# Data Pipeline with Kafka

Dr. Mole T.Y. WONG  @  HK OSCON 2018
2018 / 06 / 16 - 17

1

# whoami

## Why

深入了解用戶行為, 洞悉可行的改善方法
Understand our users.
Provide actionable insights.
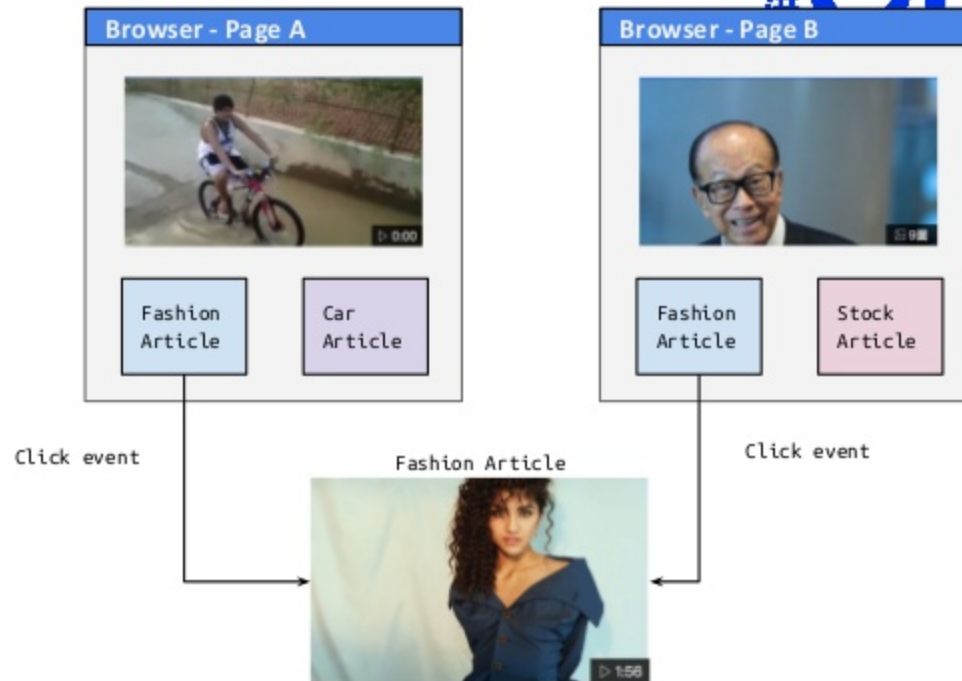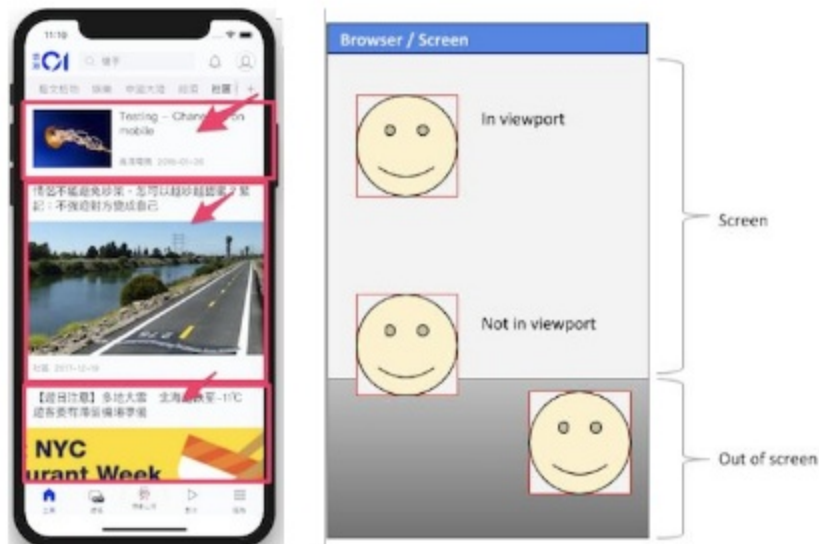
# How

以數據驅動產品方向
Data driven: steer our product direction.

# What

數據：定義、收集、處理、洞見
Data: definition, ingress, process, insight.

Click-Through Rate VS Pageview



Fashion Article

Traffic Source Analysis

Data-Driven Product Development

User Reading History

NLP Content-based Clustering

Collaborative filtering
Image source: wikipedia
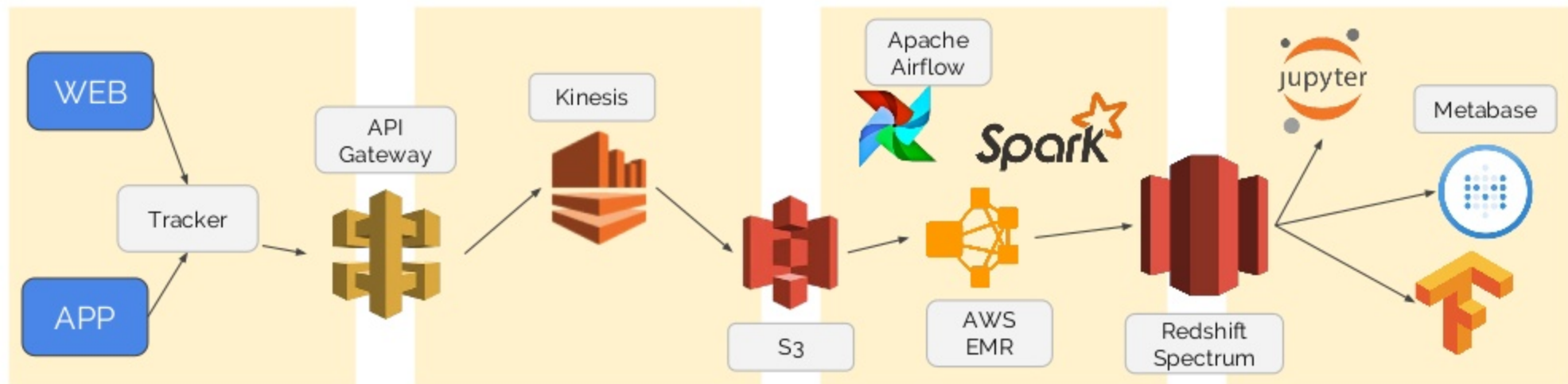
Machine Learning Products

