# Customer Sales Data Transformation with AWS Glue
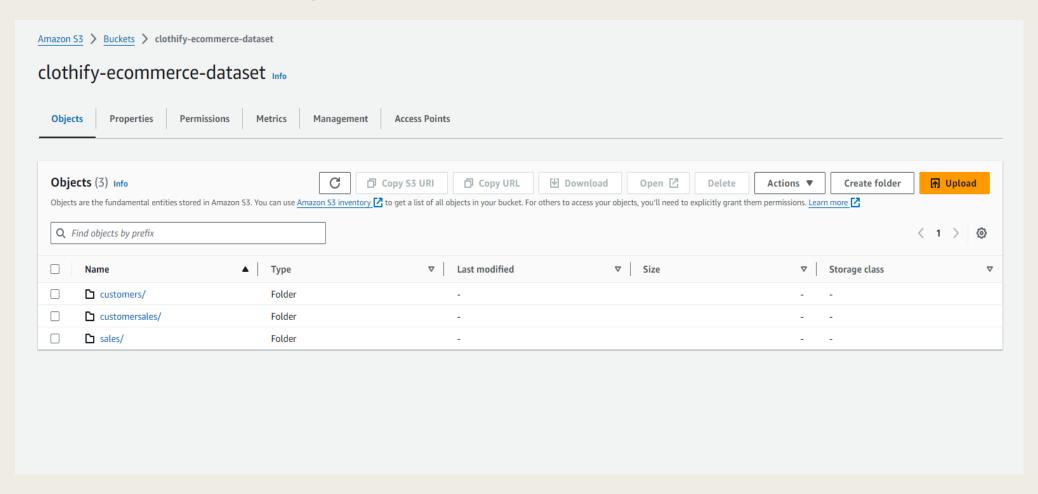
Sameer Singh

https://www.linkedin.com/in/sameer-singh-data/
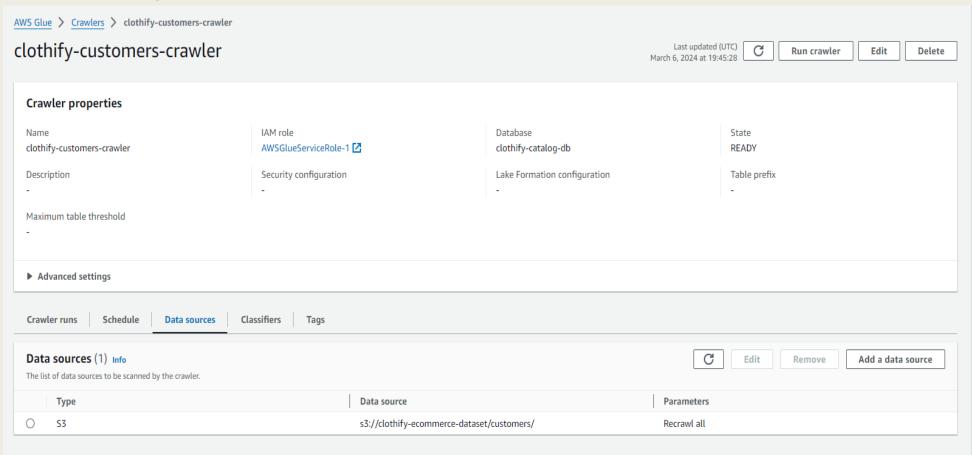
# Overview

Built an ETL pipeline using AWS Glue, seamlessly handling data from ingestion through transformation and storage in Amazon S3.

