# Collections : Lab Guide

## Streaming Simulation using Amazon Kinesis Stream

## &

## Real Time Data Analytics using Amazon Kinesis Analytics

**Developed & Tested**

**By**

**Karthigayen.Y**

**LKM, Accenture - ATCI**

# Context :

In this Lab we will be doing a Simualtion of Streaming Data from the Crime dataset using Amazon Kinesis Stream, And ingest that data into S3 and Amazon Kinesis Analytics , to perform basic Real Time Analysis of the Crime dataset.

## Steps to simulate the streaming data from chicago_crime_dataset.csv file using Kinesis Stream

- ➢ Place the **chicago_crime_dataset.csv** file in a specific directory in C: drive on your windows machine.
    - o Ex : C:\AWS_Dataset\chicago_crime_dataset.csv

- ➢ Download and install **AWS CLI V2** on windows and set the Environmental Variables as show below
    - o https://docs.aws.amazon.com/cli/latest/userguide/install-cliv2-windows.html#cliv2-windows-install
    - o

```
C:\windows\system32\cmd.exe
Microsoft Windows [Version 10.0.18363.1237]
(c) 2019 Microsoft Corporation. All rights reserved.

C:\Users\y.karthigayen>cd \

C:\>aws configure
AWS Access Key ID [****************QF4I]: AKIAXERMTLGJHD2PDAUE
AWS Secret Access Key [****************P3zr]: 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
Default region name [US_EAST_1]: US_EAST_1
Default output format [None]:

C:\>_
```
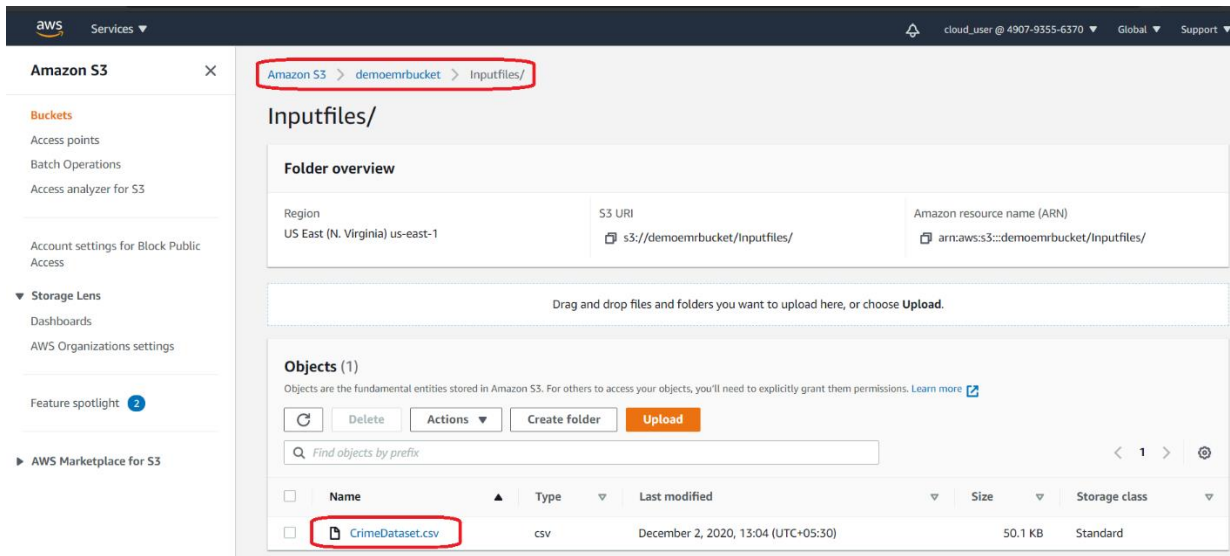
- ➢ Using the AWS console , create a Kinesis Data Stream named **"chicagoCrimeStream".**

