





Need of Messaging Systems

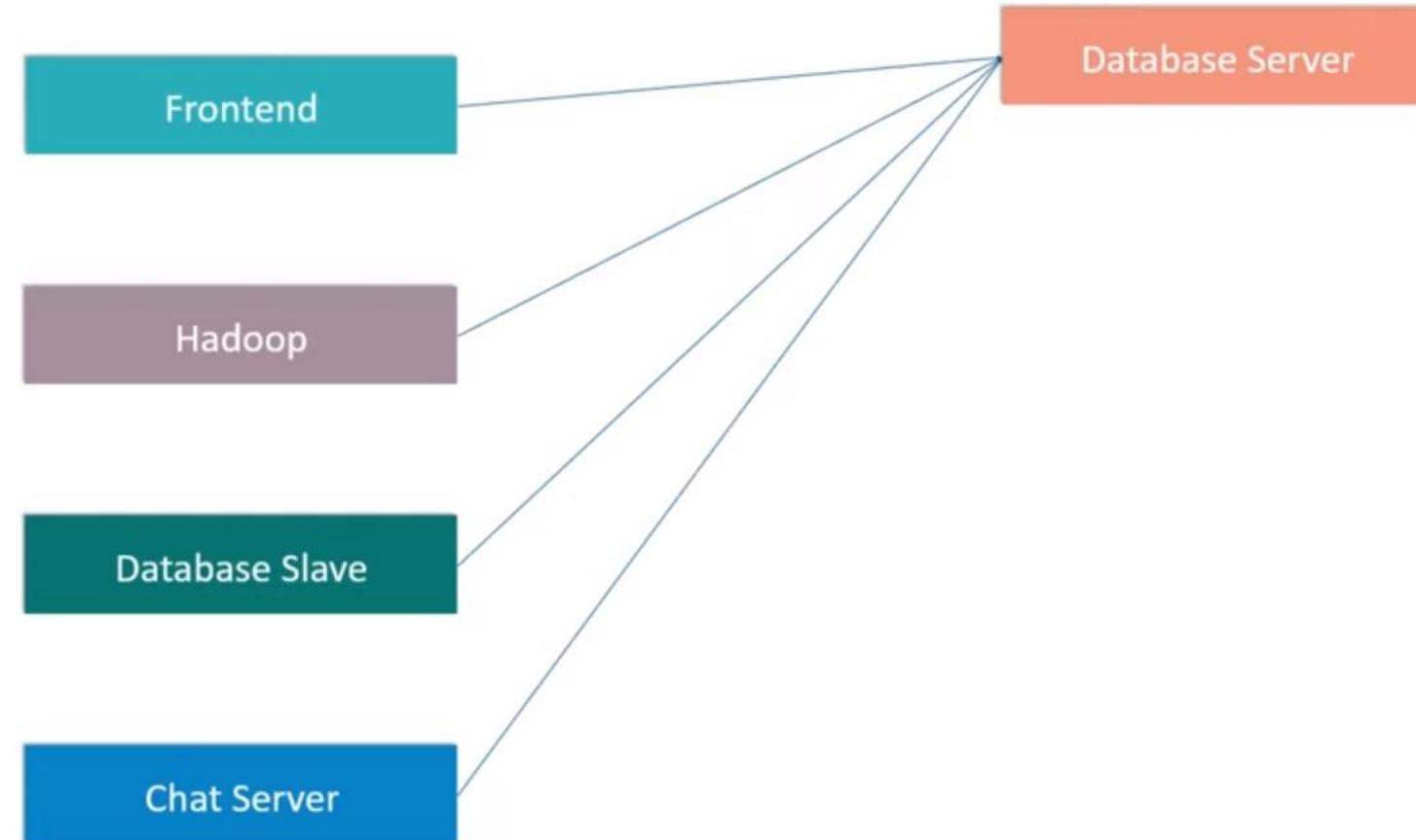
Data Pipelines

Communication is required between different systems in the real-time scenario, which is done by using data pipelines.



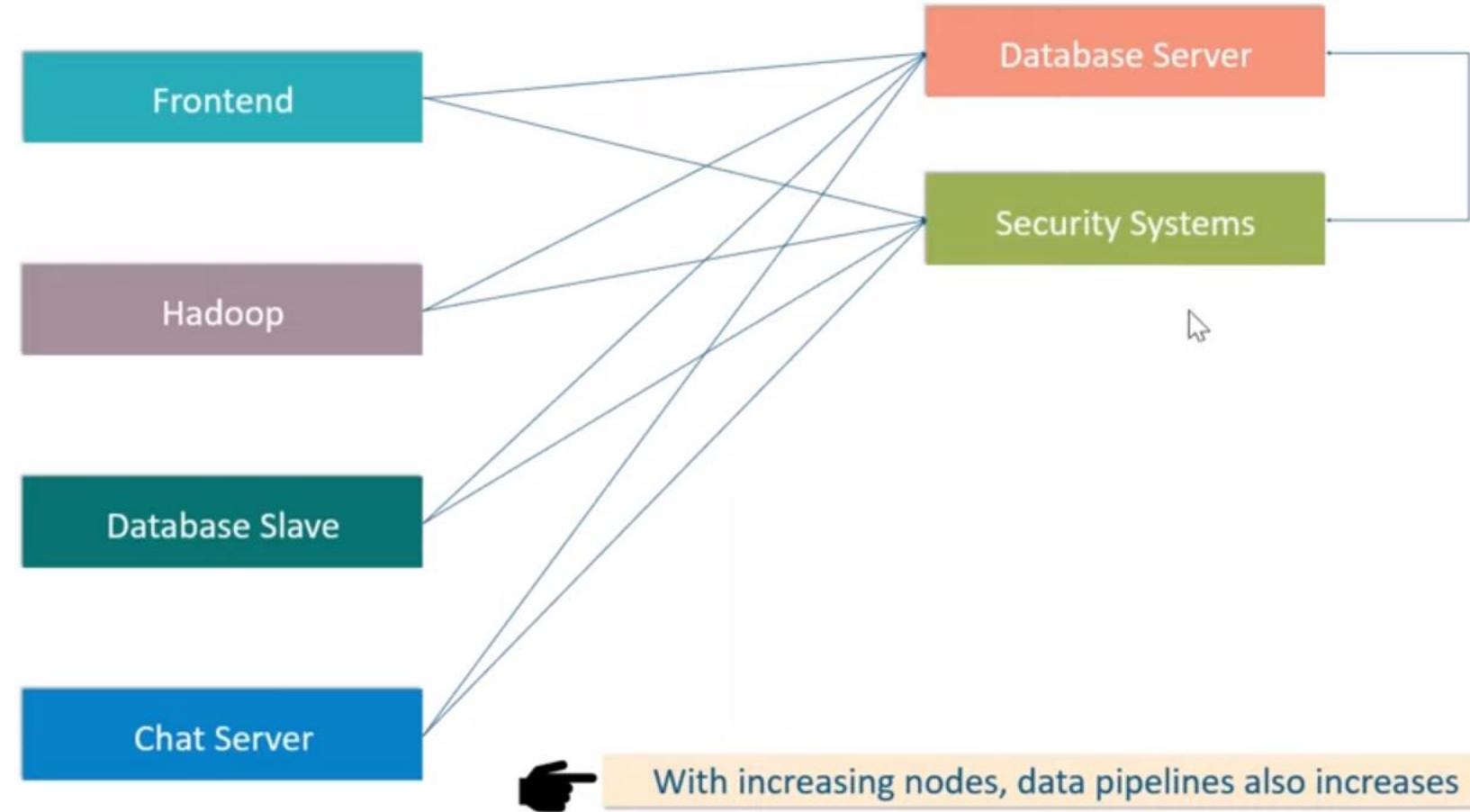
For Example: Chat Server needs to communicate with Database Server for storing messages

Increase in number of Nodes



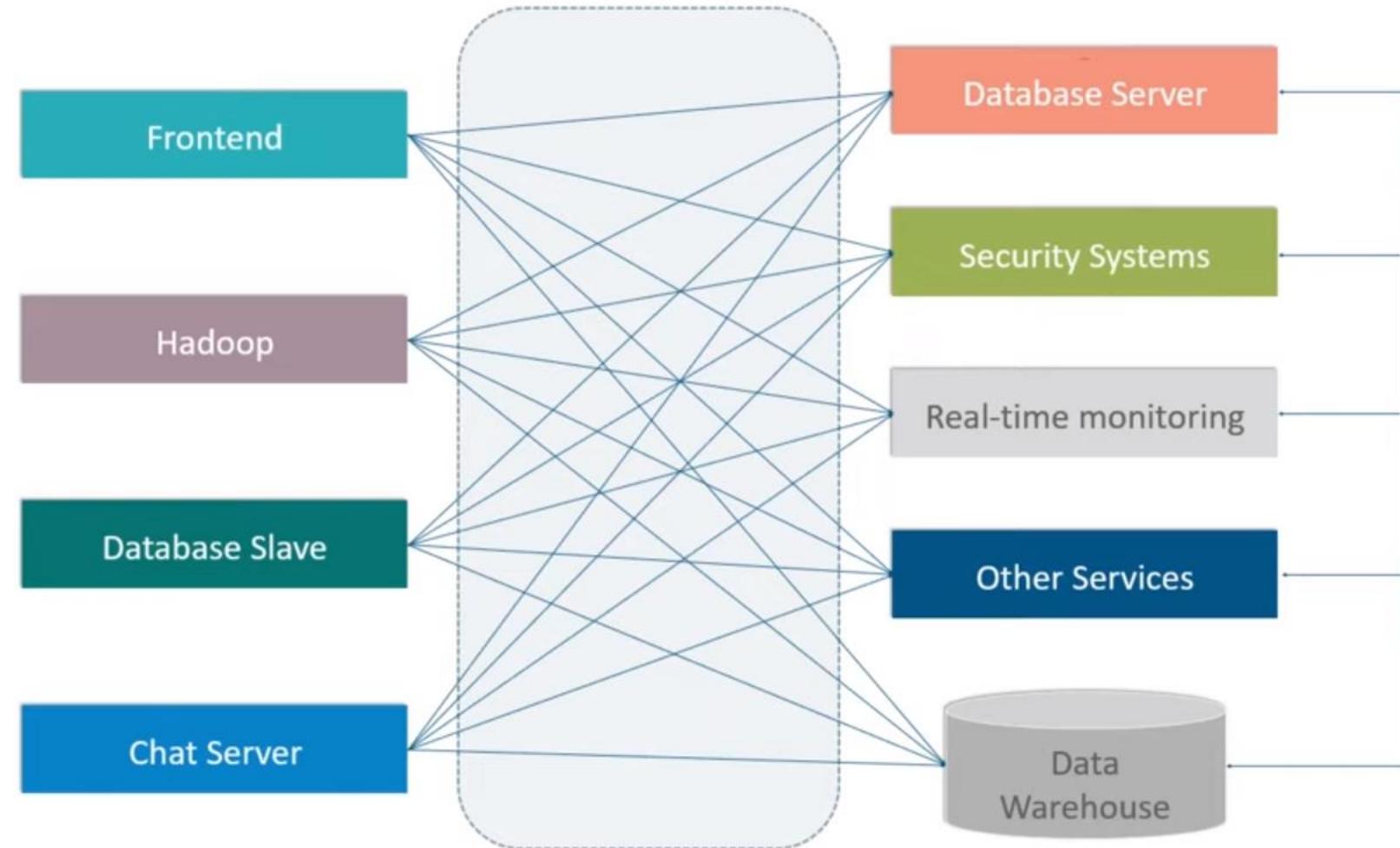
Similarly, there may be many applications wanting to access the Database Server

Increase in number of Nodes



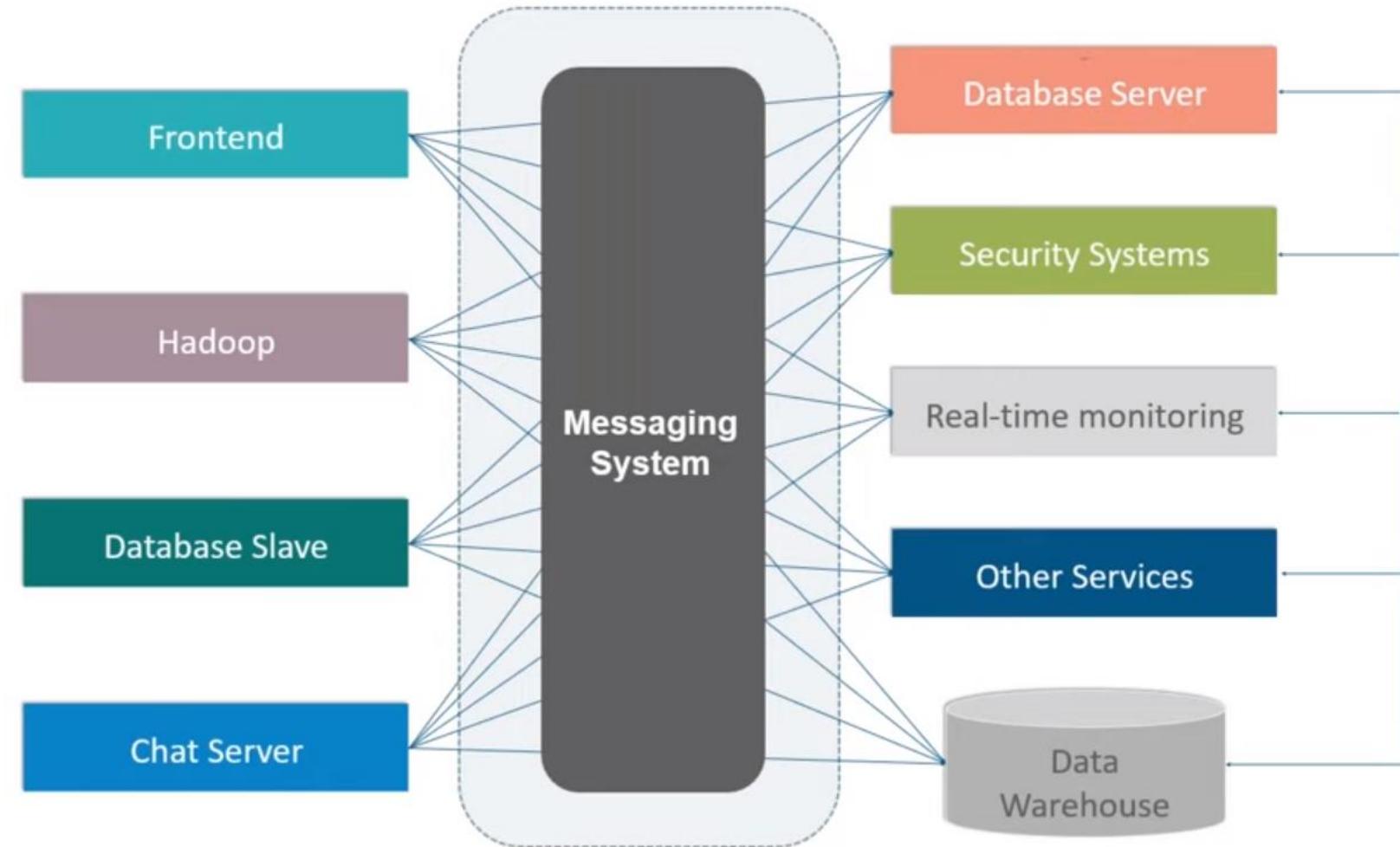
Complex Data Pipelines

Similarly, applications may also be communicating with Real-time monitoring and Other services in real-time scenario



Solution to the Complex Data Pipelines

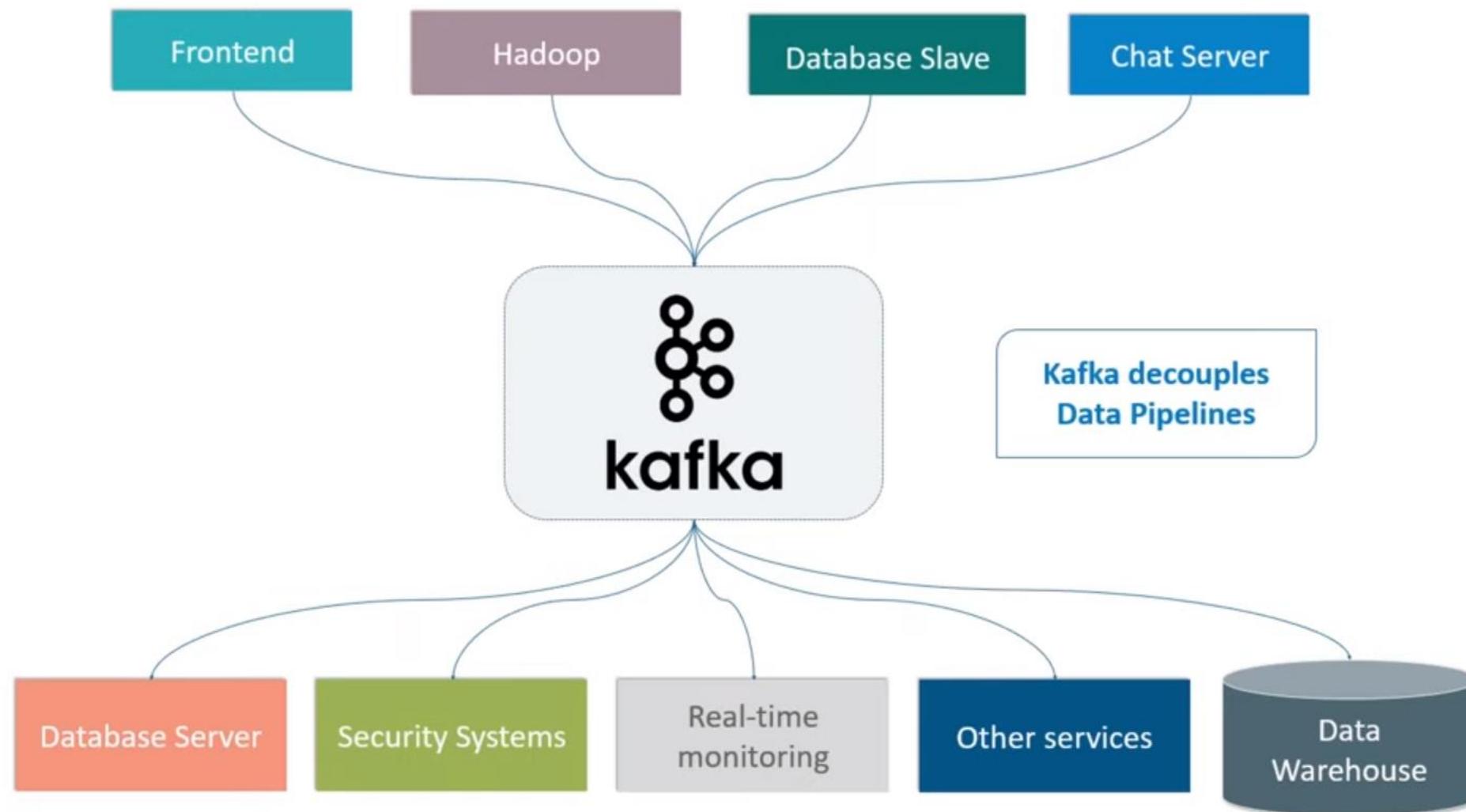
Messaging Systems helps managing the complexity of the pipelines





Let's See How Kafka Solves the Problem

Kafka Decouples Data Pipelines



What is Kafka?

- **Apache Kafka** is a distributed *publish-subscribe* messaging system
- It was originally developed at LinkedIn and later on became a part of Apache Project
- Kafka is fast, scalable, durable, fault-tolerant and distributed by design



Apache Kafka

A high-throughput distributed messaging system.

Kafka @LinkedIn

- 1100+ commodity machines
- 31,000+ topics
- 350,000+ partitions

- 675 billion messages/day
- 150 TB/day in
- 580 TB/day out

Peak Load

- 10.5 million messages/sec
- 18.5 GB/sec Inbound
- 70.5 GB/sec Outbound



Fig: A modern stream-centric data architecture built around Kafka

Kafka Growth Exploding

- More than **1/3** of all Fortune **500** companies use **Kafka**.
- These companies includes the top ten travel companies, **7** of top ten banks, **8** of top ten insurance companies, **9** of top ten telecom companies.
- **LinkedIn**, **Microsoft** and **Netflix** process billions of messages a day with Kafka (1,000,000,000,000).
- **Kafka** is used for **real-time streams** of data & used to collect big data for **real time analysis**.



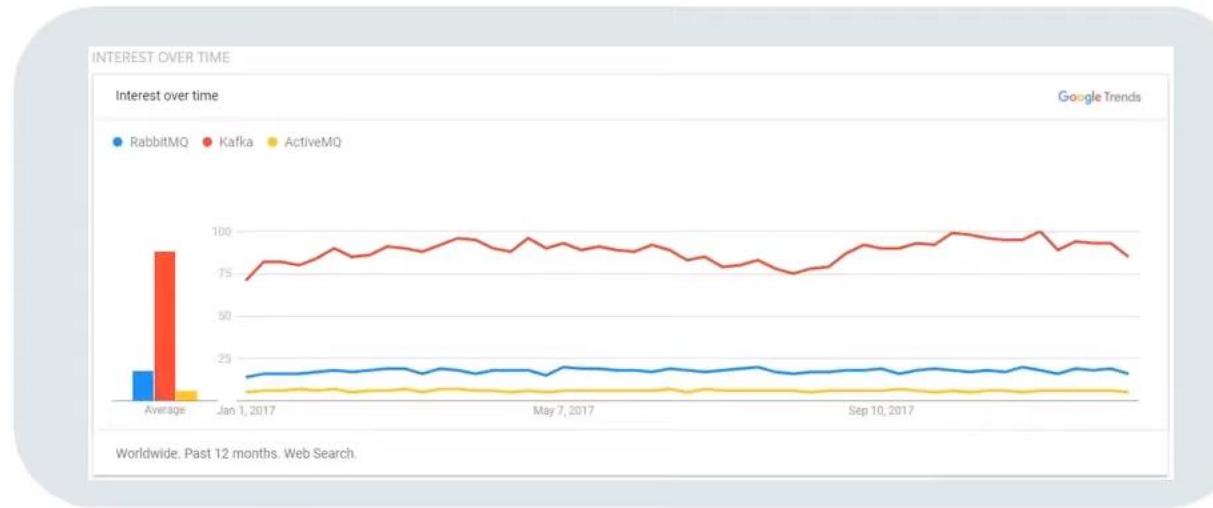
86% of respondents reported that the number of their systems that use Kafka is increasing



20% reported that the number is “growing a lot!”

