

FLASH

Generated by Doxygen 1.8.14

Contents

1	FLASH	1
2	Class Index	3
2.1	Class List	3
3	Class Documentation	5
3.1	FrequentItems Class Reference	5
3.2	greater Struct Reference	5
3.3	LSH Class Reference	6
3.3.1	Constructor & Destructor Documentation	6
3.3.1.1	LSH() [1/2]	6
3.3.1.2	LSH() [2/2]	6
3.3.1.3	~LSH()	7
3.3.2	Member Function Documentation	7
3.3.2.1	cLSH()	7
3.3.2.2	getHash() [1/4]	8
3.3.2.3	getHash() [2/4]	8
3.3.2.4	getHash() [3/4]	9
3.3.2.5	getHash() [4/4]	9
3.4	LSHReservoirSampler Class Reference	10
3.4.1	Detailed Description	11
3.4.2	Constructor & Destructor Documentation	11
3.4.2.1	LSHReservoirSampler()	11
3.4.2.2	~LSHReservoirSampler()	11
3.4.3	Member Function Documentation	12
3.4.3.1	add() [1/2]	12
3.4.3.2	add() [2/2]	12
3.4.3.3	ann() [1/2]	12
3.4.3.4	ann() [2/2]	13
3.4.3.5	checkTableMemLoad()	13
3.4.3.6	showParams()	14

Chapter 1

FLASH

FLASH (Fast [LSH](#) Algorithm for Similarity Search Accelerated with HPC) is a library for large scale approximate nearest neighbor search of sparse vectors. It is currently available in C++ for CPU parallel computing and supports OpenCL enabled GPGPU computing. See [our paper](#) for theoretical and benchmarking details.

Prerequisites

The current version of the soft is tested on 64-bit machines running Ubuntu 16.04, with CPU and at least 1 GPGPU installed. The compiler needs to support C++11 and OpenMP.

GPGPU

An active installation of OpenCL 1.1 or OpenCL 2.0 is required to support the GPGPU capability. CPUs running OpenCL is supported but not recommended as the performance will be sub-optimal. OpenCL on Nvidia graphics cards requires the installation of [CUDA](#).

Then install clinfo (using apt-get) to verify the number of OpenCL platforms and devices. These information will be used in the configurations.

Configuration and Compilation

Navigate to the FLASH directory. Configure the system by editing the following section of LSHReservoir_config.h guided by the comments:

```
// Comment out if not using GPU.
#define USE_GPU
// Comment out if using OpenCL 1.XX. Does not matter if not using GPU.
#define OPENCL_2XX

#define CL_GPU_PLATFORM 0 // Does not matter if not using OpenCL-GPU.
#define CL_CPU_PLATFORM 1 // Does not matter if not using OpenCL-CPU.
#define CL_GPU_DEVICE 0 // Does not matter if not using OpenCL-GPU.
#define CL_CPU_DEVICE 0 // Does not matter if not using OpenCL-CPU.
```

Fill in CL_GPU_PLATFORM / CL_CPU_PLATFORM according to the order that the platforms appear in the output of clinfo. If multiple devices exist, fill in CL_GPU_DEVICE / CL_CPU_DEVICE to choose the desired device according to their order in the output of clinfo.

Save and close the file. Type in terminal:

```
make
```

The compilation is complete if no errors appear.

Tutorial

The current example code in the `main()` function verifies one of the results presented in [our paper](#) on the Web-spam dataset. Download the dataset from [libsvm](#), and rename the libsvm-format file as `trigram.svm`. Download the groundtruths from this [link](#). Place the dataset and groundtruth files in the FLASH directory. Run the program from the terminal:

```
./runme
```

The test program builds multiple hash tables for the dataset and query 10,000 test vectors followed by quality evaluations. The program will run with console outputs, indicating the progress and performance. The parameters, such as L, K and R can be edited in `main.cpp`. A re-compilation is required after changing the parameters. Please note that the time for parsing the dataset from disk might take about 5-10 minutes.

A basic documentation of the API generated by [doxygen](#) is available as `doc.pdf` in the working directory.

