# Exome sequencing data analysis for diagnosing a genetic disease

## Galaxy Training! tutorial

# Galaxy trainings!

- Calling variants on diploid organism :

  https://training.galaxyproject.org/training-material/topics/variant-analysis/tutorials/dip/tutorial.html

- Calling variants on non diploid system :

  https://training.galaxyproject.org/training-material/topics/variant-analysis/tutorials/non-dip/tutorial.html

- Microbial variants calling :

  https://training.galaxyproject.org/training-material/topics/variant-analysis/tutorials/microbial-variants/tutorial.html

- Genome annotations (eukaryotes, prokaryotes, other):

  https://training.galaxyproject.org/training-material/topics/genome-annotation/

# Tutorial presentation

- Exome sequencing data from a family trio

- Boy child affected by a disease : osteopetrosis

- Parents unaffected but consanguineous

**Goal : Identify the genetic variation responsible for the disease**

# Tutorial steps

1. Perform postprocessing from premapped reads

2. Variant calling

3. Variant annotation and reporting

# Tutorial steps

1. Perform postprocessing from premapped reads

2. Variant calling

3. Variant annotation and reporting

# Premapped reads

- Data characteristics for the trio :
  - Whole exome sequencing
  - Paired-end reads

- Steps already performed :
  - Quality control (fastq)
  - Read mapping (Human Hg19 assembly)

- Format available : bam format

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

| | | | |
|---|---|---|---|
| 📁 | ☐ | PAPAA PI3K_OG:Pancancer Aberrant Pathway Activity Analysis | Summary |
| 📁 | ☐ | Proteomics | Training material for proteomics workflo ... (more) |
| 📁 | ☐ | Refining Manual Genome Annotations with Apollo | We look at how to edit Genome Annotation ... (more) |
| 📁 | ☐ | RNA interactome | RNA interactome data analysis |
| 📁 | ☐ | Sequence analysis | Analyses of sequences |
| 📁 | ☐ | Statistics and machine learning | Statistical Analyses for omics data and ... (more) |
| 📁 | ☐ | The new topic | Summary |
| 📁 | ☐ | Transcriptomics | Training material for all kinds of trans ... (more) |
| 📁 | ☐ | User Interface and Features | A collection of microtutorials explainin ... (more) |
| 📁 | ☐ | Variant Analysis | Exome sequencing means that all protein- ... (more) |

# Premapped reads upload

Libraries / GTN - Material / Variant Analysis

| | Name | | Description |
|---|---|---|---|
| ☐ | Name | | Description |
| 📁 ☐ | Calling variants in diploid systems | | |
| 📁 ☐ | Calling variants in non-diploid systems | | |
| 📁 ☐ | DOI: 10.5281/zenodo.3960260 | | latest |
| 📁 ☐ | Exome sequencing data analysis for diagnosing a genetic disease | | |
| 📁 ☐ | Identification of somatic and germline variants from tumor and normal sample pairs | | |
| 📁 ☐ | Mapping and molecular identification of phenotype-causing mutations | | |
| 📁 ☐ | Microbial Variant Calling | | |
| 📁 ☐ | Mutation calling, viral genome reconstruction and lineage/clade assignment from SARS-CoV-2 sequencing data | | |

« ‹ 1 › »    10 ⌄ per page, 8 total

# Premapped reads upload

Libraries / GTN - Material / Variant Analysis / Exome sequencing data analysis for diagnosing a genetic disease

| | | Name | Description |
|---|---|---|---|
| | ☐ | | |
| 🗀 | ☐ | DOI: 10.5281/zenodo.3054169 | latest |

« ‹ 1 › »   10 ⬍ per page, 1 total

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

**Edit Dataset Attributes**

≡ Attributes    ⚙ Convert    ≣ Datatypes    👤 Permissions

**Name**

https://zenodo.org/api/files/dd4bcd95-4412-4ac0-a7d2-23cf1c69e0bc/mapped_reads_father.bam

**Info**

uploaded bam file

**Annotation**

Add an annotation or notes to a dataset; annotations are available when a history is viewed.

**Database/Build**

unspecified (?)

💾 Save    ↻ Auto-detect

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

# Premapped reads upload

# Mapped reads postprocessing

**Warning :**

- Depends on technology
- Depends on goal
- Depends on the pipeline used (steps, software, etc.)

1. Filter reads based on characteristics :
   - Retain only forward and reverse reads mapped successfully to the reference
   - Exclude possible contaminant DNA or sequencing artefact

2. Remove/Mark duplicate reads
   - PCR-overamplification of genomic fragment during sequencing library preparation

# Mapped reads postprocessing - Filter reads

# Mapped reads processing - Filter reads

# Mapped reads postprocessing - Filter reads

**Only output alignments with all of these flag bits set**

☐ Select/Unselect all

☐ Read is paired
☐ Read is mapped in a proper pair
☐ The read is unmapped
☐ The mate is unmapped
☐ Read is mapped to the reverse strand of the reference
☐ Mate is mapped to the reverse strand of the reference
☐ Read is the first in a pair
☐ Read is the second in a pair
☐ The alignment of this read is not primary
☐ The read fails platform/vendor quality checks
☐ The read is a PCR or optical duplicate
☐ Supplementary alignment

(-f)

**Skip alignments with any of these flag bits set**

☐ Select/Unselect all

☐ Read is paired
☐ Read is mapped in a proper pair
**1** ☑ The read is unmapped
☑ The mate is unmapped
☐ Read is mapped to the reverse strand of the reference
☐ Mate is mapped to the reverse strand of the reference
☐ Read is the first in a pair
☐ Read is the second in a pair
☐ The alignment of this read is not primary
☐ The read fails platform/vendor quality checks
☐ The read is a PCR or optical duplicate
☐ Supplementary alignment

**Select alignments from Library**

(-l) Requires headers in the input SAM or BAM, otherwise no alignments will be output

**Select alignments from Read Group**

(-r) Requires headers in the input SAM or BAM, otherwise no alignments will be output

**Output alignments overlapping the regions in the BED file**

No bed dataset available.

(-L)

**Use inverse selection**

◯ No

Select the opposite of the listed chromosomes

**Select regions (only used when the input is in BAM format)**

region should be presented in one of the following formats: `chr1`, `chr2:1,000` and `chr3:1000-2,000`

**Select the output format**

BAM (-b)

**Email notification**

◯

Send an email notification when the job completes.

