

# Big Data Analytics with Scala and Apache Spark

M.Sc. course on “Technologies for Big Data Analysis” - Assignment 3

CHRISTOS BALAKTSIS (506) and VASILEIOS PAPASTERGIOS (505), Aristotle University, Greece

## 1 Introduction

The current document is a technical report for the third programming assignment in the M.Sc. course on *Technologies for Big Data Analysis*, offered by the *DWS M.Sc Program*<sup>1</sup> of the Aristotle University of Thessaloniki, Greece. The course is taught by Professor Apostolos Papadopoulos<sup>2</sup>. The authors attended the course during their first year of Ph.D. studies at the Institution.

The assignment contains 4 sub-problems and is part of a series, comprising 3 programming assignments on the following topics:

*Assignment 1* Multi-threading Programming and Inter-Process Communication

*Assignment 2* The Map-Reduce Programming Paradigm

*Assignment 3* Big Data Analytics with Scala and Apache Spark

In this document we focus on Assignment 3 and its 4 sub-problems. We refer to them as *problems* in the rest of the document for simplicity. The source code of our solution has been made available at the following public repository in the GitHub platform: <https://github.com/Bilpaster/big-data-playground>.

**Roadmap.** The rest of our work is structured as follows. We devote one section to each one of the 4 problems. That means problems 1, 2, 3 and 4 are presented in section 2, section 3, section 4 and section 5 respectively. For each problem, we first provide the problem statement, as given by the assignment. Next, we thoroughly present the reasoning and/or methodology we have adopted to approach the problem and devise a solution. Wherever applicable, we also provide insights about the source code implementation we have developed. Finally, we conclude our work in section 6. The appendix includes the evaluation results for any issues that necessitated them.

## 2 Problem 1: Word Length Analytics

We discuss here the first problem of the assignment. The main target of the assignment is to get familiar with a simple task leveraging Apache Spark and Scala programming language. This is a WordCount problem’s variation.

### 2.1 Problem Statement

Implement an **Apache Spark** (Scala) program, a variation of the **word-count** problem, to compute the average length of words that start with a specific letter (a-z). The program should sort the results based on the average length, printing first the letters with higher average length. For example:

k – 8.2

a – 5.6

b – 4.8

<sup>1</sup><https://dws.csd.auth.gr/>

<sup>2</sup><https://datalab-old.csd.auth.gr/~apostol/>

Note that you may opt for preprocessing the input text, e.g., transforming all letters to lowercase or ignoring words that start with a number.

## 2.2 Proposed approach

**2.2.1 Setting.** Our implementation is run and tested in a Linux environment with 12 cores, using the Scala programming language version 2.13.15 and Apache Spark version 3.5.3. We have used SBT as the build tool of our solution. We have used the Software Development Kit (JDK) version 11.0.11. The source code is developed in IntelliJ IDEA Community Edition 2021.1.1 and managed using SBT as the build tool. All dependencies of the project can be found in the `build.sbt` file located at the root folder of the project. The project is compiled and executed directly from the IntelliJ IDEA.

To run the project, open the `src/main/scala/Task1WordLength.scala` file in IntelliJ IDEA, and execute the main method. After successful execution of the program, the results can be viewed under the `output/task1` directory.

**2.2.2 Implementation.** The proposed approach leverages the Apache Spark computing framework to calculate the average word length for each one of the 26 letters in the Latin alphabet. In particular, the input text is read and tokenized into words. As part of a preprocessing step, all letters are transformed to lowercase and all words that start with a number are dropped.

In order to compute the average length for each initial letter, we leverage the RDD API of Apache Spark. More specifically, we use a `PairRDD` that stores composite key-value pairs in the form `(letter, (word_length, 1))`. The key of each pair is the initial letter of the respective word (`letter`). The value of each pair is a pair itself, having as (`false`) key the length of the word (`word_length`) and as (`false`) value a unary value (`1`). We use the term *false* for the key and value of the inner pair, since we do not utilize them as an actual key-value pair. As a matter of fact, we treat the inner pair as a plain tuple of values.

Subsequently, we group all pairs by key (i.e., initial letter) and compute the following two intermediate results *per initial letter*:

- (1) sum of word lengths that start with the respective letter, by adding up all `false` keys of the inner pair, and
- (2) total number of words that start with the respective letter, by adding up all unary values (`false` values) of the inner pair.

Having computed these intermediate results for each initial letter, the average length can be computed as the fraction

$$\text{average word length}_{\text{letter}} = \frac{\text{sum of word lengths}_{\text{letter}}}{\text{total number of words}_{\text{letter}}}, \forall \text{letter} \in \text{Latin Alphabet}$$

As a final step, the results are sorted in descending order w.r.t. the average length and saved to a text file. By default, after successful execution of the program, the results can be found under the `output/task1` directory.

**2.2.3 Evaluation.** The experiments were executed using the dataset `SherlockHolmes.txt`. The results are listed below in a 10-column format and graphically depicted in Figure 1.

c: 7.19	p: 7.01	v: 5.98	z: 5.74	k: 5.37	m: 5.14	w: 4.28	y: 3.72	x: 3.41
e: 7.11	r: 6.87	g: 5.93	u: 5.61	l: 5.34	n: 4.81	h: 3.81	t: 3.64	o: 3.01
q: 7.01	d: 6.37	s: 5.8	j: 5.59	f: 5.17	b: 4.54	a: 3.74	i: 3.46	

## 3 Problem 2: Movie Analytics

The second problem focuses on producing analytic insights from an IMDB dataset.

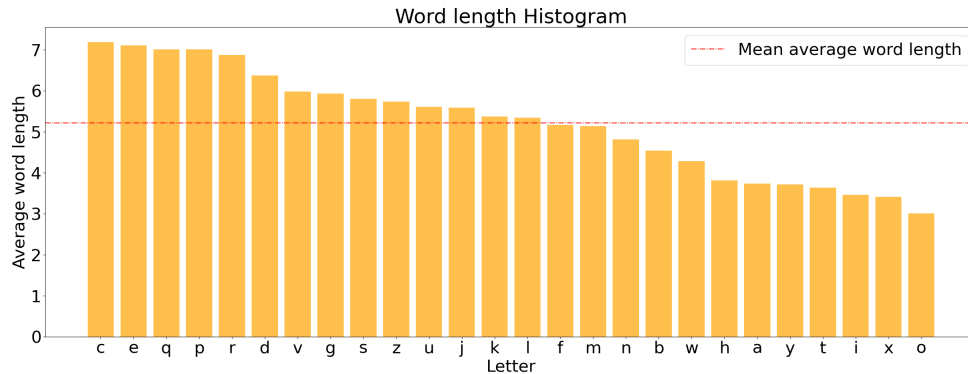


Fig. 1. A histogram depicting the average word length (y-axis) for every letter of the Latin alphabet (x-axis). The mean average length is shown as a horizontal red dashed line.

### 3.1 Problem Statement

Implement a **MapReduce** program to perform analytics tasks on an IMDB dataset about movies. The utter goal of your analysis would be to extract useful insights from the available movie data that will assist the IMDB team provide better recommendations for movies, based on their genre and/or country. In particular, the dataset (`movies.csv`) contains the following fields:

- `imdbID`: unique identifier of the movie in the IMDB database
- `title`: the movie title
- `year`: the year the movie was first released
- `duration`: the duration of the movie
- `genre(s)`: the genre or genres in which the movie is classified
- `premier date`: the date of the first showing of the movie
- `score`: the IMDB score of the movie
- `country/-ies`: the country or countries the movie was produced in

You are asked to implement Map-Reduce source code in Java programming language for the following analytics tasks:

- Calculate the total duration of all movies per country. Note that in case multiple countries are recorded for a movie, the respective duration should be counted for all of them separately.
- Calculate the total number of movies per year and genre, having IMDB score over 8. For movies that have more than one genre, the sum should be separate for each genre.

### 3.2 Proposed approach

**3.2.1 Setting.** Our implementation is run and tested in a Linux environment with 12 cores, using the Java programming language. We have used the Java Development Kit (JDK) version 11.0.11. The source code is developed in IntelliJ IDEA Community Edition 2021.1.1 and managed using Maven as the build tool.

The project's dependencies, including Hadoop libraries, are defined in the `pom.xml` file located in the root of the repository. The project is compiled and executed directly from IntelliJ IDEA.

To run the project, open the `MovieAnalyticsMaster` class in IntelliJ IDEA, and execute the main method. You will need to specify the following command-line arguments:

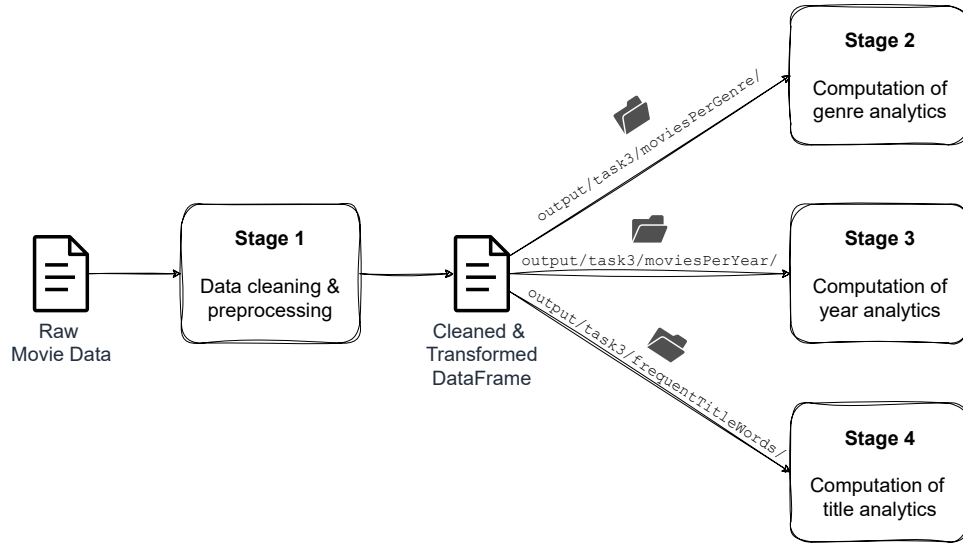


Fig. 2. The MapReduce solution architecture for the movie analytics problem. Two separate jobs are executed to solve the two tasks, namely duration per country (green) and movies per year & genre (yellow). The two tasks also create separate subdirectories for writing the results.

