

Build an AWS ETL Data Pipeline in Python on YouTube Data



Hi, I am writing this articles on how we can Build an AWS ETL Data Pipeline in Python on YouTube Data In this article I have provided how can we build the project , architecture and its components:

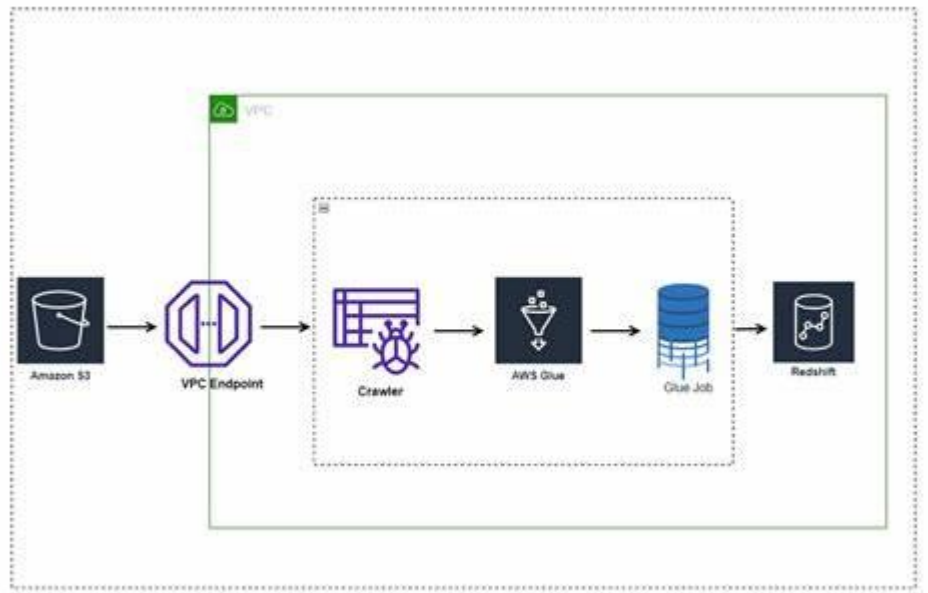
AIM: The aim is to automate the process of extracting data from Youtube Api , transform and load it into AWS Redshift using Python, and eliminate the manual effort and load it into AWS Redshift using Python, and eliminate the manual effort required. It leverage AWS services to provide a scalable, reliable, and cost-effective solutions for handling large volumes of data.

PREREQUISITE:

You should be familiar with AWS services like s3, Glue, redshift and API

PROCEDURE:

- 1) Set up AWS environment: We have to Create an AWS account, set up S3, IAM roles, Redshift, and Glue.



2) Obtain API keys: Obtain YouTube API keys to access the data.

3) Extract data: Use Python script with the Google API client library to extract data from YouTube API and save it to S3 in CSV format.

4) Set up AWS Glue: Set up a Glue job to read the CSV data from S3, transform the data, and save it back to S3 in Parquet format.

5) Load data: Use Python script to load the transformed data from S3 into Redshift using the COPY command.

6) Schedule the pipeline: Schedule the pipeline to run at regular intervals using AWS CloudWatch Events.

7) Monitor the pipeline: Use AWS CloudWatch to monitor the pipeline and capture logs.

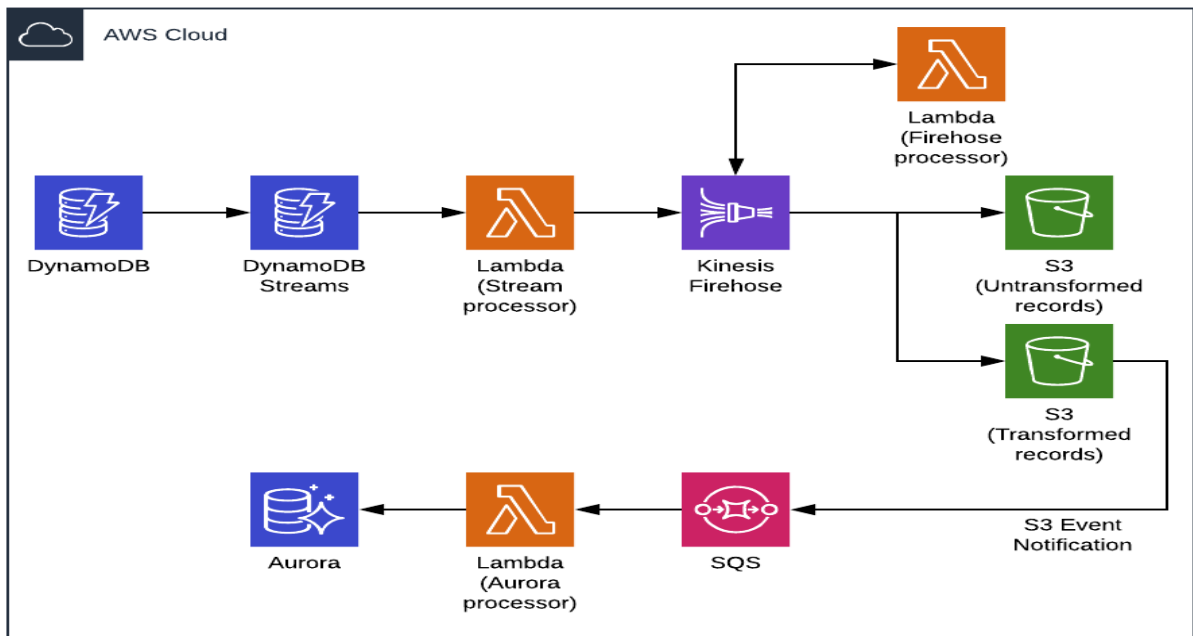
8) Visualize data: Use a BI tool such as Tableau or Power BI to visualize and analyze the data in Redshift.

