

Egor Lebedev AAI-01

Report A2

Methodology

important!!!!

the system has many points of failure during initialization (simply because there are a lot of tools), so if you see a problem, please restart the build by clearing the cache in docker, it always helps. The search itself is stable

app.sh:

Orchestrates workflow: Starts services → Prepares data → Runs indexing → Ranck for ex → Keeps container alive (tail -f /dev/null).

start-services.sh:

not changed

init-cassandra.sh:

Waits for Cassandra availability → Creates tables via `cassandra/cassandra-init.cql`.

prepare_data.sh:

Downloads a.parquet (if missing) → Transfers data to HDFS → Runs `prepare_data.py`.

prepare_data.py:

Processes Parquet → Generates .txt files → Writes metadata to HDFS `/index/data`.

index.sh:

Runs two MapReduce pipelines:

Pipeline 1: Builds inverted index and fill the tables (`mapper1.py`, `reducer1.py`).

Pipeline 2: Computes document stats (`mapper2.py`, `reducer2.py`).

search.sh:

Executes search via Spark → Passes query to `query.py` with the necessary dependencies → Outputs top 10 results.

query.py:

Implements BM25 with stability fixes → Ranks documents → Returns IDs & titles.

query.py

this script connects to cassandra and downloads the necessary data for the request, after that it calculates bm₂₅, ranks and outputs a list of files as requested in the task.

I used a different version of bm₂₅, it's just protected from division by 0 problems, we went through the formula in the IR course and I checked its performance in that course (and in general, it's just adding constants, which does not affect the ranking)

```
def calculate_bm25(row):
    term = row.term_text
    tf = row.tf
    doc_length = row.length
    df = df_bc.value.get(term, 0)
    idf = math.log((total_docs + 1) / (df + 0.5))
    score = idf * (tf * (k1 + 1)) / (tf + k1 * (1 - b + b * (doc_length / avg_length)))
    return (row.doc_id, score)
```

Demonstration

steps for start

git clone https://github.com/EgorLeb/big_data_a2.git

docker compose up

in other console

```
docker exec -it cluster-master /bin/bash
bash search.sh "book about the dog"
bash search.sh "chance of the corner if i will walk"
```

I looked at the files that came out, they contain the right words, but the meaning doesn't really fit, the problem is, I think, the number of files

screenshots:

```
Apr 15 18:43
Terminal

egorpc% git clone https://github.com/EgorLeb/big_data_a2.git
Cloning into 'big_data_a2'...
remote: Enumerating objects: 1056, done.
remote: Counting objects: 100% (1056/1056), done.
remote: Compressing objects: 100% (1038/1038), done.
remote: Total 1056 (delta 19), reused 1039 (delta 12), pack-reused 0 (from 0)
Receiving objects: 100% (1056/1056), 1.56 MiB | 1.54 MiB/s, done.
Resolving deltas: 100% (19/19), done.
egorpc% cd big_data_a2/
egorpc% docker compose up
[+] Building 0.0s (0/0)
[+] Running 3/3
 ✓ Container cluster-slave-1 Created 0.05s
 ✓ Container cassandra-server Created 0.05s
 ✓ Container cluster-master Recreated 0.45s
Attaching to cassandra-server, cluster-master, cluster-slave-1
cluster-slave-1 | * Starting OpenBSD Secure Shell server sshd [ OK ]
```

```
Apr 15 18:44
Terminal

cluster-master | 105 kB 1.3 MB/s
cluster-master | Collecting wcwidth
cluster-master | Downloading wcwidth-0.2.13-py2.py3-none-any.whl (34 kB)
cluster-master | Collecting pure-sasl
cluster-master | Downloading pure-sasl-0.6.2.tar.gz (11 kB)
cluster-master | Collecting cassandra-driver
cluster-master | Downloading cassandra_driver-3.29.2-cp38-cp38-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (4.1 MB)
cluster-master | 4.1 MB 5.1 MB/s
cluster-master | Collecting geomet<0.3,>=0.1
cluster-master | Downloading geomet-0.2.1.post1-py3-none-any.whl (18 kB)
cluster-master | Requirement already satisfied: six in /usr/lib/python3/dist-packages (from geomet<0.3,>=0.1->cassandra-driver->cqlsh) (1.14.0)
cluster-master | Collecting cqlsh
cluster-master | Downloading cqlsh-0.1.8-py3-none-any.whl (98 kB)
cluster-master | 98 kB 7.2 MB/s
cluster-master | Building wheels for collected packages: pure-sasl
cluster-master | Building wheel for pure-sasl (setup.py) ... done
cluster-master | Created wheel for pure-sasl: filename=pure_sasl-0.6.2-py3-none-any.whl size=11428 sha256=bada8294c83eeffb539566a8fe28f8354362077a9861af4ad9245aff2241e988
cluster-master | Stored in directory: /root/.cache/pip/wheels/af/5e/ca/57ff2c5801d038e3d8b227a4fb492cd84e43a535d64a86f3f2
cluster-master | Successfully built pure-sasl
cluster-master | Installing collected packages: wcwidth, pure-sasl, click, geomet, cassandra-driver, cqlsh
cluster-master | Successfully installed cassandra-driver-3.29.2 click-8.1.8 cqlsh-6.2.0 geomet-0.2.1.post1 pure-sasl-0.6.2 wcwidth-0.2.13
cluster-master | Waiting for Cassandra to start...
cluster-master | WARNING: cqlsh was built against 5.0.0, but this server is 5.0.4. All features may not work!
cluster-master | search_engine system_auth system_schema system_views
cluster-master | system system_distributed system_traces system_virtual_schema
cluster-master | Cassandra is ready. Initializing schema...
cluster-master | WARNING: cqlsh was built against 5.0.0, but this server is 5.0.4. All features may not work!
cluster-master | Starting namenodes on [cluster-master]
cluster-master | Starting datanodes
cluster-master | cluster-slave-4: ssh: Could not resolve hostname cluster-slave-4: Temporary failure in name resolution
cluster-master | cluster-slave-2: ssh: Could not resolve hostname cluster-slave-2: Temporary failure in name resolution
cluster-master | cluster-slave-5: ssh: Could not resolve hostname cluster-slave-5: Temporary failure in name resolution
cluster-master | cluster-slave-3: ssh: Could not resolve hostname cluster-slave-3: Temporary failure in name resolution
cluster-master | Starting secondary namenodes [cluster-master]
cluster-master | Starting resourcemanager
cluster-master | Starting nodemanagers
cluster-master | cluster-slave-4: ssh: Could not resolve hostname cluster-slave-4: Temporary failure in name resolution
cluster-master | cluster-slave-5: ssh: Could not resolve hostname cluster-slave-5: Temporary failure in name resolution
cluster-master | cluster-slave-2: ssh: Could not resolve hostname cluster-slave-2: Temporary failure in name resolution
cluster-master | cluster-slave-3: ssh: Could not resolve hostname cluster-slave-3: Temporary failure in name resolution
```

```
cluster-master | Present Capacity: 318056075264 (296.21 GB)
cluster-master | DFS Remaining: 316819816448 (295.06 GB)
cluster-master | DFS Used: 1236258816 (1.15 GB)
cluster-master | DFS Used%: 0.39%
cluster-master | Replicated Blocks:
cluster-master |   Under replicated blocks: 0
cluster-master |   Blocks with corrupt replicas: 0
cluster-master |   Missing blocks: 0
cluster-master |   Missing blocks (with replication factor 1): 0
cluster-master |   Low redundancy blocks with highest priority to recover: 0
cluster-master |   Pending deletion blocks: 0
cluster-master | Erasure Coded Block Groups:
cluster-master |   Low redundancy block groups: 0
cluster-master |   Block groups with corrupt internal blocks: 0
cluster-master |   Missing block groups: 0
cluster-master |   Low redundancy blocks with highest priority to recover: 0
cluster-master |   Pending deletion blocks: 0
cluster-master | -----
cluster-master | Live datanodes (1):
cluster-master | Name: 172.18.0.3:9866 (cluster-slave-1.big_data_o2_spark-cluster)
cluster-master | Hostname: cluster-slave-1
cluster-master | Decommission Status : Normal
cluster-master | Configured Capacity: 501809635328 (467.35 GB)
cluster-master | DFS Used: 1236258816 (1.15 GB)
cluster-master | Non DFS Used: 158187757568 (147.32 GB)
cluster-master | DFS Remaining: 316819816448 (295.06 GB)
cluster-master | DFS Used%: 0.25%
cluster-master | DFS Remaining%: 63.14%
cluster-master | Configured Cache Capacity: 0 (0 B)
cluster-master | Cache Used: 0 (0 B)
cluster-master | Cache Remaining: 0 (0 B)
cluster-master | Cache Used%: 100.00%
cluster-master | Cache Remaining%: 0.00%
cluster-master | Xcelfers: 0
cluster-master | Last contact: Tue Apr 15 15:44:40 GMT 2025
cluster-master | Last Block Report: Tue Apr 15 15:44:28 GMT 2025
cluster-master | Num of Blocks: 0
cluster-master | Safe mode is OFF
```

