# Overview

The search engine implementation is structured around a MapReduce pipeline, leveraging Hadoop for distributed processing and Cassandra for data storage.

Data Preparation

1. Data Ingestion: The input data is initially stored in a local directory and then uploaded to HDFS.
1. Data Formatting: Each document is processed to extract relevant metadata, such as document ID and title

MapReduce Pipeline (three reducers and mappers)

1. Document Indexing:
   - Mapper: Processes each document to extract terms and their occurrences.
   - Reducer: Aggregates term occurrences across documents to build an inverted index.
1. Document Statistics:
   - Mapper: Computes basic statistics for each document, such as term frequency.
   - Reducer: Aggregates these statistics to provide a comprehensive view of document characteristics.
1. Term Statistics:
   - Mapper: Extracts term-specific data, such as document frequency.
   - Reducer: Compiles term statistics to support relevance scoring.

Data Storage

- Cassandra Integration: The final indices and statistics are stored in Cassandra, enabling efficient retrieval and query execution. This choice supports the scalability and fault tolerance required for handling large datasets.

Query Execution

- Spark RDDs: Queries are executed using Spark, which retrieves data from Cassandra and computes relevance scores using the BM25 algorithm.

# Demonstration
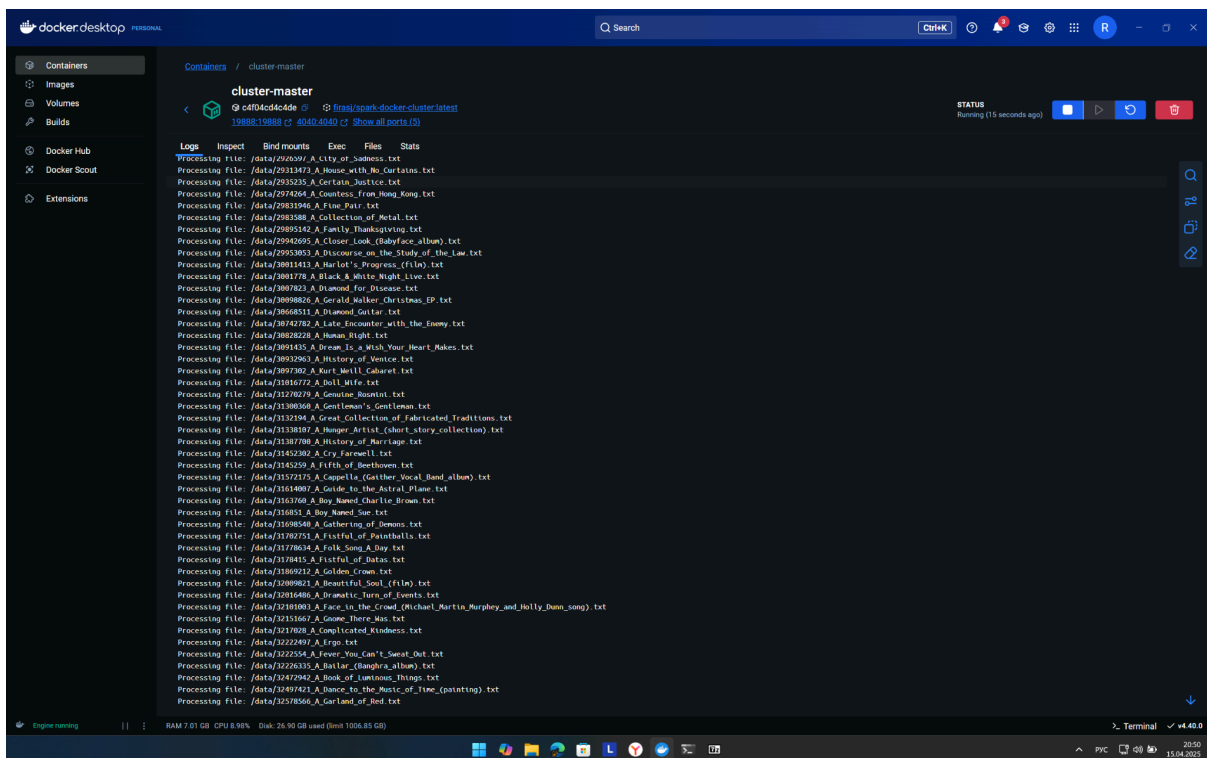
Execute the following command in your terminal:

```
docker compose up --build
```
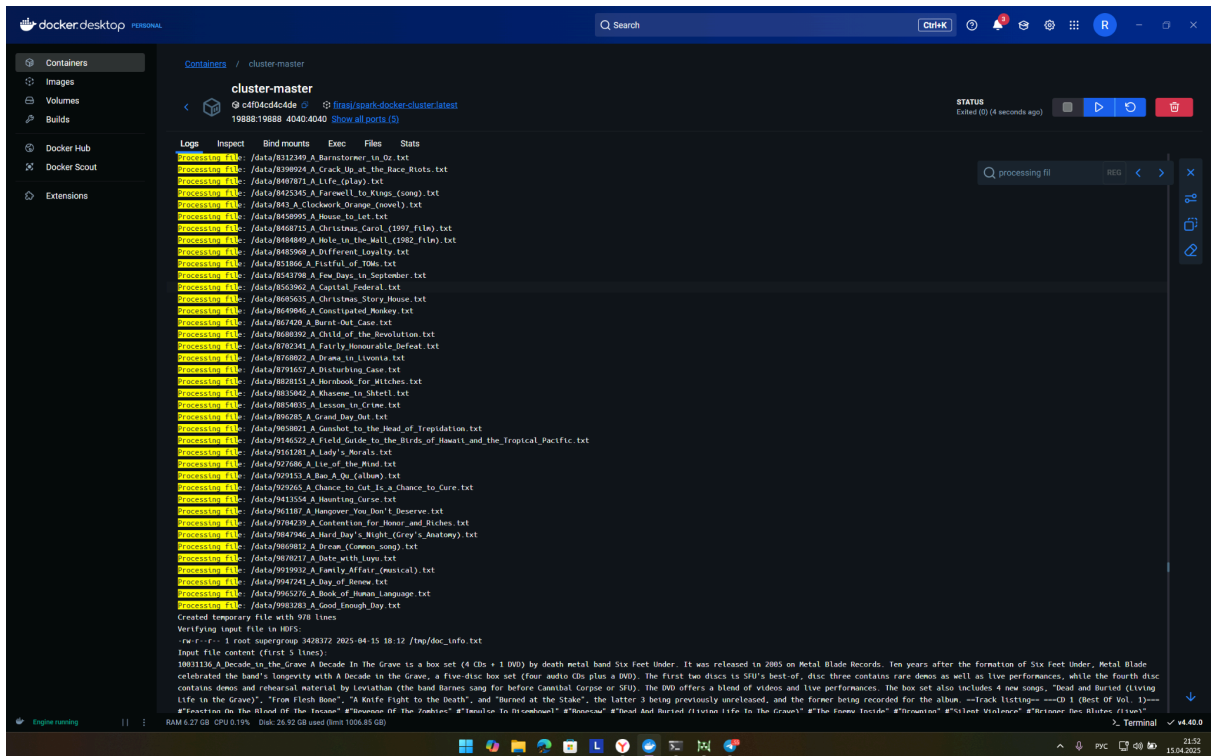
This command will build and launch all necessary services.

