

---

# **CLK hash Documentation**

*Release 0.9.0*

**N1 Analytics**

**Apr 23, 2018**



---

## Contents:

---

<b>1</b>	<b>API Documentation</b>	<b>1</b>
1.1	Bloom filter . . . . .	1
1.2	CLK . . . . .	2
1.3	identifier types . . . . .	3
1.4	IdentifierType . . . . .	3
1.5	key derivation . . . . .	4
1.6	random names . . . . .	5
1.7	schema . . . . .	6
1.8	tokenizer . . . . .	6
<b>2</b>	<b>Indices and tables</b>	<b>7</b>
	<b>Python Module Index</b>	<b>9</b>



## 1.1 Bloom filter

Generate a Bloom filter

`clkhash.bloomfilter.calculate_bloom_filters` (*dataset, schema, keys, xor\_folds=0*)

### Parameters

- **dataset** – A list of indexable records.
- **schema** – An iterable of identifier types.
- **keys** – A tuple of two lists of secret keys used in the HMAC.
- **xor\_folds** – number of XOR folds to perform

**Returns** List of bloom filters as 3-tuples, each containing bloom filter (bitarray), record first element - usually index, bitcount (int)

`clkhash.bloomfilter.crypto_bloom_filter` (*record, tokenizers, keys1, keys2, xor\_folds=0, l=1024, k=30*)

Makes a Bloom filter from a record with given tokenizers and lists of keys.

Using the method from <http://www.record-linkage.de/-download=wp-grlc-2011-02.pdf>

### Parameters

- **record** – plaintext record tuple. E.g. (index, name, dob, gender)
- **tokenizers** – A list of IdentifierType tokenizers (one for each record element)
- **keys1** – list of keys for first hash function as list of bytes
- **keys2** – list of keys for second hash function as list of bytes
- **xor\_folds** – number of XOR folds to perform
- **l** – length of the Bloom filter in number of bits
- **k** – number of hash functions to use per element

**Returns** 3-tuple: - bloom filter for record as a bytearray - first element of record (usually an index) - number of bits set in the bloomfilter

`clkhash.bloomfilter.double_hash_encode_ngrams` (*ngrams*, *key\_sha1*, *key\_md5*, *k*, *l*)  
computes the double hash encoding of the provided ngrams with the given keys.

Using the method from <http://www.record-linkage.de/-download=wp-grlc-2011-02.pdf>

**Parameters**

- **ngrams** – list of n-grams to be encoded
- **key\_sha1** – hmac secret keys for sha1 as bytes
- **key\_md5** – hmac secret keys for md5 as bytes
- **k** – number of hash functions to use per element of the ngrams
- **l** – length of the output bytearray

**Returns** bytearray of length *l* with the bits set which correspond to the encoding of the ngrams

`clkhash.bloomfilter.fold_xor` (*bloomfilter*, *folds*)  
Performs XOR folding on a Bloom filter.

If the length of the original Bloom filter is *n* and we perform *r* folds, then the length of the resulting filter is  $n / 2^{**r}$ .

**Parameters**

- **bloomfilter** – Bloom filter to fold
- **folds** – number of folds

**Returns** folded bloom filter

`clkhash.bloomfilter.serialize_bitarray` (*ba*)  
Serialize a bytearray (bloomfilter)

`clkhash.bloomfilter.stream_bloom_filters` (*dataset*, *schema\_types*, *keys*, *xor\_folds=0*)  
Yield bloom filters

**Parameters**

- **dataset** – An iterable of indexable records.
- **schema\_types** – An iterable of identifier type names.
- **keys** – A tuple of two lists of secret keys used in the HMAC.
- **xor\_folds** – number of XOR folds to perform

**Returns** Yields bloom filters as 3-tuples

## 1.2 CLK

Generate CLK from CSV file

`clkhash.clk.chunks` (*l*, *n*)  
Yield successive n-sized chunks from *l*.

`clkhash.clk.generate_clk_from_csv` (*input*, *keys*, *schema\_types*, *no\_header=False*, *progress\_bar=True*, *xor\_folds=0*)

`clkhash.clk.generate_clks` (*pii\_data*, *schema\_types*, *key\_lists*, *xor\_folds*, *callback=None*)

`clkhask.clk.hash_and_serialize_chunk` (*chunk\_pii\_data, schema\_types, keys, xor\_folds*)

Generate Bloom filters (ie hash) from chunks of PII then serialize the generated Bloom filters.

**Parameters**

- **chunk\_pii\_data** – An iterable of indexable records.
- **schema\_types** – An iterable of identifier type names.
- **keys** – A tuple of two lists of secret keys used in the HMAC.
- **xor\_folds** – Number of XOR folds to perform. Each fold halves the hash length.

