# Serverless Data Pipeline

**Data Pipeline :**

The purpose of this practical is to design and implement a data pipeline on a server that can collect data from different sources, process it, transform it, and finally load it into a storage system such as a database, data warehouse, or cloud storage. This involves using server-side scripting, automation, logging, and error-handling techniques.
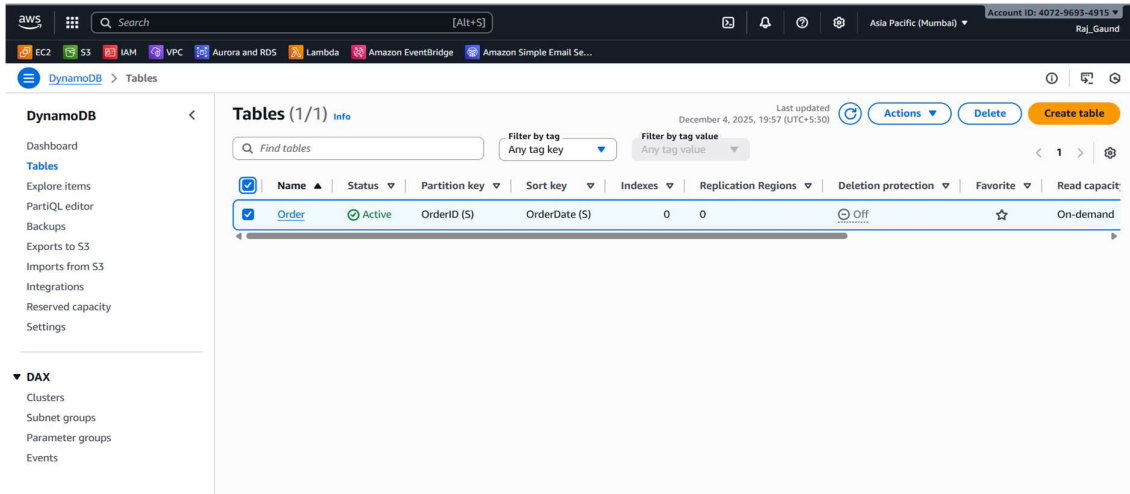
**Key Steps of the Practical -**

- Store orders in **DynamoDB**
- Use **Glue** to crawl + transform data
- Store cleaned data in **S3 (Data Lake)**
- Use **Athena** to run SQL queries on the data

<u>**Steps for Practical**</u>

**Step 1 :  Create DynamoDB Table**

- Open DynamoDB → Create Table

- Table name: Orders

- Partition Key: OrderID (String)

- Sort Key: OrderDate (String)  →  Create Table

**A new DynamoDB table is created to store order records.**

## Step 2: Insert Sample Data

- Click on Table ( Order ) → Explore Item

- Create Item → Add Attribute and Data → Create

- Add at Least five items



**DynamoDB now contains real order data that the pipeline will process.**

**Step 3: Create Glue Crawler**

- AWS Console → **Glue** → Data Catalouge

- Create Crawler → Enter Crawler Name ( OrderCrawler )

- Add Datasource → Data source : **DynamoDB**

- Table name : DynamoDB Table Name (**Order**) → Add

- Choose IAM Role : Create New IAM role

- AWS → IAM → Role → Create → Name : GlueCrawlerRole

- Trusted Entity : AWS service → Use case : **Glue**

- Add Policies / Permission : **AWSGlueServiceRole** ,
  **AmazonDynamoDBReadOnlyAccess , AmazonS3FullAccess**