Personalized Recommendation Feed

# Outline

- Data pipeline - what is it?

- Kafka - roles in a data pipeline
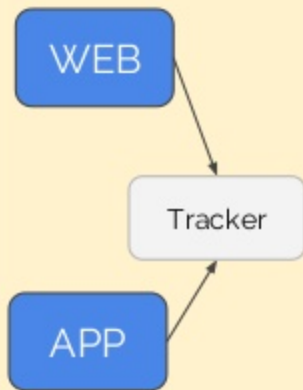
- Other use cases of Kafka

# Typical Data Pipeline Setup

## Data Ingress

JS Library（WEB）
Native Library (APP)

Google Analytics
Mixpanel
Matomo (Piwik)



## Data Tracker

- **Nature**
  - Lightweight
  - Programmable
- **Capability**
  - Page view / Screen view
  - Custom events
  - Device identification
  - Session management

Different Aspects of a Data Pipeline

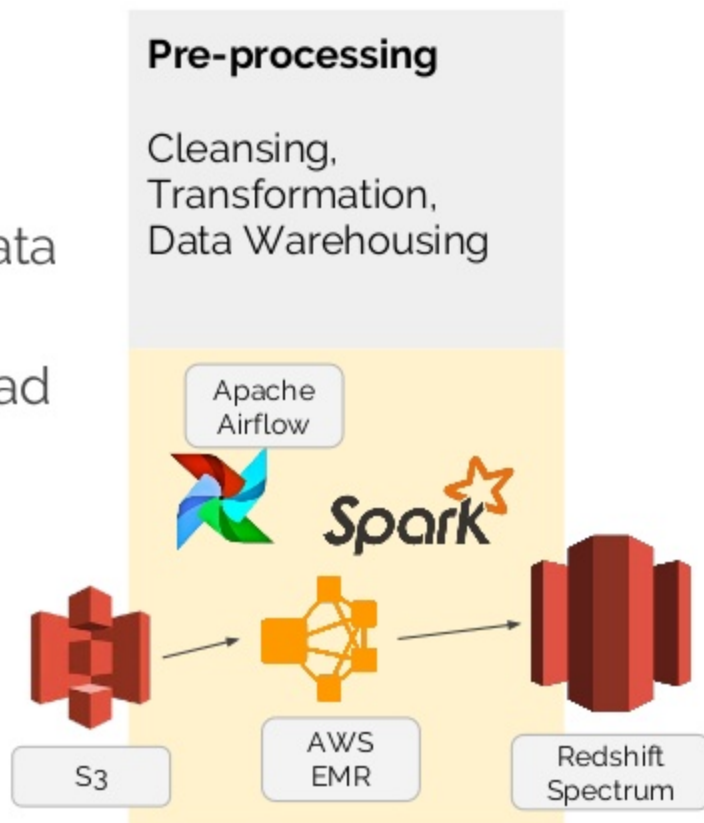Different Aspects of a Data Pipeline

# Pre-processing

- **Main Roles**
  - Avoid direct querying raw data
  - Cleansing
  - ETL - Extract, Transform, Load
  - Scheduling
- **Characteristics**
  - Defining data sets
  - Time-frame-based queries

Different Aspects of a Data Pipeline

**Pre-processing**

Cleansing,
Transformation,
Data Warehousing

Apache
Airflow

Spark

S3

AWS
EMR

Redshift
Spectrum

# Application

- **Main Roles**
  - KPI VS Exploration
  - Operators VS Data Scientists
  - Planned VS Ad-hoc queries
- **Characteristics**
  - Production-grade data
  - Fast is a must

| Pre-processing | Application |
|---|---|
| Cleansing, Transformation, Data Warehousing | Dashboard, Reporting, Recommendation Engine, etc |



S3 → AWS EMR → Redshift Spectrum → Jupyter / Metabase / etc

Different Aspects of a Data Pipeline

What is Kafka? https://kafka.apache.org/  Main Contributor: Gene NG

# Data Pipeline with Kafka

WEB

APP

Tracker

API
Gateway

Kafka
Connect
API

Kafka
Connect
API

jupyter

Metabase

What is Kafka?

Optional: data persists in S3

# Data Pipeline with Kafka



WEB

APP

Tracker

API Gateway

Kafka Connect API

Kafka Connect API

Jupyter

Metabase

What is Kafka?

