**Shaikh Ateeb Ahmed**
**08-09-2025**

## Serverless Data Pipeline

**1) Create the Orders DynamoDB table**

1. Sign in to the AWS Management Console

2. Open **Services → DynamoDB**.

3. Click **Create table**.

4. Under **Table name**, enter Orders.

5. Under **Primary key**:

**Partition key**: OrderID → select **String**.

Click **Add sort key** and enter **Sort key**: OrderDate → **String**.

6. Under **Table settings → Capacity mode**, select **On-demand (Pay per request)**.

7. Click **Create table**.



Step 2: Insert Sample Data

**1) Quick note before you start**

Your table has keys: OrderID (PK, String) and OrderDate (SK, String). **Every item must include both**.

**Using ISO-8601 datetimes (e.g. 2025-09-01T09:10:00Z) for OrderDate so lexicographic sorting works for ranges.**

**2) Sample dataset (15 orders)**

```
{
  "OrderID": { "S": "O1001" },
  "OrderDate": { "S": "2025-09-01T09:10:00Z" },
  "Customer": { "S": "Alice" },
  "Amount": { "N": "250" },
  "Status": { "S": "Shipped" }
}
{
  "OrderID": { "S": "O1002" },
  "OrderDate": { "S": "2025-09-01T10:24:00Z" },
  "Customer": { "S": "Bob" },
  "Amount": { "N": "400" },
  "Status": { "S": "Pending" }
}
{
  "OrderID": { "S": "O1003" },
  "OrderDate": { "S": "2025-09-02T11:30:00Z" },
  "Customer": { "S": "Charlie" },
  "Amount": { "N": "125" },
  "Status": { "S": "Delivered" }
}
{
  "OrderID": { "S": "O1004" },
  "OrderDate": { "S": "2025-09-03T14:05:00Z" },
  "Customer": { "S": "Diana" },
```

