

---

# **KripoDB Documentation**

*Release 3.0.0*

**Stefan Verhoeven**

**Mar 30, 2018**



---

## Contents

---

<b>1</b>	<b>Data update</b>	<b>3</b>
1.1	Baseline update . . . . .	3
1.2	Incremental update . . . . .	6
1.3	Steps . . . . .	10
1.4	Disk layout . . . . .	10
1.5	Requirements . . . . .	11
<b>2</b>	<b>DiVE visualization</b>	<b>13</b>
2.1	1. LargeVis input file from Kripo similarity matrix . . . . .	13
2.2	2. Perform embedding using LargeVis . . . . .	14
2.3	3. Generate DiVE metadata datafiles . . . . .	15
2.4	4. Create DiVE input file . . . . .	15
<b>3</b>	<b>API</b>	<b>17</b>
3.1	kripodb.canned . . . . .	17
3.2	kripodb.db . . . . .	20
3.3	kripodb.dive . . . . .	24
3.4	kripodb.frozen . . . . .	25
3.5	kripodb.hdf5 . . . . .	27
3.6	kripodb.makebits . . . . .	31
3.7	kripodb.modifiedtanimoto . . . . .	32
3.8	kripodb.pairs . . . . .	33
3.9	kripodb.pharmacophores . . . . .	35
3.10	kripodb.pdb . . . . .	38
3.11	kripodb.script . . . . .	38
3.12	kripodb.webservice . . . . .	39
<b>4</b>	<b>Command line interface</b>	<b>43</b>
4.1	Positional Arguments . . . . .	43
4.2	Named Arguments . . . . .	43
4.3	Sub-commands: . . . . .	43
<b>5</b>	<b>Indices and tables</b>	<b>59</b>
	<b>Python Module Index</b>	<b>61</b>



For installation and usage see <https://github.com/3D-e-Chem/kripodb/blob/master/README.md>



The Kripo data can be updated in 2 ways:

### 1.1 Baseline update

#### Contents

- *Baseline update*
  - *1. Create staging directory*
  - *2. Create sub-pocket pharmacophore fingerprints*
  - *3. Create fragment information*
    - \* *1. Fragment shelve*
    - \* *2. Fragment sdf*
    - \* *3. Pharmacophores*
  - *4. Add new fragment information to fragment sqlite db*
  - *5. Populate PDB metadata in fragments database*
  - *6. Check no fragments are duplicated*
  - *7. Calculate similarity scores between fingerprints*
  - *8. Convert pairs file into dense similarity matrix*
  - *9. Switch staging to current*
  - *10.0 Update web service*

The Kripo data set is generated from scratch every year or when algorithms change.

### 1.1.1 1. Create staging directory

Setup path with update scripts using:

```
export SCRIPTS=$PWD/../../kripodb/update_scripts
```

Create a new directory:

```
mkdir staging  
cd ..
```

### 1.1.2 2. Create sub-pocket pharmacophore fingerprints

Use directory listing of new pdb files as input:

```
ls $PDBS_ADDED_DIR | pdblast2fps_final_local.py
```

---

**Todo:** Too slow when run on single cpu. Chunkify input, run in parallel and merge results

---

### 1.1.3 3. Create fragment information

#### 1. Fragment shelve

Where the fragment came from is stored in a Python shelve file. It can be generated from the pharmacophore files using:

```
compiledDatabase.py
```

#### 2. Fragment sdf

The data generated thus far contains the molblocks of the ligands and atom nrs of each fragment. The fragment molblocks can be generated into a fragment sdf file with:

```
fragid2sd.py > fragments.sd
```

#### 3. Pharmacophores

The raw pharmacophores are stored in the FRAGMENT\_PPHORES sub-directory. Each pocket has a \*\_pphore.sd.gz file which contains the pharmacophore points of the whole pocket and a \*\_pphores.txt file which contains the indexes of pharmacophore points for each sub pocket or fragment. The raw pharmacophores need to be added to the pharmacophores datafile with:

```
kripodb pharmacophores add FRAGMENT_PPHORES pharmacophores.h5
```



### 1.1.4 4. Add new fragment information to fragment sqlite db

The following commands add the fragment shelve and sdf to the fragments database:

```
cp ../current/fragments.sqlite .
kripodb fragments shelve fragments.shelve fragments.sqlite
kripodb fragments sdf fragments.sd fragments.sqlite
```

Step 4 and 5 can be submitted to scheduler with:

```
jid_db=$(sbatch --parsable -n 1 -J db_append $SCRIPTS/db_append.sh)
```

### 1.1.5 5. Populate PDB metadata in fragments database

The following command will updated the PDB metadata to fragments database:

```
kripodb fragments pdb fragments.sqlite
```

### 1.1.6 6. Check no fragments are duplicated

The similarity matrix can not handle duplicates. It will result in addition of scores:

```
jid_dups=$(sbatch --parsable -n 1 -J check_dups --dependency=afterok:$jid_db $SCRIPTS/
↳baseline_duplicates.sh)
```

### 1.1.7 7. Calculate similarity scores between fingerprints

The similarities between fingerprints can be calculated with:

```
all_chunks=$(ls *fp.gz |wc -l)
jid_fpunzip=$(sbatch --parsable -n $all_chunks -J fpunzip --dependency=afterok:$jid_
↳dups $SCRIPTS/baseline_fpunzip.sh)
nr_chunks=$((($all_chunks * $all_chunks / 2 - $all_chunks))"
jid_fpneigh=$(sbatch --parsable -n $nr_chunks -J fpneigh --dependency=afterok:$jid_
↳fpunzip $SCRIPTS/baseline_similarities.sh)
jid_fpzip=$(sbatch --parsable -n $all_chunks -J fpzip --dependency=afterok:$jid_
↳fpneigh $SCRIPTS/baseline_fpzip.sh)
jid_merge_matrices=$(sbatch --parsable -n 1 -J merge_matrices --dependency=afterok:
↳$jid_fpneigh $SCRIPTS/baseline_merge_similarities.sh)
```

To prevent duplicates similarities of a chunk against itself should ignore the upper triangle.

---

**Todo:** Don't fpneigh run sequentially but submit to batch queue system and run in parallel

---

### 1.1.8 8. Convert pairs file into dense similarity matrix

---

**Tip:** Converting the pairs file into a dense matrix goes quicker with more memory.

---

The following commands converts the pairs into a compressed dense matrix:

```
jid_compress_matrix=$(sbatch --parsable -n 1 -J compress_matrix --dependency=afterok:  
↔$jid_merge_matrices $SCRIPTS/freeze_similarities.sh)
```

The output of this step is ready to be served as a webservice using the *kripodb serve* command.

### 1.1.9 9. Switch staging to current

The webserver and webservice are configured to look in the *current* directory for files.

The staging can be made current with the following commands:

```
mv current old  
mv staging current
```

### 1.1.10 10.0 Update web service

The webservice running at <http://3d-e-chem.vu-compmedchem.nl/kripodb> must be updated with the new datafiles.

The following files must be copied to the server

- fragments.sqlite
- pharmacophores.h5
- similarities.packedfrozen.h5

The webservice must be restarted.

To show how up to date the webservice is the release date of the latest PDB is stored in *version.txt* which can be reached at <http://3d-e-chem.vu-compmedchem.nl/kripodb/version.txt>. The content *version.txt* must be updated.

## 1.2 Incremental update

### Contents

- *Incremental update*
  - 1. Create staging directory
  - 2. Create sub-pocket pharmacophore fingerprints
  - 3. Create fragment information
    - \* 1. Fragment shelve
    - \* 2. Fragment sdf
    - \* 3. Pharmacophores
  - 4. Add new fragment information to fragment sqlite db
  - 5. Populate PDB metadata in fragments database
  - 6. Check no fragments are duplicated
  - 7. Calculate similarity scores between fingerprints

- 8. Convert pairs file into dense similarity matrix
- 9. Switch staging to current
  - \* 9.1 Merge fingerprint files (optional)
- 10.0 Update web service

The Kripo data set can be incrementally updated with new PDB entries.

## 1.2.1 1. Create staging directory

Setup path with update scripts using:

```
export SCRIPTS=$PWD/../../kripodb/update_scripts
```

Create a new directory:

```
mkdir staging
cd ..
```

## 1.2.2 2. Create sub-pocket pharmacophore fingerprints

The *ids.txt* file must contain a list of PDB identifiers which have not been processed before. It can be fetched from <https://www.rcsb.org/>.

Adjust the PDB save location in the *singleprocess.py* script to the staging directory.

Run the following command to generate fragments/pharmacophores/fingerprints for each PDB listed in *ids.txt*:

```
python singleprocess.py
```

## 1.2.3 3. Create fragment information

### 1. Fragment shelf

Where the fragment came from is stored in a Python shelf file. It can be generated from the pharmacophore files using:

```
compiledDatabase.py
```

### 2. Fragment sdf

The data generated thus far contains the molblocks of the ligands and atom nrs of each fragment. The fragment molblocks can be generated into a fragment sdf file with:

```
fragid2sd.py fragments.shelve > fragments.sd
```

### 3. Pharmacophores

The raw pharmacophores are stored in the FRAGMENT\_PPHORES sub-directory. Each pocket has a \*\_pphore.sd.gz file which contains the pharmacophore points of the whole pocket and a \*\_pphores.txt file which contains the indexes of pharmacophore points for each sub pocket or fragment. The raw pharmacophores of the update can be added to the existing pharmacophores datafile with:

```
cp ../current/pharmacophores.h5 .
kripodb pharmacophores add FRAGMENT_PPHORES pharmacophores.h5
```

#### 1.2.4 4. Add new fragment information to fragment sqlite db

The following commands add the fragment shelve and sdf to the fragments database:

```
cp ../current/fragments.sqlite .
kripodb fragments shelve fragments.shelve fragments.sqlite
kripodb fragments sdf fragments.sd fragments.sqlite
```

Step 4 and 5 can be submitted to scheduler with:

```
jid_db=$(sbatch --parsable -n 1 -J db_append $SCRIPTS/db_append.sh)
```

#### 1.2.5 5. Populate PDB metadata in fragments database

The following command will updated the PDB metadata to fragments database:

```
kripodb fragments pdb fragments.sqlite
```

#### 1.2.6 6. Check no fragments are duplicated

The similarity matrix can not handle duplicates. It will result in addition of scores:

```
jid_dups=$(sbatch --parsable -n 1 -J check_dups --dependency=afterok:$jid_db $SCRIPTS/
↪incremental_duplicates.sh)
```

#### 1.2.7 7. Calculate similarity scores between fingerprints

The similarities between the new and existing fingerprints and between new fingerprints themselves can be calculated with:

```
current_chunks=$(ls ../current/*fp.gz |wc -l)
all_chunks=$((current_chunks + 1))
jid_fpneigh=$(sbatch --parsable -n $all_chunks -J fpneigh --dependency=afterok:$jid_
↪dups $SCRIPTS/incremental_similarities.sh)
jid_merge_matrices=$(sbatch --parsable -n 1 -J merge_matrices --dependency=afterok:
↪$jid_fpneigh $SCRIPTS/incremental_merge_similarities.sh)
```

## 1.2.8 8. Convert pairs file into dense similarity matrix

**Note:** Converting the pairs file into a dense matrix goes quicker with more memory.

The frame size (-f) should be as big as possible, 100000000 requires 6Gb RAM.

The following commands converts the pairs into a compressed dense matrix:

```
jid_compress_matrix=$(sbatch --parsable -n 1 -J compress_matrix --dependency=afterok:
↪$jid_merge_matrices $SCRIPTS/freeze_similarities.sh)
```

The output of this step is ready used to find similar fragments, using either the webservice with the *kripodb serve* command or with the *kripodb similarities similar* command directly.

## 1.2.9 9. Switch staging to current

The webserver and webservice are configure to look in the *current* directory for files.

The current and new pharmacophores need to be combined:

```
mv staging/FRAGMENT_PPHORES staging/FRAGMENT_PPHORES.new
rsync -a current/FRAGMENT_PPHORES staging/FRAGMENT_PPHORES
rm -r staging/FRAGMENT_PPHORES.new
```

**Todo:** rsync of current/FRAGMENT\_PPHORES to destination, maybe too slow due large number of files. Switch to move old pharmacohores and rsync new pharmacophores into it when needed.