**2** ✔ Execute

# Mapped reads postprocessing - Filter reads

# Mapped reads postprocessing - Duplicate reads

# Mapped reads postprocessing - Duplicate reads

# Mapped reads postprocessing - Duplicate reads



| QNAME | FLAG | RNAME | POS | MAPQ | CIGAR | MRNM | MPOS | ISIZE | SEQ |
|---|---|---|---|---|---|---|---|---|---|
| @HD VN:1.3 SO:coordinate | | | | | | | | | |
| @SQ SN:chr8 LN:146364022 | | | | | | | | | |
| @RG ID:001 SM:father PL:ILLUMINA | | | | | | | | | |
| @PG ID:bwa PN:bwa VN:0.7.17-r1188 CL:bwa mem -t 8 -v 1 -R @RG\tID:001\tSM:father\tPL:ILLUMINA localref.fa /data/dnb02/galaxy_db/files/009/499/dataset_9499701.dat /data/dnb02/galaxy_db/files/009/499/datas |
| DCW97JN1:309:C0C42ACXX:5:2202:19629:56029 | 163 | chr8 | 11710 | 3 | 101M | = | 11865 | 256 | CCATGGCAGAGCTCCCTCCTCAGCACATGGGGAGCAGACAGGAAGT |
| DCW97JN1:309:C0C42ACXX:4:1206:10027:62829 | 163 | chr8 | 11712 | 0 | 101M | = | 11864 | 253 | ATGGCAGAGCTCCCTCCTCAGCACATGGGGAGCAGACAGGAAGTTT |
| DCW97JN1:309:C0C42ACXX:4:1115:17796:60101 | 163 | chr8 | 11712 | 15 | 101M | = | 11869 | 253 | ATGGCAGAGCTCCCTCCTCAGCACATGGGGAGCAGACAGGAAGTTT |
| DCW97JN1:309:C0C42ACXX:5:1216:6300:20909 | 99 | chr8 | 11783 | 27 | 101M | = | 11966 | 271 | AGCCACGTCTCCCCAGGTCAGTCTTAAGGACAACGAAACTCTGGGC |
| DCW97JN1:309:C0C42ACXX:4:1206:10027:62829 | 83 | chr8 | 11864 | 1 | 101M | = | 11712 | -253 | AAGCCATGGTGCCCCACCCTCGGGTGGGTCCTGAGGAGAACAAAGC |
| DCW97JN1:309:C0C42ACXX:5:2202:19629:56029 | 83 | chr8 | 11865 | 8 | 101M | = | 11710 | -256 | AGCCATGGTGACCCACCCTCGGGTGGGTCCTGAGGAGAACAAAGCT |
| DCW97JN1:309:C0C42ACXX:4:1115:17796:60101 | 83 | chr8 | 11869 | 15 | 96M5S | = | 11712 | -253 | ATGGTGACCCACCCTCGGGTGGGTCCTGAGGAGAACAAAGCTCTGG |
| DCW97JN1:309:C0C42ACXX:5:1216:6300:20909 | 147 | chr8 | 11966 | 27 | 13S88M | = | 11783 | -271 | CCAGATCCCAAACCCTGATCCCTACCCTGGATCCTAAGTCTGTCCCT |
| DCW97JN1:309:C0C42ACXX:5:2210:15831:85655 | 145 | chr8 | 98822 | 0 | 52S35M14S | = | 110566976 | 110468121 | TTTTAAATTTAAAAAAAAAAAATTGGCCAAAAAAATTTTATTTTTT |
| DCW97JN1:309:C0C42ACXX:4:2209:3455:67435 | 161 | chr8 | 98823 | 0 | 45S43M13S | = | 39494954 | 39396232 | CCCCAAAAAAATTTCGGGGTTTTGGGTTTTTTCCACCCAAAATTTT |
| DCW97JN1:309:C0C42ACXX:5:2305:4557:78030 | 2115 | chr8 | 98823 | 0 | 58H34M9H | = | 141889681 | 141790859 | TTTTTTTTTTTTTTTTTTTTTTTTTTTAAATT |
| DCW97JN1:309:C0C42ACXX:5:2111:10544:43299 | 2195 | chr8 | 98824 | 0 | 43M58H | = | 16979740 | 16880875 | TTTTTTTTTTTTTTTTTTTTTTTTAAAATTTTTTTTTTA |

- <instrument>:<run_number>:<flowcell_ID>:<lane>:<tile>:<x-pos>:<y-pos>

SO tag :

- Sorting order of alignments
- Unknown, unsorted, queryname (QNAME) or coordinate (RNAME/POS)

https://support.basespace.illumina.com/articles/descriptive/fastq-files/

# Mapped reads postprocessing - Duplicate reads

# Mapped reads postprocessing - Duplicate reads

**The maximum offset between two duplicte clusters in order to consider them optical duplicates**

| 100 |
|-----|

OPTICAL_DUPLICATE_PIXEL_DISTANCE; default=100

**Barcode Tag**

Barcode SAM tag. This tag can be utilized when you have data from an assay that includes Unique Molecular Indices. Typically 'RX'

**Select validation stringency**

| Lenient ▼ |
|-----------|

Setting stringency to SILENT can improve performance when processing a BAM file in which variable-length data (read, qualities, tags) do not otherwise need to be decoded.

**Email notification**

Send an email notification when the job completes.

✔ Execute   **6**

# Mapped reads postprocessing - Duplicate reads

# Mapped reads postprocessing - Duplicate reads

# Mapped reads postprocessing - Duplicate reads



http://broadinstitute.github.io/picard/explain-flags

# Mapped reads postprocessing - Duplicate reads

# Tutorial steps

1.  Perform postprocessing from premapped reads

2.  Variant calling

3.  Variant annotation and reporting

# Variant calling

# Variant calling

# Variant calling

**Limit variant calling to a set of regions?**

| Do not limit | ▼ |
|---|---|

Sets --targets or --region options

**Read coverage**

| Use defaults | ▼ |
|---|---|

Sets --min-coverage, --limit-coverage, and --skip-coverage

**Choose parameter selection level**                    **1**

| 2. Simple diploid calling with filtering and coverage | ▼ |
|---|---|

Select how much control over the freebayes run you need

**Email notification**

◯

Send an email notification when the job completes.

✔ Execute

**2**

**Galaxy-specific options**

Galaxy allows five levels of control over FreeBayes options, provided by the **Choose parameter selection level** menu option. These are:

1. *Simple diploid calling*: The simplest possible FreeBayes application. Equivalent to using FreeBayes with only a BAM input and no other parameter options.
2. *Simple diploid calling with filtering and coverage*: Same as #1 plus two additional options: -0 (standard filters: --min-mapping-quality 30 --min-base-quality 20 --min-supporting-allele-qsum 0 --genotype-variant-threshold 0) and --min-coverage.
3. *Frequency-based pooled calling*: This is equivalent to using FreeBayes with the following options: --haplotype-length 0 --min-alternate-count 1 --min-alternate-fraction 0 --pooled-continuous --report-monomorphic. This is the best choice for calling variants in mixtures such as viral, bacterial, or organellar genomes.
4. *Frequency-based pooled calling with filtering and coverage*: Same as #3 but adds -0 and --min-coverage like in #2.
5. *Complete list of all options*: Gives you full control by exposing all FreeBayes options as Galaxy parameters.

# Variant calling

# Variant calling - VCF

# Variant calling - VCF

```
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Float,Description="Genotype      Quality,      the      Phred-scaled
##FORMAT=<ID=GL,Number=G,Type=Float,Description="Genotype Likelihood, log10-scaled likelihoods of the
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=AD,Number=R,Type=Integer,Description="Number of observation for each allele">
##FORMAT=<ID=RO,Number=1,Type=Integer,Description="Reference allele observation count">
##FORMAT=<ID=QR,Number=1,Type=Integer,Description="Sum of quality of the reference observations">
##FORMAT=<ID=AO,Number=A,Type=Integer,Description="Alternate allele observation count">
##FORMAT=<ID=QA,Number=A,Type=Integer,Description="Sum of quality of the alternate observations">
##FORMAT=<ID=MIN_DP,Number=1,Type=Integer,Description="Minimum depth in gVCF output block.">
```

# Variant calling - VCF

**Mandatory columns**

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER |
|--------|-----|----|-----|-----|------|--------|
| chr8 | 115956 | . | A | T | 9.09784e-07 | . |
| chr8 | 116079 | . | G | A | 103.501 | . |
| chr8 | 116701 | . | A | G | 3.98084e-05 | . |
| chr8 | 116895 | . | A | G | 184.59 | . |
| chr8 | 160552 | . | G | A | 1.00485 | . |
| chr8 | 160608 | . | A | C | 722.504 | . |
| chr8 | 160609 | . | AAAAAATAAAAATAAACATAAAAATG | AAAATAAAAATAAAAATAAACATAAAAATG | 0.370623 | . |
| chr8 | 160679 | . | G | A | 5.46006e-08 | . |
| chr8 | 160719 | . | C | T | 9.28165e-15 | . |
| chr8 | 160736 | . | G | T | 530.182 | . |
| chr8 | 160760 | . | C | G | 237.975 | . |

