# Assignment 2 Report:
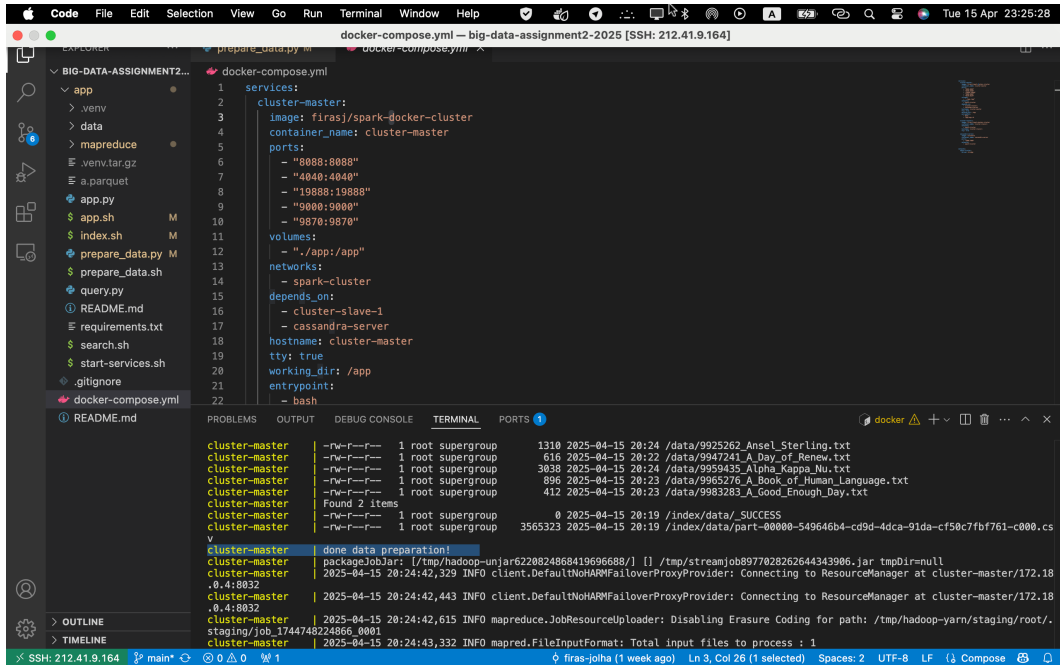# Simple Search Engine using Hadoop MapReduce

Liana Mardanova

DS-01

l.mardanova@innopolis.university

S25 - Big Data Course

April 15, 2025

# Overview

I created a simple search engine using MapReduce for indexing, Cassandra for storing statistics and Spark RDD for ranking using BM25. Initially, I prepared 1000 documents (from a.parquet) and stored them in HDFS. Afterwards, I ran MapReduce to get the info to calculate BM25 for the query and stored it in Cassandra. Next, I used PySpark to read this data and calculate BM25. Each part is described in more detail below.
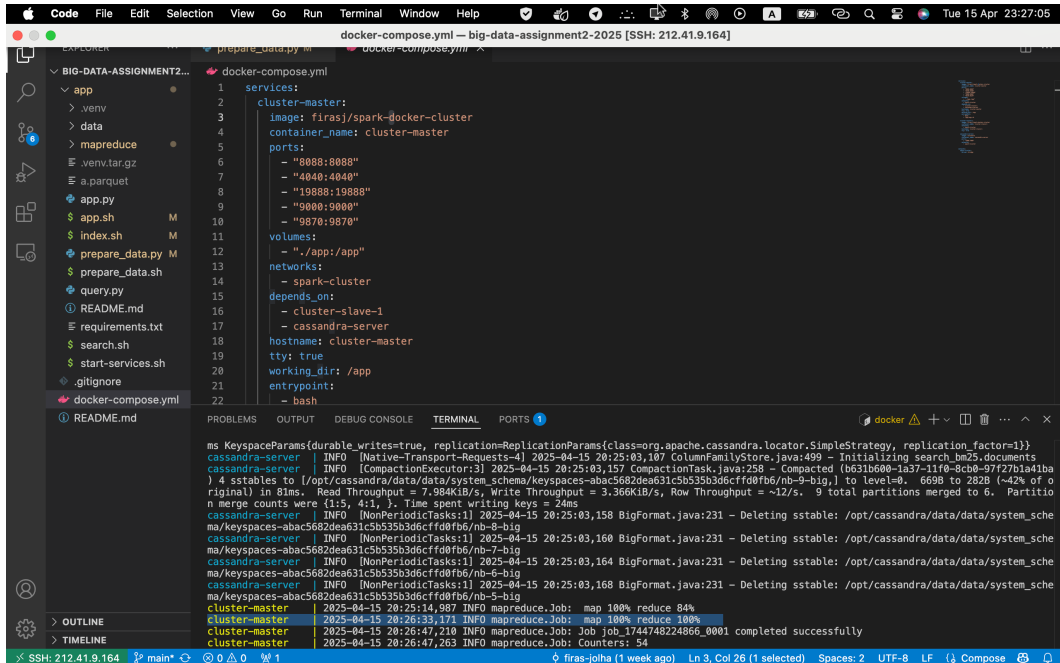
# Demo



