

How to Build a Scalable ETL Pipeline with Kafka

Alex

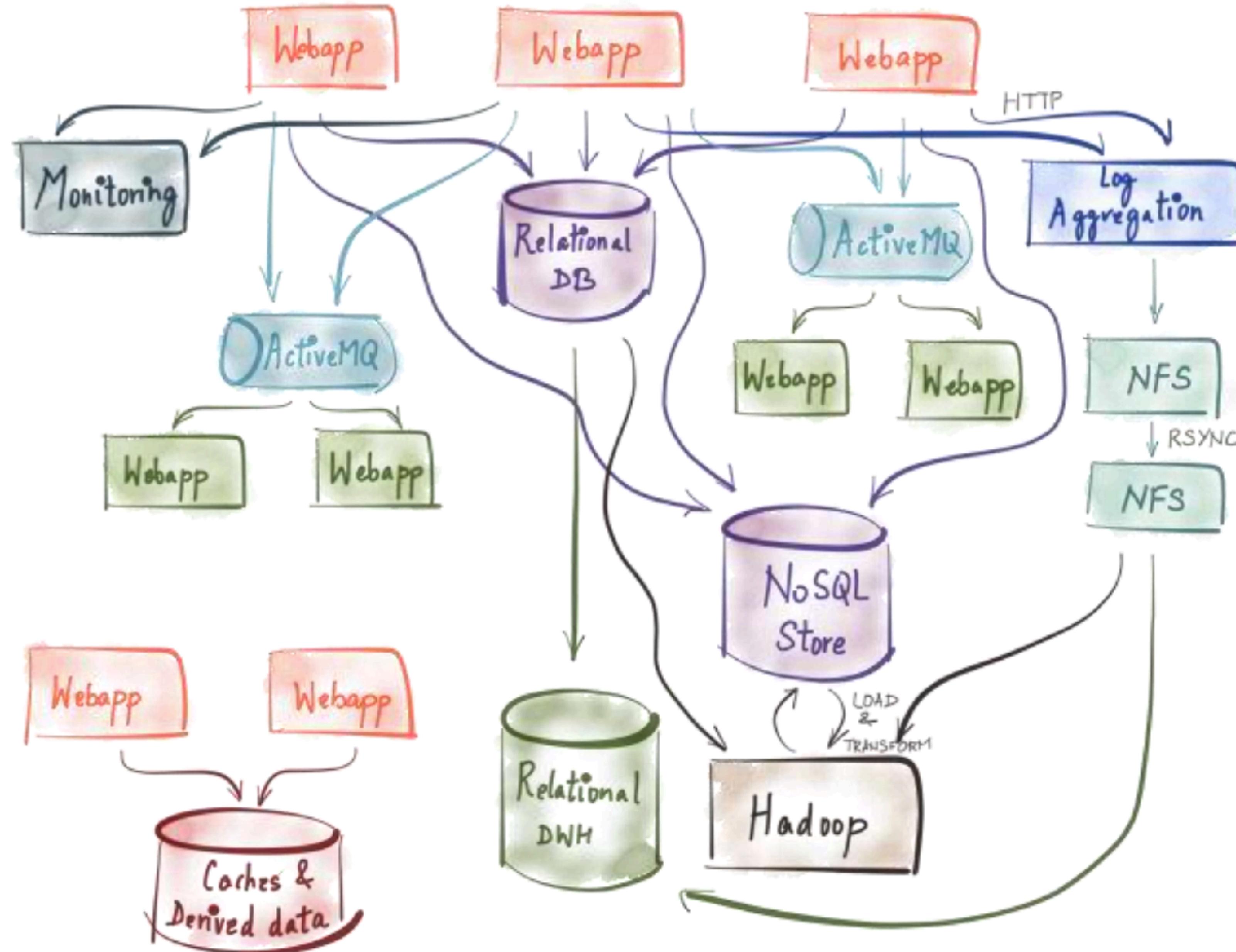
Traditional ETL



TERADATA®



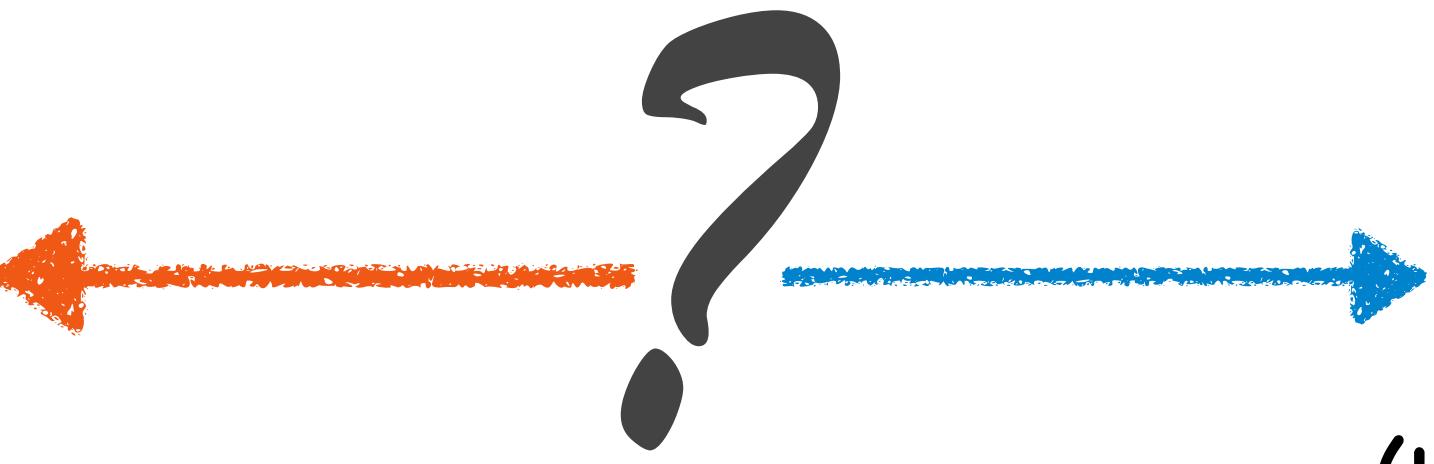




ANTI PATTERNS

1. One-off tools
2. Kitchen sink tools
3. Streaming processing frameworks

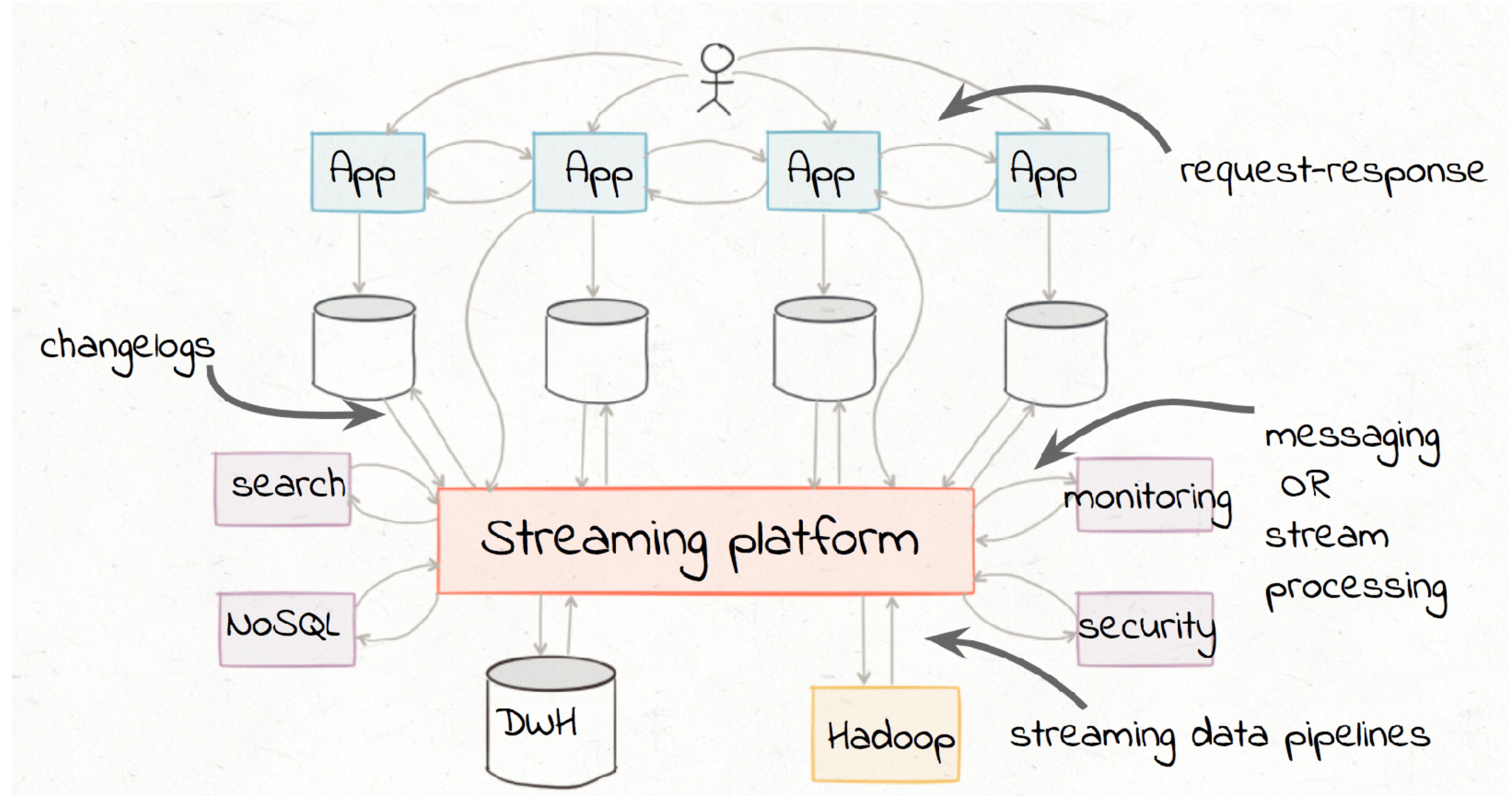
Ad-hoc
(One-off tools)

? **E,T,&L**
(Kitchen sink tools)

Data Integration

getting data to all the right places

All your data is event streams



EXAMPLE EVENT: PRODUCT VIEW

```
{  
    "time": 144704224376,  
    "User-id": 12345,  
    "product_id": 5678,  
    "Page": "product.detail",  
    ...  
}
```

