UNIVERSITÀ
DEGLI STUDI
FIRENZE

DIEF
Dipartimento di
Ingegneria Industriale

# Bigram/Trigram Generation

## Parallel Computing Course

**Antonio Castellucci**
**Michela Crulli**

In this project, we realize a program which generates **bigrams** and **trigrams** with three different approaches:

- a **sequential** implementation (Java)
- a **parallel** implementation (Java threads)
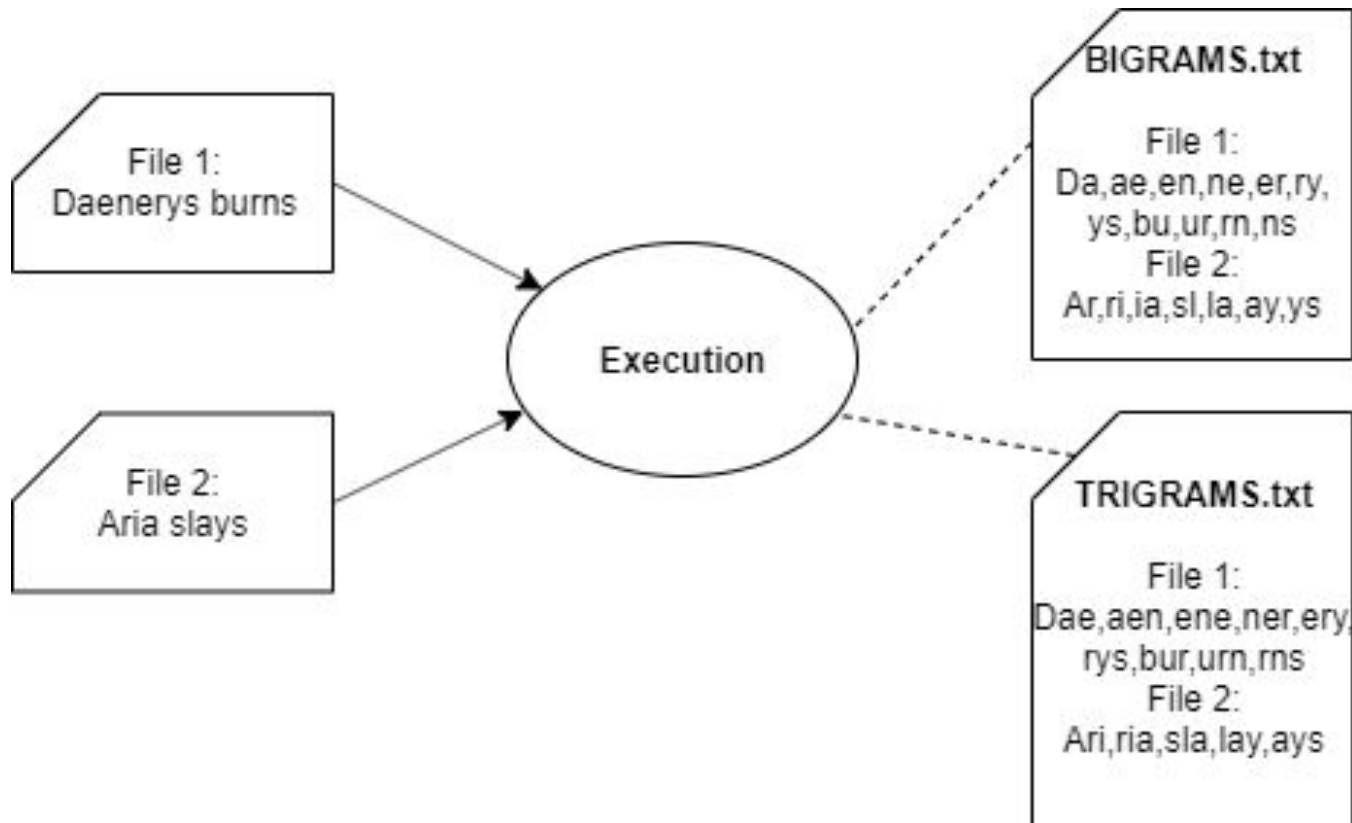- a **distributed** implementation (Hadoop)

- What are n-grams?

