**CREATION OF KAFKA CLUSTER AND EMR CLUSTER**
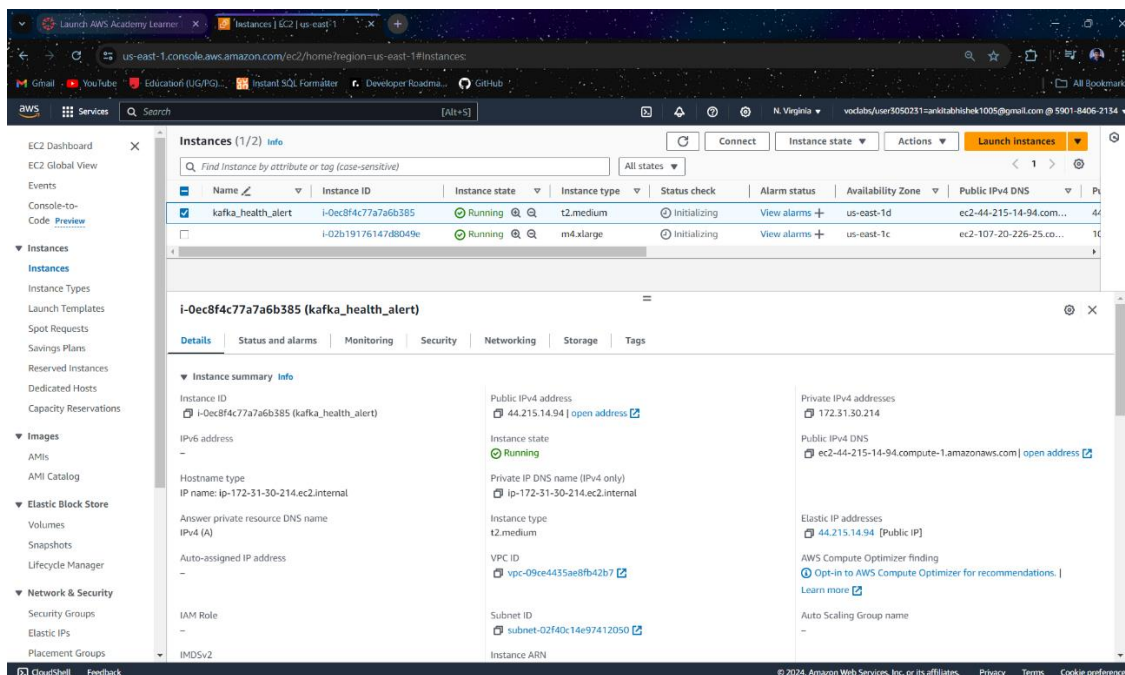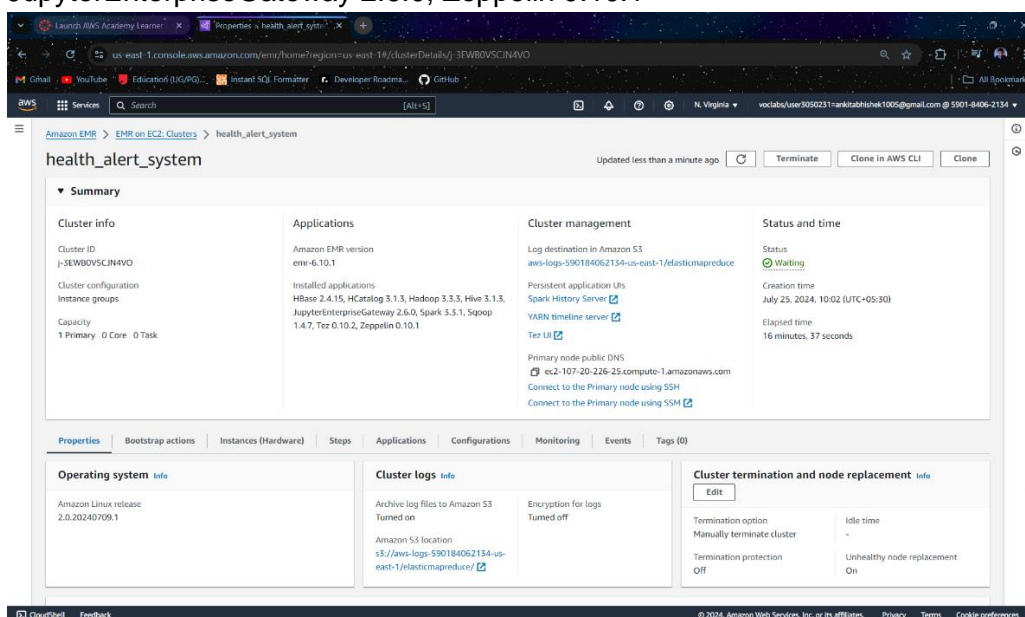
1. Create a kafka cluster with the help of ec2 instance. (Referred with PDF attached in the Apache Kafka modules). Make the required configuration needed to run kafka instance. Kafka is pre-installed on below mentioned ec2 cluster with the selection of **ami-06c41d8b5a6ddd3c2** while creating **Amazon Machine Image** as pdf within modules.