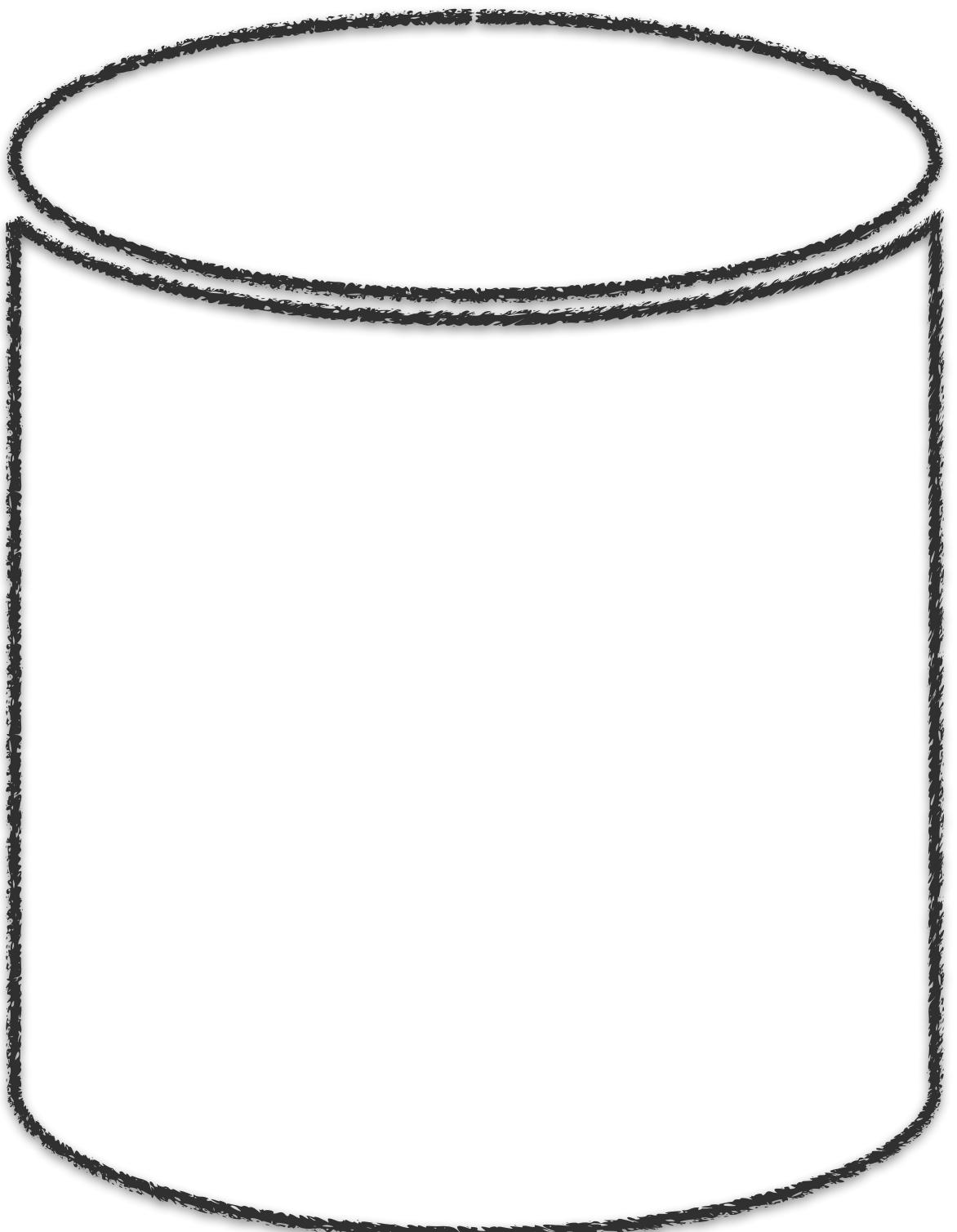
EXAMPLE EVENT: SENSORS



EXAMPLE EVENT: LOG FILES

```
jkreps-mn:~ jkreps$ tail -f -n 20 /var/log/apache2/access_log
::1 - - [23/Mar/2014:15:07:00 -0700] "GET /images/apache_feather.gif HTTP/1.1" 200 4128
::1 - - [23/Mar/2014:15:07:04 -0700] "GET /images/producer_consumer.png HTTP/1.1" 200 86
::1 - - [23/Mar/2014:15:07:04 -0700] "GET /images/log_anatomy.png HTTP/1.1" 200 19579
::1 - - [23/Mar/2014:15:07:04 -0700] "GET /images/consumer-groups.png HTTP/1.1" 200 2681
::1 - - [23/Mar/2014:15:07:04 -0700] "GET /images/log_compaction.png HTTP/1.1" 200 41414
::1 - - [23/Mar/2014:15:07:04 -0700] "GET /documentation.html HTTP/1.1" 200 189893
::1 - - [23/Mar/2014:15:07:04 -0700] "GET /images/log_cleaner_anatomy.png HTTP/1.1" 200
::1 - - [23/Mar/2014:15:07:04 -0700] "GET /images/kafka_log.png HTTP/1.1" 200 134321
::1 - - [23/Mar/2014:15:07:04 -0700] "GET /images/mirror-maker.png HTTP/1.1" 200 17054
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /documentation.html HTTP/1.1" 200 189937
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /styles.css HTTP/1.1" 304 -
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /images/kafka_logo.png HTTP/1.1" 304 -
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /images/producer_consumer.png HTTP/1.1" 304 -
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /images/log_anatomy.png HTTP/1.1" 304 -
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /images/consumer-groups.png HTTP/1.1" 304 -
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /images/log_cleaner_anatomy.png HTTP/1.1" 304
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /images/log_compaction.png HTTP/1.1" 304 -
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /images/kafka_log.png HTTP/1.1" 304 -
::1 - - [23/Mar/2014:15:08:07 -0700] "GET /images/mirror-maker.png HTTP/1.1" 304 -
::1 - - [23/Mar/2014:15:09:55 -0700] "GET /documentation.html HTTP/1.1" 200 195264
```

EXAMPLE EVENT: DATABASES



What is a table?

key	value
key1	value2
key2	value3

EXAMPLE EVENT: DATABASES

key	value
key1	value1

Time = 0

put(key1, value1)

key	value
key1	value2
key2	value2

Time = 1

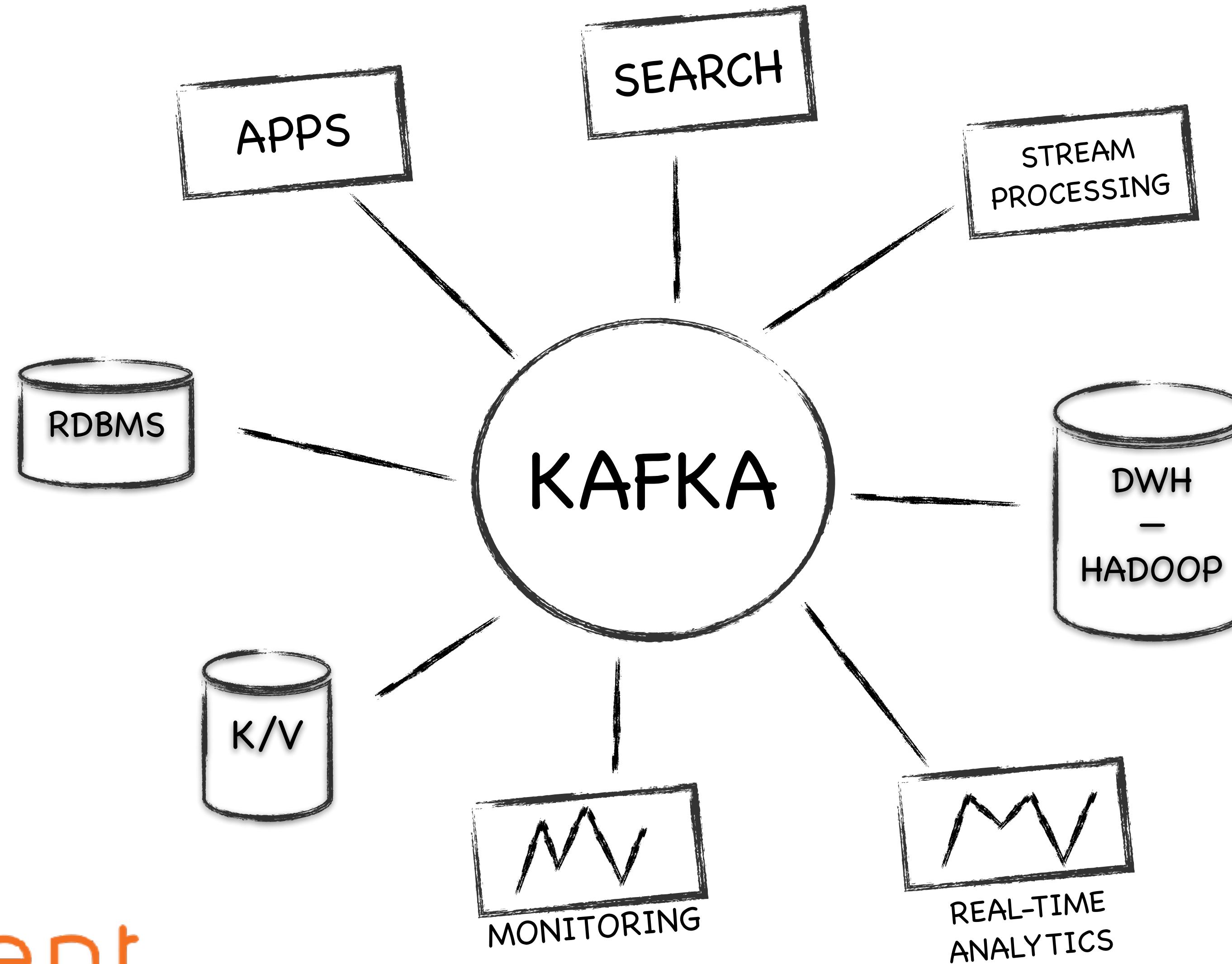
put(key2, value2)

key	value
key1	value3
key2	value2

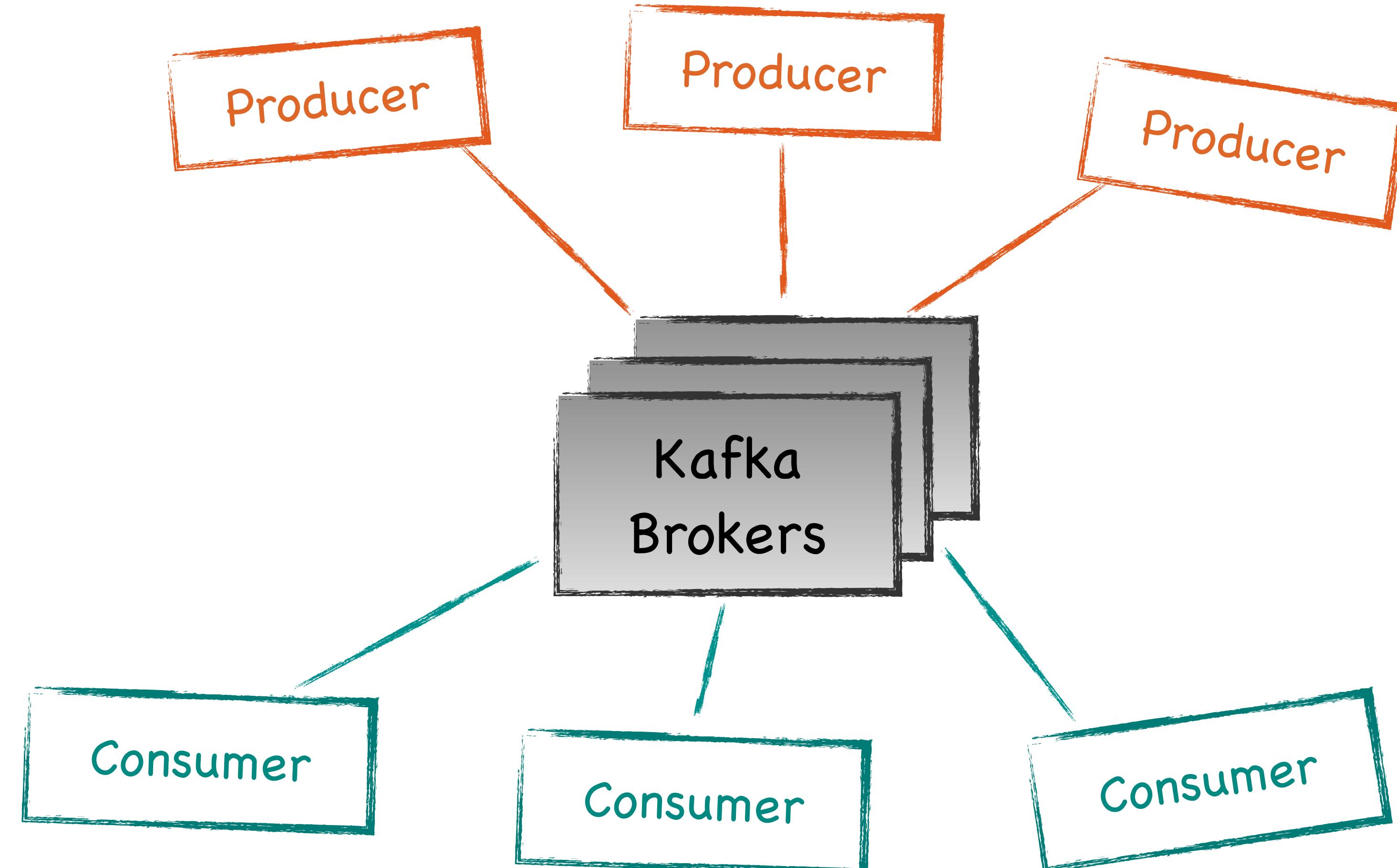
Time = 2

put(key1, value3)