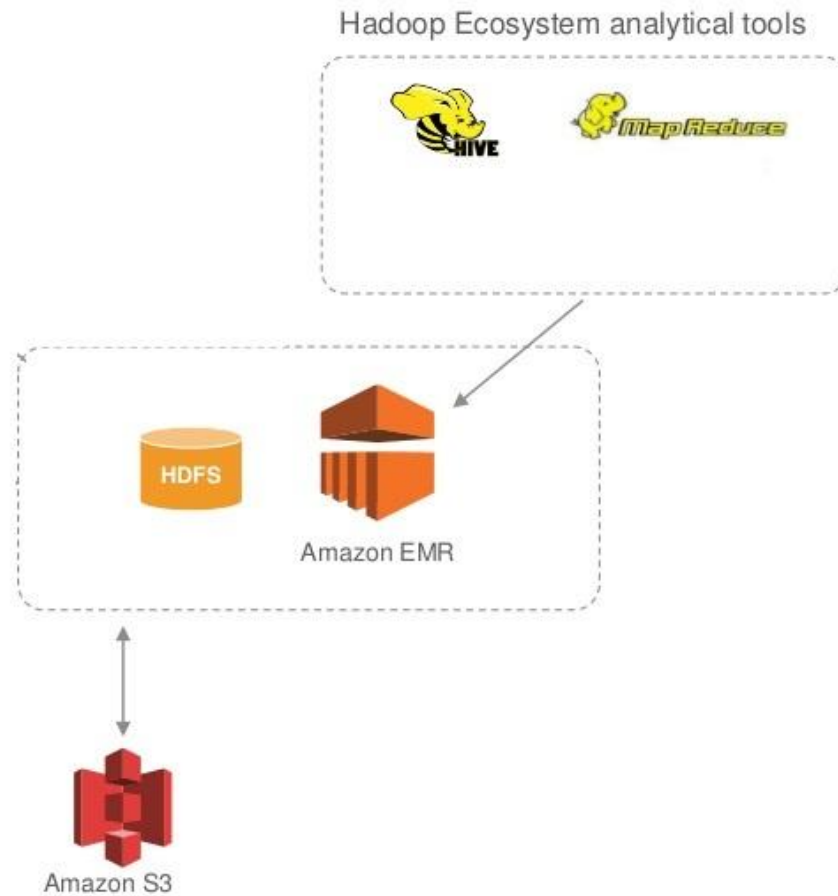An n-gram is a **contiguos sequence** of n items (letters, words, syllables or phonemes) from a given text or speech

- Why?

  - make new **word predictions**
  - make **spelling error corrections**

- Input output structure

- Hadoop & Amazon Web Service (AWS)



Hadoop Ecosystem analytical tools

HIVE  MapReduce

HDFS  Amazon EMR

Amazon S3

We use:

- a collection of **3000 books** (1.1 GB, from Gutemberg's Project) and on the collection's subsets:

  - **2000** (abt 800 megabyte)
  - **1000** (abt 400 megabyte)
  - **500** (abt 80 megabyte)
  - **150** (abt 50 megabyte)
  - **100** (abt 40 megabyte)
  - **50** (abt 20 megabyte)

The size of each book is about **360 KB**

- Sequential implementation

---
**Algorithm 1**

---
**Input:** dataset (folder with books) **Output:** bigram.txt and trigram.txt

