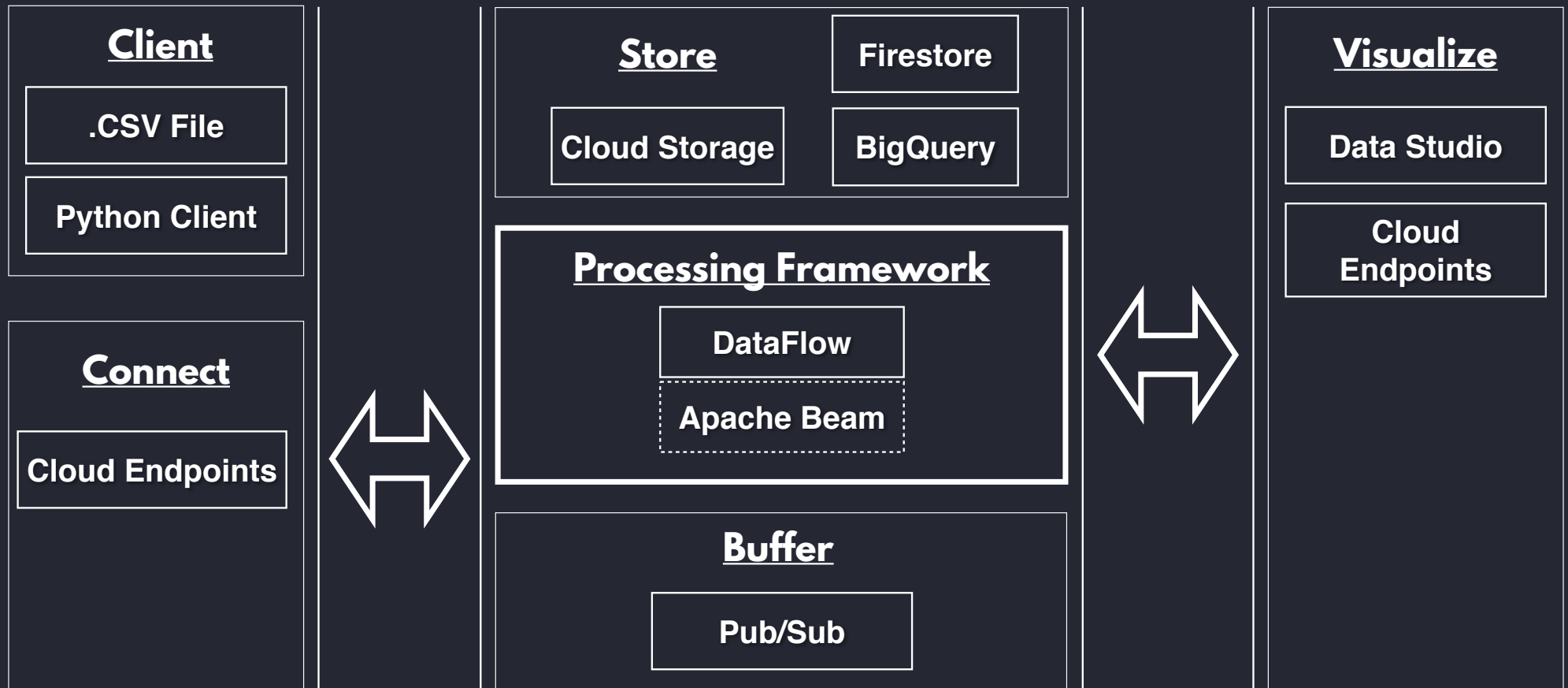


# Platform Design

# Data Engineering on GCP



# Client

- Python Client
- Reads .csv

```
data.csv
1 InvoiceNo,StockCode,Description,Quantity,InvoiceDate,UnitPrice,CustomerID,Country
2 536365,85123A,WHITE HANGING HEART T-LIGHT HOLDER,6,12/1/2010 8:26,2.55,17850,United Kingdom
3 536365,71053,WHITE METAL LANTERN,6,12/1/2010 8:26,3.39,17850,United Kingdom
4 536365,84406B,CREAM CUPID HEARTS COAT HANGER,8,12/1/2010 8:26,2.75,17850,United Kingdom
5 536365,84029G,KNITTED UNION FLAG HOT WATER BOTTLE,6,12/1/2010 8:26,3.39,17850,United Kingdom
6 536365,84029E,RED WOOLLY HOTTIE WHITE HEART.,6,12/1/2010 8:26,3.39,17850,United Kingdom
7 536365,22752,SET 7 BABUSHKA NESTING BOXES,2,12/1/2010 8:26,7.65,17850,United Kingdom
```

