# Capstone 1 Project

> Presentation

Presented by **Ahmed Alkesayer** ,**Rahaf mohammed**, and **Haneen Alghamdi**

# Agenda

## Capstone 1 Project

| | |
|---|---|
| **01** | Project Objects |
| **02** | Data Overview |
| **03** | Project Steps |
| **04** | Project Achievement |
| **05** | Future |

WeCloud**Data**

# Project Objectives

- Project Requirements

Tasks include data gathering, storage, transformation, and preparation for BI.

- Tools Used

AWS S3, AWS RDS, Lambda, Airbyte, EC2, Docker, Snowflake, Metabase.
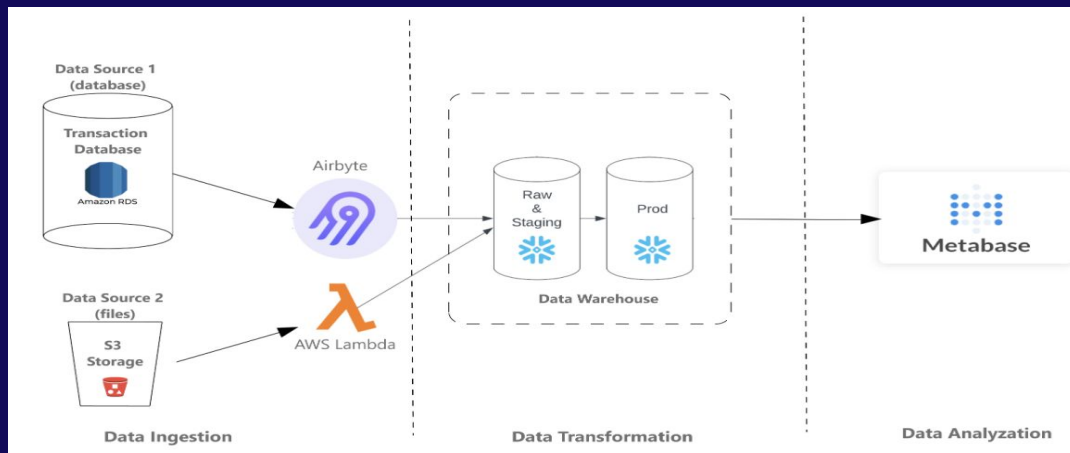


Amazon S3

AWS Lambda

Amazon EC2

snowflake

Amazon RDS

Airbyte

docker

Metabase

WeCloudData

# Project Objectives

- **Project Process Design**

Structured steps: requirements, data model, infrastructure, data processing, dashboards.



- **Project Work Split**

Independent work: everyone involved in all phases, with open communication.

# Data Overview

- The Tables Included

- Data Cleansing and Processing

- Continuous Data Loading vs. One-time Loading

| Fact tables | Dimension tables |
|---|---|
| Catalog_Sales | Date_Dim |
| Web_Sales | Customer |
| Inventory | Item |
| | Promotion |
| | Customer_Demographics |
| | Call_Center |
| | Customer_Address |
| | Catalog_Page |
| | Warehouse |
| | Time_Dim |
| | Ship_Mode |
| | Household_Demographics |
| | Income_Band |
| | Web_page |
| | Web_Site |
| | |

# Project Steps

- The processes of making our project is entirely composed of the whole ETL process for us as a data engineers and it will be explained in the coming slides