The current and new fingerprints need to be combined:

```
cp -n current/*.fp.gz staging/
```

The staging can be made current with the following commands:

```
mv current old && mv staging current
```

### 9.1 Merge fingerprint files (optional)

To keep the number of files to a minimum it is advised to merge the fingerprint files from incremental updates of a year.

The incremental fingerprint files are named like *out.<year><week>.fp.gz*, to generate *kripo\_fingerprints\_<year>\_fp.gz* run:

```
sbatch --parsable -n 1 -J merge_fp $SCRIPTS/incremental_merge_fp.sh <year>
```

## 1.2.10 10.0 Update web service

The webservice running at <http://3d-e-chem.vu-compmedchem.nl/kripodb> must be updated with the new datafiles.

The following files must copied to the server

- fragments.sqlite

- pharmacophores.h5
- similarities.packedfrozen.h5

The webservice must be restarted.

To show how up to date the webservice is the release date of the latest PDB is stored in *version.txt* which can be reached at <http://3d-e-chem.vu-compmedchem.nl/kripodb/version.txt> The content *version.txt* must be updated.

## 1.3 Steps

Overview of steps involved in updating Kripo:

1. Create staging directory
2. Create sub-pocket pharmacophore fingerprints
3. Create fragment information
4. Add new fragment information to fragment sqlite db
5. Populate PDB metadata in fragments database
6. Check no fragments are duplicated
7. Calculate similarity scores between fingerprints
8. Convert pairs file into dense similarity matrix
9. Switch staging to current
10. Update web service

---

**Note:** Steps 2 through 3 require undisclosed scripts or <https://github.com/3D-e-Chem/kripo>

---

---

**Note:** Steps 4 and 6 through 7 can be done using the KripoDB Python library.

---

---

**Todo:** Remove Kripo fragment/fingerprints of obsolete PDBs (<ftp://ftp.wwpdb.org/pub/pdb/data/status/obsolete.dat>)

---

## 1.4 Disk layout

Directories for Kripo:

- **current/**, directory which holds current dataset
- **staging/**, which is used to compute new items and combine new and old items.
- **old/**, which is used as a backup containing the previous update.

Files and directories for a data set (inside current, staging and old directories):

- **pharmacophores.h5**, pharmacophores database file
- **out.fp.sqlite**, fingerprints file
- **fragments.sqlite**, fragment information database file

- **similarities.h5**, similarities as pairs table
- **similarities.packedfrozen.h5**, similarities as dense matrix

Input directories:

- **\$PDBS\_ADDED\_DIR**, directory containing new PDB files to be processed

## 1.5 Requirements

- Slurm batch scheduler
- KripoDB and it's dependencies installed and in path
- Posix filesystem, NFS or Virtualbox share do not accept writing of hdf5 or sqlite files





DiVE homepage at <https://github.com/NLeSC/DiVE>

The Kripo similarity matrix can be embedded to 2D or 3D using largevis and then visualized using DiVE.

Steps

1. LargeVis input file from Kripo similarity matrix
2. Perform embedding using LargeVis
3. Generate DiVE metadata datafiles
4. Create DiVE input file

Input datasets

1. only fragment1 or whole unfragmented ligands
2. all fragments
3. only gpcr frag1
4. only kinase frag1
5. only gpcr and kinase frag1

Output datasets

1. 2D
2. 3D

### 2.1 1. LargeVis input file from Kripo similarity matrix

Dump the similarity matrix to csv of \*frag1 fragments:

```
kripodb similarities export --no_header --frag1 similarities.h5 similarities.frag1.txt
```

## 2.1.1 Similarities between GPCR pdb entries

Use the [GPCRDB web service](#) to fetch a list of PDB codes which contain GPCR proteins:

```
curl -X GET --header 'Accept: application/json' 'http://gpcrdb.org/services/structure/
↳' | jq -r '.[] | .pdb_code' > pdb.gpcr.txt
```

Dump the similarity matrix to csv:

```
kripodb similarities export --no_header --frag1 --pdb pdb.gpcr.txt similarities.h5_
↳similarities.frag1.gpcr.txt
```

## 2.1.2 Similarities between GPCR and Kinase pdb entries

Use the [KLIFS KNIME nodes](#) to create a file with of PDB codes of Kinases called *pdb.kinase.txt*.

Dump the similarity matrix to csv:

```
cat pdb.gpcr.txt pdb.kinase.txt > pdb.gpcr.kinase.txt
kripodb similarities export --no_header --frag1 --pdb pdb.gpcr.kinase.txt_
↳similarities.h5 similarities.frag1.gpcr.kinase.txt
```

## 2.2 2. Perform embedding using LargeVis

Get or compile LargeVis binaries from <https://github.com/lferry007/LargeVis>

Compile using miniconda:

```
conda install gsl gcc
cd LargeVis/Linux
c++ LargeVis.cpp main.cpp -o LargeVis -lm -pthread -lgsl -lgslcblas -Ofast -Wl,-rpath,
↳$CONDA_PREFIX/lib -march=native -ffast-math
cp LargeVis $CONDA_PREFIX/bin/
```

Then embed frag1 similarity matrix in 3D with:

```
LargeVis -fea 0 -outdim 3 -threads $(nproc) -input similarities.frag1.txt -output_
↳largevis.frag1.3d.txt
```

Then embed frag1 similarity matrix in 2D with:

```
LargeVis -fea 0 -outdim 2 -threads $(nproc) -input similarities.frag1.txt -output_
↳largevis.frag1.2d.txt
```

Then embed similarity matrix in 3D with:

```
LargeVis -fea 0 -outdim 3 -threads $(nproc) -input similarities.txt -output largevis.
↳3d.txt
```

Then embed similarity matrix in 2D with:

```
LargeVis -fea 0 -outdim 2 -threads $(nproc) -input similarities.txt -output largevis.
↳2d.txt
```

The *kripo export* in step 1 and the LargeVis command can be submitted to scheduler with:

```

sbatch -n 1 $SCRIPTS/dive_frag1.sh
sbatch -n 1 $SCRIPTS/dive_frag1_gpcr_kinase.sh

```

## 2.3 3. Generate DiVE metadata datafiles

Command to generate properties files:

```

wget -O uniprot.txt 'http://www.uniprot.org/uniprot/?query=database:pdb&format=tab&
↳columns=id,genes(PREFERRED),families,database(PDB) '
kripodb dive export --pdhtags pdb.gpcr.txt --pdhtags pdb.kinase.txt fragments.sqlite_
↳uniprot.txt

```

Will generate in current working directory the following files:

- kripo.props.txt
- kripo.pronames.txt

## 2.4 4. Create DiVE input file

DiVE has a script which can combine the LargeVis coordinates together with metadata. Download the `MakeVizDataWithProperMetadata.py` script from [https://github.com/NLeSC/DiVE/blob/master/scripts\\_prepareData/MakeVizDataWithProperMetadata.py](https://github.com/NLeSC/DiVE/blob/master/scripts_prepareData/MakeVizDataWithProperMetadata.py)

For more information about the script see <https://github.com/NLeSC/DiVE#from-output-of-largevis-to-input-of-dive>.

Example command to generate new DiVE input file:

```

python MakeVizDataWithProperMetadata.py -coord largevis2.similarities_frag1.gpcr.
↳kinase.txt -metadata kripo.props.txt -np kripo.pronames.txt -json largevis2.
↳similarities_frag1.gpcr.kinase.json -dir .

```

The generated file (`largevis2.similarities_frag1.gpcr.kinase.json`) can be uploaded at <https://nlesc.github.io/DiVE/> to visualize.



### 3.1 kripodb.canned

Module with functions which use pandas DataFrame as input and output.

For using Kripo data files inside KNIME (<http://www.knime.org>)

**exception** `kripodb.canned.IncompleteHits` (*absent\_identifiers, hits*)

`kripodb.canned.fragments_by_id` (*fragment\_ids, fragments\_db\_filename\_or\_url, prefix=""*)  
Retrieve fragments based on fragment identifier.

#### Parameters

- **fragment\_ids** (*List[str]*) – List of fragment identifiers
- **fragments\_db\_filename\_or\_url** (*str*) – Filename of fragments db or base url of kripodb webservice
- **prefix** (*str*) – Prefix for output columns

#### Examples

Fetch fragments of '2n2k\_MTN\_frag1' fragment identifier

```
>>> from kripodb.canned import fragments_by_id
>>> fragment_ids = pd.Series(['2n2k_MTN_frag1'])
>>> fragments = fragments_by_id(fragment_ids, 'data/fragments.sqlite')
>>> len(fragments)
1
```

Retrieved from web service instead of local fragments db file. Make sure the web service is running, for example by `kripodb serve data/similarities.h5 data/fragments.sqlite data/pharmacophores.h5`.

```
>>> fragments = fragments_by_id(fragment_ids, , 'http://localhost:8084/kripo')
>>> len(fragments)
1
```

**Returns** Data frame with fragment information

**Return type** `pandas.DataFrame`

**Raises** `IncompleteFragments` – When one or more of the identifiers could not be found.

`kripodb.canned.fragments_by_pdb_codes` (*pdb\_codes*, *fragments\_db\_filename\_or\_url*, *prefix=""*)  
Retrieve fragments based on PDB codes.

See <http://www.rcsb.org/pdb/> for PDB structures.

#### Parameters

- **pdb\_codes** (*List[str]*) – List of PDB codes
- **fragments\_db\_filename\_or\_url** (*str*) – Filename of fragments db or base url of kripodb webservice
- **prefix** (*str*) – Prefix for output columns

#### Examples

Fetch fragments of '2n2k' PDB code

```
>>> from kripodb.canned import fragments_by_pdb_codes
>>> pdb_codes = pd.Series(['2n2k'])
>>> fragments = fragments_by_pdb_codes(pdb_codes, 'data/fragments.sqlite')
>>> len(fragments)
3
```

Retrieved from web service instead of local fragments db file. Make sure the web service is running, for example by `kripodb serve data/similarities.h5 data/fragments.sqlite data/pharmacophores.h5`.

```
>>> fragments = fragments_by_pdb_codes(pdb_codes, 'http://localhost:8084/kripo')
>>> len(fragments)
3
```

**Returns** Data frame with fragment information

**Return type** `pandas.DataFrame`

**Raises** `IncompleteFragments` – When one or more of the identifiers could not be found.

`kripodb.canned.pharmacophores_by_id` (*fragment\_ids*, *pharmacophores\_db\_filename\_or\_url*)  
Fetch pharmacophore points by fragment identifiers

#### Parameters

- **fragment\_ids** (*pd.Series*) – List of fragment identifiers
- **pharmacophores\_db\_filename\_or\_url** – Filename of pharmacophores db or base url of kripodb webservice

**Returns**

**Pandas series with pharmacophores as string in phar format.** Fragment without pharmacophore will return None

**Return type** pandas.Series

## Examples

Fragments similar to '3j7u\_NDP\_frag24' fragment.

```
>>> from kripodb.canned import pharmacophores_by_id
>>> fragment_ids = pd.Series(['2n2k_MTN_frag1'], ['Row0'])
>>> pharmacophores = pharmacophores_by_id(fragment_ids, 'data/pharmacophores.h5')
>>> len(pharmacophores)
1
```

Retrieved from web service instead of local pharmacophores db file. Make sure the web service is running, for example by `kripodb serve data/similarities.h5 data/fragments.sqlite data/pharmacophores.h5`.

```
>>> pharmacophores = pharmacophores_by_id(fragment_ids, 'http://localhost:8084/
↳kripo')
>>> len(pharmacophores)
1
```

`kripodb.canned.similarities` (*queries*, *similarity\_matrix\_filename\_or\_url*, *cutoff*, *limit=1000*)

Find similar fragments to queries based on similarity matrix.

### Parameters

- **queries** (*List[str]*) – Query fragment identifiers
- **similarity\_matrix\_filename\_or\_url** (*str*) – Filename of similarity matrix file or base url of kripodb webservice
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **limit** (*int*) – Maximum number of hits for each query. Default is 1000. Use is None for no limit.