# Variant calling - VCF

**Mandatory column**

| INFO |
|------|
| AB=0;ABP=0;AC=0;AF=0;AN=6;AO=4;CIGAR=1X;DP=51;DPB=51;DPRA=2.33333;EPP=11.6962;EPPR=36.6912;GTI=0;LEN=1;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1;ODDS=15.5049;PAIRED=1;PAI |
| AB=0.276596;ABP=23.3852;AC=2;AF=0.333333;AN=6;AO=15;CIGAR=1X;DP=74;DPB=74;DPRA=0;EPP=20.5268;EPPR=29.8409;GTI=0;LEN=1;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1;ODDS=4.51 |
| AB=0.3125;ABP=7.89611;AC=1;AF=0.166667;AN=6;AO=14;CIGAR=1X;DP=240;DPB=240;DPRA=0;EPP=3.0103;EPPR=6.85361;GTI=0;LEN=1;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1;ODDS=11.6;F |
| AB=0;ABP=0;AC=6;AF=1;AN=6;AO=6;CIGAR=1X;DP=6;DPB=6;DPRA=0;EPP=8.80089;EPPR=0;GTI=0;LEN=1;MEANALT=1;MQM=60;MQMR=0;NS=3;NUMALT=1;ODDS=8.00168;PAIRED=1;PAIREDR=0;PAO=0;P |
| AB=0.25;ABP=9.52472;AC=2;AF=0.333333;AN=6;AO=3;CIGAR=1X;DP=19;DPB=19;DPRA=0.857143;EPP=3.73412;EPPR=3.55317;GTI=0;LEN=1;MEANALT=1;MQM=60;MQMR=59.5;NS=3;NUMALT=1;ODDS=1 |
| AB=0.4375;ABP=5.72464;AC=3;AF=0.5;AN=6;AO=35;CIGAR=1X;DP=80;DPB=80;DPRA=0;EPP=48.239;EPPR=49.3833;GTI=0;LEN=1;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1;ODDS=7.50894;PAIRE |
| AB=0.222222;ABP=15.074;AC=1;AF=0.166667;AN=6;AO=7;CIGAR=1M4I25M;DP=80;DPB=82.5385;DPRA=0;EPP=5.80219;EPPR=113.696;GTI=0;LEN=4;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1; |
| AB=0.130584;ABP=347.946;AC=3;AF=0.5;AN=6;AO=38;CIGAR=1X;DP=291;DPB=291;DPRA=0;EPP=54.4399;EPPR=6.10873;GTI=2;LEN=1;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1;ODDS=18.194;F |
| AB=0.101399;ABP=397.702;AC=2;AF=0.333333;AN=6;AO=29;CIGAR=1X;DP=441;DPB=441;DPRA=0.922581;EPP=3.68421;EPPR=65.7822;GTI=1;LEN=1;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1;C |
| AB=0.188995;ABP=354.186;AC=3;AF=0.5;AN=6;AO=79;CIGAR=1X;DP=418;DPB=418;DPRA=0;EPP=20.1897;EPPR=12.7531;GTI=1;LEN=1;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1;ODDS=34.1344; |
| AB=0.124567;ABP=356.825;AC=2;AF=0.333333;AN=6;AO=37;CIGAR=1X;DP=382;DPB=382;DPRA=0;EPP=7.76406;EPPR=133.565;GTI=1;LEN=1;MEANALT=1.66667;MQM=60;MQMR=60;NS=3;NUMALT=1;OD |
| AB=0.124031;ABP=161.393;AC=1;AF=0.166667;AN=6;AO=21;CIGAR=2X;DP=310;DPB=310;DPRA=0.962264;EPP=3.94093;EPPR=397.039;GTI=0;LEN=2;MEANALT=1;MQM=60;MQMR=60;NS=3;NUMALT=1;C |

# Variant calling - VCF

| FORMAT |
|---|
| GT:DP:AD:RO:QR:AO:QA:GL |
| GT:DP:AD:RO:QR:AO:QA:GL |
| GT:DP:AD:RO:QR:AO:QA:GL |
| GT:DP:AD:RO:QR:AO:QA:GL |
| GT:DP:AD:RO:QR:AO:QA:GL |
| GT:DP:AD:RO:QR:AO:QA:GL |
| GT:DP:AD:RO:QR:AO:QA:GL |
| GT:DP:AD:RO:QR:AO:QA:GL |
| GT:DP:AD:RO:QR:AO:QA:GL |

| proband |
|---|
| 0/0:30:27,3:27:891:3:92:0,-0.445657,-71.9117 |
| 0/1:24:16,8:16:644:8:260:-16.4945,0,-51.046 |
| 0/0:113:109,4:109:3436:4:144:0,-20.7059,-296.176 |
| 1/1:4:0,4:0:0:4:167:-15.4235,-1.20412,0 |
| 0/1:9:7,2:7:297:2:66:-3.55868,0,-24.3555 |
| 0/1:43:22,21:22:776:21:828:-61.8817,0,-57.2184 |
| 0/0:42:41,1:41:1422:1:34:0,-9.58258,-124.881 |
| 0/1:133:118,15:118:3976:15:509:-6.09578,0,-318.014 |
| 0/1:185:166,19:166:6862:19:635:-1.7759,0,-561.244 |

| mother |
|---|
| 0/0:12:11,1:11:353:1:33:0,-0.313225,-28.7828 |
| 0/0:27:25,2:25:1021:2:64:0,-2.04915,-86.078 |
| 0/0:111:106,5:106:3382:5:193:0,-15.6745,-286.887 |
| 1/1:1:0,1:0:0:1:36:-3.59827,-0.30103,0 |
| 0/1:3:2,1:2:85:1:35:-2.59554,0,-7.15727 |
| 0/1:17:14,3:14:502:3:114:-5.51484,0,-40.3989 |
| 0/1:18:14,4:14:477:4:132:-6.46629,0,-37.5149 |
| 0/1:59:49,10:49:1629:10:328:-12.0781,0,-129.133 |
| 0/1:101:91,10:91:3600:10:342:-0.707324,0,-293.6 |

| father |
|---|
| 0/0:9:9,0:9:286:0:0:0,-2.70927,-26.0508 |
| 0/1:23:18,5:18:728:5:166:-8.34408,0,-58.9123 |
| 0/1:16:11,5:11:364:5:178:-11.5434,0,-28.2653 |
| 1/1:1:0,1:0:0:1:33:-3.29913,-0.30103,0 |
| 0/0:7:7,0:7:271:0:0:0,-2.10721,-24.7468 |
| 0/1:20:9,11:9:307:11:421:-32.2186,0,-21.9403 |
| 0/0:20:18,2:18:614:2:64:0,-0.00155201,-49.4499 |
| 0/1:99:86,13:86:2819:13:441:-10.2124,0,-224.147 |
| 0/0:155:154,0:154:6061:0:0:0,-46.3586,-544.867 |

**Genotypes format**

**Proband genotypes information**

**Mother genotypes information**

**Father genotypes information**

# Variant calling



TP_GTN_WES_disease

14 shown

2.03 GB

14: freebayes_calling.vcf
father  mother  proband

13: markdup_proband.bam
proband

12: markdup_proband_metrics
proband

# Tutorial steps

1. Perform postprocessing from premapped reads

2. Variant calling

3. Variant annotation and reporting

# Variant normalization

# Variant normalization

# Variant normalization

# Variant normalization

# Variant normalization - Alleles

# Variant normalization - Genotypes

**Initial file**

163550 . AAGT GAGC,GAGT

1/2:169:0,61,108:0:0:61,108:2328,4362:-550.761,-359.801,-341.438,-191.22,0,-158.709

1/2:112:0,39,72:0:0:39,72:1461,2734:-343.835,-224.186,-212.446,-119.697,0,-98.023