# Project Steps: Data Infrastructure

**[Week 1, Day 1]**
Creating the database TPCDS, and within it is the RAW schema which will store the ingested data from extraction, and within the schema will be the creation of the Inventory Table that would later be used to get the data from the RDS source.
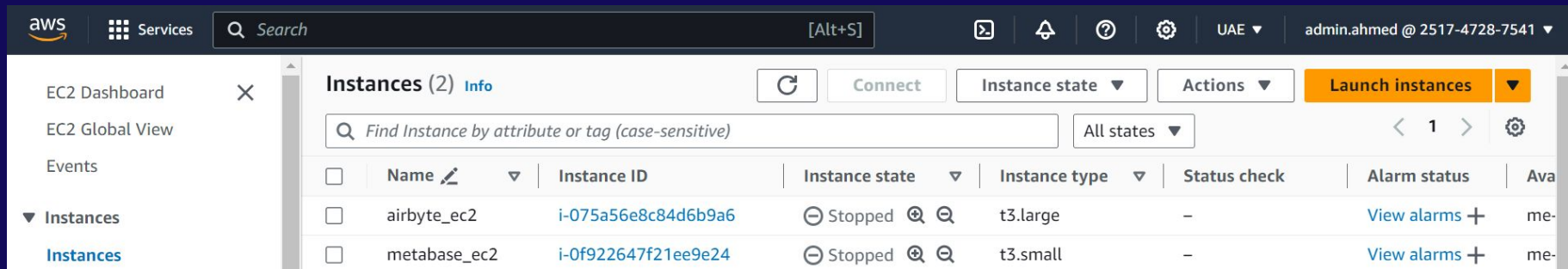
```
1  CREATE DATABASE TPCDS;
2
3  USE TPCDS;
4
5  CREATE SCHEMA RAW;
6
7  USE SCHEMA RAW;
8
9  CREATE TABLE inventory(
10     INV_DATE_SK NUMBER,
11     INV_ITEM_SK NUMBER,
12     INV_QUANTITY_ON_HAND NUMBER,
13     INV_WAREHOUSE_SK NUMBER
14 );
```

WeCloudData

# Project Steps: Data Infrastructure

***[Week 1, Day 1]***
Now, Into Configuring the EC2 instances for both **Airbyte** *(t3.large)* and **Metabase** *(t3.small)* as both are tools that are needed to be deployed on cloud for data processing (Airbyte) and production (Metabase).
Also installing + configuring **Docker** to have the ability to use these tools.

# Project Steps: Data Infrastructure

**[Week 1, Day 2]**
For this task, the main goal was to create **AWS Lambda** to extract the inventory table data from the **AWS S3** bucket that has already been provided by WCD and the AWS Lambda is triggered by the **AWS EventBridge**, the loading destination will be the Snowflake Warehouse where the TCPDS database is, and inside the RAW schema.

# Cont'd

*[Week 1, Day 2]*
Code Snippet of the **AWS Lambda** handler.

```python
import os
import boto3
import requests
import snowflake.connector as sf


def lambda_handler(event, context):

    url = 'https://de-materials-tpcds.s3.ca-central-1.amazonaws.com/inventory.csv'
    destination_folder = '/tmp'
    file_name = 'inventory.csv'
    local_file_path = '/tmp/inventory.csv'

    # Snowflake connection parameters
    account = os.environ['account']
    warehouse = os.environ['warehouse']
    database = os.environ['database']
    schema = os.environ['schema']
    table = os.environ['table']
    user = os.environ['user']
    password = os.environ['password']
    role= os.environ['role']
    stage_name = os.environ['stage_name']

    # Download the data from the API endpoint
    response = requests.get(url)
    response.raise_for_status()


    # Save the data to the destination file in /tmp directory
    file_path = os.path.join(destination_folder, file_name)
    with open(file_path, 'wb') as file:
        file.write(response.content)

    with open(file_path, 'r') as file:
        file_content = file.read()
        print("File Content:")
        print(file_content)



    # Establish Snowflake connection
    conn = sf.connect(user = user, password = password, \
                account = account, warehouse=warehouse, \
                database=database,  schema=schema,  role=role)


    cursor = conn.cursor()

    # use schema
    use_schema = f"use schema {schema};"
    cursor.execute(use_schema)

    # create CSV format
    create_csv_format = f"CREATE or REPLACE FILE FORMAT COMMA_CSV TYPE ='CSV' FIELD_DELIMITER = ',';"
    cursor.execute(create_csv_format)


    create_stage_query = f"CREATE OR REPLACE STAGE {stage_name} FILE_FORMAT =COMMA_CSV"
    cursor.execute(create_stage_query)

    # Copy the file from local to the stage
    copy_into_stage_query = f"PUT 'file://{local_file_path}' @{stage_name}"
    cursor.execute(copy_into_stage_query)

    # List the stage
    list_stage_query = f"LIST @{stage_name}"
    cursor.execute(list_stage_query)

    # truncate table
    truncate_table = f"truncate table {schema}.{table};"
    cursor.execute(truncate_table)


    # Load the data from the stage into a table (example)
    copy_into_query = f"COPY INTO {schema}.{table} FROM @{stage_name}/{file_name} FILE_FORMAT =COMMA_CSV;"
    cursor.execute(copy_into_query)


    print("File uploaded to Snowflake successfully.")


    return {
        'statusCode': 200,
        'body': 'File downloaded and uploaded to Snowflake successfully.'
    }
```