- Create IAM Role

- Add Database : **ordersdb** → OnDemand → Create Crawler
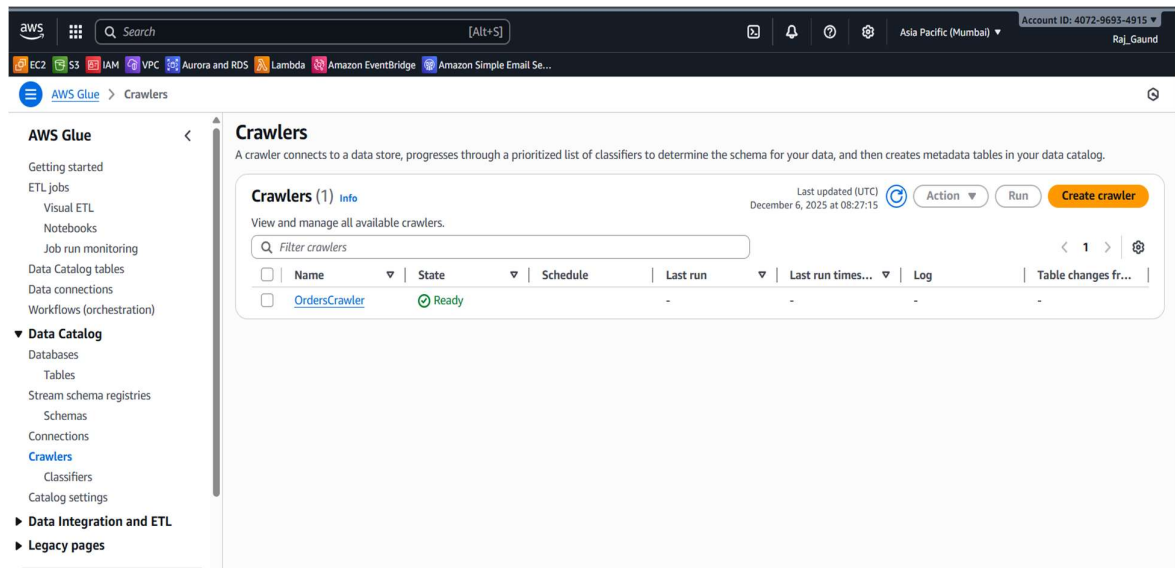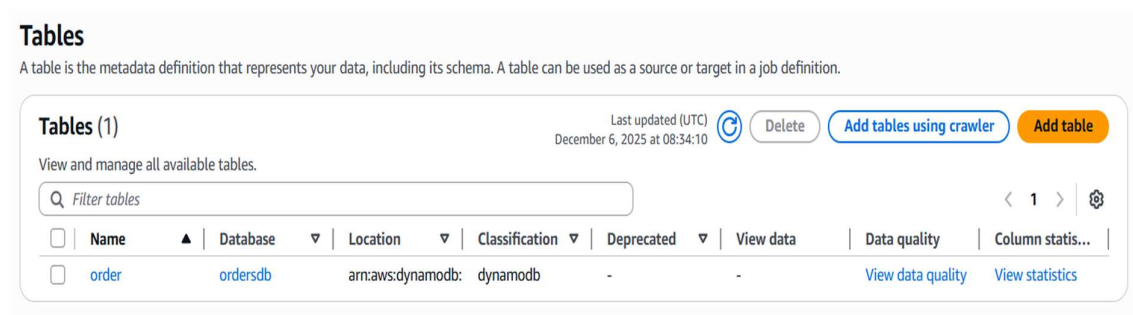


**IAM Role Created Successfully**

**Crawler create Successfully**

## Step 4 : Run The Crawler

- AWS → Glue → Crawler

- Click Crawler → **Run Crawler**

- AWS → Glue → Table



**( Table created automatic when you run the crawler )**

**Step 5 : Create S3 bucket ( for ETL output )**

- AWS Console → **S3** → General Purpose Bucket

- Bucket Name : you can enter own (my-orders-bucket123)

- Create Bucket



**Bucket used to stored Job ETL output**

**Step 6 : Glue ETL Job create**

- AWS → Glue → ETL jobs → Script editor

- Engine : Spark → Create script

- Type this code

## Code :

```python
import sys
from awsglue.transforms import Filter
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

# required argument JOB_NAME
args = getResolvedOptions(sys.argv, ['JOB_NAME'])
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)

# ----- CONFIG: change these if you used other names -----
GLUE_DATABASE = "ordersdb"          # Glue DB created by crawler
GLUE_TABLE = "order"                # Glue table name (crawler output)
S3_OUTPUT_PATH = "s3://my-orders-bucket123/shipped/"  # change to your bucket
```

```python
datasource = glueContext.create_dynamic_frame.from_catalog(
    database=GLUE_DATABASE,
    table_name=GLUE_TABLE
)

# Optional: print schema to logs (helpful)
print("Schema: ", datasource.schema())

# Filter function: keep only Status == 'Shipped'
def is_shipped(rec):
    status = rec.get('Status') or rec.get('status') or rec.get('STATUS')
    if status is None:
        return False
    try:
        return str(status).lower() == 'shipped'
    except:
        return False


filtered = Filter.apply(frame=datasource, f=is_shipped)

# Write filtered data to S3 as Parquet
glueContext.write_dynamic_frame.from_options(
    frame=filtered,
    connection_type="s3",
    connection_options={"path": S3_OUTPUT_PATH},
    format="parquet"
)

job.commit()
```

- Jobs details → Enter name : OrdersETLJob

- IAM Role : Choose existing (GlueOrdersRole)

- Glue version: Glue 3.0 (default) / Engine = Spark (keep)

- Language : Python → Save → Run Job

- Click on OrdersETLJob → job details → Run



## If succeed show means

- Script is correct
- DynamoDB read that Data
- Filter is successfully applied

**Step 7 : Open S3 bucket**

- AWS → S3 → Open Bucket

- Example : (my-orders-bucket8483)



