

**DATE: 2<sup>nd</sup> SEPTEMBER, 2024**

**NAME: HASSANA ABDULKADIR**

**EMAIL: abdulkadirhassana97@gmail.com**

**TITLE: SEARCHING FOR MUTATIONS LEADING TO ISONIAZID RESISTANCE IN TUBERCULOSIS**

**AIM: To identify and characterise genetic mutations associated with isoniazid resistance in *Mycobacterium tuberculosis*, the causative agent of TB.**

## **INTRODUCTION**

Tuberculosis (TB) remains a global health concern and a continuous threat to humanity. *Mycobacterium tuberculosis*, a large, nonmotile, rod-shaped, obligate aerobic bacterium is the pathogen that causes tuberculosis (TB). Commonly introduced into the body by inhaling droplet nuclei, which are present in the lungs' well-aerated upper lobes. In addition to the lungs, *Mycobacterium tuberculosis* can affect the liver, bones, genitourinary tract, central nervous system, lymph nodes, and gastrointestinal tract. Tuberculosis causes severe symptoms, such as a persistent cough, blood in the sputum, and weight loss. The illness claims the lives of almost 1.5 million people each year.

Those suffering from HIV/AIDS and other immune-compromised conditions are especially vulnerable to contracting tuberculosis. Four front-line medications are used to treat tuberculosis: pyrazinamide, ethambutol, rifampicin, and isoniazid.

The development of drug-resistant bacteria makes treatment considerably more difficult. An important first-line medication for tuberculosis treatment is isoniazid, a prodrug that inhibits the formation of the mycobacterial cell wall, and drug resistance to this medication can make treatment less effective.

KatG is a major gene associated with isoniazid resistance. The catalase-peroxidase enzyme, which KatG encodes, is in charge of changing isoniazid into its bactericidal form. Isoniazid resistance would arise from mutations in the KatG gene, thus decreasing the enzyme's efficiency.

The aim of this project is to identify and characterise genetic mutations associated with isoniazid resistance in *Mycobacterium tuberculosis*, the causative agent of Tuberculosis by performing variant calling and also explore associated biological pathways to understand the molecular mechanisms of isoniazid resistance.

Variant calling is a very popular analysis pipeline. Variant calling is a kind of analysis done on whole genome or whole exom sequencing data, basically to determine the variants present in the sequenced reads. Variations are usually caused by mutations.

## METHOD

The workflow for this project includes;

Sequenced reads → Trim reads (Sickle) → Index Reference Genome → Mapping sequence reads to downloaded referenced genome (bwa-mem) → SAM files → BAM files → Sorted BAM files → Variant calling(bcftools) → Variant annotation (SnpEff)

The Sequence reads were downloaded and trimmed using **sickle** (as a form of quality control). The trimmed reads were mapped or aligned to a Reference genome using BWA-Burrows-wheeler Aligner and the information was stored in a SAM file. After mapping, the SAM files were converted to BAM files and the BAM files were sorted.

VCF tools such as **mpileup** command was used to generate a coverage information for all the bases and the variance was called up using bcftools and then finally a post vcf analysis and Variant annotation was performed using SnpEff. The Variant calling algorithm applied, basically calls variance, identifies these variants and store them in a file called Variant call format file (VCF). The VCF (Variant call format) file generated stores variants present in the data. SnpEff is used to annotate variants by providing functional information about them.

## STEP 1

- The softwares necessary for this project were already installed via conda and homebrew. Dedicated directories were created for the different file types.
- The sample datasets (a total of 19), usually stored in fasta or fastq format were downloaded and the reference genome was downloaded from the NCBI database.

## STEP 2

- The sample datasets were trimmed using sickle installed via homebrew. Trimming helps to, improve the quality of the reads.

#Trimm reads

```
sickle pe -f ERR8774458_1.fastq -r ERR8774458_2.fastq -t sanger -q 20 -l 20 -g -o trimmed_R1.fastq  
-p trimmed_R2.fastq -s trimmed_S.fastq
```

```
sickle pe -f ERR8774464_1.fastq -r ERR8774464_2.fastq -t sanger -q 20 -l 20 -g -o  
trimmed2_R1.fastq -p trimmed2_R2.fastq -s trimmed2_S.fastq
```

```
sickle pe -f ERR8774480_1.fastq -r ERR8774480_2.fastq -t sanger -q 20 -l 20 -g -o  
trimmed3_R1.fastq -p trimmed3_R2.fastq -s trimmed3_S.fastq
```

```
sickle pe -f ERR8774482_1.fastq -r ERR8774482_2.fastq -t sanger -q 20 -l 20 -g -o  
trimmed4_R1.fastq -p trimmed4_R2.fastq -s trimmed4_S.fastq
```

```
sickle se -t sanger -q 20 -l 20 -g -f ERR8774487_1.fastq -o trimmed5_R1.fastq
```

```
sickle pe -f ERR8774512_1.fastq -r ERR8774512_2.fastq -t sanger -q 20 -l 20 -g -o  
trimmed6_R1.fastq -p trimmmmed6_R2.fastq -s trimmed6_S.fastq
```

```
sickle pe -f ERR8774514_1.fastq -r ERR8774514_2.fastq -t sanger -q 20 -l 20 -g -o  
trimmed7_R1.fastq -p trimmed7_R2.fastq -s trimmed7_S.fastq
```

```
sickle pe -f ERR8774522_1.fastq -r ERR8774522_2.fastq -t sanger -q 20 -l 20 -g -o  
trimmed8_R1.fastq -p trimmed8_R2.fastq -s trimmed8_S.fastq
```

```
sickle pe -f ERR8774524_1.fastq -r ERR8774524_2.fastq -t sanger -q 20 -l 20 -g -o  
trimmed9_R1.fastq -p trimmed9_R2.fastq -s trimmed9_S.fastq
```

```
sickle pe -f ERR8774534_1.fastq -r ERR8774534_2.fastq -t sanger -q 20 -l 20 -g -o  
trimmed10_R1.fastq -p trimmed10_R2.fastq -s trimmed10_S.fastq
```

## STEP 3

- The sample datasets were aligned to the downloaded reference genome. This is referred to as Genome mapping and is usually done using the BWA tool (Burrows-Wheeler Aligner). The purpose of genome mapping is to identify regions or locations

in the reference genome where the reads map uniquely to or corresponds to. The output generated by bwa-mem were stored in SAM files. The reference genome was first indexed before alignment.

*#Index the reference genome*

```
bwa index sequence.fasta
```

*#Perform genome mapping*

```
bwa mem -t 8 ref/sequence.fasta trimmed_R1.fastq trimmed_R2.fastq > output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed2_R1.fastq trimmed2_R2.fastq > 2output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed3_R1.fastq trimmed3_R2.fastq > 3output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed4_R1.fastq trimmed4_R2.fastq > 4output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed5_R1.fastq > 5output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed6_R1.fastq trimmed6_R2.fastq > 6output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed7_R1.fastq trimmed7_R2.fastq > 7output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed8_R1.fastq trimmed8_R2.fastq > 8output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed9_R1.fastq trimmed9_R2.fastq > 9output.sam  
bwa mem -t 8 ref/sequence.fasta trimmed10_R1.fastq trimmed10_R2.fastq > 10output.sam
```

## STEP 4

- The SAM files were converted into BAM files using **samtools** as they are usually large. BAM files are compressed and usually have smaller size.

*#converting SAM files to BAM files using samtools*

```
samtools view -S -b output.sam > output.bam  
samtools view -S -b 2output.sam > 2output.bam  
samtools view -S -b 3output.sam > 3output.bam  
samtools view -S -b 4output.sam > 4output.bam  
samtools view -S -b 5output.sam > 5output.bam  
samtools view -S -b 6output.sam > 6output.bam  
samtools view -S -b 7output.sam > 7output.bam  
samtools view -S -b 8output.sam > 8output.bam  
samtools view -S -b 9output.sam > 9output.bam
```

```
 samtools view -S -b 10output.sam > 10output.bam
```

## STEP 5

- The bam files were sorted using samtools sort, which was used for the variant calling.

*#sorting the bam files using samtools sort*

```
 samtools sort output.bam -o../sort/output.sort.bam  
 samtools sort 2output.bam -o 2output.sort.bam  
 samtools sort 3output.bam -o 3output.sort.bam  
 samtools sort 4output.bam -o 4output.sort.bam  
 samtools sort 5output.bam -o 5output.sort.bam  
 samtools sort 6output.bam -o 6output.sort.bam  
 samtools sort 7output.bam -o 7output.sort.bam  
 samtools sort 8output.bam -o 8output.sort.bam  
 samtools sort 9output.bam -o 9output.sort.bam  
 samtools sort 10output.bam -o 10output.sort.bam
```

## STEP 6

- The sorted bam files were indexed

*#index sorted bam files*

```
 samtools index output.sort.bam  
 samtools index 2output.sort.bam  
 samtools index 3output.sort.bam  
 samtools index 4output.sort.bam  
 samtools index 5output.sort.bam  
 samtools index 6output.sort.bam  
 samtools index 7output.sort.bam  
 samtools index 8output.sort.bam  
 samtools index 9output.sort.bam  
 samtools index 10output.sort.bam
```

## STEP 7

- The statistics was checked using flagstat on one of the reads to check how well the genome mapped;

#To view the statistics or output of the reads

-samtools flagstat output.sort.bam

793440 + 0 in total (QC-passed reads + QC-failed reads)

783766 + 0 primary

0 + 0 secondary

9674 + 0 supplementary

0 + 0 duplicates

0 + 0 primary duplicates

786124 + 0 mapped (99.08% : N/A)

776450 + 0 primary mapped (99.07% : N/A)

783766 + 0 paired in sequencing

391883 + 0 read1

391883 + 0 read2

754788 + 0 properly paired (96.30% : N/A)

774778 + 0 with itself and mate mapped

1672 + 0 singletons (0.21% : N/A)

0 + 0 with mate mapped to a different chr

0 + 0 with mate mapped to a different chr (mapQ>=5)

**This means that 99.08% of the reads for sample 1 mapped to the reference genome.**

**Majority of the reads mapped**

-samtools flagstat 2output.sort.bam

1426249 + 0 in total (QC-passed reads + QC-failed reads)

1411836 + 0 primary

0 + 0 secondary

14413 + 0 supplementary

0 + 0 duplicates

0 + 0 primary duplicates

1417810 + 0 mapped (99.41% : N/A)  
1403397 + 0 primary mapped (99.40% : N/A)  
1411836 + 0 paired in sequencing  
705918 + 0 read1  
705918 + 0 read2  
1370854 + 0 properly paired (97.10% : N/A)  
1400382 + 0 with itself and mate mapped  
3015 + 0 singletons (0.21% : N/A)  
0 + 0 with mate mapped to a different chr  
0 + 0 with mate mapped to a different chr (mapQ>=5)

## STEP 8

### VARIANT CALLING

VCF tools such as **mpileup** command was used to generate a coverage information for all the bases and the variance was called up using bcftools.

Bcftools is used to do the mpileup and this generates coverage information about the mapped bases. The mpileup sub-command of bcftools is used to count the read coverage at each sequence position.

```
bcftools mpileup -O b -o raw.bcf -f ref/sequence.fasta output.sort.bam
bcftools mpileup -O b -o 2raw.bcf -f ref/sequence.fasta 2output.sort.bam
bcftools mpileup -O b -o 3raw.bcf -f ref/sequence.fasta 3output.sort.bam
bcftools mpileup -O b -o 4raw.bcf -f ref/sequence.fasta 4output.sort.bam
bcftools mpileup -O b -o 5raw.bcf -f ref/sequence.fasta 5output.sort.bam
bcftools mpileup -O b -o 6raw.bcf -f ref/sequence.fasta 6output.sort.bam
bcftools mpileup -O b -o 7raw.bcf -f ref/sequence.fasta 7output.sort.bam
bcftools mpileup -O b -o 8raw.bcf -f ref/sequence.fasta 8output.sort.bam
bcftools mpileup -O b -o 9raw.bcf -f ref/sequence.fasta 9output.sort.bam
bcftools mpileup -O b -o 10raw.bcf -f ref/sequence.fasta 10output.sort.bam
```

#Call single nucleotide variants

```
bcftools call --ploidy 1 -m -v -o raw.vcf raw.bcf
bcftools call --ploidy 1 -m -v -o 2raw.vcf 2raw.bcf
```

```
bcftools call --ploidy 1 -m -v -o 3raw.vcf 3raw.bcf  
bcftools call --ploidy 1 -m -v -o 4raw.vcf 4raw.bcf  
bcftools call --ploidy 1 -m -v -o 5raw.vcf 5raw.bcf  
bcftools call --ploidy 1 -m -v -o 6raw.vcf 6raw.bcf  
bcftools call --ploidy 1 -m -v -o 7raw.vcf 7raw.bcf  
bcftools call --ploidy 1 -m -v -o 8raw.vcf 8raw.bcf  
bcftools call --ploidy 1 -m -v -o 9raw.vcf 9raw.bcf  
bcftools call --ploidy 1 -m -v -o 10raw.vcf 10raw.bcf
```

*#Filter the snvs for the final output in vcf format using vcftools*

```
vcfutils.pl varFilter raw.vcf > raw.filter.vcf  
vcfutils.pl varFilter 2raw.vcf > 2raw.filter.vcf  
vcfutils.pl varFilter 3raw.vcf > 3raw.filter.vcf  
vcfutils.pl varFilter 4raw.vcf > 4raw.filter.vcf  
vcfutils.pl varFilter 5raw.vcf > 5raw.filter.vcf  
vcfutils.pl varFilter 6raw.vcf > 6raw.filter.vcf  
vcfutils.pl varFilter 7raw.vcf > 7raw.filter.vcf  
vcfutils.pl varFilter 8raw.vcf > 8raw.filter.vcf  
vcfutils.pl varFilter 9raw.vcf > 9raw.filter.vcf  
vcfutils.pl varFilter 10raw.vcf > 10raw.filter.vcf
```

*#To check the number of snps and indels (variant types); the number of variants that are snps and indels*

```
grep -v -c '^#' raw.vcf  
1047  
bcftools view -v snps raw.vcf | grep -v -c '^#'  
956  
bcftools view -v indels raw.vcf | grep -v -c '^#'  
91  
  
grep -v -c '^#' 2raw.vcf
```

1025

bcftools view -v snps 2raw.vcf | grep -v -c '^#'

940

bcftools view -v indels 2raw.vcf | grep -v -c '^#'

85

grep -v -c '^#' 3raw.vcf

1067

bcftools view -v snps 3raw.vcf | grep -v -c '^#'

983

bcftools view -v indels 3raw.vcf | grep -v -c '^#'

84

grep -v -c '^#' 4raw.vcf

1049

bcftools view -v snps 4raw.vcf | grep -v -c '^#'

964

bcftools view -v indels 4raw.vcf | grep -v -c '^#'

85

grep -v -c '^#' 5raw.vcf

951

bcftools view -v snps 5raw.vcf | grep -v -c '^#'

870

bcftools view -v indels 5raw.vcf | grep -v -c '^#'

81

grep -v -c '^#' 6raw.vcf

1027

bcftools view -v snps 6raw.vcf | grep -v -c '^#'

937

bcftools view -v indels 6raw.vcf | grep -v -c '^#'

90

```
grep -v -c '^#' 7raw.vcf  
1031  
bcftools view -v snps 7raw.vcf | grep -v -c '^#'  
941  
bcftools view -v indels 7raw.vcf | grep -v -c '^#'  
90
```

```
grep -v -c '^#' 8raw.vcf  
1003  
bcftools view -v snps 8raw.vcf | grep -v -c '^#'  
918  
bcftools view -v indels 8raw.vcf | grep -v -c '^#'  
85
```

```
grep -v -c '^#' 9raw.vcf  
1027  
bcftools view -v snps 9raw.vcf | grep -v -c '^#'  
938  
bcftools view -v indels 9raw.vcf | grep -v -c '^#'  
89
```

```
grep -v -c '^#' 10raw.vcf  
1019  
bcftools view -v snps 10raw.vcf | grep -v -c '^#'  
932  
bcftools view -v indels 10raw.vcf | grep -v -c '^#'  
87
```

*#final filtered datasets*

```
grep -v "##" raw.vcf  
grep -v "##" 2raw.vcf  
grep -v "##" 3raw.vcf
```

```
grep -v "##" 4raw.vcf  
grep -v "##" 5raw.vcf  
grep -v "##" 6raw.vcf  
grep -v "##" 7raw.vcf  
grep -v "##" 8raw.vcf  
grep -v "##" 9raw.vcf  
grep -v "##" 10raw.vcf
```

*#To view results*

```
samtools tview output.sort.bam ref/sequence.fasta  
samtools tview 2output.sort.bam ref/sequence.fasta  
samtools tview 3output.sort.bam ref/sequence.fasta  
samtools tview 4output.sort.bam ref/sequence.fasta  
samtools tview 5output.sort.bam ref/sequence.fasta  
samtools tview 6output.sort.bam ref/sequence.fasta  
samtools tview 7output.sort.bam ref/sequence.fasta  
samtools tview 8output.sort.bam ref/sequence.fasta  
samtools tview 9output.sort.bam ref/sequence.fasta  
samtools tview 10output.sort.bam ref/sequence.fasta
```

## STEP 9

### VARIANT ANNOTATION

Variant annotation provides more information to the variants/mutations called from the sequence data. Synonymous or nonsynonymous substitution, start-gain codon, start-loss codon, stop-gain codon, stop-loss codon, or frameshifts are some of the variant effects that are predicted using SnpEff.

*#build database*

```
scripts/buildDbNcbi.sh NC_000962.3 Downloading genome NC_000962.3
```

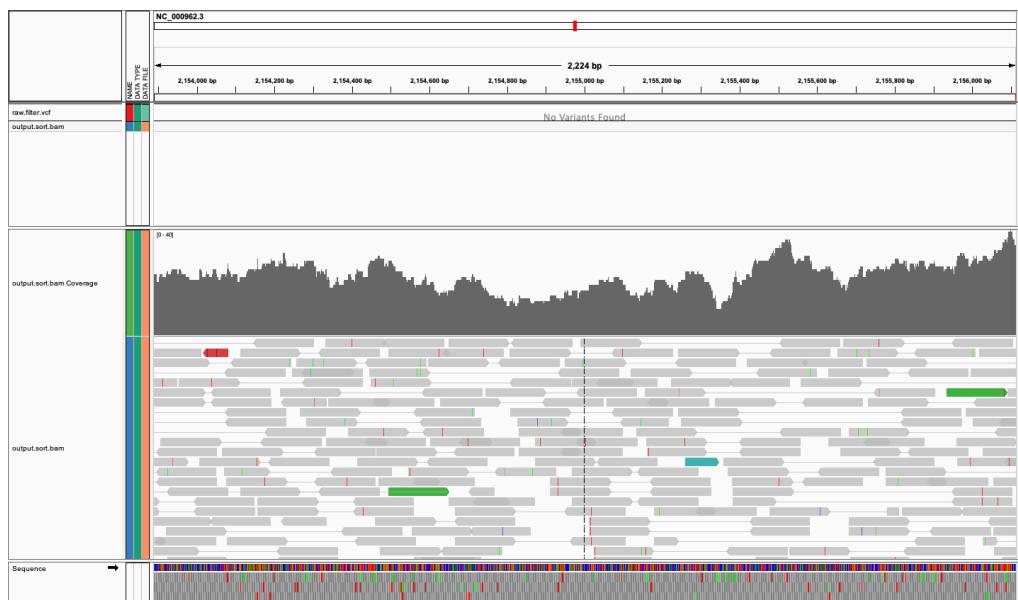
*#download database*

```
java -jar.snpEff.jar download -v NC_000962.3
```

#using snpEff to perform variant annotation

```
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 raw.filter.vcf > mmySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 2raw.filter.vcf > mySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 3raw.filter.vcf > 3mySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 4raw.filter.vcf > 4mySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 5raw.filter.vcf > 5mySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 6raw.filter.vcf > 6mySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 7raw.filter.vcf > 7mySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 8raw.filter.vcf > 8mySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 9raw.filter.vcf > 9mySNPanot.vcf  
java -Xmx8g -jar snpEff/snpEff.jar NC_000962.3 10raw.filter.vcf > 10mySNPanot.vcf
```

## RESULTS



**Fig 1.0** IGV visualisation of a *Mycobacterium tuberculosis* H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 1.

