## Key KPIs:

1. Total Revenue: Sum of the Amount or Total_Amount column.
2. Total Transactions: Count of unique Transaction_ID.
3. Average Transaction Value: Total_Amount / Total_Purchases or Total_Amount / count of Transaction_ID.
4. Average Customer Age: Average of the Age column.
5. Customer Retention Rate: Percentage of returning customers (those with multiple Transaction_ID entries) relative to all customers.
6. Product Popularity: Count of purchases by Product_Category or Product_Brand to find the most popular items.
7. Order Fulfillment Rate: Percentage of Order_Status marked as "Completed" relative to total orders.
8. Customer Satisfaction: Average of Ratings column or analysis of the Feedback column if it's qualitative.
9. Revenue by Customer Segment: Sum of Amount for each Customer_Segment.

## GitHub Link:

https://github.com/kumkumbaswal003/Cloudthat_Project.git
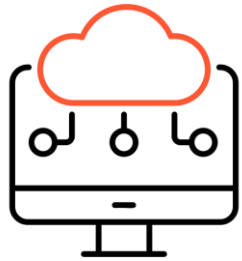
# AZURE DATA PIPELINE

### For TechRetail

**Group 4:**
- Arnab Saha
- Debadrita Acharjee
- Kumkum Baswal

# Background:

- **TechRetail**, a mid-sized retail company, wants to create a data pipeline retail data from various sources, process it using advanced analytics, and visualize the results in a dashboard. The goal is to gain insights into sales trends and improve decision-making. The company wants to leverage Azure Databricks for data processing and Microsoft Fabric for data integration and visualization.

# Objectives:



**Data Ingestion**

Azure Data Factory (ADF)

