

Introduction



GDA
International Course on
Genomic Data Analysis



PRINCIPE FELIPE
CENTRO DE INVESTIGACION

Who we are

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)

<http://bioinfo.cipf.es>

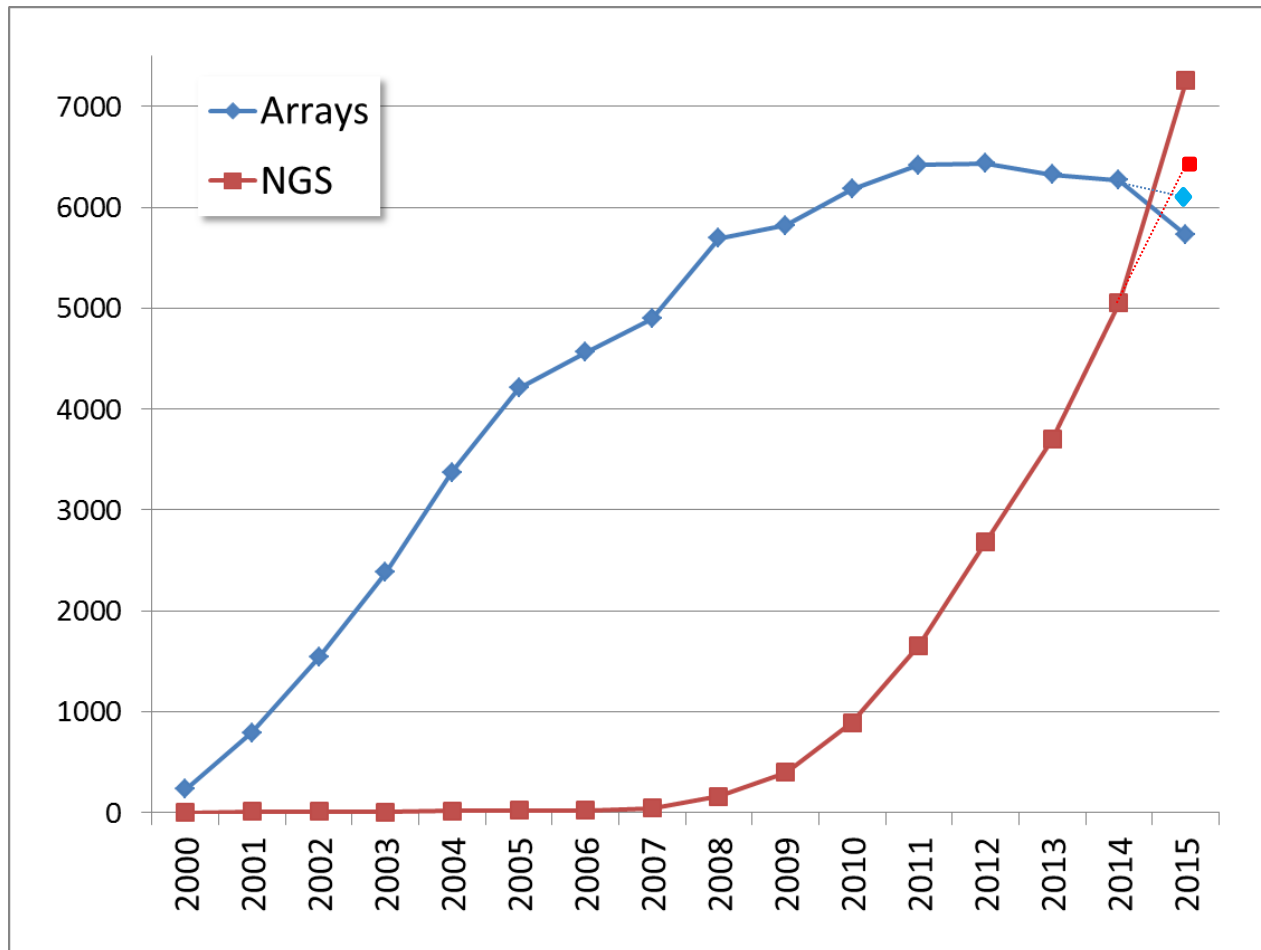
<http://www.babelomics.org>



@xdopazo @bioinfocipf



Papers published in microarrays and NGS technologies

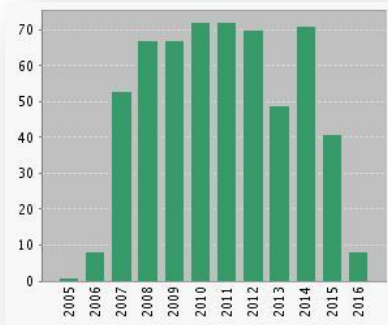


Source Pubmed. Query: "high-throughput sequencing"[Title/Abstract] OR "next generation sequencing"[Title/Abstract] OR "rna seq"[Title/Abstract] AND year[Publication Date]

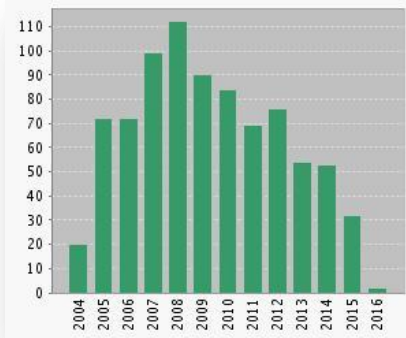
Tools used in the course

579 papers cite Babelomics (plus 835 FatiGO cites)

(source ISI Web of Knowledge, February 2016)



