

TCAG SMALL VARIANT ANNOTATION PIPELINE

Creation date: 31 August 2020.

Pipeline version: 27.4

Table of contents

- 1. General release notes**
- 2. Overall Annotation Process:**
- 3. Detailed annotation field description**
- 4. Database and tool versions**
- 5. File Descriptions**

1. GENERAL RELEASE NOTES

Following information is valid for annotations generated by TCAG annotation pipeline rev27.4 and applicable to variants called by GATK (whole genome and exome) pipelines.

2. Overall Annotation Process:

The output of many small variant calling software including GATK and HAS is a file in VCF format. The Variant Call Format is a standardized format to represent SNPs and indels. The TCAG annotation pipeline uses ANNOVAR to functionally annotate these small variants. In order to accurately add gene-based, region-based and filter-based annotations, the original VCF file generated by the variant caller needs to be converted into an ANNOVAR compatible format. The conversion involves variant decomposition (splitting of variants with two alternate alleles into multiple rows), left-aligning indels on the forward strand of the reference genome and normalization. For indels, the normalization sometimes shifts the position of the variants. The 'GT_PreNorm' captures the genotype before variant decomposition and 'Original_VCFKey' captures the position of the variant before normalization.

3. Detailed annotation field description

[coordinates]

CHROM: chromosome (autosomes 1-22 and sex chromosomes X, Y)

start: start position (1-positional system)

end: end position (1-positional system)

Original_VCFKEY: unique variant identifier for original vcf file.

MULTI_ALLELIC: Flag to represent multi-allelic sites (ie 1=multi-allelic).

[sequence and ploidy]

<Sample-name>:zygosity: heterozygous-reference (ref-alt), homozygous-alternate (hom-alt) and heterozygous-alternate(alt-alt)

ref_allele: Reference allele

alt_allele: Alternate allele

<Sample-name>: Genotype: Genotype, represented as “reference allele | alternate allele” (ref-alt) or “alternate allele | alternate allele” (hom-alt) or “alternate allele 1 | alternate allele 2” (alt-alt). This is based on pre-normalized raw genotype (eg: A/C).

<Sample-name>: GT_PostNorm: Raw genotype after left-normalization (eg : 0/1)

<Sample-name>: GT_PreNorm: Raw genotype before left-normalization (eg : 1/2)

var_type: snp (single nucleotide variation, can be ref-alt), ins (insertion), del (deletion), mnp (multiple bases substitutions), complex (could be block substitution).

SNVSB: SNV site strand bias (only for HAS variants).

SNVHPOL: SNV contextual homopolymer length (Only for HAS variants).

RU: Smallest repeating sequence unit extended or contracted in the indel allele relative to the reference. RUs are not reported if longer than 20 bases (only to HAS variants).

<Sample-name>: PhaseGT:Physical phasing haplotype information, describing how the alternate alleles are phased in relation to one another (Applicable only to GATK variants).

<Sample-name>: PhaseID:Physical phasing ID information, where each unique ID within a given sample (but not across samples) connects records within a phasing group (Applicable only to GATK variants).

<Sample-name>:NULL_CONFIG: NULL trio configuration (Applicable to trios set analysed with DeNovoGear).

<Sample-name>:DNM_CONFIG: DNM trio configuration (Applicable to trios set analysed with DeNovoGear).

[quality score and read counts]

FILTER: Filter status: PASS if this position has passed all filters.

DP: Approximate read depth for this variant (Applicable only to GATK variants).

FS: Phred-scaled p-value using Fisher's exact test to detect strand bias (Applicable only to GATK variants).

QD: Variant Confidence or Quality by Depth (Applicable only to GATK variants).

MQ: RMS Mapping Quality. This annotation provides an estimation of the overall mapping quality of reads supporting a variant call. It produces both raw data (sum of square and num of total reads) and the calculated root mean square (Applicable only to GATK variants).

MQRankSum: Z-score From Wilcoxon rank sum test of Alt vs. Ref read mapping qualities (Applicable only to GATK variants).

ReadPosRankSum: Z-score from Wilcoxon rank sum test of Alt vs. Ref read position bias (Applicable only to GATK variants).

culprit: The annotation which was the worst performing in the Gaussian mixture model, likely the reason why the variant was filtered out (Applicable only to GATK variants).

VQSLOD: Log odds of being a true variant versus being false under the trained gaussian mixture model (Applicable only to GATK variants).