  "Amount": { "N": "78.5" },

  "Status": { "S": "Processing" }

}

{

  "OrderID": { "S": "O1005" },

  "OrderDate": { "S": "2025-09-03T16:50:00Z" },

  "Customer": { "S": "Eve" },

  "Amount": { "N": "560" },

  "Status": { "S": "Shipped" }

}

{

  "OrderID": { "S": "O1006" },

  "OrderDate": { "S": "2025-09-04T08:20:00Z" },

  "Customer": { "S": "Frank" },

  "Amount": { "N": "30" },

  "Status": { "S": "Cancelled" }

}

{

  "OrderID": { "S": "O1007" },

  "OrderDate": { "S": "2025-09-04T09:15:00Z" },

  "Customer": { "S": "Grace" },

  "Amount": { "N": "210" },

  "Status": { "S": "Pending" }

}

{

  "OrderID": { "S": "O1008" },

  "OrderDate": { "S": "2025-09-05T12:00:00Z" },

  "Customer": { "S": "Heidi" },

  "Amount": { "N": "999.99" },

  "Status": { "S": "Shipped" }

}

{

  "OrderID": { "S": "O1009" },

  "OrderDate": { "S": "2025-09-05T13:45:00Z" },

  "Customer": { "S": "Ivan" },

  "Amount": { "N": "150" },

  "Status": { "S": "Delivered" }

}

{

  "OrderID": { "S": "O1010" },

  "OrderDate": { "S": "2025-09-06T15:00:00Z" },

  "Customer": { "S": "Judy" },

  "Amount": { "N": "49.99" },

  "Status": { "S": "Processing" }

}

{

  "OrderID": { "S": "O1011" },

  "OrderDate": { "S": "2025-09-06T16:30:00Z" },

  "Customer": { "S": "Ken" },

  "Amount": { "N": "320" },

  "Status": { "S": "Shipped" }

}

{

  "OrderID": { "S": "O1012" },

  "OrderDate": { "S": "2025-09-07T10:00:00Z" },

  "Customer": { "S": "Leo" },

```
  "Amount": { "N": "215.5" },

  "Status": { "S": "Returned" }

}

{

  "OrderID": { "S": "O1013" },

  "OrderDate": { "S": "2025-09-07T11:11:00Z" },

  "Customer": { "S": "Mallory" },

  "Amount": { "N": "700" },

  "Status": { "S": "Pending" }

}

{

  "OrderID": { "S": "O1014" },

  "OrderDate": { "S": "2025-09-08T18:00:00Z" },

  "Customer": { "S": "Niaj" },

  "Amount": { "N": "1200" },

  "Status": { "S": "Delivered" }

}

{

  "OrderID": { "S": "O1015" },

  "OrderDate": { "S": "2025-09-09T09:00:00Z" },

  "Customer": { "S": "Olivia" },

  "Amount": { "N": "15.75" },

  "Status": { "S": "Processing" }

}
```
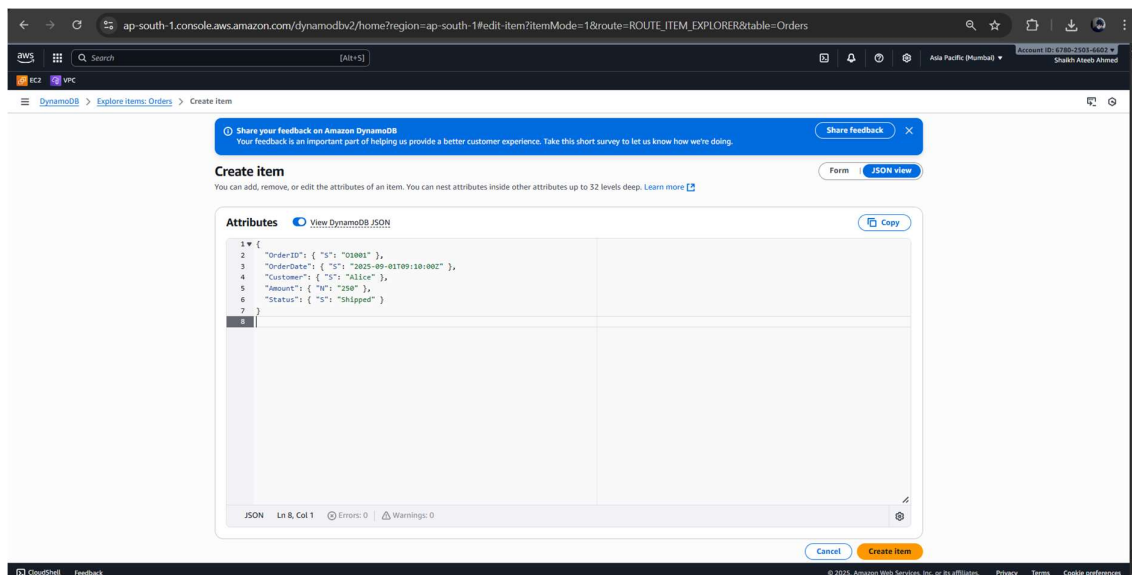
**3) Manual insert via AWS Console**

1.  Console → Services → **DynamoDB** → **Tables** → click Orders.

2.  Click **Explore items** (or **Items**), then **Create item**.

3.  Switch to **JSON** view (easier for full items).

4. Paste one item from the dataset (document format, not typed AttributeValue).

Example:

{

  "OrderID": { "S": "O1001" },

  "OrderDate": { "S": "2025-09-01T09:10:00Z" },

  "Customer": { "S": "Alice" },

  "Amount": { "N": "250" },

  "Status": { "S": "Shipped" }
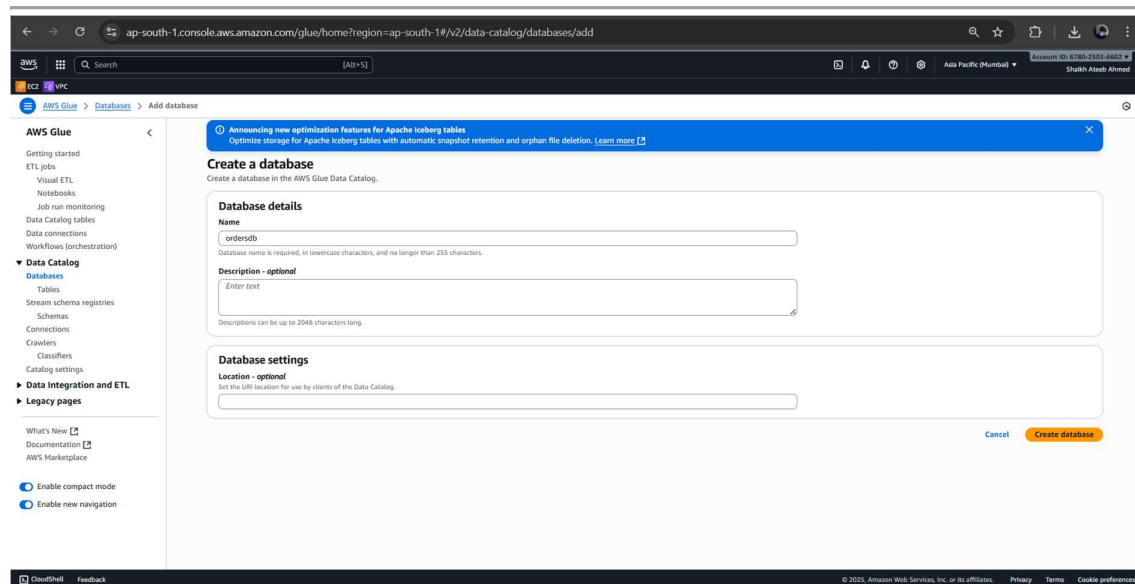
}



5. Click **Create**. Repeat for other items.

Console is manual but handy for quick checks.

---

**Step 3: Create an AWS Glue Crawler for the Orders DynamoDB table**

**1) Create Glue Database (OrdersDB)**

1. Open **AWS Console → AWS Glue → Data Catalog → Databases**.

2. Click **Add database**.

3. Name: ordersdb. (Optionally add description and location).

4. Click **Create**.

**2) Create IAM Role for Glue**

1. Open **IAM → Roles → Create role**.

2. Select **Glue** as the trusted service (choose **Glue** so the trust relationship glue.amazonaws.com is set).



3. Attach policies:

**AmazonDynamoDBReadOnlyAccess** (or a custom restricted DynamoDB policy — see example below).

---

**AmazonS3ReadOnlyAccess** *or* a custom S3 policy granting access to a specific S3 temp bucket if you use S3.

**AWSGlueServiceRole** or any AWS-managed Glue service policy if shown.

**AmazonS3FullAccess**

**DynamoDB read: DescribeTable, Scan, GetItem, Query for the Orders table ARN.**

**S3 access to your bucket: GetObject, PutObject, ListBucket for arn:aws:s3:::my-orders-analytics01.**

4.   Give the role a name such as AWSGlueOrdersRole and finish.

**Minimal custom trust policy** (if using CLI):

```
{
  "Version":"2012-10-17",
  "Statement":[
   {
     "Effect":"Allow",
     "Principal":{"Service":"glue.amazonaws.com"},
     "Action":"sts:AssumeRole"
   }
  ]
}
```

**3)  Create the Crawler**

1.  Open **AWS Console → AWS Glue → Crawlers**.

2.  Click **Add crawler** (or **Create crawler**).

3.  **Name**: OrdersCrawler → Next.



4.  **Data source**: Choose **Add a data store** → Select **DynamoDB**.

Select the table Orders from the list (region must match).

Click **Add** → Next.



5. **Choose IAM Role**: Select the role you created (AWSGlueOrdersRole).



