# Prioritization of variants and genes: BiERapp

Francisco García fgarcia@cipf.es Sep 2016





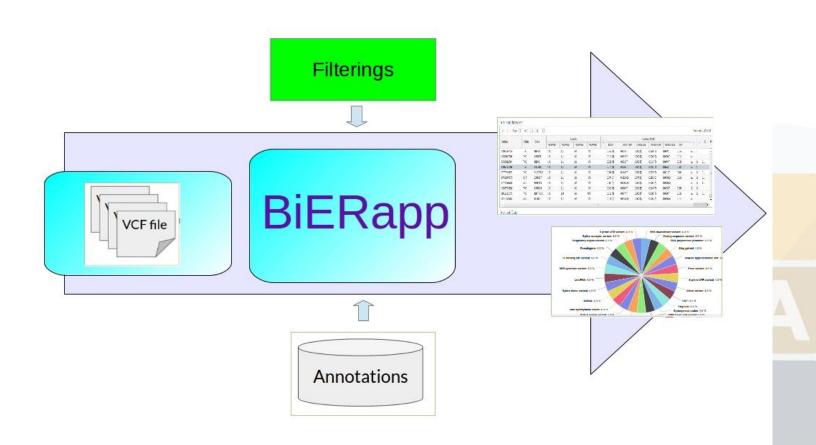


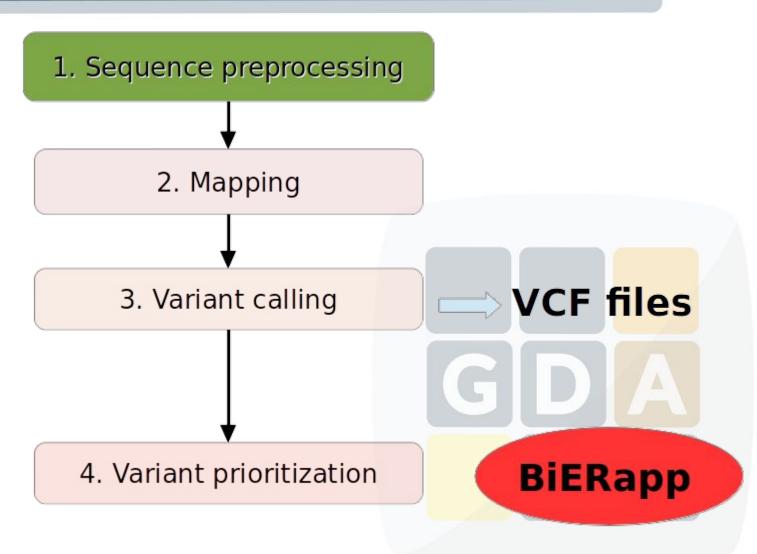


#### Introduction

- Whole-exome sequencing has become a fundamental tool for the discovery of disease-related genes of familial diseases but there are difficulties to find the causal mutation among the enormous background.
- ☐ There are different scenarios, so we need **different and immediate** strategies of prioritization.
- Vast amount of biological knowledge available in many databases.
- We need a tool to integrate this information and filter immediately to select candidate variants related to the disease

# How does BiERapp work?





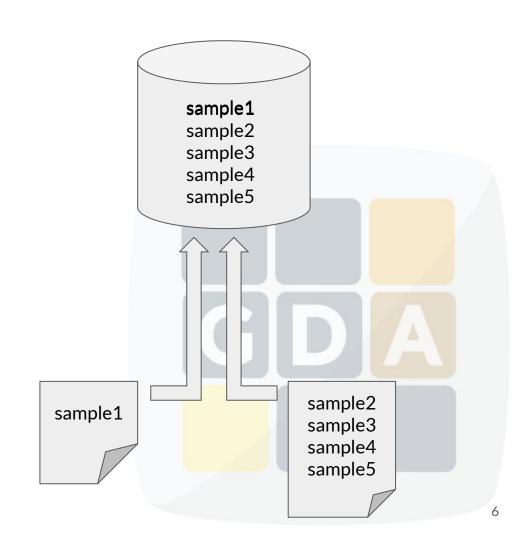
#### Input: VCF file

```
##fileformat=VCFv4.1
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seg/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS, Number=1, Type=Integer, Description="Number of Samples With Data">
##INFO=<ID=DP, Number=1, Type=Integer, Description="Total Depth">
##INFO=<ID=AF, Number=A, Type=Float, Description="Allele Frequency">
##INFO=<ID=AA, Number=1, Type=String, Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2, Number=0, Type=Flag, Description="HapMap2 membership">
##FILTER=<ID=q10, Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT, Number=1, Type=String, Description="Genotype">
##FORMAT=<ID=GQ, Number=1, Type=Integer, Description="Genotype Quality">
##FORMAT=<ID=DP, Number=1, Type=Integer, Description="Read Depth">
##FORMAT=<ID=HQ, Number=2, Type=Integer, Description="Haplotype Quality">
#CHROM POS
                                         QUAL FILTER INFO
                                                                                        FORMAT
                                                                                                    NA00001
                                                                                                                   NA00002
                                                                                                                                   NA00003
       14370
             rs6054257 G
                                             PASS
                                                    NS=3;DP=14;AF=0.5;DB;H2
                                                                                        GT:GQ:DP:HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:.,.
20
20
       17330
                                                     NS=3; DP=11; AF=0.017
                                                                                        GT:GQ:DP:HQ 0|0:49:3:58,50 0|1:3:5:65,3
                                                                                                                                   0/0:41:3
20
      1110696 rs6040355 A
                                G.T
                                              PASS
                                                    NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27 2|1:2:0:18,2
                                        67
                                                                                                                                   2/2:35:4
20
       1230237 .
                                              PASS
                                                     NS=3:DP=13:AA=T
                                                                                        GT:GQ:DP:HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
                                                     NS=3; DP=9; AA=G
20
       1234567 microsat1 GTC
                                G, GTCT 50
                                              PASS
                                                                                        GT:GO:DP
                                                                                                                                   1/1:40:3
```

- We can upload multiple VCF single/multi sample.
- You do not need to create a multi-sample file with all the samples.
- ☐ BiERapp will merge all the samples from those files in the database.

