

Overview

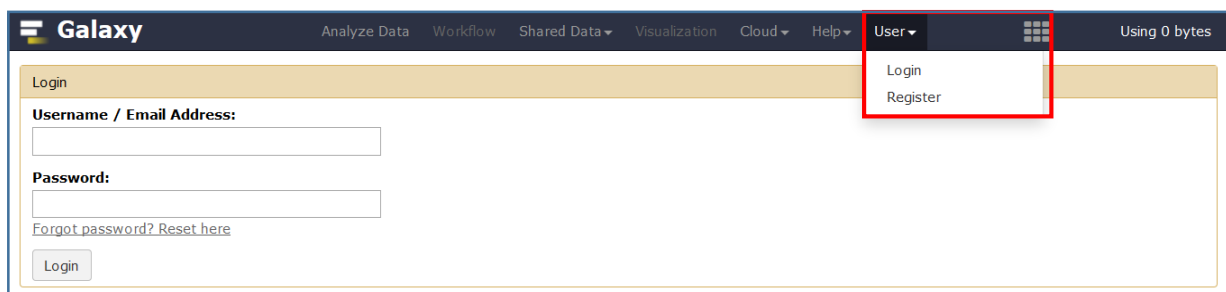
To generate a shortlist of potentially pathogenic variants we have filtered a variant call file (VCF) against public databases of known variation ([1000 genomes](#), [dbSNP](#), [Exome Sequencing Project](#)), estimates of deleteriousness such as [SIFT](#) and [PolyPhen2](#), relevant candidate genes based on the patients phenotype and mutation databases ([ClinVar](#) and [OMIM](#)). Identifying the causal mutation(s) is then a straightforward process when it occurs in a known disease-causing gene or the mutation has previously been described. In situations involving novel mutations and genes, additional filtering against related and/or unrelated individuals is often used to identify genes that are recurrently mutated in cases and are not mutated in controls. Comparison with parents may also be used to filter variants and to determine the mode of inheritance. The aim of this practical is to further reduce the number of putative causal variants by comparing VCF files between individuals. At the end of this exercise you will be able to:

1. Combine VCF files into one multisample file
2. Compare VCF files to identify shared variants and infer relatedness
2. Use GEMINI to analyse trio data and identify Mendelian errors and de-novo mutations
3. Calculate and plot alternate read percentages and use them to identify potential regions of acquired uniparental disomy (aUPD)
4. Filter VCF files using COSMIC and regions of aUPD to identify the underlying somatic driver mutation

Let's begin

For this session we will be using tools that are only available on the public version of galaxy (<https://usegalaxy.org/>).

1. Register for a public galaxy account and login



Identifying shared variants

In the first example whole exome sequencing was performed in two unrelated patients with intellectual disability, delayed speech and language development. The sequence data were analysed separately to produce two VCF files of exonic variants using BWA-MEM for alignment to hg19, samtools for variant calling and selection of variants with a minimum read depth of 4 and minimum QUAL score of phred 20. Our aim is to identify potentially causal variants that are present in both patients.

1. Rename the history for example 1 (Click on name, type 'Shared variants', press return)
2. Upload the two VCFs from the USB key in the 'FILTERING STRATEGIES' folder (**A1_exome.vcf** and **B1_exome.vcf**) to galaxy (<https://usegalaxy.org/>). Select VCF for type and Human Feb. 2009 (GRCh37/hg19)(hg19) for reference genome.

To identify shared variants we will intersect the two VCFs.

3. In **Tool Pane**: Go to **NGS: VCF Manipulation** > VCF-VCFintersect: Intersect two VCF datasets

The screenshot shows the Galaxy web interface for the VCF-VCFintersect tool. The 'Tools' pane on the left has 'NGS: VCF Manipulation' selected, with 'VCF-VCFintersect: Intersect two VCF datasets' listed below it. The main tool pane displays the configuration for 'VCF-VCFintersect: Intersect two VCF datasets (Galaxy Version 1.0.0_rc1.0)'. The first VCF dataset is '1: A1_exome.vcf' and the second is '2: B1_exome.vcf'. The reference genome is set to 'Human (Homo sapiens) (b37): hg_g1k_v37'. The 'Union or intersection' dropdown is set to 'Intersect'. The 'Execute' button is highlighted in red.

The result from VCF-VCF intersect only lists genotype calls present in sample B1. Repeat the intersection but this time compare sample B1 versus sample A1.

4. In **Tool Pane**: Go to **NGS: VCF Manipulation** > VCF-VCFintersect: Intersect two VCF datasets

The screenshot shows the Galaxy web interface for the VCF-VCFintersect tool. The 'Tools' pane on the left has 'text' selected. The main tool pane displays the configuration for 'VCF-VCFintersect: Intersect two VCF datasets (Galaxy Version 1.0.0_rc1.0)'. The first VCF dataset is '2: B1_exome.vcf' and the second is '1: A1_exome.vcf'. The reference genome is set to 'Human (Homo sapiens) (b37): hg_g1k_v37'. The 'Union or intersection' dropdown is set to 'Intersect'. The 'Execute' button is highlighted in red.

Now combine the two VCFs to make one VCF for variants shared between samples A1 and B1 and their genotype in both samples.

5. In **Tool Pane**: Go to **NGS: VCF Manipulation** > VCFcombine: Combine multiple VCF datasets

The screenshot shows the Galaxy web interface for the VCFcombine tool. The 'Tools' pane on the left has 'NGS: VCF Manipulation' selected, with 'VCFcombine: Combine multiple VCF datasets' listed below it. The main tool pane displays the configuration for 'VCFcombine: Combine multiple VCF datasets (Galaxy Version 1.0.0_rc1.0)'. The 'Select VCF Datasets' list contains '4: VCF-VCFintersect: on data 1 and data 2', '3: VCF-VCFintersect: on data 2 and data 1', '2: B1_exome.vcf', and '1: A1_exome.vcf'. The 'Execute' button is highlighted in red.

The proportion of variants with the same genotype between a pair of individuals offers a crude measure of relatedness although this will vary according to the sequencing approach (eg gene panel, whole exome, or whole genome) and capture kit (Agilent version 4 versus 5). Use text reformatting with awk to count the number of shared genotypes.

6. In **Tool Pane**: Go to **Text Manipulation** > **Text reformatting with awk**

Awk is a powerful programming language for working on large data files and is often applied to [NGS data](#). In VCF files, sample genotypes are recorded in a string of variables each separated by a colon. For example, 1/1:255,84,0:28:0:99 in the format GT:PL:DP:SP:GQ where GT is the genotype, PL is the phred scaled likelihoods, DP is the depth, SP is the phred-scaled strand bias P-value and GQ is the phred-scaled genotype quality. The awk command above splits this string into a list of separate variables for sample A1 (column 11) and sample B1 (column 10), counts every time the genotypes match and prints the number of matches when the end of the VCF file is reached.

Q1: What percentage of all variants in sample A1 have the same genotype in sample B1 and based on this value do you think these samples unrelated, at least to the 3rd degree, or not?