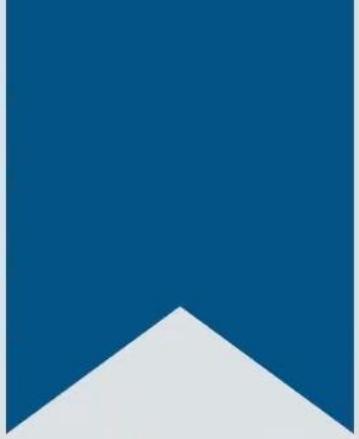


52% of organizations have at least **6** systems running Kafka



Source: Google Trends

Kafka Concepts



Kafka Terminologies

Producer

A **producer** can be any application who can publish messages to a topic

Consumer

A **consumer** can be any application that subscribes to a topic and consume the messages

Partition

Topics are *broken up into ordered commit logs called partitions*

Broker

Kafka cluster is a set of servers, each of which is called a **broker**

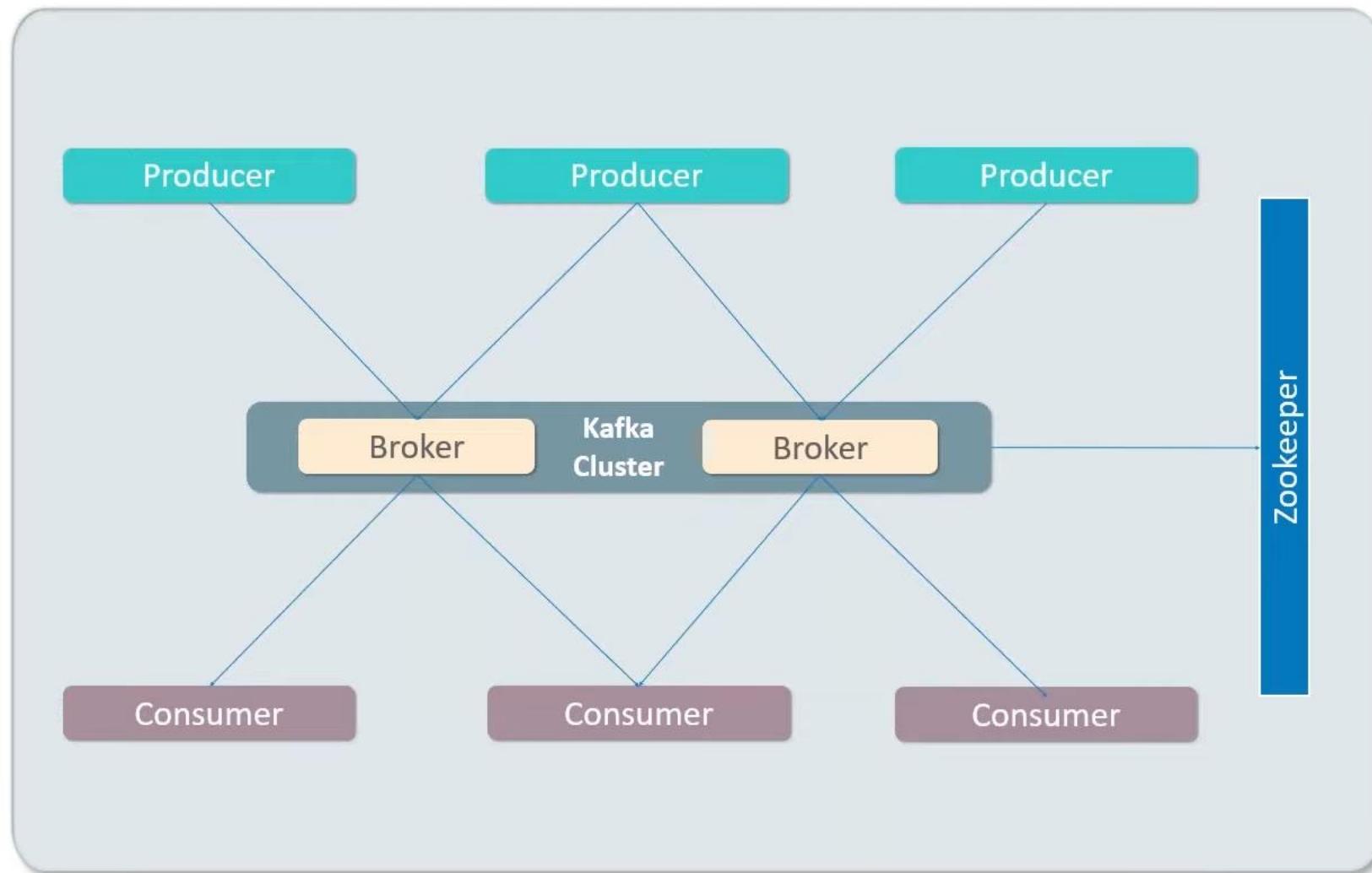
Topic

A **topic** is a category or *feed name* to which *records are published*

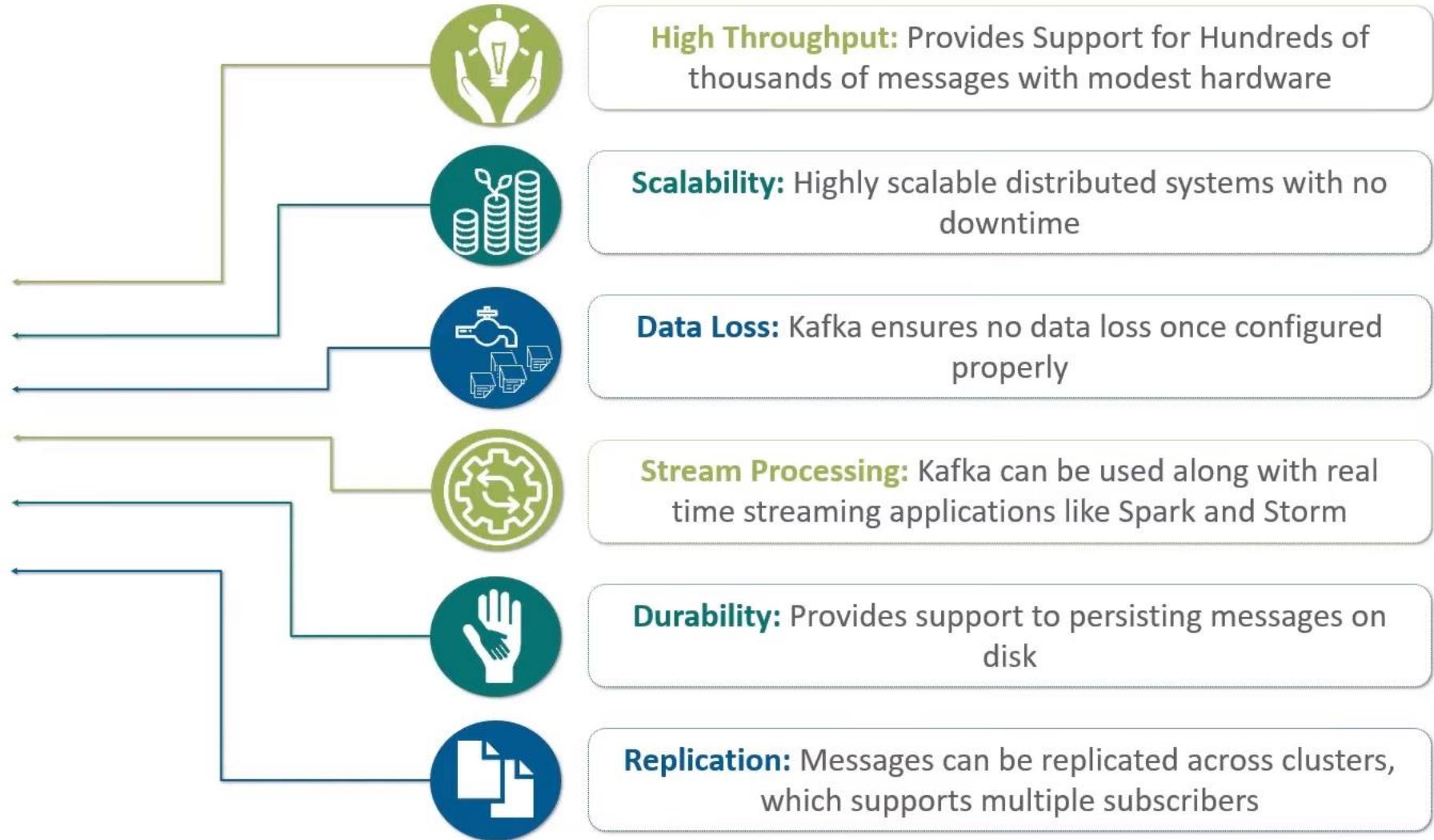
Zookeeper

ZooKeeper is used for managing and coordinating Kafka broker

Kafka Cluster



Kafka Features



Kafka Components - Topics and Partitions



A *topic* is a category or *feed name* to which *records are published*



Topics are broken up into *ordered commit logs* called *partitions*



Each *message* in a *partition* is assigned a *sequential id* called an *offset*



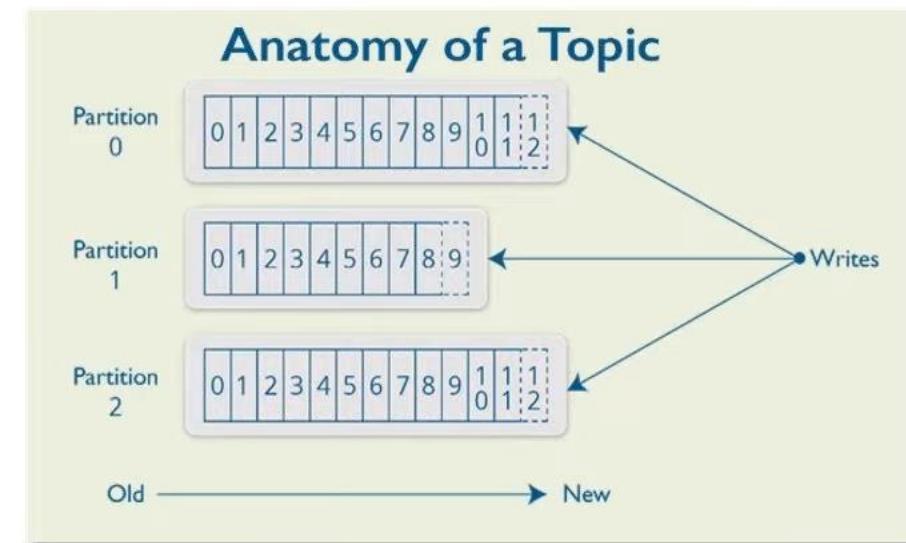
Data in a topic is retained for a *configurable period of time*



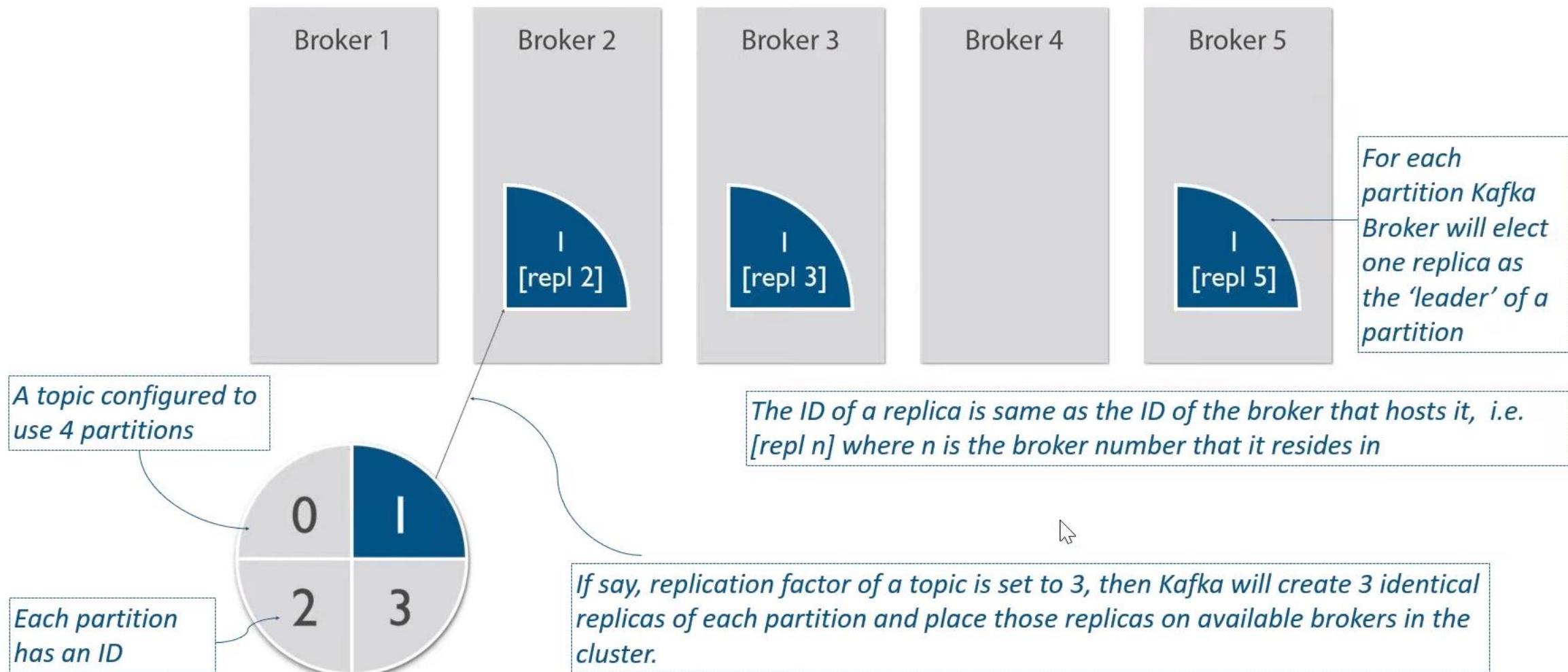
Writes to a partition are generally *sequential* thereby *reducing the number of hard disk seeks*



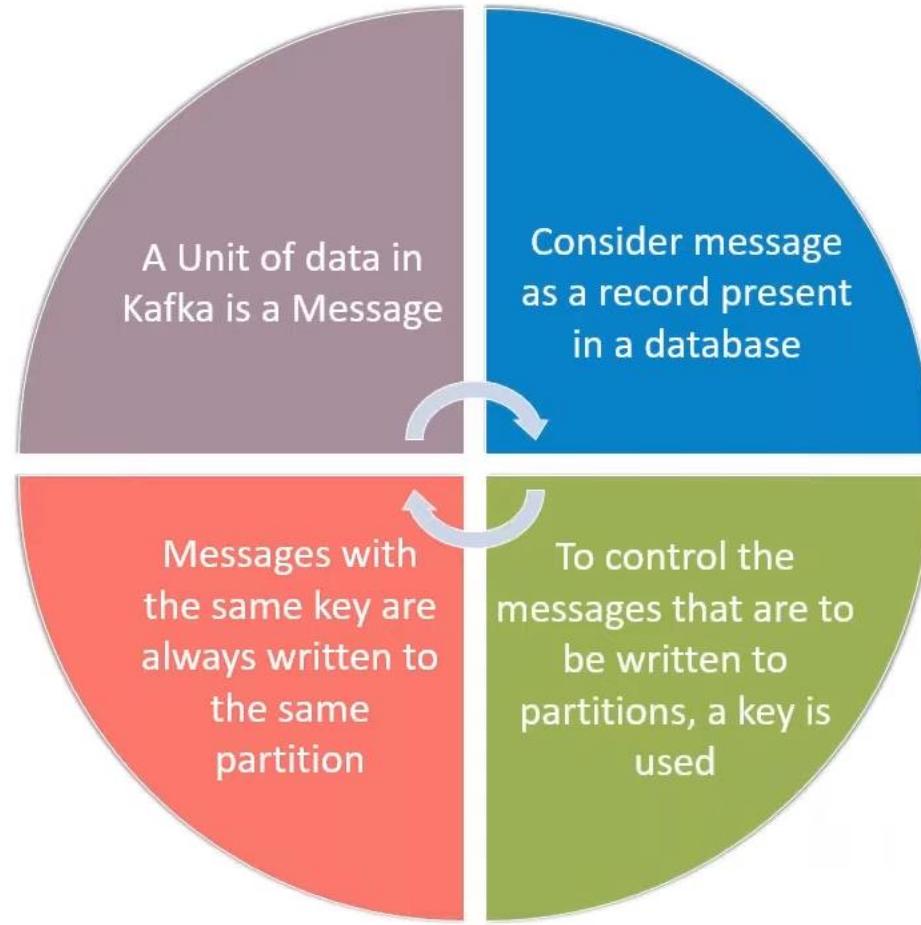
Reading messages can either be from *beginning* & also can *rewind* or skip to any point in partition by *giving an offset value*



Kafka Components - Topics, Partitions & Replicas



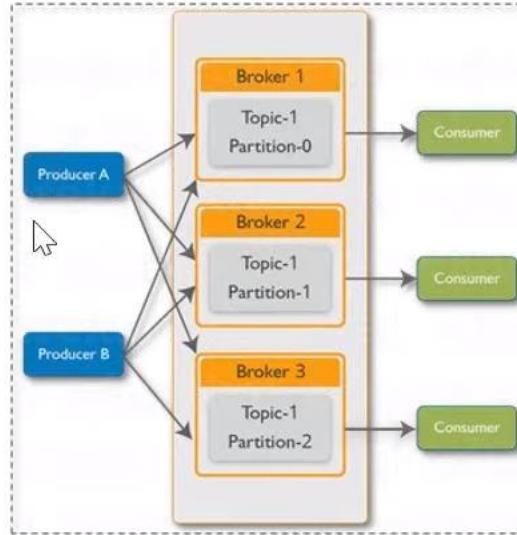
Kafka Components - Messages



Kafka Components - Producer

1

Producer (publisher or writer) publishes a new message to a **specific topic**

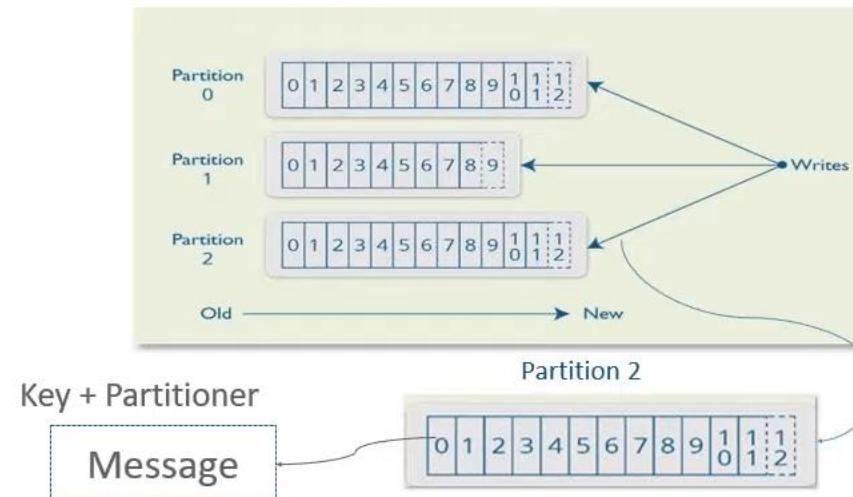


2

The producer does not care what partition a specific message is written to and will balance messages over every partition of a topic evenly

3

Directing messages to a partition is done using the **message key** and a **partitioner**, this will generate a hash of the key and map it to a partition

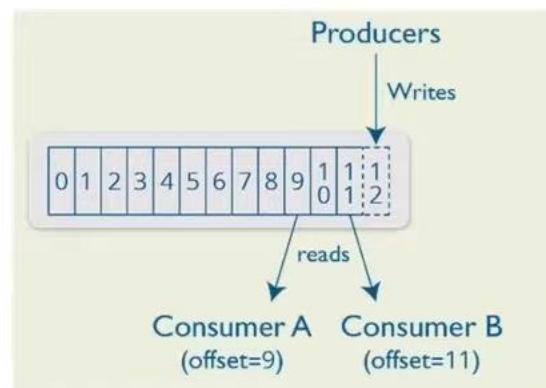


4

Every message a producer publishes in the form of a **key : value** pair

Kafka Components - Consumer

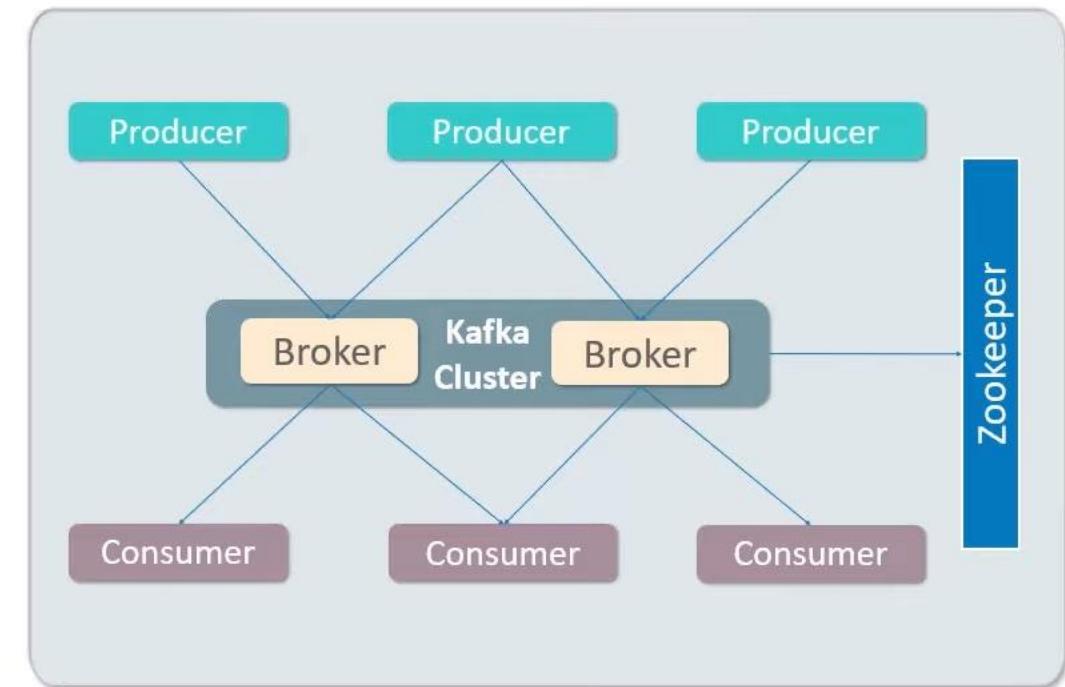
- Consumers(subscribers or readers) read messages
- The consumer subscribes to one or more topics and reads the messages sequentially
- The consumer keeps track of the messages it has consumed by keeping track on the offset of messages
- The *offset* is bit of metadata(an integer value that continually increases)that Kafka adds to each message as it is produced
- Each partition has a *unique offset* which is stored
- With the offset of the last consumed message, a consumer can *stop and restart without losing its current state*



