# Distributed Logs
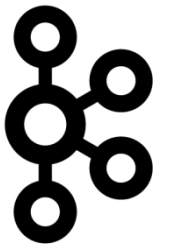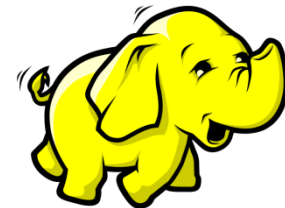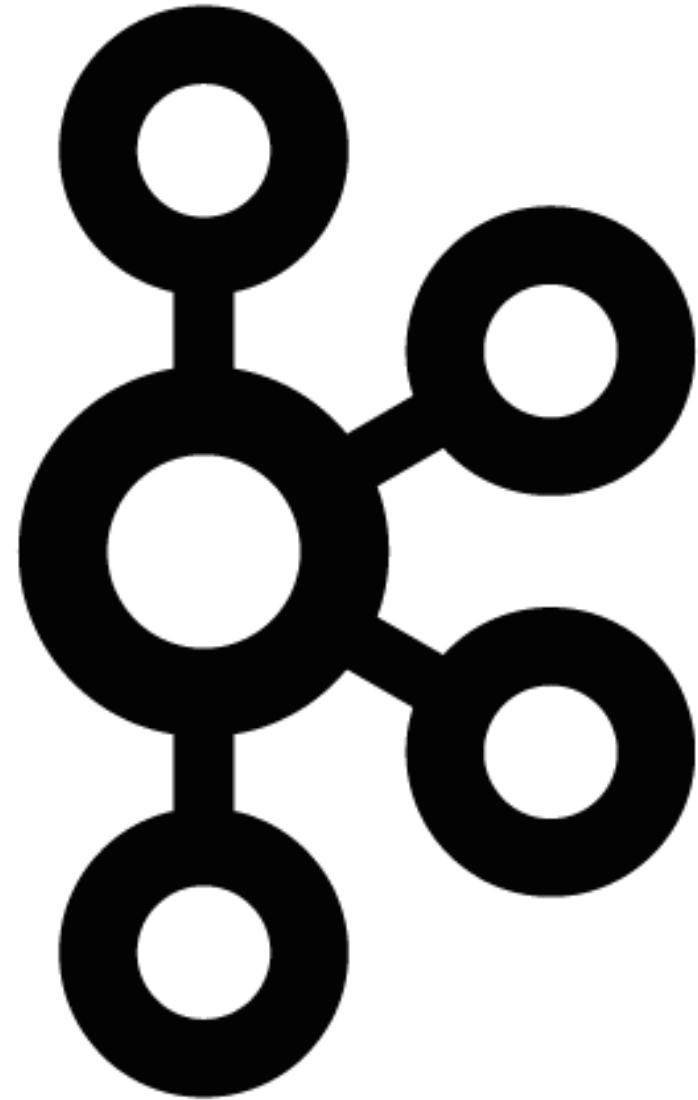
Michael Enudi

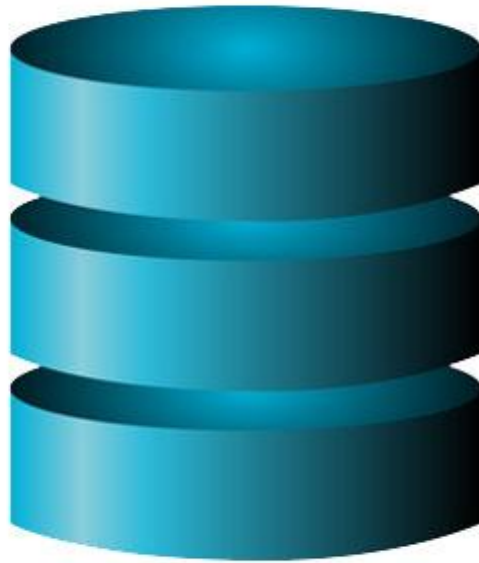*Journey through the world of databases and data engineering*

# Scope

- Data Architectures

- Introduction to Distributed Logs

- Use case: Apache Kafka

- Apache Zookeeper

- Stream processing

- Stream Processing II (Kafka Streams)

- Stream analytics

- Lab: NYC Taxi and Fares
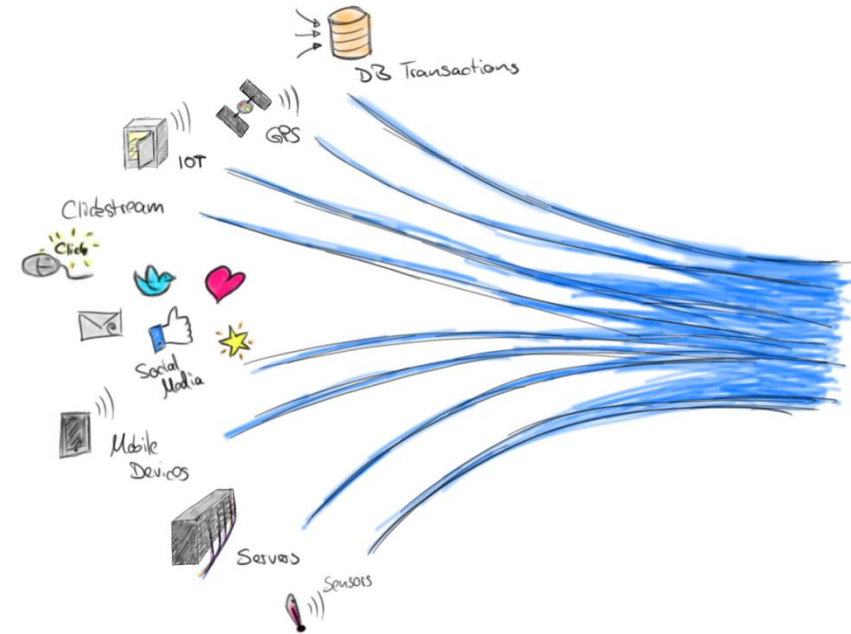
- Streaming: To be or not to be