7. Save the combined VCF and use WANNONAR to annotate and filter the variants

8. Save the annotated exome summary results as a TXT file and the 9 filtered variants from step 11 from WANNVAR.

Basic Information			
exome summary results	view	CSV file	TXT file
genome summary results	view	CSV file	TXT file
Step11:9 variants	Remove variants found in cg46		download

Unfortunately the filtered results are not annotated and Galaxy is not very good at displaying files with numerous columns. We will therefore use Excel to open and compare the annotated and filtered results.

10. Open the annotated results (query.output.exome_summary.txt) in excel

Text Import Wizard - Step 1 of 3

The Text Wizard has determined that your data is Delimited.

If this is correct, choose Next, or choose the data type that best describes your data.

Original data type

Choose the file type that best describes your data:

☒ Delimited - Characters such as commas or tabs separate each field.

☐ Fixed width - Fields are aligned in columns with spaces between each field.

Start import at row: 1 File origin: MS-DOS (PC-8)

☐ My data has headers.

Preview of file C:\Users\wgh\Documents\MSK_GENOMICS\MODULE 7 BIOINFORMATL\query.output.exome_summary.txt.

1 ChrStartEndRefAltFunc.refgeneGene.refgeneGeneDetail.refgeneExonicFunc.refgeneAACh
2 877831877831TCexonicSAMD11nononymous SNVSAMD11:NM_152486:exon10:c.T1027C:p.W3
3 878314878314GCexonicSAMD11synonymous SNVSAMD11:NM_152486:exon11:c.G1440C:p.G480G
4 881627881627GAexonicNOC2Lsynonymous SNVNOC2L:NM_015658:exon16:c.C1843T:p.L615L0
5 888639888639TCexonicNOC2Lsynonymous SNVNOC2L:NM_015658:exon9:c.A918G:p.E304E0.9
6 888659888659TCexonicNOC2Lnononymous SNVNOC2L:NM_015658:exon9:c.A998G:p.I300V

Cancel < Back **Next >** Finish

Text Import Wizard - Step 2 of 3

This screen lets you set the delimiters your data contains. You can see how your text is affected in the preview below.

Delimiters

☒ Tab

☐ Semicolon

☐ Comma

☐ Space

☐ Other:

☐ Treat consecutive delimiters as one

Text qualifier: "

Data preview

Chr	Start	End	Ref	Alt	Func.refgene	Gene.refgene	GeneDetail.refgene	ExonicFunc.refgene
1	877831	877831	T	C	exonic	SAMD11	SAMD11	nononymous SNV
1	878314	878314	G	C	exonic	SAMD11	SAMD11	synonymous SNV
1	881627	881627	G	A	exonic	NOC2L	NOC2L	synonymous SNV
1	888639	888639	T	C	exonic	NOC2L	NOC2L	synonymous SNV
1	888659	888659	T	C	exonic	NOC2L	NOC2L	nononymous SNV

Cancel < Back **Next >** Finish

Text Import Wizard - Step 3 of 3

This screen lets you select each column and set the Data Format.

Column data format

☒ General

☐ Text

☐ Date: DMY

☐ Do not import column (skip)

'General' converts numeric values to numbers, date values to dates, and all remaining values to text.

Advanced...

Data preview

Chr	Start	End	Ref	Alt	Func.refgene	Gene.refgene	GeneDetail.refgene	ExonicFunc.refgene
1	877831	877831	T	C	exonic	SAMD11	SAMD11	nononymous SNV
1	878314	878314	G	C	exonic	SAMD11	SAMD11	synonymous SNV
1	881627	881627	G	A	exonic	NOC2L	NOC2L	synonymous SNV
1	888639	888639	T	C	exonic	NOC2L	NOC2L	synonymous SNV
1	888659	888659	T	C	exonic	NOC2L	NOC2L	nononymous SNV

Cancel < Back **Next >** **Finish**

11. Create a key to compare files by merging the chromosome, location, reference and alternate alleles.

	BU	BV	BW	BX	BY	BZ	CA	CB	CC	CD	CE	CF	CG
1	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo			
2	1	877831	.	T	C	222	.	AC1=2;AF1GT:PL:DP:1/1:255,8			=CONCATENATE(BU2,BV2,BX2,BY2)		

12. Mouse over the bottom right hand corner of the cell (CE2) and double click to auto fill column CE

	BU	BV	BW	BX	BY	BZ	CA	CB	CC	CD	CE
1	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	
2	1	877831	.	T	C	222	.	AC1=2;AF1GT:PL:DP:1/1:255,84	1877831T		
3	1	878314	.	G	C	185	.	AC1=1;AF1GT:PL:DP:0/1:215,0;242:41:3:99			

13. Open the filtered results file from wANNOVAR step 11 (filter.output.step11.vcf) in Excel as a tab delimited file, create the same key and auto fill column K

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	7	79842120	.	A	G	222	.	AC1=1;AF1GT:PL:DP:0/1:253,0			=CONCATENATE(A1,B1,D1,E1)		

14. Now use the keys to compare the annotated and filtered results.

i) Go to the annotated file (query.output.exome_summary.txt) and type “=vlookup(CE2,” in cell CF2

CF2													
	BU	BV	BW	BX	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	CE	CF
1	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo	Otherinfo		
2	1	877831	.	T	C	222	.	AC1=2;AF1GT:PL:DP:1/1:255,84	1877831T				

ii) For the table_array argument, switch to the filtered results file (filter.output.step11.vcf) and highlight column K which contains the key

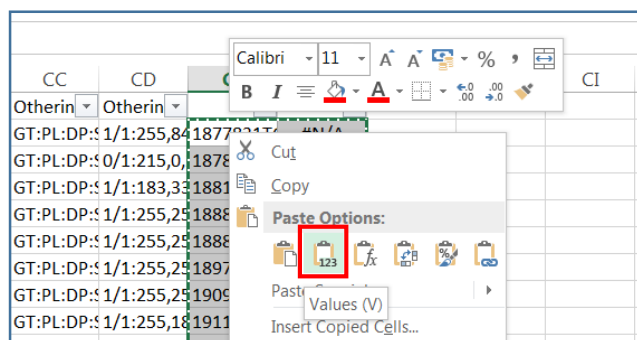
SUM													
	A	B	C	D	E	F	G	H	I	J	K	L	
1	7	79842120	.	A	G	222	.	AC1=1;AF1GT:PL:DP:0/1:253,0			779842120AG		
2	16	81242148	.	GTTT	GT	214	.	AC1=2;AF1GT:PL:DP:1/1:255,2			1681242148GTTTGT		
3	15	99511804	.	AC	ACC	217	.	AC1=1;AF1GT:PL:DP:0/1:255,0			1599511804ACACC		
4	11	1.03E+08	.	GT	GTT	126	.	AC1=2;AF1GT:PL:DP:1/1:167,2			11102738193GTGTT		