## Input: VCF multisample

- Create a Study
- Upload a new single-sample file
- BiERapp stores the sample in the created study
- Now we upload a new multisample file with 4 more samples
- BiERapp merges these samples in the study



#### Getting information: SIFT & PolyPhen

#### SIFT

- SIFT predicts whether an amino acid substitution affects protein function
- Interpretation: 1 (tolerated) to 0 (deleterious)

http://sift.jcvi.org/



#### PolyPhen

- Polymorphism Phenotyping is a tool which preditcs possible impact of an amino acid substitution on the structure and function of a human protein.
- Interpretation: 1 (probably damage) to 0 (bening)



#### Getting information: Conservation

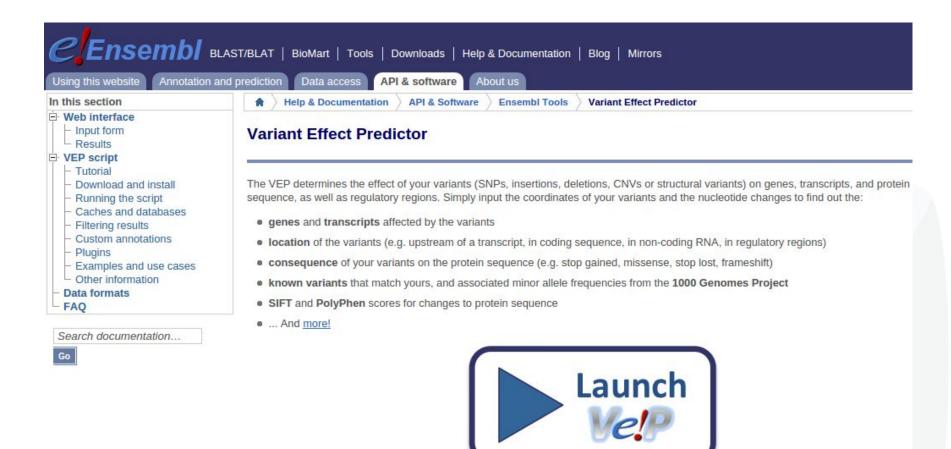
#### Phylop

- PhyloP scores measure evolutionary conservation at individual alignment sites.
   The scores are interpreted as follows compared to the evolution expected under neutral drift:
  - Positive scores -- Measure conservation, which is slower evolution than expected, at sites that are predicted to be conserved.
  - Negative scores -- Measure acceleration, which is faster evolution than expected, at sites that are predicted to be fast-evolving.

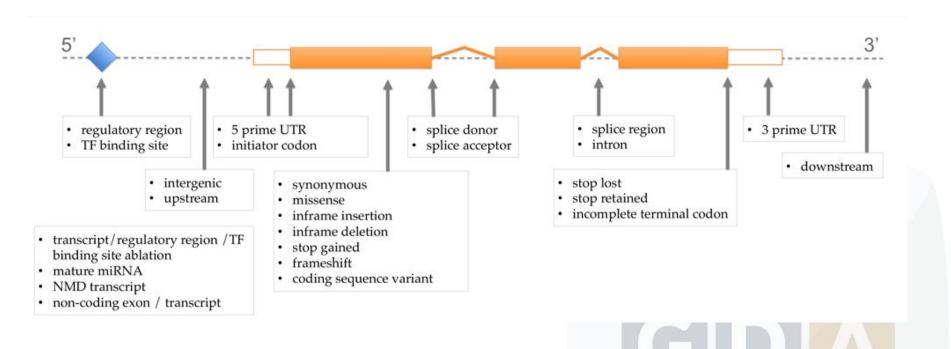
#### PhastCons

- PhastCons is a program for identifying evolutionarily conserved elements in a multiple alignment, given a phylogenetic tree.
- PhastCons essentially does three things:
  - It produces base-by-base conservation scores (as displayed in the conservation tracks in the UCSC browser)
  - It produces predictions of discrete conserved elements (as displayed in the "most conserved" tracks in the browser)
  - It estimates free parameters.

#### Getting information: Effect



#### Getting information: Effect



http://www.ensembl.org/info/genome/variation/predicted\_data.html

#### Getting information: Phenotype

#### ClinVar

ClinVar aggregates information about genomic variation and its relationship to human health.



# **GWAS Catalog**

The NHGRI-EBI Catalog of published genome-wide association studies

### Getting information: Pop. Frequencies



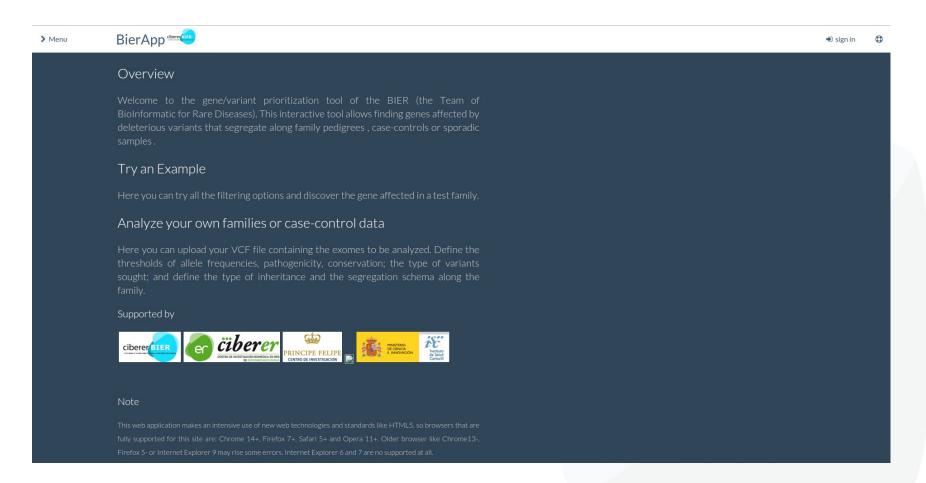


NHLBI Exome Sequencing Project (ESP)

