

---

# VAPr Documentation

*Release 3.0*

**Amanda Birmingham, Adam Mark, Carlo Mazzaferro, Guorong Xu,**

**Aug 20, 2018**



---

## Contents:

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Installation . . . . .	1
1.1.1	Ancillary Libraries . . . . .	1
1.1.2	MongoDB . . . . .	1
1.1.3	BCFtools . . . . .	1
1.1.4	Tabix . . . . .	2
1.1.5	ANNOVAR . . . . .	2
1.1.6	VAPr . . . . .	2
1.2	Annotation Quickstart using ANNOVAR . . . . .	2
1.2.1	Downloading the ANNOVAR databases . . . . .	3
<b>2</b>	<b>Downstream Analysis</b>	<b>5</b>
2.1	Filtering Variants . . . . .	5
2.1.1	1. Rare Deleterious Variants . . . . .	5
2.1.2	2. Known Disease Variants . . . . .	5
2.1.3	3. Deleterious Compound Heterozygous Variants . . . . .	5
2.1.4	4. De novo Variants . . . . .	6
2.1.5	Create your own filter . . . . .	6
2.2	Output Files . . . . .	6
2.2.1	Unfiltered Variants CSV . . . . .	6
2.2.2	Filtered Variants CSV . . . . .	7
2.2.3	Unfiltered Variants VCF . . . . .	7
2.2.4	Filtered Variants VCF . . . . .	7
<b>3</b>	<b>Core Methods</b>	<b>9</b>
<b>4</b>	<b>Tutorial</b>	<b>11</b>
<b>5</b>	<b>VAPr package</b>	<b>13</b>
5.1	Submodules . . . . .	13
5.2	VAPr.annovar_output_parsing module . . . . .	13
5.3	VAPr.annovar_running module . . . . .	14
5.4	VAPr.chunk_processing module . . . . .	14
5.5	VAPr.filtering module . . . . .	15
5.6	VAPr.validation module . . . . .	15
5.7	VAPr.vapr_core module . . . . .	16
5.8	VAPr.vcf_genotype_fields_parsing module . . . . .	20

5.9	VAPr.vcf_merging module . . . . .	22
5.10	Module contents . . . . .	22
<b>6</b>	<b>Indices and tables</b>	<b>23</b>
	<b>Python Module Index</b>	<b>25</b>

# CHAPTER 1

---

## Introduction

---

This package is aimed at providing a way of retrieving variant information using ANNOVAR and myvariant.info. In particular, it is suited for bioinformaticians interested in aggregating variant information into a single NoSQL database (MongoDB solely at the moment).

### 1.1 Installation

#### 1.1.1 Ancillary Libraries

VAPr relies on a variety of packages to function correctly. Below are packages and dependencies required to ensure that VAPr works correctly.

---

**Note:** Jupyter, Pandas, and other ancillary libraries are not installed with VAPr and must be installed separately. These can be conveniently install using [Anaconda](#):

---

```
$ conda install python=3 pandas mongodb pymongo jupyter notebook
```

#### 1.1.2 MongoDB

VAPr is written in Python and stores variant annotations in NoSQL database, using a locally-installed instance of MongoDB. [Installation instructions](#)

#### 1.1.3 BCFtools

**BCFtools** will be used for VCF file merging between samples. To download and install:

```
$ wget https://github.com/samtools/bcftools/releases/download/1.6/bcftools-1.6.tar.bz2
$ tar -vxjf bcftools-1.6.tar.bz2
$ cd bcftools-1.6
$ make
$ make install
$ export PATH=/where/to/install/bin:$PATH
```

Refer [here](#) for installation debugging.

### 1.1.4 Tabix

**Tabix** and **bzzip** binaries are available through the HTSLib project:

```
$ wget https://github.com/samtools/htslib/releases/download/1.6/htslib-1.6.tar.bz2
$ tar -vxjf htslib-1.6.tar.bz2
$ cd htslib-1.6
$ make
$ make install
$ export PATH=/where/to/install/bin:$PATH
```

Refer [here](#) for installation debugging.

### 1.1.5 ANNOVAR

(It is possible to proceed without installing ANNOVAR. Variants will only be annotated with MyVariant.info. In that case, users can skip the next steps and go straight to the section Known Variant Annotation and Storage)

Users who wish to annotate novel variants will also need to have a local installation of the popular command-line software ANNOVAR(1), which VAPr wraps with a Python interface. If you use ANNOVAR's functionality through VAPr, please remember to cite the ANNOVAR publication (see #1 in Citations)!

The base ANNOVAR program must be installed by each user individually, since its license agreement does not permit redistribution. Please visit the ANNOVAR download form [here](#), ensure that you meet the requirements for a free license, and fill out the required form. You will then receive an email providing a link to the latest ANNOVAR release file. Download this file (which will usually have a name like annoivar.latest.tar.gz) and place it in the location on your machine in which you would like the ANNOVAR program and its data to be installed—the entire disk size of the databases will be around 25 GB, so make sure you have such space available!

