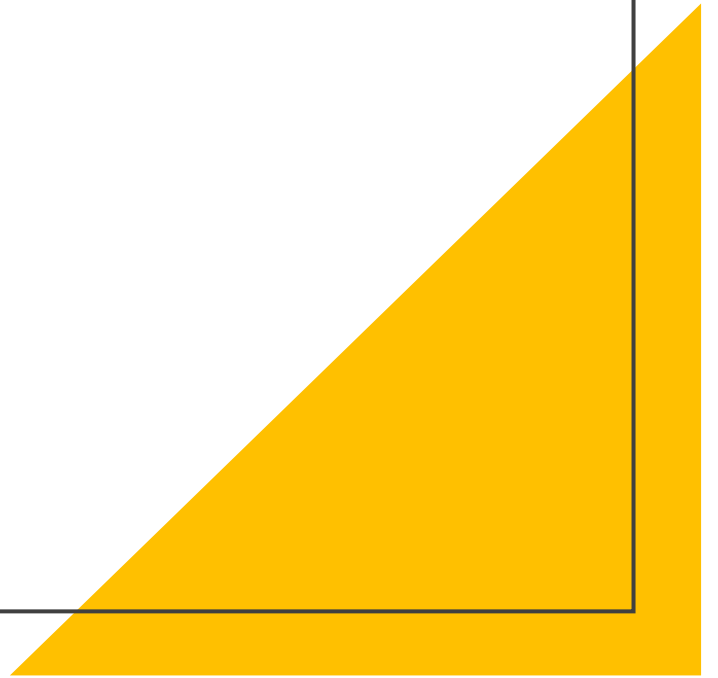# Building a Data Ecosystem in GCP for Scalable Trend Analysis

Created by: Rhythm Billore

# Introduction

**Project Background & Motivation**

- Inspired by an academic project focused on designing a database management system and conducting data analysis.

- The initial solution was constrained by limited tools, resulting in:
    - Lack of scalability and flexibility.
    - Challenges in handling multiple data sources and real-time data processing.

**Project Extension & Objectives**

- Recognized the limitations of the academic solution and reimagined it for a real-world scenario.

- Developed a scalable, end-to-end data pipeline using Google Cloud Platform (GCP) to overcome these challenges.

**Solution Overview**

- Implemented an integrated web portal for data upload.