1/1:112:0,112,0:0:0:112,0:4100,0:-368.767,-33.7154,0,-368.767,-33.7154,-368.767

**Normalized file**

163550 . AAGT GAGC

163550 . A G

1/0:169:0,61:0:0:61:2328:-550.761,-359.801,-341.438

0/1:169:0,108:0:0:108:4362:-550.761,-191.22,-158.709

1/0:112:0,39:0:0:39:1461:-343.835,-224.186,-212.446

0/1:112:0,72:0:0:72:2734:-343.835,-119.697,-98.023

1/1:112:0,112:0:0:112:4100:-368.767,-33.7154,0

0/0:112:0,0:0:0:0:0:-368.767,-368.767,-368.767

# Variant filtering

**Only Homozygous
reference**

| | | |
|---|---|---|
| 0/0:53:49,3:49:1823:3:103:0,-6.39174,-154.72 | 0/0:22:20,2:20:735:2:62:0,-0.733954,-60.5893 | 0/0:37:34,2:34:1262:2:67:0,-4.7389,-107.544 |
| 0/0:265:248,17:248:8589:17:592:0,-26.1668,-719.448 | 0/0:180:167,13:167:5745:13:447:0,-13.6264,-476.643 | 0/0:223:201,21:201:6904:21:716:0,-2.21506,-556.726 |
| 0/0:358:341,17:341:14409:17:568:0,-56.3297,-1243.15 | 0/0:250:237,13:237:9845:13:431:0,-36.1474,-845.732 | 0/0:260:238,22:238:9897:22:729:0,-12.3462,-823.558 |

**Only Homozygous
alternate**

| | | |
|---|---|---|
| 1/1:105:0,105:0:0:105:3678:-331.212,-31.6082,0 | 1/1:47:1,46:1:37:46:1559:-136.894,-10.4506,0 | 1/1:61:0,61:0:0:61:2103:-189.536,-18.3628,0 |

**Do they bring some information in our case (proband affected)
if we only consider genotypes?**

# Variant filtering

# Variant filtering



- **Metrics (INFO, FORMAT)**
- **Boolean expressions : AND (&), OR (|), NOT (!), etc.**
- **Operators : Less (<), Less or equal (<=), Equal (=), Different (!=), etc.**

https://samtools.github.io/bcftools/bcftools.html#expressions

# Variant filtering

🔧 **bcftools view** VCF/BCF conversion, view, subset and filter VCF/BCF files (Galaxy Version 1.10)  ☆  &  ▾

**VCF/BCF Data**

| 📄 | 📑 | 📁 | 15: freebayes_calling_norm.vcf | ▾ | ⬆ | 📂 |

Restrict to  👁

**Apply filters**

```

```

Skip sites where FILTER column does not contain any of the strings listed (e.g. "PASS,.") (--apply_filters)

**Regions**

Do not restrict to Regions ▾

**Targets**

Do not restrict to Targets ▾

**Include**

AF>0 & AF<1

Select sites for which the expression is true (--include)

**Exclude**

Exclude sites for which the expression is true (--exclude)

# Variant filtering

# Variant filtering

# Variant annotation

# Variant annotation

# Variant annotation

**Upstream / Downstream length**

| 5000 bases | ▼ |

(-ud)

**Set size for splice sites (donor and acceptor) in bases**

| 2 bases | ▼ |

(-ss)

**spliceRegion Settings**

| Use Defaults | ▼ |

# Variant annotation

**Annotation options**

☐ Select/Unselect all

☐ Use 'EFF' field compatible with older versions (instead of 'ANN')
☐ Use Classic Effect names and amino acid variant annotations (NON_SYNONYMOUS_CODING vs missense_variant and G180R vs p.Gly180Arg/c.538G>C)
☐ Override classic and use Sequence Ontolgy terms for effects (missense_variant vs NON_SYNONYMOUS_CODING)
☐ Override classic and use HGVS annotations for amino acid annotations (p.Gly180Arg/c.538G>C vs G180R)
☐ Old notation style notation: E.g. 'c.G123T' instead of 'c.123G>T' and 'X' instead of '*'
☐ Use one letter Amino acid codes in HGVS notation. E.g. p.R47G instead of p.Arg47Gly
☐ Use transcript ID in HGVS notation. E.g. ENST00000252100:c.914C>G instead of c.914C>G
☐ Do not shift variants according to HGVS notation (most 3prime end)
☐ Do not add HGVS annotations
☐ Only use canonical transcripts
☐ Only use protein coding transcripts
☐ Use gene ID instead of gene name (VCF output)
☐ Disable IUB code expansion in input variants
☐ Add OICR tag in VCF file
☐ Add loss of function (LOF) and nonsense mediated decay (NMD) tags
☐ Do not add LOF and NMD annotations
☐ Disable motif annotations
☐ Disable NextProt annotations
☐ Disable interaction annotations
☐ Perform 'cancer' comparisons (somatic vs. germline)

# Variant annotation

**Use custom interval file for annotation**

| 🗋 | 🗐 | 🗀 | No bed dataset available. ▾ | ⬆ | 🗁 |

(-interval)

**Only use the transcripts in this file**

| 🗋 | 🗐 | 🗀 | Nothing selected ▾ | ⬆ | 🗁 |

Format is one transcript ID per line

**Filter output**

☐ Select/Unselect all

☐ 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

**Filter out specific Effects**

| No ▾ |

# Variant annotation

**Chromosomal position**

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

**Text to prepend to chromosome name**

By default SnpEff simplifies all chromosome names. For instance 'chr1' is just '1'. You can prepend any string you want to the chromosome name (-chr)

**Produce Summary Stats**                    **1**

⬤ Yes

(-noStats)

**Suppress reporting usage statistics to server**

⬤ Yes

(-noLog)

**Email notification**

⬤

Send an email notification when the job completes.

✔ Execute    **2**

# Variant annotation - Content

## SnpEff: Variant analysis

### Contents

Summary
Variant rate by chromosome
Variants by type
Number of variants by impact
Number of variants by functional class
Number of variants by effect
Quality histogram
InDel length histogram
Base variant table
Transition vs transversions (ts/tv)
Allele frequency
Allele Count
Codon change table
Amino acid change table
Chromosome variants plots
Details by gene

**1**

18: SnpEff eff: on data 16 - HTML stats

17: SnpEff eff: on data 16

16: freebayes_calling_norm_filtered.vcf
father  mother  proband

15: freebayes_calling_norm.vcf
father  mother  proband

14: freebayes_calling.vcf
father  mother  proband

13: markdup_proband.ba

# Variant annotation - Summary

**Summary**

| | |
|---|---|
| **Genome** | hg19 |
| **Date** | 2022-03-25 11:34 |
| **SnpEff version** | SnpEff 4.3t (build 2017-11-24 10:18), by Pablo Cingolani |
| **Command line arguments** | SnpEff  -i vcf -o vcf -stats /shared/ifbstor1/galaxy/jobs/001/469/1469180/outputs/galaxy_dataset_c7e86a06-3ffe-4324-9794-c54ffaf3b4c8.dat hg19 /shared/ifbstor1/galaxy/datasets/002/674/dataset_2674023.dat |
| **Warnings** | 1,293 |
| **Errors** | 0 |
| **Number of lines (input file)** | 6,468 |
| **Number of variants (before filter)** | 6,468 |
| **Number of not variants (i.e. reference equals alternative)** | 0 |
| **Number of variants processed (i.e. after filter and non-variants)** | 6,468 |
| **Number of known variants (i.e. non-empty ID)** | 0 ( 0% ) |
| **Number of multi-allelic VCF entries (i.e. more than two alleles)** | 0 |
| **Number of effects** | 18,335 |
| **Genome total length** | 3,137,161,265 |
| **Genome effective length** | 146,364,022 |
| **Variant rate** | 1 variant every 22,628 bases |

# Variant annotation - Variants details

## Variants rate details

| Chromosome | Length | Variants | Variants rate |
|---|---|---|---|
| 8 | 146,364,022 | 6,468 | 22,628 |
| **Total** | **146,364,022** | **6,468** | **22,628** |

## Number variants by type

| Type | Total |
|---|---|
| **SNP** | 5,101 |
| **MNP** | 132 |
| **INS** | 423 |
| **DEL** | 739 |
| **MIXED** | 73 |
| **INV** | 0 |
| **DUP** | 0 |
| **BND** | 0 |
| **INTERVAL** | 0 |
| **Total** | **6,468** |

## Number of effects by impact

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| **HIGH** | 322 | 1.756% |
| **LOW** | 1,371 | 7.478% |
| **MODERATE** | 807 | 4.401% |
| **MODIFIER** | 15,835 | 86.365% |