WeCloud**Data**

# Cont'd

**[Week 1, Day 2]**
To prove the process of successful table loading to **Snowflake** after the **AWS EventBridge** trigger at every 3 AM in the morning, here's a snapshot of the **CloudWatch**.

| | | |
|---|---|---|
| ▶ | 2024-03-28T02:03:44.700+03:00 | File uploaded to Snowflake successfully. |
| ▶ | 2024-03-28T02:03:44.755+03:00 | END RequestId: 256604a4-f091-4e06-aea0-b0eaef5c76c5 |
| ▶ | 2024-03-28T02:03:44.755+03:00 | REPORT RequestId: 256604a4-f091-4e06-aea0-b0eaef5c76c5 Duration: 184974.42 ms Billed Duration: 184975 ms Memory Size: 512 MB Max Memo… |

# Project Steps: Data Infrastructure

**[Week 1, Day 3]**
As Day 3 begun, we've came to the final task in setting up the infrastructure, which is using **Airbyte** to take the data from the **AWS RDS (PostgreSQL)** which is already been provided by WCD to extract the other tables along with its skeleton which means we don't have to define those tables like we did with inventory and **Snowflake** gets synced with the RDS to take the tables there.

WeCloud**Data**

# Project Steps: Data Modeling

**[Week 2 ]**
The tasks of the 2$^{nd}$ week is all about the modeling of our data to be ready for the ETL process, and answering the business requirements for it it to be ready for analytics with Metabase.