<Sample-name>: AD_REF: Allelic depth for the ref allele.

<Sample-name>: AD_ALT: Allelic depth for the alt allele.

<Sample-name>: DP: Filtered basecall depth used for site genotyping.

<Sample-name>: DPF: Basecalls filtered from input prior to site genotyping (Applicable only to HAS variants).

<Sample-name>: DPI: Read depth associated with indel, taken from the site preceding the indel (Applicable only to HAS variants).

<Sample-name>: GQ: Genotype quality.

<Sample-name>: GQX: Empirically calibrated variant quality score for variant sites, otherwise Minimum of {Genotype quality assuming variant position, Genotype quality assuming non-variant position} (Applicable only to HAS variants).

<Sample-name>: PL: Normalized, Phred-scaled likelihoods for genotypes as defined in the VCF specification (Applicable only to GATK variants).

<Sample-name>: ML_NULL: Maximum Likelihood for the NULL configuration (Applicable to trios set analysed with DeNovoGear).

<Sample-name>: PP_NULL: Posterior probability for the NULL configuration (Applicable to trios set analysed with DeNovoGear).

<Sample-name>: ML_DNM: Maximum Likelihood for the DNM configuration (Applicable to trios set analysed with DeNovoGear).

<Sample-name>: PP_DNM: Posterior probability for the DNM configuration (Applicable to trios set analysed with DeNovoGear).

<Sample-name>: MQ: Mapping quality of the child (Applicable to trios set analysed with DeNovoGear).

<Sample-name>: RD: Read Depth of the child (Applicable to trios set analysed with DeNovoGear).

<Sample-name>: DNG_confidence: Confidence for this denovo call (High or Low). High confident entries are selected if a) for snps, PP_DNM > 95, GQ >= 99 and alt_allele fraction > 0.3 b) for indels, GQ > 90 and alt_allele fraction > 0.3 (Applicable to trios set analysed with DeNovoGear).

[allele frequency]

A1000g_all: allele frequency in the full 1000 Genome data-set

A1000g_eur: allele frequency in the caucasian-european sub-set of 1000 Genome

A1000g_amr: allele frequency in the mixed-background latin americans (e.g. Mexicans, Puerto Ricans, Peruvians) sub-set of 1000 Genome

A1000g_eas: allele frequency in the east-asian sub-set of 1000 Genome

A1000g_sas: allele frequency in the south-asian sub-set of 1000 Genome

A1000g_afr: allele frequency in the black-african sub-set of 1000 Genome (note african here indicates the ethnic background, not the actual geographical location of the population sample)

gnomAD_exome211_FILTER: gnomAD v2.1.1 Filter status: PASS if this position has passed all filters. From gnomAD exome dataset.

gnomAD_exome211_AC: gnomAD v2.1.1 alternate allele count for samples. From gnomAD exome dataset.

gnomAD_exome211_AC_female: gnomAD v2.1.1 alternate allele count for female samples. From gnomAD exome dataset.

gnomAD_exome211_AC_male: gnomAD v2.1.1 alternate allele count for male samples. From gnomAD exome dataset.

gnomAD_exome211_AN: gnomAD v2.1.1 total number of alleles in all samples. From gnomAD exome dataset.

gnomAD_exome211_AN_female: gnomAD v2.1.1 total number of alleles in female samples. From gnomAD exome dataset.

gnomAD_exome211_AN_male: gnomAD v2.1.1 total number of alleles in male samples. From gnomAD exome dataset.

gnomAD_exome211_nhomalt: gnomAD v2.1.1 count of homozygous individuals in samples. From gnomAD exome dataset.

gnomAD_exome211_nhomalt_female: gnomAD v2.1.1 count of homozygous individuals in female samples. From gnomAD exome dataset.

gnomAD_exome211_nhomalt_male: gnomAD v2.1.1 count of homozygous individuals in male samples. From gnomAD exome dataset.

gnomAD_exome211_AF: gnomAD v2.1.1 allele frequencies (exomes).

gnomAD_exome211_AF_raw: gnomAD v2.1.1 raw allele frequencies (exomes) - global.

gnomAD_exome211_AF_afr: gnomAD v2.1.1 allele frequencies (exomes) - AFR subset.

gnomAD_exome211_AF_amr: gnomAD v2.1.1 allele frequencies (exomes) - AMR subset.

