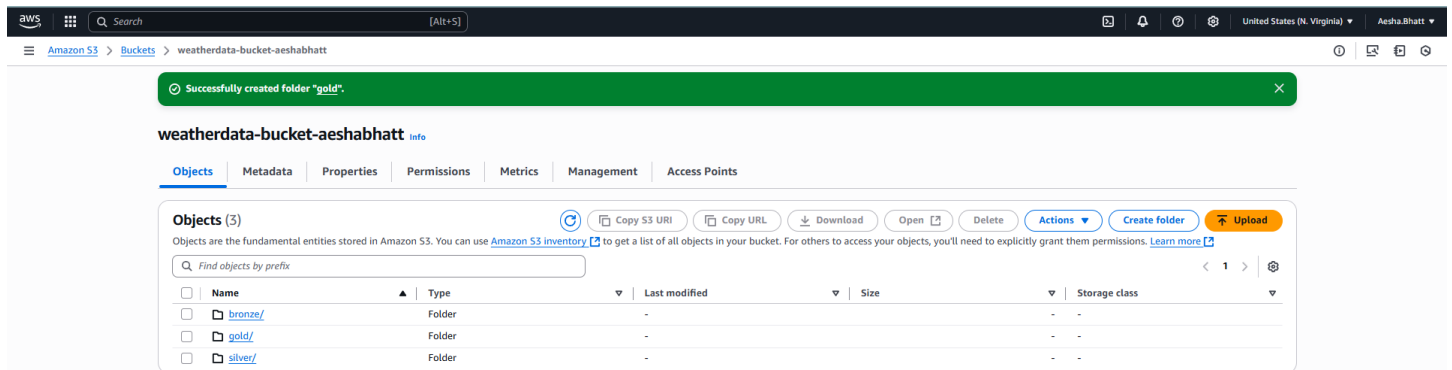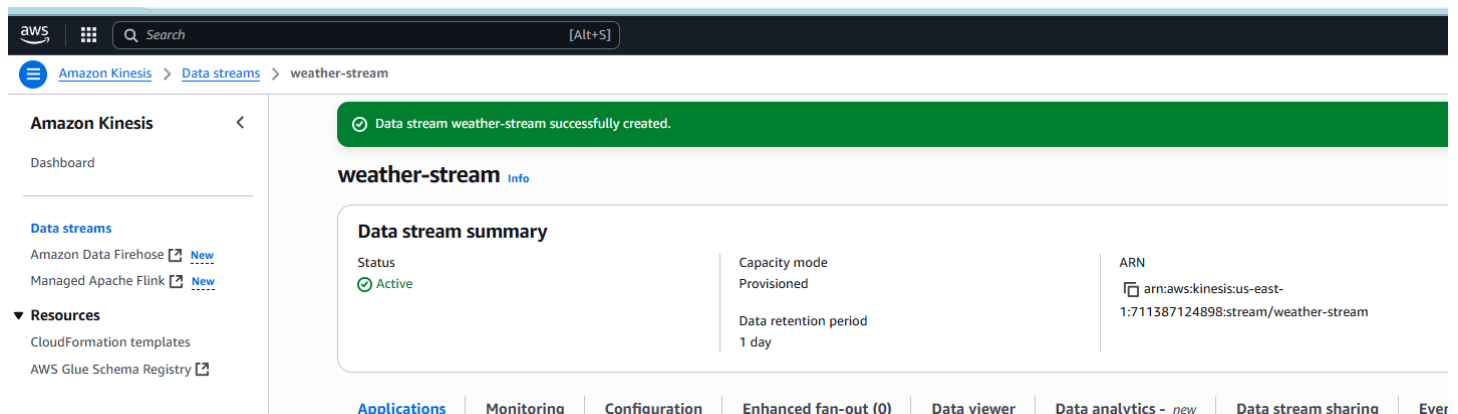# Building a Real-Time Weather Data Pipeline for Weather Analytics

Create an S3 bucket using the AWS Management Console. Once the bucket is set up, create three separate folders within it named bronze, silver, and gold to structure and manage your data efficiently.



## BRONZE LAYER:

1. Create a Kinesis Data Stream using the AWS Management Console. Specify the stream name and configure the number of shards based on the expected data throughput requirements.