Exome Variant Server

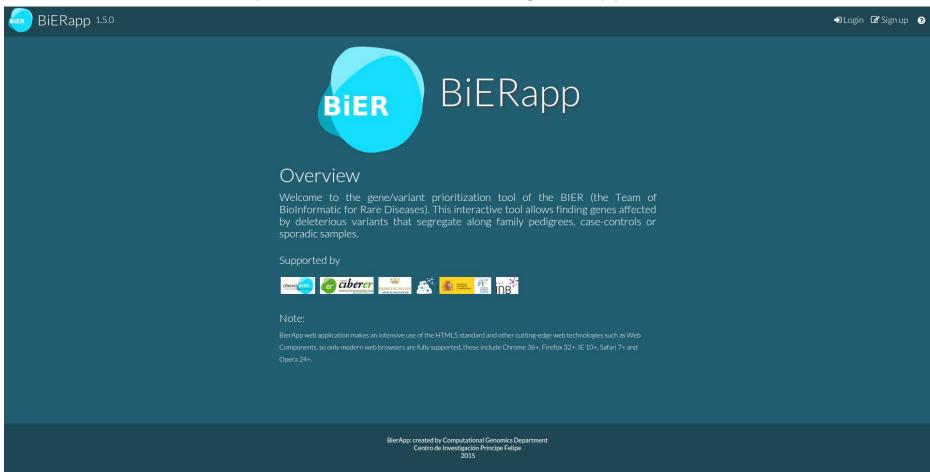
#### Tool interface: Official release

#### http://bierapp.babelomics.org/

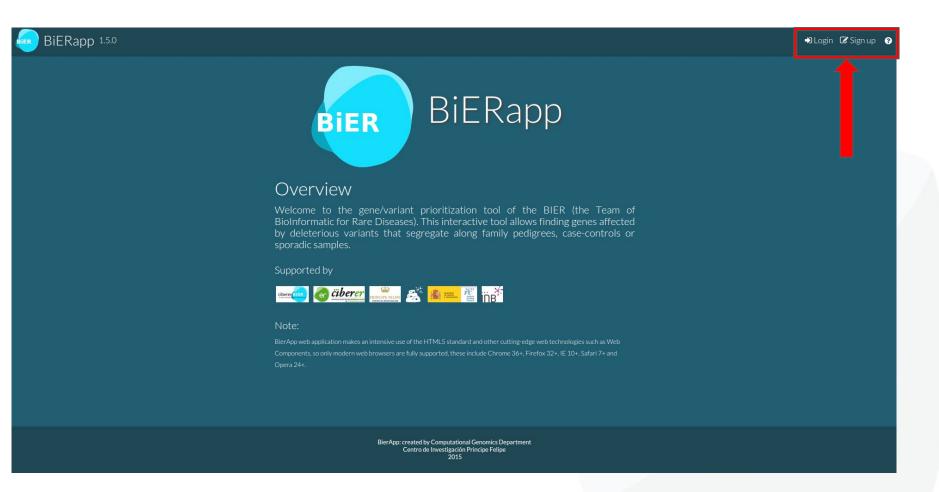


#### Tool interface:

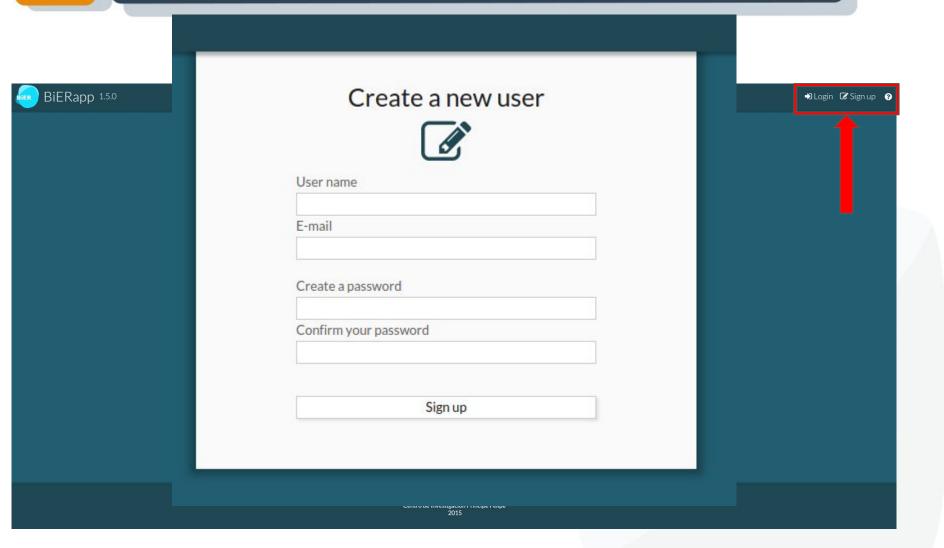
#### http://courses.babelomics.org/bierapp



