

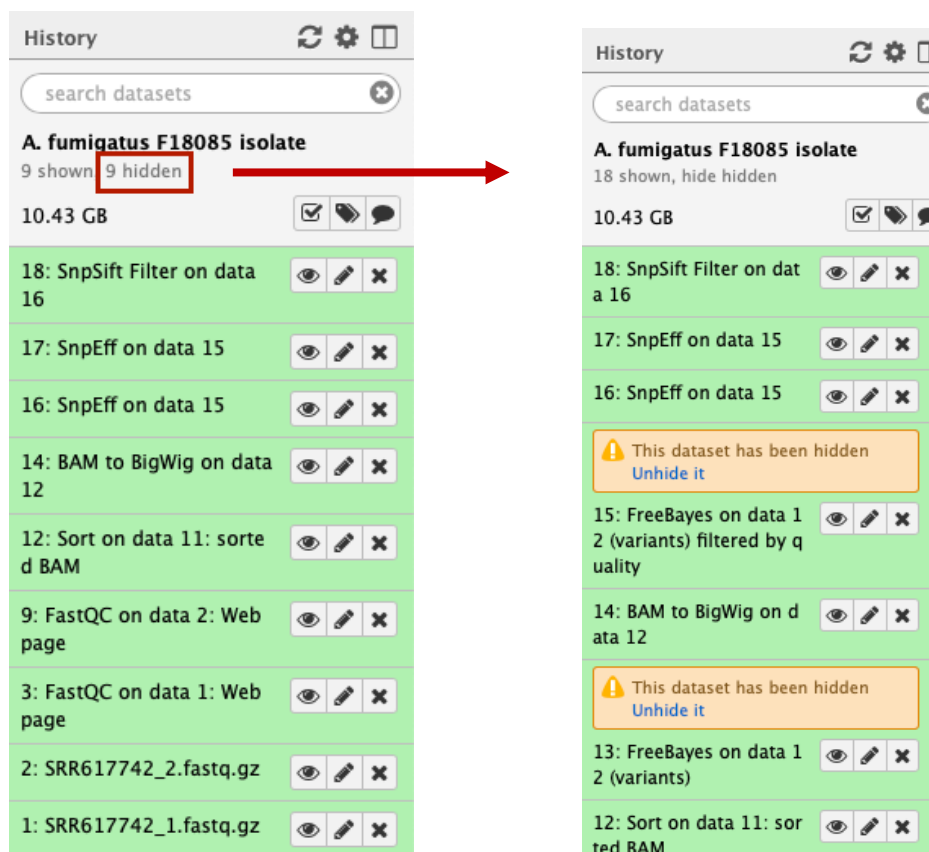
Variant calling in VEuPathDB galaxy (Part 2)

Learning objectives:

1. [Examine your results and the outputs of the workflow](#)
2. [View VCF file results in the JBrowse genome browser](#)
3. [Filter data in VCF files](#)
4. [Examine the filtered VCF file, extract Gene IDs, and create a Venny diagram](#)
5. [Share or Publish history](#)

1. Examine your results and the outputs of the workflow

- a. Examine your results:
Click on the hidden files link in the history panel to reveal all workflow output files.



- The output of Sickle is used by a program called Bowtie2.
- Bowtie generates a file called a BAM file. Whenever dealing with sequence alignment files you will likely hear of file formats called SAM or BAM. SAM stands for Sequence Alignment/Map format, and BAM is the binary version of a SAM file. Many of the downstream analysis programs that use BAM files require a sorted BAM file. This allows access to reads to be done more efficiently.
- The sorted BAM file is the input for a program called FreeBayes. This program is a Bayesian genetic variant detector designed to find small

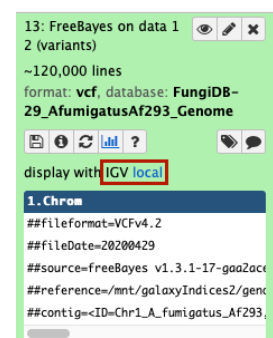
polymorphisms, specifically SNPs (single-nucleotide polymorphisms), indels (insertions and deletions), MNPs (multi-nucleotide polymorphisms), and complex events (composite insertion and substitution events) smaller than the length of a short-read sequencing alignment.