Optional: data persists in S3

# Basics: Producer-Consumer Model

Producer

Consumer

```
while(1) {
  var e = produce_event()
  producer.produce(e)
}
```

```
while(1) {
  var m = consumer.poll()
  consume_msgs(m)
}
```

What is Kafka - terminology

# Connect API

- For database / data source
- Wrapped consumer & producer code
- Nice thing: config file only!



```
Data          Kafka                  Kafka         Data
Source   →    Connect   →  ⧉  →     Connect   →    Sink
              API                    API
```

What is Kafka  -  terminology

# Connect API - common connectors

| | |
|---|---|
| JDBC - MySQL, PgSQL | S3 |
| HDFS | ElasticSearch |

Data Source → Kafka Connect API → Kafka Connect API → Data Sink
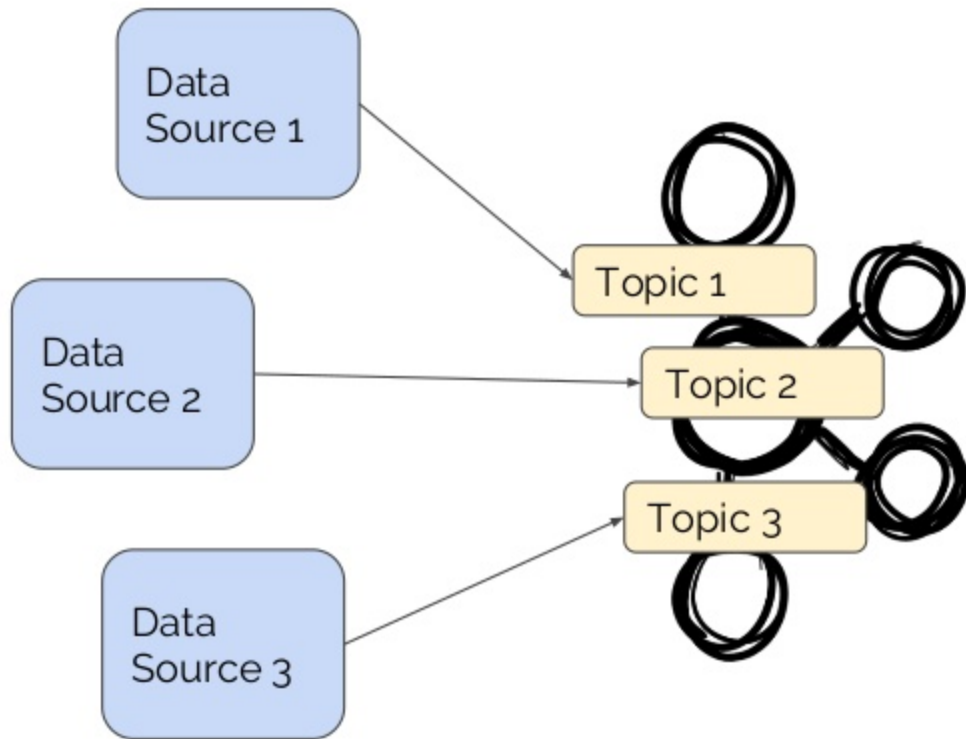
What is Kafka  -  terminology

Kafka Connect

## Data Topic Model

- One-to-one (most common)

## Feature

- Autonomous
  - Loads data from sources whenever changes occur
- Storage
  - Writes data to the hosted HDD
  - Optional: sync data to S3

```
1   name=test-source-sqlite-jdbc-autoincrement

2   connector.class=io.confluent.connect.jdbc.JdbcSourceConnector

3   tasks.max=1

4   connection.url=jdbc:sqlite:test.db

5   mode=incrementing

6   incrementing.column.name=id

7   topic.prefix=test-sqlite-jdbc-
```

# Kafka Connect - Source Property File

```
1    name=test-source-sqlite-jdbc-autoincrement

2    connector.class=io.confluent.conne

3    tasks.max=1

4    connection.url=jdbc:sqlite:test.db

5    mode=incrementing

6    incrementing.column.name=id

7    topic.prefix=test-sqlite-jdbc-
```

Topic naming convention
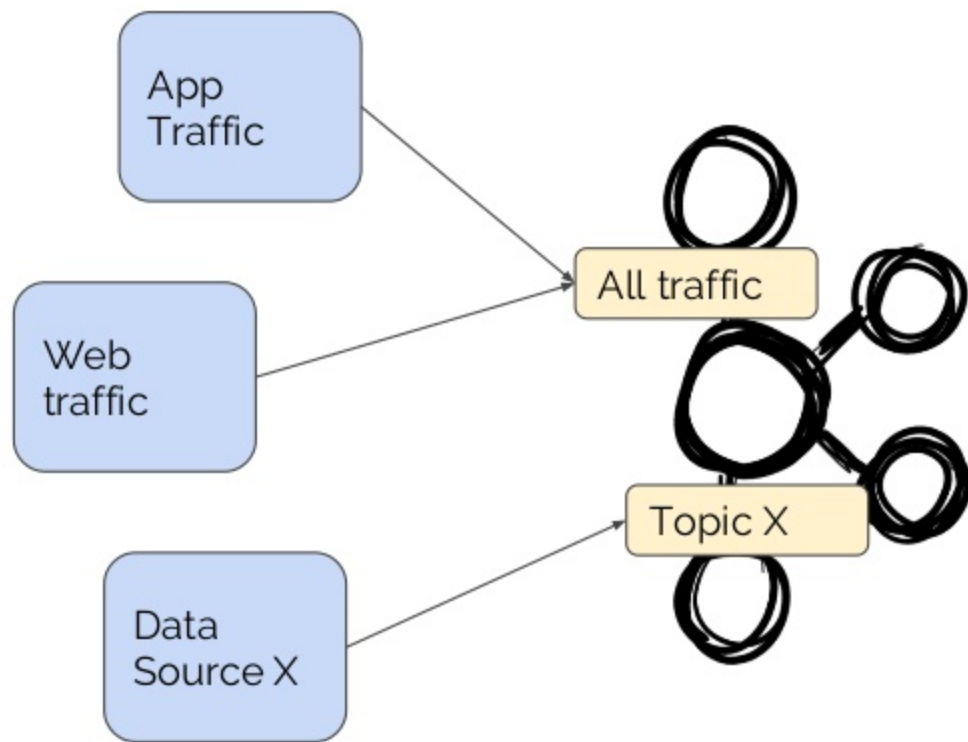
- Prefix, and
- DB table name

