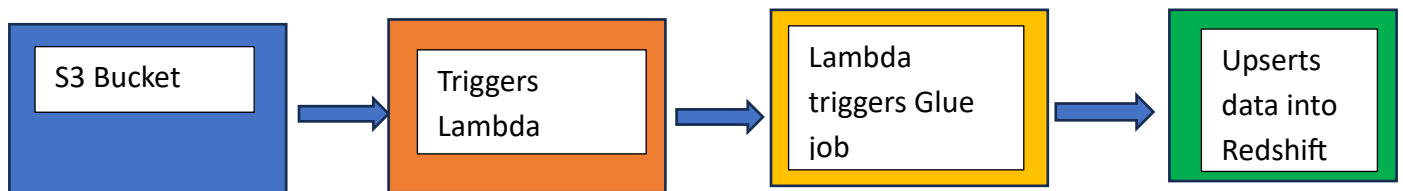


ECOMMERCE DATA PIPELINE USING S3, GLUE, LAMBDA, REDSHIFT

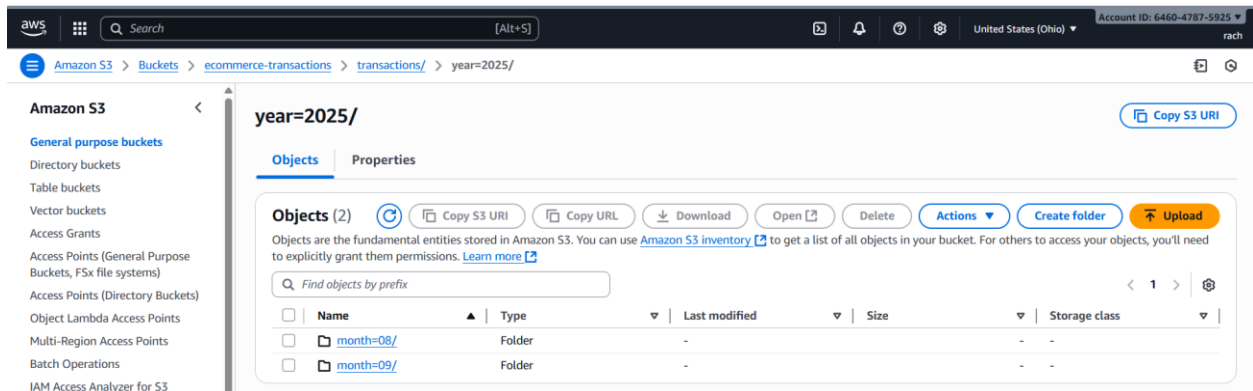
This project involves building a sophisticated event-driven data ingestion and transformation pipeline focusing on e-commerce transactional data. I have designed a system using AWS services such as S3, Lambda, Glue, Redshift, to ingest, transform, validate, and upsert data into Amazon Redshift for analytical purposes.

Architecture diagram



- **Mock Data Generation- script attached**

- Transaction Data: Generate daily transaction files in CSV format, stored using the following hive-style partitioning in S3:



- **Dimension Tables and Sample Records- Pre-load dim_customers, dim_products dimension tables into Redshift as part of the setup process. Setup Redshift cluster first.**

Amazon Redshift Serverless > Workgroup configuration > ecommerce-redshift

This workgroup is associated with a namespace
To manage database objects and users, navigate to the namespace that this workgroup is associated with.

ecom

ecommerce-redshift [info](#) Actions Query data

General information

Workgroup ecommerce-redshift	Date created August 30, 2025, 18:14 (UTC-04:00)	Endpoint ecommerce-redshift.646047875925.us-east-2.redshift-serverless.amazonaws.com:5439/dev
Namespace ecom	Status Available	JDBC URL jdbc:redshift://ecommerce-redshift.646047875925.us-east-2.redshift-serverless.amazonaws.com:5439/dev
Workgroup ARN arn:aws:redshift-serverless:us-east-2:646047875925:workgroup/efc2af71-590b-450e-8ea8-adaa465c3754	Base capacity 8 RPUs	ODBC URL Driver={Amazon Redshift (x64)}; Server=ecommerce-redshift.646047875925.us-east-2.redshift-serverless.amazonaws.com; Database=dev
Workgroup version 1.0.121035 🔗	Custom domain name -	
	Patch version Patch 192 🔗	
	Track - new	