- `<input_path>` specifies the directory containing the input text files, located at `map-reduce/movieAnalytics/input`.
- `<output_path>` specifies the directory where the MapReduce output will be written, which will be created in the out folder.

IntelliJ IDEA will handle the compilation and execution automatically when the main method is run. Make sure to configure the input and output paths as required for your specific run. Note that the command-line arguments must follow the specified order and format. If any of the arguments are missing or invalid, the program will terminate with an appropriate error message.

**3.2.2 Implementation.** The proposed approach leverages the MapReduce programming paradigm to compute the analytical insights for the two tasks. Figure 3 depicts the architecture of our solution as a diagram. We opt for executing two separate jobs; one for each task presented by the problem statement.

*Task 1: Duration per country.* One map-reduce cycle is enough to handle this task, as depicted in the top part of Figure 3 in green color. The map function parses the CSV file line by line and extracts the useful fields from each line. In this case, the useful fields are the country (or countries) the movie was produced and the movie duration, i.e., the fourth and ninth fields in the input line respectively. We employ more complex logic to handle cases where there are multiple countries in a single movie. In particular, we parse again the field and tokenize it into the separate countries, producing a key-value pair for each country-duration pair within a movie. The reduce function is, then, trivial; it just adds up the durations (value) per country (key) and outputs the results. The interested reader can refer to the `CSVProcessor.java` (mapper), `AnalyticsEngine.java` (reducer) and `MovieAnalyticsMaster.java` (driver) files for the source code implementation of our solution.

*Task 2: Movies per year & genre w.r.t. score constraint.* The task is similar to the first one, with the only difference lying in the fields extracted from the input CSV line. In this case, we are interested about the year, the genre (or

genres) and score fields, i.e., the third, fifth and seventh fields in the input CSV line. The map function produces key-value pairs in the form (`composite_key`, 1), where the composite key consists of the year and the genre of the respective movie concatenated by an underscore, e.g., `2024_action`. We handle multiple genres per movie similarly with the multiple countries in task 1. We employ the same trivial reducer that adds up the values (ones) per (composite) key, as shown in the bottom part of Figure 3 in yellow color.

**3.2.3 Evaluation.** Our solution is tested using the provided IMDB dataset, namely `movies.csv`. We list the execution results in Appendix A.

## 4 Problem 3: Movie Analytics

The third problem focuses on producing analytic insights into a movies dataset.

### 4.1 Problem Statement

In this problem, you are provided with a dataset that contains information about movies. Each record represents a movie and includes three main columns (attributes), namely `movieId`, `title` and `genres`. An example record of the dataset is the following:

```
1,Toy Story (1995),Adventure|Animation|Children|Comedy|Fantasy
```

Note that a movie can belong to multiple genres (e.g., adventure, fantasy, and so on). The genres are separated by a vertical line (|) in the respective field. Your Apache Spark (Scala) application should compute the following analytics:

- (1) **How many movies are there for each genre?** If a movie belongs to multiple genres, it should be counted to all these genres. Sort the results by the name of genre in alphabetical order.
- (2) **How many movies have been filmed per year?** Note that the year a movie was filmed is currently encoded into a composite movie title, e.g., “Toy Story (1995)”. Show the top 10 years with the most movies filmed within them.
- (3) **Which are the words that appear at least 10 times in the titles of the movies, and what is their total frequency?** You can ignore words that have less than 4 characters. Sort the results, showing first the words with the higher frequency.

### 4.2 Proposed approach

**4.2.1 Setting.** Our implementation is run and tested in a Linux environment with 12 cores, using the Scala programming language version 2.13.15 and Apache Spark version 3.5.3. We have used SBT as the build tool of our solution. We have used the Software Development Kit (JDK) version 11.0.11. The source code is developed in IntelliJ IDEA Community Edition 2021.1.1 and managed using SBT as the build tool. All dependencies of the project can be found in the `build.sbt` file located at the root folder of the project. The project is compiled and executed directly from the IntelliJ IDEA.

To run the project, open the `src/main/scala/Task3MovieAnalytics.scala` file in IntelliJ IDEA, and execute the main method. After successful execution of the program, the results can be viewed under the `output/task3` directory.

**4.2.2 Implementation.** We leverage Apache Spark *DataFrames* API to calculate the aforementioned analytical queries. Our proposed approach comprises the following four stages:

- Stage 1:* Data cleaning and preprocessing
- Stage 2:* Computation of genre analytics (query 1)
- Stage 3:* Computation of year analytics (query 2)
- Stage 4:* Computation of title analytics (query 3)

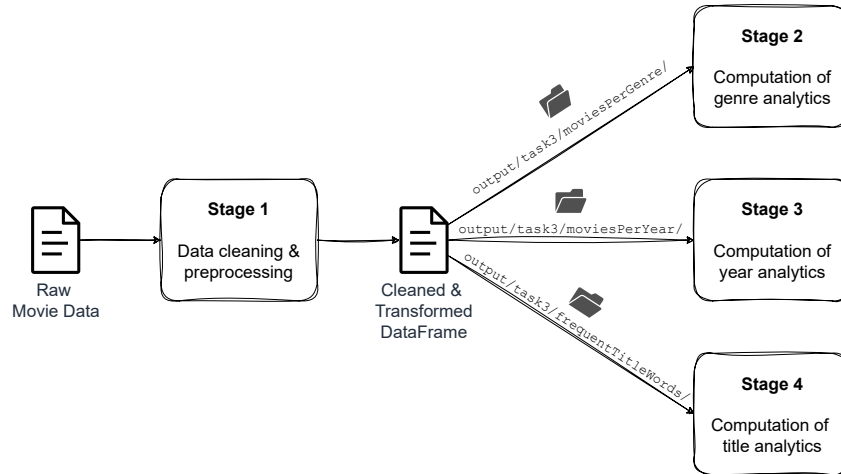


Fig. 3. The MapReduce solution architecture for the movie analytics problem. Two separate jobs are executed to solve the two tasks, namely duration per country (green) and movies per year & genre (yellow). The two tasks also create separate subdirectories for writing the results.

Figure 3 graphically depicts the four stages, as well as their assembly in a workflow. We elaborate on each stage separately in the rest of the Subsection.

*Stage1: Data Cleaning and preprocessing.* In this stage, we address data quality issues that are present in the raw movie data. Missing production year, missing genres and trailing whitespaces in the title string are the detected data quality issues. To address these issues, we explicitly mark the movie year as “N/A” (not available) & movie genre as “(no genre listed)” and trim the composite title strings, respectively. Subsequently, we transform the data representation to a more efficient form, by extracting the actual title and year to separate fields (instead of a composite one) and splitting the genres composite string field to an array of genres. Throughout the whole process, we **use solely the DataFrame API** of Apache Spark, leveraging also the pattern matching functionality of Scala. A small sample of the cleaned and transformed DataFrame shown in Figure 3 is depicted below:

movieId	title	genres	year
1	Toy Story	[Adventure, Anima...	1995
2	Jumanji	[Adventure, Child...	1995
3	Grumpier Old Men	[Comedy, Romance]	1995
4	Waiting to Exhale	[Comedy, Drama, R...	1995
5	Father of the Bri...	[Comedy]	1995

*Stage 2: Computation of genre analytics.* In this stage, we use the cleaned and transformed data to compute analytical insights into the movie genres. In particular, we use the DataFrame API (`select()`, `group_by()`, etc.) to calculate the number of movies that belong to each genre. An implementation detail of our solution lies in handling the cases where a movie belongs to multiple genres. To ensure the integrity of the computed analytics, we opt for using the `explode()` function of the DataFrame API, to transform the array of genres to separate

genre items. The analytics are, then, computed on top of this extra transformation, enabling us to include **a movie into the calculations of all the genres it belongs**.

*Stage 3: Computation of year analytics.* In this stage, we use the cleaned and transformed data to compute analytical insights into the movie production years. Similarly to Stage 2, we leverage exclusively the DataFrame API to calculate the total number of movies filmed each year. In this case, the `group_by` operation is more straightforward, since the cleaned and transformed DataFrame contains a dedicated column for the movie production year.

*Stage 4: Computation of title analytics.* Similarly to Stage 2, we use operations available from the DataFrame API to calculate the most frequently used words in the movie titles. We leverage the `split(" ")` and `explode()` functions to transform titles, first, into an array of words and, then, into separate words that are finally grouped and aggregated. We remove words that have less than 4 characters and fundamental stopwords, such as “with”, or “from”, in order to obtain analytical insights of high quality.

**4.2.3 Evaluation.** Our solution is tested using the raw movie data, namely `movies.csv`. We list the execution results in Appendix B.

*Stage 2: Computation of genre analytics.* We list here the execution results for the genre analytics query in a 6-column format.

