

# The 1000 Genomes Project: A Tutorial

Laura Clarke and Paul Flicek  
EBI Training Room  
24<sup>th</sup> January 2012



# Agenda

- Introduction
- Brief History of the 1000 Genomes Project, data and analysis
- The Raw Data and FTP site
- Finding Data
  - Exercise: Finding Data
- The Website and Browser
  - Exercise: Using the Browser
- The 1000 Genomes Tools
  - Exercise: Using the Tools



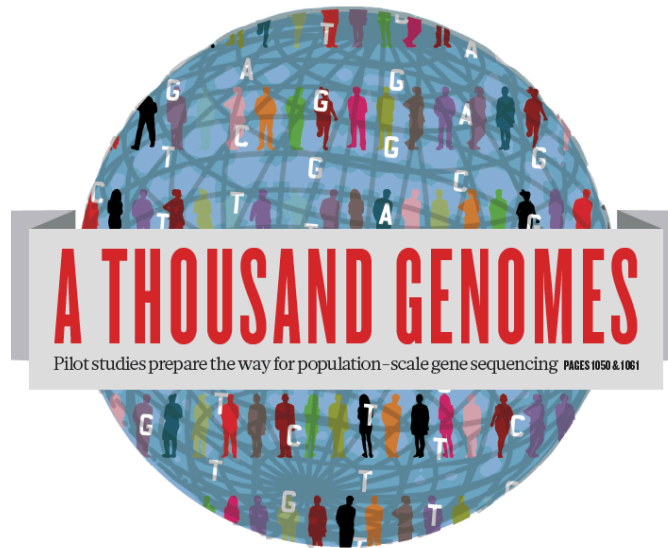
# Glossary

- **Pilot** : The 1000 Genomes project ran a pilot study between 2008 and 2010
- **Phase 1**: The initial round of exome and low coverage sequencing of 1000 individuals
- **Phase 2**: Expanded sequencing of 1700 individuals and method improvement
- **SAM/BAM**: Sequence Alignment/Map Format, an alignment format
- **VCF**: Variant Call Format, a variant format



How are you using 1000 genomes data?





# The 1000 Genomes Project: An Introduction

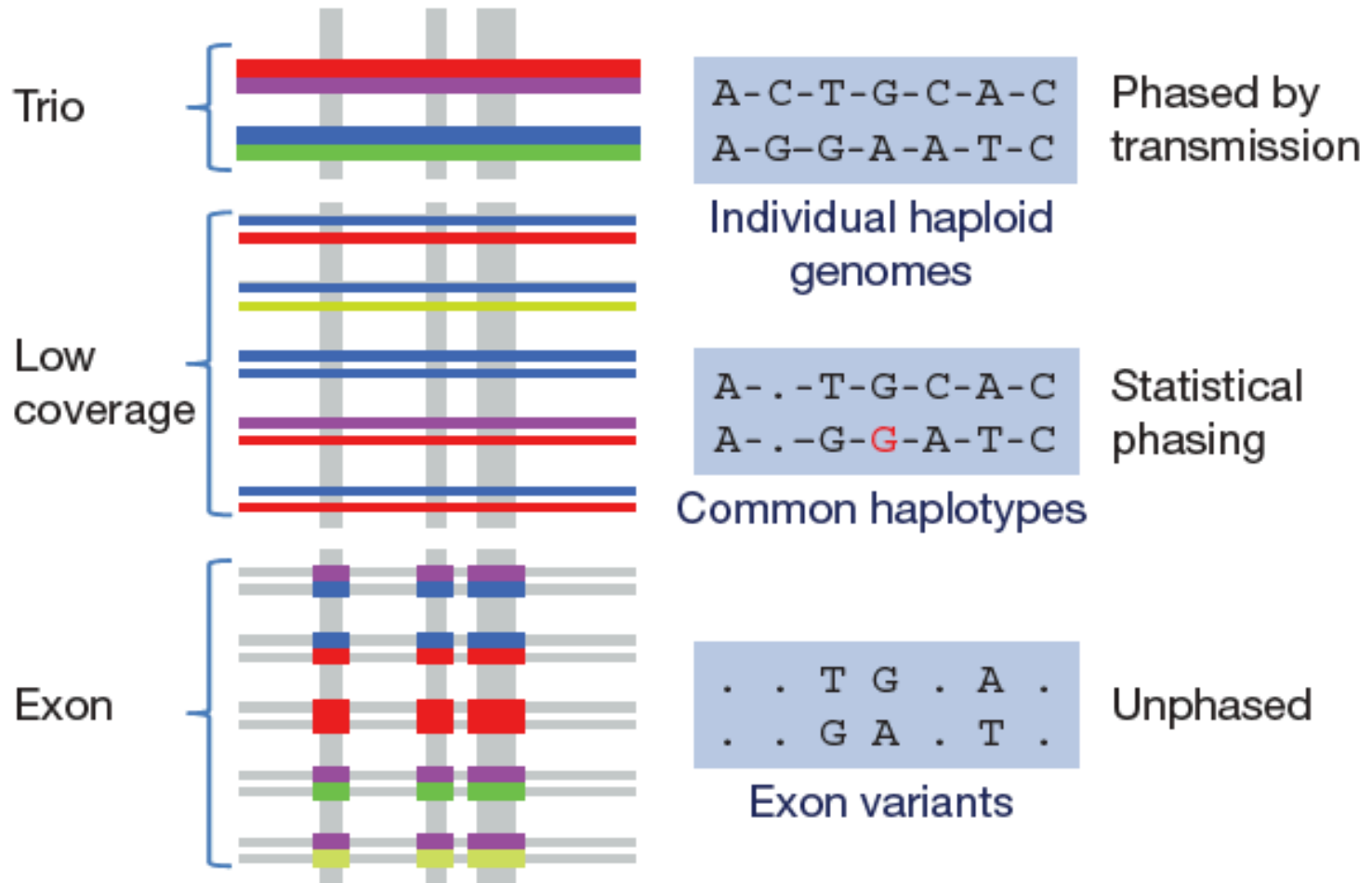


# The 1000 Genomes Project

- International project to construct a foundational data set for human genetics
  - Discover virtually all common human variations by investigating many genomes at the base pair level
  - Consortium with multiple centers, platforms, funders
- Aims
  - Discover population level human genetic variations of all types (95% of variation  $> 1\%$  frequency)
  - Define haplotype structure in the human genome
  - Develop sequence analysis methods, tools, and other reagents that can be transferred to other sequencing projects



# 3 pilot coverage strategies



# Main Project Design

- Based on the result of the pilot project, we decided to collect data on 2,500 samples from 5 continental groupings
  - Whole-genome low coverage data (>4x)
  - Full exome data at deep coverage (> 20x)
  - A number of deep coverage genomes to be sequenced, with details to be decided
  - High density genotyping at subsets of sites
- Phase 1 Release Integrated Variant Release has been made.





# Phase I (1,150)

# Phase II (1,721)

# Phase III (2,500)

CDX  
17S



CLM (70T); DNA from  
LCL




CHS (100T); DNA from  
LCL



PUR (70T); DNA from  
Blood



FIN (100S); DNA from  
LCL



GBR (96/100S); DNA from



IBS (84/100T); DNA from  
LCL



GWD



GWD



GWD



GWD (target - 100T); DNA from LCL



CDX (100S); DNA: 17 DNA from Bld, 83 from LCL

KHV (82/100) - 15 trios; DNA Bld



ACB (28/79T) - 14 trios; DNA Bld



PEL (70T); DNA from Blood



PJL (target - 100T); DNA from Blood



GIH vs. Sindhi (target - 100T)



Tamil (target -



Sri Lankan (target - 100T)



Bengalee (target - 100T)



Nigeria (target - 100T); DNA from  
Sierra Leone (target - 100T); DNA from LCL



MAB (target - 100T); DNA from  
LCL

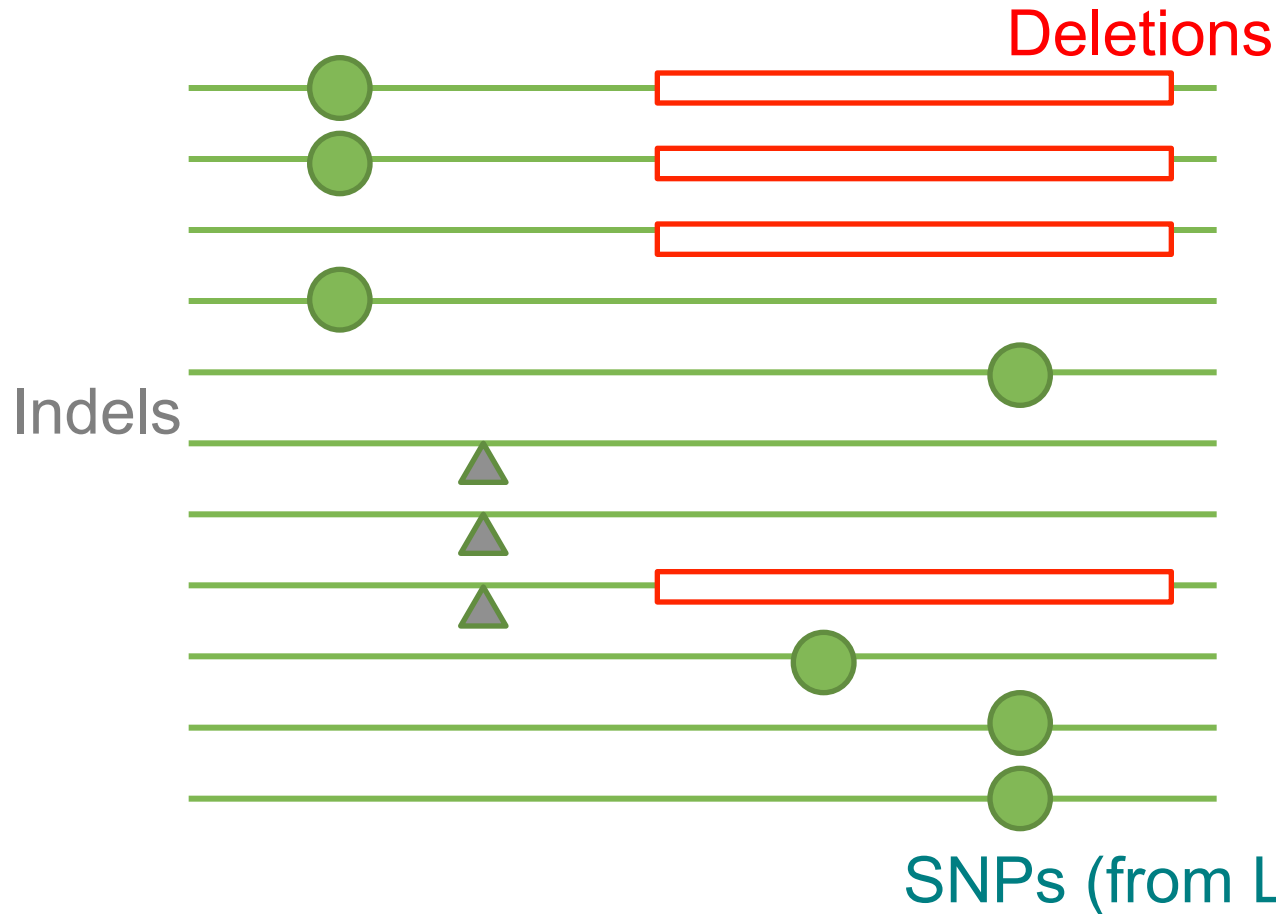


AJM (target - 80T); DNA from Bld




# Phase 1 analysis goal: an integrated view of human variations

- Reconstruct haplotypes including all variant types, using all datasets

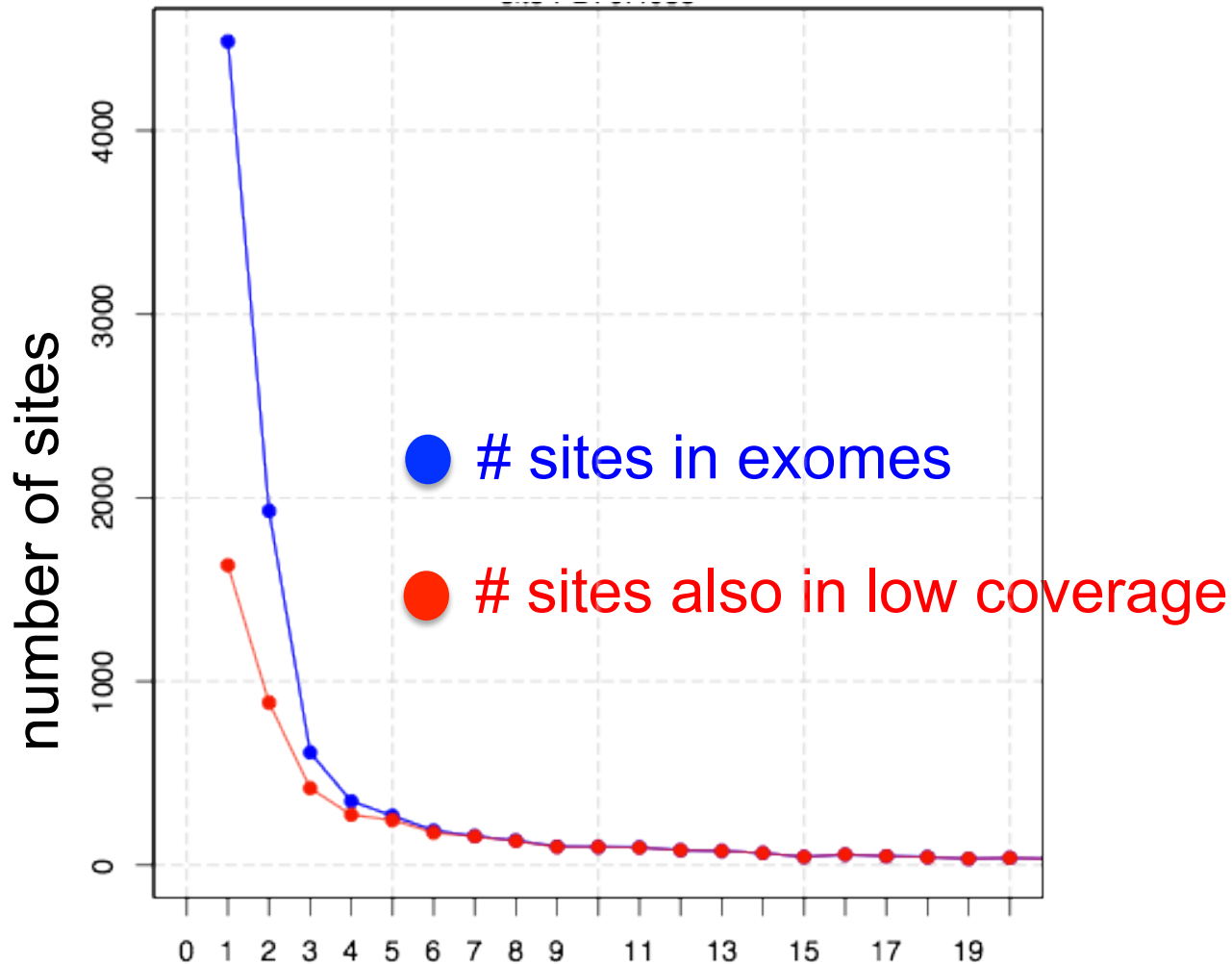


Goncalo Abecasis

SNPs (from LC, EX, OMNI)



# Deep coverage exome data is more sensitive to low-frequency variants



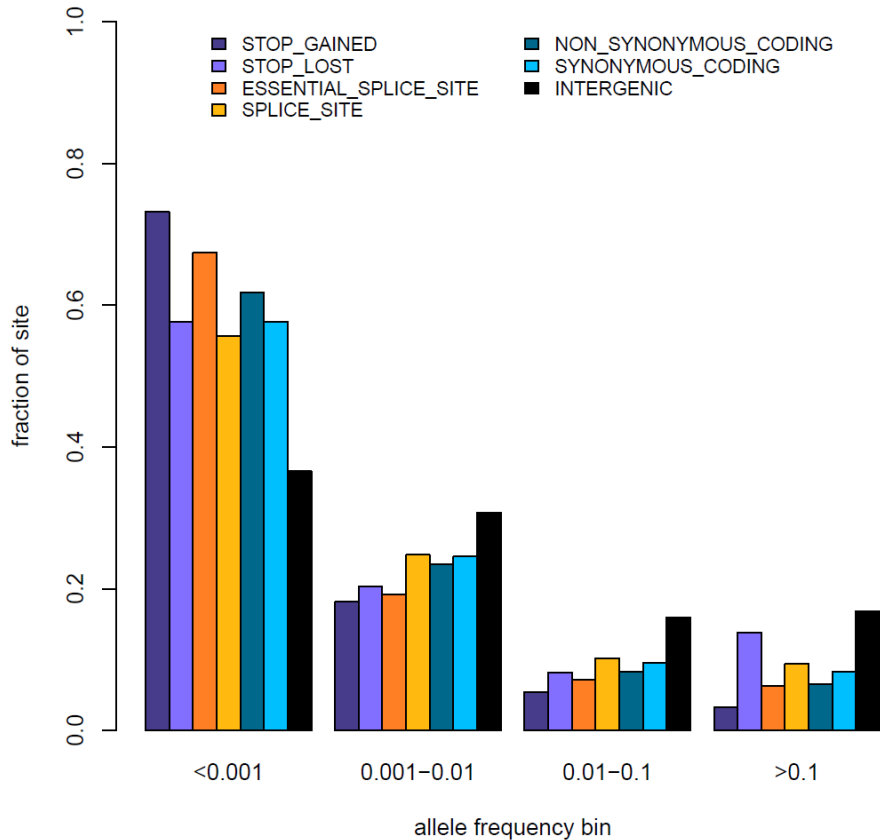
Allele count in 766 exomes (chr. 20, exons only)