- The output for many variant callers is a file called a VCF file. VCF stands for variant interchange format.
 - SnpEff is a variant annotation and effect prediction tool. It annotates and predicts the effects of variants on genes (such as amino acid changes).
- b. Examine the output files.
- The output of Sickle is used by a program called Bowtie2.
 - Bowtie generates a file called a BAM file. Whenever dealing with sequence alignment files you will likely hear of file formats called SAM or BAM. SAM stands for Sequence Alignment/Map format, and BAM is the binary version of a SAM file. Many of the downstream analysis programs that use BAM files require a sorted BAM file. This allows access to reads to be done more efficiently.
 - The sorted BAM file is the input for a program called FreeBayes. This program is a Bayesian genetic variant detector designed to find small polymorphisms, specifically SNPs (single-nucleotide polymorphisms), indels (insertions and deletions), MNPs (multi-nucleotide polymorphisms), and complex events (composite insertion and substitution events) smaller than the length of a short-read sequencing alignment.
 - The output for many variant callers is a file called a VCF file. VCF stands for variant interchange format.
 - SnpEff is a variant annotation and effect prediction tool. It annotates and predicts the effects of variants on genes (such as amino acid changes).
- c. Examine the VCF file in your results (click on the eye icon to view its contents). Detailed information about VCF file content is available here: <https://samtools.github.io/hts-specs/VCFv4.2.pdf>

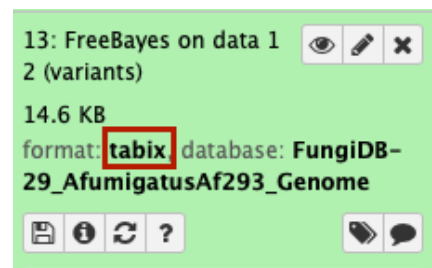
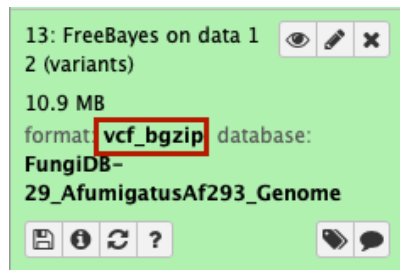
2. View VCF file results in the JBrowse genome browser:

In order to view a VCF file in JBrowse, it first has to be indexed and compressed. We will use a trick in galaxy to produce the required files:

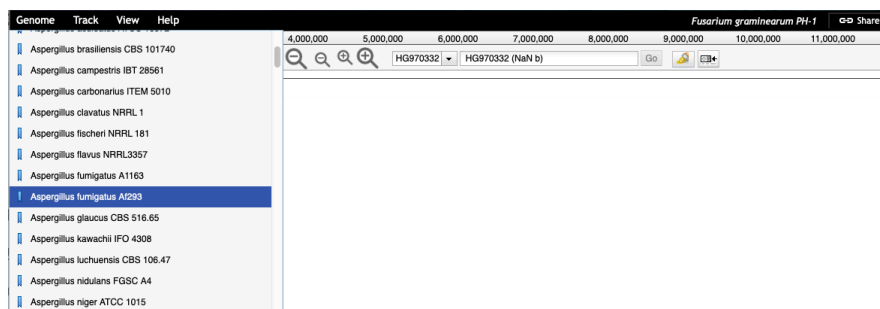
- Find the VCF file you would like to view in JBrowse. This can be the FreeBayes output or the quality filtered VCF. Click on this file to expose the available options.
- Click on “display with IGV local” (red circle).
- A new window will open up (you are not going to use this window).
- Go back to the galaxy window, you will notice that the number of hidden files has increased by 2 files.
- Show the hidden files by clicking on the word hidden.



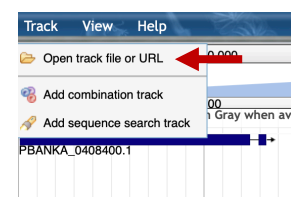
- f. The two new files you want are in the **vcf_bgzip** and **tabix** format:

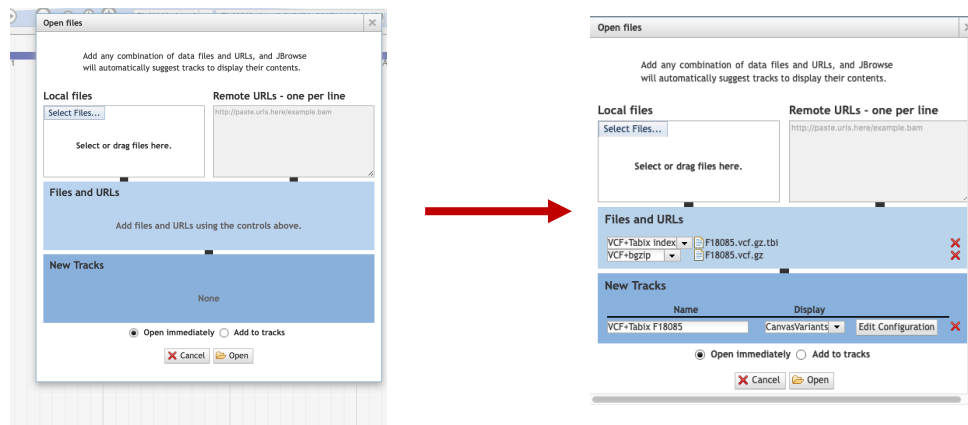


- g. Download both files by clicking on the download icon. *You will need both files.*
- h. After the files are downloaded, After the files are downloaded, rename them as follows:
- Rename the *Galaxy13-[FreeBayes_on_data_12_(variants)].vcf_bgzip* file to *sample.vcf.gz* (i.e. F18085.vcf.gz)
 - Rename the *Galaxy13-[FreeBayes_on_data_12_(variants)].tabix* file to *sample.vcf.gz.tbi* (i.e. F18085.vcf.gz.tbi)
- i. Next, navigate to the Tool section in FungiDB, click on Genome Browser link and select *Aspergillus fumigatus* Af293 from the Genome drop down list:

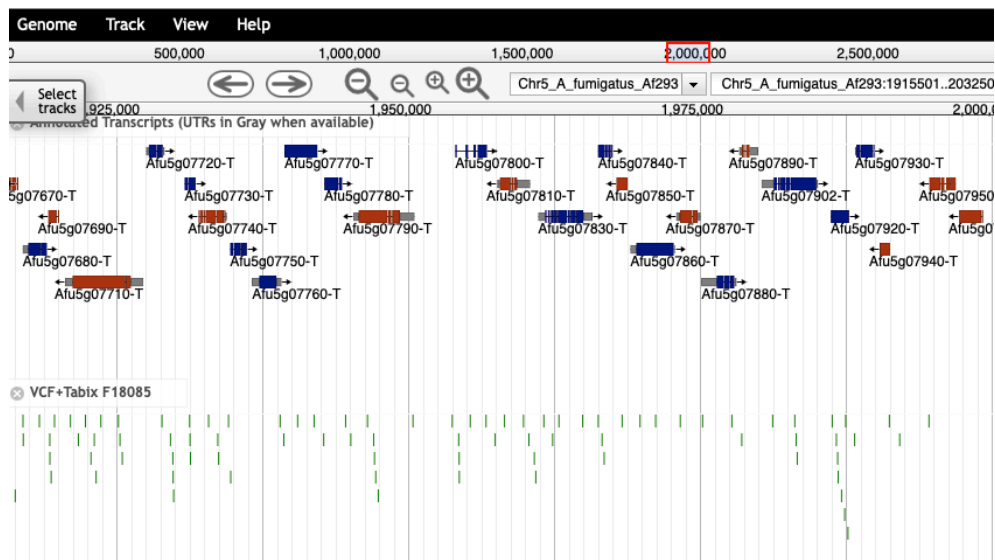


- j. Click on the Track menu, select “Open track file or URL”
- k. Drag and drop your files in the window that appears. Notice that the file formats are autodetected.





- I. Click on “Open”. You should see SNP positions displayed in a new track. You can zoom in and click on the SNPs to get more info.

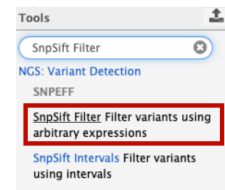


3. Filter data in VCF files:

VCF files contain a lot of data about variants and their positions.

SnpEff generates various analyses/summaries of VCF files (including GeneIDs that overlap variant positions). However, it is often necessary to filter VCF files further to obtain useful information for your specific question. For example, you may want to filter out SNP positions that have an impact on the coding sequence. One tool that can be used is called SnpSift Filter. This tool allows you to write complex expressions to filter a VCF file.