How it works:
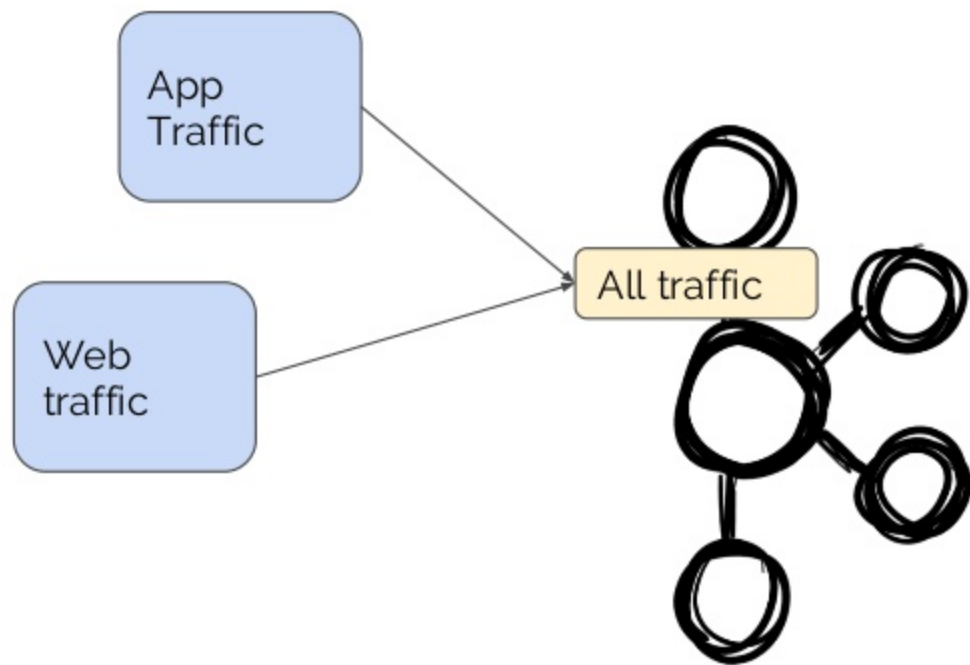- Each table implies one topic.

# Kafka Connect - Source Property File

# Data Topic Model



- One-to-one (most common)
- Many-to-one

App Traffic → All traffic

Web traffic → All traffic

Data Source X → Topic X

Kafka Connect

App
Traffic

Web
traffic

All traffic

**Schema-less**
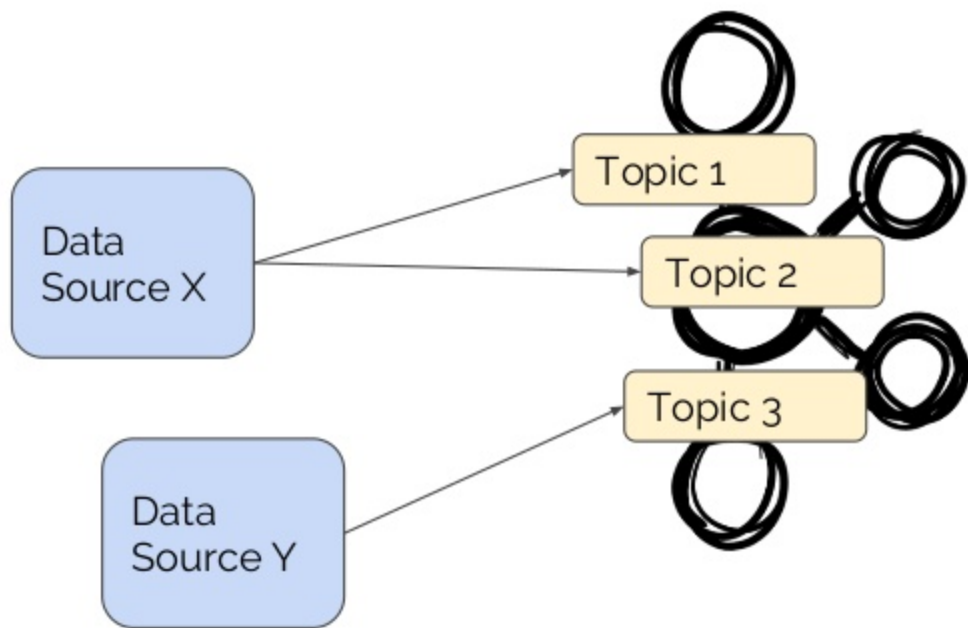
- Practically, you can write any types of data to the topic
- Most common choice is Avro