**Figure 1:** Data preparation



**Figure 2:** Indexing 1000 documents using MapReduce

# MapReduce

Mapper (`mapper1.py`) accepts input in the format `<doc_id> <doc_title> <doc_text>` separated by tabs. It tokenizes the text, counts the total number of tokens, and uses a counter to compute term frequencies. It gives three outputs: document info (`doc_id`, `title`, `doc_len`), document frequency (`term`, `doc_id`), and term frequency (`term`, `doc_id`, `tf`).

Reducer (`reducer1.py`) reads this output, does the final aggregation, and saves the results into Cassandra tables.

# Cassandra Tables

I saved the tables in this format so that it would be convenient to calculate bm25 for query later.

### term_frequency

| No. | Column Name | Type |
|-----|-------------|------|
| 1 | term | TEXT |
| 2 | doc_id | TEXT |
| 3 | tf | INT |
| **Primary Key:** (term, doc_id) | | |

**Table 1:** Schema of `term_frequency` table

### document_frequency

| No. | Column Name | Type |
|-----|-------------|------|
| 1 | term | TEXT |
| 2 | df | INT |
| **Primary Key:** (term) | | |

**Table 2:** Schema of `document_frequency` table

### documents

| No. | Column Name | Type |
|-----|-------------|------|
| 1 | doc_id | TEXT |
| 2 | title | TEXT |
| 3 | doc_len | INT |
| **Primary Key:** (doc_id) | | |

**Table 3:** Schema of `documents` table

# Query and Ranking

I wrote a script to find top 10 for a query. It connects to Cassandra and first reads all the documents to get their names and length. This is needed to calculate $N$ and the average length of the document.

It then looks up each query term and gets how many documents it occurs in (df) and how often it occurs in each document (tf). After that, it calculates BM25 for each document considering all the query terms.

At the end it shows the top 10 documents with the highest scores. The formula I used is shown below.

$$BM25(q,d) = \sum_{t \in q} \log \left[\frac{N}{df(t)}\right] \cdot \frac{(k_1 + 1) \cdot tf(t,d)}{k_1 \cdot \left[(1-b) + b \cdot \frac{dl(d)}{dl_{avg}}\right] + tf(t,d)}$$

Where

- q: the user's query
- t: a term t in the user's query
- $N$ : number of documents in the corpus
- $df(t)$ : the document frequency or number of documents containing the term t.
- $k_1, b$: model's hyperparameters (e.g. $k_1 = 1, b = 0.75$ )
- $tf(t,d)$ : the term t frequency in the document d.
- $dl(d)$ : the length of the document d.
- $dl_{avg}$: average document length.

**Figure 3:** BM25 scoring formula used for ranking

# Conclusion

This was an interesting task because it combined many tools like Hadoop, Cassandra, and Spark. I got a chance to work with all of them together and see how they connect in a real pipeline.

Since it's a real Big Data task, I faced problems related to memory. My local machine ran out of memory, so I had to run everything on a remote server. It wasn't easy, but I learned a lot during the process.