

# HOWARD Parameters Databases

---

HOWARD Parameters JSON file defines parameters for databases' downloads.

## Table of contents

- [HOWARD Parameters Databases](#)
  - [databases](#)
    - [assembly](#)
    - [genomes\\_folder](#)
    - [genome](#)
    - [genomes](#)
      - [download\\_genomes](#)
      - [download\\_genomes\\_provider](#)
      - [download\\_genomes\\_contig\\_regex](#)
    - [snpeff](#)
      - [download\\_snpeff](#)
    - [annovar](#)
      - [download\\_annovar](#)
      - [download\\_annovar\\_files](#)
      - [download\\_annovar\\_url](#)
    - [refseq](#)
      - [download\\_refseq](#)
      - [download\\_refseq\\_url](#)
      - [download\\_refseq\\_prefix](#)
      - [download\\_refseq\\_files](#)
      - [download\\_refseq\\_format\\_file](#)
      - [download\\_refseq\\_include\\_utr5](#)
      - [download\\_refseq\\_include\\_utr3](#)
      - [download\\_refseq\\_include\\_chrM](#)
      - [download\\_refseq\\_include\\_non\\_canonical\\_chr](#)
      - [download\\_refseq\\_include\\_non\\_coding\\_transcripts](#)
      - [download\\_refseq\\_include\\_transcript\\_version](#)
    - [dbnsfp](#)
      - [download\\_dbnsfp](#)
      - [download\\_dbnsfp\\_url](#)
      - [download\\_dbnsfp\\_release](#)
      - [download\\_dbnsfp\\_parquet\\_size](#)
      - [download\\_dbnsfp\\_subdatabases](#)
      - [download\\_dbnsfp\\_parquet](#)
      - [download\\_dbnsfp\\_vcf](#)
      - [download\\_dbnsfp\\_no\\_files\\_all](#)
      - [download\\_dbnsfp\\_add\\_info](#)
      - [download\\_dbnsfp\\_row\\_group\\_size](#)
    - [alphamissense](#)
      - [download\\_alphamissense](#)
      - [download\\_alphamissense\\_url](#)
    - [exomiser](#)
      - [download\\_exomiser](#)
      - [download\\_exomiser\\_application\\_properties](#)
      - [download\\_exomiser\\_url](#)
      - [download\\_exomiser\\_release](#)
      - [download\\_exomiser\\_phenotype\\_release](#)
      - [download\\_exomiser\\_remm\\_release](#)
      - [download\\_exomiser\\_remm\\_url](#)
      - [download\\_exomiser\\_cadd\\_release](#)
      - [download\\_exomiser\\_cadd\\_url](#)

- [download\\_exomiser\\_cadd\\_url\\_snv\\_file](#)
- [download\\_exomiser\\_cadd\\_url\\_indel\\_file](#)
- [dbsnp](#)
  - [download\\_dbsnp](#)
  - [download\\_dbsnp\\_releases](#)
  - [download\\_dbsnp\\_release\\_default](#)
  - [download\\_dbsnp\\_url](#)
  - [download\\_dbsnp\\_url\\_files](#)
  - [download\\_dbsnp\\_url\\_files\\_prefix](#)
  - [download\\_dbsnp\\_assemblies\\_map](#)
  - [download\\_dbsnp\\_vcf](#)
  - [download\\_dbsnp\\_parquet](#)
- [hgmd](#)
  - [convert\\_hgmd](#)
  - [convert\\_hgmd\\_file](#)
  - [convert\\_hgmd\\_basename](#)
- [from\\_Annotvar](#)
  - [input\\_annotvar](#)
  - [output\\_annotvar](#)
  - [annotvar\\_code](#)
  - [annotvar\\_to\\_parquet](#)
  - [annotvar\\_reduce\\_memory](#)
  - [annotvar\\_multi\\_variant](#)
- [from\\_extann](#)
  - [input\\_extann](#)
  - [output\\_extann](#)
  - [refgene](#)
  - [transcripts](#)
  - [param\\_extann](#)
  - [mode\\_extann](#)
- [Parameters](#)
  - [generate\\_param](#)
  - [generate\\_param\\_description](#)
  - [generate\\_param\\_releases](#)
  - [generate\\_param\\_formats](#)
  - [generate\\_param\\_bcftools](#)