A	B	C	D	E	F	G	H
REF	ALT	QUAL	FILTER	INFO			
34 NC_00962.1	1977.	A	G	225.417.	DP=35;VDB=0.5476075;GB=-0.693127;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,24;MQ=60;ANN=G upstream_gene_variant MODIFIER dnAN Rv0002 transcript Rv0002 protein_coding		
35 NC_00962.1	4013.	T	C	225.417.	DP=34;VDB=0.8290175;GB=-0.693127;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,12,20;MQ=60;ANN=C missense_variant MODERATE ref Rv0003 transcript Rv0003 protein_coding  1/1		
36 NC_00962.1	7363.	G	C	225.417.	DP=24;VDB=0.05470075;GB=0.692831;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,14,10;MQ=60;ANN=C missense_variant MODERATE gtAA Rv0006 transcript Rv0006 protein_coding  1/1		
37 NC_00962.1	7585.	G	C	225.417.	DP=28;VDB=0.0546408;GB=0.692831;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,17,15;MQ=60;ANN=C missense_variant MODERATE gtAA Rv0006 transcript Rv0006 protein_coding  1/1 C		
38 NC_00962.1	9504.	G	A	225.417.	DP=28;VDB=0.0546408;GB=0.692831;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,17,15;MQ=60;ANN=C missense_variant MODERATE gtAA Rv0006 transcript Rv0006 protein_coding  1/1 C		
39 NC_00962.1	11746.	T	G	225.417.	DP=28;VDB=0.0546408;GB=0.692831;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,7,11;MQ=60;ANN=G upstream_gene_variant MODIFIER polA Rv0009 transcript Rv0009 protein_coding		
40 NC_00962.1	11879.	A	G	225.417.	DP=21;VDB=0.462325;GB=0.630207;MQ\$BZ=0.0MQF=0AC-1AN=1;DPA=0.3,17,3;MQ=60;ANN=C missense_variant MODERATE Rv0008 transcript Rv0008 protein_coding  1/1		
41 NC_00962.1	14785.	T	C	225.417.	DP=34;VDB=0.2127945;GB=-0.693127;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,17,16;MQ=60;ANN=C missense_variant MODERATE Rv0012 Rv0121 transcript Rv0012 protein_coding  1/1		
42 NC_00962.1	18091.	G	A	225.417.	DP=48;VDB=0.8204085;GB=-0.693147;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,7,9;MQ=60;ANN=A synonymous_variant LOW pkaN Rv0015c transcript Rv0015c protein_coding  1/1 C		
43 NC_00962.1	21795.	G	A	86.415.	DP=4;VDB=0.4218035;GB=-0.556411;MQ\$BZ=0.056411;MQF=0AC-1AN=1;DPA=0.3,7,21;MQ=60;ANN=G upstream_gene_variant MODIFIER polA Rv0009 transcript Rv0009 protein_coding  1/1 C		
44 NC_00962.1	23854.	A	G	225.417.	DP=38;VDB=0.020201525;GB=-0.693139;MQ\$BZ=0.056411;MQF=0AC-1AN=1;DPA=0.3,22,7;MQ=60;ANN=G upstream_gene_variant MODIFIER phpA Rv0016c transcript Rv0016c protein_coding		
45 NC_00962.1	26959.	C	G	225.417.	DP=38;VDB=0.020201525;GB=-0.693139;MQ\$BZ=0.056411;MQF=0AC-1AN=1;DPA=0.3,22,7;MQ=60;ANN=G upstream_gene_variant MODIFIER phpA Rv0016c transcript Rv0016c protein_coding		
46 NC_00962.1	27518.	TAAAAAA	TAAAAAA	228.411.	INDEL;ID=23;MF=0.058333;OP=24;VDB=0.25565;GB=-0.692717;RPBZ=1.66531;MQ\$BZ=0.056411;MQF=0AC-1AN=1;DPA=1,0,7,16;MQ=60;ANN=TAAC		
47 NC_00962.1	27918.	G	A	225.417.	DP=17;VDB=0.9152485;GB=-0.686358;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,9,5;MQ=60;ANN=TA synonymous_variant LOW Rv0023 Rv0023 protein_coding  1/1 C 32		
48 NC_00962.1	32387.	C	T	225.417.	DP=17;VDB=0.9152485;GB=-0.686358;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,9,8;MQ=60;ANN=TA missense_variant MODERATE Rv0029 Rv0029 protein_coding  1/1		
49 NC_00962.1	33817.	C	G	225.417.	DP=43;VDB=0.2076565;GB=-0.690438;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,22,21;MQ=60;ANN=G upstream_gene_variant MODIFIER lof2 Rv0032 transcript Rv0032 protein_coding		
50 NC_00962.1	34044.	T	C	225.417.	DP=24;VDB=0.3327795;GB=-0.692717;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,16,7;MQ=60;ANN=C upstream_gene_variant MODIFIER biot2 Rv0032 transcript Rv0032 protein_coding		
51 NC_00962.1	37411.	C	G	225.417.	DP=24;VDB=0.3327795;GB=-0.692717;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,16,7;MQ=60;ANN=C upstream_gene_variant MODIFIER biot2 Rv0032 transcript Rv0032 protein_coding		
52 NC_00962.1	42967.	G	C	225.417.	DP=29;VDB=0.0778747;GB=0.6893079;MQ\$BZ=-0.0MQF=0AC-1AN=1;DPA=0.3,22,7;MQ=60;ANN=G upstream_gene_variant MODIFIER polA Rv0128c transcript Rv0128c protein_coding		
53 NC_00962.1	47138.	CTT	CTT	228.385.	INDEL;ID=23;MF=0.058333;OP=24;VDB=0.25565;GB=-0.692717;RPBZ=1.87867;CBZ=-4.79583;MQ\$BZ=0AC-1AN=1;DPA=1,0,7,16;MQ=60;ANN=TAAC		
54 NC_00962.1	49079.	C	G	225.417.	DP=32;VDB=0.0917515;GB=-0.69312;MQ\$BZ=0.056411;MQF=0AC-1AN=1;DPA=0.3,22,12;MQ=60;ANN=G synonymous_variant LOW Rv0045 Rv0045 protein_coding  1/1		
55 NC_00962.1	49233.	AGGG	AGG	225.417.	INDEL;ID=23;MF=0.0917515;GB=-0.69312;MQ\$BZ=0.056411;MQF=0AC-1AN=1;DPA=0.3,22,12;MQ=60;ANN=G synonymous_variant LOW Rv0045 Rv0045 protein_coding  1/1		
56 NC_00962.1	54303.	C	T	228.374.	DP=43;VDB=0.3767265;GB=-0.693146;MQ\$BZ=-0.056411;MQF=0AC-1AN=1;DPA=0.3,22,21;MQ=60;ANN=G upstream_gene_variant MODIFIER lof2 Rv0032 transcript Rv0032 protein_coding		
57 NC_00962.1	55553.	CCG	CGGTG	182.416.	INDEL;ID=4;MF=1;DPA=1;OP=4;VDB=0.3772756;GB=-0.556411;QBZ=1.70561;MQ\$BZ=0.056411;QBZ=1.70561;MQF=0AC-1AN=1;DPA=0.3,1,3;MQ=60;ANN=CGCTG conservative_inframe_insertion MODERATE polA Rv00		
58 NC_00962.1	56199.	G	A	197.416.	DP=8;VDB=0.3759225;GB=-0.651104;MQ\$BZ=0.056411;MQF=0AC-1AN=1;DPA=0.3,5,3;MQ=60;ANN=A synonymous_variant LOW Rv0051 Rv0051 transcript Rv0051 protein_coding  1/1 C 504		
59 NC_00962.1	62049.	A	G	225.417.	DP=17;VDB=0.00527805;GB=-0.651104;MQ\$BZ=0.056411;MQF=0AC-1AN=1;DPA=0.3,10,17;MQ=60;ANN=G missense_variant MODERATE Rv0063 Rv0063 transcript Rv0063 protein_coding  1/1		
60 NC_00962.1	67385.	G	A	225.417.	DP=31;VDB=0.9911245;GB=-0.693111;MQ\$BZ=-0.056411;MQF=0AC-1AN=1;DPA=0.3,12,19;MQ=60;ANN=A missense_variant MODERATE Rv0063 Rv0063 transcript Rv0063 protein_coding  1/1		
61 NC_00962.1	69342.	G	A	225.417.	DP=28;VDB=0.2138855;GB=-0.692976;MQ\$BZ=-0.056411;MQF=0AC-1AN=1;DPA=0.3,10,16;MQ=60;ANN=A synonymous_variant LOW Rv0064 Rv0064 transcript Rv0064 protein_coding  1/1		
62 NC_00962.1	69899.	G	A	225.417.	DP=38;VDB=0.4482855;GB=-0.693139;MQ\$BZ=-0.056411;MQF=0AC-1AN=1;DPA=0.3,20,16;MQ=60;ANN=A missense_variant MODERATE Rv0064 Rv0064 transcript Rv0064 protein_coding  1/1		
63 NC_00962.1	70816.	A	G	225.421.	DP=46;VDB=0.3279845;GB=-0.693147;MQ\$BZ=-0.056411;MQF=0AC-1AN=1;DPA=0.3,30,36;MQ=60;ANN=A missense_variant MODERATE Rv0064 Rv0064 transcript Rv0064 protein_coding  1/1		
64 NC_00962.1	71336.	G	C	225.417.	DP=24;VDB=0.3909145;GB=-0.686358;MQ\$BZ=-0.056411;MQF=0AC-1AN=1;DPA=0.3,10,16;MQ=60;ANN=C missense_variant MODERATE Rv0064 Rv0064 transcript Rv0064 protein_coding  1/1		
65 NC_00962.1	71584.	CGG	CGGAGGGCT	163.943.	INDEL;ID=6;MF=0.3;OP=20;VDB=0.0291735;GB=-0.618618;MQ\$BZ=-3.52286;MQ\$BZ=0.056411;QBZ=1.86221;CBZ=-2.4599;MQ\$BZ=0AC-1AN=1;DPA=13,1,6;MQ=60;ANN=CCAGGG		

Fig 1.1 A VCF annotated with SnpEff.

Summary	
Genome	NC_00962.3
Date	2024-08-08 17:48
SnpEff version	SnpEff 5.2c (build 2024-04-09 12:24), by Pablo Cingolani
Command line arguments	SnpEff NC_00962.3 raw.filter.vcf
Warnings	3,102
Errors	0
Number of lines (input file)	999
Number of variants (before filter)	999
Number of non-variants (i.e. reference equals alternative)	0
Number of variants processed (i.e. after filter and non-variants)	999
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 annotations	9,617
Genome total length	4,411,532
Genome effective length	4,411,532
Variant rate	1 variant every 4,415 bases

#### Variants rate details

Chromosome	Length	Variants	Variants rate
NC_00962.3	4,411,532	999	4,415
Total	4,411,532	999	4,415

Fig 1.2 SnpEff annotation summary.

#### Number variants by type

Type	Total
SNP	919
MNP	0
INS	39
DEL	50
MIXED	0
INV	0
DUP	0
BND	0
INTERVAL	0
Total	999

#### Number of effects by impact

Type (alphabetical order)	Count	Percent
HIGH	46	0.478%
LOW	320	3.327%
MODERATE	468	4.866%
MODIFIER	6,783	91.320%

#### Number of effects by functional class

Type (alphabetical order)	Count	Percent
MISSENSE	450	57.841%
NONSENSE	8	1.028%
SILENT	320	41.131%

Mis sense / Silent ratio: 1.4026

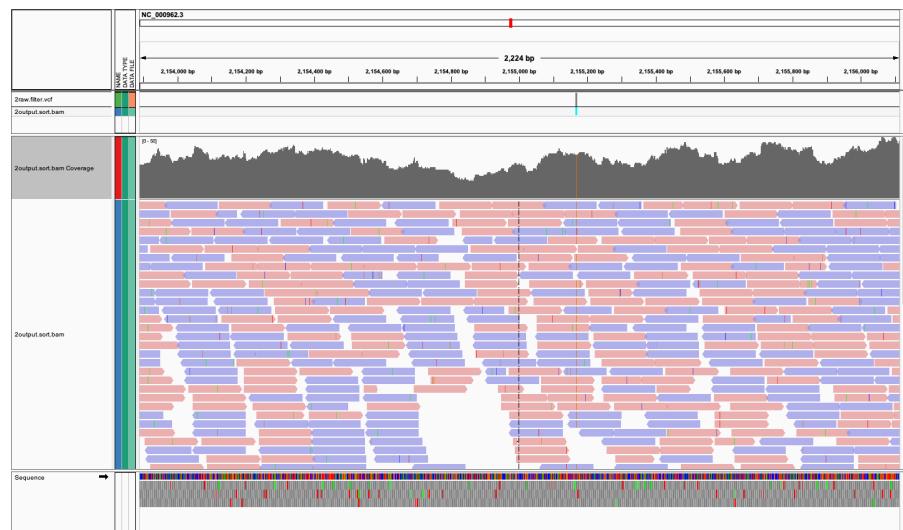
**Fig 1.3.** Variant effect summary.

Number of annotations and region counts		
Annotation	Region	
Type (alphabetical order)	Count	Percent
conservative_inframe_deletion	8	0.083%
conservative_inframe_insertion	7	0.073%
disruptive_inframe_deletion	2	0.021%
disruptive_inframe_insertion	4	0.042%
downstream_gene_variant	4,294	44.641%
frameshift_variant	35	0.364%
intergenic_region	127	1.32%
intragenic_variant	41	0.428%
missense_variant	447	4.647%
splice_region_variant	1	0.01%
start_lost	3	0.031%
stop_gained	8	0.083%
stop_lost	1	0.01%
synonymous_variant	320	3.327%
upstream_gene_variant	4,321	44.932%

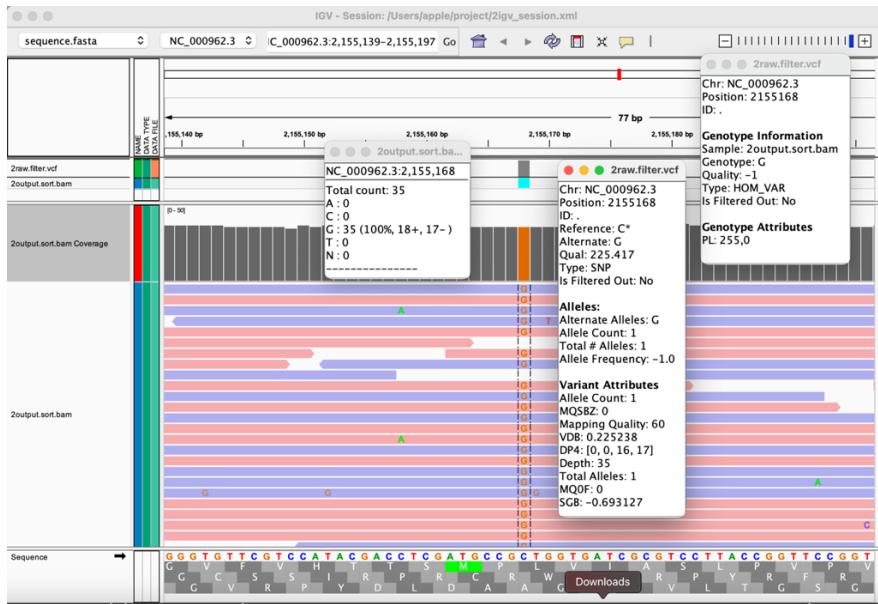
  

Type (alphabetical order)	Count	Percent
DOWNSTREAM	4,294	44.65%
EXON	834	8.672%
INTERGENIC	127	1.321%
TRANSCRIPT	41	0.426%
UPSTREAM	4,321	44.931%

**Fig 1.4** The Number of effects by type and region



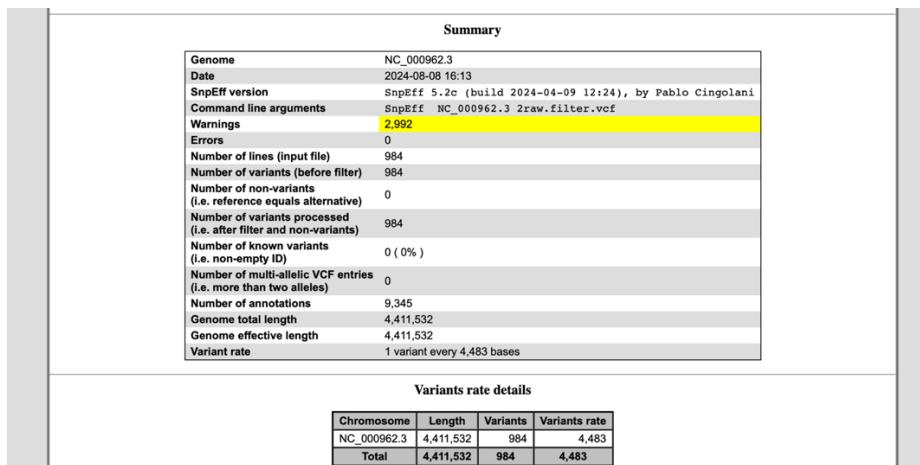
**Fig 2.0** IGV visualisation of a *Mycobacterium tuberculosis* H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 2



**Fig 2.0.1**

B	C	D	E	F	G	H	
512	G	A	225.417.	D	P=22V0B=0.503111<1.5C8B=0.02525A>C0BZ=0.240F=0AC>1AN=1.094=0.12,1.2,10,MQ=60,ANN=C upstream_gene_variant MODIFIER Rv12301 Rv12301 transcript Rv12301 protein_coding GTPL		
513	G	C	225.417.	D	P=23V0B=0.179495G8B=0.692717MGSBZ=C0M0F=0AC>1AN=1.094=0.17,6,MQ=0,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL		
514	T	T	225.417.	D	P=23V0B=0.111411NC8B=0.692811MGSBZ=C0M0F=0AC>1AN=1.094=0.17,15,MQ=0,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL		
515	2081087.	T	225.417.	D	P=26V0B=0.23540625G8B=0.693136MQSBZ=0.020F=0AC>1AN=1.094=0.19,17,MQ=60,ANN=1 transcript Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL		
516	2094911..	ACAGCGTCA	ACAGGGTCA	228.387.	INDELUDV=49.0MF=0.907407 P=54,VDB=0.01340685G8B=0.693147RPBZ=1.65926MQSBZ=0.6MQSBZ=0.80BZ=2.24935C8BZ=2.05198MQOF=0AC>1AN=1.094=1,4,3,11,MQ=0,ANN=ACA GTPL		
517	209186..	A	G	225.417..	D	P=66V0B=0.7738755G8B=0.693147MQSBZ=0.020F=0AC>1AN=1.094=0.37,29,MQ=60,ANN=C upstream_gene_variant LOW synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
518	2109523..	G	GGGG	228.413..	D	P=66V0B=0.49.0MF=0.998109 P=33,VDB=0.1071425G8B=0.693147RPBZ=1.69873MQSBZ=0.6MQSBZ=0.80BZ=1.768985C8BZ=0.020F=0AC>1AN=1.094=1,2,4,25,MQ=0,ANN=CGG upstream_gene_variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
519	2111348..	C	A	225.417..	D	P=66V0B=0.23540625G8B=0.693136MQSBZ=0.020F=0AC>1AN=1.094=0.37,29,MQ=60,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
520	2116903..	C	T	225.417..	D	P=66V0B=0.23540625G8B=0.693136MQSBZ=0.020F=0AC>1AN=1.094=0.37,29,MQ=60,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
521	2123935..	C	T	225.417..	D	P=21V0B=0.536395G8B=0.693152MQSBZ=0.020F=0AC>1AN=1.094=0.12,9,MQ=60,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
522	2123168..	T	G	225.417..	D	P=43V0B=0.7727435G8B=0.693146MQSBZ=0.020F=0AC>1AN=1.094=0.37,29,MQ=60,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
523	2123169..	T	G	225.422..	D	P=43V0B=0.7536385G8B=0.693146MQSBZ=0.020F=0AC>1AN=1.094=0.26,17,MQ=60,ANN=C upstream_gene_variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
524	2128870..	A	G	225.417..	D	P=29V0B=0.1101385G8B=0.693046MQSBZ=0.020F=0AC>1AN=1.094=0.37,29,MQ=60,ANN=C upstream_gene_variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
525	2133468..	TTCGGATGCC	TTCGGATGCC	199.8..	INDELUDV=4.4MF=0.424242 P=33,VDB=0.1393455G8B=0.688446MQSBZ=0.020F=0AC>1AN=1.094=0.37,29,MQ=60,ANN=C upstream_gene_variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL		
526	2151740..	T	C	225.417..	D	P=24V0B=0.23540625G8B=0.693136MQSBZ=0.020F=0AC>1AN=1.094=0.37,29,MQ=60,ANN=C upstream_gene_variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
527	2158245..	CTGGATTACCT	CTGGATTACCT	228.419..	D	P=53V0B=0.34MF=0.891893 P=37,VDB=0.0256995G8B=0.693117MQBZ=1.325MQSBZ=0.020F=0AC>1AN=1.094=0.10,21,MQ=60,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
528	2143228..	G	C	225.417..	D	P=31V0B=0.1133235G8B=0.693111MQSBZ=C0M0F=0AC>1AN=1.094=0.23,33,MQ=0,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
529	2143958..	C	T	228.388..	D	P=57V0B=0.1757695G8B=0.693147RPBZ=1.58194082>0AC>1AN=1.094=0.37,32,MQ=60,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 GTPL	
530	2147022..	A	C	225.422..	D	P=85V0B=0.0812575G8B=0.693147MQSBZ=0.020F=0AC>1AN=1.094=0.44,36,MQ=60,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
531	215168..	C	G	225.417..	D	P=35V0B=0.2252385G8B=0.693127MQSBZ=0.020F=0AC>1AN=1.094=0.18,17,MQ=60,ANN=C upstream_gene_variant MODERATE Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
532	2158327..	C	T	225.417..	D	P=53V0B=0.00404975G8B=0.693147MQSBZ=C0M0F=0AC>1AN=1.094=0.30,23,MQ=60,ANN=C synonymous variant LOW synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
533	2158493..	A	G	105.415..	D	P=53V0B=0.00404975G8B=0.693147MQSBZ=C0M0F=0AC>1AN=1.094=0.30,23,MQ=60,ANN=C synonymous variant LOW synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
534	2158504..	G	A	105.415..	D	P=53V0B=0.006334085G8B=0.590765MQSBZ=0.295174MQOF=0AC>1AN=1.094=0.2,2,3,MQ=39,ANN=C synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
535	2163496..	A	C	105.415..	D	P=53V0B=0.00767015G8B=0.590765MQSBZ=0.295174MQOF=0AC>1AN=1.094=0.2,2,3,MQ=39,ANN=C synonymous variant LOW synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
536	2163504..	G	A	211.417..	D	P=21V0B=0.04369225G8B=0.693147RPBZ=1.58194082>0AC>1AN=1.094=0.4,6,MQ=49,ANN=A synonymous variant LOW synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
537	2163510..	T	C	210.417..	D	P=10V0B=0.04369225G8B=0.693147RPBZ=1.58194082>0AC>1AN=1.094=0.4,6,MQ=49,ANN=A synonymous variant LOW synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
538	2163517..	A	C	201.416..	D	P=10V0B=0.02574515G8B=0.693147RPBZ=1.58194082>0AC>1AN=1.094=0.4,6,MQ=49,ANN=A synonymous variant LOW synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
539	2163520..	G	A	188.416..	D	P=10V0B=0.02209515G8B=0.670147MQSBZ=0.020F=0AC>1AN=1.094=0.4,6,MQ=49,ANN=A synonymous variant LOW synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
540	2165296..	A	C	225.417..	D	P=23V0B=0.917395G8B=0.693125MQSBZ=0.020F=0AC>1AN=1.094=0.8,14,MQ=60,ANN=C synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
541	2165503..	T	A	228.383..	D	P=41V0B=0.11318735G8B=0.693145RPBZ=1.277584MQSBZ=0.2,MQOF=0AC>1AN=1.094=0.10,21,MQ=60,ANN=C synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
542	2176101..	A	G	225.417..	D	P=30V0B=0.991415G8B=0.693147RPBZ=1.58194082>0AC>1AN=1.094=0.15,14,MQ=60,ANN=C synonymous variant MODIFIER Rv18356 Rv18356 transcript Rv18356 protein_coding 1/1 c GTPL	
543	2181805..	AG	A	228.403..	D	P=55V0B=0.9048345G8B=0.693147MQSBZ=0.020F=0AC>1AN=1.094=0.30,25,MQ=60,ANN=G synonymous variant LOW Rv1931c transcript Rv1931c protein_coding 1/1 c GTPL	
544	2183054..	T	G	225.417..	D	P=55V0B=0.9048345G8B=0.693147MQSBZ=0.020F=0AC>1AN=1.094=0.30,25,MQ=60,ANN=G synonymous variant LOW Rv1931c transcript Rv1931c protein_coding 1/1 c GTPL	

