# Using public sources of data

## H3ABioNet Data Management Workshop

# Data sources

- Genotypes
  - 1000 Genomes
  - HapMap
  - dbGaP
  - EGA
  - GEO
  - WTCCC
  - PGP

# Data sources

- Annotations
  - Chip manufacturers
  - ENSEMBL
  - RefSeq
  - dbNSFP

# Tools

- Google
- Galaxy
- PLINK ("PLINK 2")
- Scripting
  - sed, awk, grep, cut, sort, uniq, parallel
- Databases

# Data QC issues

- Genome build
- Annotation builds
- Identifiers
- Strand

# Question 1

Hi

I would like to know how many SNPs contained in the Immunochip are in the Illumina 2.5 Million SNPs chip.

The Immunochip is an Illumina Infinium genotyping chip, containing 196,524 polymorphisms (718 small insertion deletions, 195,806 SNPs), designed to perform deep replication of major autoimmune and inflammatory diseases.

The 2.5M Illumina chip promises to perform well in typing both common and rare SNP content from the 1kGP (MAF>2.5%) for diverse world populations. This array contains tagSNP data from recently released 1000 Genomes Project pilot data.

# Usenet post, 1997

Some people, when confronted with a problem, think
"I know, I'll use regular expressions."

Now they have two problems.

# Galaxy

# Visualisation in Trackster

# 1st attempt

UCSC table browser

Galaxy

# UCSC table browser

**Tools**

search tools

Get Data
Send Data
Lift-Over
Text Manipulation
Convert Formats
FASTA manipulation
Filter and Sort
Join, Subtract and Group
Extract Features
Fetch Sequences
Fetch Alignments
Get Genomic Scores
Operate on Genomic Intervals
Statistics
Graph/Display Data
Regional Variation
Multiple regression
Multivariate Analysis
Evolution
Motif Tools
Multiple Alignments
Metagenomic analyses
Genome Diversity

NGS TOOLBOX BETA

Phenotype Association
NGS: QC and manipulation
NGS: Mapping
NGS: SAM Tools
NGS: GATK Tools (beta)

https://wiki.galaxyproject.org/GalaxyIsHiring

**Galaxy** is an open source, web–based platform for data intensive biomedical research. If you are new to Galaxy start here or consult our help resources.

**Tweets**  Follow

**Galaxy Project** @galaxyproject  13h
Doing #highthroughput research? Want to save money? Register for GCC2014 by THIS FRIDAY. bit.ly/gcc2014reg #usegalaxy
Expand

**Ravi K Madduri** @madduri  14h
Our crop modeling #usegalaxy hackathon at Uchicago cc @ia nfoster pic.twitter.com/qnwcgbpYDg
↻ Retweeted by Galaxy Project
Show Photo

**Galaxy Project** @galaxyproject  16h
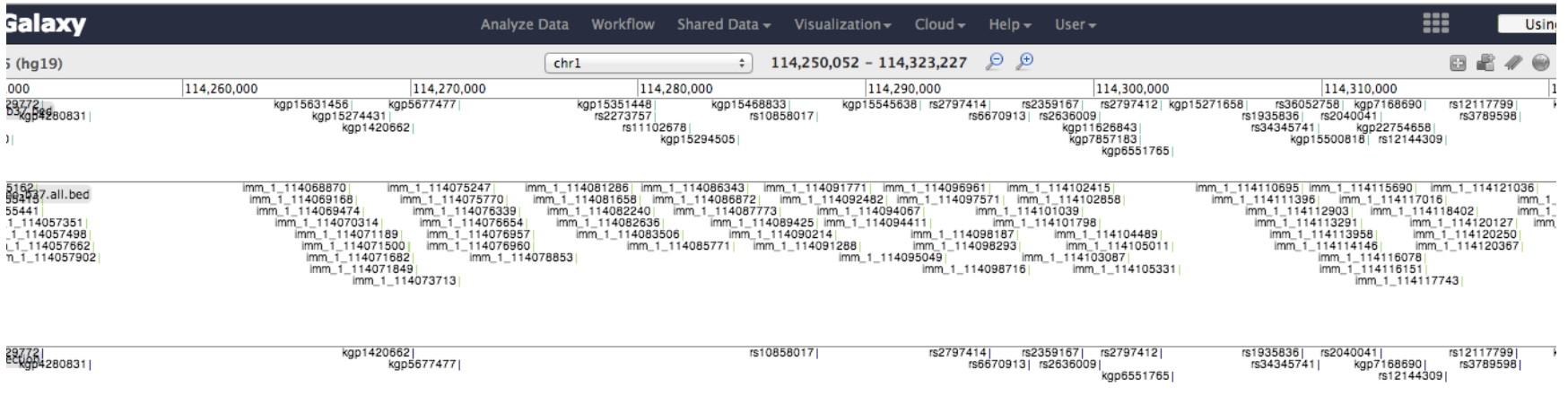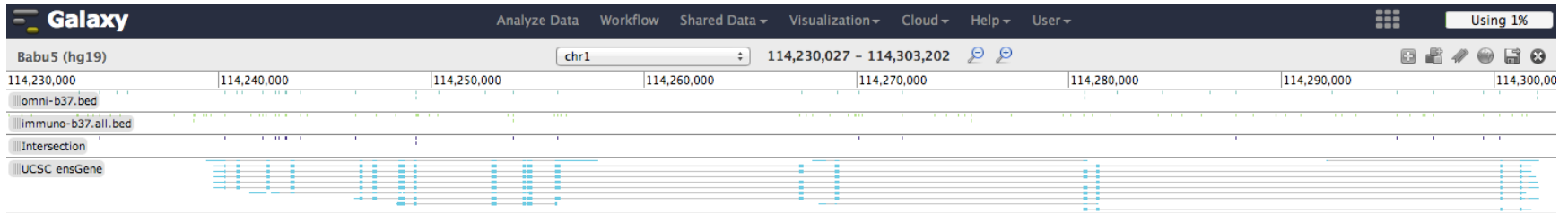@irazoqui_javier Please see the responses to a similar question at our

