# Task Name: AWS Data Pipeline Setup for ETL Processing

## Problem Statement

Organizations often struggle with managing and processing large volumes of data efficiently due to manual processes and a lack of automation. This leads to increased time for data handling and potential errors, hindering timely insights.

## Objectives

1. **Automate Data Ingestion**: Use AWS Glue to automatically retrieve and catalog data from S3.
2. **Transform Data**: Create a Glue notebook to clean and restructure the data.
3. **Load Data**: Save the processed data back to a designated S3 bucket.
4. **Catalog Data**: Maintain an organized schema using AWS Glue for easy access.
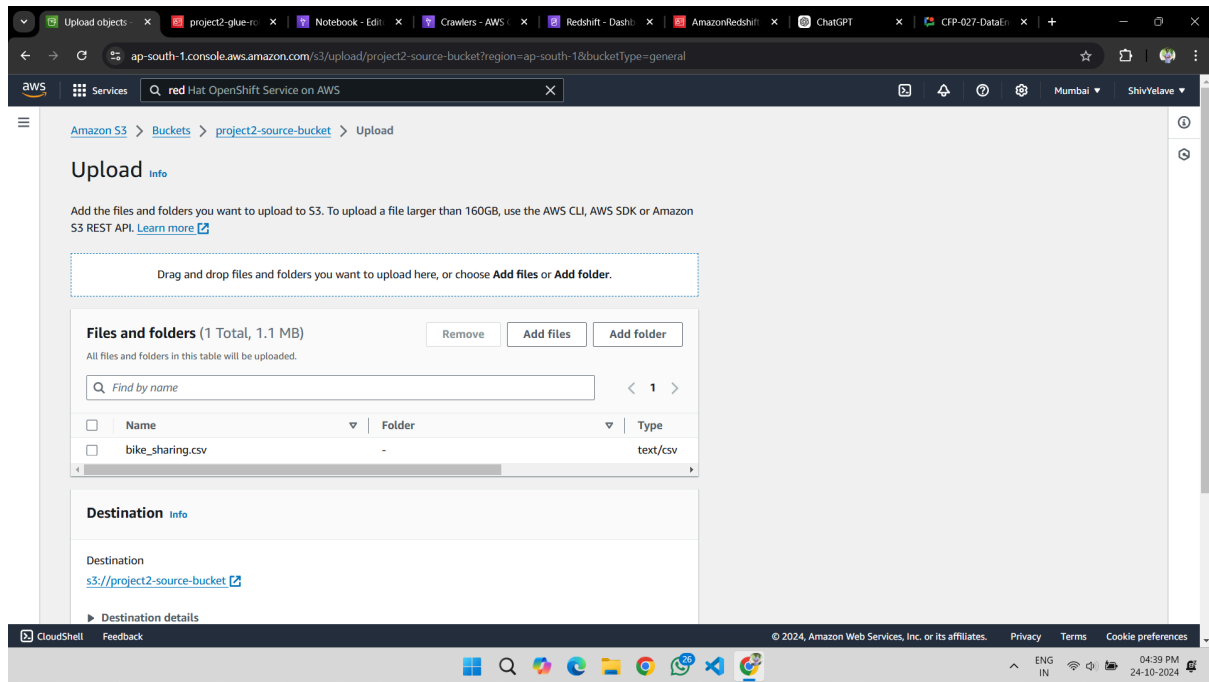
## Overview:

This task involves creating an end-to-end data pipeline in AWS that retrieves data from an S3 bucket, processes it using AWS Glue, and stores the transformed data back in another S3 bucket. The pipeline includes the following components:

1. **Source S3 Bucket**: Stores the initial raw data that needs to be processed.
2. **AWS Glue Crawler**: Automatically discovers and catalogs the schema of the data in the source bucket, making it available for querying.
3. **AWS Glue Notebook Job**: Executes the ETL (Extract, Transform, Load) logic, fetching the data from the Glue catalog, performing necessary transformations, and writing the processed data back to the destination S3 bucket.
4. **Destination S3 Bucket**: Stores the final output of the ETL process, which can be used for further analysis or reporting.

