**Week 3 case study Report**

**Name:** Balu Lakshmi sudha

**Introduction:**

This report explains how I designed and built a complete data pipeline to process, validate, transform, mask, and store customer transactions data using Spark Structured Streaming, Delta Lake, Great Expectations, and Apache NiFi.

# 1. Project Overview

## Goal:

Build an end-to-end data pipeline that:

Streams synthetic transaction data from a Kafka/Event Hub producer.
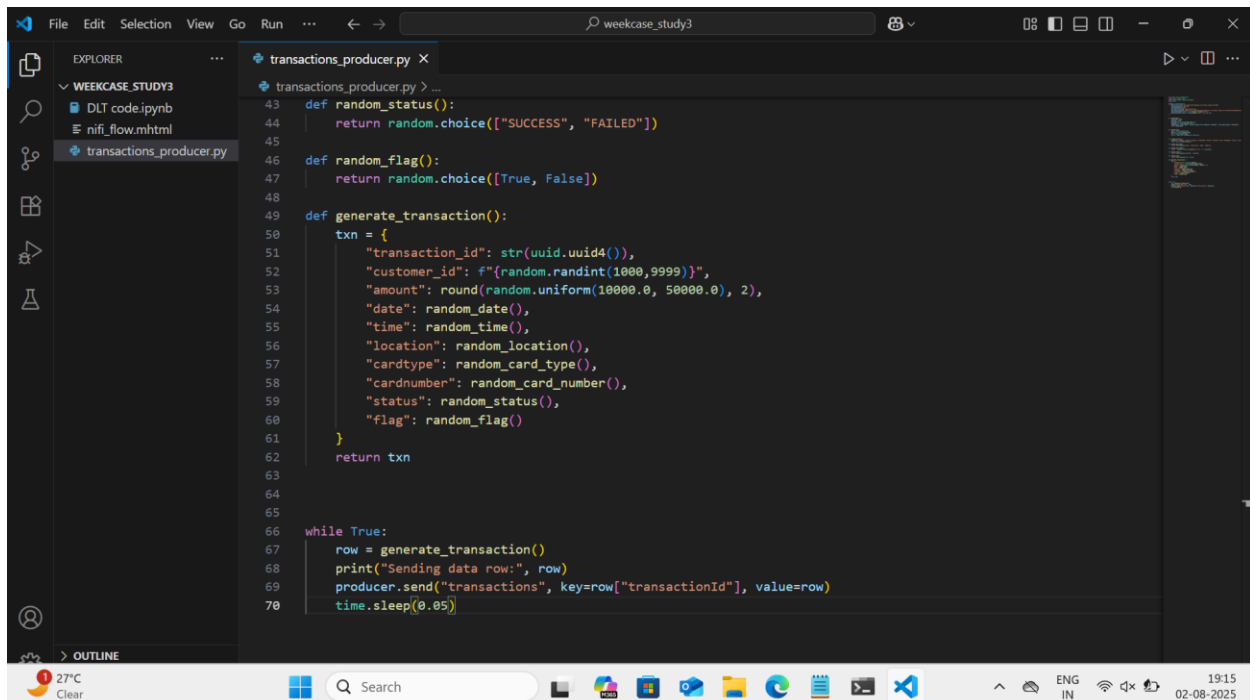
Stores raw data in Bronze Delta tables.

Validates and cleans the data using Great Expectations → Silver tables.

Applies masking, enrichment, business rules → Gold tables.

Uses Apache NiFi for additional transformation and masking.

**Data:** Simulated transactions with transaction_id, customer_id, amount, date, time, location, cardtype, cardnumber, status, flag

# 2. continuous streaming data generation



I Wrote a transactions_producer.py:

Uses uuid for transaction_id

random.randint for customer_id, amount

random_date and random_time to create date between (2015–2024) and time in (24-hour format)

random.choice for location (10 major Indian cities)