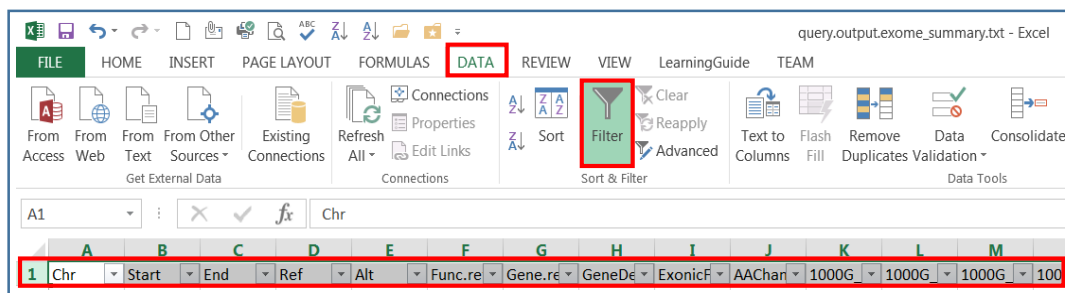
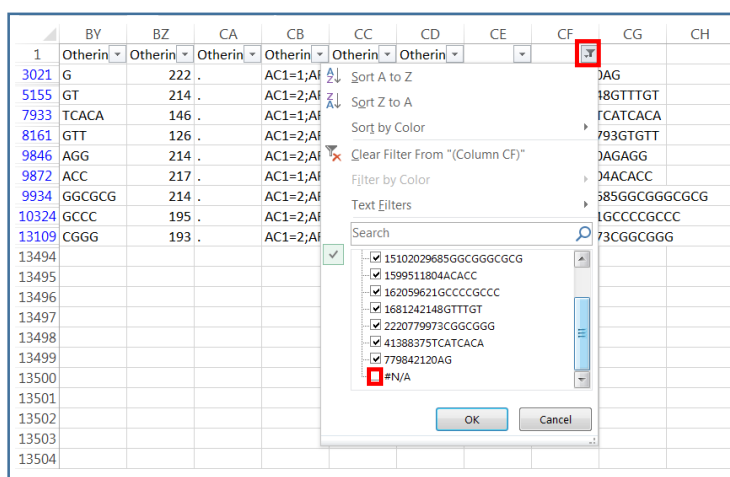
iii) Type “1” for the col_index_num, “FALSE” for [range_lookup] and press return then autofill column CF

SUM													
	A	B	C	D	E	F	G	H	I	J	K		
1	7	79842120	.	A	G	222	.	AC1=1;AF1GT:PL:DP:0/1:253,0			779842120AG		

The complete function means search column K for the value in CE2 and return the value in column K when an exact match is found.

17. Before filtering the annotated results, select cell CE2, then hold down Shift and Ctrl and use the arrow keys to highlight all cells in columns CE and CF. Once highlighted, use Ctrl c to copy the cells then right click and select paste values. This will replace all formulas with text which will speed up filtering.



18. Select the DATA tab, highlight row 1 and click Filter**19. Use the autofilters in columns BG and CF to select the 9 filtered variants and sort them in descending order of pathogenicity.**

The missense variant in GNAI1 (NM_001256414, p.K218R) is predicted to be the most deleterious shortlisted variant, although the others were not scored because they are either frameshifts or have unknown consequences. The variant is confidently called as heterozygous in both patients (GQ=99) and there are no QC issues recorded in the VCF file (high depth and no strand bias). In a recent study consisting of 7,580 whole exomes from patients with developmental disorders, de-novo GNAI1 mutations were identified in 8 patients including two with the same p.K218R mutation (Deciphering Developmental Disorders consortium 2017). Using this data, the authors demonstrated that GNAI1 was significantly enriched for deleterious de-novo mutations ($P < 5 \times 10^{-7}$) and concluded that it may be a novel cause of intellectual disability. Based on this and the observation that GNAI1 is mutated in both patients, we can conclude that this mutation is likely to cause the disease in patients A1 and B1.

Trio sequencing for de-novo mutations

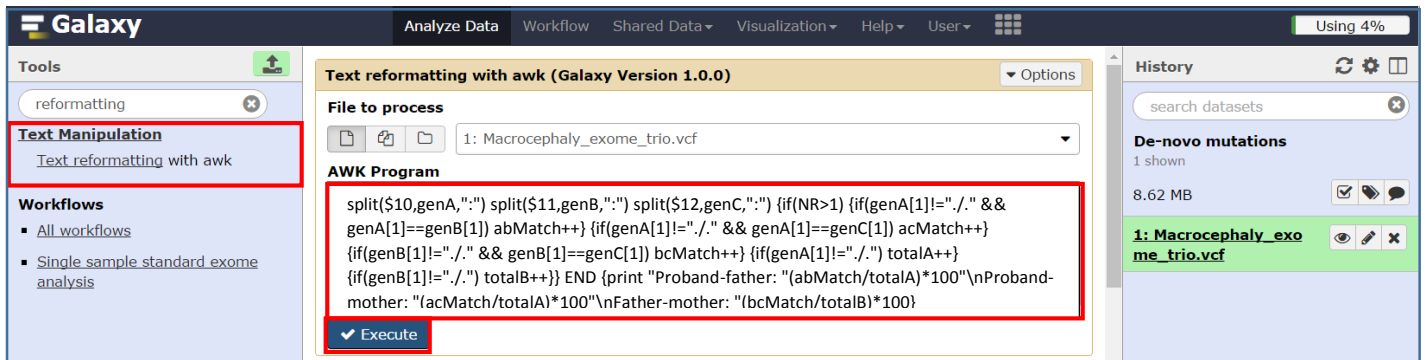
Whole exome trio sequencing of a child affected with a severe disorder and both unaffected biological parents is an effective strategy to determine diagnoses based on mutations in known genes and to generate evidence for the involvement of novel genes (Zhu et al 2015). The aim of the next example is to identify potentially causal de-novo mutations by comparing whole exome VCF files from an affected child and healthy parents. The child had macrocephaly at birth, capillary malformations on the nose and philtrum, cutis marmorata and polymicrogyria were identified by MRI. The VCF files were generated using BWA-MEM for alignment against hg19/build 37 and multi-sample GATK HaplotypeCaller for variant calling and selection of variants with a minimum QUAL score of phred 30.

1. Create a new history for example 2 (Click on name, type 'De-novo mutations', press return)

2. Upload the trio VCF from the USB key in the 'FILTERING STRATEGIES' folder (**Macrocephaly_exome_trio.vcf**) to galaxy (<https://usegalaxy.org/>), select VCF for type and Human Feb. 2009 (GRCh37/hg19)(hg19) for reference genome.

As before use text reformatting with awk to calculate the proportion of variants shared between all pairs of individuals as a simple check for relatedness.

3. In **Tool Pane**: Go to **Text Manipulation** > **Text reformatting with awk**



Text reformatting with awk (Galaxy Version 1.0.0)

File to process: 1: Macrocephaly_exome_trio.vcf

AWK Program:

```
split($10,genA,".") split($11,genB,".") split($12,genC,".") {if(NR>1) {if(genA[1]!="." && genA[1]==genB[1]) abMatch++} {if(genA[1]!="." && genA[1]==genC[1]) acMatch++} {if(genB[1]!="." && genB[1]==genC[1]) bcMatch++} {if(genA[1]!="." && genB[1]!="." && genC[1]!=".") totalA++} {if(genB[1]!="." && genC[1]!=".") totalB++} END {print "Proband-father: "(abMatch/totalA)*100"\nProband-mother: "(acMatch/totalA)*100"\nFather-mother: "(bcMatch/totalB)*100}}
```

Execute

View the results and answer these questions.

Q2: Do the percentages of shared variants agree with reported pedigree?

Q3: Is there any evidence to suggest that the parents are related?

