# Assignment №2 BigData Report

Amir Nigmatullin "am.nigmatullin@innopolis.university"

April 15, 2025

## 1 Methodology

### 1.1 Preprocessing

First of all, I preprocess data stored in HDFS. In preprocessing I make lowercasing, removing punctuation and stopwords and tokenizing. The output is stored back in HDFS for the next stages.

### 1.2 MapReduce Pair 1: Document Statistics

In first Map and Reduce pair, these jobs were processed:
  1) Computing unique terms (Vocabulary)
  2) Computing document frequencies
  3) Computing lengths and average document length

### 1.3 MapReduce Pair 2: Inverted Index Construction

In the second Map and Reduce pair I build inverted index and store for each term document IDs (where it appears) and the frequency of the term in each document.

This index is stored in the Cassandra table `inverted_index(term, doc_id, freq)`.

### 1.4 BM25 Search Engine

I deployed a search facility that did the following:

1. Preprocessed the input query the same way we handle documents.

2. For every term in the query, fetched the corresponding posting list from Cassandra.

3. Calculated the BM25 score for each candidate document.

4. Returned the top-ranked results.

I used this formula for calculating BM25. It contains from 2 parts. First:

$$\sum_{i \in q} \log \left( \frac{N - df_i + 0.5}{df_i + 0.5} + 1 \right)$$

Second:

$$\frac{(k_1 + 1) \cdot tf_{i,d}}{tf_{i,d} + k_1 \cdot (1 - b + b \cdot \frac{|d|}{avgdl})}$$

And in the result we multiply this 2 parts and get result, our score.

## 1.5 Data Storage in Cassandra

- **vocabulary**: This table stores the full set of unique terms extracted from the corpus. Each row consists of a term (the word itself) and its doc_freq.

```
CREATE TABLE IF NOT EXISTS vocabulary (
    term TEXT PRIMARY KEY,
    doc_freq INT
)
```

- **inverted_index**: This table represents the core inverted index. It maps each term to the documents it appears in (doc_id), along with the term_freq (how often the term occurs in that document), and the positions (list of word indices where the term appears in the document).

```
CREATE TABLE IF NOT EXISTS inverted_index (
    term TEXT,
    doc_id TEXT,
    term_freq INT,
    positions LIST<INT>,
    PRIMARY KEY (term, doc_id)
)
```

- **doc_stats**: This table stores metadata about each document, including its doc_id, title, total number of terms (total_terms), and number of unique terms (unique_terms).

```
CREATE TABLE IF NOT EXISTS doc_stats (
    doc_id TEXT PRIMARY KEY,
    title TEXT,
    total_terms INT,
    unique_terms INT
)
```

## 1.6 About other files

- **app.py**: The main Python script handles loading the preprocessed documents into Cassandra. It reads document-term stats like term frequencies and document lengths, then stores the inverted index, document statistics, and vocabulary into their respective Cassandra tables..

- **query.py**: This script handles BM25 ranking using the indexed data stored in Cassandra. It works with free-text queries and returns the highest-ranked documents based on their BM25 scores.

- **docker-compose.yml**: This file sets up and manages a multi-container Docker environment. It launches Hadoop for MapReduce, Spark for preprocessing, and Cassandra for storage, all in a way that's reproducible and easy to move across systems. Each service is configured with the necessary ports, volumes, and dependencies.

- **app.sh**: This shell script runs the full ingestion pipeline. It starts by launching the Docker setup, then runs preprocessing with Spark, and finally ingests the data using app.py.

- **index.sh**: This shell script handles the Hadoop-based MapReduce jobs that build the inverted index from a text corpus stored in HDFS. It takes care of compiling the Java code (if needed), runs the Map and Reduce steps, and saves the final output back into HDFS.

# 2 Demonstration

Picture 1: In this picture I run docker compose up and you can see that it starts running without errors.

Picture 2: In this picture you can see intermediate stage of working, where you can see cluster-slave 1 and 2, their info and etc.



Picture 3: Here you can see that I took 1000 documents and print the content of hdfs.

Picture 4: Here you can see my mapping and reducing jobs, they are successful and result of their job is printed.



Picture 5: Here you can see my second mapping and reducing jobs, they are successful and result of their job is also printed.

Picture 6: Here you can see that I connect to Cassandra server and indexing completed successfully and data stored in Cassandra.



Picture 7: Here you can see the result of the first query.