- Navigate to the Tools section and find SnpSift Filter.
- Click on the tool.
- Enter Expression argument



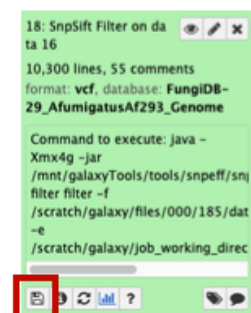
We will run the filter analysis on SNPs with high and moderate impact. Copy and paste the following expression into the SnpSift Filter tool:

```
((ANN[*].IMPACT has 'HIGH') | (ANN[*].IMPACT has 'MODERATE')) & ((na FILTER) | (FILTER = 'PASS'))
```

4. Examine the filtered VCF file, extract Gene IDs, and create a Venny diagram

Notice that the GeneIDs are buried in the file but the file has some structure which means you can extract them either programmatically or using a program like Excel. Here are some steps you can take to extract Gene IDs from two VCF files then compare them to identify genes that are in common or that distinguish the two files.

- Download the SnpSift Filter output by clicking on the save icon



- b. Open this file using Excel. Select a column to split up, choose Data>Tab to Column and make sure you select tabs and | as column delimiters

The first screenshot shows the 'Text Import Wizard - Step 1 of 3'. It indicates that the data is delimited and offers options for 'Delimited' (selected) or 'Fixed width'. The 'Delimited' option is chosen, and the 'File origin' is set to 'Macintosh'. The 'Preview of selected data' shows a VCF file header.

The second screenshot shows the 'Text Import Wizard - Step 2 of 3'. It allows setting delimiters. The 'Delimiters' section has 'Tab' selected. The 'Text qualifier' is set to '"'. The 'Preview of selected data' shows the same VCF file header.

The third screenshot shows the Excel interface with the VCF file imported. The data is pasted into cell A1. The formula bar shows the text: `##fileformat=VCFv4.1`. A red arrow points from the 'Finish' button in the wizard to the Excel interface.

- c. Now you can look for Gene IDs of interest in the excel file. For example, if this is a known drug resistant line you can find the gene(s) that might be responsible for the resistance and see what kinds of SNPs are present.

If you are comparing a mutant and a wild type or two different strains you can extract gene IDs from both VCF files and use Venny website (<https://bioinfogp.cnb.csic.es/tools/venny/>) to visualize your data.

The screenshot shows the Venny 2.1 website interface. The 'UPPERCASE' section shows two lists of gene IDs: 'List 3' (0) and 'List 4' (0). The 'Results' section shows a Venn diagram comparing the two lists.

The Venn diagram shows the overlap between the two lists. The left circle (blue) represents 'F18085' with 215 unique elements (44.3%). The right circle (yellow) represents 'AF90-PE-ILI27-01' with 262 unique elements (54%). The intersection contains 8 elements (1.6%).

*Note that in the above steps you are ultimately comparing gene IDs – do you think you might be missing some important polymorphisms using this method? Of course, the answer is yes 😊

It is quite possible that a gene with a SNP in the WT and a SNP in the mutant that will be in the intersection of the two gene lists, contains different SNPs – you will miss this by doing the above steps. Below is a description of steps you can take to create a list of unique IDs for SNPs. This list of unique IDs can then be used in Venny.

1. Start with the same excel files that you opened in the above section.
2. To create a unique ID for SNPs we will combine information from multiple columns to create something that looks like this:
chromosome:position:geneID
3. To do this you will use the concatenate function in Excel:
=concatenate(cell#1,".",cell#2,".",cell#3)

Cell#1 = cell with chromosome number

Cell#2 = cell with position

Cell#3 = cell with GeneID

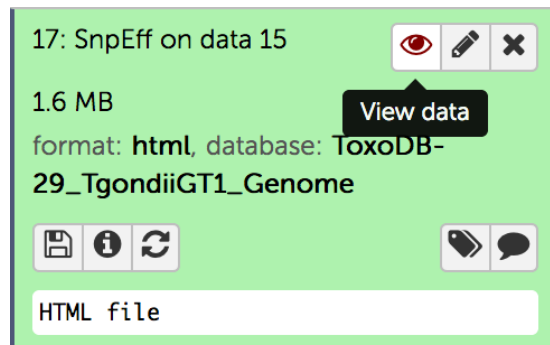
#	CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	unknown
65	TGGT1_chrla	227230		A	C	1156.55		AB=0;ABP=0; missense_va MODERATE	TGGT1_293300	
66	TGGT1_chrla	1340271		G	C	2387.77		AB=0;ABP=0; missense_va MODERATE	TGGT1_295040	
67	TGGT1_chrla	1396177		A	C	387.162		AB=0;ABP=0; missense_va MODERATE	TGGT1_295125	
68	TGGT1_chrlb	78769		A	G	1780.8		AB=0;ABP=0; missense_va MODERATE	TGGT1_207440	
69	TGGT1_chrlb	153771		T	G	1414.57		AB=0;ABP=0; missense_va MODERATE	TGGT1_207480	
70	TGGT1_chrlb	276348		T	G	2066.14		AB=0;ABP=0; missense_va MODERATE	TGGT1_207750	
71	TGGT1_chrlb	622140		G	C	2335.06		AB=0;ABP=0; missense_va MODERATE	TGGT1_208310	
72	TGGT1_chrlb	1446003		C	T	60.6579		AB=0;ABP=0; missense_va MODERATE	TGGT1_209755B	
73	TGGT1_chrlb	1446022		G	T	82.4046		AB=0;ABP=0; missense_va MODERATE	TGGT1_209755B	