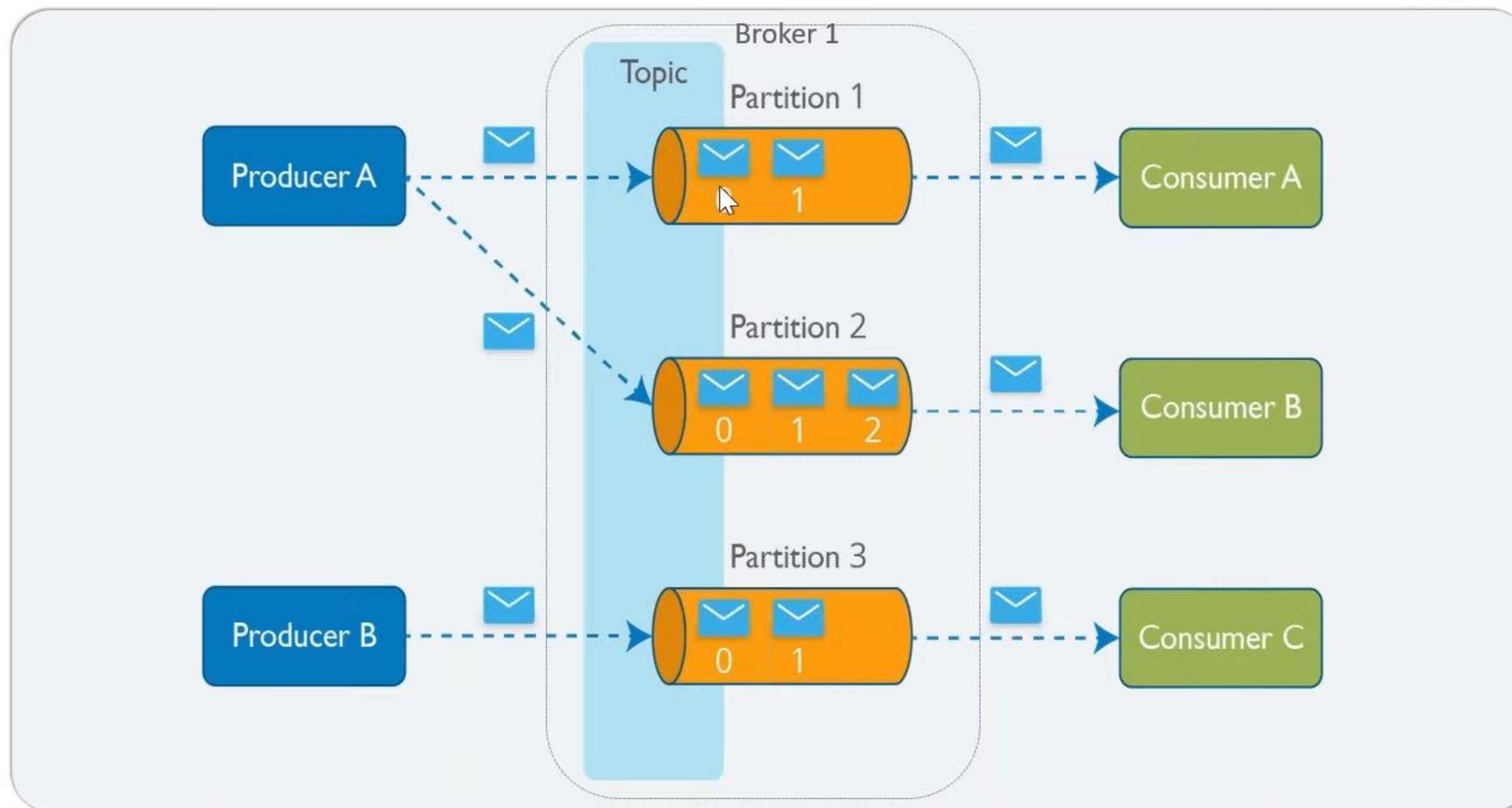
Kafka Components - ZooKeeper

ZooKeeper is used for managing and coordinating Kafka broker

- Zookeeper service is mainly used for co-ordinating between brokers in the Kafka cluster
- Kafka cluster is connected to ZooKeeper to get information about any failure nodes



Kafka Architecture



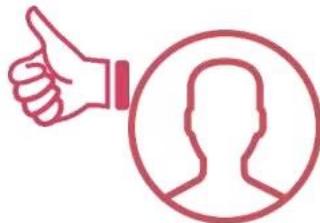
Let's see some Use Cases of Kafka

Kafka - Use Cases



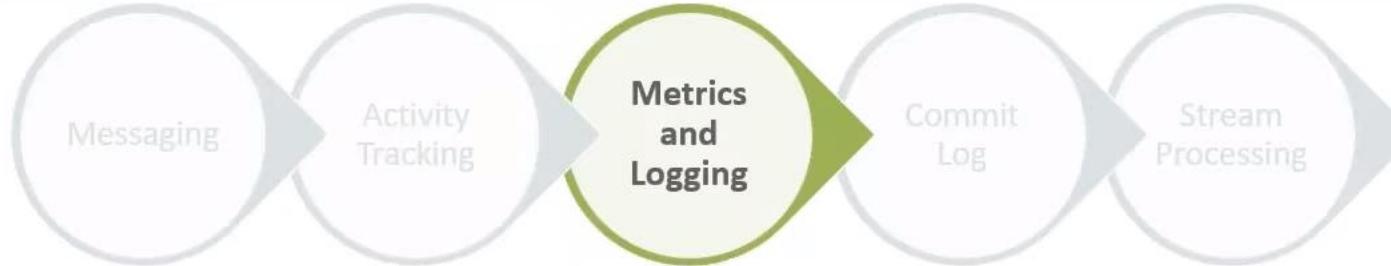
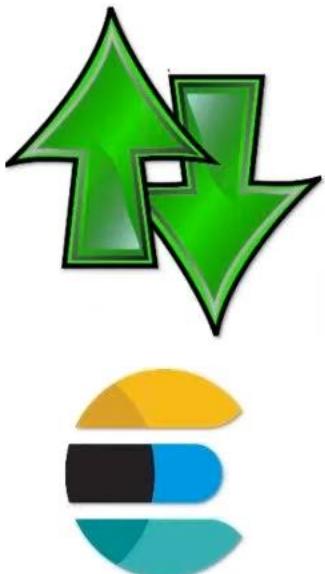
- Applications can produce messages using Kafka, without being concerned about the format of the messages
- Messages are sent and handled by a single application that can read all of them consistently, including :
 - A common formatting of messages using a common look
 - Send multiple messages in a single notification
 - Receive messages in a way that meets the users preferences

Kafka - Use Cases



- Originally Kafka was designed at LinkedIn, to track user activity
- When a user interacts with frontend applications, which generates messages regarding actions the user is taking
- Kafka keeps track of simple information like click tracking to complex information like data in a user's profile

Kafka - Use Cases



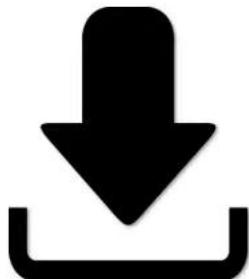
- Kafka is also ideal for collecting application's and system metrics and logs
- Applications publish metrics on a regular basis to a Kafka topic, and those metrics can be consumed by systems for monitoring and alerting
- Log messages can be published in the same way and routed to dedicated log search systems like Elasticsearch or security analysis applications

Kafka - Use Cases



- Database changes can be published to Kafka and applications can easily monitor this stream to receive live updates as they happen
- Kafka replicates database updates to a remote system for consolidating changes from multiple applications in a single database view
- Durable retention becomes useful providing a buffer for the changelog, meaning it can be replayed in the event of a failure of the consuming applications
- Log-compacted topics can be used to provide longer retention by only retaining a single change per key

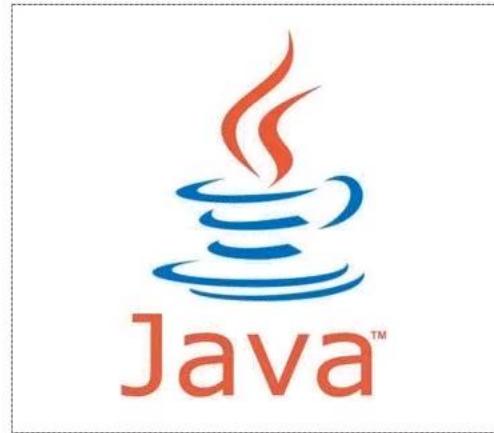
Kafka - Use Cases



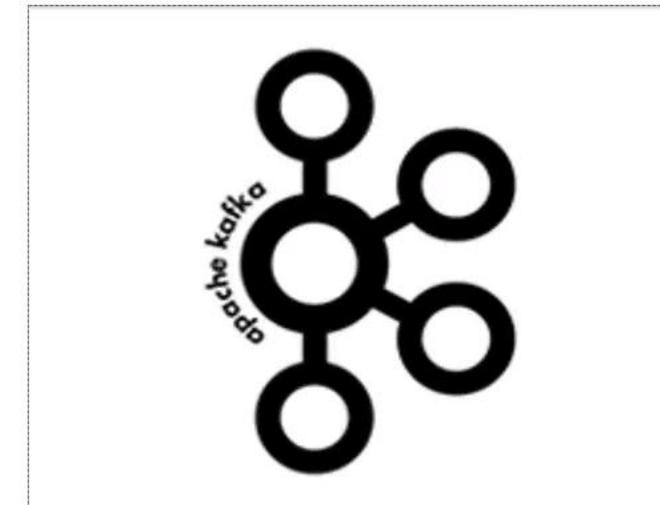
- Stream processing term is typically used to refer applications that provide similar functionality to map/reduce processing in Hadoop
- Stream processing operates on data in real-time, as quickly as messages are produced :
 - Write small applications to operate on Kafka messages,
 - Performing tasks such as counting metrics
 - Partitioning messages for efficient processing by other applications

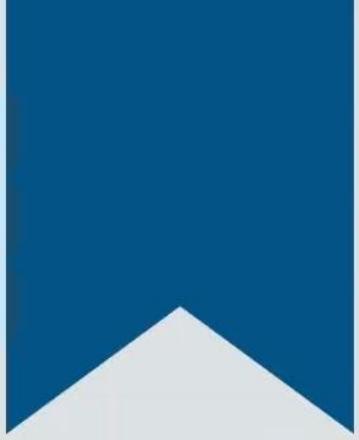
Getting Started with Kafka

- Prerequisites :



- Components :

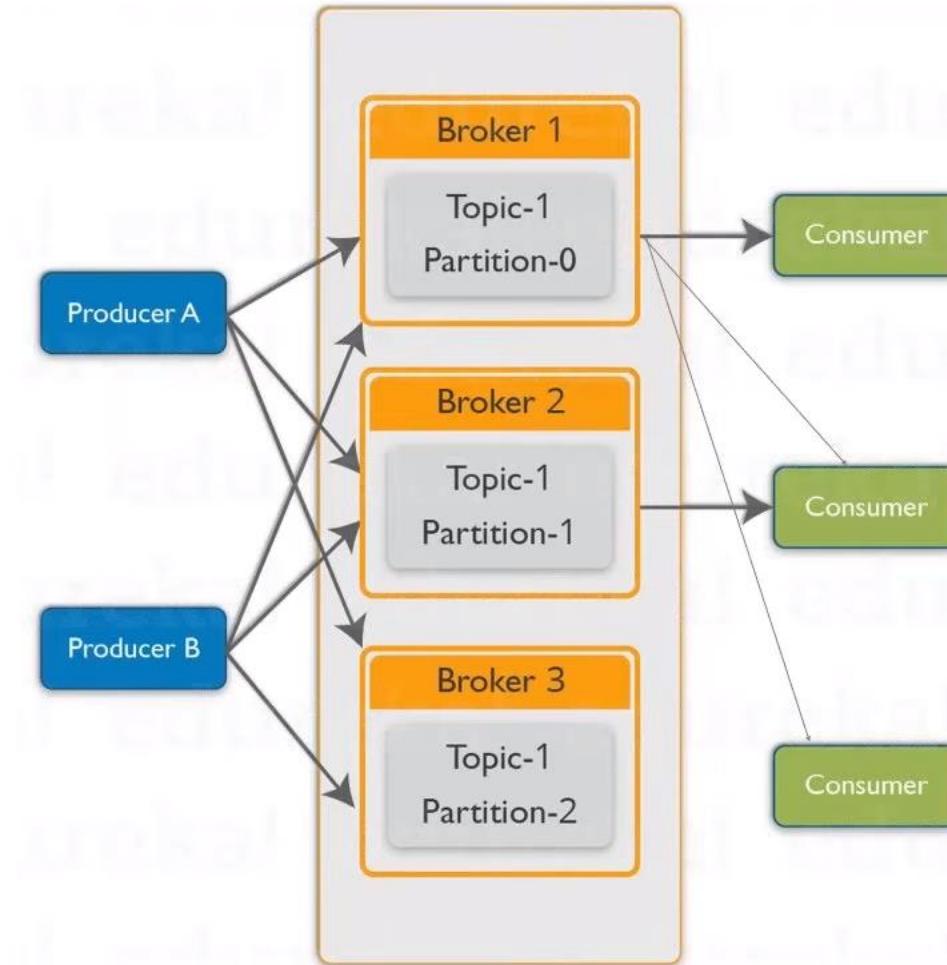




Let's Classify Different Types of Clusters in Kafka

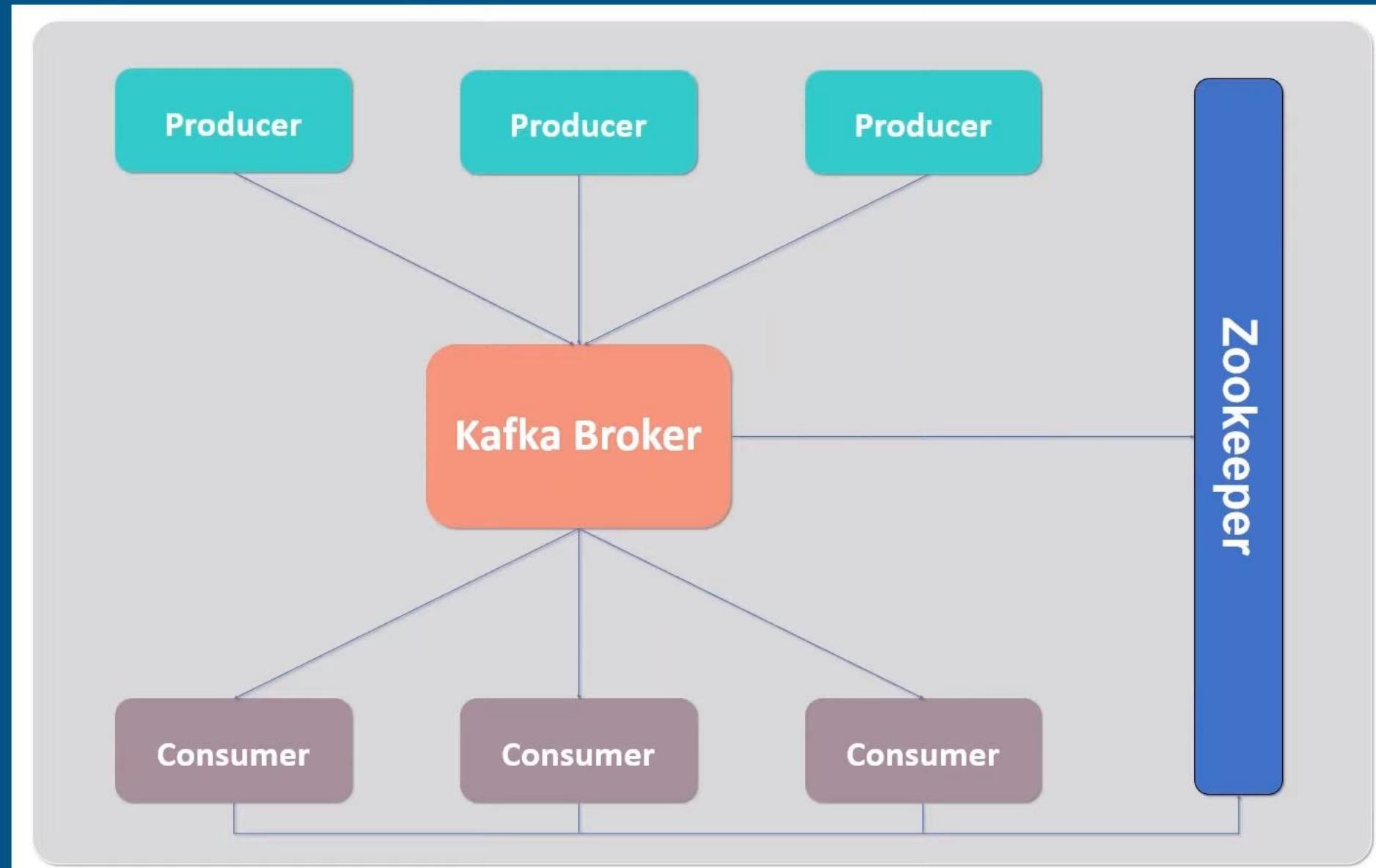
Kafka Cluster

- Kafka brokers are designed to operate as part of a cluster
- One broker will also function as the cluster controller
- Controller is responsible for administrative operations, like
 - Assigning partitions to brokers
 - Monitoring for broker failures in a cluster
- A particular partition is owned by a broker, and that broker is called the leader of the partition
- All consumers and producers operating on that partition must connect to the leader



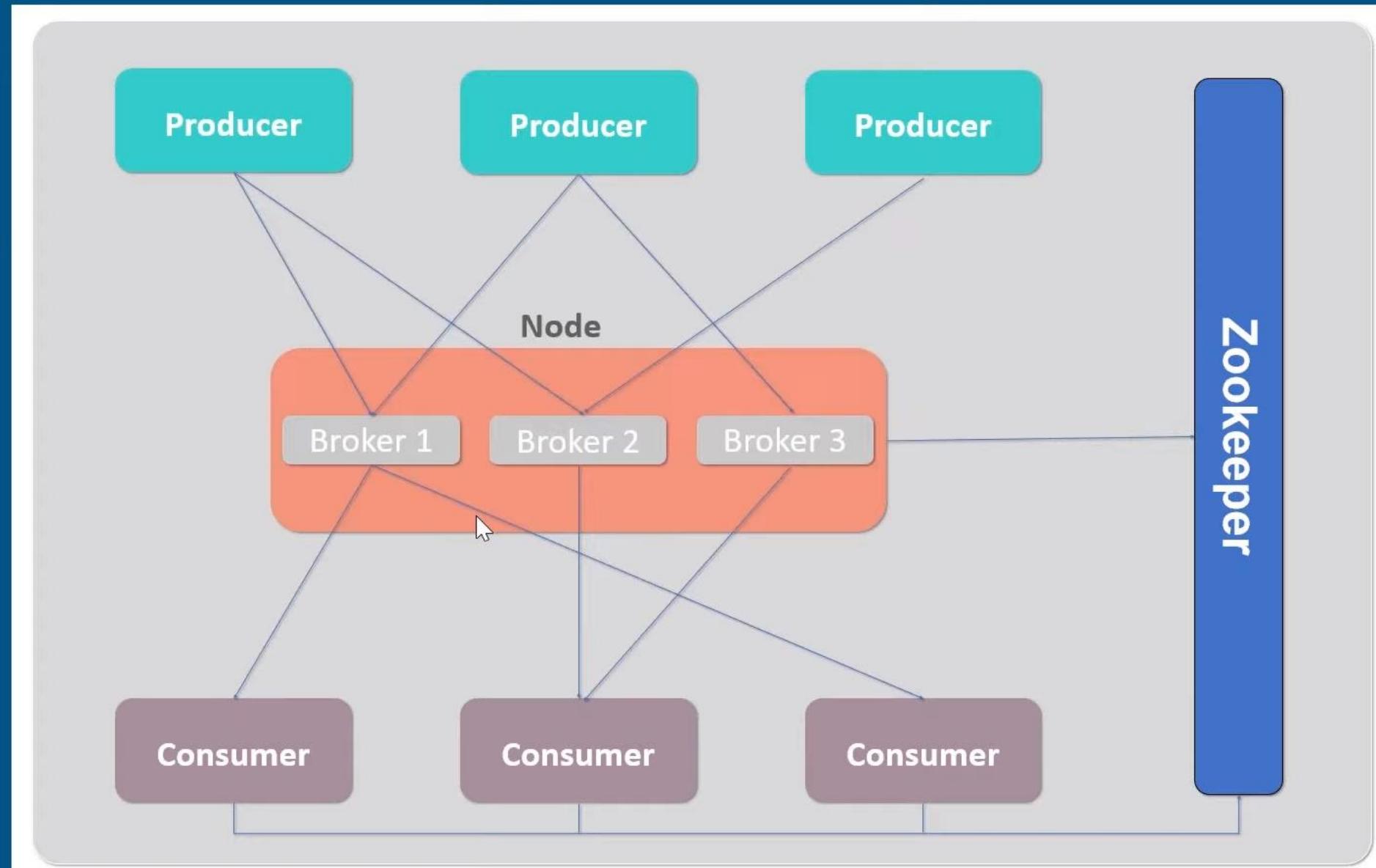
Type of Kafka Clusters

- 1 Single Node-Single Broker Cluster
- 2 Single Node-Multiple Broker Cluster
- 3 Multiple Nodes-Multiple Broker Cluster



