# Introduction to Big Data: Assignment 2

Author: Anastasiia Shvets
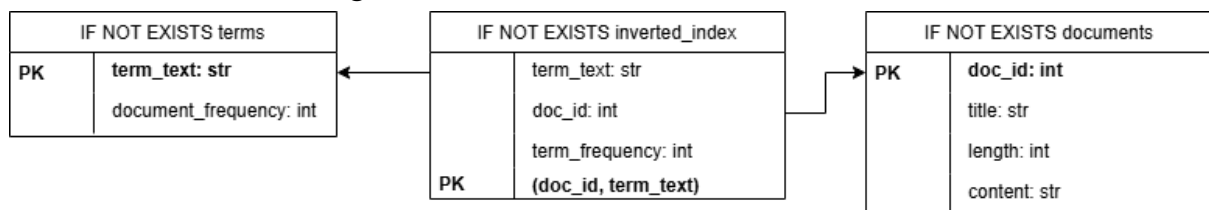Group: B22-DS-02

## Methodology

I followed the structure proposed in the github repository. However, I have made some design choices that are not present in the original solution. Here are the most import notes about my solution:

### Data preparation

Data preparation part is taken fully from the repository, nothing was changed significantly. However, I did use another parquet file from Kaggle dataset: k.parquet, since I wanted to add some twist.

### Cassandra schema

I created the following schema for the solution:



It is simple enough, but fits the queries I will use in the ranker (getting all documents, getting term frequency, getting inverted index for term and document).

### Indexer

The indexer solution is presented in the indes.sh file. First it checks if user input is provided. If it is the case, it transforms the input to one csv file in hdfs to be passed to mapreduce. Otherwise it defaults to using csv file created in /index/data during data preparation.
The mapreduce consists of two parts:
1. mapper.py: adds every document to the Cassandra database, and splits the document into (term, doc_id, 1) triples, where 1 is then used to count the number of occurrences. The solution is not too

elegant, since I do not count the number of occurrences here, however, it works anyway.
2. reducer.py: gets the triples from mapper.py and counts statistics for a term in the document. Since the indexer can be run multiple times, the reducer also gets previously collected statistics from the Cassandra database.

## Ranker

The ranker is really simple. It gets a user query, separates it into tokens and runs the BM25 as described in the homework. When it shows 10 document titles with the best scores, as well as a corresponding score.

# Demonstration

## How to run:

Before running the repository, make sure the "data" folder is created in the "app" folder, otherwise data preparation will not work!
Also make sure all files you are trying to rank exist.
Ensure you have k.parquet installed, or, if you want to use different file in data preparation, change file name both in prepare_data.py and prepare_data.sh. If you want to skip data preparation, comment out the line "bash prepare_data.sh" in app.sh

Go to app.sh file and put the file/directory you want to index after the "bash index.sh"

In the same file put queries you want to test after the "bash query.sh"

Make sure you have docker and docker compose installed.
Go to the project folder, and run

```
docker compose up
```

In the terminal.

My results:

I have used a different laptop, not the same one as for Assignment 1, so the interface may look different. This laptop belongs to a person who does not have the Big Data course, since he is on Robotics, so he could not have benefitted from my solution in any way, shape or form.
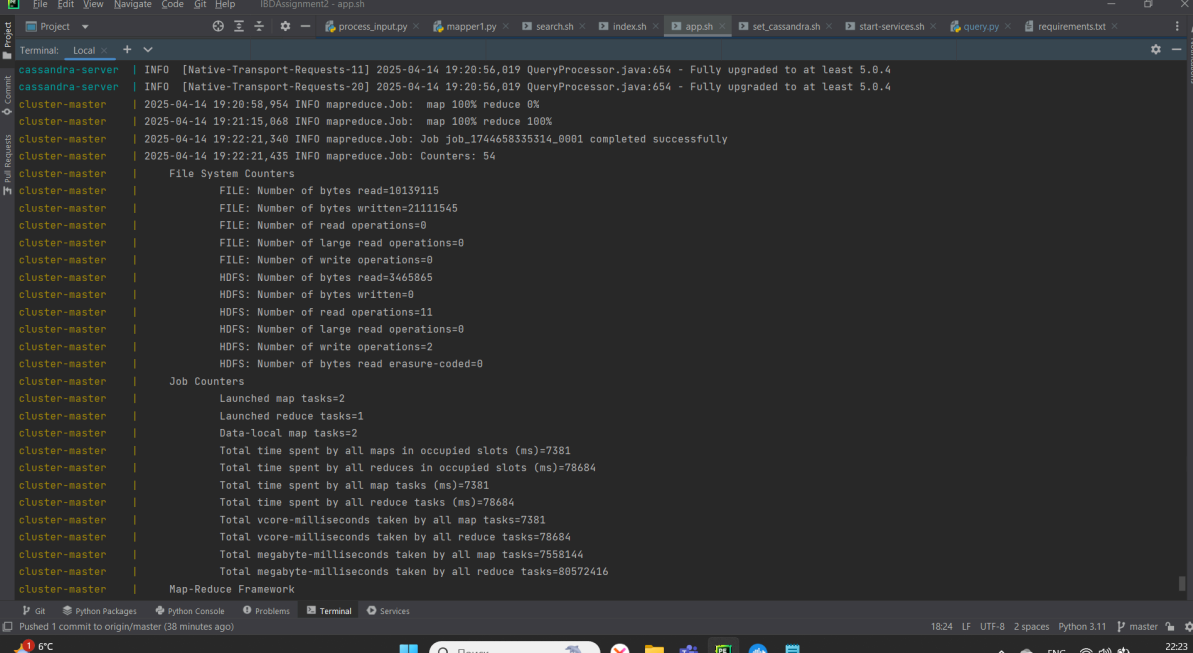
So, here is my app.sh script:



As you can see, it will index both local data directory, and local 123_cute_cats.txt file.
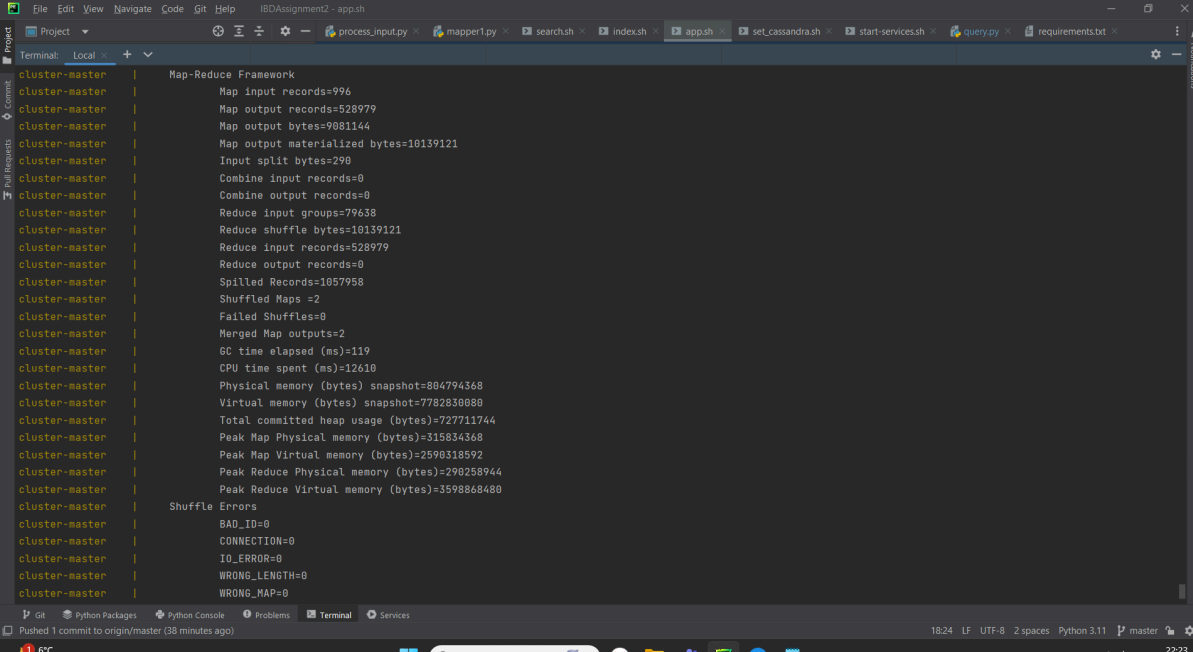So, here are the outputs for index.sh:

(As you can see, the files from the data folder have already been added to hadoop during data preparation.)

```
cluster-master |                    WRONG_MAP=0
cluster-master |                    WRONG_REDUCE=0
cluster-master |        File Input Format Counters
cluster-master |                    Bytes Read=3465575
cluster-master |        File Output Format Counters
cluster-master |                    Bytes Written=0
cluster-master | 2025-04-14 19:22:21,436 INFO streaming.StreamJob: Output directory: /tmp/index/1744658410/stage1
cluster-master | Indexing completed. Cleaning up...
cluster-master | Deleted /tmp/index/1744658410
cluster-master | Index is ready in Cassandra database.
cluster-master | User input detected. Processing...
cluster-master | Setting default log level to "WARN".
cluster-master | To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
cluster-master | 25/04/14 19:22:28 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
cluster-master | /app/.venv/lib/python3.8/site-packages/pyspark/sql/context.py:113: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
cluster-master |   warnings.warn(
User input prepared for mapreduce!
cluster-master | Starting indexing process...
cluster-master | packageJobJar: [/tmp/hadoop-unjar239807064660225895/] [] /tmp/streamjob941645282015129908.jar tmpDir=null
cluster-master | 2025-04-14 19:22:40,180 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at cluster-master/172.18.0.4:8032
cluster-master | 2025-04-14 19:22:40,327 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at cluster-master/172.18.0.4:8032
cluster-master | 2025-04-14 19:22:40,478 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/root/.staging/job_1744658335314_0002
cluster-master | 2025-04-14 19:22:51,804 INFO mapred.FileInputFormat: Total input files to process : 1
cluster-master | 2025-04-14 19:22:52,277 INFO mapreduce.JobSubmitter: number of splits:2
cluster-master | 2025-04-14 19:22:52,759 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1744658335314_0002
cluster-master | 2025-04-14 19:22:52,759 INFO mapreduce.JobSubmitter: Executing with tokens: []
cluster-master | 2025-04-14 19:22:52,878 INFO conf.Configuration: resource-types.xml not found
cluster-master | 2025-04-14 19:22:52,879 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
cluster-master | 2025-04-14 19:22:52,927 INFO impl.YarnClientImpl: Submitted application application_1744658335314_0002
cluster-master | 2025-04-14 19:22:52,953 INFO mapreduce.Job: The url to track the job: http://cluster-master:8088/proxy/application_1744658335314_0002/
cluster-master | 2025-04-14 19:22:52,955 INFO mapreduce.Job: Running job: job_1744658335314_0002
```

```
cluster-master | 2025-04-14 19:22:52,953 INFO mapreduce.Job: The url to track the job: http://cluster-master:8088/proxy/application_1744658335314_0002/
cluster-master | 2025-04-14 19:22:52,955 INFO mapreduce.Job: Running job: job_1744658335314_0002
cluster-master | 2025-04-14 19:23:05,169 INFO mapreduce.Job: Job job_1744658335314_0002 running in uber mode : false
cluster-master | 2025-04-14 19:23:05,172 INFO mapreduce.Job:  map 0% reduce 0%
cluster-master | 2025-04-14 19:23:09,229 INFO mapreduce.Job:  map 100% reduce 0%
cluster-master | 2025-04-14 19:23:25,303 INFO mapreduce.Job:  map 100% reduce 100%
cluster-master | 2025-04-14 19:24:28,542 INFO mapreduce.Job: Job job_1744658335314_0002 completed successfully
cluster-master | 2025-04-14 19:24:28,625 INFO mapreduce.Job: Counters: 54
cluster-master |        File System Counters
cluster-master |                    FILE: Number of bytes read=449
cluster-master |                    FILE: Number of bytes written=834210
cluster-master |                    FILE: Number of read operations=0
cluster-master |                    FILE: Number of large read operations=0
cluster-master |                    FILE: Number of write operations=0
cluster-master |                    HDFS: Number of bytes read=562
cluster-master |                    HDFS: Number of bytes written=0
cluster-master |                    HDFS: Number of read operations=11
cluster-master |                    HDFS: Number of large read operations=0
cluster-master |                    HDFS: Number of write operations=2
cluster-master |                    HDFS: Number of bytes read erasure-coded=0
cluster-master |        Job Counters
cluster-master |                    Launched map tasks=2
cluster-master |                    Launched reduce tasks=1
cluster-master |                    Data-local map tasks=2
cluster-master |                    Total time spent by all maps in occupied slots (ms)=5200
cluster-master |                    Total time spent by all reduces in occupied slots (ms)=75660
cluster-master |                    Total time spent by all map tasks (ms)=5200
cluster-master |                    Total time spent by all reduce tasks (ms)=75660
cluster-master |                    Total vcore-milliseconds taken by all map tasks=5200
cluster-master |                    Total vcore-milliseconds taken by all reduce tasks=75660
cluster-master |                    Total megabyte-milliseconds taken by all map tasks=5324800
```

```
cluster-master |                    Total megabyte-milliseconds taken by all map tasks=5324800
cluster-master |                    Total megabyte-milliseconds taken by all reduce tasks=77475840
cluster-master |        Map-Reduce Framework
cluster-master |                    Map input records=1
cluster-master |                    Map output records=35
cluster-master |                    Map output bytes=373
cluster-master |                    Map output materialized bytes=455
cluster-master |                    Input split bytes=290
cluster-master |                    Combine input records=0
cluster-master |                    Combine output records=0
cluster-master |                    Reduce input groups=29
cluster-master |                    Reduce shuffle bytes=455
cluster-master |                    Reduce input records=35
cluster-master |                    Reduce output records=0
cluster-master |                    Spilled Records=70
cluster-master |                    Shuffled Maps =2
cluster-master |                    Failed Shuffles=0
cluster-master |                    Merged Map outputs=2
cluster-master |                    GC time elapsed (ms)=135
cluster-master |                    CPU time spent (ms)=9310
cluster-master |                    Physical memory (bytes) snapshot=770801664
cluster-master |                    Virtual memory (bytes) snapshot=7778820096
cluster-master |                    Total committed heap usage (bytes)=706740224
cluster-master |                    Peak Map Physical memory (bytes)=291770368
cluster-master |                    Peak Map Virtual memory (bytes)=2592624640
cluster-master |                    Peak Reduce Physical memory (bytes)=272498688
cluster-master |                    Peak Reduce Virtual memory (bytes)=3588886528
cluster-master |        Shuffle Errors
cluster-master |                    BAD_ID=0
cluster-master |                    CONNECTION=0
cluster-master |                    IO_ERROR=0
```

