# DATA ENGINEERING

## Real Time YouTube Comment Sentiment Analysis

Presented by -

Rohan Mahaveer
Manoj Saligrama Harisha
Namratha Prakash
Valeska Joshna Dsouza

# INTRODUCTION

- Domain: real-time analysis of YouTube comments.

- **Goal:** Build an end-to-end pipeline from ingestion → batch processing → streaming → ML predictions → storage and monitoring

**This Pipeline Delivers:**

- Automated, serverless GCP architecture that transforms chaotic comment streams into actionable sentiment metrics
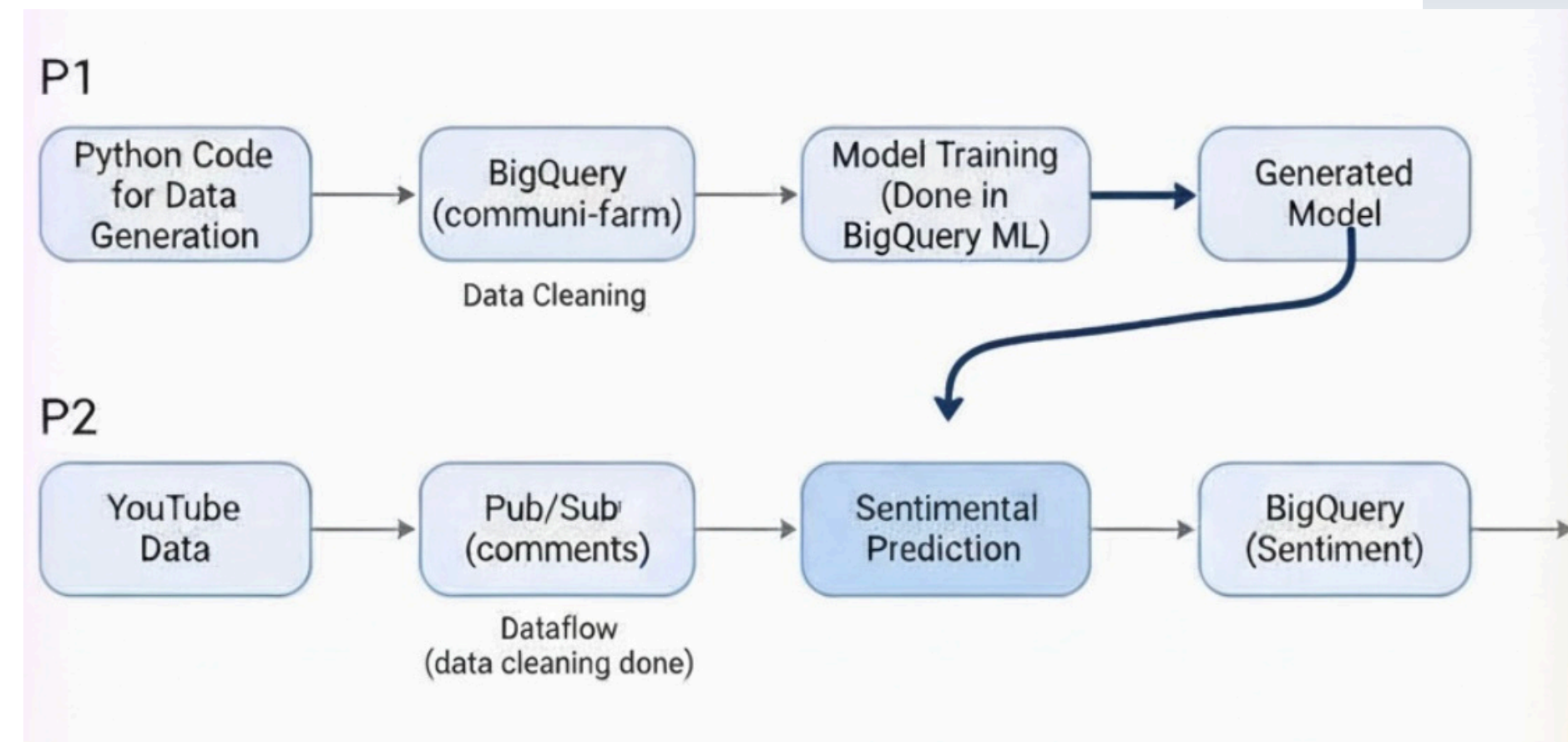
- From Raw Comments → Strategic Insights .

# PROBLEM STATEMENT

- YouTube comments contain rich signals about viewer satisfaction, content quality, and emerging issues.

- Manual review and purely batch analytics do not scale for high-volume channels or real-time moderation.

- We need a cloud-based system that can analyse comment sentiment both historically (batch) and in near real-time (streaming).