Erik Garrison

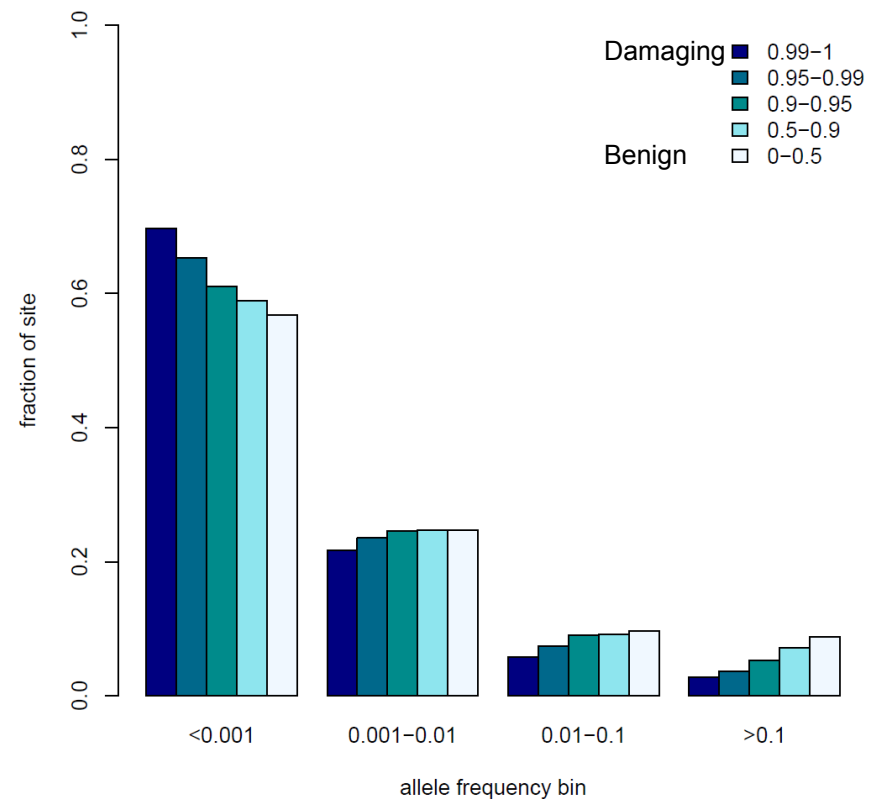


# Newly discovered SNPs are mostly at low frequency and enriched for functional variants

## Functional category

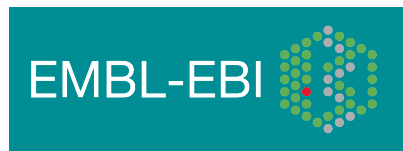


## Non-synonymous: Condel score



Presentation on using the data for GWAS by Brian Howie

Enza Colonna, Yuan Chen, Yali Xue



Fraction of variant sites present in an individual that are NOT already represented in dbSNP

Date	Fraction <u>not</u> in dbSNP
February, 2000	98%
February, 2001	80%
April, 2008	10%
February, 2011	2%
Now	<1%

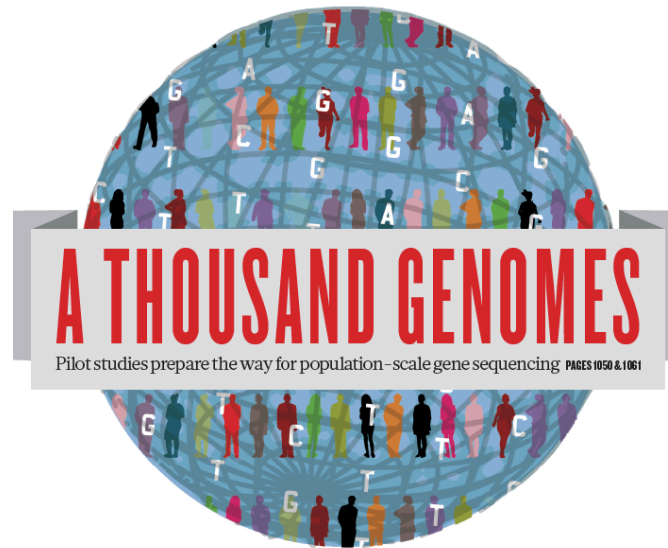
Ryan Poplin, David Altshuler



# 1000 Genomes Project: Present & Future

- First Phase 2 sequence release 14<sup>th</sup> November 2011
- First Phase 2 alignment release February 2012
- First Phase 2 variant site release Summer 2012
  
- Sample collected expected end to June 2012
- Final Phase 3 Sequence release expected December 2012
- 2013 will represent finalization of 1000 genomes analysis results and final data releases





# The 1000 Genomes Project: A Brief History of Data and Analysis Results



# Timeline

- **September 2007:** 1000 Genomes project formally proposed Cambridge, UK
- **April 2008:** First Submission of Data to the Short Read Archive.
- **May 2008:** First public data release.
- **October 2008:** SAM/BAM Format Defined.
- **December 2008:** First High Coverage Variants Released.
- **December 2008:** First 1000 genomes browser released
- **May 2009:** First Indel Calls released.
- **July 2009:** VCF Format defined
- **August 2009:** First Large Scale Deletions released.
- **December 2009:** First Main Project Sequence Data Released.
- **March 2010:** Low Coverage Pilot Variant Release made
- **July 2010:** Phased genotypes for 159 Individuals released.
- **October 2010:** A Map of Human Variation from population scale sequencing is published in Nature.
- **January 2011:** Final Phase 1 Low coverage alignments are released
- **May 2011:** @1000genomes appears on Twitter
- **May 2011:** First Variant Release made on more than 1000 individuals
- **October 2011:** Phase 1 integrated variant release made





# Sequencing Data

- The Project contains data from 3 different providers and multiple platforms

Platform	Min Read Length (bp)	Max Read Length (bp)
454 Roche GS FLX Titanium	70	400
Illumina GA	30	81
Illumina GA II	26	160
Illumina HiSeq	50	102
ABI Solid System 2.0	25	35
ABI Solid System 2.5	50	50
ABI Solid System 3.0	50	50

# Alignment Data

- The project has made more than 10 releases of Alignment Data
- Pilot Project
  - Aligned to NCBI36
  - Maq and Corona
  - Base Quality Recalibration done
- Phase 1
  - Aligned to GRCh37
  - BWA and Bfast
  - Indel Realignment
- Phase 2
  - Aligned to extended GRCh37
  - Improvements to Base Quality Recalibration



# Variant Calling

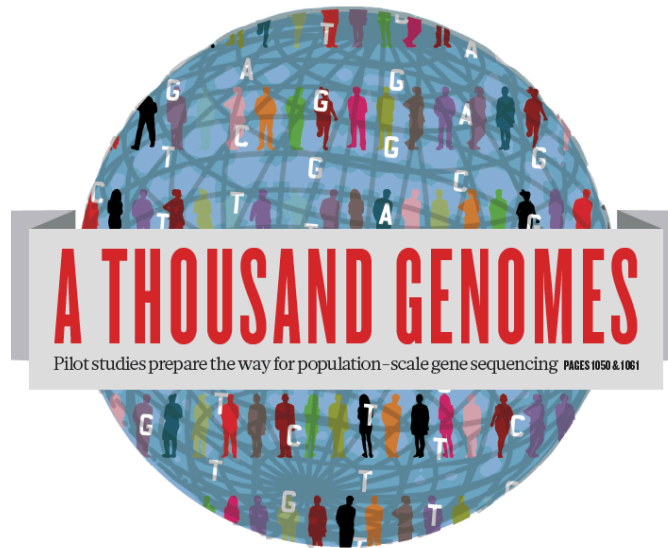
- Early call sets used a single variant caller
- Intersect approach developed during pilot
- Variant Quality Score Recalibration (VQSR) developed for Phase 1
- Genotype Likelihoods assigned to help with genotype calling
- Integrated genotype calling based on individual variant call sets
- Phase 2 looks to improve site discover and improve integration



# Data Availability

- FTP site: <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/>
  - Raw Data Files
- Web site: <http://www.1000genomes.org>
  - Release Announcements
  - Documentation
- Ensembl Style Browser: <http://browser.1000genomes.org>
  - Browse 1000 Genomes variants in Genomic Context
  - Variant Effect Predictor
  - Data Slicer
  - Other Tools





# The 1000 Genomes Project: The Raw Data











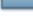







ftp://ftp.1000genomes.ebi.ac.uk

ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp

Index of ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/

 Up to higher level directory

Name	Size	Last Modified
 CHANGELOG	118 KB	05/01/2012 5/01/2012 12:40:00
 README.alignment_data	12 KB	26/01/2011 26/01/2011 12:00:00
 README.ftp_structure	9 KB	04/04/2011 4/04/2011 12:00:00
 README.pilot_data	3 KB	14/07/2011 14/07/2011 12:00:00
 README.populations	2 KB	18/02/2010 18/02/2010 12:00:00
 README.sequence_data	7 KB	23/07/2011 23/07/2011 19:03:00
 alignment_indices		14/07/2011 14/07/2011 10:53:00
 changelog_details		05/01/2012 05/01/2012 12:40:00
 current.tree	29933 KB	05/01/2012 05/01/2012 12:37:00
 data		04/07/2012 04/07/2012 18:50:00
 phase1		14/07/2011 14/07/2011 14:03:00
 pilot_data		27/07/2011 27/07/2011 12:00:00
 release		12/10/2011 12/10/2011 13:18:00
 sequence.index	27185 KB	20/12/2011 20/12/2011 12:26:00
 sequence_indices		14/11/2011 14/11/2011 10:10:00
 technical		13/12/2011 13/12/2011 10:05:00

Documentation

Raw Data

Phase 1 Data

Pilot Data

Release Data

Technical Data



# The FTP Site: Data

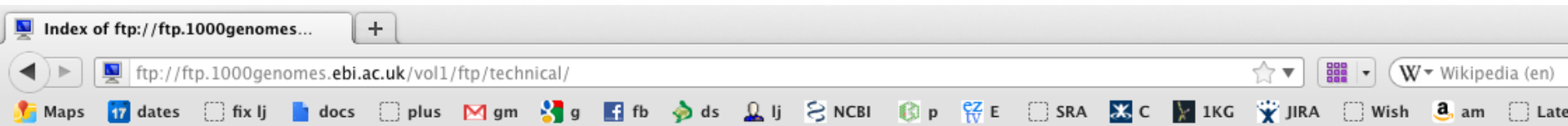
The screenshot shows the index of an FTP site at `ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data/`. The index lists folders for samples HG00104 through HG00131. Each folder is associated with a date and time. Three green boxes with red arrows point to specific folders:

- Sample Level Files** points to HG00109.
- sequence\_read** points to HG00110.
- alignment** points to HG00111.

Sample ID	Date	Time
HG00104	14/12/2011	14/12/2011 12:06:00
HG00105	13/12/2011	13/12/2011 12:45:00
HG00106	13/12/2011	13/12/2011 12:45:00
HG00107	13/12/2011	13/12/2011 12:40:00
HG00108	13/12/2011	13/12/2011 12:43:00
HG00109	13/12/2011	13/12/2011 12:43:00
HG00110	13/12/2011	13/12/2011 12:43:00
HG00111	13/12/2011	13/12/2011 12:36:00
HG00112	13/12/2011	13/12/2011 12:41:00
HG00113	13/12/2011	13/12/2011 12:41:00
HG00114	13/12/2011	13/12/2011 12:41:00
HG00115	13/12/2011	13/12/2011 12:43:00
HG00116	13/12/2011	13/12/2011 12:44:00
HG00117	13/12/2011	13/12/2011 12:38:00
HG00118	13/12/2011	13/12/2011 12:43:00
HG00119	13/12/2011	13/12/2011 12:37:00
HG00120	13/12/2011	13/12/2011 12:45:00
HG00121	13/12/2011	13/12/2011 12:43:00
HG00122	13/12/2011	13/12/2011 12:44:00
HG00123	13/12/2011	13/12/2011 12:36:00
HG00124	13/12/2011	13/12/2011 12:39:00
HG00125	13/12/2011	13/12/2011 12:39:00
HG00126	14/12/2011	14/12/2011 12:06:00
HG00127	14/12/2011	14/12/2011 12:06:00
HG00128	13/12/2011	13/12/2011 12:46:00
HG00129	13/12/2011	13/12/2011 12:44:00
HG00130	13/12/2011	13/12/2011 12:44:00
HG00131	13/12/2011	13/12/2011 12:44:00



# FTP Site: Technical



## Index of ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/

[Up to higher level directory](#)

Name	Size	Last Modified
<a href="#">README.reference</a>	1 KB	12/10/2009 12/10/2009 12 :00:00
<a href="#">browser</a>		19/12/2011 19/12/2011 3 :50:00
<a href="#">method_development</a>		06/06/2011 6/06/2011 12 :00:00
<a href="#">ncbi_varpipe_data</a>		
<a href="#">other_exome_alignments.alignment_indices</a>		20/07/2011 20/07/2011 12 :00:00
<a href="#">pilot2_high_cov_GRCh37_bams</a>		11/01/2012 11/01/2012 5 :56:00
<a href="#">pilot3_exon_targetted_GRCh37_bams</a>		
<a href="#">qc</a>		
<a href="#">reference</a>		
<a href="#">retired_reference</a>		
<a href="#">simulations</a>		04/05/2010 4/05/2010 12 :00:00
<a href="#">supporting</a>		21/12/2009 21/12/2009 12 :00:00
<a href="#">working</a>		17/01/2012 17/01/2012 4 :07:00

Alternative Alignments

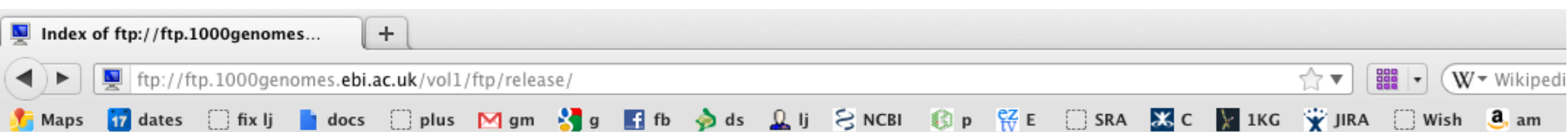
Reference Data Sets

Experimental Data





# FTP Site: Release



## Index of ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/

[Up to higher level directory](#)

Name	Size	Last Modified
<a href="#">2008_12</a>		21/02/2009 21:00:00
<a href="#">2009_02</a>		21/02/2009 21:00:00
<a href="#">2009_04</a>		07/05/2009 7:00:00
<a href="#">2009_05</a>		08/06/2009 8:00:00
<a href="#">2009_08</a>		10/08/2009 10:00:00
<a href="#">20100804</a>		
<a href="#">20101123</a>		
<a href="#">2010_11</a>		16/02/2011 16:00:00
<a href="#">20110521</a>		16/12/2011 16:09:00

Older Release Dirs

Sequence Index Dates

<a href="#">20110521.sequence.index</a>	23693 KB	19/07/2011 19:00:00
<a href="#">20110521.sequence.index.exome.stats</a>	48 KB	19/07/2011 19:00:00
<a href="#">20110521.sequence.index.low_coverage.stats</a>	53 KB	21/05/2011 21:00:00
<a href="#">20110521_20110719.exome.stats.csv</a>	2 KB	19/07/2011 19:00:00
<a href="#">20110521_20110719.low_coverage.stats.csv</a>	2 KB	19/07/2011 19:00:00
<a href="#">20110719.sequence.index</a>	23961 KB	19/07/2011 19:00:00
<a href="#">20110719.sequence.index.exome.stats</a>	52 KB	10/10/2011 10:10:00
<a href="#">20110719.sequence.index.low_coverage.stats</a>	54 KB	10/10/2011 10:13:00
<a href="#">20110719_20110920.exome.stats.csv</a>	1 KB	10/10/2011 10:45:00
<a href="#">20110719_20110920.low_coverage.stats.csv</a>	2 KB	10/10/2011 10:45:00



# FTP Site: Pilot Data



## Index of ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/pilot\_data/

[Up to higher level directory](#)

Name	Size	Last Modified
<a href="#">README.alignment.index</a>	2 KB	26/08/2009 26/08/2009 12:00:00
<a href="#">README.bas</a>	3 KB	27/08/2009 27/08/2009 12:00:00
<a href="#">README.sequence.index</a>	2 KB	22/07/2009 22/07/2009 12:00:00
<a href="#">SRP000031.sequence.index</a>	7365 KB	12/07/2010 12/07/2010 12:00:00
<a href="#">SRP000032.sequence.index</a>	2181 KB	12/07/2010 12/07/2010 12:00:00
<a href="#">SRP000033.sequence.index</a>	480 KB	12/07/2010 12/07/2010 12:00:00
<a href="#">data</a>		
<a href="#">paper_data_sets</a>		03/02/2011 3/02/2011 12:00:00
<a href="#">pilot_data.alignment.index</a>	795 KB	06/05/2010 6/05/2010 12:00:00
<a href="#">pilot_data.alignment.index.bas.gz</a>	1740 KB	14/06/2010 14/06/2010 12:00:00
<a href="#">pilot_data.sequence.index</a>	10025 KB	12/07/2010 12/07/2010 12:00:00
<a href="#">release</a>		20/07/2010 20/07/2010 12:00:00
<a href="#">technical</a>		29/07/2010 29/07/2010 12:00:00

Final Paper Data



# Data formats and key tools

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 25 no. 16 2009, pages 2078–2079  
doi:10.1093/bioinformatics/btp352

Sequence analysis

## The Sequence Alignment/Map format and SAMtools

Heng Li<sup>1,†</sup>, Bob Handsaker<sup>2,†</sup>, Alec Wysoker<sup>2</sup>, Tim Fennell<sup>2</sup>, Jue Ruan<sup>3</sup>, Nils Homer<sup>4</sup>, Gabor Marth<sup>5</sup>, Goncalo Abecasis<sup>6</sup>, Richard Durbin<sup>1,\*</sup> and 1000 Genome Project Data Processing Subgroup<sup>7</sup>

<sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge, CB10 1SA, UK, <sup>2</sup>Broad Institute of MIT and Harvard, Cambridge, MA 02141, USA, <sup>3</sup>Beijing Institute of Genomics, Chinese Academy of Science, Beijing 100029, China, <sup>4</sup>Department of Computer Science, University of California Los Angeles, Los Angeles, CA 90095, <sup>5</sup>Department of Biology, Boston College, Chestnut Hill, MA 02467, <sup>6</sup>Center for Statistical Genetics, Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA and <sup>7</sup><http://1000genomes.org>

Received on April 28, 2009; revised on May 28, 2009; accepted on May 30, 2009

Advance Access publication June 8, 2009

Associate Editor: Alfonso Valencia

## BAM alignment files

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 27 no. 15 2011, pages 2156–2158  
doi:10.1093/bioinformatics/btr330

Sequence analysis

Advance Access publication June 7, 2011

## The variant call format and VCFtools

Petr Danecek<sup>1,†</sup>, Adam Auton<sup>2,†</sup>, Goncalo Abecasis<sup>3</sup>, Cornelis A. Albers<sup>1</sup>, Eric Banks<sup>4</sup>, Mark A. DePristo<sup>4</sup>, Robert E. Handsaker<sup>4</sup>, Gerton Lunter<sup>2</sup>, Gabor T. Marth<sup>5</sup>, Stephen T. Sherry<sup>6</sup>, Gilean McVean<sup>2,7</sup>, Richard Durbin<sup>1,\*</sup> and 1000 Genomes Project Analysis Group<sup>†</sup>

<sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge CB10 1SA, <sup>2</sup>Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK, <sup>3</sup>Center for Statistical Genetics, Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, <sup>4</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA 02141, <sup>5</sup>Department of Biology, Boston College, MA 02467, <sup>6</sup>National Institutes of Health National Center for Biotechnology Information, MD 20894, USA and <sup>7</sup>Department of Statistics, University of Oxford, Oxford OX1 3TG, UK

Associate Editor: John Quackenbush

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 27 no. 5 2011, pages 718–719  
doi:10.1093/bioinformatics/btq671

Sequence analysis

Advance Access publication January 5, 2011

## Tabix: fast retrieval of sequence features from generic TAB-delimited files

Heng Li

Program in Medical Population Genetics, The Broad Institute of Harvard and MIT, Cambridge, MA 02142, USA

Associate Editor: Dmitrij Frishman

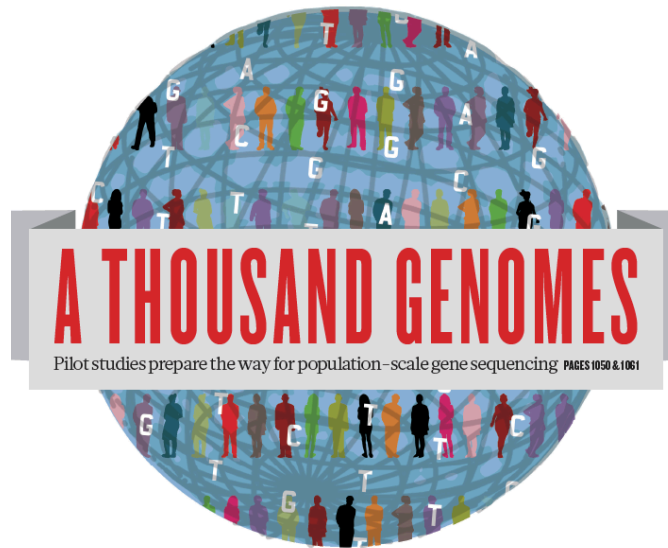
## VCF variant files

# All indexed for fast retrieval



EMBL-EBI





# The 1000 Genomes Project: Finding Data



# Finding Data

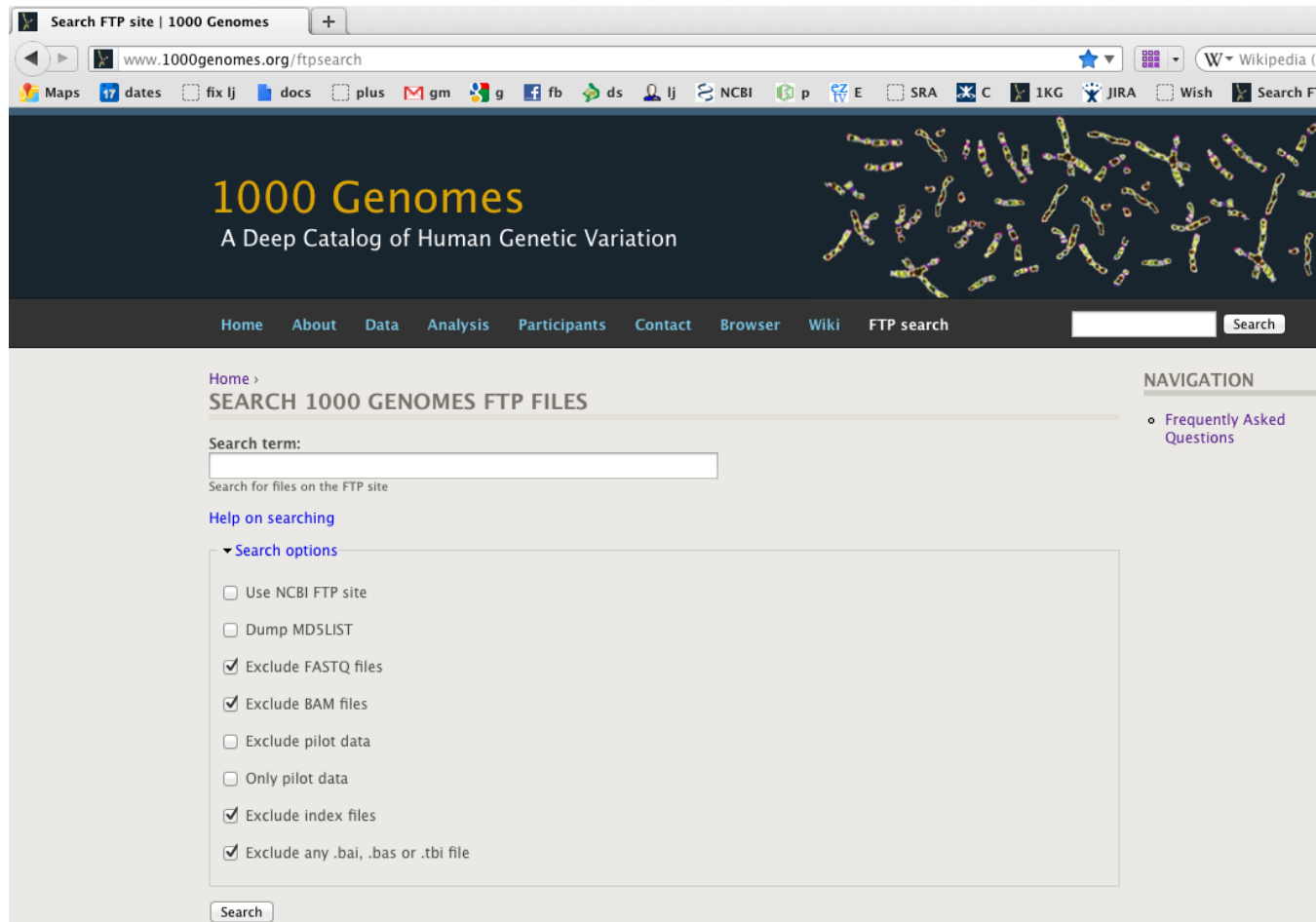
- Current.tree file
- <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/current.tree>

```
ftp://ftp.1000ge...ftp/current.tree
ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/current.tree
Wikipedia (en)
Maps 17 dates fix lj docs plus gm g fb ds lj NCBI p E SRA C 1KG JIRA Wish am Later
ftp directory 403 Tue Dec 20 16:11:25 2011
ftp/README.ftp_structure file 8408 Mon Apr 4 14:52:52 2011 2a59a3feb2540c113e10877f3ef1efe5
ftp/README.populations file 1506 Wed Jan 11 15:12:44 2012 f7c588af82396013c1737e66e58f0f05
ftp/CHANGELOG file 122151 Sat Jan 14 23:51:50 2012 ecaa9b1e0a6860cd76b1545e84ff3403
ftp/sequence.index file 27836681 Tue Dec 20 12:26:18 2011 b25557458f6c468bd13d025c17461bab
ftp/README.alignment_data file 11632 Wed Jan 26 16:22:41 2011 7528e9f4ba8c6b085e6d29c7546fc684
ftp/README.sequence_data file 6548 Sat Jul 23 22:03:54 2011 b5cfc5784ebf06998f883c629c1c0ba0
ftp/README.pilot_data file 2082 Fri Aug 14 13:58:10 2009 977fe3983de2131f9e28f6f0036b31d9
ftp/phasel directory 412 Wed Dec 14 16:03:36 2011
ftp/phasel/phasel.exome.alignment.index.HsMetrics.stats file 293 Wed Dec 14 15:53:53 2011 lebf793046daadd7ff67ececblb5361f
ftp/phasel/phasel.exome.alignment.index file 397947 Wed Dec 14 15:53:52 2011 2891d1fffe08acf3ee99c88cb42d130d
ftp/phasel/phasel.alignment.index.bas.gz file 5115518 Wed Dec 14 15:53:23 2011 2b4e1edb78f617ebfaf5087536d80f95
ftp/phasel/phasel.alignment.index file 8850348 Wed Dec 14 15:53:22 2011 ea3423858ec976af1e17839cd334c164
ftp/phasel/phasel.exome.alignment.index.bas.gz file 423691 Wed Dec 14 15:53:52 2011 7a56f22d28e860fbc65b71d1013717ae
ftp/phasel/phasel.exome.alignment.index.HsMetrics.gz file 143893 Wed Dec 14 15:53:53 2011 93ba34ab86e9c42198919d128acc13b7
ftp/phasel/phasel.exome.alignment.index_stats.csv file 715 Wed Dec 14 15:53:53 2011 376ea20314a94399cab99c723e1d974c
ftp/phasel/technical/ncbi_varpipe_data directory 137 Wed Dec 14 16:16:31 2011
ftp/phasel/technical/ncbi_varpipe_data/alignment/ncbi.20100804.alignment.summary file 39866 Wed Dec 14 16:13:58 2011 df4676c95ed2cc6f9cd4c9e24a66bbe8
ftp/phasel/technical/ncbi_varpipe_data/alignment/ncbi.20100804.alignment.index file 159169 Wed Dec 14 16:13:58 2011 a9bc22ace39cb0bcd0bf35f2ee807bbc
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA12004 directory 308 Tue Dec 13 12:16:47 2011
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA12004/NA12004.chrom20.ILLUMINA.mosaik.CEU.low_coverage.20100804.bam file 238645793 Thu Apr 14 15:24
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA12004/NA12004.ILLUMINA.mosaik.CEU.low_coverage.20100804.bam.bai file 7899352 Wed Oct 27 18:31:23 2010
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA12004.chrom20.ILLUMINA.mosaik.CEU.low_coverage.20100804.bam.bai file 166624 Thu Apr 14 15:24
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA12004/NA12004.ILLUMINA.mosaik.CEU.low_coverage.20100804.bam file 11091314322 Wed Oct 27 18:31:24 2010
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA18486 directory 308 Tue Dec 13 12:25:36 2011
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA18486/NA18486.ILLUMINA.mosaik.YRI.low_coverage.20101123.bam.bai file 8418040 Tue Jan 25 22:46:53 2011
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA18486/NA18486.ILLUMINA.mosaik.YRI.low_coverage.20101123.bam file 29068330549 Tue Jan 25 22:46:53 2011
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA18486.chrom20.ILLUMINA.mosaik.YRI.low_coverage.20101123.bam.bai file 176848 Tue Jan 25 22:47
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA18486/NA18486.chrom20.ILLUMINA.mosaik.YRI.low_coverage.20101123.bam file 685641416 Tue Jan 25 22:47
ftp/phasel/technical/ncbi_varpipe_data/alignment/NA12045 directory 604 Tue Dec 13 12:24:58 2011
```



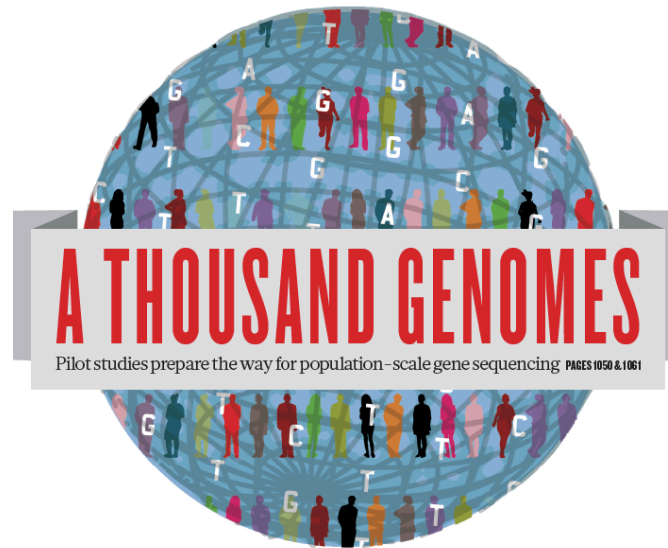
# Finding Data

- FTP search
- <http://www.1000genomes.org/ftpsearch>



The screenshot shows a web browser window with the address bar displaying [www.1000genomes.org/ftpsearch](http://www.1000genomes.org/ftpsearch). The page header features the text "1000 Genomes" in large yellow font, followed by "A Deep Catalog of Human Genetic Variation" in white. Below this is a navigation menu with links for Home, About, Data, Analysis, Participants, Contact, Browser, Wiki, and FTP search. A search bar is located to the right of the menu. The main content area is titled "SEARCH 1000 GENOMES FTP FILES" and includes a "Search term:" input field. Below the input field is a "Search" button. A "Help on searching" link is also present. A "Search options" section is expanded, showing a list of checkboxes: "Use NCBI FTP site", "Dump MDSLIST", "Exclude FASTQ files" (checked), "Exclude BAM files" (checked), "Exclude pilot data", "Only pilot data", "Exclude index files" (checked), and "Exclude any .bai, .bas or .tbi file" (checked). A "Search" button is located at the bottom of the search options section. On the right side of the page, there is a "NAVIGATION" section with a link to "Frequently Asked Questions".





## The 1000 Genomes Project:

Exercise 1: Finding Data on the 1000 genomes ftp site



# Finding Data: Exercise

1. Find what VCF files we have containing genotypes from the Illumina Omni platform.

<http://www.1000genomes.org/ftpsearch>

2. Find the FAQ question which gives you instructions on how to get a sub-section of a VCF file. The Search Box is on the top right hand corner of any website page.

<http://www.1000genomes.org/>





# Finding Data: Exercise

## 1. Finding Omni VCF Files

The screenshot shows the EBI FTP search interface. At the top, there is a navigation bar with links: Home, About, Data, Analysis, Participants, Contact, Browser, Wiki, and FTP search. A search box is located on the right side of the navigation bar. Below the navigation bar, the page title is "SEARCH 1000 GENOMES FTP FILES". The search term "omni\*vcf" is entered in the search box. Below the search box, there is a link "Help on searching" and a section titled "Search options" with a dropdown arrow. The search options are:

- Use NCBI FTP site
- Dump MD5LIST
- Exclude FASTQ files
- Exclude BAM files
- Exclude pilot data
- Only pilot data
- Exclude index files
- Exclude any .bai, .bas or .tbi file

At the bottom of the search options section, there is a "Search" button. On the right side of the page, there is a user profile for "LAURA@EBI.AC.UK" with a list of links: My account, Create content, List content, List users, Manage files, Log out, and Frequently Asked Questions.

Two red arrows are present: one pointing from the left edge of the image to the search box, and another pointing from the right edge of the image to the "Exclude pilot data" checkbox.



