VCF Miner: A Platform and Tools for Collaborative Information Mining from Next-Generation Sequence Variant Data

Zhi-Liang Hu, James E. Koltes, Eric Fritz-Waters and James M. Reecy

Department of Animal Science and Center for Integrated Animal Genomics, Iowa State University, 2255 Kildee Hall, Ames, IA 50011

Abstract

High-throughput next-generation sequencing (NGS) technology has provided powerful tools for genetics/genomics studies. It has also imposed a great challenge on researchers to determine how to efficiently handle and quickly glean information for genotype/phenotype studies. The ability to quickly process NGS data for useful information screening is important before the extensive analysis may be carried out, because the latter is usually time consuming and demands large computation power. The variant call format (VCF) that stores variations against a reference genome is one of such data formats at the heart of information crossroads of the genotype/phenotype data analysis pipelines. We propose a shared platform for collaborative VCF file storage, handling, information abstraction, querying, and data re-use. A unix-file-system based data file repository has been built, made available Plink/SEQ, BEDOPS, and BEDtools for VCF file handling as well as a number of in-house scripts to serve various purposes during the data process An initial collection of 41 VCF files from four laboratories were used as a test case. The goals of this project were to allow VCF data from multiple collaborative projects to be managed in a central location to This resource facilitate pre-computed data abstractions and query. will allow data re-use in subsequent, possibly more diverse types of data analysis in an efficient manner.

Introduction

The amount of new NGS data flowing to genetics lab is phenomenal. While researchers enjoy the increasing amount of genotype data to use, the speed and capacity to handle and possibility to re-use these data are challenging. The purpose of this project was to help researchers to meet these challenges. At this time we focused on pre-abstracting information from VCF files in order to save time for researchers

Preliminary Results

We have built a UNIX file-system based VCF file repository. The repository is equipped with existing and custom software for information abstraction, and enabled MvSQL database for information handling. We have also developed preliminary web interfaces for end users. Figures 1 and 2 show introductory information on how the web site and its front end work. The main idea behind this scheme is to pre-compute some basic but much needed statistics for users to get a summary of each VCF file to guide further analysis

There are already a few developing as well as mature software available that handle VCF files, such as vcftools, BEDOPS, and PLINK/Seq (1, 2, and 3). They are diverse in terms of functions yet practical to use (Tables 1 and 2). We have adopted these software locally in our pipeline. In Box 1 is shown a conceptual work flow in terms of data processes and jobs done at certain steps.

We were fortunate to have received 41 sample VCF files from five laboratories (Figure 1) for initial test of concept. In Tables 3, 4, and Box 2, are shown some pre-computed statistics, including number of genotypes with a minor allele, homozygotes, heterozygotes, genotyping rate, etc. It's worth pointing out that this is only our start of this work, and the preliminary data shown here is only a very small example of potential features we may facilitate with the repository when it's fully developed.

Table 3. Sample output from "pseq" tool to show combined VCF statistics. (a) "i-stats" on combined VCF file (called "projects" in PSEQ's realm); (b) "i-stats" on single VCF file with multiple animals.

 Dall D
 NAIT
 NMIN
 NHET
 NVAR
 RATE
 SING
 TITV
 PASS
 <

 NALT
 NMIN
 NHET
 NVAR
 RATE
 SING
 TIV
 PASS
 P

 7505165
 4893843
 4228099
 12.82+07
 0986531
 1138/02
 1.47731
 4496583

 730548
 475403
 3995701
 2:74+07
 983163
 105776
 1.46924
 453373

 750500
 4903176
 4327976
 1.28+07
 989388
 1102893
 1.48216
 4512845

 7596746
 5029035
 4355703
 1.28+07
 0.985318
 11027503
 1.47145
 4512845

 7596881
 493141
 4273745
 1.28+07
 0.985131
 1151957
 1.4883
 451384

naturadud ID Number of non-reference genetypes Number of genetypes with a minor allele Number of herecyclon genetypes (or individual Total number of called variants for individual Genetyping rate for individual Number of singletons individual has Number of singletons individual has a nonreference ge Number of variants PASSing for which individual has a nonreference Mean Warth OF for variants for which individual has a nonreference Mean variant DF for variants of or which individual has a nonreference Mean variant DF for variants of or which individual has a nonreference Mean variant DF for variants of providents for which individual has a nonreference Mean variant DF for variants of the individual has a nonreference Network of the second providents for which individual has a nonreference ge

