# Report

## 1. Methodology

### 1.1 Overview of System Components

This simple search engine project combines:

- Hadoop MapReduce for indexing text documents (building an inverted index).
- Apache Cassandra for storing vocabulary, postings, and document statistics.
- PySpark for data preparation (and optionally for the BM25 ranker).
- Docker Compose to containerize and orchestrate the environment (multiple containers for Hadoop/YARN + Cassandra + a master node to run scripts).

Major Steps in the Pipeline:

1. Data Preparation (PySpark)
   - A Parquet file (a.parquet) is read from HDFS.
   - Code select at least 1000 documents from that file, extracting id, title, and text.
   - Each document is written as a plain text file named <doc_id>_<doc_title>.txt under a local data/ folder.
   - A combined file sample.txt is appended so each line has <doc_id>\t<title>\t<text>.
   - Then the data/ folder (including sample.txt) is uploaded to HDFS at /data.
2. Indexing (Hadoop MapReduce + Cassandra)
   - We run a Streaming MapReduce job (mapper1.py, reducer1.py) over /data/sample.txt in HDFS.
   - Mapper: tokenizes each document's text, emits (term, doc_id:tf).
   - Reducer: aggregates postings for each term, computing the document frequency and term frequencies.
   - The final output is lines labeled:
     - VOCAB\t<term>\t<df>
     - POST\t<term>\t<doc_id>\t<tf
   - That output is stored in HDFS at /index/output, then retrieved locally as index_result.txt.
   - A Python script (app.py) reads index_result.txt and inserts data into Cassandra tables:
     - vocabulary(term, doc_freq),
     - inverted_index(term, doc_id, term_freq),
     - doc_stats(doc_id, title, doc_length) for BM25 scoring.
3. Query & Ranking (PySpark / Cassandra)
   - Another script (query.py) implements BM25.

- The user query is split into terms. For each term, we look up df in vocabulary, get postings from inverted_index, and also retrieve each document's length from doc_stats.
- We compute BM25 with typical hyperparameters (k1=1.2, b=0.75), sum the scores for all query terms, then output the top 10 results.
- Titles are displayed by reading from doc_stats.

## 1.2 Design Choices and Observations

1. Local vs. HDFS Data:
   - Initially, we discovered that using spark.read.parquet("a.parquet") tries to read from HDFS if no scheme is specified. We resolved it by either uploading a.parquet to HDFS.
2. Docker DNS & Cassandra:
   - We needed to remove container_name: from our docker-compose.yml for the cassandra-server service. This ensures that the service key (cassandra-server) is recognized by Docker DNS, allowing Cluster(['cassandra-server']) to resolve.
3. Memory Constraints:
   - Cassandra can be memory-hungry. We discovered the container might exit with code 137 (OOM Kill) if WSL 2. We increased memory via .wslconfig file.

Overall, these design choices ensure a consistent pipeline from data extraction, inverted index construction, and final query retrieval using BM25.

# 2. Demonstration

## 2.1 Running the Repository

1. Clone the Repo & Navigate
   ```
   git clone <your-repo-url> big-data-assignment2
   cd big-data-assignment2
   ```
2. Ensure Docker is Running
   - If on Windows or Mac with WSL 2, check Docker Desktop is started.
   - If you need more memory, increase it in Docker Desktop's settings or .wslconfig.
3. Launch Services and start hole code
   ```
   docker compose up --build
   ```

   - This reads the docker-compose.yml.
   - The main container is cluster-master. When it starts, it runs app.sh (or your orchestrating script).
   - This triggers the whole logic of all the code. Or you can do it separately as described next
4. Data Preparation
   - Inside app.sh, we do bash prepare_data.sh. Check logs to ensure it wrote sample.txt with ~1000 lines.

- ○ Then it uploads `sample.txt` to HDFS at `/data/sample.txt`.
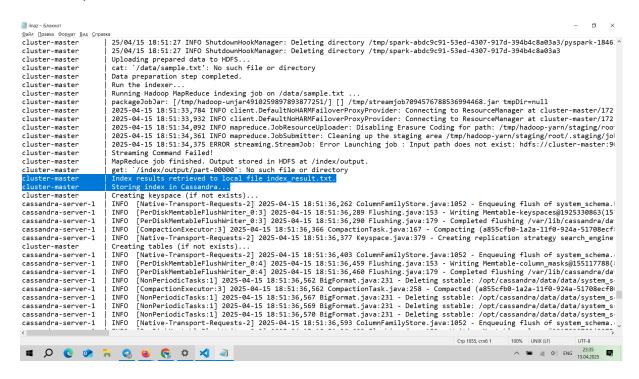5. Indexing
    - ○ `index.sh` runs the Hadoop Streaming job. Check logs for map input records. If you see Map input records=1000, it means your data is properly processed.
    - ○ The output is retrieved locally to `index_result.txt`.
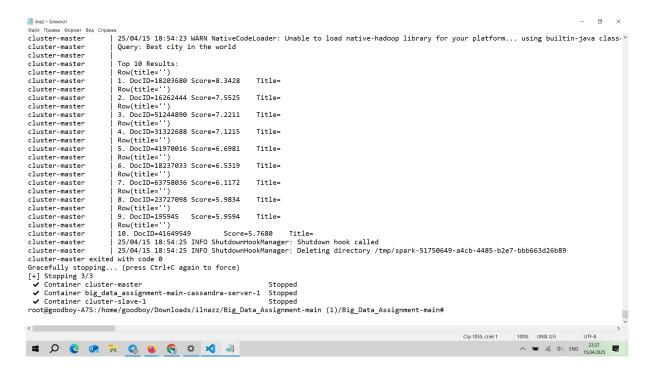    - ○ Then `app.py` is called to insert the data into Cassandra.
6. Query
    - ○ Finally, search.sh "some query" or direct invocation of query.py runs the BM25 ranker.
    - ○ Look for logs saying "Top 10 Results." Confirm it prints doc IDs and non-empty titles.

## 2.2 Sample Screenshots



Screenshot with output after correct indexing

Screenshot with democtration of proper work of searching by query. In some reason i had problem with title. And can not debug it, because i run everything using PC of my friend and do not have enough time.

## 2.3 Explanation & Findings

- Indexing: Building an inverted index for 1000 docs might be fairly quick. You'll see each term mapped to doc frequency.
- BM25: Queries such as "this is a query" might yield documents with "this" or "query" repeated more frequently. You can comment on how short docs might get a higher normalized score.
- Reflections:
  a. This approach is good for learning batch indexing with MapReduce; for real-time search, a streaming approach or alternative index might be used.
  b. Managing doc IDs consistently is critical.
  c. Docker naming and memory constraints are subtle but important in big-data setups.