# Finding Data: Exercise

Home About Data Analysis Participants Contact Browser Wiki FTP search

Home >  
**SEARCH 1000 GENOMES FTP FILES**

Search term:  
  
Search for files on the FTP site

[Help on searching](#)

– [Search options](#)

---

## RESULTS

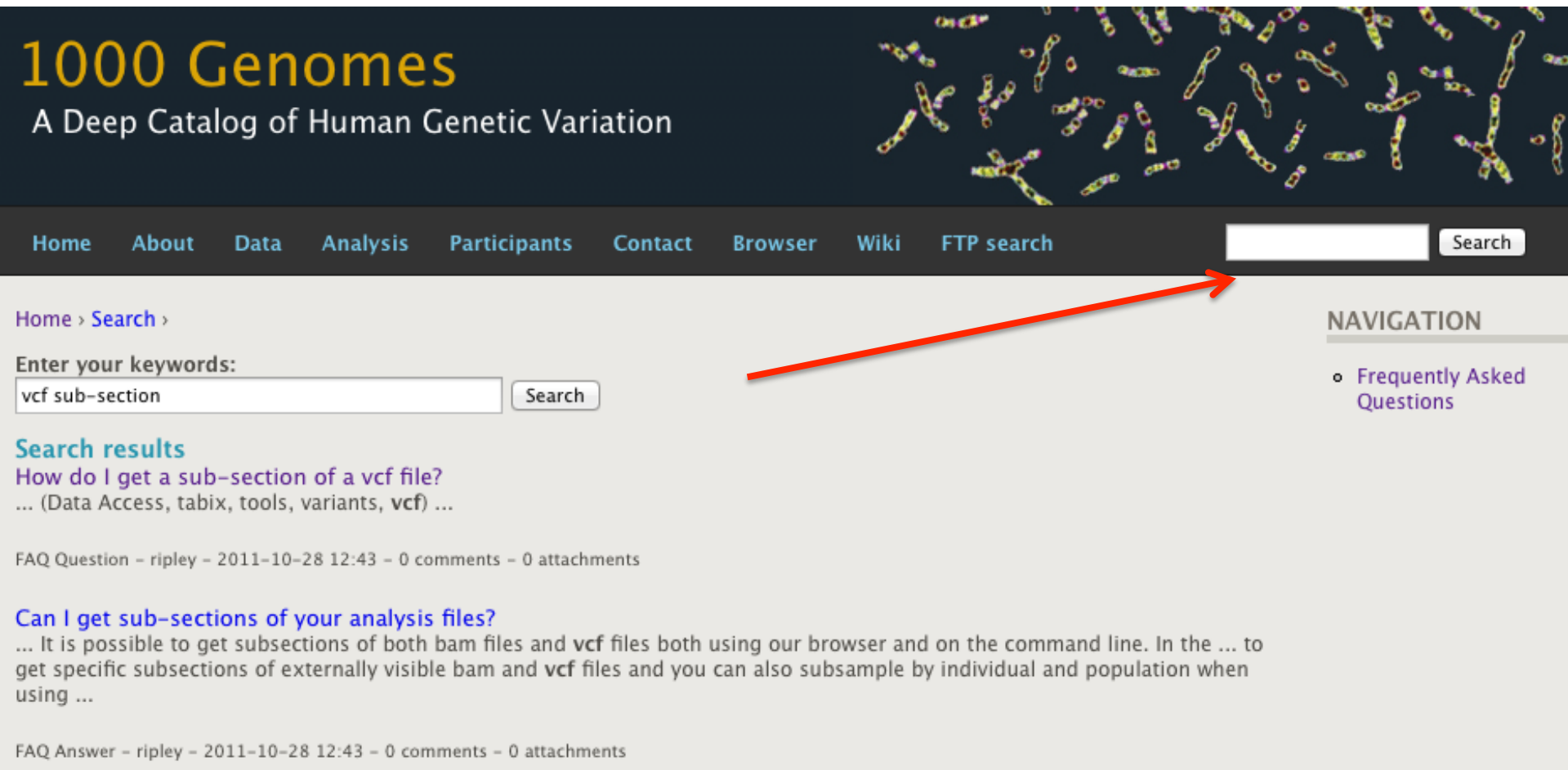
50 files found

File
<a href="ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr20.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz">ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr20.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz</a>
<a href="ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr15.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz">ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr15.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz</a>
<a href="ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr4.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz">ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr4.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz</a>
<a href="ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr9.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz">ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr9.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz</a>
<a href="ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr8.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz">ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr8.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz</a>
<a href="ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr12.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz">ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120103_omni_shapeit_haplotypes/ALL.chr12.omni_2123_samples_b37_SHAPEIT.20120103.haplotypes.vcf.gz</a>



# Finding Data: Exercise

- Finding help on getting sub-sections of VCF files



**1000 Genomes**  
A Deep Catalog of Human Genetic Variation

Home About Data Analysis Participants Contact Browser Wiki FTP search  Search

Home > Search >

Enter your keywords:  
 Search

**Search results**  
How do I get a sub-section of a vcf file?  
... (Data Access, tabix, tools, variants, vcf) ...

FAQ Question - ripley - 2011-10-28 12:43 - 0 comments - 0 attachments

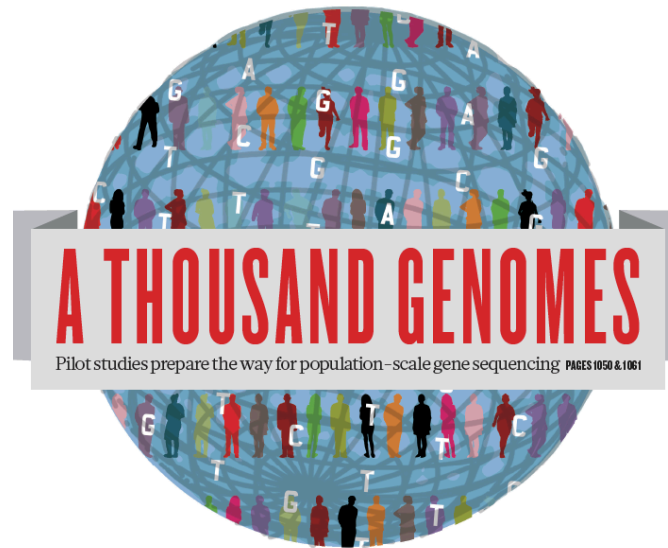
**Can I get sub-sections of your analysis files?**  
... It is possible to get subsections of both bam files and vcf files both using our browser and on the command line. In the ... to get specific subsections of externally visible bam and vcf files and you can also subsample by individual and population when using ...

FAQ Answer - ripley - 2011-10-28 12:43 - 0 comments - 0 attachments

NAVIGATION

- [Frequently Asked Questions](#)





The 1000 Genomes Project:

The 1000 Genomes Website and Ensembl-style 1000 Genomes Browser



# http://www.1000genomes.org

**1000 Genomes**  
A Deep Catalog of Human Genetic Variation

Home About Data Analysis Participants Contact **Browse** Wiki FTP search  Search

### LATEST ANNOUNCEMENTS

WEDNESDAY OCTOBER 12, 2011

#### October 2011 Integrated Variant Set release #ICHG2011

This **October 2011** release represents an integrated set of variant calls and phased genotypes including SNPS, short INDELS and Deletions based on low coverage and exome sequencing data across 1092 individuals.

Our [FAQ](#) contains instructions on how to get [smaller subsections](#) of these files

Data access links: [EBI](#) / [NCBI](#)

Link to additional information: [README file](#)

---

THURSDAY JUNE 23, 2011

#### June 2011 Data Release

Genotypes for 1094 individuals for the [May 2011 snp calls](#) from the 20101123 sequence and alignment release of the 1000 genomes project has now been made. This release is based on the GRCh37 assembly of the human genome and is released in the format [VCF 4.0](#)

Our [FAQ](#) contains instructions on how to get [smaller subsections](#) of these files

Data access links: [EBI](#) / [NCBI](#)

Link to additional information: [README file](#)

### NAVIGATION

- [Frequently Asked Questions](#)

### LINKS

- [All Project Announcements](#)
- [Sample and Project Information](#)
- [Media Archive](#)
- [Download the 1000 Genomes Pilot Paper](#)
- [Project Contacts](#)



# 1000 Genomes

A Deep Catalog of Human Genetic Variation



Tools | Help

## Search 1000 Genomes

e.g. gene BRCA2 or Chromosome 6:133098746-133108745

## Start Browsing 1000 Genomes data



[Browse Human](#) →  
GRCh37

[Protein variations](#) →  
View the consequences of sequence variation at the level of each protein in the genome.

[Individual genotypes](#) →  
Show different individual's genotype, for a variant.

## Browser update September 2011

based on interim Main project data from 20101123 for 1094 individuals and ensembl release 63. The data can be found on [the ftp site](#).

Please see [www.1000genomes.org](http://www.1000genomes.org) for more information about the data presented here and instructions for downloading the complete data set.

- [View sample data](#)

## The 1000 Genomes Browser

### Ensembl-based browser provides early access to 1000genomes data

In order to facilitate immediate analysis of the 1000 Genomes Project data by the whole scientific community, this browser (based on Ensembl) integrates the SNP calls from an [interim release 20101123](#). This data has been submitted to dbSNP, and once rsid's have been allocated, will be absorbed into the UCSC and Ensembl browsers according to their respective release cycles. Until that point any non rs SNP id's on this site are temporary and will NOT be maintained.

### Links



[1000 Genomes](#) →  
More information about the 1000 Genomes Project on the 1000 genomes main site.



[Pilot browser](#) →  
This browser is based on Ensembl release 60 and represents the variant set analysed as part of [A map of human genome variation from population-scale sequencing](#), Nature 467, 1061.1073.



[Tutorial](#) →  
The 1000 Genomes Browser Tutorial.

The 1000 Genomes Project is an international collaborative project described at [www.1000genomes.org](http://www.1000genomes.org).

The 1000 Genomes Browser is based on Ensembl web code.

Ensembl is a joint project of EMBL-EBI  and the [Wellcome Trust Sanger Institute](#)



1000 Genomes release 10 - October 2011 © [EBI](#)

[About 1000 Genomes](#) | [Contact Us](#) | [Help](#)

<http://browser.1000genomes.org>



EMBL-EBI



# 1000 Genomes

A Deep Catalog of Human Genetic Variation



Human (GRCh37)

Location: 1:114,356,433-114,414,381

Gene: PTPN22

Transcript: PTPN22-001

Tools | Help

## Location-based displays

- Whole genome
- Chromosome summary
- Region overview
- Region in detail**
- Genetic Variation
  - Resequencing (20)
  - Linkage Data
  - Markers

Configure this page

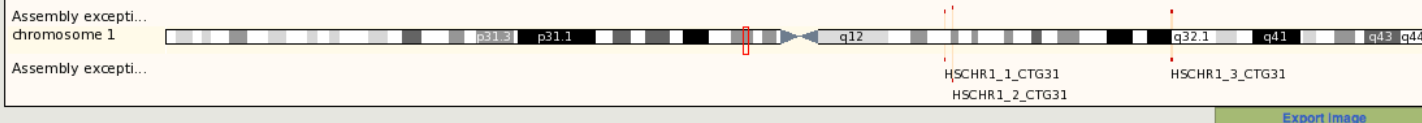
Manage your data

Export data

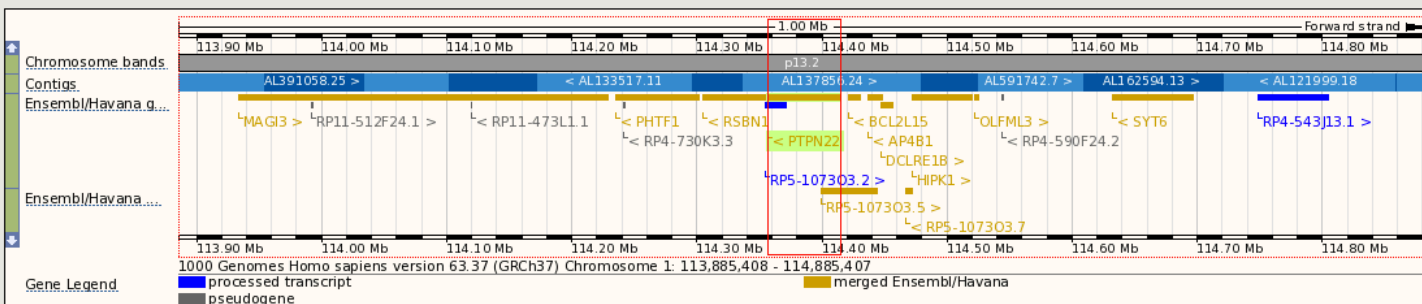
Get VCF data

Bookmark this page

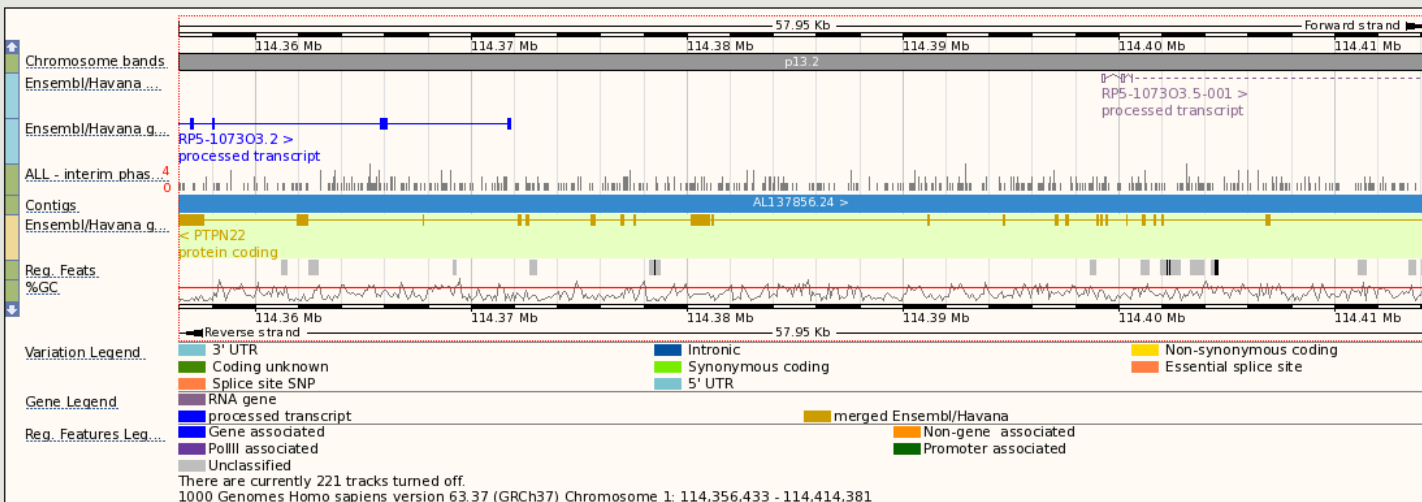
## Chromosome 1: 114,356,433-114,414,381



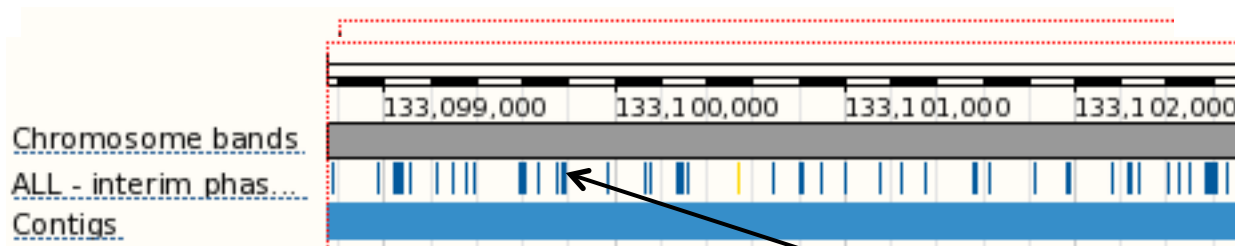
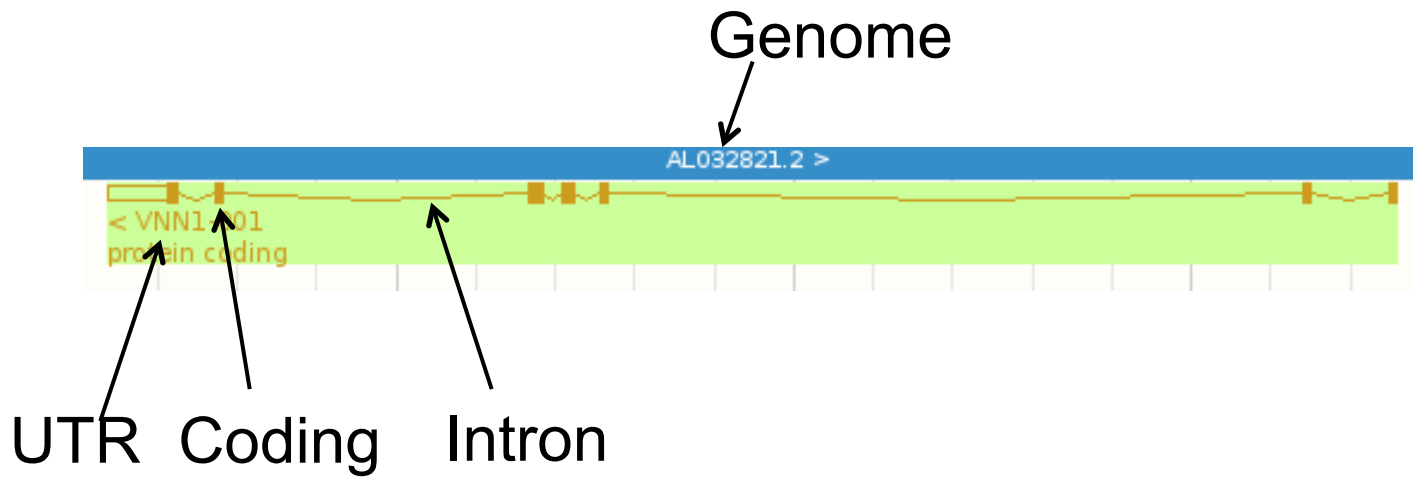
## Region in detail [help](#)



Location:  Go Gene:  Go



# What to Look For



Line indicates number of SNPS

Each Line is One SNP



Human (GRCh37)

Location: 13:32,890,598-32,890,664

Gene: BRCA2

Gene: BRCA2 (ENSG00000139618)

Gene-based displays

- Gene summary
- Splice variants (6)
- Supporting evidence
- Sequence
- External references
- Regulation
- Genetic Variation
- Variation Table
- Variation Image
- External Data
- ID History
- Gene history

Description: breast cancer 2, early onset [Source:HGNC Symbol;Acc:1101]  
Location: Chromosome 13: 32,889,811-32,973,805 forward strand.  
Transcripts: There are 6 transcripts in this gene

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
BRCA2-001	ENST00000380152	10930	ENSP00000369407	3418	Protein coding	CCDS9344
BRCA2-003	ENST00000530893	2009	ENSP00000435689	602	Protein coding	-
BRCA2-201	ENST00000544465	10984	ENSP00000439202	3418	Protein coding	CCDS9344
BRCA2-002	ENST00000470094	842	ENSP00000434988	186	Nonsense mediated decay	-
BRCA2-005	ENST00000507792	495	ENSP00000433189	64	Nonsense mediated decay	-
BRCA2-006	ENST00000537776	523	No protein product	-	Retained intron	-

- Configure this page
- Manage your data
- Export data
- Bookmark this page

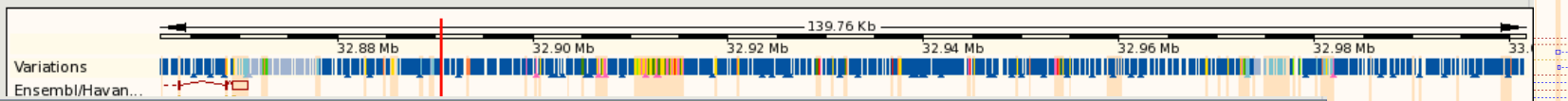
### Transcript and Gene level displays

In 1000 Genomes we provide displays at two levels:

- Transcript views which provide information specific to an individual transcript such as the cDNA and CDS sequences and protein domain annotation.
- Gene views which provide displays for data associated at the gene level such as orthologues, paralogues, regulatory regions and splice variants.