ADLS Gen2

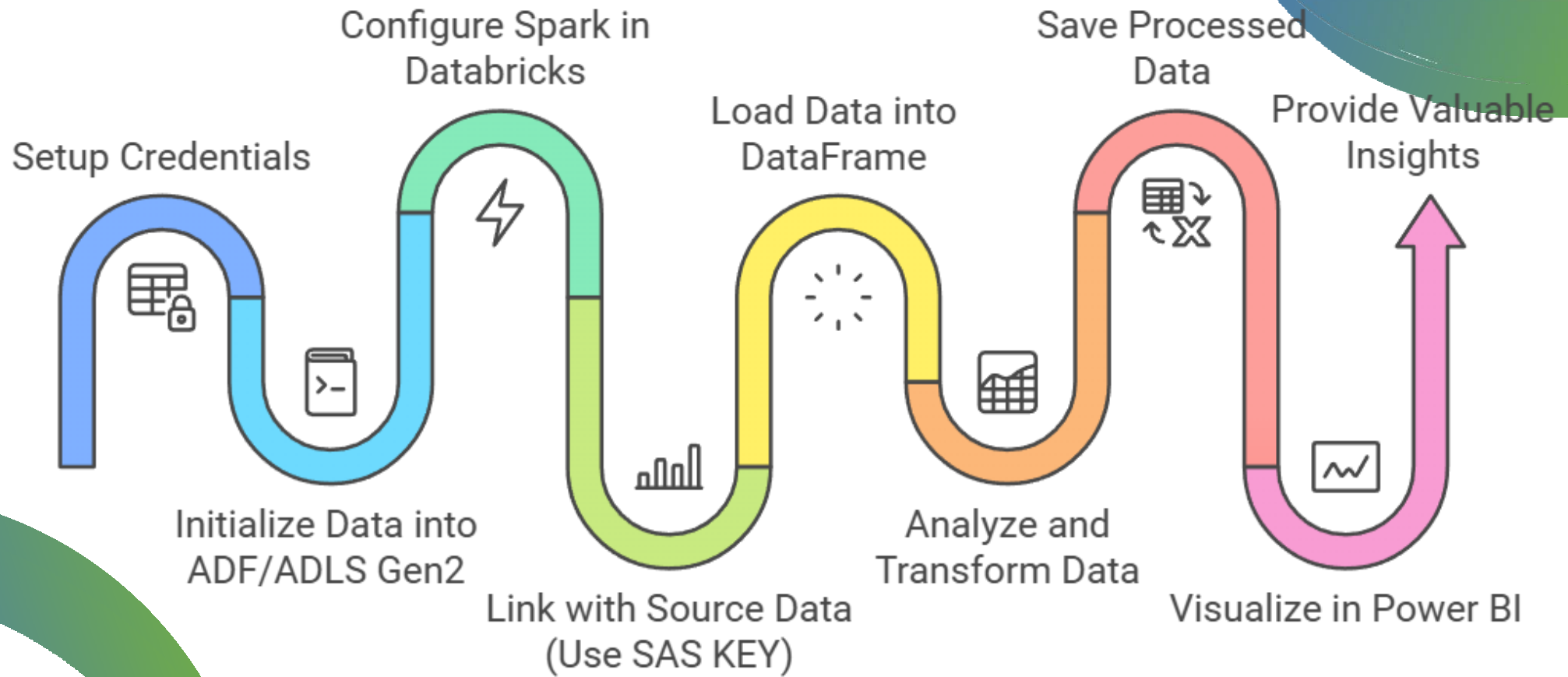**Data Processing**

Synapse Data Engineering

Databricks

**Data Storage**

ADLS Delta Table

**Data Visualization**

Microsoft PowerBI

# Architectural Process Flow:



Setup Credentials

Configure Spark in Databricks

Load Data into DataFrame

Save Processed Data

Provide Valuable Insights

Initialize Data into ADF/ADLS Gen2

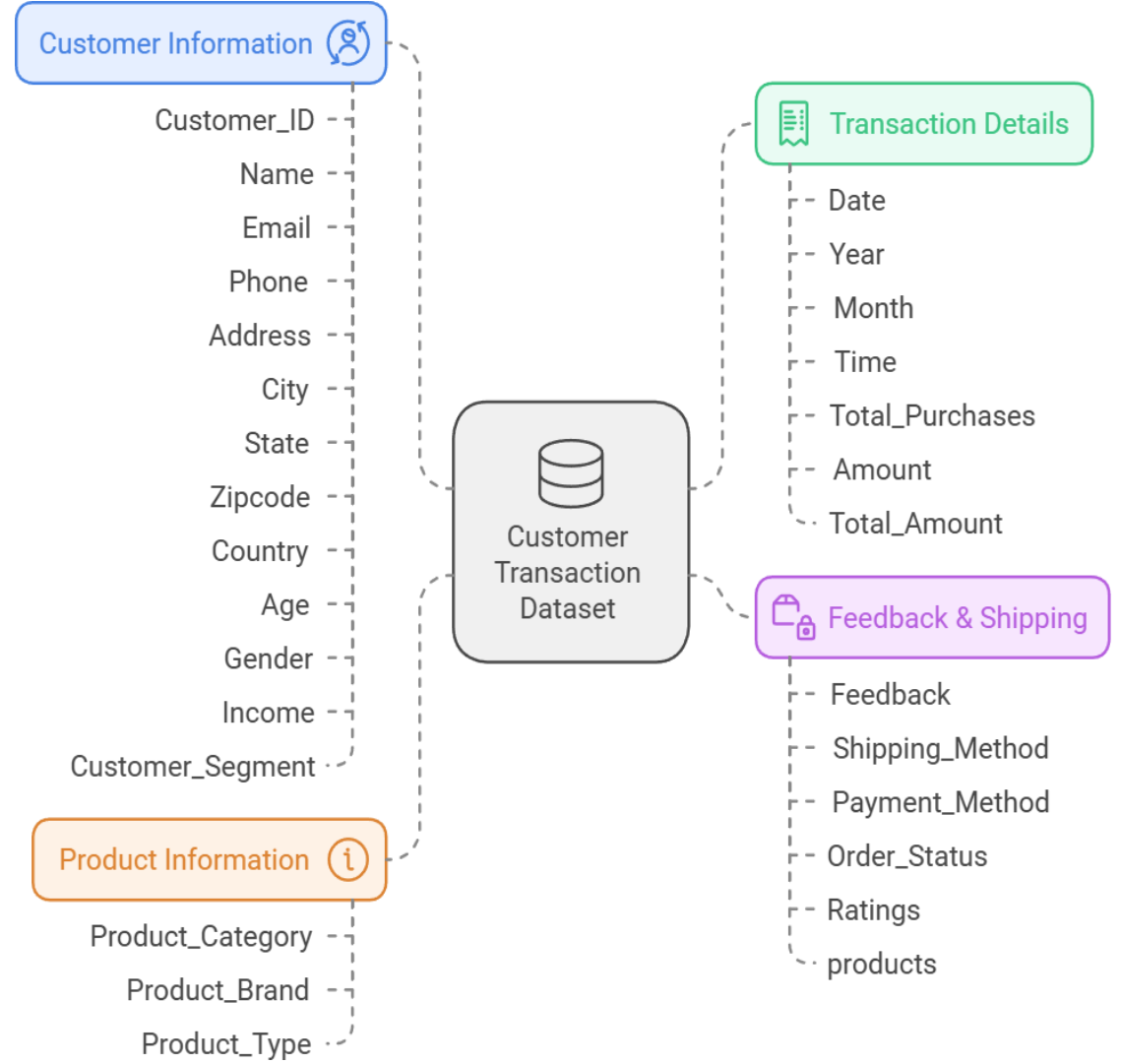Link with Source Data (Use SAS KEY)