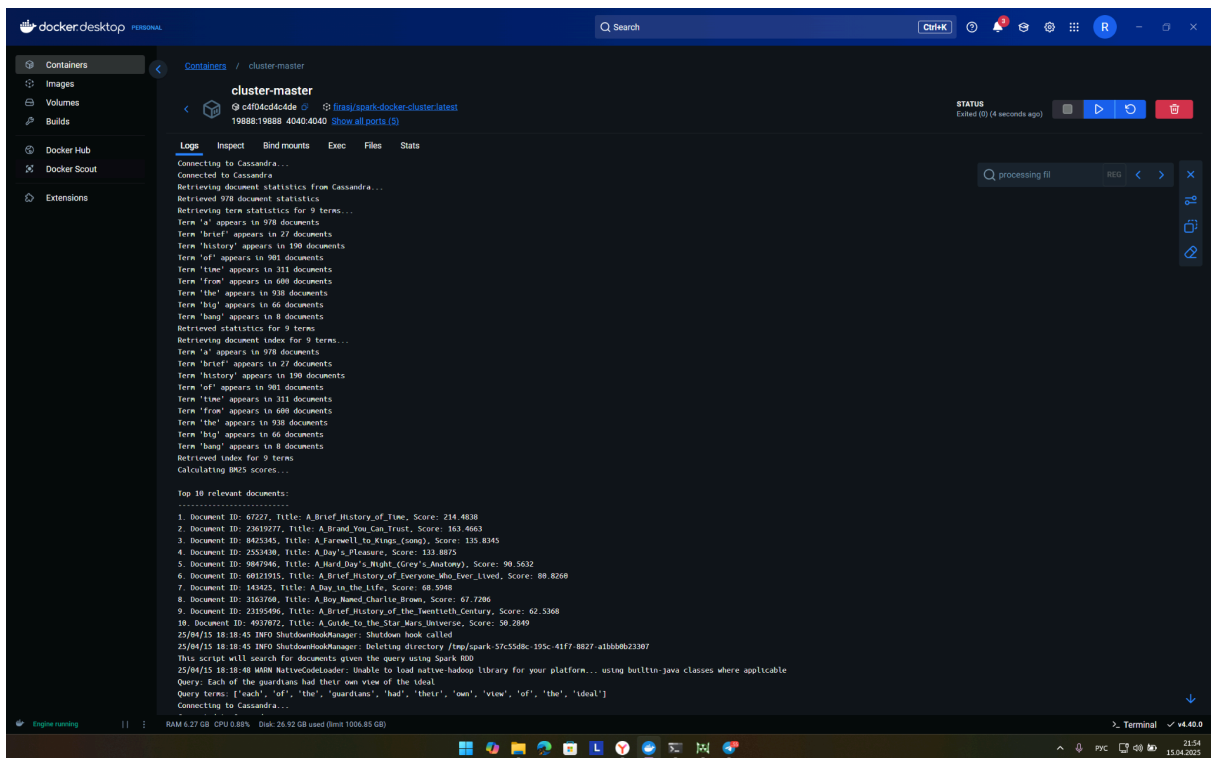All files located in the `/data` folder will be automatically processed by the system.

