# WeCloudData

## Data Engineering Bootcamp
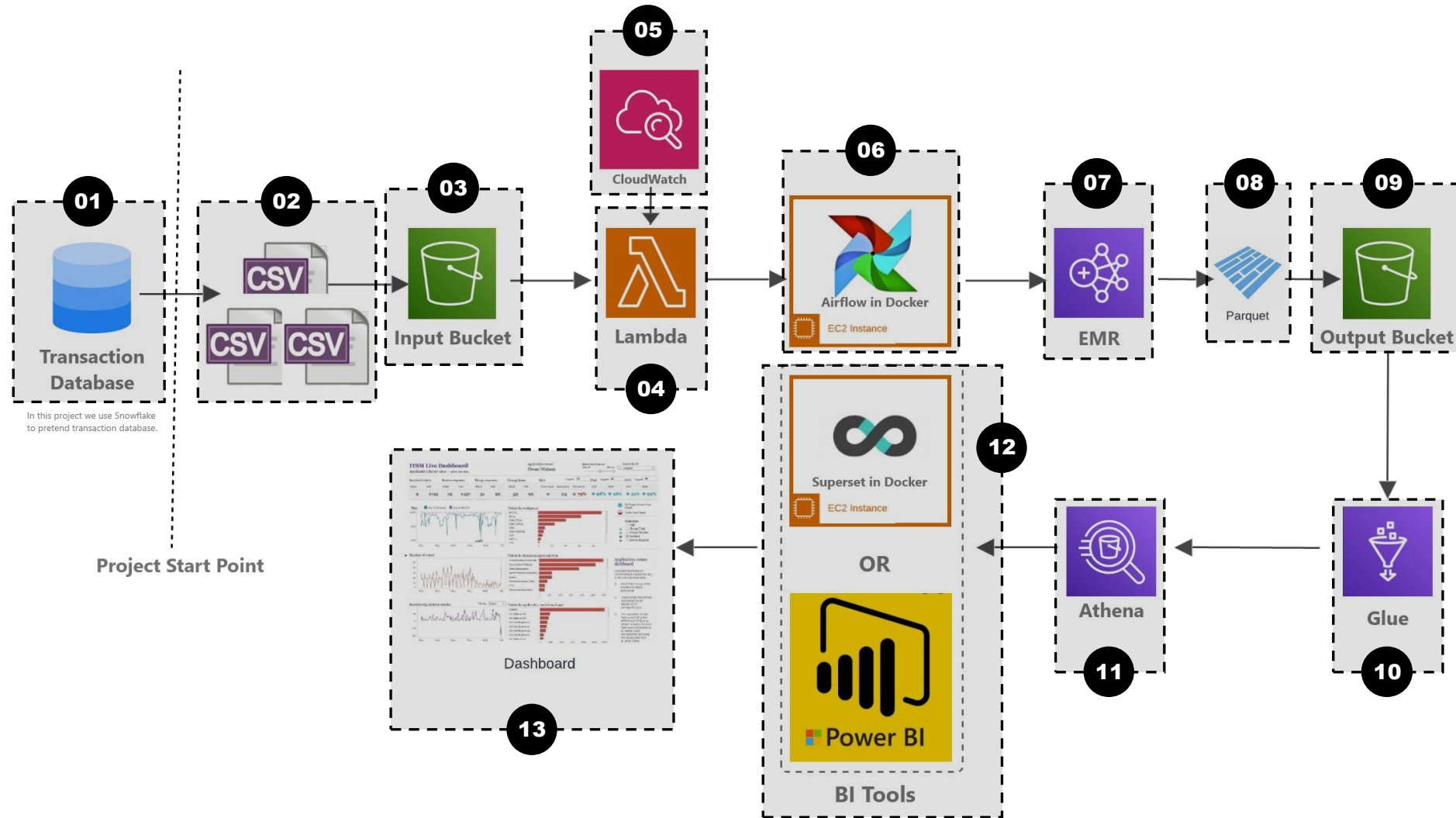
Chichi Aguoji

# Architecture 1



Transaction Database

In this project we use Snowflake to pretend transaction database.

Project Start Point

CSV

Input Bucket

CloudWatch

Lambda

Airflow in Docker
EC2 Instance

EMR

Parquet

Output Bucket

Superset in Docker
EC2 Instance

OR

Power BI

BI Tools

Athena

Glue

Dashboard

# Architecture 1



**01** Transaction Database

In this project we use Snowflake to pretend transaction database.

