

HOWARD Help Parameters

Contents

1	Introduction	3
2	hgvs	4
2.1	use_gene	4
2.2	use_exon	5
2.3	use_protein	5
2.4	add_protein	5
2.5	full_format	5
2.6	codon_type	5
2.7	refgene	6
2.8	refseqlink	6
3	annotation	6
3.1	parquet	7
3.1.1	annotations	8
3.2	bcftools	9
3.2.1	annotations	9
3.3	annovar	10
3.3.1	annotations	11
3.3.2	options	11
3.4	snpeff	11
3.4.1	options	11
3.4.2	stats	12
3.4.3	csvStats	12
3.5	snpsift	12
3.5.1	annotations	13
3.6	bigwig	13
3.6.1	annotations	14
3.7	exomiser	14
3.7.1	release	15
3.7.2	transcript_source	15
3.7.3	hpo	15
3.8	splice	15
3.8.1	split_mode	15
3.8.2	spliceai_distance	16
3.8.3	spliceai_mask	16
3.8.4	transcript	16
3.8.5	rm_snps	16
3.8.6	rm_annot	16
3.8.7	whitespace	17
3.9	options	17
3.9.1	annotations_update	17
3.9.2	annotations_append	17
4	calculation	17
4.1	calculations	18
4.2	calculation_config	18

5	prioritization	19
5.1	prioritizations	19
5.2	profiles	19
5.3	default_profile	20
5.4	pzfields	20
5.5	prioritization_score_mode	20
5.6	pzprefix	21
6	stats	21
6.1	stats_md	21
6.2	stats_json	21
7	query	22
7.1	query	22
7.2	query_limit	22
7.3	query_print_mode	22
8	export	22
8.1	order_by	22
8.2	include_header	23
8.3	parquet_partitions	23
9	explode	23
9.1	explode_infos	23
9.2	explode_infos_prefix	23
9.3	explode_infos_fields	23
10	transcripts	23
10.1	table	25
10.2	transcripts_info_field_json	25
10.3	transcripts_info_field_format	25
10.4	transcripts_info_json	25
10.5	transcripts_info_format	26
10.6	transcript_id_remove_version	26
10.7	transcript_id_mapping_file	26
10.8	Example of transcript ID mapping file	26
10.9	transcript_id_mapping_force	27
10.10	struct	27
10.10.1	from_column_format	27
10.10.2	from_columns_map	27
10.10.3	commons parameters	28
10.11	prioritization	29
10.11.1	profiles	29
10.11.2	default_profile	29
10.11.3	prioritization_config	30
10.11.4	prioritization_score_mode	30
10.11.5	pzprefix	30
10.11.6	pzfields	30
10.11.7	prioritization_transcripts_order	30
10.11.8	prioritization_transcripts	31
10.11.9	prioritization_transcripts_force	31
10.11.10	prioritization_transcripts_version_force	31
10.12	export	31
11	threads	32
12	samples	32
12.1	list	33
12.2	check	33

1 Introduction

HOWARD Parameters JSON file defines parameters to process annotations, calculations, prioritizations, conversions and queries.

Examples:

Parameters for annotation, calculation, prioritization, HGVS annotation and export

```
{
  "hgvs": {
    "full_format": true,
    "use_exon": true
  },
  "annotation": {
    "parquet": {
      "annotations": {
        "/path/to/database3.parquet": {
          "field1": null,
          "field2": "field2_renamed"
        },
        "/path/to/database4.vcf.gz": {
          "INFO": null
        },
        "/path/to/database5.bed.gz": {
          "INFO": null
        }
      }
    },
    "bcftools": {
      "annotations": {
        "/path/to/database6.vcf.gz": {
          "field1": null,
          "field2": "field2_renamed"
        },
        "/path/to/database7.bed": {
          "INFO": null
        }
      }
    },
    "annovar": {
      "annotations": {
        "annovar_keyword2": {
          "field1": null,
          "field2": "field2_renamed"
        },
        "annovar_keyword3": {
          "INFO": null
        }
      }
    },
    "snpeff": {
      "options": " -hgvs -noShiftHgvs -spliceSiteSize 3 -lof -oicr "
    },
    "exomiser": {
      "release": "2109",
      "transcript_source": "refseq",
      "hpo": ["HP:0001156", "HP:0001363", "HP:0011304", "HP:0010055"]
    }
  }
}
```

```

    },
    "options": {
      "append": true
    }
  },
  "calculation": {
    "operation1": null,
    "operation2": {
      "options": {
        "option1": "value1",
        "option2": "value2"
      }
    }
  },
  "prioritization": {
    "prioritizations": "config/prioritization_profiles.json",
    "profiles": ["GENOME", "GERMLINE"],
    "default_profile": "GERMLINE",
    "pzfields": ["PZScore", "PZFlag", "PZComment"],
    "prioritization_score_mode": "VaRank"
  },
  "export": {
    "include_header": true
  }
}

```

2 hgvs

HOWARD annotates variants with HGVS annotation using HUGO HGVS international Sequence Variant Nomenclature (<http://varnomen.hgvs.org/>). Annotation refers to refGene and genome to generate HGVS nomenclature for all available transcripts. This annotation adds 'hgvs' field into VCF INFO column of a VCF file. Several options are available, to add gene, exon and protein information, to generate a 'full format' detailed annotation, to choose codon format.

Examples:

HGVS annotation with operations for generate variant_id and variant type, extract HGVS from snpEff annotation, select NOMEN from snpEff HGVS with a list of transcripts of preference

```

{
  "hgvs": {
    "full_format": true,
    "use_exon": true
  }
}

```

2.1 use_gene

Add Gene information to generate HGVS annotation (e.g. 'NM_152232(**TAS1R2**):c.231T>C').

Default: **False**

Examples:

Use Gene in HGVS annotation

```

{
  "use_gene": true
}

```

2.2 use_exon

Add Exon information to generate HGVS annotation (e.g. 'NM_152232(exon2):c.231T>C'). Only if 'use_gene' is not enabled.

Default: **False**

Examples:

Use Exon in HGVS annotation

```
{
  "use_exon": true
}
```

2.3 use_protein

Use Protein level to generate HGVS annotation (e.g. 'NP_689418:p.Cys77Arg'). Can be used with 'use_exon' or 'use_gene'.