here need to wait

```
cluster-master | Collecting packaging
cluster-master | Downloading packaging-24.2-py3-none-any.whl (65 kB)
cluster-master | 65 kB 3.7 MB/s
cluster-master | Collecting py4j==0.10.9.7
cluster-master | Downloading py4j-0.10.9.7-py2.py3-none-any.whl (280 kB)
cluster-master | 280 kB 60.5 MB/s
cluster-master | Requirement already satisfied: click in /usr/local/lib/python3.8/dist-packages (from geomet==0.3,>=0.1->cassandra-driver->-r requirements.txt (line 1)) (8.1.8)
cluster-master | Requirement already satisfied: six in /usr/lib/python3/dist-packages (from geomet==0.3,>=0.1->cassandra-driver->-r requirements.txt (line 1)) (1.14.0)
cluster-master | Collecting pytz==2020.1
cluster-master | Downloading pytz-2020.1-py2.py3-none-any.whl (509 kB)
cluster-master | 509 kB 42.0 MB/s
cluster-master | Collecting tzdata==2025.1
cluster-master | Downloading tzdata-2025.2-py2.py3-none-any.whl (347 kB)
cluster-master | 347 kB 73.6 MB/s
cluster-master | Collecting python-dateutil==2.8.2
cluster-master | Downloading python_dateutil-2.9.0.post0-py2.py3-none-any.whl (229 kB)
cluster-master | 229 kB 53.2 MB/s
cluster-master | Building wheels for collected packages: pyspark
cluster-master | Building wheel for pyspark (setup.py) ... done
cluster-master | Created wheel for pyspark: filename=pyspark-3.5.5-py2.py3-none-any.whl size=317747881 sha256=856f3fdad2b52aba0260d83a86a1dd1abb2753d9ea77b6ceeeec57a3fa7dcd1
cluster-master | Stored in directory: /root/.cache/pip/wheels/9e/5b/b4/a3ac8d456edf8c52eb15f9eb357d961812d5f17bf203c54c18
cluster-master | Successfully built pyspark
cluster-master | Installing collected packages: venv-pack, numpy, pyarrow, pytz, tzdata, python-dateutil, pandas, cranjam, fsspec, packaging, fastparquet, pathvalidate, py4j, pyspark, tqdm
cluster-master | Successfully installed cranjam-2.10.0 fastparquet-2024.2.0 fsspec-2025.3.0 numpy-1.24.4 packaging-24.2 pandas-2.0.3 pathvalidate-3.2.1 py4j-0.10.9.7 pyarrow-17.0.0 pyspark-3.5.5 python-dateutil-2.9.0.post0 pytz-2025.2 tqdm-4.67.1 tzdata-2025.2 venv-pack-0.2.0
cluster-master | Collecting packages...
cluster-master | Packing environment at '/app/.venv' to '.venv.tar.gz'
cluster-master | [#####] | 100% Completed | 0.2s
cluster-master | Download file a.parquet...
cluster-master | --2025-04-15 15:49:36-- https://storage.googleapis.com/kaggle-data-sets/3521629/6146260/compressed/a.parquet.zip?X-Goog-Algorithm=GOOG4-RSA-SHA256&X-Goog-Credential=gcp-kaggle-com%40kaggle-16160714m.gserviceaccount.com%2F20250415%2FAutoK%2Fstorage%2Fgoog4_request&X-Goog-Date=20250415T131438Z&X-Goog-Expires=259200&X-Goog-SignedHeader=host&X-Goog-Signature=72d25f37f9779f6cad85b87ae2cceb0b1474fe6ddcfb085e08df4631ce978ce941477276d6dd4cccd24478262542f624be66719da0d4892980837fb5c25e5dab6912512b5175f01abf566ae8d54178b257b7c81686fc7cdf0773d69d734a6a2cb0a02b2ca92a69d4f869d2bb36336ae5c5391ec66ecbdf4d408be9bffa93c4f38788977ed23990955fdd265dc297ea5a7f559927e1dfa018da46f170ce1be9194432ea7c244cbdf5047548f95589a96ea93ff9b5babdfdc2cbe087d183f3d77af52eae65c0eb455ce3c72ace180236f3d894cb08362fb30a49533e3a495b9c350966ebab9b9d6d6ddc937e04d8f5c707bf329c7d79676405c56a
cluster-master | Resolving storage.googleapis.com (storage.googleapis.com)... 64.233.162.207, 64.233.163.207, 173.194.221.207, ...
cluster-master | Connecting to storage.googleapis.com (storage.googleapis.com)[64.233.162.207]:443... connected.
cluster-master | HTTP request sent, awaiting response... 200 OK
cluster-master | Length: 76835361 (733M) [application/zip]
cluster-master | Saving to: '/app/a.parquet.zip'
```