992 5151410 3081217 3081870 5.15E+06 0.222664 2656056 0.83554

NHET

Individual ID

(a)

(b) A1476 A1691

A2452 A4378

Legend: ID NALT



Figure 1

The front page of the VCF file Data Repository showing: (1) User login is required for data contributions and data meta-information management (refer to Figure 2); (2) Automated data repositry overview with some basic statistics; (3) Main features offered by the

S NRSP-8 NAGRP VCF Data Repository analysis is the study of gene it Call Format (VCF) has bee variant data from DNA serve The purpose of the NAGRP VCF Data Repository is to fact collaborative VCF file storage, handling, information abstr querying, and data re-use. Further access to the raw da authorized by the data owner. Pass Login Current data collection at a gla Species Number of files of eninels Species . Species Peported SNPs cattle : 46,858,235 chicken : 14,506,315 korse : 12,933,913 ARS, USDA : 2 ARG, USDA Iowa State University University of Hinnesota University of Hissouri Contributing P1 Number of files Carrie Finno 1 GATE Unit Van Raseil 1 2 James Recey 1 24 Gallow Hoters D. Schaberl 1 1 Naw GATE 1 Number of files Unknown 1 hws/GATE 1

- Simple statistics such as counts of each SNPs, homozyge heterozygotes, etc. (vcf-stats to jason format).
- Estimation of allele fin
- Merging multiple VCF files for combined analysis
- Build "projects" from multiple VCF files with PLINZ/BEQ tools for combined analysis (if 1 animal per VCF file, this is to bring multiple animals together ul Basic statistics (e.g. v-stats, i-stats, and g-stats) (ref)
- sic statistics (e.g. v-stats, i-stats, and g-stats) (ref) atistics across all variants (e.g. non-reference genotypes, number of genotypes with a minor allel mber of heterocrygous genotypes for an individual, total number of called variants for an individual notypen grate for an individual, etc.)
- Convert VCF files to BED format to utilize BEDOPS tools for data abstraction and analysis (ref)
- ations: extract features, match features, et tics: con stical operations by mapping overlapping features, merging files, et I File management: file starch for lexicographical speed sorting, data extraction, ci
- Pre-compute some summary statistics, including but not limited to: Filtering SNPs based on sequencing depth and frequencies
- Il Estimating allele frequencies among multiple animals Il Counting the number of homozygotes, heterozygotes, and phased/unphased SNPs (vcftools)

To participate:

- Follow these simple steps to begin using the data repository:

- Step 1: Submit general information about your VCP files.
 Step 2: Upbad your VCP files.
 Step 2: Access pre-computed data summary statistics.
 Step 4: Get in tooch with the owners of data that you are interested in to start further collaborations.

Table 1. By adopting BEDOPS tools in the VCF Data Repository, a number of functions will become available.

ping BED

e any n

File m

<u>sort-bed</u> - apply lexicographical sort to BED data

• <u>starch</u> and <u>unstarch</u> - compress and extract BED data

<u>Conversion tools</u> - convert common genomic formats to BED

assed arch

-> 5979847,

'nalt_1 ---'snp_count' -> 5% .--'shared' -> { '1' -> 597984}

-> ('9111' -> (

• starchcat - merge com

Statistics bedops - apply set o bedmap - map overlap
 elements onto target and optionally compubedextract - efficiently extract BED closest-features - matches features between BED files

NA 13.4638 NA 13.4688 NA 13.4688 NA 13.9661 NA 14.6209

NA 12 850

692419 622.393 72.270 602144 629.885 73.384

658255 621.247 72.466 777631 630.65 72.447

389198

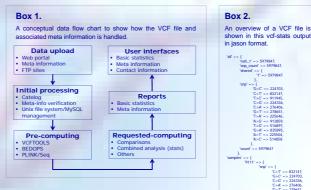


Table 2. Function sets from vcftools

Name of Tools Functions

- vcftools To analyse VCF files (multiple optoins available) vcf-stats Obtain basic statistics of a VCF file vcf-to-tab Generate tab-delimited list of genotype metrix Convert between VCF version vcf-convert
- vcf-merge Merge multiple VCF files vcf-compare Compare two or more VCF files for differences
- vcf-querv Query specific chromosome region for information

Figure 2.

A sample VCF file meta-information entry form shows the kind of information needed in order for the Repository to properly handle and process the data file. Note the VCF file contributor has access to update these information once logged in.