6. **Crawler output**:

Choose **Output to a Data Catalog database**, Database: **OrdersDB**.

7. **Configure crawler options**:
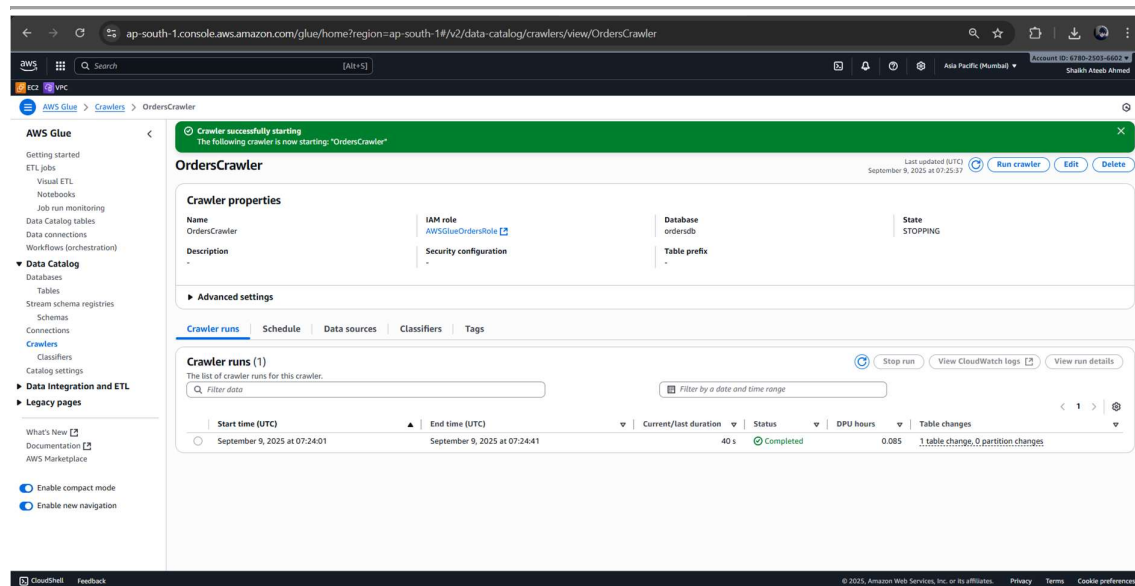
Frequency: **Run on demand**.

8.  Review and **Finish**.



9.  After creation, select the crawler and click **Run crawler**.

**Shaikh Ateeb Ahmed**
**08-09-2025**

**Step 4 — Create the AWS Glue ETL Job**

**1) Create the S3 target bucket**

**S3 → Create bucket → name my-orders-analytics (or your preferred name) → choose Region → Create.**

**Shaikh Ateeb Ahmed**
**08-09-2025**

**2) Create / verify the IAM role for Glue jobs**

**Create role AWSGlueOrdersETLRole and attach these managed policies:**

**AWSGlueServiceRole (or similar Glue-managed role)**

**AmazonDynamoDBReadOnlyAccess**

**AmazonS3FullAccess (or tightly-scoped S3 policy to your bucket)**

**DynamoDB read: DescribeTable, Scan, GetItem, Query for the Orders table ARN.**

**S3 access to your bucket: GetObject, PutObject, ListBucket for arn:aws:s3:::my-orders-analytics01.**

**Shaikh Ateeb Ahmed**
**08-09-2025**

**Minimal trust policy (Glue service principal):**

```
{

  "Version":"2012-10-17",

  "Statement":[

   {

     "Effect":"Allow",

     "Principal":{"Service":"glue.amazonaws.com"},

     "Action":"sts:AssumeRole"

   }

  ]

}
```