gnomAD_exome211_AF_asj: gnomAD v2.1.1 allele frequencies (exomes) - ASJ subset.

gnomAD_exome211_AF_eas: gnomAD v2.1.1 allele frequencies (exomes) - EAS subset.

gnomAD_exome211_AF_nfe: gnomAD v2.1.1 allele frequencies (exomes) - NFE subset.

gnomAD_exome211_AF_fin: gnomAD v2.1.1 allele frequencies (exomes) - FIN subset.

gnomAD_exome211_AF_oth: gnomAD v2.1.1 allele frequencies (exomes) - others.

gnomAD_exome211_AF_sas: gnomAD v2.1.1 allele frequencies (exomes) - SAS subset.

gnomAD_exome211_AF_female: gnomAD v2.1.1 alternate allele frequency in female samples (exomes) .

gnomAD_exome211_AF_male: gnomAD v2.1.1 alternate allele frequency in male samples (exomes) .

gnomAD_exome211_faf95_afr: Filtering allele frequency (using Poisson 95% CI) for samples of African-American/African ancestry (exomes).

gnomAD_exome211_faf95_amr: Filtering allele frequency (using Poisson 95% CI) for samples of Latino ancestry (exomes).

gnomAD_exome211_faf95_eas: Filtering allele frequency (using Poisson 95% CI) for samples of East Asian ancestry (exomes).

gnomAD_exome211_faf95_nfe: Filtering allele frequency (using Poisson 95% CI) for samples of Non-Finnish European ancestry (exomes).

gnomAD_exome211_faf95_sas: Filtering allele frequency (using Poisson 95% CI) for samples of South Asian ancestry (exomes).

gnomAD_genome211_FILTER: gnomAD v2.1.1 Filter status: PASS if this position has passed all filters. From gnomAD genome dataset.

gnomAD_genome211_AC: gnomAD v2.1.1 alternate allele count for all samples. From gnomAD genome dataset.

gnomAD_genome211_AC_female: gnomAD v2.1.1 alternate allele count for female samples. From gnomAD genome dataset.

gnomAD_genome211_AC_male: gnomAD v2.1.1 alternate allele count for male samples. From gnomAD genome dataset.

gnomAD_genome211_AN: gnomAD v2.1.1 total number of alleles in all samples. From gnomAD genome dataset.

gnomAD_genome211_AN_female: gnomAD v2.1.1 total number of alleles in female samples. From gnomAD genome dataset.

gnomAD_genome211_AN_male: gnomAD v2.1.1 total number of alleles in male samples. From gnomAD genome dataset.

gnomAD_genome211_nhomalt: gnomAD v2.1.1 count of homozygous individuals in all samples. From gnomAD genome dataset.

gnomAD_genome211_nhomalt_female: gnomAD v2.1.1 count of homozygous individuals in female samples. From gnomAD genome dataset.

gnomAD_genome211_nhomalt_male: gnomAD v2.1.1 count of homozygous individuals in male samples. From gnomAD genome dataset.

gnomAD_genome211_AF: gnomAD v2.1.1 allele frequencies (genomes).

gnomAD_genome211_AF_raw: gnomAD v2.1.1 raw allele frequencies (genomes) - global.

gnomAD_genome211_AF_afr: gnomAD v2.1.1 allele frequencies (genomes) - AFR subset.

gnomAD_genome211_AF_amr: gnomAD v2.1.1 allele frequencies (genomes) - AMR subset.

gnomAD_genome211_AF_asj: gnomAD v2.1.1 allele frequencies (genomes) - ASJ subset.

gnomAD_genome211_AF_eas: gnomAD v2.1.1 allele frequencies (genomes) - EAS subset.

gnomAD_genome211_AF_nfe: gnomAD v2.1.1 allele frequencies (genomes) - NFE subset.

gnomAD_genome211_AF_fin: gnomAD v2.1.1 allele frequencies (genomes) - FIN subset.

gnomAD_genome211_AF_oth: gnomAD v2.1.1 allele frequencies (genomes) - others.

gnomAD_genome211_AF_sas: gnomAD v2.1.1 allele frequencies (genomes) - SAS subset.

gnomAD_genome211_AF_female: gnomAD v2.1.1 alternate allele frequency in female samples (genomes) .

gnomAD_genome211_AF_male: gnomAD v2.1.1 alternate allele frequency in male samples (genomes) .

