# BigData Assignment 2 report

*Nikita Kurkulskiu*

April 14, 2025

# Methodology

## 1. Data Preparation

Begin by collecting a sample of plain text documents. Each document is stored in the `data/` folder and pushed to the Hadoop Distributed File System (HDFS). The documents follow a naming format `<doc_id>_<doc_title>.txt`, allowing us to parse `doc_id` in MapReduce jobs.

## 2. First MapReduce Pipeline (Inverted Index)

*Mapper*: For each document, the mapper reads lines of text and extracts its `doc_id` from the filename, then tokenizes the lines (using basic regex) into terms. For each term, the mapper emits `(term, doc_id, 1)`.

*Reducer*: The reducer aggregates partial term counts for each `(term, doc_id)` pair to produce the final term frequency (TF). Concurrently, it maintains the total number of terms per document (`doc_length`), so that each document's length is inserted into the `doc_length` table. The reducer stores:

- `(term, doc_id, freq)` tuples in the `inverted_index` table.
- `(doc_id, length)` in the `doc_length` table.

Both tables reside in Cassandra database under a dedicated keyspace (e.g., `search_keyspace`).

## 3. Second MapReduce Pipeline (Term Statistics)

*Mapper*: To compute document frequency (DF), each mapper outputs `(term, doc_id)` but ensures each term in a document is emitted only once. The mapper thus relies on a set-based approach to avoid duplicates within the same document.

*Reducer*: This reducer receives (`term`, `doc_id`) pairs and accumulates a set of unique documents for every term. It also looks up total document count and sums of lengths from the `doc_length` table to compute:

- Document frequency (`df`) for each term,
- Total number of documents (`n`),
- The average document length (`avg_doc_length`).

These are stored in a new Cassandra table `statistics(term, df, n, avg_doc_length)`, enabling BM25 computations later.

## 4. Retrieval (PySpark & BM25)

For query processing, developed a `query.py` Spark application that:

1. Reads a user query from the command line,
2. Loads `inverted_index`, `doc_length`, and `statistics` tables from Cassandra,
3. Prepares an RDD from the in-memory data (in a simple implementation) and filters for terms present in the user query,
4. Computes BM25 for each (`term`, `doc_id`) using aggregated statistics:

$$\text{BM25}(d, q) = \sum_{t \in q} \ln\left(\frac{n - \text{df}(t) + 0.5}{\text{df}(t) + 0.5}\right) \times \frac{\text{freq}(t, d)\,(k + 1)}{\text{freq}(t, d) + k\left(1 - b + b \cdot \frac{|d|}{|d|}\right)}$$

5. Ranks documents by their BM25 score and returns the top 10 matches.

Hence, the retrieval is orchestrated via Spark RDD operations, preserving distributed parallelism and upholding the assignment's requirement of avoiding single-machine frameworks.

## 5. Deployment on YARN

Both indexing pipelines (`index.sh`) and the BM25 retrieval job (`search.sh`) are packaged for deployment on a multi-node Hadoop cluster.

- The indexing pipelines run via Hadoop Streaming MapReduce (in either local or distributed mode).

- The `search.sh` script launches a `spark-submit` job that connects to the Cassandra cluster and processes queries across the YARN resource manager.

This architecture ensures a scalable approach: indexing leverages batch-oriented MapReduce, while real-time queries rely on Spark's low-latency RDD transformations to rank and retrieve documents.

# Demonstration

to run the repository enter the cluster-master node and run bash app.sh
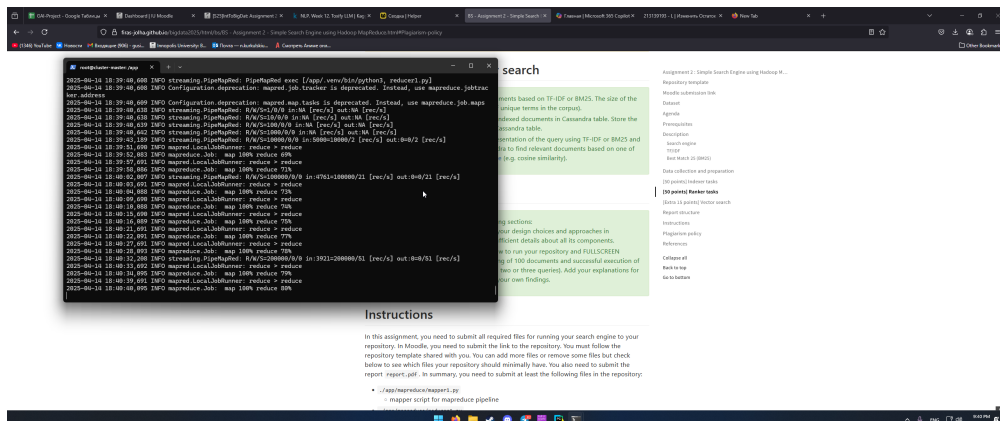


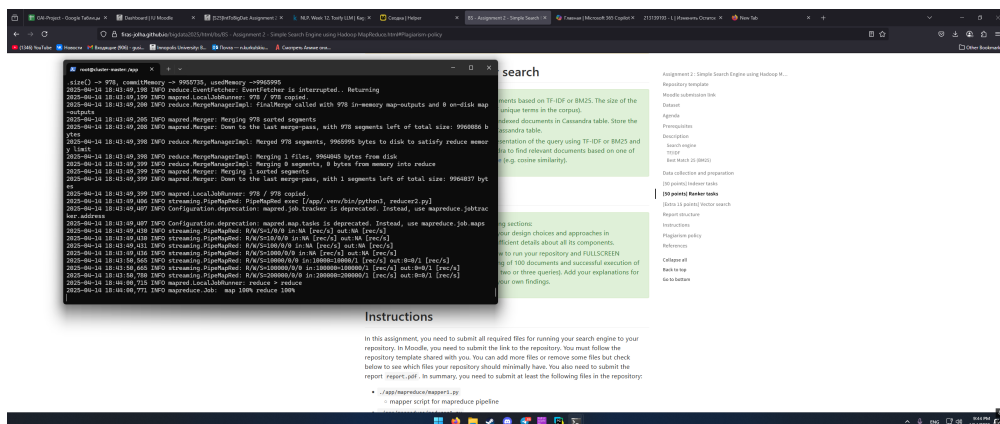Fig. 1. Cluster running MapReduce and Spark jobs



Fig. 2. Cluster ended running MapReduce and Spark jobs

Fig. 3. top 10 for query "this is a query!"



Fig. 4. top 10 for query "hello"



Fig. 5. top 10 for query "database"