**Returns** A list of serialized Bloom filters

## 1.3 identifier types

Convert PII to tokens

`clkhask.identifier_types.identifier_type_from_description` (*schema\_object*)

Convert a dictionary describing a feature into an IdentifierType

**Parameters** *schema\_object* –

**Returns** An IdentifierType

## 1.4 IdentifierType

**class** `clkhask.identifier_types.IdentifierType` (*unigram=False, weight=1, \*\*kwargs*)

Bases: `object`

Base class used for all identifier types.

Required to provide a mapping of schema to hash type uni-gram or bi-gram.

`__call__` (*entry*)

Call self as a function.

`__init__` (*unigram=False, weight=1, \*\*kwargs*)

**Parameters**

- **unigram** (*bool*) – Use uni-gram instead of using bi-grams
- **weight** (*float*) – adjusts the “importance” of this identifier in the Bloom filter. Can be set to zero to skip
- **kwargs** – Extra keyword arguments passed to the tokenizer

---

**Note:** For each n-gram of an identifier, we compute  $k$  different indices in the Bloom filter which will be set to true. There is a global  $k_{default}$  value, and the  $k$  value for each identifier is computed as

$$k = weight * k_{default},$$

rounded to the nearest integer.

Reasons why you might want to set weights:

- Long identifiers like street name will produce a lot more n-grams than small identifiers like zip code. Thus street name will flip more bits in the Bloom filter and will have a bigger influence in the overall matching score.
  - The matching might produce better results if identifiers that are stable and / or have low error rates are given higher prominence in the Bloom filter.
- 

**`__weakref__`**

list of weak references to the object (if defined)

## 1.5 key derivation

```
class clkhash.key_derivation.HKDFconfig(master_secret, salt=None, info=None,
                                         hash_algo='SHA256')
```

Bases: `object`

```
static check_is_bytes (value)
```

```
static check_is_bytes_or_none (value)
```

```
supported_hash_algos = ('SHA256', 'SHA512')
```

```
clkhash.key_derivation.generate_key_lists(master_secrets, num_identifier, key_size=64,
                                          salt=None, info=None, kdf='HKDF')
```

Generates a derived key for each identifier for each master secret using a key derivation function (KDF).

The only supported key derivation function for now is 'HKDF'.

The previous key usage can be reproduced by setting `kdf` to 'legacy'. This is highly discouraged, as this strategy will map the same n-grams in different identifier to the same bits in the Bloom filter and thus does not lead to good results.

**Parameters**

- **master\_secrets** – a list of master secrets (either as bytes or strings)
- **num\_identifier** – the number of identifiers
- **key\_size** – the size of the derived keys
- **salt** – salt for the KDF as bytes
- **info** – optional context and application specific information as bytes
- **kdf** – the key derivation function algorithm to use

**Returns** The derived keys. First dimension is the same as `master_secrets`, second dimension is of size `num_identifier`. A key is represented as bytes.

```
clkhash.key_derivation.hkdf(hkdf_config, num_keys, key_size=64)
```

Executes the HKDF key derivation function as described in rfc5869 to derive `num_keys` keys of size `key_size` from the `master_secret`.

**Parameters**

- **hkdf\_config** – an `HKDFconfig` object containing the configuration for the HKDF.
- **num\_keys** – the number of keys the kdf should produce
- **key\_size** – the size of the produced keys

**Returns** Derived keys

## 1.6 random names

Module to produce a dataset of names, genders and dates of birth and manipulate that list

Currently very simple and not realistic. Additional functions for manipulating the list of names - producing reordered and subset lists with a specific overlap

ClassList class - generate a list of length n of [id, name, dob, gender] lists

TODO: Get age distribution right by using a mortality table TODO: Get first name distributions right by using distributions TODO: Generate realistic errors TODO: Add RESTfull api to generate reasonable name data as requested

```
class clkhash.randomnames.NameList(n)
```

Bases: `object`

List of randomly generated names

```
generate_random_person(n)
```

Generator that yields details on a person with plausible name, sex and age.

**Yields** Generated data for one person tuple - (id: int, name: str('First Last'), birthdate: str('DD/MM/YYYY'), sex: str('M' | 'F'))

```
generate_subsets(sz, overlap=0.8)
```

Generate a pair of subsets of the name list with a specified overlap

**Parameters**

- **sz** – length of subsets to generate
- **overlap** – fraction of the subsets that should have the same names in them

**Returns** 2-tuple of lists of subsets

```
load_names()
```

This function loads a name database into globals firstNames and lastNames

initial version uses data files from <http://www.quietaffiliate.com/free-first-name-and-last-name-databases-csv-and-sql/>

```
schema = [{'identifier': 'INDEX'}, {'identifier': 'NAME freetext'}, {'identifier':
```

```
schema_types
```

```
clkhash.randomnames.load_csv_data(resource_name)
```

Loads a specified data file as csv and returns the first column as a Python list

```
clkhash.randomnames.random_date(start, end)
```

This function will return a random datetime between two datetime objects.