gnomAD_genome211_faf95_afr: Filtering allele frequency (using Poisson 95% CI) for samples of African-American/African ancestry (genomes).

gnomAD_genome211_faf95_amr: Filtering allele frequency (using Poisson 95% CI) for samples of Latino ancestry (genomes).

gnomAD_genome211_faf95_eas: Filtering allele frequency (using Poisson 95% CI) for samples of East Asian ancestry (genomes).

gnomAD_genome211_faf95_nfe: Filtering allele frequency (using Poisson 95% CI) for samples of Non-Finnish European ancestry (genomes).

gnomAD_genome211_faf95_sas: Filtering allele frequency (using Poisson 95% CI) for samples of South Asian ancestry (genomes).

freq_max: Maximum of GnomAD exome/GnomAD genome filtering allele frequencies(faf) .

A1000g_freq_max: Maximum of 1000 Genomes allele frequencies

gnomAD_exome_freq_max: Maximum filtering allele frequencies(faf) for gnomAD (exome).

gnomAD_genome_freq_max: Maximum filtering allele frequencies(faf) of gnomAD (genome)

.

[reference variant databases]

dbSNP: exact match (position, allele) to dbSNP.

dbSNP_common: exact match (position, allele) to common dbSNP track USCS; this is not meant as a replacement of frequency-based filtering, and is not practically used for hard filters.

dbSNP_region: overlap-based match for dbSNP; this is useful to look up variants that are split into multiple dbSNP entries, or otherwise only partially match to dbSNP entries.

dbSNP_common_region: overlap-based match for dbSNP; this can be used to hard filter (non-frameshift) substitutions or complex variants that are unlikely to be rare and/or damaging

Clinvar_SIG: Overall ClinVar significance code; "pathogenic" is the code of interest for rare disorders. Clinical significance values for all the individual submissions (SCVs) aggregated for the RCV record in ClinVar (<https://www.ncbi.nlm.nih.gov/clinvar/docs/clinsig/>)

Clinvar_CLNREF: References includes PubMed ID

Clinvar_AlleleID: ClinVar alleleID. A unique integer identifier assigned to each individual variant in ClinVar.

Clinvar_ReviewStatus: The level of review supporting the assertion of clinical significance(https://www.ncbi.nlm.nih.gov/clinvar/docs/details/#review_status)

Clinvar_SIG_Simple: expected values 0,1 or -1; 0 = no current value of pathogenic; 1 = at least one record submitted with pathogenic/likely pathogenic; -1 = no values for clinical significance at all for this variant or set of variants. Used for the "included" variants that are only in ClinVar because they are included in a haplotype or genotype with an interpretation.

cosmic: exact match (position, allele) to the Cosmic database of somatic variants; useful only for cancer projects.

Clinvar_Disease: Clinvar interpreted condition.

[gene mapping]

refseq_id: combined Annovar output on coding sequence mapping and effect, composed of: (a) for coding exonic changes ("typeset" exonic): gene official symbol : RefSeq ID : position in the coding sequence : amino acid change, [idem]; (b) for core splice site changes (typeseq "exonic"): gene official symbol (RefSeq ID : exon number : coding sequence position and change, [idem])

typeseq: type of sequence overlapped, with respect to known genes/transcripts and their coding / noncoding status: (a) "exonic" represents coding exons, (b) "splicing" represents core splicing site (by default, 15 bp on the intron side of intron-exon and exon-intron junctions), (c) "ncRNA_exonic" represents exons of non-coding RNA genes, (d) "ncRNA_splicing" represents core splicing sites of non-coding RNA genes , (e) "UTR5" represents 5' untranslated region, (f) "UTR3" represents 3' untranslated region, (g) "upstream" represents 1kb upstream of TSS, (h) "downstream" represents 1kb downstream of TSS and (i) "intergenic" represents intergenic regions (beyond upstream/downstream threshold(1kb)). For variants with multiple sequence overlaps (eg, exonic for one transcript and intronic for other), all possible typeseq values will be listed in semicolon-delimited format (eg: exonic;intronic).

typeseq_priority: Prioritized sequence overlap for multi-sequence overlap variants.