2. Create an EMR instance with required below mentioned libraries (Referred with PDF attached in the modules).
Spark 3.3.1, Sqoop 1.4.7, HBase 2.4.15, HCatalog 3.1.3, Hadoop 3.3.3, Hive 3.1.3, JupyterEnterpriseGateway 2.6.0, Zeppelin 0.10.1

# INSTALLING REQUIRED PACKAGES ON KAFKA CLUSTER

Sudo pip3 install kafka-python
Sudo pip3 install mysql-connector
Sudo pip3 install boto3

```
ec2-user@ip-172-31-30-214:~                                    —   □   ✕
[ec2-user@ip-172-31-30-214 ~]$ sudo pip3 install kafka-python
WARNING: Running pip install with root privileges is generally not a good idea.
Try `pip3 install --user` instead.
Collecting kafka-python
  Downloading kafka_python-2.0.2-py2.py3-none-any.whl (246 kB)
      |                                  | 246 kB 36.2 MB/s
Installing collected packages: kafka-python
Successfully installed kafka-python-2.0.2
[ec2-user@ip-172-31-30-214 ~]$ sudo pip3 install mysql-connector
WARNING: Running pip install with root privileges is generally not a good idea.
Try `pip3 install --user` instead.
Collecting mysql-connector
  Downloading mysql-connector-2.2.9.tar.gz (11.9 MB)
      |                                  | 11.9 MB 69 kB/s
Using legacy 'setup.py install' for mysql-connector, since package 'wheel' is no
t installed.
Installing collected packages: mysql-connector
    Running setup.py install for mysql-connector ... done
Successfully installed mysql-connector-2.2.9
[ec2-user@ip-172-31-30-214 ~]$ sudo pip3 install boto3
WARNING: Running pip install with root privileges is generally not a good idea.
Try `pip3 install --user` instead.
Collecting boto3
  Downloading boto3-1.33.13-py3-none-any.whl (139 kB)
      |                                  | 139 kB 13.2 MB/s
```

```
ec2-user@ip-172-31-30-214:~                                    —   □   ✕
Try `pip3 install --user` instead.
Collecting boto3
  Downloading boto3-1.33.13-py3-none-any.whl (139 kB)
      |                                  | 139 kB 13.2 MB/s
Collecting s3transfer<0.9.0,>=0.8.2
  Downloading s3transfer-0.8.2-py3-none-any.whl (82 kB)
      |                                  | 82 kB 122 kB/s
Collecting jmespath<2.0.0,>=0.7.1
  Downloading jmespath-1.0.1-py3-none-any.whl (20 kB)
Collecting botocore<1.34.0,>=1.33.13
  Downloading botocore-1.33.13-py3-none-any.whl (11.8 MB)
      |                                  | 11.8 MB 35 kB/s
Collecting python-dateutil<3.0.0,>=2.1
  Downloading python_dateutil-2.9.0.post0-py2.py3-none-any.whl (229 kB)
      |                                  | 229 kB 59.4 MB/s
Collecting urllib3<1.27,>=1.25.4; python_version < "3.10"
  Downloading urllib3-1.26.19-py2.py3-none-any.whl (143 kB)
      |                                  | 143 kB 63.5 MB/s
Collecting six>=1.5
  Downloading six-1.16.0-py2.py3-none-any.whl (11 kB)
Installing collected packages: six, python-dateutil, urllib3, jmespath, botocore
, s3transfer, boto3
Successfully installed boto3-1.33.13 botocore-1.33.13 jmespath-1.0.1 python-date
util-2.9.0.post0 s3transfer-0.8.2 six-1.16.0 urllib3-1.26.19
[ec2-user@ip-172-31-30-214 ~]$
```