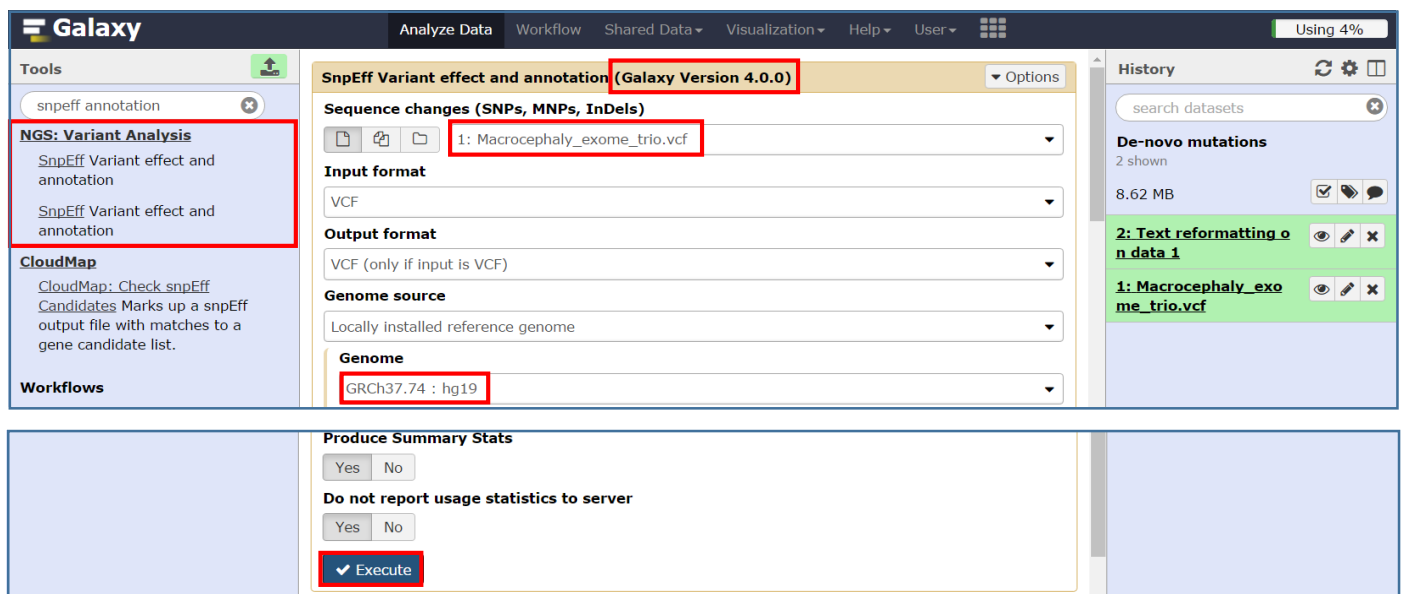
Q4: Compare the proportion of shared variants between unrelated individuals with the result from example 1 and give one reason which accounts for most of the difference.

To identify de-novo mutations we will use the **GEMINI** (GEnome MINing) software. GEMINI annotates VCF files with data from multiple sources (ENCODE tracks, UCSC tracks, OMIM, dbSNP, KEGG, and HPRD) and loads this information along with familial relationships into a SQLite database. The data can then be queried based on sample genotypes, inheritance patterns and the annotated fields.

To use GEMINI the first step is to load a VCF file and a pedigree file into the GEMINI database. The VCF file must be annotated using VEP or SnpEff before loading. We will therefore use SnpEff to annotate the VCF.

4. In **Tool Pane**: Go to **NGS: Variant Analysis** > **SnpEff Variant effect and annotation**

There are two versions of SnpEff available, make sure you choose version 4.0.0



SnpEff Variant effect and annotation (Galaxy Version 4.0.0)

Sequence changes (SNPs, MNPs, InDels): 1: Macrocephaly_exome_trio.vcf

Input format: VCF

Output format: VCF (only if input is VCF)

Genome source: Locally installed reference genome

Genome: GRCh37.74 : hg19

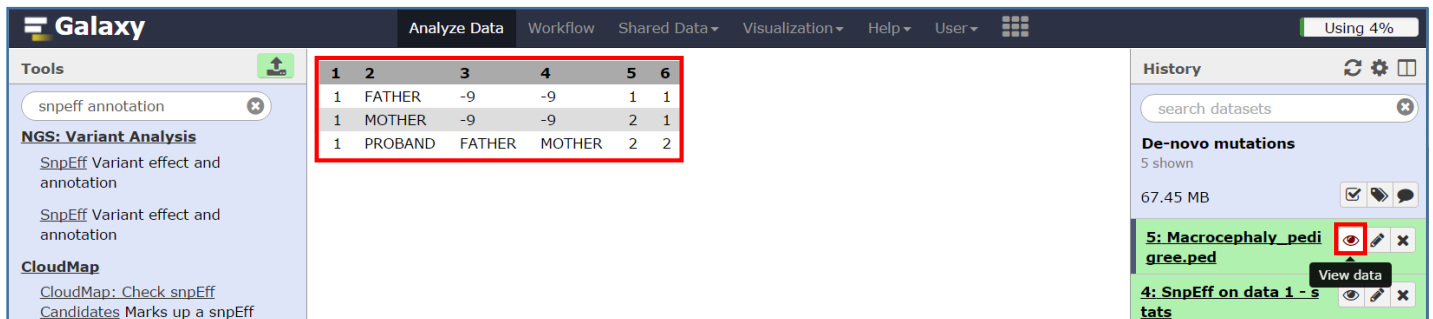
Produce Summary Stats: Yes

Do not report usage statistics to server: Yes

Execute

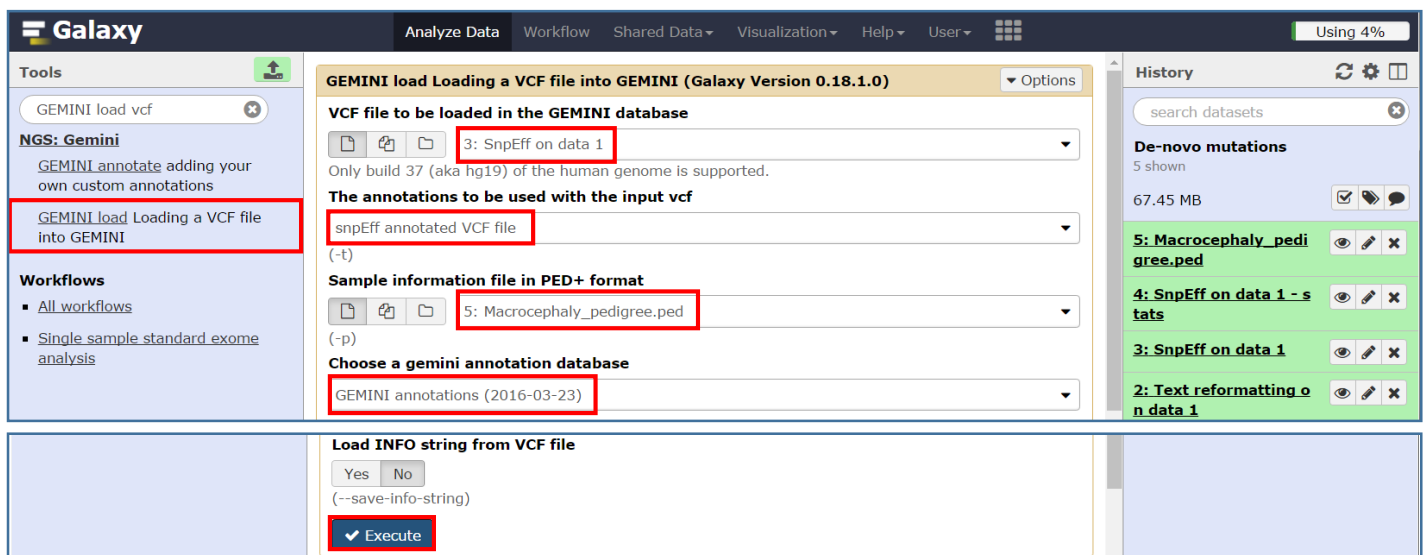
5. Now upload the trio pedigree from the USB key in the 'FILTERING STRATEGIES' folder (**Macrocephaly_pedigree.ped**) to galaxy (<https://usegalaxy.org/>), select tabular for type and Human Feb. 2009 (GRCh37/hg19)(hg19) for reference genome.