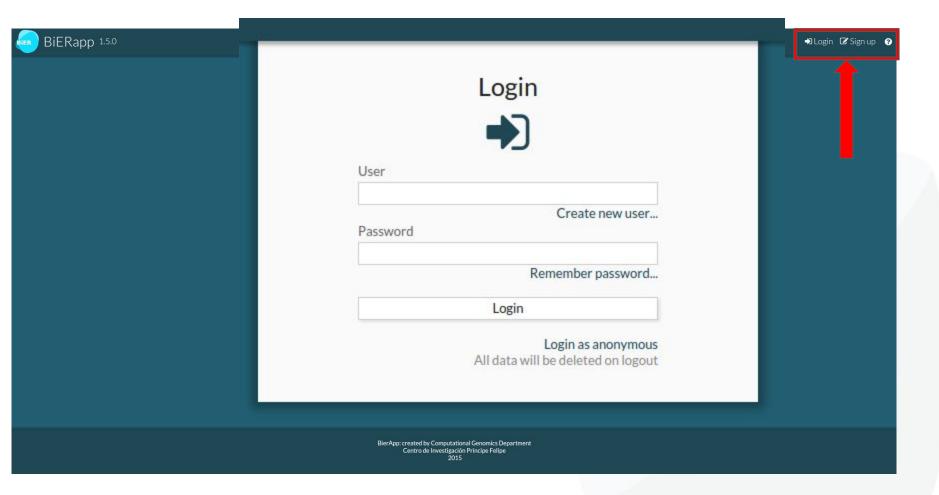
## Tool interface: Sign up /Log in

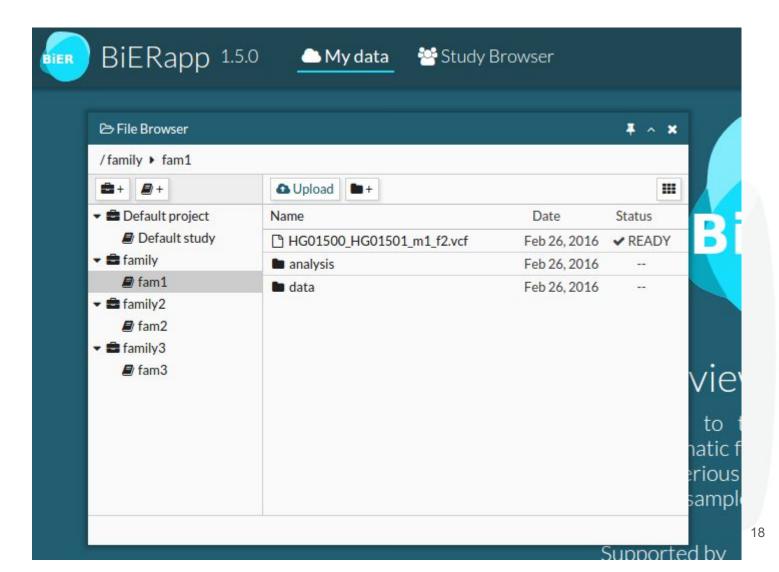


# Tool interface: Sign up /Log in

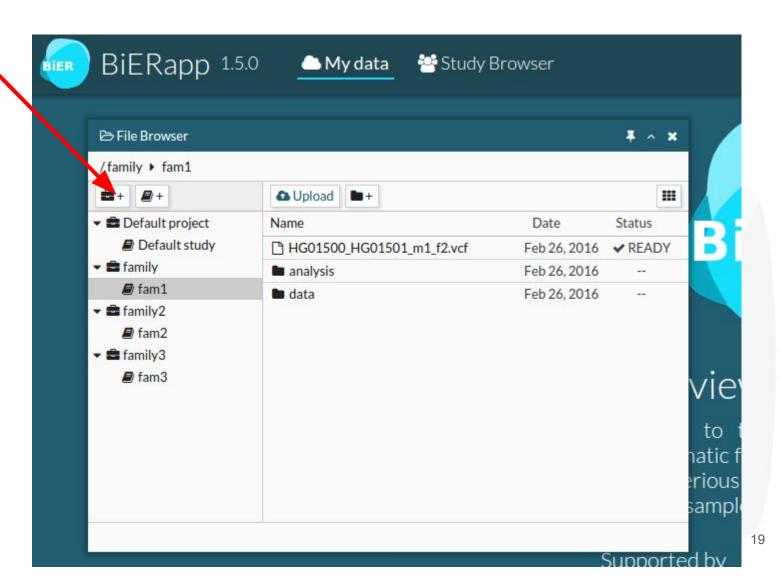


# Tool interface: Sign up /Log in

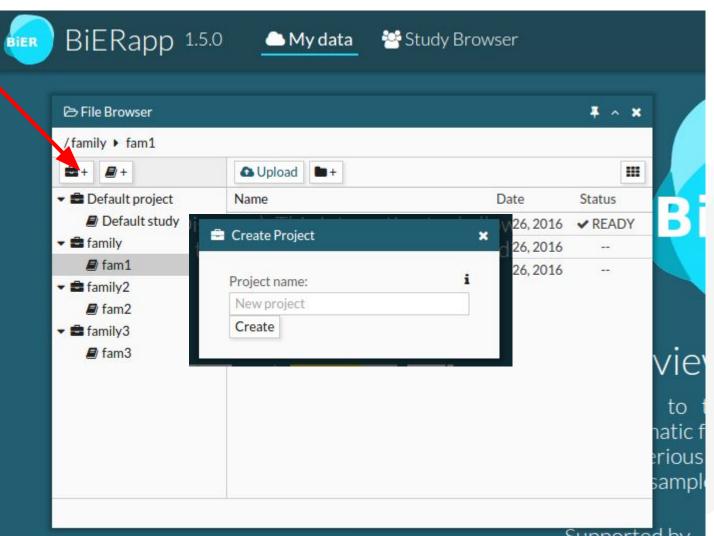


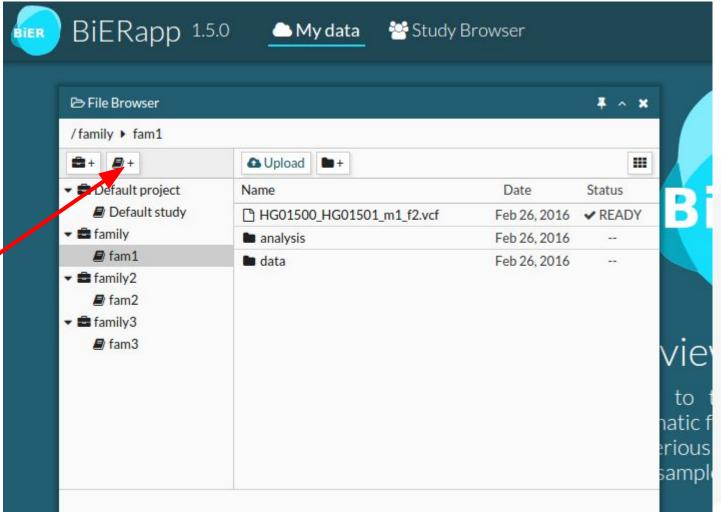


**New Project** 

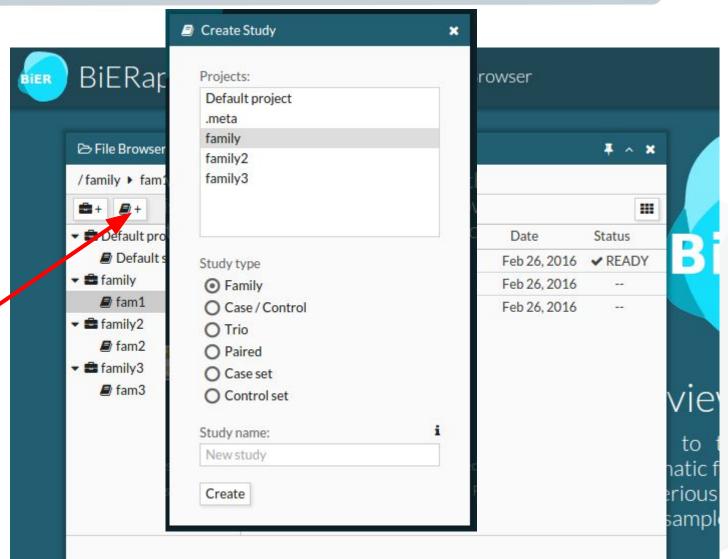


**New Project** 

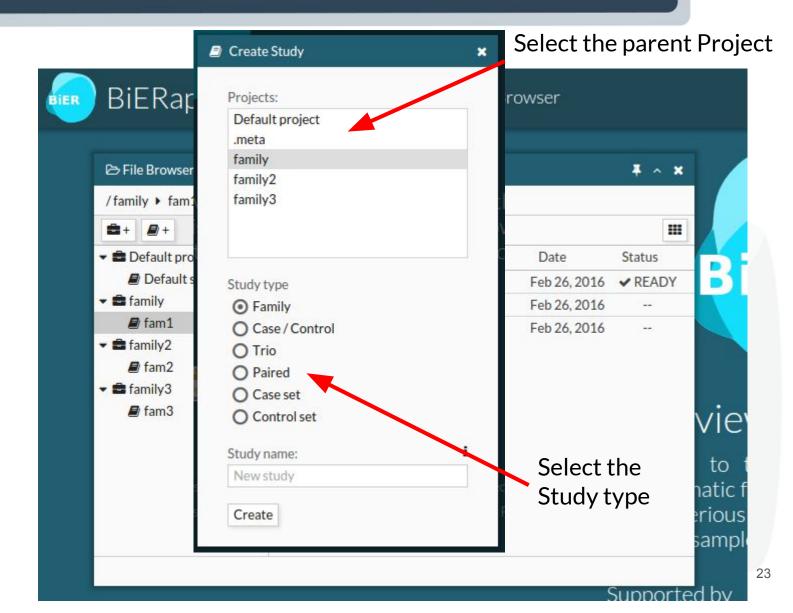




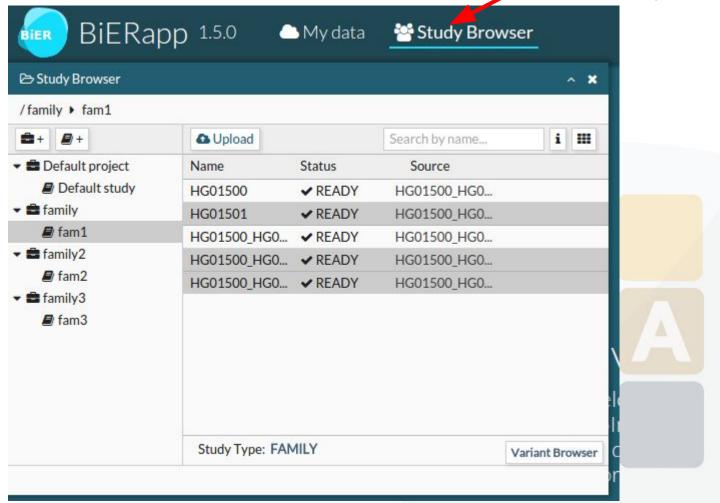
**New Study** 



**New Study** 

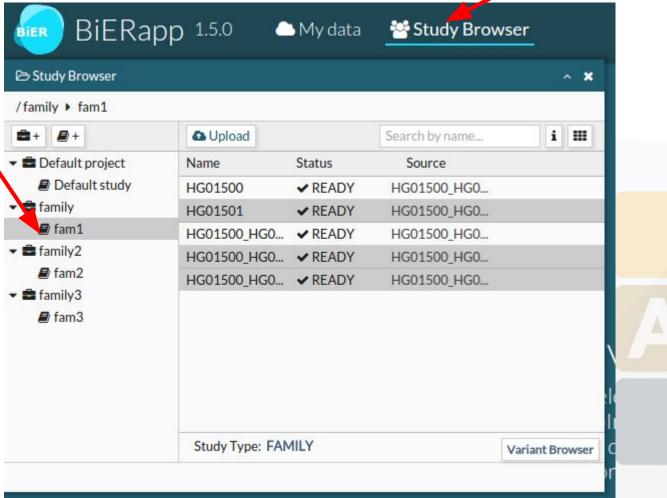


Select Study Browser