#	CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	unknown
65	TGGT1_chrla	227230		A	C	1156.55		AB=0;ABP=0; missense_va MODERATE	TGGT1_293300	
66	TGGT1_chrla	1340271		G	C	2387.77		AB=0;ABP=0; missense_va MODERATE	TGGT1_295040	
67	TGGT1_chrla	1396177		A	C	387.162		AB=0;ABP=0; missense_va MODERATE	TGGT1_295125	
68	TGGT1_chrlb	78769		A	G	1780.8		AB=0;ABP=0; missense_va MODERATE	TGGT1_207440	
69	TGGT1_chrlb	153771		T	G	1414.57		AB=0;ABP=0; missense_va MODERATE	TGGT1_207480	
70	TGGT1_chrlb	276348		T	G	2066.14		AB=0;ABP=0; missense_va MODERATE	TGGT1_207750	
71	TGGT1_chrlb	622140		G	C	2335.06		AB=0;ABP=0; missense_va MODERATE	TGGT1_208310	
72	TGGT1_chrlb	1446003		C	T	60.6579		AB=0;ABP=0; missense_va MODERATE	TGGT1_209755B	
73	TGGT1_chrlb	1446022		G	T	82.4046		AB=0;ABP=0; missense_va MODERATE	TGGT1_209755B	

4. You should get unique SNP IDs that look like this (for example):
TGGT1_chrlb:1446003:TGGT1_209755B
5. Copy this function to the rest of the column to replicate the concatenate function.
6. Copy the these newly generated unique IDs into Venny and compare the mutant and wild type.

Supporting information - Examining SnpEff summary:

- Click on the view icon (eye) in the SnpEff output file that has the html format.



- This will open the html file right in galaxy where you can view it.
 - The header contains a short summary and information about the run and it has several major components:
1. Summary table that warns about possible genomic annotation errors or inconsistencies identified in the reference genome. If there are many, use caution interpreting results and examine associated gff files for any issues (ex. missing feature values in gff files, incomplete gene sequences, more than one stop codon per gene, etc.).
 2. Summary statistics for variant types

Number variantss by type

Type	Total
SNP	114,034
MNP	12,864
INS	6,907
DEL	7,304
MIXED	2,180
INTERVAL	0
Total	143,289

Here is an example of variant calls and what they mean in terms of nucleotide changes:

Type	What is means	Example
SNP	Single-Nucleotide Polymorphism	Reference = 'A', Sample = 'C'
Ins	Insertion	Reference = 'A', Sample = 'AGT'
Del	Deletion	Reference = 'AC', Sample = 'C'
MNP	Multiple-nucleotide polymorphism	Reference = 'ATA', Sample = 'GTC'
MIXED	Multiple-nucleotide and an InDel	Reference = 'ATA', Sample = 'GTCAGT'