Babelomics



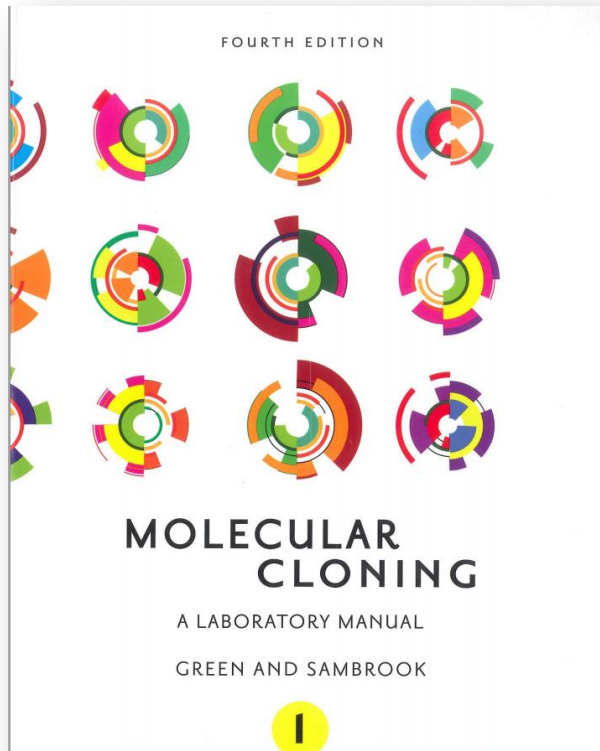
FatiGO



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

Babelomics in the textbooks

The Babelomics suite of programs becomes a classic. Now is cited as a method in the last edition of **Molecular Cloning**. 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

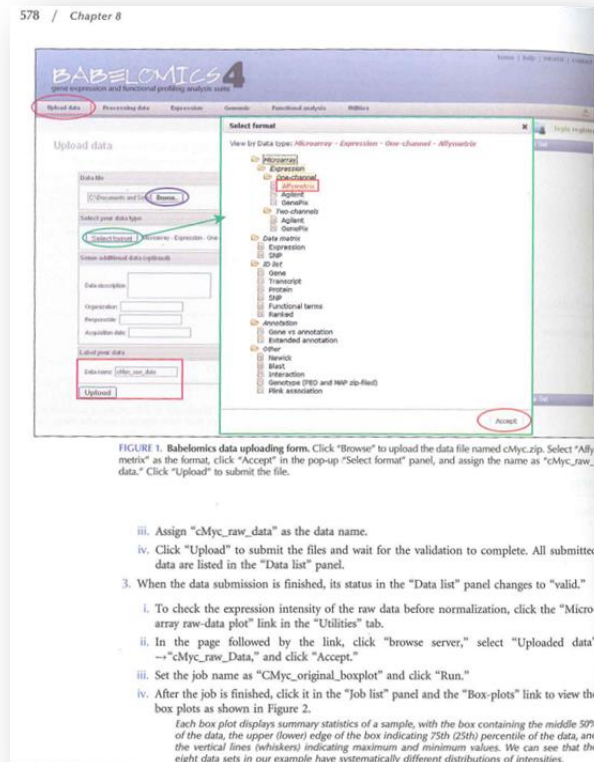


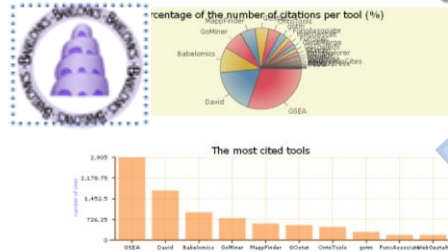
FIGURE 1. Babelomics data uploading form. Click "Browse" to upload the data file named cMyc.zip. Select "Affymetrix" as the format, click "Accept" in the pop-up "Select format" panel, and assign the name as "cMyc_raw_data." Click "Upload" to submit the file.

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

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

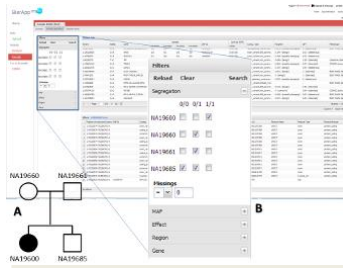
Software development

Functional analysis

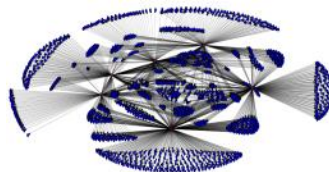


See interactive map of for the last 24h use <http://bioinfo.cipf.es/toolsusage>

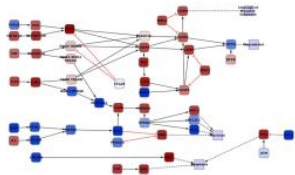
Variant prioritization



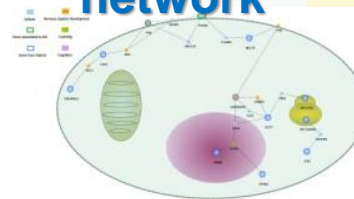
Regulatory network



Signaling network



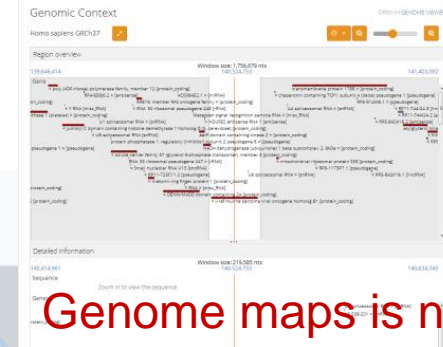
Interaction network



Mapping

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

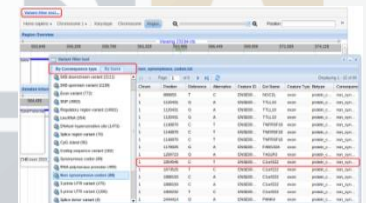
Visualization



Ultrafast
genome
viewer with
google
technology

Genome maps is now part
of the ICGC data portal

Variant annotation



CellBase



Knowledge
database

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

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.

