

Università degli studi di Firenze  
Scuola di Ingegneria

# Most Common N-Grams in the English and the Italian Languages

Elena Sesoldi, Giulia Pellegrini  
parallel18computing@gmail.com

February 26, 2018



## Introduction

Task

## Java Implementation

Sequential

Java Threads

## Experiments

Java Thread

## Results

Results

## Conclusion



- Implementation of an algorithm that computes the number of N-gram occurrences

N = 2

It was all Pride

*word  
bigrams*

it-was  
was-all  
all-pride

N = 3

It was all Pride

*word  
trigrams*

it-was-all  
was-all-pride



- Implementation of an algorithm that computes the number of N-gram occurrences

N = 2

It was all Pride

*word  
bigrams*

it-was  
was-all  
all-pride

N = 3

It was all Pride

*word  
trigrams*

it-was-all  
was-all-pride

- Programming language: [Java](#)



- Implementation of an algorithm that computes the number of N-gram occurrences

N = 2

It was all Pride

*word  
bigrams*

it-was  
was-all  
all-pride

N = 3

It was all Pride

*word  
trigrams*

it-was-all  
was-all-pride

- Programming language: **Java**
- A sequential implementation



- Implementation of an algorithm that computes the number of N-gram occurrences

N = 2

It was all Pride

*word  
bigrams*

it-was  
was-all  
all-pride

N = 3

It was all Pride

*word  
trigrams*

it-was-all  
was-all-pride

- Programming language: **Java**
- A sequential implementation
- A parallel implementation: **Java Threads**



### Steps:

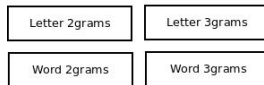
- ▶ book loading
- ▶ Hashmap creation
- ▶ Hashmap sorting  
& result sorting



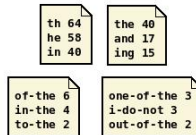
### LOADING

Lines Book1  
Lines Book2  
Lines Book3  
Lines Book4

### HASHMAP CREATION



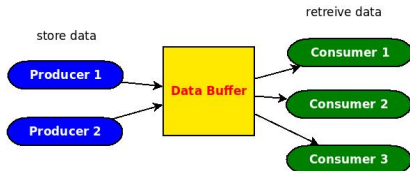
### HASHMAP SORTING & RESULT STORING



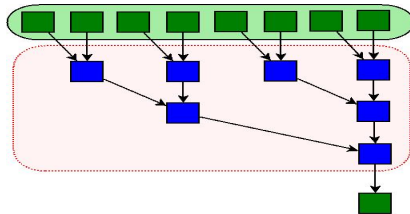


## Pattern involved

- Producer Consumer Pattern



- Reduction Pattern







- ▶ We use two different classes that implements **Runnable** interface, one named **Producer** and the other **Consumer**



- ▶ We use two different classes that implements **Runnable** interface, one named **Producer** and the other **Consumer**
- ▶ The ***Producer Thread*** loads each book's line from the directory in a **BlockingQueue**.



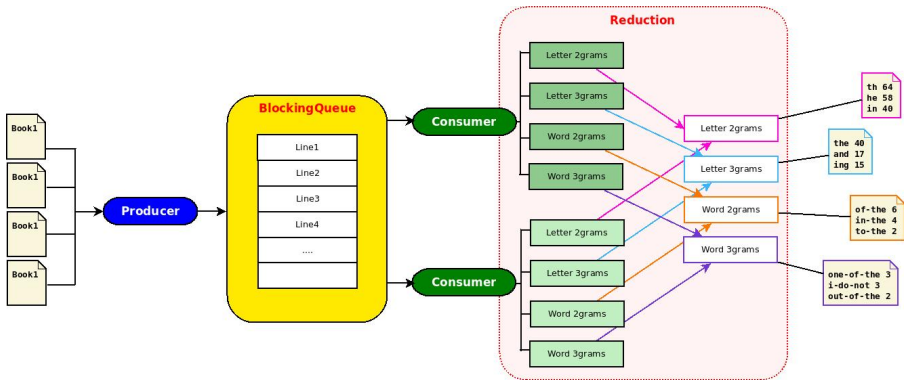
- ▶ We use two different classes that implements **Runnable** interface, one named **Producer** and the other **Consumer**
- ▶ The **Producer Thread** loads each book's line from the directory in a **BlockingQueue**.
- ▶ The **Consumer Thread** takes the lines from the same **BlockingQueue** one by one and compute the partial results of the current line, which are four **HashMap** for the N-grams of every type.



- ▶ We use two different classes that implements **Runnable** interface, one named **Producer** and the other **Consumer**
- ▶ The **Producer Thread** loads each book's line from the directory in a **BlockingQueue**.
- ▶ The **Consumer Thread** takes the lines from the same **BlockingQueue** one by one and compute the partial results of the current line, which are four **HashMap** for the N-grams of every type.
- ▶ The **Consumer** share four **ConcurrentHashMap** and as soon as the **BlockingQueue** is empty and the **Producer** has completed its task, the partial results are completed and the consumers can add them in the corresponding shared **ConcurrentHashMap**.