Authors

- [Yiqiu Wang](#) designed and implemented the CPU and GPU FLASH.
- [Anshumali Shrivastava](#) invented and is author of the [DOPH](#) hash codes. He is the main contributor to the theoretical aspect of [our work](#).
- Heejung Ryu actively participated in the process and contributed to the testing and comparison of FLASH with similar works.

License

This project is licensed under Apache License. See LICENSE for more details.

Acknowledgments

- Rice University Sketching and Hashing Lab ([RUSH Lab](#)) provided the computing platform for testing.

Chapter 2

Class Index

2.1 Class List

Here are the classes, structs, unions and interfaces with brief descriptions:

FrequentItems	5
greater	5
LSH	6
LSHReservoirSampler	10

Chapter 3

Class Documentation

3.1 FrequentItems Class Reference

Public Member Functions

- **FrequentItems** (int k)
- void **increment** (int item)
- unsigned int * **getTopk** ()
- void **getTopk** (unsigned int *outputs)

The documentation for this class was generated from the following files:

- FrequentItems.h
- FrequentItems.cpp

3.2 greater Struct Reference

Public Member Functions

- template<class T >
bool **operator()** (T const &a, T const &b) const

The documentation for this struct was generated from the following file:

- KNN_bruteforce.cpp

3.3 LSH Class Reference

Public Member Functions

- void [getHash](#) (cl_mem *hashIndices_obj, cl_mem *identity_obj, cl_mem *dataIdx_obj, cl_mem *dataVal_obj, cl_mem *dataMarker_obj, int numInputEntries, int numProbes)
- void [getHash](#) (unsigned int *hashIndices, unsigned int *identity, int *dataIdx, float *dataVal, int *dataMarker, int numInputEntries, int numProbes)
- void [getHash](#) (cl_mem *hashIndices_obj, cl_mem *identity_obj, cl_mem *input_obj, int numInputEntries, int numProbes)
- void [getHash](#) (unsigned int *hashIndices, unsigned int *identity, float *input, int numInputEntries, int numProbes)
- [LSH](#) (int hashType, int numHashPerFamily, int numHashFamilies, int dimension, int samFactor)
- [LSH](#) (int hashType, int _K_in, int _L_in, int _rangePow_in)
- void [clLSH](#) (cl_platform_id *platforms_lsh, cl_device_id *devices_lsh, cl_context context_lsh, cl_program program_lsh, cl_command_queue command_queue_lsh)
- [~LSH](#) ()

3.3.1 Constructor & Destructor Documentation

3.3.1.1 LSH() [1/2]

```
LSH::LSH (
    int hashType,
    int numHashPerFamily,
    int numHashFamilies,
    int dimension,
    int samFactor )
```

Constructor.

Construct an [LSH](#) class for signed random projection.

Parameters

<i>hashType</i>	For SRP, use 1.
<i>numHashPerFamily</i>	Number of hash (bits) per hashfamily (hash table).
<i>numHashFamilies</i>	Number of hash families (hash tables).
<i>dimension</i>	Dimensionality of input data.
<i>samFactor</i>	samFactor = dimension / samSize, have to be an integer.

3.3.1.2 LSH() [2/2]

```
LSH::LSH (
    int hashType,
```

```
int _K_in,
int _L_in,
int _rangePow_in )
```

Constructor.

Construct an [LSH](#) class for optimal densified min-hash (for more details refer to Anshumali Shrivastava, anshu@rice.edu). This hashing scheme is for very sparse and high dimensional data stored in sparse format.

Parameters

<i>hashType</i>	For optimal densified min-hash, use 2.
-----------------	--

3.3.1.3 ~LSH()

```
LSH::~~LSH ( )
```

Destructor.

Does not uninitialized the OpenCL environment (if applicable).

3.3.2 Member Function Documentation

3.3.2.1 clLSH()

```
void LSH::clLSH (
    cl_platform_id * platforms_lsh,
    cl_device_id * devices_lsh,
    cl_context context_lsh,
    cl_program program_lsh,
    cl_command_queue command_queue_lsh )
```

Initializes [LSH](#) hashing with OpenCL environment.

Takes in initialized OpenCL objects to support OpenCL hash functions, given that the particular type of hashing instantiated have an OpenCL implementation. Otherwise an error will be triggered during hashing.