- Select data e.g. from today or number of lines

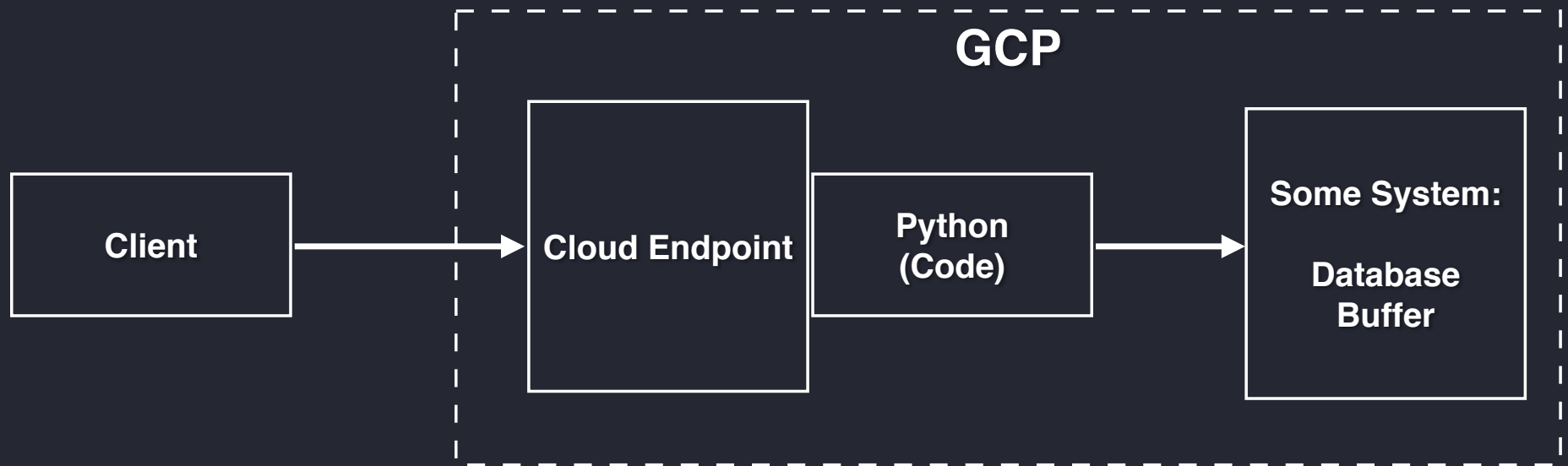
- Transforms each line into JSON string

```
1 {
2     "InvoiceNo": 536365,
3     "StockCode": "85123A",
4     "Description": "WHITE HANGING HEART T-LIGHT HOLDER",
5     "Quantity": 6,
6     "InvoiceDate": "12/1/2010 8:26",
7     "UnitPrice": 2.55,
8     "CustomerID": 17850,
9     "Country": "United Kingdom"
10 }
```

- Writes each JSON string into sink (API Gateway)