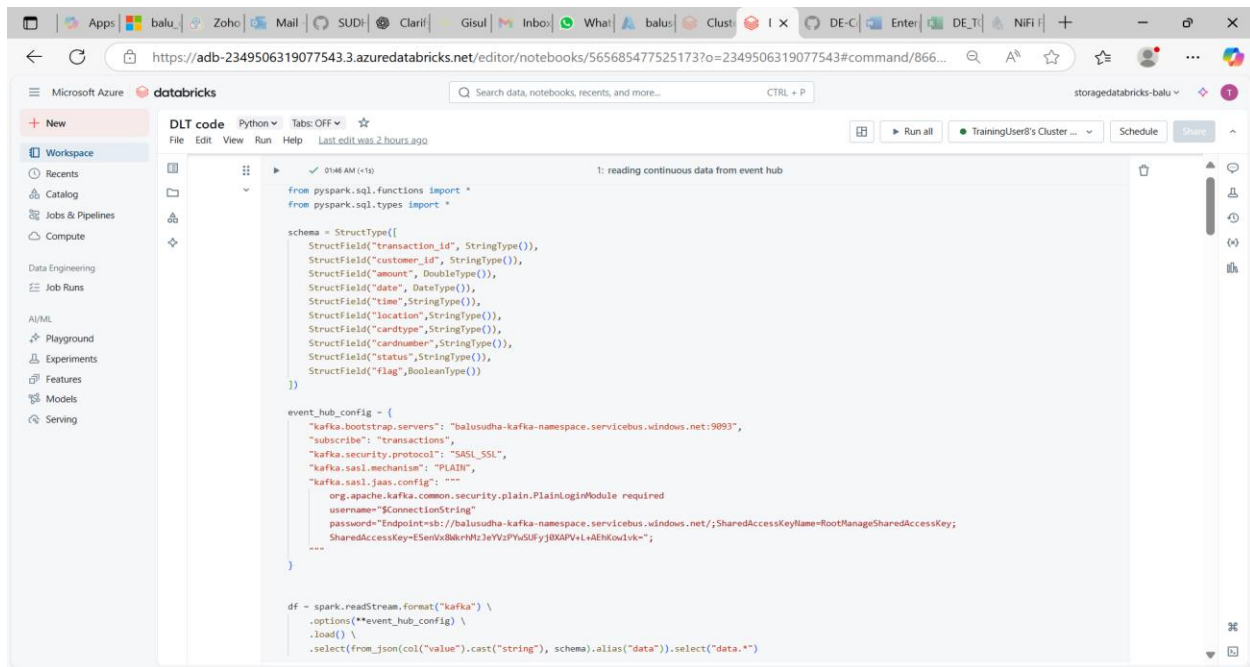random.choice for cardtype (VISA, MasterCard, AMEX, RUPAY)

cardnumber generated as random 16-digit string

status and flag set randomly.

Sent data to transactions Event Hub.

## 3.consumer

Consumes the streaming transaction data coming from the Event Hub (configured as a Kafka-compatible source) and parses it into structured format for further processing.



## 4. Data Ingestion (Bronze)

Used Spark Structured Streaming to connect to the Event Hub (Kafka-compatible).

Configured .readStream.format("kafka") with Event Hub connection settings.

Defined a schema and used from_json to parse the incoming JSON messages.

Extracted required fields from the JSON value column.

Stored the raw streaming data into a Delta Lake table: /delta_transactions.



## 5. Data Validation on Bronze Layer (Great Expectations)

Connected Bronze Delta table using Spark to load the raw streaming data.

Applied Great Expectations (GE) to check data quality on the ingested raw data.

## Defined rules:

Checked that transaction_id is not null and unique.

Verified customer_id, amount, and location are not null.

Ensured cardtype values are within an allowed set (Visa, MasterCard, Amex, RuPay).

Verified cardnumber length is exactly 16 digits.

Ran validation using SparkDFDataset and validate() method.

## 5. Data Transformation and Masking (Silver Layer)

Read Bronze data from Delta Lake.

Filtered valid records based on Great Expectations results.

Performed PII masking on the cardnumber column:

    Kept the last 4 digits visible.

    Replaced first 12 digits with ************ using a concat and substring function.

Stored the cleaned, masked data in the Delta table: /silver_transactions

# 6. Data Aggregation and Analytics (Gold Layer)

Read Silver table data for trusted, cleaned transactions.

Performed aggregation:

 Grouped data by customer_id.

 Calculated the total amount spent by each customer.

Generated customer-level insights for reporting and downstream analytics.

Stored the aggregated data in the Delta table: /gold/transactions_summary.

## 7. Data Processing using Apache NiFi

Used Apache NiFi to handle additional filtering and record-level transformations.

GetFile Processor: Picked up the source CSV file (Silver data).

QueryRecord Processor:

Applied custom SQL-like filters on the ingested records.