Action: 7348	Comedy: 16870	Drama: 25606	IMAX: 195	Sci-Fi: 3595	Western: 1399
Adventure: 4145	Crime: 5319	Fantasy: 2731	Musical: 1054	Thriller: 8654	N/A: 5062
Animation: 2929	Documentary: 5605	Film-Noir: 353	Mystery: 2925	War: 1874	
Children: 2935		Horror: 5989	Romance: 7719		

*Stage 3: Computation of year analytics.* We list here the execution results for the year analytics query in a 3-column format, i.e., the top 10 years with the most movies. We observe that these years span the time period between 2009 and 2018, excluding 2015, indicating a heavy film production process in this decade. Figure 4 depicts the distribution of movies per year for all years in chronological order. Our findings indicate a clear long-tail distribution of movies over the years, with a burst of new movies over the last decade.

2015: 2512 movies	2013: 2173 movies	2009: 1724 movies
2016: 2488 movies	2018: 2032 movies	2010: 1691 movies
2014: 2406 movies	2012: 1978 movies	
2017: 2373 movies	2011: 1838 movies	

*Stage 4: Computation of title analytics.* We list in Appendix B the execution results for the title analytics query. Also, Figure 5 depicts a word cloud with the most frequently used words in movie titles, found in the dataset used for evaluating our solution. Note that the bigger a word is, the more frequently it appears in movie titles.

## 5 Problem 4: Probabilistic graph

We discuss here the fourth problem of the assignment. The main target of the assignment is to get acquainted with performing various MapReduce phases on a pipeline frame.

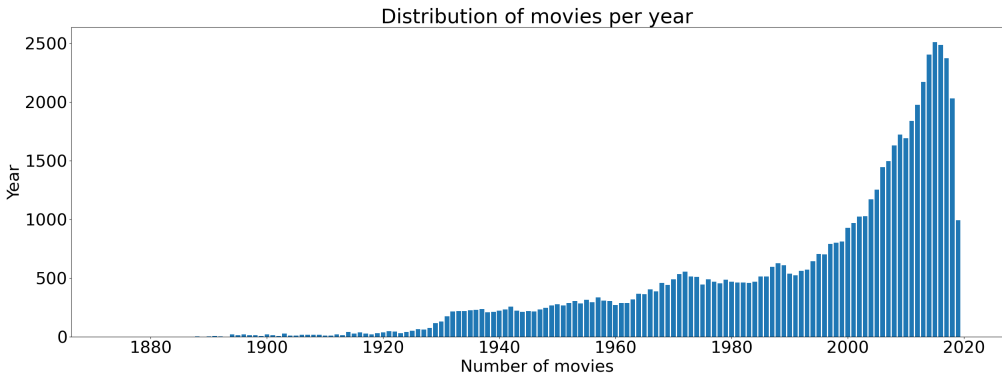


Fig. 4. The distribution of the total number of movies per year in the dataset used for evaluating our solution.



Fig. 5. A word cloud depicting the most frequent words appearing in the titles of the movies in the dataset examined. The bigger a word is, the more frequently it appears in the titles of the movies.

## 5.1 Problem Statement

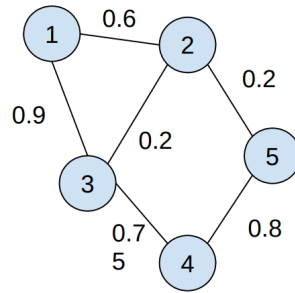
We are given an input file in text format where each line contains a connection between two vertices of a network and a probability value. For each edge  $e$ , there is a probability value  $p(e)$  which indicates the probability that the two vertices are connected by the edge. Obviously, the values of  $p(e)$  range between 0 and 1. The values in each line are separated by a space. Consider the network shown in the previous figure. The edge connecting vertices 4 and 5 has a probability of 0.8, the edge connecting vertices 2 and 3 has a probability of 0.2, etc.

The file corresponding to this graph would be:

1 2 0.6, 1 3 0.9, 2 3 0.2, 3 4 0.75, 2 5 0.2, 4 5 0.8,

The edges are generally stored in random order in the file, so we cannot assume they have a specific arrangement.





The following tasks are requested:

- (1) Write a Java program that computes the average degree for all the vertices.
- (2) The average degree is defined as the sum of the probabilities of the edges that fall on a vertex.
- (3) For example, in the previous diagram, the average degree of vertex 3 is  $0.9 + 0.2 + 0.75 = 1.85$ .
- (4) Before performing this calculation, you should ignore all edges with a probability less than a threshold  $T$ , which should be passed as a parameter to the main function.
- (5) Modify the code from task 1 so that, at the end, only the vertices with an average degree greater than the average of the degrees of all the vertices are displayed in the output.

## 5.2 Proposed approach

**5.2.1 Setting.** Our implementation is run and tested in a Linux environment with 12 cores, using the Java programming language. We have used the Java Development Kit (JDK) version 11.0.11. The source code is developed in IntelliJ IDEA Community Edition 2021.1.1 and managed using Maven as the build tool.

The project's dependencies, including Hadoop libraries, are defined in the `pom.xml` file located in the root of the repository. The project is compiled and executed directly from IntelliJ IDEA.

To run the project, open the `GraphMaster` class in IntelliJ IDEA, and execute the main method. You will need to specify the following command-line arguments:

- `<input_path>` specifies the directory containing the input text files, located at `map-reduce/ probabilisticGraph /input`.
- `<output_path>` specifies the directory where the MapReduce output will be written, which will be created in the out folder.
- `<T>` is the minimum edge-degree threshold for vertices to be included in the results.

IntelliJ IDEA will handle the compilation and execution automatically when the main method is run. Make sure to configure the input and output paths as required for your specific run.

Note that the command-line arguments must follow the specified order and format. If any of the arguments are missing or invalid, the program will terminate with an appropriate error message.

**5.2.2 Implementation.** The problem at hand involves processing a probabilistic graph, where each edge connects two vertices with a probability value. The task is to compute the average degree for each vertex, considering only edges whose probability exceeds a given threshold  $T$ . The methodology proposed here utilizes a distributed MapReduce framework, implemented using Hadoop, to efficiently process and analyze large-scale graph data. The approach is divided into three main phases, each handled by separate MapReduce jobs, with intermediate results passed between the phases.

The GraphMaster class serves as the orchestrator of the entire MapReduce process. It manages the execution of the three phases, invoking the appropriate MapReduce jobs in sequence. The main steps executed by GraphMaster are:

- (1) It retrieves the command-line arguments, including the input and output directories as well as the threshold value  $T$ , which determines which edges to consider in the graph.
- (2) It initiates the first MapReduce job to compute the degree of each vertex in the graph. The input consists of edge data, and the output is a summation of edge probabilities for each vertex.
- (3) After the first job completes, the GraphMaster starts the second MapReduce job to calculate the mean degree of the graph, using the results from the first job.
- (4) Once the mean degree is computed, GraphMaster starts the third MapReduce job, which filters out the vertices with degrees lower than the mean degree.
- (5) It cleans up intermediate directories after each phase and moves the final output to the desired location.

GraphMaster coordinates these steps by configuring and executing the MapReduce jobs, ensuring that each phase depends on the results of the previous one.

The entire process consists of three MapReduce jobs, executed sequentially:

- (1) **Phase 1:** The first job computes the degree of each vertex by summing the probabilities of the edges connected to it. It produces intermediate results for each vertex.
- (2) **Phase 2:** The second job calculates the mean degree by summing the degrees from the first phase and dividing by the total number of vertices.
- (3) **Phase 3:** The third job filters out the vertices with a degree lower than the mean degree and produces the final filtered result.

After each job completes, intermediate data is written to disk and passed as input to the subsequent job. The final output consists of the vertices whose degree exceeds or equals the mean degree.

**Phase 1: Calculating Vertex Degree** The first MapReduce job is responsible for calculating the degree of each vertex in the graph, where the degree is defined as the sum of the probabilities of the edges connected to that vertex. The degree of each vertex is computed by the GraphMapper and GraphReducer classes.

The GraphMapper class reads each edge from the input data, which consists of lines containing two vertices and the associated probability value. It splits each line into two components (the two vertices), and for each vertex, it emits the corresponding probability value as the output. The mapper also checks if the probability value exceeds the threshold  $T$ , which is provided as a configuration parameter. If the probability is greater than or equal to  $T$ , it emits the vertex along with the corresponding probability.

The GraphReducer class receives the vertex and associated probability values emitted by the mapper. For each vertex, it sums up all the probabilities of the edges connected to that vertex, and emits the pair as output.

**Phase 2: Calculating the Mean Degree** The second phase of the approach calculates the mean degree of all vertices. This is done by summing the degree values from the previous phase and dividing by the total number of vertices.

The MeanMapper class receives the degree values from the output of the first MapReduce job. It emits two key-value pairs for each input line: a count of the number of vertices and the sum of the degree values. These are emitted with fixed keys ("count" and "sum") so that they can be aggregated in the reducer.

The MeanReducer class processes the count and sum values emitted by the mapper. It calculates the total sum and count of vertices and then writes these values as output. The mean degree is calculated by dividing the sum by the count, and the result is stored for use in the next phase.

**Phase 3: Filtering Vertices by Mean Degree** In the third phase, the vertices whose degree is greater than or equal to the mean degree calculated in Phase 2 are selected. This phase uses the results from Phase 1 and Phase 2, applying the filtering criteria to retain only those vertices with a degree greater than or equal to the mean.