Tweet to @galaxyproject

PENNSTATE
1855

JOHNS HOPKINS
UNIVERSITY

TACC

iPlant Collaborative

The Galaxy Team is a part of the Center for Comparative Genomics and Bioinformatics at Penn State, and the ___ay and at Johns Hopkins University.

This instance of Galaxy is utilizing infrastructure generously provided by the iPlant Collaborative at the Texas Advanced Computing Center, with support from the

**History**

**Unnamed history**

0 bytes

**1: Homo sapiens Short Variation (SNPs and indels) (GRCh37.p13)**

# 2<sup>nd</sup> attempt

BioMart

Galaxy

set

o sapiens Short Variation
's and indels)
Ch37.p13)

ers

iation Set Name :
mina_ImmunoChip

ributes

iation Name
iation source
romosome name
sition on Chromosome (bp)

set

e Selected]

☐ Phenotype

`17-@ALPHA-HYDROXYLASE/1720-LYASE DEFICIENCY COMBINED COMPLETE`

☐ Phenotype significance [0 non significant, 1 significant]

`0  ⬍`

☑ Variation Set Name

```
Illumina_ExomeChip
Illumina_Human610_Quad
Illumina_Human660W-quad
Illumina_HumanHap550
Illumina_HumanHap650Y
Illumina_HumanOmni1-Quad
Illumina_HumanOmni2.5
Illumina_HumanOmni5
Illumina_ImmunoChip
Marjolein Kriek
```

☐ SIFT Prediction

```
tolerated
deleterious
```

☐ SIFT score <= [0 most deleterious, 1 least deleterious]

☐ PolyPhen Prediction

```
unknown
benign
possibly damaging
probably damaging
```

☐ PolyPhen score >= [1 most damaging, 0 least damaging]

☐ Global minor allele frequency <=

☐ Global minor allele frequency >=

☐ Clinical significance

```
drug-response
histocompatibility
```

et

sapiens Short Variation
s and indels)
137.p13)

rs

ation Set Name :
ina_ImmunoChip
mosome : 22

butes

ation Name
mosome name
tion on Chromosome (bp)

et

Selected]

Export  all results to    [ Galaxy ⇕ ] [ TSV ⇕ ] ☐ Unique results only   ✅ Go

Email notification to    [_____]

View    [ 10 ⇕ ] rows as [ TSV ⇕ ] ☐ Unique results only

Loading... 🔄

# 3rd attempt

Illumina website

Custom cleanup

Galaxy

# Manufacturer annotations



Support » **Downloads**                                    🖶 💬 ➕ | Follow us: 📧

## Downloads

This is the Sequencing Downloads A-Z list. Use your browser's find function (CTRL-F on PCs, or Command-F on Macs) to search by keyword.

