# AUTOMATED CLOUD BASED ETL PIPELINE WITH RDS,S3,GLUE AND SNS INTEGRATION

**COACH** : LEON YIGAL BENJAMIN

**MENTOR** : BADRI PRASAD BALAKRISHNAN

**TEAM MEMBERS** :-

RAHUL V - 227633

PRAVEEN G - 2276623

DHAMOTHARAN A S - 2276706

VISHWA R - 2276634

LIKITHA CHINTAPALLI - 2276784

**BU/SL** : Data
**Technology** : AWS

# ABSTRACT :-

This project scenario involves the implementation of an Extract, Transform, and Load (ETL) process in the AWS cloud. The project begins by creating an Amazon RDS database and establishing a connection with MySQL. A snapshot of the RDS database is exported to an Amazon S3 bucket using an IAM role for Backup and security. The exported CSV files are then uploaded to an S3 bucket. When a file is uploaded to an S3 bucket, the AWS lambda gets triggered, and it calls the Step function State machine. So the step-function state machine automatically orchestrates the glue crawler, and the AWS Glue crawler is used to scan and catalog the CSV files, storing the crawled data in a newly created database in the AWS Glue Data Catalog. The CSV files are transformed to Parquet format using the visual script editor in AWS Glue and stored in a separate S3 bucket. The Glue crawler is again automatically done by step-function, where step-function will get triggered by lambda whenever a file uploads in the S3 bucket. Then it will scan and crawl the Parquet files, storing the data automatically in another database in the Data Catalog. Data transformation is performed using the visual script editor and AWS Glue DataBrew, and the transformed data is loaded into a new S3 bucket. EventBridge is integrated to trigger the Simple Notification Service (SNS) for email notifications upon successful completion of the Glue job, providing timely updates on job status. This ETL project in the AWS cloud enables efficient data extraction, transformation, and loading, leveraging the scalability, reliability, and flexibility of AWS services.

# Introduction :-

**Purpose of the Document** :

The purpose of this document is to provide an abstract or summary of a project that involves the extraction, transformation, and loading processes using AWS services such as RDS,Glue,Databrew,S3, Eventbridge, and SNS. The document aims to highlight the key components of the project, including its objectives, the technologies involved, and the benefits it offers.

**Project Overview** :-

This project aims to create the database in the RDS and extract that database by using snapshots, then by using the AWS glue, just crawling and storing it in the data catalog, and then doing the transformation, which is concat,split,drop, join, and dateformat. Then storing the result in the S3 bucket. If the glue job is run, it will send the email to the user by using SNS.

**AWS** – **AMAZON WEB SERVICES**

Amazon Web Services (AWS) is a comprehensive and widely adopted cloud computing platform offered by Amazon. It provides a vast array of on-demand computing resources and services, delivered over the internet, to help businesses and individuals

build and run their applications, store and analyze data, and scale their operations with ease.

**This project utilizes the following AWS services** :-

**1**. **Amazon Ec2 :** Amazon Elastic Compute Cloud (Amazon EC2) is a web service provided by Amazon Web Services (AWS) that offers resizable compute capacity in the cloud. It allows you to quickly provision virtual servers, known as EC2 instances, and run applications and workloads on them.

**2**. **AWS RDS** : Amazon Relational Database Service (Amazon RDS) is a collection of managed services that makes it simple to set up, operate, and scale databases in the cloud.. process by monitoring the S3 bucket for new files and triggering the Glue job. Lambda functions are also utilized to crawl the destination bucket after the conversion is completed.

**3. Iam Roles :** AWS Identity and Access Management (IAM) roles are a fundamental component of the AWS security model. They enable you to securely delegate permissions to AWS resources, such as EC2 instances, Lambda functions, or users within your AWS account.

**4**. **Amazon S3** : Amazon Simple Storage Service (S3) is an object storage service that offers industry-leading scalability, durability, and security. It is used as the storage infrastructure for the CSV files, as well as the destination bucket for the transformed Parquet files.

**5. AWS KMS :** AWS Key Management Service (KMS) is a fully managed service provided by Amazon Web Services (AWS) that

allows you to create and control the encryption keys used to encrypt your data. KMS helps you protect sensitive information by providing a secure and scalable solution for key management.

6. **AWS Glue** : AWS Glue is a fully managed extract, transform, and load (ETL) service that simplifies the process of preparing and transforming data for analytics. It is used in this project to convert the CSV files to Parquet format and optimize the data for efficient processing and storage.