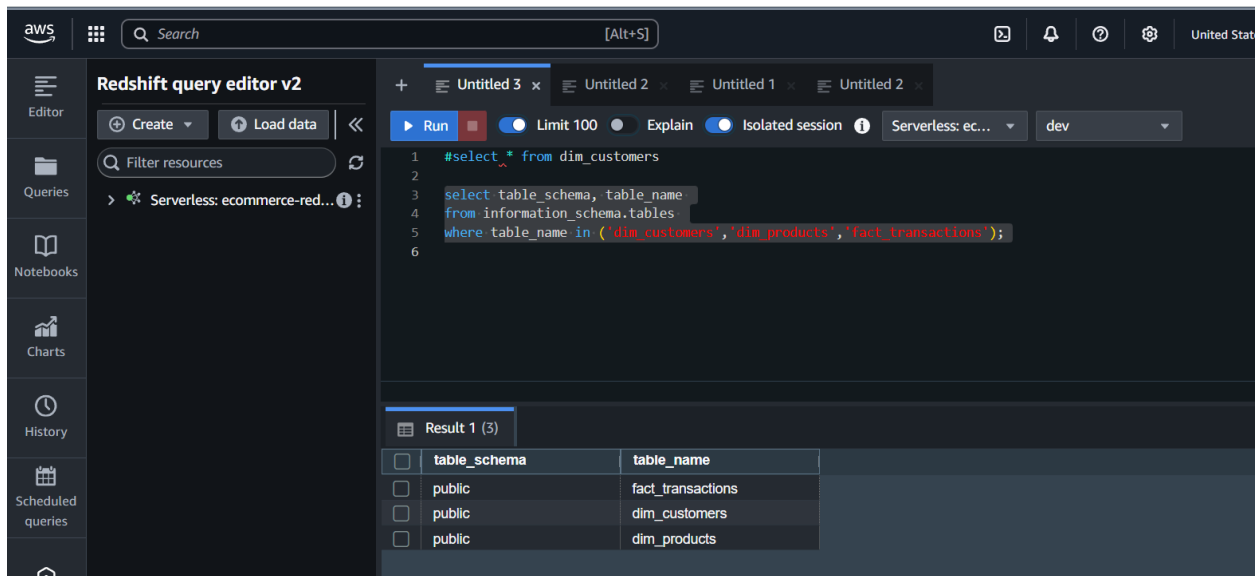
aws Search [Alt+S]

Redshift query editor v2

Create Load data Filter resources Serverless: ecommerce-red...

Run Limit 100 Explain Isolated session Serverless: ec... dev

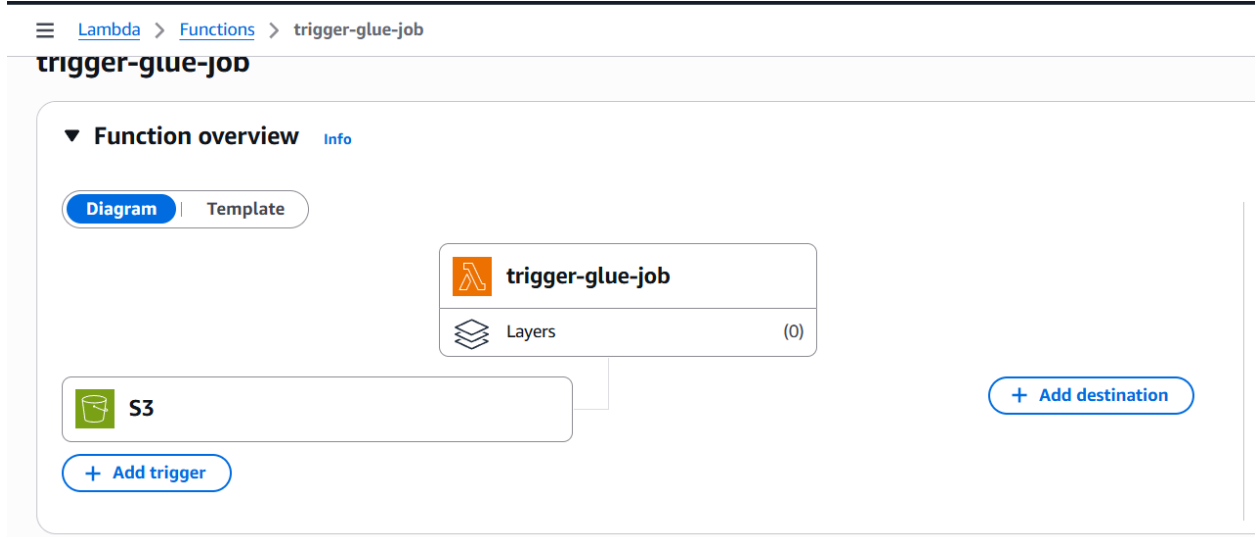
```
1 -- Customers dimension table
2 CREATE TABLE dim_customers (
3     customer_id VARCHAR(20) PRIMARY KEY,
4     first_name VARCHAR(50),
5     last_name VARCHAR(50),
6     email VARCHAR(100),
7     membership_level VARCHAR(20)
8 );
9
10 -- Products dimension table
11 CREATE TABLE dim_products (
12     product_id VARCHAR(20) PRIMARY KEY,
13     product_name VARCHAR(100),
14     category VARCHAR(50),
15     price DECIMAL(10,2),
16     supplier_id VARCHAR(20)
17 );
18
19 -- Transactions fact table
20 CREATE TABLE fact_transactions (
21     transaction_id VARCHAR(50) PRIMARY KEY,
22     customer_id VARCHAR(20),
23     product_id VARCHAR(20),
24     quantity INTEGER,
25     price DECIMAL(10,2),
26     total_amount DECIMAL(10,2),
27     transaction_date DATE,
28     payment_type VARCHAR(50),
```



