

The Variant Call Format and VCFtools

Petr Danecek¹, Adam Auton², Goncalo Abecasis³, Cornelis A. Albers¹, Eric Banks⁴, Mark A. DePristo⁴, Bob Handsaker⁴, Gerton Lunter⁵, Garbor Marth⁶, Steve Sherry⁷, Gilean McVean⁸, Richard Durbin^{1,*} and 1000 Genomes Project Analysis Group⁹

¹Wellcome Trust Sanger Institute, Cambridge, CB10 1SA, UK; ²University of Oxford, Wellcome Trust Centre for Human Genetics, Oxford, OX3 7BN, UK; ³Center for Statistical Genetics, Department of Biostatistics, University of Michigan, Ann Arbor, M48109, USA; ⁴Broad Institute of MIT and Harvard, Cambridge, MA 02141, USA; ⁵University of Oxford, Department of Physiology, Anatomy and Genetics, Oxford, OX1 3QX, UK; ⁶Boston College, Department of Biology, MA 02467, USA; ⁷National Institutes of Health National Center for Biotechnology Information, MD 20894, USA; ⁸University of Oxford Department of Statistics, Oxford, OX1 3TG, UK; ⁹<http://www.1000genomes.org>

Abstract

One of the main uses of next-generation sequencing is to discover variation amongst large populations of related samples. Recently the format for storing next-generation read alignments has been standardised by the SAM/BAM file format specification. This has significantly improved the interoperability of next-generation tools for alignment, visualisation, and variant calling. We propose the Variant Call Format (VCF) as a standardised format for storing the most prevalent types of sequence variation, including SNPs, indels and larger structural variants, together with rich annotations. VCF is usually stored in a compressed manner and can be indexed for fast data retrieval of variants from a range of positions on the reference genome. The format was developed for the 1000 Genomes Project, and has also been adopted by other projects such as UK10K, dbSNP, or the NHLBI Exome Project. VCFtools is a software suite that implements various utilities for processing VCF files, including validation, merging and comparing, and also provides a general Perl and Python API. The VCF specification and VCFtools are available from <http://vcftools.sourceforge.net>.

Example

```
##fileformat=VCFv4.0
##fileDate=20100707
##source=VCFtools
##reference=NCBI36
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality (phred score)">
##FORMAT=<ID=GL,Number=3,Type=Float,Description="Likelihoods for RR,RA,AA genotypes (R=ref,A=alt)">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##ALT=<ID=DEL,Description="Deletion">
##INFO=<ID=SVTYPE,Number=1,Type=String,Description="Type of structural variant">
##INFO=<ID=END,Number=1,Type=Integer,Description="End position of the variant">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE1 SAMPLE2
1 1 . ACG A,AT . PASS . GT:DP 1/2:13 0/0:29
1 2 rs1 C T,CT . PASS H2;AA=T GT:GQ 0|1:100 2/2:70
1 5 . A G . PASS . GT:GQ 1|0:77 1/1:95
1 100 T <DEL> . PASS SVTYPE=DEL;END=300 GT:GQ:DP 1/1:12:3 0/0:20
```

Mandatory header lines (lines starting with ##)

Optional header lines (meta-data about the annotations in the VCF body) (lines starting with #)

Reference alleles (GT=0)

Alternate alleles (GT>0 is an index to the ALT column)

Phased data (G and C above are on the same chromosome)

Deletion (row 4, ALT:)

SNP (row 2, ALT: A,AT)

Large SV (row 4, ALT: , SVTYPE=DEL;END=300)

Insertion (row 2, ALT: T,CT)

Other event (row 4, ALT:)

Types of variants

SNPs

Alignment	VCF representation
ACGT	POS REF ALT
ATGT	2 C T

Insertions

Alignment	VCF representation
AC-GT	POS REF ALT
ACTGT	2 C CT

Deletions

Alignment	VCF representation
ACGT	POS REF ALT
A--T	1 ACG A

Complex events

Alignment	VCF representation
ACGT	POS REF ALT
A-TT	1 ACG AT

Large structural variants

VCF representation
POS REF ALT INFO
100 T SVTYPE=DEL;END=300

VCF highlights

- Meta-data - flexible and extensible
- Text format - easy to generate and parse
- Stored compressed - compact size
- Indexed by **tabix** - fast random access by genomic position
- Open source implementation - VCFtools, GATK, ... (C++, Java, general Perl and Python API)

Extensible meta-data

Annotations may apply to the variant as a whole (the **INFO** column) or to each genotype (the **FORMAT** column). In addition to genotype, other commonly used annotations include genotype likelihoods, dbSNP membership, ancestral allele, read depth, mapping quality, and others.

VCFtools

- Format validation
- Annotating
- Comparing, calculating basic statistics
- Merging
- Creating intersections and complements

Examples

```
# Validate VCF files
vcf-validator file.vcf.gz
```

```
# Compare VCF files
compare-vcf A.vcf.gz B.vcf.gz C.vcf.gz
```

```
# List positions present in at least two of the files
vcf-isec -n +2 A.vcf.gz B.vcf.gz C.vcf.gz > out.vcf
```