

Final Assignment (Part 2) - Creating Streaming Data Pipelines using Kafka

Scenario

You are a data engineer at a data analytics consulting company. You have been assigned to a project that aims to de-congest the national highways by analyzing the road traffic data from different toll plazas. As a vehicle passes a toll plaza, the vehicle's data like vehicle_id, vehicle_type, toll_plaza_id and timestamp are streamed to Kafka. Your job is to create a data pipe line that collects the streaming data and loads it into a database.

Objectives

In this assignment you will create a streaming data pipe by performing these steps:

- Start a MySQL Database server.
- Create a table to hold the toll data.
- Start the Kafka server.
- Install the Kafka python driver.
- Install the MySQL python driver.
- Create a topic named toll in kafka.
- Download streaming data generator program.
- Customize the generator program to stream to toll topic.
- Download and customise streaming data consumer.
- Customize the consumer program to write into a MySQL database table.
- Verify that streamed data is being collected in the database table.

Exercise 1 - Prepare the lab environment

Step 1: Download Kafka.

- `wget https://archive.apache.org/dist/kafka/2.8.0/kafka_2.12-2.8.0.tgz`

Step 2: Extract Kafka.

- `tar -xzf kafka_2.12-2.8.0.tgz`

Step 3: Start MySQL server.

- `start_mysql`

Step 4: Connect to the mysql server, using the command below. Make sure you use the password given to you when the MySQL server starts. Please make a note or record of the password because you will need it later.

- `mysql --host=127.0.0.1 --port=3306 --user=root --password=Mjk0NDQtcnNhbm5h`

Step 5: Create a database named tolldata.

```
mysql > create database tolldata;
```

Step 6: Create a table named livetolldata with the schema to store the data generated by the traffic simulator.

Run the following command to create the table:

```
mysql > use tolldata;
```

```
mysql > create table livetolldata(timestamp datetime,vehicle_id int,vehicle_type char(15),toll_plaza_id smallint);
```

Step 7: Disconnect from MySQL server.

```
mysql > Exit
```

Step 8: Install the python module kafka-python using the pip command. This python module will help you to communicate with kafka server. It can used to send and receive messages from kafka.

- `python3 -m pip install kafka-python`

Step 9: Install the python module mysql-connector-python using the pip command. This python module will help you to interact with mysql server.

- `python3 -m pip install mysql-connector-python==8.0.31`

Exercise 2 - Start Kafka

Task 2.1 - Start Zookeeper

Start zookeeper server.

Task 2.2 - Start Kafka server

Start Kafka server

Task 2.3 - Create a topic named toll

Create a Kafka topic named toll

Task 2.4 - Download the Toll Traffic Simulator

Download the toll_traffic_generator.py from the url given below using 'wget'.

- `wget https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DB0250EN-SkillsNetwork/labs/Final%20Assignment/toll_traffic_generator.py`

Open the code using the theia editor using the "Menu -> File -> Open" option.

Task 2.5 - Configure the Toll Traffic Simulator

Open the toll_traffic_generator.py and set the topic to toll.

Task 2.6 - Run the Toll Traffic Simulator

Run the toll_traffic_generator.py.

Hint : `python3 <pythonfilename>` runs a python program on the theia lab.

Task 2.7 - Configure streaming_data_reader.py

Download the streaming_data_reader.py from the url below using 'wget'.

- `wget https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DB0250EN-SkillsNetwork/labs/Final%20Assignment/streaming_data_reader.py`

Open the streaming_data_reader.py and modify the following details so that the program can connect to your mysql server.

- TOPIC
- DATABASE
- USERNAME
- PASSWORD

Task 2.8 - Run streaming_data_reader.py

Run the streaming_data_reader.py

```
python3 streaming_data_reader.py
```

Task 2.9 - Health check of the streaming data pipeline.

If you have done all the steps till here correctly, the streaming toll data would get stored in the table livetolldata.