Parameters

<i>platforms_lsh</i>	Reference to OpenCL <code>cl_platform_id</code> .
<i>devices_lsh</i>	Reference to OpenCL <code>cl_device_id</code> .
<i>context_lsh</i>	Reference to OpenCL <code>context_lsh</code> .
<i>program_lsh</i>	Reference to OpenCL <code>program_lsh</code> .
<i>command_queue_lsh</i>	Reference to OpenCL <code>cl_command_queue</code> .

3.3.2.2 `getHash()` [1/4]

```
void LSH::getHash (
    cl_mem * hashIndices_obj,
    cl_mem * identity_obj,
    cl_mem * dataIdx_obj,
    cl_mem * dataVal_obj,
    cl_mem * dataMarker_obj,
    int numInputEntries,
    int numProbes )
```

Obtain hash indice given the (sparse) input vector, using OpenCL.

Hash indice refer to the corresponding "row number" in a hash table, in the form of unsigned integer. This function will only be valid when an OpenCL implementation exists for that type of hashing, and when OpenCL is initialized (clLSH) for that hashing. The outputs indexing is defined as `hashIndicesOutputIdx(numHashFamilies, numProbes, numInputs, inputIdx, probelIdx, tb)` (unsigned)(`numInputs * numProbes * tb + inputIdx * numProbes + probelIdx`).

Parameters

<i>hashIndices_obj</i>	Hash indice for the batch of input vectors.
<i>identity_obj</i>	Hash indice's corresponding identifications (sequential number starting 0) for this batch of inputs.
<i>dataIdx_obj</i>	Non-zero indice of the sparse format.
<i>dataVal_obj</i>	Non-zero values of the sparse format.
<i>dataMarker_obj</i>	Marks the start index of each vector in <i>dataIdx</i> and <i>dataVal</i> .
<i>numInputEntries</i>	Number of input vectors.
<i>numProbes</i>	Number of probes per input.

3.3.2.3 `getHash()` [2/4]

```
void LSH::getHash (
    unsigned int * hashIndices,
    unsigned int * identity,
    int * dataIdx,
    float * dataVal,
    int * dataMarker,
    int numInputEntries,
    int numProbes )
```

Obtain hash indice given the (sparse) input vector, using CPU.

Hash indice refer to the corresponding "row number" in a hash table, in the form of unsigned integer. This function will only be valid when an CPU implementation exists for that type of hashing. The outputs indexing is defined as `hashIndicesOutputIdx(numHashFamilies, numProbes, numInputs, inputIdx, probelIdx, tb)` (unsigned)(`numInputs * numProbes * tb + inputIdx * numProbes + probelIdx`).

Parameters

<i>hashIndices</i>	Hash indice for the batch of input vectors.
<i>identity</i>	Hash indice's corresponding identifications (sequential number starting 0) for this batch of inputs.
<i>dataIdx</i>	Non-zero indice of the sparse format.
<i>dataVal</i>	Non-zero values of the sparse format.
<i>dataMarker</i>	Marks the start index of each vector in <i>dataIdx</i> and <i>dataVal</i> .
<i>numInputEntries</i>	Number of input vectors.
<i>numProbes</i>	Number of probes per input.

3.3.2.4 `getHash()` [3/4]

```
void LSH::getHash (
    cl_mem * hashIndices_obj,
    cl_mem * identity_obj,
    cl_mem * input_obj,
    int numInputEntries,
    int numProbes )
```

Obtain hash indice given the (dense) input vector, using OpenCL.

Hash indice refer to the corresponding "row number" in a hash table, in the form of unsigned integer. This function will only be valid when an CPU implementation exists for that type of hashing. The outputs indexing is defined as $\text{hashIndicesOutputIdx}(\text{numHashFamilies}, \text{numProbes}, \text{numInputs}, \text{inputIdx}, \text{probIdx}, \text{tb}) = (\text{unsigned})(\text{numInputs} * \text{numProbes} * \text{tb} + \text{inputIdx} * \text{numProbes} + \text{probIdx})$.

Parameters

<i>hashIndices_obj</i>	Hash indice for the batch of input vectors.
<i>identity_obj</i>	Hash indice's corresponding identifications (sequential number starting 0) for this batch of inputs.
<i>input_obj</i>	Input vectors concatenated.
<i>numInputEntries</i>	Number of input vectors.
<i>numProbes</i>	Number of probes per input.

