# Big Data Assignment 2

Miasoedov Artemii

## Methodology

Most search engines consist of 2 parts:

1. Indexer - collects information about the known documents for fast retrieval in future

2. Ranker - uses the prepared index to find documents relevant to the query

## Indexer

To efficiently handle large amounts of data we use Hadoop's mapreduce pipeline. I decided to implement a pipeline consisting of 2 steps. First step handles each document separately, calculating term frequency inside each document. Imagine having 2 documents:

```
doc1: hello hello world
doc2: hello friend
```

Mapper receives the text from documents and splits in individual terms, producing output in the following format:

```
friend doc2 1
hello doc1 1
hello doc1 1
hello doc2 1
world doc1 1
```

Reducer calculates the number of term occurrences inside the document and inserts this data in Cassandra database.

```
friend doc2 1
hello doc1 2
hello doc2 1
world doc1 1

INSERT INTO tdf (term, doc, freq)
VALUES ('friend', 'doc2', 1),
       ('hello', 'doc1', 2),
       ('hello', 'doc2', 1),
       ('world', 'doc1', 1);
```

However, to build TF-IDF ranking we also need to know the length of each document and the total number of occurrences of each token. The data from the previous step might be used to

solve both of the tasks with a single reducer (counting length of each document and overall frequency of term), so mapper should produce data in a format similar to this one:

```
doc1 2
doc1 1
doc2 1
doc2 1
friend 1
hello 2
hello 1
world 1
```

Reducer will sum the values for the same key, thus finding document lengths and token frequencies:

```
doc1 3
doc2 2
friend 1
hello 3
world 1
```

However, we can't distinguish document names and terms for sure, so mapper should add some suffixes, specifying what kind of data it is:

```
doc1:dl 2
doc1:dl 1
doc2:dl 1
doc2:dl 1
friend:tf 1
hello:tf 2
hello:tf 1
world:tf 1
```

Reducer will insert obtained lengths and frequencies in the database:

```sql
INSERT INTO dl (doc, len)
VALUES ('doc1', 3),
       ('doc2', 3);

INSERT INTO tf (term, docs, total)
VALUES ('friend', 1, 1),
       ('hello', 2, 3);
```

Now, we have all the required data to do ranking.

# Ranking

To find the most relevant documents for a user query, we use a BM25 ranking algorithm implemented with PySpark and Cassandra. The process begins by tokenizing the input query, which breaks it into individual terms.

These query terms are then matched against the term frequency (tf), document frequency (tdf), and document length (dl) data previously stored in Cassandra. Each table is loaded into Spark as a DataFrame, and filters are applied to only consider terms that appear in the current query.

```python
tf = spark.read \
        .format(cassandra) \
        .options(table="tf", keyspace=keyspace) \
        .load() \
        .filter(col("term").isin(query)) \
        .alias("tf")
```

The loaded DataFrames are then joined on common fields. The goal is to bring together all relevant term-document pairs with their corresponding frequency and document statistics.

```python
index = dl \
        .join(tdf, "doc") \
        .join(tf, "term")
```

To compute BM25 scores, we first calculate the average document length and the total number of documents. In a more efficient scenario, this information should be stored, but for demonstration it is easier to be calculated.

```python
agg = dl.agg(avg("len"), count("*")).first()
avgdl, n = agg
```

Finally, obtained metrics are used to calculate BM-25 score:

```python
tf = (col("freq") * (k+1)) / (col("freq") + k * (1-b + b*n/avgdl))
idf = log(1 + (n - col("docs") + 0.5) / (col("docs") + 0.5))
index = index \
    .select(col("term"), col("doc"), (tf * idf).alias("score")) \
    .groupby("doc") \
    .agg(sum("score").alias("score")) \
    .orderBy(desc("score"))
```

# Demonstration

## Indexing



Starting processing files.
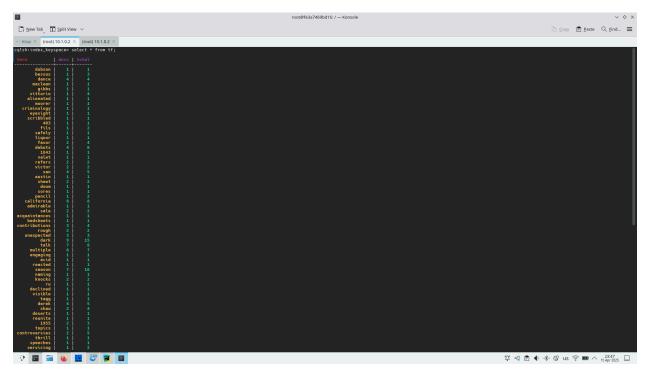
Successfully completed the first step.