Default: **False**

Examples:

Use Protein in HGVS annotation

```
{
  "use_protein": true
}
```

2.4 add_protein

Add Protein level to DNA HGVS annotation (e.g. 'NM_152232:c.231T>C,NP_689418:p.Cys77Arg').

Default: **False**

Examples:

Add Protein level to DNA HGVS annotation

```
{
  "add_protein": true
}
```

2.5 full_format

Generates HGVS annotation in a full format (non-standard, e.g. 'TAS1R2:NM_152232:NP_689418:c.231T>C:p.Cys77Arg', 'TAS1R2:NM_152232:NP_689418:exon2:c.231T>C:p.Cys77Arg'). Full format use all information to generates an exhaustive annotation. Use specifically 'use_exon' to add exon information.

Default: **False**

Examples:

Use full format for HGVS annotation

```
{
  "full_format": true
}
```

2.6 codon_type

Amino Acide Codon format type to use to generate HGVS annotation (default '3'):

- '1': codon in 1 character (e.g. 'C', 'R')
- '3': codon in 3 character (e.g. 'Cys', 'Arg')
- 'FULL': codon in full name (e.g. 'Cysteine', 'Arginine')

Type: **str**

Choices: ['1', '3', 'FULL']

Default: 3

Examples:

Amino Acide Codon format with 1 character

```
{  
  "codon_type": "1"  
}
```

Amino Acide Codon format with 3 character

```
{  
  "codon_type": "3"  
}
```

Amino Acide Codon format with full name

```
{  
  "codon_type": "FULL"  
}
```

2.7 refgene

Path to refGene annotation file (see HOWARD User Guide).

Type: **Path**

Default: **None**

Examples:

Path to refSeq file

```
{  
  "refgene": "~/howard/databases/refseq/current/hg19/ncbiRefSeq.txt"  
}
```

2.8 refseqlink

Path to refGeneLink annotation file (see HOWARD User Guide).

Type: **Path**

Default: **None**

Examples:

Path to refSeq file

```
{  
  "refseqlink": "~/howard/databases/refseq/current/hg19/ncbiRefSeqLink.txt"  
}
```

3 annotation

Annotation process using HOWARD algorithms or external tools.

For HOWARD Parquet algorithm, provide the list of database files available (formats such as Parquet, VCF, TSV, duckDB, JSON) and select fields (rename possible, 'INFO' keyword for all fields), or use 'ALL' keyword to detect available databases.

For external tools, such as Annovar, snpEff and Exomiser, specify parameters such as annotation keywords (Annovar) and options (depending on the tool), and select fields (BCFtools and Annovar, field rename available).

Examples:

Annotation with multiple tools in multiple formats with multiple options

```
{
  "annotation": {
    "parquet": {
      "annotations": {
        "/path/to/database3.parquet": {
          "field1": null,
          "field2": "field2_renamed"
        },
        "/path/to/database4.vcf.gz": {
          "INFO": null
        }
        "/path/to/database5.bed.gz": {
          "INFO": null
        }
      }
    }
    "bcftools": {
      "annotations": {
        "/path/to/database6.vcf.gz": {
          "field1": null,
          "field2": "field2_renamed"
        },
        "/path/to/database7.bed": {
          "INFO": null
        }
      }
    }
    "annovar": {
      "annotations": {
        "annovar_keyword2": {
          "field1": null,
          "field2": "field2_renamed"
        },
        "annovar_keyword3": {
          "INFO": null
        }
      }
    }
    "snpeff": {
      "options": " -hgvs -noShiftHgvs -spliceSiteSize 3 -lof -oicr "
    }
    "exomiser": {
      "release": "2109",
      "transcript_source": "refseq",
      "hpo": ["HP:0001156", "HP:0001363", "HP:0011304", "HP:0010055"]
    }
    "options": {
      "append": true
    }
  }
}
```

3.1 parquet

Annotation process using HOWARD Parquet algorithm, for the list of databases available (formats such as Parquet, VCF, TSV, duckDB, JSON).

Examples:

Annotation with multiple databases in multiple formats

```
{
  "parquet": {
    "annotations": {
      "/path/to/database3.parquet": {
        "field1": null,
        "field2": "field2_renamed",
      },
      "/path/to/database4.vcf.gz": {
        "INFO": null
      }
      "/path/to/database5.bed.gz": {
        "INFO": null
      }
    }
  }
}
```

3.1.1 annotations

Specify a list of databases files available (formats such as Parquet, VCF, TSV, duckDB, JSON). This parameter enables users to select specific database fields and optionally rename them (e.g. "field": null' to keep field name, "field": "new__name" to rename field). Use 'INFO' keyword to select all fields within the database INFO/Tags header (e.g. "INFO": null'). Use 'ALL' keyword to select all fields within the database regardless INFO/Tags header (e.g. "ALL": null').

For add all availalbe databases files, use 'ALL' keyword, with filters on type and release (e.g. 'ALL', 'ALL:parquet:current', 'ALL:parquet,vcf:current,devel').

If a full path is not provided, the system will automatically detect files within database folders (see Configuration doc) and assembly (see Parameter option).

Examples:

Annotation with dbSNP database with INFO/tags fields, and dbNSFP databases with all fields

```
{
  "annotations": {
    "tests/databases/annotations/current/hg19/avsnip150.parquet": {
      "INFO": null
    }
    "tests/databases/annotations/current/hg19/dbnsfp42a.parquet": {
      "ALL": null
    }
  }
}
```

Annotation with dbNSFP (only PolyPhen, ClinVar and REVEL score), and rename fields

```
{
  "annotations": {
    "tests/databases/annotations/current/hg19/dbnsfp42a.parquet": {
      "Polyphen2_HDIV_pred": "PolyPhen",
      "ClinPred_pred": "ClinVar",
      "REVEL_score": null
    }
  }
}
```

Annotation with refSeq as a BED file

```
{
  "annotations": {
```



```

        "tests/databases/annotations/current/hg19/refGene.bed": {
            "INFO": null
        }
    }
}

```

Annotation with dbNSFP REVEL annotation (as a VCF file) within configured annotation databases folders (default: '~/howard/databases/annotations/current') and assembly (default: 'hg19')

```

{
    "annotations": {
        "dbnsfp42a.REVEL.vcf.gz": {
            "REVEL_score": null,
            "REVEL_rankscore": null
        }
    }
}

```

Annotation with all available databases in Parquet for "current release

```

{
    "parquet": {
        "annotations": {
            "ALL": {
                "formats": ["parquet"],
                "releases": ["current"]
            }
        }
    }
}

```

3.2 bcftools

Annotation process using BCFTools. Provide a list of database files and annotation fields.