*The downloadable materials displayed on this web page are proprietary to Illumina, Inc., and are intended solely for the use of its customers and for no other purpose than use with Illumina's products or services. The downloadable materials and their contents shall not be used or distributed for any other purpose or otherwise communicated, disclosed, or reproduced in any way without the prior written consent of Illumina, Inc.

You can download current software and support files by clicking one of the links listed below. The description expands to show available downloads. Click the desired download and select Save. If you are looking for documentation, visit the Documentation page.

[ ~ Select a technology ⬍ ]  [ ~ Select a workflow ⬍ ]

| DESCRIPTION | FILE INFO | DATE |
|---|---|---|

+ ADME Plug-in Setup v1.0.1.4
This download contains the ADME Plug-in Setup v1.0.1.4 installer.

+ African American Admixture Panel Product Files
This download contains the Manifest (.opa), BeadStudio Project (.bsc), and annotation file for the GoldenGate African American Admixture Panel.

+ Amplicon Viewer Installer
The Illumina Amplicon Viewer is a desktop tool to allow users to analyze their MiSeq Amplicon data. With the Illumina Amplicon Viewer, users can aggregate samples from multiple runs for data analysis and visualization. Minimum System Requirements: Windows Vista or 7 Operating System; 32-bit system /4 GB AM or 64-bit system /8 GB RAM; Microsoft.Net Framework 4.0 or above; Microsoft Office Excel 2010 (recommended).

+ Analysis Visual Controller (AVC) v 1.7
This download contains the Installer, Installation Guide and User Guide for the Analysis Visual Controller

BED format provides a flexible way to define the data lines that are displayed in an annotation track. BED lines have three required fields and nine additional optional fields. The number of fields per line must be consistent throughout any single set of data in an annotation track. The order of the optional fields is binding: lower-numbered fields must always be populated if higher-numbered fields are used.

If your data set is BED-like, but it is very large and you would like to keep it on your own server, you should use the bigBed data format.

The first three required BED fields are:

1. **chrom** - The name of the chromosome (e.g. chr3, chrY, chr2_random) or scaffold (e.g. scaffold10671).
2. **chromStart** - The starting position of the feature in the chromosome or scaffold. The first base in a chromosome is numbered 0.
3. **chromEnd** - The ending position of the feature in the chromosome or scaffold. The *chromEnd* base is not included in the display of the feature. For example, the first 100 bases of a chromosome are defined as *chromStart=0, chromEnd=100*, and span the bases numbered 0-99.

The 9 additional optional BED fields are:

4. **name** - Defines the name of the BED line. This label is displayed to the left of the BED line in the Genome Browser window when the track is open to full display mode or directly to the left of the item in pack mode.
5. **score** - A score between 0 and 1000. If the track line *useScore* attribute is set to 1 for this annotation data set, the *score* value will determine the level of gray in which this feature is displayed (higher numbers = darker gray). This table shows the Genome Browser's translation of BED score values into shades of gray:

| shade | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| score in range | ≤ 166 | 167-277 | 278-388 | 389-499 | 500-611 | 612-722 | 723-833 | 834-944 | ≥ 945 |

6. **strand** - Defines the strand - either '+' or '-'.
7. **thickStart** - The starting position at which the feature is drawn thickly (for example, the start codon in gene displays).
8. **thickEnd** - The ending position at which the feature is drawn thickly (for example, the stop codon in gene displays).
9. **itemRgb** - An RGB value of the form R,G,B (e.g. 255,0,0). If the track line *itemRgb* attribute is set to "On", this RBG value will determine the display color of the data contained in this BED line. NOTE: It is recommended that a simple color scheme (eight colors or less) be used with this attribute to avoid overwhelming the color resources of the Genome Browser and your Internet browser.
10. **blockCount** - The number of blocks (exons) in the BED line.
11. **blockSizes** - A comma-separated list of the block sizes. The number of items in this list should correspond to *blockCount*.
12. **blockStarts** - A comma-separated list of block starts. All of the *blockStart* positions should be calculated relative to *chromStart*. The number of items in this list should correspond to *blockCount*.

Illumina, Inc.