**In software or data engineering, we encounter data in two general states**

**Data-at-rest**

**Batch**

**Data-in-motion**

**Stream/Streaming**

"Data-in-Use" / "Serving"

Data architecture (Lambda)



The Lambda Architecture

Speed Layer

Recent Data

Real-time View

Real-time View

Query

New Data

New Data

New Data

Query

Batch Layer

All Data

Serving Layer

Batch View

Batch View

Data architecture (Kappa)

Batch Layer

Batch Engine

Serving Layer

Serving Backend

Queries

Streaming Layer
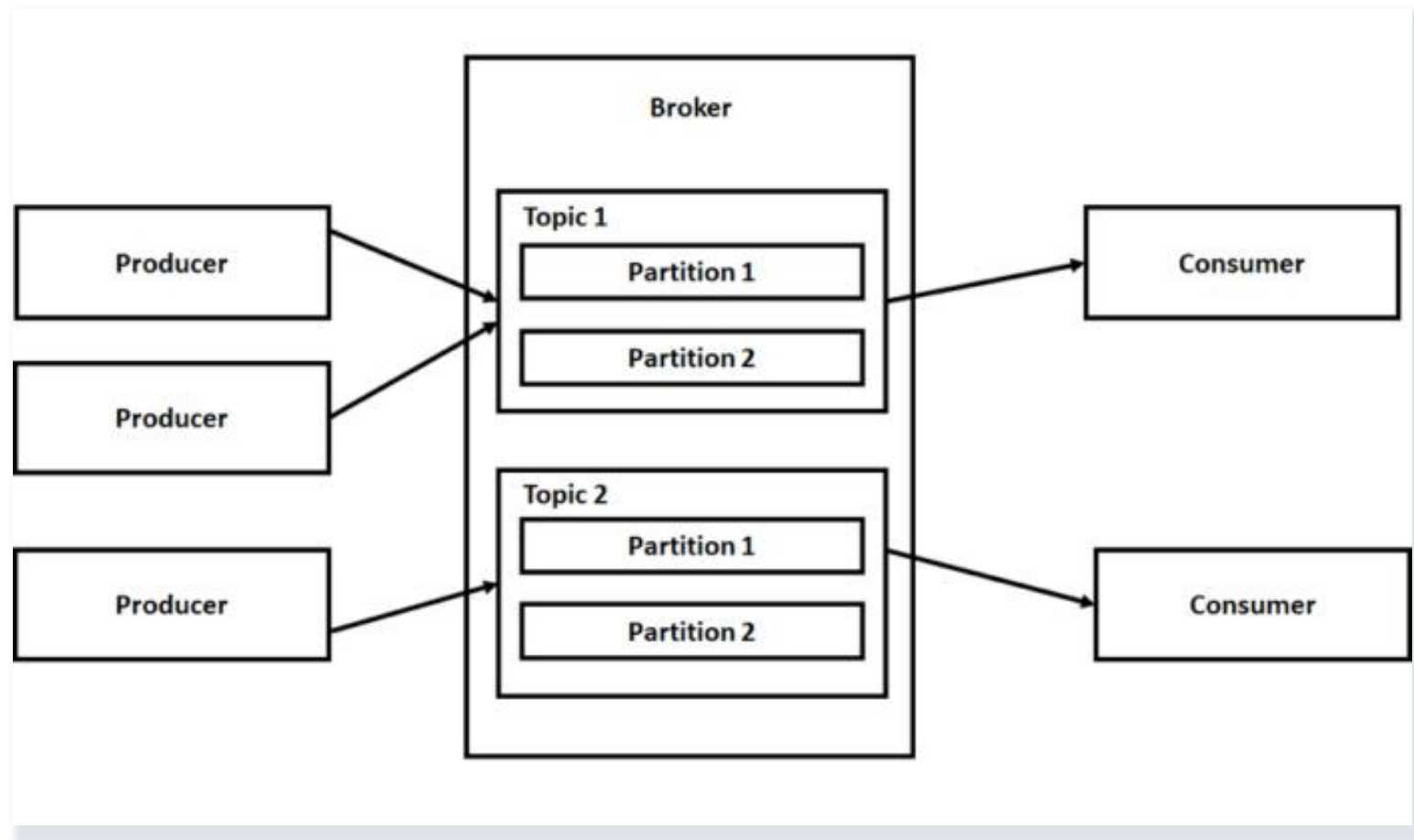
Real-time Engine

Data storage

Raw data

Results

**DISTRIBUTED COMMIT LOG AS A STORAGE PLATFORM**

❑ Used as a Message Broker or as a building real-time data pipelines for streaming applications

❑ Kafka started from within LinkedIn in January 2011

❑ Current version as at June 2019 is 2.3.0

❑ Written in Scala and Java

❑ Uses Zookeeper (just like HBase) for distributed coordination and metadata repository.

❑ Confluent currently offers Kafka as part of the Confluent Platform.

❑ Use cases include message broker, streaming application, event sourcing, commit logging, activity tracking and more.

