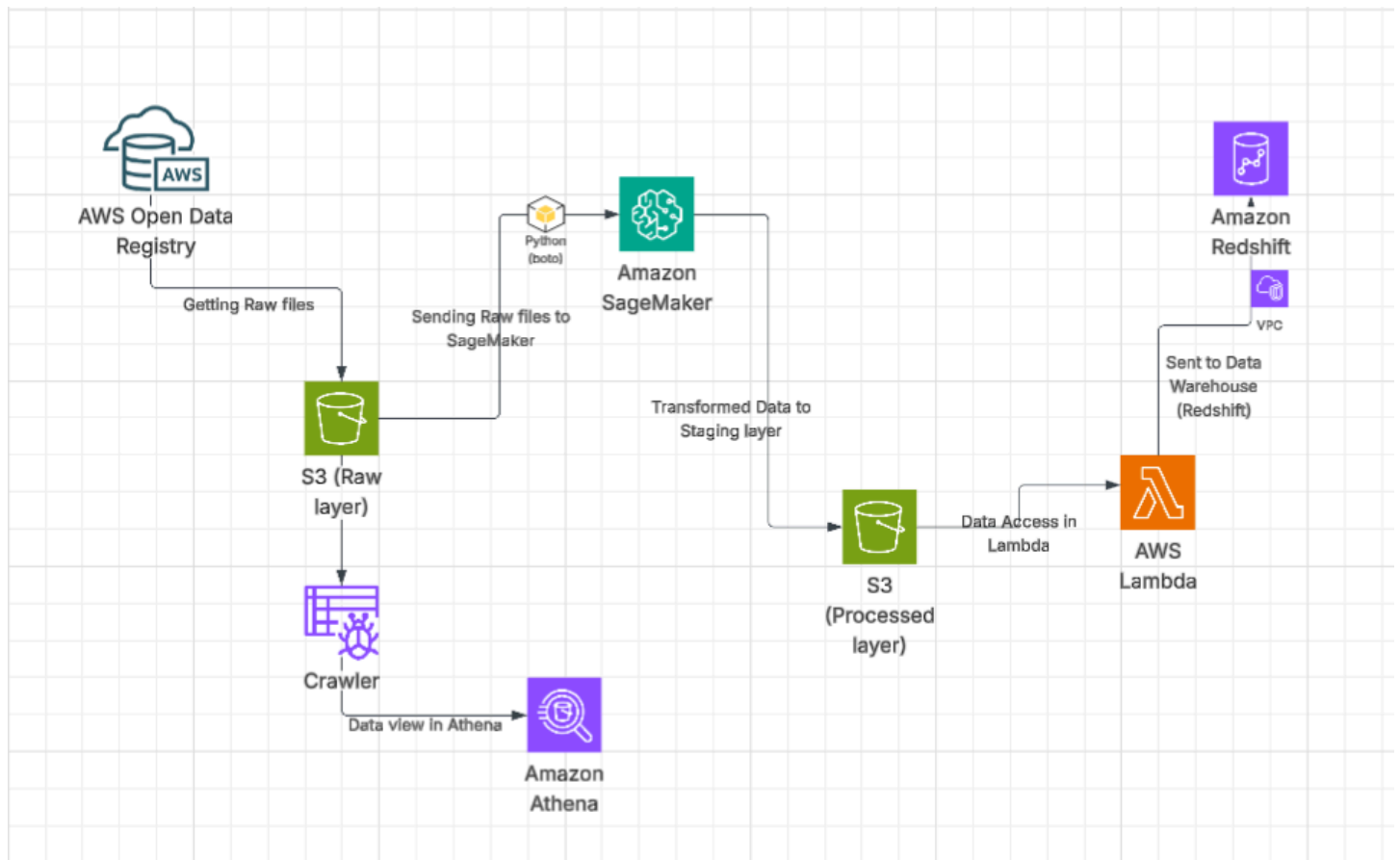# Covid Data Raw to DataWarehouse Ingestion AWS

Git: [Azure-Projects/Covid Data Ingestion from git to Data Warehouse at main · diljotrandhawaa/Azure-Projects](#)

Workflow:



**The goal of the project:**
Get the Raw files about Covid-19 data and transform the data into Dimension and Fact tables to store them in Data Warehouse.


**Tech used:**
Amazon S3
IAM
Glue Crawler
Amazon Athena
Amazon SageMaker
Amazon SageMaker Studio

Jupyter Notebook
Python boto3
S3 Event Notifications
AWS Lambda
Psycopg2 (Layer in Lambda) (to connect to Redshift directly)
Amazon VPC
Amazon Redshift


**Before** starting with Implementation:
1. Create two S3 buckets, one for Raw layer and one for Transformed/Processed layer.
2. Create a IAM role for the project, give it S3FullAccess, SageMakerFullAccess, LambdaFullAccess and RedshiftFullAccess.
3. Create a SageMaker Studio Domain if not already present.
4. Create a JupyterLab workspace and run it.
5. Create a Lambda function.
6. Download zip file for Pysopg2 layer from the internet. Used zip file is present in the git repository of the project.
7. Create a Redshift namespace and workgroup.


Project Implementation:

# Steps:

## Step 1: Raw Data Collection and Data Ingestion

In the first step, we manually download the raw files from an AWS Open Data Registry that has covid-19 Data. It's called Covid Data Lake.
We upload the files to S3 bucket (Raw layer).
Inside the Raw layer bucket, the four raw files are stored in 4 different folders (directories) so that the crawler creates one table for each raw file.
Below is the structure of the Raw bucket:

Each folder is named after the Raw file that is stored inside it.

