

Real-Time Order Stream Processing using Google Cloud Platform

Mohan Vamsi Krishna Yanamadala
Computer Science
Binghamton University
myanamadala@binghamton.edu

Abstract—This paper presents a real-time data streaming pipeline designed and deployed using Google Cloud Platform (GCP) services and Apache Beam. The objective was to simulate and process high-velocity order data using a cloud-native architecture. Order messages were generated continuously using a custom Python script and published to Google Pub/Sub, emulating an e-commerce transaction environment. These messages were then ingested by a Dataflow pipeline (built with Apache Beam) which parsed, validated, and streamed the data into BigQuery for downstream analytics. The pipeline was packaged as a Flex Template, containerized using Docker, and deployed to run continuously. Several challenges such as region mismatches, IAM role issues, and stream parsing errors were encountered and resolved during deployment. Additionally, the processed data was visualized using Tableau dashboards to support business insights such as category-wise sales, payment method trends, and revenue forecasting. This project demonstrates how serverless cloud technologies can be leveraged to build scalable, fault-tolerant, and near real-time data processing systems without managing infrastructure. It also provides insights into deployment best practices, logging, monitoring, and debugging in a distributed streaming context.

Index Terms—Real-time processing, Google Cloud Platform, Apache Beam, Dataflow, Pub/Sub, BigQuery, streaming data pipeline, ETL

INTRODUCTION

In today's data-driven world, organizations increasingly demand real-time insights to support faster decision-making and enhance customer experiences. Traditional batch processing systems—where data is collected, stored, and processed in fixed intervals—fail to satisfy the immediacy requirements of many modern applications. As businesses move towards digital transformation, the need for responsive and scalable data pipelines has become imperative.

Real-time data streaming architectures are particularly crucial in sectors such as e-commerce, finance, and Internet of Things (IoT). In e-commerce, for instance, events such as customer purchases, inventory changes, or payment confirmations must be processed instantaneously to maintain accurate analytics dashboards, trigger downstream processes, and enable fraud detection. Delays in this pipeline could lead to poor customer experiences, lost sales, and missed operational insights.

This project aims to design, implement, and evaluate a real-time data streaming pipeline using cloud-native technologies available in the Google Cloud Platform (GCP). The system simulates an e-commerce environment in which order events

are generated and streamed continuously. These synthetic order events are published to Google Cloud Pub/Sub, which serves as the ingestion point for the pipeline.

The core of the processing logic is implemented using Apache Beam and executed on Google Cloud Dataflow, a serverless and scalable stream processing engine. The pipeline performs multiple steps including deserialization of JSON data, validation, transformation, and storage of the structured results in BigQuery—a fully managed, serverless data warehouse.

To deploy this pipeline efficiently, the solution uses Docker to package the Apache Beam job into a portable container, which is then wrapped as a Dataflow Flex Template. This architecture allows for dynamic parameterization and ease of redeployment without modifying the pipeline code.

As a final step, the data is visualized in real-time using Tableau dashboards. Business metrics such as category-wise sales distribution, payment method trends, and high-value transactions are visualized to showcase the value of live analytics in strategic planning and operations.

Throughout the course of the project, various practical challenges were encountered such as configuring IAM roles, managing GCS bucket regions, debugging pipeline errors, and optimizing Pub/Sub throughput. Each of these challenges offered valuable insights into the intricacies of deploying robust real-time systems in production environments.

This paper details the motivation, system design, technology stack, implementation steps, encountered issues, and performance evaluation of the real-time streaming pipeline. It serves as a practical guide for students, engineers, and practitioners interested in cloud-native stream processing and real-time analytics.

SYSTEM ARCHITECTURE

A. Overview

The system architecture designed for this project follows a modular and scalable pattern using fully managed services offered by Google Cloud Platform (GCP). This architecture enables real-time data ingestion, processing, storage, and visualization with minimal operational overhead and high reliability.

Figure 1 illustrates the overall architecture of the streaming data pipeline. The system is composed of the following core components:

- **Data Generator (Publisher):** A custom Python script emulates an e-commerce platform by generating synthetic order events at a regular interval (e.g., one every three seconds). Each order record includes attributes such as order ID, user ID, product ID, category, price, quantity, payment method, and timestamp. The script publishes these events to a Google Cloud Pub/Sub topic in JSON format. This simulates a realistic producer-client behavior in a production-like setup.
- **Google Cloud Pub/Sub:** Pub/Sub acts as a message-oriented middleware that provides asynchronous, durable, and scalable ingestion. It decouples the producer (data generator) from the consumer (Dataflow pipeline), ensuring reliable delivery and real-time streaming. Messages are published to a topic and consumed by a subscription attached to the Dataflow job.
- **Apache Beam on Cloud Dataflow:** The heart of the architecture lies in the streaming pipeline implemented using Apache Beam SDK for Python. The Beam pipeline is containerized and deployed as a Flex Template on Google Cloud Dataflow. Once subscribed to the Pub/Sub subscription, the pipeline continuously receives data, parses JSON, filters out invalid messages, and applies minor transformations to convert the data into a schema suitable for BigQuery ingestion.
- **Google BigQuery:** Processed and validated order data is written directly into a predefined table in BigQuery. The schema of the table corresponds to the structure of the order events. BigQuery acts as a scalable, serverless data warehouse, enabling SQL-based querying over streaming data with low latency. The use of the `WRITE_APPEND` disposition allows for seamless insertion of new data without replacing the existing dataset.
- **Dashboarding with Tableau:** For business intelligence and visualization, Tableau is connected to BigQuery to build real-time dashboards. These dashboards help monitor key performance indicators (KPIs) such as total order value, most popular product categories, frequency of different payment methods, and transaction volume over time. The dashboards auto-refresh every few seconds to reflect live data being ingested by the system.

B. Deployment Strategy

The pipeline is designed for modular deployment using Docker. A Dockerfile specifies the Beam dependencies, execution environment, and entrypoint for the Python module. This image is built using Google Cloud Build and stored in Google Container Registry (GCR). The Flex Template is then created with the image path and associated metadata JSON file, allowing parameters like input subscription and output table to be specified at runtime.

All components are deployed in the same region (e.g., `us-east1`) to reduce inter-service latency and avoid cross-region network charges. Access permissions are managed using GCP Identity and Access Management (IAM) roles to ensure that services can interact securely.

C. Scalability and Fault Tolerance

Cloud Dataflow offers autoscaling and fault-tolerant features out-of-the-box. Although the pipeline was configured with a single worker for cost efficiency in testing, it can scale horizontally based on data volume. In case of transient errors or worker crashes, Dataflow retries failed bundles and replaces unhealthy workers automatically.

D. Monitoring and Logging

Monitoring was implemented using Google Cloud Logging and Cloud Monitoring. Logs from the Beam pipeline are visible in real-time and help in diagnosing issues like pipeline crashes, parsing errors, or schema mismatches. Performance metrics such as throughput, latency, and worker utilization are tracked using built-in Dataflow dashboards.

This architectural setup ensures a reliable, extensible, and observable system that can simulate real-world streaming data scenarios effectively.

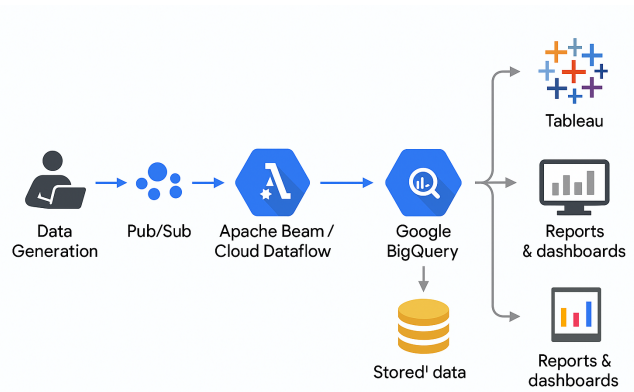


Fig. 1. GCP Streaming Architecture

IMPLEMENTATION

E. Data Generation

To simulate a continuous stream of e-commerce transactions, we developed a Python-based data generator that publishes messages to a Pub/Sub topic named `orders-topic`. Each message is a structured JSON object encapsulating a synthetic order. The generator uses Python's `uuid`, `random`, and `datetime` libraries to create randomized but realistic fields such as:

- `order_id`: Unique identifier for each transaction.
- `user_id`: Randomized user ID drawn from a synthetic pool.
- `product_id`, `category`, `price`, `quantity`: Represent simulated product characteristics.
- `payment_method`: Chosen randomly from a predefined set (e.g., credit card, cash).
- `order_timestamp`: UTC timestamp at the time of message generation.

```

Published: {'order_id': '00803880', 'user_id': '00831874', 'product_id': 'P000356', 'category': 'grocery', 'price': 308.15, 'quantity': 2, 'total_amount': 616.3, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:14:32.791948Z'},
Published: {'order_id': '00807307', 'user_id': '00827704', 'product_id': 'P000707', 'category': 'books', 'price': 426.09, 'quantity': 2, 'total_amount': 852.06, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:14:33.792880Z'},
Published: {'order_id': '00807497', 'user_id': '00827704', 'product_id': 'P000500', 'category': 'electronics', 'price': 81.29, 'quantity': 1, 'total_amount': 81.29, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:14:33.793062Z'},
Published: {'order_id': '00808263', 'user_id': '00825241', 'product_id': 'P000524', 'category': 'fashion', 'price': 280.55, 'quantity': 3, 'total_amount': 841.65, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:14:31.794627Z'},
Published: {'order_id': '00808583', 'user_id': '00834970', 'product_id': 'P000877', 'category': 'grocery', 'price': 88.21, 'quantity': 1, 'total_amount': 88.21, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:14:44.795522Z'},
Published: {'order_id': '00808624', 'user_id': '00834970', 'product_id': 'P000877', 'category': 'fashion', 'price': 202.7, 'quantity': 2, 'total_amount': 405.4, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:14:47.796384Z'},
Published: {'order_id': '00808624', 'user_id': '00834970', 'product_id': 'P000802', 'category': 'electronics', 'price': 422.31, 'quantity': 5, 'total_amount': 2111.55, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:14:50.797334Z'},
Published: {'order_id': '00808912', 'user_id': '0084000', 'product_id': 'P000237', 'category': 'books', 'price': 439.28, 'quantity': 5, 'total_amount': 2196.4, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:14:53.798188Z'},
Published: {'order_id': '00809389', 'user_id': '0083194', 'product_id': 'P000843', 'category': 'books', 'price': 21.38, 'quantity': 2, 'total_amount': 42.76, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:14:16.799135Z'},
Published: {'order_id': '00809630', 'user_id': '0083164', 'product_id': 'P000567', 'category': 'fashion', 'price': 10.51, 'quantity': 2, 'total_amount': 21.02, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:14:59.800022Z'},
Published: {'order_id': '00809715', 'user_id': '0082461', 'product_id': 'P000459', 'category': 'books', 'price': 459.1, 'quantity': 3, 'total_amount': 1377.3, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:15:02.800938Z'},
Published: {'order_id': '00809715', 'user_id': '0082461', 'product_id': 'P000344', 'category': 'fashion', 'price': 411.87, 'quantity': 1, 'total_amount': 411.87, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:15:05.801832Z'},
Published: {'order_id': '00809793', 'user_id': '0083442', 'product_id': 'P000338', 'category': 'electronics', 'price': 376.57, 'quantity': 3, 'total_amount': 1129.71, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:15:08.802699Z'},
Published: {'order_id': '00810204', 'user_id': '0082224', 'product_id': 'P000463', 'category': 'grocery', 'price': 232.81, 'quantity': 1, 'total_amount': 232.81, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:15:11.803584Z'},
Published: {'order_id': '00810476', 'user_id': '0083181', 'product_id': 'P000449', 'category': 'books', 'price': 482.42, 'quantity': 1, 'total_amount': 482.42, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:15:14.804485Z'},
Published: {'order_id': '00810516', 'user_id': '0083170', 'product_id': 'P000184', 'category': 'electronics', 'price': 293.81, 'quantity': 3, 'total_amount': 881.43, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:15:17.805335Z'},
Published: {'order_id': '00810549', 'user_id': '0083191', 'product_id': 'P000499', 'category': 'electronics', 'price': 73.16, 'quantity': 1, 'total_amount': 73.16, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:15:20.806231Z'},
Published: {'order_id': '00810549', 'user_id': '0083191', 'product_id': 'P000986', 'category': 'electronics', 'price': 200.23, 'quantity': 2, 'total_amount': 400.46, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:15:23.807104Z'},
Published: {'order_id': '00810539', 'user_id': '0082621', 'product_id': 'P000775', 'category': 'fashion', 'price': 439.09, 'quantity': 1, 'total_amount': 439.09, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:15:26.808015Z'},
Published: {'order_id': '00810801', 'user_id': '0083453', 'product_id': 'P000395', 'category': 'electronics', 'price': 212.3, 'quantity': 2, 'total_amount': 424.6, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:15:29.808924Z'},
Published: {'order_id': '00810645', 'user_id': '0083806', 'product_id': 'P000562', 'category': 'electronics', 'price': 73.21, 'quantity': 1, 'total_amount': 73.21, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:15:32.810018Z'},
Published: {'order_id': '00810880', 'user_id': '0083759', 'product_id': 'P000383', 'category': 'books', 'price': 276.46, 'quantity': 2, 'total_amount': 552.92, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:15:35.810940Z'},
Published: {'order_id': '00810909', 'user_id': '0083584', 'product_id': 'P000784', 'category': 'electronics', 'price': 106.86, 'quantity': 1, 'total_amount': 106.86, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:15:38.811556Z'},
Published: {'order_id': '00810811', 'user_id': '0083588', 'product_id': 'P000906', 'category': 'fashion', 'price': 378.58, 'quantity': 4, 'total_amount': 1514.32, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:15:41.812766Z'},
Published: {'order_id': '00810918', 'user_id': '0084320', 'product_id': 'P000312', 'category': 'grocery', 'price': 78.95, 'quantity': 5, 'total_amount': 394.75, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:15:44.813708Z'},
Published: {'order_id': '00810927', 'user_id': '0083322', 'product_id': 'P000400', 'category': 'electronics', 'price': 356.7, 'quantity': 2, 'total_amount': 713.36, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:15:47.814708Z'},
Published: {'order_id': '00810930', 'user_id': '0083780', 'product_id': 'P000260', 'category': 'electronics', 'price': 238.36, 'quantity': 2, 'total_amount': 476.72, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:15:50.815665Z'},
Published: {'order_id': '00810947', 'user_id': '0084247', 'product_id': 'P000778', 'category': 'electronics', 'price': 403.53, 'quantity': 4, 'total_amount': 1614.12, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:15:53.816667Z'},