❑ High profile users of kafka include LinkedIn, Apple, eBay, Netflix, Yelp, Oracle, Uber, Citibank, Spotify and a host of others
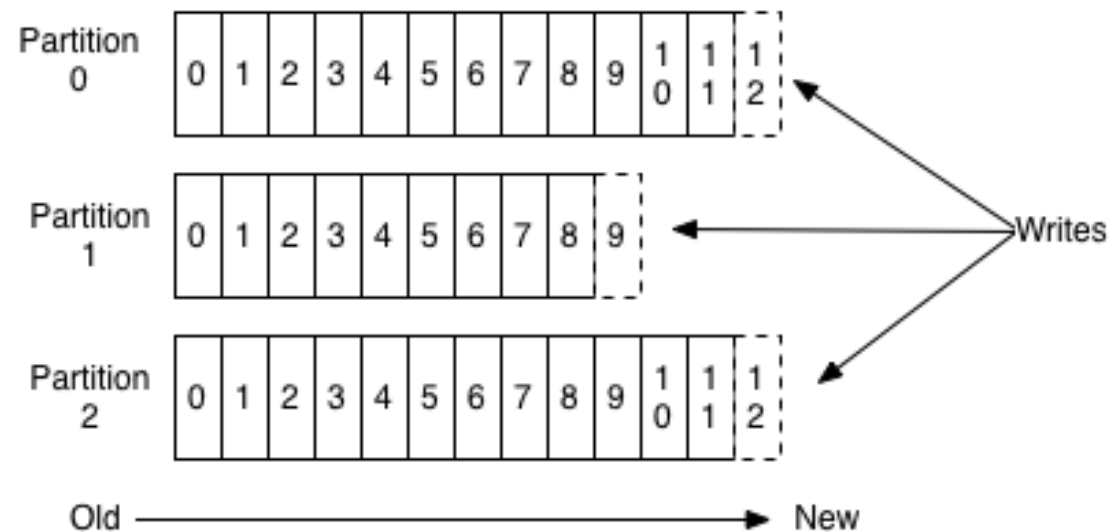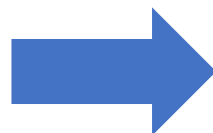
APACHE kafka™