## Step 2: Create Crawler for raw bucket and processed bucket (Data Quality Check)

We **create a Crawler for the Raw data**. The crawler points to the Raw S3 bucket. The source is directory containing four directories for raw files.



After creating, we run the crawler.
The data can be viewed in Athena for query purposes and to understand the structure of the raw data.
Below is what the "covid_data" raw file looks like

Next, we **create one Crawler for the Processed S3 bucket** to view the processed data in Athena. We don't run it yet.



## Step 3: Jupyter code implementation (Data Transformation)

Now that we understand the structure of the data better, we export the data to Jupyter notebook in SageMaker Studio's Jupyterlab to transform it further.

Below is the small snippet of the Jupyter notebook, the code file is present in the git repository.

In the notebook, we first import the raw data from S3 using boto3. We then create our Dimension tables and Fact tables using merge and other operations.
The dimension tables for Hospitals, Region and Date are created while a Fact table is created for Covid data.
(dimHospital, dimRegion, dimDate, factsCovid) are created.

Using the code in Jupyter notebook, the dim and fact tables will be exported to our Processed Layer (S3) using boto3 and StringIO buffer in the form of CSV files.

## Step 4: Test the Processed data and Create Event notification for automation (Automation)

We can test run the code to see if we receive the correct data in our Processed S3 bucket. Run the Crawler to view data in Athena. After the data is verified and validated, it is then ready to move to Data Warehouse in Redshift.

To automate the ingestion, we create a Event notification on the Processed S3 bucket and connect it to the Lambda function created earlier. The Lambda function does not have any code yet, we'll move on to the Lambda function now.
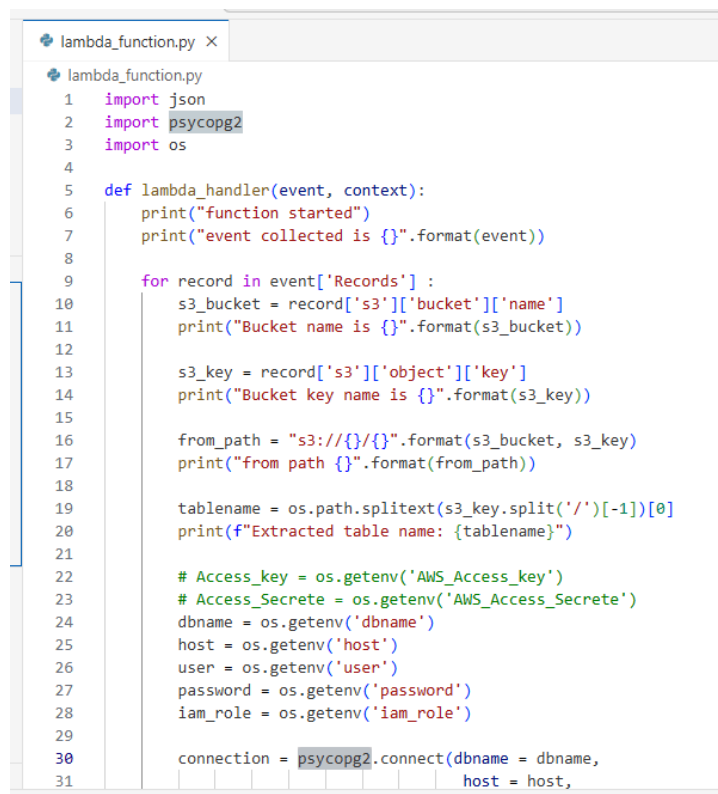The event notification created can be seen below.

**Step 5**: Configure the Lambda function (Data Storage)

We go to the Lambda function and write the code. But before that, we need to add a custom layer for Psycopg2 to connect directly to Redshift. To do that, we need a zip file for the layer, the used zip file can be found in the git repository.

After installing and attaching the layer to the Lambda function, we create some Enviornment variables for cleaner code and for safety. The variables are created for Access keys and credentials to connect to Redshift.

The code is written in the Lambda function. Code simply imports the dimension and fact tables from Processed layer. It directly connects to the Redshift and creates four tables in Redshift. It then copies the data from S3 to Redshift.

The lambda code file can be found in the git repository of the project.

```
lambda_function.py ×

lambda_function.py
1    import json
2    import psycopg2
3    import os
4
5    def lambda_handler(event, context):
6        print("function started")
7        print("event collected is {}".format(event))
8
9        for record in event['Records'] :
10           s3_bucket = record['s3']['bucket']['name']
11           print("Bucket name is {}".format(s3_bucket))
12
13           s3_key = record['s3']['object']['key']
14           print("Bucket key name is {}".format(s3_key))
15
16           from_path = "s3://{}/{}".format(s3_bucket, s3_key)
17           print("from path {}".format(from_path))
18
19           tablename = os.path.splitext(s3_key.split('/')[-1])[0]
20           print(f"Extracted table name: {tablename}")
21
22           # Access_key = os.getenv('AWS_Access_key')
23           # Access_Secrete = os.getenv('AWS_Access_Secrete')
24           dbname = os.getenv('dbname')
25           host = os.getenv('host')
26           user = os.getenv('user')
27           password = os.getenv('password')
28           iam_role = os.getenv('iam_role')
29
30           connection = psycopg2.connect(dbname = dbname,
31                                          host = host,
```

Now deploy the function.

**Step 6**: Add VPC Security group to the Lambda function

To successfully connect to Redshift, make sure your Lambda function is configured to connect to same VPC Security group as the Redshift workgroup.

Now run the code and the data will be sent to the Data warehouse Redshift.

Code run example:

So all we have to do is run the Jupyter notebook. As soon as we run the notebook and data gets sent to the Processed bucket, it will automatically trigger the Lambda function and copies the data in the Redshift warehouse.