

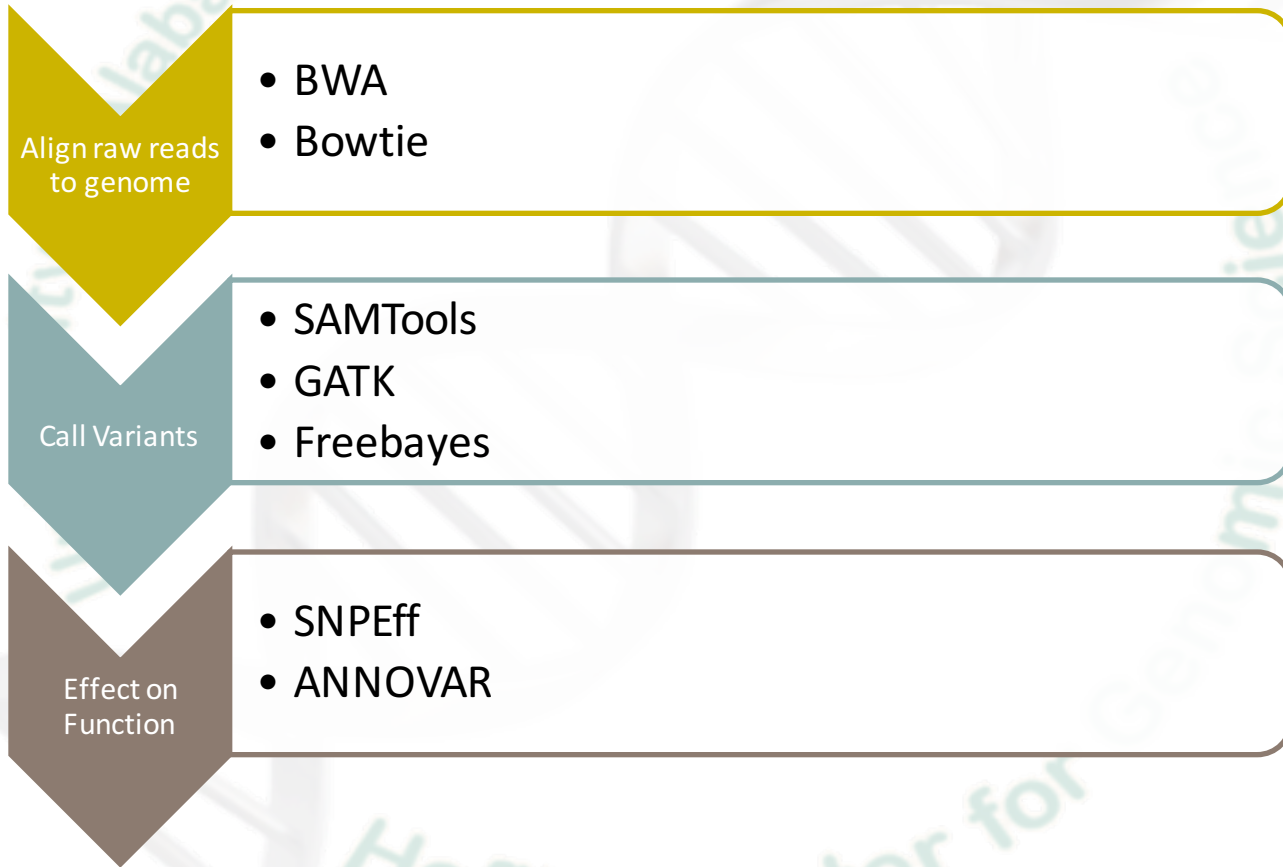


Variant Analysis of Exome/Genome Sequencing Data

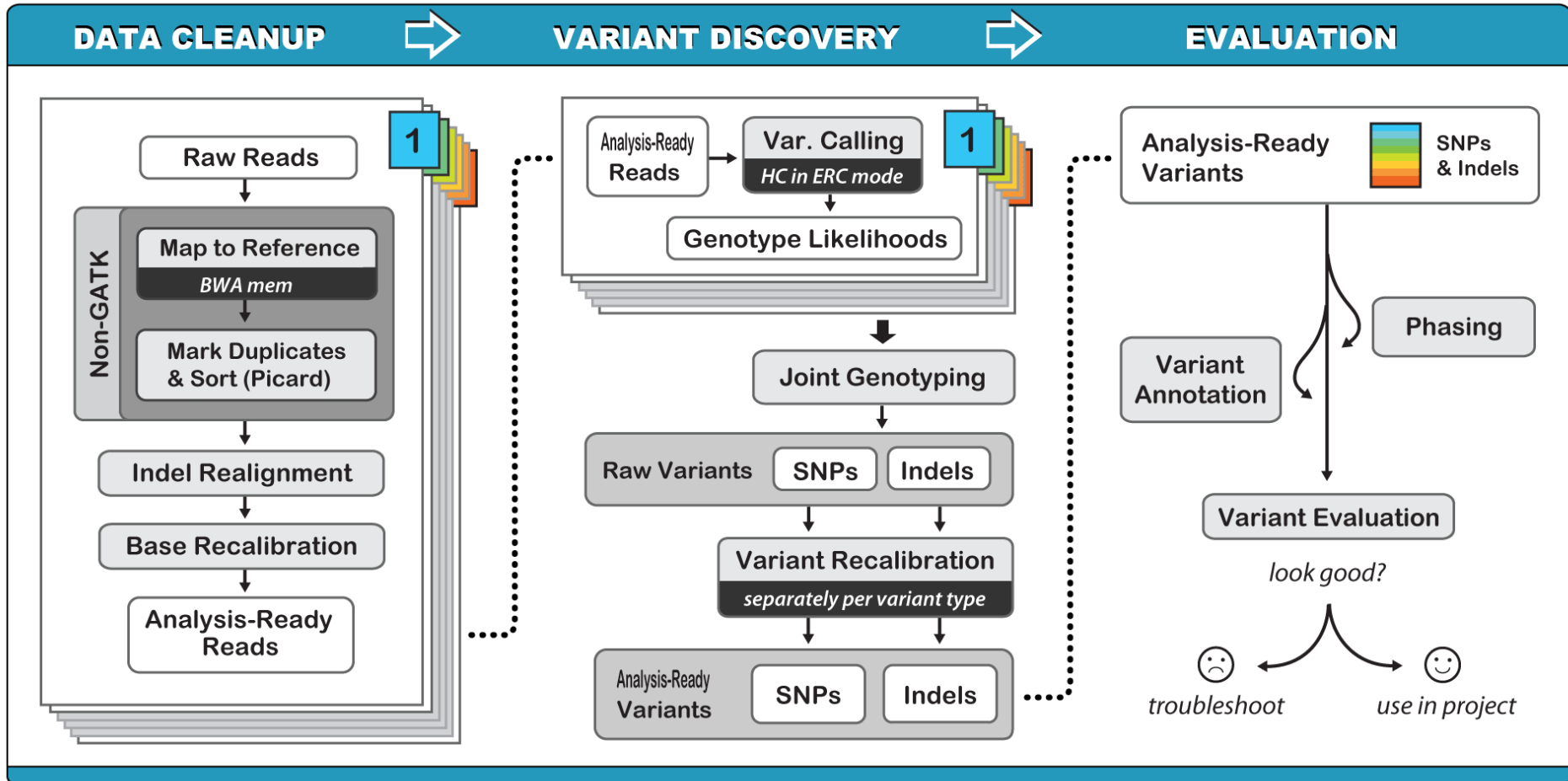
David Crossman, Ph.D.
UAB Heflin Center for Genomic Science

Immersion Course

Whole Genome/Exome (DNA-Seq) analysis pipeline



GATK pipeline



GATK Best Practices

(<http://www.broadinstitute.org/gatk/>)

Best Practice Variant Detection with the GATK v4, for release 2.0

There are 18 comments on this article. To see them or add your own, read this post on the forum →

Introduction

1. The basic workflow

Our current best practice for making SNP and indel calls is divided into four sequential steps: initial mapping, refinement of the initial reads, multi-sample indel and SNP calling, and finally variant quality score recalibration. These steps are the same for targeted resequencing, whole exomes, deep whole genomes, and low-pass whole genomes. Example commands for each tool are available on the individual tool's wiki entry. [There is also a list of which resource files to use with which tool.](#)

Note that due to the specific attributes of a project the specific values used in each of the commands may need to be selected/modified by the analyst. Care should be taken by the analyst running our tools to understand what each parameter does and to evaluate which value best fits the data and project design.