## Examples

Fragments similar to '3j7u\_NDP\_frag24' fragment.

```
>>> import pandas as pd
>>> from kripodb.canned import similarities
>>> queries = pd.Series(['3j7u_NDP_frag24'])
>>> hits = similarities(queries, 'data/similarities.h5', 0.55)
>>> len(hits)
11
```

Retrieved from web service instead of local similarity matrix file. Make sure the web service is running, for example by `kripodb serve data/similarities.h5 data/fragments.sqlite data/pharmacophores.h5`.

```
>>> hits = similarities(queries, 'http://localhost:8084/kripo', 0.55)
>>> len(hits)
11
```

**Returns** Data frame with `query_fragment_id`, `hit_frag_id` and `score` columns

**Return type** `pandas.DataFrame`

**Raises** `IncompleteHits` – When one or more of the identifiers could not be found.

## 3.2 kripodb.db

Fragments and fingerprints sqlite based data storage.

Registers `BitMap` and `molblockgz` data types in sqlite.

**class** `kripodb.db.FastInserter` (*cursor*)

Use with to make inserting faster, but less safe

By setting journal mode to WAL and turn synchronous off.

**Parameters** `cursor` (`sqlite3.Cursor`) – Sqlite cursor

### Examples

```
>>> with FastInserter(cursor):
    cursor.executemany('INSERT INTO table VALUES (?, rows)')
```

**class** `kripodb.db.FingerprintsDb` (*filename*)

Fingerprints database

**as\_dict** (*number\_of\_bits=None*)

Returns a dict-like object to query and alter fingerprints db

**Parameters** `number_of_bits` (*Optional[int]*) – Number of bits that all fingerprints have

**Returns** `BitMapDict`

**create\_tables** ()

Abstract method which is called after connecting to database so tables can be created.

Use `CREATE TABLE IF NOT EXISTS ...` in method to prevent duplicate create errors.

**class** `kripodb.db.FragmentsDb` (*filename*)

Fragments database

**add\_fragment** (*frag\_id, pdb\_code, prot\_chain, het\_code, frag\_nr, atom\_codes, hash\_code, het\_chain, het\_seq\_nr, nr\_r\_groups*)

Add fragment to database

**Parameters**

- **frag\_id** (*str*) – Fragment identifier
- **pdb\_code** (*str*) – Protein databank identifier
- **prot\_chain** (*str*) – Major chain of pdb on which pharmacophore is based
- **het\_code** (*str*) – Ligand/Hetero code
- **frag\_nr** (*int*) – Fragment number, whole ligand has number 1, fragments are >1
- **atom\_codes** (*str*) – Comma separated list of HETATOM atom names which make up the fragment (hydrogens are excluded)
- **hash\_code** (*str*) – Unique identifier for fragment



- **het\_chain** (*str*) – Chain ligand is part of
- **het\_seq\_nr** (*int*) – Residue sequence number of ligand the fragment is a part of
- **nr\_r\_groups** (*int*) – Number of R groups in fragment

**add\_fragments\_from\_shelve** (*myshelve*, *skipdups=False*)

Adds fragments from shelve to fragments table.

Also creates index on `pdb_code` column.

#### Parameters

- **myshelve** (*Dict [Fragment]*) – Dictionary with fragment identifier as key and fragment as value.
- **skipdups** (*bool*) – Skip duplicates, instead of dieing one first duplicate

**add\_molecule** (*mol*)

Adds molecule to molecules table

Uses the name of the molecule as the primary key.

**Parameters** **mol** (*rdkit.Chem.AllChem.Mol*) – the rdkit molecule

**add\_molecules** (*mols*)

Adds molecules to to molecules table.

**Parameters** **mols** (*list [rdkit.Chem.Mol]*) – List of molecules

**add\_pdbs** (*pdbs*)

Adds pdb meta data to to pdbs table.

**Parameters** **pdbs** (*Iterable [Dict]*) – List of pdb meta data

**by\_pdb\_code** (*pdb\_code*)

Retrieve fragments which are part of a PDB structure.

**Parameters** **pdb\_code** (*str*) – PDB code

**Returns** List of fragments

**Return type** List[Fragment]

**Raises** `LookupError` – When `pdb_code` could not be found

**create\_tables** ()

Create tables if they don't exist

**id2label** ()

Lookup table of fragments from an number to a label.

**Returns** `SqliteDict`

**is\_ligand\_stored** (*pdb\_code*, *het\_code*)

Check whether ligand is already in database

#### Parameters

- **pdb\_code** (*str*) – Protein databank identifier
- **het\_code** (*str*) – Ligand/hetero identifier

**Returns** `bool`

**label2id** ()

Lookup table of fragments from an label to a number.

**Returns** `SqliteDict`

**class** `kripodb.db.IntbitsetDict` (*db*, *number\_of\_bits=None*)  
Dictionary of BitMaps with sqlite3 backend.

**Parameters**

- **db** (`FingerprintsDb`) – Fingerprints db
- **number\_of\_bits** (*int*) – Number of bits

**number\_of\_bits**

*int* – Number of bits the bitsets consist of

**update** (*[E]*, *\*\*F*) → None. Update D from mapping/iterable E and F.

If E present and has a `.keys()` method, does: for k in E: D[k] = E[k] If E present and lacks `.keys()` method, does: for (k, v) in E: D[k] = v In either case, this is followed by: for k, v in F.items(): D[k] = v

**class** `kripodb.db.SqliteDb` (*filename*)  
Wrapper around a sqlite database connection

Database is created if it does not exist.

**Parameters** **filename** (*str*) – Sqlite filename

**connection**

`sqlite3.Connection` – Sqlite connection

**cursor**

`sqlite3.Cursor` – Sqlite cursor

**close** ()

Close database

**commit** ()

Commit pending changes

**create\_tables** ()

Abstract method which is called after connecting to database so tables can be created.

Use `CREATE TABLE IF NOT EXISTS ...` in method to prevent duplicate create errors.

**class** `kripodb.db.SqliteDict` (*connection*, *table\_name*, *key\_column*, *value\_column*)  
Dict-like object of 2 columns of a sqlite table.

Can be used to query and alter the table.

**Parameters**

- **connection** (`sqlite3.Connection`) – Sqlite connection
- **table\_name** (*str*) – Table name
- **key\_column** (*str*) – Column name used as key
- **value\_column** (*str*) – Column name used as value

**connection**

`sqlite3.Connection` – Sqlite connection

**cursor**

`sqlite3.Cursor` – Sqlite cursor

**items** () → list of D's (key, value) pairs, as 2-tuples

**iteritems** () → an iterator over the (key, value) items of D

**iteritems\_startswith** (*prefix*)  
item iterator over keys with prefix

**Parameters** **prefix** (*str*) – Prefix of key

### Examples

All items with key starting with letter ‘a’ are returned.

```
>>> for frag_id, fragment in fragments.iteritems_startswith('a'):
    # do something with frag_id and fragment
```

**Returns** List[Tuple[key, value]]

**intervalues** () → an iterator over the values of D

**materialize** ()

Fetches all kev/value pairs from the sqlite database.

Useful when dictionary is iterated multiple times and the cost of fetching is to high.

**Returns** Dictionary with all kev/value pairs

**Return type** Dict

**values** () → list of D’s values

kripodb.db.**adapt\_BitMap** (*ibs*)

Convert BitMap to it’s serialized format

**Parameters** **ibs** (*BitMap*) – bitset

### Examples

Serialize BitMap

```
>>> adapt_BitMap(BitMap([1, 2, 3, 4]))
'xc@\{ð\'
```

**Returns** serialized BitMap

**Return type** str

kripodb.db.**adapt\_molblockgz** (*mol*)

Convert RDKit molecule to compressed molblock

**Parameters** **mol** (*rdkit.Chem.Mol*) – molecule

**Returns** Compressed molblock

**Return type** str

kripodb.db.**convert\_BitMap** (*s*)

Convert serialized BitMap to BitMap

**Parameters** **s** (*str*) – serialized BitMap

## Examples

Deserialize BitMap

```
>>> ibs = convert_BitMap('xc@{\delta}')
BitMap([1, 2, 3, 4])
```

**Returns** bitset

**Return type** BitMap

`kripodb.db.convert_molblockgz` (*molgz*)

Convert compressed molblock to RDKit molecule

**Parameters** *molgz* – (str) zlib compressed molblock

**Returns** molecule

**Return type** rdkit.Chem.Mol

## 3.3 kripodb.dive

`kripodb.dive.dense_dump` (*inputfile, outputfile, frag1only*)

Dump dense matrix with zeros included

**Parameters**

- **inputfile** (*str*) – Filename of dense similarity matrix
- **outputfile** (*file*) – Writeable file object
- **frag1only** (*bool*) – Only dump frag1 fragments

Returns:

`kripodb.dive.dense_dump_iter` (*matrix, frag1only*)

Iterate dense matrix with zeros

**Parameters**

- **matrix** (`FrozenSimilarityMatrix`) – Dense similarity matrix
- **frag1only** (*bool*) – True to iterate over \*frag1 only

**Yields** (*str, str, float*) – Fragment label pair and score

`kripodb.dive.dive_export` (*fragmentsdb, uniprot\_annot, pdbtags, propnames, props*)

Writes metadata props for DiVE visualization

**Parameters**

- **fragmentsdb** (*str*) – Filename fo fragments db file
- **uniprot\_annot** (*file*) – Readable file object with uniprot gene and family mapping as tsv
- **pdbtags** (*list*) – List of readable file objects to tag pdb by filename
- **propnames** (*file*) – Writable file object to write prop names to
- **props** (*file*) – Writeable file object to write props to

`kripodb.dive.dive_sphere` (*inputfile*, *outputfile*, *onlyfrag1*)  
Export fragments as DiVE formatted sphere

#### Parameters

- **inputfile** (*str*) – fragments db input file
- **outputfile** (*file*) – fragments dive output file
- **onlyfrag1** (*bool*) – Only \*\_frag1

## 3.4 kripodb.frozen

Similarity matrix using pytables carray

**class** `kripodb.frozen.FrozenSimilarityMatrix` (*filename*, *mode*='r', *\*\*kwargs*)  
Frozen similarities matrix

Can retrieve whole column of a specific row fairly quickly. Store as compressed dense matrix. Due to compression the zeros use up little space.

Warning! Can not be enlarged.

Compared find performance `FrozenSimilarityMatrix` with `SimilarityMatrix`:

```
>>> from kripodb.db import FragmentsDb
>>> db = FragmentsDb('data/feb2016/Kripo20151223.sqlite')
>>> ids = [v[0] for v in db.cursor.execute('SELECT frag_id FROM fragments ORDER_
↳BY RANDOM() LIMIT 20')]
>>> from kripodb.frozen import FrozenSimilarityMatrix
>>> fdm = FrozenSimilarityMatrix('01-01_to_13-13.out.frozen.bloszlib.h5')
>>> from kripodb.hdf5 import SimilarityMatrix
>>> dm = SimilarityMatrix('data/feb2016/01-01_to_13-13.out.h5', cache_labels=True)
>>> %timeit list(dm.find(ids[0], 0.45, None))
```

```
... 1 loop, best of 3: 1.96 s per loop >>> %timeit list(fdm.find(ids[0], 0.45, None)) ... The slowest run took
6.21 times longer than the fastest. This could mean that an intermediate result is being cached. ... 10 loops, best
of 3: 19.3 ms per loop >>> ids = [v[0] for v in db.cursor.execute('SELECT frag_id FROM fragments ORDER
BY RANDOM() LIMIT 20')] >>> %timeit -n1 [list(fdm.find(v, 0.45, None)) for v in ids] ... 1 loop, best of 3:
677 ms per loop >>> %timeit -n1 [list(dm.find(v, 0.45, None)) for v in ids] ... 1 loop, best of 3: 29.7 s per loop
```

#### Parameters

- **filename** (*str*) – File name of hdf5 file to write or read similarity matrix from
- **mode** (*str*) – Can be 'r' for reading or 'w' for writing
- **\*\*kwargs** – Passed though to `tables.open_file()`

#### h5file

*tables.File* – Object representing an open hdf5 file

#### scores

*tables.CArray* – HDF5 Table that contains matrix

#### labels

*tables.CArray* – Table to look up label of fragment by id or id of fragment by label

#### close ()

Closes the hdf5file

**count** (*frame\_size=None, raw\_score=False, lower\_triangle=False*)

Count occurrences of each score

Only scores are counted of the upper triangle or lower triangle. Zero scores are skipped.