The `FilterMapper` class takes the degree values emitted by the first reducer and writes them as key-value pairs. Each key is a vertex, and the value is the degree of that vertex. The mapper does not perform any filtering; instead, it simply passes the values along to the reducer.

The `FilterReducer` class filters out the vertices whose degree is less than the mean degree. It retrieves the mean degree from the configuration, and for each vertex, it compares the degree value to the mean. If the degree is greater than or equal to the mean, it emits the vertex along with its degree. This final output contains only the vertices that satisfy the filtering criteria.

**5.2.3 Evaluation.** The experiments were executed using the dataset ‘collins.txt’ with a parameter setting of  $T = 0.8$ .

To provide a comprehensive overview, the results of the experiment, including the detailed values derived from the dataset, are presented in subsection B.1 in a 2-column format.

## 6 Conclusion

In this document we have presented our solutions and rationale for solving the second assignment of the M.Sc. course on *Technologies for Big Data Analysis*, offered by the *DWS M.Sc. Program*. For each one of the four problems of the assignment, we have presented their statement, as well as the solution approach we have adopted. All execution results with the provided datasets can be found in the appendices of our work.

## A Problem 2 execution results

We list here the execution results of our solution for the second problem (movie analytics) tested on the `movies.csv` dataset. Task 1 (duration per country) produces the following results:

Afghanistan: 815	Denmark: 44186	Kosovo: 188	Portugal: 15232
Albania: 646	Dominican Republic: 866	Kuwait: 273	Puerto Rico: 1322
Algeria: 2115	East Germany: 2095	Kyrgyzstan: 400	Qatar: 1036
American Samoa: 247	Ecuador: 893	Laos: 289	Republic of Macedonia: 1515
Angola: 357	Egypt: 3090	Latvia: 2297	Reunion: 38
Argentina: 28813	El Salvador: 167	Lebanon: 1289	Romania: 15739
Armenia: 264	Estonia: 4651	Liberia: 283	Russia: 50757
Aruba: 739	Ethiopia: 537	Libya: 520	Rwanda: 527
Australia: 66068	Faroe Islands: 97	Liechtenstein: 1011	Samoa: 110
Austria: 22925	Federal Republic of Yugoslavia: 2167	Lithuania: 3772	Saudi Arabia: 490
Bahamas: 560	Finland: 48990	Luxembourg: 11976	Senegal: 1462
Bahrain: 87	France: 467279	Macao: 175	Serbia: 4337
Bangladesh: 429	Gabon: 80	Madagascar: 204	Serbia and Montenegro: 1314
Barbados: 92	Georgia: 1846	Malaysia: 1338	Singapore: 4557
Belarus: 427	Germany: 226808	Mali: 220	Slovakia: 2525
Belgium: 54192	Ghana: 655	Malta: 885	Slovenia: 2956
Bermuda: 95	Greece: 20305	Martinique: 103	Somalia: 90
Bhutan: 291	Greenland: 90	Mauritania: 515	South Africa: 14688
Bolivia: 797	Grenada: 79	Mexico: 35560	South Korea: 56245
Bosnia and Herzegovina: 2207	Guatemala: 638	Micronesia: 85	Soviet Union: 51353
Botswana: 295	Guinea: 90	Moldova: 95	Spain: 116651
Brazil: 27333	Haiti: 498	Monaco: 189	Sri Lanka: 318
Bulgaria: 5430	Honduras: 80	Mongolia: 186	Suriname: 105
Burkina Faso: 903	Hong Kong: 66567	Montenegro: 439	Sweden: 68956
Burma: 95	Hungary: 19668	Morocco: 4525	Switzerland: 34516
Cambodia: 894	Iceland: 7839	Namibia: 105	Syria: 447
Cameroon: 615	India: 110215	Nepal: 797	Taiwan: 16618
Canada: 210826	Indonesia: 3931	Netherlands: 47054	Tajikistan: 351
Chad: 493	Iran: 11762	New Zealand: 14561	Tanzania: 441
Chile: 6178	Iraq: 1365	Nicaragua: 183	Thailand: 10189
China: 42089	Ireland: 26506	Nigeria: 494	The Democratic Republic Of Congo: 98
Colombia: 3810	Isle Of Man: 396	North Korea: 508	Trinidad and Tobago: 86
Congo: 88	Israel: 16662	Norway: 26476	Tunisia: 2236
Costa Rica: 431	Italy: 288407	Pakistan: 1728	Turkey: 20068
Croatia: 4340	Jamaica: 807	Palestine: 1530	UK: 492616
Cuba: 2899	Japan: 178558	Panama: 903	USA: 2276142
Cyprus: 948	Jordan: 1313	Papua New Guinea: 295	Uganda: 243
Czech Republic: 15317	Kazakhstan: 2145	Paraguay: 317	Ukraine: 4616
Czechoslovakia: 8434	Kenya: 523	Peru: 2436	
Côte d'Ivoire: 180	Korea: 90	Philippines: 7882	
		Poland: 35543	

United Arab Emirates: 3202	Uzbekistan: 324	Vietnam: 1261	Zaire: 80
Uruguay: 2038	Vanuatu: 100	West Germany: 58480	Zimbabwe: 210
	Venezuela: 2014	Yugoslavia: 8742	

Task 2 (movies per year and genre) produces the following results:

1973_Adventure: 1	1982_Comedy: 1	2002_Documentary: 1	2007_War: 1
1973_Drama: 1	1982_Drama: 1	2002_Music: 1	2008_Action: 1
1973_History: 1	1985_Documentary: 1	2004_Adventure: 1	2008_Crime: 1
1976_Comedy: 2	1985_History: 1	2004_Documentary: 2	2008_Drama: 1
1976_Drama: 1	1988_Comedy: 1	2004_Drama: 1	2009_Documentary: 1
1976_History: 1	1988_Drama: 2	2004_Music: 1	2015_Action: 1
1976_Mystery: 1	1988_Sci-Fi: 1	2004_Romance: 1	2015_Adventure: 1
1976_Romance: 1	1996_Documentary: 1	2006_Animation: 1	2015_Animation: 1
1979_Adventure: 1	1996_Music: 1	2006_Crime: 1	
1979_Crime: 1	2001_Documentary: 1	2006_Drama: 1	
1979_Mystery: 1	2001_Music: 1	2007_Documentary: 1	

## B Problem 3 execution results

love: 851	that: 201	hell: 146	never: 114	kiss: 99	west: 91
night: 519	blue: 200	legend: 146	beyond: 113	room: 99	case: 91
story: 495	woman: 197	killer: 145	devil: 112	magic: 99	green: 90
life: 476	good: 196	heart: 145	perfect: 111	it's: 98	crazy: 88
last: 451	three: 194	return: 144	living: 111	party: 98	water: 88
girl: 387	wild: 192	when: 142	long: 110	just: 98	once: 88
black: 352	summer: 185	like: 142	space: 108	town: 98	fall: 87
time: 350	great: 181	young: 139	monster: 108	children: 98	song: 87
little: 335	home: 181	ghost: 138	wedding: 108	man,: 98	wolf: 86
dead: 327	girls: 180	family: 137	school: 108	will: 97	paris: 85
death: 323	don't: 179	star: 137	best: 108	four: 97	deadly: 85
house: 311	live: 179	women: 135	year: 107	tale: 97	queen: 85
world: 306	lady: 169	part: 135	battle: 106	heaven: 97	america: 84
movie: 286	island: 168	boys: 134	final: 106	other: 95	miss: 84
christmas: 274	high: 166	game: 133	baby: 106	river: 95	true: 84
blood: 271	this: 165	murder: 129	come: 106	fear: 95	midnight: 84
dark: 258	road: 159	land: 129	angel: 105	over: 95	into: 83
city: 250	back: 158	fire: 125	eyes: 104	they: 94	hero: 83
american: 249	what: 156	adventures: 121	happy: 103	here: 94	snow: 82
white: 235	kill: 156	people: 121	before: 102	golden: 93	dance: 81
your: 225	street: 155	under: 120	earth: 101	sweet: 93	light: 81
king: 219	moon: 154	after: 117	without: 101	five: 93	jack: 81
lost: 209	first: 153	evil: 116	dragon: 101	years: 92	hard: 81
days: 209	about: 152	down: 114	dream: 101	beautiful: 92	goes: 81
secret: 202	the): 147	journey: 114	seven: 100	film: 91	hollywood: 80