2. Lane, Library, Sample, Cohort

There are four major organizational units for next-generation DNA sequencing processes that used throughout this documentation:

- **Lane:** The basic machine unit for sequencing. The lane reflects the basic independent run of an NGS machine. For Illumina machines, this is the physical sequencing lane.
- **Library:** A unit of DNA preparation that at some point is physically pooled together. Multiple lanes can be run from aliquots from the same library. The DNA library and its preparation is the natural unit that is being sequenced. For example, if the library has limited complexity, then many sequences are duplicated and will result in a high duplication rate across lanes.
- **Sample:** A single individual, such as human CEPH NA12878. Multiple libraries with different properties can be constructed from the original sample DNA source. Here we treat samples as independent individuals whose genome sequence we are attempting to determine. From this perspective, tumor / normal samples are different despite coming from the same individual.
- **Cohort:** A collection of samples being analyzed together. This organizational unit is the most subjective and depends intimately on the design goals of the sequencing project. For population discovery projects like the 1000 Genomes, the analysis cohort is the ~100 individual in each population. For exome projects with many samples (e.g., ESP with 800 EOMI samples) deeply sequenced we divide up the complete set of samples into cohorts of ~50 individuals for multi-sample analyses.

This document describes how to call variation within a single analysis cohort, comprised for one or many samples, each of one or many libraries that were sequenced on at least one lane of an NGS machine.

Note that many GATK commands can be run at the lane level, but will give better results seeing all of the data for a single sample, or even all of

VCF file

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | Gen011 |
|--------|---------|-----------|-----|-----|---------|--------|--|----------------|--------------------------|
| chr1 | 3311261 | rs6702992 | C | T | 663.71 | PASS | AC=2;AF=1.00;AN=2;DB;DP=3;FS=0.000;MQ=60.00;MQ0=0;POSITIVE_TRAIN_SITE;QD=31.56;VQSLOD=4.22;culprit=FS;SF=0 | GT:AD:DP:GQ:PL | 1/1:0,3:3:9:10 9,9,0 |
| chr1 | 5937168 | rs3747992 | G | A | 1245.09 | PASS | AC=1;AF=0.500;AN=2;BaseQRankSum=1.23;ClippingRankSum=-7.620e-01;DB;DP=5;FS=0.000;MQ=60.00;MQ0=0;MQRankSum=1.18;POSITIVE_TRAIN_SITE;QD=19.16;ReadPosRankSum=0.067;VQSLOD=4.25;culprit=FS;SF=0 | GT:AD:DP:GQ:PL | 0/1:3,2:5:51:5 8,0,51 |

SNPEff (effect on function of variants)

SnpEff tools

1

SnpEff Variant effect and annotation

2

- [SnpEff Download](#) Download a new database
- [SnpSift Annotate](#) Annotate SNPs from dbSnp
- [SnpSift CaseControl](#) Count samples are in 'case' and 'control' groups.
- [SnpSift Filter](#) Filter variants using arbitrary expressions
- [SnpSift Intervals](#) Filter variants using intervals

SnpEff (version 1.0)

Sequence changes (SNPs, MNPs, InDels):

8: Generate VCF with mpileup piped through bcftools view on data 7: mp

3a

Input format:

VCF *

Output format:

Tabular *

Genome:

Homo_sapiens (hg19)

3b

Upstream / Downstream length:

5000 bases

Filter homozygous / heterozygous changes: *

- No filter (analyze everything)
- Analyze homozygous sequence changes only
- Analyze heterozygous sequence changes only

Filter sequence changes: *

- No filter (analyze everything)
- Analyze deletions only
- Analyze insertions only
- Only MNPs (multiple nucleotide polymorphisms)
- Only SNPs (single nucleotide polymorphisms)

Filter output: *

Select All Unselect All

- None
- Do not show DOWNSTREAM changes
- Do not show INTERGENIC changes
- Do not show INTRON changes
- Do not show UPSTREAM changes
- Do not show 5_PRIME_UTR or 3_PRIME_UTR changes

Chromosomal position: *

- Use default (based on input type)
- Force zero-based positions (both input and output)
- Force one-based positions (both input and output)

Execute

4

1. Click on "SnpEff tools"
2. Click on "SnpEff"
3. Select options:
 - a) Choose VCF output file from previous slide
 - b) Pick genome "Homo_sapiens (hg19)"
4. Click "Execute"

* Other options to be aware of!