- ✓ Leveraged AWS services for integrated customer and sales data analysis.
- ✓ Utilized Amazon S3 for efficient data storage.
- ✓ Employed AWS Glue for metadata management, ETL processes, and data transformation.
- ✓ Combined datasets, apply filtering, and enhance data clarity with AWS Glue.
- ✓ Stored transformed data in S3, organized using Parquet files.
- ✓ Utilized Amazon Athena for efficient querying, facilitating comprehensive analysis.
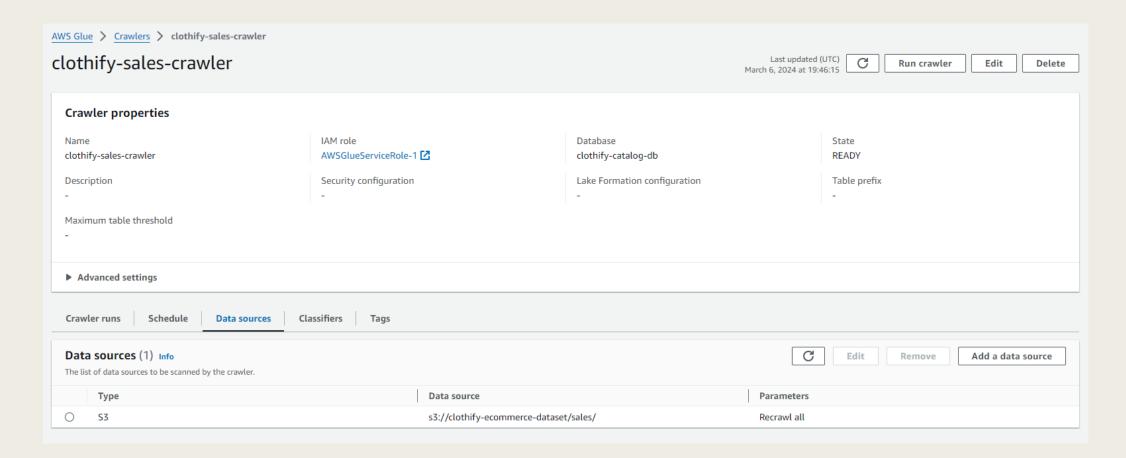
Created three S3 folders in buckets: two dedicated to customers and sales data, and one for storing Glue job results.
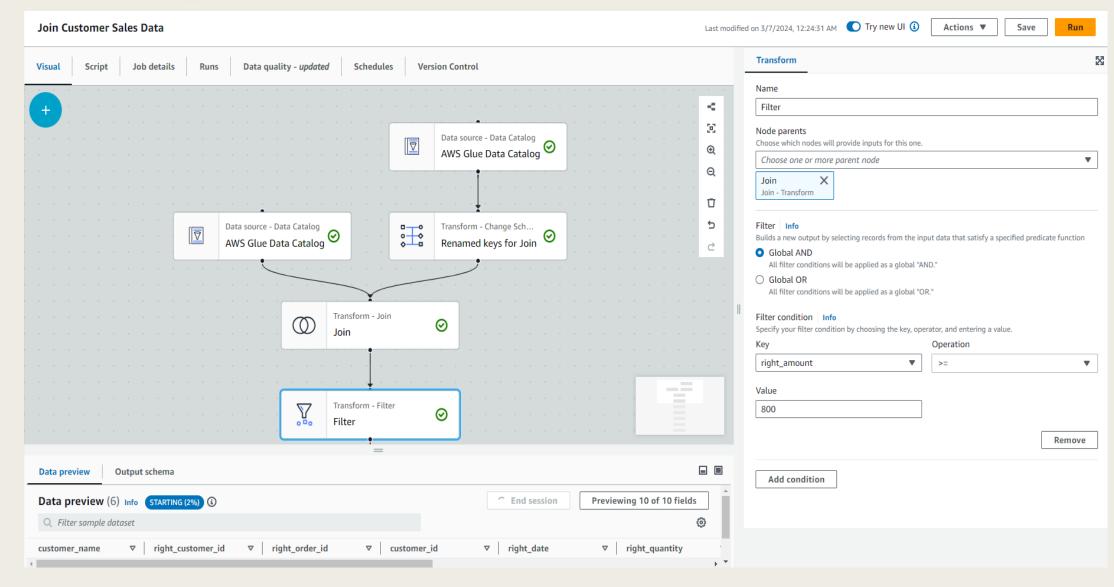
Utilized Glue Crawlers to catalog metadata for Customer Dataset and creating Data Catalog tables.