Select Study Browser

Choose your study

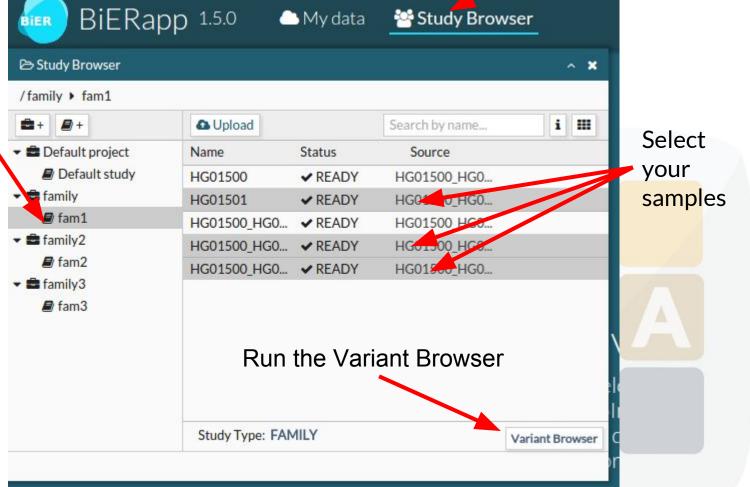


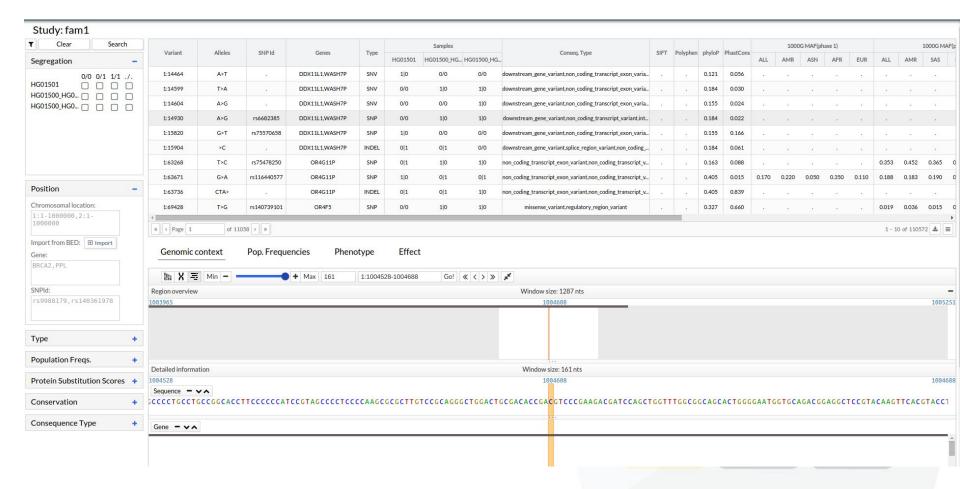
Select Study Browser BiERapp 1.5.0 My data Study Browser Study Browser ^ **X** /family > fam1 + 8+ i III **△** Upload Search by name... Select Default project Status Source Name your Default study **✓** READY HG01500 HG01500\_HG0... samples a family **✔** READY HG0" U\_HG0... HG01501 fam1 HG01500 HC0. HG01500\_HG0... ✓ READY ▼ amily2 HG01300 HG8 HG01500\_HG0... ✓ READY ■ fam2 HG01500\_HG0... ✓ READY HG01500\_HG0... ▼ amily3 ■ fam3 Study Type: FAMILY Variant Browser

