

Genetics and population analysis

Sparse Project VCF: efficient encoding of population genotype matrices

Michael F. Lin ^{1,*}, Xiaodong Bai², William J. Salerno² and Jeffrey G. Reid²

¹mclin.net LLC, San Jose, CA 95113, USA and ²Department of Regeneron Pharmaceuticals, Inc., Regeneron Genetics Center, Tarrytown, NY 10591, USA

*To whom correspondence should be addressed.
Associate Editor: Russell Schwartz

Received on September 27, 2020; revised on November 13, 2020; editorial decision on November 16, 2020; accepted on November 20, 2020

Abstract

Summary: Variant Call Format (VCF), the prevailing representation for germline genotypes in population sequencing, suffers rapid size growth as larger cohorts are sequenced and more rare variants are discovered. We present Sparse Project VCF (spVCF), an evolution of VCF with judicious entropy reduction and run-length encoding, delivering $>10\times$ size reduction for modern studies with practically minimal information loss. spVCF interoperates with VCF efficiently, including tabix-based random access. We demonstrate its effectiveness with the DiscovEHR and UK Biobank whole-exome sequencing cohorts.

Availability and implementation: Apache-licensed reference implementation: github.com/mclin/spVCF.

Contact: dna@mclin.net

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Variant Call Format (VCF) is the prevailing representation for small germline variants discovered by high-throughput sequencing (Danecek *et al.*, 2011). In addition to capturing variants sequenced in one study participant, VCF can represent the genotypes for many participants at all discovered variant loci. This ‘Project VCF’ (pVCF) form is a 2-D matrix with loci down the rows and participants across the columns, filled in with each called genotype and annotations thereof, including quality-control (QC) measures like read depth, strand ratio and genotype likelihoods.

As the number of study participants N grows (columns), more variant loci are also discovered (rows), leading to super-linear growth of the pVCF genotype matrix. And, because cohort sequencing discovers mostly rare variants, this matrix consists largely of reference-identical genotypes and their high-entropy QC measures. In recent experiments with human whole-exome sequencing (WES), doubling N from 25 000 to 50 000 also increased the pVCF locus count by 43%, and 96% of all loci had non-reference allele frequency below 0.1% (Lin *et al.*, 2018). Empirically, vcf.gz file sizes in WES and whole-genome sequencing (WGS) are growing roughly with $N^{1.5}$ in the largest studies as of this writing ($N \approx 500\,000$ WES). Unchecked, we project $N = 1\,000\,000$ WGS will yield petabytes of compressed pVCF.

2 Approach

We sought an incremental solution to these challenges for existing pVCF-based pipelines, which may be reluctant to adopt fundamentally

different formats or data models (Danek and Deorowicz, 2018; Deorowicz and Danek, 2019; Lan *et al.*, 2020; Layer *et al.*, 2015; Li, 2016; Zheng *et al.*, 2017; Supplementary Appendix S1) to minimize disruption to existing processes and users. To this end, we developed Sparse Project VCF (spVCF), which adds three simple features to VCF (Fig. 1):

1. *Squeezing: judiciously reducing QC entropy.* In those cells with zero reads supporting a variant (typically Allele Depth $AD = d, 0$ for any d) and corresponding non-variant genotype, we discard all fields except the genotype GT and the read depth DP, which we also round down to a power of two (0, 1, 2, 4, 8, 16, ...; configurable). Any cell reporting evidence of variation retains its original QC measures and other annotations.

This convention, inspired by common base quality score compression techniques, aims to preserve nearly all *useful* information, removing minor fluctuations in non-variant cells. (If required for compatibility, non-variant genotype likelihoods could be approximated from depth, albeit without read quality inputs that might subtly affect downstream calculations.)

1. *Succinct, lossless encoding for runs of reference-identical cells.* First, we replace the contents of a reference-identical (or non-called) cell with a double-quotation mark if it’s identical to the cell above it, compressing runs down the column for each sample. Then we run-length encode these quotation marks across the rows, so for example a stretch of 42 marks across a row is written `<tab>`