a.parquet installing - wait

```
Apr 15 2002

cluster-master | 25/04/15 17:01:56 INFO MemoryStore: Block broadcast_6_piece0 stored as bytes in memory (estimated size 77.8 KiB, free 2012.1 MiB)
cluster-master | 25/04/15 17:01:56 INFO BlockManagerInfo: Added broadcast_6_piece0 in memory on cluster-master:40809 (size: 77.8 KiB, free: 2012.7 MiB)
cluster-master | 25/04/15 17:01:56 INFO SparkContext: Created broadcast 6 from broadcast at DAGScheduler.scala:1585
cluster-master | 25/04/15 17:01:56 INFO DAGScheduler: Submitting 1 missing tasks from ResultStage 6 (MapPartitionsRDD[25] at save at NativeMethodAccessorImpl.java:0) (first 15 tasks are f
or partitions Vector(0))
cluster-master | 25/04/15 17:01:56 INFO TaskSchedulerImpl: Adding task set 6.0 with 1 tasks resource profile 0
cluster-master | 25/04/15 17:01:56 INFO TaskSetManager: Starting task 0.0 in stage 6.0 (TID 23) (cluster-master, executor driver, partition 0, NODE_LOCAL, 9275 bytes)
cluster-master | 25/04/15 17:01:56 INFO Executor: Running task 0.0 in stage 6.0 (TID 23)
cluster-master | 25/04/15 17:01:56 INFO FileOutputCommitter: File Output Committer Algorithm version is 1
cluster-master | 25/04/15 17:01:56 INFO FileOutputCommitter: FileOutputCommitter skip cleanup temporary folders under output directory:false, ignore cleanup failures: False
cluster-master | 25/04/15 17:01:56 INFO SQLHadoopMapReduceCommitProtocol: Using output committer class org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
cluster-master | 25/04/15 17:01:56 INFO ShuffleBlockFetcherIterator: Getting 1 (2.3 MiB) non-empty blocks including 1 (2.3 MiB) local and 0 (0.0 B) host-local and 0 (0.0 B) push-merged-lo
cal and 0 (0.0 B) remote blocks
cluster-master | 25/04/15 17:01:56 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 0 ms
cluster-master | 25/04/15 17:01:56 INFO FileOutputCommitter: Saved output of task 'attempt_202504151701566479016802585261868_0006_m_000000_23' to hdfs://cluster-master:9000/index/data/_te
mporary/0/task_202504151701566479016802585261868_0006_m_000000
cluster-master | 25/04/15 17:01:56 INFO SparkHadoopMapReduceUtil: attempt_202504151701566479016802585261868_0006_m_000000_23: Committed. Elapsed time: 17 ms.
cluster-master | 25/04/15 17:01:56 INFO TaskSetManager: Finished task 0.0 in stage 6.0 (TID 23) in 376 ms on cluster-master (executor driver) (1/1)
cluster-master | 25/04/15 17:01:56 INFO TaskSchedulerImpl: Removed TaskSet 6.0, whose tasks have all completed, from pool
cluster-master | 25/04/15 17:01:56 INFO DAGScheduler: ResultStage 6 (save at NativeMethodAccessorImpl.java:0) finished in 0.427 s
cluster-master | 25/04/15 17:01:56 INFO DAGScheduler: Job 4 is finished. Cancelling potential speculative or zombie tasks for this job
cluster-master | 25/04/15 17:01:56 INFO TaskSchedulerImpl: Killing all running tasks in stage 6: Stage finished
cluster-master | 25/04/15 17:01:56 INFO DAGScheduler: Job 4 finished: save at NativeMethodAccessorImpl.java:0, took 0.437131 s
cluster-master | 25/04/15 17:01:56 INFO FileFormatWriter: Start to commit write Job 4b84c2c6-08a1-4b58-ab45-b8f71683e336
cluster-master | 25/04/15 17:01:56 INFO FileFormatWriter: Write Job 4b84c2c6-08a1-4b58-ab45-b8f71683e336 committed. Elapsed time: 21 ms.
cluster-master | 25/04/15 17:01:56 INFO FileFormatWriter: Finished processing stats for write job 4b84c2c6-08a1-4b58-ab45-b8f71683e336.
cluster-master | 25/04/15 17:01:56 INFO SparkContext: Invoking stop() from shutdown hook
cluster-master | 25/04/15 17:01:56 INFO SparkContext: SparkContext is stopping with exitCode 0.
cluster-master | 25/04/15 17:01:56 INFO SparkUI: Stopped Spark web UI at http://cluster-master:4040
cluster-master | 25/04/15 17:01:56 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
cluster-master | 25/04/15 17:01:56 INFO MemoryStore: MemoryStore cleared
cluster-master | 25/04/15 17:01:56 INFO BlockManager: BlockManager stopped
cluster-master | 25/04/15 17:01:56 INFO BlockManagerMaster: BlockManagerMaster stopped
cluster-master | 25/04/15 17:01:56 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
cluster-master | 25/04/15 17:01:56 INFO SparkContext: Successfully stopped SparkContext
cluster-master | 25/04/15 17:01:56 INFO ShutdownHookManager: Shutdown hook called
cluster-master | 25/04/15 17:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-7b73f9c8-f09a-4b4f-954f-304deae21834
cluster-master | 25/04/15 17:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-0af59ef1-d382-4df2-aba2-e3d1e78afdc6/pyspark-8e4155fa-1d7d-4f69-892e-07ba14dd78ab
cluster-master | 25/04/15 17:01:56 INFO ShutdownHookManager: Deleting directory /tmp/spark-0af59ef1-d382-4df2-aba2-e3d1e78afdc6
cluster-master | [3/3] Finalizing...
cluster-master | Putting data to hdfs
```

```
Apr 15 2005

cluster-master | -rw-r--r-- 1 root supergroup 522 2025-04-15 17:04 /data/data/8237581_A_Bridge_Too_Far.txt
cluster-master | -rw-r--r-- 1 root supergroup 1668 2025-04-15 17:02 /data/data/8237659_A_Bridge_Too_Far_(book).txt
cluster-master | -rw-r--r-- 1 root supergroup 4925 2025-04-15 17:04 /data/data/8270511_A_Chairy_Tale.txt
cluster-master | -rw-r--r-- 1 root supergroup 3243 2025-04-15 17:02 /data/data/8312349_A_Barnstormer_in_Oz.txt
cluster-master | -rw-r--r-- 1 root supergroup 4898 2025-04-15 17:03 /data/data/8373513_A_Carol_Christmas.txt
cluster-master | -rw-r--r-- 1 root supergroup 1253 2025-04-15 17:05 /data/data/8457823_A_Change_in_the_Weather.txt
cluster-master | -rw-r--r-- 1 root supergroup 1272 2025-04-15 17:02 /data/data/8528885_A_Cheerful_Gang_Turns_the_Earth.txt
cluster-master | -rw-r--r-- 1 root supergroup 11090 2025-04-15 17:02 /data/data/8530771_A_Band_in_the_River.txt
cluster-master | -rw-r--r-- 1 root supergroup 1237 2025-04-15 17:04 /data/data/8563962_A_Capitol_Federal.txt