Example: Filtered transactions based on conditions such as status = 'SUCCESS' or specific cardtype.
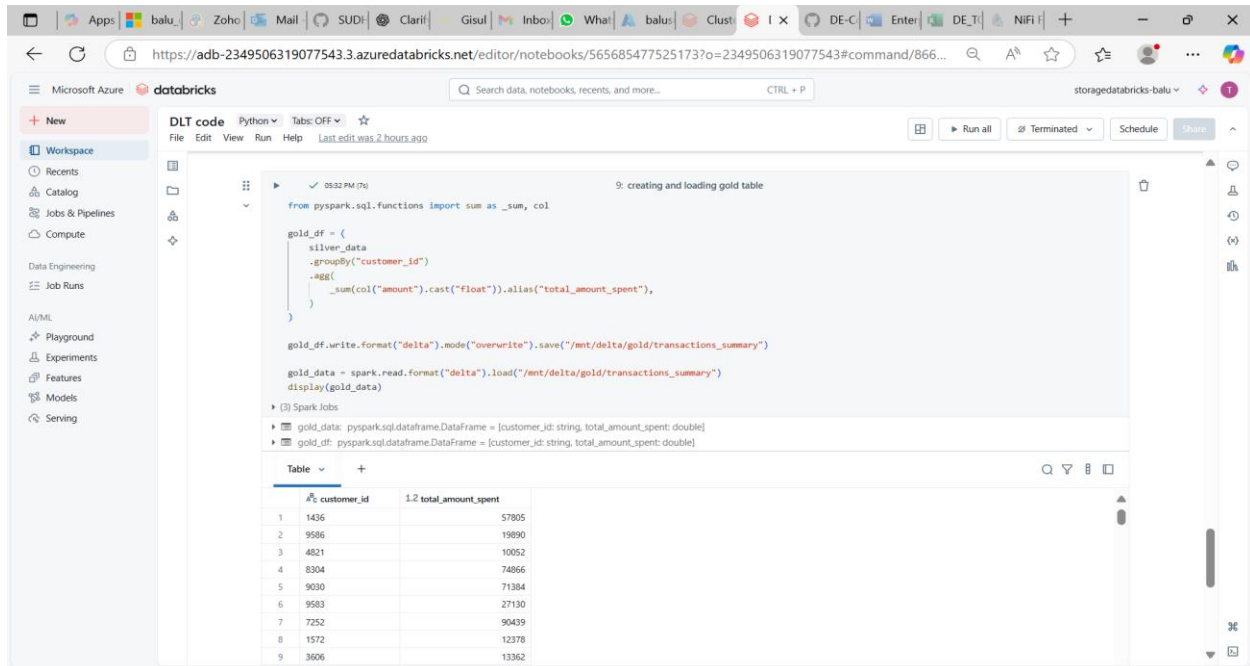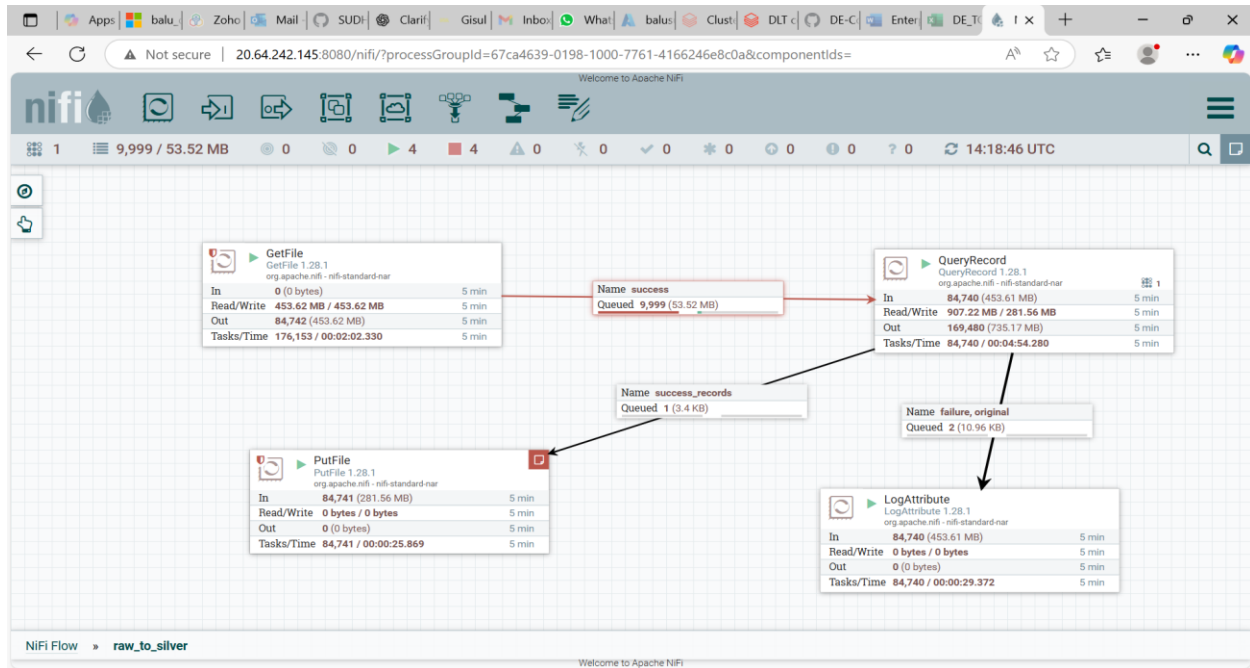
Used Record Reader/Writer with schema for structured processing.

PutFile Processor: Stored the processed (filtered) records into a target directory.

LogAttribute Processor: Captured logs for flow debug and verification — used to inspect records that failed processing.

## Conclusion:

In this project, I built a complete data pipeline from raw data to final insights.

I ingested streaming data from Event Hub using Spark Structured Streaming.

I validated the raw data with Great Expectations to ensure data quality.

I filtered and masked sensitive information and stored the clean data in Silver tables.

I aggregated the Silver data to create Gold tables for business analysis.

Finally, I used Apache NiFi to filter the processed data further using the QueryRecord processor and saved the filtered results to files for additional downstream use.

This end-to-end flow helps ensure the data is clean, secure, and ready for reporting and analytics.