- Leveraged cloud-native tools to automate data processing:
    - Google Cloud Storage for data storage.
    - Pub/Sub for event-driven data ingestion.
    - Implemented Cloud Functions, Google Cloud Run, Google Dataflow, and Composer for automated processing.
    - BigQuery for efficient data analysis.
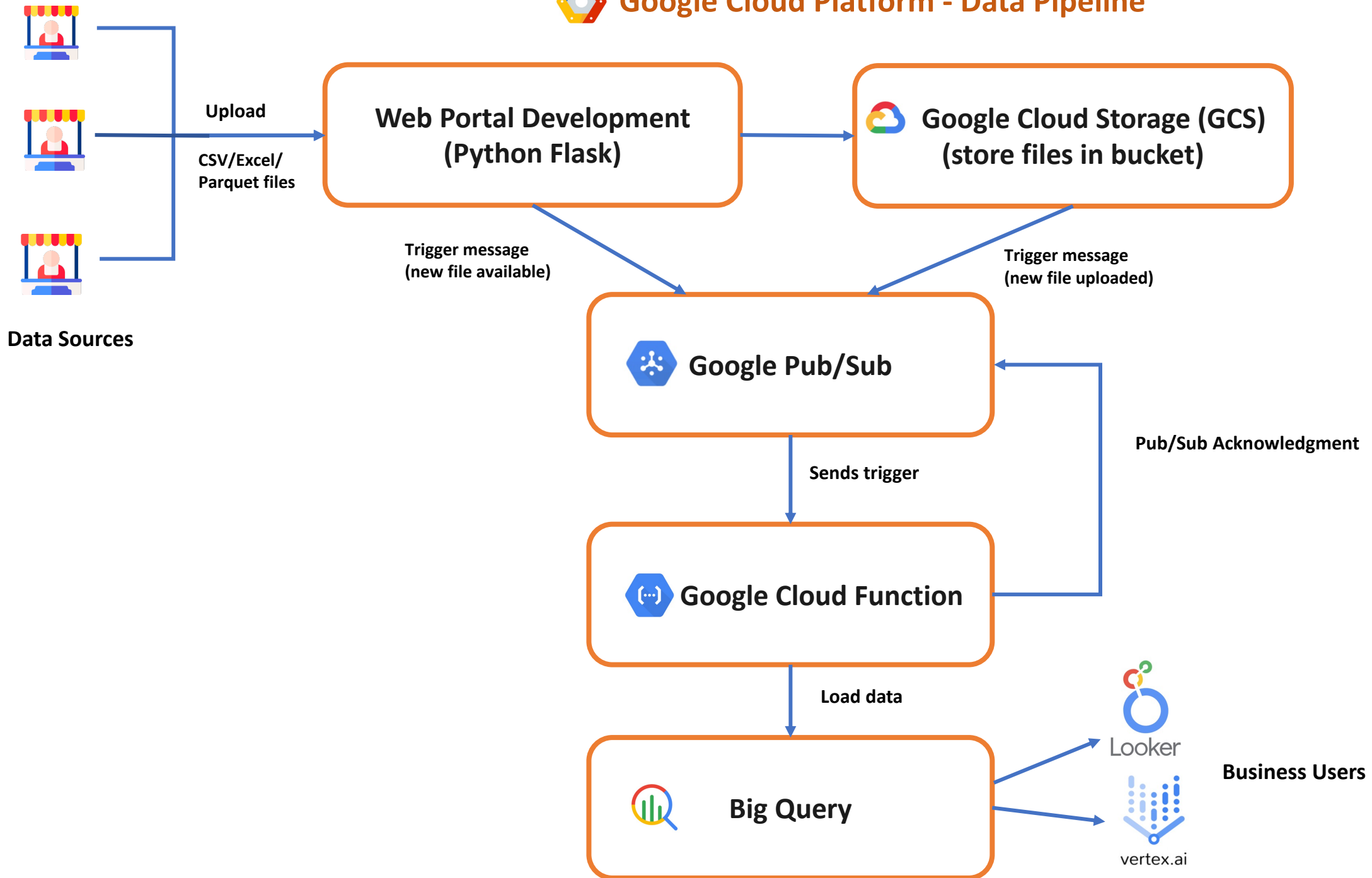
**Key Benefits**

- Created a robust data ecosystem enabling seamless data flow and real-time analysis.

- Scalable solution designed to meet the needs of multiple stakeholders.

- Provides business end-users with actionable insights through effective data visualization.

**Deliverable**

- Developed an interactive analysis dashboard to deliver in-depth insights on **Sales trends** and **customer behavior analysis**.

- Supports informed decision-making by leveraging comprehensive data analytics.

**1** Sales Data Portal

**2** GCS

**3** Pub/Sub

**4** Cloud function

**5** Big Query

4

# Sales Data Analysis - 2024

**# of Sales**
5000

**States Covered**
50

**Net Revenue**
402.2K

**# of Active Subscribers**
1,343

## Sales Across Categories

Legend: Male, Female

| Category | Male | Female |
|---|---|---|
| Clothing | 1,377 | 634 |
| Accessories | 1,030 | 491 |
| Footwear | 574 | 275 |
| Outerwear | 423 | 196 |

(Y-axis: # of Sales, from 0 to 2.5K; X-axis: Category)

## Top 10 Selling Items

Legend: Total Purchase Amount (USD)

| Item Purchased | Total Purchase Amount (USD) |
|---|---|
| Sunglasses | 19,293 |
| Shirt | 18,722 |
| Belt | 17,872 |
| Coat | 17,423 |
| Shoes | 17,320 |
| Blouse | 17,059 |
| Shorts | 16,941 |
| Dress | 16,914 |
| Scarf | 16,735 |
| Pants | 16,673 |