cluster-master | -rw-r--r-- 1 root supergroup 6682 2025-04-15 17:03 /data/data/867420_A_Burnt-Out_Case.txt
cluster-master | -rw-r--r-- 1 root supergroup 5869 2025-04-15 17:02 /data/data/8694441_A_Best_2.txt
cluster-master | -rw-r--r-- 1 root supergroup 417 2025-04-15 17:03 /data/data/8740104_A_Certain_World.txt
cluster-master | -rw-r--r-- 1 root supergroup 2639 2025-04-15 17:02 /data/data/8788071_A_Bug_and_a_Bag_of_Weed.txt
cluster-master | -rw-r--r-- 1 root supergroup 185 2025-04-15 17:03 /data/data/8783301_A_Bugged_Out_Mix_(Felix_da_Housecat_album).txt
cluster-master | -rw-r--r-- 1 root supergroup 694 2025-04-15 17:03 /data/data/8796602_A_Chiefains_Celebration.txt
cluster-master | -rw-r--r-- 1 root supergroup 1066 2025-04-15 17:05 /data/data/8922606_A_Carolina_Jubilee.txt
cluster-master | -rw-r--r-- 1 root supergroup 22618 2025-04-15 17:05 /data/data/89397_A_Bridge_Too_Far_(film).txt
cluster-master | -rw-r--r-- 1 root supergroup 12884 2025-04-15 17:02 /data/data/89541_A_Better_Tomorrow.txt
cluster-master | -rw-r--r-- 1 root supergroup 2523 2025-04-15 17:03 /data/data/8981838_A_Calendar_of_Wisdom.txt
cluster-master | -rw-r--r-- 1 root supergroup 1726 2025-04-15 17:02 /data/data/9277570_A_Break_in_the_Weather.txt
cluster-master | -rw-r--r-- 1 root supergroup 519 2025-04-15 17:02 /data/data/9279736_A_(Cass_McCombs_album).txt
cluster-master | -rw-r--r-- 1 root supergroup 1938 2025-04-15 17:03 /data/data/929153_A_Bao_Ou_(album).txt
cluster-master | -rw-r--r-- 1 root supergroup 3221 2025-04-15 17:03 /data/data/929265_A_Chance_to_Cut_Is_a_Chance_to_Cure.txt
cluster-master | -rw-r--r-- 1 root supergroup 2380 2025-04-15 17:04 /data/data/9341936_A_Cage_of_Butterflies.txt
cluster-master | -rw-r--r-- 1 root supergroup 1991 2025-04-15 17:04 /data/data/9398117_A_Biography.txt
cluster-master | -rw-r--r-- 1 root supergroup 6716 2025-04-15 17:03 /data/data/9424803_A_Breed_of_Heroes.txt
cluster-master | -rw-r--r-- 1 root supergroup 3466 2025-04-15 17:04 /data/data/9707515_A_Chaos_of_Flowers.txt
cluster-master | -rw-r--r-- 1 root supergroup 1650 2025-04-15 17:02 /data/data/9848866_A_Big_10-8_Place.txt
cluster-master | -rw-r--r-- 1 root supergroup 2368 2025-04-15 17:02 /data/data/9938275_A_Band_Called_David.txt
cluster-master | -rw-r--r-- 1 root supergroup 1770 2025-04-15 17:02 /data/data/993992_A_Beginners'_Guide_to_the_King_Crinson_Collectors'_Club.txt
cluster-master | -rw-r--r-- 1 root supergroup 896 2025-04-15 17:03 /data/data/9965276_A_Book_of_Human_Language.txt
cluster-master | -rw-r--r-- 1 root supergroup 0 2025-04-15 17:01 /index/data/_SUCCESS
cluster-master | -rw-r--r-- 1 root supergroup 3537921 2025-04-15 17:01 /index/data/part-00000-3762c4fe-2bc6-4662-9640-a11085024504-c000.csv
cluster-master | Done data preparation!
cluster-master | This script include commands to run mapreduce jobs using hadoop streaming to index documents
cluster-master | Found 6 items
cluster-master | -rw-r--r-- 1 root supergroup 873207391 2025-04-15 17:01 /a-parquet
cluster-master | drwxr-xr-x 1 root supergroup 0 2025-04-15 16:57 /apps
cluster-master | drwxr-xr-x 1 root supergroup 0 2025-04-15 17:01 /data
cluster-master | drwxr-xr-x 1 root supergroup 0 2025-04-15 17:01 /index
cluster-master | drwxrwxr-x 1 root supergroup 0 2025-04-15 16:57 /tmp
cluster-master | drwxr-xr-x 1 root supergroup 0 2025-04-15 16:58 /user
```

```
Apr 15 2020

cluster-master | -rw-r--r-- 1 root supergroup 1650 2025-04-15 17:02 /data/data/9848866_A_Big_10-8_Place.txt
cluster-master | -rw-r--r-- 1 root supergroup 2368 2025-04-15 17:02 /data/data/9938275_A_Band_called_David.txt
cluster-master | -rw-r--r-- 1 root supergroup 1778 2025-04-15 17:02 /data/data/995992_A_Beginners'_Guide_to_the_King_Crimson_Collectors'_Club.txt
cluster-master | -rw-r--r-- 1 root supergroup 896 2025-04-15 17:03 /data/data/9965276_A_Book_of_Human_Language.txt
cluster-master | -rw-r--r-- 1 root supergroup 0 2025-04-15 17:01 /index/data/_SUCCESS
cluster-master | -rw-r--r-- 1 root supergroup 3537921 2025-04-15 17:01 /index/data/part-00000-3762c4fe-2bcb-4662-9640-a11085024504-c000.csv
cluster-master | Done data preparation!
cluster-master | This script include commands to run mapreduce jobs using hadoop streaming to index documents
cluster-master | Found 6 items
cluster-master | -rw-r--r-- 1 root supergroup 873207391 2025-04-15 17:01 /a.parquet
cluster-master | drwxr-xr-x - root supergroup 0 2025-04-15 16:57 /apps
cluster-master | drwxr-xr-x - root supergroup 0 2025-04-15 17:01 /data
cluster-master | drwxr-xr-x - root supergroup 0 2025-04-15 17:01 /index
cluster-master | drwxrwx--- - root supergroup 0 2025-04-15 16:57 /tmp
cluster-master | drwxr-xr-x - root supergroup 0 2025-04-15 16:58 /user
cluster-master | packageJobJar: [/tmp/hadoop-unjar7324219116855901263/] [] /tmp/streamjob6040915109492439187.jar tmpDir=null
cluster-master | 2025-04-15 17:05:42,667 INFO client.DefaultHadoopFailoverProxyProvider: Connecting to ResourceManager at cluster-master/172.24.0.4:8032
cluster-master | 2025-04-15 17:05:42,803 INFO client.DefaultHadoopFailoverProxyProvider: Connecting to ResourceManager at cluster-master/172.24.0.4:8032
cluster-master | 2025-04-15 17:05:43,123 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/root/.staging/job_1744736240976_0001
cluster-master | 2025-04-15 17:05:44,560 INFO mapred.FileInputFormat: Total input files to process : 1
cluster-master | 2025-04-15 17:05:44,595 INFO mapreduce.JobSubmitter: number of splits:2
cluster-master | 2025-04-15 17:05:44,718 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1744736240976_0001
cluster-master | 2025-04-15 17:05:44,719 INFO mapreduce.JobSubmitter: Executing with tokens: []
cluster-master | 2025-04-15 17:05:44,875 INFO conf.Configuration: resource-types.xml not found
