

Università degli studi di Firenze
Scuola di Ingegneria

Histograms Representation of Name Entity Detector

Elena Sesoldi, Giulia Pellegrini
parallel18computing@gmail.com

February 26, 2018



Introduction

Task

Implementation C++

Sequential C++

OpenMP API

Experiments

OpenMP

Implementation Java

Apache Hadoop Framework

Results

Histograms

Conclusion



- Implementation of an algorithm that extracts Named Entities from tweet messages and computes for each different category of considered entity the corresponding histogram of occurrences



- ▶ Implementation of an algorithm that extracts Named Entities from tweet messages and computes for each different category of considered entity the corresponding histogram of occurrences
- ▶ 2 programming languages: C++ & Java



- ▶ Implementation of an algorithm that extracts Named Entities from tweet messages and computes for each different category of considered entity the corresponding histogram of occurrences
- ▶ 2 programming languages: C++ & Java
- ▶ A sequential implementation: C++



- ▶ Implementation of an algorithm that extracts Named Entities from tweet messages and computes for each different category of considered entity the corresponding histogram of occurrences
- ▶ 2 programming languages: C++ & Java
- ▶ A sequential implementation: C++
- ▶ 2 parallel implementations: OpenMP & Hadoop



Steps:

- ▶ dataset loading
- ▶ tokenization
- ▶ Named Entity prediction & map insertion (for loop)
- ▶ category splitting
- ▶ histograms' creation



LOADING



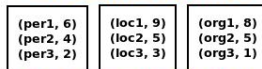
TOKENIZATION



NE PRED & MAP INSERTION



CATEGORY SPLITTING



HISTOGRAM CREATION





Steps:

- ▶ dataset loading
- ▶ tokenization
- ▶ Named Entity prediction & map insertion (for loop)
- ▶ category splitting
- ▶ histograms' creation



LOADING

Tweets1
Tweets2
Tweets3

TOKENIZATION

All Tweet
Words

NE PRED & MAP INSERTION

(entity1, 3)
(entity2, 1)
(entity3, 6)

CATEGORY SPLITTING

(per1, 6)
(per2, 4)
(per3, 2)

(loc1, 9)
(loc2, 5)
(loc3, 3)

(org1, 8)
(org2, 5)
(org3, 1)

HISTOGRAM CREATION





Parallel loop

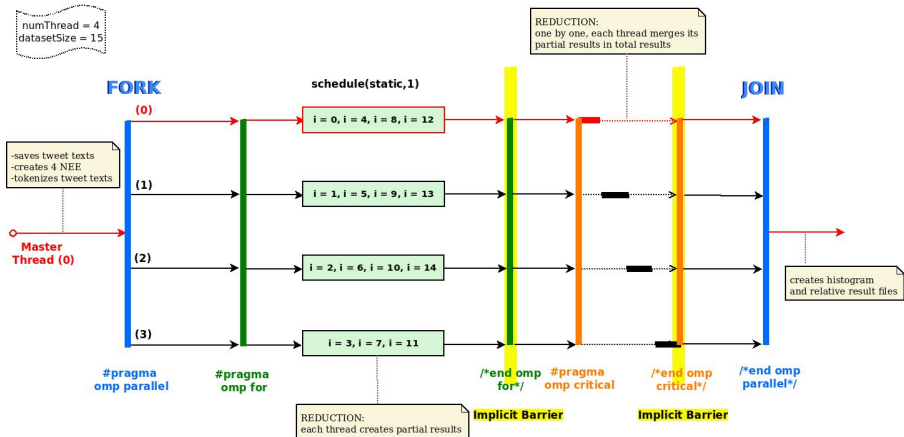
After creating the tokens as in the sequential version, the parallel part of the program starts with the directive `# pragma omp parallel`. Then with the directive `# pragma omp for` a for loop is executed in parallel and distributed among threads already existing in the parallel region.

Reduction

In the loop, a partial result map is defined and created for each thread. After the loop it is defined a critical section with the clause `# pragma omp critical` where each thread merges its partial result in a final common one.

C++ Implementation

OpenMp API





for loop scheduling

The schedule clause avails to control the distribution of the loop between threads.