2. Prepare and run a script to stream weather data. Begin by creating an IAM role with Lambda as the use case. Attach the necessary permissions, such as access to Kinesis, CloudWatch, and any other services the Lambda function will interact with. This role will be assigned to the Lambda function responsible for ingesting and streaming the weather data into the Kinesis Data Stream.



To run the script in AWS CloudShell:

1. Click the terminal icon on the bottom navigation bar of the AWS Console to open CloudShell.
2. Open a new file using the nano editor:
   a. Run the command: nano weather_stream-project-3.py
3. Paste the streaming script into the file.
4. To save and exit nano:
   a. Press Ctrl + O to write (save) the file
   b. Press Enter to confirm
   c. Press Ctrl + X to exit the editor
5. Make the script executable by running:
   a. chmod +x weather_stream-project-3.py
6. Run the script using the command:
   a. python3 weather_stream-project-3.py

```python
import boto3
import json
import time
import requests

# Create a Kinesis client
kinesis = boto3.client("kinesis", region_name="us-east-1")

# API key and list of cities
API_KEY = "f2ffdd03d1f74ac5b09fa933e24e8876"
CITIES = [
    "New York",
    "Los Angeles",
    "Mexico City",
    "Toronto",
    "Chicago",
    "Houston",
    "Miami",
    "Dallas-Fort Worth",
    "Montreal"
]

# Infinite loop to send weather data every 60 seconds
while True:
    for city in CITIES:
        url = f"http://api.openweathermap.org/data/2.5/weather?q={city}&appid={API_KEY}"
        response = requests.get(url).json()

        # Send data to Kinesis stream
```

```
kinesis.put_record(

    StreamName="weather-stream",

    Data=json.dumps(response),

    PartitionKey=city

)


    print(f"Sent weather data for {city}")


# Wait 60 seconds before the next batch

time.sleep(60)
```
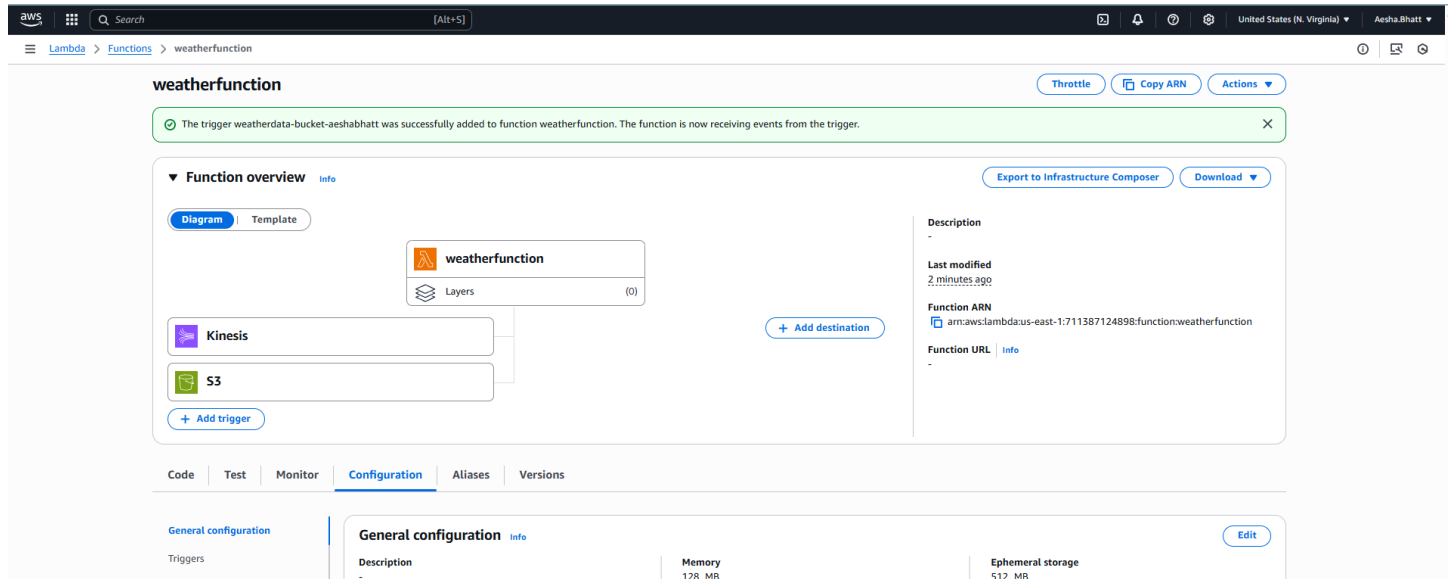
Notes:

- Make sure the stream name "weather-stream" exists in **us-east-1**.
- The IAM role associated with your environment must have permissions for:
  - kinesis:PutRecord
  - logs:* (for monitoring via CloudWatch)

# SILVER LAYER

1. Create a Lambda function to handle the weather data. Use Python 3.x as the runtime and assign the previously created IAM role with appropriate permissions (access to Kinesis, S3, and CloudWatch).
2. Add the following triggers to the Lambda function:
   - **Kinesis**: Configure the Lambda function to be triggered by the Kinesis Data Stream. This allows it to consume real-time weather data as it arrives.
   - **S3**: Integrate S3 by using the boto3 client within the Lambda function to store raw weather data in the Bronze layer and cleaned/processed data in the Silver layer of the S3 bucket.



import json

import boto3

import base64

from datetime import datetime

import pandas as pd

import io


# AWS clients

s3 = boto3.client('s3')


# bucket & prefixes

bucket = "weatherdata-bucket-aeshabhatt"

```python
bronze_prefix = "bronze/weather_data/"
silver_prefix = "silver/weather_data/"


def lambda_handler(event, context):
    print("Received event:", json.dumps(event))

    # Detect Kinesis event
    if 'Records' in event and 'kinesis' in event['Records'][0]:
        return handle_kinesis_event(event)

    # Detect S3 event
    elif 'Records' in event and 's3' in event['Records'][0]:
        return handle_s3_event(event)

    else:
        return {
            "statusCode": 400,
            "body": "Unknown event type"
        }

def handle_kinesis_event(event):
    """Triggered by Kinesis : stores raw JSON in Bronze layer (handles multiple records)"""
    try:
        # Decode all records in the batch
        records = [
            json.loads(base64.b64decode(rec['kinesis']['data']).decode('utf-8'))
            for rec in event['Records']
        ]

        now = datetime.utcnow().strftime("%Y-%m-%d-%H-%M")
```

```python
        filename = f"weather_data_{now}.json"


        # Save as a JSON array
        body = json.dumps(records, indent=2)


        # Upload to Bronze layer
        s3.put_object(
            Bucket=bucket,
            Key=f"{bronze_prefix}{filename}",
            Body=body,
            ContentType='application/json'
        )


        print(f"Uploaded {len(records)} records to s3://{bucket}/{bronze_prefix}{filename}")
        return {"statusCode": 200, "body": f"Uploaded {len(records)} records"}


    except Exception as e:
        print(f"Error in handle_kinesis_event: {e}")
        return {"statusCode": 500, "body": str(e)}


def handle_s3_event(event):
    """Triggered by S3 : reads Bronze JSON → clean
```
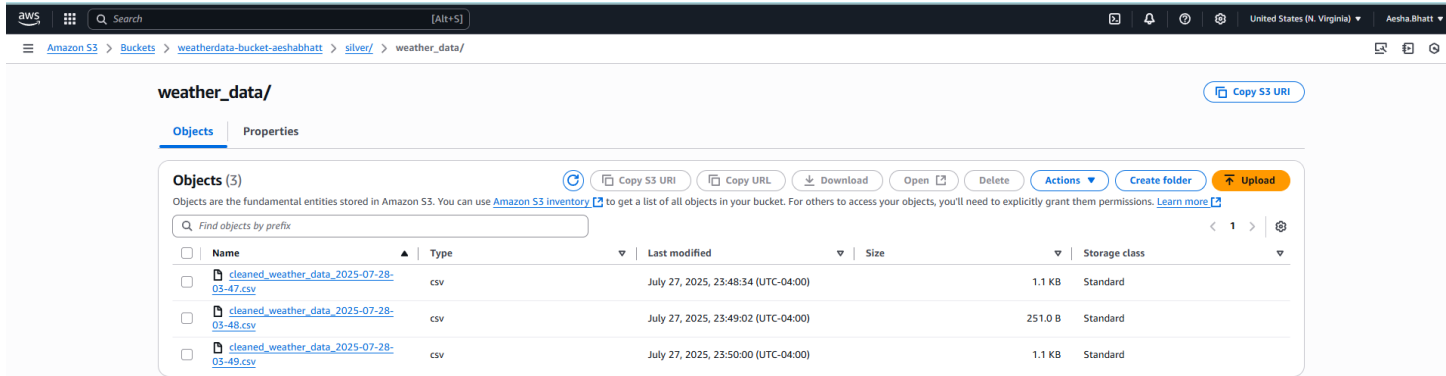
To connect your AWS Redshift cluster to an IAM role, follow these steps:

1. **Create or identify an IAM Role** with the necessary permissions for Redshift to access other AWS services (e.g., S3).