**Parameters**

- **frame\_size** (*int*) – Dummy argument to force same interface for thawed and frozen matrix
- **raw\_score** (*bool*) – When true return raw int16 score else fraction score
- **lower\_triangle** (*bool*) – When true return scores from lower triangle else return scores from upper triangle

**Returns** Score and number of occurrences

**Return type** Tuple[(str, int)]

**find** (*query, cutoff, limit=None*)

Find similar fragments to query.

**Parameters**

- **query** (*str*) – Query fragment identifier
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **limit** (*int*) – Maximum number of hits. Default is None for no limit.

**Returns** Hit fragment identifier and similarity score

**Return type** list[tuple[str,float]]

**from\_array** (*data, labels*)

Fill matrix from 2 dimensional array

**Parameters**

- **data** (*np.array*) – 2 dimensional square array with scores
- **labels** (*list*) – List of labels for each column and row index

**from\_pairs** (*similarity\_matrix, frame\_size, limit=None, single\_sided=False*)

Fills self with matrix which is stored in pairs.

Also known as COOrdinate format, the ‘ijv’ or ‘triplet’ format.

**Parameters**

- **similarity\_matrix** (*kripodb.hdf5.SimilarityMatrix*) –
- **frame\_size** (*int*) – Number of pairs to append in a single go
- **limit** (*int/None*) – Number of pairs to add, None for no limit, default is None.
- **single\_sided** (*bool*) – If false add stored direction and reverse direction. Default is False.

```
time kripodb similarities freeze -limit 200000 -f 100000 data/feb2016/01-01_to_13-13.out.h5 percell.h5
47.2s time kripodb similarities freeze -limit 200000 -f 100000 data/feb2016/01-01_to_13-13.out.h5
coo.h5 0.2m - 2m6s .4m - 2m19s .8m - 2m33s 1.6m - 2m48s 3.2m - 3m4s 6.4m - 3m50s 12.8m - 4m59s
25.6m - 7m27s
```

**to\_pairs** (*pairs*)

Copies labels and scores from self to pairs matrix.

**Parameters** **pairs** (*SimilarityMatrix*) –

**to\_pandas** ()

Pandas dataframe with labelled columns and rows.

Warning! Only use on matrices that fit in memory

**Returns** pd.DataFrame

## 3.5 kripodb.hdf5

Similarity matrix using hdf5 as storage backend.

**class** kripodb.hdf5.**AbstractSimpleTable** (*table*, *append\_chunk\_size=100000000*)

Abstract wrapper around a HDF5 table

**Parameters**

- **table** (*tables.Table*) – HDF5 table
- **append\_chunk\_size** (*int*) – Size of chunk to append in one go. Defaults to 1e8, which when table description is 10bytes will require 2Gb during append.

**Attributes** *table* (*tables.Table*): HDF5 table *append\_chunk\_size* (*int*): Number of rows to read from other table during append.

**append** (*other*)

Append rows of other table to self

**Parameters** *other* – Table of same type as self

**class** kripodb.hdf5.**LabelsLookup** (*h5file*, *expectedrows=0*)

Table to look up label of fragment by id or id of fragment by label

When table does not exist in h5file it is created.

**Parameters**

- **h5file** (*tables.File*) – Object representing an open hdf5 file
- **expectedrows** (*int*) – Expected number of pairs to be added. Required when similarity matrix is opened in write mode, helps optimize storage

**by\_id** (*frag\_id*)

Look up label of fragment by id

**Parameters** *frag\_id* (*int*) – Fragment identifier

**Raises** `IndexError` – When id of fragment is not found

**Returns** Label of fragment

**Return type** `str`

**by\_label** (*label*)

Look up id of fragment by label

**Parameters** *label* (*str*) – Fragment label

**Raises** `IndexError` – When label of fragment is not found

**Returns** Fragment identifier

**Return type** `int`

**by\_labels** (*labels*)

Look up ids of fragments by label

**Parameters** **labels** (*set [str]*) – Set of fragment labels

**Raises** `IndexError` – When label of fragment is not found

**Returns** Set of fragment identifiers

**Return type** `set[int]`

**keep** (*other, keep*)

Copy content of self to other and only keep given fragment identifiers

**Parameters**

- **other** (`LabelsLookup`) – Labels table to fill
- **keep** (*set [int]*) – Fragment identifiers to keep

**label2ids** ()

Return whole table as a dictionary

**Returns** Dictionary with label as key and frag\_id as value.

**Return type** `dict`

**merge** (*label2id*)

Merge label2id dict into self

When label does not exist an id is generated and the label/id is added. When label does exist the id of the label in self is kept.

**Parameters** **label2id** (*dict*) – Dictionary with fragment label as key and fragment identifier as value.

**Returns** Dictionary of label/id which were in label2id, but missing in self

**Return type** `dict`

**skip** (*other, skip*)

Copy content of self to other and skip given fragment identifiers

**Parameters**

- **other** (`LabelsLookup`) – Labels table to fill
- **skip** (*set [int]*) – Fragment identifiers to skip

**update** (*label2id*)

Update labels lookup by adding labels in label2id.

**Parameters** **label2id** (*dict*) – Dictionary with fragment label as key and fragment identifier as value.

**class** `kripodb.hdf5.PairsTable` (*h5file, expectedrows=0*)

Table to store similarity score of a pair of fragment fingerprints

When table does not exist in h5file it is created.

**Parameters**

- **h5file** (*tables.File*) – Object representing an open hdf5 file
- **expectedrows** (*int*) – Expected number of pairs to be added. Required when similarity matrix is opened in write mode, helps optimize storage



**score\_precision**

*int* – Similarity score is a fraction, the score is converted to an int by multiplying it with the precision

**full\_matrix**

*bool* – Matrix is filled above and below diagonal.

**append** (*other*)

Append rows of other table to self

**Parameters** *other* – Table of same type as self

**count** (*frame\_size*, *raw\_score=False*)

Count occurrences of each score

**Parameters**

- **frame\_size** (*int*) – Size of matrix loaded each time. Larger requires more memory and smaller is slower.
- **raw\_score** (*bool*) – Return raw int16 score or fraction score

**Returns** Score and number of occurrences

**Return type** Tuple[(str, int)]

**find** (*frag\_id*, *cutoff*, *limit*)

Find fragment hits which has a similarity score with frag\_id above cutoff.

**Parameters**

- **frag\_id** (*int*) – query fragment identifier
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **limit** (*int*) – Maximum number of hits. Default is None for no limit.

**Returns** Where first tuple value is hit fragment identifier and second value is similarity score

**Return type** List[Tuple]

**keep** (*other*, *keep*)

Copy pairs from self to other and keep given fragment identifiers and the identifiers they pair with.

**Parameters**

- **other** (*PairsTable*) – Pairs table to fill
- **keep** (*set [int]*) – Fragment identifiers to keep

**Returns** Fragment identifiers that have been copied to other

**Return type** set[int]

**skip** (*other*, *skip*)

Copy content from self to other and skip given fragment identifiers

**Parameters**

- **other** (*PairsTable*) – Pairs table to fill
- **skip** (*set [int]*) – Fragment identifiers to skip

**update** (*similarities\_iter*, *label2id*)

Store pairs of fragment identifier with their similarity score

**Parameters**

- **similarities\_iter** (*Iterator*) – Iterator which yields (label1, label2, similarity\_score)
- **label2id** (*Dict*) – Lookup with fragment label as key and fragment identifier as value

**class** kripodb.hdf5.**SimilarityMatrix** (*filename, mode='r', expectedpairrows=None, expected-labelrows=None, cache\_labels=False, \*\*kwargs*)

Similarity matrix

#### Parameters

- **filename** (*str*) – File name of hdf5 file to write or read similarity matrix from
- **mode** (*str*) – Can be 'r' for reading or 'w' for writing
- **expectedpairrows** (*int*) – Expected number of pairs to be added. Required when similarity matrix is opened in write mode, helps optimize storage
- **expectedlabelrows** (*int*) – Expected number of labels to be added. Required when similarity matrix is opened in write mode, helps optimize storage
- **cache\_labels** (*bool*) – Cache labels, speed up label lookups

#### h5file

*tables.File* – Object representing an open hdf5 file

#### pairs

*PairsTable* – HDF5 Table that contains pairs

#### labels

*LabelsLookup* – Table to look up label of fragment by id or id of fragment by label

#### append (*other*)

Append data from other similarity matrix to me

**Parameters** **other** (*SimilarityMatrix*) – Other similarity matrix

#### close ()

Closes the hdf5file

#### count (*frame\_size, raw\_score=False, lower\_triangle=False*)

Count occurrences of each score

#### Parameters

- **frame\_size** (*int*) – Size of matrix loaded each time. Larger requires more memory and smaller is slower.
- **raw\_score** (*bool*) – Return raw int16 score or fraction score
- **lower\_triangle** (*bool*) – Dummy argument to force same interface for thawed and frozen matrix

**Returns** Score and number of occurrences

**Return type** (*str, int*)

#### find (*query, cutoff, limit=None*)

Find similar fragments to query.

#### Parameters

- **query** (*str*) – Query fragment identifier
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **limit** (*int*) – Maximum number of hits. Default is None for no limit.

**Yields** (*str, float*) – Hit fragment identifier and similarity score

**keep** (*other, keep*)

Copy content of self to other and only keep given fragment labels and the labels they pair with

**Parameters**

- **other** (*SimilarityMatrix*) – Writable matrix to fill
- **keep** (*set [str]*) – Fragment labels to keep

**skip** (*other, skip*)

Copy content of self to other and skip all given fragment labels

**Parameters**

- **other** (*SimilarityMatrix*) – Writable matrix to fill
- **skip** (*set [str]*) – Fragment labels to skip

**update** (*similarities\_iter, label2id*)

Store pairs of fragment identifier with their similarity score and label 2 id lookup

**Parameters**

- **similarities\_iter** (*iterator*) – Iterator which yields (label1, label2, similarity\_score)
- **label2id** (*dict*) – Dictionary with fragment label as key and fragment identifier as value.

## 3.6 kripodb.makebits

Module to read/write fingerprints in Makebits file format

`kripodb.makebits.iter_file` (*infile*)

Reads Makebits formatted file Yields header first then tuples of identifier and BitMap object

**Yields** first header (format name, format version, number of bits, description), then tuples of the fingerprint identifier and an BitMap object

**Parameters** `infile` (*File*) – File object of Makebits formatted file to read

### Examples

Read a file

```
>>> f = iter_file(open('fingerprints01.fp'))
>>> read_fp_size(next(f))
4
>>> {frag_id: fp for frag_id, fp in f}
{'id1': BitMap([1, 2, 3, 4])}
```

`kripodb.makebits.write_file` (*fp\_size, bitsets, fn*)

Write makebits formatted file

**Parameters**

- **fp\_size** (*int*) – Number of bits
- **bitsets** (*dict*) – Dict with fingerprint identifier as key and BitMap object as value

- **fn** (*File*) – File object to write to

## Examples

Write a file

```
>>> write_file(4, {'id1': BitMap([1, 2, 3, 4])}, open('fingerprints01.fp', 'w'))
```

## 3.7 kripodb.modifiedtanimoto

Module to calculate modified tanimoto similarity

`kripodb.modifiedtanimoto.calc_mean_onbit_density` (*bitsets, number\_of\_bits*)

Calculate the mean density of bits that are on in bitsets collection.

### Parameters

- **bitsets** (*list* [*pyroaring.BitMap*]) – List of fingerprints
- **number\_of\_bits** – Number of bits for all fingerprints

**Returns** Mean on bit density

**Return type** *float*

`kripodb.modifiedtanimoto.corrections` (*mean\_onbit\_density*)

Calculate corrections

See *similarity()* for explanation of corrections.

**Parameters** **mean\_onbit\_density** (*float*) – Mean on bit density

**Returns**  $S_T$  correction,  $S_{T0}$  correction

**Return type** *float*

`kripodb.modifiedtanimoto.similarities` (*bitsets1, bitsets2, number\_of\_bits, corr\_st, corr\_sto, cutoff, ignore\_upper\_triangle=False*)

Calculate modified tanimoto similarity between two collections of fingerprints

Excludes similarity of the same fingerprint.