Examples:

Example of simple databases parameters JSON file for downloads

```
"databases": {
  "genomes": {
    "download_genomes": "~/howard/databases/genomes/current",
    "download_genomes_contig_regex": "chr[0-9XYM]+$"
  },
  "snpeff": {
    "download_snpeff": "~/howard/databases/snpeff/current"
  },
  "annovar": {
    "download_annovar": "~/howard/databases/annovar/current",
    "download_annovar_files": "refGene,cosmic70,nci60"
  },
  "refseq": {
    "download_refseq": "~/howard/databases/refseq/current"
  },
  "dbnsfp": {
    "download_dbnsfp": "~/howard/databases/dbnsfp/current",
    "download_dbnsfp_release": "4.4a"
  },
  "alphamissense": {
```

```

    "download_alphamissense": "~/howard/databases/alphamissense/current"
  },
  "exomiser": {
    "download_exomiser": "~/howard/databases/exomiser/current"
  },
  "dbSNP": {
    "download_dbSNP": "~/howard/databases/dbSNP/current",
    "download_dbSNP_releases": "b156",
    "download_dbSNP_vcf": true,
    "download_dbSNP_parquet": true
  },
  "assemblies": [
    "hg19"
  ]
}

```

Example of a full Databases parameters JSON file for downloads

```

"databases": {
  "genomes_folder": "~/howard/databases/genomes/current",
  "genome": "~/howard/databases/genomes/current/hg19/hg19.fa",
  "genomes": {
    "download_genomes": "~/howard/databases/genomes/current",
    "download_genomes_provider": "UCSC",
    "download_genomes_contig_regex": "chr[0-9XYM]+$"
  },
  "snpeff": {
    "download_snpeff": "~/howard/databases/snpeff/current"
  },
  "annovar": {
    "download_annovar": "~/howard/databases/annovar/current",
    "download_annovar_files": "refGene,cosmic70,nci60",
    "download_annovar_url": "http://www.openbioinformatics.org/annovar/download"
  },
  "refseq": {
    "download_refseq": "~/howard/databases/refseq/current",
    "download_refseq_url": "http://hgdownload.soe.ucsc.edu/goldenPath",
    "download_refseq_prefix": "ncbiRefSeq",
    "download_refseq_files": "ncbiRefSeq.txt,ncbiRefSeqLink.txt",
    "download_refseq_format_file": "ncbiRefSeq.txt",
    "download_refseq_include utr5": false,
    "download_refseq_include utr3": false,
    "download_refseq_include chrM": false,
    "download_refseq_include_non_canonical_chr": false,
    "download_refseq_include_non_coding_transcripts": false,
    "download_refseq_include_transcript_version": false
  },
  "dbnsfp": {
    "download_dbnsfp": "~/howard/databases/dbnsfp/current",
    "download_dbnsfp_url": "https://dbnsfp.s3.amazonaws.com",
    "download_dbnsfp_release": "4.4a",
    "download_dbnsfp_parquet_size": 100,
    "download_dbnsfp_subdatabases": true,
    "download_dbnsfp_parquet": false,
    "download_dbnsfp_vcf": false,
    "download_dbnsfp_no_files_all": false,
    "download_dbnsfp_add_info": false,
    "download_dbnsfp_row_group_size": 100000
  },
  "alphamissense": {
    "download_alphamissense": "~/howard/databases/alphamissense/current",
    "download_alphamissense_url": "https://storage.googleapis.com/dm_alphamissense"
  },
  "exomiser": {
    "download_exomiser": "~/howard/databases/exomiser/current",