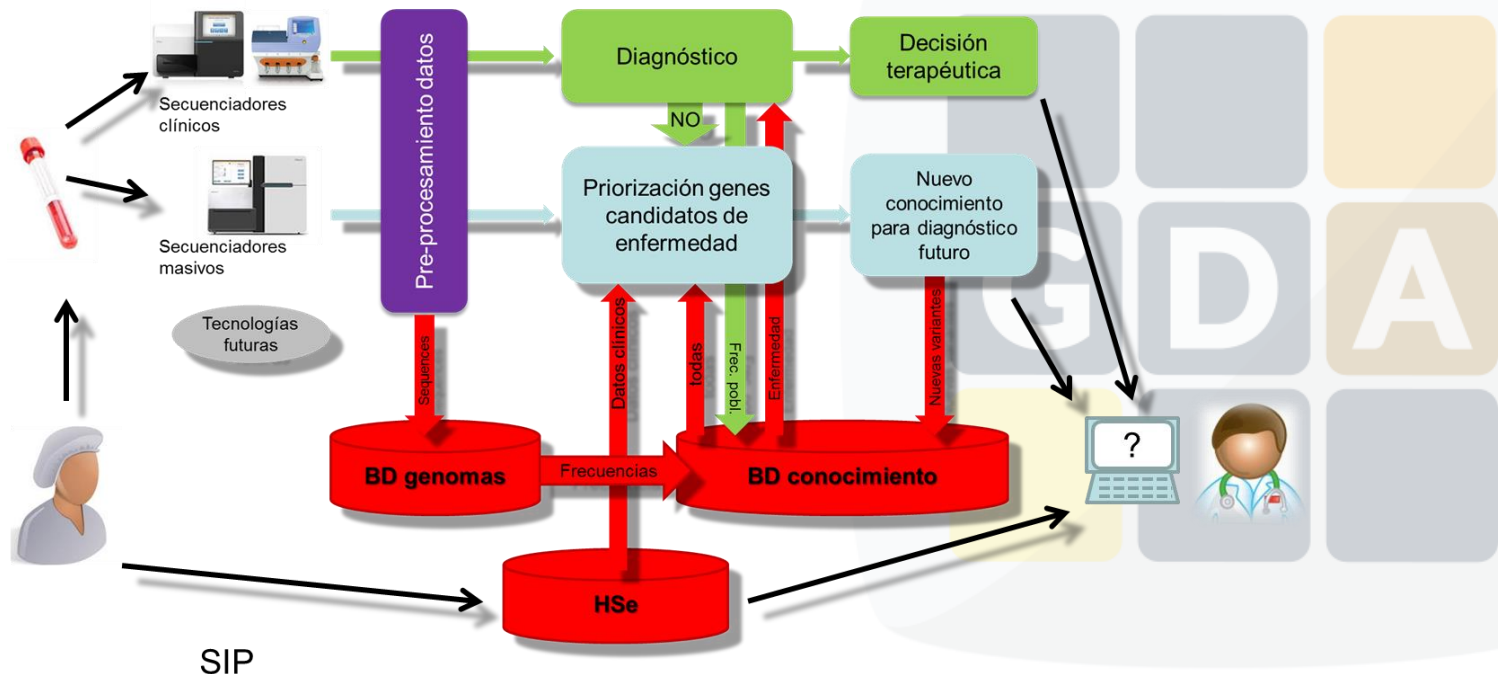
But not necessarily the way in which we really address or test them...

Here you will learn how to do so using state-of-the-art methods and software.

NGS genomic data analysis

Two main applications:

- Diagnostic
- Disease gene finding

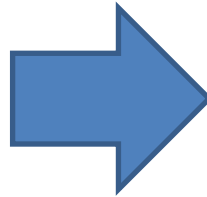


Disease gene finding

Typically, an exome renders between 40 and 60K variants (and a genome about 1 million). Only one or a few among all of them are expected to be the causative factors of the disease.



Casablanca: Round up all suspicious characters and search them for stolen documents.



40-60K suspects?
We need...



Agatha Christie

The prioritization process is like a police investigation in which suspected are discarded by their alibies

Thus, through sequential heuristic filtering steps, unlikely candidates are discarded and a final, reduced list with one or a few candidate genes is (hopefully) produced

Pipeline of data analysis



Primary processing

Initial QC
FASTQ file

Mapping
BAM file

Variant calling
VCF File

Secondary analysis (Heuristic filtering)

Variant annotation

Filtering by effect

Filtering by MAF

Filtering by family
segregation

Knowledge-based prioritization

Proximity to other
known disease genes

Functional proximity

Network proximity

Burden tests

Other prioritization
methods



Primary
analysis

Gene prioritization

Pipeline of data analysis



Primary processing

Initial QC
FASTQ file

Mapping
BAM file

Variant calling
VCF File

Secondary analysis (Heuristic filtering)

Variant annotation

Filtering by effect

Filtering by MAF

Filtering by family
segregation

Knowledge-based prioritization

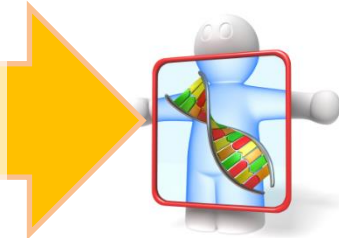
Proximity to other
known disease genes

Functional proximity

Network proximity

Burden tests

Other prioritization
methods



Primary
analysis

Gene prioritization

HPG suite

High-Performance Genomics

More info at:

<http://bioinfo.cipf.es/docs/compbio/projects/hpg/doku.php>

**Other analysis
(HPC4Genomics
consortium)**

- RNA-seq
- DNA assembly
- Methyl-seq
- Copy Number
- Structural variation
- Transcript isoform
- Etc.