## Number of effects by functional class

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| **MISSENSE** | 743 | 45.667% |
| **NONSENSE** | 4 | 0.246% |
| **SILENT** | 880 | 54.087% |

Missense / Silent ratio: 0.8443

# Variant annotation - Variants details

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| 3_prime_UTR_variant | 2,907 | 15.538% |
| 5_prime_UTR_premature_start_codon_gain_variant | 57 | 0.305% |
| 5_prime_UTR_variant | 440 | 2.352% |
| conservative_inframe_deletion | 2 | 0.011% |
| conservative_inframe_insertion | 4 | 0.021% |
| disruptive_inframe_deletion | 5 | 0.027% |
| downstream_gene_variant | 1,368 | 7.312% |
| frameshift_variant | 7 | 0.037% |
| intergenic_region | 236 | 1.261% |
| intragenic_variant | 1 | 0.005% |
| intron_variant | 9,544 | 51.013% |
| missense_variant | 766 | 4.094% |
| non_coding_transcript_exon_variant | 565 | 3.02% |
| non_coding_transcript_variant | 2 | 0.011% |
| protein_protein_contact | 6 | 0.032% |
| sequence_feature | 135 | 0.722% |
| splice_acceptor_variant | 13 | 0.069% |
| splice_donor_variant | 3 | 0.016% |
| splice_region_variant | 358 | 1.914% |
| start_lost | 2 | 0.011% |
| stop_gained | 7 | 0.037% |
| stop_lost | 3 | 0.016% |
| stop_retained_variant | 1 | 0.005% |
| structural_interaction_variant | 284 | 1.518% |
| synonymous_variant | 883 | 4.72% |
| upstream_gene_variant | 1,110 | 5.933% |

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| DOWNSTREAM | 1,368 | 7.461% |
| EXON | 2,507 | 13.673% |
| INTERGENIC | 236 | 1.287% |
| INTRON | 9,209 | 50.226% |
| SPLICE_SITE_ACCEPTOR | 11 | 0.06% |
| SPLICE_SITE_DONOR | 3 | 0.016% |
| SPLICE_SITE_REGION | 349 | 1.903% |
| TRANSCRIPT | 138 | 0.753% |
| UPSTREAM | 1,110 | 6.054% |
| UTR_3_PRIME | 2,907 | 15.855% |
| UTR_5_PRIME | 497 | 2.711% |

# Variant annotation - Variants details

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| 3_prime_UTR_variant | 2,907 | 15.538% |
| 5_prime_UTR_premature_start_codon_gain_variant | 57 | 0.305% |
| 5_prime_UTR_variant | 440 | 2.352% |
| conservative_inframe_deletion | 2 | 0.011% |
| conservative_inframe_insertion | 4 | 0.021% |
| disruptive_inframe_deletion | 5 | 0.027% |
| downstream_gene_variant | 1,368 | 7.312% |
| frameshift_variant | 7 | 0.037% |
| intergenic_region | 236 | 1.261% |
| intragenic_variant | 1 | 0.005% |
| intron_variant | 9,544 | 51.013% |
| missense_variant | 766 | 4.094% |
| non_coding_transcript_exon_variant | 565 | 3.02% |
| non_coding_transcript_variant | 2 | 0.011% |
| protein_protein_contact | 6 | 0.032% |
| sequence_feature | 135 | 0.722% |
| splice_acceptor_variant | 13 | 0.069% |
| splice_donor_variant | 3 | 0.016% |
| splice_region_variant | 358 | 1.914% |
| start_lost | 2 | 0.011% |
| stop_gained | 7 | 0.037% |
| stop_lost | 3 | 0.016% |
| stop_retained_variant | 1 | 0.005% |
| structural_interaction_variant | 284 | 1.518% |
| synonymous_variant | 883 | 4.72% |
| upstream_gene_variant | 1,110 | 5.933% |

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| DOWNSTREAM | 1,368 | 7.461% |
| EXON | 2,507 | 13.673% |
| INTERGENIC | 236 | 1.287% |
| INTRON | 9,209 | 50.226% |
| SPLICE_SITE_ACCEPTOR | 11 | 0.06% |
| SPLICE_SITE_DONOR | 3 | 0.016% |
| SPLICE_SITE_REGION | 349 | 1.903% |
| TRANSCRIPT | 138 | 0.753% |
| UPSTREAM | 1,110 | 6.054% |
| UTR_3_PRIME | 2,907 | 15.855% |
| UTR_5_PRIME | 497 | 2.711% |

# Variant annotation - Variants quality

**Quality:**

| | |
|---|---|
| **Min** | 0 |
| **Max** | 57,898 |
| **Mean** | 1,449.862 |
| **Median** | 691 |
| **Standard deviation** | 2,384.312 |
| **Values** | 0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,… |
| **Count** | 456,23,14,22,14,14,6,12,16,14,14,14,9,7,7,11,5,8,12,13,9,8,12,10,10,8,4,3,9,3,6,7,8,7,8,6,8,6,9,10,12,10… |

# Variant annotation - Insertions/Deletions

**Insertions and deletions length:**

| | |
|---|---|
| **Min** | 0 |
| **Max** | 23 |
| **Mean** | 1.104 |
| **Median** | 1 |
| **Standard deviation** | 1.693 |
| **Values** | 0,1,2,3,4,5,6,7,8,9,11,12,15,17,20,21,23 |
| **Count** | 259,797,31,35,7,11,5,4,2,1,3,2,1,1,1,1,1 |



Insertion deletion length histogram

# Variant annotation - Transitions/Transversions



**Base changes (SNPs)**

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0 | 207 | 762 | 163 |
| C | 253 | 0 | 233 | 885 |
| G | 1,014 | 255 | 0 | 219 |
| T | 140 | 763 | 207 | 0 |

Purines
Adenine — Transitions — Guanine

Transversions — Transversions

Cytosine — Transitions — Thymine
Pyrimidines

**Ts/Tv (transitions / transversions)**

**Note:** Only SNPs are used for this statistic.
**Note:** This Ts/Tv ratio is a 'raw' ratio (ratio of observed events).

| Transitions | 8,638 |
|---|---|
| Transversions | 4,186 |
| Ts/Tv ratio | 2.0635 |

**All variants:**

```
Sample  ,proband,mother,father,Total
 Transitions ,2917,2793,2928,8638
Transversions ,1437,1322,1427,4186
 Ts/Tv ,2.030,2.113,2.052,2.064
```

| Sequencing Type | # of Variants* | TiTv Ratio |
|---|---|---|
| WGS | ~4.4M | 2.0-2.1 |
| WES | ~41k | 3.0-3.3 |

*for a single sample

https://en.wikipedia.org/wiki/Transversion
https://gatk.broadinstitute.org/hc/en-us/articles/360035531572-Evaluating-the-quality-of-a-germline-short-variant-callset

# Variant annotation - Allele details

**Allele frequency**

Allele Frequency %



| | |
|---|---|
| **Min** | 16 |
| **Max** | 83 |
| **Mean** | 41.217 |
| **Median** | 33 |
| **Standard deviation** | 21.155 |
| **Values** | 16,25,33,50,66,75,83 |
| **Count** | 1665,53,1965,1229,968,45,543 |

**Allele Count**

Allele Count



| | |
|---|---|
| **Min** | 1 |
| **Max** | 5 |
| **Mean** | 2.462 |
| **Median** | 2 |
| **Standard deviation** | 1.262 |
| **Values** | 1,2,3,4,5 |
| **Count** | 1751,2029,1177,968,543 |

# Variant annotation - Genotypes details

**Hom/Het per sample**



Sample_names , proband, mother, father
Reference , 1998, 2082, 1922
Het , 3494, 3417, 3510
Hom , 957, 844, 952
Missing , 19, 125, 84