- ➢ Use the Maven Project shared with you, which has the below classes to read the Data from the specified location and provide us the simulation of streaming data and write the data into the stream which we had created.
    - o **AwsKinesisClient.java**
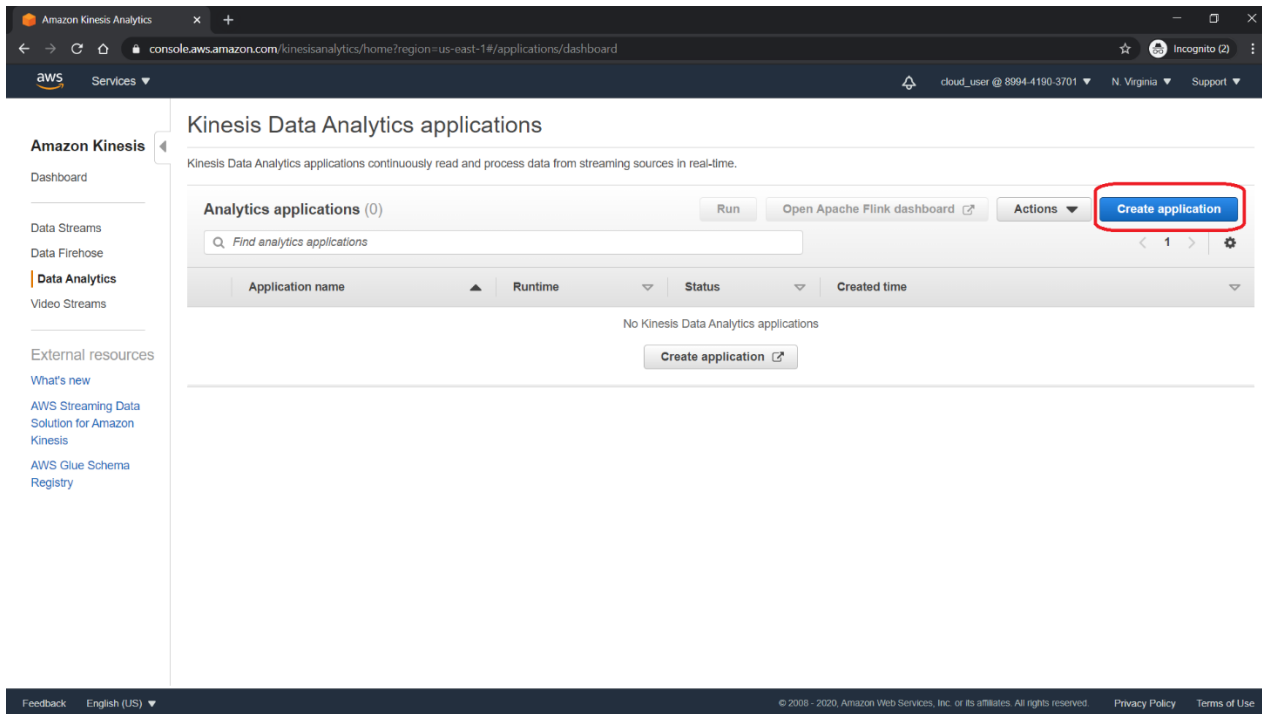    - o **CrimeDataWriter.java**

## Steps to read the Data from the Kinesis Stream and write it into a S3 bucket

- ➢ Using the AWS console, create a bucket named **"demoemrbucket"** on **S3 .**

- ➢ The below code will help us read the Data writen on to the **chicagoCrimeStream** and ingest the same into the S3 bucket which we had created
    - o **CrimeDataS3Put.java**

**Steps to create a Kinesis Data Analytics Application and fetch the data from Kinesis Stream "chicagoCrimeStream" and do the Analysis.**

**Step 1 :** Click on Create Application button

LKM                                        Collection – Lab Guide

**Step 2 :** Specify the name of the application as **CrimeDataAnalytics,** and the **SQL** as the environment to do the Analysis.