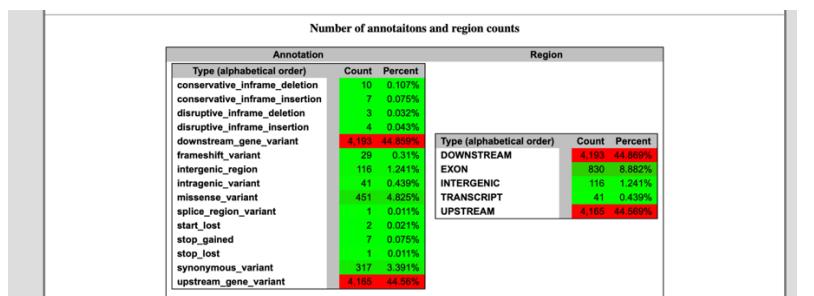
**Fig 2.1 A VCF annotated with SnpEff.**



**Fig 2.2** SnpEff annotation summary.



**Fig 2.3** Variant effect summary



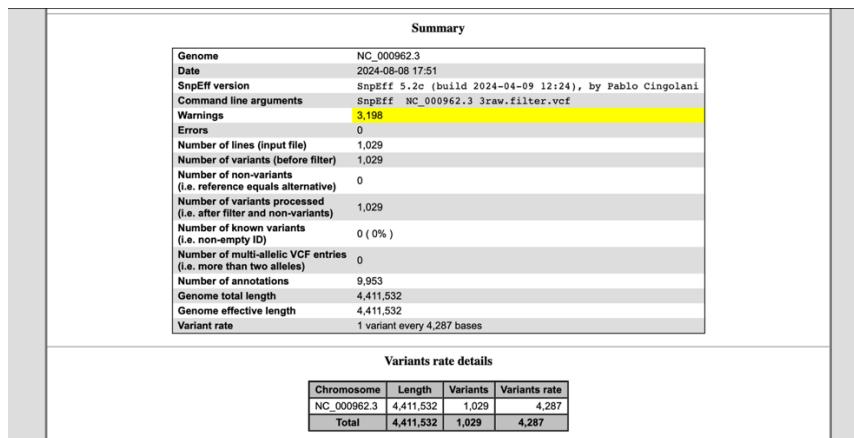
**Fig 2.4** The Number of effects by type and region



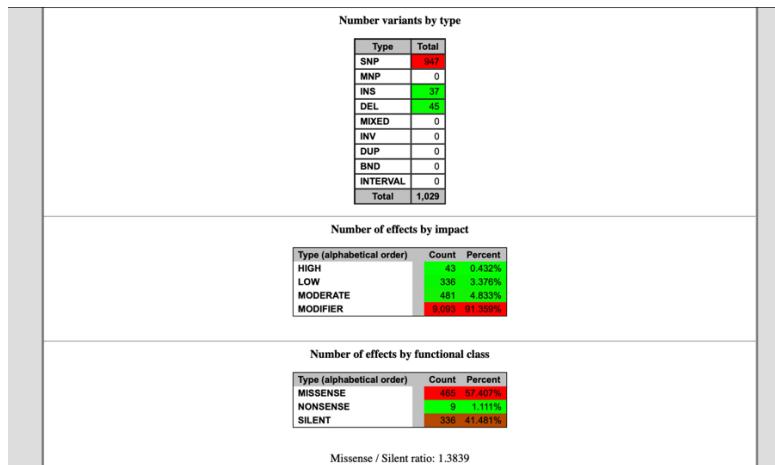
**Fig 3.0** IGV visualisation of a Mycobacterium tuberculosis H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 3

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
33	CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	3output.sort.bam										
34	NC_000962.3	1199	.	G	A	225.417	.	DP=82;VDB=GT;PL	1,255,0											
35	NC_000962.3	1977	.	A	G	225.417	.	DP=51;VDB=GT;PL	1,255,0											
36	NC_000962.3	4013	.	T	C	225.417	.	DP=70;VDB=GT;PL	1,255,0											
37	NC_000962.3	7362	.	G	C	225.417	.	DP=55;VDB=GT;PL	1,255,0											
38	NC_000962.3	7556	.	G	C	225.417	.	DP=54;VDB=GT;PL	1,255,0											
39	NC_000962.3	9304	.	G	A	225.417	.	DP=55;VDB=GT;PL	1,255,0											
40	NC_000962.3	11879	.	A	G	225.417	.	DP=56;VDB=GT;PL	1,255,0											
41	NC_000962.3	14785	.	T	C	225.417	.	DP=61;VDB=GT;PL	1,255,0											
42	NC_000962.3	18091	.	G	A	225.417	.	DP=27;VDB=GT;PL	1,255,0											
43	NC_000962.3	21795	.	G	A	95.4151	.	DP=4;VDB=0 GT;PL	1,125,0											
44	NC_000962.3	23854	.	A	G	225.417	.	DP=76;VDB=GT;PL	1,255,0											
45	NC_000962.3	26233	.	G	C	225.417	.	DP=54;VDB=GT;PL	1,255,0											
46	NC_000962.3	30695	.	C	G	225.417	.	DP=56;VDB=GT;PL	1,255,0											
47	NC_000962.3	37518	TAAAAAA	TAAAAAA		225.417	.	INDelJ0h=3 GT;PL	1,255,0											
48	NC_000962.3	27918	G	A	225.417	.	DP=64;VDB=GT;PL	1,255,0												
49	NC_000962.3	32387	C	T	225.417	.	DP=64;VDB=GT;PL	1,255,0												
50	NC_000962.3	33817	C	G	228.419	.	DP=71;VDB=GT;PL	1,255,0												
51	NC_000962.3	34044	T	C	225.417	.	DP=59;VDB=GT;PL	1,255,0												
52	NC_000962.3	37031	C	G	225.417	.	DP=37;VDB=GT;PL	1,255,0												
53	NC_000962.3	42967	G	C	225.417	.	DP=54;VDB=GT;PL	1,255,0												
54	NC_000962.3	47131	T	TTT	CTT	225.417	.	DP=50;VDB=GT;PL	1,255,0											
55	NC_000962.3	48079	C	G	225.422	.	DP=66;VDB=GT;PL	1,255,0												
56	NC_000962.3	49323	AGGG	AGG	228.419	.	INDelJ0h=6 GT;PL	1,255,0												
57	NC_000962.3	53502	G	A	225.417	.	DP=37;VDB=GT;PL	1,255,0												
58	NC_000962.3	54304	C	T	225.417	.	DP=48;VDB=GT;PL	1,255,0												
59	NC_000962.3	55553	C	T	192.416	.	DP=13;VDB=GT;PL	1,222,0												
60	NC_000962.3	60456	G	C	225.417	.	DP=43;VDB=GT;PL	1,255,0												
61	NC_000962.3	62049	A	G	225.417	.	DP=35;VDB=GT;PL	1,255,0												
62	NC_000962.3	67385	G	A	225.417	.	DP=60;VDB=GT;PL	1,255,0												
63	NC_000962.3	69042	G	A	225.417	.	DP=60;VDB=GT;PL	1,255,0												
64	NC_000962.3	69980	G	A	225.422	.	DP=58;VDB=GT;PL	1,255,0												
65	NC_000962.3	70816	A	G	225.417	.	DP=41;VDB=GT;PL	1,255,0												
66	NC_000962.3	71336	G	C	166.416	.	DP=7;VDB=0 GT;PL	1,196,0												

**Fig 3.1** A VCF annotated with SnpEff.



**Fig 3.2** SnpEff annotation summary



**Fig 3.3 .** Variant effect summary

Number of annotations and region counts		
Annotation		Region
Type (alphabetical order)	Count	Percent
conservative_inframe_deletion	7	0.07%
conservative_inframe_insertion	6	0.06%
disruptive_inframe_deletion	1	0.01%
disruptive_inframe_insertion	4	0.04%
downstream_gene_variant	4,538	48.585%
frameshift_variant	32	0.321%
intergenic_region	132	1.326%
intragenic_variant	40	0.402%
missense_variant	463	4.651%
splice_region_variant	1	0.01%
start_lost	2	0.02%
stop_gained	9	0.09%
stop_lost	1	0.01%
synonymous_variant	336	3.375%
upstream_gene_variant	4,383	44.028%
Type (alphabetical order)	Count	Percent
DOWNSTREAM	4,538	45.594%
EXON	880	8.641%
INTERGENIC	132	1.326%
TRANSCRIPT	40	0.402%
UPSTREAM	4,383	44.037%

**Fig 3.4** The Number of effects by type and region



**Fig 4.0** IGV visualisation of a Mycobacterium tuberculosis H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 4

A	B	C	D	E	F	G	H	I	J	K	L
FORMAT	output.sort.bam										
33 #CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO				
34 NC_000962:	1199	.	G	A	225.417	.	DP=93;VDB=0.458754;5GB=-0.693147;MQSBZ=1.05654;MQOF=0;AC=1;AN=1;DP4=0,0,48,43;MQ=59;ANN=A missense_variant MODERATE GT:PL	1:255,0			
35 NC_000962:	1977	.	A	G	225.417	.	DP=63;VDB=0.0874014;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,20,41;MQ=60;ANN=G upstream_gene_variant MODERATE ref GT:PL	1:255,0			
36 NC_000962:	4013	.	T	C	225.417	.	DP=73;VDB=0.941131;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,26,47;MQ=60;ANN=C missense_variant MODERATE ref GT:PL	1:255,0			
37 NC_000962:	7362	.	G	C	225.417	.	DP=52;VDB=0.250536;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,26,23;MQ=60;ANN=C missense_variant MODERATE gt GT:PL	1:255,0			
38 NC_000962:	7585	.	G	C	225.417	.	DP=46;VDB=0.097744;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,27,20;MQ=60;ANN=A missense_variant MODERATE gt GT:PL	1:255,0			
39 NC_000962:	9304	.	G	A	225.417	.	DP=48;VDB=0.093863;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,15,20;MQ=60;ANN=G missense_variant MODERATE gt GT:PL	1:255,0			
40 NC_000962:	11879	.	A	G	225.417	.	DP=48;VDB=0.093863;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,15,20;MQ=60;ANN=G missense_variant MODERATE gt GT:PL	1:255,0			
41 NC_000962:	14722	.	T	C	225.417	.	DP=78;VDB=0.097522;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,15,20;MQ=60;ANN=C missense_variant MODERATE gt GT:PL	1:255,0			
42 NC_000962:	18091	.	G	A	225.417	.	DP=38;VDB=0.917405;5GB=-0.693139;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,30,60;MQ=60;ANN=A synonymous_variant LOW phon Rv003 GT:PL	1:255,0			
43 NC_000962:	21795	.	G	A	189.416	.	DP=10;VDB=0.102407;5GB=-0.670168;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,2,8;MQ=60;ANN=A missense_variant MODERATE ptp1 GT:PL	1:219,0			
44 NC_000962:	23854	.	A	G	225.417	.	DP=57;VDB=0.511872;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,2,8;MQ=60;ANN=A missense_variant MODERATE gt GT:PL	1:255,0			
45 NC_000962:	26233	.	G	C	225.417	.	DP=38;VDB=0.259611;5GB=-0.693143;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,15,23;MQ=60;ANN=C missense_variant MODERATE Rv002 GT:PL	1:255,0			
46 NC_000962:	26959	.	C	G	225.417	.	DP=46;VDB=0.0652945;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,15,23;MQ=60;ANN=C missense_variant MODERATE gt GT:PL	1:255,0			
47 NC_000962:	27518	AAAAAA	TAAAAA		225.417	.	INDEL;DV=49;MF=1;DP=19;DB=0,0;O99999;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,24,25;MQ=60;ANN=TA synonymous_variant LOW Rv0023 gt GT:PL	1:255,0			
48 NC_000962:	27918	.	G	A	225.417	.	DP=68;VDB=0.608663;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,24,25;MQ=60;ANN=TA synonymous_variant LOW Rv0023 gt GT:PL	1:255,0			
49 NC_000962:	32387	.	C	T	225.417	.	DP=52;VDB=0.6450563;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,27,23;MQ=60;ANN=T missense_variant MODERATE Rv0023 GT:PL	1:255,0			
50 NC_000962:	33817	.	C	G	225.417	.	DP=62;VDB=0.178457;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,26,34;MQ=60;ANN=G upstream_gene_variant MODIFIER bi GT:PL	1:255,0			
51 NC_000962:	34044	.	T	C	225.417	.	DP=73;VDB=0.859975;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,42,30;MQ=60;ANN=C upstream_gene_variant MODIFIER bi GT:PL	1:255,0			
52 NC_000962:	37031	.	C	G	225.417	.	DP=58;VDB=0.658122;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,30,27;MQ=60;ANN=G synonymous_variant LOW Rv0034 gt GT:PL	1:255,0			
53 NC_000962:	42567	.	G	C	225.417	.	DP=59;VDB=0.0699934;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,28,30;MQ=60;ANN=C synonymous_variant LOW mtz28 Rv0034 gt GT:PL	1:255,0			
54 NC_000962:	47138	.	CTT	TTT	228.418	.	INDEL;DV=58;MF=0;95082;DP=61;DB=059857;5GB=-0.693147;R8PBZ=0;600747;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,31,24;MQ=60;ANN=G upstream_gene_variant MODIFIER gt GT:PL	1:255,0			
55 NC_000962:	49079	.	C	G	225.417	.	DP=82;VDB=0.655149;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,32,36;MQ=60;ANN=G synonymous_variant LOW Rv0045 gt GT:PL	1:255,0			
56 NC_000962:	49233	.	AGGG	TTTT	228.418	.	INDEL;DV=60;MF=0;983607;DP=61;DB=0;0914894;SGB=-0.693147;R8PBZ=1;70495;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,18,23;MQ=60;ANN=A upstream_gene_variant MODIFIER Rv0045 gt GT:PL	1:255,0			
57 NC_000962:	53502	.	G	A	225.417	.	DP=43;VDB=0.769179;5GB=-0.693145;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,18,23;MQ=60;ANN=A upstream_gene_variant MODIFIER Rv0045 gt GT:PL	1:255,0			
58 NC_000962:	54304	.	C	T	225.417	.	DP=54;VDB=0.489122;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,24,34;MQ=60;ANN=G upstream_gene_variant MODIFIER bi GT:PL	1:255,0			
59 NC_000962:	55553	.	C	T	169.416	.	DP=19;VDB=0.0146148;5GB=-0.691688;MQOF=0;AC=1;AN=1;DP4=0,0,19;MQ=60;ANN=T missense_variant MODERATE ponA1 Rv0050 tra GT:PL	1:199,0			
60 NC_000962:	60456	.	G	C	225.417	.	DP=77;VDB=0.348309;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,31,35;MQ=60;ANN=C missense_variant MODERATE dnab1 GT:PL	1:255,0			
61 NC_000962:	62049	.	A	G	225.417	.	DP=39;VDB=0.161600;5GB=-0.693144;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,12,27;MQ=60;ANN=G missense_variant MODERATE dnab1 GT:PL	1:255,0			
62 NC_000962:	67385	.	G	A	225.417	.	DP=53;VDB=0.999702;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,31,20;MQ=60;ANN=A missense_variant MODERATE Rv0063 GT:PL	1:255,0			
63 NC_000962:	69342	.	G	A	225.417	.	DP=80;VDB=0.196932;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,24,54;MQ=60;ANN=A synonymous_variant LOW Rv0064 gt GT:PL	1:255,0			
64 NC_000962:	69989	.	G	A	225.417	.	DP=54;VDB=0.248174;5GB=-0.693147;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,27,27;MQ=60;ANN=A missense_variant MODERATE Rv0064 GT:PL	1:255,0			
65 NC_000962:	70816	.	A	G	225.417	.	DP=43;VDB=0.110236;5GB=-0.693146;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,2,17;MQ=60;ANN=G missense_variant MODERATE Rv0064 GT:PL	1:255,0			
66 NC_000962:	71336	.	G	C	225.417	.	DP=11;VDB=0.253952;5GB=-0.670168;MQSBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,7,3;MQ=60;ANN=C missense_variant MODERATE Rv0064 GT:PL	1:255,0			