3.3.2.5 `getHash()` [4/4]

```
void LSH::getHash (
    unsigned int * hashIndices,
    unsigned int * identity,
    float * input,
    int numInputEntries,
    int numProbes )
```

Obtain hash indice given the (dense) input vector, using OpenCL.

Hash indice refer to the corresponding "row number" in a hash table, in the form of unsigned integer. This function will only be valid when an OpenCL implementation exists for that type of hashing, and when OpenCL is initialized (cLSH) for that hashing. The outputs indexing is defined as `hashIndicesOutputIdx(numHashFamilies, numProbes, numInputs, inputIdx, probelIdx, tb) (unsigned)(numInputs * numProbes * tb + inputIdx * numProbes + probelIdx)`.

Parameters

<i>hashIndices</i>	Hash indice for the batch of input vectors.
<i>identity</i>	Hash indice's corresponding identifications (sequential number starting 0) for this batch of inputs.
<i>input</i>	Input vectors concatenated.
<i>numInputEntries</i>	Number of input vectors.
<i>numProbes</i>	Number of probes per input.

The documentation for this class was generated from the following files:

- LSH.h
- LSH.cpp
- LSH_helpers.cpp
- LSH_init.cpp

3.4 LSHReservoirSampler Class Reference

```
#include <LSHReservoirSampler.h>
```

Public Member Functions

- void **restart** ([LSH](#) *hashFamIn, unsigned int numHashPerFamily, unsigned int numHashFamilies, unsigned int reservoirSize, unsigned int dimension, unsigned int numSecHash, unsigned int maxSamples, unsigned int queryProbes, unsigned int hashingProbes, float tableAllocFraction)
- [LSHReservoirSampler](#) ([LSH](#) *hashFam, unsigned int numHashPerFamily, unsigned int numHashFamilies, unsigned int reservoirSize, unsigned int dimension, unsigned int numSecHash, unsigned int maxSamples, unsigned int queryProbes, unsigned int hashingProbes, float tableAllocFraction)
- void [add](#) (int numInputEntries, int *dataIdx, float *dataVal, int *dataMarker)
- void [ann](#) (int numQueryEntries, int *dataIdx, float *dataVal, int *dataMarker, unsigned int *outputs, int k)
- void [add](#) (int numInputEntries, float *input)
- void [ann](#) (int numQueryEntries, float *queries, unsigned int *outputs, int k)
- void [showParams](#) ()
- void [checkTableMemLoad](#) ()
- [~LSHReservoirSampler](#) ()

Public Attributes

- cl_platform_id * **platforms**
- cl_device_id * **devices_gpu**
- cl_context **context_gpu**
- cl_program **program_gpu**
- cl_command_queue **command_queue_gpu**
- cl_device_id * **devices_cpu**
- cl_context **context_cpu**
- cl_program **program_cpu**
- cl_command_queue **command_queue_cpu**

3.4.1 Detailed Description

[LSHReservoirSampler](#) Class.

Providing hashtable data-structure and k-select algorithm. An [LSH](#) class instantiation is pre-required.

3.4.2 Constructor & Destructor Documentation

3.4.2.1 LSHReservoirSampler()

```
LSHReservoirSampler::LSHReservoirSampler (
    LSH * hashFam,
    unsigned int numHashPerFamily,
    unsigned int numHashFamilies,
    unsigned int reservoirSize,
    unsigned int dimension,
    unsigned int numSecHash,
    unsigned int maxSamples,
    unsigned int queryProbes,
    unsigned int hashingProbes,
    float tableAllocFraction )
```

Constructor.

Creates an instance of [LSHReservoirSampler](#).