# Same Process for Sales Dataset and creating Data Catalog tables.

# Developed a Glue ETL job using two Data Catalog tables for customers and sales datasets.
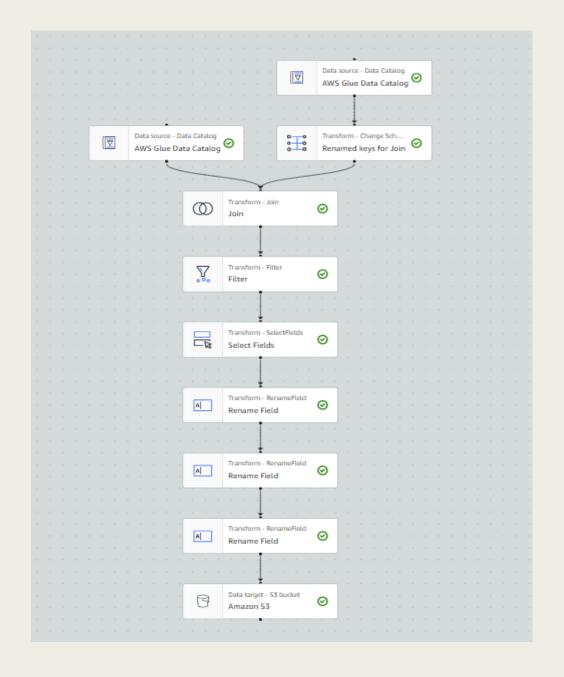
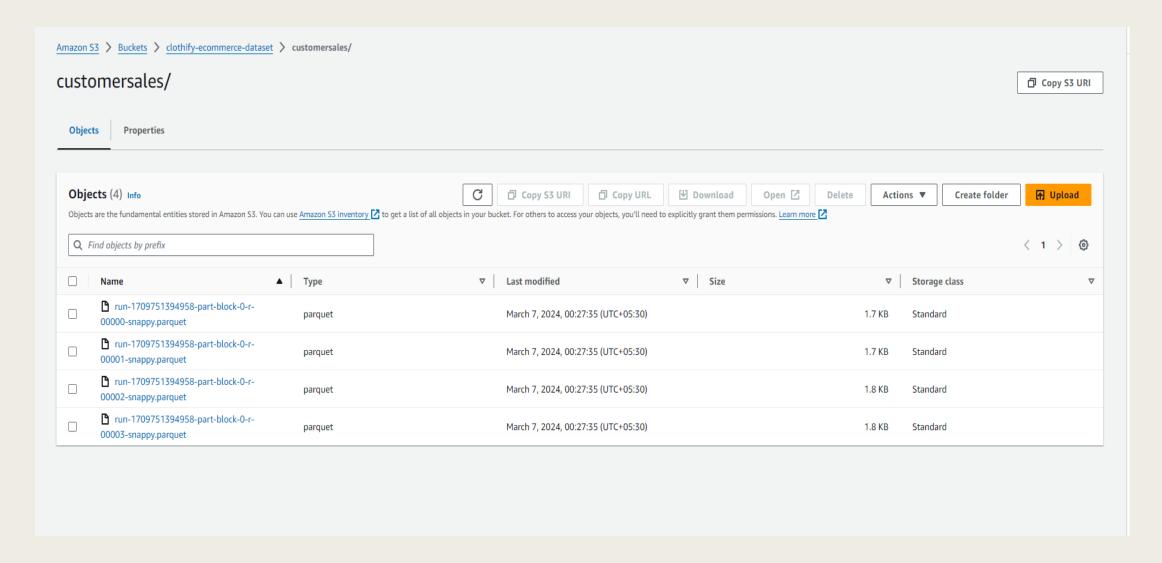Utilized AWS Glue ETL Job with metadata from Data Catalog tables.

Executed a join operation on 'customer id' and applied a filter for orders exceeding $800.

Implemented data transformation by renaming key sales fields to 'order date', 'order ID', and 'order amount.'

Organized and stored the transformed dataset in the 'customer-sales' folder within the specified S3 bucket.

# Result from Glue Job Stored in customersales folder in s3 bucket

# Leveraged Athena to query the combined dataset stored in Parquet format.