# Java Implementation

## Steps





We test our program with two datasets of books, downloaded from the Gutenberg Project <sup>1</sup>, repeating them several times. For the parallel versions the number of producer is always one for English datasets and one for the Italian ones. The number of consumers instead varies from one to four for each dataset to see the benefit of parallelism..

| Size | Seq    | Par1  | Par2  | Par3  | Par4  |
|------|--------|-------|-------|-------|-------|
| 2.7  | 2902   | 3168  | 2661  | 2483  | 2573  |
| 5.4  | 5564   | 5518  | 4249  | 4203  | 4185  |
| 10.8 | 10563  | 10430 | 7479  | 7334  | 6177  |
| 21.6 | 22269  | 21109 | 13539 | 11326 | 11121 |
| 43.2 | 54458  | 41346 | 27533 | 22001 | 21988 |
| 86.4 | 156008 | 83835 | 55353 | 43905 | 42066 |

---

<sup>1</sup><https://www.gutenberg.org/>



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

| Rep  | Sp1    | Sp2    | Sp3    | Sp4    |
|------|--------|--------|--------|--------|
| 2.7  | 0.9160 | 1.0906 | 1.1687 | 1.1279 |
| 5.4  | 1.0083 | 1.3095 | 1.3238 | 1.3295 |
| 10.8 | 1.0128 | 1.4124 | 1.4403 | 1.7101 |
| 21.6 | 1.0055 | 1.6448 | 1.9662 | 2.0024 |
| 43.2 | 1.3171 | 1.9779 | 2.4753 | 2.4767 |
| 86.4 | 1.8609 | 2.8184 | 3.5533 | 3.7086 |



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

| Rep  | Sp1    | Sp2    | Sp3    | Sp4    |
|------|--------|--------|--------|--------|
| 2.7  | 0.9160 | 1.0906 | 1.1687 | 1.1279 |
| 5.4  | 1.0083 | 1.3095 | 1.3238 | 1.3295 |
| 10.8 | 1.0128 | 1.4124 | 1.4403 | 1.7101 |
| 21.6 | 1.0055 | 1.6448 | 1.9662 | 2.0024 |
| 43.2 | 1.3171 | 1.9779 | 2.4753 | 2.4767 |
| 86.4 | 1.8609 | 2.8184 | 3.5533 | 3.7086 |

- ▶ No speedup
- ▶ Linear speedup
- ▶ Maximum speedup



# Results



| Letter2grams         | Letter3grams          | Word2grams             | Word3grams                |
|----------------------|-----------------------|------------------------|---------------------------|
| <b>th</b><br>2062144 | <b>the</b><br>1284320 | <b>of-the</b><br>99776 | <b>one-of-the</b><br>3392 |
| <b>he</b><br>1866208 | <b>and</b><br>566816  | <b>in-the</b><br>71328 | <b>i-do-not</b><br>3392   |
| <b>in</b><br>1300448 | <b>ing</b><br>510048  | <b>to-the</b><br>41952 | <b>out-of-the</b><br>3232 |

| Letter2grams         | Letter3grams         | Word2grams             | Word3grams                 |
|----------------------|----------------------|------------------------|----------------------------|
| <b>on</b><br>1063040 | <b>ent</b><br>435072 | <b>e-di</b><br>15328   | <b>per-lo-piu</b><br>1472  |
| <b>er</b><br>963808  | <b>ell</b><br>415040 | <b>per-la</b><br>12000 | <b>poco-a-poco</b><br>1472 |
| <b>en</b><br>914400  | <b>ion</b><br>332768 | <b>che-si</b><br>11840 | <b>a-poco-a</b><br>1472    |



## Conclusions

- ▶ two different implemetations of N-grams occurence calculus: bigrams and trigrams of letters and words.
- ▶ the sequential version takes one line of a book at a time and processes it.
- ▶ the parallel version with explicit *Java Threads*, *Reduction Pattern* and the *Producer and Consumer Pattern* computes the same work but the execution time is quicker.
- ▶ with the increasing of the dimension of the dataset the execution time of the parallel version is three times less of sequential one, indeed the speedup has a value of 3.7086.



## Application

- ▶ Text Prediction
- ▶ Spelling Correction
- ▶ Language Identification
- ▶ Plagiarism Detection

An abstract graphic consisting of multiple flowing, curved lines in shades of light blue and white. The lines originate from the left and curve towards the right, creating a sense of motion and fluidity. Some lines have small, glowing white dots or sparkles along their length. The overall shape is reminiscent of a stylized wave or a plume of smoke.

Thank you for the attention