Picture 8: Here you can see my tables in Cassandra. I have 3 tables, as I described in the methodology. And you can see how they look like.

```
zaurall@zaurall: ~/Documents/developer/INNO_S25/BD/amir/big-data-assignment2-2025
amir@amir:~$docker ps
CONTAINER ID   IMAGE                        COMMAND                 CREATED          STATUS         PORTS
                                                                                                    NAMES
f3ecbee615a3   firasj/spark-docker-cluster  "/bin/bash -c 'servi…"  About an hour ago  Up 32 minutes   2122/tcp, 7001-7007/tcp, 8020/tcp, 8030-8033/tcp, 8040/tcp, 8042/tcp, 8088/tcp, 8888
/tcp, 9000/tcp, 9870/tcp, 10020/tcp, 19888/tcp, 49707/tcp, 50010/tcp, 50020/tcp, 50070/tcp, 50075/tcp, 50090/tcp   cluster-slave-2
c8f6e7265455   cassandra                    "docker-entrypoint.s…"  About an hour ago  Up 32 minutes   7001/tcp, 0.0.0.0:7000->7000/tcp, 7199/tcp, 0.0.0.0:9042->9042/tcp, 9160/tcp
                                                                                                    cassandra-server
9ab9eec503c3   firasj/spark-docker-cluster  "/bin/bash -c 'servi…"  About an hour ago  Up 32 minutes   2122/tcp, 7001-7007/tcp, 8020/tcp, 8030-8033/tcp, 8040/tcp, 8042/tcp, 8088/tcp, 8888
/tcp, 9000/tcp, 9870/tcp, 10020/tcp, 19888/tcp, 49707/tcp, 50010/tcp, 50020/tcp, 50070/tcp, 50075/tcp, 50090/tcp   cluster-slave-1
amir@amir:~$docker exec -it cassandra-server cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.2.0 | Cassandra 5.0.4 | CQL spec 3.4.7 | Native protocol v5]
Use HELP for help.
cqlsh> use
index_keyspace       system                system_auth         system_distributed    system_schema      system_traces       system_views        system_virtual_schema
cqlsh> use index_keyspace ;
cqlsh:index_keyspace> select * from
doc_stats          inverted_index       system_auth.        system_schema.       system_views.      vocabulary
index_keyspace.      system.             system_distributed.  system_traces.       system_virtual_schema.
cqlsh:index_keyspace> select * from doc_stats limit 5;

 doc_id    | title                                  | total_terms | unique_terms
-----------+----------------------------------------+-------------+--------------
 10230685  |                   A Dead Sinking Story |         115 |           83
 27568194  | A Hero Ain't Nothin' but a Sandwich (film) |     567 |          297
 39710446  |                  A Little Bit of Luck |         330 |          155
 38294693  | A Change Is Gonna Come (Jack McDuff album) |     338 |          189
 51794980  |      A Family Secret (Upstairs, Downstairs) |     306 |          171

(5 rows)
cqlsh:index_keyspace> select * from inverted_index  limit 5;

 term   | doc_id   | positions                | term_freq
--------+----------+--------------------------+-----------
 dobson | 13633480 |                     [68] |         1
   sain | 14404655 |   [1840, 1339, 1108, 1568] |       4
 bessus | 12000397 |        [1017, 1000, 1029] |       3
     ix | 19789501 |                    [542] |         1
     ix | 32497421 |                     [34] |         1

(5 rows)
cqlsh:index_keyspace> select * from vocabulary   limit 5;

 term   | doc_freq
--------+----------
 dobson |        1
   sain |        1
 bessus |        1
     ix |        4
  await |        2

(5 rows)
cqlsh:index_keyspace>
```
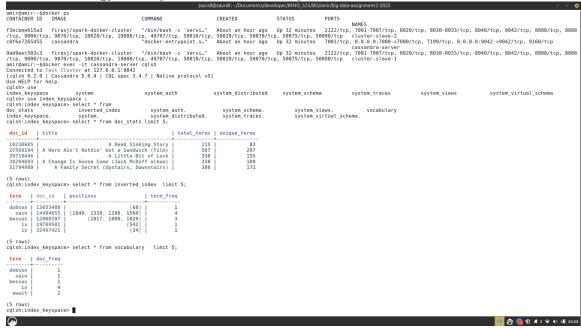
Picture 9: Here you can see the result of the second (custom) query, that I wrote by myself.

```
zaurall@zaurall: ~
cluster-master   | 25/04/15 17:41:56 INFO BlockManagerInfo: Added broadcast_4_python on disk on cluster-slave-1:45603 (size: 10.6 KiB)
cluster-master   | 25/04/15 17:41:56 INFO BlockManagerInfo: Added broadcast_0_piece0 in memory on cluster-slave-1:45603 (size: 248.0 B, free: 366.3 MiB)
cluster-master   | 25/04/15 17:41:56 INFO BlockManagerInfo: Added broadcast_0_python on disk on cluster-slave-1:45603 (size: 55.0 B)
cluster-master   | 25/04/15 17:41:56 INFO TaskSetManager: Finished task 0.0 in stage 3.0 (TID 6) in 189 ms on cluster-slave-2 (executor 2) (1/2)
cluster-master   | 25/04/15 17:41:56 INFO TaskSetManager: Finished task 1.0 in stage 3.0 (TID 7) in 209 ms on cluster-slave-1 (executor 1) (2/2)
cluster-master   | 25/04/15 17:41:56 INFO YarnScheduler: Removed TaskSet 3.0, whose tasks have all completed, from pool