```

```

    "download_exomiser_url": "http://data.monarchinitiative.org/exomiser",
    "download_exomiser_remm_url": "https://kircherlab.bihealth.org/download/ReMM",
    "download_exomiser_cadd_url": "https://kircherlab.bihealth.org/download/CADD",
    "download_exomiser_cadd_url_snv_file": "whole_genome_SNVs.tsv.gz",
    "download_exomiser_cadd_url_indel_file": "InDels.tsv.gz"
  },
  "dbSNP": {
    "download_dbSNP": "~/howard/databases/dbSNP/current",
    "download_dbSNP_releases": "b156",
    "download_dbSNP_url": "https://ftp.ncbi.nih.gov/snp/archive",
    "download_dbSNP_url_files_prefix": "GCF_000001405",
    "download_dbSNP_assemblies_map": {
      "hg19": "25",
      "hg38": "40"
    },
    "download_dbSNP_vcf": true,
    "download_dbSNP_parquet": false
  },
  "parameters": {
    "generate_param_releases": "current",
    "generate_param_formats": "parquet",
    "generate_param_bcftools": false
  },
  "assemblies": [
    "hg19"
  ]
}

```

## databases

Databases download options

### databases::assembly

Genome Assembly (e.g. 'hg19', 'hg38').

Type: **str**

Default: **hg19**

Examples:

Default assembly for all analysis tools

```
"assembly": "hg19"
```

List of assemblies for databases download tool

```
"assembly": "hg19,hg38"
```

### databases::genomes\_folder

Folder containing genomes. (e.g. '/Users/lebechea/howard/databases/genomes/current')

Type: **Path**

Default: **/Users/lebechea/howard/databases/genomes/current**

### databases::genome

Genome file in fasta format (e.g. 'hg19.fa', 'hg38.fa').

Type: **Path**

Default: **~/howard/databases/genomes/current/hg19/hg19.fa**

databases::genomes

Genomes download.

#### **databases::genomes::download\_genomes**

Path to genomes folder with Fasta files, indexes, and all files generated by pygenome module. (e.g. '/Users/lebechea/howard/databases/genomes/current').

Type: **Path**

Default: **None**

#### **databases::genomes::download\_genomes\_provider**

Download Genome from an external provider. Available: GENCODE, Ensembl, UCSC, NCBI.

Type: **str**

Choices: **['GENCODE', 'Ensembl', 'UCSC', 'NCBI']**

Default: **UCSC**

#### **databases::genomes::download\_genomes\_contig\_regex**

Regular expression to select specific chromosome (e.g 'chr[0-9XYM]+\$').

Type: **str**

Default: **None**

databases::snpeff

snpEff download.

#### **databases::snpeff::download\_snpeff**

Download snpEff databases within snpEff folder

Type: **Path**

Default: **None**

databases::annovar

Annovar download.

#### **databases::annovar::download\_annovar**

Path to Annovar databases (e.g. '/Users/lebechea/howard/databases/annovar/current').

Type: **Path**

Default: **None**

#### **databases::annovar::download\_annovar\_files**

Download Annovar databases for a list of Annovar file code (see Annovar Doc). Use None to download all available files, or Annovar keyword (e.g. 'refGene', 'cosmic70', 'clinvar\_202\*'). Note that refGene will at least be downloaded, and only files that not already exist or changed will be downloaded.

Type: **str**

Default: **None**

#### **databases::annovar::download\_annovar\_url**

Annovar databases URL (see Annovar Doc).

Type: **str**

Default: **<http://www.openbioinformatics.org/annovar/download>**

#### **databases::refseq**

refSeq download.

#### **databases::refseq::download\_refseq**

Path to refSeq databases (e.g. '/Users/lebechea/howard/databases/refseq/current').

Type: **Path**

Default: **None**

#### **databases::refseq::download\_refseq\_url**

refSeq databases URL (see refSeq WebSite) (e.g. 'http://hgdownload.soe.ucsc.edu/goldenPath')·/n

Type: **str**

Default: **<http://hgdownload.soe.ucsc.edu/goldenPath>**

#### **databases::refseq::download\_refseq\_prefix**

Check existing refSeq files in refSeq folder.

Type: **str**

Default: **[ncbiRefSeq](#)**

#### **databases::refseq::download\_refseq\_files**

List of refSeq files to download.

Type: **str**

Default: **[ncbiRefSeq.txt](#),[ncbiRefSeqLink.txt](#)**

#### **databases::refseq::download\_refseq\_format\_file**

Name of refSeq file to convert in BED format (e.g. 'ncbiRefSeq.txt'). Process only if not None.

Type: **str**

Default: **None**

#### **databases::refseq::download\_refseq\_include utr5**

Formating BED refSeq file including 5'UTR.