PI's name	Carrie Finon	
PI's email	fire0100@ump.edu	
Contact's name	Carrie Fino	
Contact's made		
and the second sec	e;finno@gmail.com	
Institute / Affiliation:	University of Minnesota College of Vet Med	
CE File information:		
Species	horse	
Project		
	2807 Q20 for minimum Phred vari- call score.	
Animals (IDs, breeds, etc):	1691, 4378, 4468, 1476, 2452	
Samples (tissue, etc):		•
Genotype platform:	Illumine HiSeq 2000	•
VCF file generator (software):	bwa/GATE	•
VCF file names:	BelgienHorses_DNASeq_FINNO.vof.gz	•
Phenotype file names:		
Genotype file names:		
Notes:	All horses are male castrated B Draft horses.	elgi

Discussions

The VCF Miner will have several potentials to facilitate genetics research. For example, to make breed and species level variant comparisons; to identify novel variants by cross breeds/species comparisons, to visualize variants in the context of QTL and other genome features, etc. Our long term goal is to build tools making the VCF information more readily accessible to bring together useful genotype, phenotype and annotation data for queries and comparisons, thus increasing the re-usability of data and chances for new discoveries.

We will look out and work with other public platforms on using and sharing VCF information that may complement our efforts. For example, we will check with NCBI Variation Database and Ensembl Variation Database for functions they make available and in ways how we may complement their works as well.

Table 4. Allele frequency estimates by vcftools, showing SNP types substitutions, insertions/deletions (number of animals=5).

POS 860 942 1237 1495	N_ALLELES 2 2 2	N_CHR 10 10	{ALLELE:FRE	G:0.4
942 1237	2			
1237	-	10	m:0.2	
	2			C:0.8
1 405		10	T:0.2	TTTTG:0.8
1495	2	10	A:0.2	AAGAT:0.8
1541	2	10	C:0.9	T:0.1
1939	2	10	C:0.6	G:0.4
2123	2	10	G:0.9	GT:0.1
2296	2	10	A:0.9	C:0.1
3348	2	10	G:0	GT:1
3688	2	10	C:0.9	T:0.1
4247	2	10	T:0.2	C:0.8
4431	2	10	TTTTA:0.6	т:0.4
4625	2	10	A:0.9	C:0.1
	1939 2123 2296 3348 3688 4247 4431	1939 2 2123 2 2296 2 3348 2 3688 2 4247 2 4431 2	1939 2 10 2123 2 10 22966 2 10 3348 2 10 3688 2 10 4247 2 10 4431 2 10	1939 2 10 C:0.6 2123 2 10 G:0.9 2296 2 10 A:0.9 3348 2 10 G:0 3688 2 10 C:0.9 4247 2 10 TTTA:0.6

References

- 1. Danecek P. Auton A. Abecasis G. Albers CA. Banks E. DePristo MA. Handsake RE, Lunter G, Marth GT, Sherry ST, McVean G, Durbin R; 1000 Genomes Project Analysis Group. (2011). The variant call format and VCFtools. Bioinformatics. 2011 Aug 1;27(15):2156-8. doi: 10.1093/bioinformatics/btr330. Foub 2011 Jun 7. Web site Epub 2011 Jun 7. W http://vcftools.sourceforge.net/ accessed on December 1, 2013. Web
- 2. Shaun Purcell et al. (2011), PLINK/Seq Analysis of genetic variation data from Analysis of genetic validon data from large-scale, population-based medical sequencing studies. Workshop at the Analytic and Translational Genetics Unit, MGH; Center for Human Genetic Research, MGH. Web site https://atgu.mgh.harvard.edu/plinkseq/ last accessed on December 21, 2013.
- З. Shane Neph, M. Scott Kuehn, Alex Reynolds, et al. BEDOPS: high-performance genomic feature operations Bioinformatics (2012) 28 (14): 1919-1920 doi: 10.1093/bioinformatics/bts277. Web site <u>https://bedops.readthedocs.org</u> last accessed on December 11, 2013.

Acknowledgements

We sincerely thank Dr. Carrie Finno from University of Minnesota, Dr. Robert D. Schnabel from University of Missouri, Dr. Curt Van Tassell from USDA-ARS, and Dr. Susan Lamont from Iowa State University for kindly sharing their VCF files for the experimental works testing various tools on the NAGRP VCF Data Repository to proof the concept. We also thank Xue-feng Zhao and James Coyle from the Baker & HPC Centers at Iowa State University for their assistance using the ISU computing facilities.