**3) Create the Glue Job**

1.  **Open AWS Console → AWS Glue → Jobs → Add job → Visual.**

2.  **Name: AWS Glue Data Catalog**

3.  **Database: ordersdb**

4. **Table: orders**

5. **IAM role: select AWSGlueOrdersETLRole (created above).**

**Create Glue job:**

**AWS Console → AWS Glue → Jobs → Add job.**

**Name: OrdersToParquet (or whatever you prefer).**

**IAM role: AWSGlueOrdersETLRole.**

**Type: Spark, Glue version: 3.0 or 4.0, Python 3.**

**This job runs: A new script authored by you OR point to an S3 script (see below).**

**TempDir: s3://my-orders-analytics01/temp/ (set as --TempDir in job properties).**

**Worker type / number: start small (Standard, 2 workers).**

**Paste the script (below) into the job script editor (or upload to S3 and supply path).**

**Run the job (Console → select job → Run) or via CLI:**

**aws glue start-job-run --job-name OrdersToParquet**

**Shaikh Ateeb Ahmed**
**08-09-2025**

aws | Q Search [Alt+S]

EC2 VPC

≡

## OrdersToParquet ✎

Script | **Job details** | Runs | Data quality | Schedules | Version Control

**Job bookmark** | Info
Specifies how AWS Glue processes job bookmark when the job runs. It can remember previously processed data (Enable), update state information (Pause), or ignore state information (Disable).

Disable ▼

**Job Run Queuing**
☐ Enable job runs to be queued to run later when they cannot run immediately due to service quotas

**Flex execution** | Info
☐ Reduce costs by running this job on spare capacity. Ideal for non-urgent workloads that don't require fast jobs start times or consistent execution times. See recommendations, limitations and pricing in the help panel by clicking on the Info link above.

**Number of retries**

0

**Job timeout (minutes)**
Set the maximum execution time. The default is 480 minutes (8 hours) for a Glue 5.0 ETL job and 2,880 minutes (48 hours) for Glue 4.0 and below. No job timeout is defaulted for a Glue Streaming job.

480

**Usage profile**

-

### ▼ Advanced properties

**Script filename**

OrdersToParquet.py

**Script path**
S3 location of the script. Path must be in the form s3://bucket/prefix/path/. It must end with a slash (/) and not include any files.

🔍 s3://my-orders-analytics01/scripts/    ✕ | View ⬈ | Browse S3

**Job metrics** | Info
☑ Enable the creation of CloudWatch metrics when this job runs.

**Job observability metrics** | Info
☑ Enable the creation of additional observability CloudWatch metrics when this job runs.

**Continuous logging** | Info
☐ Enable logs in CloudWatch.

**Spark UI** | Info
☑ Write Spark UI logs to Amazon S3.

**Spark UI logs path**

🔍 s3://aws-glue-assets-678025036602-ap-south-1/sparkHistoryLogs/    ✕ | View ⬈ | Browse S3

CloudShell    Feedback

aws | ::: | Q Search [Alt+S]

EC2 | VPC

≡

## OrdersToParquet ✎

**Script** | **Job details** | Runs | **Data quality** | Schedules | **Version Control**

**Spark UI logs path**

Q s3://my-orders-analytics01/sparkHistoryLogs/ | ✕ | View ↗ | Browse S3

**Spark UI logging and monitoring configuration**

⦿ Standard - *default*
Write logs using the Glue job run id as the filename. Enable Spark UI monitoring in Glue console.

○ Legacy
Write logs using 'spark-application-(timestamp)' as the filename. Do not enable Spark UI monitoring in Glue console.

○ Standard and legacy
Write logs to both the standard and legacy locations. Enable Spark UI monitoring in Glue console.

**Maximum concurrency**
Sets the maximum number of concurrent runs that are allowed for this job. An error is returned when this threshold is reached.

| 1 |

**Temporary path**
Working directory. Path must be in the form s3://bucket/prefix/path/. It must end with a slash (/) and not include any files.