SNPEff output (cropped)

| Chr | Pos | REF | ALT | Change_type | Homozygous | Quality | Cov. | Gene name | Transcript_ID | Effect | old_AA/new_AA | Old_codon/New_codon | Codon_Num(CDS) | Codon_Degeneracy | CDS_size |
|-----|----------|-----|-----|-------------|------------|---------|------|-----------|---------------|-----------------------|---------------|---------------------|----------------|------------------|----------|
| 21 | 10910311 | T | G | SNP | Hom | 99 | 737 | TPTE | NM_199259 | NON_SYNONYMOUS_CODING | Y/S | tAt/tCt | 464 | 1 | 1602 |
| 21 | 10970008 | C | T | SNP | Hom | 4.13 | 286 | TPTE | NM_199259 | SPLICE_SITE_DONOR | | | | | 1602 |
| 21 | 11058226 | G | C | SNP | Hom | 21 | 882 | BAGE2 | NM_182482 | NON_SYNONYMOUS_CODING | P/A | Cct/Gct | 72 | 1 | 330 |
| 21 | 15481365 | G | T | SNP | Hom | 99 | 171 | LIPI | NM_198996 | NON_SYNONYMOUS_CODING | D/E | gaC/gaA | 465 | 2 | 1446 |
| 21 | 15596772 | T | G | SNP | Hom | 99 | 63 | RBM11 | NM_144770 | NON_SYNONYMOUS_CODING | L/V | Ttg/Gtg | 116 | 2 | 846 |
| 21 | 15954528 | G | A | SNP | Hom | 99 | 405 | SAMSN1 | NM_001256370 | NON_SYNONYMOUS_CODING | H/Y | Cac/Tac | 64 | 1 | 1326 |
| 21 | 30339120 | C | A | SNP | Hom | 99 | 111 | LTN1 | NM_015565 | NON_SYNONYMOUS_CODING | G/C | Ggc/Tgc | 611 | 1 | 5439 |
| 21 | 31744127 | A | T | SNP | Hom | 99 | 110 | KRTAP13-2 | NM_181621 | STOP_GAINED | C/* | tgT/tgA | 135 | 2 | 528 |
| 21 | 34948686 | * | +A | INS | Hom | 99 | 74 | SON | NM_138927 | FRAME_SHIFT | -/? | -/A | 2413 | | 7281 |

| Tool | Link |
|--|---|
| CADD | http://cadd.gs.washington.edu/ |
| Broad ExAC | http://exac.broadinstitute.org/ |
| SeattleSeq | http://snp.gs.washington.edu/SeattleSeqAnnotation138/ |
| Ensembl VEP | http://useast.ensembl.org/info/docs/tools/vep/index.html |
| MutationAssessor | http://mutationassessor.org/ |
| MutationTaster | http://www.mutationtaster.org/ |
| OMIM | http://www.ncbi.nlm.nih.gov/omim |
| ClinVar | http://www.ncbi.nlm.nih.gov/clinvar/ |
| GeneCards | http://www.genecards.org/ |
| Wellcome Trust Sanger Institute Mouse Genomes Project | http://www.sanger.ac.uk/resources/mouse/genomes/ |



Demo

CADD

CADD

- **Combined Annotation Dependent Depletion (CADD)**
- A tool that scores deleteriousness of SNVs as well as INDELS in the human genome
- Integrates multiple annotations into one metric by contrasting variants that survived natural selection with simulated mutations.
- How to interpret scaled CADD scores:
 - 10th-% of CADD scores are assigned to CADD-10
 - Top 1% to CADD-20
 - Top 0.1% to CADD-30

Prepare SNPEff output files for downstream analysis (i.e. CADD)

1. Download SNPEff output file
2. Open in Excel
3. Keep first four (4) columns; remove all header columns
4. Create a 3rd column called ID and place a “.” for each variant. Final output in this format:

| # | CHROM | POS | ID | REF | ALT |
|----|-------|---------|----|-----|-----|
| 21 | | 1097008 | . | C | T |

1. Save as a Windows or Unix text file
2. Manually change .txt to .vcf

How to use CADD (on a small set of SNVs)

- Go to: <http://cadd.gs.washington.edu/score>
- Upload a gzip-compressed vcf file of variants
- Click on “Upload variants”
- Once analysis completes, download the tab-delimited file and open in Excel.