### Parameters

- **bitsets1** (*Dict*{*str*, *pyroaring.BitMap*}) – First dict of fingerprints with fingerprint label as key and *pyroaring.BitMap* as value
- **bitsets2** (*Dict*{*str*, *pyroaring.BitMap*}) – Second dict of fingerprints with fingerprint label as key and *pyroaring.BitMap* as value
- **number\_of\_bits** (*int*) – Number of bits for all fingerprints
- **corr\_st** (*float*) –  $S_T$  correction
- **corr\_sto** (*float*) –  $S_{T0}$  correction
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **ignore\_upper\_triangle** (*Optional* [*bool*]) – When true returns similarity where  $label1 > label2$ , when false returns all similarities

**Yields** (fingerprint label 1, fingerprint label2, similarity score)

`kripodb.modifiedtanimoto.similarity` (*bitset1, bitset2, number\_of\_bits, corr\_st, corr\_sto*)

Calculate modified Tanimoto similarity between two fingerprints

Given two fingerprints of length  $n$  with  $a$  and  $b$  bits set in each fingerprint, respectively, and  $c$  bits set in both fingerprint, selected from a data set of fingerprint with a mean bit density of  $\rho_0$ , the modified Tanimoto similarity  $S_{MT}$  is calculated as

$$S_{MT} = \left(\frac{2-\rho_0}{3}\right)S_T + \left(\frac{1+\rho_0}{3}\right)S_{T0}$$

where  $S_T$  is the standard Tanimoto coefficient

$$S_T = \frac{c}{a + b - c}$$

and  $S_{T0}$  is the inverted Tanimoto coefficient

$$S_{T0} = \frac{n - a - b + c}{n - c}$$

#### Parameters

- **bitset1** (*pyroaring.BitMap*) – First fingerprint
- **bitset2** (*pyroaring.BitMap*) – Second fingerprint
- **number\_of\_bits** (*int*) – Number of bits for all fingerprints
- **corr\_st** (*float*) –  $S_T$  correction
- **corr\_sto** (*float*) –  $S_{T0}$  correction

**Returns** modified Tanimoto similarity

**Return type** `float`

## 3.8 kripodb.pairs

Module handling generation and retrieval of similarity of fingerprint pairs

`kripodb.pairs.dump_pairs` (*bitsets1, bitsets2, out\_format, out\_file, out, number\_of\_bits, mean\_onbit\_density, cutoff, label2id, nomemory, ignore\_upper\_triangle=False*)

Dump pairs of bitset collection.

A pairs are rows of the bitset identifier of both bitsets with a similarity score.

#### Parameters

- **bitsets1** (*Dict{str, pyroaring.BitMap}*) – First dict of fingerprints with fingerprint label as key and `pyroaring.BitMap` as value
- **bitsets2** (*Dict{str, pyroaring.BitMap}*) – Second dict of fingerprints with fingerprint label as key and `pyroaring.BitMap` as value
- **out\_format** – ‘tsv’ or ‘hdf5’
- **out\_file** – Filename of output file where ‘hdf5’ format is written to.
- **out** (*File*) – File object where ‘tsv’ format is written to.
- **number\_of\_bits** (*int*) – Number of bits for all bitsets
- **mean\_onbit\_density** (*float*) – Mean on bit density

- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **label2id** – dict to translate label to id (string to int)
- **nomemory** – If true bitset2 is not loaded into memory
- **ignore\_upper\_triangle** – When true returns similarity where label1 > label2, when false returns all similarities

`kripodb.pairs.dump_pairs_hdf5` (*similarities\_iter, label2id, expectedrows, out\_file*)

Dump pairs in hdf5 file

Pro: \* very small, 10 bytes for each pair + compression Con: \* requires hdf5 library to access

#### Parameters

- **similarities\_iter** (*Iterator*) – Iterator with tuple with fingerprint 1 label, fingerprint 2 label, similarity as members
- **label2id** (*dict*) – dict to translate label to id (string to int)
- **expectedrows** –
- **out\_file** –

`kripodb.pairs.dump_pairs_tsv` (*similarities\_iter, out*)

Dump pairs in tab delimited file

Pro: \* when stored in sqlite can be used outside of Python Con: \* big, unless output is compressed

#### Parameters

- **similarities\_iter** (*Iterator*) – Iterator with tuple with fingerprint 1 label, fingerprint 2 label, similarity as members
- **out** (*File*) – Writeable file

`kripodb.pairs.merge` (*ins, out*)

Concatenate similarity matrix files into a single one.

#### Parameters

- **ins** (*list[str]*) – List of input similarity matrix filenames
- **out** (*str*) – Output similarity matrix filenames

**Raises** `AssertionError` – When nr of labels of input files is not the same

`kripodb.pairs.open_similarity_matrix` (*fn*)

Open read-only similarity matrix file.

**Parameters** **fn** (*str*) – Filename of similarity matrix

**Returns** A read-only similarity matrix object

**Return type** `SimilarityMatrix` | `FrozenSimilarityMatrix`

`kripodb.pairs.similar` (*query, similarity\_matrix, cutoff, limit=None*)

Find similar fragments to query based on similarity matrix.

#### Parameters

- **query** (*str*) – Query fragment identifier
- **similarity\_matrix** (*kripodb.db.SimilarityMatrix*) – Similarity matrix
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **limit** (*int*) – Maximum number of hits. Default is None for no limit.

**Yields** *Tuple[str, str, float]* – List of (query fragment identifier, hit fragment identifier, similarity score) sorted on similarity score

`kripodb.pairs.similar_run(query, pairsdbfn, cutoff, out)`

Find similar fragments to query based on similarity matrix and write to tab delimited file.

**Parameters**

- **query** (*str*) – Query fragment identifier
- **pairsdbfn** (*str*) – Filename of similarity matrix file or url of kripodb webservice
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **out** (*File*) – File object to write output to

`kripodb.pairs.similarity2query(bitsets2, query, out, mean_onbit_density, cutoff, memory)`

Calculate similarity of query against all fingerprints in bitsets2 and write to tab delimited file.

**Parameters**

- **bitsets2** (`kripodb.db.IntbitsetDict`) –
- **query** (*str*) – Query identifier or beginning of it
- **out** (*File*) – File object to write output to
- **mean\_onbit\_density** (*float*) – Mean on bit density
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **memory** (*Optional[bool]*) – When true will load bitset2 into memory, when false it doesn't

`kripodb.pairs.total_number_of_pairs(fingerprint_filenames)`

Count number of pairs in similarity matrix files

**Parameters** **fingerprint\_filenames** (*list[str]*) – List of file names of similarity matrices

**Returns** Total number of pairs

**Return type** `int`

## 3.9 kripodb.pharmacophores

`kripodb.pharmacophores.FEATURE_TYPES = [{'color': 'ff33cc', 'element': 'He', 'key': 'LI`

*Types of pharmacophore feature types. List of dictionaries with the following keys –*

- **key**, short identifier of type
- **label**, human readable label
- **color**, hex rrggbb color
- **element**, Element used in kripo pharmacophore sfile for this type

**class** `kripodb.pharmacophores.PharmacophorePointsTable(h5file, expectedrows=0)`

Wrapper around pytables table to store pharmacophore points

**Parameters**

- **h5file** (*tables.File*) – Pytables hdf5 file object which contains the pharmacophores table

- **expectedrows** (*int*) – Expected number of pharmacophores. Required when hdf5 file is created, helps optimize compression

Pharmacophore points of a fragment can be retrieved using:

```
points = table['frag_id1']
```

*points* is a list of points, each point is a tuple with following columns feature type key, x, y and z coordinate. The feature type key is defined in FEATURE\_TYPES.

Number of pharmacophore points can be requested using:

```
nr_points = len(table)
```

To check whether fragment identifier is contained use:

```
'frag_id1' in table
```

**add\_dir** (*startdir*)

Find \*\_pphore.sd.gz \*\_pphores.txt file pairs recursively in start directory and add them.

**Parameters** **startdir** (*str*) – Path to a start directory

**read\_phar** (*infile*)

Read phar formatted file and add pharmacophore to self

**Parameters** **infile** – File object of phar formatted file

**class** kripodb.pharmacophores.**PharmacophoresDb** (*filename, mode='r', expectedrows=0, \*\*kwargs*)

Database for pharmacophores of fragments aka sub-pockets.

**Parameters**

- **filename** (*str*) – File name of hdf5 file to write or read pharmacophores to/from
- **mode** (*str*) – Can be 'r' for reading or 'w' for writing or 'a' for appending
- **expectedrows** (*int*) – Expected number of pharmacophores. Required when hdf5 file is created, helps optimize compression
- **\*\*kwargs** – Passed to tables.open\_file

Pharmacophore points of a fragment can be retrieved using:

```
points = db['frag_id1']
```

*points* is a list of points, each point is a tuple with following columns feature type key, x, y and z coordinate. The feature type key is defined in FEATURE\_TYPES.

**h5file**

*tables.File* – Object representing an open hdf5 file

**points**

*PharmacophorePointsTable* – HDF5 table that contains pharmacophore points

**add\_dir** (*startdir*)

Find \*\_pphore.sd.gz \*\_pphores.txt file pairs recursively in start directory and add them.

**Parameters** **startdir** (*str*) – Path to a start directory

**append** (*other*)

Append pharmacophores in other db to self



**Parameters other** (`PharmacophoresDb`) – The other pharmacophores database

**close()**

Closes the hdf5file

Instead of calling close() explicitly, use context manager:

```
with PharmacophoresDb('data/pharmacophores.h5') as db:
    points = db['frag_id1']
```

**read\_phar** (*infile*)

Read phar formatted file and add pharmacophore to self

**Parameters infile** – File object of phar formatted file

**write\_phar** (*outfile*, *frag\_id=None*)

Write pharmacophore of frag\_id as phar format to outfile

**Parameters**

- **outfile** (*file*) – File object to write to
- **frag\_id** (*str*) – Fragment identifier, if None all pharmacophores are written

`kripodb.pharmacophores.as_phar` (*frag\_id*, *points*)

Return pharmacophore in \*.phar format.

See [align-it](#) for format description.

**Parameters**

- **frag\_id** (*str*) – Fragment identifier
- **points** (*list*) – List of points where each point is (key,x,y,z)

**Returns** Pharmacophore is \*.phar format

**Return type** `str`

`kripodb.pharmacophores.read_fragtxtfile` (*fragtxtfile*)

Read a fragment text file

**Parameters fragtxtfile** – Filename of fragment text file

**Returns** Dictionary where key is fragment identifier and value is a list of pharmacophore point indexes.

**Return type** `dict`

`kripodb.pharmacophores.read_fragtxtfile_as_file` (*fileobject*)

Read a fragment text file object which contains the pharmacophore point indexes for each fragment identifier.

File format is a fragment on each line, the line is space separated with fragment\_identifier followed by the pharmacophore point indexes.

**Parameters fileobject** (*file*) – File object to read

**Returns** Dictionary where key is fragment identifier and value is a list of pharmacophore point indexes.

**Return type** `dict`

`kripodb.pharmacophores.read_pphore_gzipped_sdf` (*sdf*)

Read a gzipped sdf which contains pharmacophore points as atoms

**Parameters sdf** (*string*) – Path to filename

**Returns** List of Pharmacophore points

**Return type** `list`

`kripodb.pharmacophores.read_pphore_sdf`(*sdf*)

Read a sdf file which contains pharmacophore points as atoms

**Parameters** `sdf` (*file*) – File object with sdf file contents

**Returns** List of pharmacophore points

**Return type** `list`

## 3.10 kripodb.pdb

**class** `kripodb.pdb.PdbReport` (*pdbs=None, fields=None*)

Client for the Custom Report Web Services of the RCSB PDB website

See <http://www.rcsb.org/pdb/software/wsreport.do> for more information.

**Parameters**

- **pdbs** (*List[str]*) – List of pdb identifiers to fetch. Default is ['\*'] which fetches all.
- **fields** – (*List[str]*): List of fields to fetch. Default is ['structureTitle', 'compound', 'ecNo', 'uniprotAcc', 'uniprotRecommendedName'] See <http://www.rcsb.org/pdb/results/reportField.do> for possible fields.