Examples:

Annotation with multiple databases in multiple formats

```

{
    "parquet": {
        "bcftools": {
            "/path/to/database1.vcf.gz": {
                "field1": null,
                "field2": "field2_renamed"
            },
            "database2.bed.gz": {
                "INFO": null
            }
        }
    }
}

```

3.2.1 annotations

Specify the list of database files in formats VCF or BED. Files need to be compressed and indexed.

This parameter enables users to select specific database fields and optionally rename them (e.g. "field": null' to keep field name, "field": "new_name" to rename field). Use 'INFO' or 'ALL' keyword to select all fields within the database INFO/Tags header (e.g. "INFO": null, "ALL": null').

If a full path is not provided, the system will automatically detect files within database folders (see Configuration doc) and assembly (see Parameter option).

Examples:

Annotation with dbSNP database with all fields

```
{
  "annotations": {
    "tests/databases/annotations/current/hg19/avsnp150.vcf.gz": {
      "INFO": null
    }
  }
}
```

Annotation with dbNSFP (only PolyPhen, ClinVar and REVEL score), and rename fields

```
{
  "annotations": {
    "tests/databases/annotations/current/hg19/dbnsfp42a.vcf.gz": {
      "Polyphen2_HDIV_pred": "PolyPhen",
      "ClinPred_pred": "ClinVar",
      "REVEL_score": null
    }
  }
}
```

Annotation with refSeq as a BED file

```
{
  "annotations": {
    "tests/databases/annotations/current/hg19/refGene.bed": {
      "INFO": null
    }
  }
}
```

Annotation with dbNSFP REVEL annotation (as a VCF file) within configured annotation databases folders (default: '~/howard/databases/annotations/current') and assembly (default: 'hg19')

```
{
  "annotations": {
    "dbnsfp42a.REVEL.vcf.gz": {
      "REVEL_score": null,
      "REVEL_rankscore": null
    }
  }
}
```

3.3 annovar

Annotation process using Annovar tool. Provides a list of keywords to select Annovar databases, and defines Annovar options (see Annovar documentation).

Examples:

Annotation with multiple Annovar databases, with fields selection, and Annovar options

```
{
  "annovar": {
    "annotations": {
      "annovar_keyword2": {
        "field1": null,
        "field2": "field2_renamed",
      },
      "annovar_keyword3": {
        "INFO": null
      }
    }
  }
}
```

```

    }
  }
}

```

3.3.1 annotations

List of keywords referring to Annovar databases (see Annovar Databases documentation), with a list of selected fields for each of them (rename available)

Examples:

Annotation with ClinVar (fields CLNSIG and CLNDN renamed) and Cosmic (all fields)

```

{
  "annotations": {
    "clinvar_20221231": {
      "CLNSIG": "ClinVar_class",
      "CLNDN": "ClinVar_desease",
    },
    "cosmic70": {
      "INFO": null
    },
  }
}

```

3.3.2 options

List of options available with Annovar tool (see Annovar documentation). As example, these options allows to define splicing threshold or HGVS annotation with refGene database

Examples:

HGVS Annotation with refGene (add 'refGene' to 'annotations') and a splicing threshold as 3

```

{
  "options": {
    "splicing_threshold": 3,
    "argument": "'-hgvs'"
  }
}

```

3.4 snpeff

Annotation process using snpEff tool and options (see snpEff documentation).

Examples:

Annotation with snpEff databases, with options for HGVS annotation and additional tags.

```

{
  "snpeff": {
    "options": " -hgvs -noShiftHgvs -spliceSiteSize 3 -lof -oicr "
  }
}

```

3.4.1 options

String (as command line) of options available such as:

- filters on variants (regions filter, specific changes as intronic or downstream)
- annotation (e.g. HGVS, loss of function)
- database (e.g. only protein coding transcripts, splice sites size)

Examples:

Annotation with snpEff databases, with options to generate HGVS annotation, specify to not shift variants according to HGVS notation, define splice sites size to 3, add loss of function (LOF), Nonsense mediated decay and OICR tags.

```
{
  "options": " -hgvs -noShiftHgvs -spliceSiteSize 3 -lof -oicr "
}
```

3.4.2 stats

HTML file for snpEff stats. Use keyword 'OUTPUT' to generate file according to output file.

Examples:

Annotation with snpEff databases, and generate a specific stats in HTML format.

```
{
  "stats": "/path/to/stats.html"
}
```

Annotation with snpEff databases, and generate stats in HTML format associated with output file.

```
{
  "stats": "OUTPUT.html"
}
```

3.4.3 csvStats

CSV file for snpEff stats. Use keyword 'OUTPUT' to generate file according to output file.

Examples:

Annotation with snpEff databases, and generate a specific stats in CSV format.

```
{
  "csvStats": "/path/to/stats.csv"
}
```

Annotation with snpEff databases, and generate stats in CSV format associated with output file.

```
{
  "csvStats": "OUTPUT.csv"
}
```

3.5 snpsift

Annotation process using snpSift. Provide a list of database files and annotation fields.

Examples:

Annotation with multiple databases in multiple formats

```
{
  "snpsift": {
    "annotations": {
      "/path/to/database1.vcf.gz": {
        "field1": null,
        "field2": null
      },
      "/path/to/database2.vcf.gz": {
        "field1": null,
        "field2": null
      }
    }
  }
}
```

```

    }
}

```

3.5.1 annotations

Specify the list of database files in formats VCF. Files need to be compressed and indexed.