List the top 10 rows in the table livetolldata.

```

theia@theiadocker-delacruzdan1:/home/project$ cd kafka_2.12-2.8.0
theia@theiadocker-delacruzdan1:/home/project/kafka_2.12-2.8.0$ bin/zookeeper-server-start.sh config/zookeeper.properties
JVM9VM007W Command-line option unrecognized: -Xlog:gc*:file=/home/project/kafka_2.12-2.8.0/bin/../logs/zookeeper-gc.log:time,tags:filecount=10,file
size=100M
[2023-02-24 17:32:04,181] INFO Reading configuration from: config/zookeeper.properties (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2023-02-24 17:32:04,183] WARN config/zookeeper.properties is relative. Prepend ./ to indicate that you're sure! (org.apache.zookeeper.server.quorum
.QuorumPeerConfig)
[2023-02-24 17:32:04,198] INFO clientPortAddress is 0.0.0.0:2181 (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2023-02-24 17:32:04,199] INFO secureClientPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2023-02-24 17:32:04,207] INFO autopurge.snapRetainCount set to 3 (org.apache.zookeeper.server.DataDirCleanupManager)
[2023-02-24 17:32:04,207] INFO autopurge.purgeInterval set to 0 (org.apache.zookeeper.server.DataDirCleanupManager)

```

```

theia@theiadocker-delacruzdan1:/home/project$ cd kafka_2.12-2.8.0
theia@theiadocker-delacruzdan1:/home/project/kafka_2.12-2.8.0$ bin/kafka-server-start.sh config/server.properties
JVM9VM007W Command-line option unrecognized: -Xlog:gc*:file=/home/project/kafka_2.12-2.8.0/bin/../logs/kafkaServer-gc.log:time,tags:filecount=10,fi
lesize=100M
[2023-02-24 17:33:49,556] INFO Registered kafka:type=kafka.Log4jController MBean (kafka.utils.Log4jControllerRegistration$)
[2023-02-24 17:33:49,943] INFO Setting -Djdk.tls.rejectClientInitiatedRenegotiation=true to disable client-initiated TLS renegotiation (org.apache.
zookeeper.common.X509Util)

```

```

theia@theiadocker-delacruzdan1:/home/project$ cd kafka_2.12-2.8.0
theia@theiadocker-delacruzdan1:/home/project/kafka_2.12-2.8.0$ bin/kafka-topics.sh --create --topic toll --bootstrap-server localhost:9092
Created topic toll.

```

toll_traffic_generator.py x

```

toll_traffic_generator.py > ...
1  """
2  Top Traffic Simulator
3  """
4  from time import sleep, time, ctime
5  from random import random, randint, choice
6  from kafka import KafkaProducer
7  producer = KafkaProducer(bootstrap_servers='localhost:9092')
8
9  TOPIC = 'toll'
10
11  VEHICLE_TYPES = ("car", "car", "car", "car", "car", "car", "car", "car",
12                  "car", "car", "car", "truck", "truck", "truck",
13                  "truck", "van", "van")
14  for _ in range(100000):
15      vehicle_id = randint(10000, 10000000)
16      vehicle_type = choice(VEHICLE_TYPES)
17      now = ctime(time())
18      plaza_id = randint(4000, 4010)
19      message = f"{now},{vehicle_id},{vehicle_type},{plaza_id}"
20      message = bytearray(message.encode("utf-8"))
21      print(f"A {vehicle_type} has passed by the toll plaza {plaza_id} at {now}.")
22      producer.send(TOPIC, message)
23      sleep(random() * 2)
24