# Variant annotation - Codon changes

**Codon changes**

How to read this table:
- Rows are reference codons and columns are changed codons. E.g. Row 'AAA' column 'TAA' indicates how many 'AAA' codons have been replaced by 'TAA' codons.
- Red background colors indicate that more changes happened (heat-map).
- Diagonals are indicated using grey background color
- WARNING: This table may include different translation codon tables (e.g. mamalian DNA and mitochondrial DNA).

| | - | AAA | AAC | AAG | AAT | ACA | ACC | ACG | ACT | AGA | AGC | AGG | AGT | ATA | ATC | ATG | ATT | CAA | CAC | CAG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| - | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | 3 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp |
| AAA | 1 | &nbsp | 5 | 8 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | 2 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp |
| AAC | 2 | 3 | &nbsp | 1 | 28 | &nbsp | 3 | &nbsp | &nbsp | &nbsp | 13 | &nbsp | 3 | &nbsp | 3 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp |

# Variant annotation - Amino acid changes

**Amino acid changes**

How to read this table:
- Rows are reference amino acids and columns are changed amino acids. E.g. Row 'A' column 'E' indicates how many 'A' amino acids have been replaced by 'E' amino acids.
- Red background colors indicate that more changes happened (heat-map).
- Diagonals are indicated using grey background color
- WARNING: This table may include different translation codon tables (e.g. mamalian DNA and mitochondrial DNA).

|   | * | - | ? | A | C | D | E | F | G | H | I | K | L | M | N | P | Q | R | S | T | V | W |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| * | 1 | 1 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | 2 | &nbsp | &nbsp | &nbsp | &nbsp | &nbs |
| - | &nbsp | **&nbsp** | 1 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | 3 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | 3 | &nbsp | &nbsp | &nbsp | &nbs |
| ? | &nbsp | &nbsp | **&nbsp** | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbs |
| A | &nbsp | 1 | &nbsp | **166** | &nbsp | 1 | 1 | &nbsp | 3 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | 6 | 23 | 33 | &nbs | |
| C | &nbsp | 3 | &nbsp | &nbsp | **9** | &nbsp | &nbsp | &nbsp | 3 | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | &nbsp | 5 | &nbsp | &nbsp | &nbsp | &nbs |

# Variant annotation - Chromosomes details



**Variants by chromosome**

# Variant annotation - Genes information

**1**

### Details by gene

<u>**Here**</u> you can find a tab-separated table.

```
# The following table is formatted as tab separated values.
#GeneName       GeneId  TranscriptId    BioType variants_impact_HIGH     variants_impact_LOW     variants_impact_MODERATE        variants_impact_MODIFIER
variants_effect_3_prime_UTR_variant     variants_effect_5_prime_UTR_premature_start_codon_gain_variant  variants_effect_5_prime_UTR_variant
variants_effect_conservative_inframe_deletion   variants_effect_conservative_inframe_insertion  variants_effect_disruptive_inframe_deletion
variants_effect_downstream_gene_variant variants_effect_frameshift_variant      variants_effect_intron_variant  variants_effect_missense_variant
variants_effect_non_coding_transcript_exon_variant      variants_effect_non_coding_transcript_variant   variants_effect_protein_protein_contact
variants_effect_sequence_feature        variants_effect_splice_acceptor_variant variants_effect_splice_donor_variant    variants_effect_splice_region_variant
variants_effect_start_lost      variants_effect_stop_gained     variants_effect_stop_lost       variants_effect_stop_retained_variant
variants_effect_structural_interaction_variant  variants_effect_synonymous_variant      variants_effect_upstream_gene_variant
AARD    AARD    NM_001025357.2  protein_coding  0       0       0       2       2       0       0       0       0       0       0       0       0       0       0       0
0       0       0       0       0       0       0       0       0       0       0
ABRA    ABRA    NM_139166.4     protein_coding  0       3       0       3       1       0       0       0       0       0       0       0       1       0       0       0
0       0       0       0       0       0       0       0       3       1
ADAM18  ADAM18  NM_001190956.1  protein_coding  0       2       0       2       0       0       0       0       0       0       0       0       3       0       0       0
0       0       0       1       0       0       0       0       0       1       0
ADAM18  ADAM18  NM_001320313.1  protein_coding  0       3       0       10      0       0       0       0       0       0       0       0       11      0       0       0
0       0       0       1       0       0       0       0       0       2       0
```

# Variant annotation - ANN field

# Variant annotation - Examples

**Synonymous**

ANN=G|synonymous_variant|LOW|OR4F21|OR4F21|transcript|NM_001005504.1|protein_coding|1/1|c.324T>C|p.Gly108Gly|324/939|324/939|108/312||

**Missense**

ANN=G|missense_variant|MODERATE|FBXO25|FBXO25|transcript|NM_183421.1|protein_coding|3/11|c.138C>G|p.Ile46Met|404/2441|138/1104|46/367||,

**Intronic**

ANN=G|intron_variant|MODIFIER|FBXO25|FBXO25|transcript|NM_183421.1|protein_coding|1/10|c.-7-166C>G||||||

'Allele | Annotation | Annotation_Impact | Gene_Name | Gene_ID | Feature_Type | Feature_ID | Transcript_BioType | Rank | HGVS.c | HGVS.p |
cDNA.pos / cDNA.length | CDS.pos / CDS.length | AA.pos / AA.length | Distance | ERRORS / WARNINGS / INFO' ">

# Variant reporting - Pedigree

**Individual**

**Family ID**

**Father ID**

**Mother ID**

| FAM | father | 0 | 0 | 1 | 1 |
| FAM | mother | 0 | 0 | 2 | 1 |
| FAM | proband | father | mother | 1 | 2 |

**Sex (1: male; 2: female)**    **Status (1: control; 2: case)**

4: Pedigree.txt

**1**

# Variant reporting - Database creation

# Variant reporting - Database creation

# Variant reporting - Database creation

# Variant reporting - Database creation

## Dataset Information

| | |
|---|---|
| Number | 19 |
| Name | GEMINI load on data 4 and data 17 |
| Created | Friday Mar 25th 2:37:11 2022 UTC |
| Filesize | **190.8** MB |
| Dbkey | hg19 |
| Format | gemini.sqlite |
| File contents | contents |
| History Content API ID | 319b4d6eefbba9f5 |
| History API ID | 57e9be0d003985de |
| UUID | f41f617b-fc1c-4840-9ee4-cf206a5c4555 |

## Tool Parameters

| Input Parameter | Value |
|---|---|
| VCF dataset to be loaded in the GEMINI database | 17 freebayes_calling_norm_filtered_annotated.vcf  father  mother  proband |
| The variants in this input are | annotated with snpEff |
| This input comes with genotype calls for its samples | True |
| Choose a gemini annotation source | 2022-03-23 |
| Sample and family information in PED format | 4 Pedigree.txt |
| Load the following optional content into the database | GERP scores CADD scores (non-commercial use only; see licensing note below) Gene tables Sample genotypes variant INFO field |

## Job Outputs

| Tool Outputs | Dataset |
|---|---|



search datasets

TP_GTN_WES_disease

19 shown

2.23 GB

**19: GEMINI load on data 4 and data 17**

father  mother  proband

190.8 MB

format: **gemini.sqlite**, database: **hg19**

Indexing /shared/ifbstor1/galaxy/jobs/001/474 with grabix.
Loading 6468 variants.
Breaking /shared/ifbstor1/galaxy/jobs/001/474 into 12 chunks.
Loading chunk 0.
Loading chunk 1.
Loading chunk 2.
L