**url**

*str* – Url of report, based on pdbs and fields.

**fetch** ()

Fetch report from PDB website

**Yields** *dict* – Dictionary with keys same as ['structureId', 'chainID'] + self.fields

`kripodb.pdb.parse_csv_file` (*thefile*)

Parse csv file, yielding rows as dictionary.

The csv file should have an header.

**Parameters** `thefile` (*file*) – File like object

**Yields** *dict* – Dictionary with column header name as key and cell as value

## 3.11 kripodb.script

`kripodb.script.main` (*argv=['-b', 'latex', '-D', 'language=en', '-d', '\_build/doctrees', '.', '\_build/latex']*)

Main script function.

Calls run method of selected sub commandos.

**Parameters** `argv` (*list[str]*) – List of command line arguments

`kripodb.script.make_parser` ()

Creates a parser with sub commands

**Returns** parser with sub commands

**Return type** `argparse.ArgumentParser`

`kripodb.script.fragments.make_fragments_parser(subparsers)`

Creates a parser for fragments sub commands

**Parameters** `subparsers` (*argparse.ArgumentParser*) – Parser to which to add sub commands to

`kripodb.script.fingerprints.make_fingerprints_parser(subparsers)`

Creates a parser for fingerprints sub commands

**Parameters** `subparsers` (*argparse.ArgumentParser*) – Parser to which to add sub commands to

`kripodb.script.similarities.make_similarities_parser(subparsers)`

Creates a parser for similarities sub commands

**Parameters** `subparsers` (*argparse.ArgumentParser*) – Parser to which to add sub commands to

`kripodb.script.similarities.read_fpneighpairs_file(inputfile, ignore_upper_triangle=False)`

Read fpneigh formatted similarity matrix file.

**Parameters**

- `inputfile` (*File*) – File object to read
- `ignore_upper_triangle` (*bool*) – Ignore upper triangle of input

**Yields** *Tuple((Str,Str,Float))* – List of (query fragment identifier, hit fragment identifier, similarity score)

`kripodb.script.similarities.simmatrix_export_run(simmatrixfn, outputfile, no_header, frag1, pdb)`

Export similarity matrix to tab delimited file

**Parameters**

- `simmatrixfn` (*str*) – (Compact) hdf5 similarity matrix filename
- `outputfile` (*File*) – Tab delimited output file
- `no_header` (*bool*) – Output no header
- `frag1` (*bool*) – Only output \*frag1
- `pdb` (*str*) – Filename with pdb codes inside

`kripodb.script.dive.dense_dump_sc(sc)`

Dump dense matrix with zeros

## 3.12 kripodb.webservice

Module for Client for kripo web service

**exception** `kripodb.webservice.client.Incomplete` (*message, absent\_identifiers*)

**exception** `kripodb.webservice.client.IncompleteFragments` (*absent\_identifiers, fragments*)

**exception** `kripodb.webservice.client.IncompletePharmacophores` (*absent\_identifiers, pharmacophores*)

**class** `kripodb.webservice.client.WebserviceClient` (*base\_url*)

Client for kripo web service

## Example

```
>>> client = WebserviceClient('http://localhost:8084/kripo')
>>> client.similar_fragments('3j7u_NDP_frag24', 0.85)
[{'query_frag_id': '3j7u_NDP_frag24', 'hit_frag_id': '3j7u_NDP_frag23', 'score': 0.8991}]
```

**Parameters** `base_url` (*str*) – Base url of web service. e.g. `http://localhost:8084/kripo`

**fragments\_by\_id** (*fragment\_ids*, *chunk\_size=100*)

Retrieve fragments by their identifier

### Parameters

- **fragment\_ids** (*List[str]*) – List of fragment identifiers
- **chunk\_size** (*int*) – Number of fragment to retrieve in a single http request

**Returns** List of fragment information

**Return type** `list[dict]`

**Raises** `IncompleteFragments` – When one or more of the identifiers could not be found.

**fragments\_by\_pdb\_codes** (*pdb\_codes*, *chunk\_size=450*)

Retrieve fragments by their PDB code

### Parameters

- **pdb\_codes** (*List[str]*) – List of PDB codes
- **chunk\_size** (*int*) – Number of PDB codes to retrieve in a single http request

**Returns** List of fragment information

**Return type** `list[dict]`

**Raises** `requests.HTTPError` – When one of the PDB codes could not be found.

**similar\_fragments** (*fragment\_id*, *cutoff*, *limit=1000*)

Find similar fragments to query.

### Parameters

- **fragment\_id** (*str*) – Query fragment identifier
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **limit** (*int*) – Maximum number of hits. Default is None for no limit.

**Returns** Query fragment identifier, hit fragment identifier and similarity score

**Return type** `list[dict]`

**Raises** `request.HTTPError` – When `fragment_id` could not be found

Kripo datafiles wrapped in a webservice

```
class kripodb.webservice.server.KripodbJSONEncoder (skipkeys=False, ensure_ascii=True, check_circular=True, allow_nan=True, sort_keys=False, indent=None, separators=None, encoding='utf-8', default=None)
```

JSON encoder for KripoDB object types

Copied from <http://flask.pocoo.org/snippets/119/>

**default** (*obj*)

Implement this method in a subclass such that it returns a serializable object for *o*, or calls the base implementation (to raise a `TypeError`).

For example, to support arbitrary iterators, you could implement `default` like this:

```
def default(self, o):
    try:
        iterable = iter(o)
    except TypeError:
        pass
    else:
        return list(iterable)
    return JSONEncoder.default(self, o)
```

`kripodb.webservice.server.get_fragment_phar` (*fragment\_id*)

Pharmacophore in phar format of fragment

**Parameters** `fragment_id` (*str*) – Fragment identifier

**Returns** `PharmacophoreIproblem`

**Return type** `flask.ResponseIconnexion.lifecycle.ConnexionResponse`

`kripodb.webservice.server.get_fragment_svg` (*fragment\_id*, *width*, *height*)

2D drawing of fragment in SVG format

**Parameters**

- **fragment\_id** (*str*) – Fragment identifier
- **width** (*int*) – Width of SVG in pixels
- **height** (*int*) – Height of SVG in pixels

**Returns** `SVG documentIproblem`

**Return type** `flask.ResponseIconnexion.lifecycle.ConnexionResponse`

`kripodb.webservice.server.get_fragments` (*fragment\_ids=None*, *pdb\_codes=None*)

Retrieve fragments based on their identifier or PDB code.

**Parameters**

- **fragment\_ids** (*List[str]*) – List of fragment identifiers
- **pdb\_codes** (*List[str]*) – List of PDB codes

**Returns** List of fragment information

**Return type** `list[dict]`

**Raises** `werkzeug.exceptions.NotFound` – When one of the `fragments_ids` or `pdb_code` could not be found

`kripodb.webservice.server.get_similar_fragments` (*fragment\_id*, *cutoff*, *limit*)

Find similar fragments to query.

**Parameters**

- **fragment\_id** (*str*) – Query fragment identifier
- **cutoff** (*float*) – Cutoff, similarity scores below cutoff are discarded.
- **limit** (*int*) – Maximum number of hits. Default is `None` for no limit.

**Returns** List of dict with query fragment identifier, hit fragment identifier and similarity score

**Return type** `list[dict]`

**Raises** `werkzeug.exceptions.NotFound` – When the `fragments_id` could not be found

`kripodb.webservice.server.get_version()`

**Returns** Version of web service

**Return type** `dict[version]`

`kripodb.webservice.server.serve_app` (*similarities*, *fragments*, *pharmacophores*, *internal\_port=8084*, *external\_url='http://localhost:8084/kripo'*)

Serve webservice forever

**Parameters**

- **similarities** – Filename of similarity matrix hdf5 file
- **fragments** – Filename of fragments database file
- **pharmacophores** – Filename of pharmacophores hdf5 file
- **internal\_port** – TCP port on which to listen
- **external\_url** (*str*) – URL which should be used in Swagger spec

`kripodb.webservice.server.wsgi_app` (*similarities*, *fragments*, *pharmacophores*, *external\_url='http://localhost:8084/kripo'*)

Create wsgi app

**Parameters**

- **similarities** (*SimilarityMatrix*) – Similarity matrix to use in webservice
- **fragments** (*FragmentsDb*) – Fragment database filename
- **pharmacophores** – Filename of pharmacophores hdf5 file
- **external\_url** (*str*) – URL which should be used in Swagger spec

**Returns** `connexion.App`

---

## Command line interface

---

```
usage: kripodb [-h] [--version]
             {fingerprints,fragments,similarities,dive,serve,pharmacophores}
             ...
```

### 4.1 Positional Arguments

**subcommand**      Possible choices: fingerprints, fragments, similarities, dive, serve, pharmacophores

### 4.2 Named Arguments

**--version**      show program's version number and exit

### 4.3 Sub-commands:

#### 4.3.1 fingerprints

Fingerprints

```
kripodb fingerprints [-h]
                    {import,export,meanbitdensity,similar,similarities,merge}
                    ...
```

## Sub-commands:

### import

Add Makebits file to fingerprints db

```
kripodb fingerprints import [-h] infile [infile ...] outfile
```

### Positional Arguments

<b>infile</b>	Name of makebits formatted fingerprint file (.tar.gz or not packed or - for stdin)
<b>outfile</b>	Name of fingerprints db file Default: “fingerprints.db”

### export

Dump bitsets in fingerprints db to makebits file

```
kripodb fingerprints export [-h] infile outfile
```

### Positional Arguments

<b>infile</b>	Name of fingerprints db file Default: “fingerprints.db”
<b>outfile</b>	Name of makebits formatted fingerprint file (or - for stdout)

### meanbitdensity

Compute mean bit density of fingerprints

```
kripodb fingerprints meanbitdensity [-h] [--out OUT] fingerprintsdb
```

### Positional Arguments

<b>fingerprintsdb</b>	Name of fingerprints db file (default: “fingerprints.db”) Default: “fingerprints.db”
-----------------------	---

### Named Arguments

<b>--out</b>	Output file, default is stdout (default: -) Default: -
--------------	---



## similar

Find the fragments closests to query based on fingerprints

```
kripodb fingerprints similar [-h] [--mean_onbit_density MEAN_ONBIT_DENSITY]
                             [--cutoff CUTOFF] [--memory]
                             fingerprintsdb query out
```

### Positional Arguments

<b>fingerprintsdb</b>	Name of fingerprints db file Default: "fingerprints.db"
<b>query</b>	Query identifier or beginning of it
<b>out</b>	Output file tabdelimited (query, hit, score)

### Named Arguments

<b>--mean_onbit_density</b>	Mean on bit density (default: 0.01) Default: 0.01
<b>--cutoff</b>	Set Tanimoto cutoff (default: 0.55) Default: 0.55
<b>--memory</b>	Store bitsets in memory (default: False) Default: False

## similarities

Output formats: \* tsv, tab separated id1,id2, similarity \* hdf5, hdf5 file constructed with pytables with a, b and score, but but a and b have been replaced

by numbers and similarity has been converted to scaled int

When input has been split into chunks, use *--ignore\_upper\_triangle* flag for computing similarities between same chunk. This prevents storing pair a->b also as b->a.

```
kripodb fingerprints similarities [-h] [--out_format {tsv,hdf5}]
                                  [--fragmentsdbfn FRAGMENTSDBFN]
                                  [--mean_onbit_density MEAN_ONBIT_DENSITY]
                                  [--cutoff CUTOFF] [--nomemory]
                                  [--ignore_upper_triangle]
                                  fingerprintsfn1 fingerprintsfn2 out_file
```

### Positional Arguments

<b>fingerprintsfn1</b>	Name of reference fingerprints db file
<b>fingerprintsfn2</b>	Name of query fingerprints db file
<b>out_file</b>	Name of output file (use - for stdout)

## Named Arguments

<b>--out_format</b>	Possible choices: tsv, hdf5 Format of output (default: “hdf5”) Default: “hdf5”
<b>--fragmentsdbfn</b>	Name of fragments db file (only required for hdf5 format)
<b>--mean_onbit_density</b>	Mean on bit density (default: 0.01) Default: 0.01
<b>--cutoff</b>	Set Tanimoto cutoff (default: 0.45) Default: 0.45
<b>--nomemory</b>	Do not store query fingerprints in memory (default: False) Default: False
<b>--ignore_upper_triangle</b>	Ignore upper triangle (default: False) Default: False

