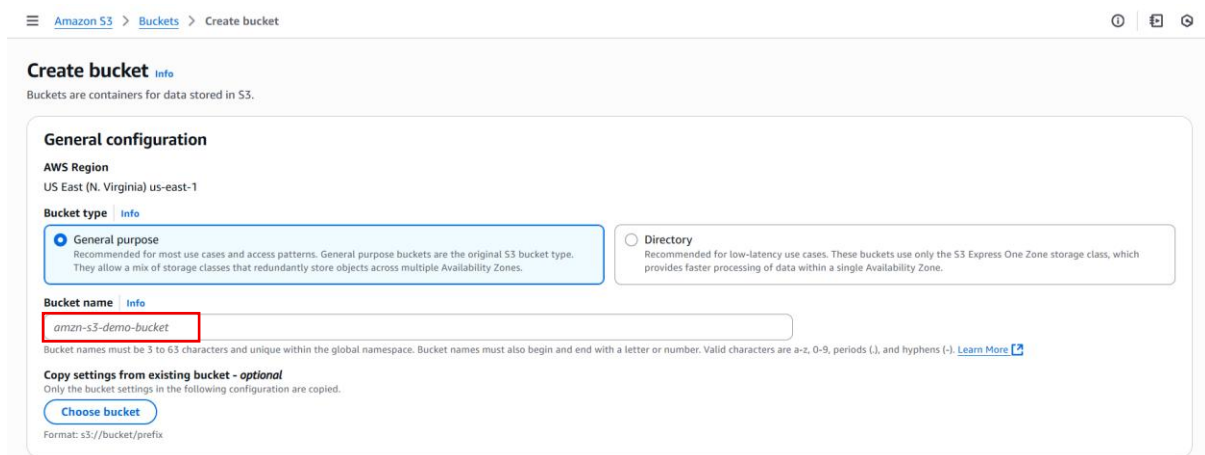# Partitioning with Glue

# To Begin with the Lab

**Summary of the Lab**

In this lab, you learn how to implement **data partitioning using AWS Glue and Athena**. You create an S3 bucket with subfolders (London, New York, Tokyo) representing partitions and upload corresponding sales data. Using AWS Glue, a crawler is configured to detect data and automatically create partitions in the Glue Data Catalog. Then, in Athena, you query the data and update metadata using the **ALTER TABLE ADD PARTITION** command to include new folders like Tokyo.

- **Create an S3 Bucket**
- Go to **AWS Management Console → S3 → Create bucket**.
- Enter a unique name.
- Keep other settings default and **create the bucket**.



- Inside the bucket, create folders for each location:
  - London/
  - NewYork/
- These folders will act as **partition keys**.



- Upload your sales data files into each folder:
- sales_london.csv → into the **London** folder.

- sales_newyork.csv → into the **New York** folder.
- Go to **AWS Glue** → **Crawlers** → **Create crawler**.
- Go to **AWS Console** → **Glue Service** (ensure you're in the same region as your S3 bucket).



- Now Click on the Crawlers and then click on Create crawler.



- Name the Crawler and click next.