Default: **False**

#### **databases::refseq::download\_refseq\_include utr3**

Formating BED refSeq file including 3'UTR.

Default: **False**

**databases::refseq::download\_refseq\_include\_chrM**

Formating BED refSeq file including Mitochondiral chromosome 'chrM' or 'chrMT'.

Default: **False**

**databases::refseq::download\_refseq\_include\_non\_canonical\_chr**

Formating BED refSeq file including non canonical chromosomes.

Default: **False**

**databases::refseq::download\_refseq\_include\_non\_coding\_transcripts**

Formating BED refSeq file including non coding transcripts.

Default: **False**

**databases::refseq::download\_refseq\_include\_transcript\_version**

Formating BED refSeq file including transcript version.

Default: **False**

**databases::dbnsfp**

dbNSFP download.

**databases::dbnsfp::download\_dbnsfp**

Download dbNSFP databases within dbNSFP folder (e.g. '/Users/lebechea/howard/databases').

Type: **Path**

Default: **None**

**databases::dbnsfp::download\_dbnsfp\_url**

Download dbNSFP databases URL (see dbNSFP website) (e.g. <https://dbnsfp.s3.amazonaws.com>).

Type: **str**

Default: **<https://dbnsfp.s3.amazonaws.com>**

**databases::dbnsfp::download\_dbnsfp\_release**

Release of dbNSFP to download (see dbNSFP website) (e.g. '4.4a').

Default: **4.4a**

**databases::dbnsfp::download\_dbnsfp\_parquet\_size**

Maximum size (Mb) of data files in Parquet folder. Parquet folder are partitioned (hive) by chromosome (sub-folder), which contain N data files.

Type: **int**

Default: **100**

**databases::dbnsfp::download\_dbnsfp\_subdatabases**

Generate dbNSFP sub-databases. dbNSFP provides multiple databases which are split onto multiple columns. This option create a Parquet folder for each sub-database (based on columns names).

Default: **False**

**databases::dbnsfp::download\_dbnsfp\_parquet**

Generate a Parquet file for each Parquet folder.

Default: **False**

**databases::dbnsfp::download\_dbnsfp\_vcf**

Generate a VCF file for each Parquet folder. Need genome FASTA file (see --download-genome).

Default: **False**

**databases::dbnsfp::download\_dbnsfp\_no\_files\_all**

Not generate database Parquet/VCF file for the entire database ('ALL'). Only sub-databases files will be generated. (see '--download-dbnsfp-subdatabases').

Default: **False**

**databases::dbnsfp::download\_dbnsfp\_add\_info**

Add INFO column (VCF format) in Parquet folder and file. Useful for speed up full annotation (all available columns). Increase memory and space during generation of files.

Default: **False**

**databases::dbnsfp::download\_dbnsfp\_row\_group\_size**

Minimum number of rows in a parquet row group (see duckDB doc). Lower can reduce memory usage and slightly increase space during generation, speed up highly selective queries, slow down whole file queries (e.g. aggregations).

Type: **int**

Default: **100000**

**databases::alphamissense**

AlphaMissense download.

**databases::alphamissense::download\_alphamissense**

Path to AlphaMissense databases

Type: **Path**

Default: **None**

**databases::alphamissense::download\_alphamissense\_url**

Download AlphaMissense databases URL (see AlphaMissense website) (e.g. 'https://storage.googleapis.com/dm\_alphamissense').

Type: **str**

Default: **https://storage.googleapis.com/dm\_alphamissense**

**databases::exomiser**

Exomiser download.

**databases::exomiser::download\_exomiser**

Path to Exomiser databases (e.g. /Users/lebechea/howard/databases/exomiser/current).

Type: **Path**

Default: **None**



**databases::exomiser::download\_exomiser\_application\_properties**

Exomiser Application Properties configuration file (see Exomiser website). This file contains configuration settings for the Exomiser tool. If this parameter is not provided, the function will attempt to locate the application properties file automatically based on the Exomiser. Configuration information will be used to download expected releases (if no other parameters). CADD and REMM will be downloaded only if 'path' are provided.

Type: **Path**

Default: **None**

**databases::exomiser::download\_exomiser\_url**

URL where Exomiser database files can be downloaded from (e.g. 'http://data.monarchinitiative.org/exomiser').