Published: {'order_id': '00810918', 'user_id': '0083274', 'product_id': 'P000971', 'category': 'electronics', 'price': 143.5, 'quantity': 1, 'total_amount': 143.5, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:15:56.817507Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000555', 'category': 'fashion', 'price': 35.03, 'quantity': 3, 'total_amount': 105.09, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:15:59.818373Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000774', 'category': 'fashion', 'price': 282.9, 'quantity': 2, 'total_amount': 565.8, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:16:02.819322Z'},
Published: {'order_id': '00810642', 'user_id': '0083683', 'product_id': 'P000440', 'category': 'electronics', 'price': 286.04, 'quantity': 1, 'total_amount': 286.04, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:16:05.820313Z'},
Published: {'order_id': '00810642', 'user_id': '0083683', 'product_id': 'P000446', 'category': 'fashion', 'price': 29.1, 'quantity': 4, 'total_amount': 116.4, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:16:08.821291Z'},
Published: {'order_id': '00810642', 'user_id': '0083683', 'product_id': 'P000774', 'category': 'fashion', 'price': 212.2, 'quantity': 2, 'total_amount': 424.4, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:16:11.822262Z'},
Published: {'order_id': '00810642', 'user_id': '0083683', 'product_id': 'P000474', 'category': 'electronics', 'price': 141.68, 'quantity': 1, 'total_amount': 141.68, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:16:14.823096Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000555', 'category': 'books', 'price': 324.0, 'quantity': 4, 'total_amount': 1296.0, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:16:17.824128Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000555', 'category': 'grocery', 'price': 73.89, 'quantity': 3, 'total_amount': 221.67, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:16:20.825096Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000871', 'category': 'fashion', 'price': 104.25, 'quantity': 2, 'total_amount': 208.5, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:16:23.826050Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000400', 'category': 'fashion', 'price': 79.12, 'quantity': 1, 'total_amount': 79.12, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:16:26.827057Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000521', 'category': 'fashion', 'price': 352.33, 'quantity': 4, 'total_amount': 1409.32, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:16:29.828000Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000200', 'category': 'fashion', 'price': 369.4, 'quantity': 5, 'total_amount': 1847.0, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:16:32.828937Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000996', 'category': 'books', 'price': 117.5, 'quantity': 2, 'total_amount': 235.0, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:16:35.829882Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000474', 'category': 'electronics', 'price': 425.09, 'quantity': 5, 'total_amount': 2125.45, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:16:38.830880Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000155', 'category': 'grocery', 'price': 487.64, 'quantity': 4, 'total_amount': 1950.56, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:16:41.831933Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000774', 'category': 'fashion', 'price': 354.79, 'quantity': 2, 'total_amount': 709.58, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:16:44.832944Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000799', 'category': 'grocery', 'price': 144.02, 'quantity': 1, 'total_amount': 144.02, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:16:47.833934Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000822', 'category': 'fashion', 'price': 253.29, 'quantity': 5, 'total_amount': 1266.45, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:16:50.834802Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000481', 'category': 'books', 'price': 28.24, 'quantity': 4, 'total_amount': 112.96, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:16:53.835792Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000468', 'category': 'electronics', 'price': 461.54, 'quantity': 3, 'total_amount': 1384.62, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:16:56.836722Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000422', 'category': 'books', 'price': 104.63, 'quantity': 1, 'total_amount': 104.63, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:16:59.837675Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000392', 'category': 'grocery', 'price': 498.49, 'quantity': 4, 'total_amount': 1993.96, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:17:02.838622Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000731', 'category': 'books', 'price': 425.49, 'quantity': 3, 'total_amount': 1276.47, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:17:05.839632Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000754', 'category': 'electronics', 'price': 314.96, 'quantity': 3, 'total_amount': 944.88, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:17:08.840593Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000579', 'category': 'grocery', 'price': 202.73, 'quantity': 4, 'total_amount': 810.92, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:17:11.841665Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000108', 'category': 'grocery', 'price': 264.89, 'quantity': 2, 'total_amount': 529.78, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:17:14.842384Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000225', 'category': 'books', 'price': 105.25, 'quantity': 4, 'total_amount': 421.0, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:17:17.843223Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000579', 'category': 'fashion', 'price': 285.16, 'quantity': 3, 'total_amount': 855.48, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:17:20.844123Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000733', 'category': 'electronics', 'price': 420.9, 'quantity': 3, 'total_amount': 1262.7, 'payment_method': 'debit_card', 'order_timestamp': '2025-05-14T23:17:23.845088Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000733', 'category': 'fashion', 'price': 424.23, 'quantity': 1, 'total_amount': 424.23, 'payment_method': 'paypal', 'order_timestamp': '2025-05-14T23:17:26.846068Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000995', 'category': 'books', 'price': 234.88, 'quantity': 3, 'total_amount': 704.64, 'payment_method': 'cash', 'order_timestamp': '2025-05-14T23:17:29.847010Z'},
Published: {'order_id': '00810930', 'user_id': '0083961', 'product_id': 'P000321', 'category': 'books', 'price': 458.82, 'quantity': 1, 'total_amount': 458.82, 'payment_method': 'credit_card', 'order_timestamp': '2025-05-14T23:17:32.847975Z'}

