# School of Information studies
# Syracuse University

# IST615 - Cloud Management

# Final Project Progress Report

Professor in-charge: Carlos E. Caicedo Bastidas
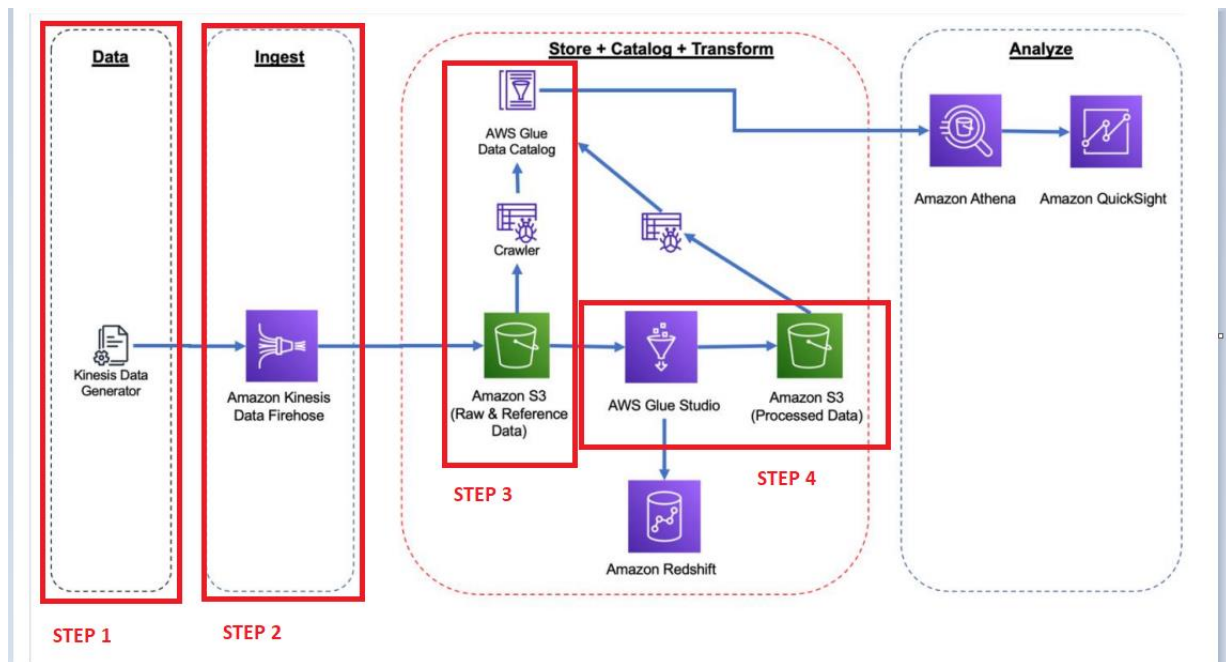
Group 6:

Names:

    1) Rahul Jadhav                 SUID: 705889151

    2) Ramazan Yener            SUID: 4604870161

    3) Francia Lizeth Ortiz Leyva    SUID: 2566771622

    4) Himamshu Chandrashekara   SUID: 254732602

Due Date – 11/10/2021

**<u>Tasks completed so far</u>** -



<u>Step 1</u> -

Generated Data from Kinesis Data generator.

<u>Step 2</u> –

Delivered the data to Amazon S3 bucket with Kinesis Firehose delivery stream. Processed and stored our data into 1$^{st}$ S3 bucket (Raw and Reference data).

<u>Step 3</u> –

Registered the datasets in the AWS Glue Data Catalog and automated the metadata captured with the help of Glue Crawlers.

<u>Step 4</u> –

Transformed Data from 1$^{st}$ AWS s3 bucket from table named raw to 2$^{nd}$ AWS S3 bucket (Processed data) from table named reference_data using AWS Glue Studio.

Description: In the above screenshot we selected our first source data from Data Catalog Table i.e. raw Table



Description: In the above screenshot we selected the another source data from Data Catalog Table i.e. reference_data Table.

Description: We used inner join transformation to join the parent s3 buckets data(raw & reference_data) using the common field as track_id



Description : In the above screenshot we apply mapping to the new target keys for the pushing the new processed dat into s3 bucket.

Description: In the above screenshot we provided the path for the target s3 bucket.

Issues Encountered

- As of now, we haven't experienced any obstacles.

Changes to the project/goals

- None

Plan for completion of project week by week plan

Week 11/15 - 11/21

- Pushed this data into an Amazon Redshift Data for easier accessibility and referencing.
- Analyse the data using AWS Athena using SQL queries.

Week 11/29 - 12/05

- Use Amazon Quicksight to build visualizations over the data collected and stored in S3.