

# SV Annotation Table

## DESCRIPTION OF EACH COLUMN

**SV\_Name:** name of the SV

**chrAnn:** Chromosome (chr1 to chr22, plus chrX, chrY, and chrM).

**startAnn:** CNV/SV start.

**endAnn:** CNV/SV end.

**variantTypeAnn:** Variant type as interpreted by the annotation pipeline; for CNVs, deletions and duplications. For CNVs, complex is used when the calling software assigns a type which cannot be ascribed as either duplication or deletion (e.g. translocations, inversions, and other events); in this case the annotation pipeline ignores the type in the database when matching it with the variant. SV types recognized by the pipeline are: duplications, deletions, inversions, and insertions. Insertions of size 1 are padded (+/- 50bp) prior to frequency calculations. BND are currently ignored.

**sizeAnn:** CNV/SV size as determined by the pipeline.

**GC\_content\_perc:** GC content based on the UCSC genome reference, hg37 or hg38.

**cytobandAnn:** Cytoband.

**numberOfGeneSymbols:** number of official gene symbol(s) for genes spanned by the CNV/SV based on the UCSC RefSeq gene definitions.

**gene\_symbol:** official gene symbol(s) for genes spanned by the CNV/SV based on the UCSC RefSeq gene definitions.

**gene\_egID:** entrez gene IDs.

**gene\_symbol\_CNVstart:** official gene symbol(s); overlap on CNV/SV start coordinate.

**gene\_symbol\_CNVend:** official gene symbol(s); overlap on CNV/SV end coordinate.

**exon\_symbol:** official gene symbol(s); exons overlap only.

**exon\_egID:** official gene entrez gene IDs, exons overlap only.

**cds\_symbol:** official gene symbol(s); coding exons overlap.

**cds\_egID:** official gene entrez IDs; coding exons overlap.

**ISCA\_haploinsufficient:** (array of) entrez gene ID, gene symbol, and score from the dosage sensitivity map, haploinsufficient phenotype defined in ClinGen (ISCA). ClinGen is a National Institutes of Health (NIH)-funded resource dedicated to building an authoritative central resource that defines the clinical relevance of genes and variants for use in precision medicine and research.

**ISCA\_triplosensitive:** (array of) entrez gene ID, gene symbol, and score for dosage sensitivity map, triplosensitive phenotype defined in ClinGen (ISCA).

**ExAC\_pLI**: probability of being loss-of-function intolerant; for more information of ExAc functional constraint scores, see Samocha *et al.* - Nature Genetics 2014 (<http://www.ncbi.nlm.nih.gov/pubmed/25086666>).

**gnomAD\_oe\_lof**: observed/expected ratio for loss of function from the genome aggregation database (gnomAD)

**gnomAD\_oe\_lof\_upper**: observed/expected upper bound loss of function from the genome aggregation database (gnomAD)

**gnomAD\_oe\_mis**: observed/expected missense ratio from the genome aggregation database (gnomAD)

**gnomAD\_oe\_mis\_upper**: observed/expected upper bound missense from the genome aggregation database (gnomAD)

**gnomAD\_pLI**: probability that a gene falls into the class of intolerant of a single LoF gene (LoF-haploinsufficient intolerant genes), from the genome aggregation database (gnomAD)

**gnomAD\_pRec**: probability that a gene falls into the class of intolerant of two LoF genes (recessive genes), from the genome aggregation database (gnomAD)

**repeatMasker\_percOverlap**: percent overlap with repeat regions (RepeatMasker annotation from UCSC).

**dirtyRegion\_percOverlap**: percent overlap with gaps (including centromeres and telomeres), and segmental duplications.

**chrRegion**: telomere/centromere tag.

**MPO\_NervousSystem**: (array of) entrez gene ID, gene symbol, and inheritance for genes in the MPO terms part of nervous system phenotype (for MPO description, see MPO\_Other).

**MPO\_Growth**: (array of) entrez gene ID, gene symbol, and inheritance for genes in the MPO terms part of growth phenotype (for MPO description, see MPO\_Other).

**MPO\_Other**: (array of) entrez gene ID, gene symbol, MPO term, and inheritance for genes in the MPO (Mammalian Phenotype Ontology) top level phenotype(s) excluding terms in nervous system or growth phenotypes, 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 confirmed on the MGI website.

**HPO\_NervousSystem**: (array of) entrez gene ID, gene symbol and inheritance for genes in the HPO terms part of nervous system phenotype (for HPO description, see HPO\_Other).

**HPO\_Growth**: (array of) entrez gene ID, gene symbol and inheritance for genes in the HPO terms part of growth system phenotype (for HPO description, see HPO\_Other).

**HPO\_Other**: (array of) entrez gene ID, gene symbol, MPO term, and inheritance for genes in the HPO (Human Phenotype Ontology) top level phenotype(s), imported from HPO, excluding nervous system and growth phenotypes. The genotype-phenotype association is typically supported by an OMIM entry.

**CGD:** (array of) entrez gene ID, gene symbol, disease name(s), and inheritance found in the Clinical Genomics Database; it 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 (AD, AR, AD/AR, XL, more complex modes); since the CGD mode of inheritance is directly added by a curator and is tied to specific genetic disorder(s), it could be considered more accurate than the mode of inheritance for top-level HPO phenotypes.