The pedigree file consists of six columns which describe the family_id, sample name, paternal_id, maternal_id, sex and phenotype where -9 is missing, 1=male, 2=female, 1=unaffected and 2=affected.

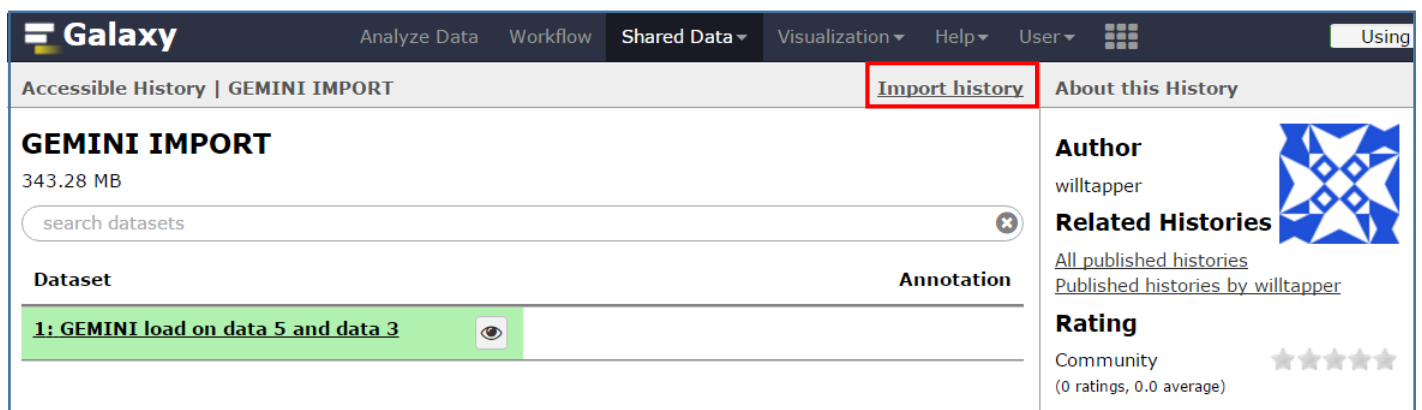


1	2	3	4	5	6
1	FATHER	-9	-9	1	1
1	MOTHER	-9	-9	2	1
1	PROBAND	FATHER	MOTHER	2	2

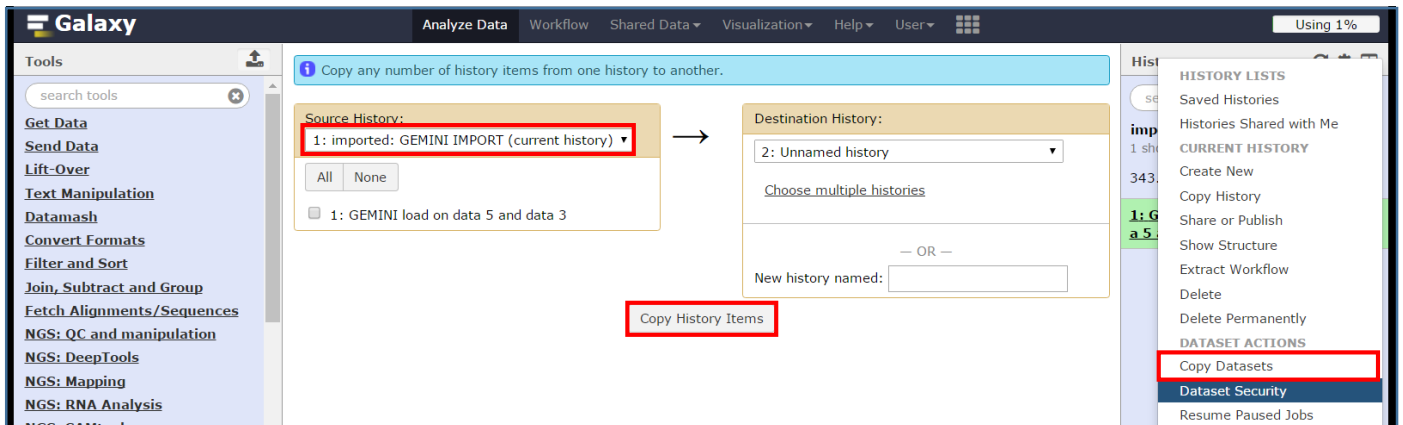
6. We're now ready to load the data into GEMINI using the options below. However, the loading step is computationally intensive and therefore very slow. To save time we will import a preloaded dataset so skip to step 7.



7. In a new browser tab enter this link <https://usegalaxy.org/u/willtapper/h/gemini-import> and click Import history. This will create a new history in your Galaxy account consisting of the GEMINI database for the macrocephaly trio.