- Our project focuses on real-time YouTube comment sentiment analysis using a two-phase big data architecture

# TWO PHASE ARCHITECTURE

**Phase 1: Training**

Batch pipeline in Airflow + BigQuery for cleaning, aggregation, feature engineering, and BigQuery ML model training.

**Phase 2: Inference**

Streaming via Pub/Sub and Dataflow, applying the trained BigQuery ML model to new comments hourly and storing predictions in BigQuery.

# OBJECTIVES

- Design a simple, scalable cloud architecture for real-time YouTube comment sentiment analysis with 4–5 core GCP components

- Implement a batch pipeline (Airflow → BigQuery) to ingest, clean, aggregate comments and train a sentiment model with BigQuery ML

- Implement a streaming pipeline (Pub/Sub → Dataflow → BigQuery) that scores new comments in micro-batches every 30–60 seconds and runs hourly sentiment analytics

- Apply Little's Law and basic queuing theory to estimate throughput, latency, and worker/consumer sizing for both batch and streaming phases

- Set up basic monitoring and logging for data volume, processing time, and error rates, and provide clear documentation, diagrams, and demo instructions

# TOOLS USED

**Orchestration Tool**

Airflow

**Streaming Platform**

Pub/Sub

**Data Store**

BigQuery

**Pipeline between Pub/Sub & BigQuery**

Dataflow

**Model Training**

BigQueryML

# PROBLEM STATEMENT

- YouTube comments contain rich signals about viewer satisfaction, content quality, and emerging issues.

- Manual review and purely batch analytics do not scale for high-volume channels or real-time moderation.

- We need a cloud-based system that can analyse comment sentiment both historically (batch) and in near real-time (streaming).

- Our project focuses on real-time YouTube comment sentiment analysis using a two-phase big data architecture

# DATA GENERATION & CLEANING

- Python script generates synthetic YouTube comments
- **Faker** is used to create synthetic YouTube comments for training
- Script continuously inserts rows into the comments_train table in **BigQuery**
- Cleaned data is stored in BigQuery and later used as the training dataset for the sentiment model

# DATA GENERATION & CLEANING

# PHASE 1: MODEL TRAINING → Generated Model

### 4-STEP PIPELINE → GENERATED MODEL

- PYTHON DATA → communi-farm dataset
  Generate labeled training dataset

- BIGQUERY CLEANING
  Transform raw text for ML

- BIGQUERY ML TRAINING  CREATE MODEL
sentiment_model (LOGISTIC_REG)
  Serverless model training.
          GENERATED MODEL READY  P2 Predictions

# SENTIMENT RESULTS (P1 PIPELINE)

# YOUTUBE COMMENT DATA AND PUB/SUB

PIPELINE ARCHITECTURE

1.**YouTube Data API v3** - Source of video data

2.**Cloud Pub/Sub** - Message broker for streaming

3.**Cloud Run** - Serverless compute for processing

4.**Push Subscription** - Event-driven integration

DATA FLOW

1. **Trigger**(Youtube webhook or API poll

2. **Publish** (Message to Pub/Sub topis)

3. **Process** (Cloud Run Handels the Event along with cleaning the data aswell)

ARCHITECTURE BENEFITS

**Real-time** - Event driven processing with minimal latency.

# SENTIMENT PREDICTICTION



- Stores processed sentiment results
- Enables analytical queries
- Can feed downstream ML systems

# DATA IN BIGQUERY

# OUTPUT EXPLANATION

| Comment | Score | Magnitude | Explanation |
|---|---|---|---|
| "Ok." | 0.0 | 0.0 | No emotion |
| "Nice video!" | 0.6 | 0.6 | Positive, moderate |
| "I absolutely love this amazing masterpiece!!!" | 0.9 | 2.3 | Strong positive |
| "I hate this so much and it ruined my day." | -0.9 | 2.1 | Strong negative |

◆ **Sentiment_score**

- Range: -1 to +1
- Measures polarity
- +1 → Very positive
- 0 → Neutral
- -1 → Very negative

◆ **Sentiment_magnitude**

- Range: 0 to ∞
- Measures emotional intensity
- Higher value = stronger emotion

# FUTURE SCOPE

- Implement an Airflow (Cloud Composer) DAG to **automate**:
    * Real-time data ingestion
    * Model training
    * Model evaluation
- Apply **performance-based model** selection:
    * Compare evaluation metrics (Accuracy, F1 score).
    * Automatically promote only the best-performing model to production.
- Replace synthetic data generation (Faker) with **real-time YouTube comment** data using YouTube Data API for realistic, production-grade inputs.

# Demo.

# Thank you