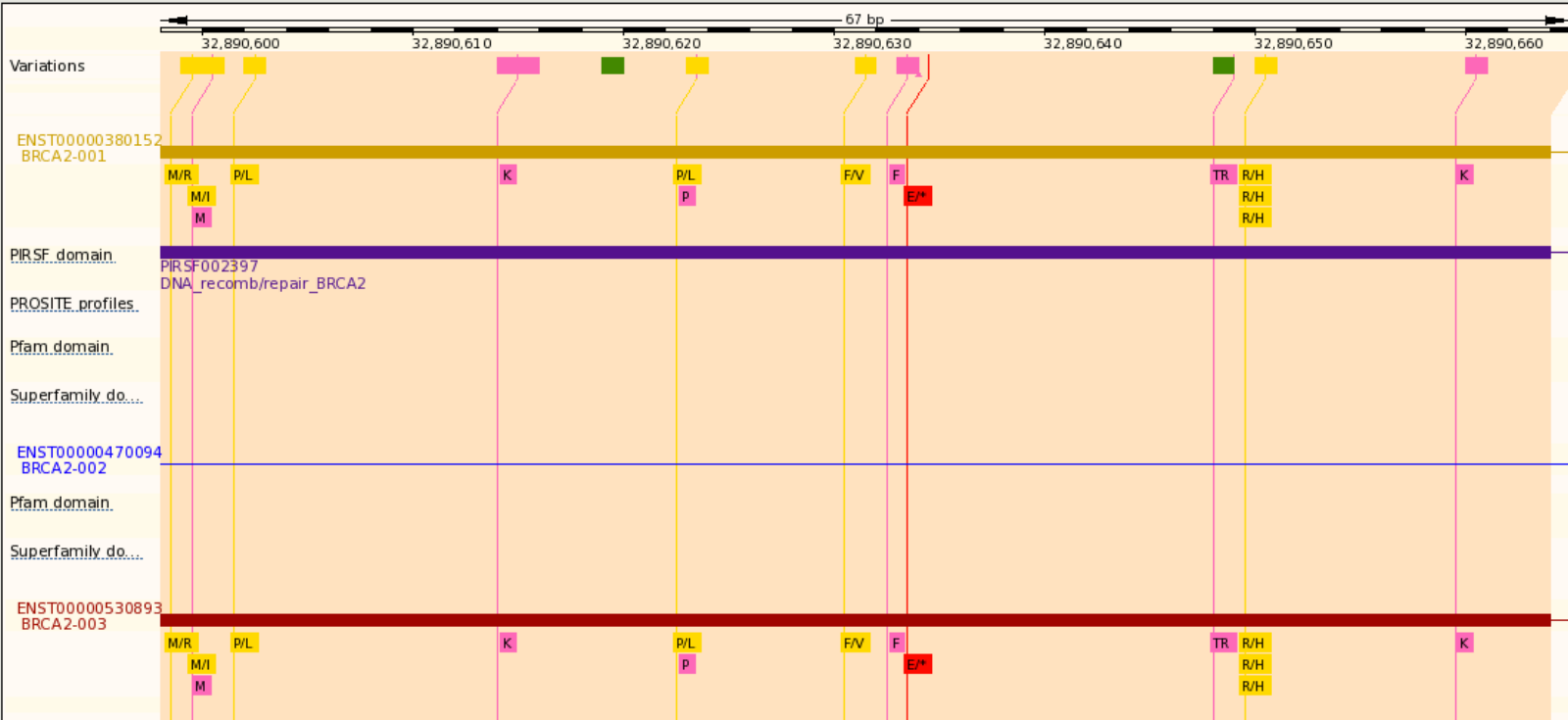
This view is a gene level view. To access the transcript level displays select a Transcript ID in the table above and then navigate to the information you want using the menu at the left hand side of the page. To return to viewing gene level information click on the Gene tab in the menu bar at the top of the page.

Variation Image [help](#)



Location: 13:32890598 - 32890664

Variation ID:



Export Image

None of the intronic variations are removed by the Context filter.

# 1000 Genomes

A Deep Catalog of Human Genetic Variation

Human (GRCh37) Location: 6:74,125,388-74,126,388 Variation: rs311685

- Variation displays**
- Flanking sequence
  - Gene/Transcript (3)
  - Population genetics (46)**
  - Individual genotypes (2769)
  - Genomic context
  - Phenotype Data
  - Phylogenetic Context
  - External Data

- Configure this page
- Manage your data
- Export data
- Get VCF data
- Bookmark this page
- Download view as CSV

Variation: rs311685

**Variation class** SNP (rs311685 source dbSNP\_132 - Variants (including SNPs and indels) imported from dbSNP [http://www.ncbi.nlm.nih.gov/projects/SNP/])

**Synonyms**  
**Affy GeneChip 100K Array** SNP\_A-1679873  
**Affy GenomeWideSNP\_6.0** AFFY\_6\_1M\_SNP\_A-8668494, SNP\_A-8668494  
**dbSNP** rs58378291, rs17756820, rs52794514, rs524803, rs3173186, rs11567000, rs17421786  
**ENSEMBL** ENSSNP9062281  
**Illumina\_Human1M-duoV3** rs311685  
**Uniprot** VAR\_057235

**Present in**  
 1000 genomes - High coverage - Trios (1000 genomes - High coverage - Trios - CEU, 1000 genomes - High coverage - Trios - YRI), 1000 genomes - Low coverage (1000 genomes - Low coverage - CEU, 1000 genomes - Low coverage - CHB+JPT, 1000 genomes - Low coverage - YRI), ALL - interim phase 1 - 1000 Genomes (AFR - interim phase 1 - 1000 Genomes, AMR - interim phase 1 - 1000 Genomes, ASN - interim phase 1 - 1000 Genomes, EUR - interim phase 1 - 1000 Genomes), ENSEMBL:Venter,HapMap

**Alleles** A/G (Ambiguity code: R)

**Ancestral allele** A

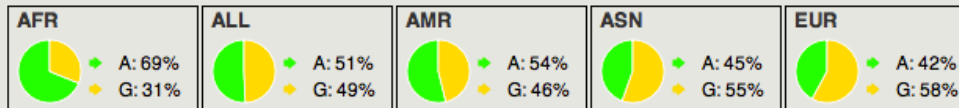
**Location** This feature maps to 6:74125888 (forward strand) | [View in location tab](#)

**Validation status** Proven by cluster, frequency, doublehit, 1000Genome HapMap variant

**HGVSN names** This feature has 4 HGVSN names - click the plus to show

## Population

### 1000 genomes alleles frequencies



### 1000 genomes

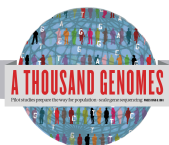
Population	Alleles A	Alleles G	Genotypes A/A	Genotypes A/G	Genotypes G/G	Count
1000GENOMES:AFR	0.689	0.311	0.463	0.451	0.085	114
1000GENOMES:ALL	0.507	0.493	0.269	0.477	0.254	294
1000GENOMES:AMR	0.539	0.461	0.293	0.492	0.215	53
1000GENOMES:ASN	0.446	0.554	0.199	0.493	0.308	57
1000GENOMES:EUR	0.421	0.579	0.184	0.475	0.341	70

### 1000 genomes pilot

Population	ssID	Submitter	Alleles A	Alleles G	Count
<a href="#">1000GENOMES:pilot_1_CEU_low_coverage_panel</a>	<a href="#">ss233534774</a>	<a href="#">1000GENOMES</a>	0.458	0.542	
<a href="#">1000GENOMES:pilot_1_CHB+JPT_low_coverage_panel</a>	<a href="#">ss240577229</a>	<a href="#">1000GENOMES</a>	0.400	0.600	
<a href="#">1000GENOMES:pilot_1_YRI_low_coverage_panel</a>	<a href="#">ss222470667</a>	<a href="#">1000GENOMES</a>	0.729	0.271	

<a href="#">CSHL-HAPMAP:HAPMAP-LWK</a>	<a href="#">ss5253350</a>	<a href="#">TSC-CSHL</a>	0.667	0.333	0.400	0.533	0.067	6
<a href="#">CSHL-HAPMAP:HAPMAP-MEX</a>	<a href="#">ss5253350</a>	<a href="#">TSC-CSHL</a>	0.490	0.510	0.245	0.490	0.285	13
<a href="#">CSHL-HAPMAP:HAPMAP-MKK</a>	<a href="#">ss5253350</a>	<a href="#">TSC-CSHL</a>	0.633	0.367	0.410	0.446	0.144	20
<a href="#">CSHL-HAPMAP:HAPMAP-TSI</a>	<a href="#">ss5253350</a>	<a href="#">TSC-CSHL</a>	0.488	0.512	0.226	0.524	0.250	21
<a href="#">CSHL-HAPMAP:HapMap-YRI</a>	<a href="#">ss5253350</a>	<a href="#">TSC-CSHL</a>	0.708	0.292	0.487	0.442	0.071	8
<a href="#">SEATTLESEQ:Eight-Hapmap-Samples</a>	<a href="#">ss159712995</a>	<a href="#">SEATTLESEQ</a>	unknown	unknown				

Other data (26)



# 1000 Genomes

A Deep Catalog of Human Genetic Variation

Human (GRCh37) Location: 1:114,356,433-114,414,381 Gene: PTPN22 Transcript: PTPN2-001

Transcript: PTPN2-001 (ENST00000359785)

**Description** protein tyrosine phosphatase, non-receptor type 22 (lymphoid) [Source:HGNC Symbol;Acc:9652]  
**Location** Chromosome 1:114,356,433-114,414,381 reverse strand.  
**Gene** This transcript is a product of gene [ENSG00000134242](#) - There are 12 transcripts in this gene

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
PTPN22-001	ENST00000359785	3654	ENSP00000352833	807	Protein coding	CCDS8883
PTPN22-002	ENST00000460620	1794	ENSP00000433141	179	Protein coding	-
PTPN22-004	ENST00000528414	3424	ENSP00000435176	752	Protein coding	-
PTPN22-006	ENST00000420377	2726	ENSP00000388229	795	Protein coding	-
PTPN22-007	ENST00000525799	2118	ENSP00000432674	668	Protein coding	-
PTPN22-201	ENST00000354605	2347	ENSP00000346621	691	Protein coding	CCDS884
PTPN22-202	ENST00000358253	2414	ENSP00000439372	563	Protein coding	-
PTPN22-008	ENST00000532224	2421	ENSP00000431249	135	Nonsense mediated decay	-
PTPN22-010	ENST00000529045	527	ENSP00000434932	92	Nonsense mediated decay	-
PTPN22-009	ENST00000534519	565	-	-	Processed transcript	-
PTPN22-003	ENST00000484147	2258	-	-	Retained intron	-
PTPN22-005	ENST00000469077	562	-	-	Retained intron	-

### Transcript and Gene level displays

Views in 1000 Genomes are separated into gene based views and transcript based views according to which level the information is more appropriately associated with. This view is a transcript level view. To flip between the two sets of views you can click on the Gene and Transcript tabs in the menu bar at the top of the page.

Variations [help](#)

Show All entries Show/hide columns Filter

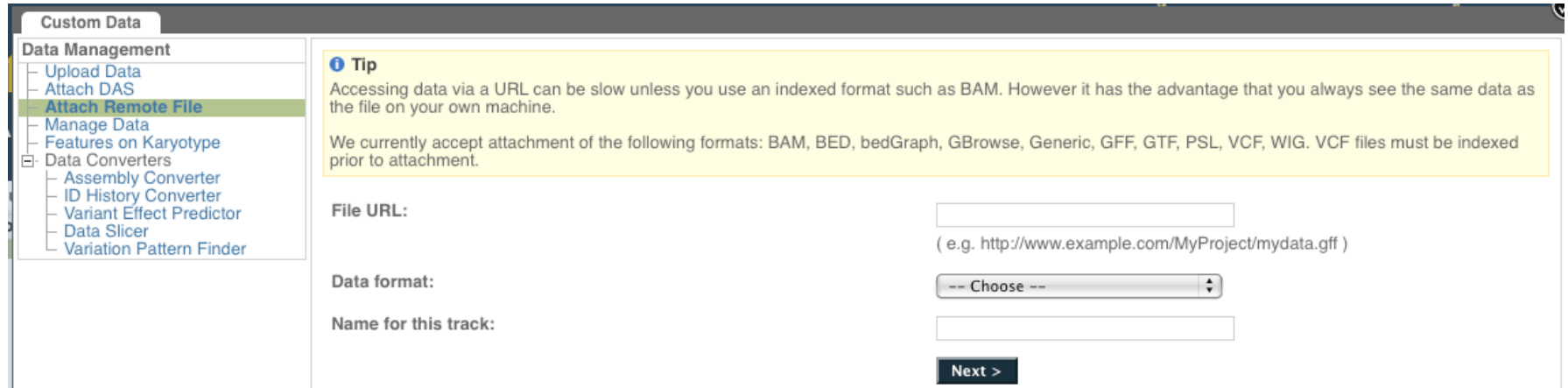
Residue	Variation ID	Variation type	Alleles	Ambiguity code	Residues	Codons	SIFT	PolyPhen
16	<a href="#">rs74163639</a>	Synonymous coding	G/A	R	S	AGC, AGT	-	-
49	<a href="#">rs61745743</a>	Synonymous coding	A/G	R	A	GCT, GCC	-	-
71	<a href="#">rs74163642</a>	Non-synonymous coding	A/G	R	V, A	GTA, GCA	deleterious	probably damaging
141	<a href="#">rs115552198</a>	Non-synonymous coding	G/A	R	R, C	CGC, TGC	deleterious	probably damaging
177	<a href="#">1KG_1_114399013</a>	Synonymous coding	C/T	Y	K	AAG, AAA	-	-
183	<a href="#">rs34590413</a>	Stop gained	G/A	R	R, *	CGA, TGA	-	-
201	<a href="#">rs74163647</a>	Non-synonymous coding	G/A	R	S, F	TCT, TTT	deleterious	probably damaging
206	<a href="#">rs61738614</a>	Non-synonymous coding	A/C	M	L, R	CTT, CGT	deleterious	probably damaging
232	<a href="#">rs78195073</a>	Synonymous coding	T/C	Y	G	GGA, GGG	-	-
247	<a href="#">rs35910094</a>	Synonymous coding	T/G	K	L	CTA, CTC	-	-
263	<a href="#">rs33996649</a>	Non-synonymous coding	C/T	Y	R, Q	CGG, CAG	tolerated	benign
266	<a href="#">rs72650670</a>	Non-synonymous coding	G/A	R	R, W	CGG, TGG	deleterious	probably damaging
277	<a href="#">rs72483511</a>	Stop gained, Splice site	C/A	M	E, *	GAA, TAA	-	-
324	<a href="#">rs113984534</a>	Synonymous coding	A/G	R	Y	TAT, TAC	-	-
366	<a href="#">rs74163654</a>	Synonymous coding	C/T	Y	E	GAG, GAA	-	-
370	<a href="#">rs72650671</a>	Non-synonymous coding	G/T	K	H, N	CAC, AAC	deleterious	possibly damaging
388	<a href="#">rs77913785</a>	Non-synonymous coding	G/T	K	D, E	GAC, GAA	deleterious	benign
413	<a href="#">1KG_1_114380784</a>	Non-synonymous coding	T/G	K	Q, P	CAA, CCA	deleterious	benign
414	<a href="#">1KG_1_114380780</a>	Synonymous coding	A/G	R	S	AGT, AGC	-	-
427	<a href="#">rs112873647</a>	Non-synonymous coding	-/ATT	-	-, N	-, AAT	-	-
444	<a href="#">rs74163655</a>	Non-synonymous coding	T/A	W	I, L	ATA, TTA	tolerated	benign
447	<a href="#">rs112191110</a>	Non-synonymous coding	G/A	R	T, I	ACC, ATC	deleterious	probably damaging
452	<a href="#">rs56174946</a>	Synonymous coding	A/G	R	F	TTT, TTC	-	-
456	<a href="#">rs72650672</a>	Non-synonymous coding	G/C	S	Q, E	CAG, GAG	deleterious	possibly damaging
477	<a href="#">rs74163656</a>	Synonymous coding	A/G	R	L	CAT, CAC	-	-

- SIFT
- PolyPhen

Download view as CSV



# File upload to view with 1000 Genomes data



The screenshot shows a web interface for uploading data. On the left is a sidebar titled 'Custom Data' with a 'Data Management' section containing a tree view of options: Upload Data, Attach DAS, Attach Remote File (highlighted), Manage Data, Features on Karyotype, Data Converters (with sub-items: Assembly Converter, ID History Converter, Variant Effect Predictor, Data Slicer, Variation Pattern Finder). The main area has a yellow tip box stating that URL access can be slow unless using indexed formats like BAM, and lists supported formats: BAM, BED, bedGraph, GBrowse, Generic, GFF, GTF, PSL, VCF, WIG. Below the tip is a form with fields for 'File URL:', 'Data format:' (a dropdown menu showing '-- Choose --'), and 'Name for this track:'. A 'Next >' button is at the bottom right.