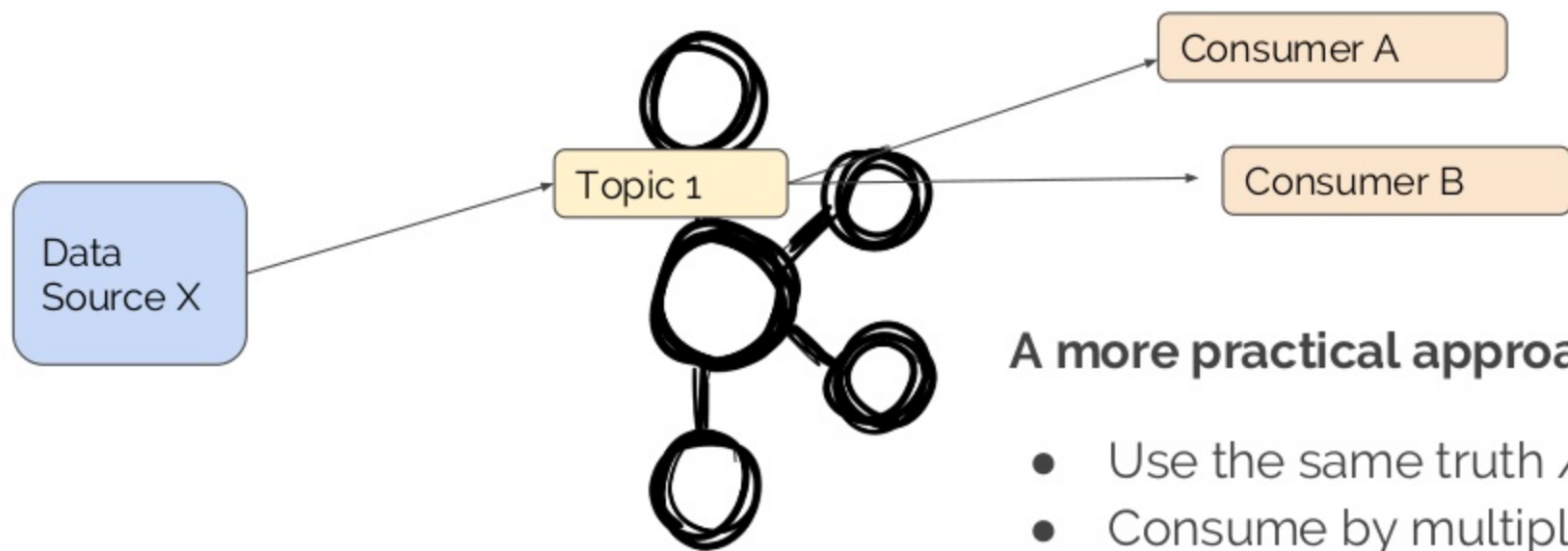Btw, Avro is an open-source library for schema specification and data serialization.

Kafka Connect

## Data Topic Model

- One-to-one (most common)
- Many-to-one
- One-to-many (most rare)

Kafka Connect

**Data Source X** → **Topic 1** → **Consumer A**, **Consumer B**

## A more practical approach

- Use the same truth / data
- Consume by multiple guys!

Kafka Connect

# Takeaway Messages

- Producers and consumers are actors
  - Push data to or pull data from Kafka


- Connect API automates the above actions
  - Work nicely with databases

# Data Pipeline Use Cases

| Kafka Internal - consumer's state | | | | |
|---|---|---|---|---|
| Consumer | Topic | Current Topic Position | Your last-read position | Lag behind by |
| hello_world | foobar | 1080 | 1000 | 80 |

Kafka keeps track on consumer's state:

- A consumer can always resume work-in-progress
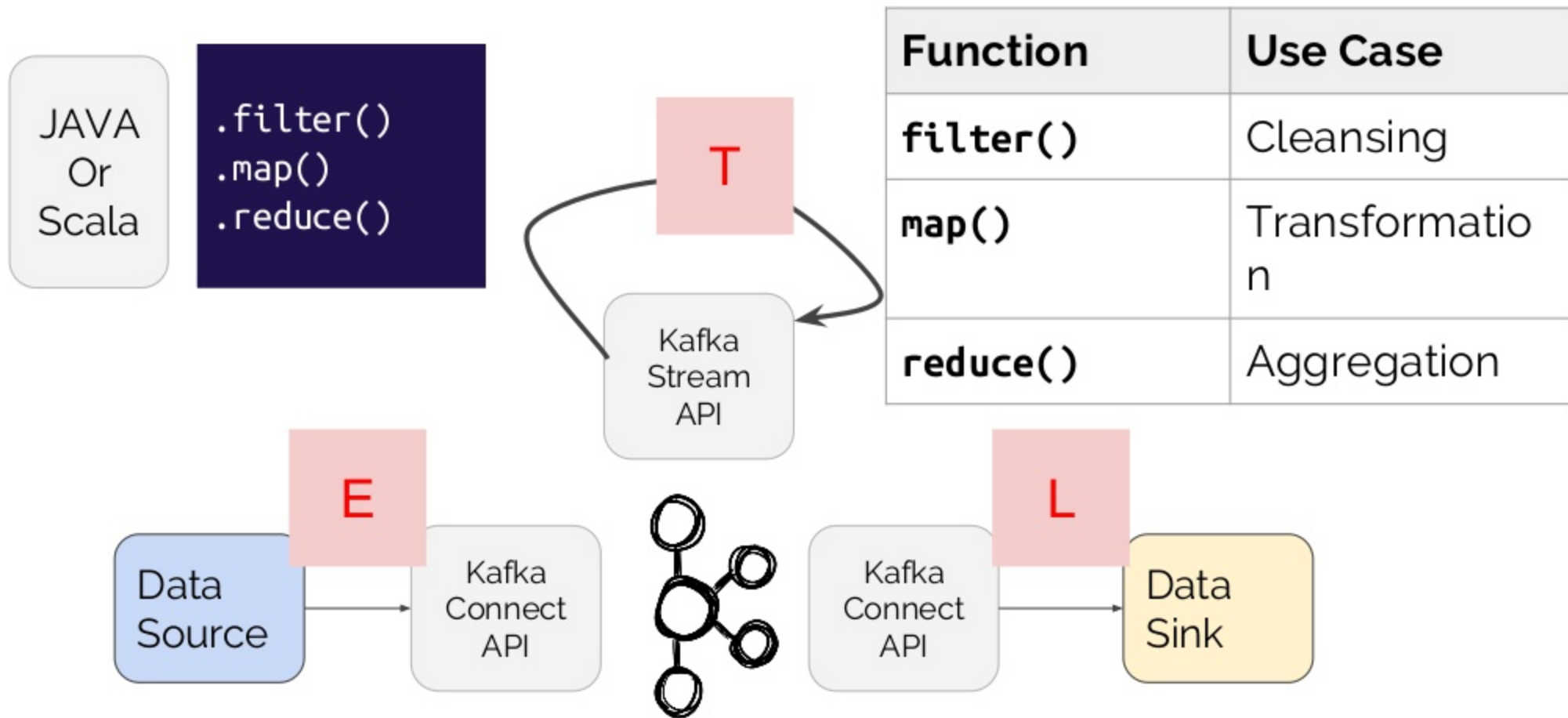- New consumer can start fresh!