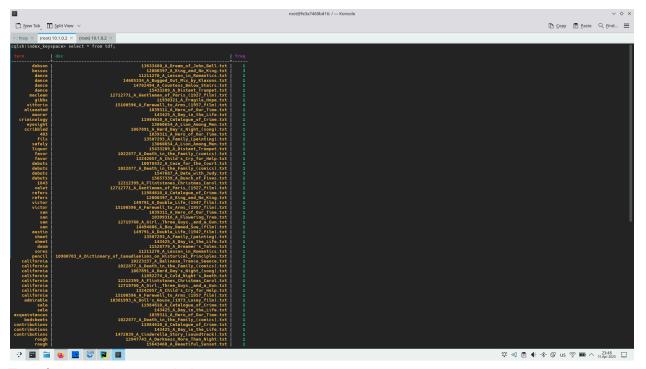Successfully completed the second step.
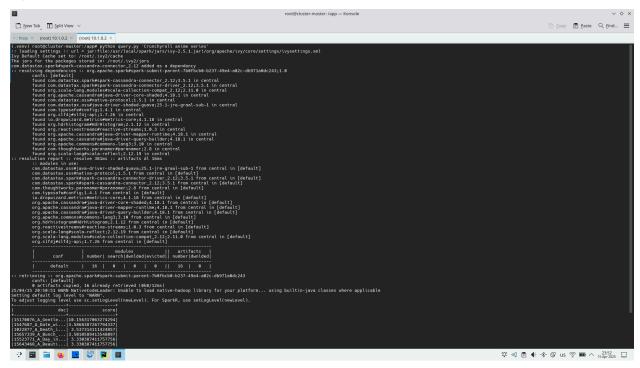
# Cassandra

Document lengths.
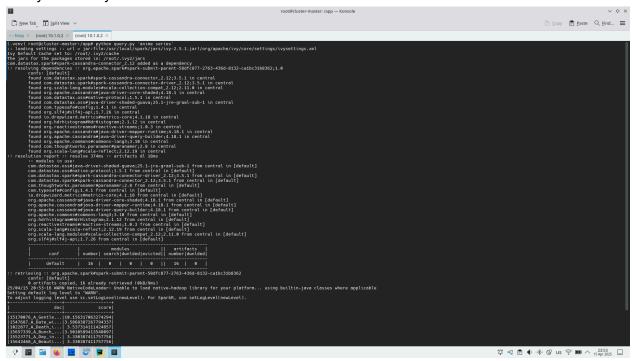


Term frequencies.

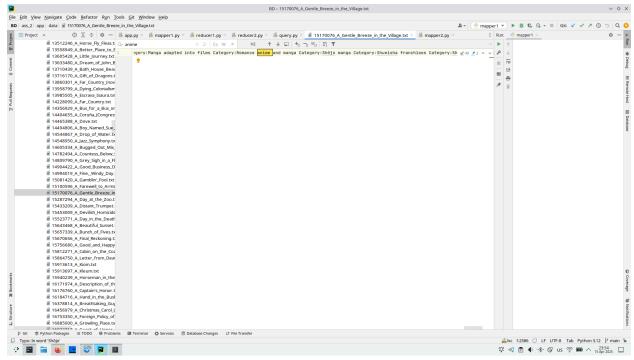

Term frequencies per each document.

# Queries



Query: "Crunchyroll anime series"



Query: "anime series"

The first document "15170076_A_Gentle..." has anime:



As it can be seen from the screenshots, similar queries lead to similar results and documents really contain keywords, and different query leads to different results.