**for** book in dataset

    **for** line in book

        $tmp \leftarrow lineWithoutSpace$

        **for** i = 0 to tmp.length

            $ngrams2.append(j + 2).append(",")$

            $ngrams3.append(j + 3).append(",")$

        **end for**
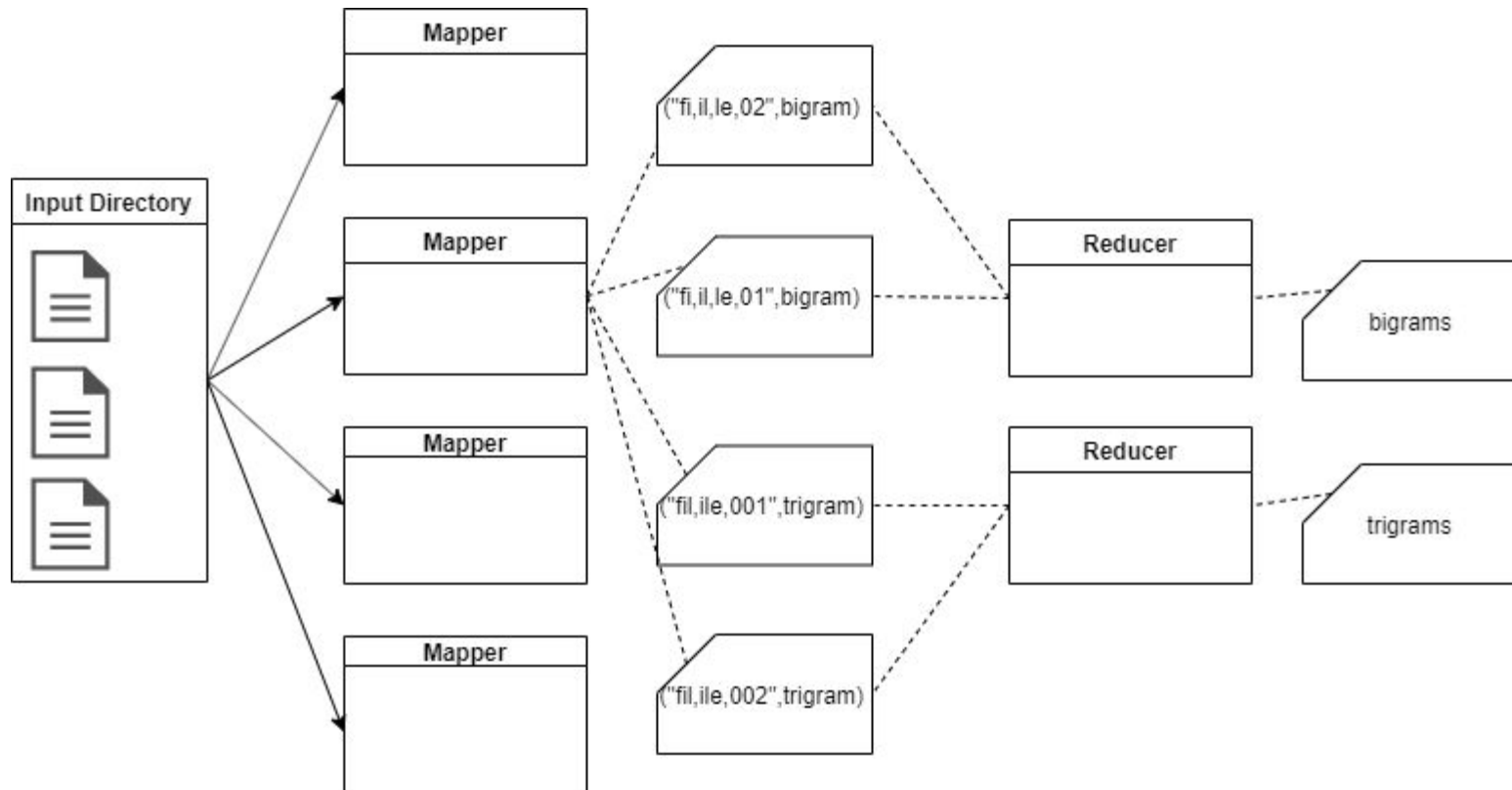
    **end for**

    write ngrams2 on ngrams2.txt

    write ngrams3 on ngrams3.txt

**end for**

---

- ## Parallel implementation

- Hadoop implementation

## Test platform

MSI GS65 with:

- i7-8750H 6 CORES 12 THREAD

- 16 GB 2400 MHz DDR4 RAM

- M2 SSD SAMSUNG MZVLW256HEHP

- Parellel implementation setup

  - 12 threads and 1 book to process each time

- AWS

  - Emr cluster:

    - m4.xlarge 1 master 4 slaves, Amazon Hadoop version 2.8.5

  - Storage:

    - S3 bucket

- All versions compared respect to the number of books

- ## AWS
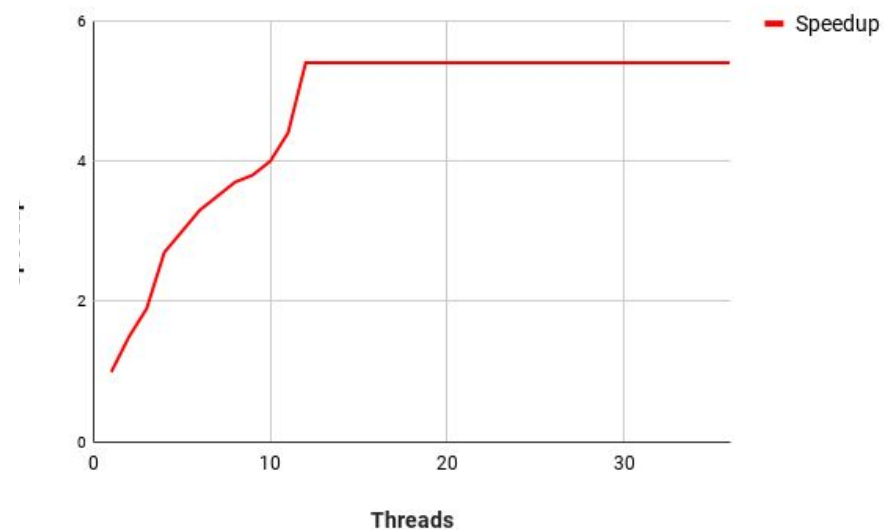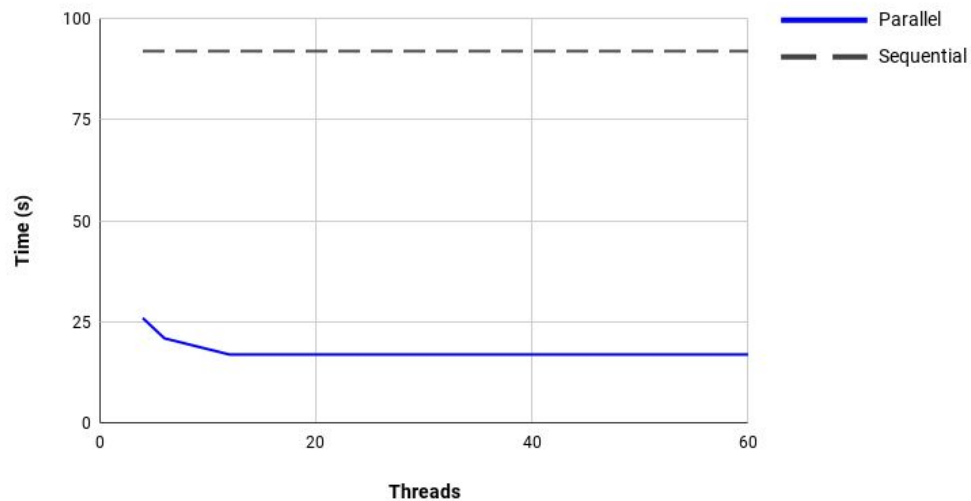
- AWS

- AWS vs Standalone respect to the number of books

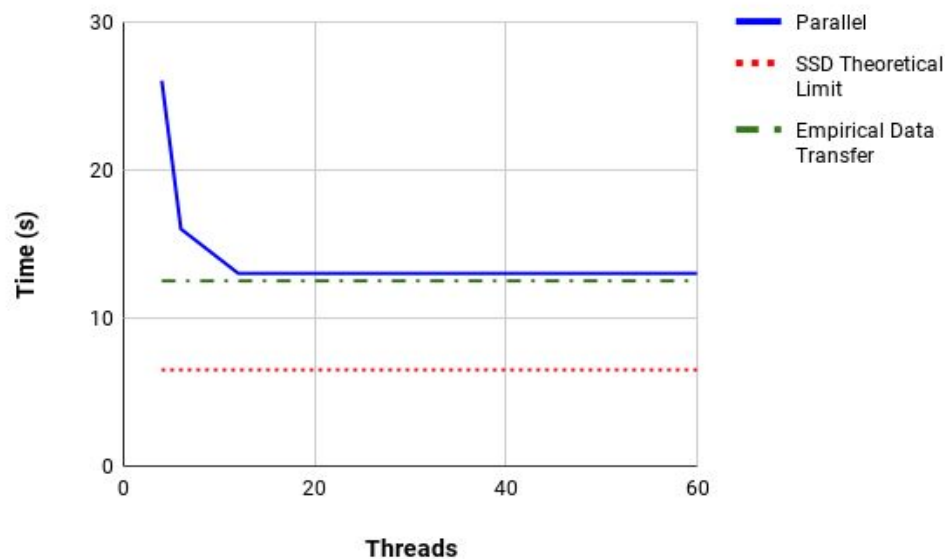- Parallel vs sequential respect to number of books

- Parallel vs sequential respect to number of threads

- A deeper analysis

- Classic case of Hadoop usage:

  - perform local computations over huge amounts of input data while returning relatively small result set, which makes our case not suitable to be implemented with Hadoop

- A program for n-grams counting rather than writing them only, probably we would have been able to see the potential of this distributed technology

# Thanks for the attention!