### 1.1.6 VAPr

VAPr is compatible with Python 2.7 or later, but it is preferred to use Python 3.5 or later to take full advantage of all functionality. The simplest way to install VAPr is from [PyPI](#) with [pip](#), Python's preferred package installer.

```
$ pip install VAPr
```

## 1.2 Annotation Quickstart using ANNOVAR

An annotation project can be started by providing the API with a small set of information and then running the core methods provided to spawn annotation jobs. This is done in the following manner:

```
# Import core module
from VAPr import vapr_core
import os

# Start by specifying the project information
IN_PATH = "/path/to/vcf"
OUT_PATH = "/path/to/out"
ANNOVAR_PATH = "/path/to/annovar"
MONGODB = 'VariantDatabase'
COLLECTION = 'Cancer'

annotator = vapr_core.VaprAnnotator(input_dir=IN_PATH,
                                      output_dir=OUT_PATH,
                                      mongo_db_name=MONGODB,
                                      mongo_collection_name=COLLECTION,
                                      build_ver='hg19',
                                      vcfs_gzipped=False,
                                      annovar_install_path=ANNOVAR_PATH)

annotator.download_databases()
dataset = annotator.annotate(num_processes=8)
```

### 1.2.1 Downloading the ANNOVAR databases

If you plan to use Annovar, the command below will download the necessary Annovar databases. The code above includes this step. When Annovar is first installed, it does not install databases by default. The `vapr_core` has a method `download_annovar_databases()` that will download the necessary annovar databases. If you do not plan on using Annovar, you should not run this command. Note: this command only needs to be run once, the first time you use VAPr.

```
annotator.download_databases()
```

This will download the required databases from ANNOVAR for annotation and will kickstart the annotation process, storing the variants in MongoDB.



# CHAPTER 2

---

## Downstream Analysis

---

For notes on how to implement these features, refer to the [Tutorial](#) and the [API Reference](#)

### 2.1 Filtering Variants

Four different pre-made filters that allow for the retrieval of specific variants have been implemented. These allow the user to query in an easy and efficient manner variants of interest

#### 2.1.1 1. Rare Deleterious Variants

- criteria 1: 1000 Genomes (ALL) allele frequency (Annovar) < 0.05 or info not available
- criteria 2: ESP6500 allele frequency (MyVariant.info - CADD) < 0.05 or info not available
- criteria 3: cosmic70 (MyVariant.info) information is present
- criteria 4: Func\_knownGene (Annovar) is exonic, splicing, or both
- criteria 5: ExonicFunc\_knownGene (Annovar) is not “synonymous SNV”

#### 2.1.2 2. Known Disease Variants

- criteria: cosmic70 (MyVariant.info) information is present or ClinVar data is present and clinical significance is not Benign or Likely Benign

#### 2.1.3 3. Deleterious Compound Heterozygous Variants

- criteria 1: genotype\_subclass\_by\_class (VAPr) is compound heterozygous
- criteria 2: CADD phred score (MyVariant.info - CADD) > 10

## 2.1.4 4. De novo Variants

- criteria 1: Variant present in proband
- criteria 2: Variant not present in either ancestor-1 or ancestor-2

## 2.1.5 Create your own filter

As long as you have a MongoDB instance running and an annotation job ran successfully, filtering can be performed through `pymongo` as shown by the code below. Running the query will return a `cursor` object, which can be iterated upon.

If instead a list is intended to be created from it, simply add: `filter2 = list(filter2)`.

**Warning:** If the number of variants in the database is large and the filtering is not set up correctly, returning a list will be probably crash your computer since lists are kept in memory. Iterating over the cursor object perform *lazy evaluations* (i.e., one item is returned at a time instead of in bulk) which are much more memory efficient.

Further, if you'd like to customize your filters, a good idea would be to look at the available fields to be filtered. Looking at the myvariant.info documentation, you can see what are all the fields available and can be used for filtering.

```
from pymongo import MongoClient

client = MongoClient()
db = getattr(client, mongodb_name)
collection = getattr(db, mongo_collection_name)

filtered = collection.find({"$and": [
    {"$or": [{"func_knowngene": "exonic"}, {"func_knowngene": "splicing"}]}, {"cosmic70": {"$exists": True}}, {"1000g2015aug_all": {"$lt": 0.05}}
]})

# filtered = list(filtered) Uncomment this if you'd like to return them as a list
for var in filtered:
    print(var)
```

## 2.2 Output Files

Although iterating over variants can be interesting for cursory analyses, we provide functionality to retrieve as well csv files for downstream analysis. A few options are available:

### 2.2.1 Unfiltered Variants CSV

```
write_unfiltered_annotated_csv(out_file_path)
```

- All variants will be written to a CSV file.

## 2.2.2 Filtered Variants CSV

```
write_filtered_annotated_csv(variant_list, out_file_path)
```

- A list of filtered variants will be written to a CSV file.

## 2.2.3 Unfiltered Variants VCF

```
write_unfiltered_annotated_vcf(vcf_out_path)
```

- All variants will be written to a VCF file.

## 2.2.4 Filtered Variants VCF

```
write_filtered_annotated_vcf(variant_list, vcf_out_path)
```

- A List of filtered variants will be written to a VCF file.



# CHAPTER 3

---

## Core Methods

---

See the API Reference for [\*VAPr:vapr\\_core module\*](#) for a detailed functionality of the core methods and classes of this package.



