

The variant call format and VCFtools

Petr Danecek^{1,†}, Adam Auton^{2,†}, Goncalo Abecasis³, Cornelis A. Albers¹, Eric Banks⁴, Mark A. DePristo⁴, Robert E. Handsaker⁴, Gerton Lunter², Gabor T. Marth⁵, Stephen T. Sherry⁶, Gilean McVean^{2,7}, Richard Durbin^{1,*} and 1000 Genomes Project Analysis Group[‡]

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

Associate Editor: John Quackenbush

ABSTRACT

Summary: The variant call format (VCF) is a generic format for storing DNA polymorphism data such as SNPs, insertions, deletions and 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 and the NHLBI Exome Project. VCFtools is a software suite that implements various utilities for processing VCF files, including validation, merging, comparing and also provides a general Perl API.

Availability: <http://vcftools.sourceforge.net>

Contact: rd@sanger.ac.uk

Received on October 28, 2010; revised on May 4, 2011; accepted on May 28, 2011

1 INTRODUCTION

One of the main uses of next-generation sequencing is to discover variation among large populations of related samples. Recently, a format for storing next-generation read alignments has been standardized by the SAM/BAM file format specification (Li *et al.*, 2009). This has significantly improved the interoperability of next-generation tools for alignment, visualization and variant calling. We propose the variant call format (VCF) as a standardized format for storing the most prevalent types of sequence variation, including SNPs, indels and larger structural variants, together with rich annotations. The format was developed with the primary intention to represent human genetic variation, but its use is not restricted to diploid genomes and can be used in different contexts as well. Its flexibility and user extensibility allows representation of a wide variety of genomic variation with respect to a single reference sequence.

Although generic feature format (GFF) has recently been extended to standardize storage of variant information in genome variant format (GVF) (Reese *et al.*, 2010), this is not tailored for storing information across many samples. We have designed the VCF format to be scalable so as to encompass millions of sites with genotype data and annotations from thousands of samples. We have adopted a textual encoding, with complementary indexing, to allow easy generation of the files while maintaining fast data access. In this article, we present an overview of the VCF and briefly introduce the companion VCFtools software package. A detailed format specification and the complete documentation of VCFtools are available at the VCFtools web site.

2 METHODS

2.1 The VCF

2.1.1 Overview of the VCF A VCF file (Fig. 1a) consists of a header section and a data section. The header contains an arbitrary number of meta-information lines, each starting with characters '##', and a TAB delimited field definition line, starting with a single '#' character. The meta-information header lines provide a standardized description of tags and annotations used in the data section. The use of meta-information allows the information stored within a VCF file to be tailored to the dataset in question. It can be also used to provide information about the means of file creation, date of creation, version of the reference sequence, software used and any other information relevant to the history of the file. The field definition line names eight mandatory columns, corresponding to data columns representing the chromosome (CHROM), a 1-based position of the start of the variant (POS), unique identifiers of the variant (ID), the reference allele (REF), a comma separated list of alternate non-reference alleles (ALT), a phred-scaled quality score (QUAL), site filtering information (FILTER) and a semicolon separated list of additional, user extensible annotation (INFO). In addition, if samples are present in the file, the mandatory header columns are followed by a FORMAT column and an arbitrary number of sample IDs that define the samples included in the VCF file. The FORMAT column is used to define the information contained within each subsequent genotype column, which consists of a colon separated list of fields. For example, the FORMAT field GT:GQ:DP in the fourth data entry of Figure 1a indicates that the subsequent entries contain information regarding the genotype, genotype quality and read depth for each sample. All data lines are TAB delimited and the number of fields in each data line must match the number of fields in the header line. It is strongly recommended that all annotation tags used are declared in the VCF header section.

*To whom correspondence should be addressed.

[†]The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint First Authors.

[‡]<http://www.1000genomes.org>

gives an overall quality score for the assertion made in ALT that the site is variant or no variant.