Fig 4.1 A VCF annotated with SnpEff

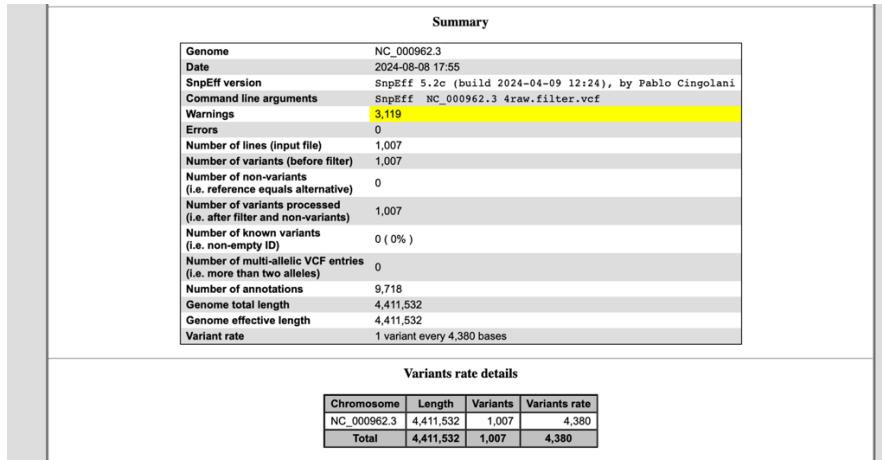


Fig 4.2 SnpEff annotation summary

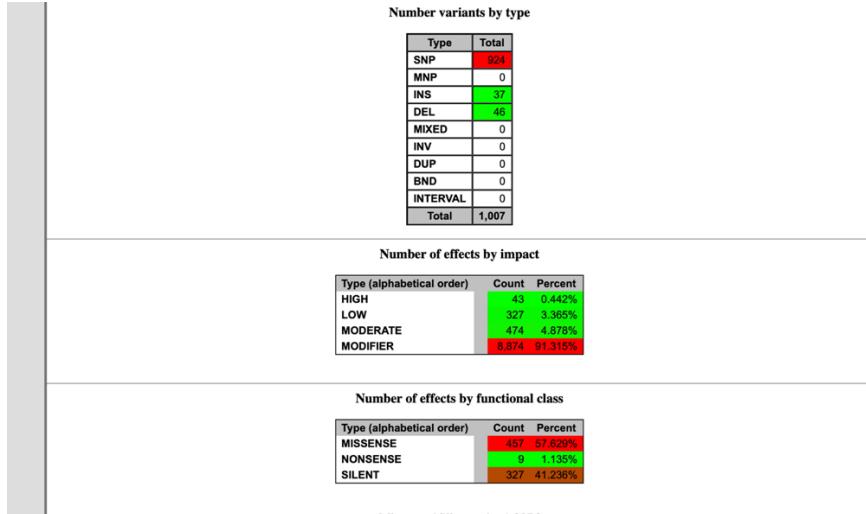


Fig 4.3 . Variant effect summary

Number of annotations and region counts		
Annotation	Region	
Type (alphabetical order)	Count	Percent
conservative_inframe_deletion	8	0.082%
conservative_inframe_insertion	6	0.062%
disruptive_inframe_deletion	1	0.01%
disruptive_inframe_insertion	4	0.041%
downstream_gene_variant	4,417	45.442%
frameshift_variant	32	0.329%
intergenic_region	126	1.296%
intragenic_variant	40	0.412%
missense_variant	455	4.681%
splice_region_variant	1	0.01%
start_lost	2	0.021%
stop_gained	9	0.093%
stop_lost	1	0.01%
synonymous_variant	327	3.364%
upstream_gene_variant	4,291	44.146%
Type (alphabetical order)	Count	Percent
DOWNSTREAM	4,417	45.452%
EXON	844	8.685%
INTERGENIC	126	1.297%
TRANSCRIPT	40	0.412%
UPSTREAM	4,291	44.155%

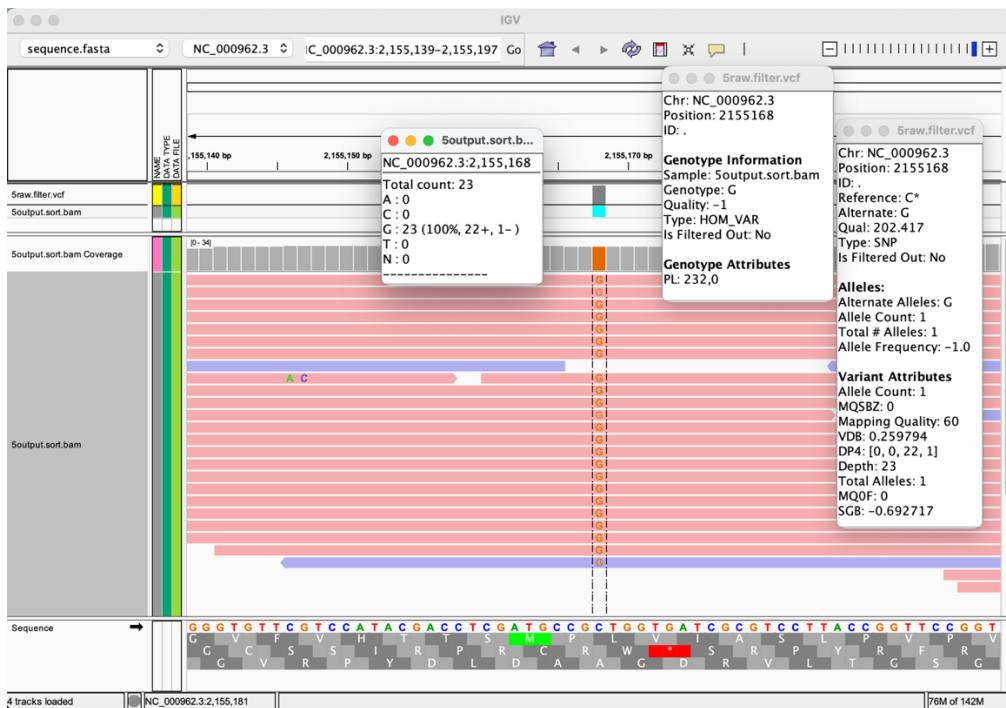
**Fig 4.4** The Number of effects by type and region



**Fig 5.0** IGV visualisation of a *Mycobacterium tuberculosis* H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 5



**Fig 5.0.1**



**Fig 5.0.2**

	B	C	D	E	F	G	H
459	2094911.		ACAGGGCTCA-AGAGGCTCA-	238_368.			
460	2096186.	A	G	225_417.			
461	2109523.	CGG	CGGG	225_417.			
462	2113058.	C	A	225_417.			
463	2115903.	C	T	225_417.			
464	2123168.	C	T	225_417.			
465	2123168.	T	G	225_417.			
466	2123169.	T	G	225_417.			
467	2123870.	A	G	225_417.			
468	2133468.	TTGCAT	TTGCATGCC	206_069.			
469	2133870.	T	C	225_417.			
470	2134212.	CT	CT	225_417.			
471	2143328.	G	C	164_416.			
472	2143958.	C	T	225_417.			
473	2147022.	A	C	225_417.			
474	215168.	C	G	202_417.			
475	2153327.	C	T	225_417.			
476	2153327.	C	T	225_417.			
477	2153444.	T	C	77_470.			
478	2163481.	G	A	12_775.			
479	2163484.	A	C	4_54011.			
480	2163493.	A	G	27_4222.			
481	2163494.	G	C	27_4222.			
482	2163510.	A	C	27_4222.			
483	2163510.	G	A	176_416.			
484	2163510.	T	C	176_416.			
485	2163517.	A	C	219_417.			
486	2163520.	G	A	225_417.			
487	2163790.	A	C	221_417.			
488	2173101.	A	G	225_417.			
489	2173101.	G	C	225_417.			
490	2181805.	AG	AG	228_417.			
491	2183054.	T	G	225_417.			
492	2207591.	T	TC	228_417.			

Fig 5.1 A VCF annotated with SnpEff

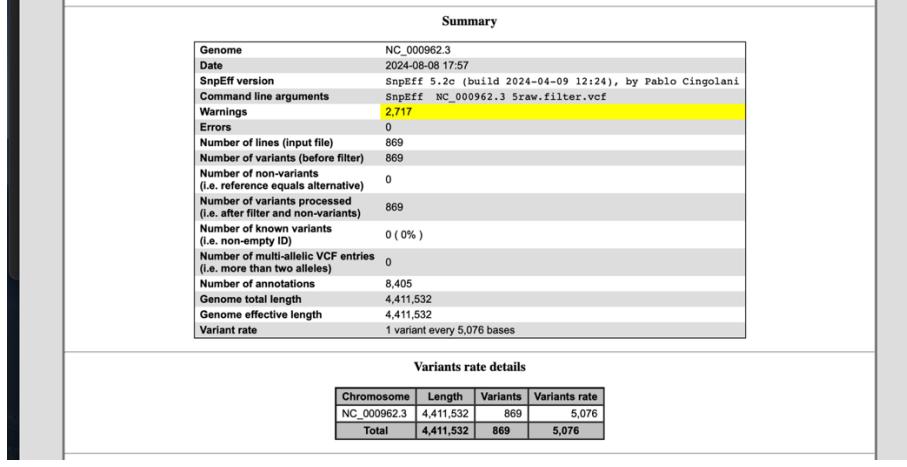


Fig 5.2 SnpEff annotation summary

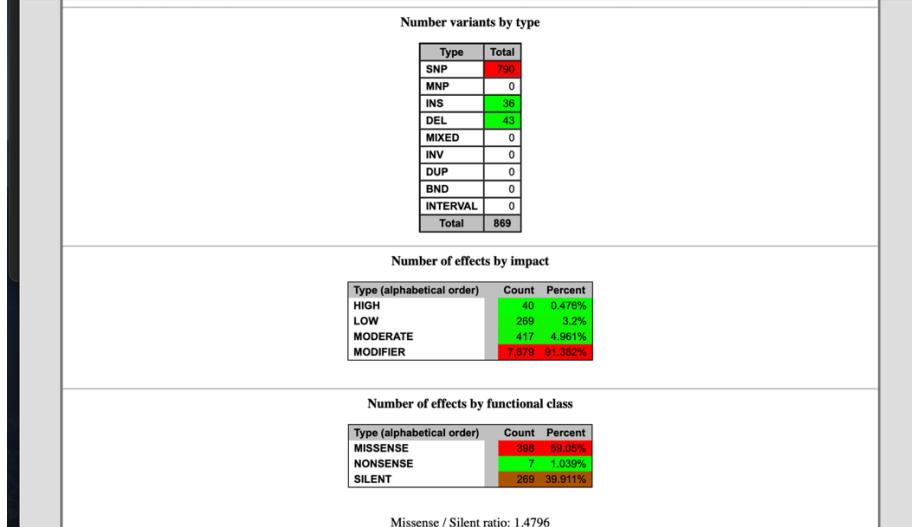


Fig 5.3 Variant effect summary.

Number of annotations and region counts		
Annotation		Region
Type (alphabetical order)	Count	Percent
conservative_inframe_deletion	6	0.071%
conservative_inframe_insertion	9	0.107%
disruptive_inframe_deletion	2	0.024%
disruptive_inframe_insertion	4	0.048%
downstream_gene_variant	3,789	45.07%
frameshift_variant	31	0.369%
intergenic_region	108	1.285%
intragenic_variant	37	0.44%
missense_variant	396	4.71%
splice_region_variant	1	0.012%
start_lost	2	0.024%
stop_gained	7	0.083%
stop_lost	1	0.012%
synonymous_variant	269	3.2%
upstream_gene_variant	3,745	44.546%
Type (alphabetical order)	Count	Percent
DOWNSTREAM	3,789	45.08%
EXON	726	8.638%
INTERGENIC	108	1.285%
TRANSCRIPT	37	0.44%
UPSTREAM	3,745	44.557%

**Fig 5.4** The Number of effects by type and region



**Fig 6.0** IGV visualisation of a *Mycobacterium tuberculosis* H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 6

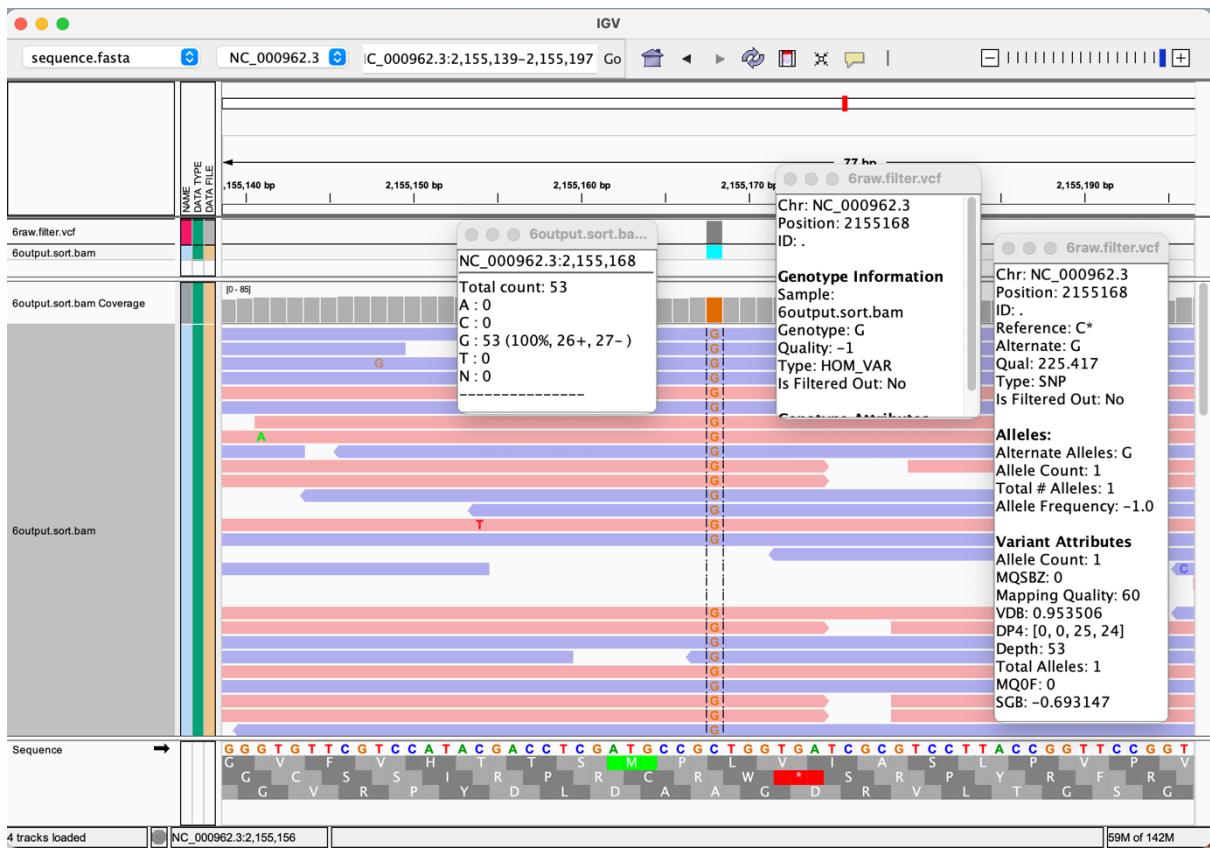
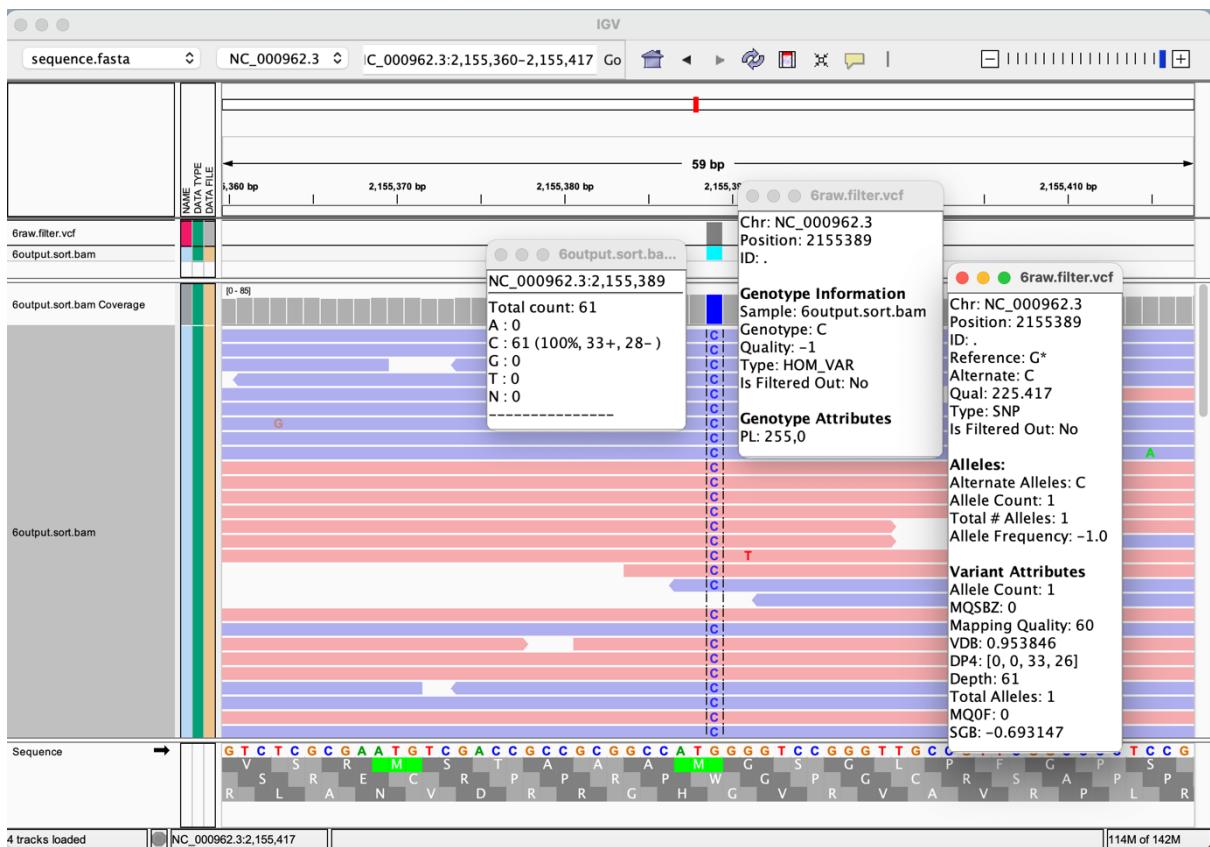