WeCloud**Data**

# Data Modeling (cont'd)

*[Week 2 ]*
Code snippet for data modeling

```
1  -- Creating this schema to process the ingested data on RAW schema and staging it for the production
2  CREATE OR REPLACE SCHEMA INTERMEDIATE;
3
4  -- Creating Customer Snapshot table which will store the data from all the customer tables
5  -- as a snapshot to be processed yo Customer Dim
6  CREATE OR REPLACE TABLE TPCDS.INTERMEDIATE.CUSTOMER_SNAPSHOT (
7      C_SALUTATION VARCHAR(16777216),
8      C_PREFERRED_CUST_FLAG VARCHAR(16777216),
9      C_FIRST_SALES_DATE_SK NUMBER(38,0),
10     C_CUSTOMER_SK NUMBER(38,0),
11     C_LOGIN VARCHAR(16777216),
12     C_CURRENT_CDEMO_SK NUMBER(38,0),
13     C_FIRST_NAME VARCHAR(16777216),
14     C_CURRENT_HDEMO_SK NUMBER(38,0),
15     C_CURRENT_ADDR_SK NUMBER(38,0),
16     C_LAST_NAME VARCHAR(16777216),
17     C_CUSTOMER_ID VARCHAR(16777216),
18     C_LAST_REVIEW_DATE_SK NUMBER(38,0),
19     C_BIRTH_MONTH NUMBER(38,0),
20     C_BIRTH_COUNTRY VARCHAR(16777216),
21     C_BIRTH_YEAR NUMBER(38,0),
22     C_BIRTH_DAY NUMBER(38,0),
23     C_EMAIL_ADDRESS VARCHAR(16777216),
24     C_FIRST_SHIPTO_DATE_SK NUMBER(38,0),
25     START_DATE TIMESTAMP_NTZ(9),
26     END_DATE TIMESTAMP_NTZ(9)
27 );
```

```
1  -- Now for this schema which will be the main engine for the METABASE as it defines the main fact table
2  CREATE OR REPLACE SCHEMA ANALYTICS;
3
4
5  -- Now this is the creation of the Customer_Dim
6  -- which will store the data aggregated and joined from the customer snapshots
7  create or replace TABLE TPCDS.ANALYTICS.CUSTOMER_DIM (
8      C_SALUTATION VARCHAR(16777216),
9      C_PREFERRED_CUST_FLAG VARCHAR(16777216),
10     C_FIRST_SALES_DATE_SK NUMBER(38,0),
11     C_CUSTOMER_SK NUMBER(38,0),
12     C_LOGIN VARCHAR(16777216),
13     C_CURRENT_CDEMO_SK NUMBER(38,0),
14     C_FIRST_NAME VARCHAR(16777216),
15     C_CURRENT_HDEMO_SK NUMBER(38,0),
16     C_CURRENT_ADDR_SK NUMBER(38,0),
17     C_LAST_NAME VARCHAR(16777216),
18     C_CUSTOMER_ID VARCHAR(16777216),
19     C_LAST_REVIEW_DATE_SK NUMBER(38,0),
20     C_BIRTH_MONTH NUMBER(38,0),
21     C_BIRTH_COUNTRY VARCHAR(16777216),
22     C_BIRTH_YEAR NUMBER(38,0),
23     C_BIRTH_DAY NUMBER(38,0),
24     C_EMAIL_ADDRESS VARCHAR(16777216),
25     C_FIRST_SHIPTO_DATE_SK NUMBER(38,0),
26     CA_STREET_NAME VARCHAR(16777216),
27     CA_SUITE_NUMBER VARCHAR(16777216),
28     CA_STATE VARCHAR(16777216),
29     CA_LOCATION_TYPE VARCHAR(16777216),
30     CA_COUNTRY VARCHAR(16777216),
31     CA_ADDRESS_ID VARCHAR(16777216),
32     CA_COUNTY VARCHAR(16777216),
33     CA_STREET_NUMBER VARCHAR(16777216),
34     CA_ZIP VARCHAR(16777216),
35     CA_CITY VARCHAR(16777216),
36     CA_GMT_OFFSET FLOAT,
37     CD_DEP_EMPLOYED_COUNT NUMBER(38,0),
38     CD_DEP_COUNT NUMBER(38,0),
39     CD_CREDIT_RATING VARCHAR(16777216),
40     CD_EDUCATION_STATUS VARCHAR(16777216),
41     CD_PURCHASE_ESTIMATE NUMBER(38,0),
42     CD_MARITAL_STATUS VARCHAR(16777216),
43     CD_DEP_COLLEGE_COUNT NUMBER(38,0),
44     CD_GENDER VARCHAR(16777216),
45     HD_BUY_POTENTIAL VARCHAR(16777216),
46     HD_DEP_COUNT NUMBER(38,0),
47     HD_VEHICLE_COUNT NUMBER(38,0),
48     HD_INCOME_BAND_SK NUMBER(38,0),
49     IB_LOWER_BOUND NUMBER(38,0),
50     IB_UPPER_BOUND NUMBER(38,0),
51     START_DATE TIMESTAMP_NTZ(9),
52     END_DATE TIMESTAMP_NTZ(9)
53 );
```

WeCloudData

# Data Modeling (cont'd)

Code snippet for data modeling

```sql
1  create or replace TABLE TPCDS.INTERMEDIATE.DAILY_AGGREGATED_SALES (
2      WAREHOUSE_SK NUMBER(38,0),
3      ITEM_SK NUMBER(38,0),
4      SOLD_DATE_SK NUMBER(38,0),
5      SOLD_WK_NUM NUMBER(38,0),
6      SOLD_YR_NUM NUMBER(38,0),
7      DAILY_QTY NUMBER(38,0),
8      DAILY_SALES_AMT FLOAT,
9      DAILY_NET_PROFIT FLOAT
10 );
```