Q s3://my-orders-analytics01/temp/ | ✕ | View ↗ | Browse S3

**Delay notification threshold (minutes)** | Info
Set a delay threshold in minutes. If the job runs longer than the specified time Glue will send a delay notification via CloudWatch.

| |

**Security configuration**
Defines the encryption options of the ETL job. The security configuration specifies how the script is encrypted using server-side encryption with AWS KMS-managed keys (SSE-KMS) or Amazon S3-managed encryption keys (SSE-S3).

| None ▼ |

**Server-side encryption**

☐ Select this option to enable S3-SSE encryption on both your data target and any data that is written to an S3 temporary directory. Ignored if a security configuration is specified **Learn More** ↗

**Use Glue data catalog as the Hive metastore**

☑ To use Glue data catalog as the Hive metastore, the IAM role used for the job should have glue:CreateDatabase permissions. A database called "default" is created in the Glue data catalog if it does not exist.

## Connections

**Additional network connections** | Info
Choose a VPC configuration to access Amazon S3 data sources located in your virtual private cloud (VPC). You can create and manage Network connections in **AWS Glue.** ↗
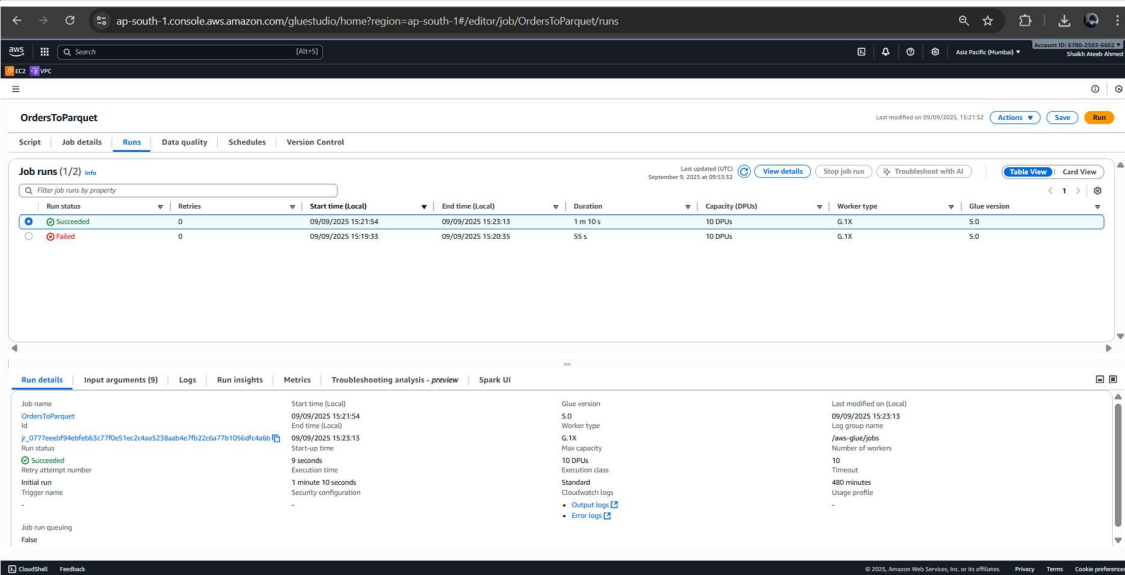
| Choose options ▼ | ↻ |

### Current connections
These are the connections currently associated with the job.

>_ CloudShell    Feedback

**Shaikh Ateeb Ahmed**
**08-09-2025**



---

**Copy-paste-ready Glue PySpark script**

```
import sys

import re

from awsglue.transforms import *

from awsglue.utils import getResolvedOptions

from pyspark.context import SparkContext

from awsglue.context import GlueContext

from awsglue.job import Job

from awsglue.dynamicframe import DynamicFrame

import pyspark.sql.functions as F


# --- JOB ARGS ---

args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()

glueContext = GlueContext(sc)

spark = glueContext.spark_session

job = Job(glueContext)
```