This parameter enables users to select specific database fields and optionally rename them (e.g. "field": null' to keep field name, "field": "new_name" to rename field). Use 'INFO' or 'ALL' keyword to select all fields within the database INFO/Tags header (e.g. "INFO": null', "ALL": null').

If a full path is not provided, the system will automatically detect files within database folders (see Configuration doc) and assembly (see Parameter option).

Examples:

Annotation with dbSNP database with all fields

```

{
  "annotations": {
    "tests/databases/annotations/current/hg19/avsnp150.vcf.gz": {
      "INFO": null
    }
  }
}

```

Annotation with dbNSFP (only PolyPhen, ClinVar and REVEL score)

```

{
  "annotations": {
    "tests/databases/annotations/current/hg19/dbnsfp42a.vcf.gz": {
      "Polyphen2_HDIV_pred": "PolyPhen",
      "ClinPred_pred": "ClinVar",
      "REVEL_score": null
    }
  }
}

```

Annotation with dbNSFP REVEL annotation (as a VCF file) within configured annotation databases folders (default: '~/howard/databases/annotations/current') and assembly (default: 'hg19')

```

{
  "annotations": {
    "dbnsfp42a.REVEL.vcf.gz": {
      "REVEL_score": null,
      "REVEL_rankscore": null
    }
  }
}

```

3.6 bigwig

Annotation process using BigWig files. Provide a list of database files in BigWig format ('.bw') and annotation fields.

Examples:

Annotation with multiple databases in BigWig format

```

{
  "bigwig": {
    "annotations": {
      "/path/to/database1.bw": {
        "field1": null,
        "field2": null
      },

```

```

        "/path/to/database2.bw": {
            "field1": null,
            "field2": null
        }
    }
}

```

3.6.1 annotations

Specify the list of database files in BigWig format.

This parameter enables users to select specific database fields and optionally rename them (e.g. "field": null' to keep field name, "field": "new_name" to rename field). Use 'INFO' or 'ALL' keyword to select all fields within the database INFO/Tags header (e.g. "INFO": null', "ALL": null').

If a full path is not provided, the system will automatically detect files within database folders (see Configuration doc) and assembly (see Parameter option).

A URL can be provided as a database file (experimental). In this case, associated header file will be automatically generated with ua uniq value as the name of the file (cleaned for avoid special characters, and '.bw' extension).

Examples:

Annotation with GERP database with all fields

```

{
  "annotations": {
    "tests/databases/annotations/current/hg19/gerp.bw": {
      "INFO": null
    }
  }
}

```

Annotation with GERP (only gerp score)

```

{
  "annotations": {
    "tests/databases/annotations/current/hg19/gerp.bw": {
      "gerp": "GERP_score"
    }
  }
}

```

Annotation with GERP from a distante database (experimental)

```

{
  "annotations": {
    "https://hgdownload.soe.ucsc.edu/gbdb/hg19/bbi/All_hg19_RS.bw": {
      "INFO": null
    }
  }
}

```

3.7 exomiser

Annotation process using Exomiser tool and options (see Exomiser website documentation).

Examples:

Annotation with Exomiser, using database release '2109', transcripts source as UCSC and a list of HPO terms.

```

{
  "exomiser": {
    "release": "2109"
    "transcript_source": "refseq"
  }
}

```

```

    "hpo": ["HP:0001156", "HP:0001363", "HP:0011304", "HP:0010055"]
  }
}

```

3.7.1 release

Release of Exomiser database. This option replace the 'release' variable in 'application.properties' file (see 'exomiser_application_properties' option). The release will be downloaded if it is not available locally.

Examples:

Annotation with release '2109' of Exomiser database.

```

{
  "release": "2109"
}

```

3.7.2 transcript_source

Transcription source of Exomiser. This option replace the 'transcript_source' variable in 'application.properties' file (see 'exomiser_application_properties' option). The release will be downloaded if it is not available locally.

Examples:

Annotation with transcription source 'refseq' of Exomiser.

```

{
  "transcript_source": "refseq"
}

```

3.7.3 hpo

List of HPO for Exomiser. This option replace the 'hpo' variable in 'application.properties' file (see 'exomiser_application_properties' option). The release will be downloaded if it is not available locally.

Examples:

Annotation with a list of 4 HPO for Exomiser.

```

{
  "hpo": ["HP:0001156", "HP:0001363", "HP:0011304", "HP:0010055"]
}

```

3.8 splice

Annotation process using Splice tool and options. This annotation will be processed only for variants that are not already annotated (i.e. without annotation like 'SpliceAI_*' and 'SPiP_*')

Examples:

Annotation with Splice, using database splice mode ('one'), spliceAI distance (500) and spliceAI mask (1).

```

{
  "splice": {
    "split_mode": "one",
    "spliceai_distance": 500,
    "spliceai_mask": 1
  }
}

```

3.8.1 split_mode

Split mode of Exomiser database (default 'one'):

- all: report all annotated transcript for one gene.
- one: keep only the transcript with the most pathogenic score (in case of identical score, take the first).

- list: keep transcript provided in transcript file, if no matching transcript in file 'one' mode is activated.
- mixed: 'one' mode, if identical score, list mode is activated.

Examples:

Split mode to report all annotated transcript for one gene.

```
{
  "split_mode": "all"
}
```

3.8.2 spliceai_distance

Maximum distance between the variant and gained/lost splice site (default: 500).

Examples:

Maximum distance of '500' between variant and splice site.

```
{
  "spliceai_distance": 500
}
```

3.8.3 spliceai_mask

Mask scores representing annotated acceptor/donor gain and unannotated acceptor/donor loss (default: 1).

Examples:

Mask score of '1' for acceptor/donor gain and loss.

```
{
  "spliceai_mask": 1
}
```

3.8.4 transcript

Path to a list of transcripts of preference (default "").

Examples:

Path to file of transcripts.

```
{
  "transcript": "tests/data/transcripts.tsv"
}
```

3.8.5 rm_snps

Do not consider SNV for the analysis, only Indels and MNV (default 'false').

Examples:

Analysing only non SNV.

```
{
  "rm_snps": "true"
}
```

3.8.6 rm_annot

Remove existing annotation before analysing (default 'true').

