# ETL Jobs

# To Begin with the Lab

**Summary of the Lab.**

In this lab, an ETL job is created in AWS Glue to move data from one S3 folder to another. The source folder **"documents"** contains CSV files, and the target folder **"target-customer"** stores the transformed data. Using **Visual ETL**, an **S3 Source** and **S3 Target** are configured, with **Parquet** as the output format. The Data Catalog is updated to create and maintain table schemas. An **IAM role** with **S3FullAccess** is assigned to allow read/write access. After saving and running the job, permissions are verified if needed. Finally, the Parquet file appears in S3, and the table is added to the **Data Catalog**.

Prerequisite for the Lab:

- o Folder created with the name "documents" inside the S3 Buckets with the files inside the folder.

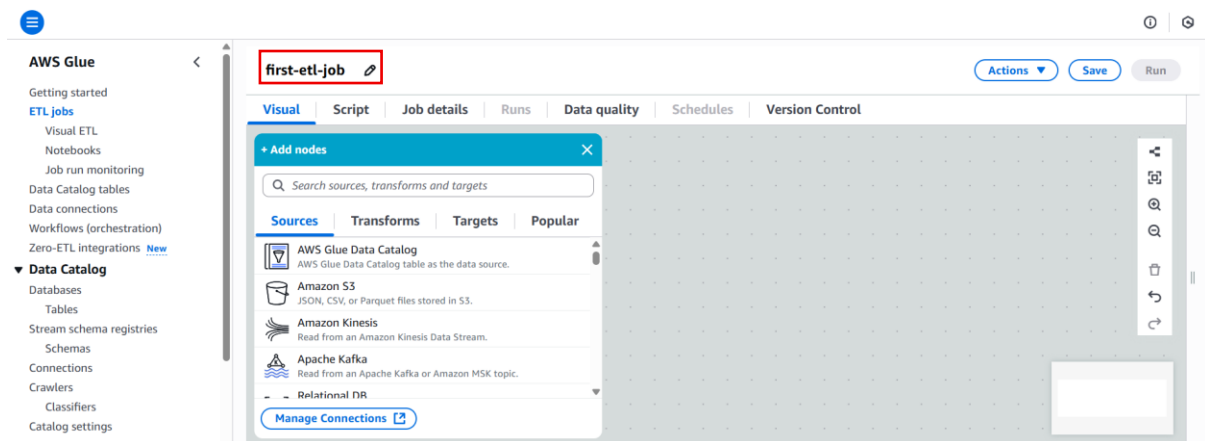- o Create another folder in the same S3 bucket directory with the name target-customer.

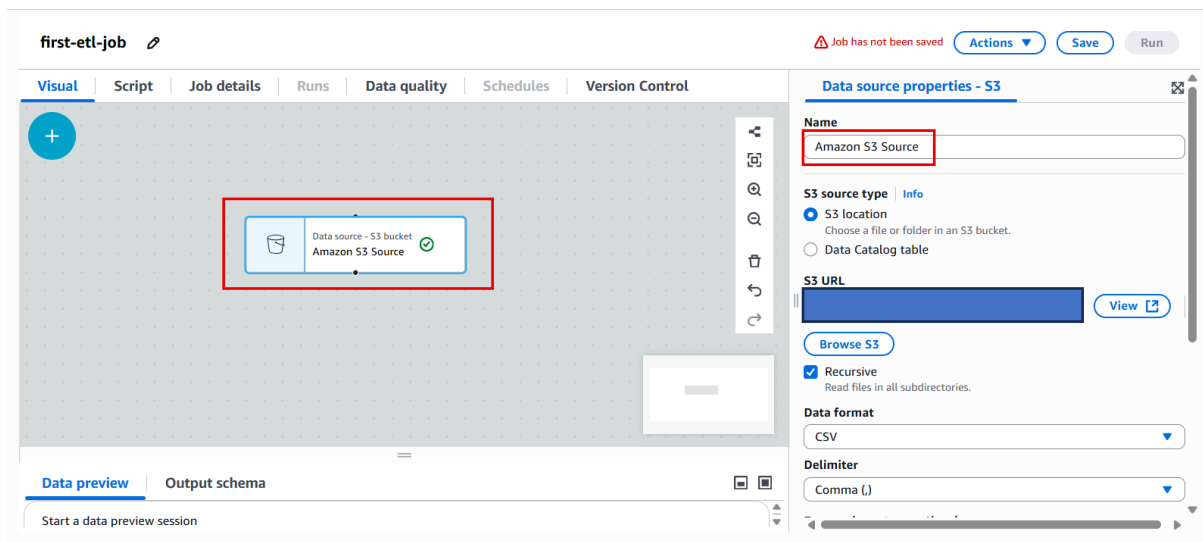- Go to **AWS Console → Glue Service** (ensure you're in the same region as your S3 bucket).