- Topic

- Records / Messages

- Producers

- Consumers

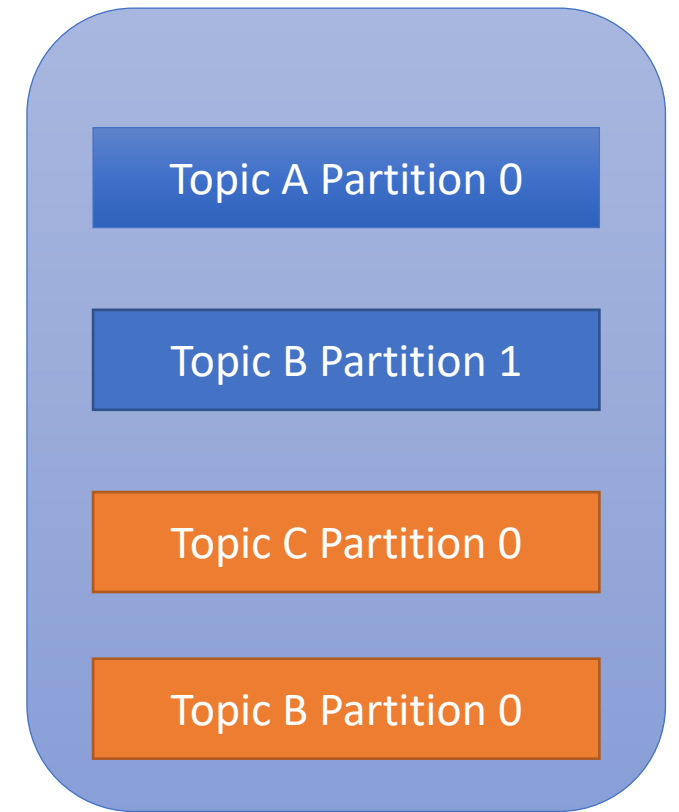# Kafka Topic

# Topic Partition Replication

Topic A - 1 Partition & Replication factor of 3
Topic B - 2 Partition & Replication factor of 3
Topic C - 1 Partition & Replication factor of 2
Topic D - 1 Partition & Replication factor of 1
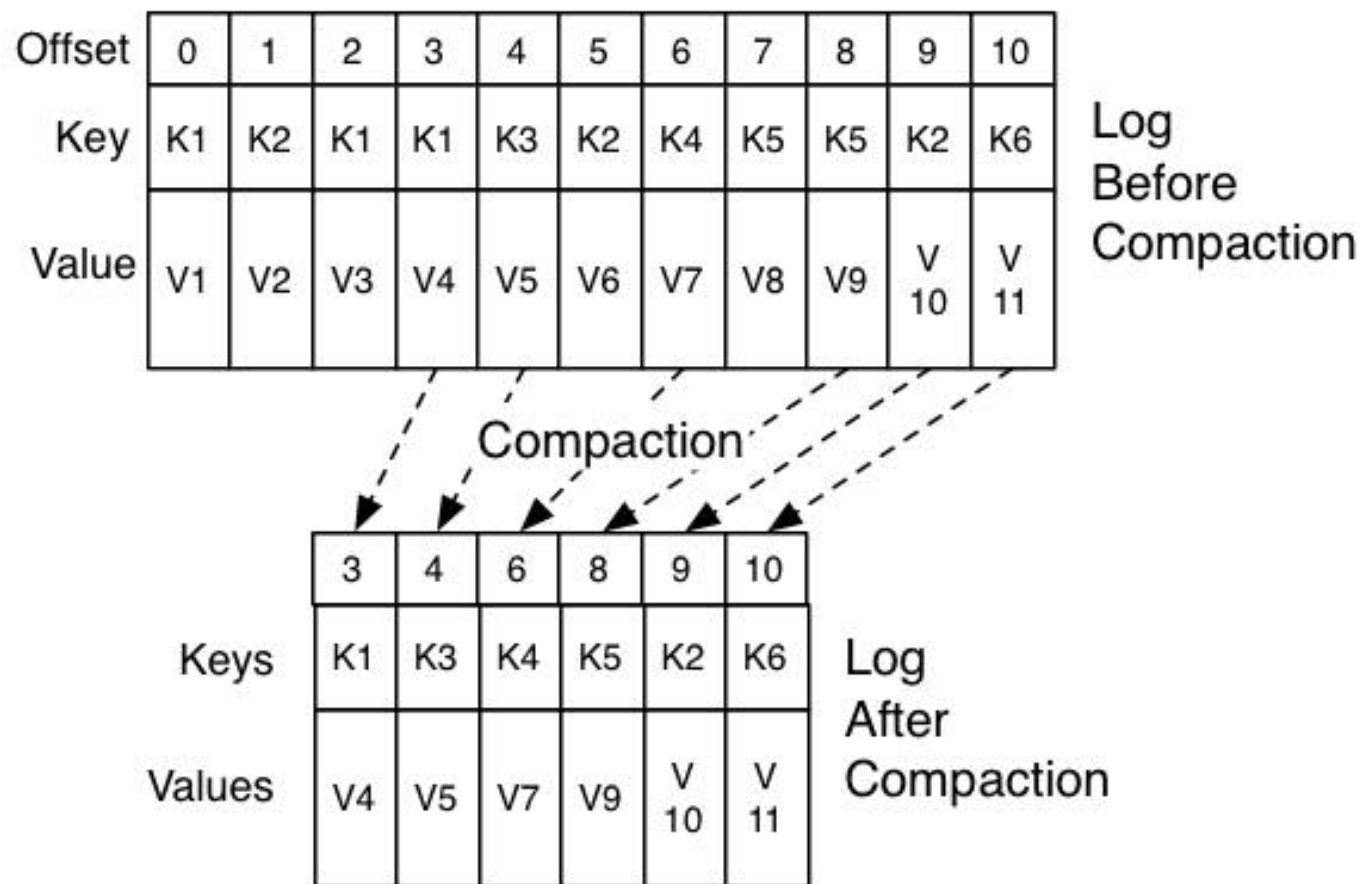
| Broker 1 | Broker 2 | Broker 3 |
|---|---|---|
| Topic A Partition 0 | Topic A Partition 0 | Topic A Partition 0 |
| Topic B Partition 1 | Topic B Partition 1 | Topic B Partition 1 |
| Topic C Partition 0 | Topic D Partition 0 | Topic C Partition 0 |
| Topic B Partition 0 | Topic B Partition 0 | Topic B Partition 0 |