- Supports popular file types:
  - BAM, BED, bedGraph, BigWig, GBrowse, Generic, GFF, GTF, PSL, VCF\*, WIG

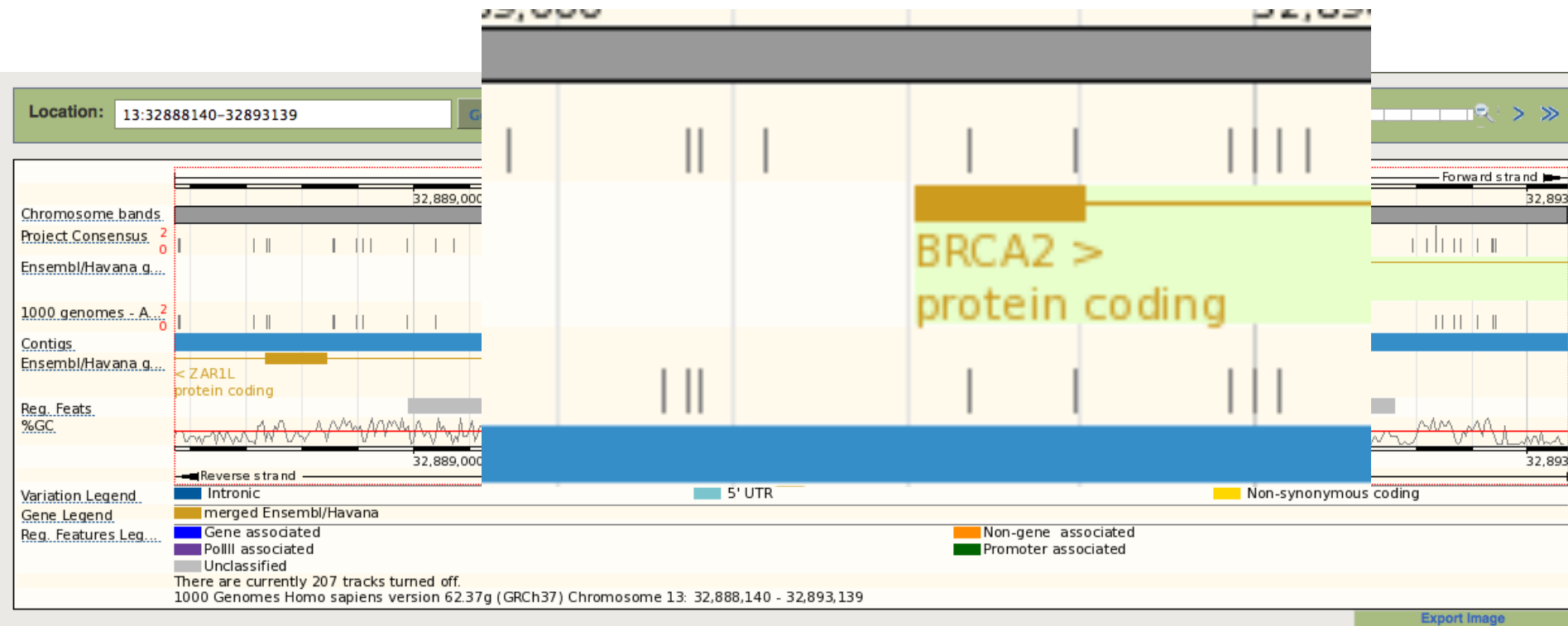
\* VCF must be indexed



# Uploaded VCF

Example:

Comparison of August calls and  
/technical/working/20110502\_vqsr\_phase1\_wgs\_snps/  
ALL.wgs.phase1.projectConsensus.snps.sites.vcf.gz



# The Browser: Coming Soon

**e!Ensembl** BLAST/BLAT | BioMart | Tools | Downloads | Help & Documentation | Blog | Mirrors Login · Register

Human (GRCh37) Location: 9:22,125,003-22,126,003 Variation: rs1333049

**Variation displays**

- Explore this variation
- Genomic context
  - Gene/Transcript (2)
- Population genetics (28)
- Individual genotypes (1737)
- Linkage disequilibrium
- Phenotype Data (8)
- Phylogenetic Context (4)
- Flanking sequence
- External Data

Configure this page | Manage your data | Export data | Bookmark this page

**rs1333049** SNP

**Source** [dbSNP 134](#) - Variants (including SNPs and indels) imported from [dbSNP](#)

**Alleles** Reference/Alternative: **G/C** | Ancestral: **C** | Ambiguity code: **S** | MAF: **0.40** (C)

**Location** Chromosome **9:22125503** (forward strand) | [View in location tab](#)

**Validation status** This variation is validated by **1000 Genomes**, **HapMap** and also cluster, doublehit, frequency, precious, submitter


**Synonyms** This feature has **7** synonyms - click the plus to show

**HGVS name** [g.22125503G>C](#)


**Explore this variation** [help](#)



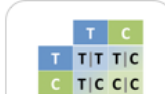
**Genomic context**




**Gene / Transcript**



**Population genetics**




**Individual genotypes**




**Linkage disequilibrium**



**Phenotype data**



**Phylogenetic context**



**Flanking sequence**

**Help with variations**

**YouTube videos**

- [SNPs and other Variations - 1 of 2](#)
- [SNPs and other Variations - 2 of 2](#)
- [Clip: Genome Variation](#)
- [BioMart: Variation IDs to HGNC Symbols](#)

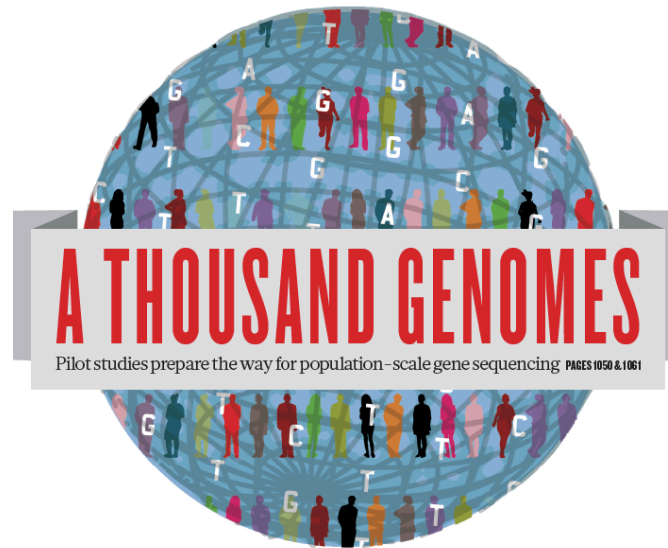
**Reference materials**

- [Ensembl variation data: background and terminology](#)
- [Variation Quick Reference card](#)

**Additional resources**

- [Accessing variation data with the Variation API](#)
- [Genomes and SNPs in Malaria](#)





# The 1000 Genomes Project:

## Exercise 2: Finding Variation Using the Browser

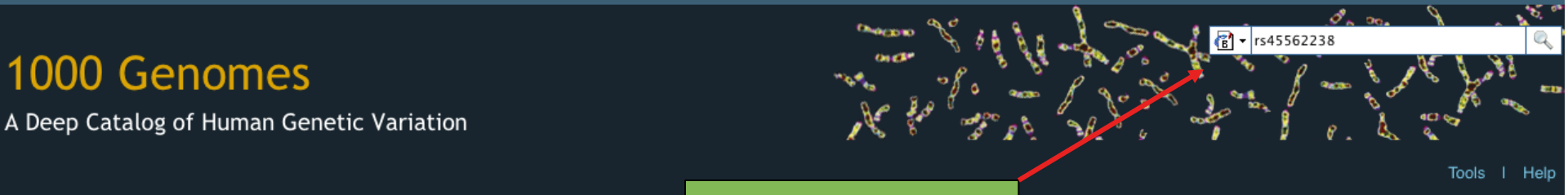


# Exercise: Finding Variation Using the Browser

- Find the variant rs45562238
  - <http://browser.1000genomes.org>
- Which 1000 Genomes Super Populations was it discovered in?
- What are it's allele frequencies?







# 1000 Genomes

A Deep Catalog of Human Genetic Variation

Search Box

## Search 1000 Genomes

Go

e.g. gene BRCA2 or Chromosome 6:133098746-133108745

## Start Browsing 1000 Genomes data



[Browse Human](#) →  
GRCh37

[Protein variations](#) →  
View the consequences of sequence variation at the level of each protein in the genome.

[Individual genotypes](#) →  
Show different individual's genotype, for a variant.

## Browser update September 2011

based on interim Main project data from 20101123 for 1094 individuals and ensembl release 63. The data can be found on [the ftp site](#).

Please see [www.1000genomes.org](#) for more information about the data presented here and instructions for downloading the complete data set.

- [View sample data](#)

## Ensembl-based browser provides early access to 1000genomes data

In order to facilitate immediate analysis of the 1000 Genomes Project data by the whole scientific community, this browser (based on Ensembl) integrates the SNP calls from an [interim release 20101123](#). This data has been submitted to dbSNP, and once rsid's have been allocated, will be absorbed into the UCSC and Ensembl browsers according to their respective release cycles. Until that point any non rs SNP id's on this site are temporary and will NOT be maintained.

## Links



[1000 Genomes](#) →  
More information about the 1000 Genomes Project on the 1000 genomes main site.




[Pilot browser](#) →  
This browser is based on Ensembl release 60 and represents the variant set analysed as part of [A map of human genome variation from population-scale sequencing](#), Nature 467, 1061-1073.



[Tutorial](#) →  
The 1000 Genomes Browser Tutorial.

The 1000 Genomes Project is an international collaborative project described at [www.1000genomes.org](#).

The 1000 Genomes Browser is based on Ensembl web code.

Ensembl is a joint project of EMBL-EBI  and the Wellcome Trust Sanger Institute 





Human (GRCh37) ▾

Search 1000 Genomes

New Search

Configure this page

Manage your data

Export data

Get VCF data

Bookmark this page

Results S

You searched for 'rs45562238'

## Gene or Gene Product

0 entrie(s) matched your search strings.

## Genetic Marker

0 entrie(s) matched your search strings.

## Array Probe Set

0 entrie(s) matched your search strings.

## SNP

1 entrie(s) matched your search strings.

1. dbSNP SNP: [rs45562238](#)

SNP Result

## Interpro Domain

0 entrie(s) matched your search strings.

## Gene Family

0 entrie(s) matched your search strings.

## Sequence Aligned to Genome, eg. EST or Protein

0 entrie(s) matched your search strings.

## Genomic Region, eg. Clone or Contig

0 entrie(s) matched your search strings.

1000 Genomes release 9 - September 2011 © [EBI](#)



- Variation displays**
- Flanking sequence
  - Gene/Transcript (1)
  - Population genetics (11)
  - Individual genotypes (2770)
  - Genomic context
  - Phenotype Data
  - Phylogenetic Context
  - External Data
- 
- Configure this page
  - Manage your data
  - Export data
  - Get VCF data
  - Bookmark this page

Variation: rs45562238

Population Genetics

<b>Variation class</b>	SNP ( <a href="#">rs45562238</a> source <a href="#">dbSNP 132</a> - Variants (including SNPs and indels) imported from dbSNP)
<b>Synonyms</b>	OMNI SNP6-133055237 Uniprot <a href="#">VAR_023973</a>
<b>Present in</b>	1000 genomes - Low coverage (1000 genomes - Low coverage - CEU),ALL - interim phase 1 - 1000 Genomes (AMR - interim phase 1 - 1000 Genomes, EUR - interim phase 1 - 1000 Genomes),ENSEMBL:Watson
<b>Alleles</b>	T/C (Ambiguity code: Y)
<b>Ancestral allele</b>	T
<b>Location</b>	This feature maps to 6:133013544 (forward strand)   <a href="#">View in location tab</a>
<b>Validation status</b>	Proven by <b>cluster, frequency, 1000Genome</b>
<b>HGVS names</b> <input type="checkbox"/>	This feature has 3 HGVS names - click the plus to show

Flanking sequence [help](#)

**Flanking Sequence (reference and dbSNP)**

```

AAAAAAAAAAAAAAAAAGATTTAGCAAATAAAATTTAGCCTGTATTTGTAGGCTTCCCAGCTA
AAGTTGAATTTCAAATAAACAACAAAATTTTAGCATACATAAAATCCCAAAAAGTGTGAT
ATGCTTATACTAAAAATCATTTCATTTATCTAAACTGTATATTTAACTGAGTGTCTT
TATTTTATCCAGTAACTCTTCTGACATCAAAATATTACCTGTAGATAAATAGCGCCCTCC
ACAGTGTGCAGTCCGTCAAATGCCCTAGAGCGTACACTTCATTTGGTATGTTCTCAGAC
ATTTGTAGCTTAAATGACAGCAGAGATCTTTCTGACAAACTGTATAATTTCCCTGCAACT
CCTGTGAGCTTCAAAAAGTGAATTCATCGAAAAAGACAGYGCCTTTAAATTCCTTGTTT
CCTGATGAGAGCGCTTCTATACTGCTGGCATAGGAAGTCCAGTTCACCACTGCAGAATGG
GATGGGTGGGAATCCAGTTGCGAGAGGAGGAGTTTCCCTCTCTGTCTTCATATCATAA
TGAAATGCCTTGAAGAATTGGGTGCATAGATGCCACTTCCGGAAGTAATAAGTTGAAA
ACATCATTTGAAATGGAAAAGTAAAAAGGGATTTTCATAATTGTTATTACTATTTAACCT
CATATTCACATTTTACTCTCTCTAAGTTATATGTTTTATTAATTTTCAGTAAAAACATGCTG
ATAAAGACAGAGAAGTGAGAGGAAATCAATTAAGTATGAAACCAAGAAGCCAAGAAGGG
AGAGTTGTATGTAAGAAGTCA
    
```

*(Variant highlighted)*



# Gene/Transcript

- Variation displays
  - Flanking sequence
  - Gene/Transcript (1)
  - Population genetics (11)**
  - Individual genotypes (2770)
  - Genomic context
  - Phenotype Data
  - Phylogenetic Context
  - External Data
- Configure this page
- Manage your data
- Export data
- Get VCF data
- Bookmark this page
- Download view as CSV

## Variation: rs45562238

**Variation class** SNP ([rs45562238](#) source [dbSNP 132](#) - Variants (including SNPs and indels) imported from dbSNP [<http://www.ncbi.nlm.nih.gov/projects/SNP/>])

**Synonyms** OMNI SNP6-133055237  
Uniprot [VAR\\_023973](#)

**Present in** 1000 genomes - Low coverage (1000 genomes - Low coverage - CEU),ALL - interim phase 1 - 1000 Genomes (AMR - interim phase 1 - 1000 Genomes, EUR - interim phase 1 - 1000 Genomes),ENSEMBL:Watson

**Alleles** T/C (Ambiguity code: Y)

**Ancestral allele** T

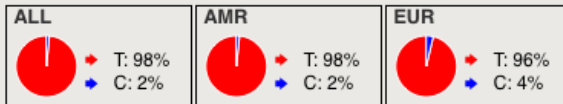
**Location** This feature maps to 6:133013544 (forward strand) | [View in location tab](#)

**Validation status** Proven by **cluster, frequency, 1000Genome**

**HGVS names** This feature has 3 HGVS names - click the plus to show

## Population genetics [help](#)

### 1000 genomes alleles frequencies



## Pie Charts

### 1000 genomes

Show/hide columns

Population	Alleles C	Alleles T	Genotypes CIT	Genotypes TIT	Allele count	Genotype count	Genotype detail
1000GENOMES:ALL	0.016	0.984	0.031	0.969	34 (C) / 2154 (T)	34 (CIT) / 1060 (TIT)	<a href="#">Show</a>
1000GENOMES:AMR	0.017	0.983	0.033	0.967		6 (CIT) / 175 (TIT)	<a href="#">Show</a>
1000GENOMES:EUR	0.037	0.963	0.073	0.927		28 (CIT) / 353 (TIT)	<a href="#">Show</a>



# 1000 Genomes

A Deep Catalog of Human Genetic Variation

Human (GRCh37)

Location: 6:133,013,044-133,014,044

Variation: rs45562238

Tools | Help

## Variation displays

- Flanking sequence
- Gene/transcript (1)**
- Population genetics (11)
- Individual genotypes (2770)
- Genomic context
- Phenotype Data
- Phylogenetic Context
- External Data

Configure this page

Manage your data

Export data

Get VCF data

Bookmark this page

Download view as CSV

## Variation: rs45562238

**Variation class** SNP ([rs45562238](#) source [dbSNP 132](#) - Variants (including SNPs and indels) imported from dbSNP [<http://www.ncbi.nlm.nih.gov/projects/SNP/>])

**Synonyms** OMNI SNP6-133055237  
Uniprot [VAR\\_023973](#)

**Present in** 1000 genomes - Low coverage (1000 genomes - Low coverage - CEU),ALL - interim phase 1 - 1000 Genomes (AMR - interim phase 1 - 1000 Genomes),ENSEMBL:Watson

**Alleles** T/C (Ambiguity code: Y)

**Ancestral allele** T

**Location** This feature maps to 6:133013544 (forward strand) | [View in location tab](#)

**Validation status** Proven by **cluster, frequency, 1000Genome**

**HGVS names** This feature has 3 HGVS names - click the plus to show

Sift/PolyPhen

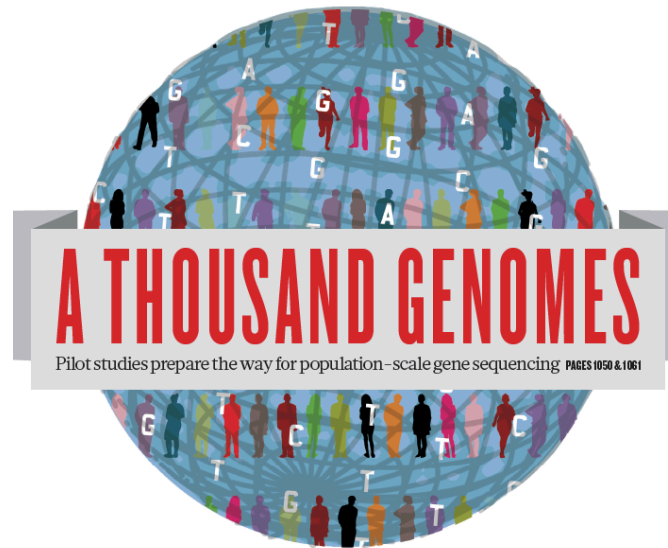
## Gene/transcript [help](#)

Gene	Transcript (strand)	Allele (transcript allele)	Type	HGVS names	Position in transcript	Position in CDS	Position in protein	Amino acid	Codons	SIFT	PolyPhen
<a href="#">ENSG00000112299</a>	<a href="#">ENST00000367928</a>	(-) C (G)	Non-synonymous coding	ENST00000367928.4:c.1006A>G ENSP00000356905.4:p.Thr336Ala	1020	1006	336	T/A	ACT/GCT	tolerated	benign

1000 Genomes release 9 - September 2011 © [EBI](#)

[About 1000 Genomes](#) | [Contact Us](#) | [Help](#)





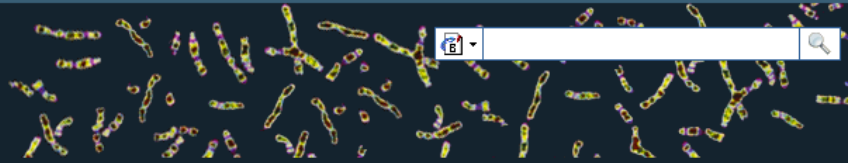
The 1000 Genomes Project:

The 1000 Genomes Tools



# 1000 Genomes

A Deep Catalog of Human Genetic Variation



[Tools](#) | [Help](#)

## Search 1000 Genomes

e.g. gene BRCA2 or Chromosome 6:133098746-133108745

## Start Browsing 1000 Genomes data



[Browse Human](#) →  
GRCh37

[Protein variations](#) →  
View the consequences of sequence variation at the level of each protein in the genome.

