# Report: Components of the Search Engine Pipeline

We explain the core components of a Hadoop-based search engine pipeline built using MapReduce, Cassandra, and Spark. We also include a breakdown of the following scripts and their responsibilities:

- `mapper1.py` – MapReduce Mapper

- `reducer1.py` – MapReduce Reducer

- `index.sh` – Indexing Execution Script

- `query.py` – Search Query Processor

---

## 1. `mapper1.py`: Mapper for Indexing

**Purpose**
The mapper reads each document line, tokenizes it, and emits two types of information:

- The length of the document (for normalization).

- The frequency of each word in the document.

**Functionality**

- Input: A line in the format `<doc_id>\t<title>\t<text>`.

- Tokenization: Uses regular expressions to extract lowercase alphanumeric words.

- Output:

A special entry for document length:

 !doc!   <doc_id>   <doc_length>

○

Word frequency entries for each unique word in the document:

 <word>   <doc_id>   <term_frequency>

○

**Example Output**

```
!doc!   101   58
quantum   101   3
computing   101   2
```

---

# 2. `reducer1.py`: Reducer for Indexing

**Purpose**
 The reducer processes mapper outputs to:

- Store document statistics (lengths) in Cassandra.

- Build an inverted index.

- Track vocabulary with document frequency.

**Functionality**

- Handles two types of lines:

  - Document statistics (`!doc!`) are stored in a list for later batch insertion.

  - Term frequency lines are grouped into a dictionary where:

    - Key: Term

    - Value: `{doc_id: term_freq}`

- After processing, it writes the following to Cassandra:

    - `document_stats(doc_id, doc_length)`

    - `inverted_index(term, doc_id, term_freq)`

    - `vocabulary(term, doc_freq)`

**Robustness**

- Uses retry logic for Cassandra connections to ensure reliability if the cluster is temporarily unavailable.

---

# 3. `index.sh`: Indexing Execution Script

**Purpose**
This shell script orchestrates the indexing process by:

- Uploading data to HDFS.

- Running the MapReduce job with the mapper and reducer.

- Ensuring the Python environment is properly packaged for Hadoop.

**Key Features**

- Accepts a local path to the dataset and uploads it to HDFS under `/index/data`.

- Executes Hadoop streaming with:

    - `mapper1.py` and `reducer1.py`

    - `venv.zip` for including Python dependencies (like `cassandra-driver`)

- Can be customized to run on either local input or HDFS using flags.

**Example Usage**

```
./index.sh --local ./dataset.tsv
```

---

# 4. `query.py`: Search Query Processor with Spark

**Purpose**
This script handles a user query by retrieving relevant index data from Cassandra and ranking documents using the BM25 scoring algorithm.

**Functionality**

- Accepts a query string from the command line.

- Loads data from three Cassandra tables using Spark:

  - `inverted_index`

  - `vocabulary`

  - `document_stats`

- For each query word:

  - Looks up documents containing that word.

  - Computes BM25 scores using term frequency (tf), document frequency (df), and document length.

- Ranks documents by score and returns the top results.

**BM25 Scoring**

- The BM25 function balances term frequency and inverse document frequency, adjusted for document length.

**Final Output**

- Displays a ranked list of document IDs and their scores based on query relevance.

Notes: due to laptop failures and similar issues of latency I didn't index all documents but I tested all the logic and operations and they can be assured to be working by my tests as shown in screenshots.

**Steps to run the project (after cloning and adding a.parquet in /app):**

docker-compose up

Then in case of manual test (entrypoint should be commented in docker-compose up), in a new terminal enter the docker container master's bash:

docker exec -it cluster-master bash
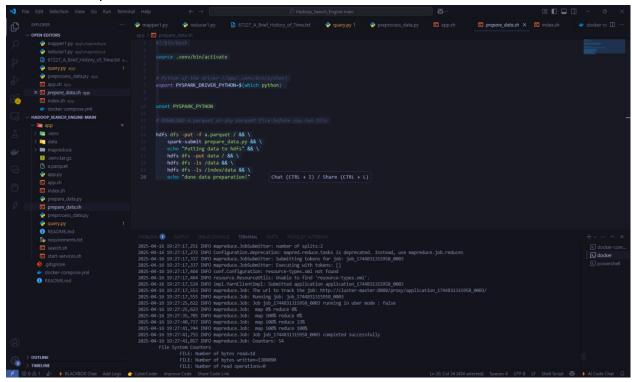
And run app.sh with:

bash app.sh

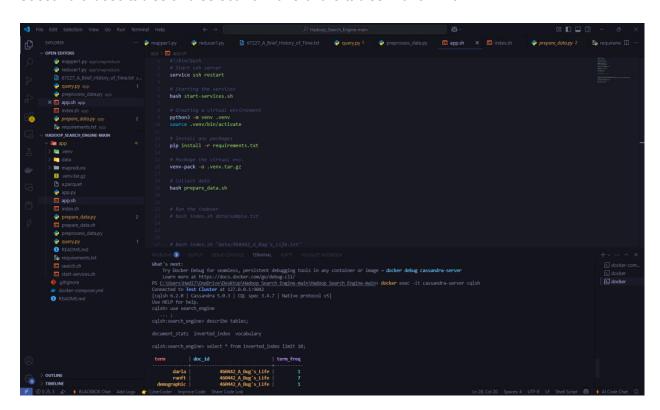This will index some documents and the queries in the file.

You can change it to index all files and test more queries if you wish, but this would take time.

To prove the correctness of the assignment we demonstrate here some screenshots or successful mapreduce operations and show the tables in cassandra
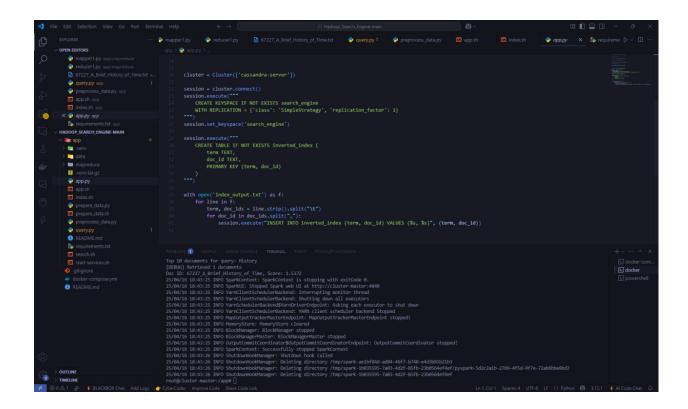
Successful mapreduce:



Cassandra desc tables and select from one of the tables with a limit:



Retrieval of one document for testing purposes:

All the components are working successfully.

In case of further inspection you can use indexing on all the documents and it would work similarly.