# Kafka Messages

## Topic

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|

## Relational Table

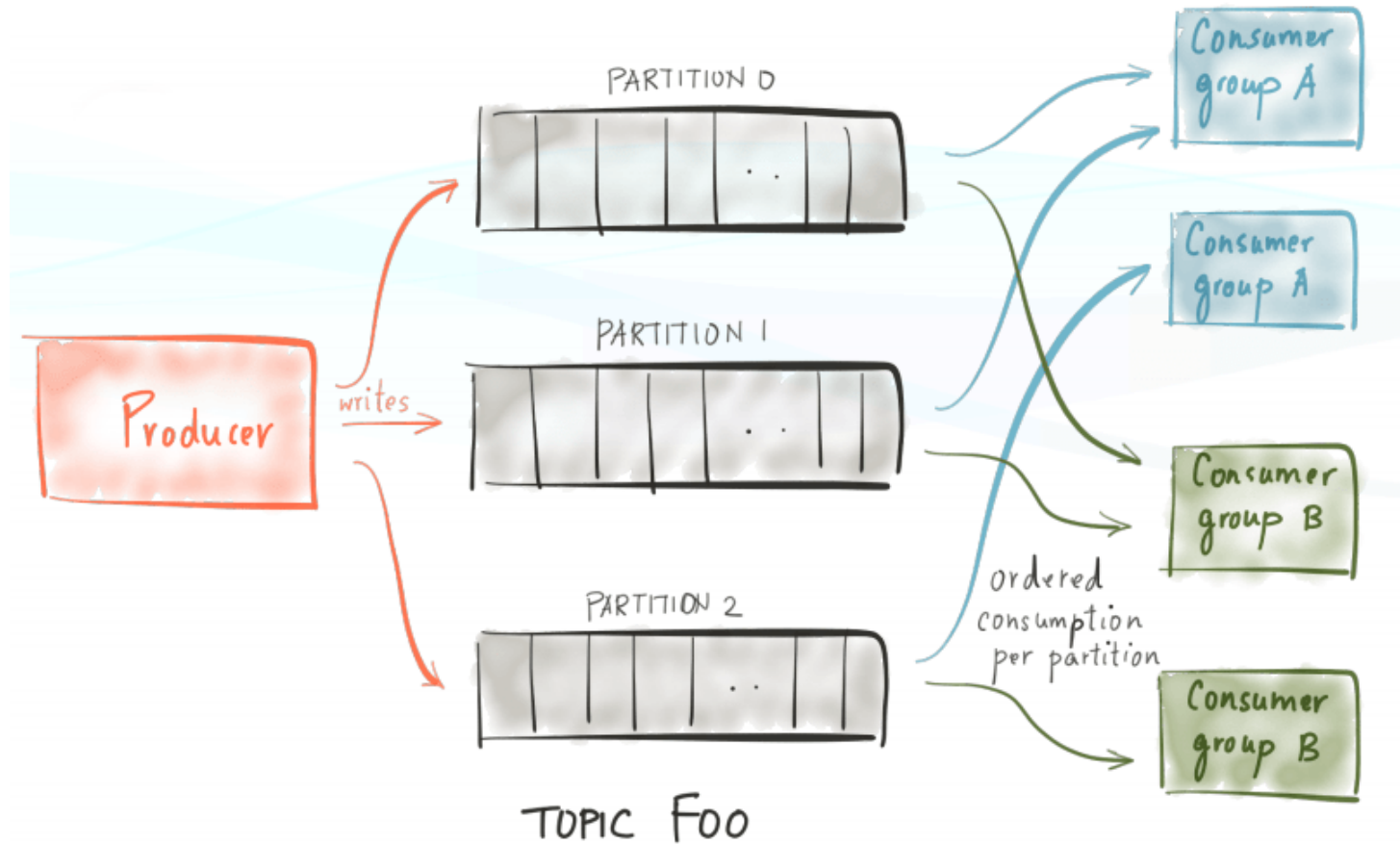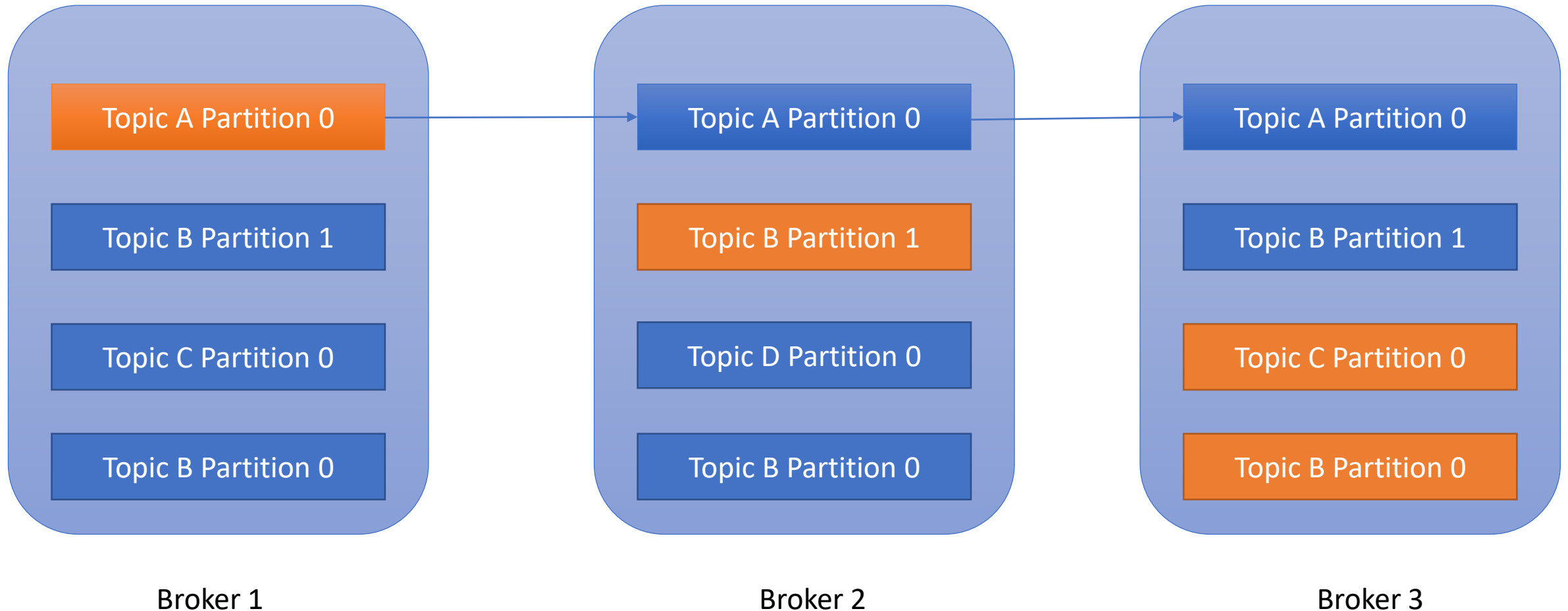| |
|---|
| 0 |
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |
| 6 |
| 4 |
| 7 |
| 8 |
| 9 |
| 10 |
| 11 |
| 12 |
| 13 |
| 14 |