7. **AWS Lambda :** AWS Lambda is a serverless compute service provided by Amazon Web Services (AWS). It allows you to run your code without provisioning or managing servers, enabling you to focus on writing and deploying your applications rather than worrying about infrastructure.

8. **AWS Step Function :** AWS Step Functions is a serverless workflow orchestration service provided by Amazon Web Services (AWS). It allows you to coordinate and visualize the execution of multiple AWS Lambda functions, AWS Batch jobs, Amazon ECS tasks, and other serverless or traditional compute resources as a series of steps, forming a workflow.

9. **Glue DataBrew** : AWS Glue DataBrew is a new visual data preparation tool that makes it easy for data analysts and data scientists to clean and normalize data to prepare it for analytics and machine learning.

10. **Amazon SNS** : Amazon Simple Notification Service (SNS) is a fully managed messaging service that enables the sending of notifications to subscribers. In this project, SNS is used to send

notifications to subscribers about the completion of the Glue job, providing timely updates on the status of the conversion process.

**11. Eventbridge :** EventBridge is a fully managed event bus service provided by Amazon Web Services (AWS). It enables the creation, routing, and processing of events between different applications or services within an AWS environment or across various software systems.

**SERVICES USED** :

AWS Cloud

| | | | | |
|---|---|---|---|---|
| Amazon RDS | Amazon S3 | Amazon EC2 | AWS KMS | IAM Roles |
| AWS Glue | Glue DataBrew | Event Bridge | Amazon SNS | AWS Lambda |
| Step Functions | | | | |

# ARCHITECTURE :

**Implementation of Project Architecture** :

**1**. **Create IAM Role** :

AWS Identity and Access Management (IAM) is a web service provided by Amazon Web Services (AWS) that allows you to manage access to AWS resources securely. IAM enables you to create and manage AWS users, groups, and permissions to control who can access specific AWS resources and what actions they can perform

**Users** : IAM users are entities within your AWS account that represent individual people or applications requiring access to AWS resources. Each user is assigned a unique set of security credentials (access key ID and secret access key) for programmatic access and a password or a password policy for AWS Management Console access.

**Groups** : IAM groups are collections of IAM users. You can manage permissions for multiple users collectively by attaching policies to a group rather than individual users.

**Roles** : IAM roles are similar to users but are not associated with a specific identity. Roles define a set of permissions that determine what actions can be performed and what resources can be accessed.

## ROLE 1 : GlueAccess

**DESCRIPTION** : The purpose of this role is to grant specific permissions that align with the needs of the user or service accessing Glue.

**PERMISSION** : Administrator Access

**PURPOSE** : The purpose of granting "Administrator access" to AWS Glue is to provide users or roles with full administrative privileges and permissions within the AWS Glue service. By assigning the Administrator access level, users can have complete control over all aspects of AWS Glue, including managing databases, tables, jobs, triggers, and connections.

**ROLE 2** : **DataBrewAcess**

**DESCRIPTION** : The "DataBrewRole" is an IAM role designed to provide specific permissions required for AWS DataBrew, a data preparation and data blending service offered by AWS. The role is tailored to meet the needs and actions associated with utilizing DataBrew effectively.

**PERMISSION** : Administrator Access

**PURPOSE** : Granting "Administrator Access" permission within AWS provides users or roles with full administrative privileges and permissions across all services and resources within the AWS account. The purpose of granting Administrator access is to allow complete control and management over the entire AWS environment.

**ROLE 3** : **RDS-S3-ExportRole:**

**DESCRIPTION** : We need to grant specific permissions that allow the role to interact with Athena database and perform necessary actions.

**PERMISSION 1** : **S3 Full Access :**Provides full access to Amazon s3 via the AWS Management Consoles3 full access.

# STEP 1: EC2 INSTANCE CREATION

Instance Name : Rds-london

Creating an instance in Amazon Ec2 (Elastic Compute Cloud) allows to provision and manage virtual servers in the cloud.

