# BUILDING AN EFFICIENT ETL PIPELINE

By:
Janani    III/CSE
Allen     III/CSE
Poonam    III/IT
Akaash    III/CSE

# /TABLE OF CONTENTS

## /AIM

To build an efficient **ETL Pipeline** which Extracts, Transforms and Loads data from a destination to another. It must consist of a web dashboard which performs **visualization** on data and helps **reporting errors** occurring in the pipeline.

# /WHY WE CHOSE THIS PS

- Interesting problem statement

- ETL, Cloud, Data Warehousing is a new concept to us - we wanted to take it up as a challenge.

- There are not many resources like other problem statements- so we learnt topics from scratch by doing a lot of research.

# INITIAL SOLUTION - what we tried performing in the beginning weeks

- Extracting data from SQL Server and loading data into Postgres
- Automating the pipeline using Airflow (using TaskFlow API) and using DAG
- Performing Incremental Data loading using Destination Change Comparison technique
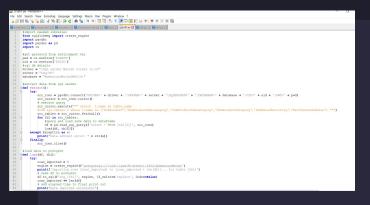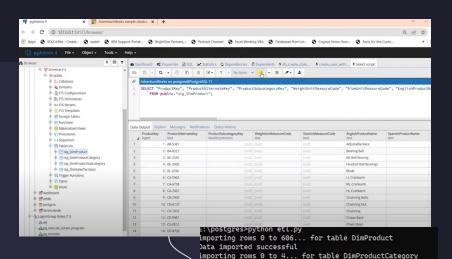- Alteryx- data visualization tool that works excellently with ETL

# INITIAL SOLUTION 1

## 1. Making a basic ETL pipeline using Python



## 2. Results were produced in Postgre from SSMS ⇒ETL successful



## 3. DAG code for automation purposes

# INITIAL SOLUTION 1 PROBLEMS

Postgres
- Loading takes a long time for big data
- Multiple data access becomes difficult
- Low reading speed

Airflow
- Airflow doesn't have versioning- if you delete a task from DAG code and redeploy it, metadata related to task will be lost

# INITIAL SOLUTION 2

# PROBLEMS WRT DOCKER

- Writing ETL pipeline in Python.
- Apply functional programming in Data Engineering.
- Use a proper object oriented code design and meta file for job control.
- Implement a pipeline in Python extracting data Using docker image with application code on Docker Hub, using Kubernetes and Argo Workflows

- Is best for building microservices for applications
- Not good if software runs in different environments or are from multiple sources
- Lesser speed
- Data storage is complicated
- Graphical applications don't work well
- In Docker files are created inside a container. It is difficult to retrieve out of the container for different kind of processes

**SO WE DIDN'T IMPLEMENT THIS**

# <OUR SOLUTION?>

## During our last few weeks?

> Using Cloud of course! <

# /WHY CLOUD?

| | | |
|---|---|---|
| Simplified data process for users | Allows organizations to combine data from multiple sources | Easy optimization |
| Handles huge volumes of data and data migration | Easy for visualization and debugging and reduce MTTR | Easy to compute data |
| No need of much security patches | Maintenance | Easy to analyze unstructured data |

# TOOLS WE ARE USING

## DBEAVER

To create and manage databases across a wide range of database management systems.

## REDSHIFT

Large scale data storage and analysis, and is frequently used to perform large database migrations.

## DATA PIPELINE

helps you reliably process and move data between different AWS compute and storage services, as well as on-premises data sources, at specified intervals.

## GLUE

Uses ETL jobs to extract data from a combination of other Amazon Web Services and incorporates it into data lakes and data warehouses.

# TOOLS WE ARE USING

## CLOUDWATCH

An AWS monitoring service you can use to monitor your applications, services, and resources.

## QUICKSIGHT

Use to deliver easy-to-understand insights to the people who you work with, wherever they are. Amazon QuickSight connects to your data in the cloud and combines data from many different sources.

