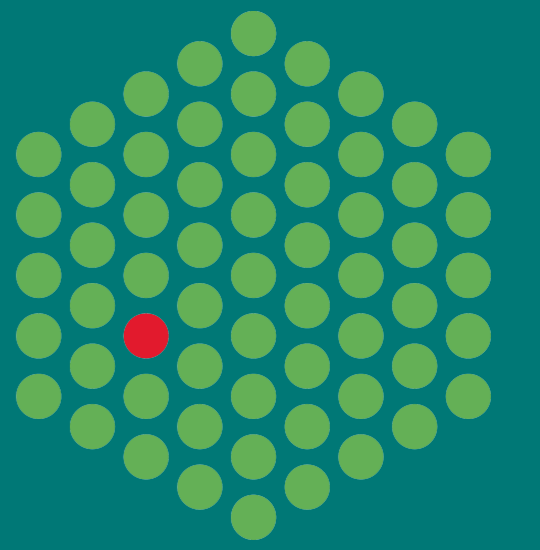


The 1000 Genomes Project Data and Tools



The main goal of the 1000 Genomes Project is to establish a comprehensive catalogue of human genome variation. We now have low coverage, whole genome and exome sequence data for more than 2500 individuals across 26 different populations.

The sequence, alignment and variant data is available from both the EBI and NCBI ftp sites

- <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp>
- <ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp>

Here we present some phase3 release stats and how to access the data and tools available from our browser (<http://browser.1000genomes.org>) to help the community utilise the data.

Phase 3 Data Release

Type	Count
biallelic_SNP	77,818,409
biallelic_del	1,930,866
biallelic_ins	1,058,174
multiallelic_SNP	259,371
multiallelic_del_ins	47,490
multiallelic_ins	34,911
biallelic_<CN0>	32,355
multiallelic_SNP_del	29,444
multiallelic_SNP_ins	28,011
multiallelic_del	13,765
biallelic_<INS:ME:ALU>	12,491
biallelic_<INS:ME:LINE1>	2,910
multiallelic_SNP_del_ins	1,248
biallelic_<INS:ME:SVA>	822
multiallelic_<CN>	169
biallelic_<INS:MT>	165
biallelic_<INV>	100
Total	81,270,701

The phase 3 release is based on 2504 unrelated individuals from 26 Populations

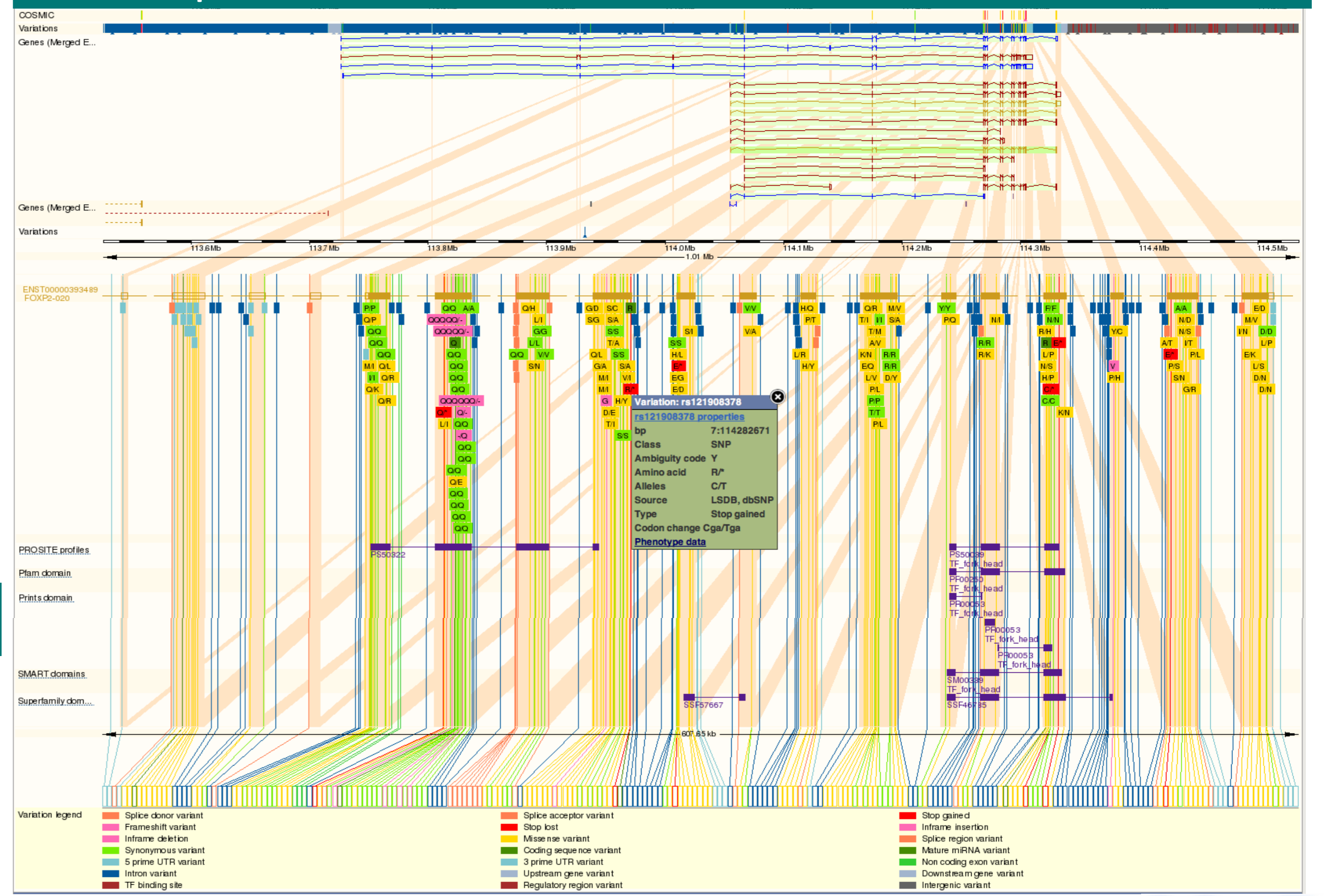
It contains phased genotypes for SNPs, short indels and large deletions like the phase 1 release. There are also new variant types including MNPs, long insertions and copy number changes.