**Shaikh Ateeb Ahmed**
**08-09-2025**

---

```python
job.init(args['JOB_NAME'], args)


# --- CONFIG ---

DATABASE = "ordersdb"

TABLE_NAME = "orders"

OUTPUT_PATH = "s3://my-orders-analytics01/shipped/"


# --- READ from Glue Data Catalog (DynamoDB table) ---

datasource_dyf = glueContext.create_dynamic_frame.from_catalog(

    database=DATABASE,

    table_name=TABLE_NAME,

    transformation_ctx="datasource_dyf"

)


# Convert to Spark DataFrame

df = datasource_dyf.toDF()


# --- FIX Amount column ---
# DynamoDB "choice" type: struct with {double, long}
df = df.withColumn(

    "Amount",

    F.when(F.col("Amount.double").isNotNull(), F.col("Amount.double"))

     .when(F.col("Amount.long").isNotNull(), F.col("Amount.long").cast("double"))

     .otherwise(F.lit(None))

)


# --- FILTER shipped orders (case-insensitive) ---

filtered_df = df.filter(F.lower(F.col("Status")) == "shipped")
```
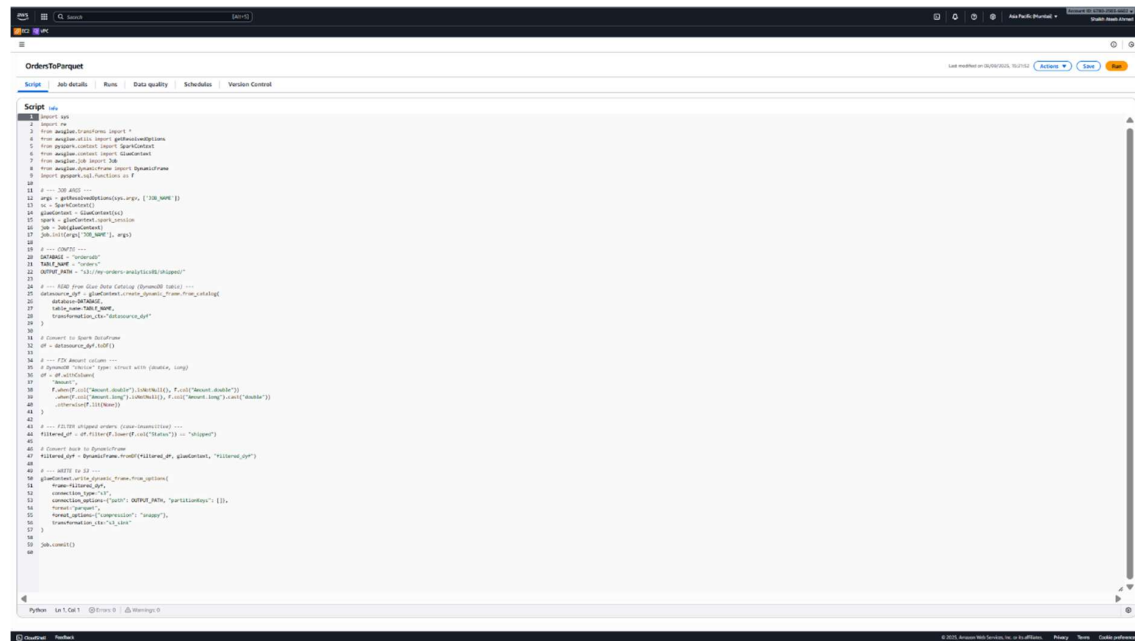
**Shaikh Ateeb Ahmed**
**08-09-2025**

---

# Convert back to DynamicFrame

filtered_dyf = DynamicFrame.fromDF(filtered_df, glueContext, "filtered_dyf")