only: 79	vampire: 68	diary: 62	zombie: 54	stone: 48	hands: 44
princess: 79	stranger: 68	nothing: 62	blind: 54	grace: 48	gone: 44
rock: 79	dawn: 68	point: 62	glory: 54	horror: 48	soldier: 44
another: 79	child: 68	beauty: 61	force: 54	iii:: 48	very: 44
dreams: 78	behind: 68	spring: 61	horse: 53	invisible: 48	sister: 44
prince: 78	name: 67	history: 61	power: 53	storm: 48	hope: 43
angels: 78	mother: 67	broken: 61	welcome: 53	walk: 48	voyage: 43
john: 77	dangerous: 67	were: 61	castle: 53	sword: 48	david: 43
escape: 77	planet: 67	play: 61	head: 53	line: 48	hidden: 43
terror: 77	trouble: 67	nightmare: 60	music: 53	mind: 48	alice: 43
cold: 77	born: 67	comedy: 60	massacre: 52	friend: 47	dust: 43
shadow: 77	crime: 66	comes: 59	attack: 52	mrs.: 47	grand: 43
zero: 76	special: 66	master: 59	real: 52	samurai: 47	louis: 43
brothers: 76	there: 66	place: 59	camp: 52	you're: 47	demon: 43
call: 76	winter: 66	naked: 58	going: 52	bill: 47	missing: 43
strange: 76	park: 66	super: 58	fight: 52	express: 47	mission: 43
take: 75	blues: 66	johnny: 58	kids: 52	against: 47	rising: 43
inside: 75	human: 66	made: 58	shadows: 52	most: 47	guns: 43
money: 74	lives: 66	killing: 58	affair: 52	came: 47	state: 43
iron: 74	lake: 66	show: 57	miracle: 51	harry: 46	being: 43
charlie: 74	meet: 66	stars: 57	hunter: 51	texas: 46	romance: 43
country: 74	north: 66	(les: 57	flying: 51	enemy: 46	making: 43
wind: 74	friends: 65	witch: 57	soul: 51	hunt: 46	called: 43
side: 74	beach: 65	santa: 57	small: 51	silence: 46	trail: 42
wife: 73	silent: 65	holiday: 57	alive: 51	hearts: 46	wall: 42
mystery: 73	free: 65	rose: 57	future: 50	ninja: 46	thief: 42
paradise: 73	next: 65	still: 56	heroes: 50	george: 46	code: 42
tales: 73	times: 65	nights: 56	london: 50	company: 46	yellow: 42
where: 73	captain: 65	ride: 56	dogs: 50	open: 46	tell: 42
bride: 73	have: 65	hill: 56	through: 50	stories: 46	than: 42
father: 72	jungle: 65	things: 56	something: 50	haunted: 45	tree: 42
again: 72	tomorrow: 64	flight: 56	vengeance: 50	full: 45	thing: 42
gold: 72	forever: 64	upon: 56	tiger: 50	desert: 45	pink: 42
mountain: 71	dirty: 64	fury: 56	operation: 49	lies: 44	sisters: 42
truth: 71	detective: 64	rise: 55	more: 49	tokyo: 44	holy: 42
revenge: 71	away: 64	rain: 55	wars: 49	always: 44	story: 42
book: 71	curse: 64	garden: 55	sunday: 49	door: 44	sleep: 42
beast: 69	alien: 64	club: 55	mary: 49	mine: 44	bloody: 42
face: 69	body: 63	brother: 55	flesh: 49	les): 44	killed: 42
lucky: 69	train: 63	right: 55	between: 49	left: 44	can't: 41
adventure: 69	double: 63	know: 55	gang: 49	project: 44	forbidden: 41
edge: 69	private: 63	doctor: 55	secrets: 49	glass: 44	bridge: 41
york: 69	second: 63	valley: 54	broadway: 49	lovers: 44	mars: 41
deep: 69	alone: 63	phantom: 54	china: 49	make: 44	chapter: 41
hotel: 68	darkness: 63	treasure: 54	skin: 48	dear: 44	morning: 41
devil's: 68	daughter: 62	machine: 54	love: 48	hour: 44	class: 41

michael: 41	every: 38	century: 35	ring: 32	dancing: 30	rome: 28
breaking: 41	look: 38	quiet: 35	paul: 32	tears: 30	creek: 28
warrior: 41	chance: 38	south: 35	richard: 32	danger: 30	brooklyn: 28
goodbye: 41	running: 38	within: 34	looking: 32	county: 30	wonderland: 28
justice: 41	flower: 37	greatest: 34	lord: 32	lonely: 30	waiting: 28
frankenstein: 41	burning: 37	agent: 34	carry: 32	elephant: 30	pirates: 28
kingdom: 40	unknown: 37	hand: 34	beat: 32	keep: 30	della: 28
faces: 40	savage: 37	woman,: 34	some: 31	almost: 30	cinema: 28
touch: 40	bear: 37	kings: 34	henry: 31	flowers: 29	jones: 28
heat: 40	date: 37	central: 34	untold: 31	quest: 29	those: 27
anna: 40	let's: 37	beginning: 34	stand: 31	returns: 29	hercules: 27
fast: 40	fever: 37	zombies: 34	action: 31	movie,: 29	falls: 27
kung: 40	search: 37	season: 34	berlin: 31	meets: 29	holmes: 27
late: 40	forest: 37	ghosts: 34	past: 31	guys: 29	scream: 27
games: 40	sherlock: 37	rage: 34	tour: 31	circus: 29	uncle: 27
teenage: 40	life,: 37	frank: 34	lover: 31	manhattan: 29	assassin: 27
work: 40	funny: 37	empire: 34	wonder: 31	world,: 29	luck: 27
dracula: 40	pretty: 37	happiness: 34	generation: 31	sing: 29	adam: 27
freedom: 40	ever: 37	billy: 34	jesus: 31	nowhere: 29	fine: 27
vacation: 40	circle: 36	hood: 34	them: 31	band: 29	seventh: 27
man's: 40	saint: 36	easy: 34	talk: 31	beneath: 29	border: 27
trip: 40	thunder: 36	legacy: 34	woods: 31	remember: 29	female: 27
cool: 40	revolution: 36	east: 34	monkey: 31	finding: 29	team: 27
i'll: 39	takes: 36	half: 34	brown: 31	rabbit: 29	playing: 27
must: 39	race: 36	weekend: 34	movie:: 31	musical: 29	dollar: 27
everything: 39	cinderella: 36	bullet: 34	jimmy: 31	lights: 29	control: 27
brave: 39	minutes: 36	among: 33	wanted: 31	nine: 29	barbie: 27
kind: 39	national: 36	blonde: 33	strangers: 30	getting: 29	duck: 27
giant: 39	police: 36	chasing: 33	honey: 30	picture: 29	sons: 27
shark: 39	coming: 36	spirit: 33	butterfly: 30	gangster: 29	birds: 27
third: 39	hours: 36	amazing: 33	eagle: 30	loves: 29	violence: 27
want: 39	friday: 36	hate: 33	marriage: 30	farm: 28	till: 27
sound: 39	better: 36	silver: 33	falling: 30	magnificent: 28	order: 27
bird: 39	walking: 36	married: 33	hills: 30	sleeping: 28	moment: 27
confessions: 39	gods: 36	bang: 33	letter: 30	below: 28	halloween: 27
honor: 39	lone: 36	shoot: 33	steel: 30	chan: 28	village: 27
monsters: 39	wonderful: 36	souls: 33	killers: 30	happened: 28	change: 27
business: 39	brain: 36	around: 33	scooby-doo!: 30	belle: 28	werewolf: 27
fish: 39	thousand: 36	worlds: 33	perry: 30	roll: 28	witness: 27
lion: 38	army: 35	shot: 33	wings: 30	merry: 28	give: 27
passion: 38	above: 35	number: 33	public: 30	doll: 28	ladies: 27
portrait: 38	french: 35	blade: 32	warriors: 30	ball: 28	apocalypse: 27
wrong: 38	royal: 35	trial: 32	break: 30	you,: 28	(zatôichi: 27
youth: 38	wish: 35	robin: 32	snake: 30	chinese: 28	stop: 27
across: 38	jane: 35	their: 32	mark: 30	grave: 28	watch: 26
million: 38	short: 35	madness: 32	color: 30	zone: 28	bunny: 26

rich: 26	damned: 25	desire: 23	hole: 22	vegas: 21	guide: 20
peter: 26	bell: 25	rescue: 23	paper: 22	courage: 21	safe: 20
haunting: 26	nobody: 25	murders: 23	hide: 22	close: 21	viva: 20
tarzan: 26	scarlet: 25	food: 23	talking: 22	moonlight: 21	rush: 20
jeff: 26	fatal: 25	united: 23	tower: 22	move: 21	underground: 20
sunshine: 26	cross: 25	wives: 23	criminal: 22	creature: 21	twelve: 20
please: 26	cowboy: 25	mysterious: 23	assassination: 22	lane: 21	prisoner: 20
highway: 26	resurrection: 25	chaos: 23	riding: 22	tough: 21	wants: 20
radio: 26	animals: 25	classic: 23	poor: 22	cobra: 21	volume: 20
(die: 26	shaolin: 25	crossing: 23	wake: 22	twist: 21	pain: 20
these: 26	memory: 25	eternal: 23	kong: 22	care: 21	purple: 20
dinner: 26	twilight: 25	general: 23	peace: 22	theroux:: 21	die): 20
student: 26	wolves: 25	ground: 23	squad: 22	breakfast: 21	september: 20
dick: 26	straight: 25	sunset: 23	hollow: 22	mister: 21	trust: 20
august: 26	leave: 24	who's: 23	much: 22	africa: 21	turn: 20
tall: 26	eight: 24	birthday: 23	impossible: 22	fool: 21	runner: 20
virgin: 26	godzilla: 24	knight: 23	alley: 22	theatre: 21	gets: 20
mask: 26	genius: 24	series: 23	odyssey: 22	chicken: 21	front: 20
voice: 26	hello: 24	square: 23	daddy: 22	strikes: 21	fairy: 20
price: 26	angry: 24	hall: 23	style: 22	mike: 21	clown: 20
step: 26	sugar: 24	plan: 23	spider: 22	think: 21	darling: 20
fist: 26	yours: 24	girl,: 23	fort: 22	prey: 21	dave: 20
bachelor: 26	reunion: 24	wife,: 23	requiem: 22	dolls: 21	diamond: 20
violent: 26	nature: 24	single: 23	what's: 22	smart: 21	boat: 20
eden: 26	target: 24	california: 23	opera: 22	nice: 21	raiders: 20
demons: 26	animal: 24	teen: 23	miles: 22	birth: 21	find: 20
burn: 25	catch: 24	found: 23	maria: 22	then: 21	mercy: 20
satan: 25	trap: 24	palace: 23	destiny: 22	society: 21	donald's: 20
mason:: 25	forgotten: 24	saving: 23	hunters: 22	lightning: 21	nick: 20
batman: 25	drive: 24	lego: 23	one,: 22	camera: 21	league: 20
conspiracy: 25	moving: 24	ones: 23	electric: 22	chronicles: 21	hunting: 20
rainbow: 25	she's: 24	shanghai: 23	pass: 22	wizard: 21	god's: 20
ballad: 25	scared: 24	house,: 23	bitter: 22	speed: 21	need: 20
wicked: 25	veggietales:: 24	shock: 23	pride: 21	taking: 20	yesterday: 20
harvest: 25	incident: 24	simple: 23	martin: 21	soldiers: 20	western: 20
loved: 25	invasion: 24	rules: 23	daughters: 21	could: 20	woman's: 20
streets: 25	eddie: 24	factory: 23	annie: 21	sailor: 20	we're: 20
autumn: 25	modern: 24	thieves: 22	universe: 21	documentary: 20	gift: 20
match: 25	bank: 24	innocent: 22	forget: 21	robert: 20	cats: 20
nation: 25	fair: 24	conan:: 22	boss: 21	fantastic: 20	thin: 20
fighting: 25	james: 24	words: 22	report: 21	honeymoon: 20	center: 20
together: 25	bright: 24	hare: 22	2000: 21	duel: 20	fighter: 20
april: 25	dying: 23	ultimate: 22	frozen: 21	letters: 20	tango: 20
madame: 25	loving: 23	bobby: 22	normal: 21	station: 20	venus: 20
taxi: 25	clouds: 23	cage: 22	sand: 21	teacher: 20	delta: 19
psycho: 25	prison: 23	night,: 22	gate: 21	presents:: 20	promise: 19