Analyze and Transform Data

Visualize in Power BI

# KNOW YOUR DATA

Dataset Insights

# DATASET OVERVIEW

**Customer Information**

- Customer_ID
- Name
- Email
- Phone
- Address
- City
- State
- Zipcode
- Country
- Age
- Gender
- Income
- Customer_Segment

**Product Information**

- Product_Category
- Product_Brand
- Product_Type

**Customer Transaction Dataset**

**Transaction Details**

- Date
- Year
- Month
- Time
- Total_Purchases
- Amount
- Total_Amount

**Feedback & Shipping**

- Feedback
- Shipping_Method
- Payment_Method
- Order_Status
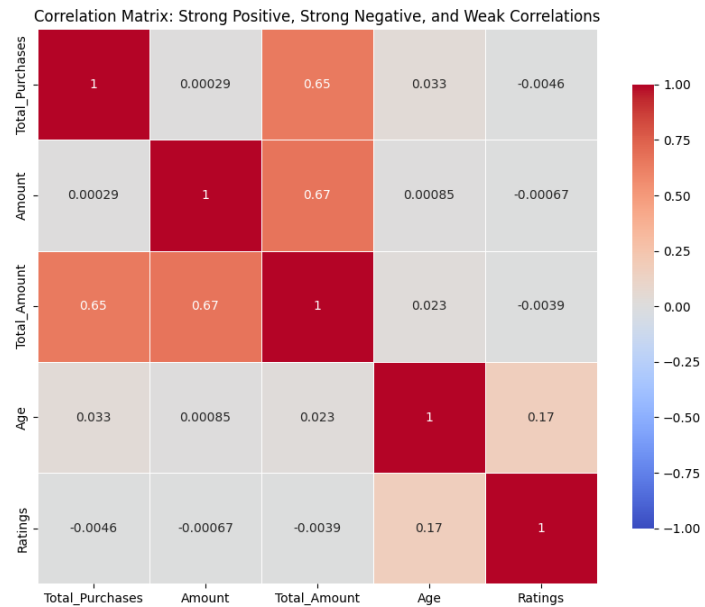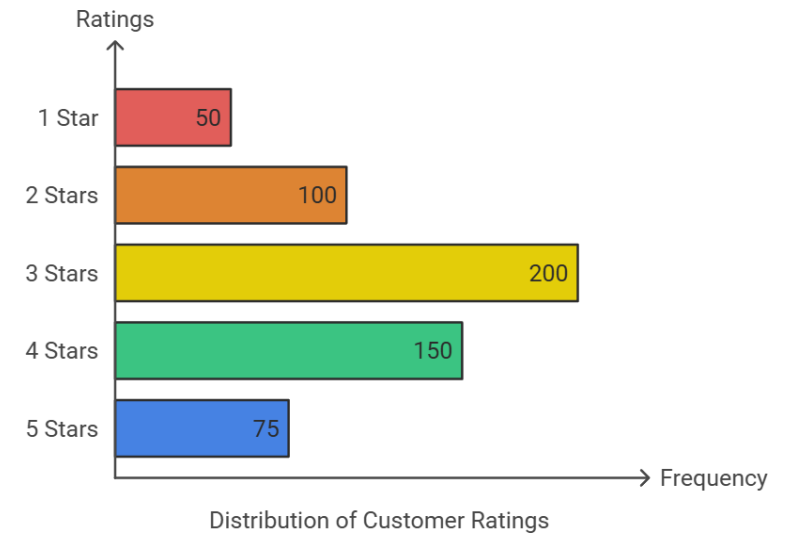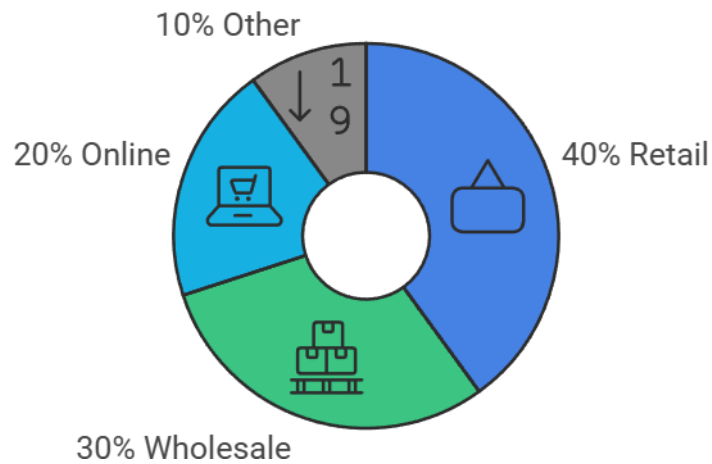- Ratings
- products