### Step 1: Set Up AWS S3 Buckets

1. **Create Source S3 Bucket**
   - Navigate to the S3 service in the AWS Management Console.
   - Click on "Create bucket."
   - Provide a unique bucket name (e.g., `source-bucket-name`) and configure settings as needed (e.g., versioning, encryption).

- Click "Create bucket."

2. **Create Destination S3 Bucket**
   - Repeat the steps above to create another bucket (e.g., `destination-bucket-name`).

## Step 2: Configure IAM Roles and Permissions

1. **Create a New Policy**:
   - In the IAM console, click on **Policies** in the left sidebar.
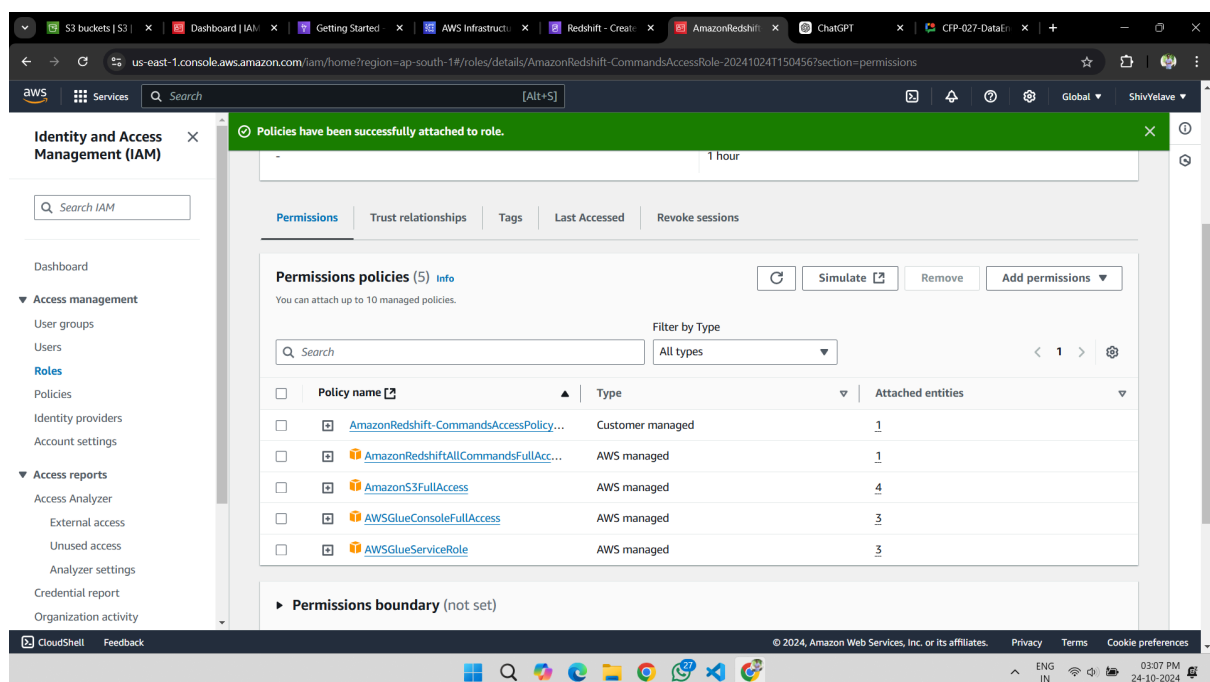   - Click on **Create policy**.