This new release also contains multi allelic variants.

This is thanks to the new haplotype refinement pipeline which utilised Shapelt2 and MVNCall to create this dataset

Tool	Description
The Browser	Browse our variations in genomic context alongside other annotation
The Data Slicer	Enabling download of genomic subsections of our alignments and variant files
The Ensembl Variant Effect Predictor (VEP)	Providing functional consequence annotation of given variants
The VCF to PED Converter	Converting VCF files to PED files
The Variation Pattern Finder	Finding patterns of shared inheritance between different individuals in the same VCF file
FORGE	Performs functional overlap analysis of GWAS SNPs
The Allele Frequency Calculator	Provides population specific allele frequencies for a region in a VCF file

Transcript View of Variants



The Ensembl Variant Effect Predictor

The Ensembl Variant Effect Predictor allows you to provide variants and discover

- The genes affected by the variants.
- The location of the variants relative to other annotation (e.g. upstream of a transcript, in coding sequence, in non-coding RNA, in regulatory regions)
- The consequence your variants on the protein sequence (e.g. stop gained, missense, frameshift)
- Known variants that match yours, and the associated minor allele frequencies
- SIFT and PolyPhen scores for changes protein sequence

Subscribed Variation	Location	Allele	Gene	Feature	Feature type	Consequence	Position in cDNA	Position in CDS	Position in protein	Amino acid change	Codon change	Co-located Variation	Extra
1.881907_<C>	1.881906-881907	C	ENSG0000013763	ENST0000021465	Transcript	downstream_gene_variant	-	-	-	-	-	-	DISTANCE=3502

Variation Pattern Finder

The Variation Pattern Finder calculates patterns of shared inheritance between individuals. In a genomic region, some individuals will share distinct combinations of inheritance. The finder calculates:

- The distinct patterns of inheritance
- Which individuals have a particular pattern.
- The frequency of those patterns in each population.
- The pattern finder focuses on variants with functional consequences.

Here is example output:

ASW	CEU	CHB	CHS	CLM	FIN	GBR	JPT	LWK	MLL	PUR	TSI	YRI	Freq	r11885371-G/A	r14977987-T/C	r20024870-C/T	r9399628-C/T
NA19908	NA18815	NA19058	NA19058	NA19058	NA19058	NA19058	NA19058	NA19058	NA19058	NA19058	NA19058	NA19058	0.082	-	-	-	-

Allele Frequency Calculator

CHR	POS	ID	REF	ALT	TOTAL_CNT	ALT_CNT	FRQ
22	17004085	rs182269758	A	G	170	9	0.05
22	17004141	rs192917218	A	G	170	2	0.011

The Allele Frequency Calculator takes:

- A VCF file
- A sample to population mapping file
- A chromosomal region
- A population name

It calculates population-wide allele frequency for sites within the defined region. Here is example output.

The International Genome Sample Resource Beyond 1000 Genomes Project

The International Genome Sample Resource (IGSR), a Wellcome Trust funded project that will be built on the foundation of the 1000 Genomes Project starts in 2015.

IGSR plans to:

- Maintain the existing 1000 genomes data and move to GRCh38
- Collect other data sets generated on the Coriell Cell Lines including Geuvadis
- Add new populations to explain the global diversity of the variant catalog

Acknowledgements

We would like to thank the Ensembl variation team for all their help and the Wellcome Trust and European Molecular Biology Laboratory for their funding.