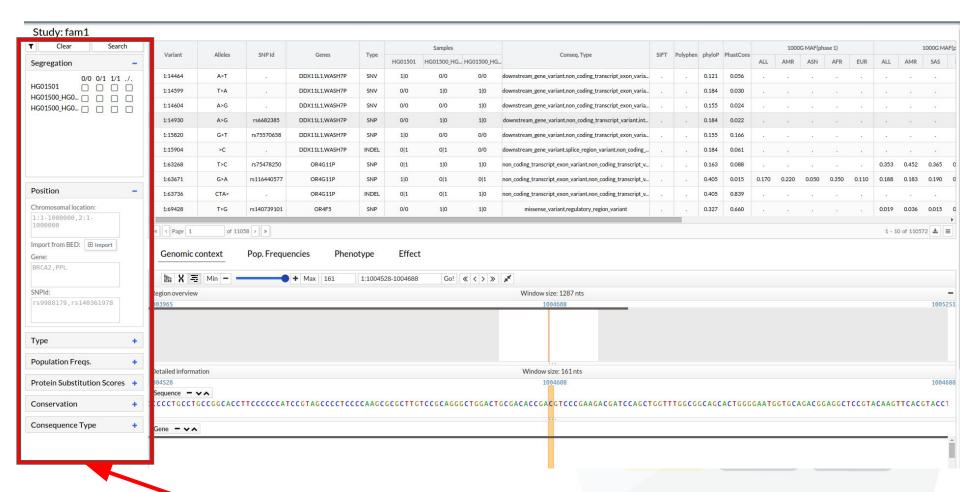
Choose your study

Select Study Browser ^ **X** i III Select your samples

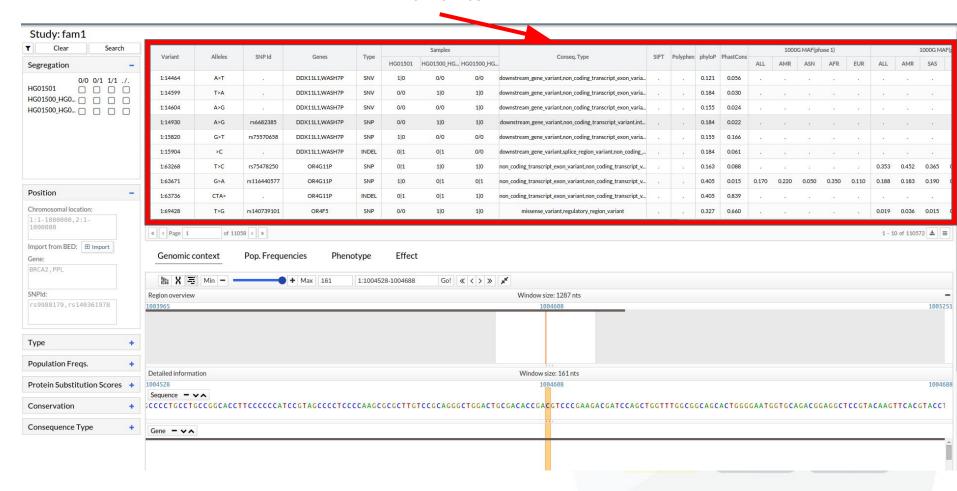
Choose your study

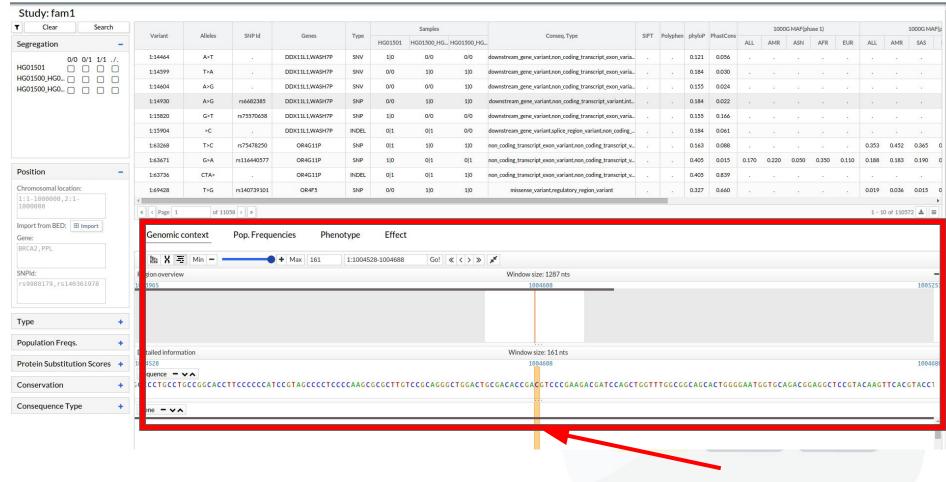




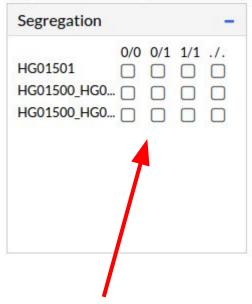


#### **Variants**

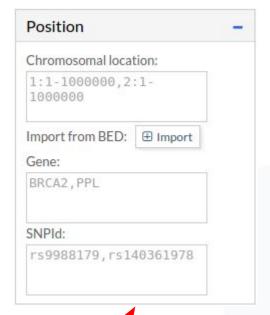




#### Tool interface: Filters



Choose your genotypes.

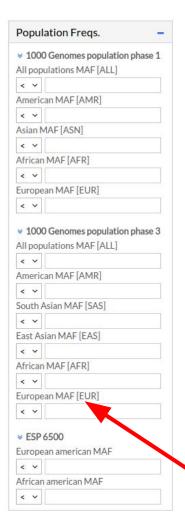


You can filter by region, Chromosome, Gene and SNPid. You can also import regions from a BED file