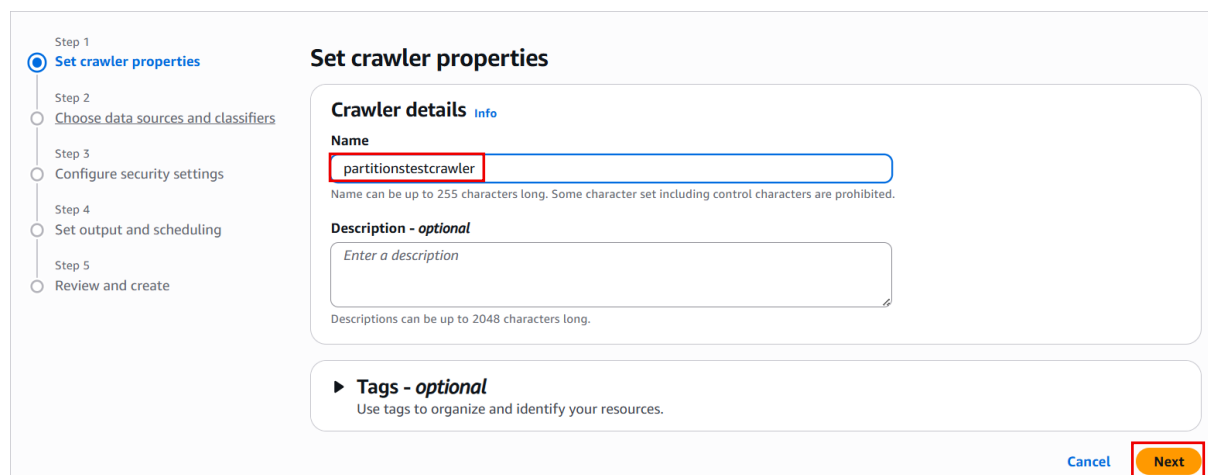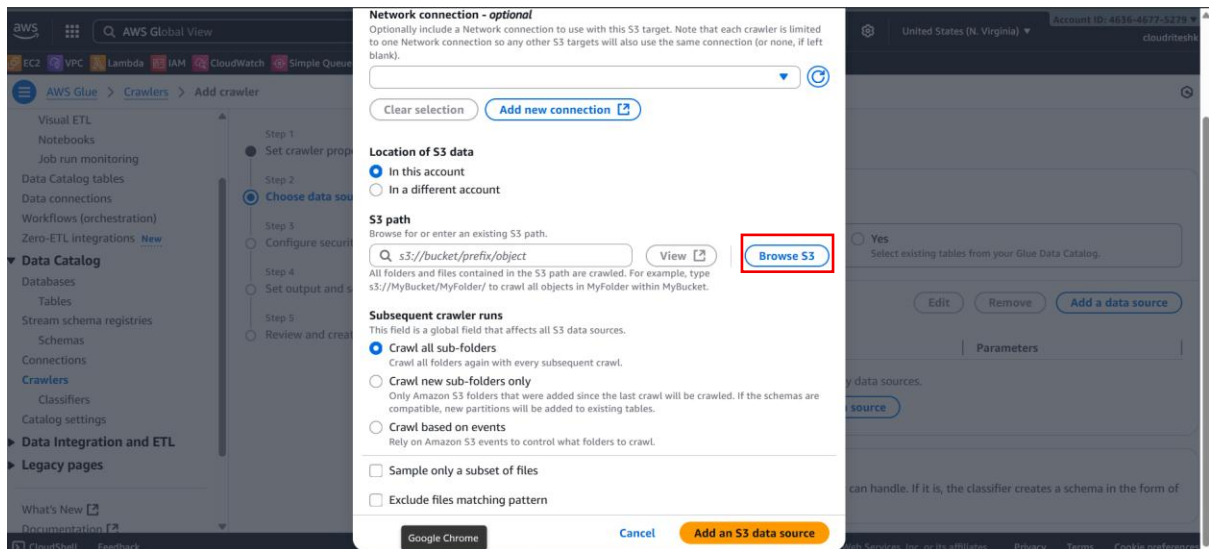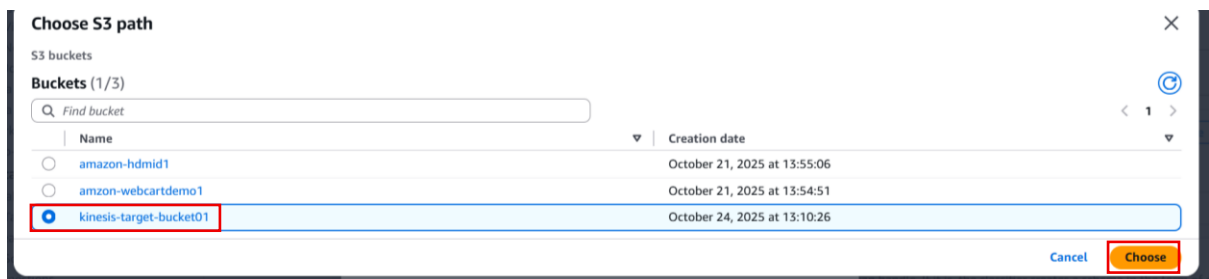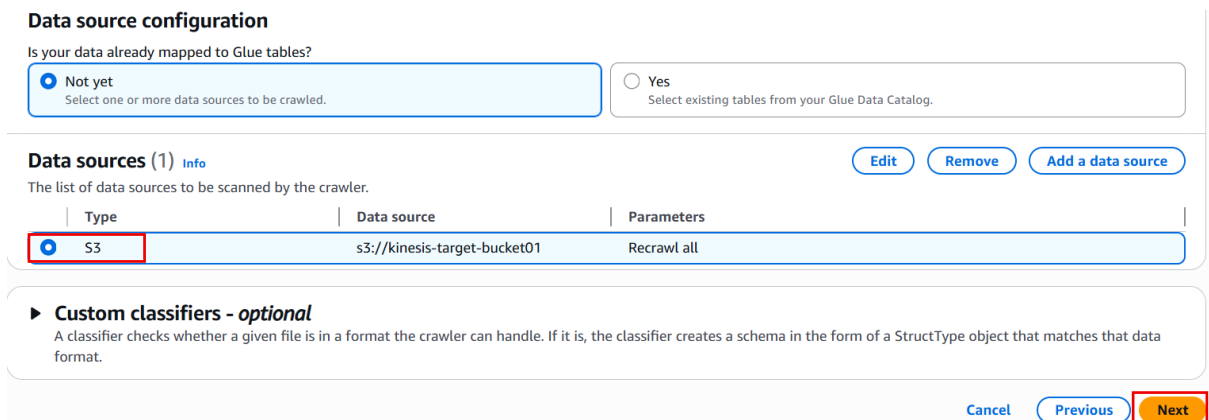- Now click on the Add a data source.
- Click on the Browse S3