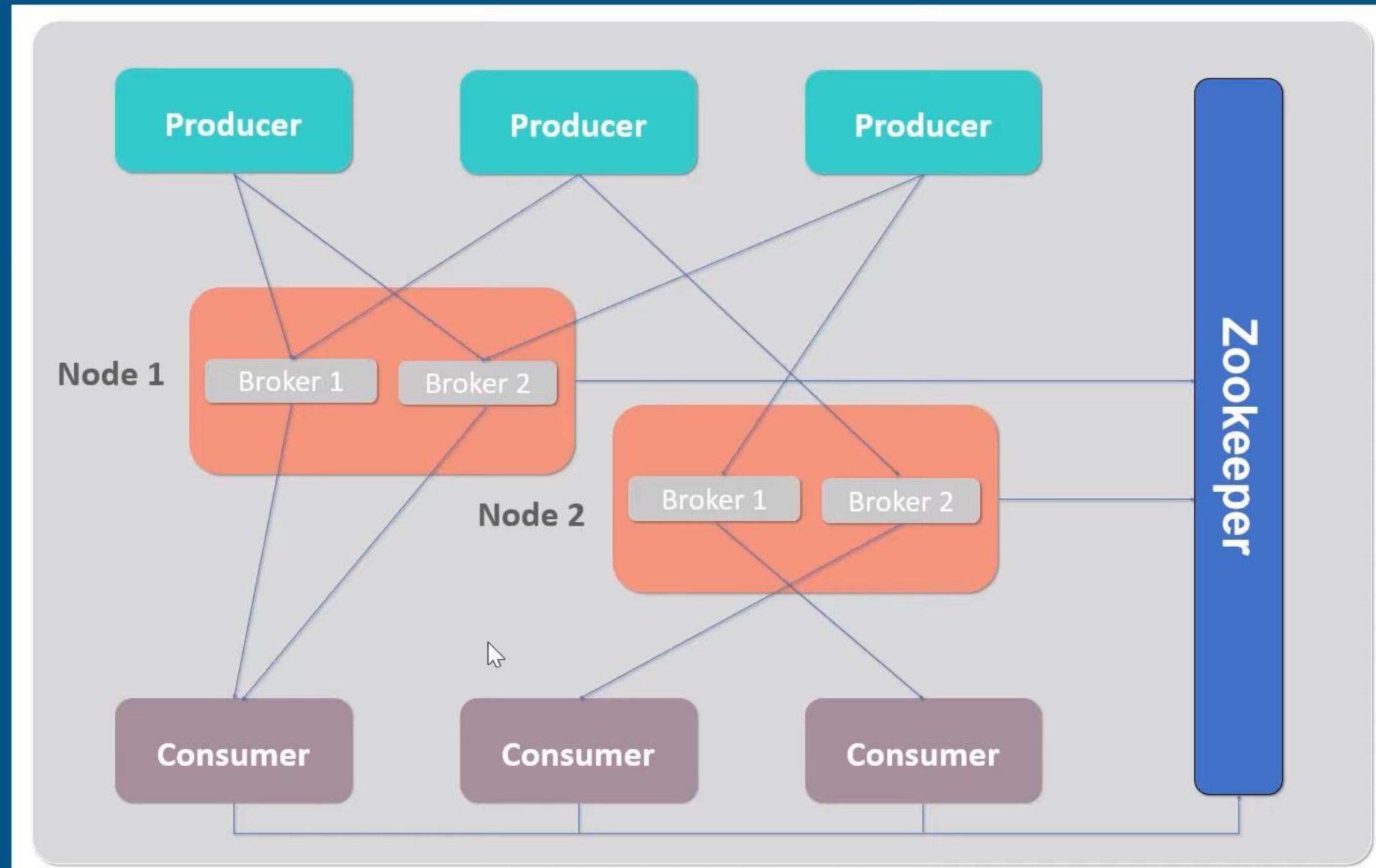
Type of Kafka Clusters

- 1 Single Node-Single Broker Cluster
- 2 Single Node-Multiple Broker Cluster
- 3 Multiple Nodes-Multiple Broker Cluster



Type of Kafka Clusters

- 1 Single Node-Single Broker Cluster
- 2 Single Node-Multiple Broker Cluster
- 3 Multiple Nodes-Multiple Broker Cluster



“

Getting Started with KAFKA & Spring Boot

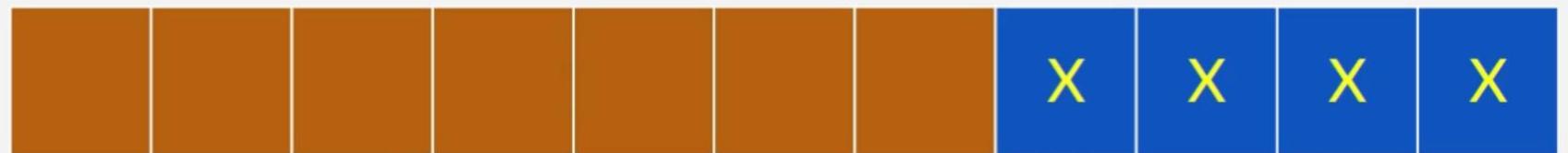
”

Producer & Consumer

Consumer Offset on First Run

auto.offset.reset =
latest
(default value)

Producer
start sending



Consumer started

auto.offset.reset =
earliest

Producer
start sending



Consumer started

Consumer
not started

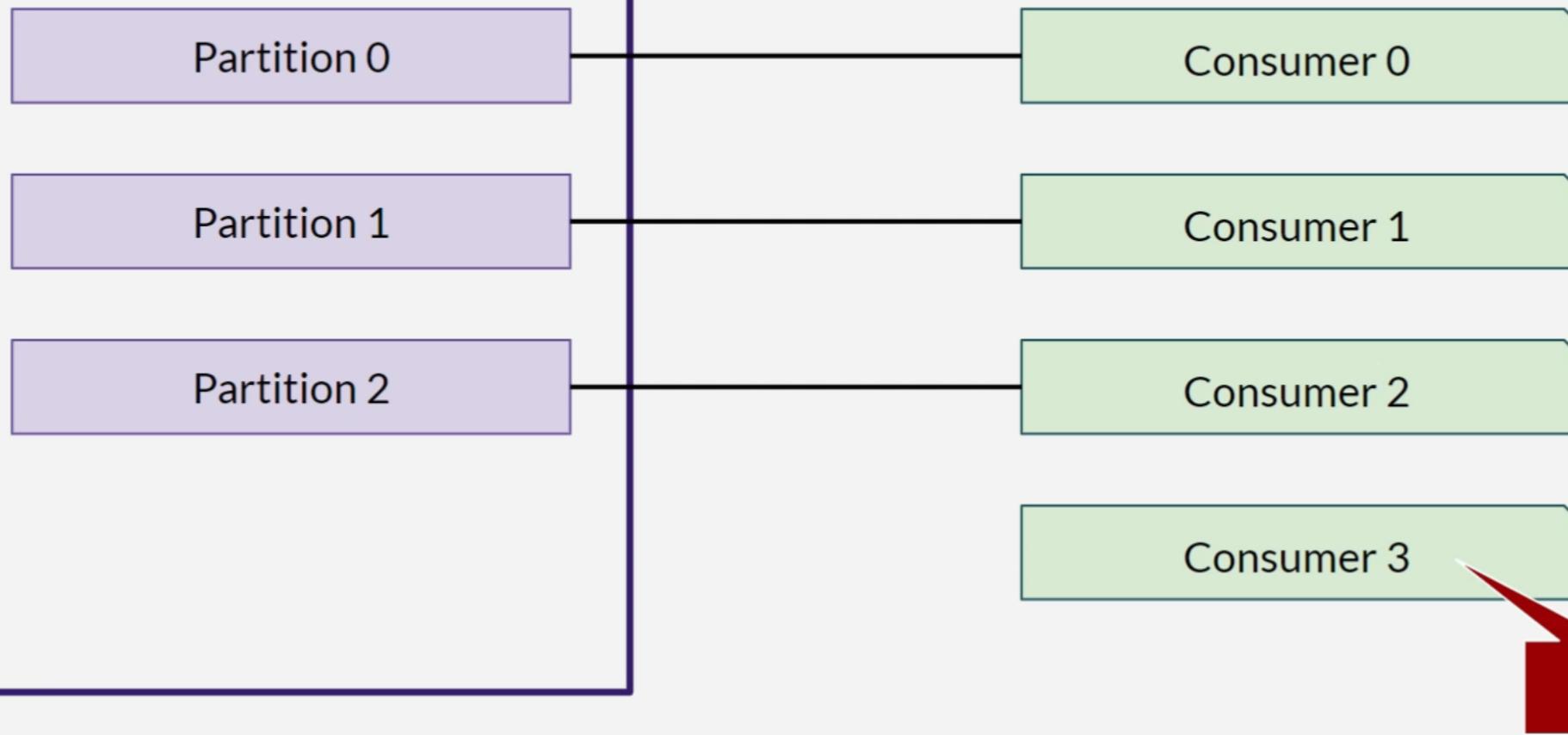
Message with key

- ▶ Same key goes to same partition
- ▶ As long as we don't change the partition number

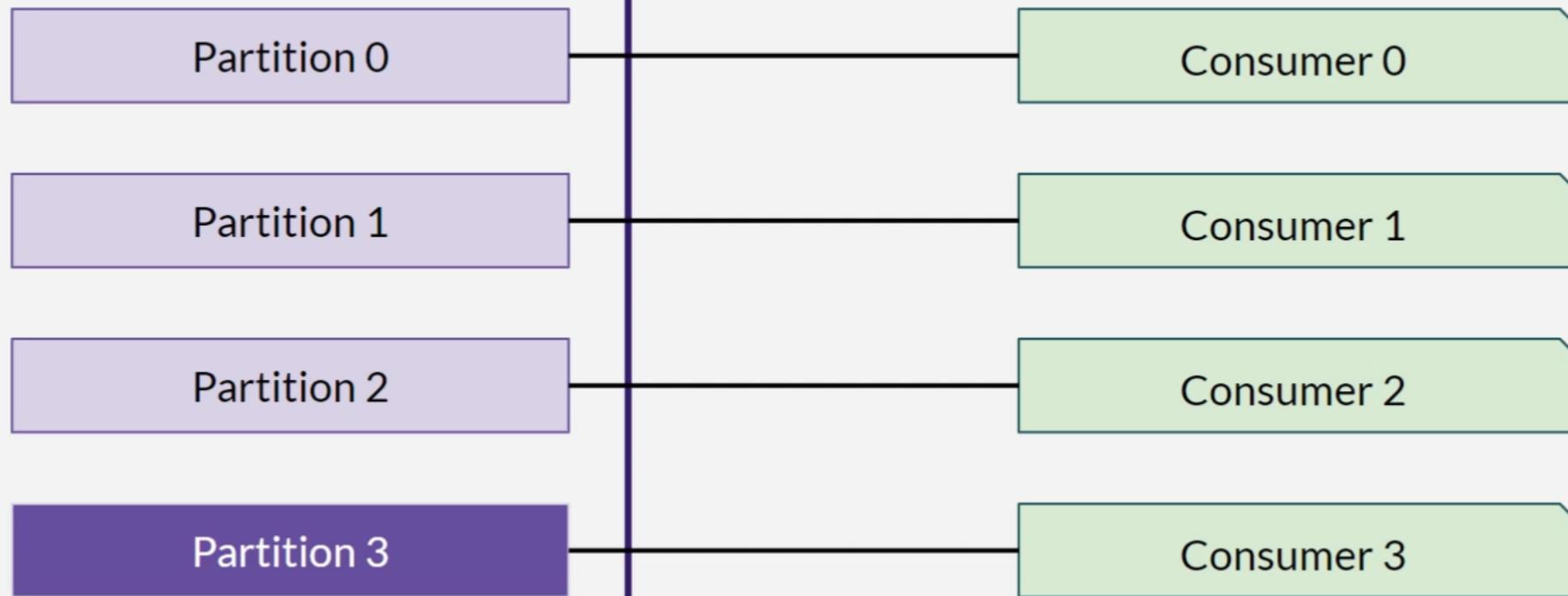
Multiple Consumers, Do we need it?

- ▶ Publisher (producer) works faster than consumer
- ▶ Subscriber (consumer) bottleneck

t-multi-partitions



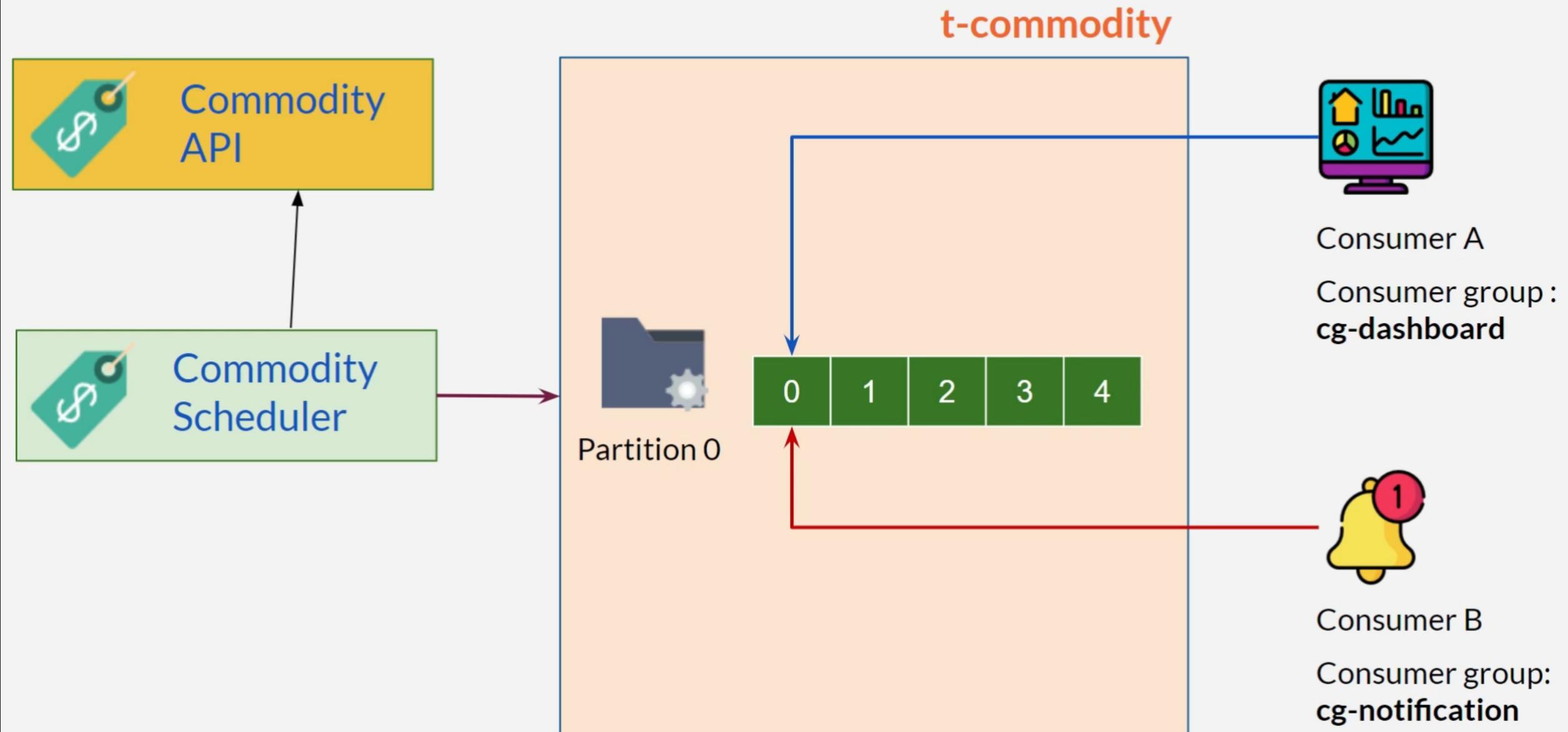
t-multi-partitions



What about deleting partition?

- ▶ Delete topic is OK
- ▶ Can't delete partition
- ▶ Can cause data loss
- ▶ Wrong key distribution
- ▶ Decrease partition > delete & recreate topic
- ▶ Real life usually use kafka native linux

Kafka Schema



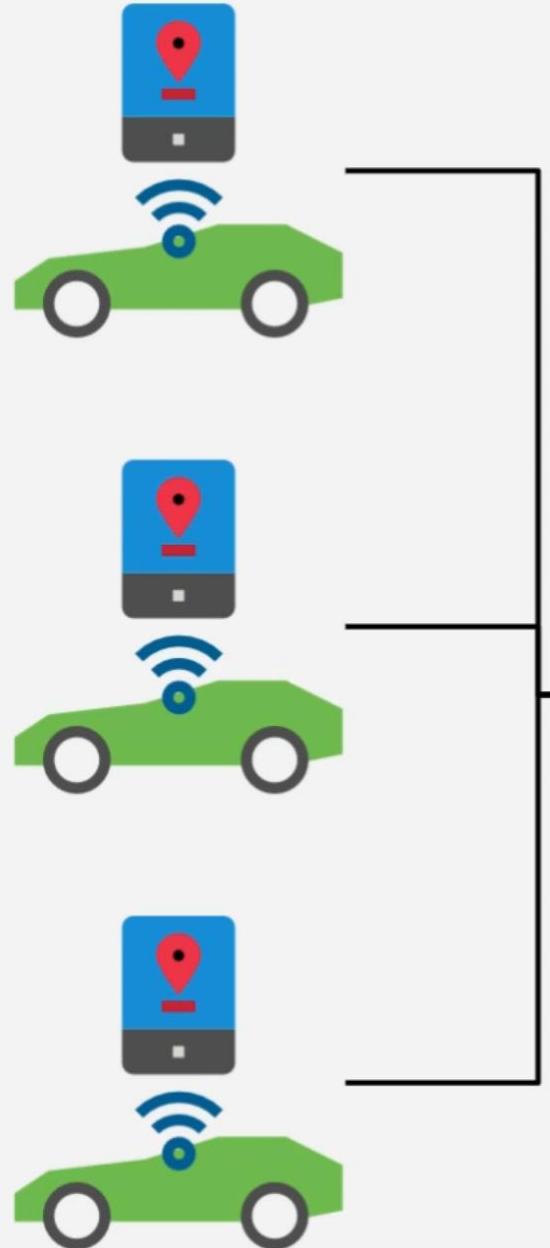
Commodity.java

name : String Example : oil, gold, coffee

price : double Example : 185.89

measurement : String Example : barrel, ounce, tonne

timestamp : long UNIX timestamp



t.location

```
{  
    "car_id" : "CX1580",  
    "timestamp" : 1577750400,  
    "distance" : 74  
}
```



Idempotent Consumer

- ▶ Duplicate message is OK
 - ▶ Outcome of processing message always same even for duplicate messages
 - ▶ ex:- update search engine index
- ▶ duplicate Message is dangerous
 - ▶ Duplicate transaction
 - ▶ ex:- create (duplicate) payment
 - ▶ Filter out duplicate messages

How to duplicate?