Please upload a VCF file containing up to 100,000 variants

Please provide a (preferentially gzip-compressed) VCF file of your variants. For information on the VCF format see <http://vcftools.sourceforge.net/specs.html>. It is sufficient to provide the first 5 columns of a VCF file without header, as all other information than CHROM, POS, REF, ALT will be ignored anyway. The maximum accepted file size is set at 2MB (>100,000 variants for 5 column compressed VCF). If you try to upload files larger than 2MB, you will receive an error ("Connection reset"). You will be able to retrieve your variants faster, if you upload them in smaller sets. The file that will be provided for download is a gzip-compressed tab-separated text file. Make sure that your browser does not alter the file extension (.tsv.gz) during download; otherwise your operating system will not be able to automatically pick the right programs for opening the output. If you need more variants, we suggest [downloading](#) the full set of variants. To learn about differences between versions, please check the [release notes](#).

Choose File No file chosen

v1.1 

Include underlying annotation in output (not only the scores)

Upload variants

University of Alabama at Birmingham

Heflin Center for Genomic Science

Demo

BROAD EXAC

Broad's ExAC Browser

- **Exome Aggregation Consortium (ExAC)**
- Consists of 60,706 unrelated individuals that were exome sequenced
- Various disease-specific and population genetic studies
- Search for any human gene or variant or region

How to use ExAC

- Go to: <http://exac.broadinstitute.org/>
- Input either gene or variant or region

ExAC Browser (Beta) | Exome Aggregation Consortium

Search for a gene or variant or region

Examples - Gene: [PCSK9](#), Transcript: [ENST00000407236](#), Variant: [22-46615880-T-C](#), Multi-allelic variant: [rs1800234](#), Region: [22:46615715-46615880](#)

Gene: PRDM1

PRDM1 PR domain containing 1, with ZNF domain
Number of variants 509 (Including filtered: 555)
UCSC Browser [6:106534195-106557814](#)
GeneCards [PRDM1](#)
OMIM [PRDM1](#)
Other

Transcripts ▾

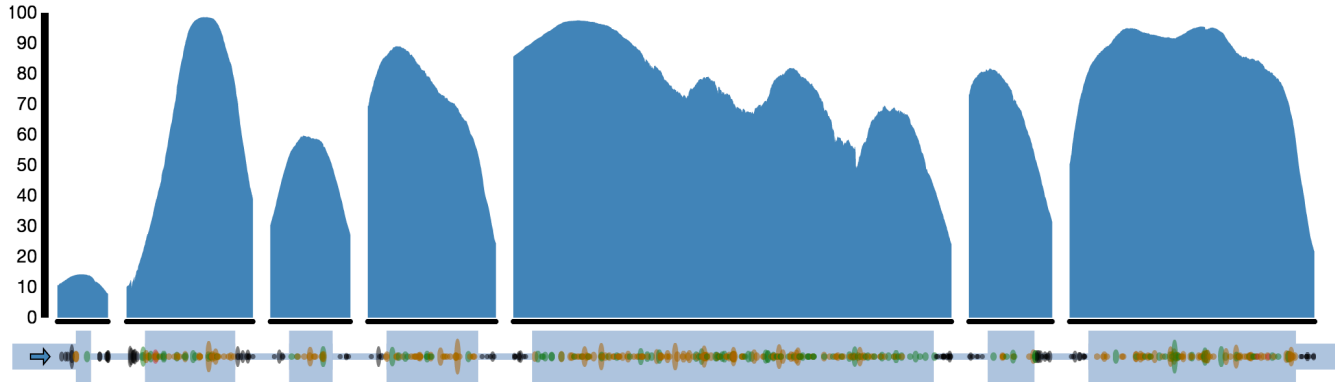
Gene summary

(Coverage shown for **canonical transcript**: ENST00000369096)