## PYSPARK

A standard ETL tool like PySpark, supports all basic data transformation features like sorting, mapping, joins, operations, etc.

## ATHENA

Performs interactive queries in the web-based cloud storage service, Amazon Simple Storage Service (S3). Athena is used with large-scale data sets. Amazon S3 is designed for online backup and archiving of data

# /TECHNICAL EXPLANATION

OUR SOLUTION EXPLANATION →

# TECHNICAL EXPLANATION

- We first made a database connection between DBeaver and RDS console. DBeaver will have the entire database.
- Then we set up Redshift cluster by creating roles, editing IAM roles and VPC and configured it as per the requirement to build the ETL
- Created a S3 bucket which has an initial load and incremental load
- In Data Pipeline we configured DataNodes, table structure and created an EC2 instance

# DBEAVER

# REDSHIFT CLUSTER

Go to query editor v2

Amazon Redshift > Clusters > dwhredshift

## dwhredshift ▾

Actions ▼   Edit   Add partner integration   Query data ▼

### General information

| | | | |
|---|---|---|---|
| **Cluster identifier** | **Status** | **Node type** | **Endpoint** |
| dwhredshift | ✓ Available | dc2.large | ▣ dwhredshift.c7bqmdwjxjsz.eu-west-1.red... |
| **Cluster namespace** | **Date created** | **Number of nodes** | **JDBC URL** |
| 6063fba5-785d-43e0-8194-6245813dee75 | May 02, 2022, 22:54 (UTC+05:30) | 1 | ▣ jdbc:redshift://dwhredshift.c7bqmdwjxjs... |
| | **Storage used** | **AQUA** | **ODBC URL** |
| | 0.23% (0.37 of 160 GB used) | Not available | ▣ Driver={Amazon Redshift (x64)}; Server=... |

# /COMPARISON

| | /BIGQUERY | /REDSHIFT | |
|---|---|---|---|
| Service | Serverless/Managed | Not Managed | |
| What are they | fully managed cloud-based data warehouse which is designed for handling large-scale data set storage. | fully managed and serverless data warehouse. | |
| Regions Available | Americas – 4<br>Europe – 3<br>Asia Pacific – 7 | Americas – 8<br>Europe – 5<br>Asia Pacific – 9 | |
| Storage format | Columnar & uncompressed storage | Columnar & compressed storage | |
| Billing | Not Predictable | Predictable | |
| | | | |

# EXPLANATION

- AWS Redshift connection was established with DBeaver
- Initial data load was performed into Redshift for data load for a base object
- Incremental data load was performed into Redshift for subsequent load process requirements
- Incremental Data Pipeline from MySQL to S3 was performed by altering DataNodes and EC2 Resource
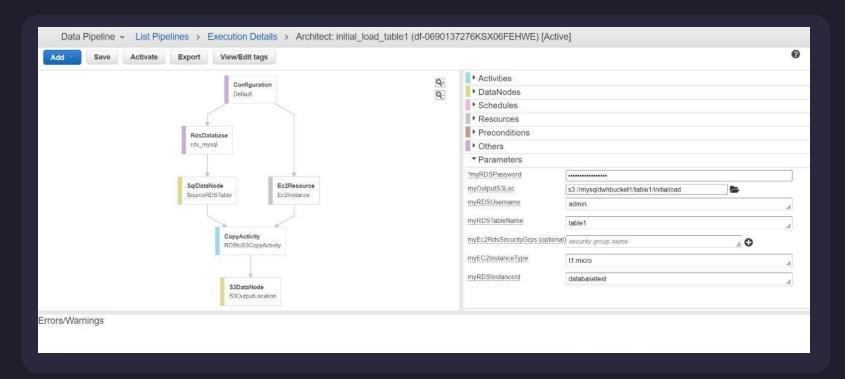- AWS Glue Job from S3 was then configured with Redshift

# DATA PIPELINE

# EXPLANATION

- Lambda Function for AWS Glue Job from S3 was configured to AWS Redshift
- Redshift Snapshot and Restore process performed
- Our ETL was ready