- ▶ Use database for permanent unique values
- ▶ Use cache for temporary unique values
- ▶ Cache
 - ▶ Better performance
 - ▶ Automatically remove data after certain time
 - ▶ Example: Redis
- ▶ Database
 - ▶ Might publish duplicate after longer period
 - ▶ Virtually unlimited storage



id : primary key database
(surrogate)

po number : natural key
(business)

Kafka record key using ID



Partition 0

Offset 0 :
5551

Offset 1 :
5552



Partition 1

Offset 0 :
5553



5551, 5553, 5552
5553, 5551, 5552

Event	PR Event ID	PR number
Budget reserve	5551	PR-One
Approval workflow	5552	PR-One
Push notification	5553	PR-One

Kafka record key using PR Number



Partition 0

Offset 0 :
PR-One
5551

Offset 1 :
PR-One
5552

Offset 2:
PR-One
5553



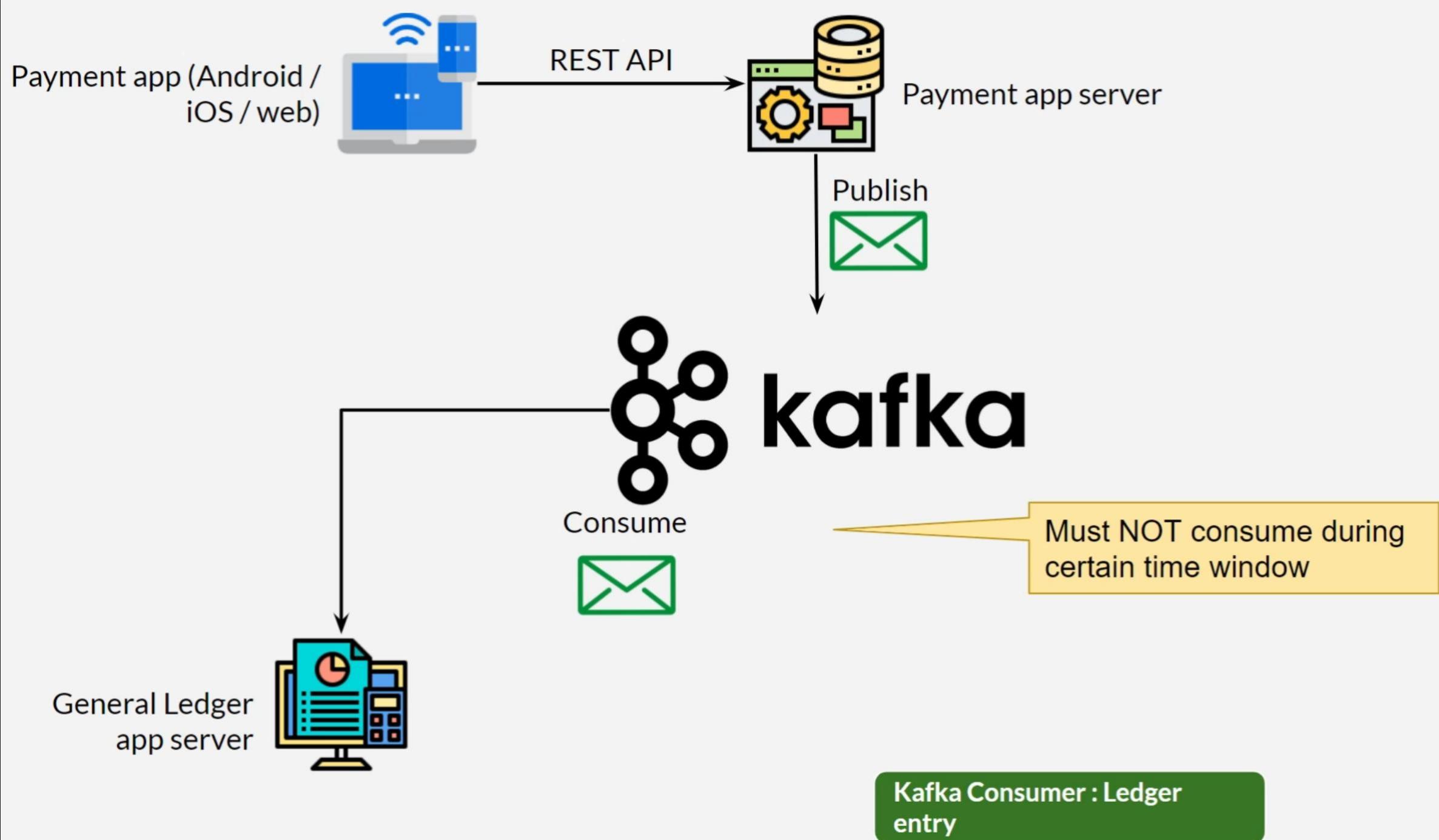
Partition 1

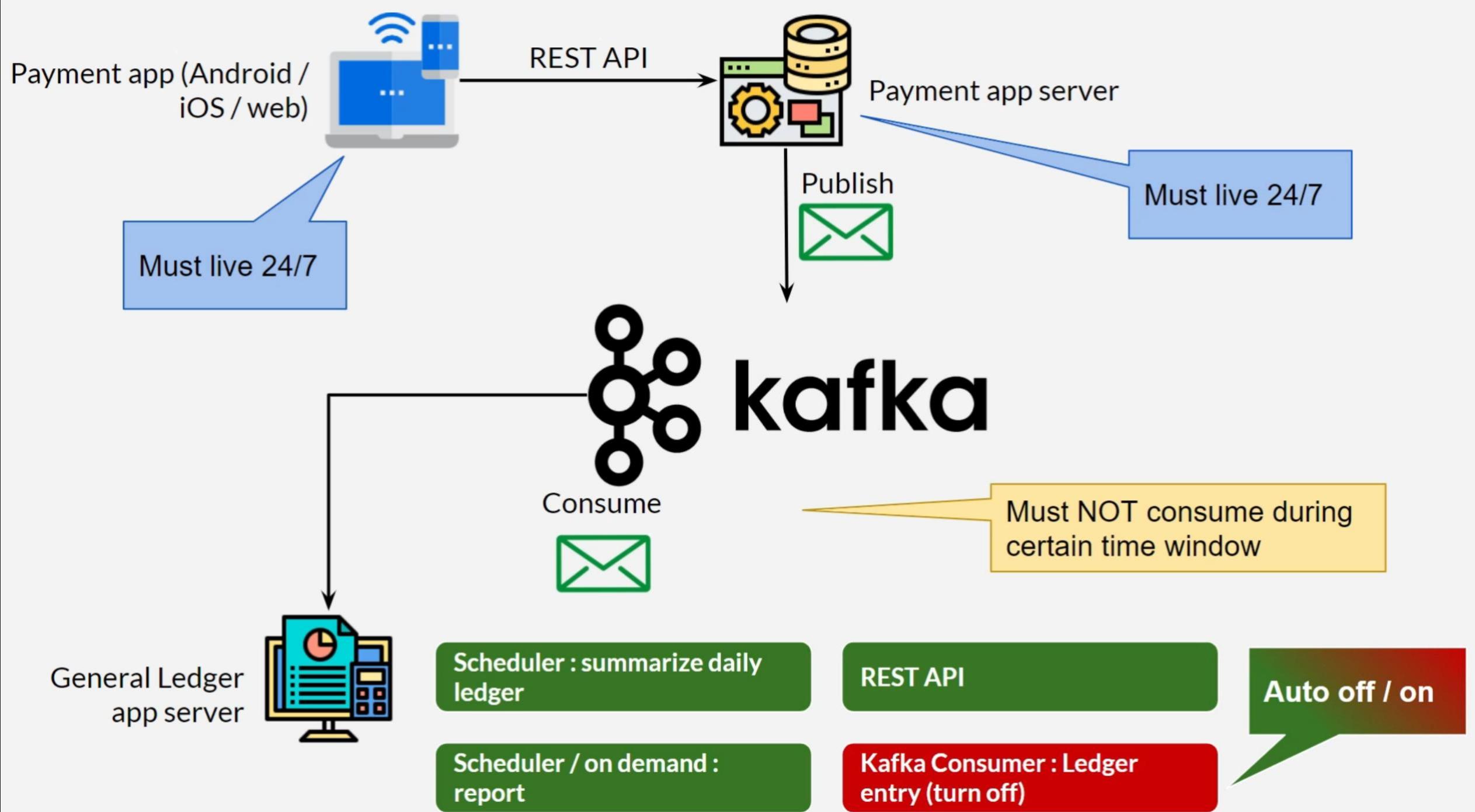


“

Scheduling Consumer

”

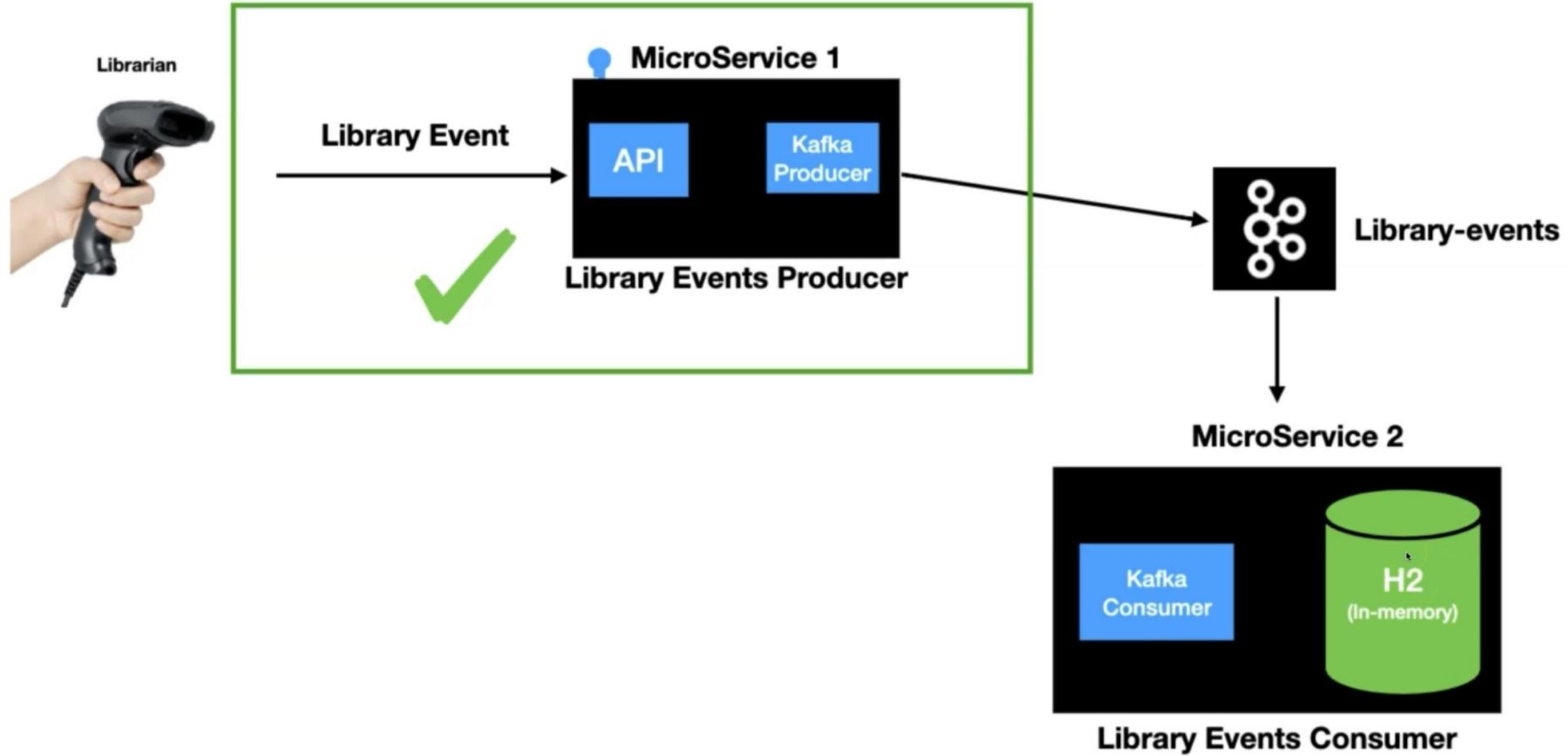




“

KAFKA PROJECT

”





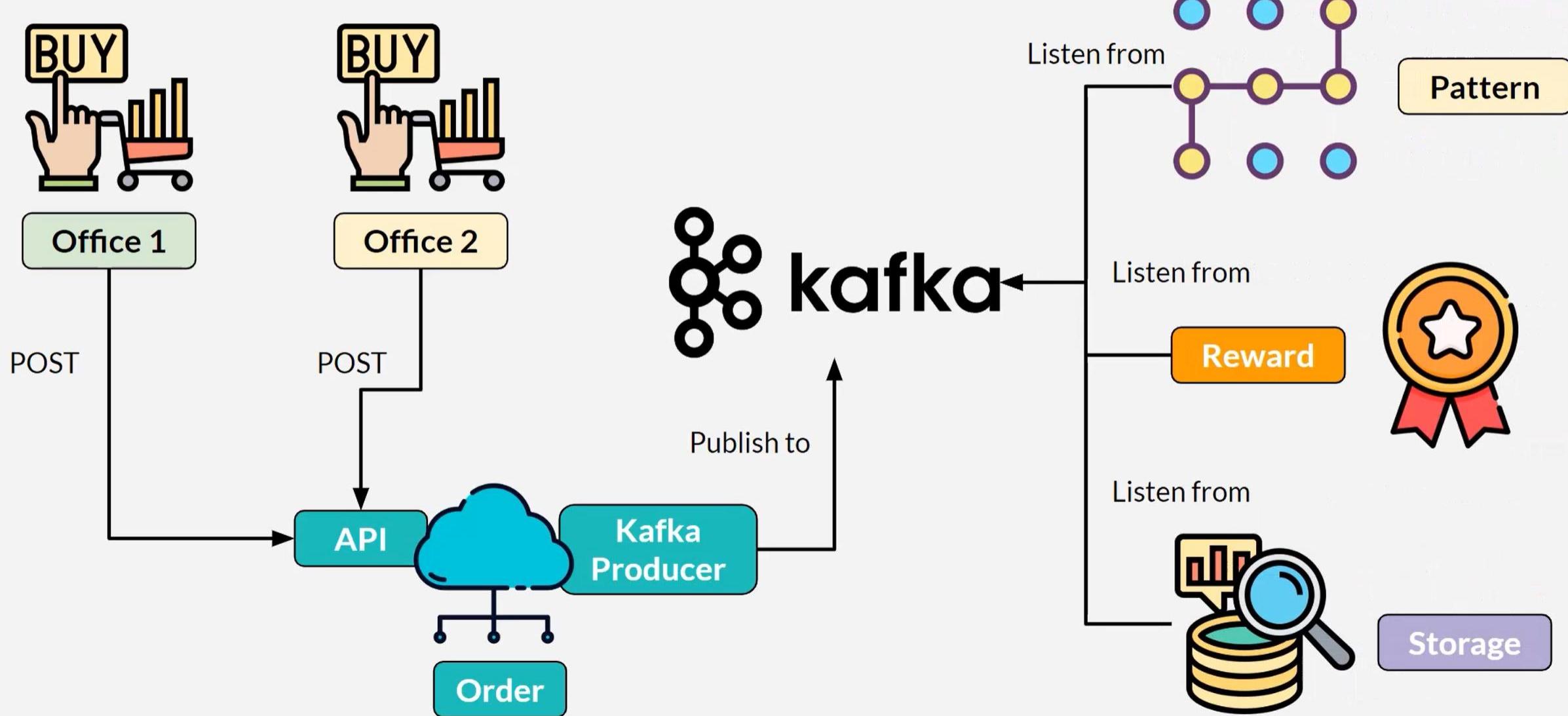
Application Overview

- ▶ Oversimplified application
- ▶ Kafka usage on real life
- ▶ Microservice architecture & pattern

Use Case

- ▶ Commodity trading company with multiple branches
- ▶ Branch Submit purchase order to head office
- ▶ After Process
 - ▶ Pattern Analysis
 - ▶ Reward
 - ▶ Bigdata storage/Persist data
- ▶ Speed is the key
- ▶ Branch office submit each order once
- ▶ Head Office process using kafka

Application Architecture



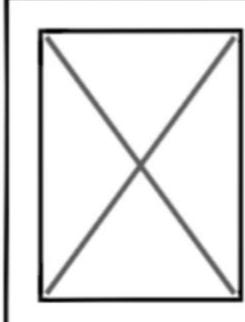


Order location

Credit Card Number

1958 2850 6094 3758

Items

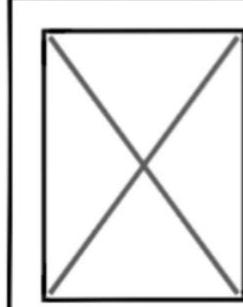


A beautiful book

\$14



Quantity

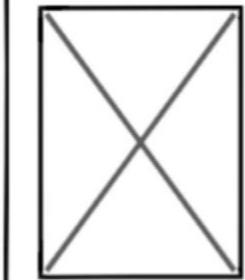
 

An exotic fruit

\$3



Quantity

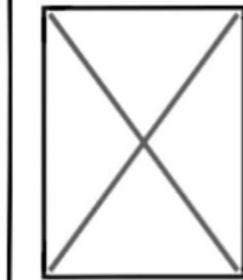
 

A mysterious box

\$14



Quantity

A luxury dress

\$26



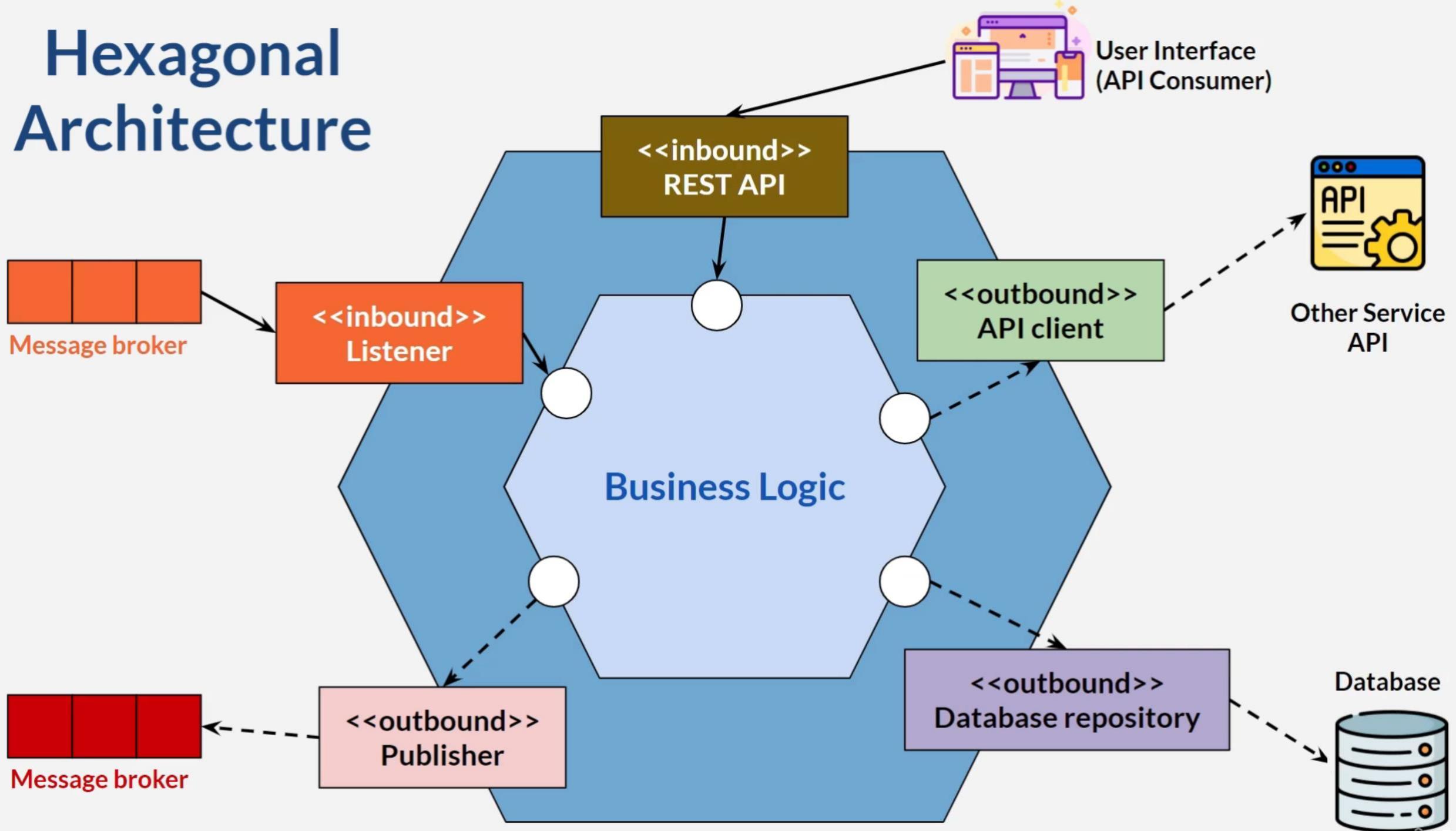
Quantity

Submit Order

The need of organizing code

- ▶ Many codebase
- ▶ Applications
- ▶ Easy to work with
- ▶ Easy to future changes
- ▶ Easy transfer knowledge & employee onboarding
- ▶ Multiple coding & source code organizations
- ▶ patterns for code structure & organization

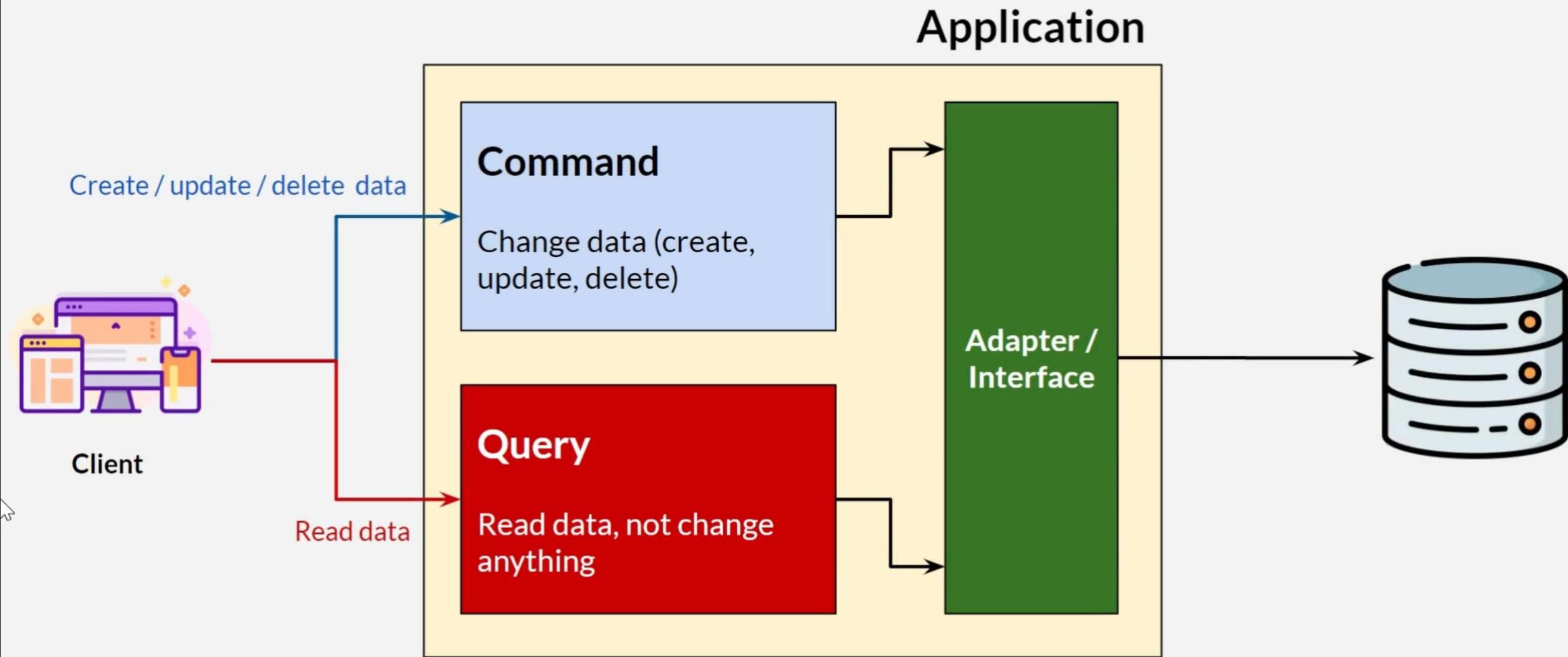
Hexagonal Architecture



Hexagonal Architecture

- ▶ Benefit: decouple business logic with data access
- ▶ Easier to test or change
- ▶ Communication using multiple adapters