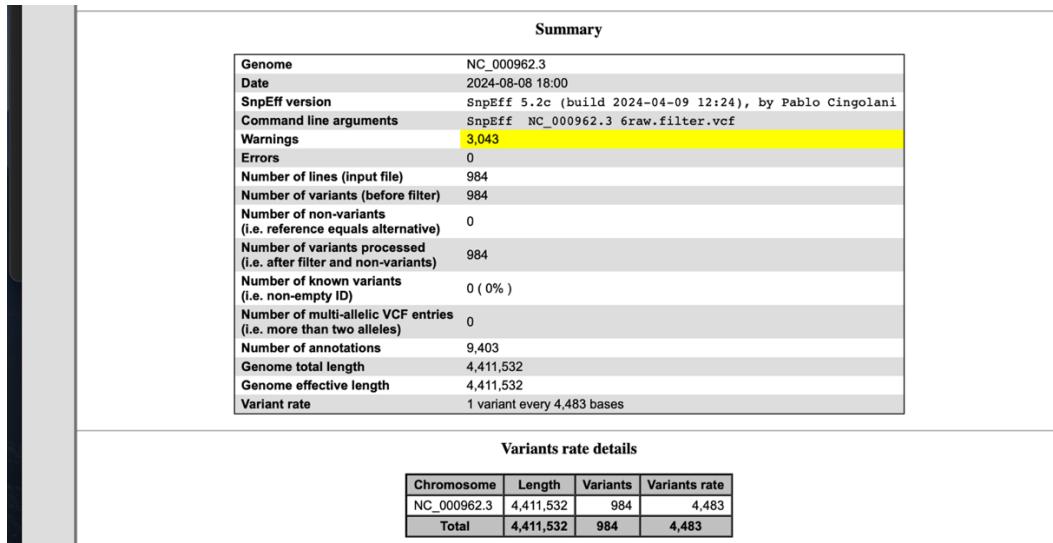
Fig 6.0.1



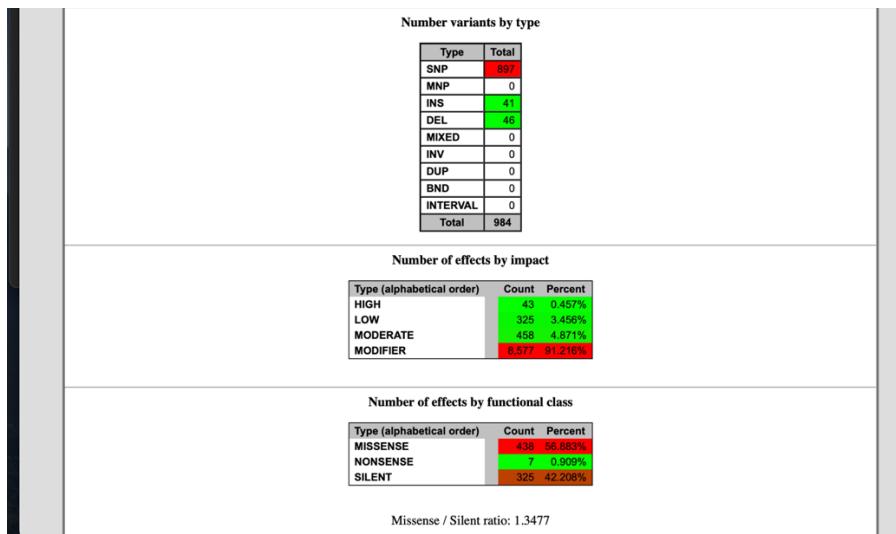
**Fig 6.0.2**

A	B	C	D	E	F	G	H
524 00962.3	2133168	T	G	225.417	.	DP=56;VD=0.992288;5;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,33,22;MQ=60;ANN=C upstream_gene_variant MODIFIER Rv1860c transcript Rv1860c protein_coding	
525 00962.3	2133169	T	G	225.417	.	DP=56;VD=0.992288;5;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,31,32;MQ=60;ANN=C upstream_gene_variant MODIFIER Rv1860c transcript Rv1860c protein_coding	
526 00962.3	2128870	A	G	225.417	.	DP=39;VD=0.430526;5;GB=-0.693139;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,1,26;MQ=60;ANN=C synonymous_variant LOW gfa3 Rv1878 transcript Rv1878 protein_coding 1/1 c-8494	
527 00962.3	2133468	TTGGCAT	TTGGCATGCC	177.916	.	INDEL;DV=22;IM=0,4488;DF=-49;VD=0,000870;6;GB=-0.692552;RPB2=4,1364;MQ=2;0;ANS2;0;QBZ=2,1,85501;SCBZ=2,4,0657;MQF=0;AC<1;AN>1;DP4=0,12,13,5,17;MQ=60;ANN=TTG	
528 00962.3	2135870	T	C	225.417	.	DP=57;VD=0.999922;6;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,21,34;MQ=60;ANN=C upstream_gene_variant MODIFIER ip transcript Rv1880c protein_coding	
529 00962.3	2138245	CTGGTAATC	CT	228.335	.	INDEL;DV=43;IM=0,82693;DP=52;VD=0,010759;2;5;GB=-0.693143;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,3,6,17,26;MQ=60;ANN=CT u	
530 00962.3	2143328	G	C	225.417	.	DP=34;VD=0.996281;5;GB=-0.693132;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,1,18;MQ=60;ANN=C missense_variant MODERATE Rv1895 Rv1895 protein_coding 1/1 c	
531 00962.3	2143958	C	T	225.417	.	DP=70;VD=0.073712;8;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,1,10;MQ=60;ANN=I strand_gained HIGH Rv1896c transcript Rv1896c protein_coding 1/1 c-489564	
532 00962.3	2147022	A	C	228.419	.	DP=72;VD=0.605654;6;GB=-0.693147;RPB2=0,0247615;MQB2=0,02872;0;QBZ=2,1,7278;SCBZ=2,6,6275;MQF=0;AC<1;AN>1;DP4=1,0,40,29;MQ=60;ANN=C missense_variant MODERATE ip	
533 00962.3	2155168	C	G	225.417	.	DP=53;VD=0.953506;5;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,25,24;MQ=60;ANN=G missense_variant MODERATE ip transcript Rv1908c protein_coding 1/1 c	
534 00962.3	2155389	G	C	225.417	.	DP=61;VD=0.953846;5;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,33,26;MQ=60;ANN=C synonymous_variant LOW katG Rv1908c transcript Rv1908c protein_coding 1/1 c c-230	
535 00962.3	2158327	C	T	225.417	.	DP=79;VD=0.180001;6;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,42,35;MQ=60;ANN=T synonymous_variant LOW fadB5 Rv1912c protein_coding 1/1 c 765	
536 00962.3	2163493	A	G	192.416	.	DP=17;VD=0,142431;6;GB=-0.690438;MQ\$B2=1,6902;MQF=0,058823;5;AC<1;AN>1;DP4=0,0,3,14;MQ=49;ANN=C synonymous_variant LOW PPE34 Rv1917c protein_coding	
537 00962.3	2163494	G	C	191.416	.	DP=17;VD=0,960242;9;5;GB=-0.690438;MQ\$B2=1,6902;MQF=0,058823;5;AC<1;AN>1;DP4=0,0,10,8;MQ=57;ANN=C synonymous_variant LOW PPE34 Rv1917c protein_coding	
538 00962.3	2163496	A	C	186.416	.	DP=17;VD=0,129926;9;8;GB=-0.690438;MQ\$B2=1,7602;MQF=0,058823;5;AC<1;AN>1;DP4=0,0,1,18;MQ=60;ANN=C missense_variant MODERATE Rv1895 Rv1895 protein_coding 1/1 c	
539 00962.3	2163504	G	A	225.417	.	DP=25;VD=0,467187;6;5;GB=-0.692914;MQ\$B2=1,8616;MQF=0,04;AC<1;AN>1;DP4=0,0,4,21;MQ=52;ANN=A synonsense_variant LOW PPE34 Rv1917c transcript Rv1917c protein_coding	
540 00962.3	2163510	T	C	225.417	.	DP=26;VD=0,000110;6;5;GB=-0.692976;MQ\$B2=1,32156;MQF=0,0384615;AC<1;AN>1;DP4=0,0,5,21;MQ=53;ANN=C missense_variant MODERATE Rv1917c transcript Rv1917c protein_coding	
541 00962.3	2163517	A	C	225.417	.	DP=26;VD=0,829224;6;5;GB=-0.692976;MQ\$B2=1,32156;MQF=0,0384615;AC<1;AN>1;DP4=0,0,5,21;MQ=53;ANN=C synonymous_variant LOW PPE34 Rv1917c transcript Rv1917c protein_coding	
542 00962.3	2163520	G	A	225.417	.	DP=26;VD=0,829223;6;5;GB=-0.692976;MQ\$B2=1,32156;MQF=0,0384615;AC<1;AN>1;DP4=0,0,5,21;MQ=53;ANN=C synonymous_variant LOW PPE34 Rv1917c transcript Rv1917c protein_coding	
543 00962.3	2163790	A	C	225.417	.	DP=18;VD=0,180001;6;5;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,42,35;MQ=60;ANN=T synonymous_variant LOW PPE34 Rv1917c protein_coding	
544 00962.3	2165286	A	C	225.417	.	DP=45;VD=0,934707;7;5;GB=-0.693144;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,12,27;MQ=60;ANN=C missense_variant MODERATE Rv1917c Rv1917c protein_coding	
545 00962.3	2165503	T	A	162.081	.	DP=50;VD=0,423516;5;GB=-0.693141;RPB2=1,36686;MQB2=1,99062;MQF=0,04;AC<1;AN>1;DP4=0,4,23,14;MQ=50;ANN=A synonymous_variant LOW PPE34 Rv1917c transcript Rv1917c protein_coding	
546 00962.3	2167101	A	G	225.417	.	DP=53;VD=0,930469;5;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,22,27;MQ=60;ANN=G missense_variant MODERATE ip transcript Rv1923 Rv1923 protein_coding 1/1 c 92	
547 00962.3	2181805	AG	A	228.397	.	INDEL;DV=67;IM=0,917808;D9=73;VD=0,999927;5;GB=-0.693147;RPB2=0,0412;MQB2=0;MQF=0,04;AC<1;AN>1;DP4=0,4,21;MQ=52;ANN=A frameshift_variant LOW PPE34 Rv1917c transcript Rv1917c protein_coding	
548 00962.3	2183054	T	G	225.422	.	DP=90;VD=0,781921;5;GB=-0.693145;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,36;MQ=60;ANN=G synonymous_variant LOW Rv1931c Rv1931c transcript Rv1931c protein_coding 1/1 c 1	
549 00962.3	2207591	G	A	225.409	.	DP=34;VD=0,818886;5;GB=-0.693145;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,38,-0,04;AC<1;AN>1;DP4=0,0,44,44;MQ=60;ANN=T synonymous_variant LOW PPE34 Rv1917c protein_coding 1/1 c	
550 00962.3	2211809	T	C	225.409	.	INDEL;DV=65;IM=0,20239;D9=5;VD=0,000232;5;GB=-0,693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,31,48;MQ=60;ANN=C synonymous_variant LOW PPE34 Rv1917c protein_coding 1/1 c	
551 00962.3	2211826	A	G	225.417	.	DP=68;VD=0,576215;5;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,31,36;MQ=60;ANN=C synonymous_variant LOW imc2 Rv1968 protein_coding 1/1 c	
552 00962.3	2216443	C	A	225.417	.	DP=71;VD=0,818886;5;GB=-0.693143;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,36,33;MQ=60;ANN=A missense_variant MODERATE imc2f Rv1971 Rv1971 protein_coding 1/1 c	
553 00962.3	2220512	T	G	225.417	.	DP=34;VD=0,424561;5;GB=-0.693127;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,22,11;MQ=60;ANN=G synonymous_variant LOW Rv1977 Rv1977 transcript Rv1977 Rv1977 protein_coding 1/1 c 755	
554 00962.3	2221746	CGGGCG	CGGGCGGCC	228.379	.	INDEL;DV=45;IM=0,918367;D9=49;VD=0,000132;5;GB=-0.693147;RPB2=2,44942;MQB2=0;MQF=0,08;QBZ=2,149091;SCBZ=2,3,3543;MQF=0;AC<1;AN>1;DP4=1,2,12,23;MQ=60;ANN=CGGG	
555 00962.3	2223293	T	C	225.417	.	DP=43;VD=0,859494;5;GB=-0.693145;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,23,18;MQ=60;ANN=C upstream_gene_variant MODIFIER Rv1976c Rv1976c transcript Rv1976c protein_coding	
556 00962.3	2228067	A	G	225.417	.	DP=63;VD=0,350528;5;GB=-0.693147;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,22,39;MQ=60;ANN=G upstream_gene_variant MODIFIER mpf64 Rv1980c transcript Rv1980c protein_coding	
557 00962.3	2233947	T	C	225.417	.	DP=43;VD=0,109514;5;GB=-0.693145;MQ\$B2=0;MQF=0;AC<1;AN>1;DP4=0,0,23,18;MQ=60;ANN=C upstream_gene_variant MODIFIER Rv1985c Rv1985c transcript Rv1985c protein_coding	

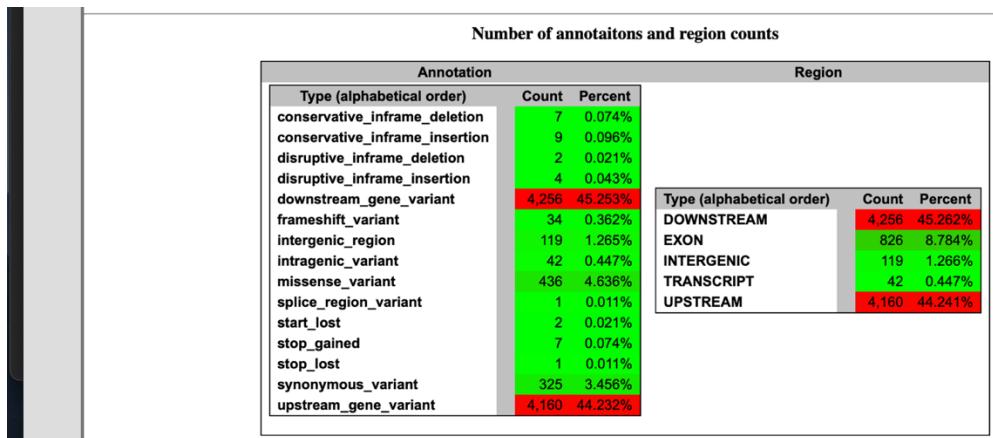
**Fig 6.1 A VCF annotated with SnpEff.**



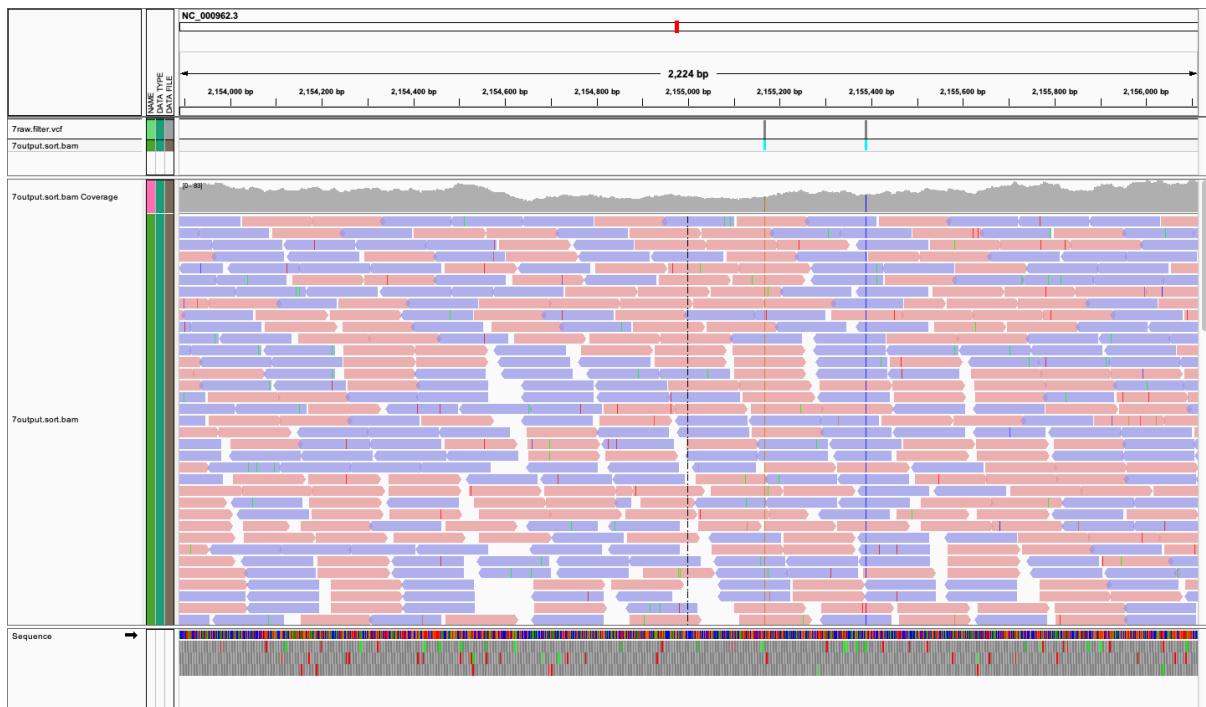
**Fig 6.2 SnpEff annotation summary**



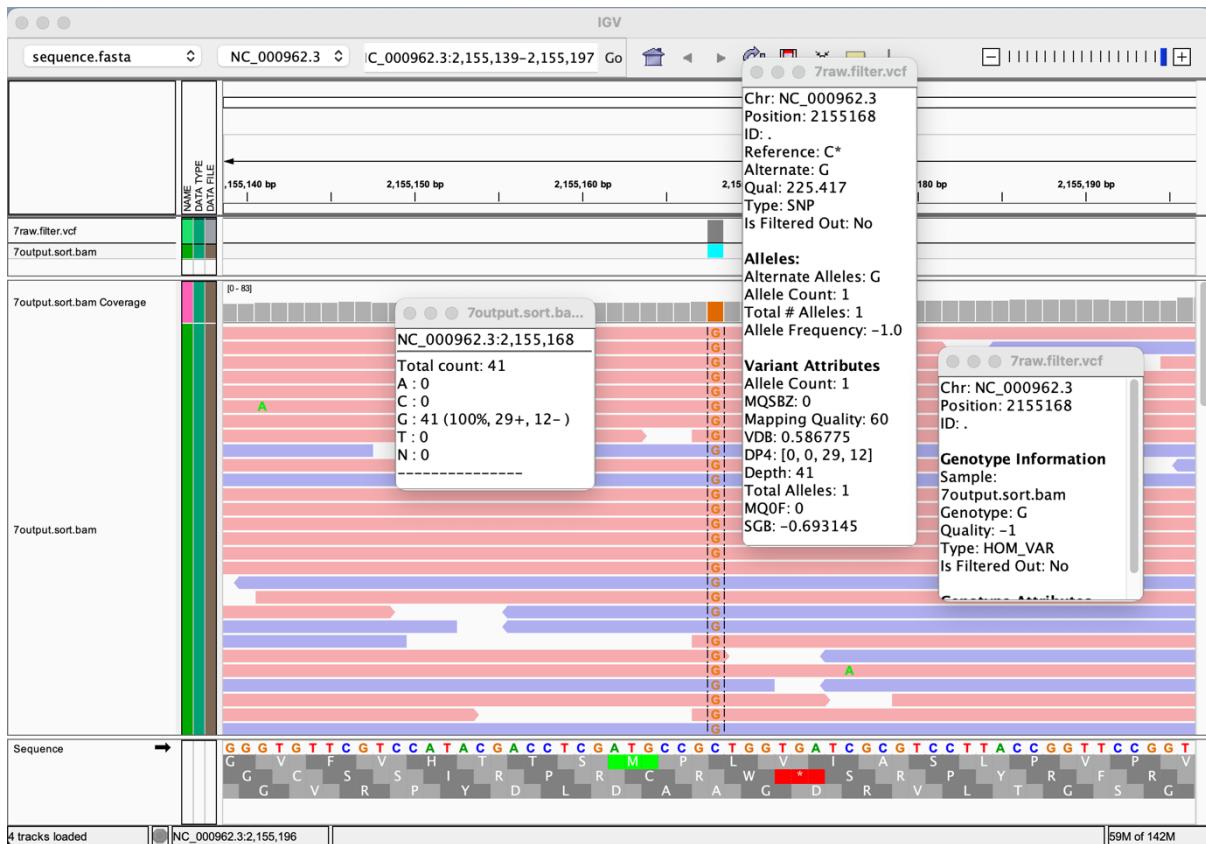
**Fig 6.3** Variant effect summary.



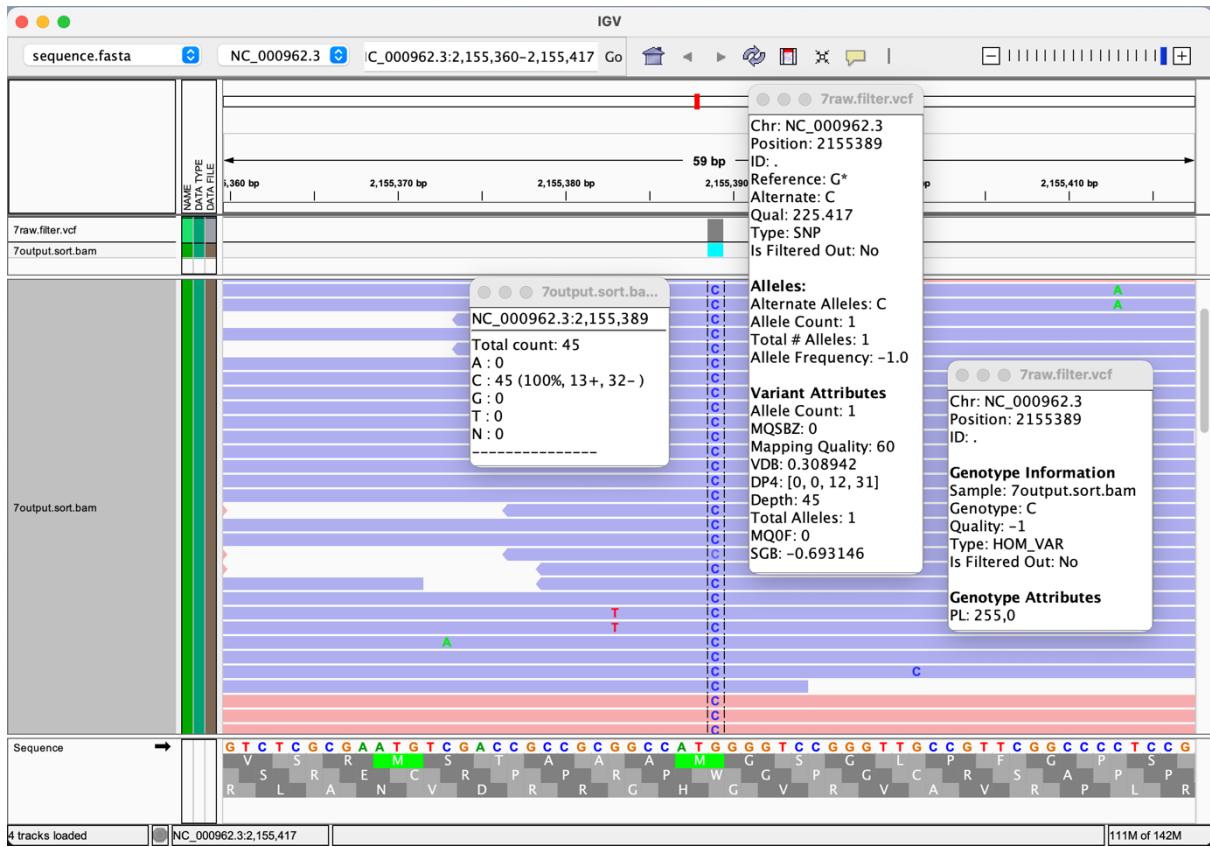
**Fig 6.4** The Number of effects by type and region