Fastq file, up to hundreds of GB per run

QC and preprocessing

QC stats, filtering and preprocessing options

HPG Aligner, short read aligner

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



SAM/BAM file

QC and preprocessing

QC stats, filtering and preprocessing options

Variant calling analysis

GATK and SAM mPileup HPC
Implementation.
Statistics genomic tests



VCF file

QC and preprocessing

QC stats, filtering and preprocessing options

Variant VCF viewer

HTML5+SVG Web based viewer



HPG Variant, Variant analysis

Consequence type, GWAS, regulatory
variants and system biology information

Pipeline of data analysis



Primary processing

Initial QC
FASTQ file

Mapping
BAM file

Variant calling
VCF File

Secondary analysis (Heuristic filtering)

Variant annotation

Filtering by effect

Filtering by MAF

Filtering by family
segregation

Knowledge-based prioritization

Proximity to other
known disease genes

Functional proximity

Network proximity

Burden tests

Other prioritization
methods

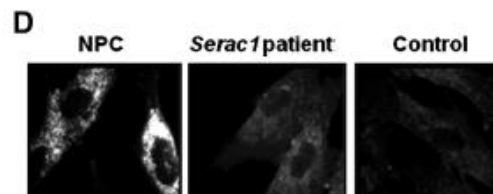
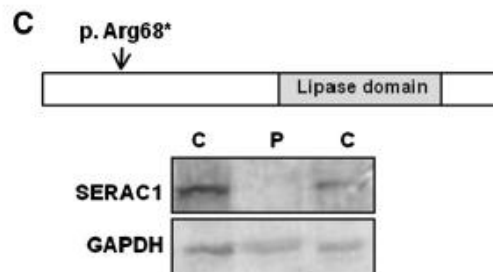
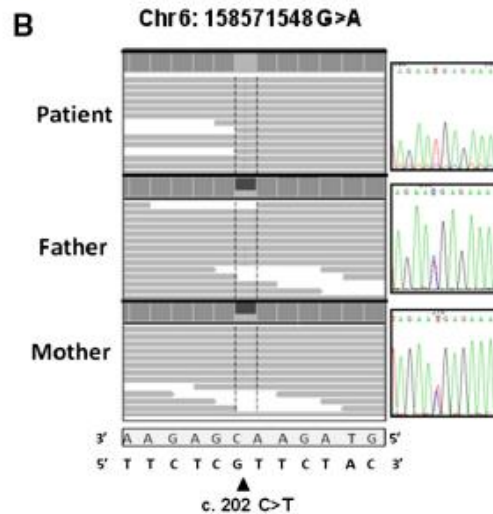
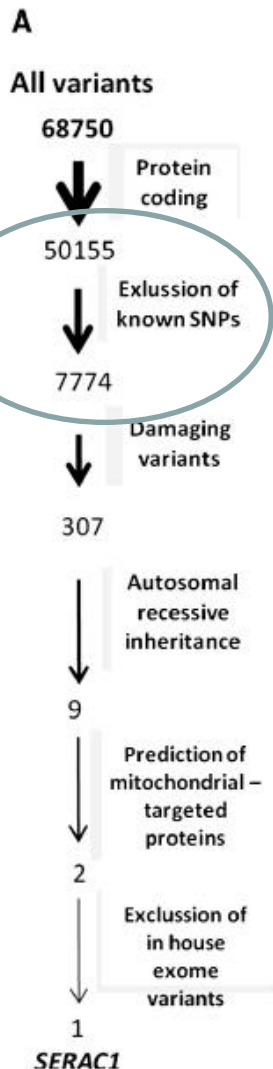


Primary
analysis

Gene prioritization

Heuristic filtering

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



3-Methylglutaconic aciduria (3-MGA-uria) is a heterogeneous group of syndromes characterized by an increased excretion of 3-methylglutaconic and 3-methylglutaric acids.

WES with a consecutive filter approach is enough to detect the new mutation in this case.

Contents lists available at SciVerse ScienceDirect

Molecular Genetics and Metabolism

ELSEVIER

journal homepage: www.elsevier.com/locate/ymgme

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}, Maria 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 Mitocondriales-Enfermedades Metabólicas Hereditarias, Servicio de Pediatría, Hospital 12 de Octubre, Madrid, Spain

^d IBER, CIBERER, Centro de Investigación Príncipe Felipe, Valencia, Spain

^e Computational Medicine Institute, Centro de Investigación Príncipe Felipe (CIPIF), Valencia, Spain

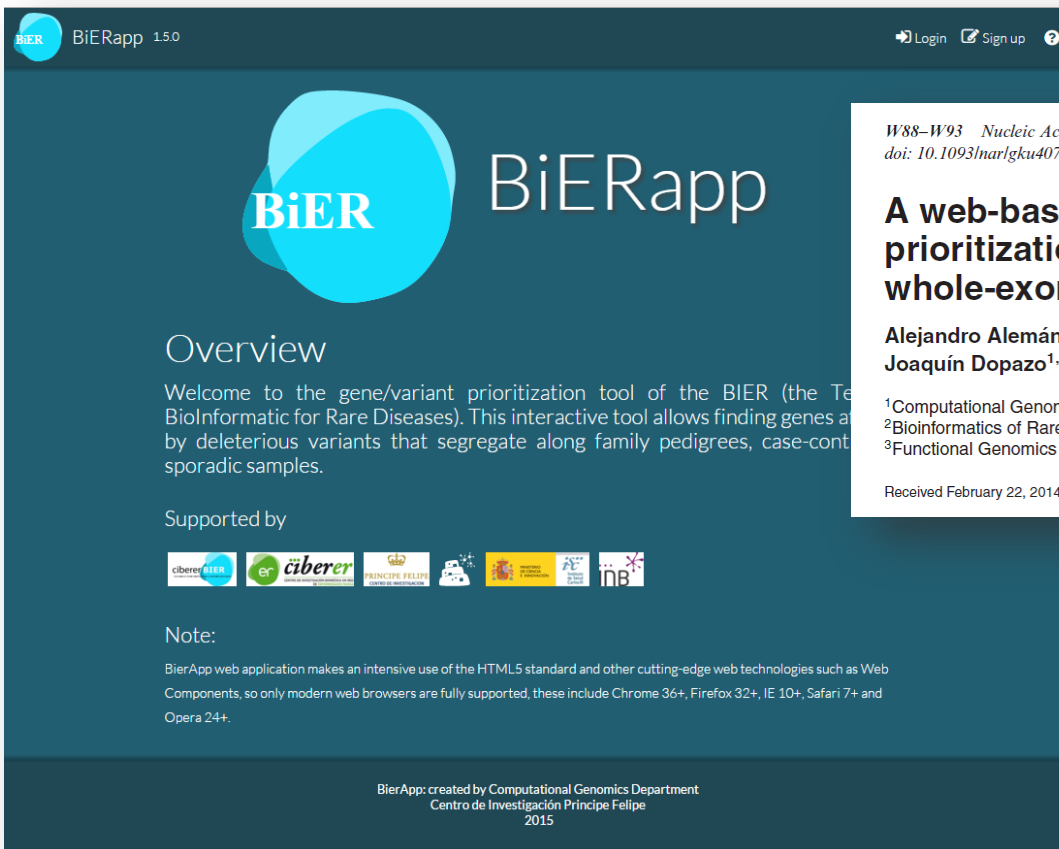
^f Functional Genomics Node, (INB) at CIPIF, Valencia, Spain