   **Define the Policy**:

   - Switch to the **JSON** tab and paste the following policy, which allows the Glue service to assume roles:

```json
{
    "Version": "2012-10-17",
    "Statement": [
        {
            "Effect": "Allow",
            "Action": "iam:PassRole",
            "Resource": "arn:aws:iam::<account-id>:role/<target-role>"
        }
    ]
}
```

2. **Create an IAM Role for AWS Glue**
   - Go to the IAM service in the AWS Management Console.
   - Click on "Roles" > "Create role."
   - Select "Glue" as the trusted entity.
   - Attach the following policies:
     - `AmazonS3FullAccess` (or a more restrictive policy specific to your buckets)
     - `AWSGlueServiceRole`
     - `Above Policy that we created.`
   - Click "Next: Tags," then "Next: Review."
   - Provide a name for the role (e.g., `GlueServiceRole`) and click "Create role."



**Step 3: Create an AWS Glue Crawler**
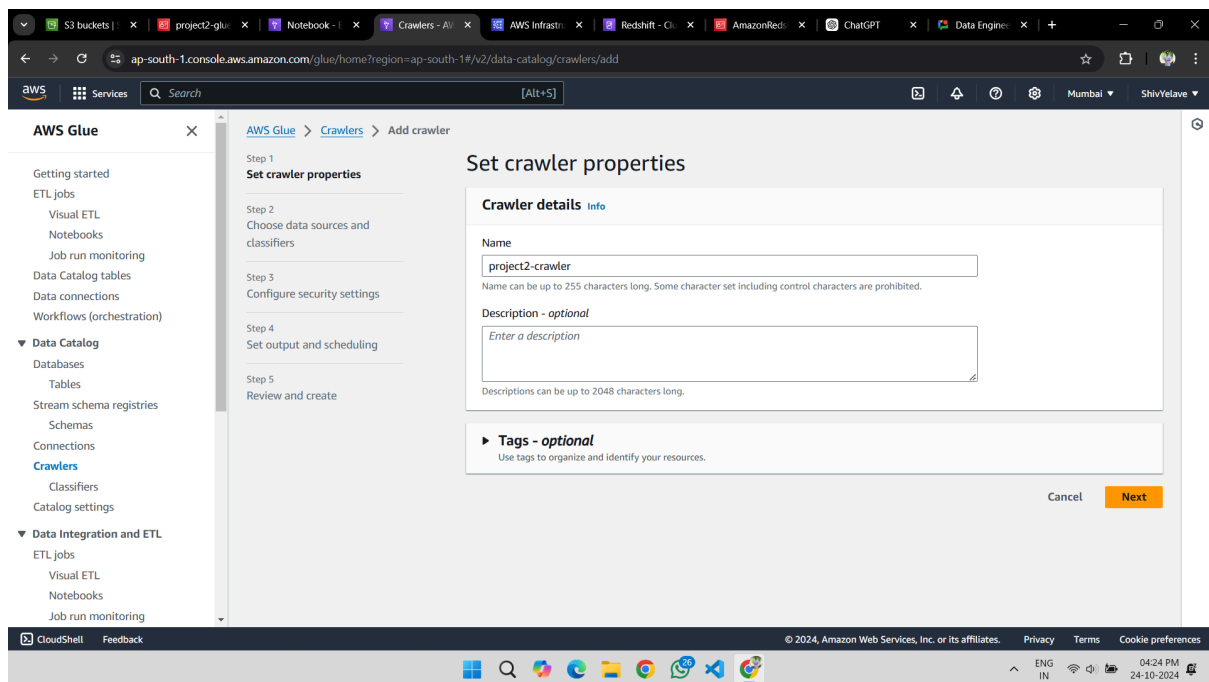
1. **Navigate to AWS Glue Service**
   - In the AWS Management Console, go to the Glue service.
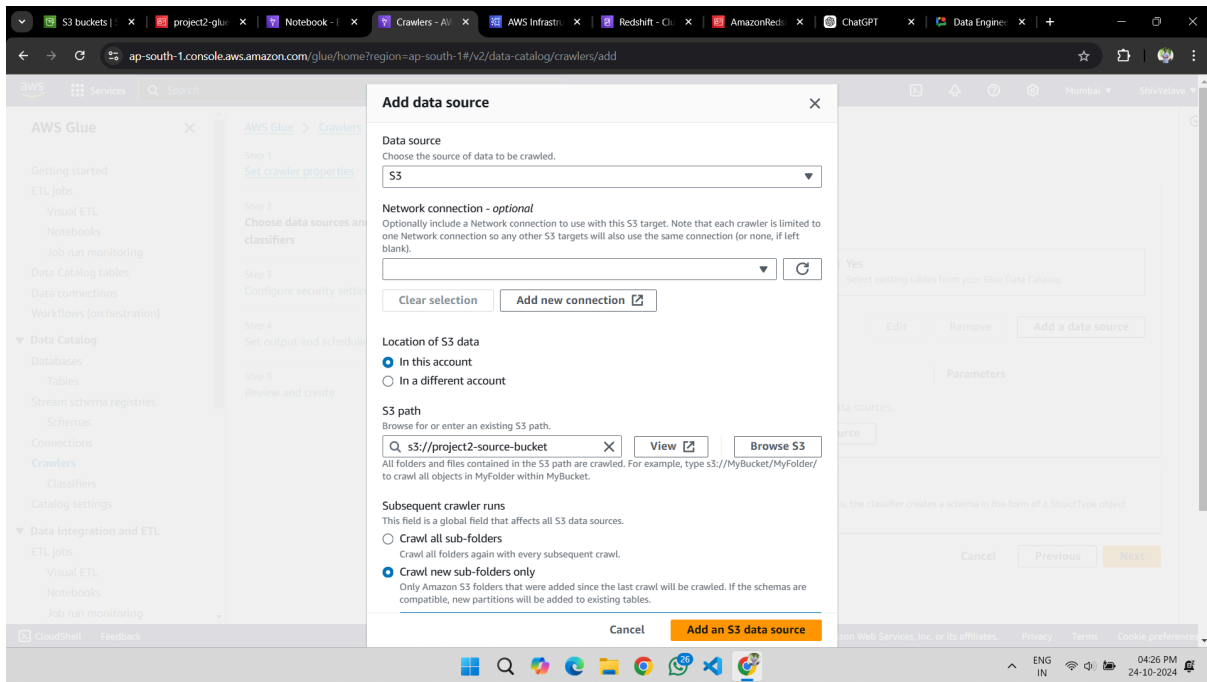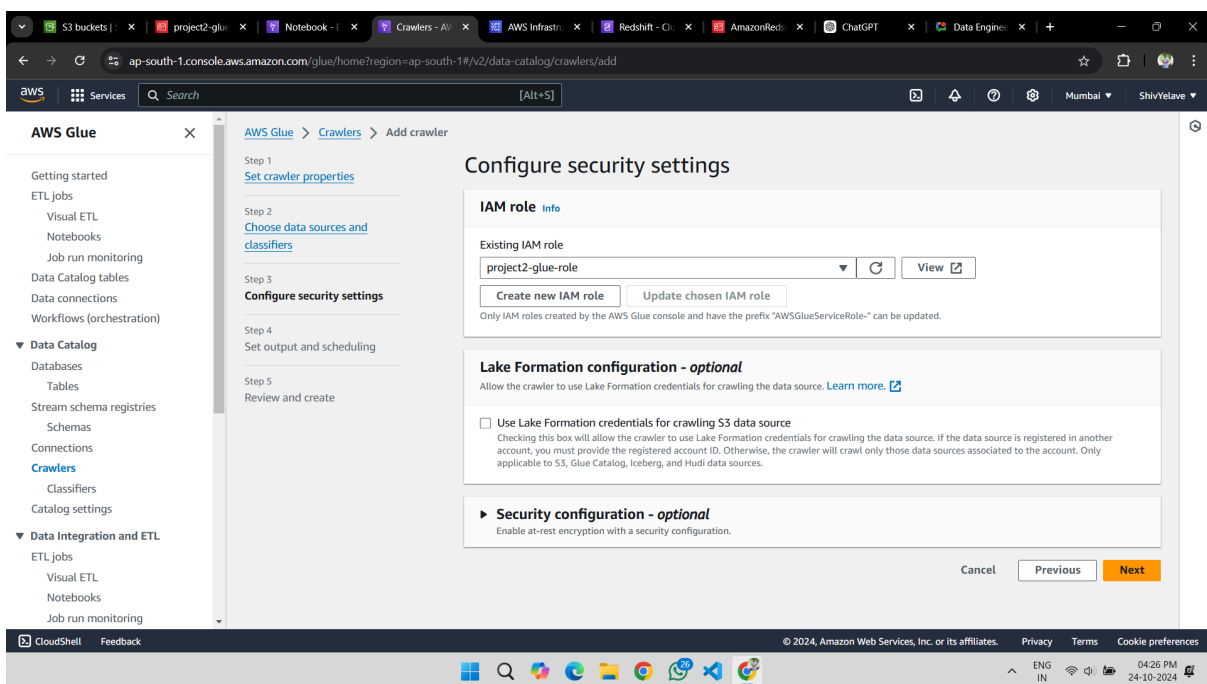2. **Create a Crawler**
   - Click on "Crawlers" > "Add crawler."
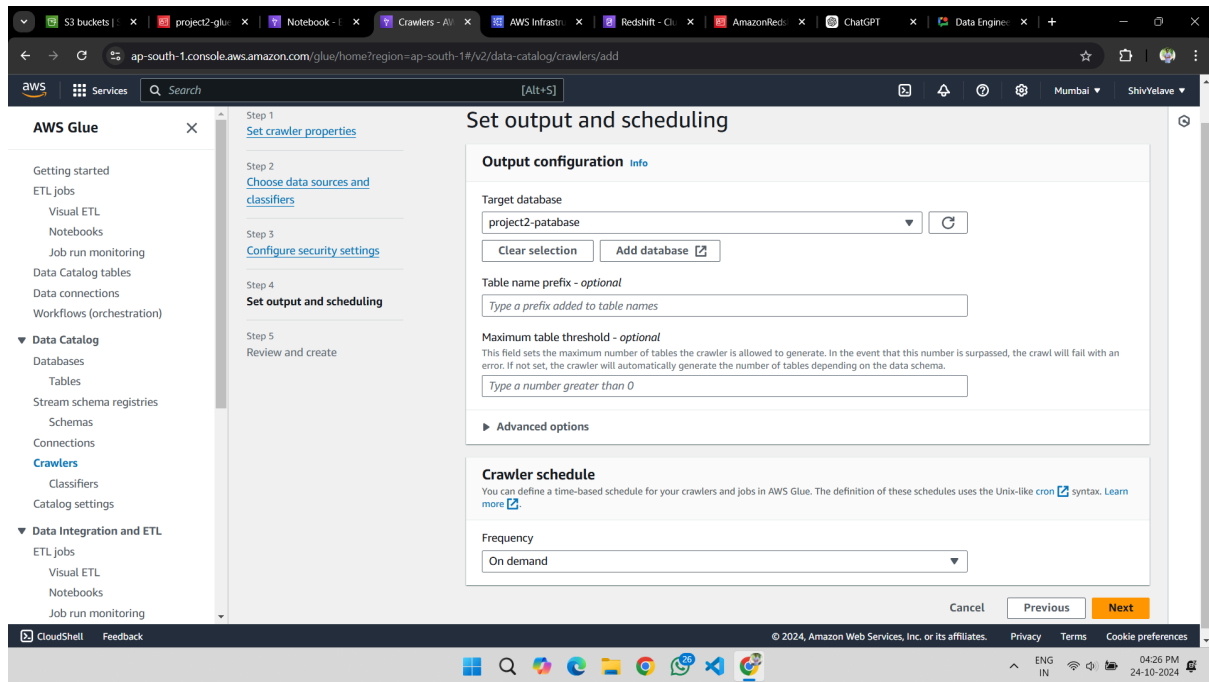
○ Name your crawler (e.g., `my-data-crawler`).