incredible: 19	connection: 19	sheep: 18	mouth: 17	devils: 17	tapes: 16
chicago: 19	crimes: 19	albums:: 18	pursuit: 17	well: 17	does: 16
shoes: 19	went: 19	beloved: 18	noon: 17	tony: 17	jerry: 16
cousin: 19	frontier: 19	emperor: 18	johnson: 17	inc.: 17	chronicle: 16
outer: 19	sinners: 19	andy: 18	desperate: 17	wars:: 17	wrath: 16
cave: 19	dans: 19	outlaw: 18	velvet: 17	fortune: 17	rider: 16
shop: 19	scandal: 19	mighty: 18	afternoon: 17	many: 17	rolling: 16
pleasure: 19	smoke: 19	heights: 18	enchanted: 17	hold: 17	sale: 16
choice: 19	strike: 19	satan's: 18	begins: 17	high:: 16	time:, 16
lesson: 19	oscar: 18	pacific: 18	serial: 17	bone: 16	anything: 16
everybody: 19	academy: 18	saturday: 18	taste: 17	parts: 16	mondo: 16
funeral: 19	falcon: 18	arms: 18	indian: 17	laughing: 16	panther: 16
along: 19	kevin: 18	knock: 18	bigfoot: 17	sense: 16	twenty: 16
witchcraft: 19	october: 18	citizen: 18	(das: 17	carnival: 16	endless: 16
mirror: 19	superman: 18	chase: 18	sight: 17	vice: 16	mouse: 16
monsieur: 19	awakening: 18	kid,: 18	path: 17	office: 16	crooked: 16
count: 19	field: 18	guest: 18	somewhere: 17	avenue: 16	video: 16
candy: 19	sarah: 18	captive: 18	innocence: 17	twin: 16	poison: 16
bruce: 19	ship: 18	trapped: 18	monk: 17	saved: 16	hurricane: 16
siege: 19	fallen: 18	that's: 18	hound: 17	table: 16	personal: 16
divine: 19	lovely: 18	blondie: 18	round: 17	main: 16	romeo: 16
wave: 19	swan: 18	early: 18	gray: 17	india: 16	buddy: 16
italian: 19	russell: 18	raid: 18	knows: 17	songs: 16	tender: 16
stolen: 19	parents: 18	hell's: 18	different: 17	nest: 16	store: 16
science: 19	cannibal: 18	hamlet: 18	skies: 17	parade: 16	perry's: 16
bullets: 19	roast: 18	stage: 18	apache: 17	lust: 16	rebel: 16
exit: 19	ticket: 18	robbery: 18	states: 17	inspector: 16	knights: 16
robot: 19	enemies: 18	believe: 18	mother's: 17	tyler: 16	(der: 16
chocolate: 19	hong: 18	panic: 18	cruel: 17	moscow: 16	caught: 16
farewell: 19	cherry: 18	mutant: 18	galaxy: 17	seduction: 16	voices: 16
word: 19	grass: 18	rough: 18	wise: 17	pigs: 16	service: 16
encounters: 19	der): 18	monte: 17	alex: 17	u.s.a.: 16	13th: 16
save: 19	claus: 18	suicide: 17	feet: 17	kitchen: 16	sign: 16
thank: 19	lola: 18	para: 17	fools: 17	boys:: 16	crush: 16
follow: 19	pour: 18	even: 17	coast: 17	later: 16	weapon: 16
bears: 19	joan: 18	college: 17	week: 17	hawk: 16	boots: 16
evening: 19	conquest: 18	mermaid: 17	murderer: 17	extraordinary: 16	bound: 16
faith: 19	crimson: 18	dead,: 17	list: 17	16	affairs: 16
buffalo: 19	confession: 18	knife: 17	movies: 17	labyrinth: 16	murder,: 16
heist: 19	encounter: 18	simon: 17	canyon: 17	steal: 16	smile: 16
experiment: 19	steps: 18	russian: 17	guilty: 17	william: 16	named: 16
fate: 19	passage: 18	paranormal: 17	ends: 17	atomic: 16	laugh: 16
matter: 19	inferno: 18	babylon: 17	phoenix: 17	dynamite: 16	mickey: 16
persian: 19	episode: 18	grey: 17	extreme: 17	lisa: 16	django: 16
husband: 19	chain: 18	ways: 17	rocky: 17	sea,: 16	giants: 16
prime: 19	ashes: 18	lampoon's: 17	someone: 17	challenge: 16	reality: 16

bells: 16	strong: 15	musketeers: 14	wilderness: 14	jurassic: 13	mickey's: 13
flash: 16	memories: 15	(gojira: 14	there's: 14	losing: 13	salt: 13
d'un: 16	films: 15	wwe:: 14	lily: 14	gypsy: 13	pure: 13
horses: 16	pokémon: 15	fourth: 14	lethal: 14	brian: 13	heaven's: 13
donald: 15	ugly: 15	lupin: 14	crisis: 14	boy,: 13	dinosaur: 13
cell: 15	clear: 15	news: 14	buck: 14	eternity: 13	artist: 13
bulldog: 15	kelly: 15	williams:: 14	isle: 14	widow: 13	father's: 13
shooting: 15	metal: 15	education: 14	ricky: 14	nous: 13	cops: 13
christ: 15	loose: 15	tomb: 14	heart,: 14	page: 13	bomb: 13
confidential: 15	destination: 15	ordinary: 14	monogatari): 14	hunger: 13	trees: 13
feast: 15	male: 15	war,: 14	places: 14	valentine: 13	said: 13
deal: 15	street,: 15	jeekyll: 14	kino-pravda: 14	fearless: 13	vanishing: 13
stay: 15	arrow: 15	suit: 14	flame: 14	daisy: 13	bliss: 13
famous: 15	anne: 15	beware: 14	driver: 14	afraid: 13	theory: 13
liberty: 15	middle: 15	outside: 14	temple: 14	concert: 13	swamp: 13
roman: 15	keeper: 15	same: 14	been: 14	apple: 13	whale: 13
twisted: 15	beverly: 15	brief: 14	command: 14	ties: 13	hyde: 13
foreign: 15	bikini: 15	alla: 14	degrees: 14	game,: 13	miami: 13
robinson: 15	patrol: 15	elvis: 14	november: 14	sexy: 13	paris,: 13
major: 15	chris: 15	mistress: 14	floor: 14	neighbor: 13	undercover: 13
according: 15	fingers: 15	calls: 14	bull: 14	possession: 13	julia: 13
obsession: 15	walls: 15	tail: 14	warning: 14	standing: 13	sorority: 13
scenes: 15	mountains: 15	fright: 14	bandit: 14	bangkok: 13	prayer: 13
steve: 15	showdown: 15	trick: 14	judge: 14	molly: 13	bastards: 13
mile: 15	assault: 15	arthur: 14	heavy: 14	area: 13	president: 13
model: 15	nanny: 15	shift: 14	club,: 14	the: 13	aliens: 13
plus: 15	spies: 15	cleopatra: 14	alexander: 14	pool: 13	invincible: 13
hair: 15	calling: 15	roses: 14	marvel: 14	boys,: 13	navy: 13
seconds: 15	image: 15	symphony: 14	african: 14	tunnel: 13	stones: 13
natural: 15	kills: 15	others: 14	russia: 14	heads: 13	coffee: 13
america's: 15	dragons: 15	assassins: 14	english: 14	heaven,: 13	carol: 13
wait: 15	plastic: 15	king,: 14	crash: 14	cowboys: 13	july: 13
diaries: 15	femme: 15	triple: 14	seas: 14	tale,: 13	liar: 13
scorpion: 15	waters: 14	solo: 14	dancer: 14	live:: 13	marie: 13
thirst: 15	mafia: 14	bread: 14	countdown: 14	melody: 13	crystal: 13
amityville: 15	traffic: 14	fifth: 14	crown: 14	puppet: 13	victory: 13
nude: 15	havana: 14	worst: 14	grande: 14	sharpe's: 13	rhapsody: 13
stella: 15	priest: 14	asylum: 14	corner: 14	enough: 13	shall: 13
l.a.: 15	furious: 14	tonight: 14	julie: 14	apes: 13	destruction: 13
men,: 15	orange: 14	whole: 14	march: 14	francisco: 13	home,: 13
weeks: 15	fred: 14	miracles: 14	trailer: 14	corpse: 13	near: 13
arizona: 15	polar: 14	rest: 14	swing: 14	boyfriend: 13	wood: 13
tide: 15	affair,: 14	mean: 14	disaster: 13	foot: 13	pieces: 13
day,: 15	redemption: 14	note: 14	slow: 13	world's: 13	drop: 13
while: 15	lose: 14	temptation: 14	motion: 13	knew: 13	winnie: 13
enter: 15	mama: 14	window: 14	jackson: 13	runs: 13	saints: 13

