# INFSCI 2750 - Mini Project 03

Yichi Zhang (yiz141@pitt.edu)
Quan Zhou (quz3@pitt.edu)

## Part 1: Setting up Cassandra

Setting up Cassandra on a cluster with Debian package installation is relatively easy compared to manually setting up Hadoop or Spark. Here is how we did it.

```
# ssh on to master(159.89.43.89) and slave(159.89.43.152) and run the following commands:
echo "deb http://www.apache.org/dist/cassandra/debian 311x main" \
 | sudo tee -a /etc/apt/sources.list.d/cassandra.sources.list
curl https://www.apache.org/dist/cassandra/KEYS | sudo apt-key add -
sudo apt-get update
sudo apt-get install cassandra
sudo service cassandra stop
```

Modify the configuration file `cassandra.yaml` on each node:

```
vim /etc/cassandra/cassandra.yaml
# on both nodes, set
 - seeds: "master,slave"
 # setting the read timeout to a larger number to make sure UDF wouldn't timeout
 read_request_timeout_in_ms: 600000
# on master node, set
 listen_address: master
 rpc_address: master
# on slave node, set
 listen_address: slave
 rpc_address: slave
```

With all the preparation steps done, we can start the Cassandra cluster by running the following command on each node:

```
cassandra -Rf
```

On a new ssh session, run `nodetool status` to check the status of the Cassandra cluster.

Screenshot of the cluster up and running

To start a Cassandra CQL Shell on the cluster, simply run the following command. Note we specified the request timeout (in seconds) for the shell to match the timeout we set in the configuration file.

```
cqlsh master --request-timeout=600
```



Screenshot of CQL Shell and result of a test CQL

The project's source code was written in JAVA and used Maven as dependency management and build tool. The JAVA code and `pom.xml` file for the project is in the `source_code` folder. The `mini-project-03-1.0.0-jar-with-dependencies.jar` file was built locally with maven to include all the source code provided and was uploaded to the VM server for running. Note here we need to build the jar file with dependencies in order to successfully run it.

```
mvn package
scp target/mini-project-03-1.0.0-jar-with-dependencies.jar root@159.89.43.89:~/
```

Other than the JAVA code and jar file, we also provided all the CQL commands used in this project in `source_code/cql.txt`

## Part 2 : Import Data into Cassandra

The `ImportData.java` file is the source code for importing the `access_log` file into Cassandra. The import process can be launched by:

```
java -cp mini-project-03-1.0.0-jar-with-dependencies.jar log.ImportData
```

The program first run setup CQL for the `project_03` keyspace synchronously:

```
# drop existing data
DROP TABLE IF EXISTS project_03.log;
DROP TABLE IF EXISTS project_03.ip;
DROP TABLE IF EXISTS project_03.path;
DROP KEYSPACE IF EXISTS project_03;

# create keyspace
CREATE KEYSPACE project_03
 WITH replication = {
   'class': 'SimpleStrategy',
   'replication_factor' : 2};

# create log table and index on columns ip and path
CREATE TABLE project_03.log (
 id int,
 ip text,
 identity text,
 username text,
 time text,
 method text,
 path text,
 protocol text,
 status int,
 size int,
 PRIMARY KEY ((id), ip, path));
CREATE INDEX ip_index ON project_03.log (ip);
CREATE INDEX path_index ON project_03.log (path);

# create counter table for ip and path
CREATE TABLE project_03.ip (ip text PRIMARY KEY, count counter);
CREATE TABLE project_03.path (path text PRIMARY KEY, count counter);

# create UDF for retrieving max group count on ip and path columns directly from log table
CREATE FUNCTION project_03.state_group_and_count(state map<text, int>, type text)
 CALLED ON NULL INPUT
 RETURNS map<text, int>
 LANGUAGE java
 AS $$
```

```
    Integer count = (Integer) state.get(type);
  if (count == null)
    count = 1;
  else
    count++;
  state.put(type, count);
  return state;
$$;

CREATE FUNCTION project_03.ccmapmax(input map<text, int>)
    RETURNS NULL ON NULL INPUT
    RETURNS map<text, int>
    LANGUAGE java
    AS $$
     Integer max = Integer.MIN_VALUE;
     String data="";
     for (String k : input.keySet()) {
         Integer tmp = input.get(k);
         if (tmp > max) { max = tmp; data = k; }
     }
     Map<String,Integer> mm = new HashMap<String,Integer>();
     mm.put(data,max);
     return mm;
$$;

CREATE OR REPLACE AGGREGATE group_and_count_q34(text)
    SFUNC state_group_and_count
    STYPE map<text, int>
    FINALFUNC ccmapmax
    INITCOND {};
```