Examples:

Remove annotation before analysing.


```
{
  "rm_annot": "true"
}
```

3.8.7 whitespace

Remove spaces in INFO field, 'true' to remove (default 'true').

Examples:

Remove spaces in INFO field.

```
{
  "whitespace": "true"
}
```

3.9 options

Options for annotations, such as annotation strategy (skip if exists, update, append)

Examples:

Annotation with Parquet databases, with update annotation strategy.

```
{
  "options": {
    "update": true
  }
}
```

3.9.1 annotations__update

Update option for annotation (only for Parquet annotation). If True, annotation fields will be removed and re-annotated. These options will be applied to all annotation databases.

Default: **False**

Examples:

Apply update on all annotation fields for all databases.

```
{
  "update": true
}
```

3.9.2 annotations__append

Append option for annotation (only for Parquet annotation). If True, annotation fields will be annotated only if not annotation exists for the variant. These options will be applied to all annotation databases.

Default: **False**

Examples:

Apply append on all annotation fields for all databases.

```
{
  "append": true
}
```

4 calculation

Calculation process operations that are defined in a Calculation Configuration JSON file. List available calculation operations with possible options (see Calculation JSON file help).

Examples:

Calculation of operations 'operation1' and 'operation2' (with options) defined in 'calculation_config.json' file

```
{
  "calculation": {
    "calculations": {
      "operation1": null,
      "operation2": {
        "options": {
          "option1": "value1",
          "option2": "value2"
        }
      }
    },
    "calculation_config": "calculation_config.json"
  }
}
```

4.1 calculations

List of operations to process with possible options (see Calculation JSON file help).

Examples:

Calculation with operations for generate variant_id and variant type, extract HGVS from snpEff annotation, select NOMEN from snpEff HGVS with a prioritized transcript (from prioritization transcript calculation) and list of transcripts of preference, with a specific NOMEN pattern

```
{
  "calculations": {
    "variant_id": null,
    "vartype": null,
    "snpeff_hgvs": null,
    "NOMEN": {
      "options": {
        "hgvs_field": "snpeff_hgvs",
        "transcripts": "tests/data/transcripts.tsv",
        "transcripts_table": "variants",
        "transcripts_column": "PZTTranscript",
        "transcripts_order": ["column", "file"],
        "pattern": "GNOMEN:TNOMEN:ENOMEN:CNOMEN:RNOMEN:NNOMEN:PNOMEN"
      }
    }
  }
}
```

4.2 calculation_config

Calculation configuration JSON file.

Type: Path

Default: None

Examples:

Calculation configuration JSON file as an option

```
{
  "calculation_config": "calculation_config.json"
}
```

5 prioritization

Prioritization process use a JSON configuration file that defines all profiles that can be used. By default, all profiles will be calculated from the JSON configuration file, and the first profile will be considered as default. Prioritization annotations (INFO/tags) will be generated depending of a input list (default 'PZScore' and 'PZFlag'), for all profiles (e.g. 'PZScore_GERMLINE' for 'GERMLINE' profile) and for default profile (e.g. 'PZScore' for default). Prioritization score mode is 'HOWARD' by default.

Examples:

Prioritization with 'GENOME' and 'GERMLINE' (default) profiles, from a list of configured profiles, only 3 prioritization fields returned, and score calculated in 'VaRank' mode

```
{
  "prioritization": {
    "prioritizations": "config/prioritization_profiles.json",
    "profiles": ["GENOME", "GERMLINE"],
    "default_profile": "GERMLINE",
    "pzfields": ["PZScore", "PZFlag", "PZComment"],
    "prioritization_score_mode": "VaRank"
  }
}
```

5.1 prioritizations

Prioritization configuration profiles JSON file defining profiles to calculate. All configured profiles will be calculated by default (see 'profiles' parameter). First profile will be considered as 'default' if none are provided (see 'default_profile' parameter). Default score calculation mode is 'HOWARD'. This option refers to the quick prioritization command line parameter `--prioritizations`.

Type: `str`

Default: `None`

Examples:

Prioritization configuration profile JSON file

```
{
  "prioritizations": "config/prioritization_profiles.json"
}
```

5.2 profiles

Prioritization profiles to consider, from the list of configured profiles. If empty, all configured profiles will be calculated. First profile will be considered as 'default' if none are provided (see 'default_profile' parameter). Prioritization annotations (INFO/tags) will be generated for all these profiles (e.g. 'PZScore_GERMLINE' for 'GERMLINE' profile).

Type: `str`

Default: `None`

Examples:

Prioritization with 'GERMLINE' profile only

```
{
  "profiles": ["GERMLINE"]
}
```

Prioritization with 'GENOME' and 'GERMLINE' profiles

```
{
  "profiles": ["GENOME", "GERMLINE"]
}
```

5.3 default_profile

Prioritization default profile from the list of processed profiles. Prioritization annotations (INFO/tags) will be generated for this default profile (e.g. 'PZScore', 'PZFlags').

Type: **str**

Default: **None**

Examples:

```
Prioritization default profile 'GERMLINE'
{
  "default_profile": "GERMLINE"
}
```

5.4 pzfields

Prioritization annotations (INFO/tags) to generate. By default 'PZScore', 'PZFlags'.

Prioritization fields can be selected from:

- PZScore: calculated score from all passing filters, depending of the mode
- PZFlag: final flag ('PASS' or 'FILTERED'), with strategy that consider a variant is filtered as soon as at least one filter do not pass. By default, the variant is considered as 'PASS' (no filter pass)
- PZComment: concatenation of all passing filter comments
- PZTags: combinason of score, flags and comments in a tags format (e.g. 'PZFlag#PASS|PZScore#15|PZComment#Described on ...')
- PZInfos: information about passing filter criteria

Type: **str**

Default: **PZScore,PZFlag**

Examples:

```
Prioritization annotations (INFO/tags) list
{
  "pzfields": ["PZScore", "PZFlag", "PZComment"]
}
```

5.5 prioritization_score_mode

Prioritization score can be calculated following multiple mode. The HOWARD mode will increment scores of all passing filters (default). The VaRank mode will select the maximum score from all passing filters.