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

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

* Corresponding author. Tel.: +34 91 410 85 00; fax: +34 91 410 85 00. E-mail address: anton.ribe@ciqpf.es (A. Ribes).

The BiERapp



BiERapp 1.5.0

Login Sign up

BiERapp

Overview

Welcome to the gene/variant prioritization tool of the BIER (the Bioinformatics for Rare Diseases). This interactive tool allows finding genes associated by deleterious variants that segregate along family pedigrees, case-control studies, and sporadic samples.

Supported by

Note:

BierApp web application makes an intensive use of the HTML5 standard and other cutting-edge web technologies such as Web Components, so only modern web browsers are fully supported, these include Chrome 36+, Firefox 32+, IE 10+, Safari 7+ and Opera 24+.

BierApp: created by Computational Genomics Department
Centro de Investigación Príncipe Felipe
2015

W88–W93 *Nucleic Acids Research*, 2014, Vol. 12, Web Server issue
doi: 10.1093/nar/gku407

Published online 06 May 2014

A web-based interactive framework to assist in the prioritization of disease candidate genes in whole-exome sequencing studies

Alejandro Alemán^{1,2}, Francisco Garcia-Garcia¹, Francisco Salavert^{1,2}, Ignacio Medina¹ and Joaquín Dopazo^{1,2,3,*}

¹Computational Genomics Department, Centro de Investigación Príncipe Felipe (CIPF), Valencia 46012, Spain,

²Bioinformatics of Rare Diseases (BIER), CIBER de Enfermedades Raras (CIBERER), Valencia 46010, Spain and

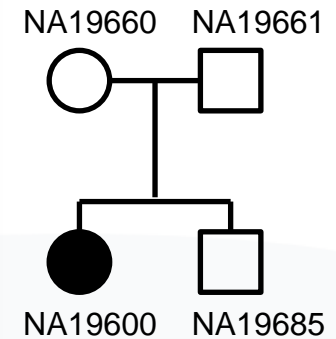
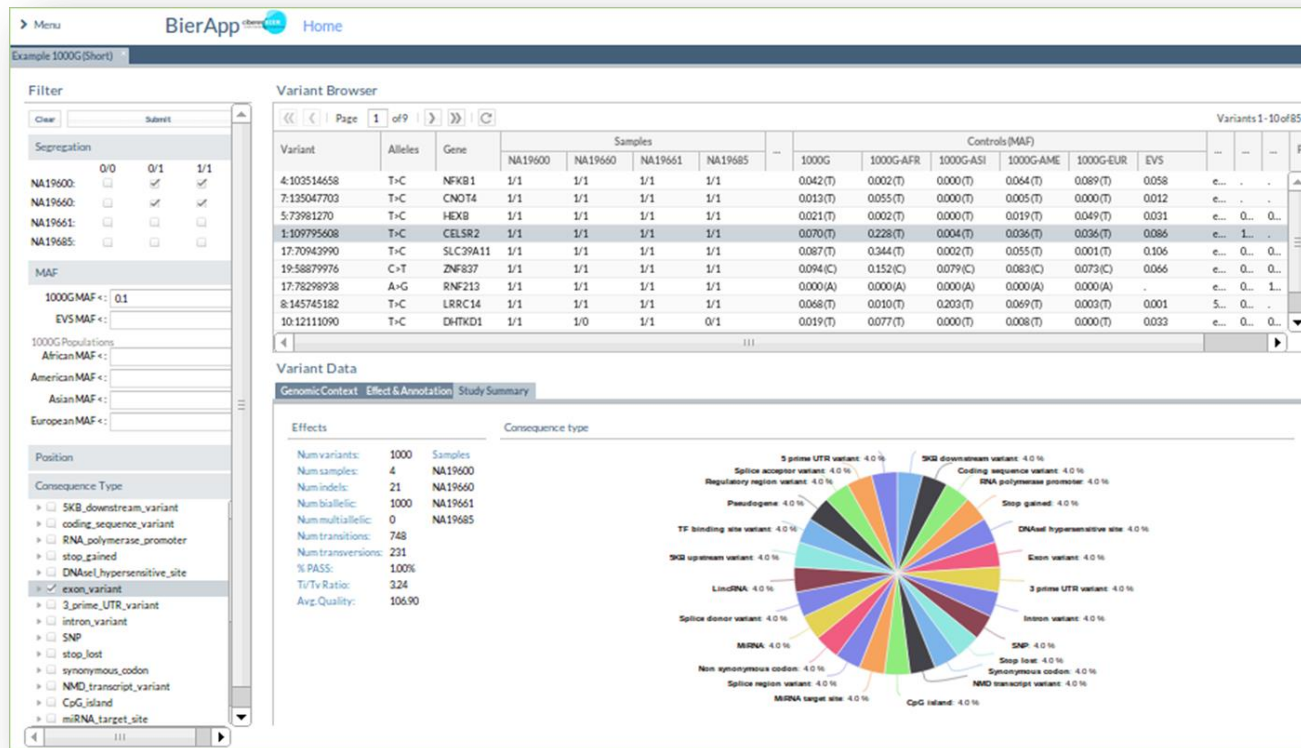
³Functional Genomics Node, (INB) at CIPF, Valencia 46012, Spain