Type: **str**

Default: **<http://data.monarchinitiative.org/exomiser>**

**databases::exomiser::download\_exomiser\_release**

Release of Exomiser data to download. If "default", "auto", or "config", retrieve from Application Properties file. If not provided (None), from Application Properties file (Exomiser data-version) or default '2109'.

Type: **str**

Default: **None**

**databases::exomiser::download\_exomiser\_phenotype\_release**

Release of Exomiser phenotype to download. If not provided (None), from Application Properties file (Exomiser Phenotype data-version) or Exomiser release.

Type: **str**

Default: **None**

**databases::exomiser::download\_exomiser\_remm\_release**

Release of ReMM (Regulatory Mendelian Mutation) database to download. If "default", "auto", or "config", retrieve from Application Properties file.

Type: **str**

Default: **None**

**databases::exomiser::download\_exomiser\_remm\_url**

URL where ReMM (Regulatory Mendelian Mutation) database files can be downloaded from (e.g. 'https://kircherlab.bihealth.org/download/ReMM').

Type: **str**

Default: **<https://kircherlab.bihealth.org/download/ReMM>**

**databases::exomiser::download\_exomiser\_cadd\_release**

Release of CADD (Combined Annotation Dependent Depletion) database to download. If "default", "auto", or "config", retrieve from Application Properties file.

Type: **str**

Default: **None**

**databases::exomiser::download\_exomiser\_cadd\_url**

URL where CADD (Combined Annotation Dependent Depletion) database files can be downloaded from (e.g. 'https://kircherlab.bihealth.org/download/CADD').

Type: `str`

Default: `https://kircherlab.bihealth.org/download/CADD`

#### **databases::exomiser::download\_exomiser\_cadd\_url\_snv\_file**

Name of the file containing the SNV (Single Nucleotide Variant) data for the CADD (Combined Annotation Dependent Depletion) database.

Type: `str`

Default: `whole_genome_SNVs.tsv.gz`

#### **databases::exomiser::download\_exomiser\_cadd\_url\_indel\_file**

Name of the file containing the INDEL (Insertion-Deletion) data for the CADD (Combined Annotation Dependent Depletion) database.

Type: `str`

Default: `InDels.tsv.gz`

#### **databases::dbsnp**

dbSNP download.

#### **databases::dbsnp::download\_dbsnp**

Path to dbSNP databases (e.g. '/Users/lebechea/howard/databases/exomiser/dbsnp').

Type: `Path`

Default: `None`

#### **databases::dbsnp::download\_dbsnp\_releases**

Release of dbSNP to download (e.g. 'b152', 'b152,b156').

Type: `str`

Default: `b156`

#### **databases::dbsnp::download\_dbsnp\_release\_default**

Default Release of dbSNP ('default' symlink) (e.g. 'b156'). If None, first release to download will be assigned as default only if it does not exists.

Type: `str`

Default: `None`

#### **databases::dbsnp::download\_dbsnp\_url**

URL where dbSNP database files can be downloaded from. (e.g. 'https://ftp.ncbi.nih.gov/snp/archive').

Type: `str`

Default: `https://ftp.ncbi.nih.gov/snp/archive`

#### **databases::dbsnp::download\_dbsnp\_url\_files**

Dictionary that maps assembly names to specific dbSNP URL files. It allows you to provide custom dbSNP URL files for specific assemblies instead of using the default file naming convention.

Type: **str**

Default: **None**

#### **databases::dbsnp::download\_dbsnp\_url\_files\_prefix**

String that represents the prefix of the dbSNP file name for a specific assembly. It is used to construct the full URL of the dbSNP file to be downloaded.

Type: **str**

Default: **GCF\_000001405**

#### **databases::dbsnp::download\_dbsnp\_assemblies\_map**

dictionary that maps assembly names to their corresponding dbSNP versions. It is used to construct the dbSNP file name based on the assembly name.

Type: **str**

Default: **{'hg19': '25', 'hg38': '40'}**

#### **databases::dbsnp::download\_dbsnp\_vcf**

Generate well-formatted VCF from downloaded file:

- Add and filter contigs associated to assembly
- Normalize by splitting multiallelics
- Need genome (see --download-genome)

Default: **False**

#### **databases::dbsnp::download\_dbsnp\_parquet**

Generate Parquet file from VCF.