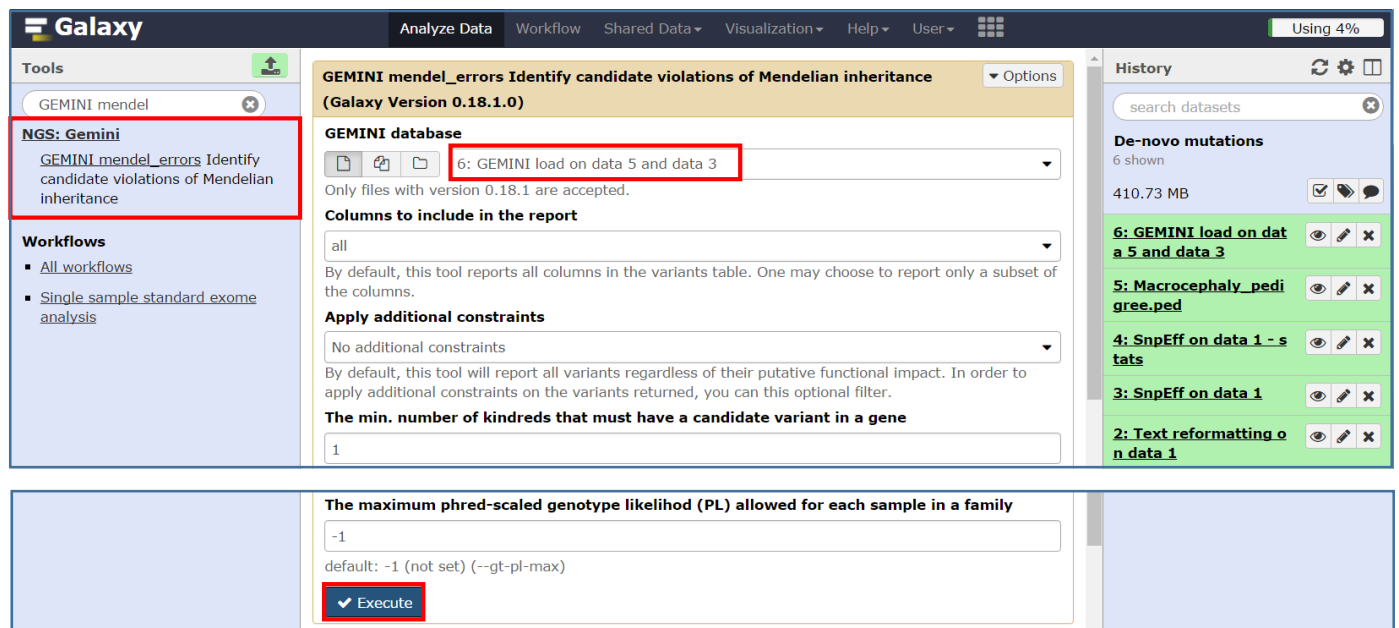


8. Click the cog icon in the history pane, select 'Copy Datasets', choose 'GEMINI IMPORT' as the source history and 'De-novo mutations' as the destination then click 'Copy History Items'.



We will now use GEMINI to identify all mendelian errors in the trio whereby an allele in the child could not have been received from either of its biological parents by Mendelian inheritance.

9. In **Tool Pane**: Go to **NGS: Gemini** > GEMINI mendel_errors Identify candidate violations of Mendelian Inheritance



Categories of Mendelian error:

Category	Description	Child	Parent A	Parent B
Loss of heterozygosity (LOH)	Child and one parent are opposite homozygotes and the other parent is heterozygous	AA	BB	AB
Uniparental disomy (UPD)	Child is homozygous and the parents are opposite homozygotes	AA	AA	BB
Implausible de-novo mutations	Child is homozygous and both parents are homozygous for the opposite genotype to child	AA	BB	BB
Plausible de-novo	Child is heterozygous and parents have the same homozygous genotype as each other	AB	AA	AA

Use text reformatting with awk to count the number of variants in each category of Mendelian errors.

10. In **Tool Pane**: Go to **Text Manipulation** > Text reformatting with awk

Galaxy Analyze Data Workflow Shared Data Visualization Help User Using 4%

Tools text reformatting

Text Manipulation
Text reformatting with awk

Workflows
All workflows
Single sample standard exome analysis

Text reformatting with awk (Galaxy Version 1.0.0) Options

File to process
7: GEMINI mendel_errors on data 6

AWK Program

```
{if($156=="implausible de novo") idn++; if($156=="plausible de novo") pdn++;  
if($156=="loss of heterozygosity") loh++; if($156=="uniparental disomy") upd++;}  
END {print "implausible de novo="idn" plausible de novo="pdn" loss of  
heterozygosity="loh" uniparental disomy="upd"}
```

Execute

History search datasets

De-novo mutations
7 shown
411.31 MB

7: GEMINI mendel_err
ors on data 6

6: GEMINI load on dat
a 5 and data 3

5: Macrocephaly_pedi
gree.ped

View the results and answer these questions.

Q5. How many variants are identified in each type of Mendelian error?

Q6. How does the number of de-novo mutations compare with expectation?

Based on the human mutation rate ($\sim 1.1 \times 10^{-8}$ per bp per generation) and genome size (3.1×10^9) there should be approximately 31 de-novo mutations per haploid genome and 0-2 mutations per exome since the exome is approximately 2% of the genome. Given the low mutation rate, the number of Mendelian errors can be used to estimate the genotyping error rate. Errors occur due to factors such as PCR artefacts, incorrect alignment and low read depth and are unevenly distributed throughout the genome. In a study by Patel et al 2014, approximately 95% of Mendelian errors were removed while retaining 80% of the called variants by applying filters based on read depth (DP<15 removed), genotype quality (GQ<20 removed) and alternate allele ratio (homozygous reference variants with alternate ratio >0.15 were removed, homozygous alternate variants with alternate ratio <0.85 were removed, and heterozygous variants with alternate ratio <0.3 or >0.7 were removed).

Use GEMINI to select all plausible de-novo mutations with: a maximum alternate allele frequency in public databases of ≤ 0.01 , an estimated impact severity of medium or high, a minimum read depth of 4 in all samples and a minimum genotype quality of 20 in all samples.

11. In **Tool Pane**: Go to **NGS: Gemini** > Gemini de_novo Identifying potential de novo mutations

Galaxy Analyze Data Workflow Shared Data Visualization Help User Using 4%

Tools gemini de novo

NGS: Gemini
GEMINI de_novo Identifying potential de novo mutations

Workflows
All workflows
Single sample standard exome analysis

GEMINI de_novo Identifying potential de novo mutations (Galaxy Version 0.18.1.0) Options

GEMINI database
6: GEMINI load on data 5 and data 3

Only files with version 0.18.1 are accepted.

Columns to include in the report
all

By default, this tool reports all columns in the variants table. One may choose to report only a subset of the columns.

Apply additional constraints
Apply additional constraints

By default, this tool will report all variants regardless of their putative functional impact. In order to apply additional constraints on the variants returned, you can this optional filter.

Constraints in SQL syntax
impact_severity != 'LOW' and max_aaf_all <= 0.01

Conditions applied here will become WHERE clauses in the query issued to the GEMINI database. E.g. alt='G' or impact_severity = 'HIGH'. (--filter)

The min. number of kindreds that must have a candidate variant in a gene
1

History search datasets

De-novo mutations
8 shown
411.31 MB

8: Text reformatting o
n data 7

7: GEMINI mendel_err
ors on data 6

6: GEMINI load on dat
a 5 and data 3

5: Macrocephaly_pedi
gree.ped

4: SnpEff on data 1 - s
tats

3: SnpEff on data 1

2: Text reformatting o
n data 1

<p>The minimum aligned sequence depth (genotype DP) required for each sample</p> <input type="text" value="4"/> <p>default: 0 (-d)</p> <p>the minimum genotype quality required for each sample in a family</p> <input type="text" value="20"/> <p>default: 0 (--min-gq)</p> <p>The maximum phred-scaled genotype likelihood (PL) allowed for each sample in a family</p> <input type="text" value="-1"/> <p>default: -1 (not set) (--gt-pl-max)</p> <p><input checked="" type="button" value="Execute"/></p>	<p>6: GEMINI load on data 5 and data 3</p> <p>5: Macrocephaly_pedigree.ped</p> <p>4: SnpEff on data 1 - stats</p> <p>3: SnpEff on data 1</p> <p>2: Text reformatting on data 1</p> <p>1: Macrocephaly_exome trio.vcf</p>
--	--

View the results and answer the question below.

Q7. How many plausible de-novo variants meet the search criteria?

A list of candidate has been generated using Phenomizer to assess the patients phenotype. Upload the

12. Upload the list of candidate genes from the USB key in the 'FILTERING STRATEGIES' folder (**candidate_gene_list.txt**) to galaxy, select tabular for type and leave the genome blank.