Type: **str**

Choices: ['HOWARD', 'VaRank']

Default: **HOWARD**

Examples:

```
Prioritization score calculation mode 'HOWARD'
{
  "prioritization_score_mode": "HOWARD"
}

Prioritization score calculation mode 'VaRank'
{
  "prioritization_score_mode": "VaRank"
}
```

5.6 pzprefix

Prioritization prefix for all annotations generated by prioritization.

Type: `str`

Default: `PZ`

Examples:

Prioritization prefix by default ('PZ'):

```
{  
  "pzprefix": "PZ"  
}
```

Specific prioritization prefix:

```
{  
  "pzprefix": "Prioritization_"  
}
```

Specific prioritization prefix for transcript (see below):

```
{  
  "pzprefix": "PZT"  
}
```

6 stats

Statistics on loaded variants.

6.1 stats_md

Stats Output file in Markdown format.

Type: `Path`

Default: `None`

Examples:

Export statistics in Markdown format

```
{  
  "stats_md": "/tmp/stats.pdf"  
}
```

6.2 stats_json

Stats Output file in JSON format.

Type: `Path`

Default: `None`

Examples:

Export statistics in JSON format

```
{  
  "stats_json": "/tmp/stats.json"  
}
```

7 query

Query options tools. Mainly a SQL query, based on 'variants' table corresponding on input file data, or a independant query. Print options for 'query' tool allow limiting number of lines and choose printing mode.

Type: `str`

Default: `None`

7.1 query

Query in SQL format (e.g. 'SELECT * FROM variants LIMIT 50').

Type: `str`

Default: `None`

Examples:

Simple query to show all variants file

```
SELECT "#CHROM", POS, REF, ALT, INFO
FROM variants
```

7.2 query_limit

Limit of number of row for query (only for print result, not output).

Type: `int`

Default: `10`

7.3 query_print_mode

Print mode of query result (only for print result, not output). Either `None` (native), 'markdown', 'tabulate' or disabled.

Type: `str`

Choices: `[None, 'markdown', 'tabulate', 'disabled']`

Default: `None`

8 export

Export options for output files, such as data order, include header in output and hive partitioning.

8.1 order_by

List of columns to sort the result-set in ascending or descending order. Use SQL format, and keywords `ASC` (ascending) and `DESC` (descending). If a column is not available, order will not be considered. Order is enable only for compatible format (e.g. TSV, CSV, JSON). Examples: 'ACMG_score DESC', 'PZFlag DESC, PZScore DESC'.

Type: `str`

Default:

Examples:

Order by ACMG score in descending order

```
{
  "order_by": "ACMG_score DESC"
}
```

Order by PZFlag and PZScore in descending order

```
{
  "order_by": "PZFlag DESC, PZScore DESC"
}
```

8.2 include_header

Include header (in VCF format) in output file. Only for compatible formats (tab-delimiter format as TSV or BED).

Default: `False`

8.3 parquet_partitions

Parquet partitioning using hive (available for any format). This option is faster parallel writing, but memory consuming. Use 'None' (string) for NO partition but split parquet files into a folder. Examples: '#CHROM', '#CHROM,REF', 'None'.

Type: `str`

Default: `None`

9 explode

Explode options for INFO/tags annotations within VCF files.

9.1 explode_infos

Explode VCF INFO/Tag into 'variants' table columns.

Default: `False`

9.2 explode_infos_prefix

Explode VCF INFO/Tag with a specific prefix.

Type: `str`

Default:

9.3 explode_infos_fields

Explode VCF INFO/Tag specific fields/tags. Keyword `*` specify all available fields, except those already specified. Pattern (regex) can be used, such as `.*_score` for fields named with `'_score'` at the end. Examples:

- 'HGVS,SIFT,Clinvar' (list of fields)
- 'HGVS,*,Clinvar' (list of fields with all other fields at the end)
- 'HGVS,.*_score,Clinvar' (list of 2 fields with all scores in the middle)
- 'HGVS,.*_score,*' (1 field, scores, all other fields)
- 'HGVS,._score' (1 field, all other fields, all scores)

Type: `str`

Default: `*`

10 transcripts

Transcripts information to create transcript view. Useful to add transcripts annotations in INFO field, to calculate transcripts specific scores (see Calculation JSON file help), to merge and map transcript IDs (e.g. from Ensembl to refSeq), or prioritize transcripts (see Priorization JSON file help).

Type: `dict`

Default: `{}`

Examples:

Transcripts information from snpEff and dbNSFP annotation

```

{
  "transcripts": {
    "table": "transcripts",
    "transcripts_info_field_json": "transcripts_json",
    "transcripts_info_field_format": "transcripts_ann",
    "transcripts_info_json": "transcripts_json",
    "transcripts_info_format": "transcripts_format",
    "transcript_id_remove_version": true,
    "transcript_id_mapping_file": "transcripts.for_mapping.tsv",
    "transcript_id_mapping_force": false,
    "struct": {
      "from_column_format": [
        {
          "transcripts_column": "ANN",
          "transcripts_infos_column": "Feature_ID",
          "column_clean": true,
          "column_case": "lower"
        }
      ],
      "from_columns_map": [
        {
          "transcripts_column": "Ensembl_transcriptid",
          "transcripts_infos_columns": [
            "genename",
            "Ensembl_geneid",
            "LIST_S2_score",
            "LIST_S2_pred"
          ],
          "column_rename": {
            "LIST_S2_score": "LISTScore",
            "LIST_S2_pred": "LISTPred"
          },
        },
        {
          "transcripts_column": "Ensembl_transcriptid",
          "transcripts_infos_columns": [
            "genename",
            "VARITY_R_score",
            "Aloft_pred"
          ]
        }
      ]
    },
    "prioritization": {
      "profiles": ["transcripts"],
      "prioritization_config": "config/prioritization_transcripts_profiles.json",
      "pzprefix": "PZT",
      "pzfields": ["Score", "Flag", "LISTScore", "LISTPred"],
      "prioritization_transcripts_order": {
        "PZTFlag": "ASC",
        "PZTScore": "DESC"
      },
      "prioritization_score_mode": "HOWARD",
      "prioritization_transcripts": null,
      "prioritization_transcripts_force": false,
      "prioritization_transcripts_version_force": false
    },
    "export": {
      "output": "/tmp/output.tsv.gz"
    }
  }
}

```



```

    }
  }
}

```

10.1 table

Transcripts table name to create.

