# Batch Time Analysis of Transactional Data

## Description

Lenodo is a multinational e-commerce organization that sells products directly to consumers. The database administrator exports the data every night in a CSV file, but this export functionality is unused. Lenodo wants to use this data to uncover insights about the most-sold item and the countries where customers have bought this item.
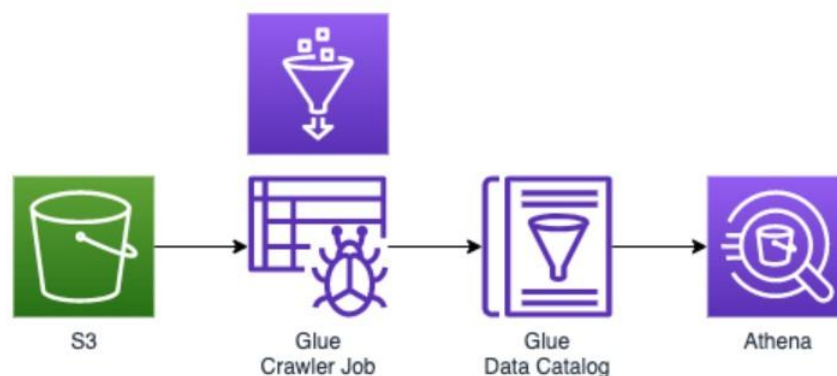
You are a data analytics consultant, and you're asked to provide valuable insights and statistics across products, brands, categories, segments to the marketing, product, sales, and procurement teams and inform them about which product has the highest amount of sales and which product and its marketing needs the most improvement. These statistics will help to run effective digital marketing campaigns. The scope of this project is limited to data engineering and analysis.

## Objective:

To use AWS Big Data stack for data engineering to analyze transactions, uncover patterns, and share actionable insights
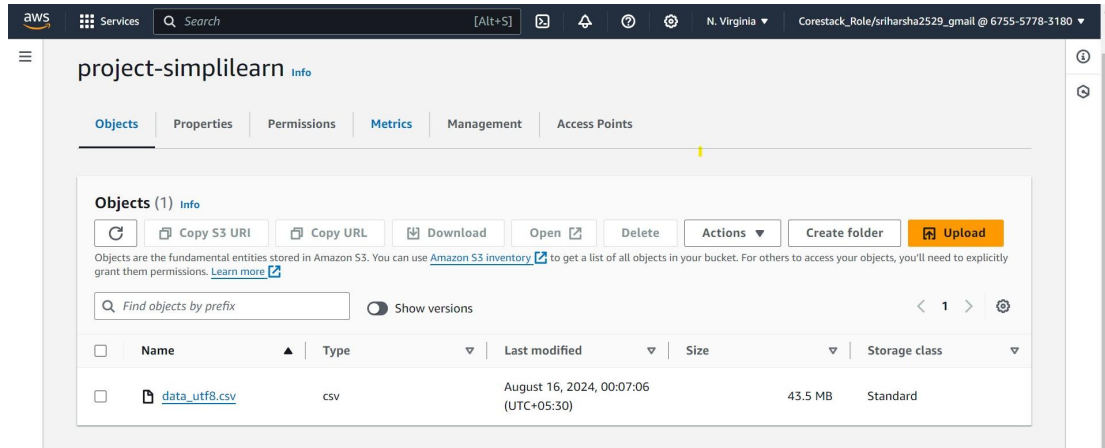
## Steps to perform

1. Create an S3 bucket with a unique name and upload the CSV file to the S3 bucket

2. Create a crawler to crawl the CSV data and generate a metadata catalog

3. Create a Glue job to transform the data into the Parquet format as CSV is not optimal for data warehouse queries

4. Add another crawler to crawl the Parquet data files to generate the metadata catalog of the Parquet file in order to query it with Athena