# CHAPTER 4

---

## Tutorial

---

A brief, although comprehensive tour of the functionality offered by this package can be found in this [Jupyter Notebook](#). To run it interactively, download the github repo (or just the Notebook), install the required dependencies (see [Installation](#))



# CHAPTER 5

---

VAPr package

---

## 5.1 Submodules

### 5.2 VAPr.annovar\_output\_parsing module

```
class VAPr.annovar_output_parsing.AnnovarAnnotatedVariant
Bases: object

    ALLELE_DEPTH_KEY = 'AD'
    FILTER_PASSING_READS_COUNT_KEY = 'filter_passing_reads_count'
    GENOTYPE_KEY = 'genotype'
    GENOTYPE_LIKELIHOODS_KEY = 'genotype_likelihoods'
    GENOTYPE_SUBCLASS_BY_CLASS_KEY = 'genotype_subclass_by_class'
    HGVS_ID_KEY = 'hgvs_id'
    SAMPLES_KEY = 'samples'
    SAMPLE_ID_KEY = 'sample_id'

    @classmethod make_per_variant_annotation_dict(fields_by_annovar_header,
                                                   hgvs_id, format_string, genotype_field_strings_by_sample_name)

class VAPr.annovar_output_parsing.AnnovarTxtParser
Bases: object

    Class that processes an Annovar-created tab-delimited text file.

    ALT_HEADER = 'alt'
    CHR_HEADER = 'chr'
    CYTOBAND_HEADER = 'cytoband'
```

```
END_HEADER = 'end'
ESP6500_ALL_HEADER = 'esp6500siv2_all'
EXONICFUNC_KNOWNGENE_HEADER = 'exonicfunc_knowngene'
FUNC_KNOWNGENE_HEADER = 'func_knowngene'
GENEDETAIL_KNOWNGENE_HEADER = 'genedetail_knowngene'
GENE_KNOWNGENE_HEADER = 'gene_knowngene'
GENOMIC_SUPERDUPS_HEADER = 'genomicsuperdups'
NCI60_HEADER = 'nci60'
OTHERINFO_HEADER = 'otherinfo'
RAW_CHR_MT_SUFFIX_VAL = 'M'
RAW_CHR_MT_VAL = 'chrM'
REF_HEADER = 'ref'
SCORE_KEY = 'Score'
STANDARDIZED_CHR_MT_SUFFIX_VAL = 'MT'
STANDARDIZED_CHR_MT_VAL = 'chrMT'
START_HEADER = 'start'
TFBS_CONS_SITES_HEADER = 'tfbsconsites'
THOU_G_2015_ALL_HEADER = '1000g2015aug_all'
classmethod read_chunk_of_annotations_to_dicts_list(annovar_txt_file_like_obj,
                                                    sample_names_list,
                                                    chunk_index, chunk_size)
```

## 5.3 VAPr.annovar\_running module

```
class VAPr.annovar_running.AnnovarWrapper(annovar_install_path, genome_build_version,
                                            custom_annovar_dbs_to_use=None)
Bases: object
Wrapper around ANNOVAR download and annotation functions
download_databases()
hg_19_databases = {'1000g2015aug': 'f', 'knownGene': 'g'}
hg_38_databases = {'1000g2015aug': 'f', 'knownGene': 'g'}
run_annotation(single_vcf_path, output_basename, output_dir)
```

## 5.4 VAPr.chunk\_processing module

```
class VAPr.chunk_processing.AnnotationJobParamsIndices

CHUNK_INDEX_INDEX = 0
CHUNK_SIZE_INDEX = 2
```

```

COLLECTION_NAME_INDEX = 4
DB_NAME_INDEX = 3
FILE_PATH_INDEX = 1
GENOME_BUILD_VERSION_INDEX = 5
SAMPLE_LIST_INDEX = 7
VERBOSE_LEVEL_INDEX = 6
classmethod get_num_possible_indices()

VAPr.chunk_processing.collect_chunk_annotations_and_store(job_params_tuple)

```

## 5.5 VAPr.filtering module

VAPr.filtering.get\_any\_of\_sample\_ids\_filter(sample\_names\_list)

VAPr.filtering.get\_sample\_id\_filter(sample\_name)

VAPr.filtering.make\_de\_novo\_variants\_filter(proband, ancestor1, ancestor2)

Function for de novo variant analysis. Can be performed on multisample files or on data coming from a collection of files. In the former case, every sample contains the same variants, although they have differences in their allele frequency and read values. A de novo variant is defined as a variant that occurs only in the specified sample (sample1) and not on the other two (sample2, sample3). Occurrence is defined as having allele frequencies greater than [0, 0] ([REF, ALT]).

VAPr.filtering.make\_deleterious\_compound\_heterozygous\_variants\_filter(sample\_ids\_list=None)

VAPr.filtering.make\_known\_disease\_variants\_filter(sample\_ids\_list=None)

Function for retrieving known disease variants by presence in Clinvar and Cosmic.

VAPr.filtering.make\_rare\_deleterious\_variants\_filter(sample\_ids\_list=None)

Function for retrieving rare, deleterious variants

## 5.6 VAPr.validation module

This module exposes utility functions to validate user inputs

By convention, validation functions in this module raise an appropriate Error if validation is unsuccessful. If it is successful, they return either nothing or the appropriately converted input value.