**OMIM\_MorbidMap:** (array of) entrez gene ID, gene symbol, and disorder/disease name(s) found in OMIM.

**ASDgenes:** list of genes involved in ASD/DD, with their evidence source.

**ISCA\_region:** Genomic disease region from ClinGen (ISCA).

**CNV\_ISCA\_percOverlap:** % length of CNV/SV overlapped by ClinGen (ISCA) region(s).

**ISCA\_CNV\_percOverlap:** (array of) % length of ClinGen (ISCA) region(s) overlapped by CNV/SV.

**ISCA\_CNV\_percOverlap\_max:** highest % length of ClinGen (ISCA) region overlapped by CNV/SV.

**ISCA\_matchCNVmax\_percOverlap:** % length of CNV/SV overlapped by largest ClinGen (ISCA) region.

**exon\_symbol\_ISCA:** gene symbol(s) in the region covered by ClinGen (ISCA).

**decipher\_region:** Genomic disease region from Decipher; the Database of genomic variation and Phenotype in Humans using Ensembl Resources is an interactive web-based database which incorporates a suite of tools designed to aid the interpretation of genomic variants. For more information, see: <https://decipher.sanger.ac.uk/>.

**CNV\_decipher\_percOverlap:** % length of CNV/SV overlapped by Decipher region(s).

**decipher\_CNV\_percOverlap:** % length of Decipher region(s) overlapped by CNV.

**decipher\_CNV\_percOverlap\_max:** highest % length of Decipher region(s) overlapped by CNV.

**decipher\_matchCNVmax\_percOverlap:** % length of CNV/SV overlapped by largest Decipher region.

**exon\_symbols\_Decipher:** gene symbol(s) in the region covered by Decipher.

**gnomAD\_commonSV:** % overlap with the common (greater than 1%) features in the structural variants genome aggregation database (gnomAD)

**gnomAD\_rareSV:** % overlap with the rare (less than or equal to 1%) features in the structural variants genome aggregation database (gnomAD)

**DGV\_N\_studies\_50percRecipOverlap:** Number of studies in DGV where at least one subject in the study has a variant overlapping the CNV, restricted to 50% reciprocal overlap.

**DGV\_N\_subjects\_50percRecipOverlap:** Number of subjects in DGV where the variant overlaps the CNV, restricted to 50% reciprocal overlap.

**DGVpercFreq\_subjects\_allStudies\_50percRecipOverlap:** % frequency in DGV with at least 50% reciprocal overlap; all studies combined.

**DGVpercFreq\_subjects\_coverageStudies\_50percRecipOverlap:** % frequency in DGV with at least 50% reciprocal overlap; only studies where at least one of the subjects had coverage.

**DGV\_percOverlap\_any:** (array of) % length of DGV region(s) overlapped by CNV (no cutoff used). The Database of Genomic Variants provides a comprehensive summary of structural variation in the human genome. For more information: <http://dgv.tcag.ca/dgv/app/about?ref=GRCh37/hg19>. The DVG was lifted over to obtain the corresponding intervals in the GRCh38 reference genome.

**DGV\_50percRecipOverlap:** % length covered by merged variants in DGV, restricted to those with at least 50% reciprocal overlap.

**DGV\_commonPerc:** % overlap with CNVs with a frequency higher than 1% in DGV.

**CGparentalPercFreq\_90percRecipOverlap:** frequency based on internal database - parents sequenced by Complete Genomics, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**otgMantaPercFreq\_90percRecipOverlap:** frequency based on the 1000G+ collection – parents sequenced by NovaSeq6000, called by Manta, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**otgDellyPercFreq\_90percRecipOverlap:** frequency based on the 1000G+ collection – parents sequenced by NovaSeq6000, called by Delly, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**svMantaXPercFreq\_90percRecipOverlap:** frequency based on internal database - parents sequenced by Illumina HiSeqX called by Manta, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**svManta2PercFreq\_90percRecipOverlap:** frequency based on internal database - parents sequenced by Illumina HiSeq2000/2500 called by Manta, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**svDellyXPercFreq\_90percRecipOverlap:** frequency based on internal database - parents sequenced by Illumina HiSeqX called by Delly, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**svDelly2PercFreq\_90percRecipOverlap:** frequency based on internal database - parents sequenced by Illumina HiSeq2000/2500 called by Delly, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**hsDragenPercFreq\_90percRecipOverlap:** frequency based on internal database – unrelated samples sequenced by Illumina NovaSeq 6000 with average coverage of 35x, called by Dragen 3.8.4, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**pacBioPercFreq\_90percRecipOverlap:** frequency based on internal database – unrelated 18 samples called by pbsv 2.6.2, with at least 90% reciprocal overlap, matched by variant type; the type is ignored when complex.

**nearestLeftExonBoundary:** gene with exon/intron junction closest to the left CNV/SV boundary.

**nearestLeftExonDistance:** distance of closest exon/intron junction to the left boundary; negative values indicate that the CNV/SV is upstream of the junction.

**nearestRightExonBoundary:** gene with exon/intron junction closest to the right CNV/SV boundary.

**nearestRightExonDistance:** distance of the closest exon/intron junction to the right boundary; negative values indicate that the CNV/SV is upstream of the junction.