Kafka Connect API

Data Sink

Kafka as a data pipeline - data resiliency

```
$ /usr/bin/kafka-consumer-groups --zookeeper zk01.example.com:2181 --describe --group

GROUP     TOPIC     PARTITION   CURRENT-OFFSET   LOG-END-OFFSET   LAG     OWNER
flume     t1        0           1                3                2       test-consumer-group_
```

Source:
https://www.cloudera.com/documentation/kafka/latest/topics/kafka_command_line.html

JAVA
Or
Scala

```
.filter()
.map()
.reduce()
```

T

Kafka
Stream
API

| Function   | Use Case        |
|------------|-----------------|
| filter()   | Cleansing       |
| map()      | Transformation  |
| reduce()   | Aggregation     |

E

Data
Source

Kafka
Connect
API

L

Kafka
Connect
API

Data
Sink

Kafka as a data pipeline - Replace ETL

```
map, filter, and reduce
explained with emoji 😂

map([🐄, 🍠, 🐔, 🌽], cook)
=> [🍔, 🍟, 🍗, 🍿]


filter([🍔, 🍟, 🍗, 🍿], isVegetarian)
=> [🍟, 🍿]


reduce([🍔, 🍟, 🍗, 🍿], eat)
=> 💩
```

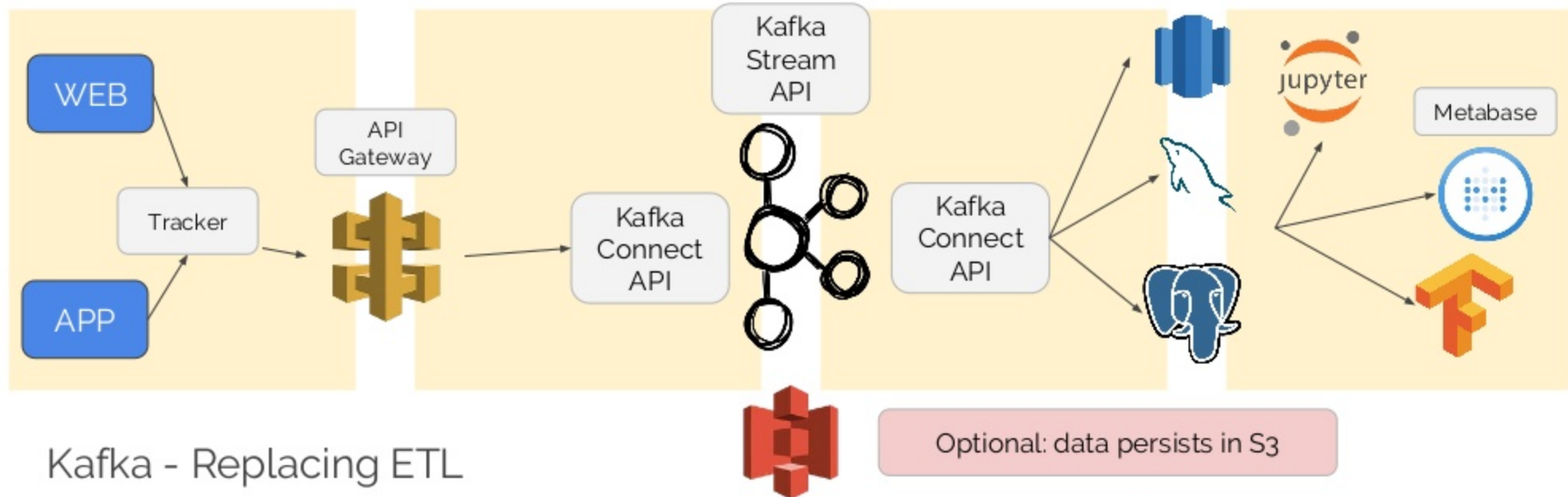Source: https://i.redd.it/yf7rw3pjiapx.jpg

```
KStream<String, String> source = builder.stream("streams-plaintext-input");
source.flatMapValues(value -> Arrays.asList(value.split("\\W+")))
      .to("streams-linesplit-output");
```
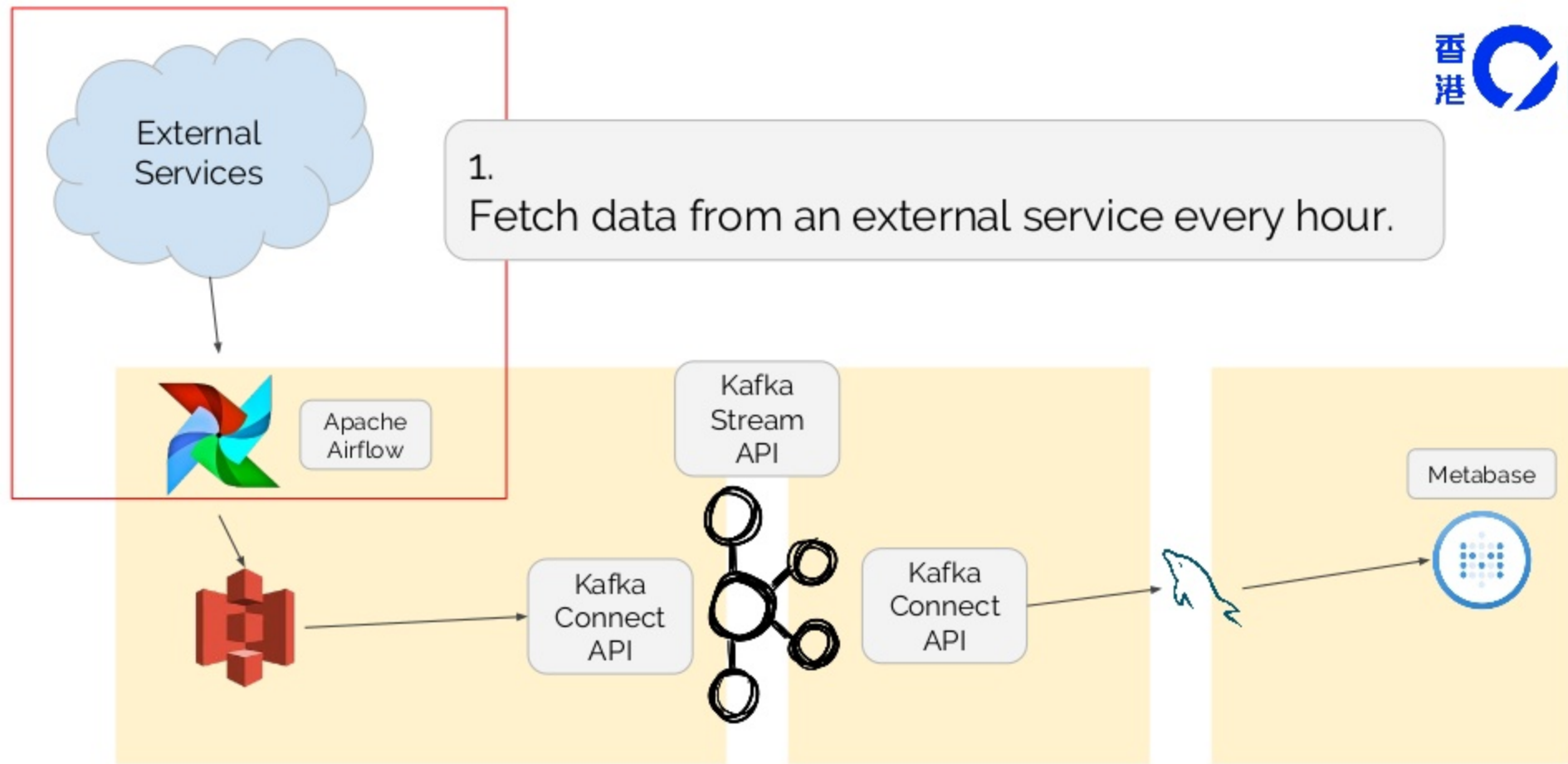
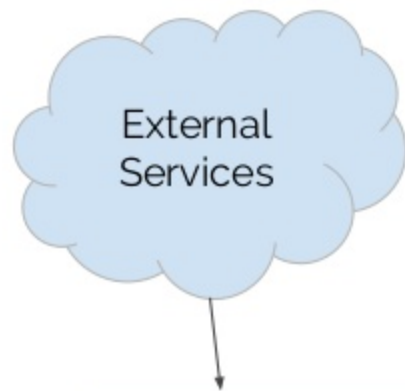A New Topic is Created!

Kafka - Streaming Example Code
Source: https://kafka.apache.org/11/documentation/streams/tutorial