VAPr.validation.convert\_to\_nonneg\_int(input\_val, nullable=False)

For non-null input\_val, cast to a non-negative integer and return result; for null input\_val, return None.

### Parameters

- **input\_val** (*Any*) – The value to attempt to convert to either a non-negative integer or a None (if nullable). The recognized null values are ‘.’, None, ‘’, and ‘NULL’
- **nullable** (*Optional[bool]*) – True if the input value may be null, false otherwise. Defaults to False.

**Returns** None if nullable=True and the input is a null value. The appropriately cast non-negative integer if input is not null and the cast is successful.

**Raises** `ValueError` – if the input cannot be successfully converted to a non-negative integer or, if allowed, None

VAPr.validation.convert\_to\_nullable(*input\_val*, *cast\_function*)

For non-null *input\_val*, apply *cast\_function* and return result if successful; for null *input\_val*, return None.

#### Parameters

- **input\_val** (*Any*) – The value to attempt to convert to either a None or the type specified by *cast\_function*. The recognized null values are ‘.’, None, ‘’, and ‘NULL’
- **cast\_function** (*Callable[[Any], Any]*) – A function to cast the *input\_val* to some specified type; should raise an error if this cast fails.

**Returns** None if input is the null value. An appropriately cast value if input is not null and the cast is successful.

**Raises** `Error` – whatever error is provided by *cast\_function* if the cast fails.

## 5.7 VAPr.vapr\_core module

```
class VAPr.vapr_core.VaprAnnotator(input_dir,          output_dir,          mongo_db_name,
                                      mongo_collection_name, annovar_install_path=None,
                                      design_file=None, build_ver=None, vcfs_gzipped=False)
```

Bases: `object`

Class in charge of gathering requirements, finding files, downloading databases required to run the annotation

#### Parameters

- **input\_dir** (*str*) – Input directory to vcf files
- **output\_dir** (*str*) – Output directory to annotated vcf files
- **mongo\_db\_name** (*str*) – Name of the database to which you’ll store the collection of variants
- **mongo\_collection\_name** (*str*) – Name of the collection to which you’d store the annotated variants
- **annoivar\_install\_path** (*str*) – Path to locally installed annovar scripts
- **design\_file** (*str*) – path to csv design file
- **build\_ver** (*str*) – genome build version to which annotation will be done against. Either `hg19` or `hg38`
- **vcfs\_gzipped** (*bool*) – if the vcf files are gzipped, set to True

Returns:

```
DEFAULT_GENOME_VERSION = 'hg19'
HG19_VERSION = 'hg19'
HG38_VERSION = 'hg38'
SAMPLE_NAMES_KEY = 'Sample_Names'
SUPPORTED_GENOME_BUILD VERSIONS = ['hg19', 'hg38']
```

**annotate** (*num\_processes=4*, *chunk\_size=2000*, *verbose\_level=1*, *allow\_adds=False*)

This is the main function of the package. It will run Annovar beforehand, and will kick-start the full annotation functionality. Namely, it will collect all the variant data from Annovar annotations, combine it with data coming from MyVariant.info, and parse it to MongoDB, in the database and collection specified in *project\_data*.

It will return the class `VaprDataset`, which can then be used for downstream filtering and analysis.

#### Parameters

- `num_processes` (`int`, *optional*) – number of parallel processes. Defaults to 8
- `chunk_size` (`int`, *optional*) – int number of variants to be processed at once. Defaults to 2000
- `verbose_level` (`int`, *optional*) – int higher verbosity will give more feedback, raise to 2 or 3 when debugging. Defaults to 1
- `allow_adds` (`bool`, *optional*) – bool Allow adding new variants to a pre-existing Mongo collection, or overwrite it (Default value = False)

**Returns** class:~`VAPr.vapr_core.VaprDataset`

**Return type** class

`annotate_lite` (`num_processes=8, chunk_size=2000, verbose_level=1, allow_adds=False`)

‘Lite’ Annotation: it will query `myvariant.info` only, without generating annotations from Annovar. It requires solely VAPr to be installed. The execution will grab the HGVS ids from the vcf files and query the variant data from MyVariant.info.

and inability to run native VAPr queries on the data.

It will return the class `VaprDataset`, which can then be used for downstream filtering and analysis.

#### Parameters

- `num_processes` (`int`, *optional*) – number of parallel processes. Defaults to 8
- `chunk_size` (`int`, *optional*) – int number of variants to be processed at once. Defaults to 2000
- `verbose_level` (`int`, *optional*) – int higher verbosity will give more feedback, raise to 2 or 3 when debugging. Defaults to 1
- `allow_adds` (`bool`, *optional*) – bool Allow adding new variants to a pre-existing Mongo collection, or overwrite it (Default value = False)

**Returns** ~`VAPr.vapr_core.VaprDataset`

**Return type** class

`download_annovar_databases()`

Needed for ANNOVAR to run, it will download the required databases

Args:

Returns:

```
class VAPr.vapr_core.VaprDataset(mongo_db_name, mongo_collection_name,
                                 merged_vcf_path=None)
```

Bases: `object`

`full_name`

Full name of database and collection

Args:

**Returns** Full name of database and collection

**Return type** str

**get\_all\_variants()**  
Self-explanatory

Args:

**Returns** list of variants

**Return type** list

**get\_custom\_filtered\_variants(filter\_dictionary)**

See [Create your own filter](#) for more information on how to implement

**Parameters** `filter_dictionary(dictionary – dict):` mongodb custom filter

**Returns** list of variants

**Return type** list

**get\_de\_novo\_variants(proband, ancestor1, ancestor2)**

See [4. De novo Variants](#) for more information on how this is implemented

**Parameters**

- `proband(str)` – proband variant
- `ancestor1(str)` – ancestor #1 variant
- `ancestor2(str)` – ancestor #2 variant

**Returns** list of variants

**Return type** list

**get\_deleterious\_compound\_heterozygous\_variants(sample\_names\_list=None)**

See [3. Deleterious Compound Heterozygous Variants](#) for more information on how this is implemented

**Parameters** `sample_names_list(list – list, optional):` list of samples to draw variants from (Default value = None)

**Returns** list of variants

**Return type** list

**get\_distinct\_sample\_ids()**

Self-explanatory

Args:

**Returns** list of sample ids

**Return type** list

**get\_known\_disease\_variants(sample\_names\_list=None)**

See [2. Known Disease Variants](#) for more information on how this is implemented

**Parameters** `sample_names_list(list – list, optional):` list of samples to draw variants from (Default value = None)

**Returns** list of variants

**Return type** list

**get\_rare\_deleterious\_variants(sample\_names\_list=None)**

See [1. Rare Deleterious Variants](#) for more information on how this is implemented

**Parameters** `sample_names_list(list – list, optional):` list of samples to draw variants from (Default value = None)

**Returns** list of variants

**Return type** list

**get\_variants\_as\_dataframe** (*filtered\_variants=None*)  
Utility to get a dataframe from variants, either all of them or a filtered subset

**Parameters** **filtered\_variants** – a list of variants (Default value = None)

**Returns** pandas.DataFrame

**get\_variants\_for\_sample** (*sample\_name*)  
Return variants for a specific sample

**Parameters** **sample\_name** (*str*) – name of sample

**Returns** list of variants

**Return type** list

**get\_variants\_for\_samples** (*specific\_sample\_names*)  
Return variants from multiple samples

**Parameters** **specific\_sample\_names** (*list*) – name of samples

**Returns** list of variants

**Return type** list

**is\_empty**  
If there are no records in the collection, returns True

Args:

**Returns** if there are no records in the collection, returns True

**Return type** bool

**num\_records**  
Number of records in MongoDB collection

Args:

**Returns** Number of records in MongoDB collection

**Return type** int

**write\_filtered\_annotated\_csv** (*filtered\_variants, output\_fp*)  
Filtered csv file containing annotations from a list passed to it, coming from MongoDB

**Parameters**

- **filtered\_variants** (*list*) – variants coming from MongoDB
- **output\_fp** (*str*) – Output file path

**Returns** None

**write\_filtered\_annotated\_vcf** (*filtered\_variants, vcf\_output\_path, info\_out=True*)

**Parameters**

- **filtered\_variants** (*list*) – variants coming from MongoDB
- **vcf\_output\_path** (*str*) – Output file path
- **info\_out** – if True, extra annotation information will be written to the vcf file (Default value = True)

- **info\_out** – bool (Default value = True)

**Returns** None

**write\_unfiltered\_annotated\_csv**(*output\_fp*)

Full csv file containing annotations from both annovar and myvariant.info

**Parameters** **output\_fp** (*str*) – Output file path

**Returns** None

**write\_unfiltered\_annotated\_csvs\_per\_sample**(*output\_dir*)

**Parameters** **output\_dir** – return: None

**Returns** None

**write\_unfiltered\_annotated\_vcf**(*vcf\_output\_path*, *info\_out=True*)

Filtered vcf file containing annotations from a list passed to it, coming from MongoDB

**Parameters**

- **vcf\_output\_path** (*str*) – Output file path
- **info\_out** – if True, extra annotation information will be written to the vcf file (Default value = True)
- **info\_out** – bool (Default value = True)

**Returns** None

## 5.8 VAPr.vcf\_genotype\_fields\_parsing module

**class** VAPr.vcf\_genotype\_fields\_parsing.**Allele**(*unfiltered\_read\_counts=None*)  
Bases: *object*

Store unfiltered read counts, if any, for a particular allele.

**unfiltered\_read\_counts**

*int or None* – Number of unfiltered reads counts for this sample at this site, from AD field.

**class** VAPr.vcf\_genotype\_fields\_parsing.**GenotypeLikelihood**(*allele1\_number*, *allele2\_number*, *likelihood\_neg\_exponent*)  
Bases: *object*

Store parsed info from VCF genotype likelihood field for a single sample.

**allele1\_number**

*int* – The allele identifier for the left-hand allele inferred for this genotype likelihood.

**allele2\_number**

*int* – The allele identifier for the right-hand allele inferred for this genotype likelihood.

**likelihood\_neg\_exponent**

*float* – The “normalized” Phred-scaled likelihood of the genotype represented by allele1 and allele2.

**class** VAPr.vcf\_genotype\_fields\_parsing.**VCFGenotypeInfo**(*raw\_string*)  
Bases: *object*

Store parsed info from VCF genotype fields for a single sample.