# --- WRITE to S3 ---

glueContext.write_dynamic_frame.from_options(

   frame=filtered_dyf,

   connection_type="s3",

   connection_options={"path": OUTPUT_PATH, "partitionKeys": []},

   format="parquet",

   format_options={"compression": "snappy"},
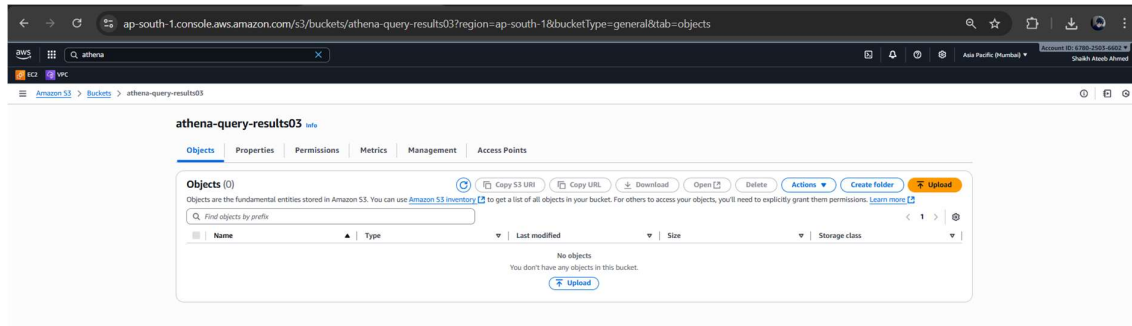
   transformation_ctx="s3_sink"

)

job.commit()

**Shaikh Ateeb Ahmed**
**08-09-2025**

---

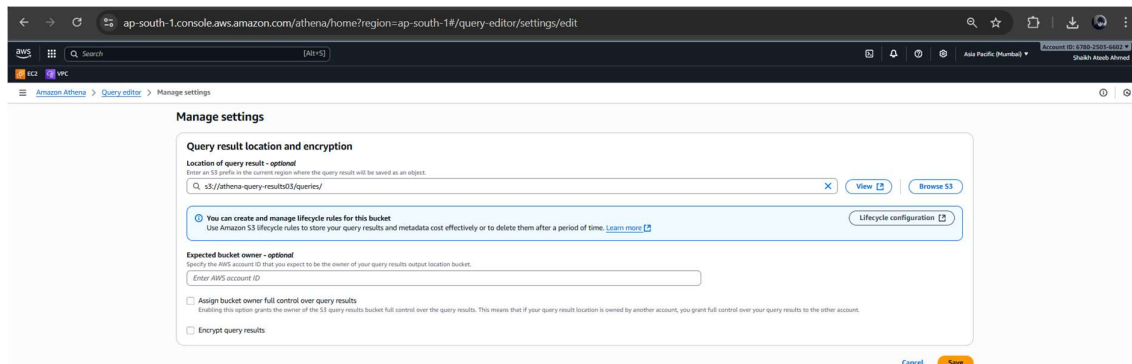**Step 5: Query with Athena**

**Create an S3 bucket for Athena query results**

1.  **Open S3 → Create bucket.**

2.  **Name it e.g. athena-query-results03 (bucket names must be globally unique).**

3.  **Choose same Region as your Glue/Athena region (ap-south-1).**

4.  **Create the bucket.**



---

**Configure Athena to use that S3 location**

1.  **Open AWS Console → Athena (ensure region = same region as Glue & S3).**

2.  **Click the edit Settings.**

3.  **Under Query result location, paste:**

    **s3://athena-query-results03/queries/**

4.  **Save.**



---

**Example queries (use your database/table name)**

**If your table name in Glue/Athena is orders in ordersdb, qualify it as shipped**

**Shaikh Ateeb Ahmed**
**08-09-2025**

---

**Top customers by total spend (Shipped)**

**SELECT Customer, SUM(Amount) AS TotalSpent**

**FROM orders_shipped**

**WHERE lower(Status) = 'shipped'**

**GROUP BY Customer**

**ORDER BY TotalSpent DESC**

**LIMIT 20;**