: e1000350.

Exomes are only covering a small fraction of the potential functionality of the genome (2.4%).

Is the **missing heritability** hidden in the remaining 78%?

If so, what type of variant should be expect to discover? SNVs? SVs?

Future prospects

We need to efficiently query all the information contained in the genome, including all the epigenomic signatures.

This means **data integration** and “**epistatic**” queries

We need to prepare our **health systems** to deal with all the genomic data flood

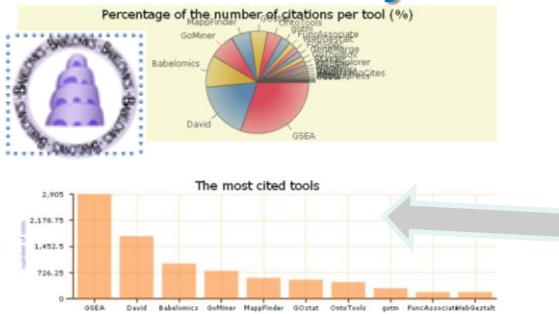
Information about variations	Processed	Raw
Genome variant information (VCF)	150 MB	250 GB
Epigenome	150 MB	250 GB
Each transcriptome	20 MB	80 GB
Individual complete variability	400 MB	525 GB
Hospital (100.000 patients)	40 TB	50 PB

There are **technical** problems and **conceptual** problems on how genomic information is managed that must be addressed in the near future.

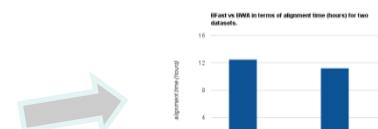


Software development

Functional analysis



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



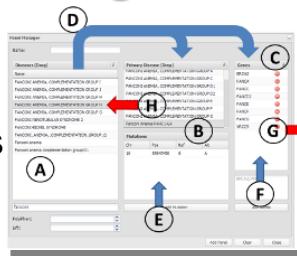
Mapping

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

Visualization

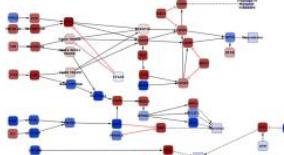


Diagnostic

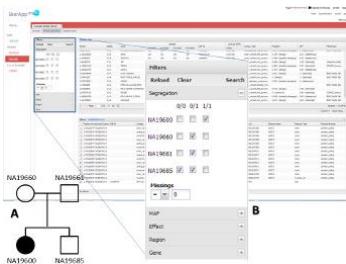


NGS
panels

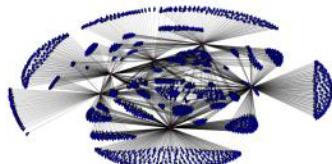
Signaling network



Variant prioritization



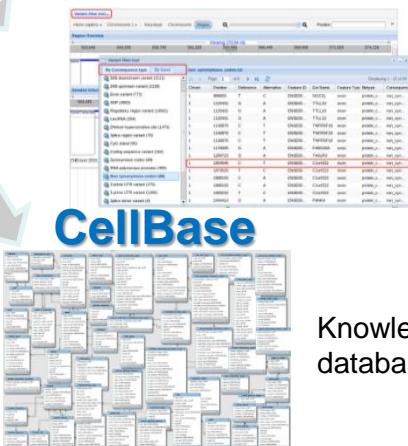
Regulatory network



Interaction network



Variant annotation



Knowledge
database

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