Default: **False**

#### **databases::hgmd**

HGMD convert.

#### **databases::hgmd::convert\_hgmd**

Convert HGMD databases. Folder where the HGMD databases will be stored. Fields in VCF, Parquet and TSV will be generated. If the folder does not exist, it will be created.

Type: **Path**

Default: **None**

#### **databases::hgmd::convert\_hgmd\_file**

File from HGMD. Name format 'HGMD\_Pro\_\_\_.vcf.gz'.

Type: **Path**

Default: **None**

#### **databases::hgmd::convert\_hgmd\_basename**

File output basename. Generated files will be prefixed by basename (e.g. 'HGMD\_Pro\_MY\_RELEASE') By default (None), input file name without '.vcf.gz'.

Type: `str`

Default: `None`

`databases::from_Annovar`

Annovar convert.

**`databases::from_Annovar::input_annovar`**

Input Annovar file path. Format file must be a Annovar TXT file, associated with '.idx'.

Type: `Path`

Default: `None`

**`databases::from_Annovar::output_annovar`**

Output Annovar file path. Format file must be either VCF compressed file '.vcf.gz'.

Type: `Path`

Default: `None`

**`databases::from_Annovar::annovar_code`**

Annovar code, or database name. Usefull to name databases columns.

Type: `str`

Default: `None`

**`databases::from_Annovar::annovar_to_parquet`**

Parquet file conversion.

Type: `Path`

Default: `None`

**`databases::from_Annovar::annovar_reduce_memory`**

Reduce memory option for Annovar convert, either 'auto' (auto-detection), 'enable' or 'disable'.

Type: `str`

Choices: [`'auto'`, `'enable'`, `'disable'`]

Default: `auto`

**`databases::from_Annovar::annovar_multi_variant`**

Variant with multiple annotation lines on Annovar file. Either 'auto' (auto-detection), 'enable' or 'disable'.

Type: `str`

Choices: [`'auto'`, `'enable'`, `'disable'`]

Default: `auto`

`databases::from_extann`

Extann convert (gene annotation).

**`databases::from_extann::input_extann`**

Input Extann file path. Format file must be a Extann TXT file or TSV file. File need to have at least the genes column.

Type: **Path**

Default: **None**

**databases::from\_extann::output\_extann**

Output Extann file path. Output extann file, should be BED or BED.gz.

Type: **Path**

Default: **None**

**databases::from\_extann::refgene**

Path to refGene annotation file.

Type: **Path**

Default: **None**

**databases::from\_extann::transcripts**

Transcripts TSV file, with Transcript in first column, optional Gene in second column.

Type: **Path**

Default: **None**

**databases::from\_extann::param\_extann**

Param extann file path. Param containing configuration, options to replace chars and bedlike header description, conf vcf specs.

Type: **Path**

Default: **~/howard/config/param.extann.json**

**databases::from\_extann::mode\_extann**

Mode extann selection. How to pick transcript from ncbi, keep all, keep the longest, or keep the chosen one (transcript\_extann).

Type: **str**

Choices: **['all', 'longest', 'chosen']**

Default: **longest**

**databases::Parameters**

Parameters generation.

**databases::Parameters::generate\_param**

Parameter file (JSON) with all databases found. Databases folders scanned are defined in config file. Structure of databases follow this structure (see doc): **.../[\*].[parquet|vcf.gz|...]**

Type: **Path**

Default: **None**

**databases::Parameters::generate\_param\_description**

Description file (JSON) with all databases found. Contains all databases with description of format, assembly, fields...

Type: **Path**

Default: **None**

**databases::Parameters::generate\_param\_releases**

List of database folder releases to check (e.g. 'current', 'latest').

Type: `str`

Default: `current`

**databases::Parameters::generate\_param\_formats**

List of database formats to check (e.g. 'parquet', 'parquet,vcf,bed,tsv').

Type: `str`

Default: `parquet`

**databases::Parameters::generate\_param\_bcftools**

Generate parameter JSON file with BCFTools annotation for allowed formats (i.e. 'vcf', 'bed').

Default: `False`