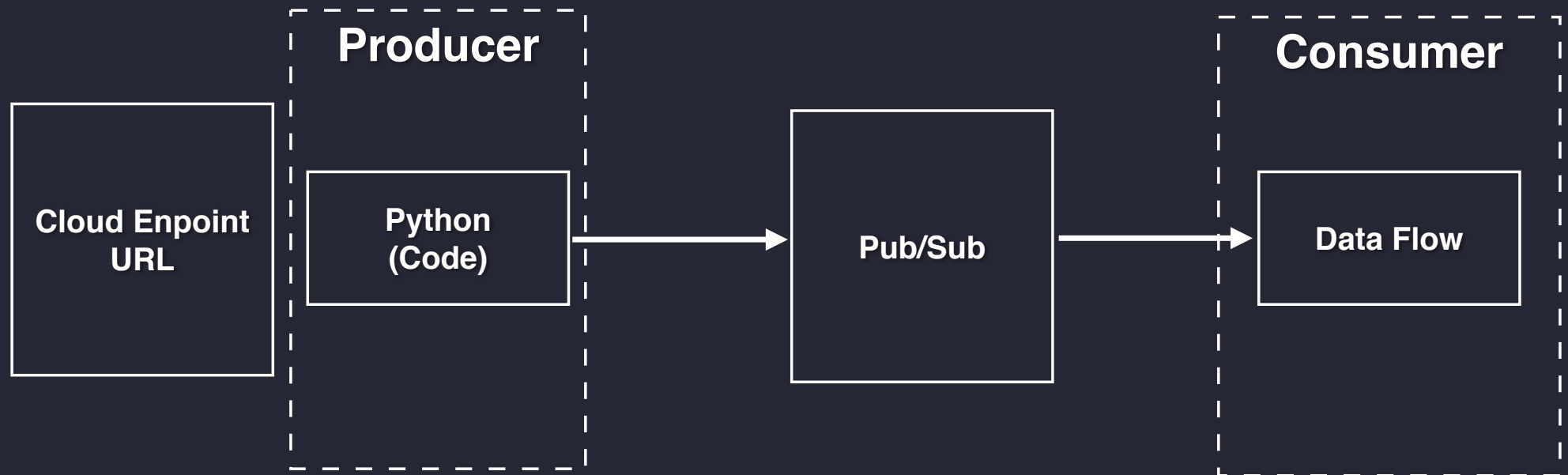
# Connect

- Cloud Endpoint
- Python



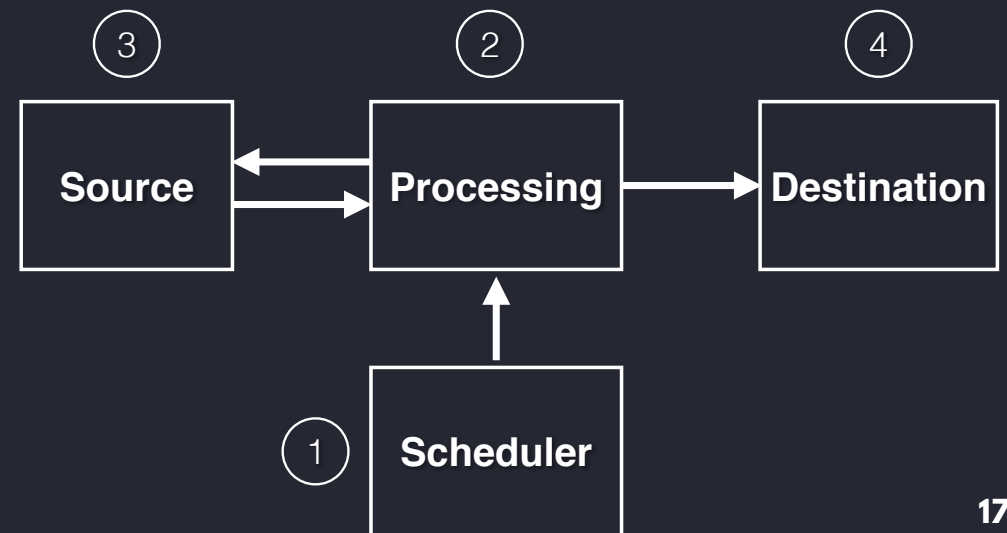
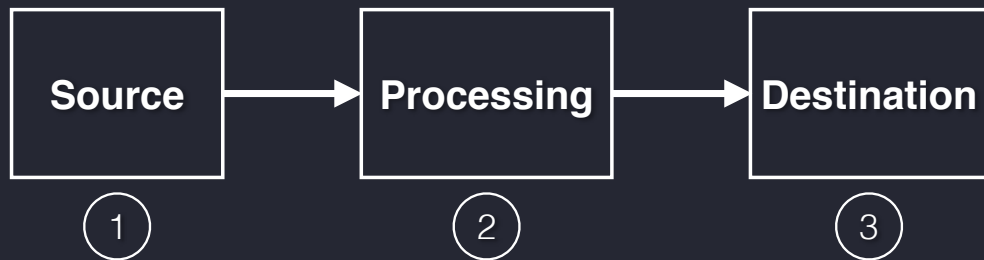
# Buffer

- Pub/Sub



# Process

- Streaming Processing
  - DataFlow with triggers on Source
  - Continuous Process
- Batch Processing
  - DataFlow
  - Cloud Composer for Scheduling



# Store

- Cloud Storage
- Firestore NoSQL
  - Document based Store like MongoDB
  - Transactions
- BigQuery Data Warehouse
  - Analytics Layer
  - Distributed Storage And Processing

# Visualize

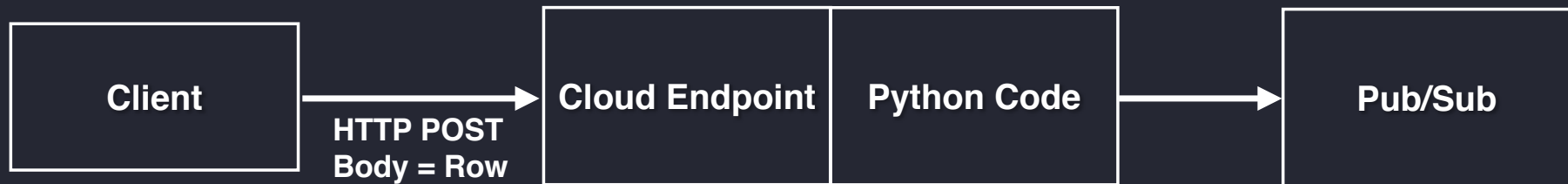
- APIs
  - Access for Apps, UIs..
  - Execute Queries and Transactions
  - Simple, Stateless
- Tableau
  - Business Intelligence Tool
  - Installed On Your PC
  - Connects to Redshift



