

# Summary Report

## 1 Introduction

This document provides summary for processing and filtering one raw VCF file (/home/brb/testdata/GSE48215subset/output/bt20\_raw.vcf) as well as annotating the filtered VCF file through the Somatic Mutation Annotator through ANNOVAR in BRB-SeqTools. We generate the following files in the variant annotation process:

- A gene list (/home/brb/testdata/GSE48215subset/output/annovar/bt20\_raw\_genelist.txt) containing nonsynonymous and splicing variants which are not known polymorphisms unless in COSMIC.
- An annotation table (/home/brb/testdata/GSE48215subset/output/annovar/bt20\_raw\_annoTable.txt) for the detected variants.
- An annotated VCF file (/home/brb/testdata/GSE48215subset/output/annovar/bt20\_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  $QUAL \geq 20$ , the read depth  $DP \geq 5$  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 ANNOVAR.
6. A gene list is retrieved for the variants through ANNOVAR, which may be a potential list related with the data of interest.

## 3 Summary Statistics

Table 1 summarizes the stastics related with the variant annotation process via ANNOVAR.

Table 1: Statistics summary associated witht the variant annotation via ANNOVAR.

Statistics	Count
Total number of variants in the raw VCF file	1610
Number of variants left after the filter $QUAL \geq 20$ , $DP \geq 5$ , $MQ \geq 1$	331
Number of variants remaining after removing variants reported in dbSNP while keeping variants in COSMIC	139
Number of variants (out of 139 variants) that are nonsynonymous or splicing ones	20
Number of variants (out of 20 variants) that are reported in COSMIC	18
Number of genes associated with 20 variants	17

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 deletion	0
Exonic	Frameshift insertion	0
Exonic	Stoploss	0
Exonic	Stopgain	0
Exonic	Mis-sense	20
Splicing	/	0
Total	/	20

## 4 Charts

We summarize here statistics of gene annotations for 139 variants that pass the quality, read depth and mapping quality filtering criteria. These variants are annotated by RefSeq, UCSC Known Gene and Ensembl gene annotation sources. 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.