## merge

Combine fingerprints databases into a single new one

```
kripodb fingerprints merge [-h] ins [ins ...] out
```

## Positional Arguments

<b>ins</b>	Input fingerprints database files
<b>out</b>	Output fingerprints database file

## 4.3.2 fragments

Fragments

```
kripodb fragments [-h] {shelve,sdf,pdb,filter,merge,export_sd} ...
```

### Sub-commands:

#### shelve

Add fragments from shelve to sqlite

```
kripodb fragments shelve [-h] [--skipdups] shelvefn fragmentsdb
```

## Positional Arguments

**shelvefn**

**fragmentsdb** Name of fragments db file (default: "fragments.db")  
Default: "fragments.db"

## Named Arguments

**--skipdups** Skip duplicates, instead of dieing one first duplicate  
Default: False

## sdf

Add fragments sdf to sqlite

```
kripodb fragments sdf [-h] sdffns [sdffns ...] fragmentsdb
```

## Positional Arguments

**sdffns** SDF filename

**fragmentsdb** Name of fragments db file (default: "fragments.db")  
Default: "fragments.db"

## pdb

Add pdb metadata from RCSB PDB website to fragment sqlite db

```
kripodb fragments pdb [-h] fragmentsdb
```

## Positional Arguments

**fragmentsdb** Name of fragments db file (default: "fragments.db")  
Default: "fragments.db"

## filter

Filter fragments database

```
kripodb fragments filter [-h] [--pdbs PDBS] [--matrix MATRIX] input output
```

## Positional Arguments

<b>input</b>	Name of fragments db input file
<b>output</b>	Name of fragments db output file, will overwrite file if it exists

## Named Arguments

<b>--pdbs</b>	Keep fragments from any of the supplied pdb codes, one pdb code per line, use - for stdin
<b>--matrix</b>	Keep fragments which are in similarity matrix file

## merge

Combine fragments databases into a single new one

```
kripodb fragments merge [-h] ins [ins ...] out
```

## Positional Arguments

<b>ins</b>	Input fragments database files
<b>out</b>	Output fragments database file

## export\_sd

Export molblocks of all fragments as SDF file

```
kripodb fragments export_sd [-h] fragmentsdb sdfile
```

## Positional Arguments

<b>fragmentsdb</b>	Input fragments database file
<b>sdfile</b>	Output SDF file

## 4.3.3 similarities

Similarity matrix

```
kripodb similarities [-h]
                        {similar,merge,export,import,filter,freeze,thaw,fpneigh2tsv,
↪ histogram}
                        ...
```

**Sub-commands:****similar**

Find the fragments closets to query based on similarity matrix

```
kripodb similarities similar [-h] [--out OUT] [--cutoff CUTOFF]
                             pairsdbfn query
```

**Positional Arguments**

<b>pairsdbfn</b>	hdf5 similarity matrix file or base url of kripodb webservice
<b>query</b>	Query fragment identifier

**Named Arguments**

<b>--out</b>	Output file tab delimited (query, hit, similarity score) Default: -
<b>--cutoff</b>	Similarity cutoff (default: 0.55) Default: 0.55

**merge**

Combine pairs files into a new file

```
kripodb similarities merge [-h] ins [ins ...] out
```

**Positional Arguments**

<b>ins</b>	Input pair file in hdf5_compact format
<b>out</b>	Output pair file in hdf5_compact format

**export**

Export similarity matrix to tab delimited file

```
kripodb similarities export [-h] [--no_header] [--frag1] [--pdb PDB]
                             simmatrixfn outputfile
```

**Positional Arguments**

<b>simmatrixfn</b>	Compact hdf5 similarity matrix filename
<b>outputfile</b>	Tab delimited output file, use - for stdout

## Named Arguments

<b>--no_header</b>	Output no header (default: False) Default: False
<b>--frag1</b>	Only output *frag1 fragments (default: False) Default: False
<b>--pdb</b>	Only output fragments which are from pdb code in file, one pdb code per line (default: None)

## import

When input has been split into chunks, use `--ignore_upper_triangle` flag for similarities between same chunk. This prevents storing pair a->b also as b->a.

```
kripodb similarities import [-h] [--inputformat {tsv,fpneigh}]
                           [--nrrows NRROWS] [--ignore_upper_triangle]
                           inputfile fragmentsdb simmatrixfn
```

## Positional Arguments

<b>inputfile</b>	Input file, use - for stdin
<b>fragmentsdb</b>	Name of fragments db file (default: "fragments.db") Default: "fragments.db"
<b>simmatrixfn</b>	Compact hdf5 similarity matrix file, will overwrite file if it exists

## Named Arguments

<b>--inputformat</b>	Possible choices: tsv, fpneigh tab delimited (tsv) or fpneigh formatted input (default: "fpneigh") Default: "fpneigh"
<b>--nrrows</b>	Number of rows in inputfile (default: 65536) Default: 65536
<b>--ignore_upper_triangle</b>	Ignore upper triangle (default: False) Default: False

## filter

Filter similarity matrix

```
kripodb similarities filter [-h] [--fragmentsdb FRAGMENTSDDB | --skip SKIP]
                           input output
```

## Positional Arguments

<b>input</b>	Input hdf5 similarity matrix file
<b>output</b>	Output hdf5 similarity matrix file, will overwrite file if it exists

## Named Arguments

<b>--fragmentsdb</b>	Name of fragments db file, fragments in it will be kept as well as their pair counter parts.
<b>--skip</b>	File with fragment identifiers on each line to skip

## freeze

Optimize similarity matrix for reading

```
kripodb similarities freeze [-h] [-f FRAME_SIZE] [-m MEMORY] [-l LIMIT] [-s]
                           in_fn out_fn
```

## Positional Arguments

<b>in_fn</b>	Input pairs file
<b>out_fn</b>	Output array file, file is overwritten

## Named Arguments

<b>-f, --frame_size</b>	Size of frame (default: 100000000) Default: 100000000
<b>-m, --memory</b>	Memory cache in Gigabytes (default: 1) Default: 1
<b>-l, --limit</b>	Number of pairs to copy, None for no limit (default: None)
<b>-s, --single_sided</b>	Store half matrix (default: False) Default: False

## thaw

Optimize similarity matrix for writing

```
kripodb similarities thaw [-h] [--nonzero_fraction NONZERO_FRACTION]
                           in_fn out_fn
```

## Positional Arguments

<b>in_fn</b>	Input packed frozen matrix file
<b>out_fn</b>	Output pairs file, file is overwritten

## Named Arguments

<b>--nonzero_fraction</b>	Fraction of pairs which have score above threshold (default: 0.012) Default: 0.012
---------------------------	---

## fpneigh2tsv

Convert fpneigh formatted file to tab delimited file

```
kripodb similarities fpneigh2tsv [-h] inputfile outputfile
```

## Positional Arguments

<b>inputfile</b>	Input file, use - for stdin
<b>outputfile</b>	Tab delimited output file, use - for stdout

## histogram

Distribution of similarity scores

```
kripodb similarities histogram [-h] [-f FRAME_SIZE] [-r] [-l]  
inputfile outputfile
```

## Positional Arguments

<b>inputfile</b>	Filename of similarity matrix hdf5 file
<b>outputfile</b>	Tab delimited output file, use - for stdout

## Named Arguments

<b>-f, --frame_size</b>	Size of frame (default: 100000000) Default: 100000000
<b>-r, --raw_score</b>	Return raw score (16 bit integer) instead of fraction score Default: False
<b>-l, --lower_triangle</b>	Return scores from lower triangle else return scores from upper triangle Default: False



### 4.3.4 dive

DiVE visualization utils

```
kripodb dive [-h] {fragments,dump,export} ...
```

#### Sub-commands:

##### fragments

Export fragments as DiVE formatted sphere

```
kripodb dive fragments [-h] [--onlyfrag1] inputfile outputfile
```

#### Positional Arguments

<b>inputfile</b>	Name of fragments db input file
<b>outputfile</b>	Name of fragments dive output file, use - for stdout

#### Named Arguments

<b>--onlyfrag1</b>	Only *_frag1 (default: False) Default: False
--------------------	---

##### dump

Dump dense matrix with zeros

```
kripodb dive dump [-h] [--frag1only] inputfile outputfile
```

#### Positional Arguments

<b>inputfile</b>	Name of dense similarity matrix
<b>outputfile</b>	Name of output file, use - for stdout

#### Named Arguments

<b>--frag1only</b>	Only *frag1 (default: False) Default: False
--------------------	--

## export

Writes props for DiVE visualization

```
kripodb dive export [-h] [--propnames PROPNames] [--props PROPS]
                    [--pdbtags PDBTAGS]
                    fragmentsdb uniprot_annot
```

### Positional Arguments

<b>fragmentsdb</b>	Name of fragments db input file
<b>uniprot_annot</b>	<b>Uniprot download accession 2 gene symbol, family mapping.</b> Fetch “ <a href="http://www.uniprot.org/uniprot/?query=database:pdb&amp;format=tab&amp;columns=id,genes(PREFERRED),families,database(PDB)">http://www.uniprot.org/uniprot/?query=database:pdb&amp;format=tab&amp;columns=id,genes(PREFERRED),families,database(PDB)</a> ”

### Named Arguments

<b>--propnames</b>	Name of prop names file Default: kripo.propnames.txt
<b>--props</b>	Name of props file Default: kripo.props.txt
<b>--pdbtags</b>	Tag pdb in file by filename

## 4.3.5 serve

Serve similarity matrix, fragments db and pharmacophores db as webservice

```
kripodb serve [-h] [--internal_port INTERNAL_PORT]
              [--external_url EXTERNAL_URL]
              similarities fragments pharmacophores
```

### Positional Arguments

<b>similarities</b>	Filename of similarity matrix hdf5 file
<b>fragments</b>	Filename of fragments sqlite database file
<b>pharmacophores</b>	Filename of pharmacophores hdf5 file

### Named Arguments

<b>--internal_port</b>	TCP port on which to listen (default: 8084) Default: 8084
<b>--external_url</b>	URL which should be used in Swagger spec (default: “ <a href="http://localhost:8084/kripo">http://localhost:8084/kripo</a> ”) Default: “ <a href="http://localhost:8084/kripo">http://localhost:8084/kripo</a> ”

## 4.3.6 pharmacophores

Pharmacophores

```
kripodb pharmacophores [-h] {add,get,filter,merge,import,sd2phar} ...
```

### Sub-commands:

#### add

Add pharmacophores from directory to database

```
kripodb pharmacophores add [-h] [--nrrows NRROWS] startdir pharmacophoresdb
```

### Positional Arguments

**startdir** Directory to start finding \*.pphores.sd.gz and \*.pphores.txt files in  
**pharmacophoresdb** Name of pharmacophore db file

### Named Arguments

**--nrrows** **Number of expected pharmacophores**, only used when database is created  
 (default: 65536)  
 Default: 65536

#### get

Retrieve pharmacophore of a fragment

```
kripodb pharmacophores get [-h] [--query QUERY] [--output OUTPUT]  

  pharmacophoresdb
```

### Positional Arguments

**pharmacophoresdb** Name of pharmacophore db file

### Named Arguments

**--query** Query fragment identifier  
**--output** Phar formatted text file  
 Default: -

## filter

Filter pharmacophores

```
kripodb pharmacophores filter [-h] [--fragmentsdb FRAGMENTSDB]
                               inputfn outputfn
```

### Positional Arguments

<b>inputfn</b>	Name of input pharmacophore db file
<b>outputfn</b>	Name of output pharmacophore db file

### Named Arguments

<b>--fragmentsdb</b>	Name of fragments db file, fragments present in db are passed (default: “fragments.db”) Default: “fragments.db”
----------------------	--

## merge

Merge pharmacophore database files into new one

```
kripodb pharmacophores merge [-h] ins [ins ...] out
```

### Positional Arguments

<b>ins</b>	Input pharmacophore database files
<b>out</b>	Output pharmacophore database file