#### Tool interface: Filters

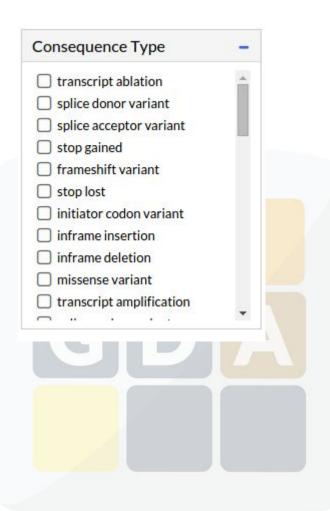






#### Filter by MAF

- 1000G phase1
- 1000G phase3
- ESP 6500



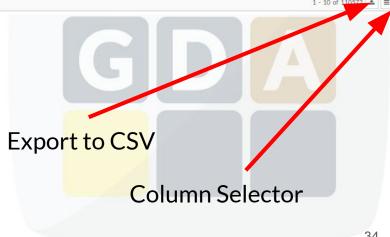
# Tool interface: Variant grid

#### Resizable columns

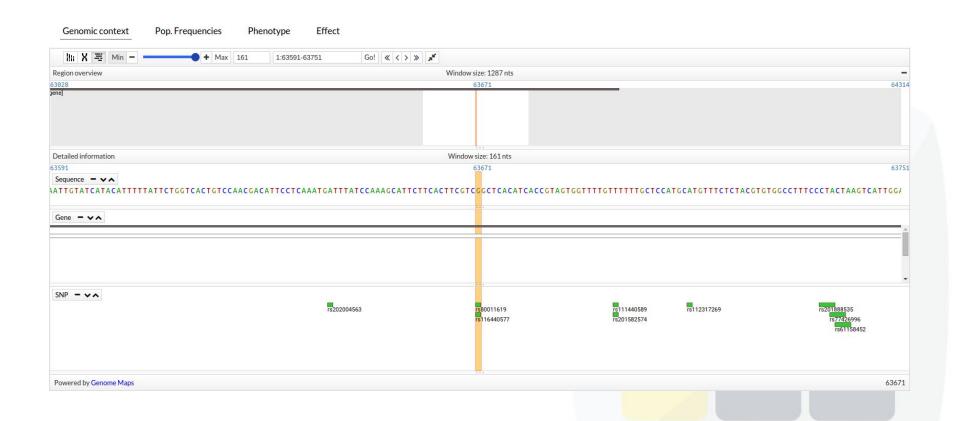
Variant	Alleles	SNP Id	Genes	Type		Samples		Conseq. Type	SIFT	Polyphen	phyloD	PhastCons		10000	G MAF(ph	ase 1)				1000G MAF
variant	Alleles	SIMP IQ	Genes	rype	HG01501	HG01500_HG	HG01500_HG		SIFI	roiypnen	phylop	PilasiCons	ALL	AMR	ASN	AFR	EUR	ALL	AMR	SAS
1:14464	A>T	20	DDX11L1,WASH7P	SNV	1 0	0/0	0/0	downstream_gene_variant,non_coding_transcript_exon_varia	10	1000	0.121	0.056			200	*			48	
1:14599	T>A		DDX11L1,WASH7P	SNV	0/0	1 0	1 0	downstream_gene_variant,non_coding_transcript_exon_varia			0.184	0.030		083		*				
1:14604	A>G	8	DDX11L1,WASH7P	SNV	0/0	0/0	1 0	downstream_gene_variant,non_coding_transcript_exon_varia	15	0.50	0.155	0.024	15		10				10	
1:14930	A>G	rs6682385	DDX11L1,WASH7P	SNP	0/0	1 0	1 0	downstream_gene_variant,non_coding_transcript_variant,int	- 1		0.184	0.022	12	7.0	7	7	- 1		10	
1:15820	G>T	rs75570658	DDX11L1,WASH7P	SNP	1 0	0/0	0/0	downstream_gene_variant,non_coding_transcript_exon_varia	-	10.63	0.155	0.166		1960	20	*			43	÷
1:15904	>C		DDX11L1,WASH7P	INDEL	0 1	0 1	0/0	downstream_gene_variant,splice_region_variant,non_coding			0.184	0.061			-6	-6	18			
1:63268	T>C	rs75478250	OR4G11P	SNP	0 1	1 0	1 0	non_coding_transcript_exon_variant,non_coding_transcript_v	100	0.00	0.163	0.088						0.353	0.452	0.365
1:63671	G>A	rs116440577	OR4G11P	SNP	1 0	0 1	0 1	non_coding_transcript_exon_variant,non_coding_transcript_v	1		0.405	0.015	0.170	0.220	0.050	0.350	0.110	0.188	0.183	0.190
1:63736	CTA>	27	OR4G11P	INDEL	0 1	0 1	1 0	non_coding_transcript_exon_variant,non_coding_transcript_v	10.	10.00	0.405	0.839	- 12		20	2			43	i.
1:69428	T>G	rs140739101	OR4F5	SNP	0/0	1 0	1 0	missense_variant,regulatory_region_variant			0.327	0.660			40			0.019	0.036	0.015



« < Page 1

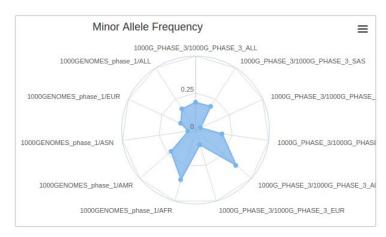


## Tool interface: Genomic Context



# Tool interface: Pop. Frequencies

Study	Population	SuperPopulation	Ref. Allele	Alt. Allele	Ref. Allele Fr	Alt. Allele Fr	MAF	0/0	0/1	
1000G_PHASE_3	1000G_PHASE_3_ALL	1000G_PHASE_3_ALL	G	Α	0.813	0.188	0.188	0	0	(
1000G_PHASE_3	1000G_PHASE_3_SAS	1000G_PHASE_3_SAS	G	Α	0.810	0.190	0.190	0	0	(
1000G_PHASE_3	1000G_PHASE_3_EAS	1000G_PHASE_3_EAS	G	Α	0.962	0.038	0.038	0	0	(
1000G_PHASE_3	1000G_PHASE_3_AMR	1000G_PHASE_3_AMR	G	Α	0.817	0.183	0.183	0	0	,
1000G_PHASE_3	1000G_PHASE_3_AFR	1000G_PHASE_3_AFR	G	Α	0.635	0.365	0.365	0	0	1
1000G_PHASE_3	1000G_PHASE_3_EUR	1000G_PHASE_3_EUR	G	A	0.896	0.104	0.104	0	0	1
000GENOMES_phase_1	AFR	AFR	G	A	0.650	0.350	0.350	0	0	
000GENOMES_phase_1	AMR	AMR	G	Α	0.780	0.220	0.220	0	0	1
000GENOMES_phase_1	ASN	ASN	G	Α	0.950	0.050	0.050	0	0	0
000GENOMES_phase_1	EUR	EUR	G	Α	0.890	0.110	0.110	0	0	
.000GENOMES_phase_1	ALL	ALL	G	Α	0.830	0.170	0.170	0	0	1