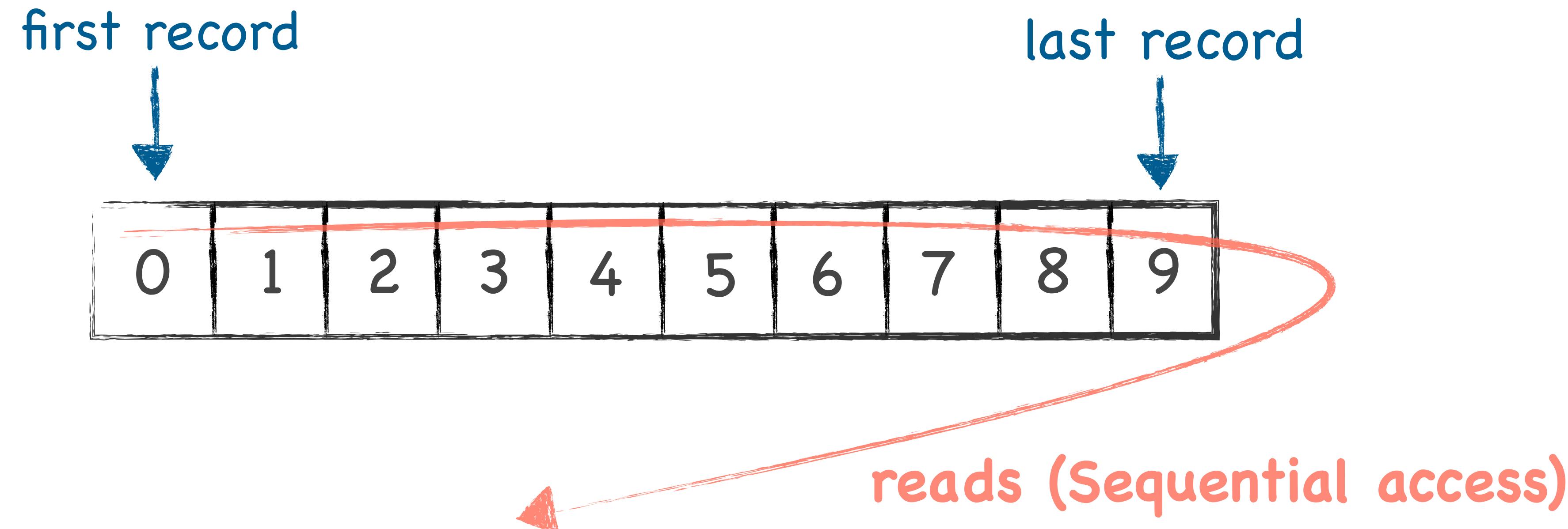
STREAMING PLATFORM



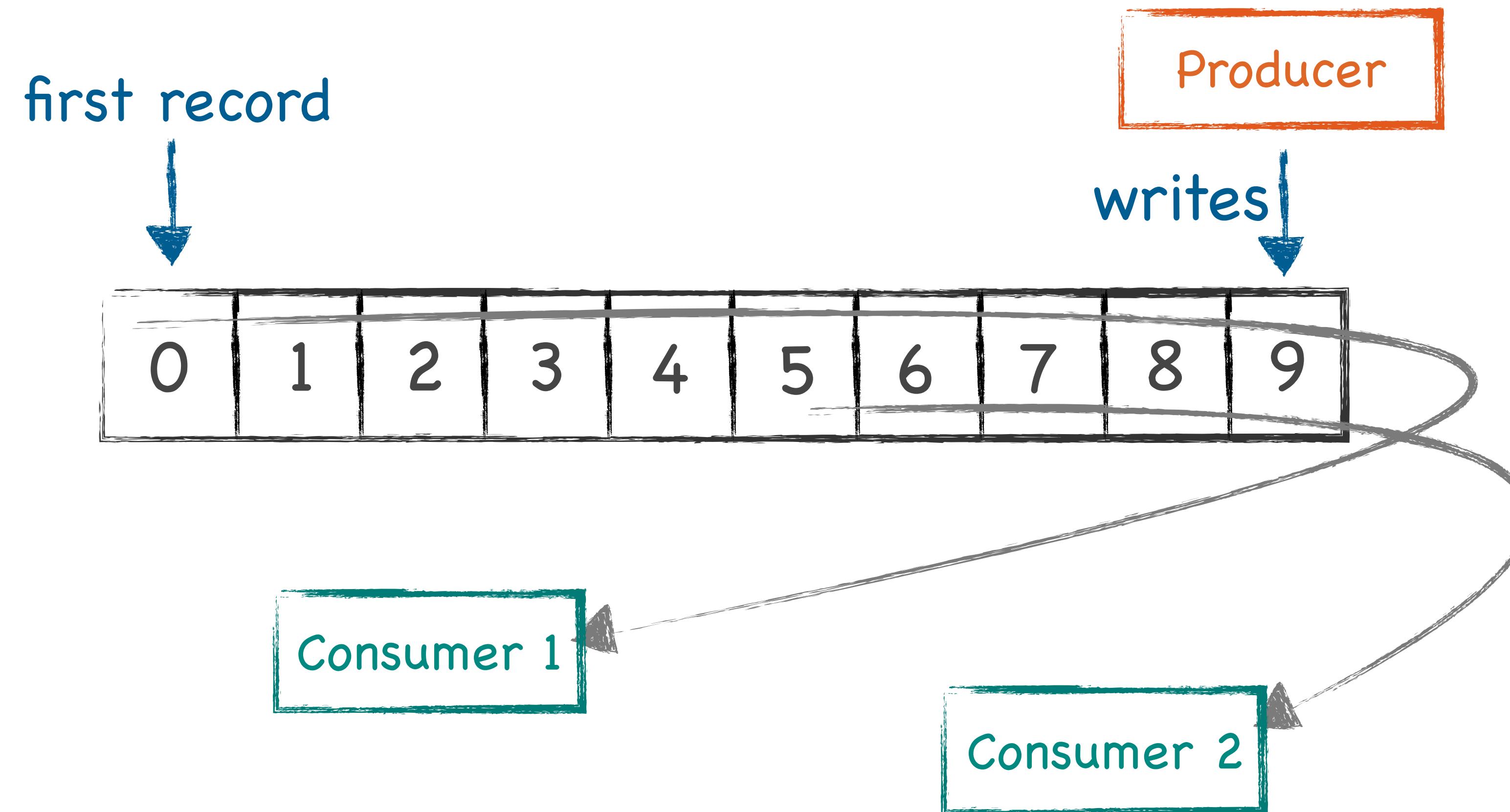
APACHE KAFKA:



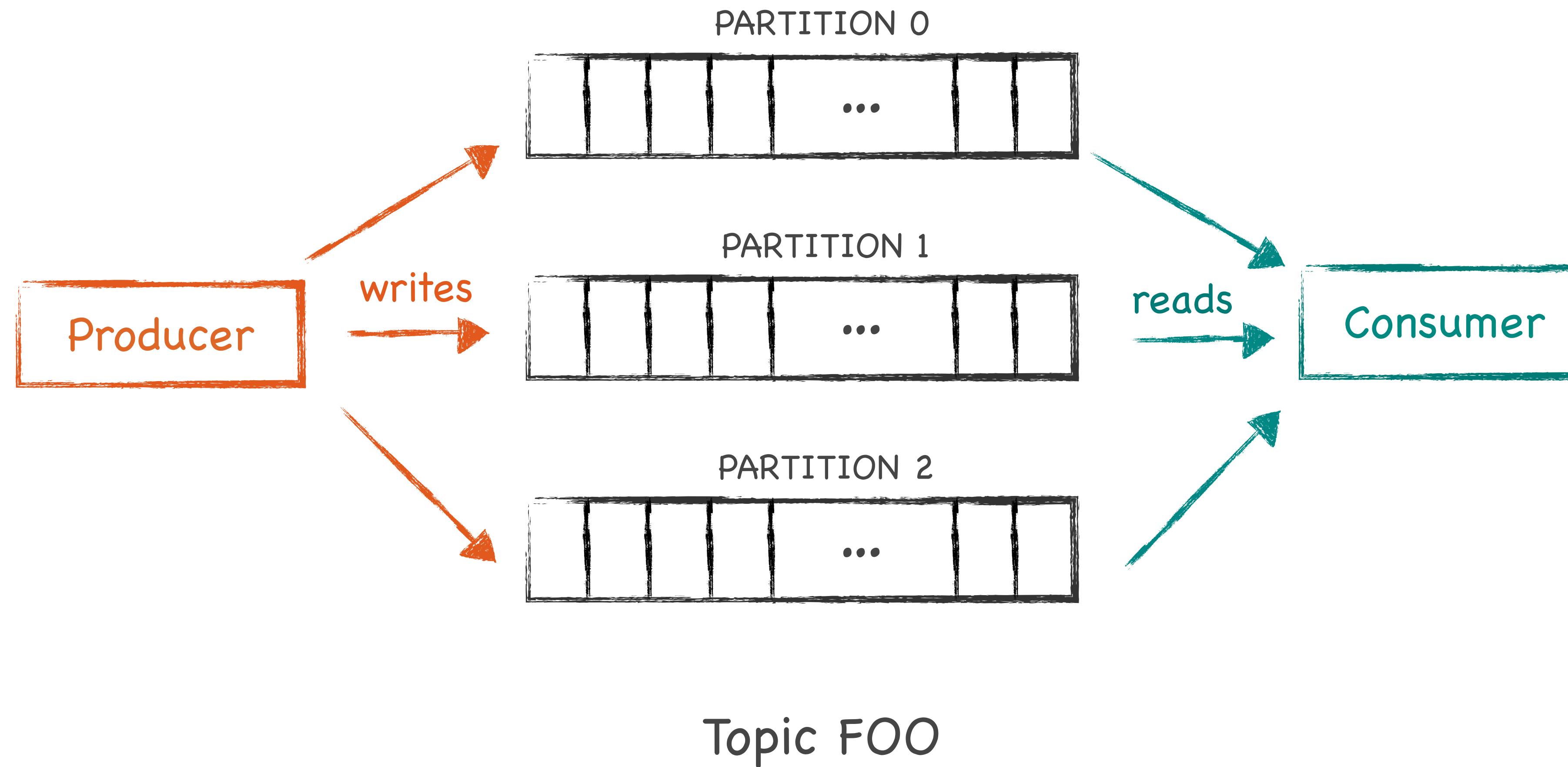
APACHE KAFKA: Key idea(Log)