**Fig 7.0** IGV visualisation of a Mycobacterium tuberculosis H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 7



**Fig 7.0.1**



**Fig 7.0.2**

Formatting as Table Styles															
Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To preserve these features, save it in an Excel file format.															
Office Update To keep up to date with security updates, fixes and improvements, choose Check for Updates.															
A53	B	C	D	E	F	G		H							
534	T	G	225.417	D <small>P=57 VDB=-0.99844 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.21,34 MQ=60 ANN=G upstream_gene_variant MODIFIER </small>	I <small>V=1869c transcript </small>	R <small>v=1869c protein_coding </small>	C <small>-24 </small>								
535	A	G	225.417	D <small>P=7 VDB=-0.117704 G SGB=0.693139 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.18,18 MQ=60 ANN=G synonymous_variant </small>	I <small>VDB=0 transcript </small>	R <small>v=1870c protein_coding </small>	C <small>1 1 c-849 A </small>	L <small>u </small>							
536	TTCGGAT	TTCGGATGCC	211.067	D <small>DEL DV=28 IMF=0.529302 NP=53 VDB=-0.032079 G SGB=-0.693054 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 2.142,2.200 MQ=60 ANN=TTCGGAT </small>	I <small>VDB=0 transcript </small>	R <small>v=1871c protein_coding </small>	C <small>-2,1983 SCBZ=-2,1983 MQ=60 ANN=TTCGGAT </small>								
537	T	C	225.417	D <small>P=31 VDB=-0.896255 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.25,27 MQ=60 ANN=G upstream_gene_variant MODIFIER </small>	I <small>V=140 transcript </small>	R <small>v=1880c protein_coding </small>	C <small> 1 c-401</small>								
538	CTGTAATCCT	228.325	D <small>DEL DV=42 IMF=0.823529 NP=51 VDB=-0.038216 G SGB=0.693146 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 1.20583 SCBZ=6,55022 MQ=60 ANN=CT upstream_gene_variant MODIFIER </small>	I <small>V=140 transcript </small>	R <small>v=1884c protein_coding </small>	C <small>1,20583 SCBZ=6,55022 MQ=60 ANN=CT upstream_gene_variant MODIFIER </small>									
539	G	C	225.417	D <small>P=52 VDB=-0.057074 G SGB=0.693145 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.13,28 MQ=60 ANN=G missense_variant MODERATE </small>	I <small>V=1895c transcript </small>	R <small>v=1895c protein_coding </small>	C <small> 1 c-808 G </small>								
540	C	T	225.417	D <small>P=64 VDB=-0.106659 G SGB=0.693146 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.35,28 MQ=60 ANN=G missense_variant MODERATE </small>	I <small>V=1896c transcript </small>	R <small>v=1896c protein_coding </small>	C <small> 1 c-489 G p Trp16 </small>								
541	A	C	228.398	D <small>P=82 VDB=-0.107693 G SGB=0.693147 MQ582=1 MQOF=0 AC=1 AN=1 DP4=0 1.12,14 MQ=60 ANN=C missense_variant MODERATE </small>	I <small>V=1900c transcript </small>	R <small>v=1900c protein_coding </small>	C <small> 1 c-109 G p Trp17 </small>								
542	C	G	225.417	D <small>P=41 VDB=-0.586745 G SGB=0.693147 MQ582=1 MQOF=0 AC=1 AN=1 DP4=0 0.29,12 MQ=60 ANN=C missense_variant MODERATE </small>	I <small>V=1906c transcript </small>	R <small>v=1906c protein_coding </small>	C <small> 1 c-944 G p Trp18 </small>								
543	G	C	225.417	D <small>P=45 VDB=-0.878942 G SGB=0.693146 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.12,31 MQ=60 ANN=C missense_variant MODERATE </small>	I <small>V=1908c transcript </small>	R <small>v=1908c protein_coding </small>	C <small> 1 c-732 G p Trp19 </small>								
544	C	T	225.417	D <small>P=19 VDB=-0.692376 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.35,28 MQ=60 ANN=C missense_variant MODERATE </small>	I <small>V=1912c transcript </small>	R <small>v=1912c protein_coding </small>	C <small> 1 c-765 G p Trp20 </small>								
545	A	G	143.416	D <small>P=5 VDB=-0.118298 G SGB=0.693097 MQ582=1 MQOF=0 AC=1 AN=1 DP4=0 0.63,63 MQ=37 ANN=G synonymous_variant </small>	I <small>V=1914c transcript </small>	R <small>v=1914c protein_coding </small>	C <small> 1 c-10 G </small>								
546	G	C	152.416	D <small>P=9 VDB=-0.124207 G SGB=0.662043 MQ582=1 MQOF=0 AC=1 AN=1 DP4=0 0.63,63 MQ=37 ANN=C missense_variant MODERATE </small>	I <small>V=1915c transcript </small>	R <small>v=1915c protein_coding </small>	C <small> 1 c-10 G </small>								
547	A	C	155.416	D <small>P=5 VDB=-0.120245 G SGB=0.662043 MQ582=1 MQOF=0 AC=1 AN=1 DP4=0 0.63,63 MQ=37 ANN=C missense_variant MODERATE </small>	I <small>V=1916c transcript </small>	R <small>v=1916c protein_coding </small>	C <small> 1 c-10 G </small>								
548	G	A	225.417	D <small>P=24 VDB=-0.096133 G SGB=0.692831 MQ582=2 2.2287 MQOF=0 AC=1 AN=1 DP4=0 0.42,49 MQ=60 ANN=C synonymous_variant LOW </small>	I <small>V=1917c transcript </small>	R <small>v=1917c protein_coding </small>	C <small> 1 c-99 G p Trp21 </small>								
549	T	C	225.417	D <small>P=26 VDB=-0.064417 G SGB=0.692831 MQ582=2 1.9267 MQOF=0 AC=1 AN=1 DP4=0 0.11,15 MQ=52 ANN=C missense_variant MODERATE </small>	I <small>V=1917c transcript </small>	R <small>v=1917c protein_coding </small>	C <small> 1 c-99 G p Trp22 </small>								
550	A	C	225.417	D <small>P=28 VDB=-0.043741 G SGB=0.693054 MQ582=2 1.69186 MQOF=0 AC=1 AN=1 DP4=0 0.13,15 MQ=52 ANN=C missense_variant MODERATE </small>	I <small>V=1917c transcript </small>	R <small>v=1917c protein_coding </small>	C <small> 1 c-765 G p Trp23 </small>								
551	G	A	225.417	D <small>P=30 VDB=-0.025514 G SGB=0.693097 MQ582=2 1.83889 MQOF=0 AC=1 AN=1 DP4=0 0.13,17 MQ=53 ANN=G synonymous_variant LOW </small>	I <small>V=1917c transcript </small>	R <small>v=1917c protein_coding </small>	C <small> 1 c-10 G </small>								
552	A	C	225.417	D <small>P=20 VDB=-0.010609 G SGB=0.69168 MQ582=1 MQOF=0 AC=1 AN=1 DP4=0 0.63,63 MQ=37 ANN=C missense_variant MODERATE </small>	I <small>V=1917c transcript </small>	R <small>v=1917c protein_coding </small>	C <small> 1 c-352 G </small>								
553	A	C	225.417	D <small>P=41 VDB=-0.785384 G SGB=0.693145 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.16,24 MQ=60 ANN=C missense_variant MODERATE </small>	I <small>V=1917c transcript </small>	R <small>v=1917c protein_coding </small>	C <small> 1 c-2026 G </small>								
554	T	A	139.032	D <small>P=70 VDB=-0.806051 G SGB=0.693147 MQ582=1 MQOF=0 AC=1 AN=1 DP4=0 0.63,63 MQ=37 ANN=C missense_variant MODERATE </small>	I <small>V=1917c transcript </small>	R <small>v=1917c protein_coding </small>	C <small> 1 c-27,23 G </small>								
555	A	G	225.417	D <small>P=64 VDB=-0.875835 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.22,35 MQ=60 ANN=G missense_variant MODERATE </small>	I <small>V=1923 transcript </small>	R <small>v=1923 protein_coding </small>	C <small> 1 c-29 G </small>								
556	AG	A	228.415	D <small>INDEL DV=IMF=0.967033 NP=91 VDB=-0.542119 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.34,42 MQ=60 ANN=A frameshift </small>	I <small>V=1924 transcript </small>	R <small>v=1924 protein_coding </small>	C <small> 1 c-10 G </small>								
557	T	G	225.417	D <small>P=80 VDB=-0.632202 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.44,32 MQ=60 ANN=G synonymous_variant LOW </small>	I <small>V=1932 transcript </small>	R <small>v=1932 protein_coding </small>	C <small> 1 c-186 G </small>								
558	G	A	225.417	D <small>P=70 VDB=-0.803183 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.31,33 MQ=60 ANN=A missense_variant MODERATE </small>	I <small>V=1941 transcript </small>	R <small>v=1941 protein_coding </small>	C <small> 1 c-241 G </small>								
559	T	TC	228.422	D <small>P=47 VDB=-0.301283 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.31,36 MQ=60 ANN=G synonymous_variant LOW </small>	I <small>V=1941 transcript </small>	R <small>v=1941 protein_coding </small>	C <small> 1 c-149,26 MQ=60 ANN=G </small>								
560	A	G	225.417	D <small>P=46 VDB=-0.091211 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.22,22 MQ=60 ANN=A missense_variant MODERATE </small>	I <small>V=1941 transcript </small>	R <small>v=1941 protein_coding </small>	C <small> 1 c-118 G </small>								
561	C	A	225.417	D <small>P=29 VDB=-0.222118 G SGB=0.693054 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.16,12 MQ=60 ANN=G synonymous_variant LOW </small>	I <small>V=1947 transcript </small>	R <small>v=1947 protein_coding </small>	C <small> 1 c-759 G </small>								
562	T	G	225.421	D <small>P=22 VDB=-0.222107 G SGB=0.693054 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.16,12 MQ=60 ANN=G synonymous_variant LOW </small>	I <small>V=1947 transcript </small>	R <small>v=1947 protein_coding </small>	C <small> 1 c-759 G </small>								
563	CCGGCG	CCGGCGGGCG	228.417	D <small>INDEL DV=43 IMF=0.95556 NP=45 VDB=-0.10916 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 1.1,11.32 MQ=60 ANN=A CGGGCG CGGGCG </small>	I <small>V=1947 transcript </small>	R <small>v=1947 protein_coding </small>	C <small> 1 c-40 G </small>								
564	T	G	228.417	D <small>P=64 VDB=-0.066495 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.35,26 MQ=60 ANN=G upstream_gene_variant MODIFIER </small>	I <small>V=1976 transcript </small>	R <small>v=1976 protein_coding </small>	C <small> 1 c-40 G </small>								
565	A	G	225.421	D <small>P=66 VDB=-0.628435 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.24,40 MQ=60 ANN=G upstream_gene_variant MODIFIER </small>	I <small>V=1980 transcript </small>	R <small>v=1980 protein_coding </small>	C <small> 1 c-4938 G </small>								
566	T	C	225.417	D <small>P=40 VDB=-0.622336 G SGB=0.693141 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.24,13 MQ=60 ANN=G upstream_gene_variant MODIFIER </small>	I <small>V=1985 transcript </small>	R <small>v=1985 protein_coding </small>	C <small> 1 c-40 G </small>								
567	A	G	225.417	D <small>P=63 VDB=-0.433889 G SGB=0.693147 MQ582=0 MQOF=0 AC=1 AN=1 DP4=0 0.28,31 MQ=60 ANN=G upstream_gene_variant MODIFIER </small>	I <small>V=2003 transcript </small>	R <small>v=2003 protein_coding </small>	C <small> 1 c-25 G </small>								

**Fig 7.1 A VCF annotated with SnpEff**

Summary	
Genome	NC_000962.3
Date	2024-08-08 18:01
SnpEff version	SnpEff 5.2c (build 2024-04-09 12:24), by Pablo Cingolani
Command line arguments	SnpEff NC_000962.3 7raw.filter.vcf
Warnings	3,059
Errors	0
Number of lines (input file)	995
Number of variants (before filter)	995
Number of non-variants (i.e. reference equals alternative)	0
Number of variants processed (i.e. after filter and non-variants)	995
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 annotations	9,496
Genome total length	4,411,532
Genome effective length	4,411,532
Variant rate	1 variant every 4,433 bases

Variants rate details	
Chromosome	Length
NC_000962.3	4,411,532
Total	4,411,532
	995
	4,433

Fig 7.2 SnpEff annotation summary



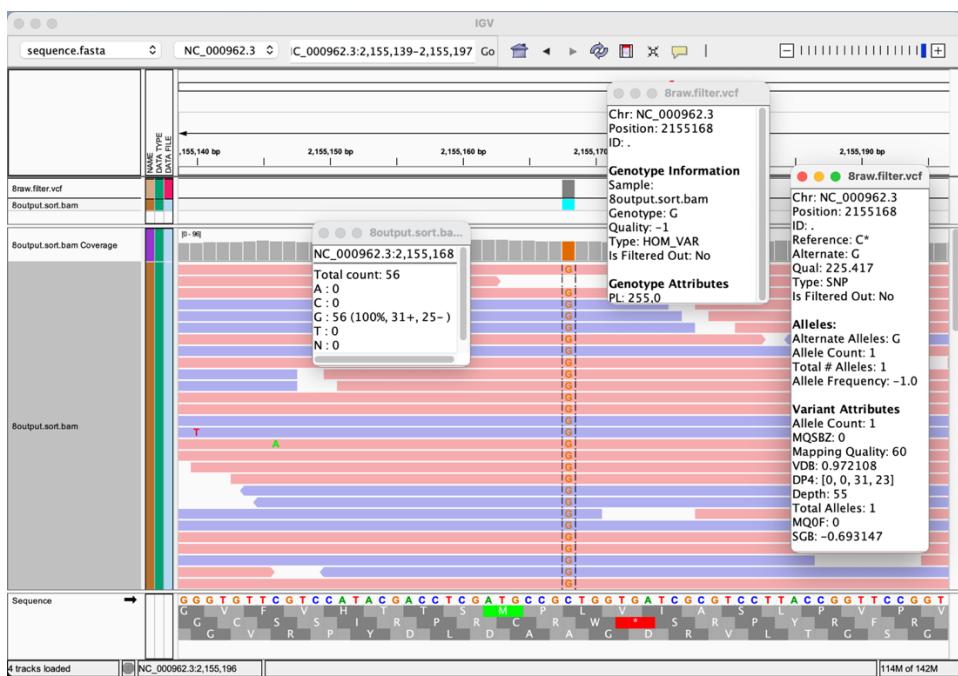
Fig 7.3 Variant effect summary

Number of annotations and region counts		
Annotation	Region	
Type (alphabetical order)	Count	Percent
conservative_inframe_deletion	6	0.063%
conservative_inframe_insertion	10	0.105%
disruptive_inframe_deletion	2	0.021%
disruptive_inframe_insertion	4	0.042%
downstream_gene_variant	4,302	45.294%
frameshift_variant	34	0.358%
intergenic_region	119	1.253%
intragenic_variant	42	0.442%
missense_variant	443	4.664%
splice_region_variant	1	0.011%
start_lost	2	0.021%
stop_gained	7	0.074%
stop_lost	1	0.011%
synonymous_variant	329	3.464%
upstream_gene_variant	4,196	44.178%
Type (alphabetical order)	Count	Percent
DOWNSTREAM	4,302	45.303%
EXON	837	8.814%
INTERGENIC	119	1.253%
TRANSCRIPT	42	0.442%
UPSTREAM	4,196	44.187%