[Heading]

Descriptor File Name,Immuno_BeadChip_11419691_B.bpm

Assay Format,Infinium HD Ultra

Date Manufactured,11/01/2010

Loci Count ,196524

[Assay]

IlmnID,Name,IlmnStrand,SNP,AddressA_ID,AlleleA_ProbeSeq,AddressB_ID,AlleleB_ProbeSeq,GenomeBuild,Chr,MapInfo,Ploidy,Species,Source,SourceVersion,SourceStrand,SourceSeq,TopGenomicSeq,BeadSetID

1-159076491-G-DELETION-1_P_F_1767851002,1-159076491-G-DELETION,PLUS,[D/I],0049656396,ACAGCAATCCTGTGAGGTACTTATTATCACCCCCATTTTACTCAAGGGGG,,,36,1,159076491,diploid,Homo sapiens,WTCCCseq,1,PLUS,TTTACTCTTAACAGCAATCCTGTGAGGTACTTATTATCACCCCCATTTTACTCAAGGGGG[-/G]AAGAAAATTGAGGCTCAGAGAGGTTAATGAATCTGCCAGAGATCACAGAGCTTCTTTTTT,TTTACTCTTAACAGCAATCCTGTGAGGTACTTATTATCACCCCCATTTTACTCAAGGGGG[-/G]AAGAAAATTGAGGCTCAGAGAGGTTAATGAATCTGCCAGAGATCACAGAGCTTCTTTTTT,285

1-159093319-A-DELETION-1_M_R_1767851004,1-159093319-A-DELETION,MINUS,[I/D],0043648403,TCTATTCTGCATATTAGTTGCCTGTAGGATTCATAGTTTGCAATTTTTTT,,,36,1,159093319,diploid,Homo sapiens,WTCCCseq,1,PLUS,ATGTACAGTAAAGGAAACAATTCACAGAGTGAAAAGGCAACCAATAGAATAGGAAAAAAA[-/A]TTGCAAACTATGAATCCTACAGGCAACTAATATGCAGAATAGACAAGAAATTCAAACGTC,ATGTACAGTAAAGGAAACAATTCACAGAGTGAAAAGGCAACCAATAGAATAGGAAAAAAA[-/A]TTGCAAACTATGAATCCTACAGGCAACTAATATGCAGAATAGACAAGAAATTCAAACGTC,285

# Steps

- Extract columns wanted
- Add end position
- Split by build numbers
- Upload to Galaxy
- Concatenate B36 annotations
- Liftover to B37

```bash
#!/bin/bash

# for some reason, the immuno annotation file has both build 36 and 37.1 annotations
# need to split them

cat input/annot/immuno_beadchip_11419691_b.csv \
    | grep -e ".*\[[ACGT]/[ACGT]" \
    | awk -F',' '{if ($9=="36"  )print "chr"$10"\t"$11"\t"$11+1"\t"$2}' \
    | grep -v -f skip > tmp/immuno-b36.bed
cat input/annot/immuno_beadchip_11419691_b.csv \
    | grep -e ".*\[[ACGT]/[ACGT]" \
    | awk -F',' '{if ($9=="36.2")print "chr"$10"\t"$11"\t"$11+1"\t"$2}' \
    | grep -v -f skip > tmp/immuno-b36.2.bed
cat input/annot/immuno_beadchip_11419691_b.csv \
    | grep -e ".*\[[ACGT]/[ACGT]" \
    | awk -F',' '{if ($9=="37.1")print "chr"$10"\t"$11"\t"$11+1"\t"$2}' \
    | grep -v -f skip > tmp/immuno-b37.bed
cat input/annot/humanomni25m-8v1-1_b.csv \
    | grep -e ".*\[[ACGT]/[ACGT]" \
    | awk -F',' '{if ($9=="37.1")print "chr"$10"\t"$11"\t"$11+1"\t"$2}' \
    | grep -v -f skip > tmp/omni-b37.bed
```

```
chrY    21762685 21762686 200610-147
chrY    21779251 21779252 200610-148
chrY    21867854 21867855 200610-149
chrY    21751440 21751441 200610-150
chrY    21888865 21888866 200610-151
chrY    21730357 21730358 200610-152
chrY    21740450 21740451 200610-153
chrY    21753199 21753200 200610-155
chrY    21868776 21868777 200610-156
chrY    17286006 17286007 200610-158
```

chr6   31,500,371 – 31,511,110

rs6915692|  rs2516474|  rs9267482|  rs10456396|  rs2734583|  rs3093974|  rs933208|  rs2239709|  rs11757236|  rs2523508|  kgp9277093|  rs2239527|
85|  rs3131628|  rs2471826|  kgp8005916|  rs3130058|  rs1129640|  rs12665501|  rs7738430|  rs3130059|  rs2523506|
38|  rs3093976|  rs2734582|  rs2071596|  rs12178599|  rs2239525|  rs2523505|
055388|  kgp1273477|  rs2075581|  rs2516393|  kgp10338335|  rs2239526|  rs22395
rs2523512|
rs2523511|

