

# GeneRollup Documentation

Chris Gates and Jessica Bene  
UM BRCF Bioinformatics Core  
January 2016

# Overview

The variant lines in the input file are grouped by gene. Each output line aggregates the mutation, SnpEff, and dbNSFP values, as described in the below sections.

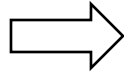
The number of output lines is equal to the number of distinct genes. The output lines are sorted such that the most impacted genes are listed first (based on *dbNSFP overall damaging rank* and *SnpEff overall impact rank*).

A sample is considered to be '**passed**' if it contains a value other than ".", "0", or blank. Otherwise, it is considered to be 'failed'.

A **sample-variant** is defined to be the intersection of a locus and a sample column.

# Overview

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



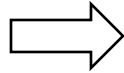
gene symbol	total impacted samples	distinct loci	total mutations	dbNSFP   overall damaging rank	dbNSFP   damaging votes   Sample1	dbNSFP   damaging votes   Sample2	dbNSFP   damaging votes   Sample3	SnpEff   overall impact rank	SnpEff   impact category   HIGH	SnpEff   impact category   MODERATE	SnpEff   impact category   LOW	SnpEff   impact category   MODIFIER	SnpEff   impact   Sample1	SnpEff   impact   Sample2	SnpEff   impact   Sample3
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	hl	hhl	hl
BRCA1	2	3	4	2	11	2		2	1	0	3	0	hl	ll	

# *gene symbol*

Variant lines are merged by gene. In this example, the *SNPEFF\_TOP\_EFFECT\_GENE\_SYMBOL* column lists the genes.

# gene symbol

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



gene symbol	total impacted samples	distinct loci	total mutations	dbNSFP   overall damaging rank	dbNSFP   damaging votes   Sample1	dbNSFP   damaging votes   Sample2	dbNSFP   damaging votes   Sample3	SnpEff   overall impact rank	SnpEff   impact category   HIGH	SnpEff   impact category   MODERATE	SnpEff   impact category   LOW	SnpEff   impact category   MODIFIER	SnpEff   impact   Sample1	SnpEff   impact   Sample2	SnpEff   impact   Sample3
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	hl	hhl	hl
BRCA1	2	3	4	2	11	2		2	1	0	3	0	hl	ll	

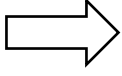
# *total impacted samples*

For each gene, the count of passed samples is calculated. In this example, the *JQ\_SUMMARY\_SOM\_COUNT* prefix depicts the columns for each sample.

*total impacted samples*

gene symbol	total impacted samples	distinct loci	total mutations	dbNSFP   overall damaging rank	dbNSFP   damaging votes   Sample1	dbNSFP   damaging votes   Sample2	dbNSFP   damaging votes   Sample3	SnpEff   overall impact rank	SnpEff   impact category   HIGH	SnpEff   impact category   MODERATE	SnpEff   impact category   LOW	SnpEff   impact category   MODIFIER	SnpEff   impact   Sample1	SnpEff   impact   Sample2	SnpEff   impact   Sample3
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	hI	hhl	hl
BRCA1	2	3	4	2	11	2		2	1	0	3	0	hl	ll	

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	✓	✓		HIGH	9
V2	ERBB2		✓	✓	HIGH	7
V3	ERBB2	✓	✓	✓	LOW	3
V4	BRCA1	✓	✓	✓	HIGH	9
V5	BRCA1	✓	✓	✓	LOW	.
V6	BRCA1	✓	✓	✓	LOW	2



✓ ✓ ✗  
1 + 1 + 0 = 2

# *distinct loci*

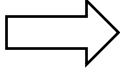
Each line in the input file corresponds to a locus. For each gene, a locus is counted if there is at least one passed sample-variant at that locus. Duplicate loci are only counted once.



# distinct loci

gene symbol																
	total impacted samples		distinct loci		total mutations											
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	h	h	h	h
BRCA1	2	3	4	2	11	2		2	1	0	3	0	h	h	h	h

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	✓	✓		HIGH	9
V2	ERBB2		✓	✓	HIGH	7
V3	ERBB2	✓	✓	✓	LOW	3
V4	BRCA1	✓			HIGH	9
V5	BRCA1		✓		LOW	.
V6	BRCA1	✓	✓		LOW	2



$$1 + 1 + 1 = 3$$

# *total mutations*

For each gene, the total mutations is determined by the number of passed sample-variants.