○ Choose "Data stores" as the crawler source type.
○ Select "S3" and specify your source bucket (e.g., `source-bucket-name`).

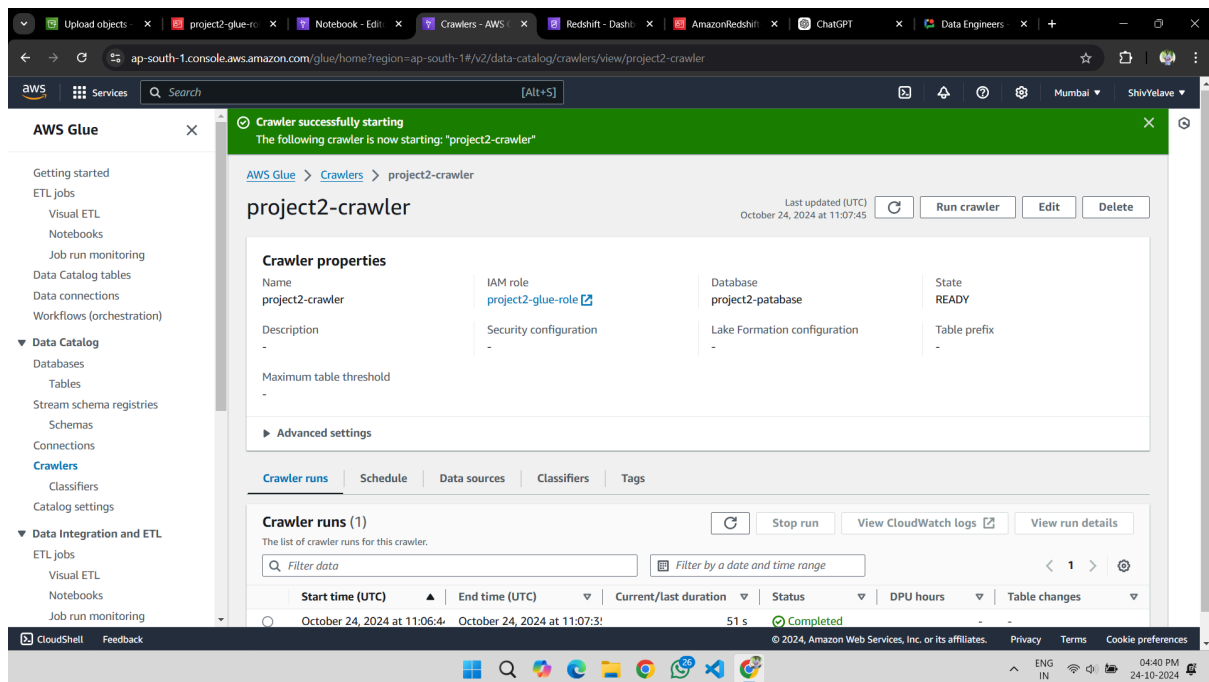- ○ Specify the IAM role you created earlier (e.g., `GlueServiceRole`).



- ○ Click "Next" and configure options for output (databases, tables).