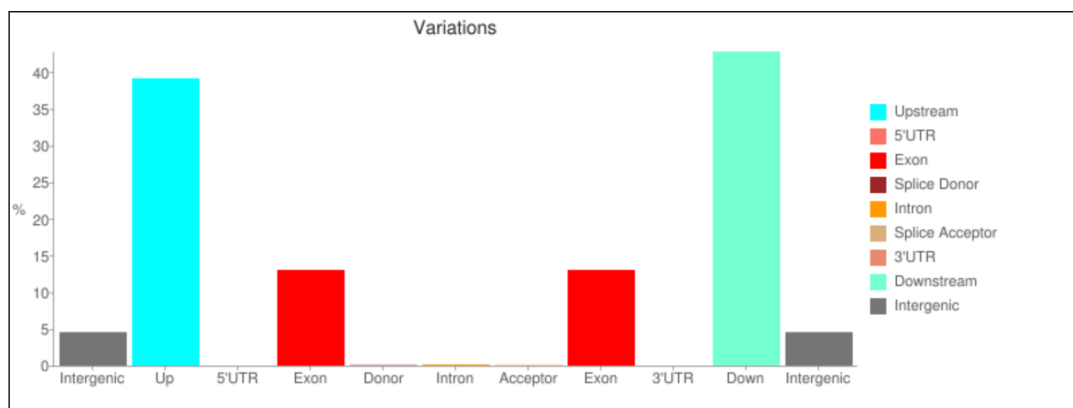
3. Statistics for the variant effects and impacts:

Number of effects by functional class

Type (alphabetical order)	Count	Percent
MISSENSE	21,588	35.949%
NONSENSE	131	0.218%
SILENT	38,332	63.832%

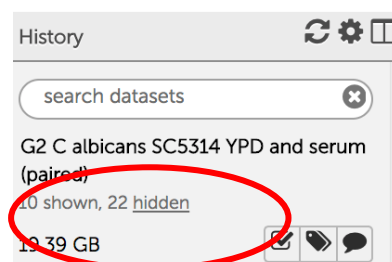
Type (alphabetical order)	Count	Percent
DOWNSTREAM	321,858	40.292%
EXON	67,505	8.451%
INTERGENIC	74,749	9.358%
INTRON	1,064	0.133%
NONE	1	0%
SPLICE_SITE_ACCEPTOR	5	0.001%
SPLICE_SITE_DONOR	4	0.001%
SPLICE_SITE_REGION	176	0.022%
TRANSCRIPT	12	0.002%
UPSTREAM	333,432	41.741%

Base changes summary. SnpEff html files provides a breakdown of SNPs across gene features:



The SNP workflow you are using is set up to generate certain files that will provide you with the information you can export and use further in your analysis (yellow stars).

If you select certain options they will be shown in your history. If you do not select to display these files, you can view the output by clicking on displaying the hidden files from the history menu:



Now, let's take a look at the files generated by the workflow and steps that you can take to further evaluate them.

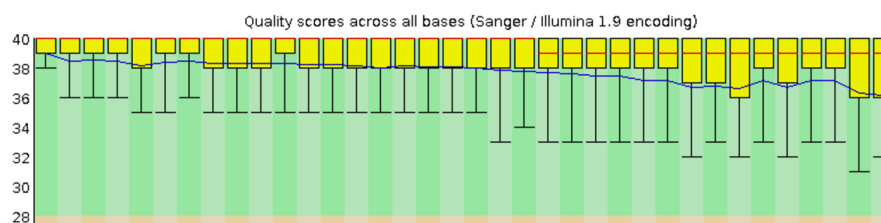
1. Examine sequence quality based on FastQC quality scores.

FastQC provides an easy-to-navigate visual representation sequencing data quality and distribution of nucleotides per read position.

Basic Statistics

Measure	Value
Filename	SRR298691.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	4887868
Sequences flagged as poor quality	0
Sequence length	36
%GC	58

Per base sequence quality



2. Download vcf files and evaluate workflow results.

The vcf file generated by SnpEff contains information about SNPs and the genomic location.