# Kafka Messages
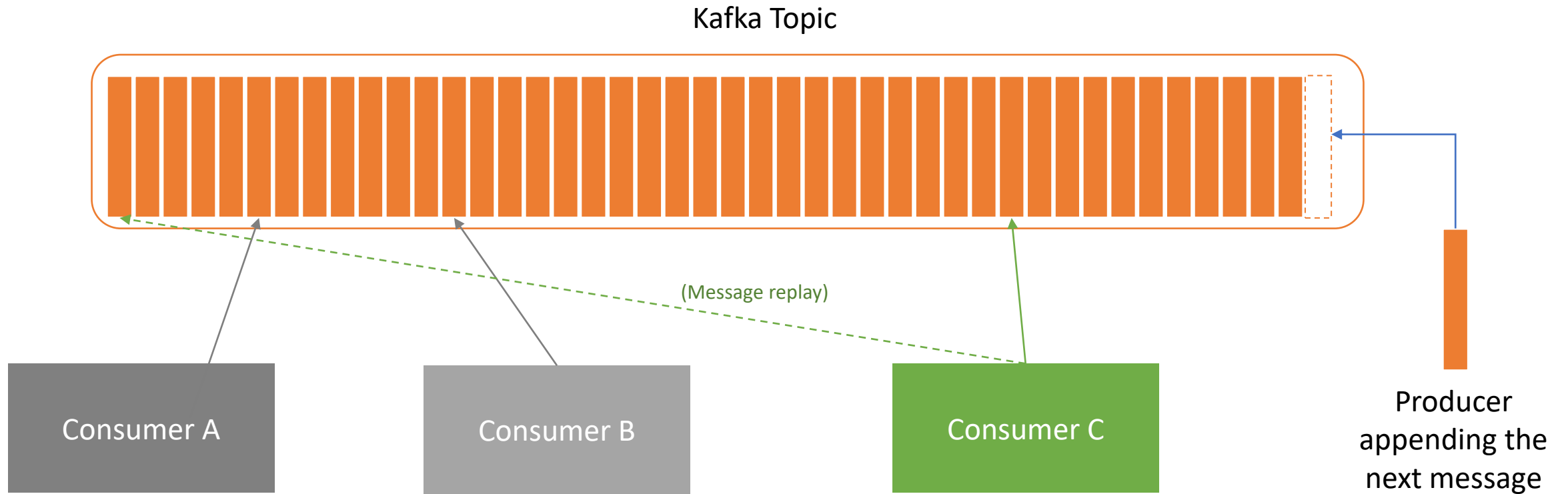# Compaction Clean up policy

# Kafka Producer

# Kafka Producer Acknowledgement

Topic A - 1 Partition & Replication factor of 3
Topic B - 2 Partition & Replication factor of 3
Topic C - 1 Partition & Replication factor of 2
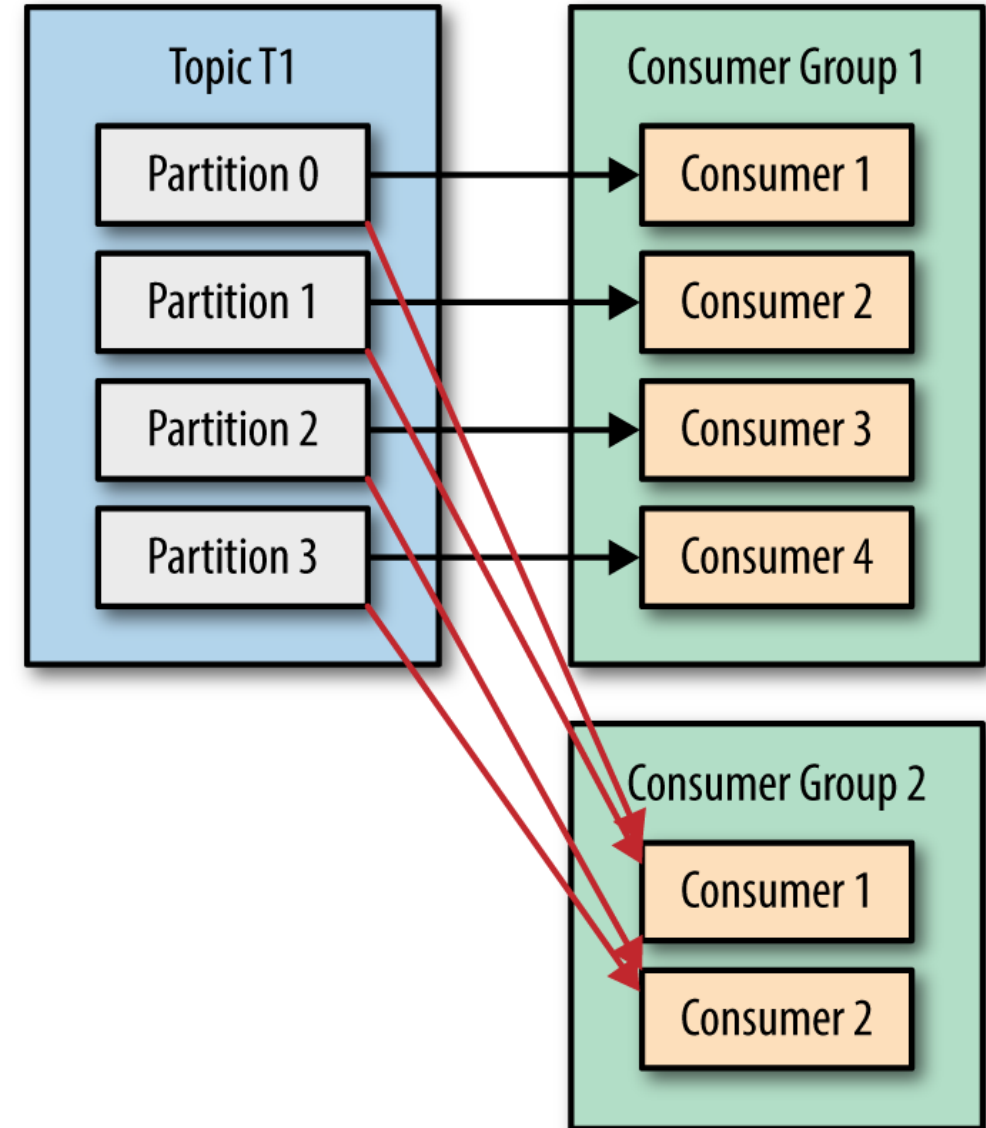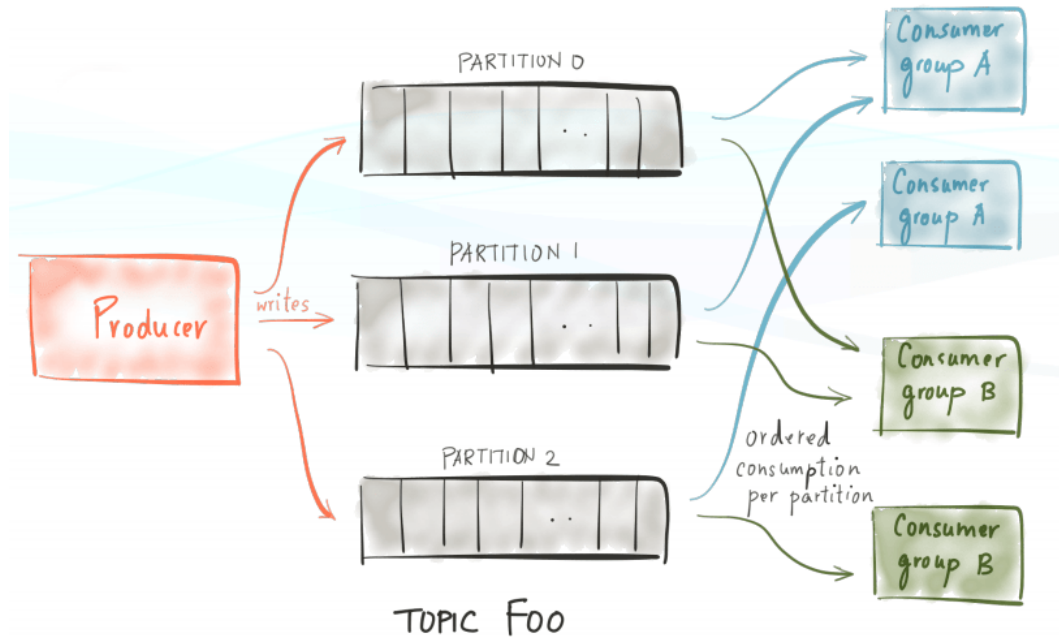Topic D - 1 Partition & Replication factor of 1