5. Query the data to identify the best-selling item and countries where customers have bought the most-sold item using Athena.
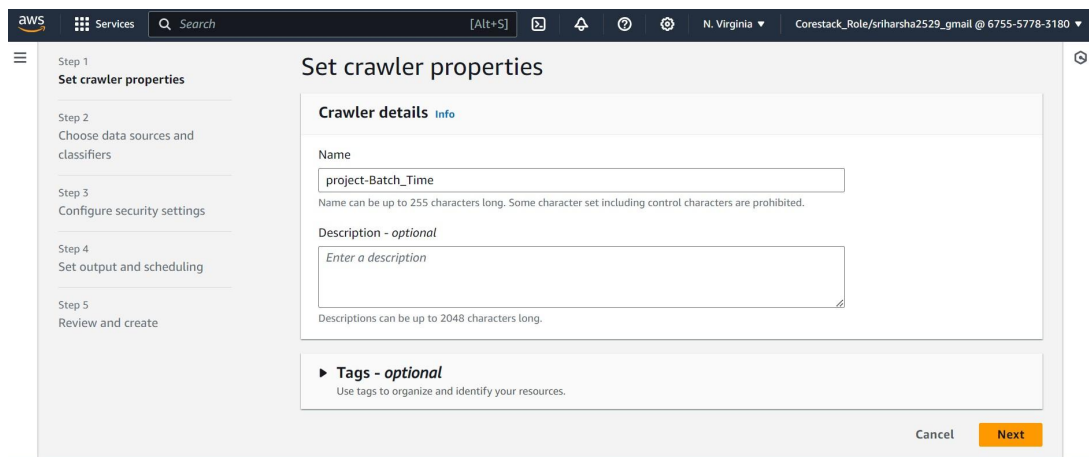
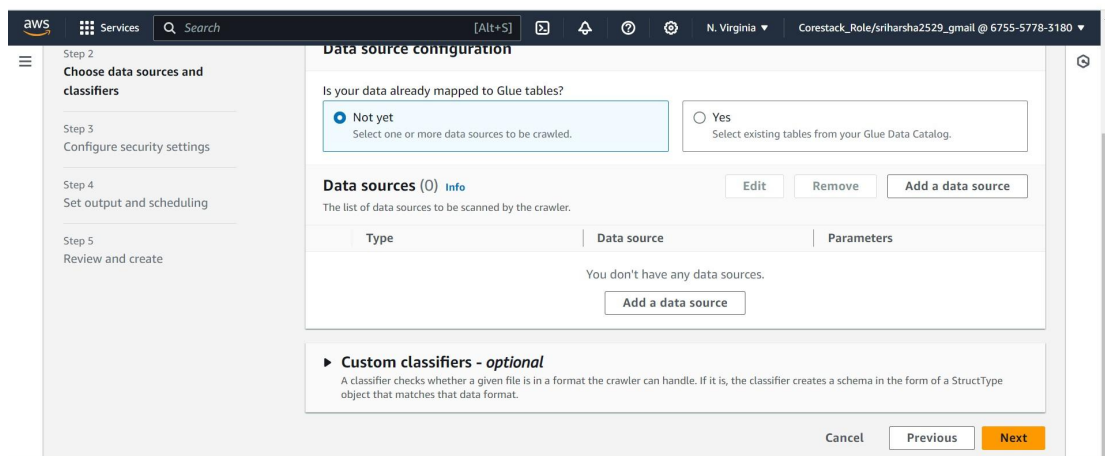# Batch Time Analysis of Transactional Data

Create a S3 Bucket and loaded data



Step 2 : Created A crawler



Step 3: Add a Data source from S3

# Batch Time Analysis of Transactional Data

## Step 4 : Data Source Added from S3



## Step 5 Data Base name



## Step 6 database created

# Batch Time Analysis of Transactional Data

## Step7schedule



## Step 8 Iam Rule creation

# Batch Time Analysis of Transactional Data

Step 9 Preview and create



Step 10 AThena

# Batch Time Analysis of Transactional Data

Step 11 Quering



Step 12 result

# Batch Time Analysis of Transactional Data

# Batch Time Analysis of Transactional Data





Step 11 Quering

# Batch Time Analysis of Transactional Data