Mean coverage 70.23

Display: **Overview** **Detail** Include UTRs in plot

Coverage metric: **Average** **Individuals over X**
 Metric: mean ▾



All **Missense + LoF** **LoF** Include filtered (non-PASS) variants

Invert (highlight rare variants)

| Variant | Chrom | Position | Protein Consequence | Filter | Annotation | Allele Count | Allele Number | Number of Homozygotes | Allele Frequency |
|------------------------------------|-------|-----------|---------------------|--------|-----------------|--------------|---------------|-----------------------|------------------|
| 6:106534391 A / G | 6 | 106534391 | | PASS | 5' UTR | 1 | 23570 | 0 | 0.00004243 |
| 6:106534402 G / C | 6 | 106534402 | | PASS | 5' UTR | 1 | 23606 | 0 | 0.00004236 |
| 6:106534408 G / A | 6 | 106534408 | | PASS | 5' UTR | 1 | 23640 | 0 | 0.00004230 |
| 6:106534419 C / G | 6 | 106534419 | | PASS | 5' UTR | 162 | 23696 | 1 | 0.006837 |
| 6:106534419 CT / C | 6 | 106534419 | | PASS | 5' UTR | 3 | 23696 | 0 | 0.0001266 |
| 6:106534430 T / C | 6 | 106534430 | p.Met1? | PASS | initiator codon | 1 | 23752 | 0 | 0.00004210 |
| 6:106534431 G / T | 6 | 106534431 | p.Met1? | PASS | initiator codon | 1 | 23756 | 0 | 0.00004209 |

OMIM

*609252

LIPASE I; **LIPI**

Alternative titles; symbols

LPD LIPASE; LPDL

PRED5

HGNC Approved Gene Symbol: **LIPI**

Cytogenetic location: [21q11.2](#) *Genomic coordinates (GRCh37):* [21:15,480,783-15,583,360](#) (from NCBI)

Gene-Phenotype Relationships

| Location | Phenotype | Phenotype MIM number | Phenotype mapping key |
|-------------------------|---|------------------------|-----------------------|
| 21q11.2 | {Hypertriglyceridemia, susceptibility to} | 145750 | 3 |

TEXT

Cloning and Expression

By screening a testis cDNA library using mouse Lpdl as probe, [Wen et al. \(2003\)](#) cloned human **LIPI**, which they called LPDL. The deduced 460-amino acid protein contains a hydrophobic leader sequence with a putative cleavage site after amino acid 15, a central lipase consensus sequence (GxSxG), a central 12-amino acid lipase lid sequence,