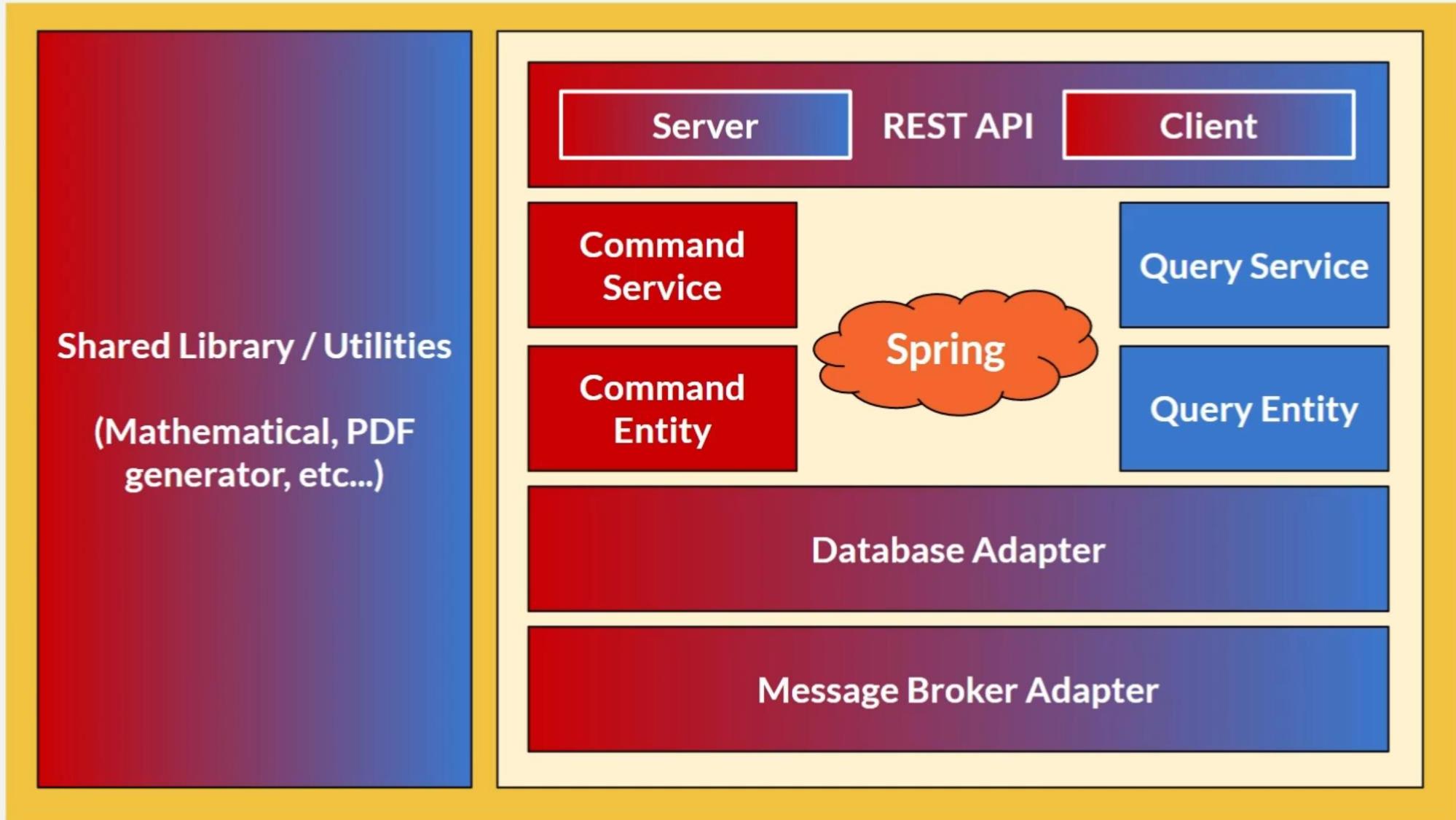
Command Query Separation (CQS)



Command Query Separation

- ▶ Command / transaction / modifier / mutator
- ▶ Query / view
- ▶ Separate, don't mix
- ▶ Easier maintenance and change

Application source code



Generate Projects

Use Spring initializr / IDE

- ▶ Create 4 projects
 - ▶ Group: com.virtusa.kafka
 - ▶ Artifact:
 - ▶ kafka-ms-order
 - ▶ kafka-ms-pattern
 - ▶ kafka-ms-reward
 - ▶ kafka-ms-storage
 - ▶ Package name: com.virtusa.kafka (remove any suffix)
- ▶ Spring boot 2.X + Java 11/17

Dependencies

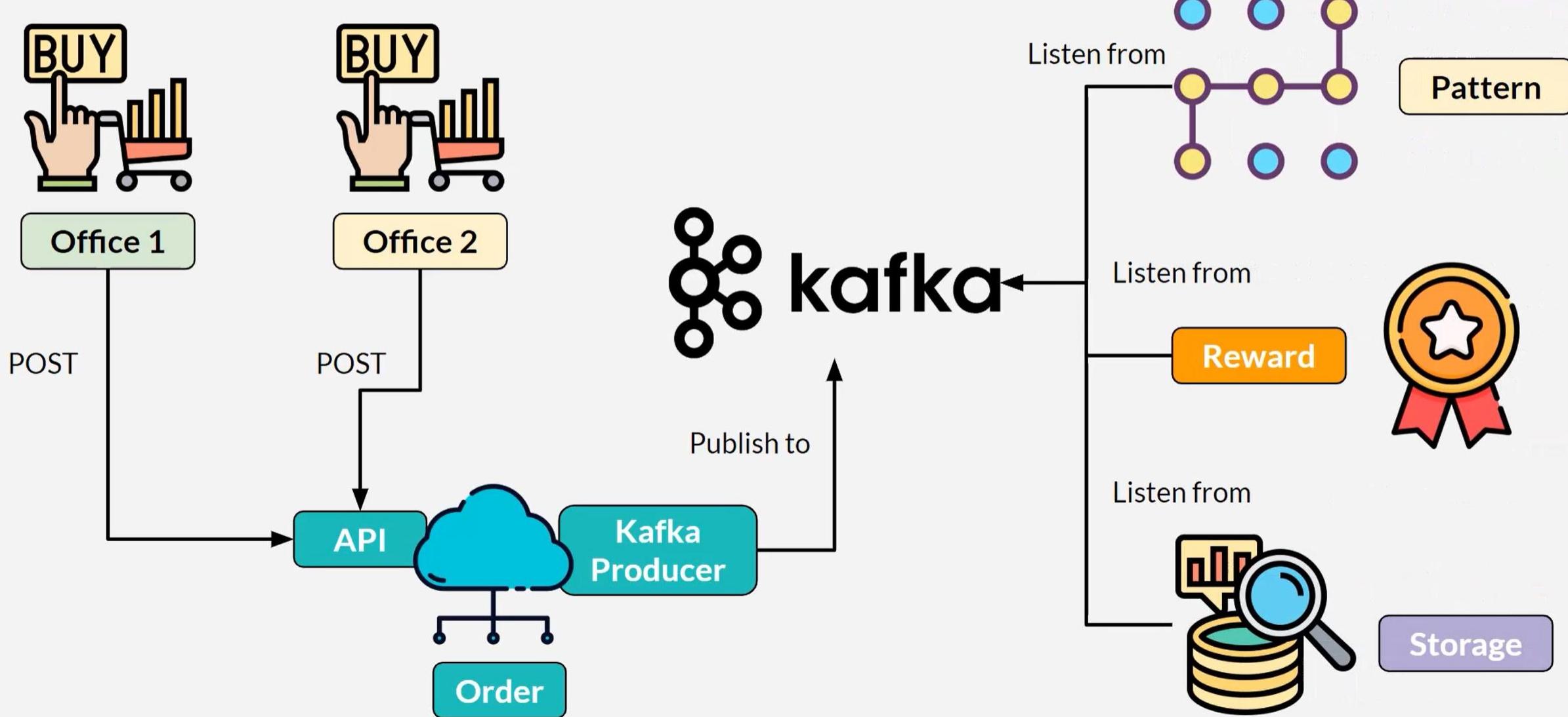
Kafka-ms-order

- ▶ Web
- ▶ Spring Kafka
- ▶ Dev tools
- ▶ JPA
- ▶ H2 database

Other 3 projects

- ▶ Spring Kafka
- ▶ Dev tools

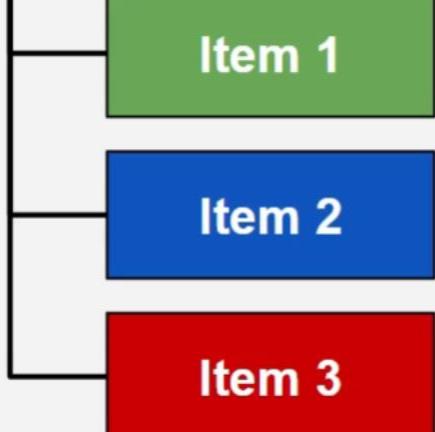
Application Architecture



Order - Kafka Message

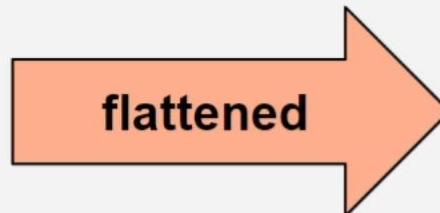


- ★ order number
- ★ location
- ★ order date & time
- ★ credit card number



Each item has

- name
- price
- quantity

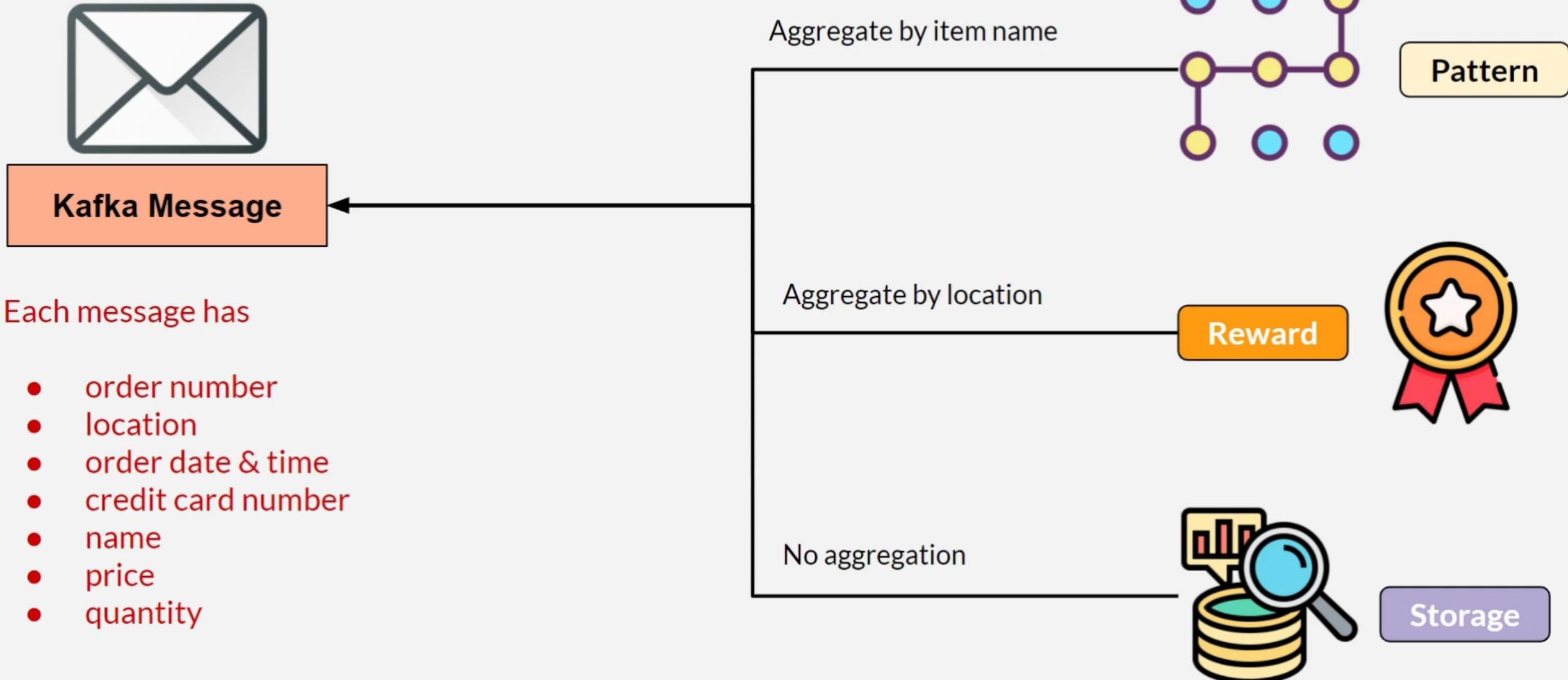


Kafka Message

Each message has

- order number
- location
- order date & time
- credit card number
- name
- price
- quantity

Order - Kafka Message

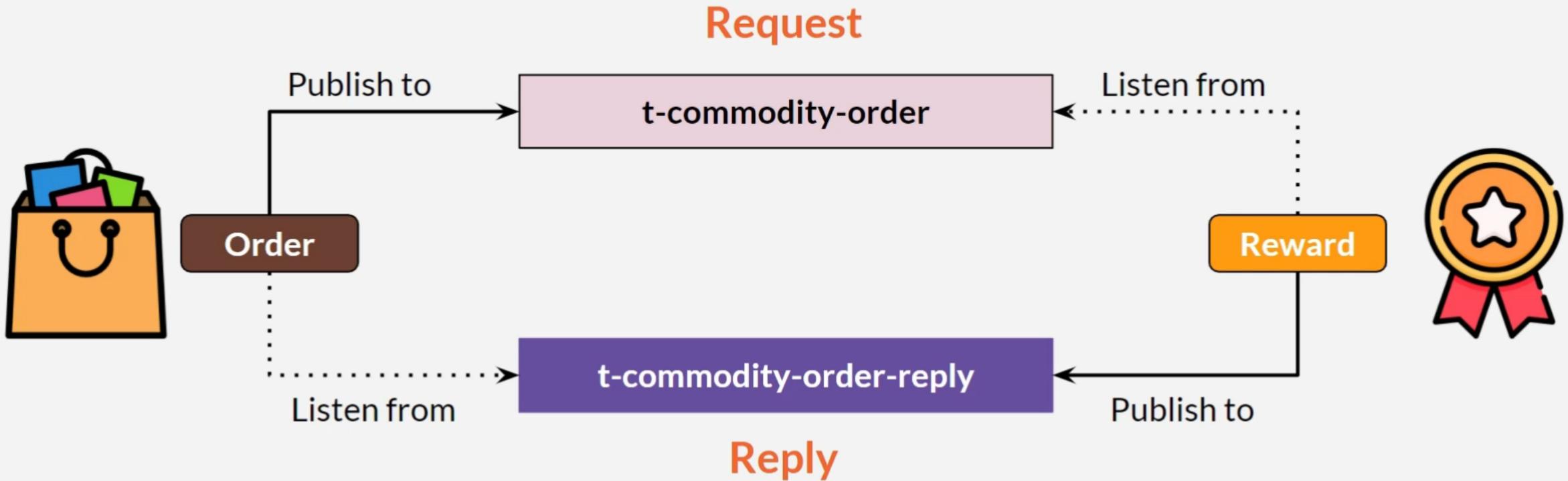


Order Team Says



- Reward team, please confirm order received
- We provide kafka topic for confirmation

Asynchronous Request / Reply

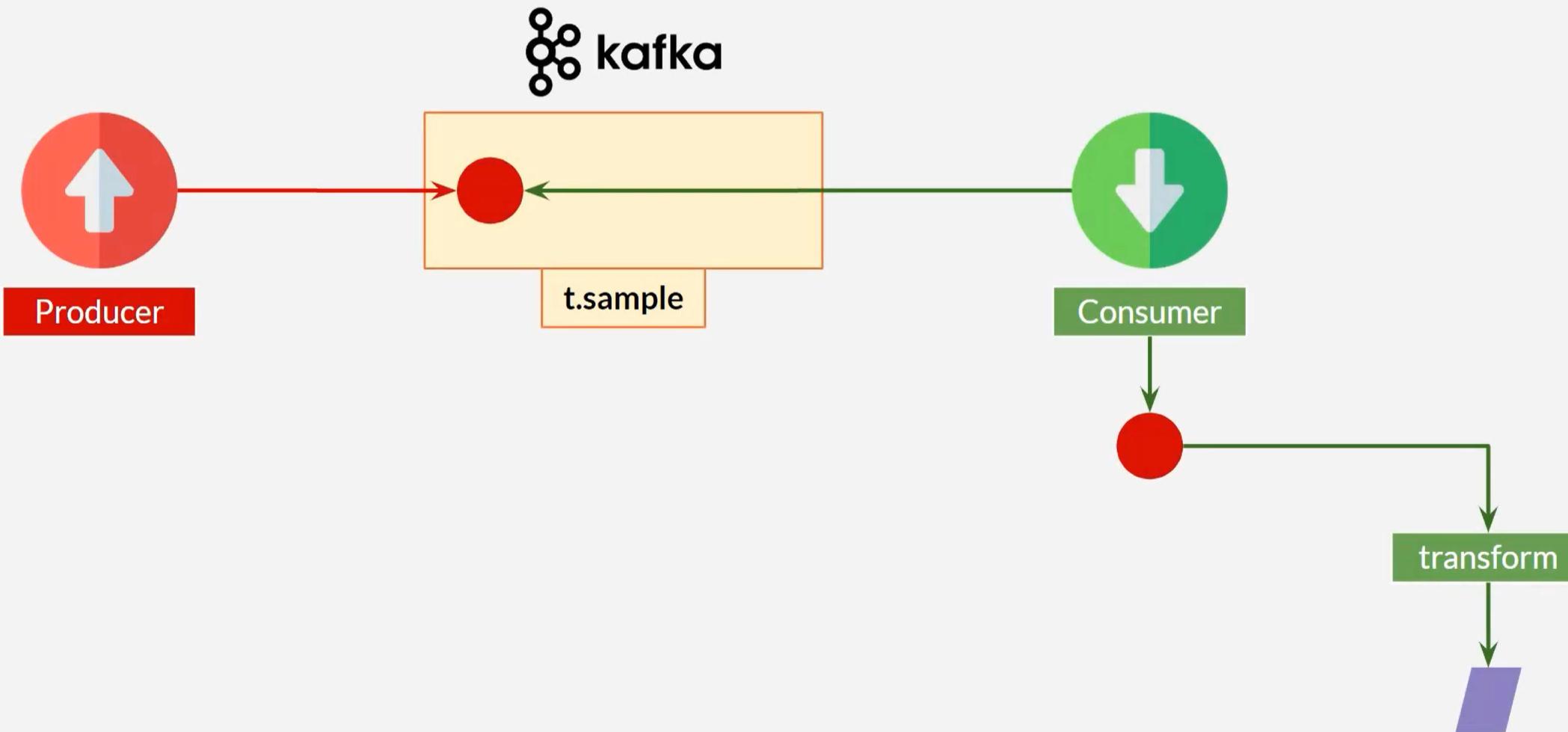


Asynchronous Request / Reply

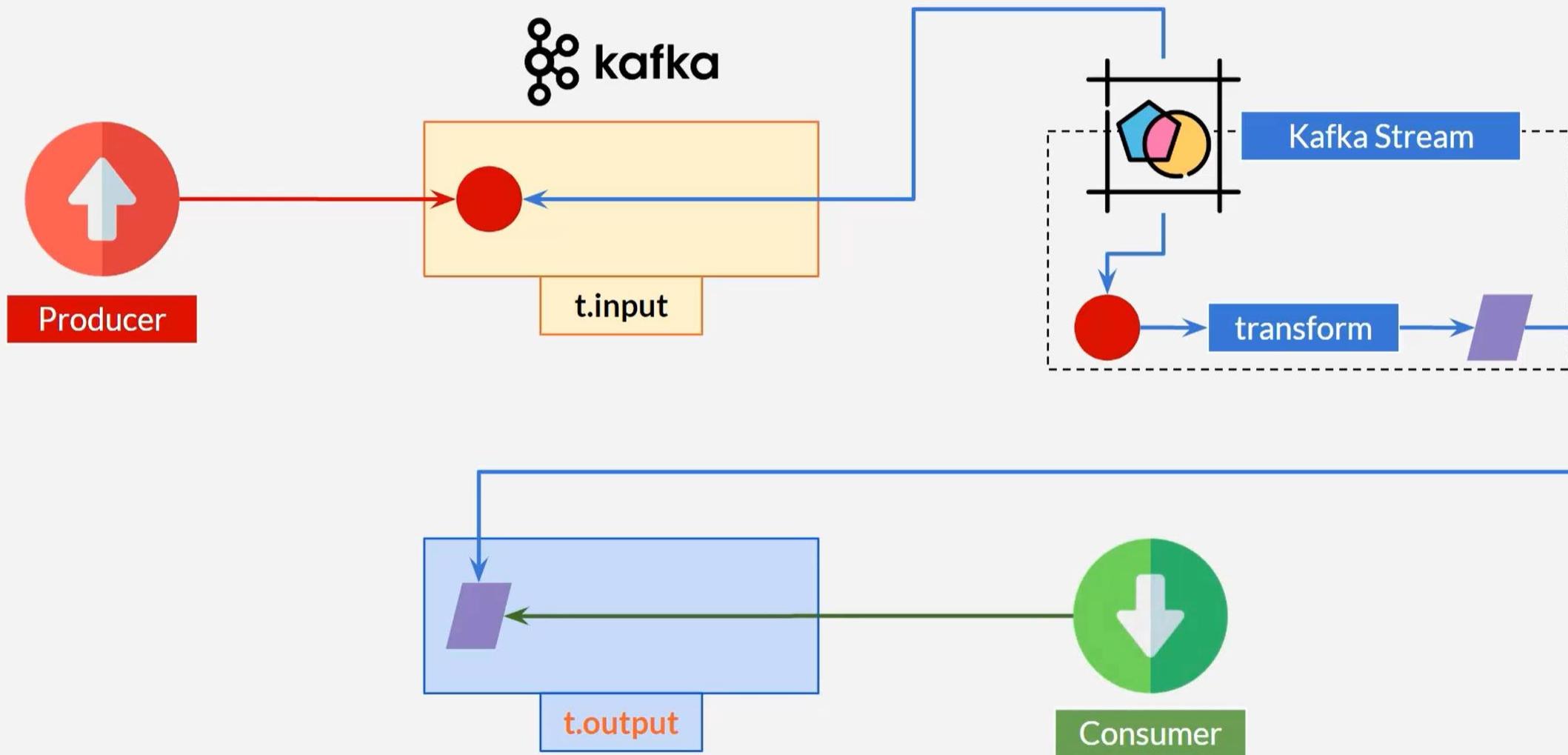
```
public class ConsumerOnReward {  
  
    @KafkaListener(topics = "t-commodity-order")  
    @SendTo("t-commodity-order-reply")  
    public OrderReplyMessage listen(OrderMessage requestMessage) {  
        // do some process for request  
        // ...  
  
        OrderReplyMessage replyMessage = new OrderReplyMessage(...);  
        return replyMessage;  
    }  
}
```