After the setup is done, we will read the log file row by row and insert it into the database. Here, for each row of the file, we insert the the raw data pre-processed by regular expression into the `log` table.

```
INSERT INTO project_03.log
 (id, ip, identity, username, time, method, path, protocol, status, size)
 VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?, ?)
```

For the accessing IP address and the resource path being accessed, we use the `UPDATE` CQL to increment the `count` column by 1 in the corresponding tables, namely `ip` and `path`. This is feasible by using the `UPDATE` CQL alone since the `count` columns in the tables have the data type of counter. Upon invocation of a `UPDATE` CQL on a non-existing key, Cassandra automatically creates a row with the key and set the counter to 0. By setting `count = count + 1`, the value of the `count` column will be 1 for the key's appearance.

```
UPDATE project_03.ip SET count = count + 1 WHERE ip = ?
UPDATE project_03.path SET count = count + 1 WHERE path = ?
```

We used asynchronous execution on inserts. The JAVA class `Semaphore` was used to limit the number of

existing asynchronous requests. No more than 256 requests can exist in the request pool, matching the configuration of the Cassandra database.





The two figures above shows the process of inserting the log file. A total of `4477813 * 3 = 13433439` `INSERT`s and `UPDATE`s was performed in 1433 seconds, resulting in an average of about `9400 iops`. We can view the resulting tables in CQLSH. Results are shown in figures below.
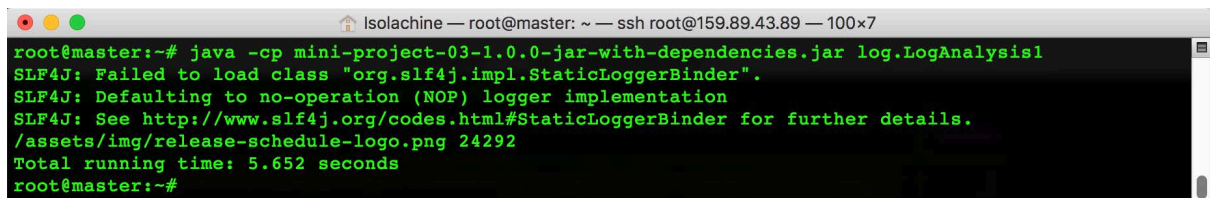
## Part 3: Operate Data in Cassandra

### Problem 1

The `LogAnalysis1.java` file is the source code for problem 1.
The program can be launched by the following command at any where, either on the cluster or a local machine.
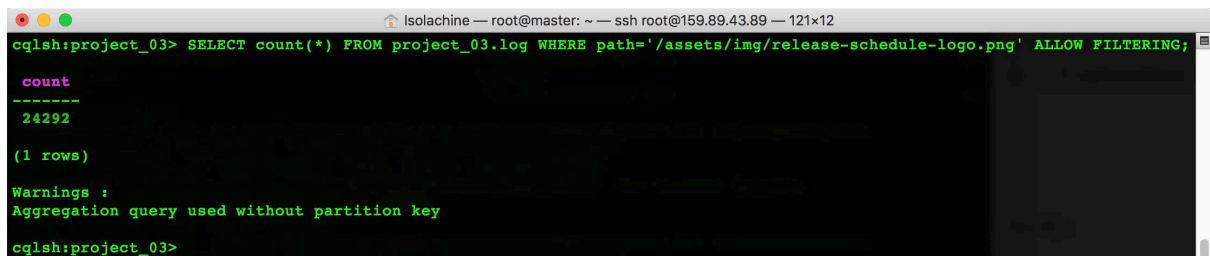
```
java -cp mini-project-03-1.0.0-jar-with-dependencies.jar log.LogAnalysis1
```