# STEP 2 : RDS DATABASE CREATION

Database Name : finalproject-shop-db

Created An Database Using Amazon RDS

**STEP 3** : **CONNECTING THE DATABASE IN MYSQL WORKBENCH USING END-POINT AND CREDENTIALS**

TABLES :

1. CUSTOMER  - CustomerID,FirstName,LastName,Address

2. ORDER  -  OrderID,OrderDate,CustomerID

3. ORDER DETAILS  - OrderItemID,OrderID,ProductID,Quantity

4. PRODUCT -  ProductID,ProductName,Price

# STEP 4 : CREATED A SNAPSHOT USING RDS

## NAME : finalproject-shop-snap



# STEP 5 : EXPORTED RDS SNAPSHOT TO S3 BUCKET IN PARQUET FORMAT USING IAM ROLE

## Target Name : finalproject-rds-snapshot

## STEP 6 : **Uploading the database in S3 - CSV format**

Target Name : finalproject-shop-csv

FOLDER :

1. Customer-csv
2. Orders-csv
3. Orderdetail-csv
4. Product-csv

## STEP 7 : TRIGGER THE S3 BUCKET USING LAMBDA FUNCTION

Function Name 1 : lambda-stepfunction

Function Name 2 : lambda-stepfunc-transformation

**CODE** :

```python
import boto3

import time

def lambda_handler(event, context):

    # Sample bucket name

    bucket_name = 'finalproject-shop'

    # Instantiate the AWS Step Functions client

    sf_client = boto3.client('stepfunctions')

    # Generate a unique execution name with a timestamp

    execution_name = f"execution_name_{int(time.time())}"

    # Start the execution of the Step Function state machine

    response = sf_client.start_execution(

stateMachineArn='arn:aws:states:eu-west-2:631707933068:stateMachine:fileformat-csv-parquet',

        name=execution_name,

        input='{"bucket": "' + bucket_name + '"}'

    )
```

# Return the execution ARN as the output of the Lambda function

return {

'statusCode': 200,

'body': f"Step Function execution started: {response['executionArn']}"