Received February 22, 2014; Revised April 27, 2014; Accepted April 28, 2014

An interactive web tool
that implements different
heuristic filters for disease
variant/gene prioritization



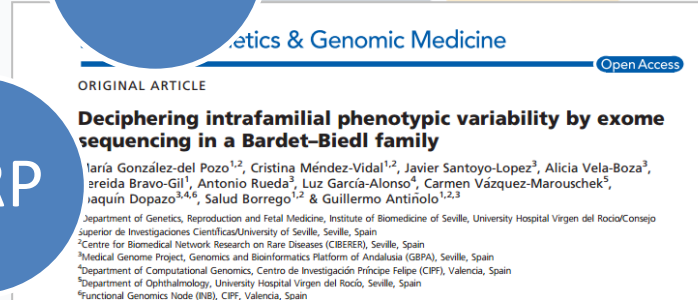
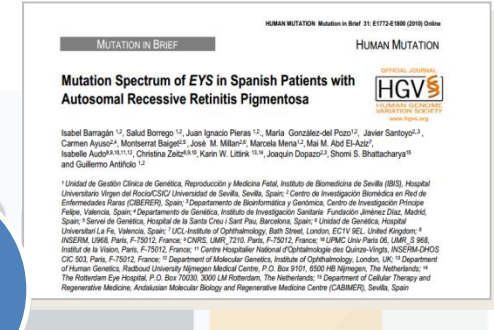
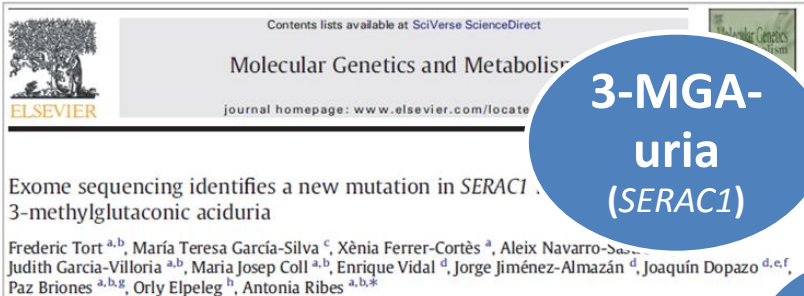
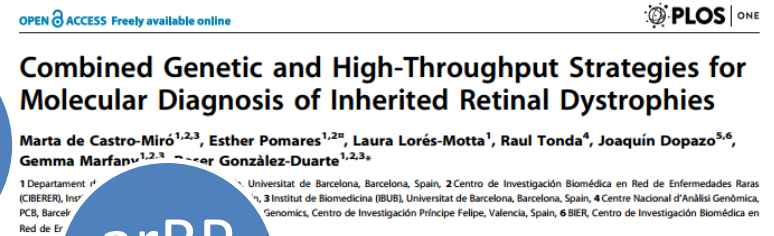
The BiERapp



| Filters | | | |
|-----------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Reload | Clear | Search | |
| Segregation | | | <input type="text" value="-"/> |
| 0/0 0/1 1/1 | | | |
| NA19600 | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| NA19660 | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| NA19661 | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| NA19685 | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| Missings | | | |
| = | <input type="text" value="v"/> | <input type="text" value="0"/> | |
| MAF | | | <input type="text" value="+"/> |
| Effect | | | <input type="text" value="+"/> |
| Region | | | <input type="text" value="+"/> |
| Gene | | | <input type="text" value="+"/> |

Filters include family pedigree segregation, population frequencies, pathogenic indexes, etc.

Successful use of BiERapp



Pipeline of data analysis



Primary processing

Initial QC
FASTQ file

Mapping
BAM file

Variant calling
VCF File

Secondary analysis (Heuristic filtering)

Variant annotation

Filtering by effect

Filtering by MAF

Filtering by family
segregation

Knowledge-based prioritization

Proximity to other
known disease genes

Functional proximity

Network proximity

Burden tests

Other prioritization
methods



Primary
analysis

Gene prioritization

Knowledge-based prioritization

Network analysis

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^{2,4}, 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 Antiñolo^{1,2}, Joaquín Dopazo^{2,3,4,*} and Salud Borrego^{1,2,*}

* Corresponding authors: Joaquín Dopazo jdopazo@cipf.es - Salud Borrego salud.borrego.sspa@iuntadeandalucia.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

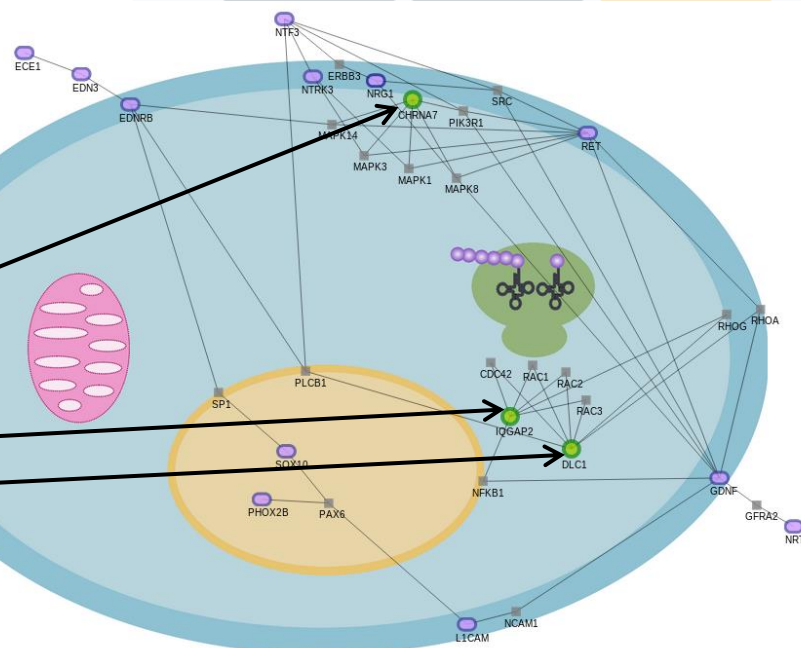
Nucleic Acids Research, 2009, 1–6
doi:10.1093/nar/gkp402