Fig 7.4 The Number of effects by type and region



**Fig 8.0** IGV visualisation of a *Mycobacterium tuberculosis* H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 8



**Fig 8.0.1**

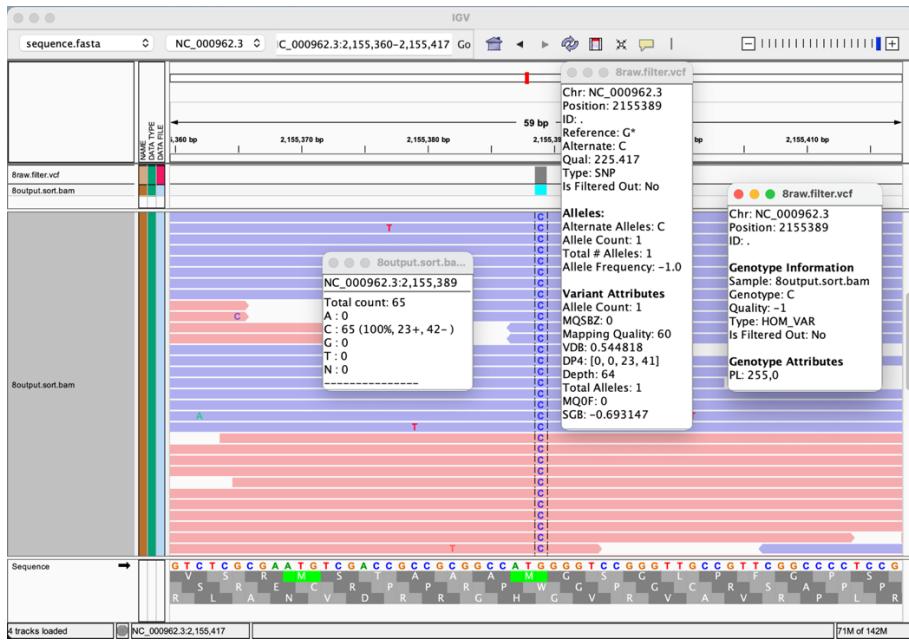


Fig 8.0.2

	B	C	D	E	F	G	H
505.	2096186.	A	G	225.417..	DP=58;VDB=0.560957;MQ0F=0.693147;MQ5B2=0.693147;MQ0B=0.AC>1;AN>1;DP4=0.33,21;MQ=60;ANN>1 synonymous_variant LOW bl Rv1846c transcript Rv1846c protein_coding 1/1;c.414T>c;p.Thr138Ile upstream_gene_variant MODERATE		
506.	2105923.	CGG	CGGG	228.391..	INDEL;DV=51;IMF=0.910714;DP=56;VD8=0.0883392;SGB=0.693147;RPB2=2.20004;MQB2=0;MQ5B2=0;QBZ=1.76025;SCB2=6.56409;MQ0F=0;AC>1;AN>1;DP4=0,22,29;MQ=60;ANN>1 upstream_gene_variant MODERATE		
507.	C	A	225.417..	DP=54;VDB=0.992609;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,30,32;MQ=60;ANN>1 upstream_gene_variant MODERATE			
508.	2109903.	C	T	225.417..	DP=54;VDB=0.992609;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,30,32;MQ=60;ANN>1 upstream_gene_variant MODERATE		
509.	2122390.	G	T	225.417..	DP=63;VDB=0.181996;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,31,35;MQ=60;ANN>1 upstream_gene_variant MODERATE		
510.	2123168.	T	G	225.417..	DP=70;VDB=0.041056;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,35,35;MQ=60;ANN>1 upstream_gene_variant MODERATE		
511.	2123169.	T	G	225.417..	DP=70;VDB=0.044287;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,35,35;MQ=60;ANN>1 upstream_gene_variant MODERATE		
512.	2128870.	A	G	225.417..	DP=51;VDB=0.016473;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,22,29;MQ=60;ANN>1 synonymous_variant LOW gfa3 Rv1878 transcript Rv1878 protein_coding 1/1;c.8494A>c p.L284Q		
513.	2133468.	TTGGCAT	TTGGCATGCC	226.15..	INDEL;DV=37;IMF=0.587302;DP=63;VD8=2.1095;MQB2=0;MQ5B2=1.85861;MQ0F=0;AC>1;AN>1;DP4=0,17,19,18;MQ=60;ANN>1 TTGGCAT		
514.	2135870.	T	C	225.417..	DP=62;VDB=1;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,24,36;MQ=60;ANN>1 upstream_gene_variant MODERATE		
515.	2138245.	CTGTGATATC	CT	228.411..	INDEL;DV=6;IMF=0.00142482;SGB=0.693147;RPB2=6.21095;MQB2=0;MQ5B2=0;QBZ=2.7568;SCB2=7.87298;MQ0F=0;AC>1;AN>1;DP4=0,11,28,33;MQ=60;ANN>1 upstream_gene_variant MODERATE		
516.	2143328.	G	C	225.417..	DP=49;VDB=0.782512;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,13,32;MQ=60;ANN>1 missense_variant MODERATE		
517.	2147022.	T	C	225.417..	DP=80;VDB=0.048497;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,13,32;MQ=60;ANN>1 missense_variant MODERATE		
518.	2147022.	A	C	225.417..	DP=74;VDB=0.174046;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,31,40;MQ=60;ANN>1 missense_variant MODERATE		
519.	2155168.	C	G	225.417..	DP=55;VDB=0.977108;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,31,32;MQ=60;ANN>1 missense_variant MODERATE		
520.	2155389.	G	C	225.417..	DP=64;VDB=0.544818;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,23,41;MQ=60;ANN>1 synonymous_variant LOW gfa3 Rv1860c transcript Rv1860c protein_coding 1/1;c.723C>c p.P233L		
521.	2158327.	C	T	225.417..	DP=72;VDB=0.312828;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,31,40;MQ=60;ANN>1 synonymous_variant LOW fad3 Rv1912c transcript Rv1912c protein_coding 1/1;c.705G>a p.L205A		
522.	2163493.	A	G	199.416..	DP=14;VDB=0.161999..06;SGB=0.686358;MQS2=1.41505;MQ0F=0;AC>1;AN>1;DP4=0,31,31;MQ=44;ANN>1 synonymous_variant LOW fad3 Rv1917c transcript Rv1917c protein_coding 1/1;c.384G>t p.R128Q		
523.	2163494.	G	C	198.416..	DP=14;VDB=0.368446..06;SGB=0.686358;MQS2=1.41505;MQ0F=0;AC>1;AN>1;DP4=0,31,31;MQ=44;ANN>1 missense_variant MODERATE		
524.	2163496.	A	C	199.416..	DP=15;VDB=0.710116..07;SGB=0.686358;MQS2=1.44885;MQ0F=0;AC>1;AN>1;DP4=0,32,32;MQ=45;ANN>1 synonymous_variant LOW fad3 Rv1917c transcript Rv1917c protein_coding 1/1;c.384G>t p.R128Q		
525.	2165510.	G	A	225.417..	DP=41;VDB=0.48549..07;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,15,26;MQ=54;ANN>1 missense_variant MODERATE		
526.	2165510.	T	C	225.417..	DP=41;VDB=0.48549..07;SGB=0.693147;MQS2=0;MQ0F=0;AC>1;AN>1;DP4=0,15,26;MQ=54;ANN>1 missense_variant MODERATE		
527.	2165517.	A	C	225.417..	DP=44;VDB=0.306559..07;SGB=0.693146;MQS2=0;0.708003;MQ0F=0;AC>1;AN>1;DP4=0,18,25;MQ=56;ANN>1 missense_variant LOW PPE34 Rv1917c transcript Rv1917c protein_coding 1/1		
528.	2165520.	G	A	225.417..	DP=47;VDB=0.327283..07;SGB=0.693147;MQS2=0;0.635047;MQ0F=0;AC>1;AN>1;DP4=0,18,29;MQ=55;ANN>1 missense_variant LOW PPE34 Rv1917c transcript Rv1917c protein_coding 1/1		
529.	2165790.	A	C	225.417..	DP=36;VDB=0.217732..06;SGB=0.693149;MQS2=0;0.160422;MQ0F=0;AC>1;AN>1;DP4=0,20,26;MQ=58;ANN>1 missense_variant LOW PPE34 Rv1917c transcript Rv1917c protein_coding 1/1		
530.	2165286.	A	C	225.417..	DP=54;VDB=0.622543..06;SGB=0.693147;MQS2=0;0.635047;MQ0F=0;AC>1;AN>1;DP4=0,25,28;MQ=60;ANN>1 missense_variant MODERATE		
531.	2165503.	T	A	211.135..	DP=119;VDB=0.162658..05;SGB=0.693147;MQS2=0;1.41505;MQ0F=0;AC>1;AN>1;DP4=0,31,31;MQ=44;ANN>1 missense_variant MODERATE		
532.	2176101.	A	G	225.417..	DP=74;VDB=0.783844..06;SGB=0.693147;MQS2=0;1.41505;MQ0F=0;AC>1;AN>1;DP4=0,31,31;MQ=44;ANN>1 missense_variant MODERATE		
533.	2181805.	AG	A	228.407..	INDEL;DV=68;IMF=0.944444..07;SGB=0.693147;MQS2=0;1.0303;MQ0F=2.93247;MQB2=0.75289;MQ0B=0.00840336;AC>1;AN>1;DP4=0,11,44,46;MQ=44;ANN>1 synonymous_variant LOW mce3F Rv1923 transcript Rv1923 protein_coding 1/1;c.923A>a p.T126M		
534.	2193904.	T	G	225.417..	DP=60;VDB=0.909827..05;SGB=0.693147;MQS2=0;0.635047;MQ0F=0;AC>1;AN>1;DP4=0,35,35;MQ=60;ANN>1 missense_variant MODERATE		
535.	2193904.	G	A	225.417..	DP=60;VDB=0.909827..05;SGB=0.693147;MQS2=0;0.635047;MQ0F=0;AC>1;AN>1;DP4=0,35,35;MQ=60;ANN>1 missense_variant MODERATE		
536.	2207591.	T	TC	228.419..	INDEL;DV=80;IMF=0.97561;DP=82;D0=0.0339507;SGB=0.693147;RPB2=0;MQB2=0;MQ5B2=0;QBZ=2.36087;SCB2=0.158114;MQ0F=0;AC>1;AN>1;DP4=0,2,52,28;MQ=60;ANN>1 TC upstream_gene_variant LOW mec3G Rv1968 transcript Rv1968 protein_coding 1/1;c.201A>c p.L101P		
537.	2211826.	A	G	225.417..	DP=60;VDB=0.819904..05;SGB=0.693147;MQS2=0;0.635047;MQ0F=0;AC>1;AN>1;DP4=0,25,28;MQ=60;ANN>1 missense_variant MODERATE		
538.	2216443.	C	A	225.417..	DP=53;VDB=0.492568..05;SGB=0.693147;MQS2=0;0.635047;MQ0F=0;AC>1;AN>1;DP4=0,4,32,36;MQ=60;ANN>1 frameshift		

Fig 8.1 A VCF annotated with SnpEff

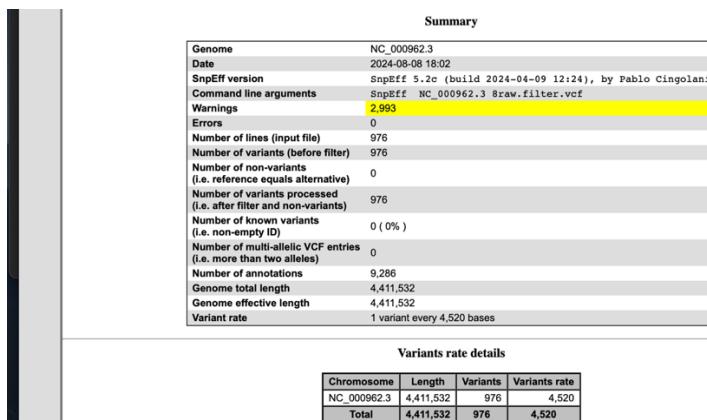
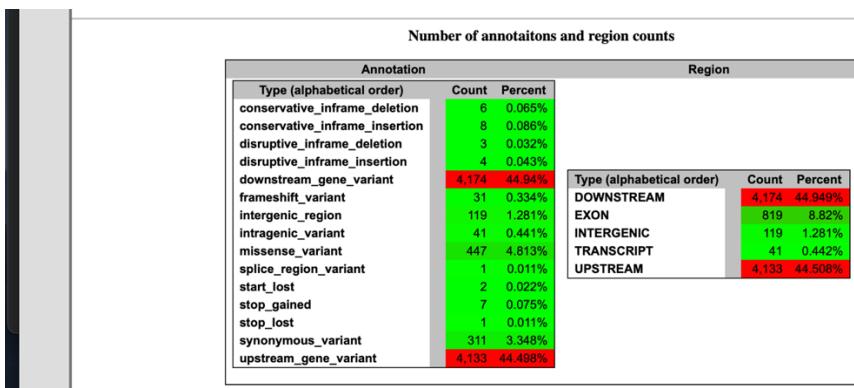


Fig 8.2 SnpEff annotation summary



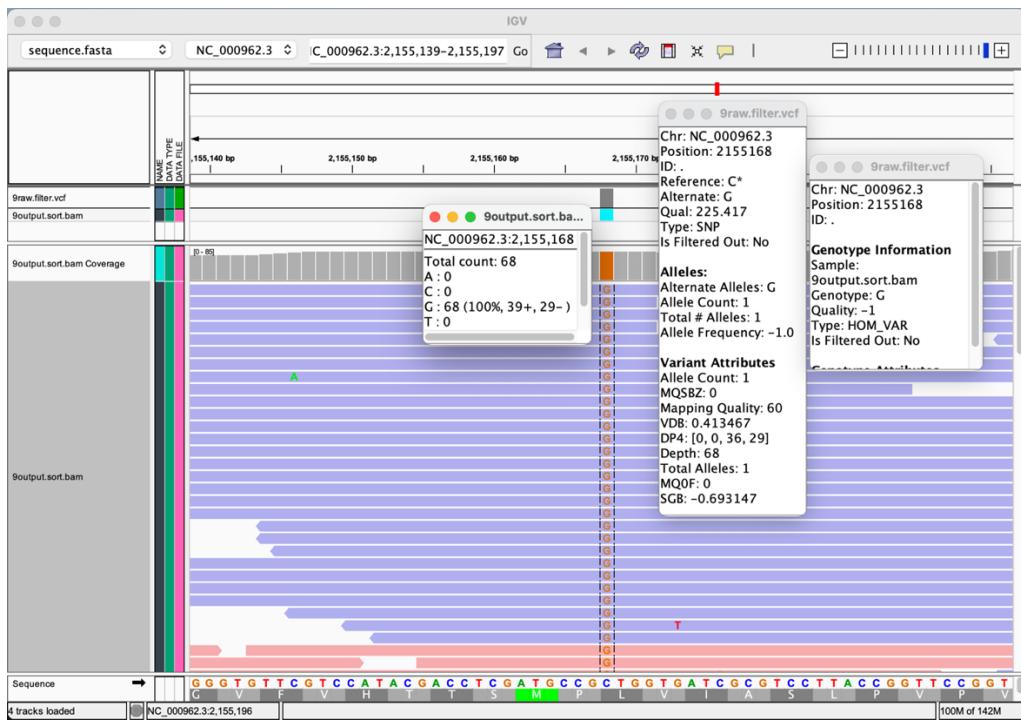
**Fig 8.3** Variant effect summary.



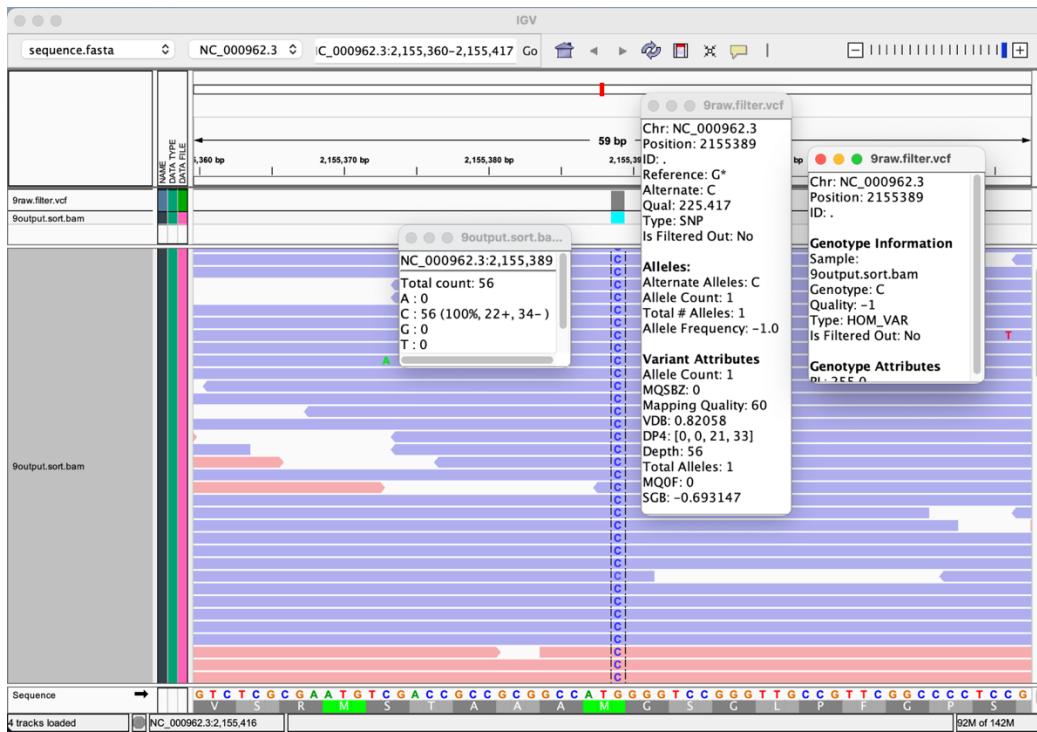
**Fig 8.4** The Number of effects by type and region



**Fig 9.0** IGV visualisation of a *Mycobacterium tuberculosis* H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 9



**Fig 9.0.1**



**Fig 9.0.2**

A	B	C	D	E	F	G	H
S20	NC_000962.3	2109523	.	CGG	CGGG	228.397	INDEL;ID=64;IM=0.941176;DP=68;VDB=0.379609;SGB=0.693147;RPBZ=-1.79936;MQBZ=0;MQSBZ=0;8QBZ=0.569672;5CBZ=-5.69891;MQOF=0;AC=1;AN=1;DP4=3,1,3,6,28;MQ=60;ANN=A upstream_gene_variant MODIFIER R
S21	NC_000962.	2113058	.	C	A	228.414	DP=65;VDB=0.992354;SGB=0.693147;RPBZ=-1.34237;MQBZ=0;MQSBZ=0;8QBZ=0.8082-1.91148;CBZ=0;MQOF=0;AC=1;AN=1;DP4=1,20,28,33;MQ=60;ANN=A upstream_gene_variant MODIFIER R
S22	NC_000962.	2116903	.	C	T	225.417	DP=76;VDB=0.767867;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,31,43;MQ=60;ANN=A synonymous_variant LOW RV1867 RV1867 protein_coding 1/c-
S23	NC_000962.	2122395	.	C	T	225.417	DP=54;VDB=0.512289;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,29,22;MQ=60;ANN=A missense_variant MODERATE I D2 RV1872 transcript RV1872 protein_coding 1/c-
S24	NC_000962.	2123168	.	T	G	228.398	DP=78;VDB=0.341513;SGB=0.693147;RPBZ=-1.68682;MQBZ=0;MQSBZ=0;8QBZ=1.514285;CBZ=0;MQOF=0;AC=1;AN=1;DP4=1,0,35,39;MQ=60;ANN=A upstream_gene_variant MODIFIER R
S25	NC_000962.	2123169	.	T	G	225.417	DP=78;VDB=0.414479;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,36,39;MQ=60;ANN=A upstream_gene_variant MODIFIER R RV1869c RV1869c transcript RV1869c protein_coding 1/c-
S26	NC_000962.	2128870	.	A	G	225.417	DP=78;VDB=0.032513;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,39,39;MQ=60;ANN=A synonymous_variant LOW pinA3 RV1878 transcript RV1878 protein_coding 1/c-
S27	NC_000962.	2133468	.	TTGGCAT	TTGGCATGCC	218.1	INDEL;ID=34;IM=0.557377;DP=61;VDB=0.0033375;MQBZ=0;8QBZ=0.31065;MQOF=0;AC=1;AN=1;DP4=15,12,25,9;MQ=60;ANN=A upstream_gene_variant MODIFIER R RV1880 transcript RV1880 protein_coding 1/c-
S28	NC_000962.	2135870	.	T	C	225.417	DP=88;VDB=0.996416;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,43,44;MQ=60;ANN=C upstream_gene_variant MODIFIER R cyt40 RV1880c transcript RV1880c protein_coding 1/c-
S29	NC_000962.	2138245	.	CTGGTAATC	CT	228.405	INDEL;ID=61;IM=0.968254;DP=63;VDB=0.80702;05;5GB=0.693147;RPBZ=-2.39329;MQBZ=0;MQSBZ=0;8QBZ=-3.02995;5CBZ=-5.52268;MQOF=0;AC=1;AN=1;DP4=1,24,37;MQ=60;ANN=A upstream_gene_variant MODIFIER R RV1881 transcript RV1881 protein_coding 1/c-
S30	NC_000962.	2143328	.	G	C	225.417	DP=67;VDB=0.840117;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,21,46;MQ=60;ANN=C missense_variant MODERATE I RV1895 RV1895 transcript RV1895 protein_coding 1/c-
S31	NC_000962.	2143958	.	C	T	225.422	DP=87;VDB=0.030396;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,44,43;MQ=60;ANN=A upstream_gene_variant MODIFIER R RV1896c RV1896c transcript RV1896c protein_coding 1/c-
S32	NC_000962.	2147022	.	A	C	225.417	DP=83;VDB=0.557001;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,42,37;MQ=60;ANN=C missense_variant MODERATE I RV1900c RV1900c transcript RV1900c protein_coding 1/c-
S33	NC_000962.	2155168	.	G	C	225.437	DP=68;VDB=0.413467;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,39,39;MQ=60;ANN=C missense_variant MODERATE R RV1808c RV1808c transcript RV1808c protein_coding 1/c-
S34	NC_000962.	2158327	.	C	T	225.437	DP=86;VDB=0.030521;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,33,33;MQ=60;ANN=C missense_variant MODERATE R RV1912c RV1912c transcript RV1912c protein_coding 1/c-
S35	NC_000962.	2163491	.	A	G	225.417	DP=17;VDB=0.0071793;SGB=0.688148;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,1,1;MQ=0;ANN=C synonymous_variant LOW pinA3 RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S36	NC_000962.	2163494	.	G	C	225.417	DP=16;VDB=0.0077709;SGB=0.688148;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,5,10;MQ=26;ANN=C missense_variant MODERATE R RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S37	NC_000962.	2163496	.	A	C	225.417	DP=18;VDB=0.0037093;SGB=0.690436;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,38,38;MQ=60;ANN=C missense_variant LOW pinA3 RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S38	NC_000962.	2163504	.	G	A	225.417	DP=46;VDB=0.0039414;SGB=0.693147;MQBZ=0;1.57889;MQOF=0;AC=1;AN=1;DP4=0,0,5,12;MQ=29;ANN=C synonymous_variant LOW pinA3 RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S39	NC_000962.	2163510	.	T	C	225.417	DP=46;VDB=0.118835;SGB=0.693147;MQBZ=0;0.313789;MQOF=0;AC=1;AN=1;DP4=0,0,20,25;MQ=52;ANN=C missense_variant MODERATE R RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S40	NC_000962.	2163517	.	A	C	225.417	DP=47;VDB=0.629476;SGB=0.693147;MQBZ=0;0.397747;MQOF=0;AC=1;AN=1;DP4=0,0,21,25;MQ=52;ANN=C synonymous_variant LOW PRE34 RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S41	NC_000962.	2163520	.	G	A	225.417	DP=48;VDB=0.400467;SGB=0.693147;MQBZ=0;0.328887;MQOF=0;AC=1;AN=1;DP4=0,0,21,26;MQ=52;ANN=A synonymous_variant LOW PRE34 RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S42	NC_000962.	2163790	.	A	C	228.404	DP=37;VDB=0.29081;0.85;SGB=0.693139;RPBZ=-0.375325;MQBZ=0.23895;MQOF=0.23895;MQSBZ=0.2346488;8QBZ=2.34126;5CBZ=0;MQOF=0;AC=1;AN=1;DP4=1,15,21;MQ=58;ANN=C synonymous_variant LOW PRE34 RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S43	NC_000962.	2165286	.	A	C	225.417	DP=68;VDB=0.524285;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,37,30;MQ=60;ANN=C missense_variant MODERATE R PE34 RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S44	NC_000962.	2165503	.	T	A	228.183	DP=141;VDB=0.145195;SGB=0.693147;RPBZ=-3.84247;MQBZ=2.83564;MQOF=0.348127;8QBZ=2.81704;5CBZ=0;MQOF=0.0143844;8QBZ=0;MQOF=0.145195;8QBZ=0;MQOF=0.145195;5CBZ=0;MQOF=0.145195;ANN=A synonymous_variant LOW PRE34 RV1917c RV1917c transcript RV1917c protein_coding 1/c-
S45	NC_000962.	2176101	.	A	G	225.417	DP=72;VDB=0.882796;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,28,39;MQ=60;ANN=G missense_variant MODERATE R RV1923 RV1923 transcript RV1923 protein_coding 1/c-
S46	NC_000962.	2181805	.	AG	A	228.381	INDEL;ID=56;IM=0.933333;DP=60;VDB=0.633592;5CBZ=0;1.55872;8QBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,4,40,36;MQ=60;ANN=A frameshift LOW PRE34 RV1931c RV1931c transcript RV1931c protein_coding 1/c-
S47	NC_000962.	2183054	.	T	G	225.422	DP=81;VDB=0.610442;SGB=0.693147;MQBZ=0;0.95505;8QBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,20,25;MQ=52;ANN=C missense_variant MODERATE R RV1931c RV1931c transcript RV1931c protein_coding 1/c-
S48	NC_000962.	2193904	.	G	A	225.417	DP=75;VDB=0.999359;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,38,38;MQ=60;ANN=G missense_variant MODERATE R RV1941 RV1941 transcript RV1941 protein_coding 1/c-
S49	NC_000962.	2207591	.	T	TC	228.413	INDEL;ID=85;IM=0.95505;DP=89;VDB=0.0279376;SGB=0.693147;RPBZ=-0.308524;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,33,33;MQ=60;ANN=G synonymous_variant LOW mce3 RV1968 RV1968 transcript RV1968 protein_coding 1/c-
S50	NC_000962.	2211826	.	A	G	225.417	DP=66;VDB=0.887205;SGB=0.693147;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,33,33;MQ=60;ANN=G synonymous_variant LOW mce3 RV1968 RV1968 transcript RV1968 protein_coding 1/c-
S51	NC_000962.	2216443	.	C	A	228.413	DP=60;VDB=0.831311;SGB=0.693147;RPBZ=-0.29233;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,1,20,23,34;MQ=60;ANN=A missense_variant MODERATE R RV1977 RV1977 transcript RV1977 protein_coding 1/c-
S52	NC_000962.	2220512	.	T	G	225.417	DP=39;VDB=0.958925;SGB=0.693144;MQBZ=0;MQOF=0;AC=1;AN=1;DP4=0,0,33,33;MQ=60;ANN=G missense_variant LOW RV1977 RV1977 transcript RV1977 protein_coding 1/c-

Fig 9.1 A VCF annotated with SnpEff.