The `app.sh` script will execute a predefined query.
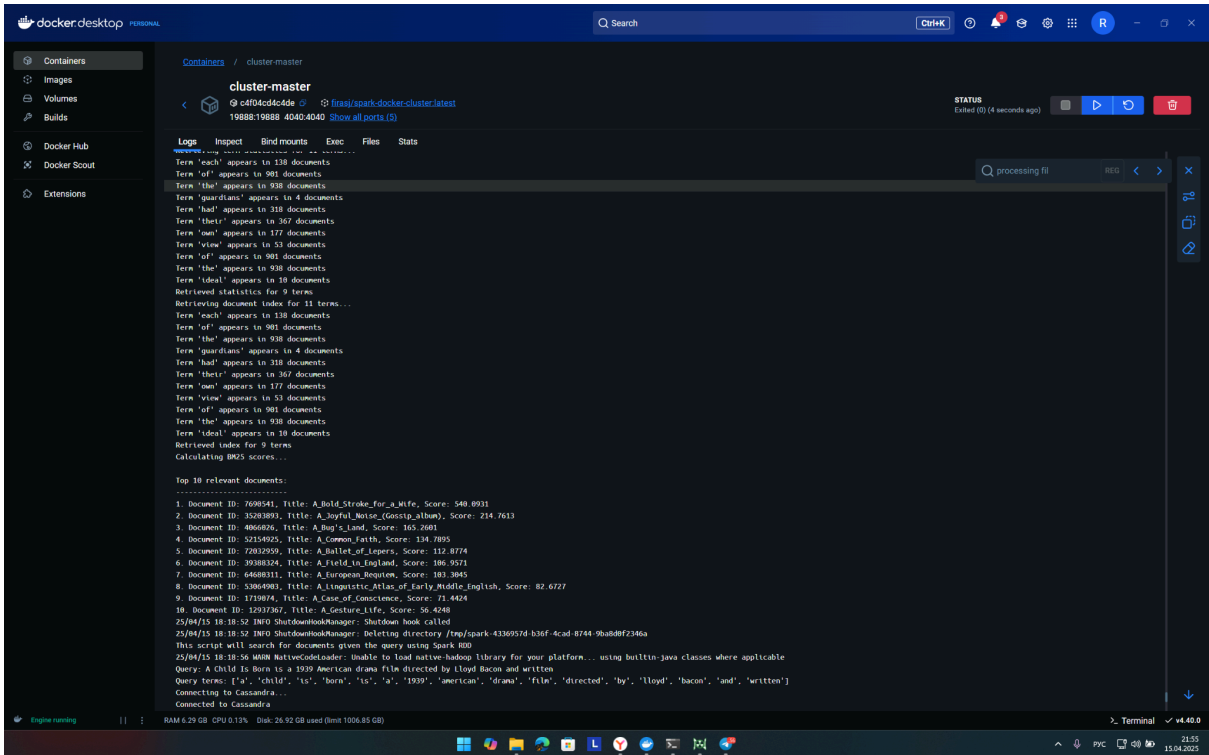
All 1000 documents are indexed successfully:

The first query is "A Brief History of Time: From the Big Bang" from book "67227_A_Brief_History_of_Time", as we can see stats are calculated correctly and the expected document ranked 1st:

Second query from random book "Each of the guardians had their own view of the ideal", the same, stats are calculated:



Third query "A Child Is Born is a 1939 American drama film directed by Lloyd Bacon and written" from the book named the smae ""A Child Is Born…", and its being ranked 1st.