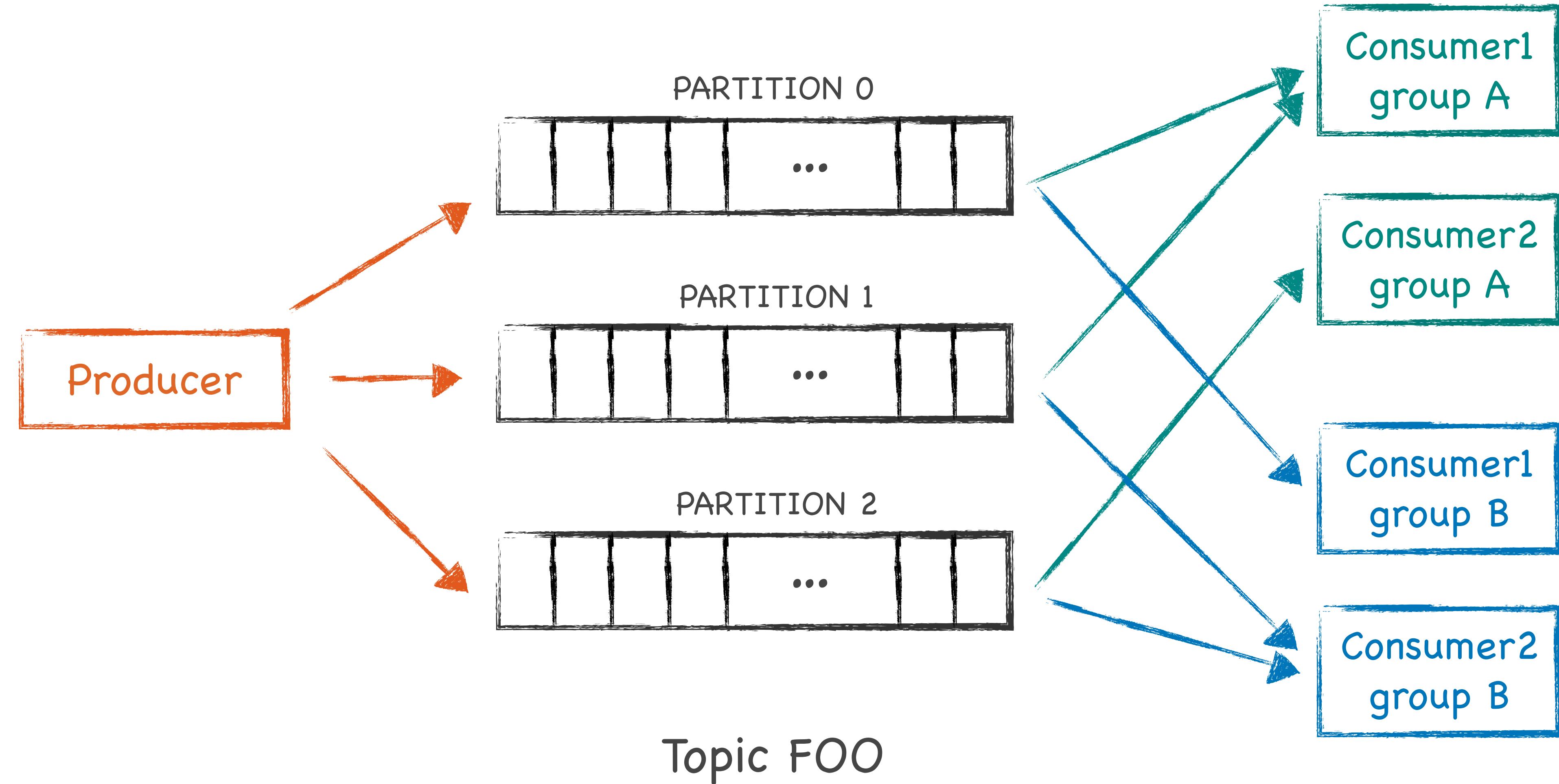
APACHE KAFKA: Logs & Pub-Sub



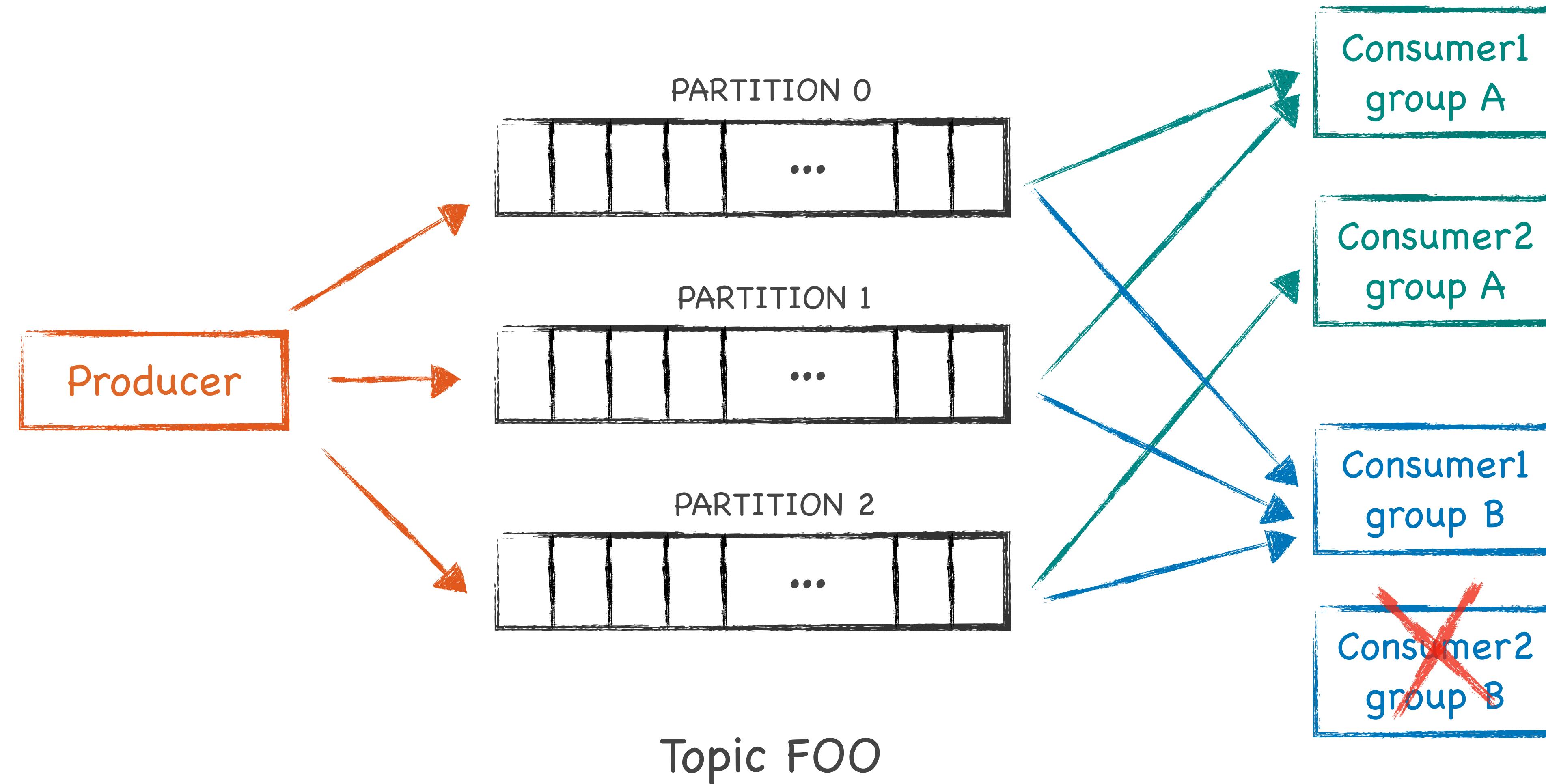
APACHE KAFKA: Topic = partitioned log



APACHE KAFKA: Scalable Consumption

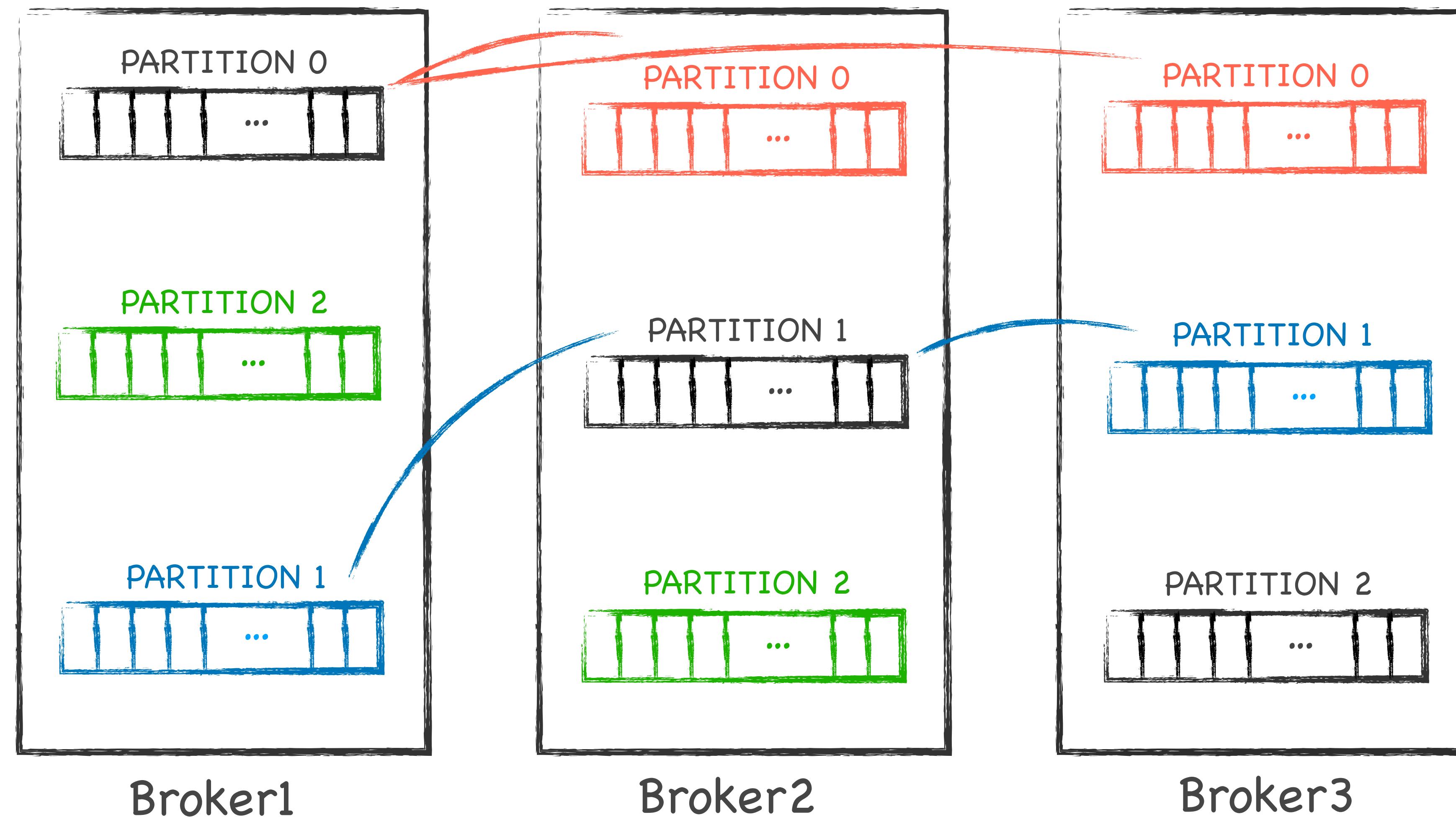


APACHE KAFKA: Fault Tolerant Consumption



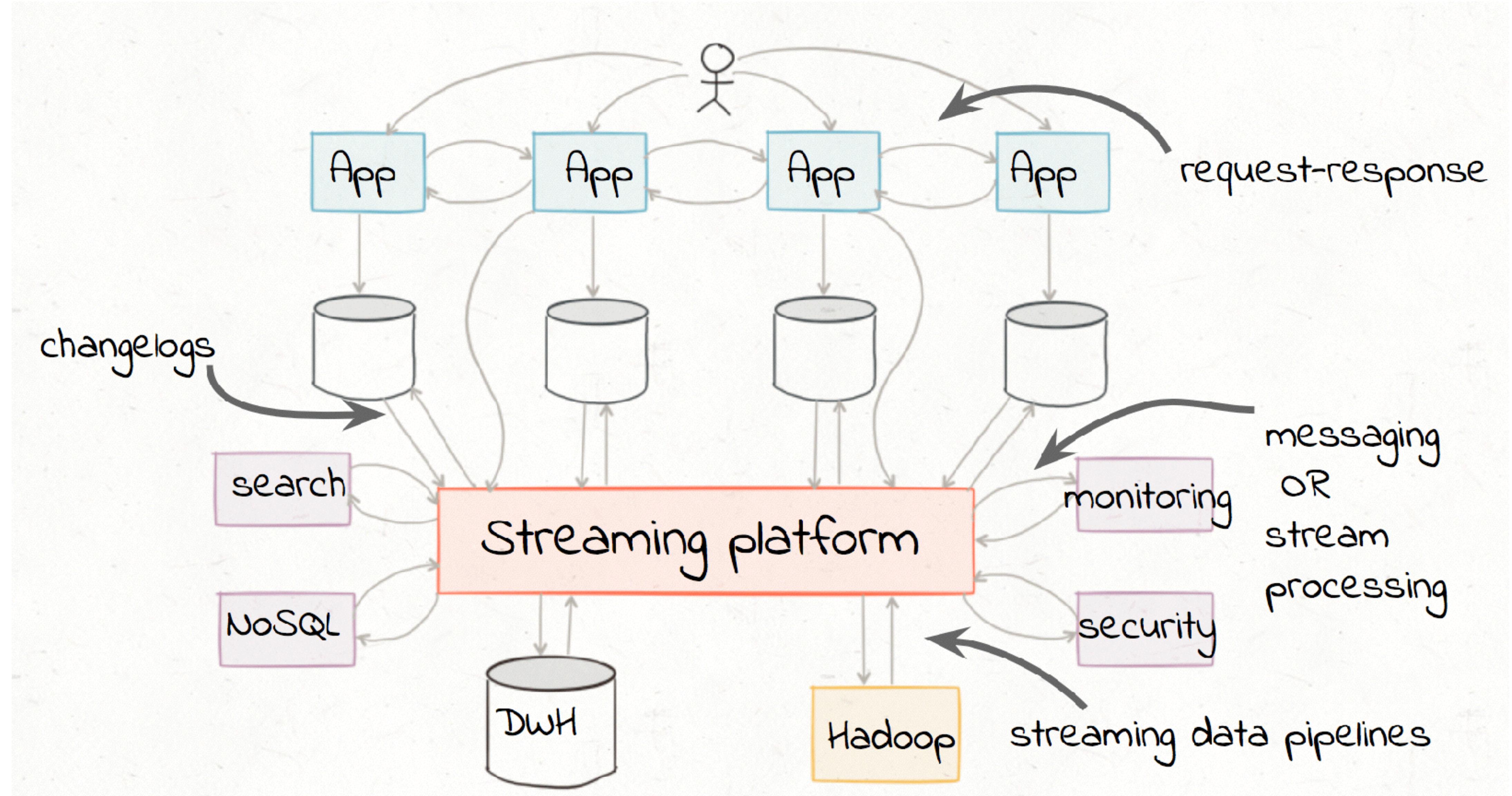
APACHE KAFKA: Partition & Replication

Topic FOO



KAFKA CONNECT

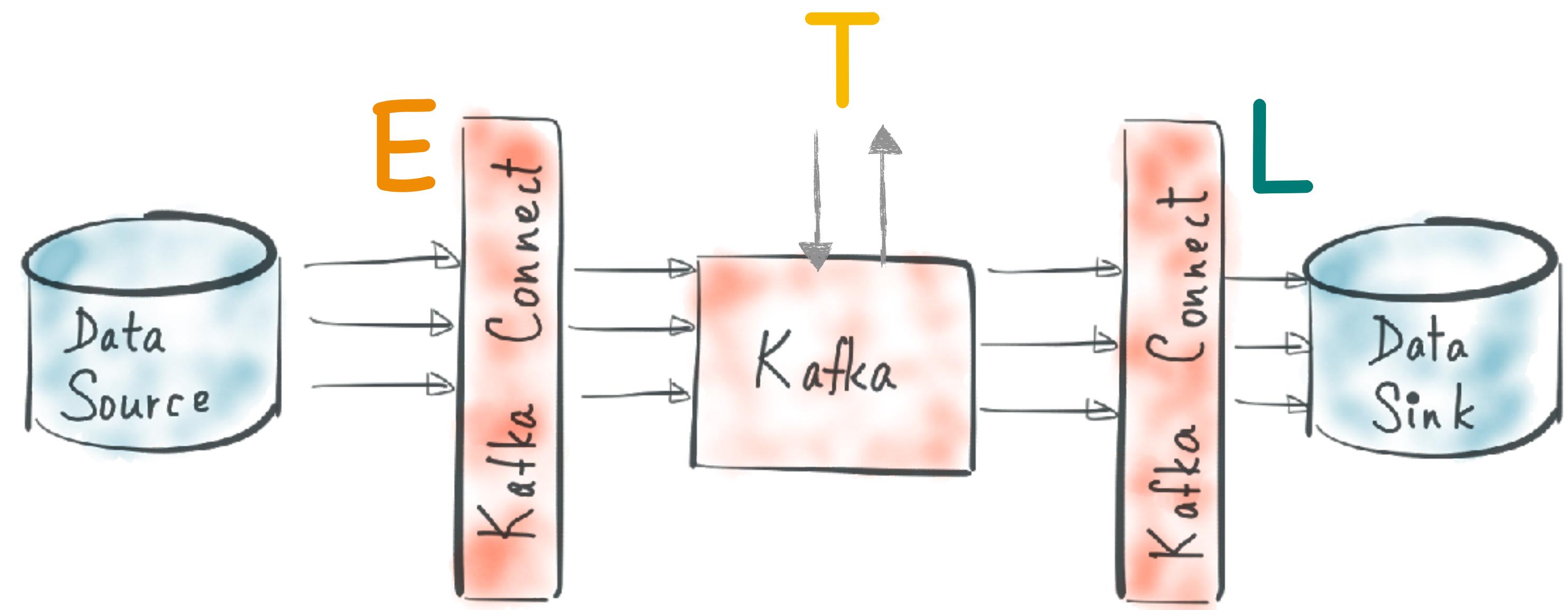
Large-scale streaming data import/export for Kafka



