



Batch Annotation output description

snpXplorer_input_XXXX.txt: Raw input as pasted in snpXplorer website including SNPs to annotate.

invalid_queries.txt: Any invalid SNP part of the raw input that could not be parsed by snpXplorer.

variant_annotation_combined.tsv: the final annotation of the target SNPs. Includes the following columns:

Column	Description
<i>Query Type</i>	Query type identified by snpXplorer
<i>RsID</i>	RS identifier of the variant
<i>Chromosome</i>	Chromosome
<i>Position (hg38)</i>	Position with respect to GRCh38
<i>Position (hg19)</i>	Position with respect to GRCh37
<i>Reference</i>	Reference Allele (from TOPMed)
<i>Alternative</i>	Alternative Allele (from TOPMed)
<i>Alternative Allele Frequency</i>	Alternative Allele Frequency (from TOPMed)
<i>Locus ID</i>	Locus identifier
<i>CADD Reference</i>	Reference Allele (CADD)
<i>CADD Alternative</i>	Alternative Allele (CADD)
<i>CADD Score</i>	CADD Phred Score (CADD)
<i>CADD Gene</i>	Likely affected gene (CADD)
<i>CADD Annotation Types</i>	Type of annotation (CADD)
<i>CADD Consequences</i>	SNP Consequences (CADD)
<i>CADD CpG Max</i>	Maximum CpG context score (CADD)

<i>CADD GC Max</i>	Maximum GC content score (CADD)
<i>CADD H3K27ac Max</i>	Overlap with H3K27ac histone marks (CADD)
<i>CADD H3K4me3 Max</i>	Overlap with H3K4me3 marks (CADD)
<i>CADD DNAase Max</i>	DNAse I hypersensitivity (CADD)
<i>CADD SpliceAI Acc Loss Max</i>	Predicted loss of acceptor splice site (CADD)
<i>CADD SpliceAI Don Loss Max</i>	Predicted loss of donor splice site (CADD)
<i>CADD SIFT Max</i>	SIFT Score (CADD)
<i>CADD PolyPhen Max</i>	PolyPhen Score (CADD)
<i>CADD GERPRS Max</i>	Evolutionary context (CADD)
<i>CADD PhyloP Max</i>	Nucleotide-level conservation (CADD)
<i>eQTL Reference</i>	Reference Allele for eQTL (GTEx)
<i>eQTL Alternative</i>	Alternative Allele for eQTL (GTEx)
<i>eQTL ensemble</i>	Ensemble gene ID (GTEx)
<i>eQTL Gene</i>	eQTL Gene name (GTEx)
<i>eQTL Tissue</i>	eQTL Tissue (GTEx)
<i>eQTL TSS Distance</i>	Distance of SNP-TSS (GTEx)
<i>eQTL P-value</i>	eQTL P-value (GTEx)
<i>eQTL Slope</i>	eQTL effect size of Alternative allele (GTEx)
<i>eQTL MAF</i>	Minor allele frequency (GTEx)
<i>sQTL Reference</i>	Reference Allele for sQTL (GTEx)
<i>sQTL Alternative</i>	Alternative Allele for sQTL (GTEx)
<i>sQTL ensemble</i>	Ensemble gene ID (GTEx)
<i>sQTL Gene</i>	sQTL Gene name (GTEx)
<i>sQTL Tissue</i>	sQTL Tissue (GTEx)
<i>sQTL TSS Distance</i>	Distance of SNP-TSS (GTEx)
<i>sQTL P-value</i>	sQTL P-value (GTEx)
<i>sQTL Slope</i>	sQTL effect size of Alternative allele (GTEx)
<i>sQTL MAF</i>	Minor allele frequency (GTEx)
<i>Most Likely Genes</i>	Most likely gene(s) affecting variant (snpXplorer)
<i>Source</i>	Source(s) used to identify genes (snpXplorer)

target_annotations/cadd_annotation_target.tsv: Complete CADD annotation for each variant. Columns are the same as in *variant_annotation_combined.tsv*, however, multiple

annotations per variant may exist. In the *variant_annotation_combined.tsv*, the annotation with the largest CADD Score is prioritized and reported.