Join the list of plausible de-novo variants with the candidate gene list.

13. In **Tool Pane**: Go to **Text Manipulation** > **Join** two files

<p>Galaxy</p> <p>Tools</p> <p>join two files</p> <p>Text Manipulation</p> <p>Join two files</p> <p>Join, Subtract and Group</p> <p>Subtract Whole Dataset from another dataset</p> <p>Join two Datasets side by side on a specified field</p> <p>Compare two Datasets to find common or distinct rows</p> <p>NGS: QC and manipulation</p> <p>FASTQ joiner on paired end reads</p> <p>Operate on Genomic Intervals</p> <p>Join the intervals of two datasets side-by-side</p> <p>Workflows</p> <ul style="list-style-type: none"> All workflows Single sample standard exome analysis 	<p>Join two files (Galaxy Version 1.0.0)</p> <p>1st file</p> <p>9: GEMINI de_novo on data 6</p> <p>Column to use from 1st file</p> <p>Column: 58</p> <p>2nd File</p> <p>10: candidate_gene_list.txt</p> <p>Column to use from 2nd file</p> <p>Column: 1</p> <p>Output lines appearing in</p> <p>Both 1st & 2nd file.</p> <p>First line is a header line</p> <p>Yes No</p> <p>Use if first line contains column headers. It will not be sorted.</p> <p>Ignore case</p> <p>Yes No</p> <p>Sort and Join key column values regardless of upper/lower case letters.</p> <p>Value to put in unpaired (empty) fields</p> <p>0</p> <p><input checked="" type="button" value="Execute"/></p>	<p>History</p> <p>search datasets</p> <p>De-novo mutations</p> <p>10 shown</p> <p>411.38 MB</p> <p>10: candidate_gene_list.txt</p> <p>9: GEMINI de_novo on data 6</p> <p>8: Text reformatting on data 7</p> <p>7: GEMINI mendel errors on data 6</p> <p>6: GEMINI load on data 5 and data 3</p> <p>5: Macrocephaly_pedigree.ped</p> <p>4: SnpEff on data 1 - stats</p> <p>3: SnpEff on data 1</p> <p>2: Text reformatting</p>
--	--	--

Download the joined file and open it in Excel to view the results.

Q8. Which variant is most likely to cause the disease in the child?

Q9. What proportion of reads have the mutant allele (mutant allele frequency)?

Q10. Is there anything unusual about the mutant allele frequency and if so how could it be related to the disease?

Identifying genetic targets of acquired uniparental disomy

The concept of uniparental disomy (UPD) was introduced in the de-novo mutation section as a type of Mendelian error whereby the child is homozygous (AA) and the parents are opposite homozygotes (AA and BB). Inherited UPD can occur when a person receives both copies of a chromosome pair, or parts of chromosomes, from one parent. Constitutional UPD is associated with developmental disorders caused by the abnormal expression of imprinted genes, i.e. genes that are differentially expressed depending on whether they have been maternally or paternally inherited. For example, Prader-Willi syndrome and Angelman syndrome can be caused by UPD or other errors in imprinting involving genes on the long arm of chromosome 15.

By contrast, somatically acquired UPD (aUPD) in cancer is a mechanism by which a pre-existing driver mutation (usually somatically acquired) is converted to homozygosity, thereby providing an additional clonal advantage. aUPD may involve whole chromosomes as a result of non-disjunction or, more commonly, whole chromosome arms or terminal segments as a consequence of mitotic recombination. Regions of aUPD can be identified by NGS as long tracts of allelic imbalance (ie deviation from the expected percentage of alternate reads for heterozygous variants which should be approximately 50%) and many genes involved in a wide range of cancers have been identified by detecting minimal regions of recurrent aUPD and searching these regions for functionally relevant genes.

The aim of this example is to identify a region of aUPD and the underlying somatic driver mutation from whole exome sequencing of tumour DNA in a patient with polycythemia vera. The VCF file was generated using BWA-MEM for alignment to hg19, samtools for variant calling and selection of variants with a minimum read depth of 4 and minimum QUAL score of phred 20.

1. Create a new history for example 3 (Click on name, type 'aUPD and somatic driver', press return)
2. Upload the tumour VCF from the USB key in the 'FILTERING STRATEGIES' folder (**UPD_exome.vcf**) to galaxy (<https://usegalaxy.org/>), select VCF for type and Human Feb. 2009 (GRCh37/hg19)(hg19) for reference genome.

Use awk to determine the percentage of alternate reads for all variants on chromosome 1 with a read depth greater than 50.

3. In **Tool Pane**: Go to **Text Manipulation** > **Text reformatting with awk**

The screenshot displays the Galaxy web interface. On the left, the 'Tools' panel shows 'Text Manipulation' selected, with 'Text reformatting with awk' highlighted. The main workspace shows the tool configuration for 'Text reformatting with awk (Galaxy Version 1.0.0)'. The 'File to process' dropdown is set to '1: UPD_exome.vcf'. The 'AWK Program' field contains the following script:

```
split($8,info,",") {if($1 !~ /^#/) {if(info[1] !~ /^INDEL/ && $1==1) {split(info[5],dp,"="); split(dp[2],dp4,","); if(dp4[1]+dp4[2]+dp4[3]+dp4[4]>50) print $1,$2/1000000,((dp4[3]+dp4[4])/(dp4[1]+dp4[2]+dp4[3]+dp4[4]))*100;}} {if(info[1] ~ /^INDEL/ && $1==1) {split(info[6],dp,"="); split(dp[2],dp4,","); if(dp4[1]+dp4[2]+dp4[3]+dp4[4]>50) print $1,$2/1000000,((dp4[3]+dp4[4])/(dp4[1]+dp4[2]+dp4[3]+dp4[4]))*100;}}}
```

The 'Execute' button is highlighted with a red box. The right sidebar shows the 'History' panel with a search bar and a list of datasets, including 'aUPD and somatic driver' and '1: UPD_exome.vcf'.

The output produced by this awk code consists of three columns of data for chromosome, sequence location in megabases and percentage of alternate reads for each variant as rows.

The screenshot shows the Galaxy web interface. On the left, the 'Tools' panel is open, showing 'Text Manipulation' > 'Text reformatting with awk'. The main panel displays a table of genomic data with columns: Chrom, Pos, ID, Ref, Alt, Qual, Filter, Info, Format, data. The table contains 14 rows of data for chromosome 1. On the right, the 'History' panel shows a list of datasets, including 'aUPD and somatic driver' and '3: Text reformatting on data 1'.

Chrom	Pos	ID	Ref	Alt	Qual	Filter	Info	Format	data
1	0.016977	32.0755							
1	0.017538	23.4375							
1	0.069511	100							
1	0.808922	100							
1	0.808928	99.0566							
1	0.809178	30.5263							
1	0.897325	98.4848							
1	0.909768	100							
1	0.981087	61.4035							
1	0.982941	53.5211							
1	0.982994	51.2195							
1	1.246	98.5714							
1	1.24919	92.5926							

Repeat the text reformatting for chromosome 9 by changing the two instances of “\$1==1” to “\$1==9”.