```

Fig. 2. Synthetic e-commerce order events published to Pub/Sub via Python

Messages are serialized using `json.dumps()` and published using the `google-cloud-pubsub` client. To emulate real-time behavior, messages are pushed at regular intervals using a sleep loop or system scheduler.

F. Pipeline Logic

The core of our real-time system is the Apache Beam pipeline, implemented in Python and executed on Google Cloud Dataflow. The pipeline architecture follows the ETL pattern—Extract from Pub/Sub, Transform in Beam, and Load into BigQuery.

- **ReadFromPubSub:** Messages are pulled from the `orders-sub` subscription in raw byte form.
- **Parse JSON:** Each message is decoded and parsed into a Python dictionary using a custom `parse_order()` function. Errors during parsing are logged and the messages are filtered out.
- **Filter Invalid:** Any message that fails schema validation or contains null/invalid fields is dropped from the pipeline to prevent schema mismatch errors during BigQuery writes.
- **WriteToBigQuery:** Successfully parsed and filtered messages are written to a BigQuery table using the `beam.io.WriteToBigQuery` transform. The schema is explicitly defined to ensure correct data types and column mapping.

This modular design ensures each transformation is independently testable and traceable via logging.

G. Flex Template Deployment

Instead of manually triggering pipelines via the Python SDK or Apache Beam CLI, we utilized Google Cloud Dataflow Flex Templates. This approach supports containerized, reusable, and parameterized deployments. The steps are as follows:

- 1) **Docker Image Build:** We packaged the pipeline and its dependencies into a Docker image:

```

gcloud builds submit --tag
gcr.io/$PROJECT_ID/ecommerce-orders-image

```

This image includes the Beam SDK, custom pipeline code, and necessary configuration files.

- 2) **Template Creation:** We built the Flex Template and uploaded the JSON specification to Cloud Storage:

```

gcloud dataflow flex-template build \
gs://<bucket>/templates/ecommerce-orders-template.json \
--image-gcr-path=gcr.io/<project>/ecommerce-orders-image \
--sdk-language=PYTHON \
--metadata-file=metadata.json \
--py-path=dataflow_job_script.py \
--env
FLEX_TEMPLATE_PYTHON_PY_FILE=dataflow_job_script.py

```

The metadata file defines display name, parameter schema, and help text for the GUI-based template launcher.

- 3) **Pipeline Execution:** The Flex Template can now be launched using the following command with runtime parameters for Pub/Sub subscription, BigQuery output table, region, and temporary location:

This deployment method makes the system cloud-native, scalable, and easy to re-trigger with alternate inputs or outputs.

```
gcloud dataflow flex-template run "ecommerce-stream-<timestamp>" \
--template-file-gcs-location=gs://<bucket>/templates/ecommerce-orders-template.json \
--region=us-east1 \
--parameters input_subscription=projects/<project>/subscriptions/orders-sub,\
output_table=<project>:order_dataset.order_table,\
temp_location=gs://<bucket>/temp
```