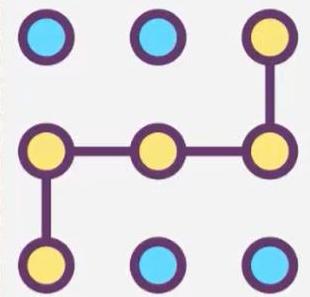


Data Transformation



Kafka Stream





Pattern

Spring Kafka

Kafka Stream

Reward



Spring Kafka

Kafka Stream



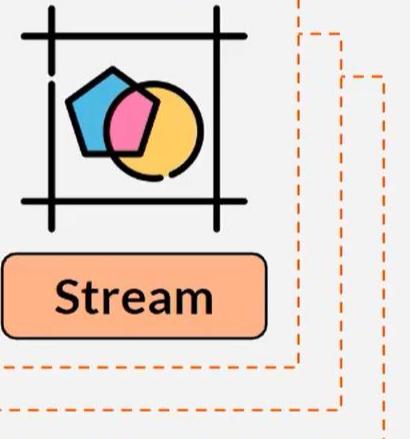
Storage

Spring Kafka

Kafka Stream

Kafka Stream

Spring Kafka

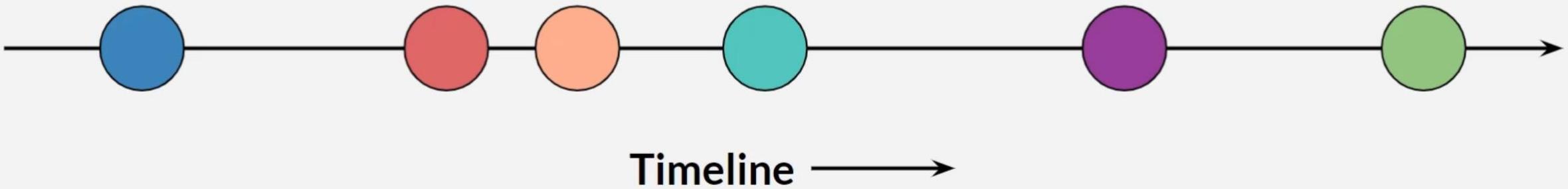


Stream

Kafka Stream

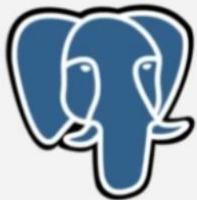
- ▶ Stream processing framework
- ▶ Released on 2017
- ▶ Alternative for Apache Spark, NIFI or Flink
- ▶ Stream & stream processing?

Data Stream

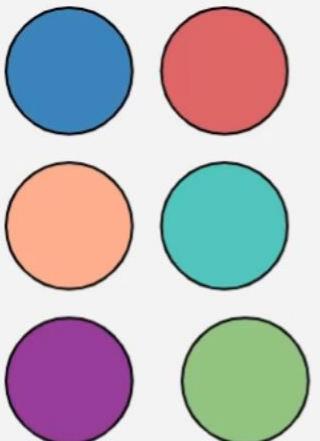


Each circle represents a data
Endless
Data (event) is immutable
Can be replayed

Data Processing



PostgreSQL

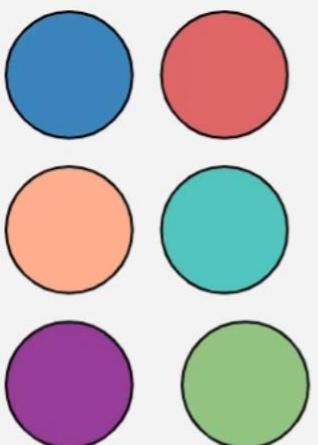


elasticsearch

ORACLE
DATABASE



Data Processing



Transformation

Aggregation

Filter

Calculation

Combination

etc

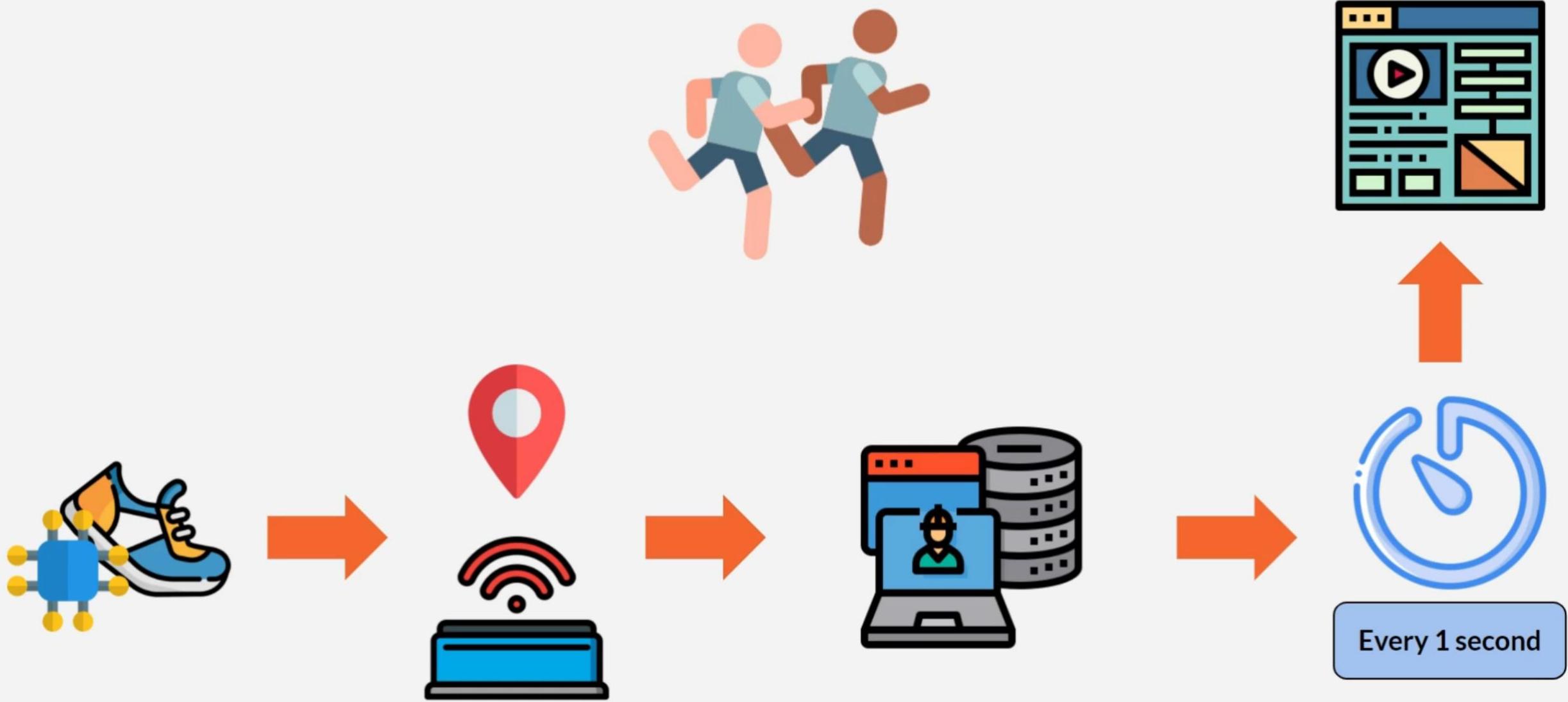
Everyday at
23:00

Every 30 min

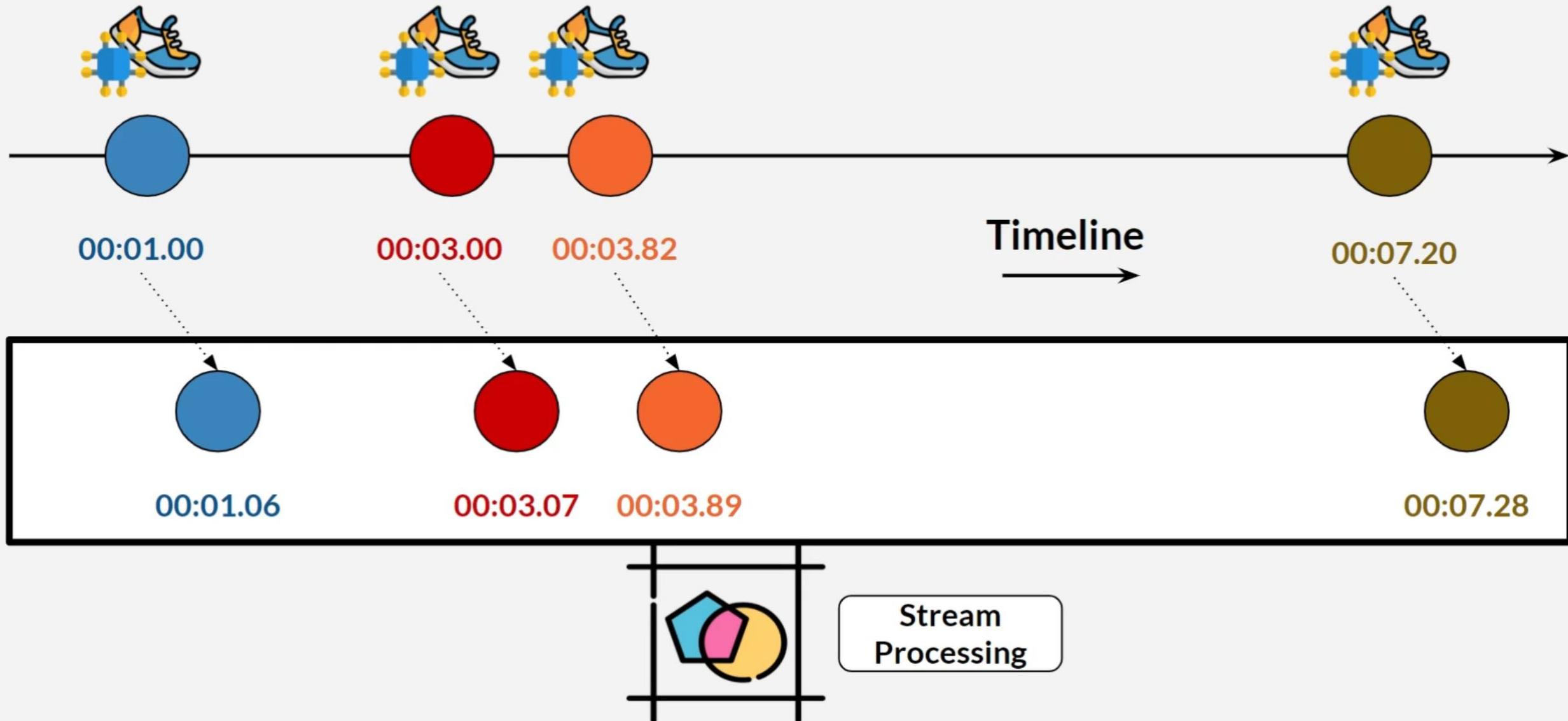


Every 1 second

Micro Batching



Stream Processing



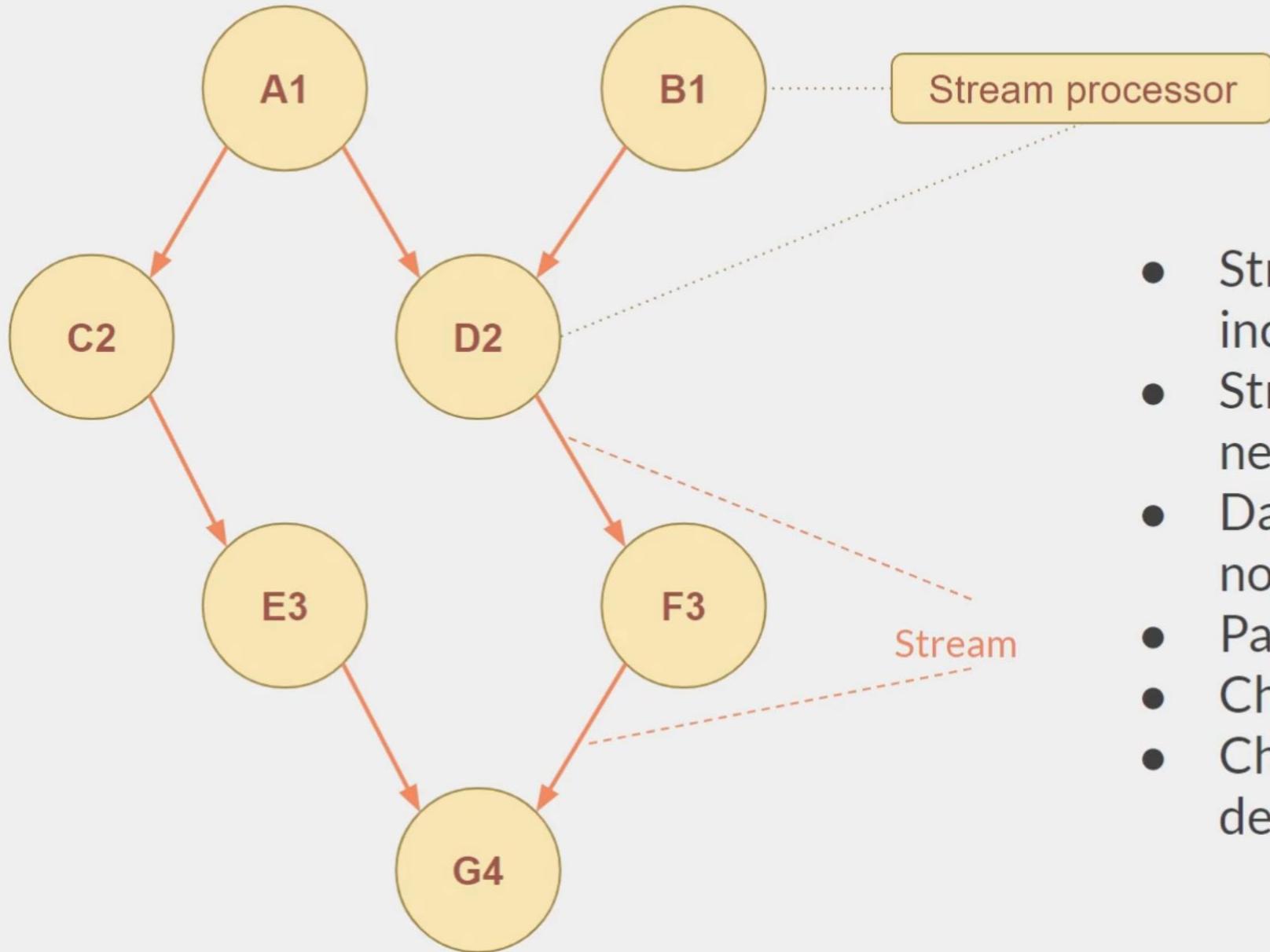
Yes To Stream Processing

- ▶ Yes, When:
 - (relatively) fast data flow
 - Application need to response quick to most recent data
- ▶ Example
 - Marathon
 - Credit card fraud
 - Stock trading
 - Log analysis

No To Stream Processing

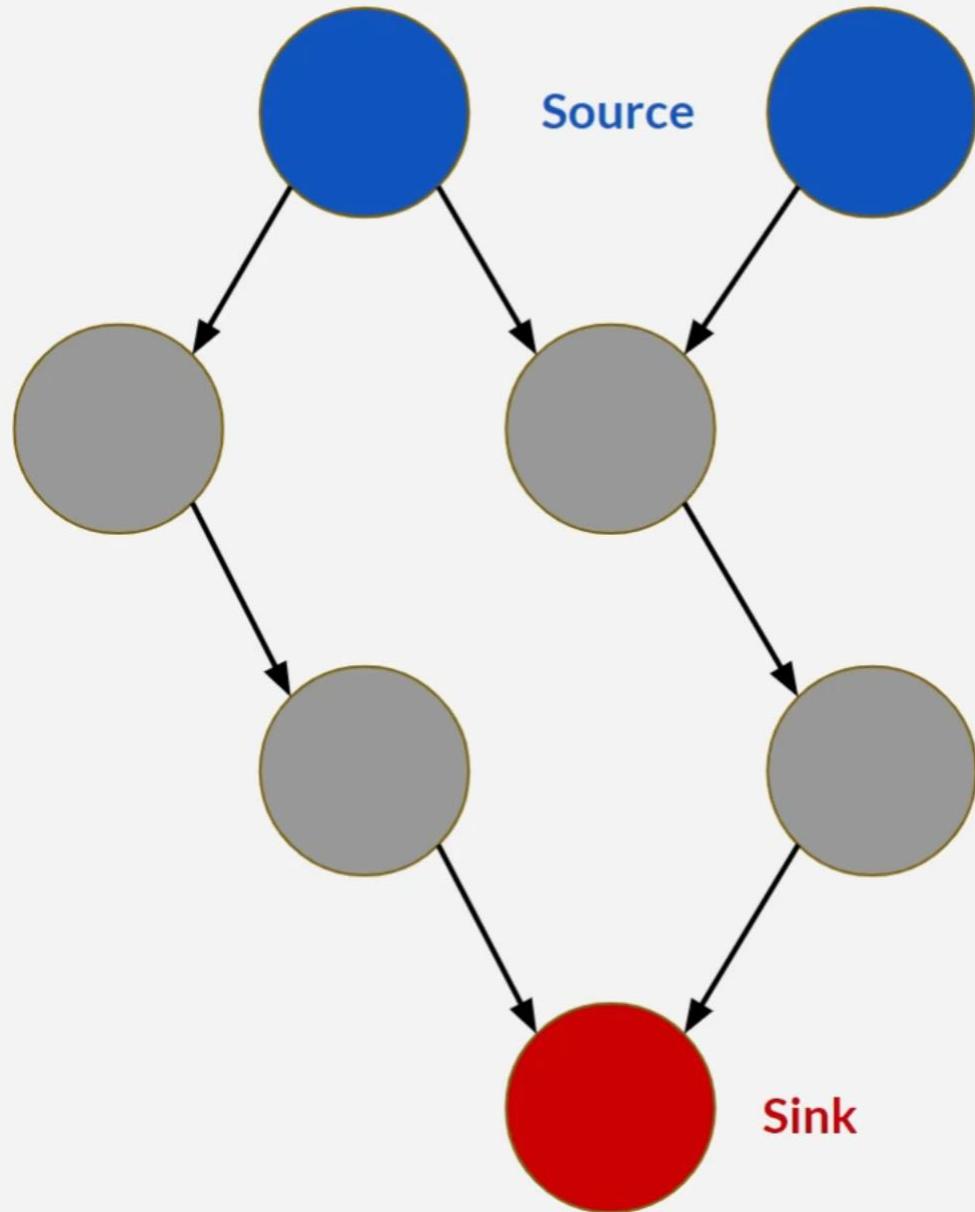
- ▶ Example
 - Daily Interest
 - Forecasting

Topology / DAG



- Stream processor process incoming data stream
- Stream processor can create new output stream
- Data flows from parent to child, not vice versa
- Parent = upstream
- Child = downstream
- Child stream processor can define another child(ren)

Kafka Stream Topology



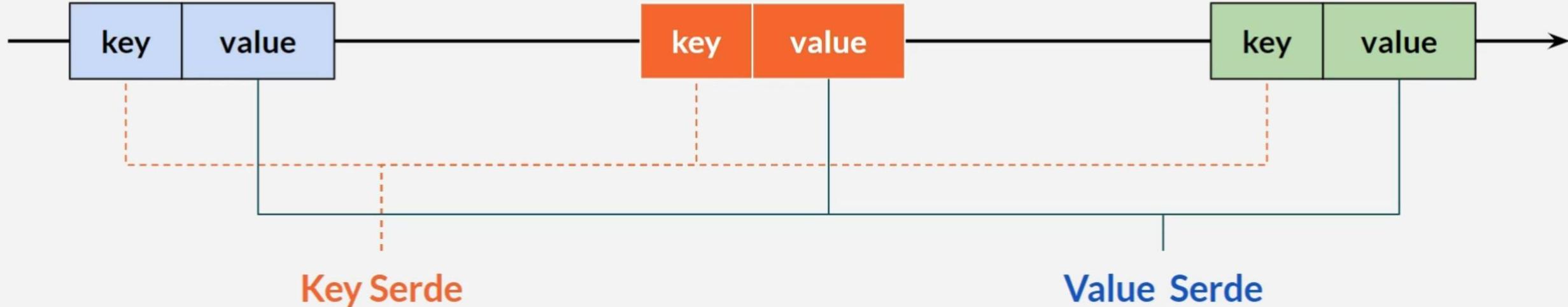
Source Processor

- Does not have upstream
- Consumes from one or more kafka topics
- Forwarding data to downstream