cluster-master | 2025-04-15 17:05:45,005 INFO resource.ResourceTypes: Unable to find 'resource-types.xml'.
cluster-master | 2025-04-15 17:05:45,324 INFO ImplVarClientImpl: Submitted application application_1744736240976_0001
cluster-master | 2025-04-15 17:05:45,384 INFO mapreduce.Job: The url to track the job: http://cluster-master:8088/proxy/application_1744736240976_0001/
cluster-master | 2025-04-15 17:05:45,386 INFO mapreduce.Job: Running job: job_1744736240976_0001
cluster-master | 2025-04-15 17:05:57,549 INFO mapreduce.Job: Job job_1744736240976_0001 running in uber mode : false
cluster-master | 2025-04-15 17:05:57,551 INFO mapreduce.Job: map 0% reduce 0%
cluster-master | 2025-04-15 17:06:03,660 INFO mapreduce.Job: map 100% reduce 0%
cluster-master | 2025-04-15 17:06:19,780 INFO mapreduce.Job: map 100% reduce 68%
cluster-master | 2025-04-15 17:06:37,915 INFO mapreduce.Job: map 100% reduce 70%
cluster-master | 2025-04-15 17:06:56,048 INFO mapreduce.Job: map 100% reduce 72%
cluster-master | 2025-04-15 17:07:14,159 INFO mapreduce.Job: map 100% reduce 73%
cluster-master | 2025-04-15 17:07:20,195 INFO mapreduce.Job: map 100% reduce 74%
cluster-master | 2025-04-15 17:07:38,300 INFO mapreduce.Job: map 100% reduce 76%
cluster-master | 2025-04-15 17:07:56,394 INFO mapreduce.Job: map 100% reduce 78%
cluster-master | 2025-04-15 17:08:02,428 INFO mapreduce.Job: map 100% reduce 79%
cluster-master | 2025-04-15 17:08:20,515 INFO mapreduce.Job: map 100% reduce 80%
cluster-master | 2025-04-15 17:08:26,542 INFO mapreduce.Job: map 100% reduce 81%
cluster-master | 2025-04-15 17:08:32,566 INFO mapreduce.Job: map 100% reduce 82%
```

```
Apr 15 2013

cluster-master | 25/04/15 17:13:28 INFO TaskSetManager: Starting task 4.0 in stage 12.0 (TID 25) (cluster-slave-1, executor 1, partition 4, RACK_LOCAL, 11160 bytes)
cluster-master | 25/04/15 17:13:28 INFO TaskSetManager: Finished task 2.0 in stage 12.0 (TID 23) in 98 ms on cluster-slave-1 (executor 1) (3/6)
cluster-master | 25/04/15 17:13:28 INFO TaskSetManager: Starting task 5.0 in stage 12.0 (TID 26) (cluster-slave-1, executor 2, partition 5, RACK_LOCAL, 11036 bytes)
cluster-master | 25/04/15 17:13:28 INFO TaskSetManager: Finished task 3.0 in stage 12.0 (TID 24) in 98 ms on cluster-slave-1 (executor 2) (4/6)
cluster-master | 25/04/15 17:13:28 INFO TaskSetManager: Finished task 5.0 in stage 12.0 (TID 26) in 72 ms on cluster-slave-1 (executor 2) (5/6)
cluster-master | 25/04/15 17:13:28 INFO TaskSetManager: Finished task 4.0 in stage 12.0 (TID 25) in 123 ms on cluster-slave-1 (executor 1) (6/6)
cluster-master | 25/04/15 17:13:28 INFO YarnScheduler: Removed TaskSet 12.0, whose tasks have all completed, from pool.
cluster-master | 25/04/15 17:13:28 INFO DAGScheduler: ResultStage 12 (collectAsMap at /app/query.py:69) finished in 0.379 s
cluster-master | 25/04/15 17:13:28 INFO DAGScheduler: Job 8 is finished. Cancelling potential speculative or zombie tasks for this job
cluster-master | 25/04/15 17:13:28 INFO YarnScheduler: Killing all running tasks in stage 12: Stage finished
cluster-master | 25/04/15 17:13:28 INFO DAGScheduler: Job 8 finished: collectAsMap at /app/query.py:69, took 0.385522 s
cluster-master |
cluster-master | Top 10 Results:
cluster-master | 1. 2828410 A (musical note) [Score: 33.31]
cluster-master | 2. 58726751 A (Los Angeles Railway) [Score: 9.90]
cluster-master | 3. 15547032 A G G Price [Score: 8.92]
cluster-master | 4. 31831074 A Cambridge Mass [Score: 7.49]
cluster-master | 5. 14810863 A Carnival Christmas [Score: 7.15]
cluster-master | 6. 4697001 A Car-Tune Portrait [Score: 6.47]
cluster-master | 7. 46860526 A Beautiful Soul (song) [Score: 6.21]
cluster-master | 8. 72523905 A Bu [Score: 6.04]
cluster-master | 9. 44464104 A Battle of Nerves [Score: 5.89]
cluster-master | 10. 3561416 A Baha [Score: 5.29]
cluster-master | 25/04/15 17:13:28 INFO SparkContext: SparkContext is stopping with exitCode 0.
cluster-master | 25/04/15 17:13:28 INFO SparkUI: Stopped Spark web UI at http://cluster-master:4040
cluster-master | 25/04/15 17:13:28 INFO YarnClientSchedulerBackend: Interrupting monitor thread
cluster-master | 25/04/15 17:13:28 INFO YarnClientSchedulerBackend: Shutting down all executors
cluster-master | 25/04/15 17:13:28 INFO YarnSchedulerBackend$YarnDriverEndpoint: Asking each executor to shut down
cluster-master | 25/04/15 17:13:28 INFO YarnClientSchedulerBackend: YARN client scheduler backend Stopped
cluster-master | 25/04/15 17:13:28 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
cluster-master | 25/04/15 17:13:28 INFO MemoryStore: MemoryStore cleared
cluster-master | 25/04/15 17:13:28 INFO BlockManager: BlockManager stopped
cluster-master | 25/04/15 17:13:28 INFO BlockManagerMaster: BlockManagerMaster stopped
cluster-master | 25/04/15 17:13:28 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
cluster-master | 25/04/15 17:13:28 INFO SparkContext: Successfully stopped SparkContext
cluster-master | 25/04/15 17:13:29 INFO ShutdownHookManager: Shutdown hook called
cluster-master | 25/04/15 17:13:29 INFO ShutdownHookManager: Deleting directory /tmp/spark-b92dbec-84a0-4d07-8d2c-2eb2fa63204e
cluster-master | 25/04/15 17:13:29 INFO ShutdownHookManager: Deleting directory /tmp/spark-7677ed75-82dc-4a17-bee6-41e5c612d0b6
cluster-master | 25/04/15 17:13:29 INFO ShutdownHookManager: Deleting directory /tmp/spark-7677ed75-82dc-4a17-bee6-41e5c612d0b6/pyspark-48f12f5e-f144-42ab-9959-57e0d82c564e
cluster-master | 25/04/15 17:13:29 INFO CassandraConnector: Disconnected from Cassandra cluster.
cluster-master | 25/04/15 17:13:29 INFO SerialShutdownHooks: Successfully executed shutdown hook: Clearing session cache for C* connector
cluster-master | Services started. Holding the container open...
```