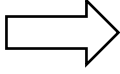
# total mutations

gene symbol	total impacted samples	distinct loci	total mutations	dbNSFP   overall damaging rank	dbNSFP   damaging votes   Sample1	dbNSFP   damaging votes   Sample2	dbNSFP   damaging votes   Sample3	SnpEff   overall impact rank	SnpEff   impact category   HIGH	SnpEff   impact category   MODERATE	SnpEff   impact category   LOW	SnpEff   impact category   MODIFIER	SnpEff   impact   Sample1	SnpEff   impact   Sample2	SnpEff   impact   Sample3
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	h	hh	h
BRCA1	2	3	4	2	11	2		2	1	0	3	0	h	h	



$$1 + 1 + 1 + 1 = 4$$

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	✓	✓		HIGH	9
V2	ERBB2		✓	✓	HIGH	7
V3	ERBB2	✓	✓	✓	LOW	3
V4	BRCA1	✓			HIGH	9
V5	BRCA1		✓		LOW	.
V6	BRCA1	✓	✓		LOW	2

# *dbNSFP / damaging votes / SampleX*

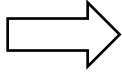
For each gene, the *dbNSFP\_rollup\_damaging* values are summed within each sample if their corresponding sample-variant is passed.

In other words, for a particular passed sample-variant within a gene, the *dbNSFP\_rollup\_damaging* value which corresponds to that locus is added to the total damaging votes for that sample.

# dbNSFP / damaging votes / SampleX

gene symbol	total impacted samples			distinct loci	total mutations			dbNSFP   overall damaging rank	dbNSFP   damaging votes   Sample1	dbNSFP   damaging votes   Sample2	dbNSFP   damaging votes   Sample3	SnpEff   overall impact rank	SnpEff   impact category   HIGH	SnpEff   impact category   MODERATE	SnpEff   impact category   LOW	SnpEff   impact category   MODIFIER	SnpEff   impact   Sample1	SnpEff   impact   Sample2	SnpEff   impact   Sample3
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	hl	hl	hl	hl	hl	hl	hl
BRCA1	2	3	4	2	11	2		2	1	0	3	0	hl	hl	ll				

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	9	9		HIGH	9
V2	ERBB2		7	7	HIGH	7
V3	ERBB2	3	3	3	LOW	3
V4	BRCA1	9	.	.	HIGH	9
V5	BRCA1	.	.	.	LOW	.
V6	BRCA1	2	2	.	LOW	2



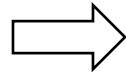
$9 + 0 + 2 = 11$   
 $0 + 0 + 2 = 2$   
 $0 + 0 + 0 = 0$

# *dbNSFP / overall damaging rank*

For each gene, the sum of the dbNSFP votes across all samples is calculated. These sums are then ranked such that the highest sum corresponds to a rank of 1.

# dbNSFP / overall damaging rank

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



gene symbol	total impacted samples	distinct loci	total mutations	dbNSFP   overall damaging rank	dbNSFP   damaging votes   Sample1	dbNSFP   damaging votes   Sample2	dbNSFP   damaging votes   Sample3	SnpEff   overall impact rank	SnpEff   impact category   HIGH	SnpEff   impact category   MODERATE	SnpEff   impact category   LOW	SnpEff   impact category   MODIFIER	SnpEff   impact   Sample1	SnpEff   impact   Sample2	SnpEff   impact   Sample3
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	hl	hhl	hl
BRCA1	2	3	4	2	11	2		2	1	0	3	0	hl	ll	

$$12 + 19 + 10 = 41$$

$$11 + 2 + 0 = 13$$

$$41 > 13$$

1. 41
2. 13

# *SnpEff / impact category / X*

In this case, the *SNPEFF\_TOP\_EFFECT\_IMPACT* column contains impact categories. For each impact category within each gene, the total number of mutations is stored in its corresponding *SnpEff / impact category* column.

In other words, for each gene, the data is pivoted on *SNPEFF\_TOP\_EFFECT\_IMPACT*, and the count of passed samples-variants is listed for each impact category.



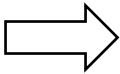
# SnpEff / impact category / X

gene symbol	total impacted samples	distinct loci	total mutations	dbNSFP   overall damaging rank	dbNSFP   damaging votes   Sample1	dbNSFP   damaging votes   Sample2	dbNSFP   damaging votes   Sample3	SnpEff   overall impact rank	SnpEff   impact category   HIGH	SnpEff   impact category   MODERATE	SnpEff   impact category   LOW	SnpEff   impact category   MODIFIER	SnpEff   impact   Sample1	SnpEff   impact   Sample2	SnpEff   impact   Sample3
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	h	h	h
BRCA1	2	3	4	2	11	2		2	1	0	3	0	h		



h → 1  
 m → 0  
 l → 1 + 1 + 1 = 3  
 x → 0

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL					JQ_SUMMARY_SOM_COUNT Sample1			JQ_SUMMARY_SOM_COUNT Sample2			JQ_SUMMARY_SOM_COUNT Sample3			SNPEFF_TOP_EFFECT_IMPACT		dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9			1	1	0	HIGH	9				
V2	ERBB2		1	1	HIGH	7											
V3	ERBB2	1	1	1	LOW	3											
V4	BRCA1	1	0	.	HIGH	9											
V5	BRCA1	.	1	.	LOW	.											
V6	BRCA1	1	1	.	LOW	2											