055388|  rs3131628|  rs929138|  imm_6_31612416|  rs2734583|  rs3093974|  rs2071596|  rs2239709|  rs12665489|  rs7738430|  rs3130059|  rs22395
rs3093976|  rs2075580|  rs10456396|  rs9380261|  rs1129640|  rs2239525|
84|  rs9267482|  rs3130057|  rs3130058|  rs2523512|  rs2239526|

rs3131628|  rs9267482|  rs10456396|  rs2734583|  rs3093974|  rs2071596|  rs2239709|  rs7738430|  rs3130059|  rs22395
36|  rs3093976|  rs3130058|  rs1129640|  rs2239525|
055388|  rs2523512|  rs2239526|

ENSG000002340
ENSG00000198563
ENSG00000201785
ENSG00000198563
ENSG00000254870

ENSG00000198563

ENSG00000198563
ENSG00000198563
ENSG00000234
ENSG00000198563
ENSG00000198563
ENSG00000198563
ENSG00000198563
ENSG00000198563
ENSG00000198563
ENSG00000254870
ENSG00000265236

# Question 2

I would like to view the allele frequencies of my study populations and compare them to the frequencies found in the HapMap populations, for the following genes

# BioMart 0.8

# Data marts

# Databases

# BioMart

- Integration
- Not "live"
- R interface (biomaRt)

# SQL query

- SELECT main.rsid_106, main.chromosome_106, main.position_106, afmap_mart.marker__allele_frequency_KHS__dm.a_freq_101,
afmap_mart.marker__allele_frequency_HER__dm.a_freq_101, afmap_mart.marker__allele_frequency_STS__dm.a_freq_101,
afmap_mart.marker__allele_frequency_CON__dm.a_freq_101, afmap_mart.marker__allele_frequency_ZUL__dm.a_freq_101,
afmap_mart.marker__allele_frequency_ASW__dm.a_freq_101, afmap_mart.marker__allele_frequency_CEU__dm.a_freq_101,
afmap_mart.marker__allele_frequency_CHB__dm.a_freq_101, afmap_mart.marker__allele_frequency_CHD__dm.a_freq_101,
afmap_mart.marker__allele_frequency_GIH__dm.a_freq_101, afmap_mart.marker__allele_frequency_JPT__dm.a_freq_101,
afmap_mart.marker__allele_frequency_LWK__dm.a_freq_101, afmap_mart.marker__allele_frequency_MEX__dm.a_freq_101,
afmap_mart.marker__allele_frequency_MKK__dm.a_freq_101, afmap_mart.marker__allele_frequency_TSI__dm.a_freq_101,
afmap_mart.marker__allele_frequency_YRI__dm.a_freq_101 FROM afmap_mart.marker__allele_frequency_TSI__dm,
afmap_mart.marker__allele_frequency_CHD__dm, afmap_mart.marker__allele_frequency_ZUL__dm,
afmap_mart.marker__allele_frequency_KHS__dm, afmap_mart.marker__allele_frequency_LWK__dm,
afmap_mart.marker__allele_frequency_CHB__dm, afmap_mart.marker__allele_frequency_ASW__dm,
afmap_mart.marker__allele_frequency_JPT__dm, afmap_mart.marker__allele_frequency_STS__dm,
afmap_mart.marker__allele_frequency_HER__dm, afmap_mart.marker__allele_frequency_CEU__dm,
afmap_mart.marker__allele_frequency_MEX__dm, afmap_mart.marker__allele_frequency_MKK__dm,
afmap_mart.marker__allele_frequency_CON__dm, afmap_mart.marker__allele_frequency_YRI__dm,
afmap_mart.marker__allele_frequency_GIH__dm, afmap_mart.marker__association__main main WHERE (main.genesymbol_103 like '%BRCA2%')
AND main.id_106_key=afmap_mart.marker__allele_frequency_TSI__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_KHS__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_ASW__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_CEU__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_YRI__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_ZUL__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_CHD__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_JPT__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_GIH__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_STS__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_CON__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_MKK__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_MEX__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_LWK__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_HER__dm.id_106_key AND
main.id_106_key=afmap_mart.marker__allele_frequency_CHB__dm.id_106_key