**\_raw\_string**

*str* – The genotype fields values string from a VCF file (e.g., ‘0/1:173,141:282:99:255,0,255’).

**genotype**

`Optional[str]` – The type of each of the sample's two alleles, such as 0/0, 0/1, etc.

**alleles**

`List[Allele]` – One Allele object for each allele detected for this variant (this can be across samples, so there can be more than 2 alleles).

**genotype\_likelihoods**

`List[GenotypeLikelihood]` – The GenotypeLikelihood object for each allele.

**unprocessed\_info**

`Dict[str, Any]` – Dictionary of field tag and value(s) for any fields not stored in dedicated attributes of VCFGenotypeInfo. Values are parsed to lists and/or floats if possible.

**genotype\_subclass\_by\_class**

`Dict[str, str]` – Genotype subclass (reference, alt, compound) keyed by genotype class (homozygous/heterozygous).

**filter\_passing\_reads\_count**

`int or None` – Filtered depth of coverage of this sample at this site from the DP field.

**genotype\_confidence**

`str` – Genotype quality (confidence) of this sample at this site, from the GQ field.

**class** VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeParser

Bases: `object`

Mine format string and genotype fields string to create a filled VCFGenotypeInfo object.

```
FILTERED_ALLELE_DEPTH_TAG = 'DP'
GENOTYPE_QUALITY_TAG = 'GQ'
GENOTYPE_TAG = 'GT'
NORMALIZED_SCALED_LIKELIHOODS_TAG = 'PL'
UNFILTERED_ALLELE_DEPTH_TAG = 'AD'
```

**static is\_valid\_genotype\_fields\_string(genotype\_fields\_string)**

Return true if input has any real genotype fields content, false if is just periods, zeroes, and delimiters.

**Parameters** `genotype_fields_string(str)` – A VCF-style genotype fields string, such as 1/1:0,2:2:6:89,6,0 or ./:::..

**Returns** bool: true if input has any real genotype fields content, false if is just periods, zeroes, and delimiters.

**classmethod parse(format\_key\_string, format\_value\_string)**

Parse the input format string and genotype fields string into a filled VCFGenotypeInfo object.

**Parameters**

- `format_key_string(str)` – The VCF format string (e.g., ‘GT:AD:DP:GQ:PL’) for this sample at this site.
- `format_value_string(str)` – The VCF genotype fields values string (e.g., ‘1/1:0,34:34:99:1187,2,101,0’) corresponding to the format\_key\_string for this sample at this site.

**Returns**

A filled VCFGenotypeInfo for this sample at this site unless an error was encountered, in which case `None` is returned. encountered, in which case `None` is returned.

Return type `VCFGenotypeInfo` or `None`

## 5.9 VAPr.vcf\_merging module

`VAPr.vcf_merging.bgzip_and_index_vcf(vcf_path)`  
bgzip and index each vcf so it can be merged with bcftools.

`VAPr.vcf_merging.merge_vcfs(input_dir, output_dir, project_name, raw_vcf_path_list=None, vcfs_gzipped=False)`  
Merge vcf files into single multisample vcf, bgzip and index merged vcf file.

## 5.10 Module contents

# CHAPTER 6

---

## Indices and tables

---

- genindex
- modindex
- search



---

## Python Module Index

---

### V

VAPr, 22  
VAPr.annovar\_output\_parsing, 13  
VAPr.annovar\_running, 14  
VAPr.chunk\_processing, 14  
VAPr.filtering, 15  
VAPr.validation, 15  
VAPr.vapr\_core, 16  
VAPr.vcf\_genotype\_fields\_parsing, 20  
VAPr.vcf\_merging, 22



### Symbols

_raw_string (VAPr.vcf_genotype_fields_parsing.VCFGenotypeInfo attribute), 20	CHUNK_SIZE_INDEX (VAPr.chunk_processing.AnnotationJobParamsIndices attribute), 14
	collect_chunk_annotations_and_store() (in module VAPr.chunk_processing), 15
	COLLECTION_NAME_INDEX (VAPr.chunk_processing.AnnotationJobParamsIndices attribute), 15
	convert_to_nonneg_int() (in module VAPr.validation), 15
	convert_to_nullable() (in module VAPr.validation), 15
	CYTOBAND_HEADER (VAPr.annovar_output_parsing.AnnovarTxtParser attribute), 13
A	Allele (class in VAPr.vcf_genotype_fields_parsing), 20
	allele1_number (VAPr.vcf_genotype_fields_parsing.GenotypeLikelihood attribute), 20
	allele2_number (VAPr.vcf_genotype_fields_parsing.GenotypeLikelihood attribute), 20
	ALT_HEADER (VAPr.annovar_output_parsing.AnnovarTxtParser attribute), 13
	ALLELE_DEPTH_KEY (VAPr.annovar_output_parsing.AnnovarAnnotatedVariant attribute), 13
	annotation() (VAPr.vapr_core.VaprAnnotator method), 16
	annotate_lite() (VAPr.vapr_core.VaprAnnotator method), 17
	AnnotationJobParamsIndices (class in VAPr.chunk_processing), 14
	AnnovarAnnotatedVariant (class in VAPr.annovar_output_parsing), 13
	AnnovarTxtParser (class in VAPr.annovar_output_parsing), 13
	AnnovarWrapper (class in VAPr.annovar_running), 14
B	bgzip_and_index_vcf() (in module VAPr.vcf_merging), 22
C	CHR_HEADER (VAPr.annovar_output_parsing.AnnovarTxtParser attribute), 13
	CHUNK_INDEX_INDEX (VAPr.chunk_processing.AnnotationJobParamsIndices attribute), 14
D	DB_NAME_INDEX (VAPr.chunk_processing.AnnotationJobParamsIndices attribute), 15
	DEFAULT_GENOME_VERSION (VAPr.vapr_core.VaprAnnotator attribute), 16
	download_annoivar_databases() (VAPr.vapr_core.VaprAnnotator method), 17
	download_databases() (VAPr.annovar_running.AnnovarWrapper method), 14
E	END_HEADER (VAPr.annovar_output_parsing.AnnovarTxtParser attribute), 13
	ESP6500_ALL_HEADER (VAPr.annovar_output_parsing.AnnovarTxtParser attribute), 14
	EXONICFUNC_KNOWNGENE_HEADER (VAPr.annovar_output_parsing.AnnovarTxtParser attribute), 14
F	FILE_PATH_INDEX (VAPr.chunk_processing.AnnotationJobParamsIndices attribute), 15