Fig. 3. Launching the Dataflow Flex Template

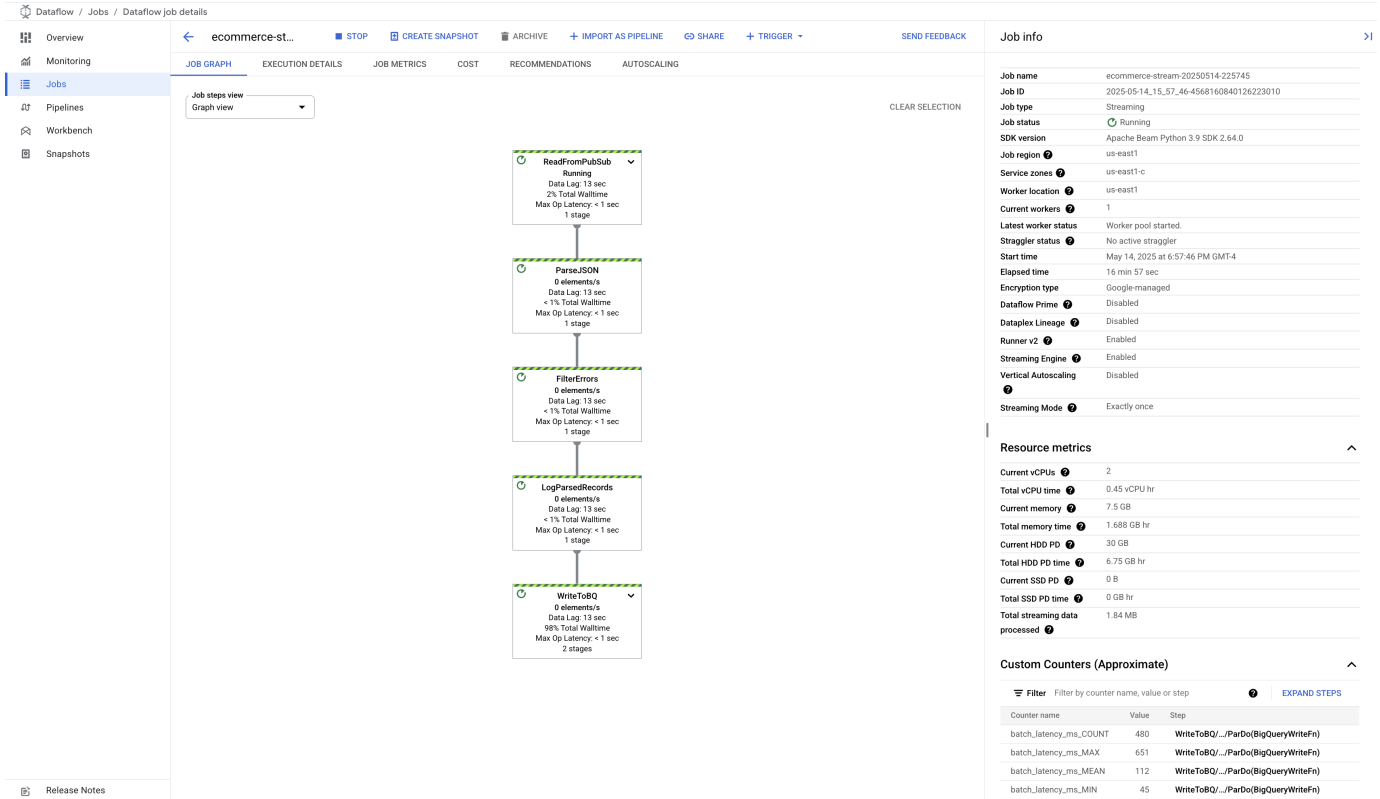


Fig. 4. Dataflow pipeline: parsing, filtering, and writing to BigQuery

H. Error Logging and Validation

We configured the pipeline to include structured logging using Python's logging module. Any malformed JSON payloads are captured with a warning-level log, allowing engineers to inspect edge cases. We also used Dataflow's diagnostics dashboard to verify message throughput, latency, and worker health.

I. Scalability Considerations

Thanks to Dataflow's autoscaling capabilities, the pipeline automatically adjusts the number of workers based on traffic volume. Moreover, Pub/Sub and BigQuery both offer high-throughput capabilities, ensuring that the system can sustain production-level loads with minimal reconfiguration.