Implementation of the Annovar prioritization scheme

(<http://annovar.openbioinformatics.org/en/latest/user-guide/gene/>).

typeseq_RefseqSelect: Type of sequence overlapped, with respect to RefSeq Select transcript

(https://www.ncbi.nlm.nih.gov/refseq/refseq_select/).

effect: type of effect on the coding sequence: (a) "synonymous SNV", (b) "nonsynonymous SNV", (c) "stopgain SNV", (d) "frameshift deletion", (e) "frameshift insertion", (f) "frameshift substitution", (g) "nonframeshift deletion", (h) "nonframeshift insertion", (i) "nonframeshift substitution", (j) "stoploss SNV". For variants with multiple effects, all possible values will be represented in comma-separated fashion.

effect_priority: Prioritized effects for coding variants with multiple effects

(<http://annovar.openbioinformatics.org/en/latest/user-guide/gene/>).

effect_RefseqSelect: Effect on the coding sequence with respect to RefSeq Select transcript

(https://www.ncbi.nlm.nih.gov/refseq/refseq_select/).

aa_flag: this flag is set to 1 if more than one distinct amino acid change is reported in the "refseq-id" field.

distance_spliceJunction: A positive integer represents the distance from the nearest exon boundary, in genomic level.

gene_symbol: official gene symbol.

entrez.id: entrez-gene id.

gene_desc: full gene name.

omim_id: OMIM gene accession id.

omim_Phenotype: OMIM disorder/disease description when available for the corresponding omim gene accession.

MPO: (array of) MPO (Mammalian Phenotype Ontology) top level phenotype(s), imported from MGI and mapped from an orthologous mouse gene; the genotype-phenotype association is typically supported by a heterozygous/homozygous knock-out or other transgenic experiment, sometimes involving more than one gene: these details are exported by TCAG from MGI, but not included in this annotation field, so they should be looked up on the mgi website.

HPO: (array of) HPO (Human Phenotype Ontology) top level phenotype(s), imported from HPO; the genotype-phenotype association is typically supported by an OMIM entry; modes of inheritance are also exported by TCAG from HPO, but not included in this annotation field, so they should be looked up on the HPO website.

CGD_disease: the Clinical Genomics Database is compiled by curators and maintained by the NHGRI (National Human Genome Research Institute); for every gene in the database, the CGD

provides a list of one or more genetic disorders and a mode of inheritance; this field reports the genetic disorder(s) .

CGD_inheritance: the Clinical Genomics Database is compiled by curators and maintained by the NHGRI (National Human Genome Research Institute); for every gene in the database, the CGD provides a list of one or more genetic disorders and a mode of inheritance; this field reports the mode of inheritance (AD, AR, AD/AR, XL, more complex modes); since the CGD mode of inheritance is directly added by a curator and it's tied to specific genetic disorder(s), it could be considered more accurate than the mode of inheritance for top-level HPO phenotypes.

gnomAD_oe_lof: LOF observed/expected (oe) metric from gnomAD constraint matrix.

gnomAD_oe_lof_upper: LOF observed/expected (oe) metric - CI upper bound

gnomAD_oe_mis: Missense observed/expected (oe) metric from gnomAD constraint matrix.

gnomAD_oe_mis_upper: missense observed/expected (oe) metric - CI upper bound.

gnomAD_pLI: Probability of being loss-of-function intolerant.

gnomAD_pRec: Probability of being intolerant of homozygous but not heterozygous LOF variants.

gnomAD_mis_z: GnomAD missense Z score.

ACMG_disease: Any (exonic, intronic or splice) variants in genes in ACMG published recommendations for reporting incidental findings
(<https://www.ncbi.nlm.nih.gov/clinvar/docs/acmg/>) .

[conservation and predicted impact]

sift_score: SIFT score for predicted protein impact, values ≤ 0.05 correspond to damaging (can be interpreted as a p-value); note that SIFT is based on amino acid conservation / substitution rates inferred from protein sequence alignments.

polyphen_score: Polyphen2 scores for predicted protein impact, values ≥ 0.95 correspond to damaging; note that Polyphen2 is based on a set of sequence and structural attributes, combined in a machine classifier trained on a data-set of positive cases, and thus it partially is complementary and partially correlated to SIFT and MA.

ma_score: mutation assessor scores for predicted protein impact, values ≥ 2 correspond to damaging; note that MA is based on amino acid conservation / substitution rates inferred from protein sequence alignments, additionally modeling protein family groupings, thus can be regarded as an improved version of SIFT (improved performance has been shown particularly for somatic variants), although to this date SIFT remains more popular / broadly used.

mt_score: Mutationtaster prediction scores.