Gemini SQLite Database, version 0.20.1

# Variant reporting - Database content

# Variant reporting - Database content

# Variant reporting - Database content

| table_name | column_name | type |
|---|---|---|
| variants | chrom | VARCHAR(20) |
| variants | start | INTEGER |
| variants | end | INTEGER |
| variants | vcf_id | TEXT |
| variants | variant_id | INTEGER |
| variants | anno_id | INTEGER |
| variants | ref | TEXT |
| variants | alt | TEXT |
| variants | qual | FLOAT |
| variants | filter | TEXT |
| variants | type | VARCHAR(20) |
| variants | sub_type | TEXT |
| variants | gts | BLOB |
| variants | gt_types | BLOB |
| variants | gt_phases | BLOB |
| variants | gt_depths | BLOB |
| variants | gt_ref_depths | BLOB |
| variants | gt_alt_depths | BLOB |
| variants | gt_alt_freqs | BLOB |
| variants | gt_quals | BLOB |
| variants | gt_copy_numbers | BLOB |
| variants | call_rate | FLOAT |
| variants | max_aaf_all | FLOAT |
| variants | in_dbsnp | BOOLEAN |
| variants | rs_ids | TEXT |

| | | |
|---|---|---|
| variant_impacts | variant_id | INTEGER |
| variant_impacts | anno_id | INTEGER |
| variant_impacts | gene | VARCHAR(60) |
| variant_impacts | transcript | VARCHAR(60) |
| variant_impacts | is_exonic | BOOLEAN |
| variant_impacts | is_coding | BOOLEAN |
| variant_impacts | is_lof | BOOLEAN |
| variant_impacts | exon | TEXT |
| variant_impacts | codon_change | TEXT |
| variant_impacts | aa_change | TEXT |
| variant_impacts | aa_length | TEXT |
| variant_impacts | biotype | TEXT |
| variant_impacts | impact | VARCHAR(60) |
| variant_impacts | impact_so | TEXT |
| variant_impacts | impact_severity | VARCHAR(20) |
| variant_impacts | polyphen_pred | TEXT |
| variant_impacts | polyphen_score | FLOAT |
| variant_impacts | sift_pred | TEXT |
| variant_impacts | sift_score | FLOAT |

# Variant reporting - Database content

| | | |
|---|---|---|
| samples | sample_id | INTEGER |
| samples | family_id | TEXT |
| samples | name | TEXT |
| samples | paternal_id | TEXT |
| samples | maternal_id | TEXT |
| samples | sex | TEXT |
| samples | phenotype | TEXT |

| | | |
|---|---|---|
| gene_detailed | uid | INTEGER |
| gene_detailed | chrom | VARCHAR(60) |
| gene_detailed | gene | VARCHAR(60) |
| gene_detailed | is_hgnc | BOOLEAN |
| gene_detailed | ensembl_gene_id | TEXT |
| gene_detailed | transcript | VARCHAR(60) |
| gene_detailed | biotype | TEXT |
| gene_detailed | transcript_status | TEXT |

| | | |
|---|---|---|
| gene_summary | uid | INTEGER |
| gene_summary | chrom | VARCHAR(60) |
| gene_summary | gene | VARCHAR(60) |
| gene_summary | is_hgnc | BOOLEAN |
| gene_summary | ensembl_gene_id | TEXT |
| gene_summary | hgnc_id | TEXT |
| gene_summary | transcript_min_start | INTEGER |
| gene_summary | transcript_max_end | INTEGER |
| gene_summary | strand | TEXT |
| gene_summary | synonym | TEXT |

# Variant reporting - Querying

# Variant reporting - Querying



**Which inheritance pattern to select ?**

# Variant reporting - Inheritance pattern

# Variant reporting - Inheritance pattern



X-linked dominant

Note: some X-linked dominant disorders are embryonic lethal in males, and most affect females less severely.

# Variant reporting - Inheritance pattern



X-linked recessive

# Variant reporting - Inheritance pattern

- Autosomal de-novo : mutation on autosomes (chr1-22), mutation not present in parents
- X-linked de-novo : mutation on the sex chromosome X, mutation not present in parents
- Compound heterozygous : 2 or more recessive alleles at a particular locus
- Violation of mendelian laws :
  - LOH : Loss of Heterozygosity, cross chromosomal event resulting in in loss of an entire gene and the surrounding chromosomal region
  - Plausible de-novo : parents are homozygous reference, offspring is heterozygous
  - Implausible de-novo : parents are homozygous reference, offspring is homozygous alternate
  - Uniparental disomy : one parent and the offspring are homozygous reference, the other parent is homozygous alternate OR one parent and the offspring are homozygous alternate and the other parent is homozygous reference

# Variant reporting - Inheritance pattern

- Autosomal recessive
- Autosomal dominant
- X-linked recessive
- X-linked dominant
- Autosomal de-novo
- X-linked de-novo
- Compound heterozygous
- Violation of mendelian laws

# Variant reporting - Inheritance pattern

- Autosomal recessive
- Autosomal dominant
- X-linked recessive
- X-linked dominant
- Autosomal de-novo
- X-linked de-novo
- Compound heterozygous
- Violation of mendelian laws

**Parents are unaffected**

# Variant reporting - Inheritance pattern

- Autosomal recessive
- Autosomal dominant
- X-linked recessive
- X-linked dominant
- Autosomal de-novo
- X-linked de-novo
- Compound heterozygous
- Violation of mendelian laws

**Parents are unaffected**

**Parents are consiguineous**

# Variant reporting - Inheritance pattern

- Autosomal recessive
- Autosomal dominant
- X-linked recessive
- X-linked dominant
- Autosomal de-novo
- X-linked de-novo
- Compound heterozygous
- Violation of mendelian laws

**Parents are unaffected**

**Parents are consiguineous**

**Chromosome 8**

# Variant reporting - Inheritance pattern

- Autosomal recessive
- ~~Autosomal dominant~~
- X-linked recessive
- ~~X-linked dominant~~
- Autosomal de-novo
- X-linked de-novo
- Compound heterozygous
- Violation of mendelian laws

**Parents are unaffected**

**Parents are consiguineous**

**Chromosome 8**

# Variant reporting - Inheritance pattern

- Autosomal recessive
- ~~Autosomal dominant~~
- ~~X-linked recessive~~
- ~~X-linked dominant~~
- Autosomal de-novo
- ~~X-linked de-novo~~
- Compound heterozygous
- Violation of mendelian laws

**Parents are unaffected**

**Parents are consanguineous**

**Chromosome 8**

# Variant reporting - Inheritance pattern

**1** • Autosomal recessive
• ~~Autosomal dominant~~
• ~~X-linked recessive~~
• ~~X-linked dominant~~
• Autosomal de-novo
• ~~X-linked de-novo~~
• Compound heterozygous
• Violation of mendelian laws

**Parents are unaffected**

**Parents are consanguineous**

**Chromosome 8**

# Variant reporting - Inheritance pattern

**1** ● Autosomal recessive
● ~~Autosomal dominant~~
● ~~X-linked recessive~~
● ~~X-linked dominant~~
**2** ● Autosomal de-novo
● ~~X-linked de-novo~~
● Compound heterozygous
● Violation of mendelian laws

**Parents are unaffected**

**Parents are consanguineous**

**Chromosome 8**

# Variant reporting - Inheritance pattern

**1** • Autosomal recessive

• ~~Autosomal dominant~~

• ~~X-linked recessive~~

• ~~X-linked dominant~~

**2** • Autosomal de-novo

• ~~X-linked de-novo~~

**3** • Compound heterozygous

• Violation of mendelian laws

**Parents are unaffected**

**Parents are consanguineous**

**Chromosome 8**

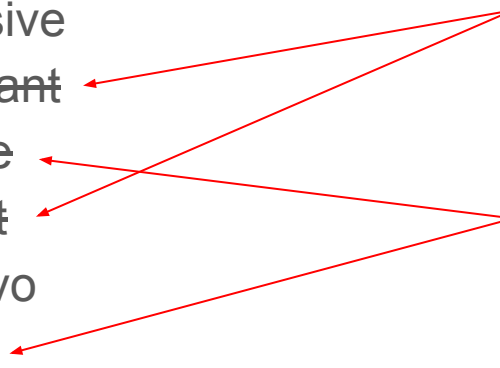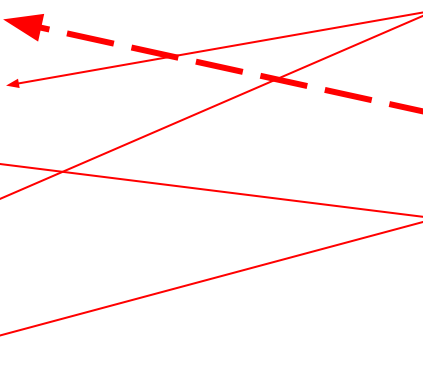# Variant reporting - Inheritance pattern