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	unknown		
CM001231	189057	.	AG	CT	787.449	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:143:0:0:143:5341:-207.887,-43.0473,0				
CM001231	483825	.	G	A	64.8756	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:4:0:0:4:146:-10.0999,-1.20412,0				
CM001231	518226	.	G	C	51.7908	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:8:0:0:7:276:-11.5007,-2.10721,0				
CM001231	574021	.	C	G	237.265	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:17:0:0:17:583:-39.079,-5.11751,0				
CM001231	609879	.	GAA	CAG	55.2785	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:32:8:277:22:861:-18.1711,-0.694735,0				
CM001231	1090073	.	G	T	79.4156	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:8:2:75:6:238:-11.5539,-1.36362,0				
CM001231	1090104	.	A	T	70.961	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:6:0:0:6:220:-12.5146,-1.80618,0				
CM001231	1153611	.	CCTC	GCTG	111.123	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:8:5:188:3:97:-9.30616,-6.1461,0				
CM001231	1159150	.	CT	GC	126.126	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:31:0:0:19:741:-29.7713,-5.71957,0				
CM001231	1159438	.	C	G	82.3312	.	AB=0;ABP=0;GT:DP:RO:QF 0/0:47:30:1092:17:640:0,-9.53002,-3.50705				
CM001231	1159465	.	G	C	249.656	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:126:47:1770:79:3013:-53.8644,-25.2134,0				
CM001231	1159499	.	T	C	124.95	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:143:32:1167:111:4248:-76.1575,-33.4865,0				
CM001231	1181576	.	CC	TG	191.675	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:27:0:0:25:924:-41.7448,-7.52575,0				
CM001231	1293309	.	C	G	51.22	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:2:0:0:2:78:-6.92763,-0.60206,0				
CM001231	1323058	.	TT	GC	71.3001	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:6:0:0:6:223:-12.5485,-1.80618,0				
CM001231	1485397	.	A	G	3558.42	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:499:0:0:497:18671:-804.678,-149.612,0				
CM001231	1485429	.	G	A	3783.33	.	AB=0;ABP=0;GT:DP:RO:QF 1/1:517:1:38:516:20010:-843.425,-151.978,0				

Post-processing of SNP data is normally required to make sense of thousands of SNPs and to decide which ones have biological and functional importance. Data processing can help you to extract SNP distribution and parse associated data including Gene IDs, protein-coding annotations, and effects in sequence ontology terms such as missense or synonymous variants, stop codon gain, etc. and also link changes to the genome model.

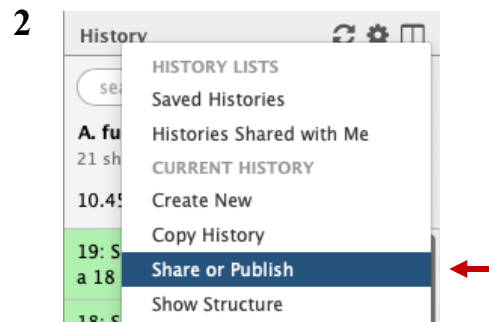
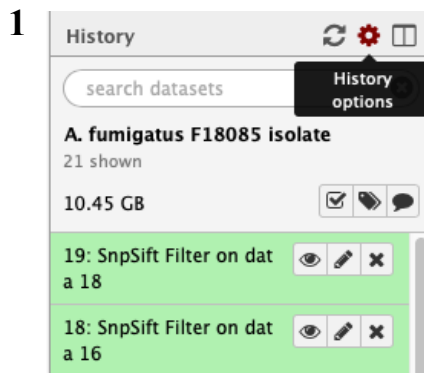
Summary

Genome	ToxoDB-29_TgondiiGT1_Genome
Date	2017-06-17 05:56
SnEff version	SnEff 4.11 (build 2015-10-03), by Pablo Cingolani
Command line arguments	SnEff -i vcf -o vcf -stats /scratch/galaxy/files/008/dataset_8107.dat ToxoDB-29_TgondiiGT1_Genome /scratch/galaxy/files/008/dataset_8105.dat
Warnings	3,941
Errors	0
Number of lines (input file)	8,411
Number of variants (before filter)	8,483
Number of not variants (i.e. reference equals alternative)	0
Number of variants processed (i.e. after filter and non-variants)	8,483
Number of known variants (i.e. non-empty ID)	0 (0%)
Number of multi-allelic VCF entries (i.e. more than two alleles)	72
Number of effects	14,149
Genome total length	63,945,332
Genome effective	

4. Share or publish workflow results (history).

There are three ways to share your history. Click on the on the “History options”, “Share or Publish” to access these options.

- Select “Make History Accessible via Link” if you would like to share a link with your colleague who also has an account with VEuPathDB.
- Select “Make History Accessible and Publish” if you would like to make the history visible in the Shared section of the VEuPathDB Galaxy.
- Click “Share with a user” to share the history directly with one or several users.



3 Share or Publish History 'A. fumigatus F18085 isolate'

Make History Accessible via Link and Publish It

This history is currently restricted so that only you and the users listed below can access it. You can:

[Make History Accessible via Link](#)

Generates a web link that you can share with other people so that they can view and import the history.

[Make History Accessible and Publish](#)

Makes the history accessible via link (see above) and publishes the history to Galaxy's Published Histories section, where it is publicly listed and searchable.

Share History with Individual Users

You have not shared this history with any users.

[Share with a user](#)