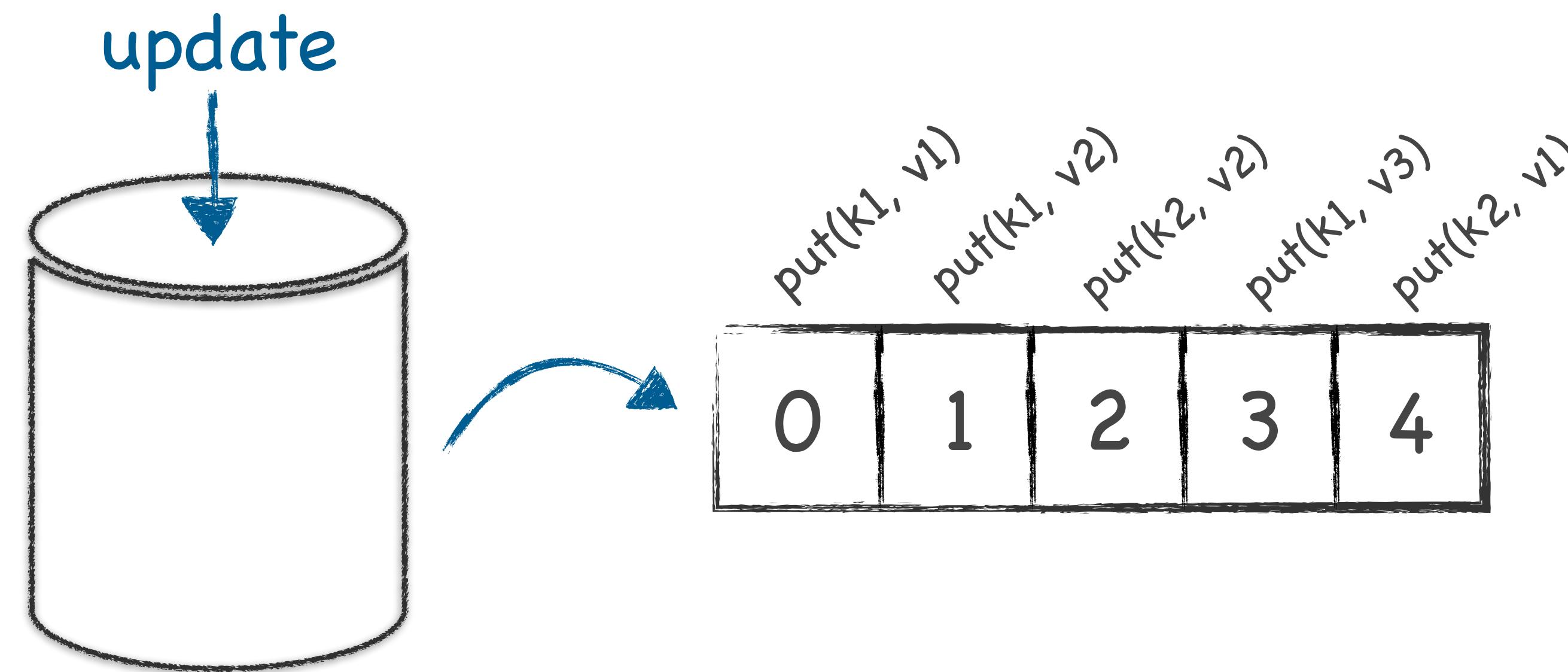
GOAL

1. Focus on copying
2. Batteries included
3. Parallelism
4. Scale

KAFKA CONNECT



Database change events/logs



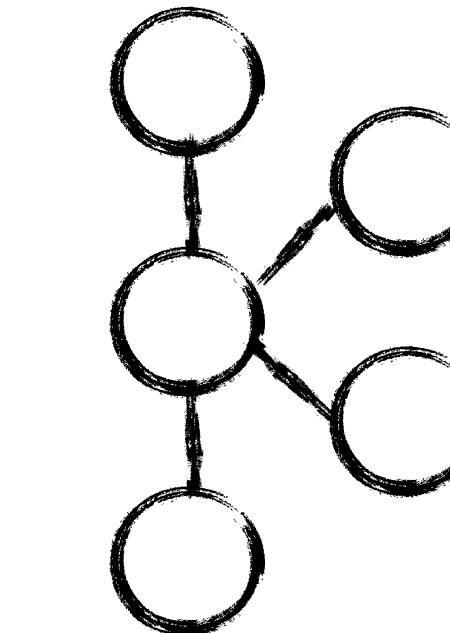
Table

TS	DATA
12:00	A
12:20	B
12:30	C
13:00	D
13:05	E

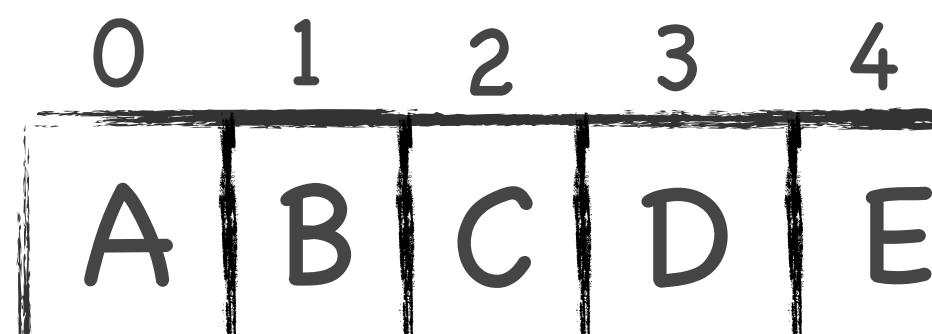
Database Source (connector)



offset = timestamp

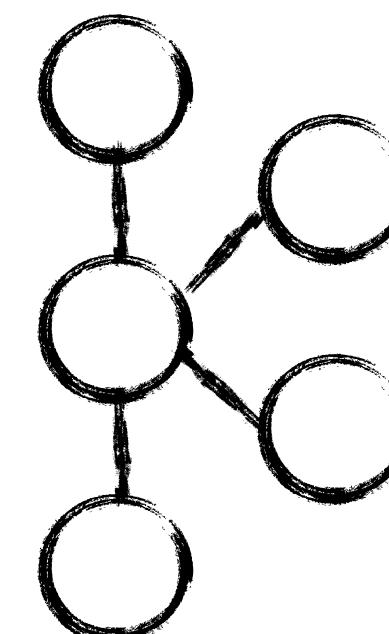
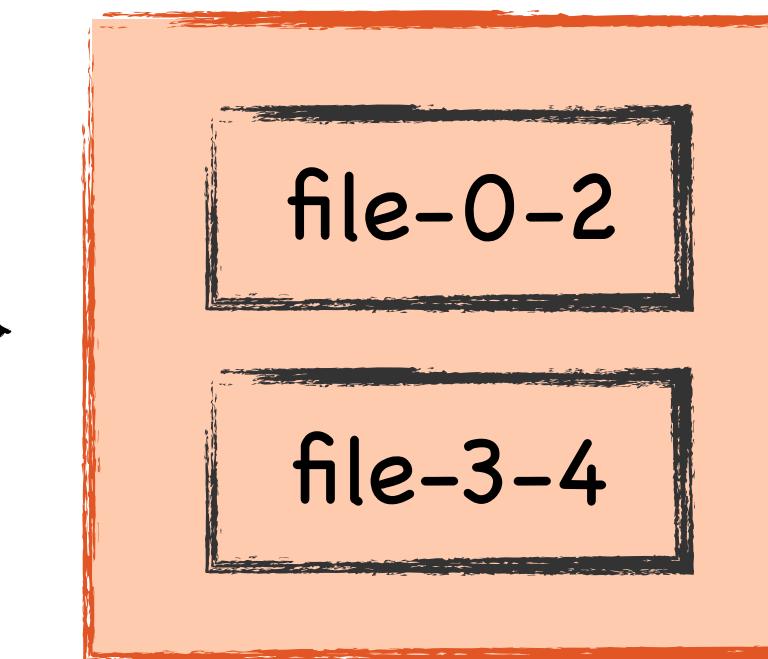


HDFS Sink (connector)