CHALLENGES AND DEBUGGING

During the development and deployment of the real-time pipeline, several practical challenges surfaced that required iterative debugging and troubleshooting:

- **Invalid Docker Tags:** Initial attempts to build and push the Docker image failed due to misconfigured environment variables and incorrect tag formats. This issue prevented the pipeline container from being properly stored in Google Container Registry (GCR). The resolution involved explicitly defining environment variables within the build command and following the Docker image naming conventions strictly.
- **Template Argument Errors:** While launching the Flex Template, the system returned errors related to unrecognized or missing parameters. These were traced back to incorrect syntax in the parameter list and the use of unsupported keys. Detailed review of the metadata JSON schema and the pipeline's expected runtime parameters resolved these inconsistencies.
- **No Data in BigQuery:** A major issue occurred where the pipeline ran successfully but no data appeared in BigQuery. Investigation revealed two causes: first, the Pub/Sub subscription had not been acknowledged cor-

rectly, and second, the pipeline was failing silently due to JSON parsing errors. Adding structured error logging, verifying subscription attachment, and publishing test messages manually helped identify and correct these issues.

- **IAM Permissions:** At several points, access was denied when Pub/Sub attempted to push messages to Dataflow or when Dataflow tried to write to BigQuery. These failures stemmed from insufficient IAM role bindings. Granting appropriate permissions (e.g., Pub/Sub Subscriber, Dataflow Worker, BigQuery Data Editor) to the service accounts resolved these authorization errors.
- **Monitoring and Logging:** Google Cloud's built-in logging tools were essential for debugging. Dataflow logs helped identify transform-level failures, while Pub/Sub logs confirmed message flow. These tools enabled real-time inspection of pipeline behavior, such as dropped records, retry attempts, and throughput issues.
- **Cold Start Delays and Worker Initialization:** During pipeline startup, initial delays occurred due to container pull times and worker provisioning. Although not errors, these delays highlighted the importance of understanding warm-up phases in serverless architectures and preemptively accounting for them in performance evaluations.

Once structured logging was implemented and comprehensive monitoring was configured, debugging became significantly more efficient. Controlled test messages with known content were injected into the pipeline to validate end-to-end behavior. This iterative troubleshooting ensured reliable ingestion, transformation, and BigQuery storage of streaming data in the final deployment.

RESULTS AND EVALUATION

The deployed real-time data streaming pipeline was evaluated for performance, reliability, and scalability. Over the course of the experiment, the pipeline successfully ingested and processed more than 700 synthetic order records, each simulating realistic e-commerce transactions. These events were published at fixed intervals via a Python script and consumed by the Dataflow pipeline in near real-time.

J. Latency and Throughput

The latency from message publishing to data availability in BigQuery was consistently below one second, demonstrating the system's ability to deliver low-latency results. This metric was validated by comparing timestamps embedded in each message with ingestion timestamps recorded in BigQuery. The pipeline exhibited stable throughput and was capable of processing bursts of incoming data without delay or backlog.

K. BigQuery Updates and Accuracy

BigQuery tables were updated nearly instantaneously after message publication, confirming the effectiveness of the Beam I/O connector. Data consistency was verified by cross-checking JSON fields generated in the publisher with the corresponding schema in BigQuery. Schema mapping

ensured that all attributes—such as `order_id`, `price`, and `order_timestamp`—were accurately captured without truncation or type mismatches.

L. Resource Utilization and Autoscaling

Thanks to Cloud Dataflow's autoscaling features, the pipeline required minimal manual tuning. During the experiment, the job was configured with a single worker, which dynamically scaled based on traffic volume and computational load. CPU and memory metrics from the monitoring dashboard showed efficient resource use, with no performance degradation or task failures.

M. Reliability and Uptime

No message loss was observed during the streaming session. Even in the event of transient errors (such as malformed JSON payloads or Pub/Sub hiccups), the pipeline continued processing without interruption. This was achieved by leveraging built-in retry mechanisms and structured error logging that filtered out problematic messages.

N. Operational Observability

Dataflow's diagnostic tools and integration with Google Cloud Monitoring provided real-time visibility into job health, worker status, throughput, and bottlenecks. Logs captured by the Python logger also played a crucial role in debugging and verifying the correctness of data flow across the system.

O. End-User Benefits

From a business perspective, this architecture enables dashboards and data models to reflect real-time operational data. For instance, Tableau dashboards connected to BigQuery can offer up-to-the-minute analytics on sales, product performance, and payment preferences—empowering faster decision-making.

In conclusion, the pipeline exhibited high performance, robustness, and operational efficiency in handling streaming workloads. The use of GCP's managed services reduced the overhead of infrastructure management and offered scalability, fault tolerance, and ease of monitoring—all critical for production-grade stream processing systems.

CONCLUSION

This project successfully demonstrates the power and flexibility of building real-time data streaming pipelines using Google Cloud Platform (GCP) and Apache Beam. By leveraging a combination of managed services such as Pub/Sub, Dataflow, BigQuery, and Container Registry, the system achieves seamless integration, scalability, and fault-tolerance without the need to provision or manage any underlying infrastructure.