Type: **str**

Default: **transcripts**

Examples:

Name of transcript table:

```

{
  "table": "transcripts"
}

```

10.2 transcripts_info_field_json

Transcripts INFO field name to add in VCF INFO field in JSON format.

Type: **str**

Default: **None**

Examples:

Transcripts INFO field name:

```

{
  "transcripts_info_field_json": "transcripts_json"
}

```

10.3 transcripts_info_field_format

Transcripts INFO field name to add in VCF INFO field in structured format.

Type: **str**

Default: **None**

Examples:

Transcripts INFO field name:

```

{
  "transcripts_info_field_format": "transcripts_ann"
}

```

10.4 transcripts_info_json

Transcripts column name to add to transcripts table in JSON format.

Type: **str**

Default: **None**

Examples:

Transcripts column name:

```

{
  "transcripts_info_json": "transcripts_json"
}

```

10.5 transcripts_info_format

Transcripts column name to add to transcripts table in structured format.

Type: `str`

Default: `None`

Examples:

Transcripts column name:

```
{
  "transcripts_info_format": "transcripts_format"
}
```

10.6 transcript_id_remove_version

When merging and mapping transcript IDs, remove possible version of transcript (e.g 'NM_123456.2' to 'NM_123456').

Type: `Boolean`

Default: `False`

Examples:

Remove transcript version when merging and mapping:

```
{
  "transcript_id_remove_version": true
}
```

10.7 transcript_id_mapping_file

When merging and mapping transcript IDs, indicate a transcript mapping file that provides mapping between transcripts IDs (useful to map refSeq and Ensembl transcript IDs).

Type: `Path`

Default: `None`

Examples:

Transcript IDs mapping file:

```
{
  "transcript_id_mapping_file": "My_transcripts_mapping_file.tsv.gz"
}
```

10.8 Example of transcript ID mapping file

Transcript IDs mapping file is a tab-delimited file with 2 columns:

- first column corresponds to the reference transcript ID to use
- second column correspond to an alias of the reference transcript

Second column can be empty (no alias is provided). Transcript IDs can include version or not (see `transcript_id_remove_version` section)

Examples:

Example of transcripts mapping file:

```
NM_001005484    ENST00000641515.1
NM_005228.5    ENST00000275493.7
NM_001346897    ENSG00000146648.21
NM_001346941
NM_005228
```

10.9 transcript_id_mapping_force

When merging and mapping transcript IDs, allows to filter transcript IDs only if they are present in first column of the transcript mapping file (see transcript_id_mapping_file section).

Type: `Boolean`

Default: `False`

Examples:

Filter transcripts IDs only if present in mapping file:

```
{
  "transcript_id_mapping_force": true
}
```

10.10 struct

Structure of transcripts information, corresponding to columns dedicated to transcripts, such as:

- 'from_column_format': a uniq annotation field with a specific format, like snpEff annotation,
- 'from_columns_map': a list of annotation fields corresponding to transcripts in another specific field, like dbNSFP annotation.

Some parameters are commons between these structure (e.g. 'column_rename', 'column_clean' and 'column_case').

Type: `dict`

Default: `{}`

10.10.1 from_column_format

Structure of transcripts information from a uniq annotation field with a specific format (such as snpEff annotation):

- 'transcripts_column' correspond to INFO field with annotations
- 'transcripts_infos_column' correspond to transcription ID annotations field within INFO field

Column can be renamed, cleaned and/or case changed (see below).

Type: `dict`

Default: `{}`

Examples:

Structure from snpEff annotation (columns names must be clean or changed for standard snpEff annotations):

```
{
  "from_column_format": [
    {
      "transcripts_column": "ANN",
      "transcripts_infos_column": "Feature_ID",
      "column_rename": null,
      "column_clean": true,
      "column_case": null
    }
  ]
}
```

10.10.2 from_columns_map

list of annotation fields corresponding to transcripts in another specific field (such as dbNSFP annotation):

- 'transcripts_column' correspond to INFO field with transcription ID
- 'transcripts_infos_columns' correspond to a list of INFO fields with transcript information

Column can be renamed, cleaned and/or case changed (see below).

Type: dict

Default: {}

Examples:

Structure from dbNSFP annotations (with 2 columns renamed):

```
{
  "from_columns_map": [
    {
      "transcripts_column": "Ensembl_transcriptid",
      "transcripts_infos_columns": [
        "genename",
        "Ensembl_geneid",
        "LIST_S2_score",
        "LIST_S2_pred"
      ],
      "column_rename": {
        "LIST_S2_score": "LISTScore",
        "LIST_S2_pred": "LISTPred"
      },
      "column_clean": false,
      "column_case": null
    }
  ]
}
```

10.10.3 commons parameters

Some parameters are commons between these structure:

- 'column_rename': dict defining mapping of column name changes
- 'column_clean': if true, clean column name to remove not alphanum characters not allowed in VCF (e.g. space, dash)
- 'column_case': rename column into lowercase ('lower') or uppercase ('upper')

Combining 'column_clean' and 'column_case' ensure well formed VCF field name and merging identical columns (e.g. same field 'Gene_Name', 'Gene name' and 'genename' from multiple sources). However, controlling column names through 'column_rename' is much more efficient.

Examples:

Commons parameters by default:

```
{
  "column_rename": null,
  "column_clean": false,
  "column_case": null
}
```

Parameters to change 2 column names:

```
{
  "column_rename": {
    "LIST_S2_score": "LISTScore",
    "LIST_S2_pred": "LISTPred"
  },
  "column_clean": false,
  "column_case": null
}
```

Parameters to ensure well-named column from format annotations field (such as snpEff):

```
{
  "column_rename": null,
  "column_clean": true,
  "column_case": null
}
```

10.11 prioritization

Prioritization parameters for transcripts (see Prioritization section and Priorization JSON file help), defining prioritization criteria with a configuration file of all available profiles and which profiles to use, which prioritization method to use (e.g. 'HOWARD', 'VaRank').