2. **Attach the IAM Role to your Redshift cluster:**

   o Open the AWS Management Console.

   o Navigate to **Amazon Redshift > Clusters**.

   o Select your Redshift cluster.

   o Choose **Actions > Manage IAM roles**.

   o Attach the IAM role you want Redshift to use.

3. **Use the IAM role for authentication or to access resources** such as S3 for COPY/UNLOAD commands without using AWS keys.

This enables Redshift to securely access AWS resources using the assigned IAM role.

## Project 3 – Aesha Bhatt

Here's a clear step-by-step for your document:

1. Create a new VPC in the AWS Management Console.
2. Navigate to **Security Groups** under the VPC dashboard.
3. Select the security group associated with your VPC (or create a new one).
4. Edit the **Inbound Rules** of the security group.
5. Add a rule with the following settings:
   - Protocol: **TCP**
   - Port Range: **5439**
   - Source: specify the IP range or security group allowed to connect (e.g., your IP or 0.0.0.0/0 for all)
6. Save the inbound rule to allow Redshift traffic on port 5439.

This opens the default Redshift port so clients can connect.

To create a Redshift workgroup configuration, follow these steps:

1. Open the AWS Management Console and navigate to **Amazon Redshift**.

2. In the left navigation pane, select **Workgroups**.

3. Click on **Create workgroup**.

4. Enter a **workgroup name**.

5.  Select the **VPC** where your Redshift cluster and resources reside.

6.  Choose the **subnets** (preferably private subnets in multiple availability zones) that the workgroup will use.

7.  Assign **security groups** that control access to the workgroup.

8.  Configure additional settings such as:

    o  Encryption options.

    o  Enhanced VPC routing.

    o  Logging options.

9.  Review the configuration and click **Create workgroup** to finalize.

This workgroup configuration helps manage the networking, security, and access settings for your Redshift Serverless environment.

# GOLD LAYER

To create an AWS Glue Connection, follow these steps:

1. Open the AWS Management Console and navigate to **AWS Glue**.
2. In the left navigation pane, select **Connections**.
3. Click **Add connection**.
4. Enter a **name** for the connection.
5. Choose the **connection type** (e.g., JDBC for databases).
6. Configure the connection properties such as:
    o JDBC URL or endpoint
    o Username and password (if required)
    o VPC, subnet, and security groups (to enable Glue to access resources inside your VPC)
7. Review the details and click **Create connection** to save it.

This connection allows AWS Glue jobs and crawlers to securely access data stores inside your VPC or external databases.

To create an AWS Glue Visual ETL job, follow these steps:

1. Open the AWS Management Console and go to **AWS Glue**.
2. In the left navigation pane, select **Jobs**.
3. Click **Add job**.
4. Enter a **name** for your job.
5. Choose **Visual with a source and target** as the job type.
6. Select the **IAM role** that has the necessary permissions for Glue to access your data sources and targets.
7. Under **This job runs**, choose the type of data processing (e.g., Spark).
8. Choose the **data source** for your ETL job (e.g., a table in the Glue Data Catalog or a connection).
9. Choose the **data target** where you want to write the processed data (e.g., S3 bucket).
10. Use the **Visual editor** to drag and drop transforms, apply mappings, filters, and other transformations as needed.
11. Review the job settings and click **Save**.

12. To run the job, select it and click **Run job**.

This creates and executes a Glue ETL job using the visual interface, simplifying ETL development without writing code manually.



To check data in a Redshift table, follow these steps:

1. Connect to your Redshift cluster using a SQL client (such as **Redshift Query Editor**, **SQL Workbench/J**, or any other SQL client).
2. Select the appropriate **database** and **schema** where your table resides.
3. Review the query results to verify the data in the table.

This lets you quickly inspect the contents of your Redshift table.

**Visualisation:**



Sum of Temp_celsius, Sum of Humidity, and Sum of Wind_speed by City