[Individual genotypes](#) →  
Show different individual's genotype, for a variant.

## Browser update September 2011

based on interim Main project data from 20101123 for 1094 individuals and ensembl release 63. The data can be found on [the ftp site](#).

Please see [www.1000genomes.org](http://www.1000genomes.org) for more information about the data presented here and instructions for downloading the complete data set.

- [View sample data](#)

## The 1000 Genomes Browser

### Ensembl-based browser provides early access to 1000genomes data

In order to facilitate immediate analysis of the 1000 Genomes Project data by the whole scientific community, this browser (based on Ensembl) integrates the SNP calls from an [interim release 20101123](#). This data has been submitted to dbSNP, and once rsid's have been allocated, will be absorbed into the UCSC and Ensembl browsers according to their respective release cycles. Until that point any non rs SNP id's on this site are temporary and will NOT be maintained.

### Links



[1000 Genomes](#) →  
More information about the 1000 Genomes Project on the 1000 genomes main site.



[Pilot browser](#) →  
This browser is based on Ensembl release 60 and represents the variant set analysed as part of [A map of human genome variation from population-scale sequencing](#), Nature 467, 1061.1073.



[Tutorial](#) →  
The 1000 Genomes Browser Tutorial.

The 1000 Genomes Project is an international collaborative project described at [www.1000genomes.org](http://www.1000genomes.org).

The 1000 Genomes Browser is based on Ensembl web code.

[Ensembl](#) is a joint project of EMBL-EBI  and the [Wellcome Trust Sanger Institute](#)



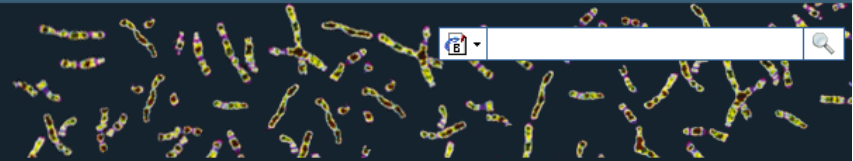
<http://browser.1000genomes.org>



# Tools page

## 1000 Genomes

A Deep Catalog of Human Genetic Variation



[Tools](#) | [Help](#)

We provide a number of ready-made tools for processing your data. At the moment, small datasets can be uploaded to our servers and processed online; for larger datasets, we provide an API script that can be downloaded (you will also need to [install our Perl API](#) to use these).

In the near future we aim to offer an intermediate service, whereby medium-to-large data sets can be submitted to a queue, similar to BLAST.

Currently available:

Tool	Description		
Assembly converter	Map your data to the current assembly. Accepted file formats: <a href="#">GFE</a> , <a href="#">GTE</a> , <a href="#">BED</a> , <a href="#">PSL</a> . N.B. Export is currently in GFF only	<a href="#">Online version</a>	<a href="#">API script</a>
ID History converter	Convert a set of Ensembl IDs from a previous release into their current equivalents.	<a href="#">Online version</a> (max 30 ids)	<a href="#">API script</a>
Variant Effect Predictor	(Formerly SNP Effect Predictor). Upload a set of SNPs in our <a href="#">standard format</a> and export a file containing consequence types. Uploaded tracks can also be viewed on Location pages.	<a href="#">Online version</a> (max 750 SNPs)	<a href="#">API script</a>
Data Slicer	Get a subset of data from a BAM or VCF file.	<a href="#">Online version</a> (max 10K region)	
Variation Pattern Finder	Identify variation patterns in a chromosomal region of interest for different individuals. Only variations with functional significance such non-synonymous coding, splice site will be reported by the tool.	<a href="#">Online version</a>	

1000 Genomes release 10 - October 2011 © [EBI](#)

[About 1000 Genomes](#) | [Contact Us](#) | [Help](#)





# Data slicer for subsets of the data

**1000 Genomes**  
A Deep (

**Custom Data**

**Data Management**

- Upload Data
- Attach DAS
- Attach Remote File
- Manage Data
- Features on Karyotype
- Data Converters
  - Assembly Converter
  - ID History Converter
  - Variant Effect Predictor
  - Data Slicer**
  - Variation Pattern Finder

We provide (use these).  
In the near  
Currently a  
**Tool**  
Assembly  
ID History  
Variant Ef  
Data Slice  
Variation F  
1000 Geno

**Tip**  
When slicing a VCF or BAM file, both the data file and its index file should be present on the web server and named correctly. The VCF file should have a ".vcf.gz" extension, and the index file should have a ".vcf.gz.tbi" extension, E.g: MyData.vcf.gz, MyData.vcf.gz.tbi  
The BAM file should have a ".bam" extension, and the index file should have a ".bam.bai" extension, E.g: MyData.bam, MyData.bam.bai  
Click [here](#) for more extensive documentation.

**VCF / BAM File URL:**   
e.g.  
ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\_phase1\_release/ALL.chr6.phase1

**Region:**   
( e.g. 1:1-50000 )

**Use VCF filters (this doesn't apply to BAM files):**

- None
- By individual(s)
- By population(s) \*

(to filter by populations please provide URL to a Sample-Population Mapping File in the box below)

**Sample-Population Mapping File URL:**   
e.g.  
ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\_phase1\_release/interim\_phase1.2

# Get VCF Button

Human (GRCh37) Location: 6:31,830,969-31,846,823 Gene: SLC44A4 Tools | Help

**Gene-based displays**

- Gene summary
- Splice variants (9)
- Supporting evidence
- Sequence
- External references
- Regulation
- Genetic Variation
  - Variation Table
  - Structural Variation
  - Variation Image
- External Data
- ID History
  - Gene history

**Gene: SLC44A4 (ENSG00000204385)**

**Description** solute carrier family 44, member 4 [Source:HGNC Symbol;Acc:13941]

**Location** [Chromosome 6: 31,830,969-31,846,823](#) reverse strand.

**Transcripts**  There are 9 transcripts in this gene

Show/hide columns

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
SLC44A4-001	<a href="#">ENST00000229729</a>	2589	<a href="#">ENSP00000229729</a>	710	Protein coding	<a href="#">CCDS4724</a>
SLC44A4-004	<a href="#">ENST00000414427</a>	1233	<a href="#">ENSP00000398901</a>	411	Protein coding	-
SLC44A4-201	<a href="#">ENST00000375562</a>	2505	<a href="#">ENSP00000364712</a>	668	Protein coding	-
SLC44A4-202	<a href="#">ENST00000544672</a>	2634	<a href="#">ENSP00000444109</a>	634	Protein coding	-
SLC44A4-002	<a href="#">ENST00000465707</a>	681	No protein product	-	Processed transcript	-
SLC44A4-003	<a href="#">ENST00000462671</a>	426	No protein product	-	Processed transcript	-
SLC44A4-007	<a href="#">ENST00000487680</a>	392	No protein product	-	Processed transcript	-
SLC44A4-005	<a href="#">ENST00000475563</a>	575	No protein product	-	Retained intron	-
SLC44A4-006	<a href="#">ENST00000479777</a>	655	No protein product	-	Retained intron	-

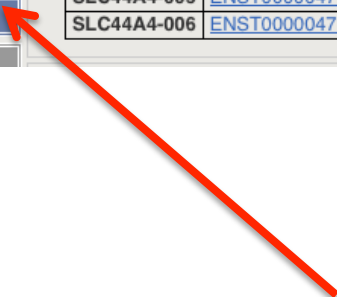
Configure this page

Manage your data

Export data

**Get VCF data**

Bookmark this page



# Ensembl Variant Effector Predictor (VEP)

- Takes list of variation and annotates with respect to Ensembl features
- Returns whether the SNP has been seen in the 1000 Genomes and if it has an rs number (if one has been assigned)
- Returns SIFT, PolyPhen and Condel scores
- Extensive filtering options by MAF and populations
- Web and command line versions



Data Management

- Upload Data
- Attach DAS
- Attach Remote File
- Manage Data
- Features on Karyotype
- Data Converters
  - Assembly Converter
  - ID History Converter
  - Variant Effect Predictor**
  - Data Slicer
  - Variation Pattern Finder

**Variant Effect Predictor:**

This tool takes a list of variant positions and alleles, and predicts the effects of each of these on overlapping transcripts and regulatory regions annotated in Ensembl. The tool accepts substitutions, insertions and deletions as input, uploaded as a list of [tab separated values](#), [VCF](#) or Pileup format input.

Upload is limited to 750 variants; lines after the limit will be ignored. Users with more than 750 variations can split files into smaller chunks, use the standalone [perl script](#) or the [variation API](#). See also [full documentation](#)

**Input file**

Species:

Human (Homo sapiens): GRCh37

Name for this upload (optional):

Paste file:

Upload file:

Choose File no file selected

or provide file URL:

Input file format:

Ensembl default

**Options**

Get regulatory region consequences:

Type of consequences to display:

Ensembl terms

Check for existing co-located variants:

Yes

Return results for variants in coding regions only:

Show HGNC identifier for genes where available:

Show Ensembl protein identifiers where available:

Show HGVS identifiers for variants where available:

No

**Non-synonymous SNP predictions (human only)**

SIFT predictions:

No

PolyPhen predictions:

No

Condel consensus (SIFT/PolyPhen) predictions:

No

**Frequency filtering of existing variants (human only)**

Filter variants by frequency:

NB: Enabling frequency filtering may be very slow for large datasets

Filter: Exclude variants with MAF greater than 0.1 in any 1KG low coverage population

Next >

# Variation Pattern Finder

- [http://browser.1000genomes.org/Homo\\_sapiens/UserData/VariationsMapVCF](http://browser.1000genomes.org/Homo_sapiens/UserData/VariationsMapVCF)
- VCF input
- Discovers patterns of Shared Inheritance
- Variants with functional consequences considered
- Web output with CSV and Excel downloads



Custom Data

Data Management

- Upload Data
- Attach DAS
- Attach Remote File
- Manage Data
- Features on Karyotype
- Data Converters
  - Assembly Converter
  - ID History Converter
  - Variant Effect Predictor
  - Data Slicer
- Variation Pattern Finder**

### Variation Pattern Finder

Export data: [CSV](#) [Excel](#)

**Go to collapsed view**

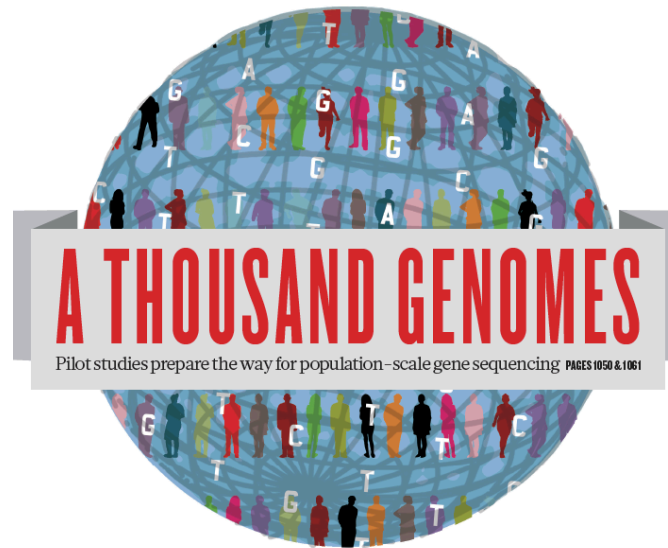
Population ASW	CEU	Freq	Variation info rs9369628:C/T	rs61361828:C/T	rs12192544:C/G	rs599
			6:46620135 ENST00000275016 SPLICE_SITE	6:46620240 ENST00000275016 NON_SYNONYMOUS_CODING:R/H	6:46620252 ENST00000275016 NON_SYNONYMOUS_CODING:R/P	6:466 ENST0 NON_S
NA20314, NA20322	NA12348, N	0.095	CIC	CIC	GIC	GIG
NA20356, NA19625 and 1 other(s)	NA11919, N	0.092	CIC	CIC	CIG	GIG
NA20291, NA19985 and 5 other(s)		0.069	CIT	CIC	CIC	GIG
NA20289, NA20294 and 4 other(s)		0.057	TIC	CIC	CIC	GIG
	NA12546, N	0.026	CIC	CIC	GIG	GIG
NA19819		0.012	TIT	CIC	CIC	GIG
	NA12283	0.011	TIC	CIC	CIG	GIG
NA19908, NA20278		0.011	CIT	CIC	GIC	GIG
NA19703		0.008	CIC	CIC	CIC	GIG
NA20351		0.007	CIC	CIC	CIC	GIG
		0.006	CIC	CIC	CIG	GIG
NA19712		0.004	CIC	CIC	CIC	CIG
		0.003	CIC	CIC	GIC	GIG
		0.003	TIC	CIC	CIC	GIG
		0.002	CIC	CIC	CIC	GIG



# Access to backend Ensembl databases

- Public MySQL database at
  - [mysql-db.1000genomes.org](http://mysql-db.1000genomes.org) port 4272
- Full programmatic access with Ensembl API
  - The 1000 Genomes Pilot uses Ensembl v60 databases and the NCBI36 assembly (this is frozen)
  - The 1000 Genomes main project currently uses Ensembl v63 databases





# The 1000 Genomes Project:

## Exercise 3: Using 1000 Genomes Tools





# Exercise: Using 1000 Genomes Tools

- Find the gene SLC44A4 using the search box on <http://browser.1000genomes.org>
- Get a VCF file for this Gene using the Get VCF button.
- Uncompress this file
  - You can get a copy at [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120120\\_1000genomes\\_tutorial/6.31830969-31846823.ALL.chr6.phase1.projectConsensus.genotypes.vcf](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120120_1000genomes_tutorial/6.31830969-31846823.ALL.chr6.phase1.projectConsensus.genotypes.vcf)
- Use this file with the Variant Effect Predictor
  - [http://browser.1000genomes.org/Homo\\_sapiens/UserData/UploadVariations](http://browser.1000genomes.org/Homo_sapiens/UserData/UploadVariations)
- Do any of the variants have deleterious effects according to SIFT or PolyPhen
- Use the example url on the page and the coordinates 6:31830700-31840700 with the Variation Pattern Finder
  - [http://browser.1000genomes.org/Homo\\_sapiens/UserData/VariationsMapVCF](http://browser.1000genomes.org/Homo_sapiens/UserData/VariationsMapVCF)



# 1000 Genomes

A Deep Catalog of Human Genetic Variation



SLC44A4

Tools | Help

## Search 1000 Genomes

e.g. gene **BRCA2** or **Chromosome 6:133098746-133108745**

## The 1000 Genomes Browser

### Ensembl-based browser provides early access to 1000genomes data

In order to facilitate immediate analysis of the 1000 Genomes Project data by the whole scientific community, this browser (based on Ensembl) integrates the SNP calls from an [interim release 20101123](#). This data has been submitted to dbSNP, and once rsid's have been allocated, will be absorbed into the UCSC and Ensembl browsers according to their respective release cycles. Until that point **any non rs SNP's on this site are temporary and will NOT be maintained.**

### Links



#### [1000 Genomes](#) →

More information about the 1000 Genomes Project on the 1000 genomes main site.



#### [Pilot browser](#) →

This browser is based on Ensembl release 60 and represents the variant set analysed as part of [A map of human genome variation from population-scale sequencing](#), Nature 467, 1061.1073.

## Start Browsing 1000 Genomes data



#### [Browse Human](#) →

GRCh37

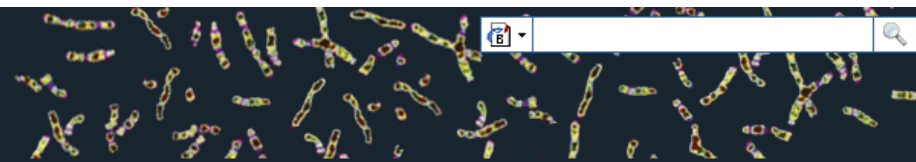
#### [Protein variations](#) →

View the consequences of sequence variation at the level of each protein in the genome.

#### [Individual genotypes](#) →

Show different individual's genotype, for a variant.





Human (GRCh37) ▾

Search 1000 Genomes

New Search

Configure this page

Manage your data

Export data

Get VCF data

Bookmark this page

## Results Summary

You searched for 'SLC44A4'

### Gene or Gene Product

10 entrie(s) matched your search strings.