oblivion: 13	serpent: 12	really: 12	ransom: 11	magical: 11	block: 11
betty: 13	macbeth: 12	washington: 12	wheels: 11	avenger: 11	limits: 11
rocks: 13	blow: 12	sexual: 12	rebirth: 11	dollars: 11	tiny: 11
munro: 13	lines: 12	brother's: 12	port: 11	disappearance: 11	magician: 11
creatures: 13	today: 12	homme: 12	pooh: 11	campus: 11	slave: 11
king's: 13	vampires: 12	divorce: 12	childhood: 11	sartana: 11	thing,: 11
distance: 13	civil: 12	johan: 12	surviving: 11	blank: 11	hello,: 11
boxer: 13	ruby: 12	fell: 12	shame: 11	eagles: 11	astérix: 11
effect: 13	should: 12	legion: 12	mummy: 11	invitation: 11	emmanuelle: 11
caesar: 13	letter,: 12	sorrow: 12	pictures: 11	farmer: 11	shell: 11
fashion: 13	thursday: 12	contract: 12	suite: 11	thirteen: 11	fortress: 11
empty: 13	pearl: 12	pyaar: 12	revelation: 11	noise: 11	khan: 11
painted: 12	fields: 12	cargo: 12	vida: 11	possessed: 11	delivery: 11
christmas,: 12	reason: 12	property: 12	son,: 11	junior: 11	ants: 11
carter: 12	contact: 12	dare: 12	one:: 11	pray: 11	vision: 11
ahead: 12	night's: 12	stealing: 12	goose: 11	immortal: 11	seed: 11
kisses: 12	watching: 12	experience: 12	mass: 11	stray: 11	wore: 11
survival: 12	smith: 12	nancy: 12	genesis: 11	lincoln: 11	bottle: 11
becoming: 12	mississippi: 12	milk: 12	bowery: 11	ivan: 11	fishing: 11
stallion: 12	leaving: 12	asterix: 12	barefoot: 11	kitty: 11	kick: 11
goddess: 12	wishes: 12	vita: 12	museum: 11	exile: 11	scene: 11
intruder: 12	man:: 12	quick: 12	infinity: 11	carmen: 11	brides: 11
underworld: 12	amor: 12	chainsaw: 12	villa: 11	rosa: 11	dean: 11
dies: 12	gentlemen: 12	closed: 12	graveyard: 11	masters: 11	help: 11
bronx: 12	room,: 12	cabin: 12	vincent: 11	maid: 11	titanic: 11
pirate: 12	corn: 12	amour: 12	montana: 11	each: 11	died: 11
lessons: 12	festival: 12	marry: 12	chamber: 11	trash: 11	fruit: 11
rides: 12	favorite: 12	toys: 12	rites: 11	dreaming: 11	dead:: 11
heavenly: 12	rouge: 12	larry: 12	apart: 11	amore: 11	emma: 11
doors: 12	loser: 12	volcano: 12	nurse: 11	fiction: 11	olsen: 11
ecstasy: 12	intimate: 12	horrors: 12	billion: 11	adult: 11	thomas: 11
wide: 12	buried: 12	poker: 12	delle: 11	rebellion: 11	viking: 11
crocodile: 12	zatoichi: 12	atlantis: 12	maniac: 11	anatomy: 11	fugitive: 11
turning: 12	uuno: 12	apartment: 12	porn: 11	holocaust: 11	twins: 11
tigers: 12	mail: 12	tattoo: 12	monty: 11	view: 11	notte: 11
homecoming: 12	won't: 12	party,: 12	jackie: 11	angel,: 11	remains: 11
butcher: 12	hundred: 12	potter: 12	death,: 11	player: 11	buddha: 11
range: 12	love's: 12	pluto: 12	emperor's: 11	mexican: 11	samson: 11
pete: 12	curious: 12	punk: 12	hart:: 11	wine: 11	swim: 11
sick: 12	rogue: 12	amazons: 12	ernest: 11	superstar: 11	feeling: 11
falk:: 12	bugs: 12	curtain: 12	goodbye,: 11	sleeps: 11	notorious: 11
terminal: 12	twice: 12	aurora: 12	gentleman: 11	wrestling: 11	riders: 11
impact: 12	sometimes: 12	june: 12	myself: 11	boom: 11	gospel: 11
plague: 12	done: 12	sex,: 12	building: 11	america:: 11	vendetta: 11
'the: 12	jazz: 12	motel: 12	beats: 11	pale: 11	universal: 11
acts: 12	danny: 12	hart: 12	holidays: 11		restless: 11

café: 11	surf: 10	attic: 10	doug: 10	sorry: 10	would: 10
juliet: 11	shut: 10	evil:: 10	turtles: 10	stupid: 10	graves: 10
accidental: 11	terra: 10	deadline: 10	alias: 10	juan: 10	polish: 10
england: 11	grow: 10	runaway: 10	leben: 10	drummond: 10	women's: 10
commando: 11	beau: 10	boston: 10	somebody: 10	searching: 10	queens: 10
santo: 11	country,: 10	batman:: 10	gladiators: 10	pony: 10	siberia: 10
jesse: 11	meat: 10	mobile: 10	cost: 10	unholy: 10	roof: 10
roger: 11	spin: 10	shelter: 10	turhapuro: 10	jean: 10	weather: 10
slaves: 11	rhythm: 10	cannibals: 10	cocaine: 10	election: 10	ninjas: 10
gates: 11	contre: 10	marine: 10	babes: 10	flag: 10	triumph: 10
waltz: 11	leather: 10	professor: 10	track: 10	germany: 10	carr:: 10
manhunt: 11	dove: 10	heroes:: 10	gamera: 10	bait: 10	husbands: 10
beer: 11	plain: 10	costello: 10	mortal: 10	brigade: 10	punch: 10
ocean: 11	fantasy: 10	victoria: 10	screaming: 10	clean: 10	says: 10
balls: 11	aces: 10	gorilla: 10	tribe: 10	travels: 10	court: 10
marathon: 11	jimi: 10	evidence: 10	fiend: 10	pokémon:: 10	back,: 10
factor: 11	wing: 10	dating: 10	pocket: 10	drums: 10	anthony: 10
maker: 11	eye,: 10	prairie: 10	flies: 10	chair: 10	chez: 10
girls,: 11	undead: 10	bare: 10	father,: 10	truth,: 10	sun,: 10
dog's: 11	phone: 10	wonders: 10	monday: 10	sands: 10	large: 10
carlin:: 11	success: 10	reich: 10	madagascar: 10	prom: 10	walks: 10
myth: 11	couple: 10	minds: 10	soft: 10	pack: 10	bong: 10
where's: 11	arabian: 10	beijing: 10	moves: 10	naughty: 10	lonesome: 10
eating: 11	morgan: 10	rocket: 10	dolly: 10	rose,: 10	harold: 10
prophecy: 11	amazon: 10	abbott: 10	nazi: 10	longest: 10	pitch: 10
chosen: 11	die,: 10	lawless: 10	dumb: 10	hannah: 10	exorcism: 10
tramp: 11	blackout: 10	ginger: 10	everybody's: 10	dark:: 10	judgment: 10
which: 11	doesn't: 10	testament: 10	hurt: 10	pope: 10	cemetery: 10
bounty: 11	pinocchio: 10	arctic: 10	i've: 10	accident: 10	voodoo: 10
nobody's: 11	level: 10	hostage: 10	venice: 10	jonathan: 10	union: 10
strawberry: 10	interview: 10	piranha: 10	way,: 10	winning: 10	pants: 10
freaks: 10	soup: 10	ears: 10	suspect: 10	casa: 10	suspicion: 10
president's: 10	dallas: 10	god,: 10	senior: 10	barbie:: 10	minute: 10
hardy: 10	he's: 10	crowd: 10	bleeding: 10	spirits: 10	thoughts: 10
cook: 10	hear: 10	seventeen: 10	jump: 10	waves: 10	piano: 10
boogie: 10	cane: 10	tragedy: 10	hansel: 10	hunter,: 10	daniel: 10
avengers: 10	emanuelle: 10	execution: 10	anniversary: 10	eleven: 10	

### B.1 Problem 4 execution results

In the following 4-columned pages, the experimental results of section's 5 problem are presented, as the raw output of the 3rd MapReduce phase execution.