**Parameters**

- **start** – datetime of start
- **end** – datetime of end

**Returns** random datetime between start and end

```
clkhash.randomnames.save_csv(data, schema, file)
```

Output generated data as csv with header.

**Parameters**

- **data** – An iterable of tuples containing raw data.

- **schema** – Iterable of schema definition dicts
- **file** – A writeable stream in which to write the csv

## 1.7 schema

`clckhash.schema.get_schema_types` (*schema*)

`clckhash.schema.load_schema` (*schema\_file*)

## 1.8 tokenizer

Functions to tokenize words (PII)

`clckhash.tokenizer.bigramlist` (*word, toremove=None*)

Make bigrams from word with pre- and ap-pended spaces

*s* -> [' ' + *s*0, *s*0 + *s*1, *s*1 + *s*2, .. *s*N + ' ']

### Parameters

- **word** – string to make bigrams from
- **toremove** – List of strings to remove before construction

**Returns** list of bigrams as strings

`clckhash.tokenizer.positional_unigrams` (*instr*)

Make positional unigrams from a word.

E.g. 1987 -> ["1 1", "2 9", "3 8", "4 7"]

**Parameters** *instr* – input string

**Returns** list of strings with unigrams

`clckhash.tokenizer.unigramlist` (*instr, toremove=None, positional=False*)

Make 1-grams (unigrams) from a word, possibly excluding particular substrings

### Parameters

- **instr** – input string
- **toremove** – Iterable of strings to remove

**Returns** list of strings with unigrams

## CHAPTER 2

---

### Indices and tables

---

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



**C**

`clckhash.bloomfilter`, 1  
`clckhash.clk`, 2  
`clckhash.identifier_types`, 3  
`clckhash.key_derivation`, 4  
`clckhash.randomnames`, 5  
`clckhash.schema`, 6  
`clckhash.tokenizer`, 6



## Symbols

`__call__()` (clkhash.identifier\_types.IdentifierType method), 3  
`__init__()` (clkhash.identifier\_types.IdentifierType method), 3  
`__weakref__` (clkhash.identifier\_types.IdentifierType attribute), 4

## B

`bigramlist()` (in module clkhash.tokenizer), 6

## C

`calculate_bloom_filters()` (in module clkhash.bloomfilter), 1  
`check_is_bytes()` (clkhash.key\_derivation.HKDFconfig static method), 4  
`check_is_bytes_or_none()` (clkhash.key\_derivation.HKDFconfig static method), 4  
`chunks()` (in module clkhash.clk), 2  
`clkhash.bloomfilter` (module), 1  
`clkhash.clk` (module), 2  
`clkhash.identifier_types` (module), 3  
`clkhash.key_derivation` (module), 4  
`clkhash.randomnames` (module), 5  
`clkhash.schema` (module), 6  
`clkhash.tokenizer` (module), 6  
`crypto_bloom_filter()` (in module clkhash.bloomfilter), 1

## D

`double_hash_encode_ngrams()` (in module clkhash.bloomfilter), 2

## F

`fold_xor()` (in module clkhash.bloomfilter), 2

## G

`generate_clk_from_csv()` (in module clkhash.clk), 2  
`generate_clks()` (in module clkhash.clk), 2

`generate_key_lists()` (in module clkhash.key\_derivation), 4  
`generate_random_person()` (clkhash.randomnames.NameList method), 5  
`generate_subsets()` (clkhash.randomnames.NameList method), 5  
`get_schema_types()` (in module clkhash.schema), 6

## H

`hash_and_serialize_chunk()` (in module clkhash.clk), 2  
`hkdf()` (in module clkhash.key\_derivation), 4  
`HKDFconfig` (class in clkhash.key\_derivation), 4

## I

`identifier_type_from_description()` (in module clkhash.identifier\_types), 3  
`IdentifierType` (class in clkhash.identifier\_types), 3

## L

`load_csv_data()` (in module clkhash.randomnames), 5  
`load_names()` (clkhash.randomnames.NameList method), 5  
`load_schema()` (in module clkhash.schema), 6

## N

`NameList` (class in clkhash.randomnames), 5

## P

`positional_unigrams()` (in module clkhash.tokenizer), 6

## R

`random_date()` (in module clkhash.randomnames), 5

## S

`save_csv()` (in module clkhash.randomnames), 5  
`schema` (clkhash.randomnames.NameList attribute), 5  
`schema_types` (clkhash.randomnames.NameList attribute), 5

`serialize_bitarray()` (in module `clkhask.bloomfilter`), 2  
`stream_bloom_filters()` (in module `clkhask.bloomfilter`), 2  
`supported_hash_algos` (`clkhask.key_derivation.HKDFconfig`  
attribute), 4

## U

`unigramlist()` (in module `clkhask.tokenizer`), 6