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 Joaquín 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



NGS for diagnostic

Published online 26 May 2014

Nucleic Acids Research, 2014, Vol. 42, Web Server issue W83-W87
doi: 10.1093/nar/gku472

A web tool for the design and management of panels of genes for targeted enrichment and massive sequencing for clinical applications

Alejandro Alemán^{1,2}, Francisco Garcia-García¹, Ignacio Medina¹ and Joaquín Dopazo^{1,2,3,*}

¹Computational Genomics Department, Centro de Investigación Príncipe Felipe (CIPF), Valencia, 46012, Spain,
²Bioinformatics of Rare Diseases (BIER), CIBER de Enfermedades Raras (CIBERER), Valencia, 46012, Spain and
³Functional Genomics Node, (INB) at CIPF, Valencia, 46012, Spain

Received March 1, 2014; Revised May 01, 2014; Accepted May 13, 2014

TEAM

Example Data

Search

Panel:

VCF File:

Results

Diagnostic Secondary findings

| | Chromosome | Position | SNP Id | Ref | Alt | Gene | Conseq. Type | Phenotype | Source | SIFT | PolyPhen |
|---------------------|------------|-----------|--------|-----|-----|------|------------------------|-------------------------|---------------|------|----------|
| gene: (1 Item) | | | | | | | | | | | |
| 1 | 3 | 129247734 | . | T | C | . | exon_variant, codin... | RETINITIS PIGMENT... | dbSNP_ClinVar | . | . |
| gene: RHO (3 Items) | | | | | | | | | | | |
| 2 | 3 | 129247734 | . | T | C | RHO | exon_variant, codin... | RETINITIS PIGMENT... | OMIM | . | . |
| 3 | 3 | 129247734 | . | T | C | RHO | exon_variant, codin... | RETINITIS PIGMENT... | Unprot | . | . |
| 4 | 3 | 129247734 | . | T | C | RHO | exon_variant, codin... | Retinitis pigmentosa... | Unprot | . | . |

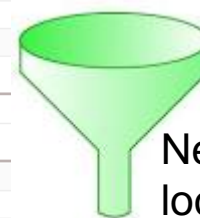
Diagnostic mutations

If no diagnostic variants appear, then secondary findings are studied

Results

Diagnostic Secondary findings

| | Chromosome | Position | SNP Id | Ref | Alt | Gene | Conseq. Type | Phenotype | Source | SIFT | PolyPhen |
|------------------------|------------|-----------|--------|-----|-----|-------|------------------------|-----------|--------|------|----------|
| 5 | 2 | 182413259 | . | A | G | CERKL | intron_variant, NMD... | . | . | . | . |
| 6 | 2 | 182413602 | . | A | T | CERKL | intron_variant, NMD... | . | . | . | . |
| 7 | 2 | 182521578 | . | G | A | CERKL | intron_variant, NMD... | . | . | . | . |
| 8 | 2 | 182543455 | . | T | C | CERKL | intron_variant, 5KB... | . | . | 0.76 | 0.003 |
| gene: CNGA1 (1 Item) | | | | | | | | | | | |
| 9 | 4 | 47953515 | . | A | T | CNGA1 | intron_variant, SNP | . | . | . | . |
| gene: CNGB1 (12 Items) | | | | | | | | | | | |
| 10 | 16 | 57937788 | . | T | C | CNGB1 | 5KB_upstream_vari... | . | . | 0.74 | 0.002 |
| 11 | 16 | 57937856 | . | G | C | CNGB1 | 5KB_upstream_vari... | . | . | . | . |
| 12 | 16 | 57937895 | . | G | A | CNGB1 | intron_variant, 5KB... | . | . | . | . |
| 13 | 16 | 57949251 | . | G | A | CNGB1 | intron_variant, 5KB... | . | . | . | . |



New filter based on local population variant frequencies

documentation tutorial

Show P

Panel Manager

Name:

Disorders (Drag)

Primary Disease (Drop)

Genes

Panel Manager interface showing a list of disorders and genes, with a search bar and a 'Generate Report' button.

Tool for defining panels

Uses of BiERapp and TEAM

Implementation of tools in the IT4I Supercomputing Center (Czech Republic)

The pipelines of primary and secondary analysis developed by the Computational Genomics Department of the CIPF in close collaboration with the Bull Chair of computational genomics has proven its efficiency in the analysis of more than 1000 exomes in a joint collaborative project of the CIBERER and the MGP. A first pilot implementation has been done in the IT4I supercomputing center, which aims to centralize the analysis of genomics data in the country.

IT4Innovations
národní
superpočítačové
centrum

BULL

OMICS MASTER 24.7.2014

PI 332 Budova Podnikatelského inkubátoru v areálu VŠB – Technické univerzity Ostrava 9:00–16:00

Obsahem kurzu je praktické seznámení se sadou nástrojů OMICS MASTER, které budou nasazeny v IT4Innovations národním superpočítačovém centru. Cílem instalace je vytvoření standardního prostředí, nástrojů a postupů pipeline pro analýzu dat – NGS (Next Generation Sequencing). Zaměřením připravované pipeline je především řešení genomický výzkum, zejména pak sestavení genomu, identifikace genů a anotace jejich variant pro diagnostické účely.

Seznámení s pipeline proběhne prostřednictvím případové studie. První část bude zaměřena na analýzu reálného vzorku primárních dat. Druhá část se soustředí na práci s diagnostickými moduly. Kurz přináší zručnosti potřebné pro efektivní zpracování a analýzu NGS dat na infrastruktuře IT4Innovations, od primárních dat až po anotaci variant včetně identifikace jejich biologické significance a diagnostické hodnoty.