ARCHITECTURE:

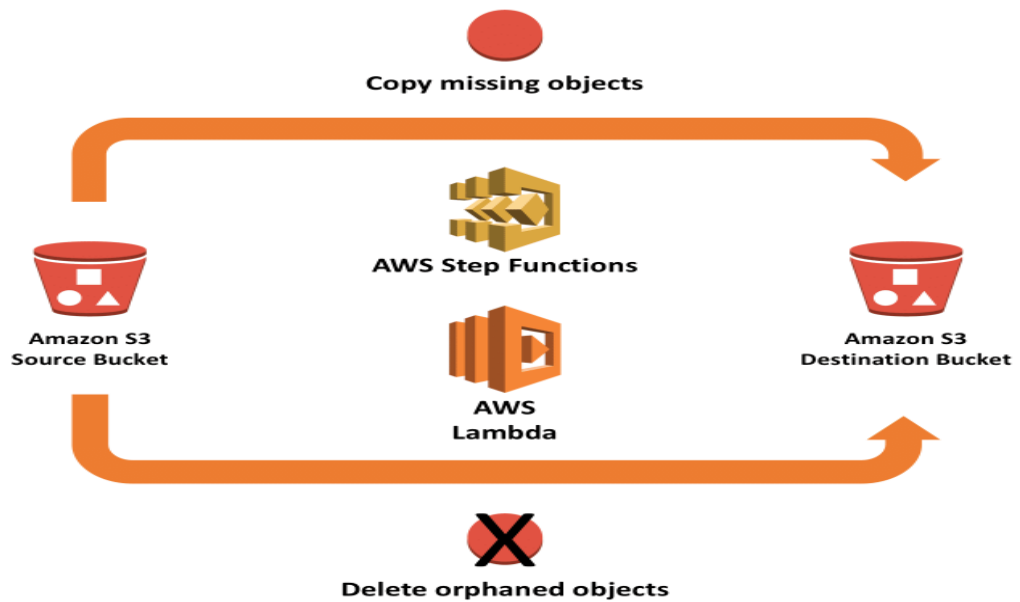
The architecture of the AWS ETL data pipeline for YouTube data involves several components:



1) YouTube API: The data source that provides the raw data.



- 2) AWS S3: The storage component where the raw data extracted from the YouTube API is stored in a CSV file.



- 3) Python: Python is the programming language used to extract, transform, and load data in the pipeline.

4) AWS Redshift: This is the data warehousing component where the transformed data is stored.

5) AWS Glue: This is the ETL (extract, transform, load) service used to transform the data.

6) AWS CloudWatch: This is the monitoring and logging service used to monitor the pipeline and capture logs.

7) IAM Roles: This is the security component used to manage access to the AWS resources used in the pipeline.

The pipeline involves the following steps:

1) Extraction: The Python script extracts data from the YouTube API and saves it as a CSV file in S3.

2) Transformation: The AWS Glue service reads the CSV file from S3 and transforms the data using Python scripts. The transformed data is then saved back to S3.

3) Loading: The Python script loads the transformed data from S3 into Redshift using the COPY command.

4) Monitoring: The AWS CloudWatch service is used to monitor the pipeline and capture logs.

Overall, this architecture provides a scalable and cost-effective solution for processing large amounts of YouTube data.

CONCLUSION:

In conclusion, building an AWS ETL data pipeline in Python on YouTube data provides a powerful solution for processing large amounts of data from the YouTube API. By leveraging AWS services such as S3, Glue, Redshift, and CloudWatch, we can automate the ETL process, eliminate manual effort, reduce the risk of errors, and increase the efficiency of the data processing workflow.