## import

Convert phar formatted file to pharmacophore database file

```
kripodb pharmacophores import [-h] [--nrrows NRROWS] infile outfile
```

### Positional Arguments

<b>infile</b>	Input phar formatted file
<b>outfile</b>	Output pharmacophore database file

### Named Arguments

**--nrrows**            **Number of expected pharmacophores**, only used when database is created  
                          (default: 65536)  
                          Default: 65536

### **sd2phar**

Convert sd formatted pharmacophore file to phar formatted file

```
kripodb pharmacophores sd2phar [-h] [--frag_id FRAG_ID] infile outfile
```

### Positional Arguments

**infile**              Input sd formatted file  
**outfile**             Output phar formatted file

### Named Arguments

**--frag\_id**            Fragment identifier  
                          Default: "frag"



## CHAPTER 5

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`





### k

- kripodb.canned, 17
- kripodb.db, 20
- kripodb.dive, 24
- kripodb.frozen, 25
- kripodb.hdf5, 27
- kripodb.makebits, 31
- kripodb.modifiedtanimoto, 32
- kripodb.pairs, 33
- kripodb.pdb, 38
- kripodb.pharmacophores, 35
- kripodb.script, 38
  - kripodb.script.dive, 39
  - kripodb.script.fingerprints, 39
  - kripodb.script.fragments, 38
  - kripodb.script.similarities, 39
- kripodb.webservice, 39
  - kripodb.webservice.client, 39
  - kripodb.webservice.server, 40



## A

AbstractSimpleTable (class in kripodb.hdf5), 27  
 adapt\_BitMap() (in module kripodb.db), 23  
 adapt\_molblockgz() (in module kripodb.db), 23  
 add\_dir() (kripodb.pharmacophores.PharmacophorePointsTable method), 36  
 add\_dir() (kripodb.pharmacophores.PharmacophoresDb method), 36  
 add\_fragment() (kripodb.db.FragmentsDb method), 20  
 add\_fragments\_from\_shelve() (kripodb.db.FragmentsDb method), 21  
 add\_molecule() (kripodb.db.FragmentsDb method), 21  
 add\_molecules() (kripodb.db.FragmentsDb method), 21  
 add\_pdb() (kripodb.db.FragmentsDb method), 21  
 append() (kripodb.hdf5.AbstractSimpleTable method), 27  
 append() (kripodb.hdf5.PairsTable method), 29  
 append() (kripodb.hdf5.SimilarityMatrix method), 30  
 append() (kripodb.pharmacophores.PharmacophoresDb method), 36  
 as\_dict() (kripodb.db.FingerprintsDb method), 20  
 as\_phar() (in module kripodb.pharmacophores), 37

## B

by\_id() (kripodb.hdf5.LabelsLookup method), 27  
 by\_label() (kripodb.hdf5.LabelsLookup method), 27  
 by\_labels() (kripodb.hdf5.LabelsLookup method), 27  
 by\_pdb\_code() (kripodb.db.FragmentsDb method), 21

## C

calc\_mean\_onbit\_density() (in module kripodb.modifiedtanimoto), 32  
 close() (kripodb.db.SqliteDb method), 22  
 close() (kripodb.frozen.FrozenSimilarityMatrix method), 25  
 close() (kripodb.hdf5.SimilarityMatrix method), 30  
 close() (kripodb.pharmacophores.PharmacophoresDb method), 37  
 commit() (kripodb.db.SqliteDb method), 22  
 connection (kripodb.db.SqliteDb attribute), 22

connection (kripodb.db.SqliteDict attribute), 22  
 convert\_BitMap() (in module kripodb.db), 23  
 convert\_molblockgz() (in module kripodb.db), 24  
 corrections() (in module kripodb.modifiedtanimoto), 32  
 count() (kripodb.frozen.FrozenSimilarityMatrix method), 25  
 count() (kripodb.hdf5.PairsTable method), 29  
 count() (kripodb.hdf5.SimilarityMatrix method), 30  
 create\_tables() (kripodb.db.FingerprintsDb method), 20  
 create\_tables() (kripodb.db.FragmentsDb method), 21  
 create\_tables() (kripodb.db.SqliteDb method), 22  
 cursor (kripodb.db.SqliteDb attribute), 22  
 cursor (kripodb.db.SqliteDict attribute), 22

## D

default() (kripodb.webservice.server.KripodbJSONEncoder method), 41  
 dense\_dump() (in module kripodb.dive), 24  
 dense\_dump\_iter() (in module kripodb.dive), 24  
 dense\_dump\_sc() (in module kripodb.script.dive), 39  
 dive\_export() (in module kripodb.dive), 24  
 dive\_sphere() (in module kripodb.dive), 24  
 dump\_pairs() (in module kripodb.pairs), 33  
 dump\_pairs\_hdf5() (in module kripodb.pairs), 34  
 dump\_pairs\_tsv() (in module kripodb.pairs), 34

## F

FastInserter (class in kripodb.db), 20  
 FEATURE\_TYPES (in module kripodb.pharmacophores), 35  
 fetch() (kripodb.pdb.PdbReport method), 38  
 find() (kripodb.frozen.FrozenSimilarityMatrix method), 26  
 find() (kripodb.hdf5.PairsTable method), 29  
 find() (kripodb.hdf5.SimilarityMatrix method), 30  
 FingerprintsDb (class in kripodb.db), 20  
 fragments\_by\_id() (in module kripodb.canned), 17  
 fragments\_by\_id() (kripodb.webservice.client.WebserviceClient method), 40

fragments\_by\_pdb\_codes() (in module kripodb.canned), 18

fragments\_by\_pdb\_codes() (kripodb.webservice.client.WebserviceClient method), 40

FragmentsDb (class in kripodb.db), 20

from\_array() (kripodb.frozen.FrozenSimilarityMatrix method), 26

from\_pairs() (kripodb.frozen.FrozenSimilarityMatrix method), 26

FrozenSimilarityMatrix (class in kripodb.frozen), 25

full\_matrix (kripodb.hdf5.PairsTable attribute), 29

## G

get\_fragment\_phar() (in module kripodb.webservice.server), 41

get\_fragment\_svg() (in module kripodb.webservice.server), 41

get\_fragments() (in module kripodb.webservice.server), 41

get\_similar\_fragments() (in module kripodb.webservice.server), 41

get\_version() (in module kripodb.webservice.server), 42

## H

h5file (kripodb.frozen.FrozenSimilarityMatrix attribute), 25

h5file (kripodb.hdf5.SimilarityMatrix attribute), 30

h5file (kripodb.pharmacophores.PharmacophoresDb attribute), 36

## I

id2label() (kripodb.db.FragmentsDb method), 21

Incomplete, 39

IncompleteFragments, 39

IncompleteHits, 17

IncompletePharmacophores, 39

IntbitsetDict (class in kripodb.db), 22

is\_ligand\_stored() (kripodb.db.FragmentsDb method), 21

items() (kripodb.db.SQLiteDict method), 22

iter\_file() (in module kripodb.makebits), 31

iteritems() (kripodb.db.SQLiteDict method), 22

iteritems\_startswith() (kripodb.db.SQLiteDict method), 22

itervalues() (kripodb.db.SQLiteDict method), 23

## K

keep() (kripodb.hdf5.LabelsLookup method), 28

keep() (kripodb.hdf5.PairsTable method), 29

keep() (kripodb.hdf5.SimilarityMatrix method), 31

kripodb.canned (module), 17

kripodb.db (module), 20

kripodb.dive (module), 24

kripodb.frozen (module), 25

kripodb.hdf5 (module), 27

kripodb.makebits (module), 31

kripodb.modifiedanimoto (module), 32

kripodb.pairs (module), 33

kripodb.pdb (module), 38

kripodb.pharmacophores (module), 35

kripodb.script (module), 38

kripodb.script.dive (module), 39

kripodb.script.fingerprints (module), 39

kripodb.script.fragments (module), 38

kripodb.script.similarities (module), 39

kripodb.webservice (module), 39

kripodb.webservice.client (module), 39

kripodb.webservice.server (module), 40

KripodbJSONEncoder (class in kripodb.webservice.server), 40

## L

label2id() (kripodb.db.FragmentsDb method), 21

label2ids() (kripodb.hdf5.LabelsLookup method), 28

labels (kripodb.frozen.FrozenSimilarityMatrix attribute), 25

labels (kripodb.hdf5.SimilarityMatrix attribute), 30

LabelsLookup (class in kripodb.hdf5), 27

## M

main() (in module kripodb.script), 38

make\_fingerprints\_parser() (in module kripodb.script.fingerprints), 39

make\_fragments\_parser() (in module kripodb.script.fragments), 38

make\_parser() (in module kripodb.script), 38

make\_similarities\_parser() (in module kripodb.script.similarities), 39

materialize() (kripodb.db.SQLiteDict method), 23

merge() (in module kripodb.pairs), 34

merge() (kripodb.hdf5.LabelsLookup method), 28

## N

number\_of\_bits (kripodb.db.IntbitsetDict attribute), 22

## O

open\_similarity\_matrix() (in module kripodb.pairs), 34

## P

pairs (kripodb.hdf5.SimilarityMatrix attribute), 30

PairsTable (class in kripodb.hdf5), 28

parse\_csv\_file() (in module kripodb.pdb), 38

PdbReport (class in kripodb.pdb), 38

PharmacophorePointsTable (class in kripodb.pharmacophores), 35

pharmacophores\_by\_id() (in module kripodb.canned), 18

PharmacophoresDb (class in kripodb.pharmacophores), 36  
 points (kripodb.pharmacophores.PharmacophoresDb attribute), 36

## R

read\_fpneighpairs\_file() (in module kripodb.script.similarities), 39  
 read\_fragtxtfile() (in module kripodb.pharmacophores), 37  
 read\_fragtxtfile\_as\_file() (in module kripodb.pharmacophores), 37  
 read\_phar() (kripodb.pharmacophores.PharmacophorePointsTable method), 36  
 read\_phar() (kripodb.pharmacophores.PharmacophoresDb method), 37  
 read\_pphore\_gzipped\_sdfile() (in module kripodb.pharmacophores), 37  
 read\_pphore\_sdfile() (in module kripodb.pharmacophores), 38

## S

score\_precision (kripodb.hdf5.PairsTable attribute), 28  
 scores (kripodb.frozen.FrozenSimilarityMatrix attribute), 25  
 serve\_app() (in module kripodb.webservice.server), 42  
 similar() (in module kripodb.pairs), 34  
 similar\_fragments() (kripodb.webservice.client.WebserviceClient method), 40  
 similar\_run() (in module kripodb.pairs), 35  
 similarities() (in module kripodb.canned), 19  
 similarities() (in module kripodb.modifiedtanimoto), 32  
 similarity() (in module kripodb.modifiedtanimoto), 32  
 similarity2query() (in module kripodb.pairs), 35  
 SimilarityMatrix (class in kripodb.hdf5), 30  
 simmatrix\_export\_run() (in module kripodb.script.similarities), 39  
 skip() (kripodb.hdf5.LabelsLookup method), 28  
 skip() (kripodb.hdf5.PairsTable method), 29  
 skip() (kripodb.hdf5.SimilarityMatrix method), 31  
 SQLiteDatabase (class in kripodb.db), 22  
 SQLiteDict (class in kripodb.db), 22

## T

to\_pairs() (kripodb.frozen.FrozenSimilarityMatrix method), 26  
 to\_pandas() (kripodb.frozen.FrozenSimilarityMatrix method), 26  
 total\_number\_of\_pairs() (in module kripodb.pairs), 35

## U

update() (kripodb.db.IntbitsetDict method), 22  
 update() (kripodb.hdf5.LabelsLookup method), 28

update() (kripodb.hdf5.PairsTable method), 29  
 update() (kripodb.hdf5.SimilarityMatrix method), 31  
 url (kripodb.pdb.PdbReport attribute), 38

## V

values() (kripodb.db.SQLiteDict method), 23

## W

WebserviceClient (class in kripodb.webservice.client), 39  
 write\_file() (in module kripodb.makebits), 31  
 write\_phar() (kripodb.pharmacophores.PharmacophoresDb method), 37  
 wsgi\_app() (in module kripodb.webservice.server), 42