# Analyzing Retail Sales of Electricity

## Milestone 6

# Group 15

Manaswini Kamtam
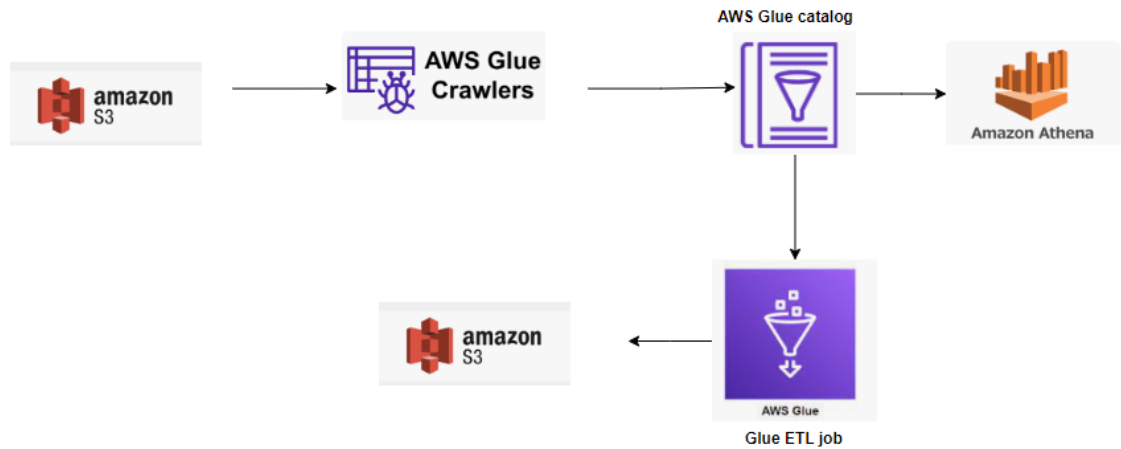
Aditya Bharadwaj Shivapura Guruprasad

kamtam.m@northeastern.edu

shivapuraguruprasa.a@northeastern.edu

**Submission Date:_____12/01/2023_____**

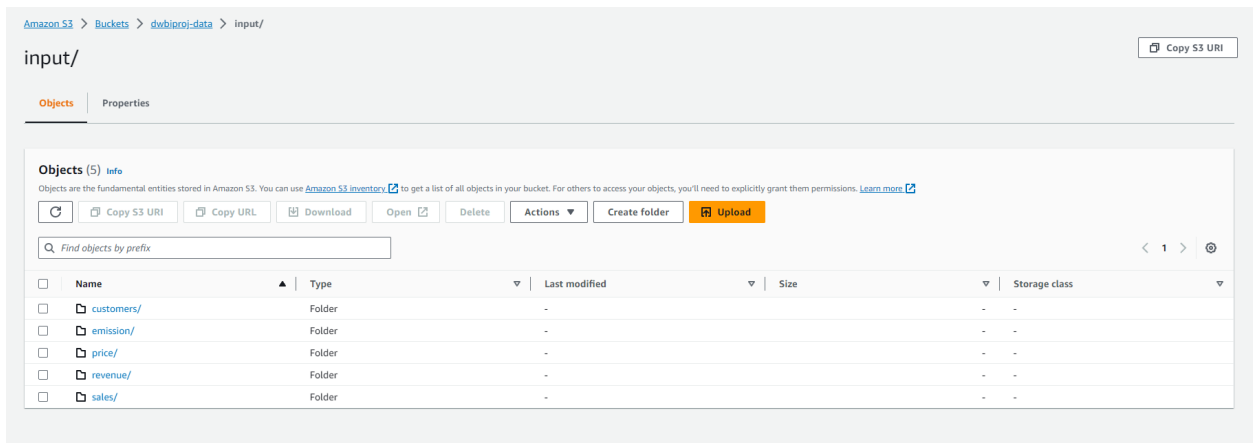**Architecture Diagram of the Data Pipeline**



**Steps involved in building the pipeline:**

1. Loaded data in S3 bucket

At first, we loaded the csv data to the S3 bucket by manual upload. Separate folders are created for each file

2. Created Glue crawlers for each data table

Then AWS Glue crawler is created for each of the csv file and a new database is created in the catalog to store the data



3. Ran the crawler and loaded the data in Glue catalog

Crawler is run to read the schema of the data present in the s3 bucket and load the data in the catalog DB



Made the following change to table schema since input data was hasquotes in the csv file

Data shifting was happening without the below configuration. Delimeters were introduced in the below steps

AWS Glue > Tables > Edit table

## Edit table

### ▼ Table details

Name

sales

Input format

org.apache.hadoop.mapred.TextInputFormat

Output format

org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat

Serde name

Serialization lib

org.apache.hadoop.hive.serde2.OpenCSVSerde

Description

### ▼ Serde parameters

| Key | Value | |
|---|---|---|
| escapeChar | \ | Remove |
| quoteChar | " | Remove |
| separatorChar | \| | Remove |

Add

### ▼ Table properties

4. The data is queried using Aws Athena to read the data from the catalog DB



5. Created a new Glue ETL visual based job with the below configuration

6. Joining two input tables (sales and revenue) and loading data to S3 bucket - as a csv



Job by default writes to multiple csv parallely. This feature is disabled by updating the job script and only one csv is obtained

7. ETL job output as seen in s3 bucket