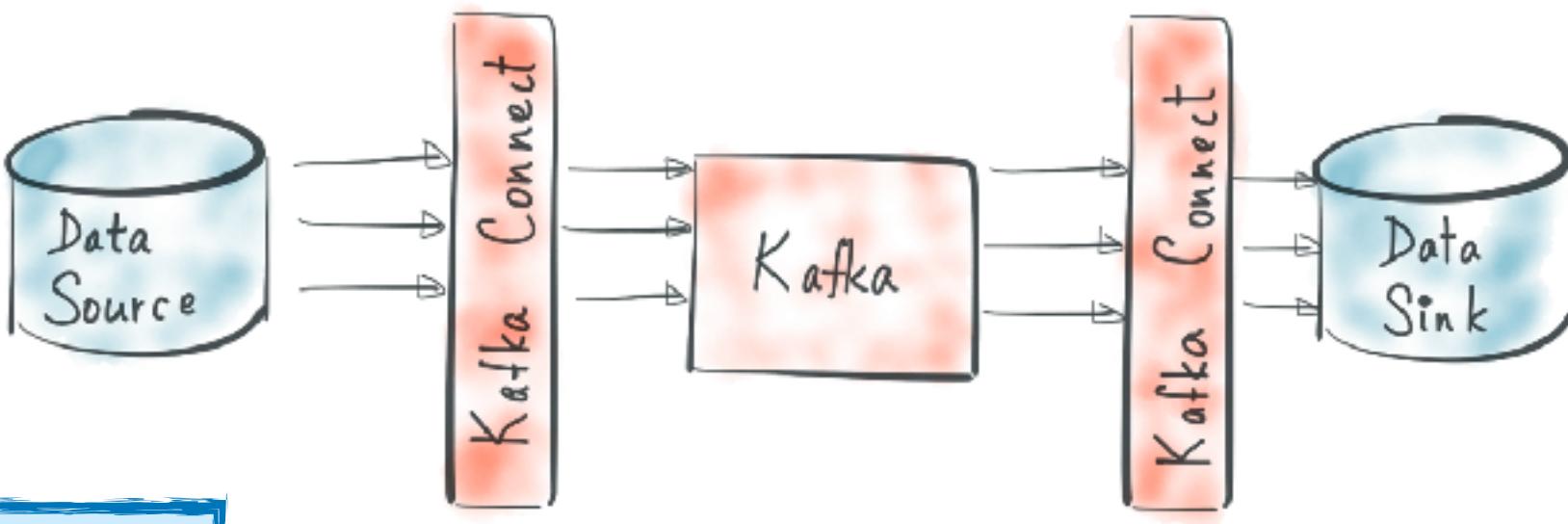


offset = kafka offset

HDFS Directory



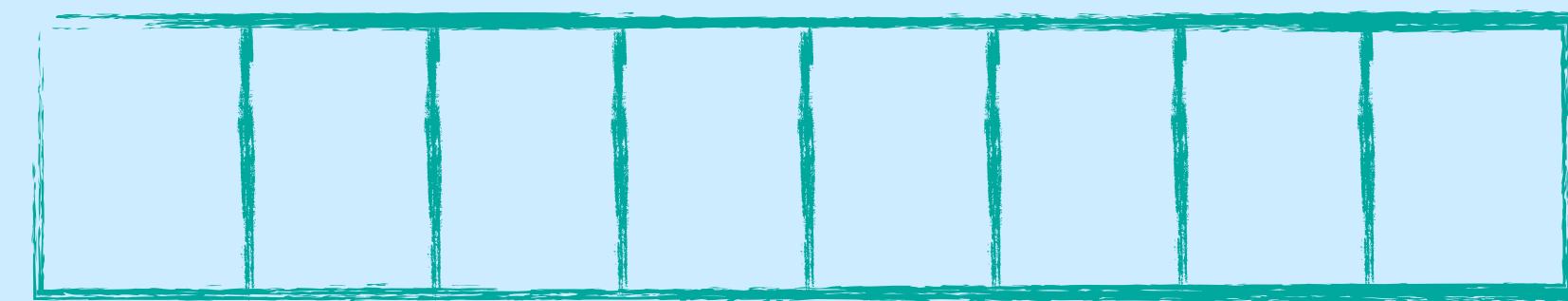
Data Model



Partitioned Stream - DataBase

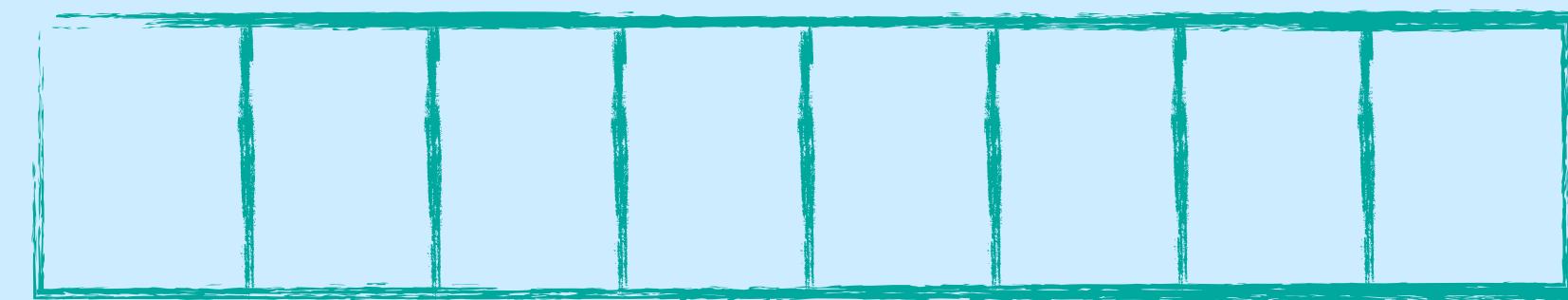
Partition 1

(Table 1)



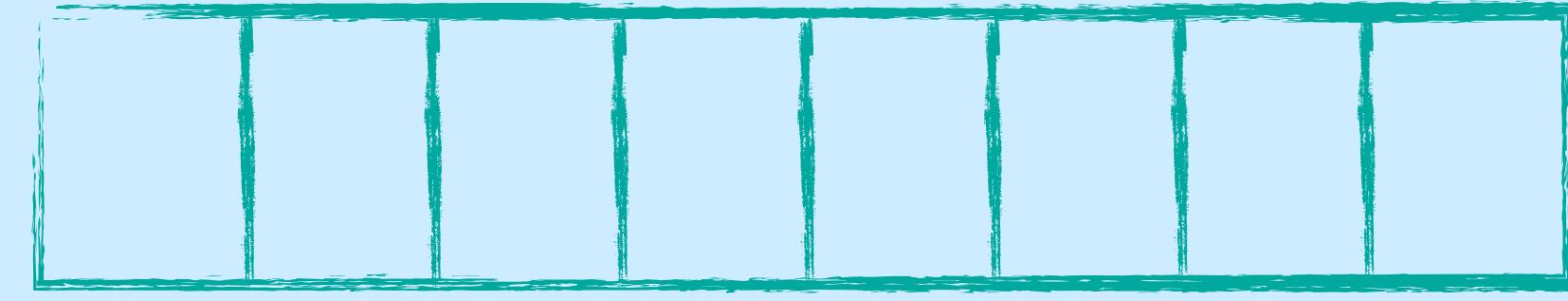
Partition 2

(Table 2)



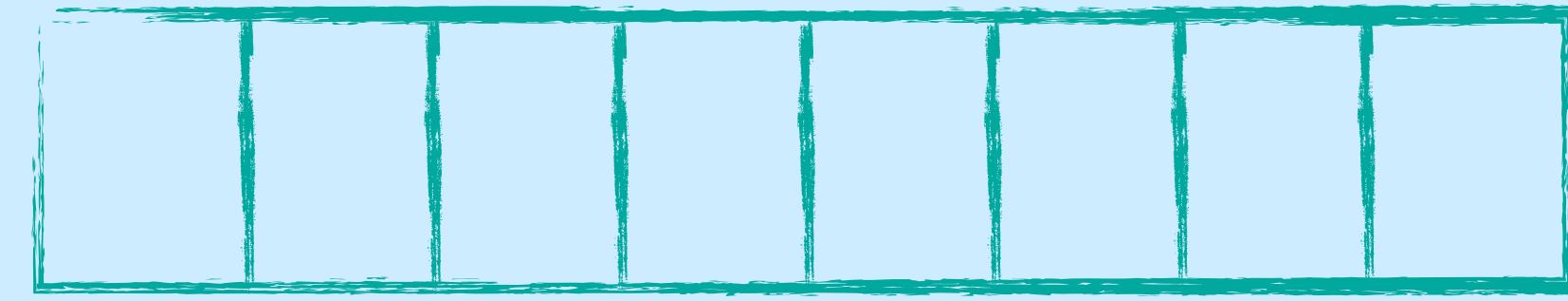
Partition 3

(Table 3)



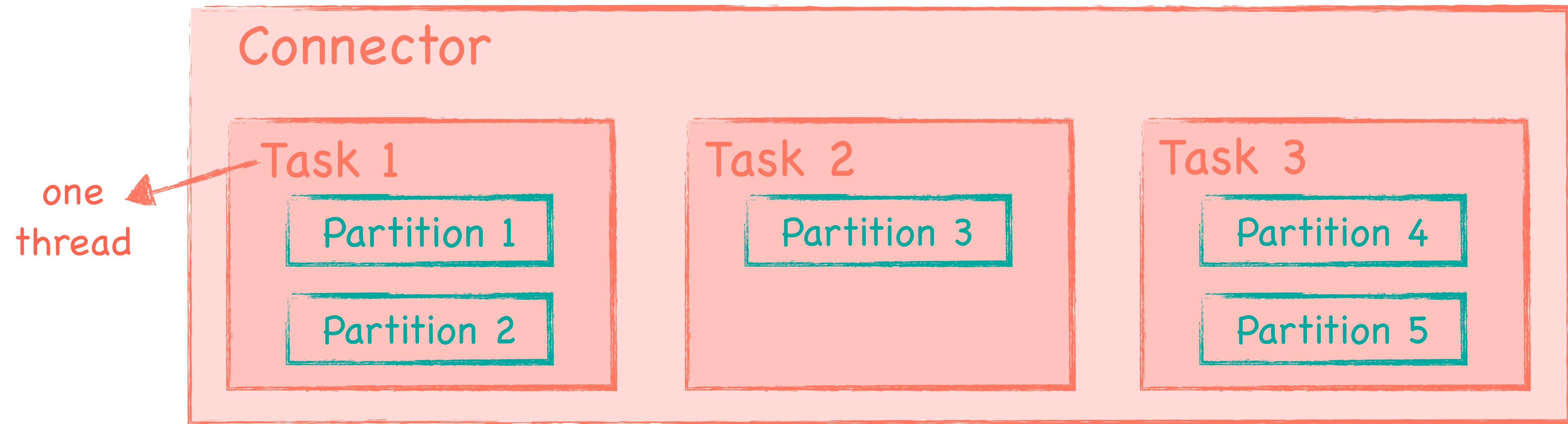
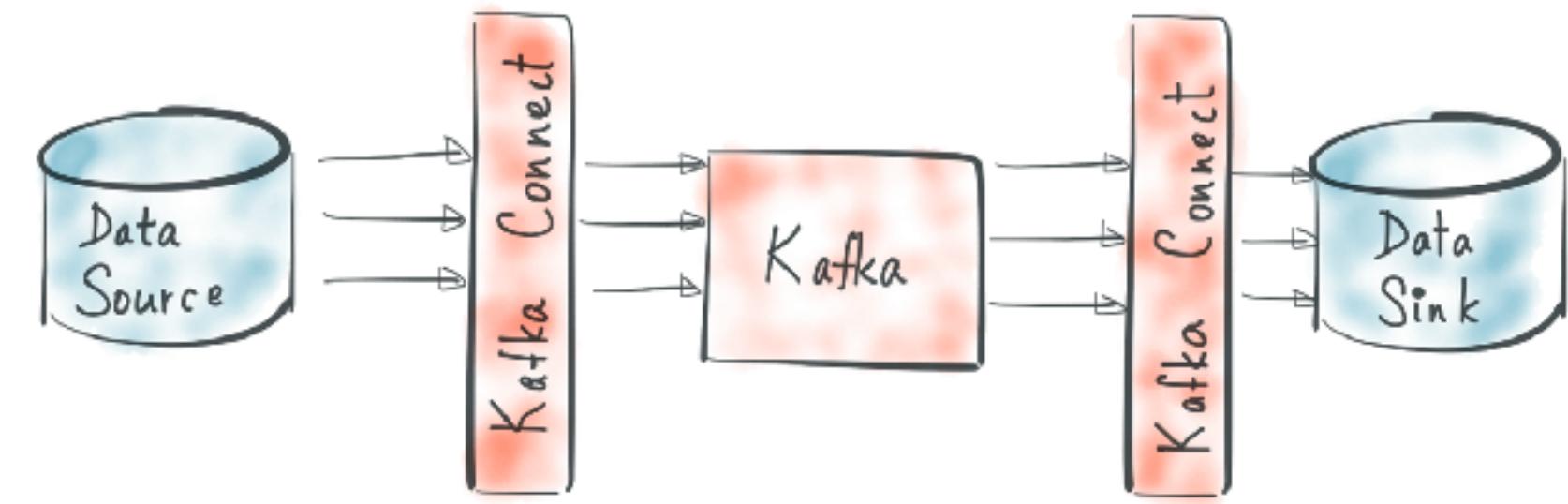
Partition 4

(Table 4)