cluster-master   | 25/04/15 17:41:56 INFO DAGScheduler: ResultStage 3 (top at /app/query.py:114) finished in 0.221 s
cluster-master   | 25/04/15 17:41:56 INFO DAGScheduler: Job 1 is finished. Cancelling potential speculative or zombie tasks for this job
cluster-master   | 25/04/15 17:41:56 INFO YarnScheduler: Killing all running tasks in stage 3: Stage finished
cluster-master   | 25/04/15 17:41:56 INFO DAGScheduler: Job 1 finished: top at /app/query.py:114, took 0.380639 s
cluster-master   |
cluster-master   | Top 10 documents for query: How to learn to do backflip?
cluster-master   | ------------------------------
cluster-master   | Document ID: 46835946
cluster-master   | Title: A Damaged Mirror
cluster-master   | Score: 9.6316
cluster-master   | ------------------------------
cluster-master   | Document ID: 41951836
cluster-master   | Title: A Arma Escarlate
cluster-master   | Score: 7.6061
cluster-master   | ------------------------------
cluster-master   | Document ID: 58622515
cluster-master   | Title: A Boy and a Priest
cluster-master   | Score: 7.4477
cluster-master   | ------------------------------
cluster-master   | Document ID: 41086718
cluster-master   | Title: A Fair to Remember (Modern Family)
cluster-master   | Score: 6.8893
cluster-master   | ------------------------------
cluster-master   | Document ID: 42403439
cluster-master   | Title: A Fairly Odd Summer
cluster-master   | Score: 6.8518
cluster-master   | ------------------------------
cluster-master   | Document ID: 5446583
cluster-master   | Title: A House on a Street in a Town I'm From
cluster-master   | Score: 6.3815
cluster-master   | ------------------------------
cluster-master   | Document ID: 50242777
cluster-master   | Title: A Couple of Poor, Polish-Speaking Romanians
cluster-master   | Score: 6.2015
cluster-master   | ------------------------------
cluster-master   | Document ID: 37916885
cluster-master   | Title: A Day Late and a Dollar Short (novel)
cluster-master   | Score: 6.1461
cluster-master   | ------------------------------
cluster-master   | Document ID: 56610847
cluster-master   | Title: A Legacy of Spies
cluster-master   | Score: 6.0178
cluster-master   | ------------------------------
cluster-master   | Document ID: 60121915
cluster-master   | Title: A Brief History of Everyone Who Ever Lived
cluster-master   | Score: 5.6352
cluster-master   | ------------------------------
cluster-master   | 25/04/15 17:41:56 INFO SparkContext: SparkContext is stopping with exitCode 0.
cluster-master   | 25/04/15 17:41:56 INFO SparkUI: Stopped Spark web UI at http://cluster-master:4040
cluster-master   | 25/04/15 17:41:56 INFO YarnClientSchedulerBackend: Interrupting monitor thread
cluster-master   | 25/04/15 17:41:56 INFO YarnClientSchedulerBackend: Shutting down all executors
cluster-master   | 25/04/15 17:41:56 INFO YarnSchedulerBackend$YarnDriverEndpoint: Asking each executor to shut down
cluster-master   | 25/04/15 17:41:57 INFO YarnClientSchedulerBackend: YARN client scheduler backend Stopped
cluster-master   | 25/04/15 17:41:57 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
cluster-master   | 25/04/15 17:41:57 INFO MemoryStore: MemoryStore cleared
cluster-master   | 25/04/15 17:41:57 INFO BlockManager: BlockManager stopped
cluster-master   | 25/04/15 17:41:57 INFO BlockManagerMaster: BlockManagerMaster stopped
cluster-master   | 25/04/15 17:41:57 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
cluster-master   | 25/04/15 17:41:57 INFO SparkContext: Successfully stopped SparkContext
cluster-master   | 25/04/15 17:41:57 INFO ShutdownHookManager: Shutdown hook called
cluster-master   | 25/04/15 17:41:58 INFO ShutdownHookManager: Deleting directory /tmp/spark-b011e977-2852-4e3f-98f2-343341e8c836
cluster-master   | 25/04/15 17:41:58 INFO ShutdownHookManager: Deleting directory /tmp/spark-07992f42-4f9b-476e-8179-87acdede8933
cluster-master   | 25/04/15 17:41:58 INFO ShutdownHookManager: Deleting directory /tmp/spark-07992f42-4f9b-476e-8179-87acdede8933/pyspark-58dbe212-68ae-46d3-82d1-51c4ee6dd22f
cluster-master exited with code 0
■
w Enable Watch
```

Picture 10: Here you can see the result of the third (custom) query, that I wrote by myself.



Picture 11: I forgot to check it before, but here you can see the size of the my tables in Cassandra.