**STATEMENT FOR STARTING KAFKA SERVER**

1. **STARTING ZOOKEEPER SERVER**:
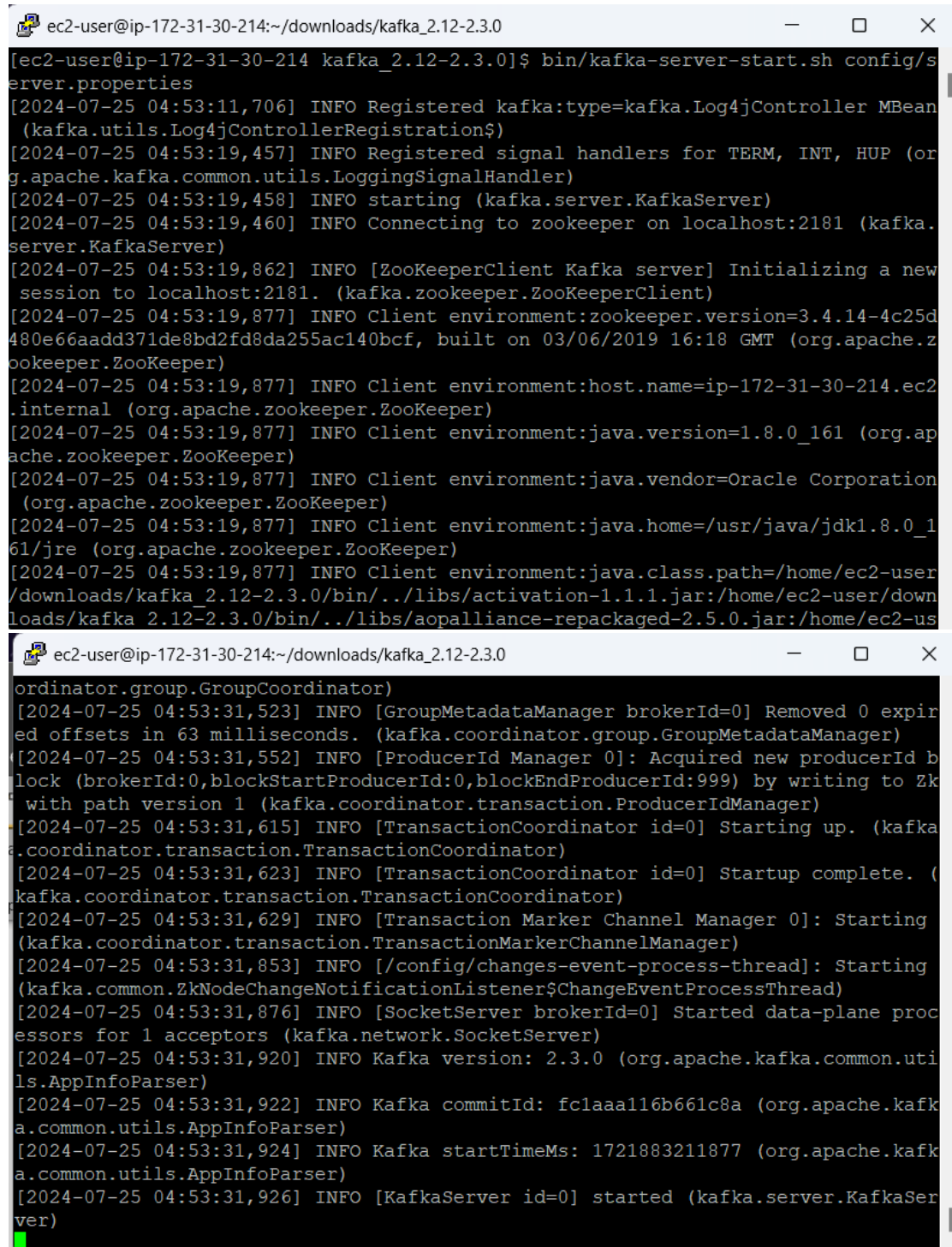
   Inside cd downloads/kafka_2.12-2.3.0 run

