2015-11-27

# 5. SNP Annotation

- **VCF format**
- **SnpEff tool**

# Difference between Ref. versions

http://www.ncbi.nlm.nih.gov/nuccore/299033929/

# Variant Call Format(VCF)



PER cultivar

Detected variants

Reference:
1) https://en.wikipedia.org/wiki/Variant_Call_Format
2) https://samtools.github.io/hts-specs/VCFv4.2.pdf
(current version)

# Variant Effect Annotation

## Predicts coding effects of genomic variants

### ■List of variant annotation tools

(https://en.wikipedia.org/wiki/SNP_annotation)

**List of available SNP annotation tools** [edit]

To annotate large number of available NGS data, currently a large number of SNPs annotation tools is available. Some of them are specific to some specific annotation. Some of the available SNPs annotation tools are as follows SNPeff, VEP, ANNOVAR, FATHMM, PhD-SNP, PolyPhen-2, SuSPect, F-SNP, AnnT SeattleSeq, SNPit, SCAN, Snap, SNPs&GO, LS-SNP, Snat, TREAT, TRAMS, Maviant, SNPdat, Snpranker, NGS - SNP, SVA, VARIANT, SIFT, PhD-SNP and Function and approach used in SNPs annotation tools are listed below

| Tools | Description | External resources use | WebsiteURL |
|---|---|---|---|
| SNPeff | SnpEff annotates variants based on their genomic locations and predicts coding effects. Use an interval forest approach | ENSEMBL, UCSC and organism based e.g. FlyBase, WormBase and TAIR | http://snpeff.sourceforge.net/SnpEff_manual.html |
| VEP | Provides the location of specific variants in individuals. Variants are calculated using sanger-style resequencing data | dbSNP, Ensembl, UCSC and NCBI | http://www.ensembl.org/ |
| ANNOVAR | This tools is suitable for pinpoint a small subset of functionally important variant. Use mutation prediction approach for annotation | UCSC, RefSe and Ensembl | http://www.openbioinformatics.org/annovar/ |

### ■SnpEff tool (example outputs※)
※http://snpeff.sourceforge.net/SnpEff_manual.html

**Type**

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| DOWNSTREAM | 2,093 | 1.766% |
| INTERGENIC | 26,314 | 22.204% |
| INTRAGENIC | 78 | 0.066% |
| INTRON | 54,238 | 45.767% |
| NON_SYNONYMOUS_CODING | 237 | 0.2% |
| NON_SYNONYMOUS_START | 1 | 0.001% |
| SPLICE_SITE_DONOR | 4 | 0.003% |
| START_GAINED | 57 | 0.048% |
| STOP_GAINED | 3 | 0.003% |
| STOP_LOST | 1 | 0.001% |
| SYNONYMOUS_CODING | 378 | 0.319% |
| TRANSCRIPT | 32,163 | 27.14% |
| UPSTREAM | 2,102 | 1.774% |
| UTR_3_PRIME | 690 | 0.582% |
| UTR_5_PRIME | 149 | 0.126% |

**Region**

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| DOWNSTREAM | 2,093 | 1.766% |
| EXON | 620 | 0.523% |
| INTERGENIC | 26,314 | 22.204% |
| INTRON | 54,238 | 45.767% |
| NONE | 32,241 | 27.206% |
| SPLICE_SITE_DONOR | 4 | 0.003% |
| UPSTREAM | 2,102 | 1.774% |
| UTR_3_PRIME | 690 | 0.582% |
| UTR_5_PRIME | 206 | 0.174% |

## ■Effect type of the SnoEff tool※

| Effect Type | Region |
|---|---|
| NONE CHROMOSOME CUSTOM CDS | NONE |
| INTERGENIC INTERGENIC_CONSERVED | INTERGENIC |
| UPSTREAM | UPSTREAM |
| UTR_5_PRIME UTR_5_DELETED START_GAINED | UTR_5_PRIME |
| SPLICE_SITE_ACCEPTOR | SPLICE_SITE_ACCEPTOR |
| SPLICE_SITE_DONOR | SPLICE_SITE_DONOR |
| SPLICE_SITE_REGION | SPLICE_SITE_REGION |
| INTRAGENIC START_LOST SYNONYMOUS_START NON_SYNONYMOUS_START GENE TRANSCRIPT | EXON or NONE |
| EXON EXON_DELETED NON_SYNONYMOUS_CODING SYNONYMOUS_CODING FRAME_SHIFT CODON_CHANGE CODON_INSERTION CODON_CHANGE_PLUS_CODON_INSERTION CODON_DELETION CODON_CHANGE_PLUS_CODON_DELETION STOP_GAINED SYNONYMOUS_STOP STOP_LOST RARE_AMINO_ACID | EXON |
| INTRON INTRON_CONSERVED | INTRON |
| UTR_3_PRIME UTR_3_DELETED | UTR_3_PRIME |
| DOWNSTREAM | DOWNSTREAM |
| REGULATION | REGULATION |