- Choose the folder and click choose.



- Then, Click on Add an S3 data source.
- Check the Data Sources and Click Next.



- Select an existing role with S3 and Glue permissions
- Click on the Add database

## Set output and scheduling

### Output configuration Info

**Target database**

Choose a database ▼ 🔄

Clear selection    Add database 🔗

- Give the database name and click Create.

**Database details**

**Name**

retail-data

Database name is required, in lowercase characters, and no longer than 255 characters.

**Description - optional**

Enter text

Descriptions can be up to 2048 characters long.

**Database settings**

**Location - optional**
Set the URI location for use by clients of the Data Catalog.

An S3 location is required for managed tables and Zero-ETL integrations.

Cancel    **Create database**

- Now go back to the clawer setup and refresh, you will see the database.
- Optional add a **table prefix** and click Next.
- Then Click on Create crawler.

**Step 4: Set output and scheduling**    Edit

**Set output and scheduling**

| Database | Table prefix - optional | Maximum table threshold - optional | Schedule |
|---|---|---|---|
| customers | table- | - | On demand |

Cancel    Previous    **Create crawler**

- After the crawler is create, Click on the Run Crawler.

Last updated (UTC)
October 16, 2025 at 18:03:25    🔄    Run crawler    Edit    Delete

- Go to **AWS Glue → Tables**.
- Open your new table and verify:
  - Columns detected from CSV files.
  - **Partition key** (e.g., partition_0) added automatically.

- Open **Athena** and select your Glue database.



- Add a New Partition Folder

In S3, create a new folder (e.g., Tokyo/) and upload sales_tokyo.csv.



- Run the same Athena query — **Tokyo data won't appear yet**.
- Re-run the crawler and data will appear.
- **Update Metadata for New Partition**
- Go to **Athena** and run:
  ```
  o ALTER TABLE your_table_name
  o ADD PARTITION (partition_0='Tokyo')
  o LOCATION 's3://your-bucket-name/Tokyo/';
  ```
- Re-run your query → **Tokyo data now appears**.