docker exec -it cluster-master /bin/bash
bash search.sh "book about the dog"

```
Apr 15 20:19
root@cluster-master: /app

Terminal
root@cluster-master: /app

egorpc% docker exec -it cluster-master /bin/bash
root@cluster-master:/app# bash search.sh "book about the dog"
This script will include commands to search for documents given the query using Spark RDD
:: loading settings :: url = jar:fille:/usr/local/spark/jars/ivy-2.5.1.jar!/org/apache/ivy/core/settings/ivysettings.xml
ivy Default Cache set to: /root/.ivy2/cache
The jars for the packages stored in: /root/.ivy2/jars
com.datastax.spark#spark-cassandra-connector_2.12 added as a dependency
com.github.jnr#jnr-posix added as a dependency
:: resolving dependencies :: org.apache.spark#spark-submit-parent-91d0fcd1-dd0b-4199-9dc5-f30a97da0ed9;1.0
  confs: [default]
    found com.datastax.spark#spark-cassandra-connector_2.12;3.2.0 in central
    found com.datastax.spark#spark-cassandra-connector-driver_2.12;3.2.0 in central
    found com.datastax.oss#java-driver-core-shaded;4.13.0 in central
    found com.datastax.oss#native-protocol;1.5.0 in central
    found com.datastax.oss#java-driver-shaded-guava;25.1-jre-graal-sub-1 in central
    found com.typesafe#config;1.4.1 in central
    found org.slf4j#slf4j-api;1.7.26 in central
    found io.dropwizard.metrics#metrics-core;4.1.18 in central
    found org.hdrhistogram#HdrHistogram;2.1.12 in central
    found org.reactivestreams#reactive-streams;1.0.3 in central
    found com.github.stephenc.jcip#jcip-annotations;1.0.1 in central
    found com.github.spotbugs#spotbugs-annotations;3.1.12 in central
    found com.google.code.findbugs#jsr305;3.0.2 in central
    found com.datastax.oss#java-driver-mapper-runtime;4.13.0 in central
    found com.datastax.oss#java-driver-query-builder;4.13.0 in central
    found org.apache.commons#commons-lang3;3.10 in central
    found com.thoughtworks.paranamer#paranamer;2.6 in central
    found org.scala-lang#scala-reflect;2.12.11 in central
    found com.github.jnr#jnr-posix;3.1.15 in central
    found com.github.jnr#jnr-ffi;2.2.11 in central
    found com.github.jnr#jffi;1.3.9 in central
    found org.ow2.asm#asm;9.2 in central
    found org.ow2.asm#asm-commons;9.2 in central
    found org.ow2.asm#asm-tree;9.2 in central
    found org.ow2.asm#asm-analysis;9.2 in central
    found org.ow2.asm#asm-util;9.2 in central
    found com.github.jnr#jnr-a64asm;1.0.0 in central
    found com.github.jnr#jnr-x86asm;1.0.2 in central
    found com.github.jnr#jnr-constants;0.10.3 in central
:: resolution report :: resolve 492ms :: artifacts dl 26ms
  modules in use:
    com.datastax.oss#java-driver-core-shaded;4.13.0 from central in [default]
    com.datastax.oss#java-driver-mapper-runtime;4.13.0 from central in [default]
    com.datastax.oss#java-driver-query-builder;4.13.0 from central in [default]
```

```
Apr 15 20:20
root@cluster-master: /app

Terminal
root@cluster-master: /app

25/04/15 17:20:18 INFO TaskSetManager: Finished task 1.0 in stage 12.0 (TID 22) in 115 ms on cluster-slave-1 (executor 1) (2/6)
25/04/15 17:20:18 INFO TaskSetManager: Starting task 4.0 in stage 12.0 (TID 25) (cluster-slave-1, executor 2, partition 4, RACK_LOCAL, 11036 bytes)
25/04/15 17:20:18 INFO TaskSetManager: Finished task 2.0 in stage 12.0 (TID 23) in 99 ms on cluster-slave-1 (executor 2) (3/6)
25/04/15 17:20:18 INFO TaskSetManager: Starting task 5.0 in stage 12.0 (TID 26) (cluster-slave-1, executor 1, partition 5, RACK_LOCAL, 11160 bytes)
25/04/15 17:20:18 INFO TaskSetManager: Finished task 3.0 in stage 12.0 (TID 24) in 96 ms on cluster-slave-1 (executor 1) (4/6)
25/04/15 17:20:18 INFO TaskSetManager: Finished task 4.0 in stage 12.0 (TID 25) in 98 ms on cluster-slave-1 (executor 2) (5/6)
25/04/15 17:20:18 INFO TaskSetManager: Finished task 5.0 in stage 12.0 (TID 26) in 90 ms on cluster-slave-1 (executor 1) (6/6)
25/04/15 17:20:18 INFO YarnScheduler: Removed TaskSet 12.0, whose tasks have all completed, from pool
25/04/15 17:20:18 INFO DAGScheduler: ResultStage 12 (collectAsMap at /app/query.py:69) finished in 0.308 s
25/04/15 17:20:18 INFO DAGScheduler: Job 8 is finished. Cancelling potential speculative or zombie tasks for this job
25/04/15 17:20:18 INFO YarnScheduler: Killing all running tasks in stage 12: Stage finished
25/04/15 17:20:18 INFO DAGScheduler: Job 8 finished: collectAsMap at /app/query.py:69, took 0.312398 s

Top 10 Results:
1. 34488106 A Ball for Daisy [Score: 11.26]
2. 17488265 A Boy and His Dog (1946 film) [Score: 7.81]
3. 2008010 A Boy and His Dog [Score: 7.71]
4. 41649549 A Boy and His Dog (1975 film) [Score: 7.27]
5. 23013668 A Bone for a Bone [Score: 7.07]
6. 5677646 A Boy's Best Friend [Score: 7.01]