# Traditional Relational DB Schema

| User | |
|------|------|
| UserId | Name |
| 1 | Bob |
| 2 | Chris |
| 3 | Fred |

| Blog | | |
|------|------|------|
| BlogId | Name | DatePosted |
| 3 | NoSQL vs RDBMS | 1/30/2012 |

| Comment | | | | |
|---------|--------|--------|-----------------------------|----------|
| CommentId | BlogId | UserId | Value | Date |
| 1 | 3 | 1 | This blog rocks | 2/1/2012 |
| 2 | 3 | 1 | Exactly what I was looking for! | 2/2/2012 |
| 3 | 3 | 2 | I'm a hater, too generalized | 2/2/2012 |

# Relational

# Galaxy

# Showing 62 kbp from chr18, positions 148,000 to 210,000

**+ Instructions**
[Bookmark this] [Upload your own data] [Show banner] [Share these tracks] [Link to Image] [SNP genotype data] [HapMap LD Data] [High-res Image] [Help] [Reset]
**+ Search**
**− Overview**

chr18

0M          10M          20M          30M          40M          50M          60M          70M

**− 🔊 ❓ Ideogram**

**− Region**

chr18

0M   0.1M  0.2M  0.3M  0.4M  0.5M  0.6M  0.7M  0.8M  0.9M   1M   1.1M  1.2M  1.3M  1.4M  1.5M  1.6M  1.7M  1.8M  1.9M   2M

**− Details**

150k          160k          170k          180k          190k          200k          210k

**− 🔊 ❓ Genotyped SNPs**

rs551835(+)                                                        rs553212(+)                                    rs592537(+
1 ▌▌▌▌▌▌▌                                                          1 ▌▌▌▌▌▌▌                                        1 ▌▌▌▌▌▌▌
CHJYHKSXZ                                                          CHJYHKSXZ                                        CHJYHKSX

rs16944691(+)                                                      rs655781(+)                                      rs600220
1 ▌▌▌▌▌▌▌                                                          1 ▌▌▌▌▌▌▌                                        1 ▌▌▌▌▌▌▌
CHJYHKSXZ                                                          CHJYHKSXZ                                        CHJYHK

            rs12960837(+)                                                      rs16951450(+)                        rs621782
            1 ▌▌▌▌▌▌▌                                                          1 ▌▌▌▌▌▌▌                            1 ▌▌▌▌▌▌▌
            CHJYHKSXZ                                                          CHJYHKSXZ                            CHJYHK

**− 🔊 ❓ Ensembl genes**

Clear highlighting                                                                                                 Update Image

**Please restrict your query using criteria below**

**Dataset**
Allele frequencies

**Filters**
Genesymbol : %BRCA2%

**Attributes**
RSID
Chrm
Pos
KHS
HER
STS
XHS
ZUL
ASW
CEU
CHB
CHD
GIH
JPT
LWK
MEX
MKK
TSI
YRI

⊞ Marker

⊞ Location

⊟ Associated Genes
☑ Identifiers
    Association count
    Accession
    Gene id
    Genesymbol        %BRCA2%
    Unigeneid
    Description
    Ncbigeneid

☐ Location
    Location        3
    Distance

# AFRIVARIANTS

## African Mutations and Polymorphisms Database

**Dataset**

allele_frequencies

**Filters**

[None selected]

**Attributes**

Chrm
Pos
KHS
HER
STS
XHS
ZUL
ASW
CEU
CHB
CHD
GIH
JPT
LWK
MEX
MKK
TSI
YRI

Export all results to   [File ▾] [TSV ▾] ☐ Unique results only   [Go]

Email notification to [_____]

