

Summary Report

1 Introduction

This document provides summary for processing and filtering one raw VCF file (/home/brb/SeqTestdata/RNASeqFibroblast/outputhg38/LFB_scramble_repA_raw.vcf) as well as annotating the filtered VCF file through the Somatic Mutation Annotator through SnpEff in BRB-SeqTools. We generate the following files in the variant annotation process:

- A gene list (/home/brb/SeqTestdata/RNASeqFibroblast/outputhg38_cli_snpeff/LFB_scramble_repA_raw_genelist.txt) containing nonsynonymous and splicing variants which are not known polymorphisms unless in COSMIC.
- An annotation table (/home/brb/SeqTestdata/RNASeqFibroblast/outputhg38_cli_snpeff/LFB_scramble_repA_raw_annoTable.txt) for the detected variants.
- An annotated VCF file (/home/brb/SeqTestdata/RNASeqFibroblast/outputhg38_cli_snpeff/LFB_scramble_repA_raw_annotated.vcf) associated with the annotation table.

2 Variant Annotation Process

The raw VCF file is processed and filtered in the following steps:

1. We keep those variants that pass the criterion that the variant call quality $QUAL \geq 1$, the read depth $DP \geq 1$ and the mapping quality $MQ \geq 1$.
2. We decompose and left normalize the remaining variants.
3. We remove those variants reported in dbSNP database but keep those variants reported in COSMIC database.
4. Nonsynonymous and splicing variants are identified from the remaining variants for further analyses.
5. The remaining variants are annotated through SnpEff.
6. A gene list is retrieved for the variants through SnpEff, which may be a potential list related with the data of interest.

3 Summary Statistics

Table 1 summarizes the statistics related with the variant annotation process via SnpEff.

Table 1: Statistics summary associated with the variant annotation via SnpEff.

Statistics	Count
Total number of variants in the raw VCF file	451
Number of variants left after the filter by $QUAL \geq 1$, $DP \geq 1$, $MQ \geq 1$	451
Number of variants after decomposing and left normalization	451
Number of variants reported in dbSNP database	195
Number of variants reported in COSMIC database	61

Statistics	Count
Number of variants reported in both dbSNP and COSMIC database	50
Number of variants remaining after removing variants reported in dbSNP while keeping variants in COSMIC	306
Number of variants (out of 306 variants) that are nonsynonymous or splicing ones	124
Number of variants (out of 124 variants) that are reported in COSMIC	22
Number of genes associated with 124 variants	114

We also provide a statistics table for the nonsynonymous and splicing variants kept for annotation. Table 2 summarizes the effects the nonsynonymous variants have.

Table 2: Nonsynonymous and splicing variants after filtering.

Region	Effect	Count
Exonic	Frameshift	4
Exonic	Stop loss	0
Exonic	Stop gain	9
Exonic	Start loss	1
Exonic	Mis-sense	110
Splicing	/	0
Total	/	124

4 Charts

We summarize here statistics of gene annotations for 451 variants that pass the quality, read depth and mapping quality filtering criteria. We draw figures for the proportion of variants that hit different regions such as exonic and intronic regions as shown in Figure 1, and for the proportion of exonic with different functional effects (e.g., synonymous, nonsynonymous) as shown in Figure 2.