Topic A Partition 0

Topic B Partition 1

Topic C Partition 0

Topic B Partition 0

Topic A Partition 0

Topic B Partition 1

Topic D Partition 0

Topic B Partition 0

Topic A Partition 0

Topic B Partition 1

Topic C Partition 0

Topic B Partition 0

Broker 1

Broker 2

Broker 3

# Kafka Consumers



Kafka Topic

(Message replay)

Consumer A

Consumer B

Consumer C

Producer appending the next message

Consumers subscribe to Topics and read from them at their various offsets

# Kafka Consumer Groups

Use case

Toll payment collection and monitoring platform
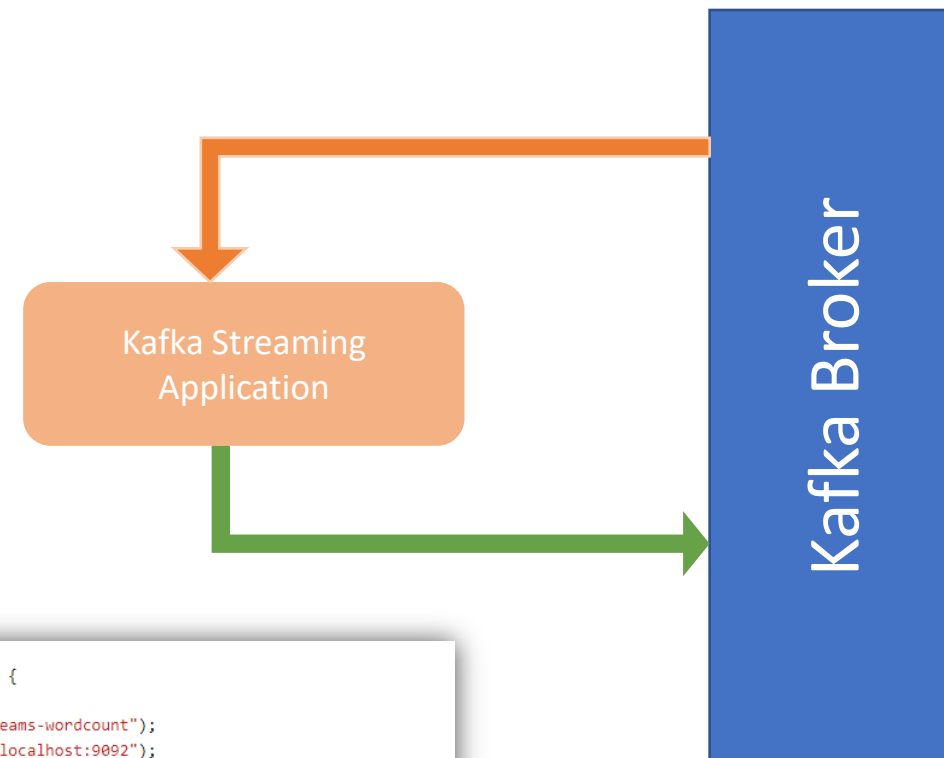
# Kafka Stream

A library for implementing scalable stream processing applications that react to events in real-time in Scala and Java.

It is more than just executing a function between a consumer and a producer.