# Data Cleaning and Transformation Process

Raw Dataset



Handle Missing Values

Convert Data Types

Remove Unnecessary Columns

Encode and Scale Data (Optional)

Cleaned and Transformed Dataset

# Dataset Insights Visualization



Correlation Matrix: Strong Positive, Strong Negative, and Weak Correlations

**Distribution of Customer Segments**

10% Other
20% Online
30% Wholesale
40% Retail

Distribution of Customer Ratings

# TechRetail Sales Analytics Dashboard

# TechRetaiI  Sales  Analysis

Key Metrics for Business Insights

| Total Revenue | Total Transactions |
|---|---|
| 124.70M | 90K |

| ATV | Average Age |
|---|---|
| J.38K | 35.83 |

## Monthly  5ale Tzends

Total Sales — 11. / 10. / 10. DM — Month (0, 5, 10)

## Revenue Per Quarter

Total enue — Z0E >4 / 0M — Quarter (2, 4)

## Customer Retention Rat:e

0.42

## Total Products 5old by Category

Product_Category

- Grocery — 21K
- 20K
- 17K
- Clothing — 17K
- 17K

Total Products Sold (0K, 10K, 20K)

## Sales Distribution

Total Sales — 40M / 30M / 20M / 10M

shots

UK

## OzzJer Status Breakdown

● Delivered  G Shipped  G Processing  G Pending

- 16.73K (9.17%)
- 33K (17.96%)
- 17K (9.17%)
- 19. (10.79%)
- 32.75K (17.96%)
- 20K (10.79%)
- 22K (12.08%)
- 22.03K (12.08%)

| Date |  |
|---|---|
| All |  |

| Customer_Seg. |  |
|---|---|
| New |  |

| All |  |
|---|---|

| Country, City, St... |  |
|---|---|

S
n
a
p

```python
# Load the data from DBFS
df = spark.read.csv("dbfs:/FileStore/tables/data-1.csv", header=True, inferSchema=True)

# Clean and process the data (e.g., removing null values)
df_cleaned = df.dropna()

# Save the processed data as a Parquet file (or any other format)
df_cleaned.write.mode('overwrite').parquet("dbfs:/FileStore/tables/processed_data.parquet")
```

▶ (3) Spark Jobs

▶ 🗔 df: pyspark.sql.dataframe.DataFrame = [Transaction_ID: integer, Customer_ID: integer ... 28 more fields]

▶ 🗔 df_cleaned: pyspark.sql.dataframe.DataFrame = [Transaction_ID: integer, Customer_ID: integer ... 28 more fields]

```python
# Load the CSV file into a DataFrame
df = spark.read.csv("dbfs:/FileStore/tables/data-1.csv", header=True, inferSchema=True)
```

---

**Untitled Notebook 2024-11-07 12:38:09**  Python ▾  ☆

▶ Run all    ● Connect ▾    Schedule

File  Edit  View  Run  Help    Last edit was 3 days ago

```python
# Total Products Sold per City
products_per_city = df_clean.groupBy("City").agg({"products": "sum"}).withColumnRenamed("sum(products)",
"Total_Products_Sold")
products_per_city.show()
```

▶ (3) Spark Jobs

▶ 🗔 products_per_city: pyspark.sql.dataframe.DataFrame = [City: string, Total_Products_Sold: double]

```
|   Winnipeg|        NULL|
|     Cairns|        NULL|
|    Kelowna|        NULL|
|   Brighton|        NULL|
|      Omaha|        NULL|
|    Bendigo|        NULL|
|   Canberra|        NULL|
|     Ottawa|        NULL|
|  Edinburgh|        NULL|
|     Dallas|        NULL|
| Manchester|        NULL|
|    Oakland|        NULL|
|   Adelaide|        NULL|
|   Frankfurt|       NULL|
|       Hull|        NULL|
```

# THANK YOU