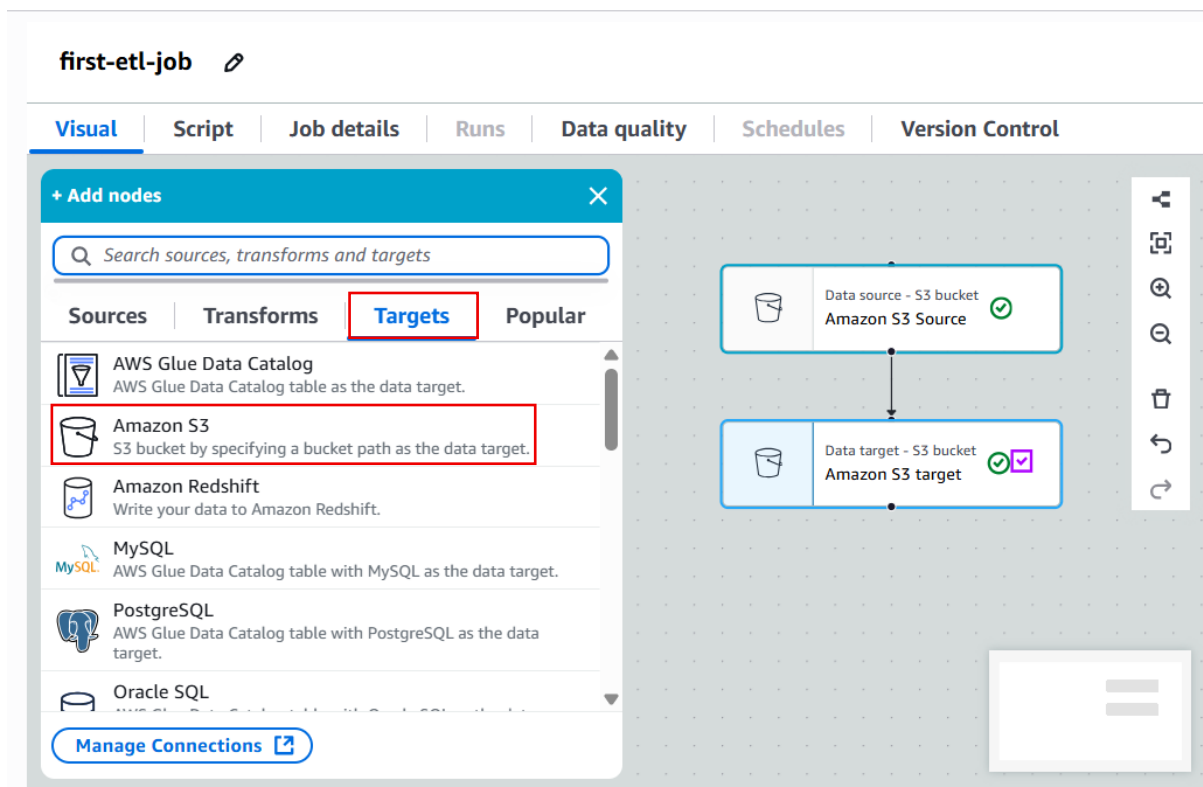- In the **AWS Glue Console**, navigate to **ETL Jobs** and create a new job using **Visual ETL**.



- Name the job



- In the **Visual ETL editor**, add an **S3 Source** node and connect it to the folder containing the CSV file.

- Click on the Add node button and choose Amazon S3.

- We can rename the node and choose the path of the folder from the S3 bucket.

- Click Infer Schema

- Click on the Targets and again select the Amazon S3 bucket.



- Rename it.

- Choose the format as Parquet.

- Choose the S3 Target Location

- In Data Catalog update options, choose the "Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions"

- Choose the Database and give the table name.

- Now go to Job details, Assign an **IAM role** with permissions to read and write from S3 (e.g., attach AmazonS3FullAccess for this lab).
- Change the Requested number of workers.
- Then Click Save and Run.

- Now we can go to Run details to check the status of the job.



- If the job fails with "Access Denied," update the IAM role to include proper S3 permissions and rerun.



- Now Go to IAM role and give the permissions



- Under Permissions, we have to attach policies for the S3 full access.

- Go back to the Job and Rerun it.



- Go back to your target-customer folder and check for the file created.



- Also, we can go the Data Catalog and in database, under tables we can see the table.