**02** CSV CSV CSV

**03** Input Bucket

**04** Lambda

**05** CloudWatch

**06** Airflow in Docker — EC2 Instance

**07** EMR

**08** Parquet

**09** Output Bucket

**10** Glue

**11** Athena

**12** BI Tools — Superset in Docker (EC2 Instance) OR Power BI

**13** Dashboard — ITSM Live Dashboard

Project Start Point

# Transaction Database

## Challenges

- Zero padded dates in Javascript

## Modifications

- Adjusted the stored procedure to include:

```
${"0"+ n.getDate().slice(-2) }
```

- This adds a "0" in the front of the date then uses the slice function to get a substring of the last 2 digits in the date.

# Data Pull

## Challenges

- None

## Modifications

- None

Input Bucket

![Input Bucket]

## Challenges

- Daily data being appended into S3 bucket

## Modifications

- Created another S3 bucket to store historical data.

- Used Lambda function to move data into this separate bucket.

```
lambda_function ×      Execution results ×      ⊕

1  import boto3
2  from datetime import datetime, timedelta
3
4  DESTINATION_BUCKET = 'wcd-midterm-raw-archive'
5  SOURCE_BUCKET = 'wcd-midterm-raw-chichi'
6
7  s3_client=boto3.client('s3')
8  s3_resource = boto3.resource('s3')
9
0
1  def lambda_handler(event,content):
2
3      for object in s3_client.list_objects_v2(Bucket='wcd-midterm-raw-chichi')['Contents']:
4
5          s3_client.copy_object(
6              Bucket=DESTINATION_BUCKET,
7              Key=object['Key'],
8              CopySource={'Bucket':SOURCE_BUCKET, 'Key':object['Key']}
9          )
0
1          s3_client.delete_object(Bucket=SOURCE_BUCKET, Key=object['Key'])
2
3      print("objects copied to archive folder and deleted")
4
```

# Lambda

**Lambda**

## Challenges

- Issues with the run date being 1 day after the file date

- Call to subprocess kept failing

## Modifications

- Attempted to use datetime library to use timedelta to get previous day's date, but it wasn't working so ended up using substrings. Substrings logic broke during testing when new month began so used the following:

```python
today = date.today()
yesterday = today-timedelta(1)
```

- Didn't realize the content-type was entered incorrectly, entered 'application/json/' instead of 'application/json'

**CloudWatch**

# Cloud Watch

## Challenges

- None

## Modifications

- None

# Airflow in Docker

Airflow in Docker
EC2 Instance

## Challenges

- Updating EC2 IP address after every start

- Unhelpful Logs

- Having to leave EC2 instance running in order to test

## Modifications

- Set up Elastic IP Address for EC2 instance

- Manually Timed Errors

- Set up 2 additional Lambda functions to start and stop EC2 instance each night to test the data drop

```
lambda_function ×          Execution results ×

import boto3
region = 'us-east-2'
instances = ['i-009aaa8206ad6988f']
ec2 = boto3.client('ec2', region_name=region)

def lambda_handler(event, context):
    ec2.start_instances(InstanceIds=instances)
    print('started your instances: ' + str(instances))
```

```
   lambda_function ×          Execution results ×

1  import boto3
2  region = 'us-east-2'
3  instances = ['i-009aaa8206ad6988f']
4  ec2 = boto3.client('ec2', region_name=region)
5
6  def lambda_handler(event, context):
7      ec2.stop_instances(InstanceIds=instances)
8      print('stopped your instances: ' + str(instances))
9
```

**EMR**

EMR

## Challenges

- Manual Start

- Configuring access to Airflow

- Testing PySpark code in Databricks w/o a local run

## Modifications

- Added a step to start EMR programmatically

- Trial/Error, Google, Rewatching Lectures

- Ended up having to run code on EMR in chunks until found errors

# Parquet

Parquet

## Challenges

- Understanding partitions labeling

- Partition amount and number of partitions not aligning

## Modifications

- None

- None

**Output Bucket**

Output Bucket

## Challenges

- Using folder structure to facilitate BI analysis

## Modifications

- None

**Glue**

**Glue**

## Challenges

- None

## Modifications

- None

**Athena**

## Challenges

- None

## Modifications

- None

Power BI

Superset in Docker
EC2 Instance

# BI Tools

## Challenges

- Connecting PowerBI to Athena

## Modifications

- Leveraged ODBC Connection and AWS Secret Access Key

# Dashboard

Dashboard

## Challenges

- Having questions about the data that the data couldn't answer

- No Automatic Refresh of data

- Wanting to use visualizations not supported by the data

## Modifications

- Only visualized questions the data could answer

- None

- Created supporting dummy data