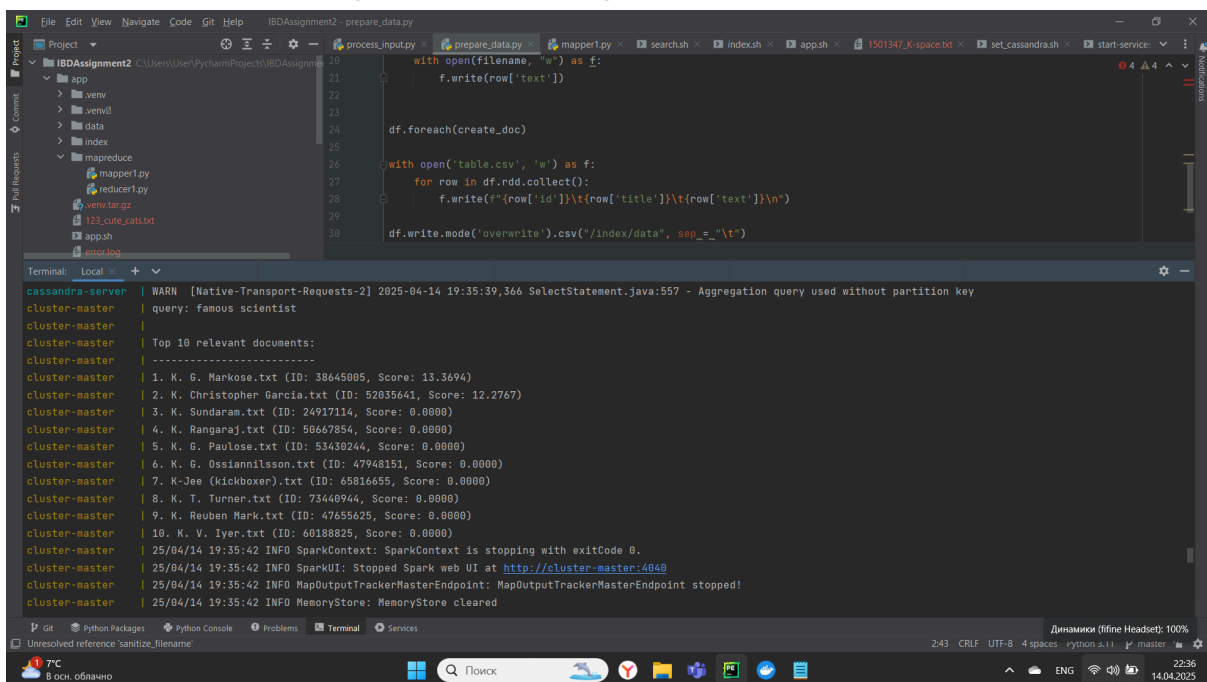As you can see, both files were successfully indexed!

Now for the ranker part. I have created three test queries:

`"famous scientist"`

`"cute cats"`

`"a spatial Fourier transform"`

I expected the first query to output some scientist, the second to rank the cute_cats file at the top, and the last one to rank K-space file high, since it is a quote from it. I got the following results:

The results for queries 2 and 3 were as expected, however for 1 the top ranked documents are K. G. Markose (an Indian singer) and K. Christofer Garcia (scientist, expectedly). I think the reason for the singer appearing in the results is the word "famous".

Overall, I am pleased with the results. Most of the queries work as expected, and even unexpected results have logical reasons.