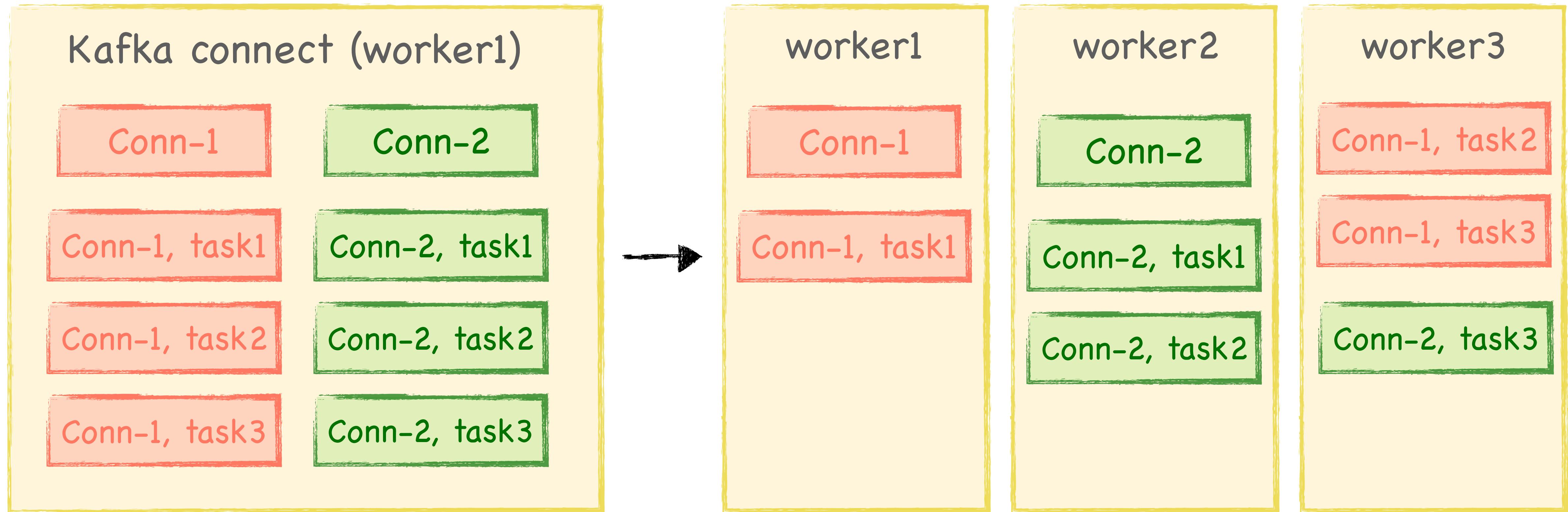
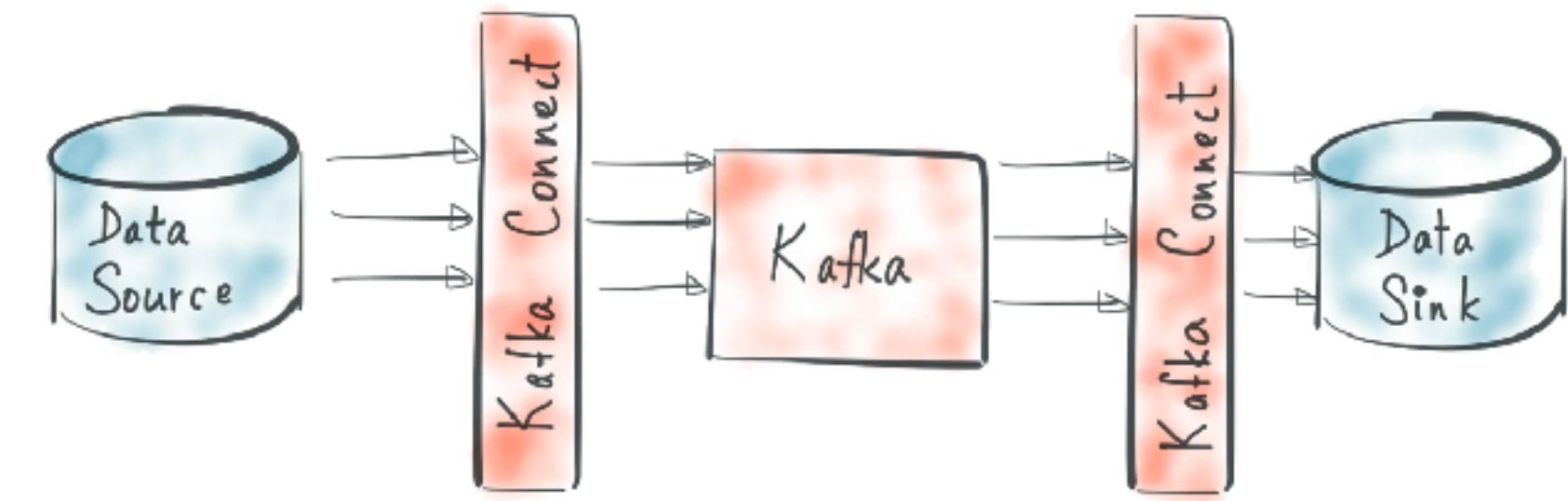


id=1 id=2 id=3 id=4 id=5 id=6 id=7 id=8

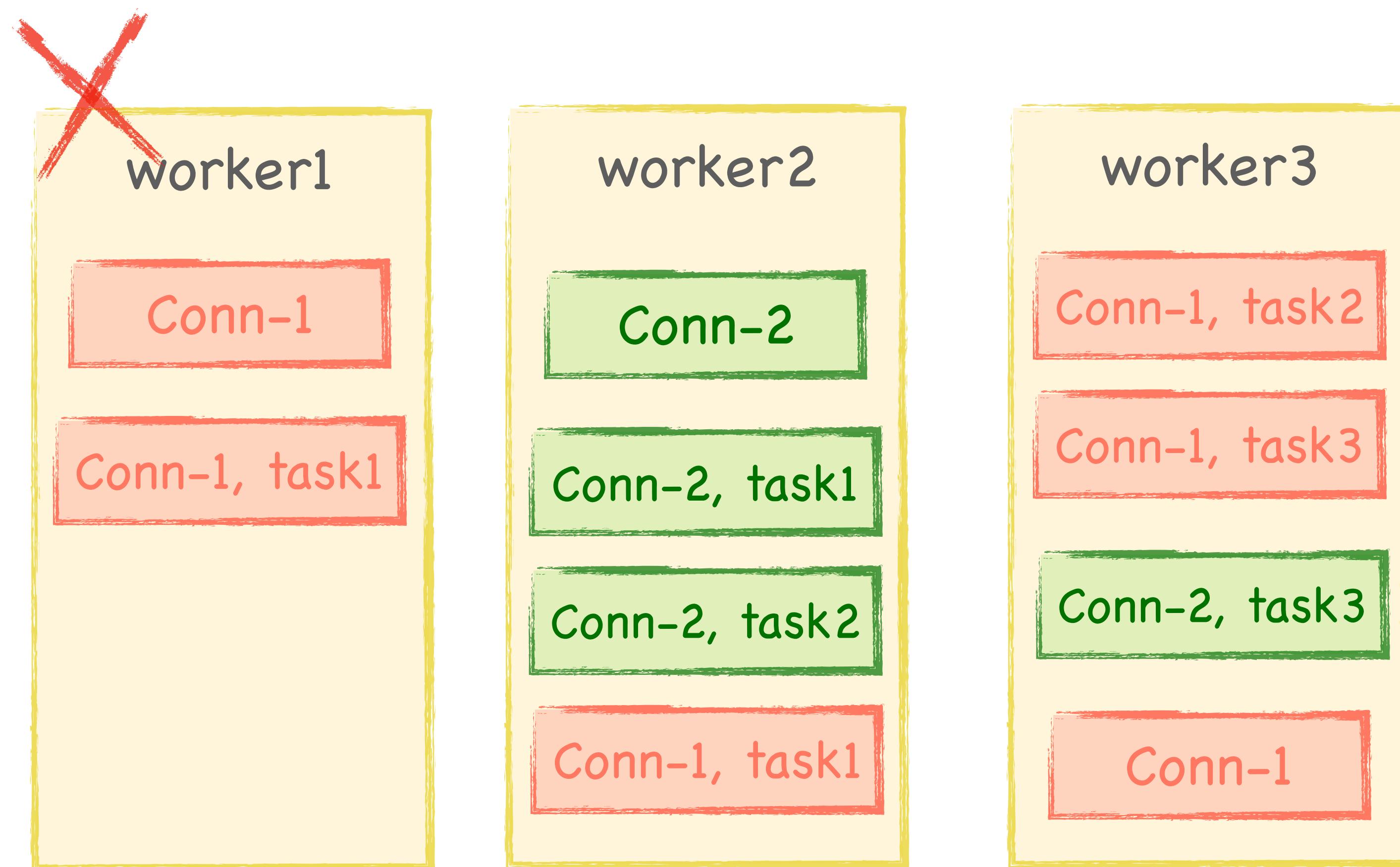
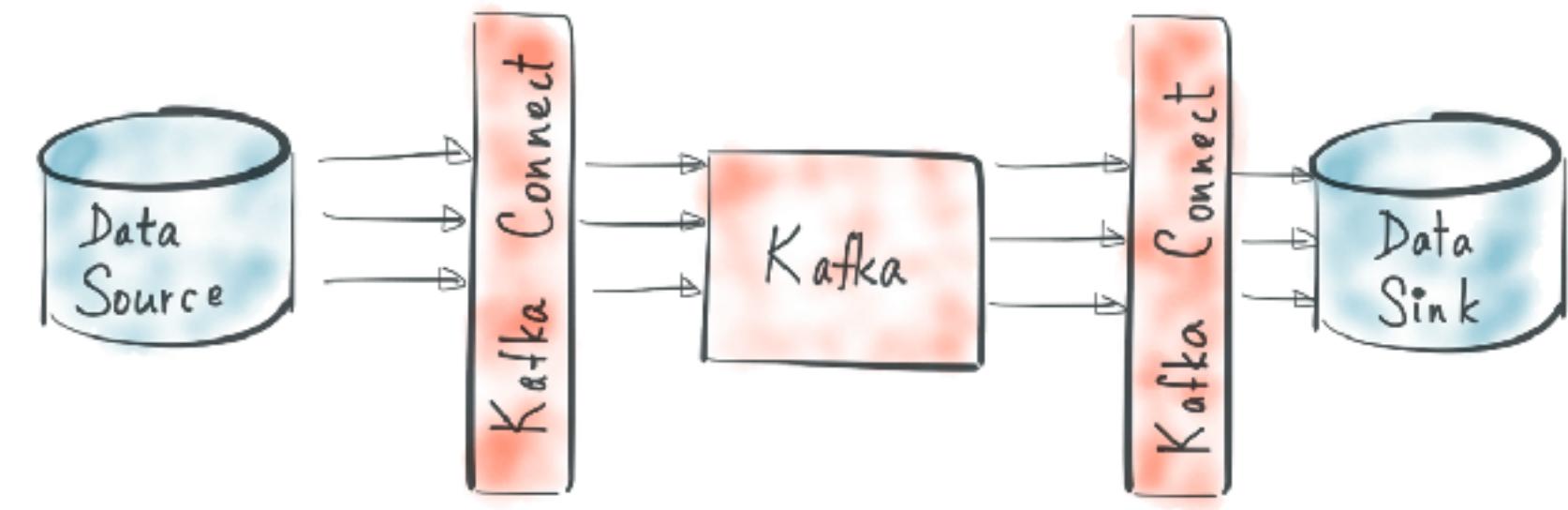
Parallelism Model