INFO column:

- DB, dbSNP membership;
- H3, membership in HapMap3;
- VALIDATED, validated by follow-up experiment;
- AN, total number of alleles in called genotypes;
- AC, allele count in genotypes, for each ALT allele, in the same order as listed;
- SVTYPE, type of structural variant (DEL for deletion, DUP for duplication, INV for inversion, etc. as described in the specification);
- END, end position of the variant;
- IMPRECISE, indicates that the position of the variant is not known accurately; and
- CIPOS/CIEND, confidence interval around POS and END positions for imprecise variants.

Missing values are represented with a dot. For practical reasons, the VCF specification requires that the data lines appear in their chromosomal order. The full format specification is available at the VCFtools web site.

2.1.3 Variation types VCF is flexible and allows to express virtually any type of variation by listing both the reference haplotype (the REF column) and the alternate haplotypes (the ALT column). This permits redundancy such that the same event can be expressed in multiple ways by including different numbers of reference bases or by combining two adjacent SNPs into one haplotype (Fig. 1g). Users are advised to follow recommended practice whenever possible: one reference base for SNPs and insertions, and one alternate base for deletions. The lowest possible coordinate should be used in cases where the position is ambiguous. When comparing or merging indel variants, the variant haplotypes should be reconstructed and reconciled, such as in the Figure 1g example, although the exact nature of the reconciliation can be arbitrary. For larger, more complex, variants, quoting large sequences becomes impractical, and in these cases the annotations in the INFO column can be used to describe the variant (Fig. 1f). The full VCF specification also includes a set of recommended practices for describing complex variants.

2.1.4 Compression and indexing Given the large number of variant sites in the human genome and the number of individuals the 1000 Genomes Project aims to sequence (Durbin *et al.*, 2010), VCF files are usually stored in a compact binary form, compressed by bgzip, a program which utilizes the zlib-compatible BGZF library (Li *et al.*, 2009). Files compressed by bgzip can be decompressed by the standard gunzip and zcat utilities. Fast random access can be achieved by indexing genomic position using tabix, a generic

indexer for TAB-delimited files. Both programs, bgzip and tabix, are part of the samtools software package and can be downloaded from the SAMtools web site (<http://samtools.sourceforge.net>).

2.2 VCFtools software package

VCFtools is an open-source software package for parsing, analyzing and manipulating VCF files. The software suite is broadly split into two modules. The first module provides a general Perl API, and allows various operations to be performed on VCF files, including format validation, merging, comparing, intersecting, making complements and basic overall statistics. The second module consists of C++ executable primarily used to analyze SNP data in VCF format, allowing the user to estimate allele frequencies, levels of linkage disequilibrium and various Quality Control metrics. Further details of VCFtools can be found on the web site (<http://vcftools.sourceforge.net/>), where the reader can also find links to alternative tools for VCF generation and manipulation, such as the GATK toolkit (McKenna *et al.*, 2010).

3 CONCLUSIONS

We describe a generic format for storing the most prevalent types of sequence variation. The format is highly flexible, and can be adapted to store a wide variety of information. It has already been adopted by a number of large-scale projects, and is supported by an increasing number of software tools.

Funding: Medical Research Council, UK; British Heart Foundation (grant RG/09/012/28096); Wellcome Trust (grants 090532/Z/09/Z and 075491/Z/04); National Human Genome Research Institute (grants 54 HG003067, R01 HG004719 and U01 HG005208); Intramural Research Program of the National Institutes of Health, the National Library of Medicine.

Conflict of Interest: none declared.

REFERENCES

- Durbin, R.M. *et al.* (2010) A map of human genome variation from population-scale sequencing. *Nature*, **467**, 1061–1073.
- Li, H. *et al.* (2009) The sequence alignment/map format and SAMtools. *Bioinformatics*, **25**, 2078–2079.
- McKenna, A.H. *et al.* (2010) The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.*, **20**, 1297–1303.
- Reese, M.G. *et al.* (2010) A standard variation file format for human genome sequences. *Genome Biol.*, **11**, 20796305.