7. 35717559 A Beuk o' Newcassell Sangs [Score: 6.67]
8. 63087399 A Candle in Her Room [Score: 6.55]
9. 61110799 A Boy, a Girl and a Dog [Score: 6.48]
10. 47515595 A Canine Sherlock Holmes [Score: 6.29]
25/04/15 17:20:18 INFO SparkContext: SparkContext is stopping with exitCode 0.
25/04/15 17:20:18 INFO SparkUI: Stopped Spark web UI at http://cluster-master:4040
25/04/15 17:20:18 INFO YarnClientSchedulerBackend: Interrupting monitor thread
25/04/15 17:20:18 INFO YarnClientSchedulerBackend: Shutting down all executors
25/04/15 17:20:18 INFO YarnSchedulerBackend$YarnDriverEndpoint: Asking each executor to shut down
25/04/15 17:20:18 INFO YarnClientSchedulerBackend: YARN client scheduler backend Stopped
25/04/15 17:20:18 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
25/04/15 17:20:18 INFO MemoryStore: MemoryStore cleared
25/04/15 17:20:18 INFO BlockManager: BlockManager stopped
25/04/15 17:20:18 INFO BlockManagerMaster: BlockManagerMaster stopped
25/04/15 17:20:18 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
25/04/15 17:20:18 INFO SparkContext: Successfully stopped SparkContext
25/04/15 17:20:18 INFO ShutdownHookManager: Shutdown hook called
25/04/15 17:20:19 INFO ShutdownHookManager: Deleting directory /tmp/spark-45c77951-8fa0-4dbd-a4d5-3bcc46c5013b
25/04/15 17:20:19 INFO ShutdownHookManager: Deleting directory /tmp/spark-f0f686e5-d731-42dc-8dfb-1bb4c583e775
25/04/15 17:20:19 INFO ShutdownHookManager: Deleting directory /tmp/spark-f0f686e5-d731-42dc-8dfb-1bb4c583e775/pyspark-cff0573b-8222-4917-89ea-d3ee3fb422e3
25/04/15 17:20:19 INFO CassandraConnector: Disconnected from Cassandra cluster.
25/04/15 17:20:19 INFO SerialShutdownHooks: Successfully executed shutdown hook: Clearing session cache for C* connector
root@cluster-master:/app#
```



```
bash search.sh "chance of the corner if i will walk"
```

```
root@cluster-master:/app
25/04/15 17:21:55 INFO ShutdownHookManager: Deleting directory /tmp/spark-f35c2a2e-cbab-4432-9199-9c1df5a56c19/pyspark-a79298df-0d41-4b58-93da-2e4ec8b99bc3
25/04/15 17:21:55 INFO CassandraConnector: Disconnected from Cassandra cluster.
25/04/15 17:21:55 INFO SerialShutdownHooks: Successfully executed shutdown hooks: Clearing session cache for C* connector
root@cluster-master:/app bash search --name-of-the-corner-12 I will walk
This script will include commands to search for documents given the query using Spark RDD
:: Loading settings :: url = jar:file:/usr/local/spark/jars/ivy-2.5.1.jar!/org/apache/ivy/core/settings/ivysettings.xml
Ivy Default Cache set to: /root/.ivy2/cache
The jars for the packages stored in: /root/.ivy2/jars
com.datastax.spark#spark-cassandra-connector_2.12 added as a dependency
com.github.jnr#jnr-posix added as a dependency
:: resolving dependencies :: org.apache.spark#spark-submt-parent-240b2b4f-adc7-4b68-99f8-6e531a6ac1d;1.0
  confs: [default]
    found com.datastax.spark#spark-cassandra-connector_2.12;3.2.0 in central
    found com.datastax.spark#spark-cassandra-connector-driver_2.12;3.2.0 in central
    found com.datastax.oss#java-driver-core-shaded;4.13.0 in central
    found com.datastax.oss#native-protocol;1.5.0 in central
    found com.datastax.oss#java-driver-shaded-guava;25.1-jre-graal-sub-1 in central
    found com.typesafe#config;1.4.1 in central
    found org.slf4j#slf4j-api;1.7.26 in central
    found io.dropwizard.metrics#metrics-core;4.1.18 in central
    found org.hdrhistogram#HdrHistogram;2.1.12 in central
    found org.reactivestreams#reactive-streams;1.0.3 in central
    found com.github.stephenc.jcip#jcip-annotations;1.0.1 in central
    found com.github.spotbugs#spotbugs-annotations;3.1.12 in central
    found com.google.code.findbugs#jsr305;3.0.2 in central
    found com.datastax.oss#java-driver-mapper-runtime;4.13.0 in central
    found com.datastax.oss#java-driver-query-builder;4.13.0 in central
    found org.apache.commons#commons-lang3;3.10 in central
    found com.thoughtworks.paranamer#paranamer;2.8 in central
    found org.scala-lang#scala-reflect;2.12.11 in central
    found com.github.jnr#jnr-posix;3.1.15 in central
    found com.github.jnr#jnr-ffi;2.2.11 in central
    found com.github.jnr#jffi;1.3.9 in central
    found org.ow2.asm#asm;9.2 in central
    found org.ow2.asm#asm-commons;9.2 in central
    found org.ow2.asm#asm-tree;9.2 in central
    found org.ow2.asm#asm-analysis;9.2 in central
    found org.ow2.asm#asm-util;9.2 in central
    found com.github.jnr#jnr-a64asm;1.0.8 in central
    found com.github.jnr#jnr-x86asm;1.0.2 in central
    found com.github.jnr#jnr-constants;0.10.3 in central
:: resolution report :: resolve 514ms :: artifacts dl 30ms
:: modules in use:
```