}

## STEP 8 : AUTOMATING THE STEP FUNCTION USING AWS LAMBDA

State Machine 1 : lambda-stepfunc-transformation

State Machine 2 : transform-parquet

## STEP 9 : CRAWL THE DATA FROM S3 AND STORED IN DATA CATALOG

Name : Extract-csv



## STEP 10 : CHANGING THE FILE FORMAT OF CSV TO PARQUET USING AWS GLUE AND UPLOADING TO S3

**STEP 11 : AWS GLUE DATA CATALOG CREATION**

Database Name : transform-parquet

**STEP 12 : CRAWL THE DATA FROM S3 AND STORED IN DATA CATALOG**

Crawl Name : transformation-crawl

**STEP 13 : CHANGING THE DATE FORMAT USING DATABREW-SOURCE DATA CATALOG And UPDATING IN SAME DATA CATALOG**

Job : date-format-change

Step name : Change format-of-orderdate-to dd/mm/yyyy

# STEP 14 : TRANSFORMATION USING AWS GLUE -VISUALLY STORING IN S3

Job Name : transformation-job

Transformation : Date Format Change ,Split ,Drop ,Concat ,Join

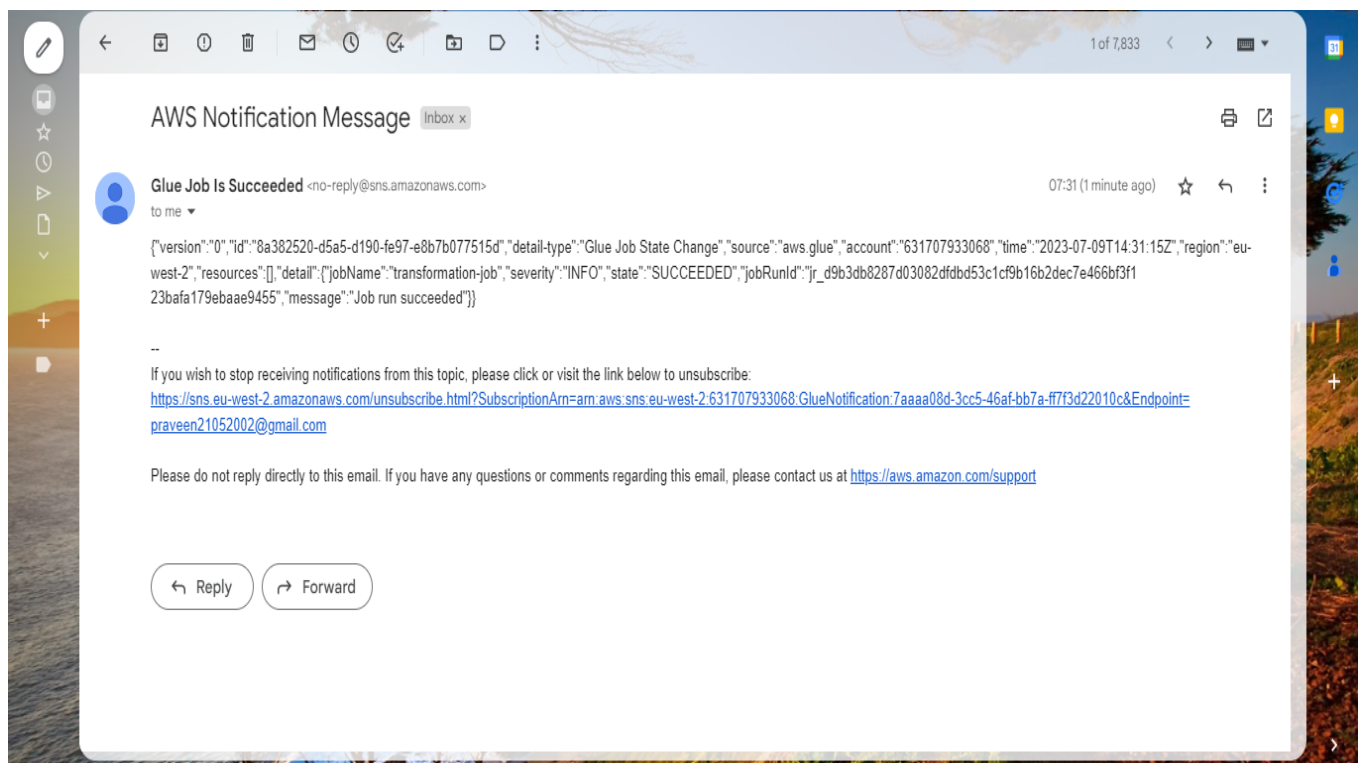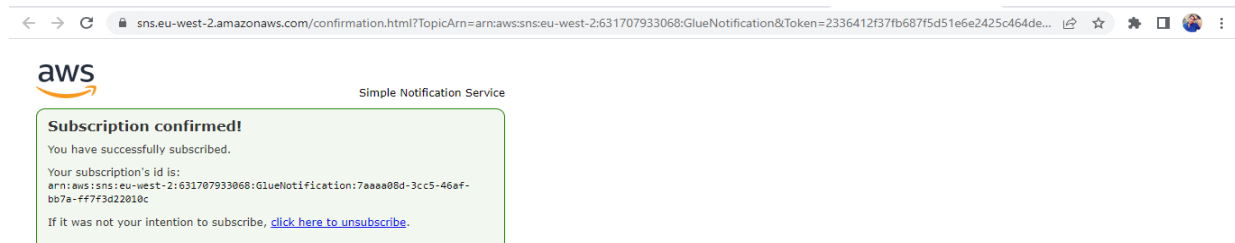Target Bucket : finalproject-transformed-parquet

## STEP 15 : SNS CREATION (SIMPLE NOTIFICATION SERVICE)

Created a topic named as SNS Glue Notification and added a EMAIL subscription of our personal mail's.So whenever SNS gets triggered it will send a notification.

Topic Name : Glue Notification

Subscription : dhamotharan596@gmail.com

praveen21052002@gmail.com

**STEP 16** : **EVENT BRIDGE RULE CREATION**

Rule Name : gluesuccess-event

In the event pattern,we chose the source service as glue and wrote the code for event-type gluejob-statechange. If gluejob-statechange succeeds, we set the target to trigger the SNS topic. So it will send a notification to EMAIL as "GLUE JOB IS SUCCEEDED.


**JSON Code** :-

```
{

  "source": ["aws.glue"],

  "detail-type": ["Glue Job State Change"],

  "detail": {

    "jobName": ["transformation-job"],

    "state": ["SUCCEEDED"]

  }

}
```

Description:

The provided code snippet is an event pattern designed for monitoring the state change of an AWS Glue job.

# Conclusion :-

In conclusion, this project showcases the successful implementation of an Extract, Transform, and Load (ETL) process in the AWS cloud. It leverages a range of AWS services including Amazon RDS, S3, AWS Glue, DataBrew, and EventBridge to achieve efficient data extraction, transformation, and loading capabilities.

By establishing a connection with MySQL and exporting a snapshot of the RDS database to an S3 bucket, data backup and security are ensured. The use of the AWS Glue crawler facilitates the scanning, cataloging, and transformation of CSV files into Parquet format, providing organized and refined data. Furthermore, the integration of EventBridge and Simple Notification Service (SNS) allows for timely email notifications upon job completion, ensuring stakeholders are informed.

Overall, this ETL project in the AWS cloud harnesses the scalability, reliability, and flexibility of AWS services, enabling organizations to efficiently manage and process their data.