# Variant Effect Annotation(DatePalm vcf)

■ SNPeff manual = http://snpeff.sourceforge.net/SnpEff_manual.html

■Install at NIG supercomputer
 wget http://sourceforge.net/projects/snpeff/files/snpEff_latest_core.zip  (url in manual)
 unzip snpEff_latest_core.zip
 cd snpeff

■Generate target database
 cd data
 mkdir  GU811709
 cp SRA100551/snpeff/db/*  ./data/GU811709/
 (add snpeff.config) ──────────────────────→
 java -Xmx400M  -jar snpEff.jar build  -genbank -v GU811709

```
#Phoenix_dactylifera
GU811709.genome : Phoenix_dactylifera
```

■Apply query vcf files to target database
  cp SRA100551/snpeff/query/*.vcf  ./
  edit 1st column from GU811709|GU811709.2 to GU811709 for all *.vcf
  remove all comments lines (start #) for all *.vcf
  java -Xmx400M -jar snpEff.jar  -ud 200 GU811709 query/AJW.vcf > AJW_ann.vcf

# Variant Effect Annotation(DatePalm vcf)

## ■Check variant annotations

GU811709    4853    .    GTTTTTTTTTTTTT    "GTTTTTTTTTTTTTTTTT,GTTTTTTTTTTTTTTTTT,GTTTTTTTTTTT"
90.5    .    "INDEL;DP=33;VDB=0.0298;AF1=1;AC1=2;DP4=0,1,10,22;MQ=44;FQ=-59.5;PV4=1,1,0.27,1";ANN="GTTTTTTTTTTTTTTTTTTT|downstream_gene_variant|MODIFIER|rps16|Gene_4869_5988|transcript|ADD63156.2|Coding||c.*4_*17delAAAAAAAAAAAAACinsAAAAAAAAAAAAAAAAAAC
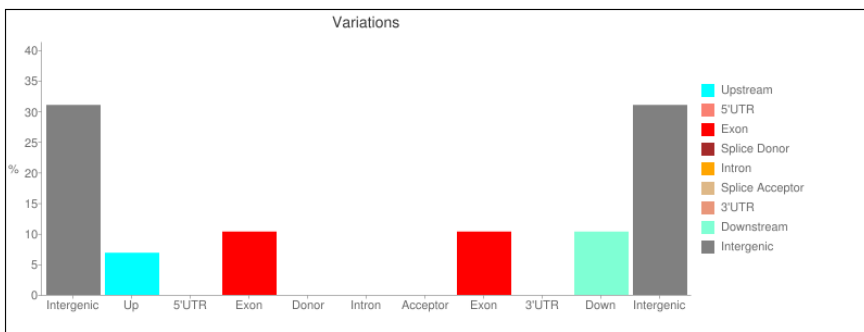
GU811709    21746    .    G    T    225.0    .
"DP=229;VDB=0.0487;AF1=0.5;AC1=1;DP4=29,27,87,86;MQ=42;FQ=141;PV4=0.88,1,1,0.074";ANN=T|synonymous_variant|LOW|rpoC1|Gene_21631_24418|transcript|ADD63235.1|Coding|2...

GU811709    21750    .    T    C    70.0    .
"DP=231;VDB=0.0403;AF1=0.5;AC1=1;DP4=90,95,25,21;MQ=42;FQ=73;PV4=0.51,0.37,1,1";ANN=C|missense_variant|MODERATE|rpoC1|...

## ■Check SNP supEff_summary.html



Number of effects by type and region

## ■Count unique annotations for AJW_ann.vcf

| AnnType | Count | Percent |
|---|---|---|
| downstream gene | | |
| intergenic region | | |
| intron variant | | |
| missense variant | | |
| synonymous variant | | |
| upstream gene | | |

# THANK YOU!