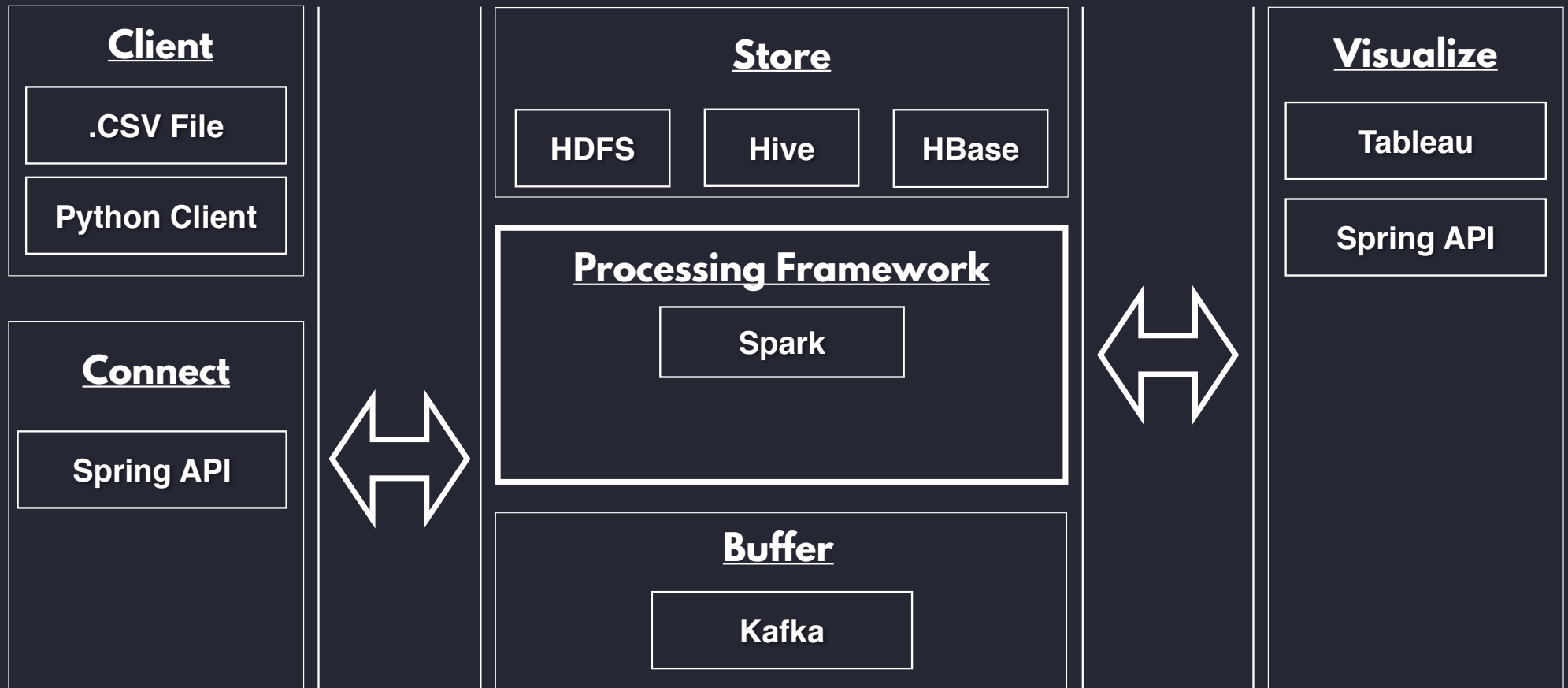


Platform Design

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

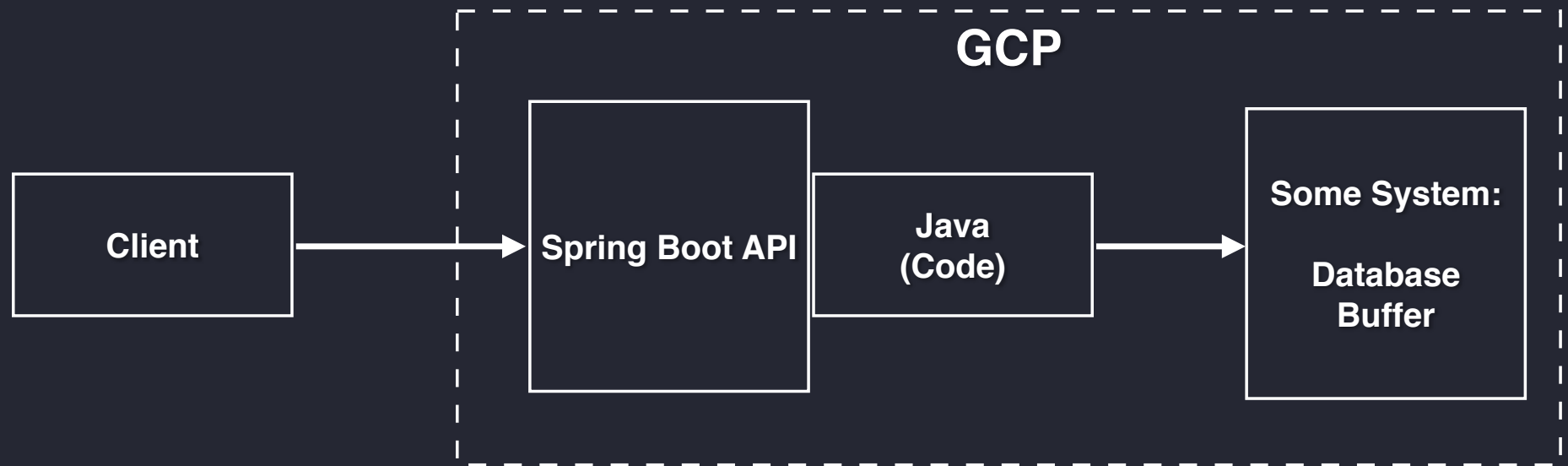
- Transforms each line into JSON string

```
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)

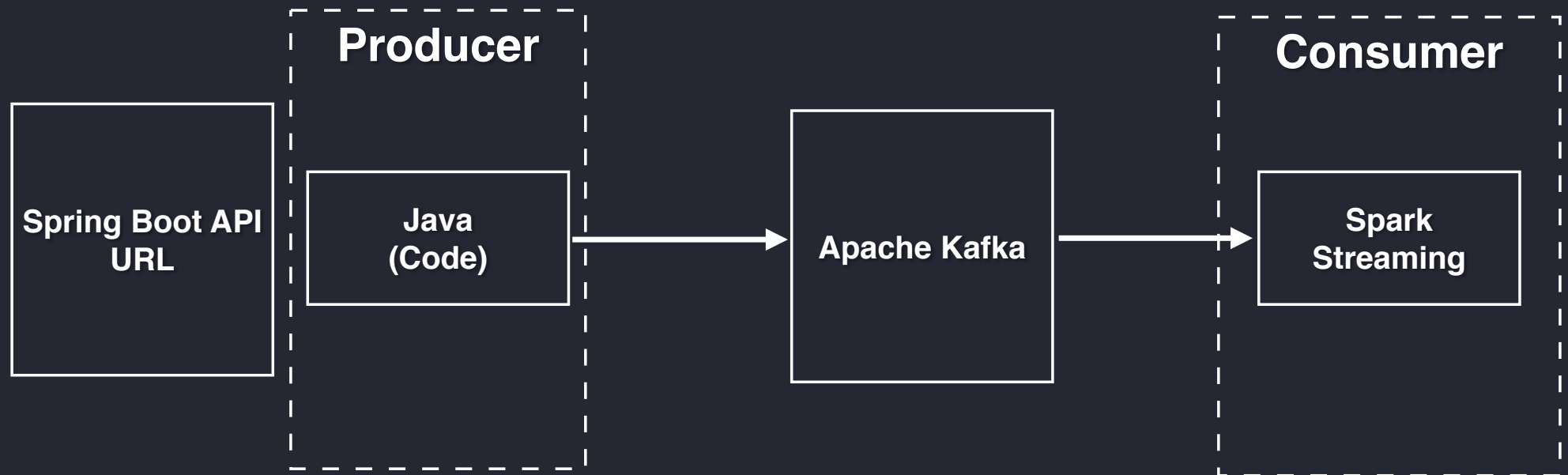
Connect

- Spring Boot
- Java



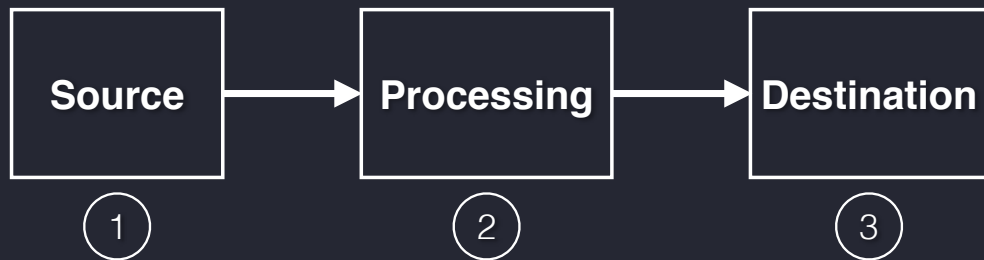
Buffer

- Apache Kafka

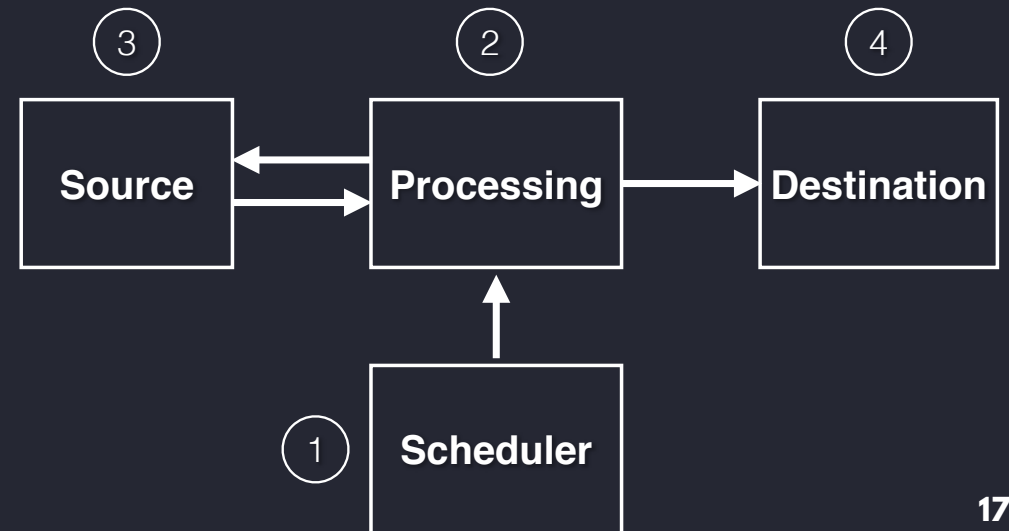


Process

- Streaming Processing
 - Spark Streaming with triggers on Source