1. Gene: [ENSG00000204385](#) [Region in detail]  
SLC44A4 - solute carrier fam 44, member 4 [Source:HGNC Symbol;Acc:13941]
2. Transcript: [ENST00000229729](#) [Region in detail]
3. Peptide: [ENSP00000398764](#) [Region in detail]  
SLC44A4
4. Peptide: [ENSP00000392054](#) [Region in detail]  
SLC44A4
5. Peptide: [ENSP00000404572](#) [Region in detail]  
SLC44A4
6. Peptide: [ENSP00000398901](#) [Region in detail]  
SLC44A4
7. Peptide: [ENSP00000415708](#) [Region in detail]  
SLC44A4
8. Peptide: [ENSP00000400263](#) [Region in detail]  
SLC44A4
9. Peptide: [ENSP00000414296](#) [Region in detail]  
SLC44A4
10. Peptide: [ENSP00000399161](#) [Region in detail]  
SLC44A4

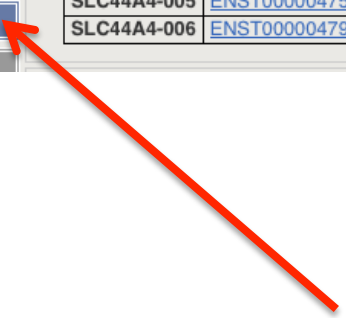


- Gene-based displays**
- Gene summary
  - Splice variants (9)
  - Supporting evidence
  - Sequence
  - External references
  - Regulation
  - Genetic Variation
    - Variation Table
    - Structural Variation
    - Variation Image
  - External Data
  - ID History
    - Gene history
- Configure this page
- Manage your data
- Export data
- Get VCF data
- Bookmark this page

**Gene: SLC44A4 (ENSG00000204385)**

**Description** solute carrier family 44, member 4 [Source:HGNC Symbol;Acc:13941]  
**Location** [Chromosome 6: 31,830,969-31,846,823](#) reverse strand.  
**Transcripts**  There are 9 transcripts in this gene

Show/hide columns		Filter				
Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
SLC44A4-001	<a href="#">ENST00000229729</a>	2589	<a href="#">ENSP00000229729</a>	710	Protein coding	<a href="#">CCDS4724</a>
SLC44A4-004	<a href="#">ENST00000414427</a>	1233	<a href="#">ENSP00000398901</a>	411	Protein coding	-
SLC44A4-201	<a href="#">ENST00000375562</a>	2505	<a href="#">ENSP00000364712</a>	668	Protein coding	-
SLC44A4-202	<a href="#">ENST00000544672</a>	2634	<a href="#">ENSP00000444109</a>	634	Protein coding	-
SLC44A4-002	<a href="#">ENST00000465707</a>	681	No protein product	-	Processed transcript	-
SLC44A4-003	<a href="#">ENST00000462671</a>	426	No protein product	-	Processed transcript	-
SLC44A4-007	<a href="#">ENST00000487680</a>	392	No protein product	-	Processed transcript	-
SLC44A4-005	<a href="#">ENST00000475563</a>	575	No protein product	-	Retained intron	-
SLC44A4-006	<a href="#">ENST00000479777</a>	655	No protein product	-	Retained intron	-



Data Management

- Upload Data
- Attach DAS
- Attach Remote File
- Manage Data
- Features on Karyotype
- Data Converters
  - Assembly Converter
  - ID History Converter
  - Variant Effect Predictor
  - Data Slicer**
  - Variation Pattern Finder

**Tip**

When slicing a VCF or BAM file, both the data file and its index file should be present on the web server and named correctly. The VCF file should have a ".vcf.gz" extension, and the index file should have a ".vcf.gz.tbi" extension, E.g: MyData.vcf.gz, MyData.vcf.gz.tbi. The BAM file should have a ".bam" extension, and the index file should have a ".bam.bai" extension, E.g: MyData.bam, MyData.bam.bai

Click [here](#) for more extensive documentation.

VCF / BAM File URL:

e.g. ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\_phase1\_release/ALL.chr1.phase1.projectConsensus.genotypes.vcf.gz

Region:

( e.g. 1:1-50000 )

Use VCF filters (this doesn't apply to BAM files):

- None
- By individual(s)
- By population(s) \*

(to filter by populations please provide URL to a Sample-Population Mapping File in the box below)

Sample-Population Mapping File URL:

e.g. ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\_phase1\_release/interim\_phase1.20101123.ALL.panel



Data Management

- Upload Data
- Attach DAS
- Attach Remote File
- Manage Data
- Features on Karyotype
- Data Converters
  - Assembly Converter
  - ID History Converter
  - Variant Effect Predictor
  - Data Slicer
  - Variation Pattern Finder

Thank you - your VCF file [\[6.31830969-31846823.ALL.chr6.phase1.projectConsensus.genotypes.vcf.gz\]](#) [Size: 83436] has been generated. Right click on the file name and choose "Save link as ..." from the menu

Preview

```
##fileformat=VCFv4.0
##source=BCM:SNPTools:hapfuse
##reference=1000Genomes-NCBI37
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=AP,Number=2,Type=Float,Description="Allelic Probability, P(Allele=1
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT HG00096 1
6 31831159 rs3869144 C T 100 PASS . GT:AP
6 31831167 . T C 100 PASS . GT:AP
```



### Input file

Species: Human (Homo sapiens): GRCh37 ▾

Name for this upload (optional): SLC44A4

Paste file: 

Upload file: /Users/laura/Downloads/6.

or provide file URL:

Input file format: VCF ▾

### Options

Get regulatory region consequences:

Type of consequences to display: Ensembl terms ▾

Check for existing co-located variants: Yes ▾

Return results for variants in coding regions only:

Show HGNC identifier for genes where available:

Show Ensembl protein identifiers where available:

Show HGVS identifiers for variants where available: No ▾

### Non-synonymous SNP predictions (human only)

SIFT predictions: Prediction only ▾

PolyPhen predictions: Prediction only ▾

Condel consensus (SIFT/PolyPhen) predictions: No ▾

6_31833249_A/G	<a href="#">6:31833249</a>	G	<a href="#">ENSG00000204385</a>	<a href="#">ENST00000487680</a>	Transcript	UPSTREAM	-	-	-	-
6_31833249_A/G	<a href="#">6:31833249</a>	G	<a href="#">ENSG00000204385</a>	<a href="#">ENST00000414427</a>	Transcript	DOWNSTREAM	-	-	-	-
6_31833249_A/G	<a href="#">6:31833249</a>	G	<a href="#">ENSG00000204385</a>	<a href="#">ENST00000479777</a>	Transcript	DOWNSTREAM	-	-	-	-
6_31833249_A/G	<a href="#">6:31833249</a>	G	<a href="#">ENSG00000204385</a>	<a href="#">ENST00000475563</a>	Transcript	DOWNSTREAM	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	-	<a href="#">ENSR00000487922</a>	RegulatoryFeature	REGULATORY_REGION	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">ENSG00000204386</a>	<a href="#">ENST00000495807</a>	Transcript	UPSTREAM	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">ENSG00000204386</a>	<a href="#">ENST00000480384</a>	Transcript	UPSTREAM	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">ENSG00000204386</a>	<a href="#">ENST00000491768</a>	Transcript	UPSTREAM	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">ENSG00000204386</a>	<a href="#">ENST00000375631</a>	Transcript	UPSTREAM	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">ENSG00000204386</a>	<a href="#">ENST00000479533</a>	Transcript	UPSTREAM	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">ENSG00000204385</a>	<a href="#">ENST00000229729</a>	Transcript	NON_SYNONYMOUS_CODING	1625	1604	535	R/H
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">ENSG00000204385</a>	<a href="#">ENST00000375562</a>	Transcript	NON_SYNONYMOUS_CODING	1544	1478	493	R/H
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">ENSG00000204385</a>	<a href="#">ENST00000544672</a>	Transcript	NON_SYNONYMOUS_CODING	1673	1376	459	R/H
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">Ei</a>	-	-	<a href="#">1KG 6 31833357</a>	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">Ei</a>	-	-	<a href="#">1KG 6 31833357</a>	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">Ei</a>	-	-	<a href="#">1KG 6 31833357</a>	-	-	-	-
6_31833357_C/T	<a href="#">6:31833357</a>	T	<a href="#">Ei</a>	535	R/H	cGc/cAc	<a href="#">1KG 6 31833357</a>	<b>SIFT=deleterious;</b>	<b>PolyPhen=probably_damaging;</b>	<b>Condel=deleterious</b>
6_31833612_C/G	<a href="#">6:31833612</a>	G	-	-	-	-	-	-	-	-
6_31833612_C/G	<a href="#">6:31833612</a>	G	<a href="#">Ei</a>	493	R/H	cGc/cAc	<a href="#">1KG 6 31833357</a>	<b>SIFT=deleterious;</b>	<b>PolyPhen=possibly_damaging;</b>	<b>Condel=deleterious</b>
6_31833612_C/G	<a href="#">6:31833612</a>	G	<a href="#">Ei</a>	459	R/H	cGc/cAc	<a href="#">1KG 6 31833357</a>	<b>SIFT=deleterious;</b>	<b>PolyPhen=probably_damaging;</b>	<b>Condel=deleterious</b>
6_31833612_C/G	<a href="#">6:31833612</a>	G	<a href="#">Ei</a>	-	-	-	<a href="#">1KG 6 31833357</a>	-	-	-
				-	-	-	<a href="#">1KG 6 31833357</a>	-	-	-
				-	-	-	<a href="#">1KG 6 31833357</a>	-	-	-





## Data Management

- Upload Data
- Attach DAS
- Attach Remote File
- Manage Data
- Features on Karyotype
- Data Converters
  - Assembly Converter
  - ID History Converter
  - Variant Effect Predictor
  - Data Slicer
  - Variation Pattern Finder

**i Variation Pattern Finder:**

The Variation Pattern Finder allows one to look for patterns of shared variation between individuals in the same vcf file. The finder looks for distinct variation combinations within the region, as well as individuals associated with each variation combination pattern. Only variants which have potentially functional consequences are considered, both intergenic and intronic snps are excluded. Click [here](#) for more extensive documentation.

The search will be performed on any VCF file you provided. It should be a URL for the file location. Please refer to <http://vcftools.sourceforge.net/specs.html> for VCF format specification. A URL for the latest VCF file for variation calls and genotypes released by the 1000 Genomes Project is displayed as an example below the input box. A mapping file between individual sample and population is required as well. The latest mapping file between individual sample and population released by the 1000 Genomes Project is displayed as well below the input box.

**Upload files****VCF File URL:**

[ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\\_phase1\\_release/ALL.chr6.phase1.projectConsensus.genotypes.vcf.gz](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim_phase1_release/ALL.chr6.phase1.projectConsensus.genotypes.vcf.gz)

[Clear box](#)

e.g. [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\\_phase1\\_release/ALL.chr6.phase1.projectConsensus.genotypes.vcf.gz](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim_phase1_release/ALL.chr6.phase1.projectConsensus.genotypes.vcf.gz)

**Sample-Population Mapping File URL:**

[ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\\_phase1\\_release/interim\\_phase1.20101123.ALL.panel](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim_phase1_release/interim_phase1.20101123.ALL.panel)

[Clear box](#)

e.g. [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\\_phase1\\_release/interim\\_phase1.20101123.ALL.panel](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim_phase1_release/interim_phase1.20101123.ALL.panel)

**Region:**

6:31830700-31840700

e.g. 6:46620015-46620998

[Next >](#)

Custom Data

Data Management

- Upload Data
- Attach DAS
- Attach Remote File
- Manage Data
- Features on Karyotype
- Data Converters
- Assembly Converter
- ID History Converter
- Variant Effect Predictor
- Data Slicer
- Variation Pattern Finder

### Variation Pattern Finder

Export data: [CSV](#) [Excel](#)

30016714.CEVW1AAACT.C.XIS

100%

Home Layout Tables Charts SmartArt Formulas Data Review

Edit Font Alignment Number Format

Calibri (Body) 12

General

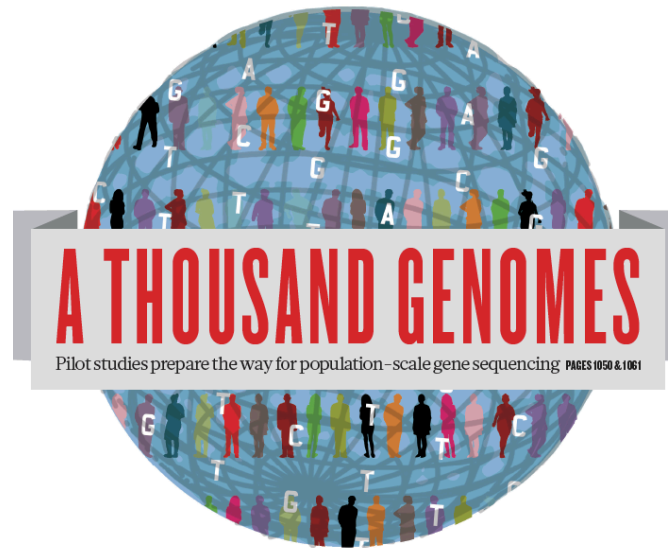
Normal

Bad

	A1	K	L	M	N	O	P	Q	R	S	T	U	V	W	X
1	MXL	PUR	TSI	YRI	Freq	rs3869144:C	6:31831167:	6:31831206:	6:31831301:	6:31831478:	6:31833357:	rs6915800:G	rs116706632	rs11712	
2						6:31831159	6:31831167	6:31831206	6:31831301	6:31831478	6:31833357	6:31833660	6:31836976	6:31837	
3						ENST000002	ENST000002	ENST000002	ENST000002	ENST000002	ENST000002	ENST000002	ENST000002	ENST000002	ENST000002
4						ENST000003	ENST000003	ENST000003	ENST000003	ENST000003	ENST000003	ENST000003	ENST000003	ENST000003	ENST000003
5						ENST000005	ENST000005	ENST000005	ENST000005	ENST000005	ENST000005	ENST000005	ENST000005	ENST000005	ENST000005
6						ENST000004	ENST000004	ENST000004	ENST000004	ENST000004	ENST000004	ENST000004	ENST000004	ENST000004	ENST000004
7	NA19758; NA	HG01048; H	NA20787; NA	NA19116; NA	0.293	C C	T T	G G	C C	T T	C C	G G	G G	G G	G G
8	NA19741; NA	HG00640; H	NA20801; NA	NA19189; NA	0.203	C C	T T	G G	C C	T T	C C	G G	G G	G G	G G
9	NA19654; NA	HG01085; H	NA20588; NA	NA19236; NA	0.195	C C	T T	G G	C C	T T	C C	G G	G G	G G	G G
10	NA19747; NA	HG01102	NA20522; NA	NA19119; NA	0.032	C T	T T	G G	C C	T T	C C	G G	G G	G G	G G
11	NA19746		NA20826; NA	NA19172; NA	0.026	T C	T T	G G	C C	T T	C C	G G	G G	G G	G G
12	NA19756		NA19160; NA	NA19160; NA	0.016	C T	T T	G G	C C	T T	C C	G G	G G	G G	G G
13			NA20752	NA19209	0.013	T C	T T	G G	C C	T T	C C	G G	G G	G G	G G
14			NA20510; NA		0.008	T C	T T	G G	C C	T T	C C	G G	G G	G G	G G
15		HG01105	NA20783; NA		0.005	C T	T T	G G	C C	T T	C C	G G	G G	G G	G G
16	NA19682; NA				0.005	C T	T T	G G	C C	T T	C C	G G	G G	G G	G G
17					0.004	T C	T T	G G	C C	T T	C C	A G	G G	G G	G G
18				NA18933	0.003	C T	T T	G G	C C	T T	C C	G A	G G	G G	G G
19				NA19098; NA	0.003	T C	T T	G G	C C	T T	C C	A G	G G	G G	G G
20	NA19740				0.003	C C	T T	G G	C C	T T	C C	G G	G A	G G	G G
21					0.003	T C	T T	G G	C C	T T	C C	G G	G G	G G	G G
22					0.002	T T	T T	G G	C C	T T	C C	A A	G G	G G	G G
23					0.002	C C	T T	G G	C C	T T	C C	G G	G G	G G	G G
24					0.002	C C	T T	G G	C C	T T	C C	G G	G G	G G	G G
25	NA19675				0.002	C C	T T	G G	C C	C T	C C	G G	G G	G G	G G
26				NA19093	0.002	C C	C T	G G	C C	T T	C C	G G	G G	G G	G G
27		HG00734		NA18486	0.002	C C	T T	G A	C C	T T	C C	G G	G G	G G	G G
28					0.001	C C	T T	G A	C C	T T	C C	G G	G G	G G	G G
29					0.001	T T	T T	G G	C C	T T	C C	G G	G G	G G	C C
30	NA19678				0.001	C C	T T	G G	C C	T C	C C	G G	G G	G G	G G
31					0.001	C C	T T	A G	C C	T T	C C	G G	G G	G G	G G
32			NA20818		0.001	C C	T T	G G	C C	T T	C C	G G	G G	G G	G G
33					0.001	T T	T T	G G	C C	T T	C C	G G	G G	G G	G G
34				NA19095	0.001	T C	T T	G G	C C	T T	C C	G G	G G	G G	G G
35					0.001	T C	T T	G G	C C	T T	C C	G G	A G	G G	G G
36				NA19152	0.001	T T	T T	G G	C C	T T	C C	G G	G G	G G	G G
37					0.001	C C	T T	G G	C G	T T	C C	G G	G G	G G	G G
38				NA19146	0.001	C T	T T	G G	C C	T T	C C	G A	G G	G G	G G
39				NA19099	0.001	C C	T C	G G	C C	T T	C C	G G	G G	G G	G G
40	NA19752				0.001	C T	T T	G G	C C	T T	C C	G G	G A	G G	G G
41					0.001	C C	T T	G G	C C	T T	C C	G G	G G	G G	G G
42					0.001	C C	T T	G G	C C	T T	C T	G G	G G	G G	G G
43		HG01083			0.001	C C	T T	G G	C C	T T	C C	G G	A G	G G	G G
44															
45															
46															

	7493:G/C	rs644827:T/C
	009	6:31838441
00229729	NON_SYNONYMOUS_CODING:Q/E	ENST00000229729 NON_SYNONYMOUS
00375562	NON_SYNONYMOUS_CODING:Q/E	ENST00000375562 NON_SYNONYMOUS
00544672	NON_SYNONYMOUS_CODING:Q/E	ENST00000544672 NON_SYNONYMOUS
00414427	NON_SYNONYMOUS_CODING:Q/E	ENST00000414427 NON_SYNONYMOUS





## The 1000 Genomes Project:

Finding out about New Data and using Data on Campus



# Announcements

- <http://1000genomes.org>
- [1000announce@1000genomes.org](mailto:1000announce@1000genomes.org)
- <http://www.1000genomes.org/1000-genomes-announcement-mailing-list>
- <http://www.1000genomes.org/announcements/rss.xml>
- <http://twitter.com/#!/1000genomes>
- [info@1000genomes.org](mailto:info@1000genomes.org)



# Thanks

- The 1000 Genomes Project Consortium
- Paul Flicek, Laura Clarke
- Richard Smith and Holly Zheng Bradley
- Giulietta Spudich and Denise Carvalho-Silva