filter\_passing\_reads\_count  
 (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeInfo  
 attribute), 21

FILTER\_PASSING\_READS\_COUNT\_KEY  
 (VAPr.annovar\_output\_parsing.AnnovarAnnotatedVariant)  
 attribute), 13

FILTERED\_ALLELE\_DEPTH\_TAG  
 (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeParser  
 attribute), 21

full\_name (VAPr.vapr\_core.VaprDataset attribute), 17

FUNC\_KNOWNGENE\_HEADER  
 (VAPr.annovar\_output\_parsing.AnnovarTxtParser  
 attribute), 14

**G**

GENE\_KNOWNGENE\_HEADER  
 (VAPr.annovar\_output\_parsing.AnnovarTxtParser  
 attribute), 14

GENEDETAIL\_KNOWNGENE\_HEADER  
 (VAPr.annovar\_output\_parsing.AnnovarTxtParser  
 attribute), 14

GENOME\_BUILD\_VERSION\_INDEX  
 (VAPr.chunk\_processing.AnnotationJobParamsIndices  
 attribute), 15

GENOMIC\_SUPERDUPS\_HEADER  
 (VAPr.annovar\_output\_parsing.AnnovarTxtParser  
 attribute), 14

genotype (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeInfo  
 attribute), 20

genotype\_confidence (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeInfo), 16

attribute), 21

GENOTYPE\_KEY (VAPr.annovar\_output\_parsing.AnnovarAnnotatedVariant), 14

attribute), 13

genotype\_likelihoods (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeInfo), 14

attribute), 21

GENOTYPE\_LIKELIHOODS\_KEY  
 (VAPr.annovar\_output\_parsing.AnnovarAnnotatedVariant  
 attribute), 13

GENOTYPE\_QUALITY\_TAG  
 (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeParser  
 attribute), 21

genotype\_subclass\_by\_class  
 (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeInfo  
 attribute), 21

GENOTYPE\_SUBCLASS\_BY\_CLASS\_KEY  
 (VAPr.annovar\_output\_parsing.AnnovarAnnotatedVariant  
 attribute), 13

GENOTYPE\_TAG (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeParser  
 attribute), 21

GenotypeLikelihood  
 (class  
 VAPr.vcf\_genotype\_fields\_parsing), 20

get\_all\_variants() (VAPr.vapr\_core.VaprDataset method), 17

get\_any\_of\_sample\_ids\_filter() (in module VAPr.filtering), 15

get\_custom\_filtered\_variants()  
 (VAPr.vapr\_core.VaprDataset method), 18

get\_deleterious\_compound\_heterozygous\_variants()  
 (VAPr.vapr\_core.VaprDataset method), 18

get\_distinct\_sample\_ids() (VAPr.vapr\_core.VaprDataset  
 method), 18

get\_known\_disease\_variants()  
 (VAPr.vapr\_core.VaprDataset method), 18

get\_num\_possible\_indices()  
 (VAPr.chunk\_processing.AnnotationJobParamsIndices  
 class method), 15

get\_rare\_deleterious\_variants()  
 (VAPr.vapr\_core.VaprDataset method), 18

get\_sample\_id\_filter() (in module VAPr.filtering), 15

get\_variants\_as\_dataframe()  
 (VAPr.vapr\_core.VaprDataset method), 19

get\_variants\_for\_sample() (VAPr.vapr\_core.VaprDataset  
 method), 19

get\_variants\_for\_samples() (VAPr.vapr\_core.VaprDataset  
 method), 19

**H**

HG19\_VERSION (VAPr.vapr\_core.VaprAnnotator  
 attribute), 16

HG38\_VERSION (VAPr.vapr\_core.VaprAnnotator  
 attribute), 16

hg\_19\_databases (VAPr.annovar\_running.AnnovarWrapper  
 attribute), 14

hg\_38\_databases (VAPr.annovar\_running.AnnovarWrapper  
 attribute), 14

HGV\_ID\_KEY (VAPr.annovar\_output\_parsing.AnnovarAnnotatedVariant  
 attribute), 13