Prioritized transcripts fields can be defined to provide VCF fields specific to the choosen transcript, using a specific prefix (e.g. 'PZTScore', 'PZTFlag'). The selected 'best' transcript ID is always provided (e.g. 'PZTTranscript'). Extra annotation fields can also be defined (e.g. 'LISTScore', 'LISTPred'), as well as prioritizations informations from multiple profiles (e.g. 'PZTScore_myprofile', 'PZTScore_myotherprofile' in 'pzfields' section with 2 profiles 'myprofile' and 'myotherprofile' in 'profiles' section).

In order to choose the 'best' transcript, parameteres can define order of transcripts (by annotation columns or a preference transcripts file), by dealing with transcript versions.

Type: dict

Default: {}

Examples:

Prioritization of transcripts in 'HOWARD' mode with 'transcripts' profiles available in a configuration JSON file, with 'PZT' as prefix:

```
{
  "prioritization": {
    "profiles": ["transcripts"],
    "default_profile": "transcripts",
    "prioritization_config": "config/prioritization_transcripts_profiles.json",
    "prioritization_score_mode": "HOWARD",
    "pzprefix": "PZT",
    "pzfields": ["Score", "Flag", "LISTScore", "LISTPred"],
    "prioritization_transcripts_order": {
      "PZTFlag": "ASC",
      "PZTScore": "DESC"
    },
    "prioritization_transcripts": null,
    "prioritization_transcripts_force": false,
    "prioritization_transcripts_version_force": false
  }
}
```

10.11.1 profiles

See Prioritization section

Type: str

Default: None

10.11.2 default__profile

See Prioritization section

Type: str

Default: None

10.11.3 prioritization_config

See Prioritization section

Type: Path

Default: None

Examples:

Prioritization configuration JSON file as an option

```
{
  "prioritization_config": "prioritization_config.json"
}
```

10.11.4 prioritization_score_mode

See Prioritization section

Type: str

Choices: ['HOWARD', 'VaRank']

Default: HOWARD

10.11.5 pzprefix

See Prioritization section

10.11.6 pzfields

See Prioritization section

Type: str

Default: PZScore,PZFlag

10.11.7 prioritization_transcripts_order

Defines the order of transcripts to determine which one is chosen (by default PZTFlag and PZTScore). All available annotation can be used (e.g. scores, length of transcript, predictions...). The first transcript will be chosen in case of equal order.

Type: dict

Default: {}

Examples:

Default order of transcript using Flag and Score:

```
{
  "prioritization_transcripts_order": {
    "PZTFlag": "ASC",
    "PZTScore": "DESC"
  }
}
```

Order of transcript using Flag, Score and additional specific spliceAI_score:

```
{
  "prioritization_transcripts_order": {
    "PZTFlag": "ASC",
    "PZTScore": "DESC",
    "spliceAI_score": "DESC"
  }
}
```

10.11.8 prioritization_transcripts

Defines a file with an ordered list of transcripts of preference. The first transcript in this list will be chosen, if no order is define (see above), or if this list is not forced (see below).

Type: Path

Default: None

Examples:

File with transcripts of preference:

```
{
  "prioritization_transcripts": "transcripts.tsv"
}
```

10.11.9 prioritization_transcripts_force

Force to use the list of transcripts of preference define in the provided file (see above).

Type: Boolean

Default: False

Examples:

Force using transcripts of preference in file:

```
{
  "prioritization_transcripts_force": true
}
```

10.11.10 prioritization_transcripts_version_force

By default, versions of transcripts are not considered when comparison is needed (e.g. transcript ID and list of transcript of preference). If true, all transcript ID will be considered with their version.

Type: Boolean

Default: False

Examples:

Force using transcripts version:

```
{
  "prioritization_transcripts_version_force": true
}
```

10.12 export

Options to export transcripts view/table into a file ('output' parameter). All HOWARD format are available (e.g. VCF, Parquet, TSV). For VCF format, all columns will be concatenate into INFO column, otherwise each column will be exported.

Type: Dict

Default: {}

Examples:

Export as compressed TSV:

```
{
  "export": {
    "output": "/tmp/output.tsv.gz"
  }
}
```

Export as VCF:

```
{
  "export": {
    "output": "/tmp/output.vcf"
  }
}
```

Export as Parquet:

```
{
  "export": {
    "output": "/tmp/output.parquet"
  }
}
```

11 threads

Specify the number of threads to use for processing HOWARD. It determines the level of parallelism, either on python scripts, duckdb engine and external tools. It can help speed up the process/tool. Use -1 to use all available CPU/cores. Either non valid value is 1 CPU/core.

Type: int

Default: -1

Examples:

Automatically detect all available CPU/cores

```
{
  "threads": -1
}
```

Define 8 CPU/cores

```
{
  "threads": 8
}
```

12 samples

Samples parameters to defined a list of samples or use options to check samples. Only for export in VCF format. By default, if no samples are listed, all existing samples are checked if they contain well-formed genotype annotations (based on 'FORMAT' VCF column).

Type: dict

Examples:

Export only a list of samples:

```
{
  "samples": {
    "list": ["sample1", "sample2"]
  }
}
```

Do not check existing samples (all VCF columns after FORMAT column):

```
{
  "samples": {
    "check": false
  }
}
```

Default configuration, with all samples are considered (null) and checked (true):


```
{
  "samples": {
    "list": null,
    "check": true
  }
}
```

12.1 list

List of columns that correspond to samples (with well formed genotype, based on 'FORMAT' VCF column). Only for export in VCF format. Only these samples are exported in VCF format file.

Type: dict

Examples:

Export only a list of samples:

```
{
  "list": ["sample1", "sample2"]
}
```

12.2 check

Check if samples (either provided in 'list' parameters, or all existing column after 'FORMAT' column) according to 'FORMAT' VCF column. Only for export in VCF format. By default, samples are checked (beware of format if check is disabled) and removed if they are not well-formed.

Type: dict

Examples:

Do not check existing samples:

```
{
  "check": false
}
```

13 databases

HOWARD Parameters Databases JSON describes configuration JSON file for databases download and convert.