The image shows a terminal window with a dark background. At the top, there's a title bar with window controls and the text 'root@cluster-master: /app'. The terminal content is divided into two main sections. The first section contains a series of log messages from the Spark ecosystem, including TaskSetManager, DAGScheduler, and YarnScheduler, detailing the execution of a job with 8 tasks. The second section, titled 'Top 10 Results:', lists ten items with their scores, such as 'A Bug's Life (video game) [Score: 9.57]'. The third section continues with more logs, showing the shutdown of the SparkContext and the deletion of directories and hooks. The terminal ends with the prompt 'root@cluster-master: /app#'.

```
Apr 15 2024
root@cluster-master: /app

Terminal
root@cluster-master: /app

25/04/15 17:24:00 INFO TaskSetManager: Finished task 0.0 in stage 12.0 (TID 21) in 138 ms on cluster-slave-1 (executor 1) (2/6)
25/04/15 17:24:00 INFO TaskSetManager: Starting task 4.0 in stage 12.0 (TID 25) (cluster-slave-1, executor 2, partition 4, RACK_LOCAL, 11036 bytes)
25/04/15 17:24:00 INFO TaskSetManager: Finished task 2.0 in stage 12.0 (TID 23) in 87 ms on cluster-slave-1 (executor 2) (3/6)
25/04/15 17:24:00 INFO TaskSetManager: Starting task 5.0 in stage 12.0 (TID 26) (cluster-slave-1, executor 1, partition 5, RACK_LOCAL, 11040 bytes)
25/04/15 17:24:00 INFO TaskSetManager: Finished task 3.0 in stage 12.0 (TID 24) in 69 ms on cluster-slave-1 (executor 1) (4/6)
25/04/15 17:24:00 INFO TaskSetManager: Finished task 5.0 in stage 12.0 (TID 26) in 36 ms on cluster-slave-1 (executor 1) (5/6)
25/04/15 17:24:00 INFO TaskSetManager: Finished task 4.0 in stage 12.0 (TID 25) in 89 ms on cluster-slave-1 (executor 2) (6/6)
25/04/15 17:24:00 INFO YarnScheduler: Removed TaskSet 12.0, whose tasks have all completed, from pool
25/04/15 17:24:00 INFO DAGScheduler: ResultStage 12 (collectAsMap at /app/query.py:69) finished in 0.271 s
25/04/15 17:24:00 INFO DAGScheduler: Job 8 is finished. Cancelling potential speculative or zombie tasks for this job
25/04/15 17:24:00 INFO YarnScheduler: Killing all running tasks in stage 12: Stage finished
25/04/15 17:24:00 INFO DAGScheduler: Job 8 finished: collectAsMap at /app/query.py:69, took 0.275089 s

Top 10 Results:
1. 2354739 A Bug's Life (video game) [Score: 9.57]
2. 63007399 A Candle in Her Room [Score: 9.57]
3. 12187025 A Bride of the Plains [Score: 9.32]
4. 31318265 A Chance to Make History [Score: 8.97]
5. 62095113 A Beautifully Foolish Endeavour [Score: 8.33]
6. 7095820 A Brush with the Law [Score: 7.95]
7. 38966582 A & R Recording [Score: 7.71]
8. 32341396 A Bird came down the Walk [Score: 7.48]
9. 20574008 A Billion Hands Concert [Score: 7.33]
10. 8373513 A Carol Christmas [Score: 7.26]

25/04/15 17:24:00 INFO SparkContext: SparkContext is stopping with exitCode 0.
25/04/15 17:24:00 INFO SparkUI: Stopped Spark web UI at http://cluster-master:4040
25/04/15 17:24:00 INFO YarnClientSchedulerBackend: Interrupting monitor thread
25/04/15 17:24:00 INFO YarnClientSchedulerBackend: Shutting down all executors
25/04/15 17:24:00 INFO YarnSchedulerBackend$YarnDriverEndpoint: Asking each executor to shut down
25/04/15 17:24:00 INFO YarnClientSchedulerBackend: YARN client scheduler backend Stopped
25/04/15 17:24:00 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
25/04/15 17:24:00 INFO MemoryStore: MemoryStore cleared
25/04/15 17:24:00 INFO BlockManager: BlockManager stopped
25/04/15 17:24:00 INFO BlockManagerMaster: BlockManagerMaster stopped
25/04/15 17:24:00 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
25/04/15 17:24:00 INFO SparkContext: Successfully stopped SparkContext
25/04/15 17:24:09 INFO ShutdownHookManager: Shutdown hook called
25/04/15 17:24:09 INFO ShutdownHookManager: Deleting directory /tmp/spark-e4fe49e9-312a-4b5d-ace-30fbf18bb34f
25/04/15 17:24:09 INFO ShutdownHookManager: Deleting directory /tmp/spark-e4f325cf-de13-4c2d-9087-b55e123bc080
25/04/15 17:24:09 INFO ShutdownHookManager: Deleting directory /tmp/spark-e4f49e9-312a-4b5d-ace-30fbf18bb34f/pyspark-66bbf9f4-c48a-490f-b711-b3f352f8b080
25/04/15 17:24:09 INFO CassandraConnector: Disconnected from Cassandra cluster.
25/04/15 17:24:09 INFO SerialShutdownHooks: Successfully executed shutdown hook: Clearing session cache for C* connector

root@cluster-master: /app#
```