 - Continuous Process



- Batch Processing
 - Spark Batch Processing
 - Oozie for Scheduling



Store

- HDFS (Hadoop File System)
- HBase NoSQL
 - Key Value Store
 - Transactions (possible but manual coding needed)
- Hive Data Warehouse
 - Analytics Layer processing with Spark
 - Can access data from HDFS or HBase

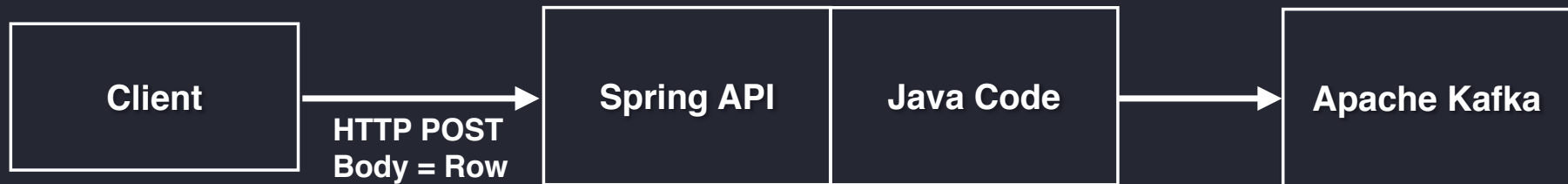
Visualize

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

Data Pipelines

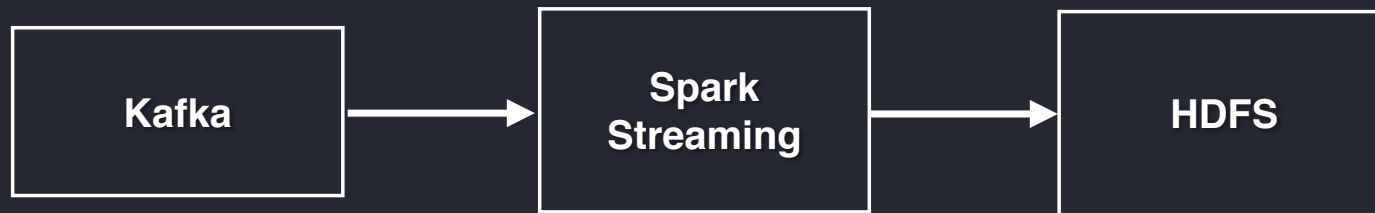
Data Ingestion Pipeline

- Client
- Simulates Streaming