PROVEAN_score: Amino acid substitution or indel prediction score from Provean software (<http://provean.jcvi.org/index.php>). Values < -2.5 corresponds to a damaging variant.

phylopMam: value array of PhyloP nucleotide-level conservation inferred from the Placental Mammal genome group; values ≥ 1 indicate moderate conservation, values ≥ 2.5 indicate strong conservation; note that PhyloP is based uniquely on nucleotide substitutions and does not factor structural variation.

phylopMam_avg: Average value of PhyloP nucleotide-level conservation inferred from the Placental Mammal genome group; values ≥ 1 indicate moderate conservation, values ≥ 2.5 indicate strong conservation; this field is suitable to hard filter, especially for missense variants and non-frameshift substitutions; note that PhyloP is based uniquely on nucleotide substitutions and does not factor structural variation.

phylopVert100: Value array of PhyloP nucleotide-level conservation inferred from the 100 Vertebrate genome group; values ≥ 1.5 indicate moderate conservation, values ≥ 4 indicate strong conservation; note that PhyloP is based uniquely on nucleotide substitutions and does not factor structural variation.

phylopVert100_avg: Average value of PhyloP nucleotide-level conservation inferred from the 100 Vertebrate genome group; this field is suitable to hard filter, especially for missense variants and non-frameshift substitutions; values ≥ 1.5 indicate moderate conservation, values ≥ 4 indicate strong conservation; note that PhyloP is based uniquely on nucleotide substitutions and does not factor structural variation.

phastCons_placental: PhastCons score for the Placental Mammal genome group, based on UCSC track (calculated by HMM scan of PhyloP values to infer conserved status); this is useful to assess conservation at the *regional level* (rather than PhyloP's *nucleotide level*); this field can be used as an evidence *suggestive* of functional relevance of a sequence (especially if non-coding), and thus could be used to hard filter non-coding variation.

gerp_elem: Rejected substitution(RS) score for GERP++ elements.

gerp_wgs: whole-genome GERP++ RS scores greater than 2.

pfam_annotar: overlap with coding sequence matching to a PFAM protein domain; can help further assess protein impact, but not meant for hard-filter.

per_cds_affected: percentage of coding exonic sequence affected, should be used to prioritize loss-of-function variation (stopgain, frameshift, splicing), but not meant for hard filters.

per_transcripts_affected: percentage of transcripts with variant overlapping them, should be used to prioritize variation, but not meant for hard filters.

CADD_Raw: Raw score from CADD(Combined Annotation Dependent Depletion).

CADD_phred: PHRED-like c-score ranking from CADD.

EPDnew: EPD Promoter sequence overlap (<https://epd.epfl.ch/index.php>)

spx_gene: The RefSeq gene affected by the variant.

spx_transcript: The RefSeq transcript affected by the variant

spx_exon_number: The exon for which percent inclusion is predicted

spx_ss_dist: The distance from the variant to the splice site

spx_wt_psi: The wild-type percentage exon inclusion

spx_sequence_event_type: Classification of effect of the variant on sequence

spx_dpsi: Splice site prediction score from SPIDEX. The change in percentage exon inclusion reported as the maximum across tissues (<https://www.ncbi.nlm.nih.gov/pubmed/25525159>)

dbscSNV_ADA_SCORE, dbscSNV_RF_SCORE: Splice site prediction scores from dbscSNV (<https://www.ncbi.nlm.nih.gov/pubmed/26555599>)

spliceAI_SYMBOL: Gene symbol from SpliceAI (<https://pypi.org/project/spliceai/>)

spliceAI_STRAND: Strand from SpliceAI.

spliceAI_TYPE: Sequence overlap for this variant (I=intronic; E=exonic) from SpliceAI.

spliceAI_DIST: Distance from exon/intron junction from SpliceAI.

spliceAI_DS_AG: SpliceAI Delta score, the probability of the variant being splice-altering (acceptor gain). To filter use 0.5 (recommended), for higher recall 0.2 can be used.

spliceAI_DS_AL: SpliceAI Delta score, the probability of the variant being splice-altering Delta score (acceptor loss). To filter use 0.5 (recommended), for higher recall 0.2 can be used.

spliceAI_DS_DG: SpliceAI Delta score, the probability of the variant being splice-altering Delta score (donor gain). To filter use 0.5 (recommended), for higher recall 0.2 can be used.

spliceAI_DS_DL: SpliceAI Delta score, the probability of the variant being splice-altering Delta score (donor loss). To filter use 0.5 (recommended), for higher recall 0.2 can be used.

spliceAI_DP_AG: SpliceAI Delta position (acceptor gain).

spliceAI_DP_AL: SpliceAI Delta position (acceptor loss).

spliceAI_DP_DG: SpliceAI Delta position (donor gain).

spliceAI_DP_DL: SpliceAI Delta position (donor loss).