- ***schedule(static, 1)***

In this case it divides the cycle into pre-defined blocks of given size (one) and the blocks are statically assigned to the threads sequentially following the identifier of the individual thread. Once the available threads are finished, the blocks are reassigned from the first thread.



for loop scheduling

The schedule clause avails to control the distribution of the loop between threads.

- ▶ ***schedule(static, 1)***

In this case it divides the cycle into pre-defined blocks of given size (one) and the blocks are statically assigned to the threads sequentially following the identifier of the individual thread. Once the available threads are finished, the blocks are reassigned from the first thread.

- ▶ ***schedule(dynamic, 1)***

A block is assigned to the first available thread to execute it. At the end of its execution that same thread can be used again to calculate another block.

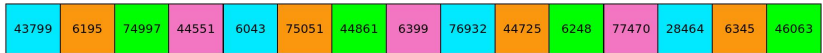


for loop scheduling

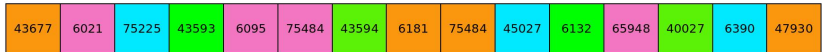


0 *Loop iteration index* 14 →

static, 1



dynamic, 1



Implementation

OpenMp Schedule



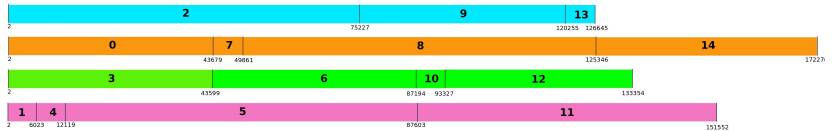
for loop scheduling



static, 1 time = 172169 ms



dynamic, 1 time = 172276 ms





We test our program with the same three datasets, repeating them several times, through setting a variable named *numRepetition*. In each execution we use a number of datasets equal to three multiplied by *numRepetition*.

We tried more parallel versions with different number of threads using *omp_set_num_threads(int num_threads)*, to see how execution times changes.

Rep	Seq	Omp1	Omp2	Omp3	Omp4
1	66533	66727	63783	44663	44529
2	133256	134112	72604	84968	66612
3	199776	200994	135204	127669	100748
4	266349	267882	147127	142258	123365
5	336083	334772	211378	184322	164915



Once the execution times have been detected we compute the relative **Speedup**

Rep	Sp1	Sp2	Sp3	Sp4
1	0.9970	1.0431	1.4896	1.4941
2	0.9936	1.8353	1.5683	2.0004
3	0.9939	1.4775	1.5647	1.9829
4	0.9942	1.8103	1.8722	2.1598
5	1.0039	1.5899	1.9861	2.2198



Once the execution times have been detected we compute the relative **Speedup**

Rep	Sp1	Sp2	Sp3	Sp4
1	0.9970	1.0431	1.4896	1.4941
2	0.9936	1.8353	1.5683	2.0004
3	0.9939	1.4775	1.5647	1.9829
4	0.9942	1.8103	1.8722	2.1598
5	1.0039	1.5899	1.9861	2.2198

Steps:

- ▶ dataset loading
- ▶ tweet texts' extraction
- ▶ tweet files' creation
- ▶ Map Reduce Process
- ▶ category splitting
- ▶ histograms' creation



LOADING

Dataset1
Dataset2
Dataset3

TWEET EXTRACTION

Tweet1
Tweet2
Tweet3

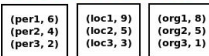
TWEET FILE CREATION



MAP REDUCE PROCESS

(entity1, 3)
(entity2, 1)
(entity3, 6)

CATEGORY SPLITTING



HISTOGRAM CREATION





MapReduce

The task of the *MapReduce* framework is counting the number of named entities in a collection of *.xml* files containing tweet messages.

- ▶ **Mapper**

Processes line by line the input files, looking for the entities. Once obtained the chunk set of entities, it loops on it creating the key-value pairs. The key is a string of the entity combined with its category and the value is an integer always setting at one.

- ▶ **Reducer**

Takes all the key-value pairs for a given real-word object and adds up the values, generating a single key-value pair for each entity, that is written on the context.