We can also simply get the result by running the CQL

```sql
SELECT count(*)
 FROM project_03.log
 WHERE path='/assets/img/release-schedule-logo.png'
 ALLOW FILTERING;
```



As the screenshots show, `/assets/img/release-schedule-logo.png` was accessed 24292 times.

### Problem 2

The `LogAnalysis2.java` file is the source code for problem 2.
The program can be launched by the following command at any where, either on the cluster or a local machine.

```
java -cp mini-project-03-1.0.0-jar-with-dependencies.jar log.LogAnalysis2
```

```
● ● ●                    Isolachine — root@master: ~ — ssh root@159.89.43.89 — 100×7
[root@master:~# java -cp mini-project-03-1.0.0-jar-with-dependencies.jar log.LogAnalysis2
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
10.207.188.188 398
Total running time: 3.446 seconds
root@master:~#
```

We can also simply get the result by running the CQL

```
SELECT count(*)
 FROM project_03.log
 WHERE ip='10.207.188.188'
 ALLOW FILTERING;
```

```
● ● ●                    Isolachine — root@master: ~ — ssh root@159.89.43.89 — 121×12
cqlsh:project_03> SELECT count(*) FROM project_03.log WHERE ip='10.207.188.188' ALLOW FILTERING;

 count
-------
   398

(1 rows)

Warnings :
Aggregation query used without partition key

cqlsh:project_03>
```

As the screenshots show, the IP address `10.207.188.188` accessed the website 398 times.


**Problem 3**

The `LogAnalysis3.java` file is the source code for problem 3.
The program can be launched by the following command at any where, either on the cluster or a local machine.

```
java -cp mini-project-03-1.0.0-jar-with-dependencies.jar log.LogAnalysis3
```

```
● ● ●                    Isolachine — root@master: ~ — ssh root@159.89.43.89 — 100×7
[root@master:~# java -cp mini-project-03-1.0.0-jar-with-dependencies.jar log.LogAnalysis3
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
/assets/css/combined.css 117348
Total running time: 3.975 seconds
root@master:~#
```

The JAVA program actually retrieves the whole `path` table and find the max count row internally. We can also get the same result from running the two following CQL in CQLSH.

Since Cassandra doesn't support subqueries in the latest version, we have to do it with two separate queries instead of a single nested query.

It's also possible to get the result from the `log` table with the UDF we defined earlier. Due to the fact the system has to query and aggregate over the big table of 4.7 million rows, this query may take a few minutes to run.



As the screenshots show, `/assets/css/combined.css` was the most accessed resource, with 117348 times.

**Problem 4**

The `LogAnalysis4.java` file is the source code for problem 4.
The program can be launched by the following command at any where, either on the cluster or a local machine.

```
java -cp mini-project-03-1.0.0-jar-with-dependencies.jar log.LogAnalysis4
```



The JAVA program actually retrieves the whole `ip` table and find the max count row internally. We can also get the same result from running the two following CQL in CQLSH.

Since Cassandra doesn't support subqueries in the latest version, we have to do it with two separate queries instead of a single nested query.

It's also possible to get the result from the `log` table with the UDF we defined earlier. Due to the fact the system has to query and aggregate over the big table of 4.7 million rows, this query may take a few minutes to run.



As the screenshots show, `10.216.113.172` accessed the website the most, with 158614 times.