# Data Pipeline with Kafka v2



Kafka - Replacing ETL

# Experimenting Kafka in HK01

External Services

1.
Fetch data from an external service every hour.

Apache Airflow

Kafka Stream API

Kafka Connect API

Kafka Connect API

Metabase

Experimenting Kafka in HK01

External Services

2.
When data arrives at S3, Kafka takes it in.

Apache Airflow

Kafka Stream API

Metabase

Kafka Connect API

Kafka Connect API

Experimenting Kafka in HK01

External
Services

3.
Stream API counts the number of new users
using certain services.

Apache
Airflow

Kafka
Stream
API

Metabase

Kafka
Connect
API

Kafka
Connect
API

Experimenting Kafka in HK01

External
Services

4.
Connect API automatically updates the MySQL table. Metabase can display the updates.

Apache Airflow

Kafka Stream API

Metabase

Kafka Connect API
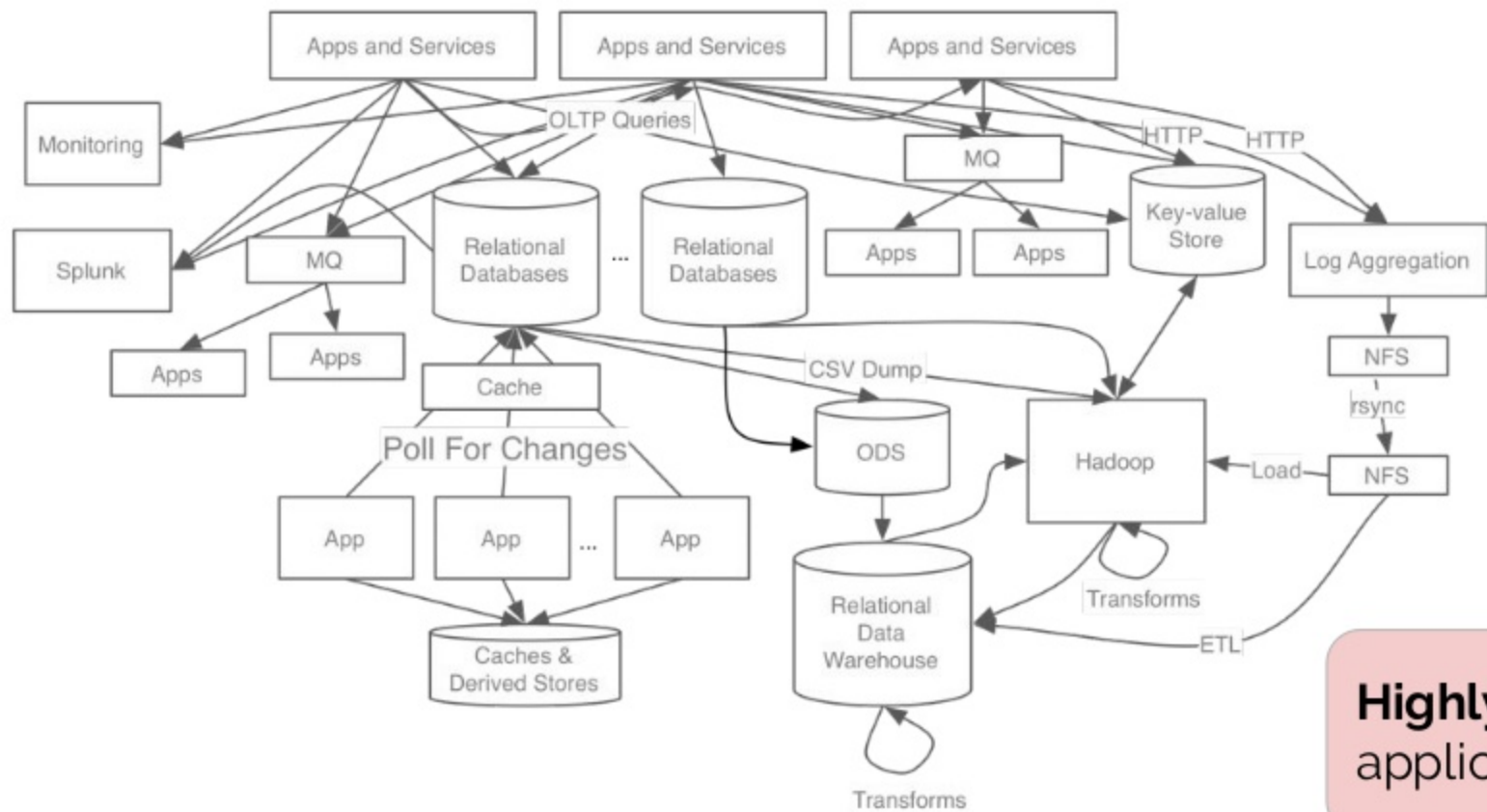
Kafka Connect API

Experimenting Kafka in HK01

# Will display live dashboard during the talk



**Experimenting Kafka in HK01**

# Other Use Cases

**Highly-coupled**: application & storage

Message Queue | Source: https://www.confluent.io/blog/stream-data-platform-1/

**As a message queue (MQ):**
- Pub/Sub
- Transformation
- Roles; clear that who are the sources and the sinks, respectively

Message Queue | Source: https://www.confluent.io/blog/stream-data-platform-1/

Things that we didn't explore

- Logs aggregation

- Database log compaction

- Event sourcing

Other Use Cases | Source: https://kafka.apache.org/uses

# Key Takeaways

Pros

1. Kafka simplifies your ETL tasks.

2. Kafka unitifies your data storage.

3. Kafka gives your other possibilities.

# Key Takeaways

Cons

1. Ops problems – scalability, HA, Zookeeper, etc.

2. Learning curve is *STEEP*.

# We Love to Share

**Mole Wong**
Data Pipeline with Apache Kafka

Day 1 17:40
Conference Hall 4-5

**Sunday Ku**
Video.js with HLS

Day 2 12:30
Conference Hall 4-5

**Ivan Ha**
React Async Rendering - Paradigm Shift After React Fiber

Day 2 15:10
Conference Hall 6

# HK01 Engineering

Build something that makes us proud!

✏️ APPLY NOW

https://goo.gl/j74Ztt