Java Implementation

Apache Hadoop Framework



15

Input

Map

Shuffle

Reduce

Output

`("name entity _CATEGORY", 1)`
key: Text value: IntWritable

`("name entity _CATEGORY", 1)`
key: Text value: IntWritable

Australia's competition
US, Ginnifer Goodwin
John Brumby in Victoria

Map

`("australia _LOCATION", 1)`
`("us _LOCATION", 1)`
`("ginnifer goodwin _PERSON", 1)`
`("john brumby _PERSON", 1)`
`("victoria _LOCATION", 1)`

Intel focuses on the
Australia vs England live
NEW YORK -- The Dow end

Map

`("intel _ORGANIZATION", 1)`
`("australia _LOCATION", 1)`
`("england _LOCATION", 1)`
`("new york _LOCATION", 1)`

Reduce

Thanks Harry Potter for
Anne Hataway to co-host
SYDNEY - police treated

Map

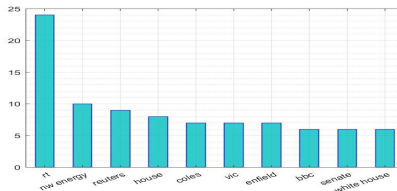
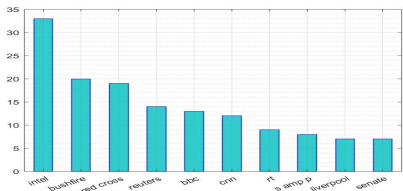
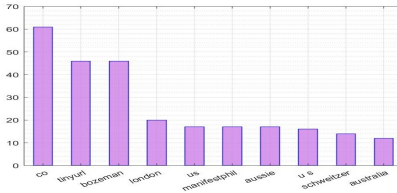
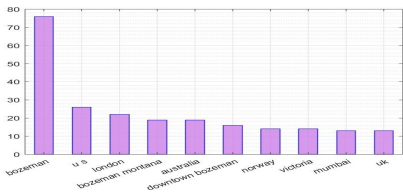
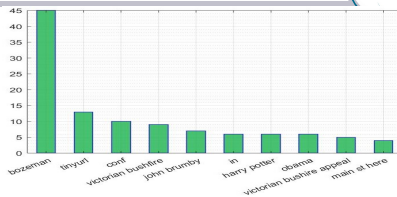
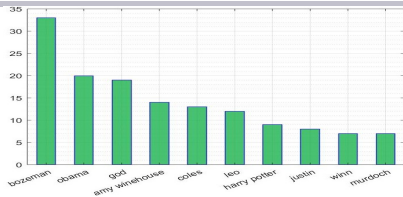
`("harry potter _PERSON", 1)`
`("anne hataway _PERSON", 1)`
`("sydney _LOCATION", 1)`

`("australia _LOCATION", 2)`
`("us _LOCATION", 1)`
`("ginnifer goodwin _PERSON", 1)`
`("john brumby _PERSON", 1)`
`("victoria _LOCATION", 1)`
`("intel _ORGANIZATION", 1)`
`("england _LOCATION", 1)`
`("new york _LOCATION", 1)`
`("harry potter _PERSON", 1)`
`("anne hataway _PERSON", 1)`
`("sydney _LOCATION", 1)`

Histograms



16





Conclusions

- ▶ three different implementations of the creation of histograms that represent the frequency of named entity in tweets.
- ▶ two versions written in C++, sequential and with the shared memory OpenMP API, and one in Java with the distributed memory Apache Hadoop tool.
- ▶ for the extraction of named entities we benefit of a library for each programming language; the histograms are different because these have distinct NER methodologies, although both are based on Machine Learning.
- ▶ we test our programs to see the benefits of the parallelization. For C++ we can evaluate the parallelism computing the speed up, and we obtained the best value equals to 2.2198 with four threads and fifteen total datasets



Applications

- ▶ Information Extraction
- ▶ Marketing Initiatives
- ▶ Information Security

A decorative graphic consisting of multiple overlapping, flowing lines in shades of light blue and white. The lines curve from the upper left towards the lower right, creating a sense of movement and depth. Some lines have a slight gradient, and there are small, faint white dots scattered along the curves.

Thank you for the attention