View   [10 ▾] rows as [HTML ▾] ☐ Unique results only

| Chrm | Pos | KHS | HER | STS | XHS | ZUL | ASW | CEU | CHB | CHD | GIH | JPT | LWK | MEX | MKK | TSI | YRI |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 84647761 | 0.375 | 0.540 | 0.480 | 0.446 | 0.471 | 0.480 | 0.549 | 0.590 | 0.641 | 0.500 | 0.583 | 0.506 | 0.580 | 0.500 | 0.557 | 0.518 |
| 5 | 156323558 | 1.000 | 0.920 | 0.840 | 0.893 | 0.912 | 0.847 | 1.000 | 0.964 | 0.941 | 1.000 | 1.000 | 0.883 | 0.980 | 0.969 | 1.000 | 0.827 |
| 5 | 158662525 | 0.458 | 0.160 | 0.240 | 0.268 | 0.176 | 0.250 | 0.259 | 0.247 | 0.282 | 0.153 | 0.190 | 0.228 | 0.300 | 0.378 | 0.364 | 0.195 |
| 9 | 22966592 | 0.917 | 0.960 | 0.920 | 0.929 | 0.941 | 0.939 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 0.922 | 1.000 | 0.982 | 1.000 | 0.903 |
| 11 | 26257545 | 0.479 | 0.620 | 0.540 | 0.589 | 0.588 | 0.643 | 0.833 | 0.886 | 0.882 | 0.920 | 0.917 | 0.528 | 0.888 | 0.631 | 0.837 | 0.545 |
| 13 | 76334809 | 0.188 | 0.380 | 0.340 | 0.446 | 0.353 | 0.378 | 0.281 | 0.223 | 0.271 | 0.432 | 0.161 | 0.356 | 0.350 | 0.427 | 0.312 | 0.385 |
| 2 | 224934756 | 0.042 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.247 | 0.212 | 0.000 | 0.286 | 0.000 | 0.110 | 0.017 | 0.000 | 0.009 |
| 8 | 119481632 | 0.188 | 0.080 | 0.080 | 0.125 | 0.088 | 0.133 | 0.196 | 0.000 | 0.000 | 0.051 | 0.000 | 0.133 | 0.120 | 0.101 | 0.176 | 0.076 |
| 6 | 169694703 | 1.000 | 1.000 | 1.000 | 0.982 | 1.000 | 0.949 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 0.978 | 1.000 | 0.989 | 1.000 | 0.965 |
| 7 | 12953585 | 0.021 | 0.120 | 0.080 | 0.143 | 0.000 | 0.245 | 0.558 | 0.476 | 0.512 | 0.453 | 0.572 | 0.100 | 0.490 | 0.112 | 0.528 | 0.116 |

# AFRIVARIANTS

African Mutations and Polymorphisms Database

Home    Biomart    Gbrowse    Upload data    Documentation    Publications

## AfMap - dev version

### Showing 500 kbp from chr18, positions 411,072 to 911,070

**Instructions**

**Searching:** Search using a sequence name, gene name, locus, or other landmark. The wildcard character * is allowed.
**Navigation:** Click one of the rulers to center on a location, or click and drag to select a region. Use the Scroll/Zoom buttons to change magnification and position.

**Examples:** Chr18:411072..911070, SNP:rs6870660, NM_153254, BRCA2, 5q31, ENm010, gwa*, PARK3.

[Bookmark this] [Upload your own data] [Hide banner] [Share these tracks] [Link to Image] [SNP genotype data] [HapMap LD Data] [High-res Image] [Help] [Reset]

**Search**

**Landmark or Region:**

Chr18:411072..911070    [Search]

**Data Source**

AfMap - dev version

**Reports & Analysis:**

Annotate LD Plot    [Configure...]    [Go]

**Scroll/Zoom:** ◄◄ ◄ ─ Show 500 kbp ─ ► ►►  ☐ Flip

**Overview**

chr18
0M    10M    20M    30M    40M    50M    60M    70M

**Region**

chr18
0M  0.1M 0.2M 0.3M 0.4M 0.5M 0.6M 0.7M 0.8M 0.9M 1M  1.1M 1.2M 1.3M 1.4M 1.5M 1.6M 1.7M 1.8M 1.9M 2M

**Details**

500k    600k    700k    800k    900k

Genotyped SNPs

funded by ...

# What's next?

- more samples

- more variants

- more data

# BioQ (http://bioq.saclab.net/)

**Change Database:** [ 1000 Genomes Phase 1 Analysis ⟷ ]

# 1000 Genomes Phase 1 Analysis

Frequency and QC data from the Phase 1 Analysis (July 2012) release of the 1000 Genomes project. This release contains an integrated set of variant calls and phased genotypes including SNPS and short insertions and deletions based on low coverage and exome sequencing. See the 1000 Genomes Announcements Page for more information.

Documentation | Downloads

| Simple Query | Advanced Query |

1000 Genomes Sample [ ASW: African American/Southwest USA ⟷ ]

| | | | |
|---|---|---|---|
| Populations | Sample Counts | ASW Samples | ASW Pedigrees |
| ASW DNA | VCF Sites | ASW Allele Freqs | ASW HWE |
| Downloads | | | |

[ Execute Query ]

Max Rows [ 1000 ⟷ ]

**Enter Genomic Features** ⊙

You may enter SNPs, genes, regions and other types of genomic features. When possible, the selected query will pertain only to these features - see *Related genomic features* in the information box for the selected query. Regions must use GRCh37 coordinates. Click *Get/Configure Features* or see the documentation for more information.