- Clicked shipped named folder

- Under shipped folder list of four items can show



## Step 8 : Athena Setup to run query

- AWS → Athena

- Select : Query your data in Athena console

- Launch Query editor

-  Click on setting → edit your S3 location

- **s3://my-orders-bucket123/athena-results/**  → Save

- perform query to create stored and view

1) Query :

```sql
SHOW DATABASES;
```

 Its shows  : **ordersdb**

- left side Database → select **ordersdb** into dropdown

2) Run this query to create table :

```sql
CREATE EXTERNAL TABLE IF NOT EXISTS ordersdb.orders_shipped (
  OrderID string,
  OrderDate string,
  Customer string,
  Amount bigint,
  Status string
)
STORED AS PARQUET
LOCATION 's3://my-orders-bucket8483/shipped/';
```

| Query results | Query stats |
| --- | --- |

| ✓ Completed | Time in queue: 37 ms | Run time: 387 ms | Data scanned: - |
| --- | --- | --- | --- |

Query successful.

3) Query To view or fetch data :

```sql
SELECT * FROM ordersdb.orders_shipped
LIMIT 50;
```

**Query results**    **Query stats**

⊘ Completed                         Time in queue: 90 ms    Run time: 426 ms    Data scanned: 1.03 KB

**Results** (3)                                    ⎘ Copy      **Download results CSV**

🔍 *Search rows*                                                    ‹ 1 ›   ⚙

| # ▽ | orderid ▽ | orderdate ▽ | customer ▽ | amount ▽ | status ▽ |
|-----|-----------|-------------|------------|----------|----------|
| 1 | 1 | 20-04-2004 | deep | 1000 | shipped |
| 2 | 3 | 22-04-2004 | azad | 3000 | shipped |
| 3 | 5 | 25-04-2004 | gopi | 5000 | shipped |

4) Query To rows count

```sql
SELECT COUNT(*) AS total_orders FROM ordersdb.orders_shipped;
```

⊘ Completed                         Time in queue: 57 ms    Run time: 356 ms    Data scanned: -

**Results** (1)                                    ⎘ Copy      **Download results CSV**

🔍 *Search rows*                                                    ‹ 1 ›   ⚙

| # ▽ | total_orders ▽ |
|-----|----------------|
| 1 | 3 |

5) Query to Total amount (sum) overall

```sql
SELECT SUM(Amount) AS total_revenue FROM ordersdb.orders_shipped;
```

| ⊘ Completed | Time in queue: 55 ms | Run time: 445 ms | Data scanned: 0.23 KB |

**Results** (1)                                      Copy    Download results CSV

🔍 Search rows                                                      < 1 >  ⚙

| # | ▽ | total_revenue | ▽ |
|---|---|---------------|---|
| 1 |   | 9000          |   |

6) Revenue by customer (top customers)

```sql
SELECT Customer, SUM(Amount) AS total_spent
FROM ordersdb.orders_shipped
GROUP BY Customer
ORDER BY total_spent DESC
LIMIT 20;
```

| ⊘ Completed | Time in queue: 55 ms | Run time: 350 ms | Data scanned: 0.40 KB |

**Results** (3)                                      Copy    Download results CSV

🔍 Search rows                                                      < 1 >  ⚙

| # | ▽ | Customer | ▽ | total_spent | ▽ |
|---|---|----------|---|-------------|---|
| 1 |   | gopi     |   | 5000        |   |
| 2 |   | azad     |   | 3000        |   |
| 3 |   | deep     |   | 1000        |   |

7) Orders per day (if OrderDate is YYYY-MM-DD)

```sql
SELECT OrderDate, COUNT(*) AS orders_count, SUM(Amount) AS daily_revenue
FROM ordersdb.orders_shipped
GROUP BY OrderDate
ORDER BY OrderDate;
```

| # | OrderDate | orders_count | daily_revenue |
|---|-----------|--------------|---------------|
| 1 | 20-04-2004 | 1 | 1000 |
| 2 | 22-04-2004 | 1 | 3000 |
| 3 | 25-04-2004 | 1 | 5000 |

Final output :



- **Purpose :**

  The purpose of this project is to build a fully automated and serverless data processing pipeline using AWS DynamoDB, AWS Glue, Amazon S3, and Amazon Athena. The aim is to extract raw order data from DynamoDB, automatically detect and catalog its schema using a Glue Crawler, transform and clean the data through a Glue ETL job, and store the optimized output in S3 using the Parquet format.

- **Conclusion :**

In conclusion, this project successfully demonstrates how AWS Glue and Amazon Athena can work together to create an end-to-end serverless data analytics solution. The pipeline extracted data from DynamoDB, transformed it using a Glue ETL job, stored it in an optimized Parquet format in S3, and allowed quick querying using Athena.

❖ Architecture Diagram :

**Serverless Data Pipeline (DynamoDB , Glue, S3 , Athena)**