target_annotations/eqtl_annotation_target.tsv: Complete eQTL annotation for each variant. Columns are the same as in *variant_annotation_combined.tsv*.

target_annotations/sqtl_annotation_target.tsv: Complete sQTL annotation for each variant. Columns are the same as in *variant_annotation_combined.tsv*.

target_annotations/gwas_annotation_target.tsv: GWAS hits identified for each variant. Note that not all GWAS traits from OpenGWAS are queried (>10,000 traits). To redundancy between similar terms, trait names were embedded using a combination of MiniLD and SapBERT language models to capture semantic similarity. Pairwise cosine similarity between embeddings was then used to cluster related traits. Columns are as follows:

Column	Description
<i>GWAS Trait</i>	Trait name
<i>Effect Allele</i>	Effect Allele
<i>Non-effect Allele</i>	Other Allele
<i>EAF</i>	Effect Allele Frequency
<i>GWAS Beta</i>	Effect size of Effect allele
<i>GWAS SE</i>	Standard error of effect size
<i>GWAS P</i>	P-value of association
<i>Sample Size</i>	Sample size of GWAS
<i>GWAS ID</i>	GWAS identifier (OpenGWAS)
<i>RsID</i>	Variant Identifier

target_annotations/sv_annotation_target.tsv: Structural variants in the vicinity of each variant. The set of structural variants was derived from [Tesi et al, 2024](#). Columns are as follows:

Column	Description
<i>Repeat Class</i>	Class of the repeat
<i>Chromosome</i>	Chromosome
<i>Start Position (hg38)</i>	Start Position in GRCh38
<i>End Position (hg38)</i>	End Position in GRCh38

<i>Length</i>	Length of the structural variant
<i>Repeat Name</i>	Repeat name
<i>Repeat Family</i>	Repeat family
<i>Color</i>	Color as shown in snpXplorer Exploration section

plots_variant_annotation/gwas_annotation_top_traits.png: Summary image of the GWAS traits (OpenGWAS) associated with each variant.

plots_variant_annotation/variant_annotation_summary.pdf: Summary image of each variant, frequency, type and chromosomal distribution.

Id_partner_annotations/ld_annotation_target.tsv: Linkage disequilibrium (LD) between input variants and all their LD partners. For LD, we used TOPMed maps. Columns are as it follows:

Column	Description
<i>Query Position</i>	Position of query variant (input variant) in GRCh38
<i>LD Partner Position (hg38)</i>	Position of the partner variant in GRCh38
<i>LD Partner RsID</i>	RsID of the partner variant
<i>LD R2</i>	LD R2 value between variants (TOPMed)
<i>LD Distance (bp)</i>	Distance between query and partner
<i>Chromosome</i>	Chromosome
<i>Position (hg38)</i>	NA
<i>Partner unique ID</i>	Unique ID of the partner variant

Id_partner_annotations/cadd_annotation_LD_partners.tsv: CADD annotation of all LD partners (i.e. variants in LD with the input variants). Columns are the same as in *cadd_annotation_target.tsv*.

Id_partner_annotations/eqtl_annotation_LD_partners.tsv: eQTL annotation of all LD partners (i.e. variants in LD with the input variants). Columns are the same as in *eqtl_annotation_target.tsv*.

ld_partner_annotations/sql_annotation_LD_partners.tsv: sQTL annotation of all LD partners (i.e. variants in LD with the input variants). Columns are the same as in *sql_annotation_target.tsv*.