[other]

SegDup: overlap with segmental duplications, often useful to explain why ploidy is not 2 outside male sex chromosomes.

Repeat: Repeatmasker annotation from UCSC.

Comment: Comments related to annotation database issues (ie, ambiguous liftover or incomplete ORF for gene transcript).

4 .Database and tool versions

Genome build : hg19

Annovar : October 2019 version

Gene-type : refGene

Gene_Updated: 08-JUL-2020

OMIM: version 08-JUL-2020

HPO: version 08-JUL-2020

MPO: version 08-JUL-2020

gnomAD Constraint matrices: Downloaded August 2019

CGD: version 08-JUL-2020

dbsnp: dbsnp150, Downloaded from annovar, updated on 29-SEP-2017

dbsnpCommon: dbsnp151 common, Download from NCBI on 10-APR-2020

clinvar: Downloaded from NCBI and curated , 11-AUG-2020 version

cosmic: v91, Download from cosmic , APR-2020 version

1000g: Downloaded from annovar, 24-AUG-2015 version

gnomAD : Version 2.1.1, Downloaded from gnomAD, 10-MAR-2019

pfam: Downloaded from UCSC on June 2017

SPIDEX: splicing predictions v2.2 , Updated July2016

segdup:genomicSuperDups, Downloaded from UCSC , Oct 2011 version

phastcon: Downloaded from UCSC, April 2014 version

phylop_placental: phyloP46way/placentalMammals , Downloaded from UCSC Nov 2009 version

phylop_100way: phyloP100way, Downloaded from UCSC on Dec 2014, UCSC update : Feb 2014

gerp++elem: Downloaded from annovar, 23-FEB-2014 version

gerp_wgs : Downloaded from annovar, 21-JUN-2012 version

CADD: v1.6, Updated AUG-2020

dbnsfp: Downloaded from annovar, 21-FEB-2019 version

dbcsnv: Downloaded from annovar, 18-DEC-2015 version

SpliceAI: Downloaded from Illumina, 08-OCT-2019

5. File Descriptions

*_FINAL_rev27.1.tsv.gz: Complete set of annotated variants.

*_FINAL_rev27.1.tsv.gz.tbi: Tabix index file for above.

*_SUBSET_rev27.1.tsv: Subset of annotated variants from above file. This subset consists of variants that satisfy following conditions:

rare (<5% for all allele frequency columns (A1000g_all , A1000g_eur , A1000g_amr , A1000g_eas , A1000g_afr , A1000g_sas ,
gnomAD_exome_controls_AF_popmax, gnomAD_genome_controls_AF_popmax))

AND

typeseq column contains at least one of the following terms; exonic, splicing, downstream,upstream or UTR

OR

have significant splicing prediction score ($2 < \text{spx_dpsi} < -2$, $\text{dbscSNV_ADA_SCORE} > 0.6$, $\text{dbscSNV_RF_SCORE} > 0.6$, $(\text{spliceAI_DS_AG} > 0.2 \text{ and } \text{abs}(\text{spliceAI_DP_AG}) \leq 50)$, $(\text{spliceAI_DS_AL} > 0.2 \text{ and } \text{abs}(\text{spliceAI_DP_AL}) \leq 50)$, $(\text{spliceAI_DS_DG} > 0.2 \text{ and } \text{abs}(\text{spliceAI_DP_DG}) \leq 50)$ or $(\text{spliceAI_DS_DL} > 0.2 \text{ \& } \text{abs}(\text{spliceAI_DP_DL}) \leq 50)$)

OR

$\text{distance_spliceJunction} \leq 100$.

OR

$\text{DNG_confidence} = \text{'High'}$ (for data set analysed with DeNovoGear) .

Use the column named "SUBSET_filter" to identify the criteria satisfied for each variant in SUBSET file.