Závěrečná pipeline je postavena na otevřeném softwaru a zabezpečí efektivní provedení výpočetně intenzivních úloh spojených se zpracováním primárních dat ze NGS přístrojů. Mezi výstupy analýzy je přímo VCF formát a nástroje pro jeho další analýzu.

Pipeline je postavená na těchto nástrojích:

- FASTQC
- HPG-aligner
- Samtools
- Picard
- GATK
- hpg-variant
- variant (annotation)
- OpenCGA indexer
- BiERapp (for gene discovery)
- TEAM (for diagnosis)

Diagram illustrating the OMICS MASTER pipeline workflow:

- Raw data (FASTQ) is processed by FASTQC and HPG-aligner.
- The aligned data is stored in the Sequence DB.
- The Sequence DB feeds into the Diagnostic module and the Discovery module.
- The Discovery module feeds into the Knowledge DB.
- The Knowledge DB feeds into the Diagnostic module.
- The Diagnostic module outputs results to the Discovery module.

Na seminář se můžete přihlásit do pátku 18. července 2014 na emailové adrese zuzana.kosarikova@vstb.cz

WWW.IT4I.CZ

Evropská unie
Investice do budoucnosti
OP Vzdělávání a odborná příprava
www.bull.cz

Course schedule

GDA2016: 29 Feb - 4 Mar

| | Monday 29 | Tuesday 1 | Wednesday 2 | Thursday 3 | Friday 4 |
|-------------|--|---|--|--|-----------------------------------|
| 9:00-9:30 | Registration | | | | |
| 9:30-10:00 | Course presentation | Biological and Clinical Databases. CSVS | Panel of genes: design and analysis for clinical applications. TEAM | Differential Expression Analysis. Babelomics 5 | Pathways Analysis |
| 10:00-10:30 | Introduction to NGS Technologies for Genomic Analysis | | | | |
| 10:30-11:00 | | | | | |
| 11:00-11:30 | Coffee break | | | | |
| 11:30-12:00 | Introduction to linux | Prioritization of variants and genes: BiERapp | Panel of genes: design and analysis for clinical applications. TEAM | Clustering and Supervised Classification Analysis. Babelomics 5 | Pathways Analysis |
| 12:00-12:30 | Primary Analysis: Quality control for raw data | | | | |
| 12:30-13:00 | | | Closing | | |
| 13:00-13:30 | Lunch | | | | |
| 13:30-14:00 | | | | | |
| 14:00-14:30 | | | | | |
| 14:30-15:00 | Primary Analysis: Mapping NGS Reads and visualization for Genomics Studies | Prioritization of variants and genes: BiERapp and Network tools | Course presentation | Functional Profiling from Babelomics 5: FatiGO and Gene Set Analysis | Enrichment Analysis for microRNAs |
| 15:00-15:30 | | | Introduction to NGS Technologies for Expression Analysis | | |
| 15:30-16:00 | | | | | |
| 16:00-16:30 | Coffee break | | | | |
| 16:30-17:00 | Primary Analysis: Variant Calling SNPs and INDELs. Variant Annotation | Bring your own data! | Primary Analysis for RNA-Seq data: quality control, mapping and quantification | Functional Profiling from Babelomics 5: Network Analysis | Bring your own data! |
| 17:00-17:30 | | | | | Closing |
| 17:30-18:00 | | | | | |
| 21:00 | | SOCIAL DINNER | | SOCIAL DINNER | |

GENOMICS: 29 February, 1 and 2 March

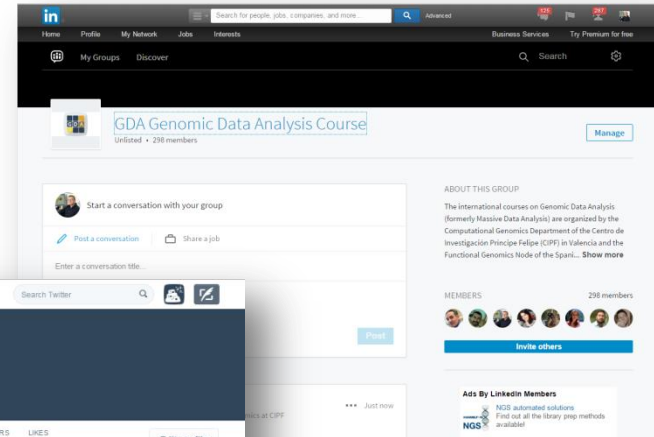
TRANSCRIPTOMICS: 2, 3 and 4 March

Social

GDA group in Linked-in



<https://www.linkedin.com/groups/1934338>



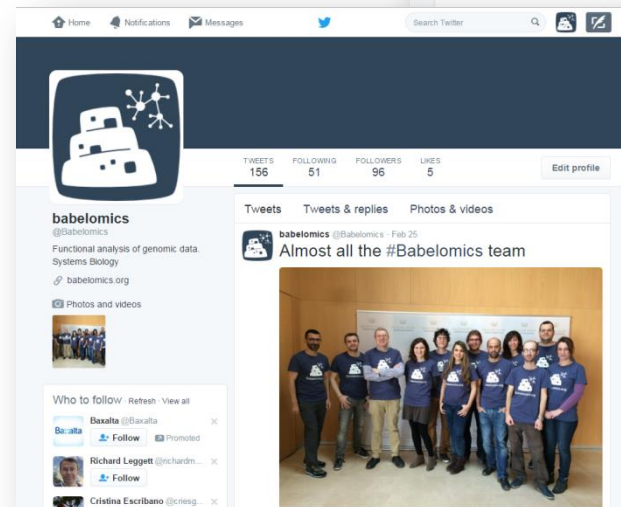
Babelomics group in twitter



@xdopazo

@bioinfocipf

@babelomics



<http://bioinfo.cipf.es>

<http://www.babelomics.org>

And the social dinner (and *mascletás*)... Networking is very important for your career. Keep in touch with fellows and instructors