```

```

theia@theiadocker-delacruzdan1:/home/project$ python3 toll_traffic_generator.py
A truck has passed by the toll plaza 4001 at Fri Feb 24 17:40:02 2023.
A van has passed by the toll plaza 4008 at Fri Feb 24 17:40:04 2023.
A truck has passed by the toll plaza 4008 at Fri Feb 24 17:40:05 2023.
A van has passed by the toll plaza 4000 at Fri Feb 24 17:40:07 2023.
A car has passed by the toll plaza 4007 at Fri Feb 24 17:40:08 2023.
A car has passed by the toll plaza 4008 at Fri Feb 24 17:40:09 2023.
A car has passed by the toll plaza 4005 at Fri Feb 24 17:40:11 2023.
A truck has passed by the toll plaza 4000 at Fri Feb 24 17:40:12 2023.
A car has passed by the toll plaza 4010 at Fri Feb 24 17:40:14 2023.
A truck has passed by the toll plaza 4007 at Fri Feb 24 17:40:14 2023.
A van has passed by the toll plaza 4006 at Fri Feb 24 17:40:16 2023.
A truck has passed by the toll plaza 4007 at Fri Feb 24 17:40:17 2023.
A truck has passed by the toll plaza 4010 at Fri Feb 24 17:40:18 2023.
A car has passed by the toll plaza 4001 at Fri Feb 24 17:40:20 2023.
A truck has passed by the toll plaza 4009 at Fri Feb 24 17:40:21 2023.
A car has passed by the toll plaza 4008 at Fri Feb 24 17:40:21 2023.
A truck has passed by the toll plaza 4001 at Fri Feb 24 17:40:22 2023.
A car has passed by the toll plaza 4006 at Fri Feb 24 17:40:23 2023.
A car has passed by the toll plaza 4008 at Fri Feb 24 17:40:23 2023.
A car has passed by the toll plaza 4001 at Fri Feb 24 17:40:24 2023.

```

toll_traffic_generator.py streaming_data_reader.py ×

```
streaming_data_reader.py > ...
1  """
2  Streaming data consumer
3  """
4  from datetime import datetime
5  from kafka import KafkaConsumer
6  import mysql.connector
7
8  TOPIC='toll'
9  DATABASE = 'mysql'
10 USERNAME = 'delacruzdan1'
11 PASSWORD = 'MTA1NzgtZGVsYWMy'
12
13 print("Connecting to the database")
14 try:
15     connection = mysql.connector.connect(host='localhost', database=DATABASE, user=USERNAME, password=PASSWORD)
16 except Exception:
17     print("Could not connect to database. Please check credentials")
18 else:
19     print("Connected to database")
20     cursor = connection.cursor()
21
22 print("Connecting to Kafka")
23 consumer = KafkaConsumer(TOPIC)
24 print("Connected to Kafka")
25 print(f"Reading messages from the topic {TOPIC}")
26 for msg in consumer:
27
28     # Extract information from kafka
29
30     message = msg.value.decode("utf-8")
31
32     # Transform the date format to suit the database schema
33     (timestamp, vehcile_id, vehicle_type, plaza_id) = message.split(",")
34
35     dateobj = datetime.strptime(timestamp, '%a %b %d %H:%M:%S %Y')
36     timestamp = dateobj.strftime("%Y-%m-%d %H:%M:%S")
37
38     # Loading data into the database table
39
40     sql = "insert into livetolldata values(%s,%s,%s,%s)"
41     result = cursor.execute(sql, (timestamp, vehcile_id, vehicle_type, plaza_id))
42     print(f"A {vehicle_type} was inserted into the database")
43     connection.commit()
44 connection.close()
```

theia@theiadocker-delacruzdan1:/home/project\$ python3 streaming_data_reader.py

```
Connecting to the database
Connected to database
Connecting to Kafka
Connected to Kafka
Reading messages from the topic toll
A truck was inserted into the database
A car was inserted into the database
A car was inserted into the database
A car was inserted into the database
A van was inserted into the database
A truck was inserted into the database
A car was inserted into the database
A truck was inserted into the database
A car was inserted into the database
A car was inserted into the database
A car was inserted into the database
A truck was inserted into the database
A truck was inserted into the database
A car was inserted into the database
A van was inserted into the database
A car was inserted into the database
A car was inserted into the database
A van was inserted into the database
```

```
mysql> SELECT * FROM livetolldata LIMIT 10;
```

timestamp	vehicle_id	vehicle_type	toll_plaza_id
2023-02-24 17:52:01	5117515	truck	4005
2023-02-24 17:52:02	3986772	car	4009
2023-02-24 17:52:03	3252296	car	4007
2023-02-24 17:52:05	611507	car	4002
2023-02-24 17:52:06	1320647	van	4010
2023-02-24 17:52:07	8055309	truck	4010
2023-02-24 17:52:08	1439727	car	4002
2023-02-24 17:52:08	9523750	truck	4004
2023-02-24 17:52:10	9543838	car	4005
2023-02-24 17:52:12	3795626	car	4003

```
10 rows in set (0.00 sec)
```