   **bin/zookeeper-server-start.sh config/zookeeper.properties**

2. **STARTING KAFKA SERVER:**

Into another putty Session of kafka cluster inside cd downloads/kafka_2.12-2.3.0 run
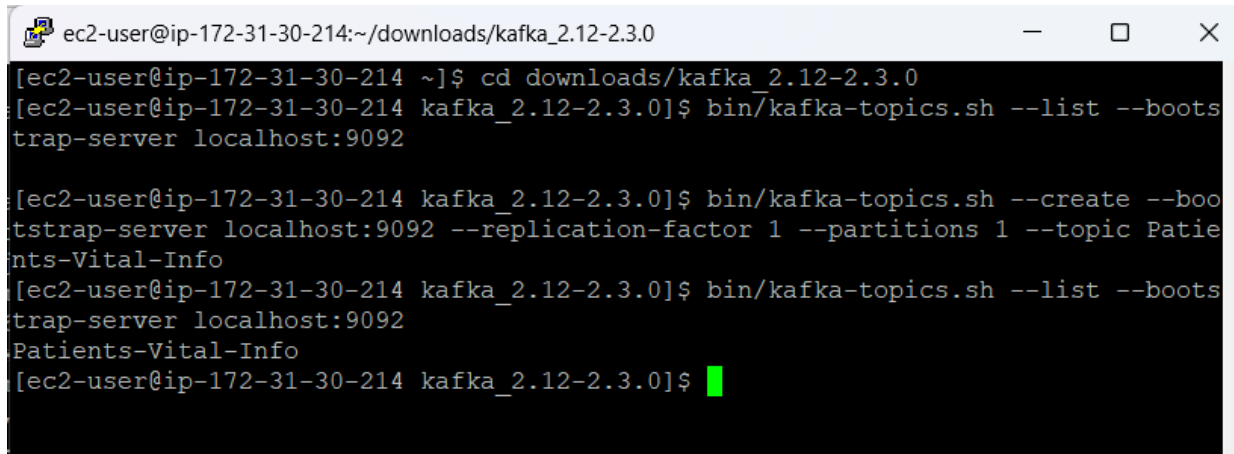**bin/kafka-server-start.sh config/server.properties**

## STATEMENT TO CREATE TOPICS

To create topic in kafka server, the command used is

**bin/kafka-topics.sh --create --bootstrap-server localhost:9092 --replication-factor 1 --partitions 1 --topic Patients-Vital-Info**

## STATEMENT TO LIST TOPICS

To list the created topic inside cd downloads/kafka_2.12-2.3.0, the command used is

**bin/kafka-topics.sh --list --bootstrap-server localhost:9092**

```
ec2-user@ip-172-31-30-214:~/downloads/kafka_2.12-2.3.0                    —    □    ×

[ec2-user@ip-172-31-30-214 ~]$ cd downloads/kafka_2.12-2.3.0
[ec2-user@ip-172-31-30-214 kafka_2.12-2.3.0]$ bin/kafka-topics.sh --list --boots
trap-server localhost:9092

[ec2-user@ip-172-31-30-214 kafka_2.12-2.3.0]$ bin/kafka-topics.sh --create --boo
tstrap-server localhost:9092 --replication-factor 1 --partitions 1 --topic Patie
nts-Vital-Info
[ec2-user@ip-172-31-30-214 kafka_2.12-2.3.0]$ bin/kafka-topics.sh --list --boots
trap-server localhost:9092
Patients-Vital-Info
[ec2-user@ip-172-31-30-214 kafka_2.12-2.3.0]$ █
```

## EXECUTING PRODUCER APPLICATION AND CONSUMER APPLICATION:

Producer application which is file named as **kafka_produce_patient_vitals.py** is built on the **python language** which will consume data residing on rds with below mentioned credentials:

Hostname = "upgraddetest.cyaielc9bmnf.us-east-1.rds.amazonaws.com"
username = "student"
password = "STUDENT123"
dbname =  "testdatabase".

Consumer Application which is file named as **kafka_spark_patient_vitals.py** is built on the **Apache PySpark** language which will consume data being produced with the help of above mentioned producer application

**NOTE: Run the producer application on ec2 Kafka cluster after starting the consumer application on EMR cluster created with Spark, Hive and another libraries**

**STATEMENT FOR EXECUTING PRODUCER APPLICATION AND CONSUMER APPLICATION**

**Spark Submitting Job to Consume Message from The Topic Patients-Vital-Info And Stored To HDFS Location**

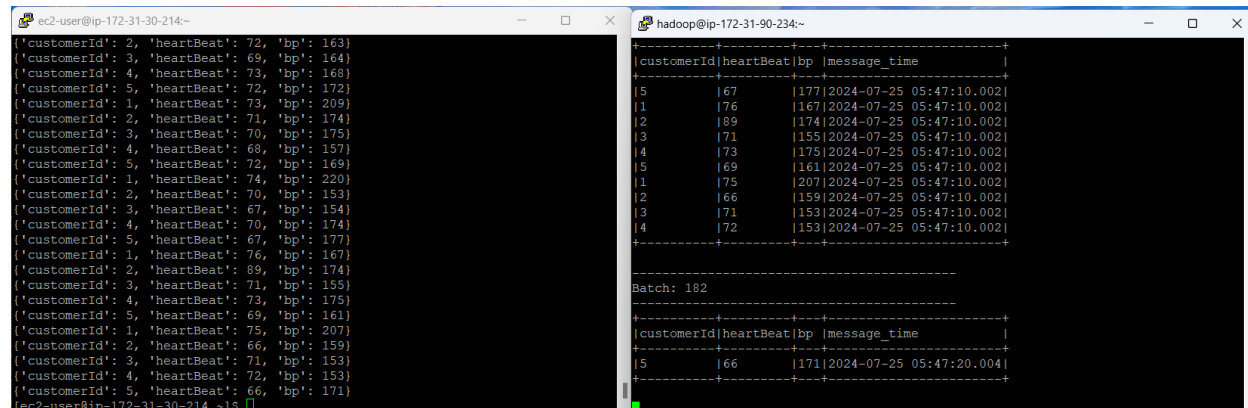For Producer application: **python3 kafka_produce_patients_vitals.py**