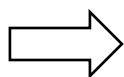
Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL					JQ_SUMMARY_SOM_COUNT Sample1			JQ_SUMMARY_SOM_COUNT Sample2			JQ_SUMMARY_SOM_COUNT Sample3			SNPEFF_TOP_EFFECT_IMPACT		dbNSFP_rollup_damaging
V1	ERBB2	h	h		HIGH	9											
V2	ERBB2		h	h	HIGH	7											
V3	ERBB2	l	l	l	LOW	3											
V4	BRCA1	h	l		HIGH	9											
V5	BRCA1	l	l		LOW	.											
V6	BRCA1	l	l		LOW	2											

# *SnpEff / overall impact rank*

Impact categories are weighted such that HIGH has the most significance and MODIFIER has the least significance. The *SnpEff / impact category* values are multiplied by their corresponding weight to obtain an impact score. These scores then ranked such that the highest score corresponds to a rank of 1.

# *SnpEff / overall impact rank*

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



gene symbol	total impacted samples	distinct loci	total mutations	dbNSFP   overall damaging rank	dbNSFP   damaging votes   Sample1	dbNSFP   damaging votes   Sample2	dbNSFP   damaging votes   Sample3	SnpEff   overall impact rank	SnpEff   impact category   HIGH	SnpEff   impact category   MODERATE	SnpEff   impact category   LOW	SnpEff   impact category   MODIFIER	SnpEff   impact   Sample1	SnpEff   impact   Sample2	SnpEff   impact   Sample3
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	hl	hhl	hl
BRCA1	2	3	4	2	11	2		2	1	0	3	0	hl	ll	

$$4w_1 + 0w_2 + 3w_3 + 0w_4 = 4000003000$$

$$1w_1 + 0w_2 + 3w_3 + 0w_4 = 1000003000$$

$$4000003000 > 1000003000$$



- 4000003000
- 1000003000

$$\begin{aligned} w_1 &= 10^9 \\ w_2 &= 10^6 \\ w_3 &= 10^3 \\ w_4 &= 10^0 \end{aligned}$$

# *SnpEff / impact / SampleX*

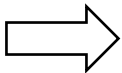
The *SNPEFF\_TOP\_EFFECT\_IMPACT* values are mapped to a single letter such that: HIGH = h, MODERATE = m, LOW = l, MODIFIER = x. For each gene, its single letters are concatenated within each sample if their corresponding sample-variant is passed.

In other words, for a particular passed sample-variant within a gene, the single letter for the *SNPEFF\_TOP\_EFFECT\_IMPACT* value which corresponds to that locus is appended to the SnpEff impact.

# SnpEff / impact / SampleX

gene symbol	total impacted samples			distinct loci	total mutations		dbNSFP   overall damaging rank		dbNSFP   damaging votes   Sample1		dbNSFP   damaging votes   Sample2		dbNSFP   damaging votes   Sample3		SnpEff   overall impact rank		SnpEff   impact category   HIGH		SnpEff   impact category   MODERATE		SnpEff   impact category   LOW		SnpEff   impact category   MODIFIER		SnpEff   impact   Sample1		SnpEff   impact   Sample2		SnpEff   impact   Sample3	
ERBB2	3	3	7	1	12	19	10	1	4	0	3	0	hl	hhl	hl	4	0	3	0	hl	hhl	hl	4	0	3	0	hl	hhl	hl	
BRCA1	2	3	4	2	11	2		2	1	0	3	0	hl	ll		2	1	0	3	0	hl	hhl	hl	4	0	3	0	hl	hhl	hl

Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	1	1	0	HIGH	9
V2	ERBB2		1	1	HIGH	7
V3	ERBB2	1	1	1	LOW	3
V4	BRCA1	1	0	.	HIGH	9
V5	BRCA1	.	1	.	LOW	.
V6	BRCA1	1	1	.	LOW	2



Variant	SNPEFF_TOP_EFFECT_GENE_SYMBOL	JQ_SUMMARY_SOM_COUNT Sample1	JQ_SUMMARY_SOM_COUNT Sample2	JQ_SUMMARY_SOM_COUNT Sample3	SNPEFF_TOP_EFFECT_IMPACT	dbNSFP_rollup_damaging
V1	ERBB2	h	h		HIGH	9
V2	ERBB2		h	h	HIGH	7
V3	ERBB2	l	l	l	LOW	3
V4	BRCA1	h	l	l	HIGH	9
V5	BRCA1	l	l	l	LOW	.
V6	BRCA1	l	l	l	LOW	2

hl ll blank