❑ Stream abstraction

❑ Streams-Tables duality

❑ Time

  ❑ Windowing

❑ Aggregations

❑ State

**Kafka Streaming Application**

**Kafka Broker**

```java
public static void main(String[] args) throws Exception {
    Properties props = new Properties();
    props.put(StreamsConfig.APPLICATION_ID_CONFIG, "streams-wordcount");
    props.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, "localhost:9092");
    props.put(StreamsConfig.DEFAULT_KEY_SERDE_CLASS_CONFIG, Serdes.String().getClass());
    props.put(StreamsConfig.DEFAULT_VALUE_SERDE_CLASS_CONFIG, Serdes.String().getClass());

    final StreamsBuilder builder = new StreamsBuilder();

    KStream<String, String> source = builder.stream("streams-plaintext-input");
    source.flatMapValues(value -> Arrays.asList(value.toLowerCase(Locale.getDefault()).split("\\W+")))
          .groupBy((key, value) -> value)
          .count(Materialized.<String, Long, KeyValueStore<Bytes, byte[]>>as("counts-store"))
          .toStream()
          .to("streams-wordcount-output", Produced.with(Serdes.String(), Serdes.Long()));

    final Topology topology = builder.build();
    final KafkaStreams streams = new KafkaStreams(topology, props);
    final CountDownLatch latch = new CountDownLatch(1);

    // ... same as Pipe.java above
}
```
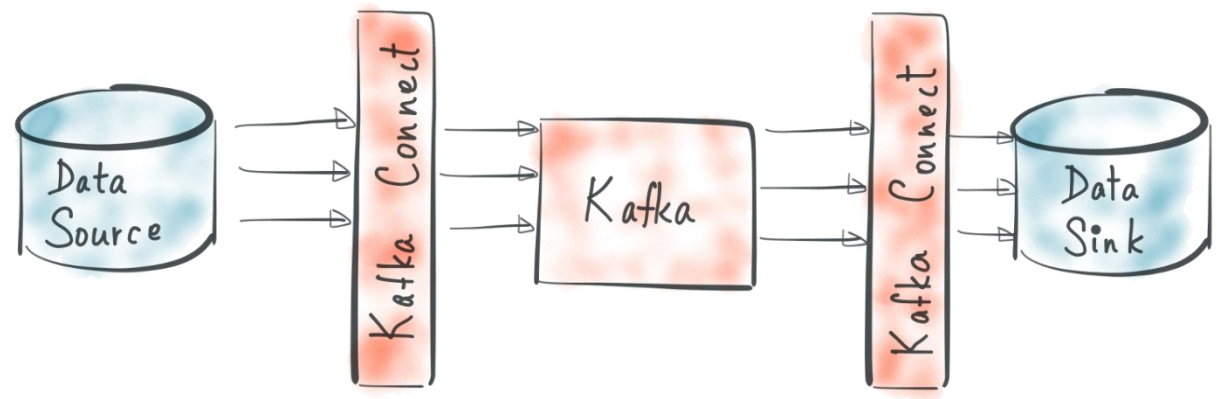
A library for reusable component from export data from a Kafka topic to an external destination or for import data from an external source into Kafka.
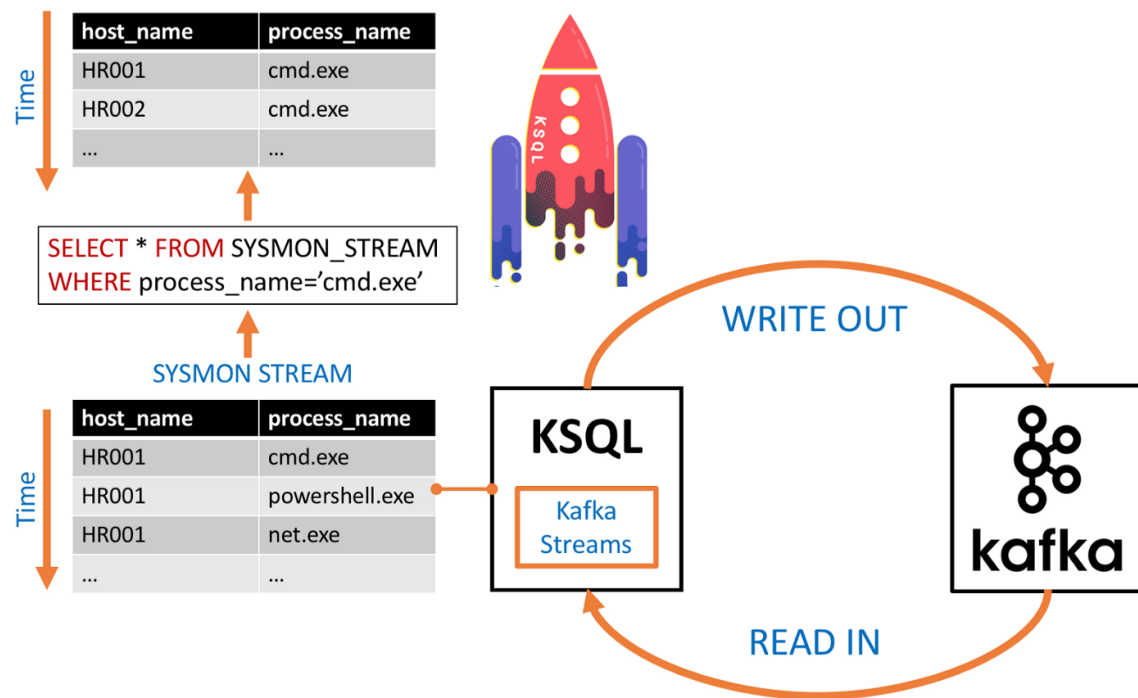
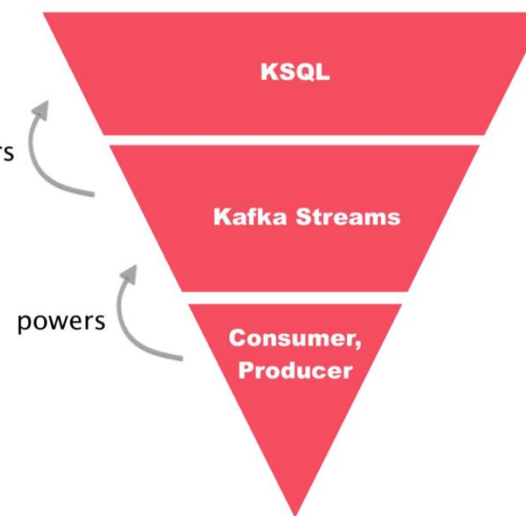The name for this reusable component is called a connector.

# KSQL

# DISTRIBUTED LOG
# WRAP-UP