For Consumer Application: **spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.3.1 kafka_spark_patient_vitals.py**

```
hadoop@ip-172-31-90-129:~                                          —    □    ×
EE:::::EEEEEEEE:::::E M:::::M              M:::::M   R:::R        R::::R
E:::::::::::::::::::E M:::::M              M:::::M RR:::R         R:::R
EEEEEEEEEEEEEEEEEEEE MMMMMMM              MMMMMMM RRRRRRR        RRRRRR

[hadoop@ip-172-31-90-129 ~]$ ls
[hadoop@ip-172-31-90-129 ~]$ ls
kafka_spark_patient_vitals.py
[hadoop@ip-172-31-90-129 ~]$ spark-submit --packages org.apache.spark:spark-sql-
kafka-0-10_2.12:3.3.1 kafka_spark_patient_vitals.py
:: loading settings :: url = jar:file:/usr/lib/spark/jars/ivy-2.5.0.jar!/org/apa
che/ivy/core/settings/ivysettings.xml
Ivy Default Cache set to: /home/hadoop/.ivy2/cache
The jars for the packages stored in: /home/hadoop/.ivy2/jars
org.apache.spark#spark-sql-kafka-0-10_2.12 added as a dependency
:: resolving dependencies :: org.apache.spark#spark-submit-parent-3252d621-c6ec-
4b42-ba7b-10c52d78769d;1.0
        confs: [default]
        found org.apache.spark#spark-sql-kafka-0-10_2.12;3.3.1 in central
        found org.apache.spark#spark-token-provider-kafka-0-10_2.12;3.3.1 in cen
tral
        found org.apache.kafka#kafka-clients;2.8.1 in central
        found org.lz4#lz4-java;1.8.0 in central
        found org.xerial.snappy#snappy-java;1.1.8.4 in central
```

```
hadoop@ip-172-31-90-129:~                                          —    □    ×
-129.ec2.internal:45885
24/07/21 10:19:50 INFO BlockManager: Using org.apache.spark.storage.RandomBlockR
eplicationPolicy for block replication policy
24/07/21 10:19:50 INFO BlockManager: external shuffle service port = 7337
24/07/21 10:19:50 INFO BlockManagerMaster: Registering BlockManager BlockManager
Id(driver, ip-172-31-90-129.ec2.internal, 45885, None)
24/07/21 10:19:50 INFO BlockManagerMasterEndpoint: Registering block manager ip-
172-31-90-129.ec2.internal:45885 with 912.3 MiB RAM, BlockManagerId(driver, ip-1
72-31-90-129.ec2.internal, 45885, None)
24/07/21 10:19:50 INFO BlockManagerMaster: Registered BlockManager BlockManagerI
d(driver, ip-172-31-90-129.ec2.internal, 45885, None)
24/07/21 10:19:50 INFO BlockManager: Initialized BlockManager: BlockManagerId(dr
iver, ip-172-31-90-129.ec2.internal, 45885, None)
24/07/21 10:19:51 INFO SingleEventLogFileWriter: Logging events to hdfs:/var/log
/spark/apps/local-1721557189930.inprogress
-------------------------------------------
Batch: 0
-------------------------------------------
+----------+---------+---+------------+
|customerId|heartBeat|bp |message_time|
+----------+---------+---+------------+
+----------+---------+---+------------+
```

After 30 minutes when all 1800 data being streamed and saved to Parquet file of the required HDFS location



**STATEMENT TO CHECK DATA STORED IN HDFS LOCATION:**
  hadoop fs -ls /user/hadoop/health-alert/patients-vital-info/

**STATEMENT TO READ ONE OF THE FILES USING '-CAT'**
    hadoop fs -cat  /user/hadoop/health-alert/patients-vital-info/part-00000-ffcd45dc-ef2e-4219-b11e-f26f0e051d73-c000.snappy.parquet