- Sends CSV Rows as JSON
- Spring API
- Java Code For Write Into Kafka
- Kafka buffer for streaming



Stream to HDFS Pipeline

- Kafka insert gets fetched by Spark Streaming job
- Spark Streaming micro batches run in small intervals
- Microbatch writes all messages in Kafka to HDFS as file



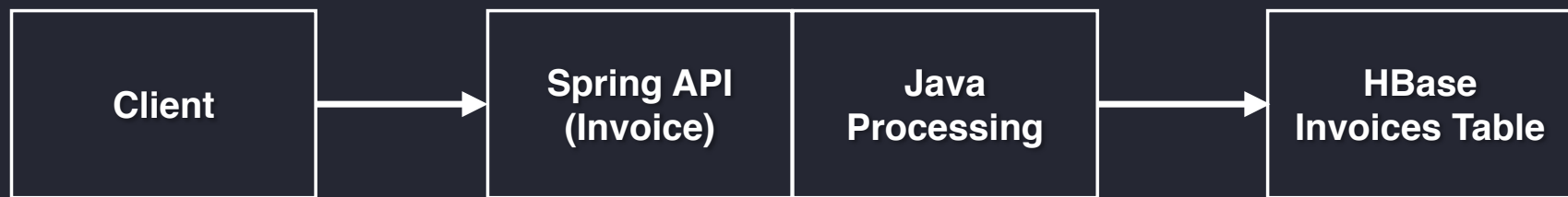
Stream to HBase Pipeline

- Kafka messages get processed by Spark
- Spark reformats/preprocesses messages
- Spark writes customer data (customer + invoices)
- Spark writes invoice data (invoice + stockcode)



Visualization Pipeline API

- APIs for UI (Items in Invoice)
- Data rests in Hbase table Invoices
- Client requests Items for InvoiceNo (Request parameter)
- Java code of Spring API queries HBase with InvoiceNo



Visualization Pipeline Hive Data Warehouse

- Messages are in Kafka
- Spark Streaming writes messages into Hive
- Visualize with Tableau



Batch Processing Pipeline

- Bulk import Pipeline
- Triggered through Oozie
- Spark reads from HDFS
- Writes data into HBase
- Writes into Hive