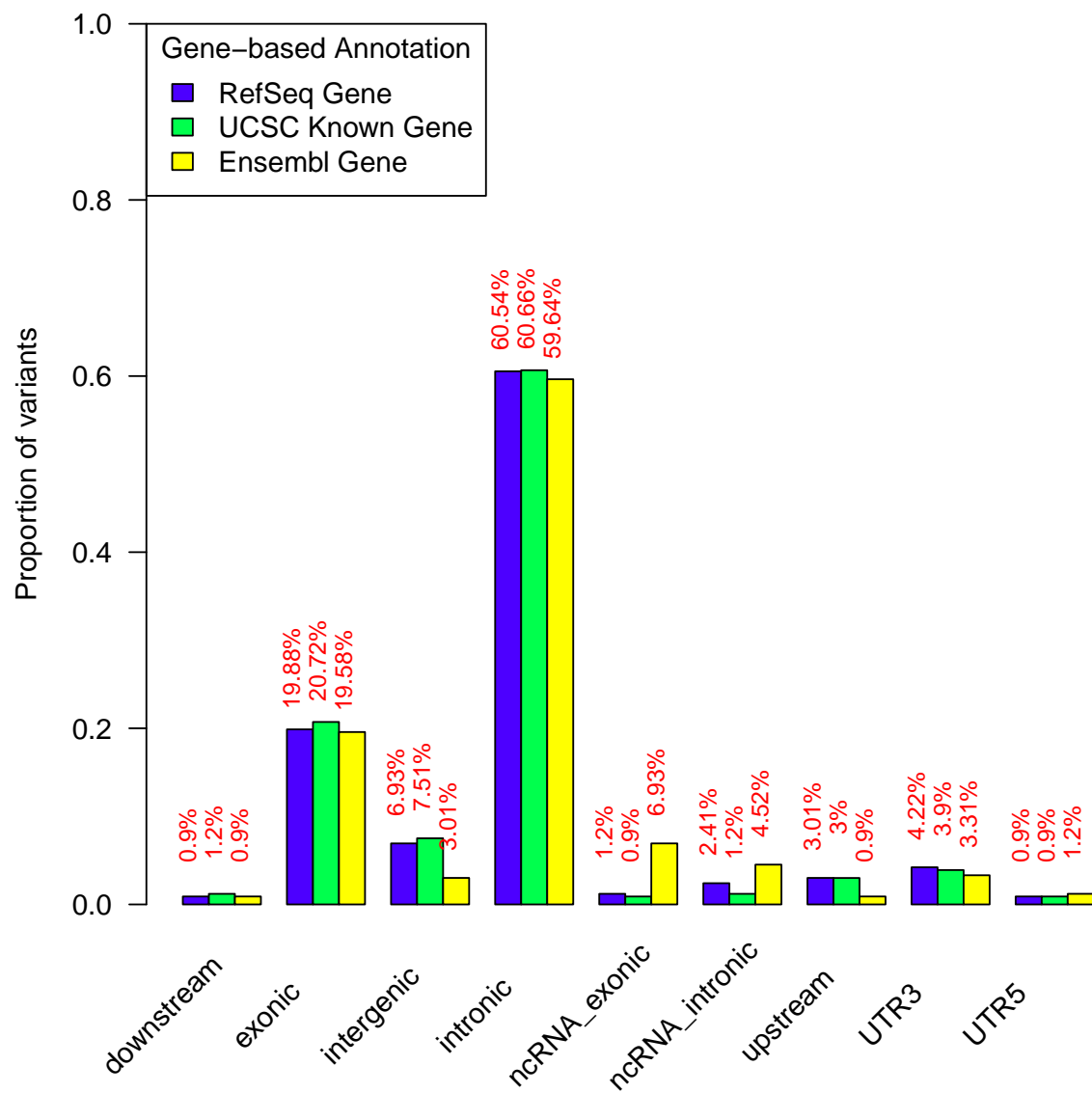


Figure 1: Proportion of variants that hit different regions based on RefSeq, UCSC Known Gene and Ensembl gene annotation sources.

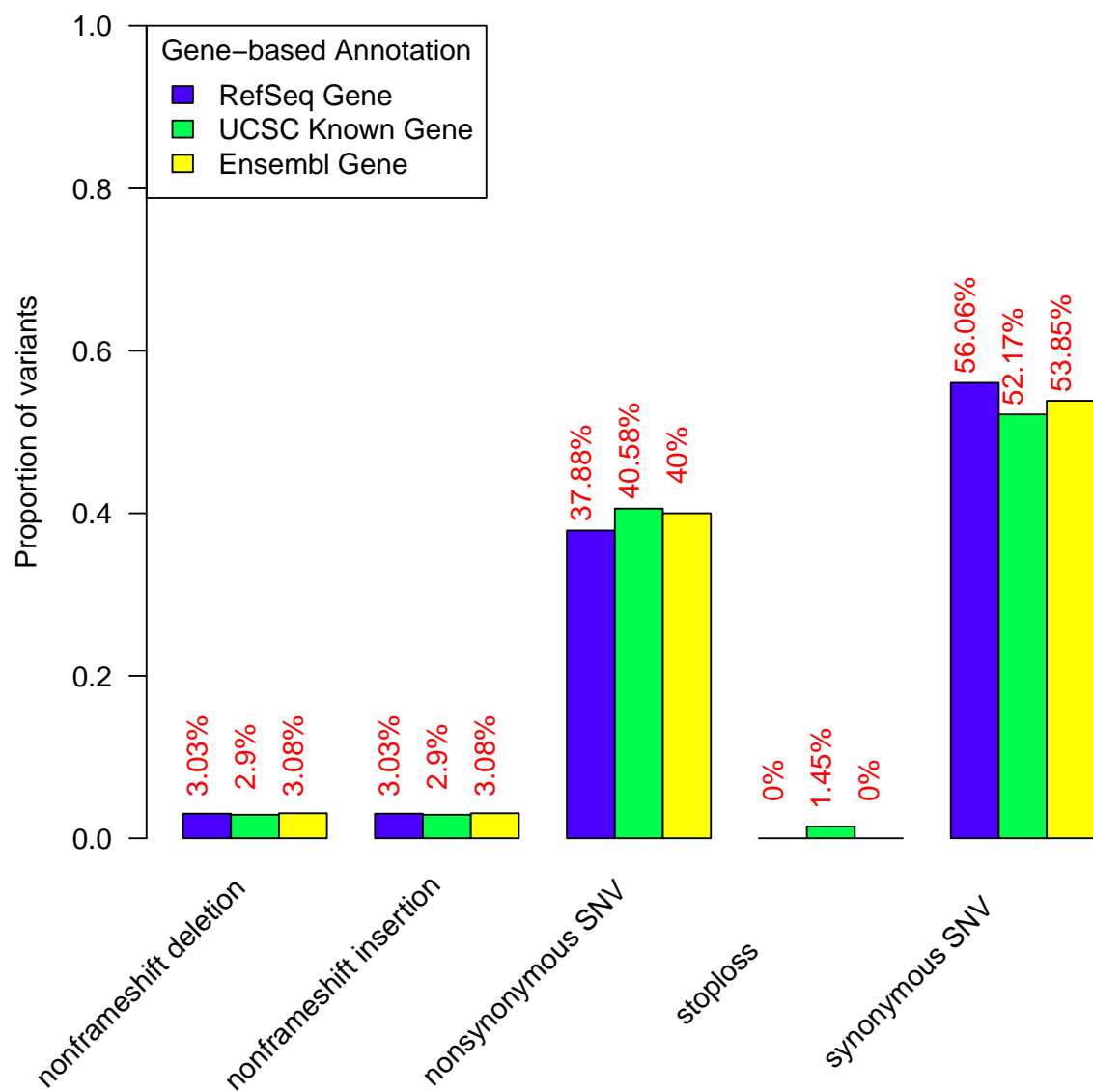


Figure 2: Proportion of exonic variants with their functional effects based on RefSeq, UCSC Known Gene and Ensembl gene annotation sources.