**1** ● Autosomal recessive

● ~~Autosomal dominant~~

● ~~X-linked recessive~~

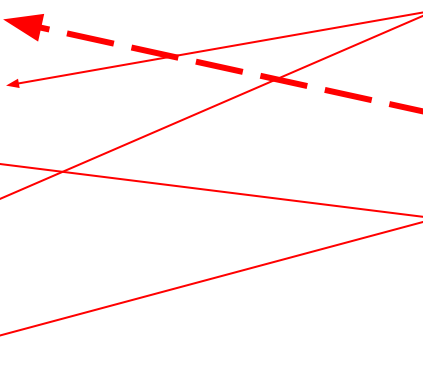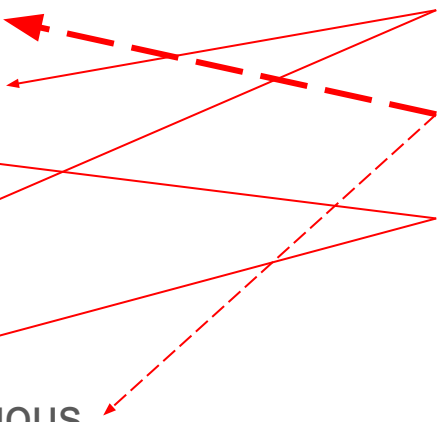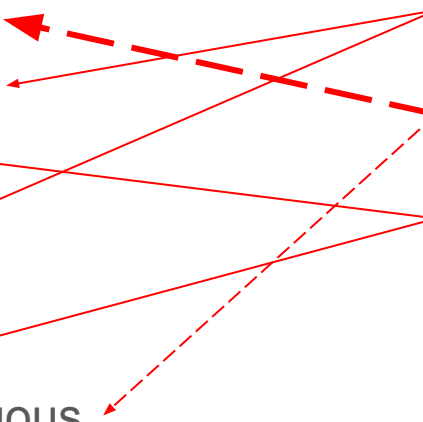● ~~X-linked dominant~~

**2** ● Autosomal de-novo

● ~~X-linked de-novo~~

**3** ● Compound heterozygous

**4** ● Violation of mendelian laws

**Parents are unaffected**

**Parents are consanguineous**

**Chromosome 8**

# Variant reporting - Querying

# Variant reporting - Querying

**Include hits with less convincing inheritance patterns**

⬤ No

The exact consequence of this setting depends on the type of inheritance pattern you are looking for (see the tool help below). (--lenient)

**Report candidates shared by unaffected samples**

⬤ No

Activating this option will enable the reporting of variants as candidate causative even if they are shared by unaffected samples in the family tree. The default will only report variants that are unique to affected samples. (--allow-unaffected)

Family-wise criteria for variant selection                                                                                        👁

**Minimum number of families with a candidate variant for a gene to be reported**

| 1 |
|---|

This is the number of families required to have a variant fitting the inheritance model in the same gene in order for the gene and its variants to be reported. For example, we may only be interested in candidates where at least 4 families have a variant (with a fitting inheritance pattern) in that gene. (--min-kindreds)

**List of families to restrict the analysis to (comma-separated)**

|  |
|---|

Leave empty for an analysis including all families (--families)

**Specify additional criteria to exclude families on a per-variant basis**

| No, analyze all variants from all included families                                                                    ▼ |
|---|

# Variant reporting - Querying

# Variant reporting - Querying

**Additional columns (comma-separated)**

chrom,start,ref,alt,impact,gene,clinvar_sig,clinvar_disease_name,clinvar_gene_phenotype,rs_ids

Column must be specified by the exact name they have in the GEMINI database, e.g., is_exonic or num_hom_alt, but, for genotype columns, GEMINI wildcard syntax is supported. The order of columns in the list is maintained in the output.

**Email notification**

Send an email notification when the job completes.

✔ Execute **8**

# Variant reporting - Results

| max_aaf_all | chrom | start | ref | alt | impact | gene | clinvar_sig | clinvar_disease_name |
|---|---|---|---|---|---|---|---|---|
| 0.6831 | chr8 | 2048830 | A | G | missense_variant | MYOM2 | None | None |
| 0.6716 | chr8 | 6479041 | C | T | missense_variant | MCPH1 | benign | Primary_autosomal_recessive_microcephaly_1\|not_specified\|Primary_Microcepha |
| 0.935555555556 | chr8 | 6681255 | A | C | splice_region_variant | XKR5 | None | None |
| -1.0 | chr8 | 11666217 | GTCCCAC | G | conservative_inframe_deletion | FDFT1 | None | None |
| 0.7798 | chr8 | 12878806 | T | G | missense_variant | KIAA1456 | None | None |
| 0.8221 | chr8 | 12879098 | G | A | missense_variant | KIAA1456 | None | None |
| 0.8221 | chr8 | 12879538 | A | G | missense_variant | KIAA1456 | None | None |
| 0.8313 | chr8 | 17434640 | G | C | splice_region_variant | PDGFRL | None | None |
| 0.847026781661 | chr8 | 17743019 | G | A | missense_variant | FGL1 | None | None |
| -1.0 | chr8 | 17796381 | AC | GT | missense_variant | PCM1 | None | None |
| 0.842472840145 | chr8 | 17814914 | A | G | missense_variant | PCM1 | None | None |

History

search datasets

TP_GTN_WES_disease

21 shown

2.23 GB

21: GEMINI autosomal_recessive pattern on data 19

father  mother  proband

1

# Variant reporting - Results

| clinvar_gene_phenotype |
| --- |
| None |
| primary_microcephaly\x2c_recessive\|primary_autosomal_recessive_microcephaly_1 |
| None |
| None |
| None |
| None |
| None |
| carcinoma_of_colon |

# Variant reporting - Results

| rs_ids | variant_id | family_id | family_members | family_genotypes | samples | family_count |
|--------|-----------|-----------|----------------|------------------|---------|--------------|
| rs968381 | 228 | FAM | proband(proband;affected;male),mother(mother;unaffected;female),father(father;unaffected;male) | G/G,A/G,A/G | proband | 1 |
| rs1057090 | 462 | FAM | proband(proband;affected;male),mother(mother;unaffected;female),father(father;unaffected;male) | T/T,C/T,C/T | proband | 1 |
| rs9772979 | 490 | FAM | proband(proband;affected;male),mother(mother;unaffected;female),father(father;unaffected;male) | C/C,A/C,A/C | proband | 1 |
| rs71711801 | 862 | FAM | proband(proband;affected;male),mother(mother;unaffected;female),father(father;unaffected;male) | G/G,GTCCCAC/G,GTCCCAC/G | proband | 1 |
| rs3739310 | 936 | FAM | proband(proband;affected;male),mother(mother;unaffected;female),father(father;unaffected;male) | G/G,T/G,T/G | proband | 1 |
| rs545589847,rs502882 | 939 | FAM | proband(proband;affected;male),mother(mother;unaffected;female),father(father;unaffected;male) | A/A,G/A,G/A | proband | 1 |

# Variant reporting - Results

**Most likely variant candidate for child's disease ?**

# Variant reporting - Results

| max_aaf_all | chrom | start | ref | alt | impact | gene | clinvar_sig | clinvar_disease_name |
|---|---|---|---|---|---|---|---|---|
| 3.24886289799e-05 | chr8 | 86385979 | G | A | stop_gained | CA2 | None | None |

clinvar_gene_phenotype

carbonic_anhydrase_ii_variant|osteopetrosis_with_renal_tubular_acidosis

| rs_ids | variant_id | family_id | family_members | family_genotypes |
|---|---|---|---|---|
| None | 3883 | FAM | proband(proband;affected;male),mother(mother;unaffected;female),father(father;unaffected;male) | A/A,G/A,G/A |