## Sales Trend by Size

Legend: M, L, S, XL

| Item Purchased | M | L | S | XL |
|---|---|---|---|---|
| Shirt | 105 | 61 | 28 | 29 |
| Pants | 103 | 54 | 34 | 27 |
| Sunglasses | 105 | 52 | 45 | 16 |
| Blouse | 97 | 55 | 36 | 25 |
| Jewelry | 89 | 52 | 46 | 25 |
| Dress | 98 | 57 | 38 | 19 |
| Belt | 88 | 55 | 44 | 25 |
| Sweater | 103 | 53 | 33 | 21 |
| Shorts | 88 | 62 | 40 | 19 |
| Socks | 96 | 57 | 30 | 24 |

(X-axis: # of Sales)

## Season-wise Sales Trend

Legend: Winter, Fall, Spring, Summer

| Item Purchased | Winter | Fall | Spring | Summer |
|---|---|---|---|---|
| Sunglasses | 5,059 | 5,255 | 4,187 | 4,792 |
| Shirt | 5,736 | 3,940 | 4,778 | 4,268 |
| Belt | 3,798 | 6,005 | 3,135 | 4,934 |
| Coat | 5,059 | 4,443 | 4,331 | 3,590 |
| Shoes | 4,124 | 4,755 | 4,999 | 3,442 |
| Blouse | 4,222 | 4,548 | 4,345 | 3,944 |
| Shorts | 5,448 | 4,398 | 3,156 | 3,939 |
| Dress | 4,519 | 3,702 | 4,352 | 4,341 |
| Scarf | 3,737 | 4,473 | 3,528 | 4,997 |
| Pants | 4,263 | 3,980 | 3,271 | 5,159 |

(X-axis: Total Purchase Amount (USD))

# Customer Behavior analysis - 2024

## Month-on-Month Purchase Trends by Category (2024)

Legend: Clothing — Accessories — Footwear — Outerwear

Y-axis: Total Purchase Amount (USD) — 0, 10K, 20K

X-axis: Jan 2024, Feb 2024, Mar 2024, Apr 2024, May 2024, Jun 2024, Jul 2024, Aug 2024, Sep 2024, Oct 2024

DateOfPurchase (Date)

## Transaction Share by Payment Method

- Credit Card — 40%
- PayPal — 29%
- Venmo — 20.4%
- Cash
- Bank Transfer

## Promo Code Usage vs. Purchase Frequency

(Promo codes applied on 85 / 295 days)

Legend: Average frequency of purchase (days) — # products sold

| PromoCodeUsed | Average frequency of purchase (days) | # products sold |
| --- | --- | --- |
| false | 1.65 | 2,848 |
| true | 1.63 | 2,152 |

Reduce Delivery Times for High CLV Customers to Prevent Loss

Scatter plot: CustomerID — Y-axis: DeliveryTimeDays (0, 5, 10, 15) — X-axis: CustomerLifetimeValue (0, 200, 400, 600)

| | Loyalty Segment | SubscriptionStatus | Current Status | Conversion Target | No. of Customers ▾ | Frequency of purchase ( in days) | Avg Purchase Amount (USD) | CustomerLifetimeValue |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 1. | Low Loyalty Points | false | One-Off Buyers | High Loyality | 3,135 | 1.62 | 59.9 | 68.2 |
| 2. | Low Loyalty Points | true | | | 1,141 | 1.68 | 59.1 | 65.07 |
| 3. | Medium Loyalty Points | false | | | 273 | 1.63 | 160.9 | 236.83 |
| 4. | High Loyalty Points | false | Loyal Non-Subscribers | Subscribers | 247 | 1.7 | 256.2 | 333.24 |
| 5. | Medium Loyalty Points | true | | | 105 | 1.76 | 154.8 | 241.54 |
| 6. | High Loyalty Points | true | Subscribers | Retain | 95 | 1.39 | 246.9 | 335.02 |

1 - 8 / 8  < >