# Data Pipelines

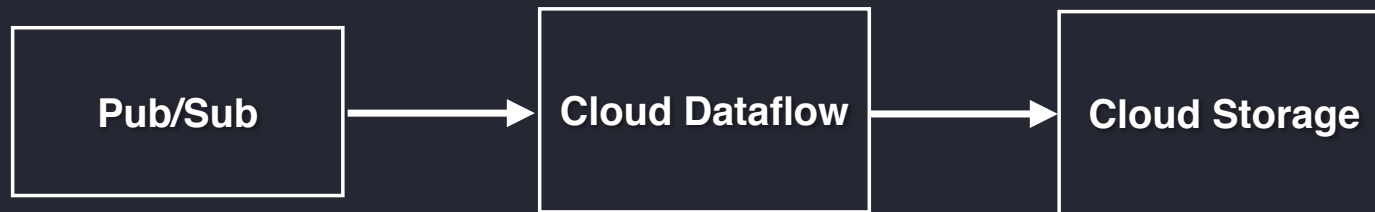
# Data Ingestion Pipeline

- Client
- Simulates Streaming
- Sends CSV Rows as JSON
- Cloud Endpoint API
- Python Code For Write Into Pub/Sub
- Pub/Sub buffer for streaming



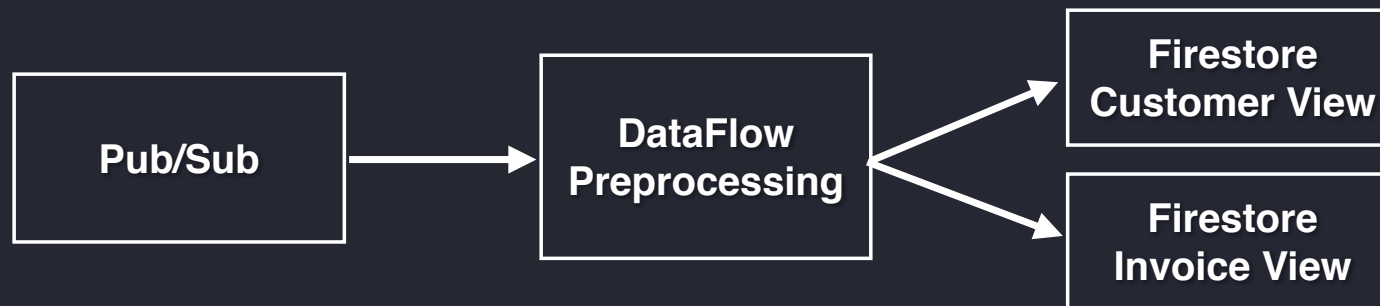
# Stream to Cloud Storage Pipeline

- Pub/Sub insert triggers Dataflow
- Dataflow waits for some time
- Writes all messages in queue to Cloud Storage as file



# Stream to Firestore Pipeline

- Pub/Sub Insert triggers Dataflow insert for Firestore
- Dataflow reformats/preprocesses messages
- Dataflow writes customer data (customer + invoices)
- Dataflow writes invoice data (invoice + stockcode)



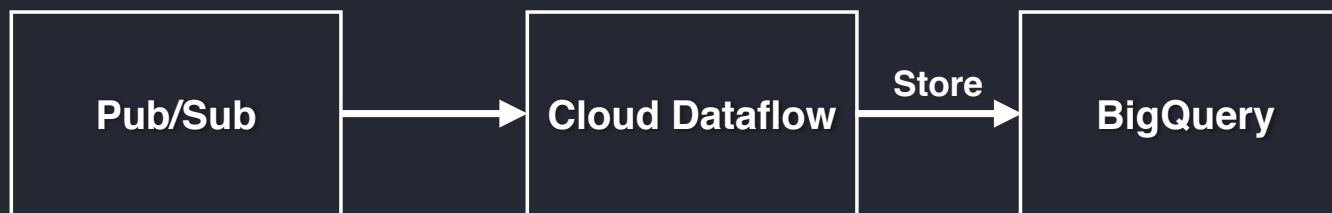
# Visualization Pipeline API

- APIs for UI (Items in Invoice)
- Data rests in Firestore table Invoices
- Client requests Items for InvoiceNo (Request parameter)
- Python of Cloud Endpoint triggered by API queries Firestore with InvoiceNo



# Visualization Pipeline BigQuery Data Warehouse

- Pub/Sub triggers Data Flow
- Data Flow writes messages into BigQuery



# Batch Processing Pipeline

- Bulk import Pipeline
- Triggered through Cloud Composer
- Dataflow reads from Cloud Storage
- Writes data into Firestore
- Writes into BigQuery