Parameters

<i>hashFam</i>	An LSH class, a family of hash functions.
<i>numHashPerFamily</i>	Number of hashes (bits) per hash table, have to be the same as that of the hashFam.
<i>numHashFamilies</i>	Number of hash families (tables), have to be the same as that of the hashFam.
<i>reservoirSize</i>	Size of each hash rows (reservoir).
<i>dimension</i>	For dense vectors, this is the dimensionality of each vector. For sparse format data, this number is not used. (TBD)
<i>numSecHash</i>	The number of secondary hash bits. A secondary (universal) hashing is used to shrink the original range of the LSH for better table occupancy. Only a number \leq numHashPerFamily makes sense.
<i>maxSamples</i>	The maximum number incoming data points to be hashed and added.
<i>queryProbes</i>	Number of probes per query per table.
<i>hashingProbes</i>	Number of probes per data point per table.
<i>tableAllocFraction</i>	Fraction of reservoirs to allocate for each table, will share with other table if overflows.

3.4.2.2 ~LSHReservoirSampler()

```
LSHReservoirSampler::~~LSHReservoirSampler ( )
```

Destructor. Frees memory allocations and OpenCL environments.

3.4.3 Member Function Documentation

3.4.3.1 `add()` [1/2]

```
void LSHReservoirSampler::add (
    int numInputEntries,
    int * dataIdx,
    float * dataVal,
    int * dataMarker )
```

Adds input vectors (in sparse format) to the hash table.

Each vector is assigned ascending identification starting 0. For `numInputEntries > 1`, simply concatenate data vectors.

Parameters

<i>numInputEntries</i>	Number of input vectors.
<i>dataIdx</i>	Non-zero indice of the sparse format.
<i>dataVal</i>	Non-zero values of the sparse format.
<i>dataMarker</i>	Marks the start index of each vector in <i>dataIdx</i> and <i>dataVal</i> . Have an additional marker at the end to mark the (end+1) index.

3.4.3.2 `add()` [2/2]

```
void LSHReservoirSampler::add (
    int numInputEntries,
    float * input )
```

Adds input vectors (in dense format) to the hash table.

Each vector is assigned ascending identification starting 0. For `numInputEntries > 1`, simply concatenate data vectors.

Parameters

<i>numInputEntries</i>	Number of input vectors.
<i>input</i>	Concatenated data vectors (fixed dimension).

3.4.3.3 `ann()` [1/2]

```
void LSHReservoirSampler::ann (
    int numQueryEntries,
```



```

int * dataIdx,
float * dataVal,
int * dataMarker,
unsigned int * outputs,
int k )

```

Query vectors (in sparse format) and return top k neighbors for each.

Near-neighbors for each query will be returned in descending similarity. For numQueryEntries > 1, simply concatenate data vectors.

Parameters

<i>numQueryEntries</i>	Number of query vectors.
<i>dataIdx</i>	Non-zero indice of the sparse format.
<i>dataVal</i>	Non-zero values of the sparse format.
<i>dataMarker</i>	Marks the start index of each vector in dataIdx and dataVal. Have an additional marker at the end to mark the (end+1) index.
<i>outputs</i>	Near-neighbor identifications. The i_th neighbor of the q_th query is outputs[q * k + i]
<i>k</i>	number of near-neighbors to query for each query vector.

3.4.3.4 ann() [2/2]

```

void LSHReservoirSampler::ann (
    int numQueryEntries,
    float * queries,
    unsigned int * outputs,
    int k )

```

Query vectors (in dense format) and return top k neighbors for each.

Near-neighbors for each query will be returned in descending similarity. For numQueryEntries > 1, simply concatenate data vectors.

Parameters

<i>numQueryEntries</i>	Number of query vectors.
<i>queries</i>	Concatenated data vectors (fixed dimension).
<i>outputs</i>	Near-neighbor identifications. The i_th neighbor of the q_th query is outputs[q * k + i]
<i>k</i>	number of near-neighbors to query for each query vector.

3.4.3.5 checkTableMemLoad()

```

void LSHReservoirSampler::checkTableMemLoad ( )

```

Check the memory load of the hash table.

3.4.3.6 showParams()

```
void LSHReservoirSampler::showParams ( )
```

Print current parameter settings to the console.

The documentation for this class was generated from the following files:

- LSHReservoirSampler.h
- LSHReservoirSampler.cpp
- LSHReservoirSampler_helpers.cpp
- LSHReservoirSampler_init.cpp
- LSHReservoirSampler_misc.cpp
- LSHReservoirSampler_routines.cpp
- LSHReservoirSampler_segsort.cpp