The modular nature of Apache Beam allowed us to define a clear and maintainable ETL pipeline, where each transformation—such as parsing, validation, and schema enforcement—was independently testable and traceable. The pipeline's deployment using Dataflow Flex Templates proved

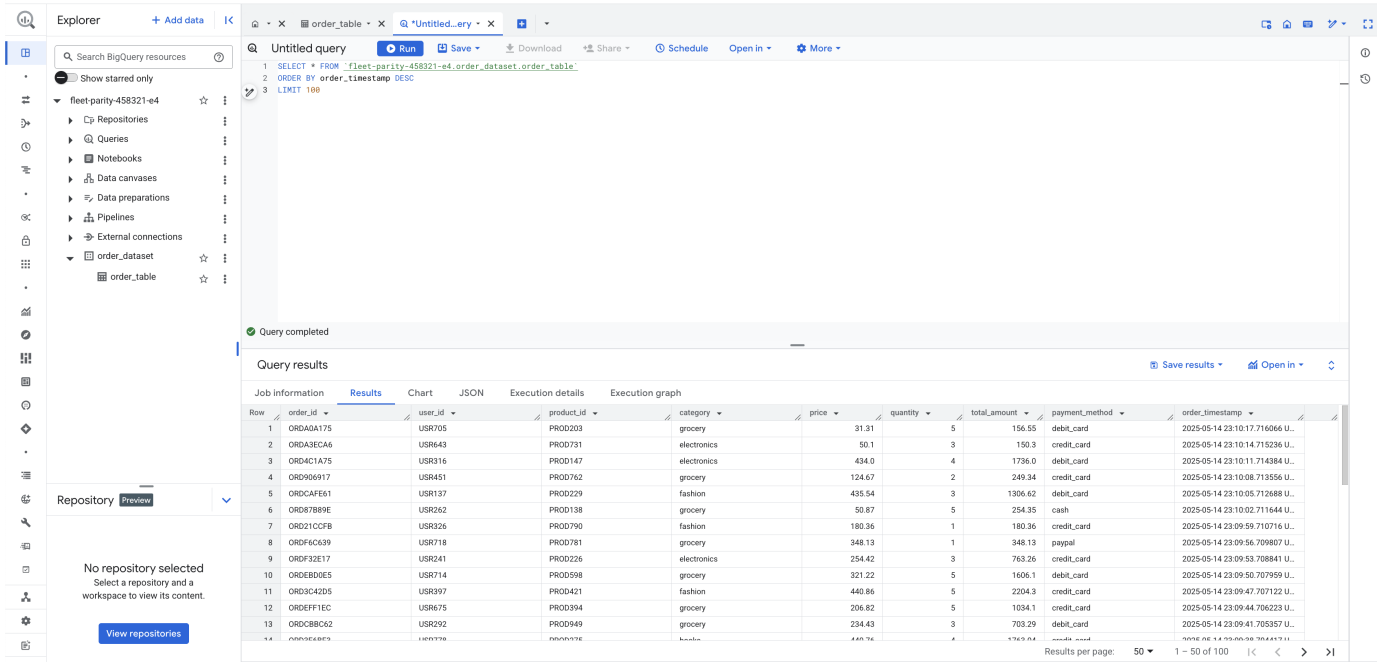


Fig. 5. Order Table in BigQuery



Fig. 6. Real-time Tableau dashboard showing aggregated metrics

highly advantageous, enabling containerized, reusable, and parameterized executions with minimal manual intervention.

The pipeline’s performance validated the architectural design, with ingestion latencies consistently below one second and high reliability in message delivery. Additionally, the use of real-time dashboards through BigQuery and Tableau highlighted the practical utility of such systems in business analytics and decision support. The ability to monitor key metrics like order value, category distribution, and payment trends in near real-time can empower organizations to act proactively rather than reactively.

Beyond technical implementation, this project also provided

valuable experience in debugging cloud-native systems, configuring IAM roles, managing Docker workflows, and setting up cross-service communication. These skills are critical when deploying production-ready data infrastructure in enterprise environments.

Looking forward, several enhancements could be introduced to elevate this system further. These include integrating machine learning models for predictive analytics (e.g., forecasting sales trends or detecting anomalies in transactions), extending the pipeline to support multiple data sources, and incorporating windowing and triggering logic to support session-based aggregations. Additionally, connecting to alerting systems such as Cloud Functions or Slack webhooks could automate responses to specific events or thresholds.

In summary, this project not only serves as a blueprint for building low-latency, cloud-native ETL systems but also demonstrates how real-time data processing can unlock operational insights and competitive advantages for data-driven organizations.

ACKNOWLEDGMENTS

Thanks to Professor Yiming Zeng for course support and Google for GCP credits.