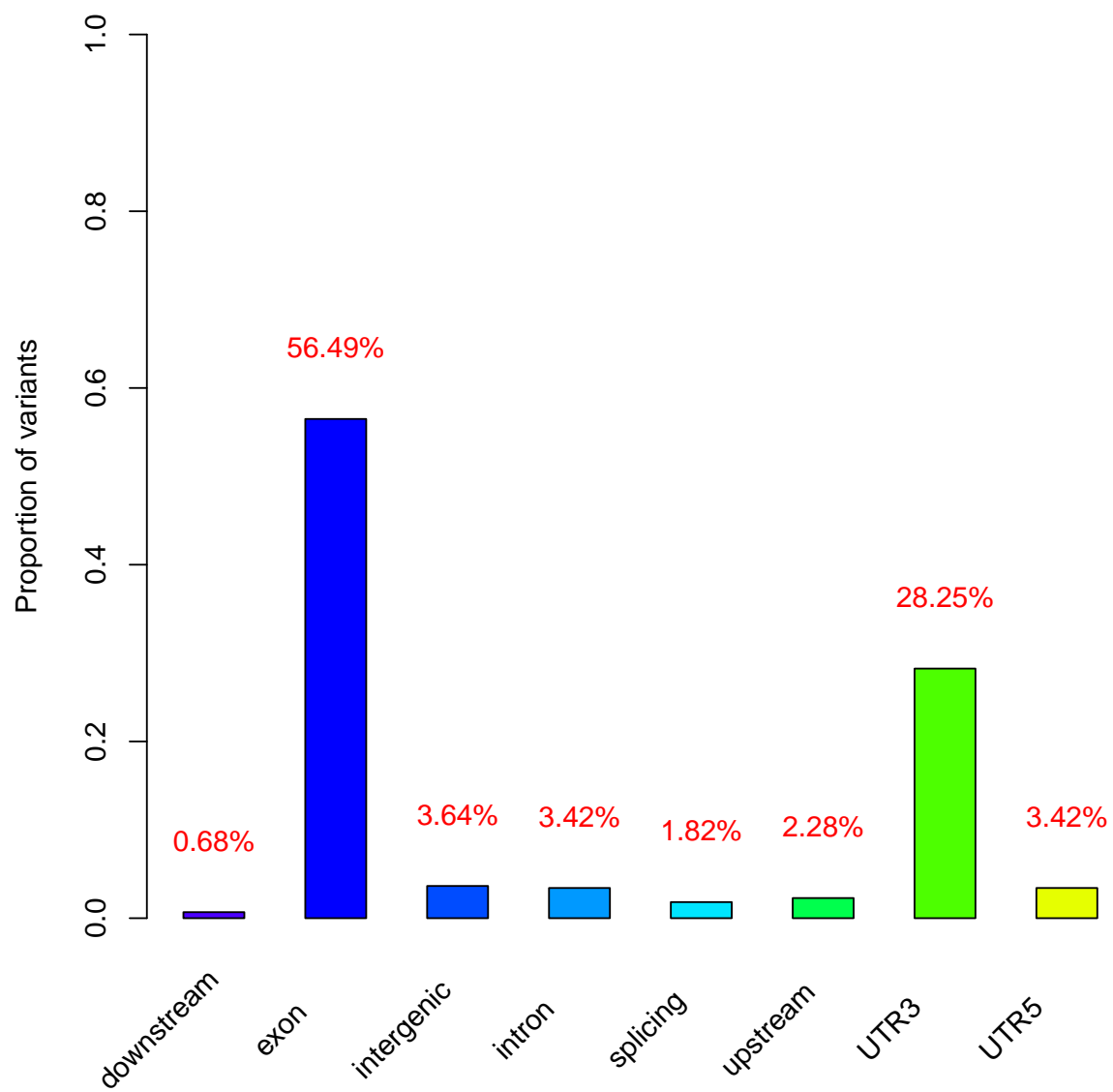


Figure 1: Proportion of variants that hit different regions.

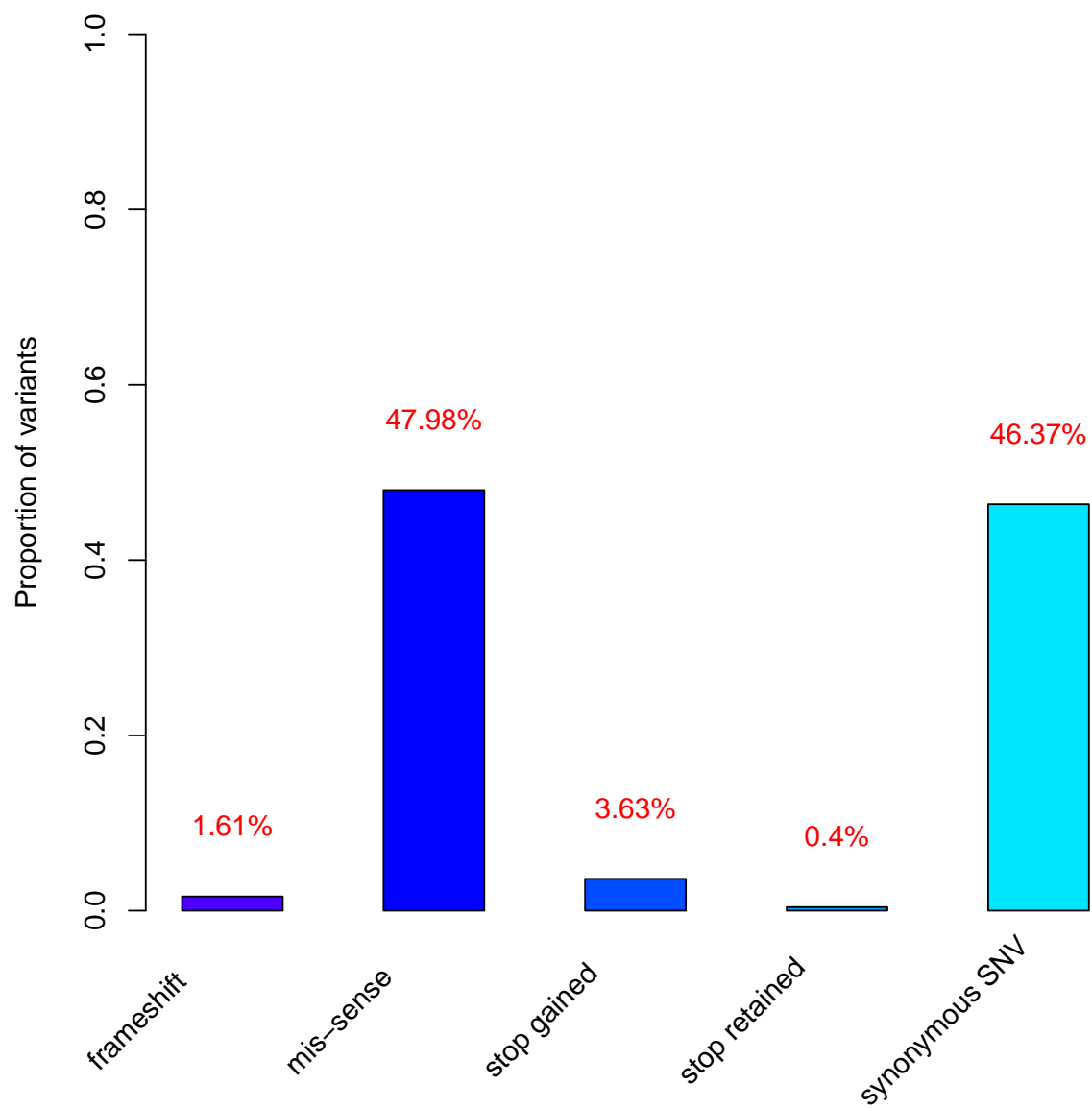


Figure 2: Proportion of exonic variants with their functional effects.