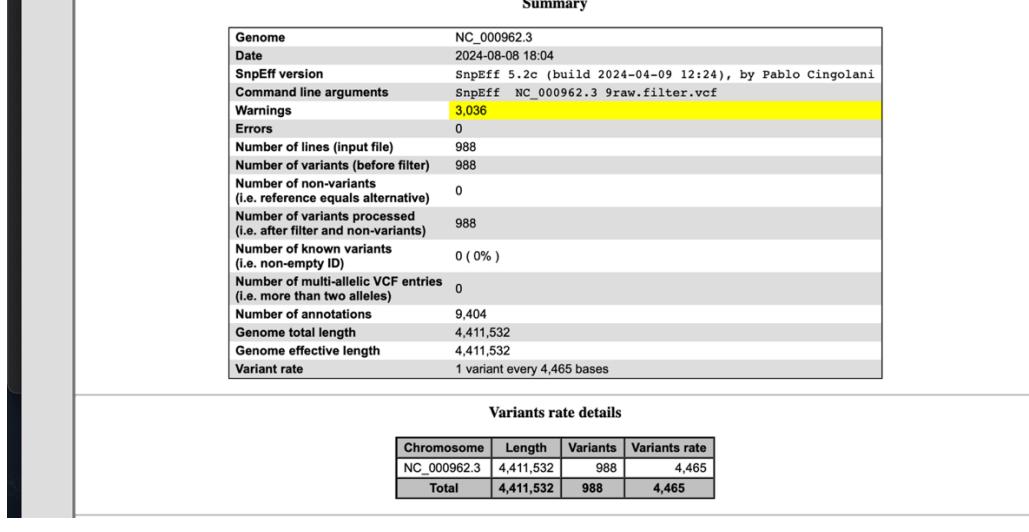
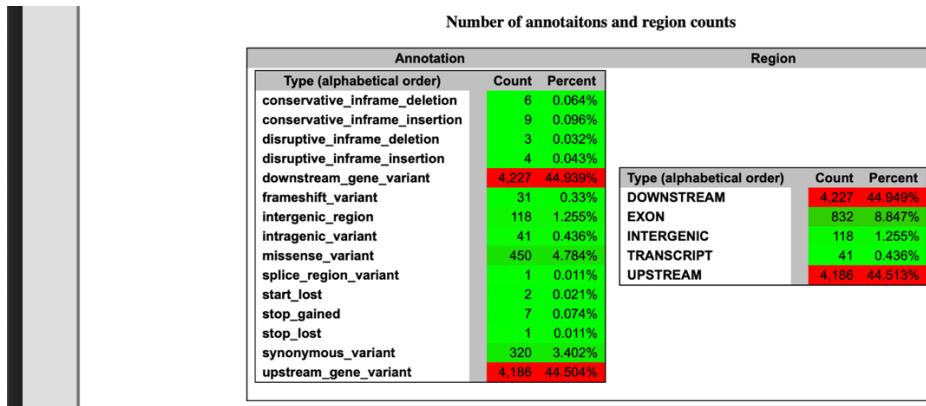


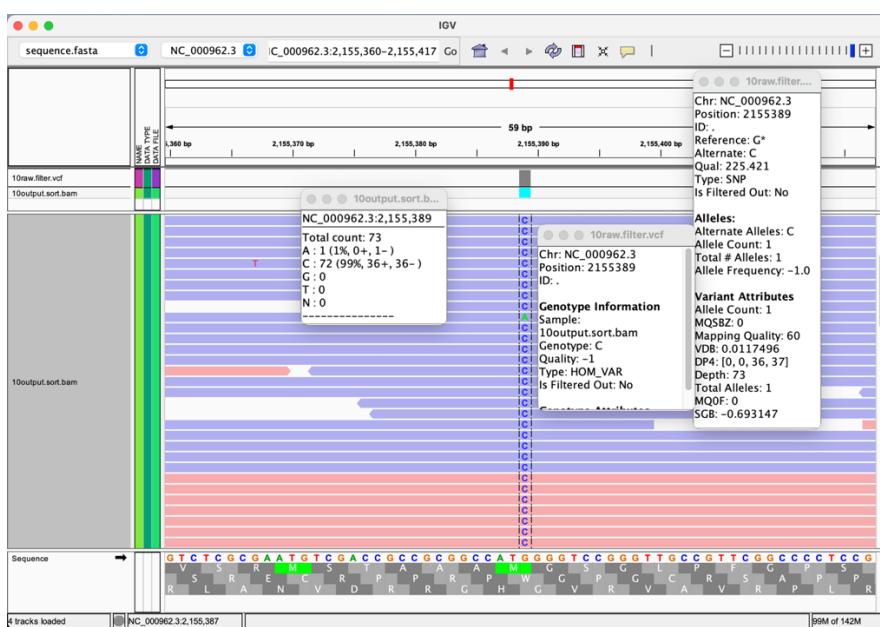
Fig 9.3 Variant effect summary.



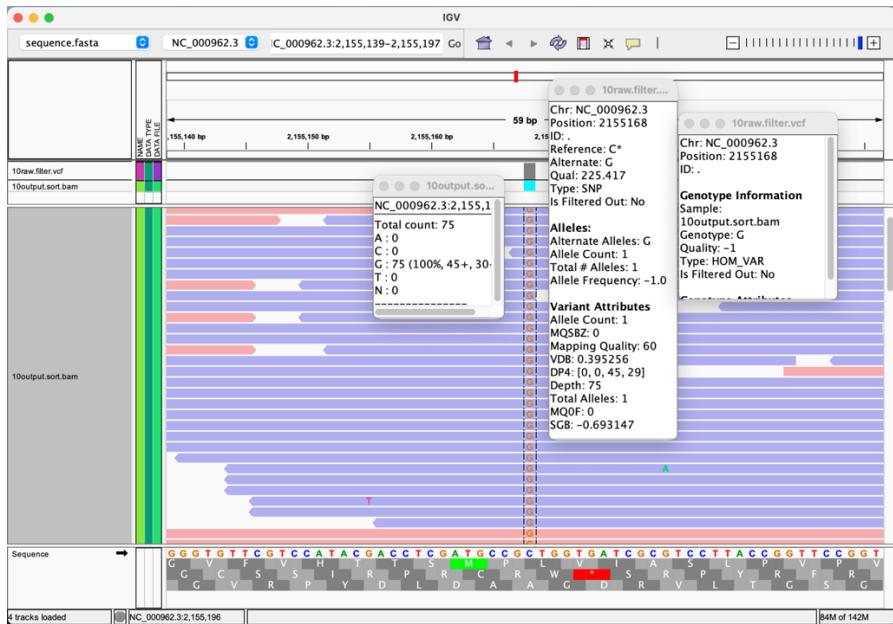
**Fig 9.4** The Number of effects by type and region



**Fig 10.0** IGV visualisation of a *Mycobacterium tuberculosis* H37RV variant site (with respect to the NC\_000962.3 reference) in a simulated sample 10



**Fig 10.0.1**



## 10.0.2

	A	C	D	E	G	
521	T	C	225.417.		DP=80 VDB=-0.99980 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,31,49 MQD=60 ANN=C ligand_receptor_gene, transcript RV2805 protein, coding [- c-40]	
522	T	G	225.417.		DP=80 VDB=-0.899608 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,24,34 MQD=60 ANN=N lncRNA, gene, variant MODIFIER RV2802 transcript RV2802 protein, coding [- c-45]	
523			CTGTAATGCT		228.376.	INDEL DVY-73 0.888889 DP=81 VDB=-1.02894 SGB=-0.693147 MQB=+0.87582 MQB2=+0.86523 -0.328753 MQBZ=0 MQOF=0 AC=1 AN=1 DP=0,34,38,34 MQD=60 ANN=N CT upstream_gene_variant MODIFIER RV2802 transcript RV2802 protein, coding [- c-45]
524	T	C	225.417.		DP=52 VDB=0.522584 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,21,31 MQD=60 ANN=C missense, variant MODERATE RV2895 protein, coding [- c c-208 c-6]	
525	T	C	225.417.		DP=85 VDB=0.452211 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,2,41,49 MQD=60 ANN=N stop gained HIGH RV2896 protein, coding [- c c-489 c-6 t p16]	
526	A	C	225.417.		DP=11 VDB=0.574046 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,47,46 MQD=60 ANN=C missense, variant MODERATE lip RV1900c transcript RV1900c protein, coding [- c c-6217	
527	C	G	225.417.		DP=75 VDB=0.395256 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,45,29 MQD=60 ANN=N stop gained HIGH RV2896 protein, coding [- c c-489 c-6 t p16]	
528	G	C	225.421.		DP=73 VDB=0.011749 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,43,37,37 MQD=60 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
529	T	C	225.417.		DP=59 VDB=0.452211 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,41,49 MQD=60 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
530	A	G	225.417.		DP=17 VDB=1.202, -0.56 SGB=-0.690438 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,6,15 MQD=41 ANN=4 stop gained MODIFIER RV2812 transcript RV2812 protein, coding [- c c-177 c-6 p16]	
531	G	C	225.417.		DP=17 VDB=0.488888 SGB=-0.75 MQD=1 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,117647 AC=1 AN=1 DP=0,6,11 MQD=41 ANN=4 stop gained MODIFIER RV2812 transcript RV2812 protein, coding [- c c-177 c-6 p16]	
532	A	C	225.417.		DP=18 VDB=0.484742 SGB=-0.65 MQD=0 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,2,11,32 MQD=60 ANN=C synonymous, variant MODERATE RV2895 protein, coding [- c c-208 c-6]	
533	G	A	225.417.		DP=46 VDB=0.116896 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,2,11,31 MQD=60 ANN=C synonymous, variant MODERATE RV2895 protein, coding [- c c-208 c-6]	
534	T	C	225.417.		DP=50 VDB=0.294079 SGB=-0.65 MQD=0 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,2,42,41 MQD=60 ANN=N stop gained HIGH RV2896 protein, coding [- c c-489 c-6 t p16]	
535	A	C	225.417.		DP=51 VDB=0.839356 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,17,33 MQD=53 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
536	G	A	225.417.		DP=52 VDB=0.839356 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,18,33 MQD=53 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
537	A	C	225.417.		DP=46 VDB=0.430846 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,15,31 MQD=53 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
538	T	C	225.417.		DP=59 VDB=0.452211 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,17,33 MQD=53 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
539	T	A	228.183.		DP=11 VDB=0.893065 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,6,15 MQD=41 ANN=4 stop gained MODIFIER RV2812 transcript RV2812 protein, coding [- c c-177 c-6 p16]	
540	A	G	225.417.		DP=84 VDB=0.879467 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,31,52 MQD=60 ANN=N missense, variant MODERATE lip RV2893 transcript RV2893 protein, coding [- c c-929 c-6 p16]	
541	AG	A	228.399.		INDEL DVY-106 IMF=0.921739 DP=0.97 VDB=-0.571175 SGB=-0.693147 RPB2=+0.497227 MQB=0 MQB2=0 MQOF=0 AC=1 AN=1 DP=0,1,8,29,53 MQD=60 ANN=C frame-shift, variant LOW RV2931 transcript RV2931 protein, coding [- c c-1386 c-6]	
542	T	G	225.417.		DP=79 VDB=0.932828 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,36,43 MQD=60 ANN=C synonymous, variant LOW RV2931 transcript RV2931 protein, coding [- c c-1386 c-6]	
543	G	A	225.417.		DP=90 VDB=0.972182 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,2,2443 MQD=0 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,1,8,38 MQD=53 ANN=C synonymous, variant LOW RV2931 transcript RV2931 protein, coding [- c c-1386 c-6]	
544	T	TC	225.417.		DP=51 VDB=0.839356 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,118,27 MQD=53 ANN=C synonymous, variant LOW RV2931 transcript RV2931 protein, coding [- c c-1386 c-6]	
545	A	G	225.417.		DP=52 VDB=0.839356 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,18,33 MQD=53 ANN=C synonymous, variant LOW RV2931 transcript RV2931 protein, coding [- c c-1386 c-6]	
546	T	TC	225.417.		DP=46 VDB=0.430846 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,18,33 MQD=53 ANN=C synonymous, variant LOW RV2931 transcript RV2931 protein, coding [- c c-1386 c-6]	
547	A	G	225.417.		DP=59 VDB=0.452211 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,17,33 MQD=53 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
548	G	A	225.417.		DP=49 VDB=0.891238 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,3,43 MQD=60 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
549	T	G	225.417.		DP=41 VDB=0.554365 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,17,23 MQD=60 ANN=C synonymous, variant MODERATE LOW kappa RV2808c transcript RV2808c protein, coding [- c c-9446 c-6 p16]	
550	A	G	228.42.		INDEL DVY-58 IMF=0.983051 DP=0.59 VDB=-0.050558 SGB=-0.693147 RPB2=+0.70776 MQB=0 MQB2=0 MQOF=0 AC=1 AN=1 DP=0,1,10,7,72,26 MQD=43 ANN=A synonymous, variant LOW RV2976 transcript RV2976 protein, coding [- c c-407 c-6 p16]	
551	T	C	225.417.		DP=10 VDB=0.400519 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,4,42,52 MQD=60 ANN=C upstream_gene_variant MODIFIER RV1976c transcript RV1976c protein, coding [- c c-407 c-6 p16]	
552	A	G	225.422.		DP=94 VDB=0.971441 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,4,42,52 MQD=60 ANN=C upstream_gene_variant MODIFIER RV1980c transcript RV1980c protein, coding [- c c-493 c-6 p16]	
553	T	C	225.417.		DP=68 VDB=0.922804 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,37,29 MQD=60 ANN=C upstream_gene_variant MODIFIER RV1985c transcript RV1985c protein, coding [- c c-493 c-6 p16]	
554	A	G	225.417.		DP=92 VDB=0.649636 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,42,44 MQD=60 ANN=C upstream_gene_variant MODIFIER RV2003c transcript RV2003c protein, coding [- c c-25 c-6 p16]	
555	C	A	225.417.		DP=80 VDB=0.661101 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,37,43 MQD=60 ANN=C missense, variant MODERATE stop RV2006 transcript RV2006 protein, coding [- c c-407 c-6 p16]	
556	T	C	225.417.		DP=77 VDB=0.996203 SGB=-0.693147 MQSB2=0 MQOF=0 AC=1 AN=1 DP=0,3,28,48 MQD=60 ANN=C upstream_gene_variant MODIFIER RV2007 transcript RV2007 protein, coding [- c c-407 c-6 p16]	

**Fig 10.1** A VCF annotated with SnpEff.

Summary						
Genome	NC_000962.3	Date	2024-08-08 18:04	SnpEff version	SnpEff 5.2c (build 2024-04-09 12:24), by Pablo Cingolani	Command line arguments
Warnings	3,065	Errors	0	Number of lines (input file)	995	Number of variants (before filter)
Number of non-variants	0	Number of variants processed (i.e. after filtering non-variants)	995	Number of known variants (i.e. non-empty ID)	0 (0 %)	Number of multi-allelic VCF entries (i.e. more than two alleles)
Number of annotations	9,486	Genome total length	4,411,532	Genome effective length	4,411,532	Variant rate
Chromosome	NC_000962.3	Length	4,411,532	Variants	995	Variants rate
Total	NC_000962.3		4,411,532		995	

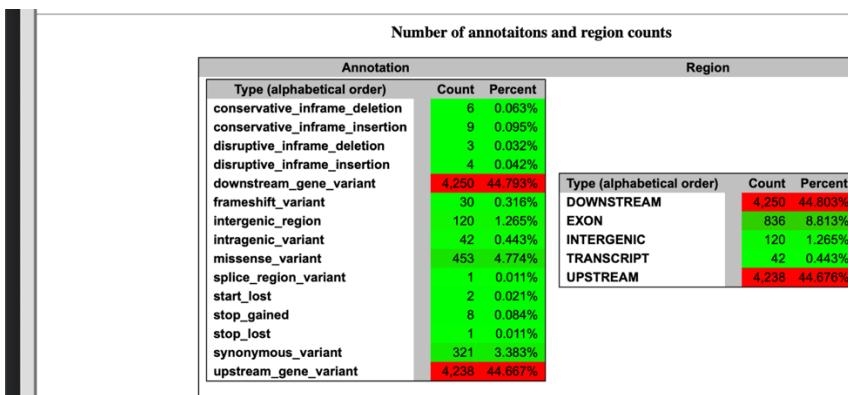
  

Variants rate details			
Chromosome	NC_000962.3	Length	4,411,532

**Fig 10.2** SnpEff annotation summary



**Fig 10.3** Variant effect summary.



**Fig 10.4** The Number of effects by type and region

## INTERPRETATIONS/DISCUSSION

*katG* mutations are prevalent in isoniazid resistance. Across all samples; IGV visualisation of mutations in *katG* sequences (Position: 2153889 - 2156111) and IGV browser gives colours to interesting events. Variants that affect a gene's function are the most harmful and

detrimental, yet they can occur anywhere along the genome sequence. A variation or mutation that impacts the gene regulatory region has the potential to repress, inhibit, or activate or deactivate a gene.

At Position: 2153889 – 2156111; Sample 1,3 and 4 showed no variants, While Variants were detected in Sample 2, 5,6,7,8,9, and 10, particularly homozygous variants, where NC\_000962.3 is the reference sequence in Genbank

For Sample 2 and Sample 5; The variant NC\_000962.3:c.944G>C p.(Ser315Thr), indicating a single nucleotide polymorphism(SNP), where a guanine (G) nucleotide at position 944 has been replaced by a cytosine (C) nucleotide. The p.(Ser315Thr) indicates the amino acid change caused by the SNP, In this case, a serine (Ser) amino acid at position 315 has been replaced by a threonine (Thr) amino acid. This is a missense mutation, as it results in a change in a single amino acid , causing Isoniazid resistance.

Sample 6, 7,8,9 and 10 produced two variants; The variant NC\_000962.3:c.944G>C p.(Ser315Thr) and The variant NC\_000962.3:c.723C>G p.(Pro241Pro). The former being a missense variation and the Latter being a synonymous variant (the variation doesn't alter the encoded amino acid).

The mutation occurs at a coding region causing a structural and conformational change and preventing it from binding to the target resulting hence resulting in Isoniazid deactivation and resistance. [Larsen et al, 2002].

Isoniazid resistance in *Mycobacterium tuberculosis* can arise from mutations in various genes and alterations in biological pathways. Understanding these pathways is crucial for developing new strategies to combat drug-resistant Tuberculosis. Some key pathways associated with Isoniazid resistance include; Pharamacokinetics/Pharmacodynamics factors, Drug efflux systems, Mycobacterial Cell wall, Drug misuse, generation or use of alternative biochemical pathways.

## CONCLUSION

Mutations were identified in 7 samples; with 5 samples having two variants. Isoniazid resistance is due to the mutation of katG gene --Ser315-Thr--. These mutations reduce its activity, making it difficult to treat Tuberculosis.

## REFERENCES

1. James T Robinson, Helga Thorvaldsdottir, Douglass Turner, Jill P Mesirov, igv.js: an embeddable JavaScript implementation of the Integrative Genomics Viewer (IGV), Bioinformatics, Volume 39, Issue 1, January 2023, btac830,  
<https://doi.org/10.1093/bioinformatics/btac830>
2. HackBio. (2023). Project Section: Searching for Mutations Leading to Isoniazid Resistance in TB (HackBio) [Data set]. Zenodo.  
<https://doi.org/10.5281/zenodo.10426436>
3. Hamid D. Ismail. Bioinformatics, A Practical Guide to Next Generation Sequencing Data Analysis. First edition published 2023 by CRC Press. Pages 110-160
4. Larsen, M. H., Vilchèze, C., Kremer, L., Besra, G. S., Parsons, L., Salfinger, M., Heifets, L., Hazbon, M. H., Alland, D., Sacchettini, J. C., & Jacobs, W. R., Jr (2002). Overexpression of inhA, but not kasA, confers resistance to isoniazid and ethionamide in *Mycobacterium smegmatis*, *M. bovis* BCG and *M. tuberculosis*. Molecular microbiology, 46(2), 453–466. <https://doi.org/10.1046/j.1365-2958.2002.03162.x>