ClinVar

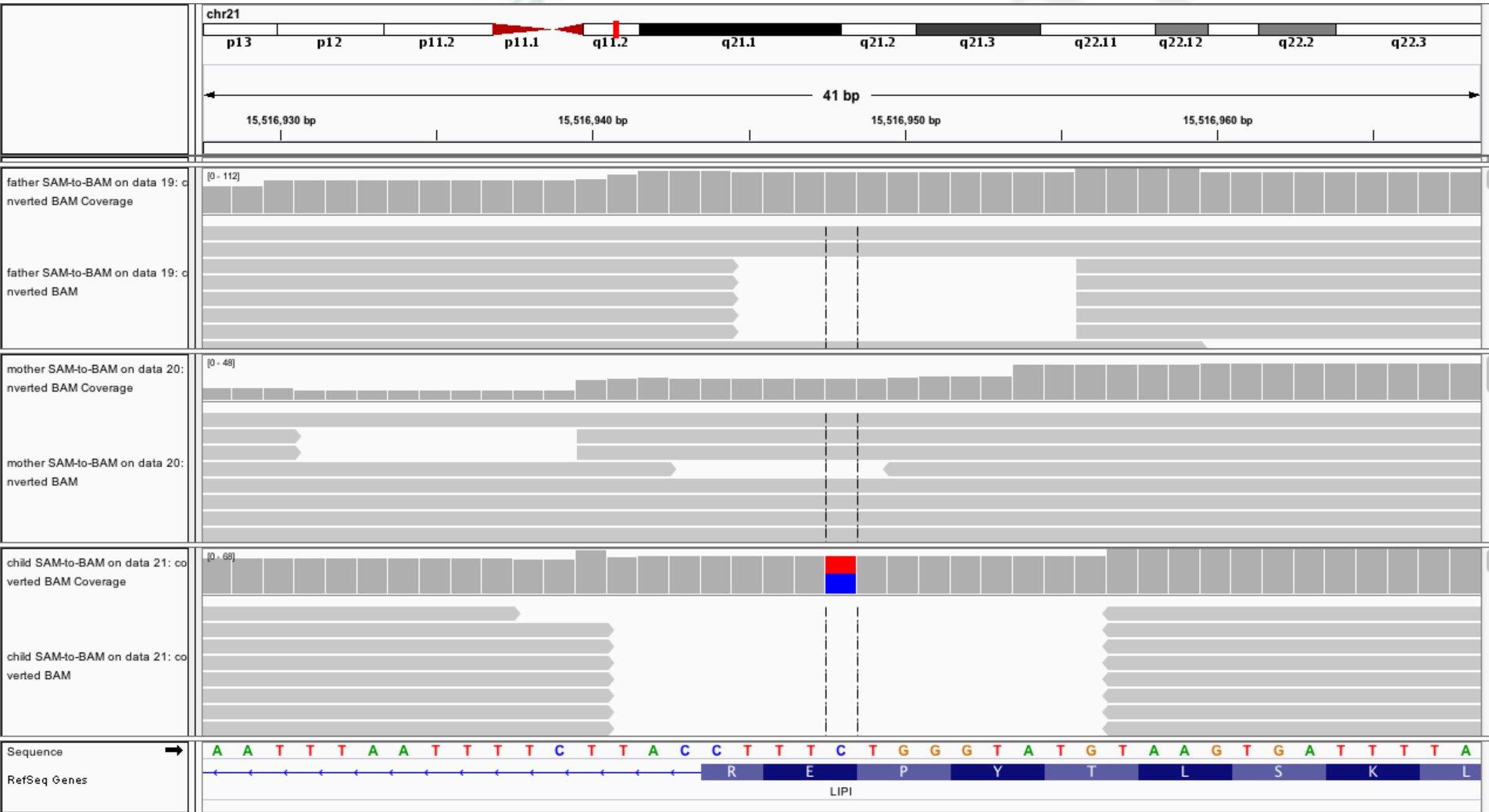
Results: 1 to 20 of 49

<< First < Prev Page 1 of 3 Next > Last >>

| | Gene(s) | Condition(s) | Frequency | Clinical significance (Last reviewed) | Review status | C |
|--|----------------------|----------------------|------------------|---------------------------------------|--------------------------------|---|
| <input type="checkbox"/> 1. NM_198996.2(LIPI):c.1358+1266T>C | LIPI | Lung cancer | | Uncertain significance | classified by single submitter | 2 |
| <input type="checkbox"/> 2. NM_198996.3(LIPI):c.322G>A (p.Gly108Ser) | LIPI | Malignant melanoma | | not provided | not classified by submitter | 2 |
| <input type="checkbox"/> 3. NM_198996.3(LIPI):c.691C>T (p.His231Tyr) | LIPI | Malignant melanoma | | not provided | not classified by submitter | 2 |
| <input type="checkbox"/> 4. NM_198996.3(LIPI):c.749C>T (p.Pro250Leu) | LIPI | Malignant melanoma | | not provided | not classified by submitter | 2 |
| <input type="checkbox"/> 5. NM_198996.3(LIPI):c.1158G>A (p.Met386Ile) | LIPI | Malignant melanoma | | not provided | not classified by submitter | 2 |
| <input type="checkbox"/> NM_198996.3(LIPI):c.227G>A | LIPI | Hypertriglyceridemia | GO-FSP:000969(T) | risk factor | classified | 2 |

IGV

at Birmingham



References and web links

- Galaxy
 - PSU Public website: <https://usegalaxy.org/>
 - UAB: <https://www.uab.edu/galaxy>
- [Bowtie](#)
- [GATK](#)
- [SAMTools](#)
- [Picard Tools](#)
- [FreeBayes](#)
- [IGV](#)
- [SNPEff](#)
- [CADD](#)
- [Broad Exac](#)
- [dbSNP](#)
- [OMIM](#)
- [ClinVar](#)

University of Alabama at Birmingham

Heflin Center for Genomic Science

Thanks! Questions?

Contact info:

David K. Crossman, Ph.D.

Bioinformatics Director

Heflin Center for Genomic Science

University of Alabama at Birmingham

<http://www.heflingenetics.uab.edu>

dkcrossm@uab.edu