4. In **Tool Pane**: Go to **Text Manipulation** > Text reformatting with awk

The screenshot shows the Galaxy web interface with the 'Text reformatting with awk' tool configuration. The 'File to process' field is set to '1: UPD_exome.vcf'. The 'AWK Program' field contains the following code:

```
split($8,info,",") {if($1 !~ /^#/){if(info[1] !~ /^INDEL/ && $1==9){split(info[5],dp,"="); split(dp[2],dp4,","); if(dp4[1]+dp4[2]+dp4[3]+dp4[4]>50) print $1,$2/1000000,((dp4[3]+dp4[4])/(dp4[1]+dp4[2]+dp4[3]+dp4[4]))*100;}} {if(info[1] ~ /^INDEL/ && $1==9){split(info[6],dp,"="); split(dp[2],dp4,","); if(dp4[1]+dp4[2]+dp4[3]+dp4[4]>50) print $1,$2/1000000,((dp4[3]+dp4[4])/(dp4[1]+dp4[2]+dp4[3]+dp4[4]))*100;}}}
```

The 'Execute' button is highlighted.

Make a scatter plot for chromosome 1 and chromosome 9 with the alternate read percentages on the Y-axis and megabase location on the X-axis.

5. In **Tool Pane**: Go to **Graph/Display Data** > Scatterplot of two numeric columns

The screenshot shows the Galaxy web interface with the 'Scatterplot of two numeric columns' tool configuration. The 'Dataset' field is set to '2: Text reformatting on data 1'. The 'Numerical column for x axis' is set to 'Column: 2'. The 'Numerical column for y axis' is set to 'Column: 3'. The 'Plot title' is set to 'Chromosome 1'. The 'Label for x axis' is set to 'Chromosome position (Mb)'. The 'Label for y axis' is set to 'Alternate read percentage'. The 'Execute' button is highlighted.

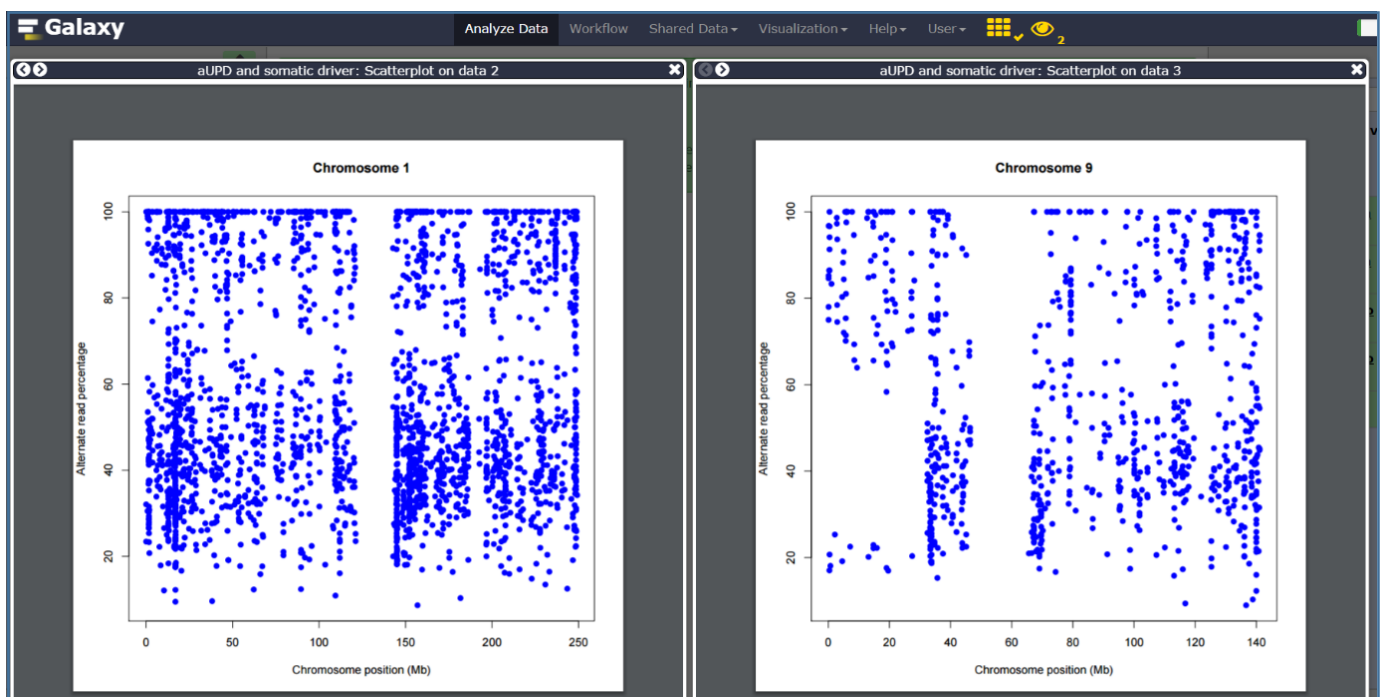
6. In **Tool Pane**: Go to **Graph/Display Data** > **Scatterplot** of two numeric columns

The screenshot shows the Galaxy web interface. In the 'Tools' pane on the left, the 'Graph/Display Data' section is selected, and the 'Scatterplot of two numeric columns' tool is chosen. The main tool pane shows the configuration for '3: Text reformatting on data'. The 'Dataset' is '3: Text reformatting on data'. The 'Numerical column for x axis' is 'Column: 2'. The 'Numerical column for y axis' is 'Column: 3'. The 'Plot title' is 'Chromosome 9'. The 'Label for x axis' is 'Chromosome position (Mb)'. The 'Label for y axis' is 'Alternate read percentage'. The 'Execute' button is highlighted.

Enable the Scratchbook and view both alternate read frequency plots side by side for comparison by clicking the view icon for steps 4 and 5.

The screenshot shows a success message in the center of the Galaxy interface: '1 job has been successfully added to the queue - resulting in the following datasets: 5: Scatterplot on data 3'. The message also states: 'You can check the status of queued jobs and view the resulting data by refreshing the History pane. When the job has been run the status will change from 'running' to 'finished' if completed successfully or 'error' if problems were encountered.' The 'Enable/Disable Scratchbook' button is highlighted. The 'History' pane on the right shows the job history.

The plots should look like this.



Q11. Does either chromosome have any evidence for aUPD and if so what are the approximate genomic coordinates in megabases?

Annotate the VCF file using ANNOVAR.

7. In Tool Pane: Go to **NGS: Variant Analysis** > ANNOVAR Annotate VCF with functional information using ANNOVAR

Use awk to filter the annotated VCF file. Replacing the question marks with coordinated for the potential aUPD region (\$1=="?" && \$2>?bp and \$2<?bp). Also select coding variants (\$6 == "exonic") that are either rare or absent from databases of known variation such as [1000 genomes](#) (\$12 <0.01 || \$12=="") && \$14=="") && (\$15<0.01 || \$15=="") && \$17!="") print \$0; if(NR==1) print \$0;}

8. In Tool Pane: Go to **NGS: Variant Analysis** > ANNOVAR Annotate VCF with functional information using ANNOVAR

View the filtered results and answer question 12.

Q12. Which variant is most likely to cause the clonal proliferation of erythrocytes (polycythemia vera) in this patient?

Well done Bioinformaticians you finished all exercises!