- ○ Review and finish creating the crawler.
3. **Run the Crawler**
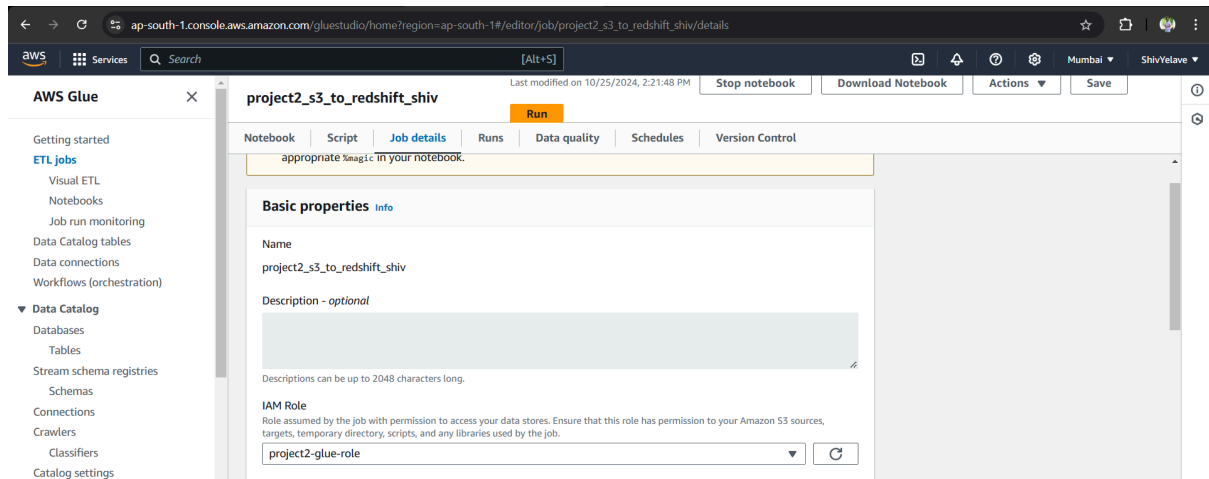   - ○ Go to the Crawlers page, select your crawler, and click "Run crawler."



## Step 4: Create a Glue Notebook Job

1. **Navigate to Glue Jobs**
   - ○ In the Glue service, click on "Jobs" .
   - ○ Click on **Add notebook**.
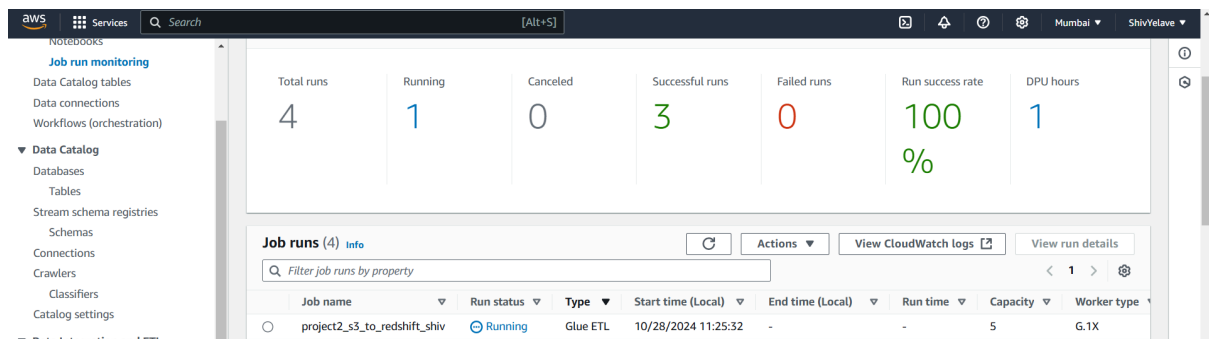   - ○ Name your job (e.g., `my-glue-notebook-job`).

○ Choose the IAM role you created earlier (e.g., `GlueServiceRole`).



○ Write and Save Glue Notebook Code.
○ Click "Next" and configure any additional options (like job bookmarks).
○ Review and create the job.
2. **Save the Notebook Job**
   ○ Save your notebook and ensure it's in the correct format.

## Step 5: Monitor and Maintain the Pipeline

● Regularly check logs in CloudWatch for the Glue job and crawler.



● Monitor data in the destination S3 bucket to ensure proper ETL execution.