standalone & Distributed

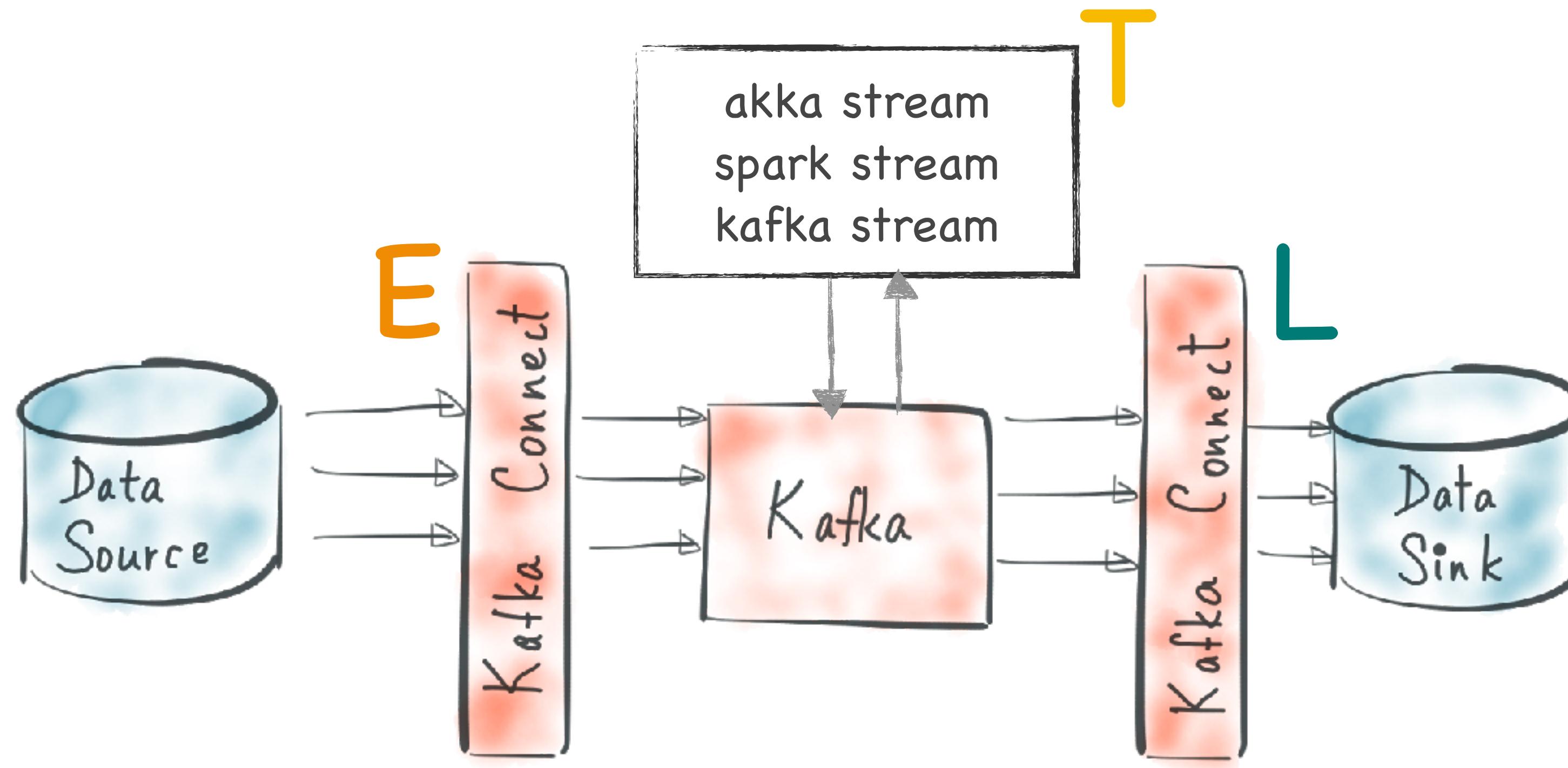


standalone & Distributed





KAFKA CONNECT



Kafka Connect Quick Start

- Confluent Platform



Confluent Platform Install

- Start all available services:

```
confluent start
```

Your output should resemble:

```
Starting zookeeper
zookeeper is [UP]
Starting kafka
kafka is [UP]
Starting schema-registry
schema-registry is [UP]
Starting kafka-rest
kafka-rest is [UP]
Starting connect
connect is [UP]
```

Confluent Connectors

[Product](#)[Solutions](#)[Developers](#)[Blog](#)[Docs](#)[Download](#)

Additional Connectors Available

Other notable Connectors that have been developed utilizing the Kafka Connect framework.

CONNECTOR	TAGS	DEVELOPER/SUPPORT	DOWNLOAD
Amazon Kinesis (Sink)	messaging	Community	Community
Apache Ignite (Source)	File System	Community	Community
Apache Ignite (Sink)	File System	Community	Community
ArangoDB (Sink)	NoSQL	Community	Community
AWS Lambda (Sink)	Cloud, AWS, Lambda	Community	Community
Azure DocumentDB (Sink)	DocumentDB, Azure, NoSQL	Community	Community
Blockchain (Source)	Bitcoin, Blockchain	Community	Community

Want to build a connector?

Interested Open Source Developers and Vendors can get started with the following resources:

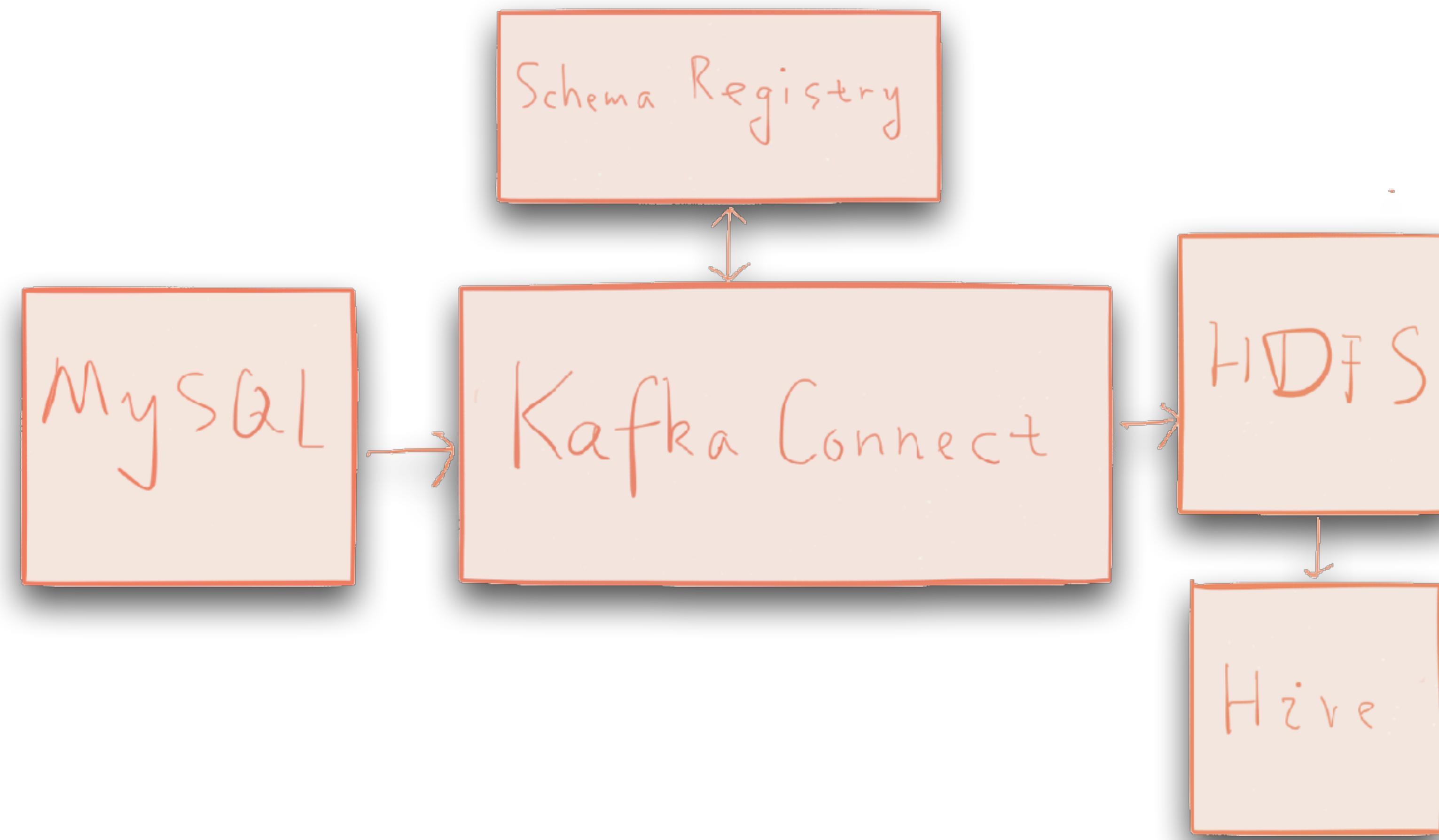
- > [Kafka Connect Overview](#)
- > [Kafka Connect Developers Guide](#)
- > [Partner Development Guide for Kafka Connect](#)

Already built a connector?

If you have a connector that you'd like to see added to our list, let us know at:
confluent-platform@googlegroups.com

If your connector meets our criteria, we'll

To demonstrate Kafka Connect,
A simple data pipeline: MySQL → Kafka → HDFS → Hive



Simple Application: transform a batch pipeline into a real-time one



類別

搜尋課程



成為講師

我的課程



電子商務

內容行銷

通訊工程與軟體

辦公室提升效率

個人成長

設計

市場行銷

生活品味

攝影

健康和保健

教師培訓

將該課程作為禮物贈送

Apache Kafka Series - Learn Apache Kafka for Beginners

Tutorial: Learn the Apache Kafka Ecosystem, Core Concepts, Operations, Kafka API, Build Your Own Producers and Consumers

暢銷課程

★★★★★ 4.6 (2072 個評等) 9323 名學生註冊

建立者：Stephane Maarek 上次更新 2/2018 語言 英語



預覽此課程

\$300 \$1,800 折扣 83%

立即購買

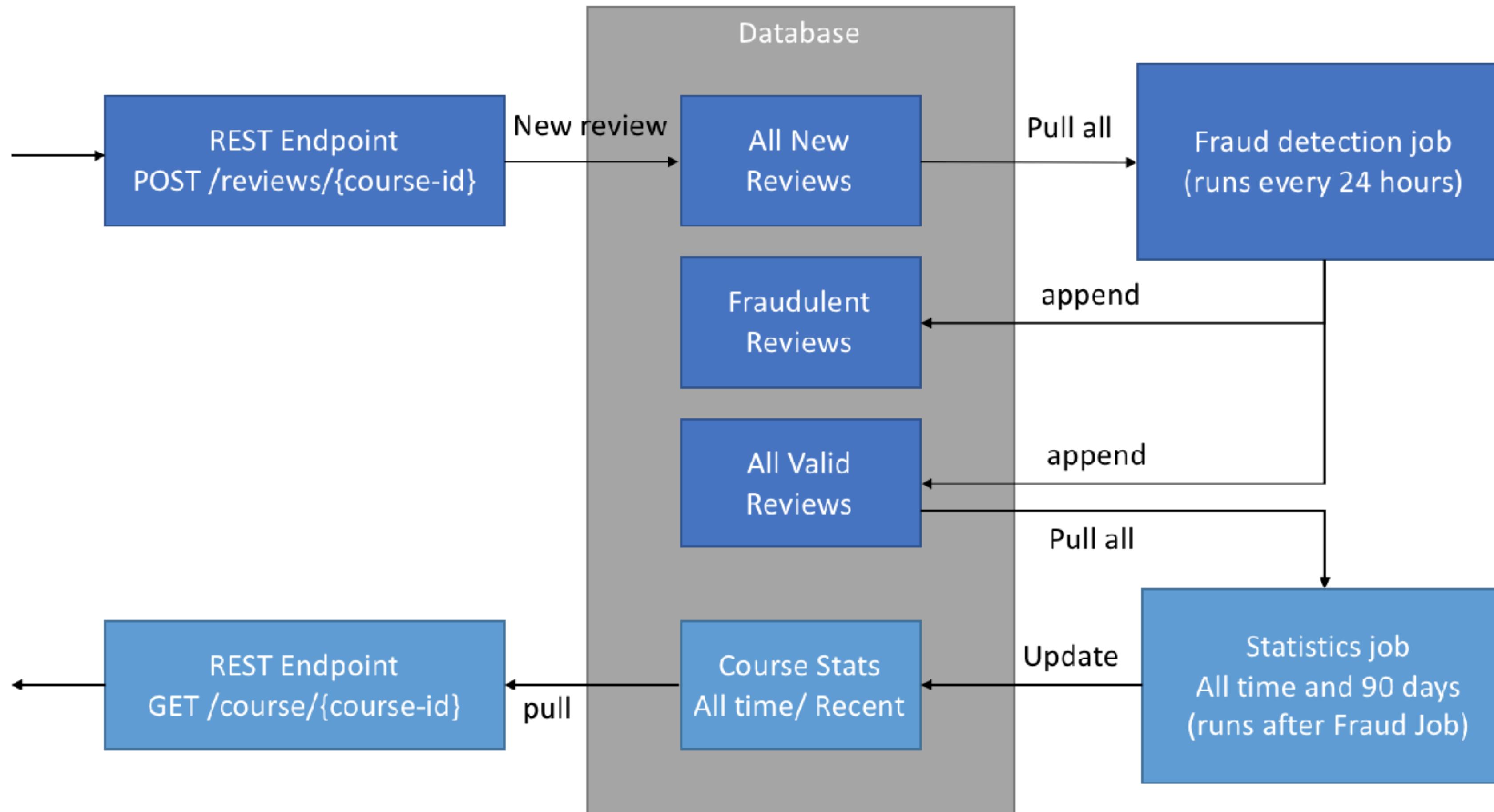
我會學些什麼呢？

✓ Learn about Apache Kafka Ecosystem, Architecture, Core Concepts and Operations

✓ Understand Fundamental Concepts behind Apache Kafka Like Topics, Partitions, Brokers, Producers.

新增至購物車

Simple Application: transform a batch pipeline into a real-time one



Simple Application: transform a batch pipeline into a real-time one