102: 20.772281	107: 14.8008850000000003	111: 42.771938	114: 17.573652999999997
104: 29.609879	109: 27.163701999999994	112: 9.6700650000000001	115: 33.936378
105: 27.598046999999998	110: 16.304225	113: 28.5879290000000003	116: 25.7654580000000002

117: 30.000875	211: 26.206232	277: 13.899377000000001	34: 8.902532
1180: 9.566474	213: 8.21817	278: 45.286856999999998	341: 13.7120190000000002
1181: 8.554887	214: 19.8097530000000004	279: 36.933221999999994	342: 12.7771680000000001
1182: 10.499468	215: 22.834134999999996	280: 36.452380999999999	343: 21.602388999999999
1183: 8.2721589999999998	216: 22.450407999999996	283: 21.347036	344: 13.763916
1186: 8.551493	217: 12.86093	285: 39.845339	345: 9.50445
120: 8.667856	219: 13.857943	287: 15.863291	346: 13.803495
121: 24.678666000000003	220: 26.279680999999997	288: 24.174710999999995	347: 13.770425000000001
122: 62.881582	221: 15.366022000000003	289: 20.507460000000005	348: 13.6963880000000002
124: 23.14891	222: 10.174368	29: 10.500267	349: 12.777321
128: 13.860000000000001	223: 31.783763000000004	290: 15.802384000000002	35: 8.831461000000001
129: 13.828262	226: 11.464125999999998	292: 34.518677	350: 13.700356000000001
130: 12.778032	227: 9.910114	293: 42.414617	351: 12.82703
131: 14.548629	228: 26.034533	294: 39.387104	352: 18.913772999999996
132: 13.780112000000003	229: 17.984762	295: 33.319982999999986	353: 13.840302000000001
133: 12.777426	23: 8.068862	296: 26.437148	354: 12.026345
134: 17.525341	231: 11.600534	30: 13.124654000000001	355: 13.838998000000002
135: 12.858882000000001	232: 21.044421999999994	300: 16.612497	356: 13.709090000000002
136: 13.721149	24: 15.912360000000001	301: 35.508855000000004	36: 7.751157000000001
137: 13.660803000000001	244: 20.739924000000002	302: 40.46924	361: 13.740885
138: 13.738254000000001	245: 26.425946999999999	303: 30.324057999999997	37: 9.843972
140: 11.448863	246: 12.422665	304: 34.186174	378: 7.955433
141: 12.826131	247: 12.198811000000001	306: 12.180608	379: 11.753690999999998
143: 12.803862	249: 19.911965999999996	308: 49.780783000000014	38: 10.698117
154: 13.723290000000002	25: 9.950251999999999	309: 21.140049	380: 9.03648
170: 35.833753999999999	250: 20.231632999999995	310: 16.898876	381: 31.255191
174: 39.935224000000005	251: 35.698739999999999	311: 9.782459	385: 16.697326
175: 42.116396999999999	252: 13.855876000000002	318: 11.424189000000002	386: 13.088604999999998
176: 33.780936999999994	254: 8.187458	32: 7.906732000000001	389: 43.11891100000001
180: 9.111828999999998	255: 27.296618	320: 15.995923000000003	39: 10.792126000000001
187: 10.061923	257: 20.167998999999995	322: 12.928547	391: 7.855742000000001
191: 21.271617	258: 13.343784000000001	323: 12.95803	40: 7.842008000000001
192: 22.415023999999995	260: 33.483572999999986	324: 14.145467000000002	405: 18.636048
193: 24.815881999999995	261: 32.508742999999996	325: 16.820611	406: 30.941025999999999
194: 15.367017	264: 12.784288	326: 9.210550999999999	408: 16.025723
197: 19.721797	265: 17.312394	327: 8.77006	409: 21.124420999999998
198: 20.10406	266: 16.947997	328: 11.502478	410: 19.690642
200: 26.785185	267: 26.829961999999999	329: 9.477534	412: 12.448818000000001
202: 20.554610999999998	268: 13.143828000000001	330: 7.819681999999999	413: 21.007549999999995
203: 26.437709999999992	269: 9.54012	331: 12.199909	414: 19.201738999999996
204: 21.193861000000002	27: 18.650638999999998	333: 17.355712	415: 20.822085
205: 18.054205999999997	270: 12.200784000000002	334: 16.773555	416: 20.244446999999994
206: 12.843962	272: 27.8627	335: 12.816709	417: 20.336876999999998
208: 18.285626	273: 20.449326000000003	336: 9.472820000000002	418: 20.905703999999997
209: 9.478037	274: 21.781969000000004	337: 16.717836	419: 16.571078
210: 18.463774999999995	275: 18.921694000000002	338: 14.435681000000002	420: 14.833722000000003

421: 22.127405999999997	529: 16.491923000000003	660: 14.484713000000001	779: 8.382993
422: 20.207241999999997	530: 20.546952999999999	661: 14.486255000000002	78: 35.074965
423: 19.837138999999997	531: 13.581469	663: 14.393031	79: 18.829224
424: 22.204759999999993	532: 18.63035	664: 15.657084000000003	792: 9.650819
425: 20.3222	533: 16.137291	665: 12.470333000000002	793: 21.405161999999994
427: 18.573423000000002	534: 15.585291000000002	666: 15.570483000000001	794: 7.747058000000001
428: 18.402874999999998	535: 13.489995	667: 8.476292	796: 7.746258000000001
429: 20.996679999999998	536: 15.663219000000002	668: 8.400375	797: 20.176615999999996
431: 17.247738000000005	537: 15.558095000000002	669: 12.719472000000001	80: 19.708188999999997
432: 7.730989	538: 20.604023999999995	670: 12.759070000000001	81: 27.764969999999987
433: 20.990592000000003	539: 15.567038000000002	671: 18.786868000000002	82: 8.954957
437: 13.467927000000001	540: 10.592666000000001	672: 12.737939000000003	822: 7.824752000000001
438: 22.538098999999995	541: 13.301682	673: 11.354866	827: 26.284926
439: 20.257774999999995	542: 13.708515000000002	674: 12.464726000000002	828: 31.257229999999999
440: 18.640948999999996	544: 8.796812000000001	675: 14.512069	83: 16.828261
441: 22.565751999999993	545: 7.992185	676: 12.633483	830: 12.777338
442: 22.468430999999992	546: 8.689068	677: 15.30606	84: 17.39505
443: 22.236632999999994	547: 17.435207	678: 13.630705	845: 19.799999999999997
444: 12.668758	548: 14.626078000000001	679: 22.209649	846: 23.372872999999995
445: 21.391153999999997	549: 8.902671000000002	680: 26.495342999999999	847: 21.634598999999994
446: 20.229419	550: 15.475738000000002	681: 11.586703	848: 21.599697999999993
447: 23.098830999999993	563: 12.416060000000002	682: 11.780345000000002	849: 20.692736999999997
448: 23.389274999999994	604: 14.693641000000001	683: 10.837639	85: 17.810776
449: 21.425622999999995	605: 14.628854	684: 11.369396	850: 21.486940999999995
450: 20.304112999999994	606: 13.673755000000003	685: 20.680481999999994	851: 20.637819999999998
451: 21.309653999999995	607: 13.823655000000002	686: 28.644553999999992	852: 20.758516999999994
452: 21.574356999999992	608: 13.666609	687: 13.618457000000001	853: 19.714644999999994
453: 21.216737999999992	609: 14.749676000000001	70: 12.576885000000003	854: 20.674989999999998
454: 21.305387999999994	610: 14.726731000000003	702: 15.804968	858: 10.127550999999999
455: 22.443621999999994	611: 12.480203	703: 18.989590999999997	859: 13.744827
456: 18.427684	612: 21.530917999999996	711: 8.769808000000001	86: 9.620828000000001
457: 23.612626999999996	639: 16.345097000000003	72: 13.937429	860: 13.392533000000002
458: 20.137554999999995	64: 11.651773	729: 11.920043	861: 12.716919
468: 16.434241	640: 26.795459999999999	73: 33.317902999999994	862: 14.757012000000001
473: 15.624302000000002	641: 27.955416999999999	739: 14.833940000000002	863: 14.737124000000001
490: 15.118589000000002	642: 22.714243999999994	74: 9.418914	870: 12.684738000000001
492: 29.037201999999994	643: 28.925386	741: 17.572777	871: 12.674369
494: 14.964478999999999	644: 27.951106	743: 16.54965	872: 11.814398
497: 7.779587000000001	645: 9.891883	744: 13.711973000000002	873: 12.763509
505: 22.493396999999998	646: 22.266540999999997	745: 14.759685000000001	874: 10.89
507: 16.502023	647: 17.500691000000003	746: 14.843160000000001	875: 12.642457000000002
508: 13.205226000000001	653: 31.201086999999998	75: 8.17184	876: 12.794634000000002
509: 10.297368000000002	654: 11.212689	76: 35.164774999999998	877: 9.48321
510: 17.989523	655: 15.470210000000002	765: 8.972965	878: 12.762451
511: 14.096020000000003	656: 13.714576000000001	773: 9.301499	879: 13.635698000000001
512: 7.9094589999999998	659: 14.530663000000002	778: 9.388625000000001	880: 12.768133

881: 12.801868000000002	915: 9.631426000000001	929: 20.752049999999997	941: 25.234674999999996
882: 10.89	917: 9.638277000000002	930: 13.043696999999998	948: 16.893155
893: 7.887907	918: 10.553774	931: 18.611175	950: 20.797809999999995
894: 9.918199	919: 8.431745	932: 20.685094999999997	951: 19.461768999999997
895: 16.402685000000005	922: 10.480677000000002	933: 20.786459999999995	953: 18.305716999999998
896: 15.014125	926: 20.516713999999997	935: 20.676130999999998	
901: 10.018733	927: 18.237036999999997	936: 21.667931999999997	
902: 18.710283	928: 18.419195	937: 19.520974999999996	