# Query Results

Select a row for a detailed view.

LD proxies can be shown for simple queries when a *snp_id* column (dbSNP ID) is present. Unless "Merge dbSNP" is checked, mapping data for proxies will be limited to the original query results.

Add LD Proxies ☐    $r^2$ Threshold [ 0.8 ⇕ ]    HapMap Sample [ ASW: African American (Southwest USA) ⇕ ]    Merge dbSNP ☐

## ASW Allele Freqs

| sites_id | chr | pos_bp | pos_global | snp_id | gene_function_list | ref | alt | maf | het | pvalue_hwe | filter |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 10583 | 1000010583 | 58108140 | DDX11L1/nearGene-5 | G | A | 0.0983607 | 0.196721 | 1 | PASS |
| 2 | 1 | 10611 | 1000010611 | 189107123 | DDX11L1/nearGene-5 | C | G | 0.0163934 | 0.0327869 | 1 | PASS |
| 3 | 1 | 13302 | 1000013302 | 180734498 | DDX11L1/ncRNA | C | T | 0.180328 | 0.295082 | 1 | PASS |
| 4 | 1 | 13327 | 1000013327 | 144762171 | DDX11L1/ncRNA | G | C | 0.0163934 | 0.0327869 | 1 | PASS |
| 5 | 1 | 13957 | 1000013957 | | | TC | T | 0.00819672 | 0.0163934 | 1 | PASS |
| 6 | 1 | 13980 | 1000013980 | 151276478 | DDX11L1/ncRNA,WASH7P/intron | T | C | 0.00819672 | 0.0163934 | 1 | PASS |
| 7 | 1 | 30923 | 1000030923 | 140337953 | MIR1302-2/nearGene-3,WASH7P/nearGene-5 | G | T | 0.418033 | 0.409836 | 0.290336 | PASS |
| 8 | 1 | 46402 | 1000046402 | | | C | CTGT | 0.0163934 | 0.0327869 | 1 | PASS |
| 9 | 1 | 47190 | 1000047190 | | | G | GA | 0.0491803 | 0.0983607 | 1 | PASS |
| 10 | 1 | 51476 | 1000051476 | 187298206 | | T | C | 0.00819672 | 0.0163934 | 1 | PASS |
| 11 | 1 | 51479 | 1000051479 | 116400033 | | T | A | 0.0901639 | 0.180328 | 1 | PASS |
| 12 | 1 | 51914 | 1000051914 | 190452223 | | T | G | 0 | 0 | 1 | PASS |
| 13 | 1 | 51935 | 1000051935 | 181754315 | | C | T | 0 | 0 | 1 | PASS |
| 14 | 1 | 51954 | 1000051954 | 185832753 | | G | C | 0 | 0 | 1 | PASS |
| 15 | 1 | 52058 | 1000052058 | 62637813 | | G | C | 0.0491803 | 0.0983607 | 1 | PASS |
| 16 | 1 | 52144 | 1000052144 | 190291950 | | T | A | 0.00819672 | 0.0163934 | 1 | PASS |
| 17 | 1 | 52185 | 1000052185 | | | TTAA | T | 0 | 0 | 1 | PASS |
| 18 | 1 | 52238 | 1000052238 | 150021059 | | T | G | 0.262295 | 0.360656 | 0.527478 | PASS |
| 19 | 1 | 53234 | 1000053234 | | | CAT | C | 0.0245902 | 0.0491803 | 1 | PASS |
| 20 | 1 | 54353 | 1000054353 | 140052487 | | C | A | 0.00819672 | 0.0163934 | 1 | PASS |

# NoSQL based storage



```
mongoDB
{
    "name": "NoSQL vs RDBMS",
    "datePosted" : Date("2012-01-30T12:00:00.3Z"),
    "comments" : [
        {
            "userName" : "Bob",
            "value" : "This blog rocks",
            "datePosted" : Date("2012-02-01T12:00:00.3Z")
        },
        {
            "userName" : "Bob",
            "value" : "Exactly what I was looking for!",
            "datePosted" : Date("2012-02-02T12:00:00.3Z")
        },
        {
            "userName" : "Chris",
            "value" : "I'm a hater, too generalized",
            "datePosted" : Date("2012-02-02T12:00:00.3Z"),
            "email" : "chris@hater.com"
        }
    ]
}
```

# Advice

- Know the datasets
- Understand what is going on behind the scenes
  - database queries
  - where the annotations are from
  - how measures were calculated
- Make it reproducable