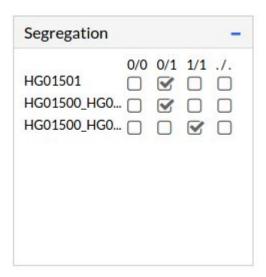
# Tool interface: Phenotype & Effect

mic:									
Gene name	Histology subtype	Mutation ID	Mutation somatic status	Primary histology	Primary site	Sample source	Site subtype	Tumour orig	
AGRN	adenocarcinoma	1126908	Confirmed somatic variant	carcinoma	prostate	fresh/frozen - NOS	NS	primary	
AGRN	adenocarcinoma	1126908	Confirmed somatic variant	carcinoma	large_intestine	NS	colon	NS	
AGRN	neoplasm	1126908	Confirmed somatic variant	other	thyroid	NS	NS	NS	
AGRN	neoplasm	1126908	Confirmed somatic variant	other	thyroid	NS	NS	NS	
AGRN	adenocarcinoma	1126908	Confirmed somatic variant	carcinoma	large_intestine	NS	colon	NS	

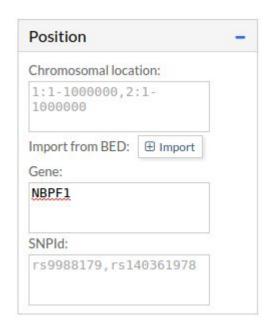
Clinvar:				
Accession	Clinical significance	Gene name	Review status	Traits
RCV000116259	Benign	AGRN	CLASSIFIED_BY_SINGLE_SUBMITTER	not specified,AllHighlyPenetrant,Not Specified
4				

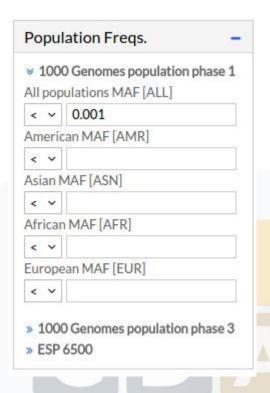
Gene Name	Ensembl Gene Id	Ensembl Transcript Id	Conseq. type	Relative Position	Codon	Strand	Biotype	cDna Position	cds Position	AA Position	AA Change	Sift	Polyphen
AGRN	ENSG00000188157	ENST00000379370	synonymous_variant		tcA/tcG	+	protein_coding	3116	3066	1022	SER/SER		,
AGRN	ENSG00000188157	ENST00000479707	2KB_downstream_gene_variant			+	retained_intron						
AGRN	ENSG00000188157	ENST00000466223	2KB_upstream_gene_variant			+	retained_intron						
AGRN	ENSG00000188157	ENST00000478677	2KB_upstream_gene_variant			+	retained_intron						
AGRN	ENSG00000188157	ENST00000492947	2KB_upstream_gene_variant			+	retained_intron						
AGRN	ENSG00000188157	ENST00000419249	upstream_gene_variant			+	protein_coding						
			regulatory_region_variant										

#### Results



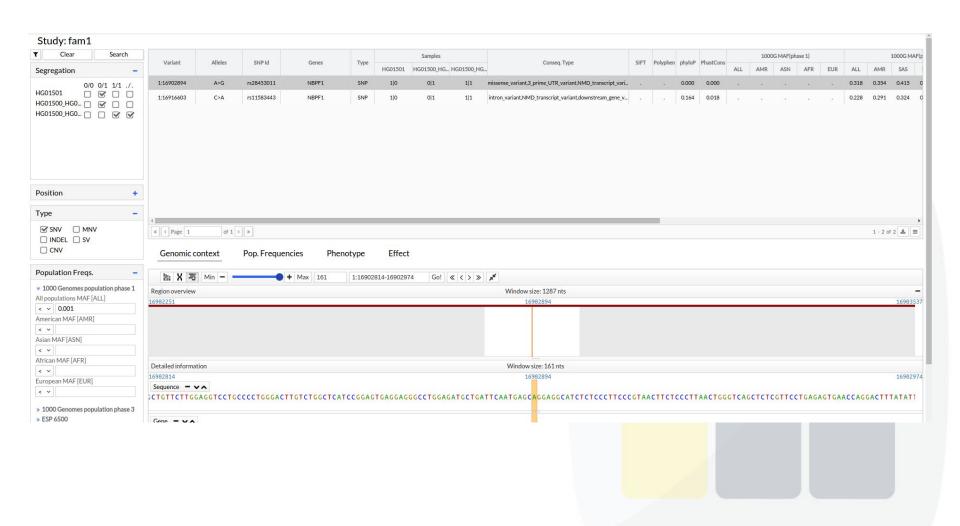






Click on "Search" and view the results

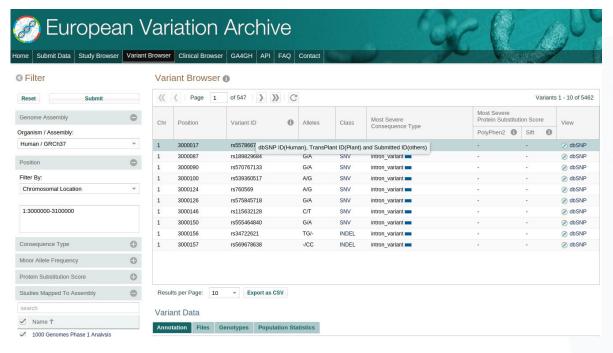
#### Results



#### Who is using BierApp?



Centro de Investigación Biomédica en Red Enfermedades Raras



IT4Innovations
national \$11€0
supercomputing
center1001\$1\$0



The EVA has its own customized version of BiERApp.

#### Conclusions

The proposed web-based interactive framework has **great potential to detect disease-related variants** in familial diseases as demonstrated by its successful use in several studies.

The use of the filters is interactive and the results are almost instantaneously displayed in a panel that includes the genes affected, the variants and specific information for them.

Candidate variants are new knowledge useful for future diagnostic.

#### More info: publication

Nucleic Acids Research Advance Access published May 6, 2014

Nucleic Acids Research, 2014 1 doi: 10.1093/nar/gku407

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

Alejandro Alemán<sup>1,2</sup>, Francisco Garcia-Garcia<sup>1</sup>, Francisco Salavert<sup>1,2</sup>, Ignacio Medina<sup>1</sup> and Joaquín Dopazo<sup>1,2,3,\*</sup>

<sup>&</sup>lt;sup>1</sup>Computational Genomics Department, Centro de Investigación Príncipe Felipe (CIPF), Valencia 46012, Spain,

<sup>&</sup>lt;sup>2</sup>Bioinformatics of Rare Diseases (BIER), CIBER de Enfermedades Raras (CIBERER), Valencia 46010, Spain and

<sup>&</sup>lt;sup>3</sup>Functional Genomics Node, (INB) at CIPF, Valencia 46012, Spain

# GDA CIBERER

# http://bioinfo.cipf.es/gda16ciberer/