- We then used various tools to visualize and debug our ETL

# LAMBDA FUNCTION in Glue

# /QUICKSIGHT

## /VISUALIZATION

- Interactive dashboards
- Combine a variety of data into a single analysis.and publish dashboards
- Less expensive than Tableau
- Any data can be used

## /SPICE

- Super-fast parallel in-memory calculation engine
- It's faster and easier to retrieve than having to query the original data source each time
- Rapidly perform advanced calculations and serve the data

# QUICKSIGHT

# /CLOUDWATCH

## /LOGS

- Collects compute performance metrics such as CPU, memory, and network as performance events

## /METRICS

- CloudWatch Metric Streams enables you to create continuous, near-real-time streams of metrics to a destination of your choice

## /MONITOR

- Create reusable graphs and visualize our cloud resources and applications in a unified view

## /INSIGHTS

- Provides automated setup of observability for enterprise applications to get visibility into their health

# CLOUDWATCH



Metrics

Logs

# REDSHIFT SPECTRUM ⟶

Users can query data without having to load it.Multiple users can access same S3 bucket at same time.



```
<databasetest> Script-1    *<dev> Script-2    dev    *<spectrum> Script-3  ×

create external schema dwh_external_data_spectrum from data catalog
database 'ecommerce_external_data'
iam_role 'YOUR_ARN'
region 'eu-west-1';

create external table dwh_external_data_spectrum.user_data_file(
    session_id varchar(255),
    event_timestamp varchar(150),
    event_type varchar(50),
    userid varchar(255),
    product_id varchar(255),
    category_id varchar(255)
)
row format delimited
fields terminated by ','
stored as textfile
location 's3://BUCKET_NAME/user_behaviour/'
table properties ('skip.header.line.count'='1')
```
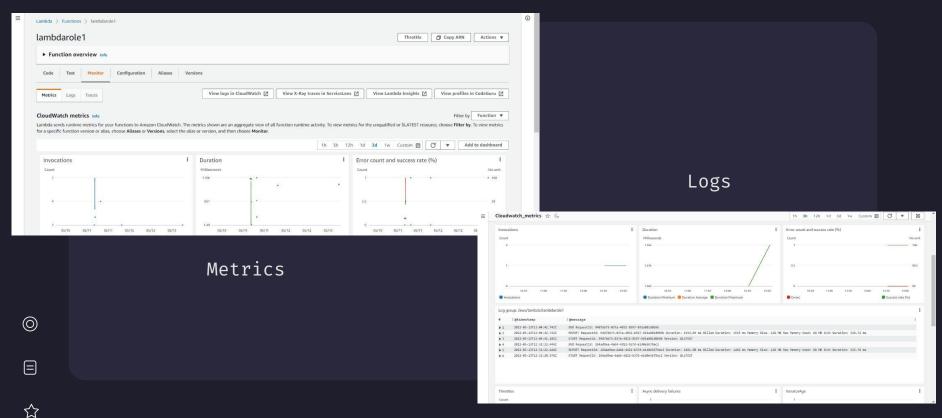
```
select
    b.product_category_name_english ,
    a.year,
    count(1) as view_count
from
    dwh_external_data_spectrum.parquet_output a
join
    mysql_dwh.product_category_name_translation b
on
    a.category_id = b.product_category_name
where a.event_type = 'view'
group by 1,2
order by a.year
```

## /Front-end

1. Build using ReactJS
2. Deployed in Netlify
3. Configured to a new domain
4. Added the domain in Domain and embedding section of the AWS page
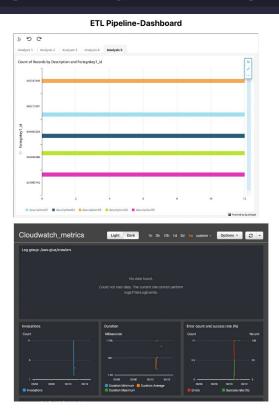
# /CURRENTLY WORKING APPLICATION

# OPTIMIZATION TECHNIQUES FOR DATA ●

## Compound sort keys

Most frequently used
columns put at front
for speed

## Redshift DIST keys

Data distribution
across slices by Leader
node matching values of
a specific column

## Redshift Vacuuming

Reclaims disk space occupied
by rows marked for deletion in
previous operations⇒ compacts
table

# /INNOVATION

- Using Cloud based service - Amazon Web Services (AWS) for efficiency of the ETL Pipeline.
- AWS optimizes the software in terms of time, space and complexity.