gene_set_enrichment/gene_set_enrichment_results.tsv: Only available when gene-set enrichment analysis was requested. This file contains enrichment results. Note that our SNP-gene annotation procedure can link a single SNP with multiple genes, depending on annotation uncertainty. To take this uncertainty into account, gene-set enrichment analysis is performed in a sampling-based framework (100 iterations). Each iteration, a single gene per SNP is taken. The resulting gene list is used for gene-set enrichment analysis. Enrichment results are corrected for multiple tests (false discovery rate, FDR). At the end of all iterations, enrichment p-value of all terms is averaged across iterations. The final Enrichment P-value is therefore an average enrichment across all iterations, and does not need to be further corrected for multiple tests. Columns are as it follows:

Column	Description
<i>Enrichment Term ID</i>	Term ID (gProfiler2)
<i>Enrichment Term Name</i>	Term Name (gProfiler2)
<i>Enrichment Source</i>	Enrichment source (gProfiler2)
<i>Enrichment P-value</i>	P-value of enrichment (gProfiler2)
<i>Enrichment Term Size</i>	Size of the Term (gProfiler2)
<i>Enrichment Query Size</i>	Query Size (gProfiler2)
<i>Enrichment Precision</i>	Precision (gProfiler2)
<i>Enrichment Recall</i>	Recall (gProfiler2)
<i>Enrichment Term Description</i>	Term Description (gProfiler2)
<i>Enrichment Intersections</i>	Intersection between Term and Query (gProfiler2)

gene_set_enrichment/semantic_similarity_matrix.tsv: To improve interpretation of gene-set enrichment analysis, we use semantic similarity between Terms. Note that this analysis is restricted to significant terms (P-value<0.05) from GeneOntology, as it follows a tree-like structure that improves similarity measurement. We use *Lin* as semantic similarity distance. This matrix shows pairwise similarity between all significant GO terms after gene-set enrichment analysis.

gene_set_enrichment/pheatmap_lin_distance.pdf: Image showing the similarity between GO terms.

gene_set_enrichment/clustering_max_{5,8,10,15}_clusters.tsv: The semantic similarity goes through a dynamic cut tree algorithm to identify clusters. We use different parameters to allow a maximum number of clusters ranging 5-15, where 5 gives less clusters (max 5) and 15 gives more fine-tuned clustering (max 15 clusters). Depending on the use, the user might prefer one clustering over another one. Note that in case two clustering parameters (e.g., max 8 clusters and max 10 clusters) produced the same output, only one is reported. Columns are as it follows:

Column	Description
<i>Term</i>	GO term ID of the term
<i>Cluster</i>	Cluster number
<i>Cluster in dendrogram</i>	Cluster as it appears on the dendrogram image
<i>Native</i>	GO term ID of the term
<i>Name</i>	Description of the term

gene_set_enrichment/clustering_max_{5,8,10,15}_dendrogram_clusters.png: Images of the dynamic cut tree algorithm showing the dendrogram of the significant terms, colored by their cluster assignment. One image for each clustering parameter.

gene_set_enrichment/wordclouds_{5,8,10,15}_clusters: To facilitate interpretation of the enrichment clusters, we use wordclouds of the term underlying each cluster. These folders contain an image per cluster, specifying the order in which they occur in the dendrogram (left to right). There is one folder for each clustering parameter (max 5, 8, 10, 15 clusters).

pathway_prs_weights/pathway_prs_weights_{number of cluster}_clusters.tsv: These files, one for each successful clustering run (similar and aligned to the clustering above and wordclouds), report weights to be used for pathway-specific polygenic risk scores (PRS). Specifically, for each input variant, depending on what genes it was associated with (SNP-gene annotation), and depending on which significant terms the gene was involved into (gene-set enrichment analysis), it is reported the effect of that variant on each of the pathways (in proportion, summing to 1 across clusters of pathways). Pathway-PRS

constructed in this way were published multiple times ([Tesi et al., 2021](#), [Lorenzini et al., 2023](#), [Tranfa et al., 2025](#)). The file shows one variant per line, and each column represent a cluster, with the relative proportion of effect of each variant on that cluster. Using this file (with some minimal edits), users can calculate pathway-specific PRS in their cohort. We provide a tool for that, [jordan](#), which can be run both as a command-line and with a user interface.