- **Data Ingestion and Transformation with AWS Glue**

- **Event-Driven Ingestion: Configure an AWS Lambda function to trigger AWS Glue**

jobs upon detecting new files in the S3 transactions folder.



- **Data Transformation and Validation using Glue ETL**

Establish Glue connection to Redshift by creating VPC endpoint

Redshift-managed VPC endpoints [Info](#)
(1)
Endpoints to access serverless endpoints that are in another VPC or subnet.

< 1 > ⚙️

Endpoint name	Endpoint ARN	Status	Account ID
redshift-glue-endpoint	arn:aws:redshift-serverless:us-east-2:646047875925:redshift-glue-endpoint-endpoint-eyijoybjkinrns2znbau.646047875925.us-east-2:redshift-serverless.amazonaws.com	Active	646047875925

redshift-glue-endpoint [Edit](#) [Delete](#)

Settings for Redshift-managed VPC endpoint [Refresh](#)

Account ID 646047875925	Virtual private cloud (VPC) vpc-09c253a589e6bb474	Network interface ID eni-03647d79555a736f1	Endpoint URL redshift-glue-endpoint-endpoint-eyijoybjkinrns2znbau.646047875925.us-east-2:redshift-serverless.amazonaws.com
Date created September 01, 2025, 00:30 (UTC-04:00)	VPC endpoint ID vpce-088513fb26919429a	Subnet ID subnet-07ea4c18bcd4773f1	JDBC URL jdbc:redshift://redshift-glue-endpoint-endpoint-eyijoybjkinrns2znbau.646047875925.us-east-2:redshift-serverless.amazonaws.com:5439/
Status Active	VPC security group sg-08c406610285a9cab	Private IP address 172.31.37.143	ODBC URL Driver={Amazon Redshift (x64)}; Server=redshift-glue-endpoint-endpoint-eyijoybjkinrns2znbau.646047875925.us-east-2:redshift-serverless.amazonaws.com;
Endpoint ARN arn:aws:redshift-serverless:us-east-2:646047875925:managedvpcendpoint/ae2f66a3-49ac-43f4-918c-0b3950ebbf5c	Subnet group subnet-07ea4c18bcd4773f1, subnet-084ff7bb17173d5a4, subnet-0eec3d82b6e636dd6	Availability Zone us-east-2c	

- **Join Operations:** Enrich transactional data by joining with dim_products and dim_customers based on product_id and customer_id.
- **Data Validation:** Include validation logic in the Glue job to filter out transactions with invalid customer_id or product_id (e.g., missing in dimension tables).
- **Additional Transformations:** Calculate the total transaction amount (quantity * price) and categorize transactions into different classes based on the amount (e.g., "Small", "Medium", "Large").
- **Upsert Operation in Amazon Redshift**
 - Design the Glue job to perform an upsert operation into the fact_transactions table in Redshift, using transaction_id as the key. Consider transaction date and status when determining if an existing record should be updated.



ecommerce-ettl-job

Last modified on 9/1/2025, 5:07:08

[Script](#) | [Job details](#) | [Runs](#) | [Data quality](#) | [Schedules](#) | [Version Control](#)

Script [Info](#)

```
1 import sys
2 from awsglue.transforms import *
3 from awsglue.utils import getResolvedOptions
4 from pyspark.context import SparkContext
5 from awsglue.context import GlueContext
6 from awsglue.job import Job
7 from pyspark.sql import functions as F
8 from awsglue.dynamicframe import DynamicFrame
9
10 # Get job parameters
11 args = getResolvedOptions(sys.argv, ['JOB_NAME'])
12 sc = SparkContext()
13 glueContext = GlueContext(sc)
14 spark = glueContext.spark_session
15 job = Job(glueContext)
16 job.init(args['JOB_NAME'], args)
17
18 print("Starting ETL job...")
```

Python | Ln 1, Col 1 | Errors: 0 | Warnings: 0