Implemented using AWS Free Tier.

# /VIABILITY - INITIAL

**<0.114s>**

> TIME IN QUEUE

**<0.627s>**

> RUN TIME

**<2.52KB>**

> DATA SCANNED

# /VIABILITY - FINAL

## <0.108s>
> TIME IN QUEUE

## <0.446s>
> RUN TIME

## <2.52KB>
> DATA SCANNED

# HOW REALISTIC IN THE MARKET?

## HUGE VOLUMES OF DATA

Great for large scale businesses to work on large data sets

## CENTRALIZED LOCATION

Diverse data such as unstructured data could also be used

## BUSINESS INTELLIGENCE

Easy to analyze using Quick Shift and Cloudwatch

## EVERYTHING IS CLOUD NOWADAYS

Safety of all data sets, easy optimization and editing

# /ROUGH ESTIMATE ($USD)

## My Estimate  Edit ✏

### Estimate summary  Info

| Upfront cost | Monthly cost | Total 12 months cost |
|---|---|---|
| 0.00 USD | 9,110.49 USD | **109,325.83 USD**<br>Includes upfront cost |

# /ROUGH ESTIMATE - SERVICES USED

| | Service Name | | Upfront cost | | Monthly cost | | Description | | Region | | Config Summary | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | AWS Data Pipeline | ✎ | 0.00 USD | | 117.00 USD | | Data Pipeline | | EU (Ireland) | | Number of high ... | |
| ☐ | AWS Data Transfer | ✎ | 0.00 USD | | 40.96 USD | | Data Transfer | | EU (Ireland) | | DT Inbound: All ... | |
| ☐ | AWS Glue | ✎ | 0.00 USD | | 3.25 USD | | Glue | | EU (Ireland) | | Number of DPU... | |
| ☐ | AWS Key Management Service | ✎ | 0.00 USD | | 11.00 USD | | KMS | | EU (Ireland) | | Number of cust... | |
| ☐ | AWS Lambda | ✎ | 0.00 USD | | 0.00 USD | | Lambda | | EU (Ireland) | | Architecture (x8... | |
| ☐ | Amazon Athena | ✎ | 0.00 USD | | 74.27 USD | | Athena | | EU (Ireland) | | Total number of... | |
| ☐ | Amazon CloudWatch | ✎ | 0.00 USD | | 1.50 USD | | Cloudwatch | | EU (Ireland) | | Number of Metri... | |
| ☐ | Amazon EC2 | ✎ | 0.00 USD | | 71.12 USD | | EC2 | | EU (Ireland) | | Operating syste... | |
| ☐ | Amazon QuickSight | ✎ | 0.00 USD | | 1,618.00 USD | | Quicksight | | EU (Ireland) | | Number of work... | |
| ☐ | Amazon RDS Custom for SQL Server | ✎ | 0.00 USD | | 2,412.36 USD | | RDS | | EU (Ireland) | | Storage for each... | |
| ☐ | Amazon Redshift | ✎ | 0.00 USD | | 4,094.81 USD | | Redshift | | EU (Ireland) | | Nodes ( 1 instan... | |
| ☐ | Amazon Simple Notification Service (SNS) | ✎ | 0.00 USD | | 0.00 USD | | SNS | | EU (Ireland) | | | |
| ☐ | Amazon Simple Storage Service (S3) | ✎ | 0.00 USD | | 21.76 USD | | S3 | | EU (Ireland) | | S3 Standard sto... | |
| ☐ | Amazon Virtual Private Cloud (VPC) | ✎ | 0.00 USD | | 644.46 USD | | VPC | | EU (Ireland) | | Working days pe... | |

## The Exterminators

## ETL Pipelining + Dashboard + Error Logging

**UPLOAD DATA**

CSV File
SQL Query

- •
- •

# The Exterminators

## ETL Pipelining + Dashboard + Error Logging

# <THANK YOU>