**Step 3 :** Click on Create Application button

LKM                                        Collection – Lab Guide

**Step 4 :** Click on the Connect Streaming Data button and choose Kinesis Data Stream as the Source.

**Step 5 :** Click on the Discover Schema button to Automatically detect the schema for Streaming Data



**Step 6 :** If the schema is not as expected , then click on Edit Schema button and edit the schema as required.

LKM                                        Collection – Lab Guide

**Step 7 :** Add the below mentioned Columns for Crime Dataset.

| Column order | Column name | Column type | |
|---|---|---|---|
| 1 | ID | INT | |
| 2 | CASENO | VARCHAR | Length: 20 |
| 3 | DATEOFCRIME | DATE | |
| 4 | BLOCK | VARCHAR | Length: 50 |
| 5 | IUCR_CODE | VARCHAR | Length: 10 |
| 6 | LOCATION_DESC | VARCHAR | Length: 50 |
| 7 | ARREST | VARCHAR | Length: 7 |
| 8 | DOMESTIC | VARCHAR | Length: 7 |
| 9 | BEAT_NUM | INT | |
| 10 | DISTRICT_CODE | INT | |
| 11 | WARD_NO | INT | |
| 12 | COMMUNITY_CODE | INT | |
| 13 | FBI_CODE | VARCHAR | Length: 7 |
| 14 | X_COORD | INT | |

**Step 8 :** Once all the columns are add , click on Save Schema button

| | | | |
|---|---|---|---|
| 15 | Y_COORD | INT | |
| 16 | CASEYEAR | INT | |
| 17 | DATE_OF_UPDATE | DATE | |
| 18 | LATITUDE | DECIMAL | |
| 19 | LONGITUDE | DECIMAL | |
| 20 | LOCATION | VARCHAR | Length: 30 |

Cancel     **Save schema and update stream samples**

Application status:  READY

Raw     Lambda output     Formatted     Error stream

**Step 9 :** Now click on Exit(Done) link to go to Analytics page.



**Step 10 :** Click on Go to SQL editor , to create the Analytical Application code



## Steps to Write Analytical Queries on Kinesis- Analytics:

  ➢  Create a Destination Stream
  ➢  Create a Stream Pump to push the resultant data into the Destination stream
  ➢  With the Analytical query triggered on the Source Data Stream

**Use Case 1 : To list all the Case numbers where Arrest has happened**

CREATE OR REPLACE STREAM "DESTINATION_SQL_STREAM" (CASENO VARCHAR(20) );

CREATE OR REPLACE PUMP "STREAM_PUMP" AS INSERT INTO "DESTINATION_SQL_STREAM"

SELECT STREAM CASENO

FROM "SOURCE_SQL_STREAM_001"

WHERE ARREST SIMILAR TO '%TRUE%';

**Destination Stream Output :**

## Use Case 2 : To check if the Case-ID is unique :

CREATE OR REPLACE STREAM "DESTINATION_SQL_STREAM" (ID INT, ID_count INTEGER);

CREATE OR REPLACE  PUMP "STREAM_PUMP" AS INSERT INTO "DESTINATION_SQL_STREAM"

SELECT STREAM ID, COUNT(*) AS ID_count

FROM "SOURCE_SQL_STREAM_001"

GROUP BY ID, FLOOR(("SOURCE_SQL_STREAM_001".ROWTIME - TIMESTAMP '2020-12-01 11:30:00') SECOND / 10 TO SECOND);

**Output :**



## Use Case 3 : Create report on total number of crime cases on each day from crimes dataset

CREATE OR REPLACE STREAM "DESTINATION_SQL_STREAM" (DATEOFCRIME DATE, COUNT_OF_CASES INT );

CREATE OR REPLACE PUMP "STREAM_PUMP" AS INSERT INTO "DESTINATION_SQL_STREAM"

SELECT STREAM DATEOFCRIME, COUNT(CASENO) AS COUNT_OF_CASES

FROM "SOURCE_SQL_STREAM_001"

GROUP BY DATEOFCRIME, FLOOR(("SOURCE_SQL_STREAM_001".ROWTIME - TIMESTAMP '2020-12-01 11:30:00') SECOND / 10 TO SECOND);

**Output :**  Note - the aggregation is for specific duration (2 to 10 Sec), hence the groups may repeat.

LKM                                    Collection – Lab Guide