WeCloudData

# Data Modeling (cont'd)

## [Week 2 ]
Code snippet for data modeling

```sql
1  -- This is the main fact table which will be the analytical engine
2  create or replace TABLE TPCDS.ANALYTICS.WEEKLY_SALES_INVENTORY (
3      WAREHOUSE_SK NUMBER(38,0),
4      ITEM_SK NUMBER(38,0),
5      SOLD_WK_SK NUMBER(38,0),
6      SOLD_WK_NUM NUMBER(38,0),
7      SOLD_YR_NUM NUMBER(38,0),
8      SUM_QTY_WK NUMBER(38,0),
9      SUM_AMT_WK FLOAT,
10     SUM_PROFIT_WK FLOAT,
11     AVG_QTY_DY NUMBER(38,6),
12     INV_QTY_WK NUMBER(38,0),
13     WKS_SPLY NUMBER(38,6),
14     LOW_STOCK_FLG_WK BOOLEAN
15 );
```

WeCloudData

# Project Steps: ETL

*[Week 2 ]*
The ETL is achieved through creating **stored procedures** and applying tasks to those procedures so data gets updated every time new data added to the Snowflake.

WeCloud**Data**

# ETL (cont'd)

*[Week 2 ]*
Code snippet for creating a stored procedure and that one for the **weekly sales** that is crucial and the engine of the analytics



WeCloudData

# ETL (cont'd)

**[Week 2]**
Code snippet for the tasks which will be the key to automating the performing of ETL by calling the stored procedures and updating the tables based on the CRON time given

```sql
1  CREATE OR REPLACE TASK tpcds.intermediate.creating_daily_aggregated_sales_incrementally
2      WAREHOUSE = COMPUTE_WH
3      SCHEDULE = 'USING CRON * 8 * * * UTC'
4      AS
5  CALL populating_daily_aggregated_sales_incrementally();
6
7  CREATE OR REPLACE TASK tpcds.analytics.creating_weekly_aggregated_sales_incrementally
8      WAREHOUSE = COMPUTE_WH
9      SCHEDULE = 'USING CRON * 8 * * * UTC'
10     AS
11 CALL populating_weekly_aggregated_sales_incrementally();
12
13 ALTER TASK tpcds.intermediate.creating_daily_aggregated_sales_incrementally RESUME;
14
15 EXECUTE TASK tpcds.intermediate.creating_daily_aggregated_sales_incrementally;
16
17
18 ALTER TASK tpcds.analytics.creating_weekly_aggregated_sales_incrementally RESUME;
19
20 EXECUTE TASK tpcds.analytics.creating_weekly_aggregated_sales_incrementally;
21
```

WeCloudData

# Project Steps:

# Final Schema

# Project Steps: Metabase

*[Week 2 ]*

For the **Metabase** and the final task of the capstone project, the EC2 instance needed for the analytics tool that would allow the visual representation of the data finding, we used the created **t3.small** instance, as **Docker** is installed and configured to run *Metabase* as a container, and then connecting it to our Snowflake warehouse to sync the data and start answering the analytical tool.

# Project Achievement

**DASHBOARD**

# Dashboard

## Detect items experiencing low sto...

| ^ ITEM_SK | LOW_STOCK_FLAG ^ |
|---|---|
| 10,240 | true |
| 13,033 | true |
| 13,966 | true |
| 2,612 | true |
| 16,288 | true |

Rows 1-5 of first 2000  ‹  ›

## the amount of items which has lo...

| ^ COUNT(ITEM_SK) | LOW_STOCK_FLAG |
|---|---|
| 1,214,866 | true |

## Displaying items with low supply levels for each week.



## Identify the highest and lowest performing items of the week by analy...

● TOTAL_SALES_AMOUNT   ● TOTAL_QUANTITY_SOLD

# Future

- If we had more time we will clean data to find perfect result(especially customer information)
- For further develop the project we are thinking about making Predictive Analytics Model , include predictive analytics by building machine learning models
- For further research that we want to do on the project in the future is Data Governance compliance , such as data protection mechanisms

Thank you!

WeCloudData