**I**

is\_empty (VAPr.vapr\_core.VaprDataset attribute), 19

is\_valid\_genotype\_fields\_string()  
 (VAPr.vcf\_genotype\_fields\_parsing.VCFGenotypeParser  
 static method), 21

**L**

likelihood\_neg\_exponent  
 (VAPr.vcf\_genotype\_fields\_parsing.GenotypeLikelihood  
 attribute), 20

**M**

make\_de\_novo\_variants\_filter() (in module VAPr.filtering), 15

make\_deleterious\_compound\_heterozygous\_variants\_filter()  
 (in module VAPr.filtering), 15

**N**

- make\_known\_disease\_variants\_filter() (in module `STANDARDIZED_CHR_MT_SUFFIX_VAL` (`VAPr.filtering`), 15)
  - (`VAPr.annovar_output_parsing.AnnovarAnnotatedSTANDARDIZED_CHR_MT_VAL` class method), 13
- make\_rare\_deleterious\_variants\_filter() (in module `STANDARDIZED_CHR_MT_VAL` (`VAPr.filtering`), 15)
  - (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14
- merge\_vcfs() (in module `VAPr.vcf_merging`), 22
  - (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14

**O**

- NCI60\_HEADER (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 16
- NORMALIZED\_SCALED\_LIKELIHOODS\_TAG (`VAPr.vcf_genotype_fields_parsing.VCFGenotypeParser` attribute), 21
- num\_records (`VAPr.vapr_core.VaprDataset` attribute), 19

**P**

- parse() (`VAPr.vcf_genotype_fields_parsing.VCFGenotypeParser` class method), 21

**R**

- RAW\_CHR\_MT\_SUFFIX\_VAL (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14
- RAW\_CHR\_MT\_VAL (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14
- read\_chunk\_of\_annotations\_to\_dicts\_list() (`VAPr.annovar_output_parsing.AnnovarTxtParser` class method), 14
- REF\_HEADER (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14
- run\_annotation() (`VAPr.annovar_running.AnnovarWrapper` method), 14

**S**

- SAMPLE\_ID\_KEY (`VAPr.annovar_output_parsing.AnnovarAnnotatedVariant` attribute), 13
- SAMPLE\_LIST\_INDEX (`VAPr.chunk_processing.AnnotationJobParamsIndices` attribute), 15
- SAMPLE\_NAMES\_KEY (`VAPr.vapr_core.VaprAnnotator` attribute), 16

**T**

- TEBS\_CONS\_SITES\_HEADER (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14
- THOU\_G\_2015\_ALL\_HEADER (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14

**U**

- UNFILTERED\_ALLELE\_DEPTH\_TAG (`VAPr.vcf_genotype_fields_parsing.VCFGenotypeParser` attribute), 21
- unfiltered\_read\_counts (`VAPr.vcf_genotype_fields_parsing.Allele` attribute), 20
- unprocessed\_info (`VAPr.vcf_genotype_fields_parsing.VCFGenotypeInfo` attribute), 21

**V**

- VAPr (module), 22
- VAPr.annovar\_output\_parsing (module), 13
- VAPr.annovar\_running (module), 14
- VAPr.chunk\_processing (module), 14
- VAPr.filtering (module), 15
- VAPr.validation (module), 15
- VAPr.vapr\_core (module), 16
- VAPr.vcf\_genotype\_fields\_parsing (module), 20
- VAPr.vcf\_merging (module), 22
- VaprAnnotator (class in `VAPr.vapr_core`), 16
- VaprDataset (class in `VAPr.vapr_core`), 17
- VCFGenotypeInfo (class in `VAPr.vcf_genotype_fields_parsing`), 20
- VCFGenotypeParser (class in `VAPr.vcf_genotype_fields_parsing`), 21

**W**

- VERBOSE\_LEVEL\_INDEX (`VAPr.chunk_processing.AnnotationJobParamsIndices` attribute), 15
- SAMPLES\_KEY (`VAPr.annovar_output_parsing.AnnovarAnnotatedVariant` attribute), 13
  - (`VAPr.annovar_output_parsing.AnnovarAnnotatedVariant` attribute), 13
  - (`VAPr.annovar_output_parsing.AnnovarAnnotatedVariant` attribute), 13
- SCORE\_KEY (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14
  - (`VAPr.annovar_output_parsing.AnnovarTxtParser` attribute), 14
  - (`VAPr.vapr_core.VaprDataset` method), 19
  - (`VAPr.vapr_core.VaprDataset` method), 19
- SCORE\_KEY (`VAPr.vapr_core.VaprDataset` method), 19
  - (`VAPr.vapr_core.VaprDataset` method), 19
- SUPPORTED\_GENOME\_BUILD\_VERSIONS (`VAPr.vapr_core.VaprAnnotator` attribute), 14

```
write_unfiltered_annotated_csv()
    (VAPr.vapr_core.VaprDataset method), 20
write_unfiltered_annotated_csvs_per_sample()
    (VAPr.vapr_core.VaprDataset method), 20
write_unfiltered_annotated_vcf()
    (VAPr.vapr_core.VaprDataset method), 20
```