Sink Processor

- Does not have downstream
- Receive data from upstream
- Send data to specified kafka topic

Serde (Serializer / Deserializer)

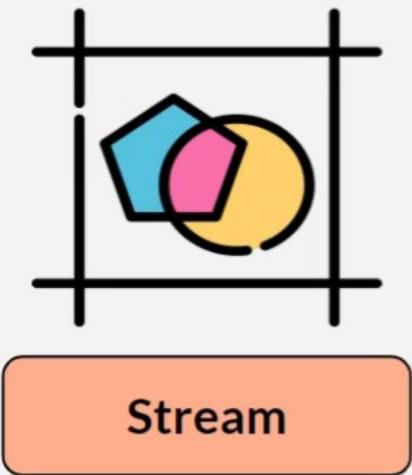


```
Serdes.String()  
Serdes.Long()  
Serdes.ByteArray()  
...
```

```
new JSONSerde<T>()
```

```
class CsvSerde<T> implements Serde<T>
```

What We Will Have

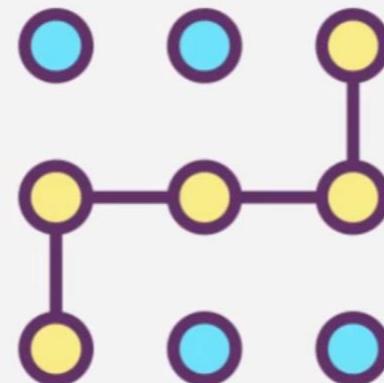


Order



Storage

Reward



Pattern

Spring Initializr

- × Generate 1 java / gradle project from start.spring.io
 - × Group: **com.course.kafka**
 - × Artifact: **kafka-stream-sample**
 - × Package name: **com.course.kafka** (remove any suffix)
 - × Dependency: **Spring for Apache Kafka, Spring for Apache Kafka Streams, Spring Boot Devtools**
- × Spring boot 2.x

Stream & Table

KStream

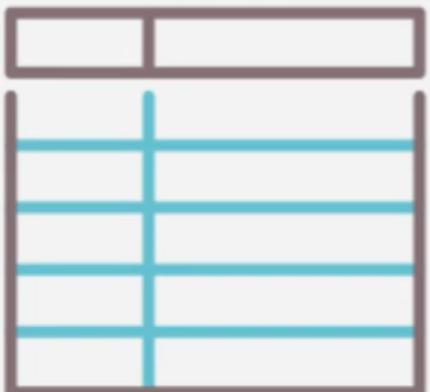
Ordered sequence of messages

Unbounded

Inserts data



Table



KTable

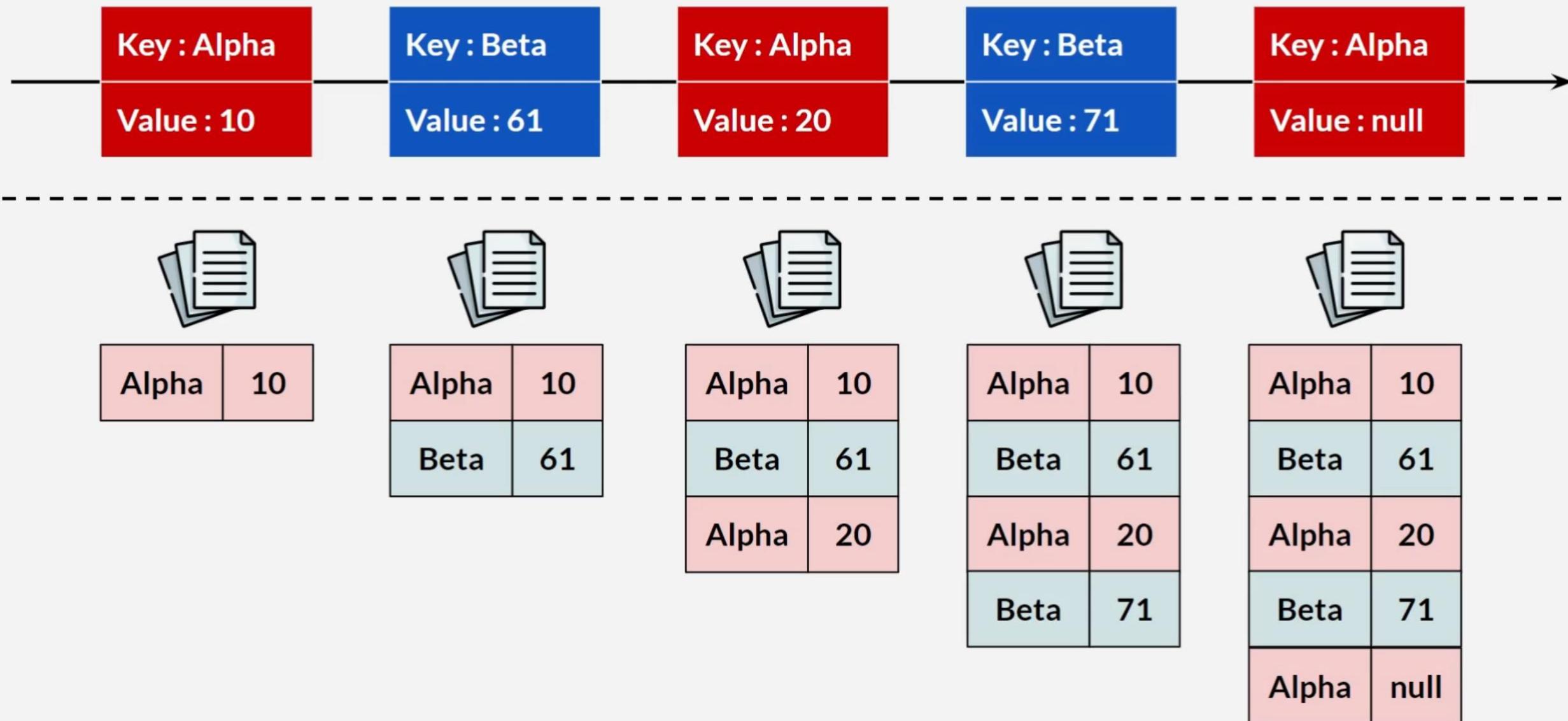
Unbounded

Upserts data : insert or update based on key

Delete on null value

Analogy : database table

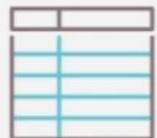
KStream : Inserts Data



KTable : Upserts / Delete Data

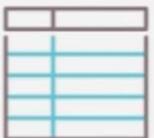


Insert Alpha



Alpha	10
-------	----

Insert Beta



Alpha	10
Beta	61

Update Alpha



Alpha	20
Beta	61

Update Beta



Alpha	20
Beta	71

Delete Alpha



Beta	71
------	----

When to Use KStream / KTable?

- ▶ KStream
 - ▶ Topic not log-compacted
 - ▶ Data is partial information
- ▶ KTable
 - ▶ Topic is log-compacted
 - ▶ Data is self sufficient

Log Compaction

- ▶ Kafka admin process
- ▶ Keep at least latest value & delete the older
- ▶ Based on record key
- ▶ Useful if we need latest snapshot
- ▶ Configure when creating topic

Log Compaction

Topic : T, partition 0

Offset	0	1	2	3	4	5	6	7	8	9	10
Key	Alpha	Sigma	Beta	Alpha	Omega	Alpha	Delta	Delta	Beta	Omega	Omega
Value	10	180	20	11	240	12	40	40	21	241	242

Log compaction

Topic : T, partition 0

Offset	1	5	7	8	10
Key	Sigma	Alpha	Delta	Beta	Omega
Value	180	12	41	21	242

Log Compaction

Topic : T, partition 0

Offset	0	1	2	3	4	5	6	7	8	9	10
Key	Alpha	Sigma	Beta	Alpha	Omega	Alpha	Delta	Delta	Beta	Omega	Omega
Value	10	180	20	11	240	12	40	40	21	241	242

Log compaction

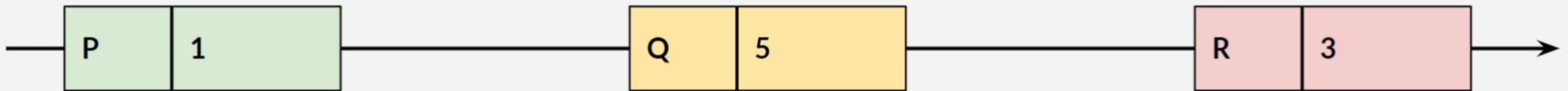
Topic : T, partition 0

Offset	1	5	7	8	9	10
Key	Sigma	Alpha	Delta	Beta	Omega	Omega
Value	180	12	40	21	241	242

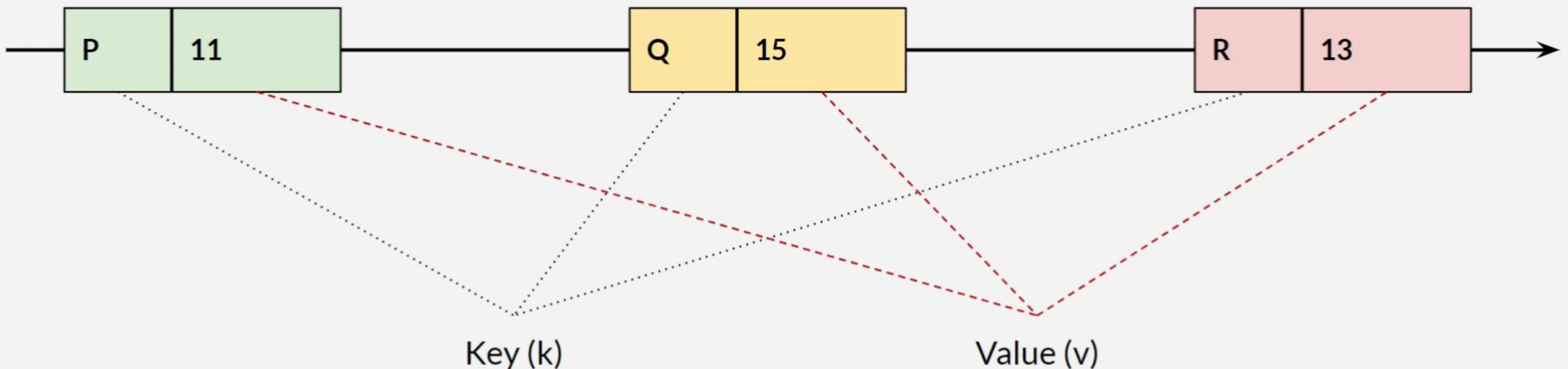
Log Compaction

- ▶ Keep the order
- ▶ Not change offset
- ▶ Not duplication validator
- ▶ Can fail

Diagram



$v = v + 10$



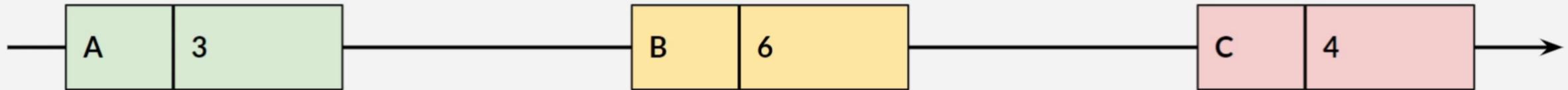
Intermediate & Terminal Operation

- × Intermediate
 - × KStream -> KStream
 - × KTable -> Ktable
- × Terminal
 - × KStream -> void
 - × KTable -> void
 - × “Final” operation

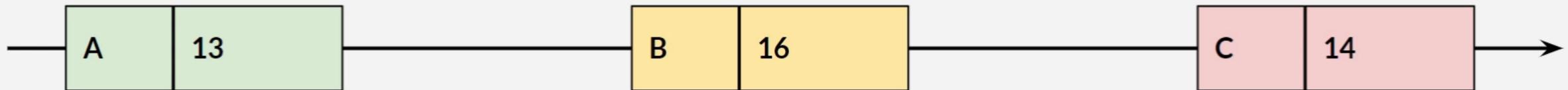
Reminder : Key & Partition

- × Key & partition is related
- × Partition according to key
- × Repartition
 - × From partition A to partition X

mapValues

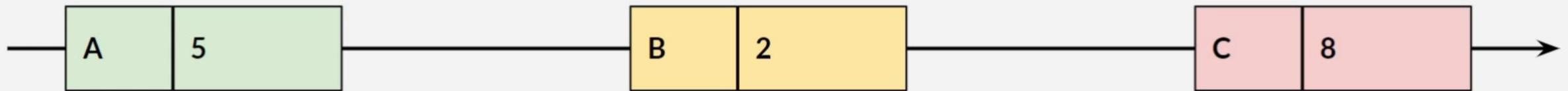


```
stream.mapValues(v -> v + 10)
```



- Takes one record, produces one record
- Does not change key
- Affect only value
- Not trigger repartition
- Intermediate operation
- KStream & Ktable

map

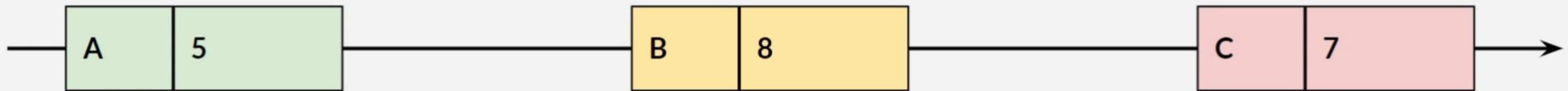


```
stream.map( (k, v) -> KeyValue.pair("X" + k, v * 5) )
```



- Takes one record, produces one record
- Change key
- Change value
- Trigger repartition
- Intermediate operation
- KStream

filter

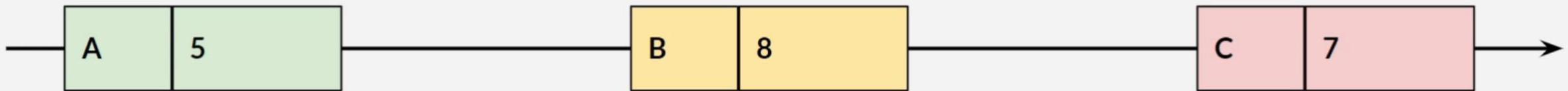


```
stream.filter((k, v) -> v % 2 == 0)
```

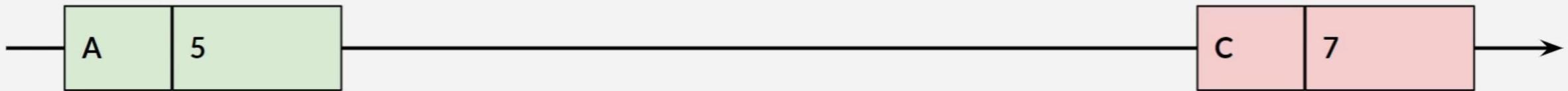


- Takes one record, produces one or zero record
- Produce record that match condition
- Does not change key or value
- Not trigger repartition
- Intermediate operation
- KStream & Ktable

filterNot

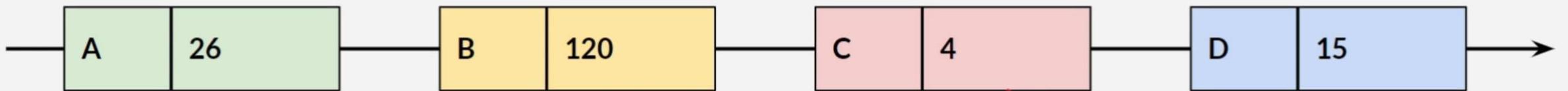


```
stream.filterNot((k, v) -> v % 2 == 0)
```



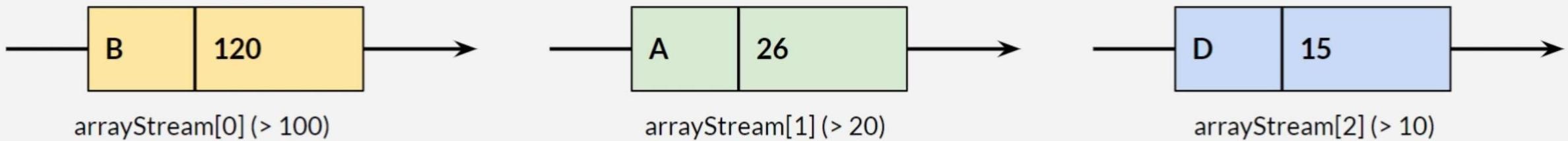
- Takes one record, produces one or zero record
- Produce record that NOT match condition
- Does not change key or value
- Not trigger repartition
- Intermediate operation
- KStream & KTable

branch



```
var arrayStream = stream.branch(  
    (k, v) -> v > 100,  
    (k, v) -> v > 20,  
    (k, v) -> v > 10  
)
```

Dropped, no match



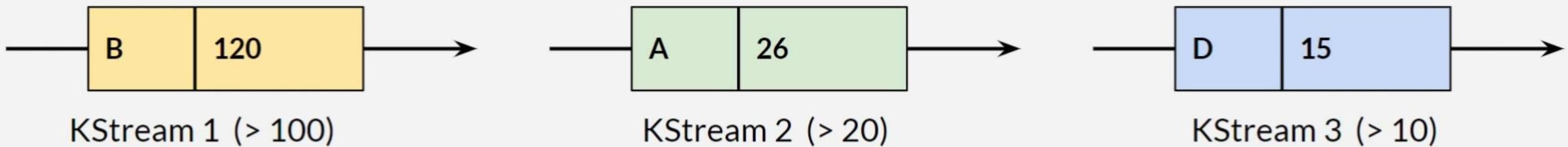
- Split stream based on predicates
- Evaluate predicate in order
- Record only placed once on first match, drop unmatched record
- Returns array of stream
- Intermediate operation
- KStream

split & branch



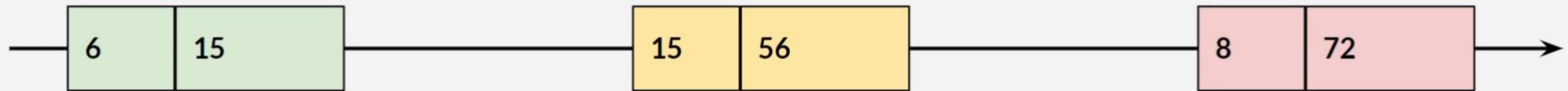
```
stream.split().  
    branch((k, v) -> v > 100, Branched.withConsumer(ks -> ks.to("t-x"))  
    .branch((k, v) -> v > 20, Branched.withConsumer(ks -> ks.to("t-y"))  
    .branch((k, v) -> v > 10, Branched.withConsumer(ks -> ks.to("t-z"))  
)
```

Dropped, no match

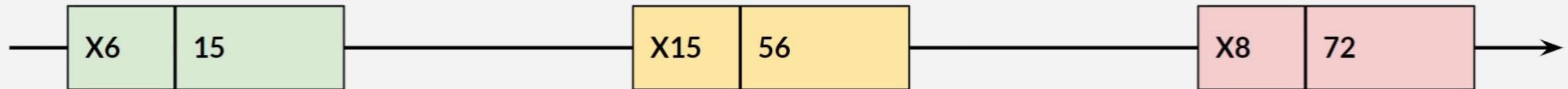


- Split stream based on predicates
- Evaluate predicate in order
- Record only placed once on first match, drop unmatched record
- Get KStream for each branch
- `split()` returns final `BranchedKStream`
- Each branch returns KStream to be processed further

selectKey

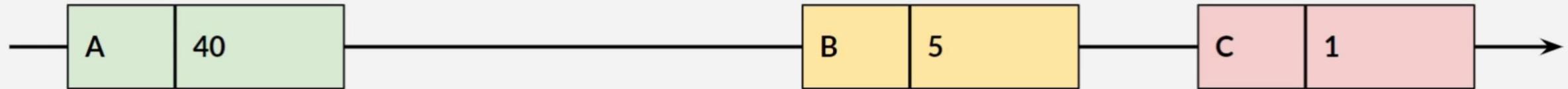


```
stream.selectKey((k, v) -> "X" + k)
```

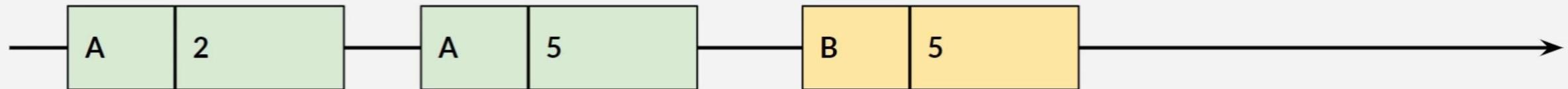


- Takes one record, produces one record
- Set / replace record key
- Possible to change key data type
- Trigger repartitioning
- Value not change
- Intermediate operation
- KStream

flatMapValues



```
stream.flatMapValues(listPrimeFactors())
```

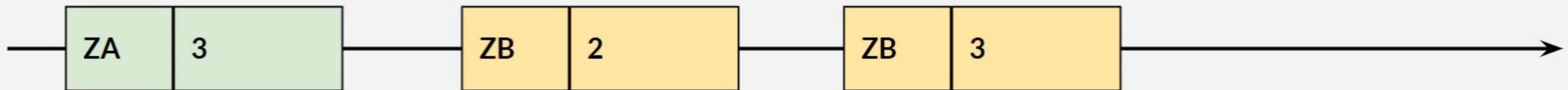


- Takes one record, produces zero or more record
- Does not change key
- Affect only value
- Not trigger repartition
- Intermediate operation
- KStream

flatMap

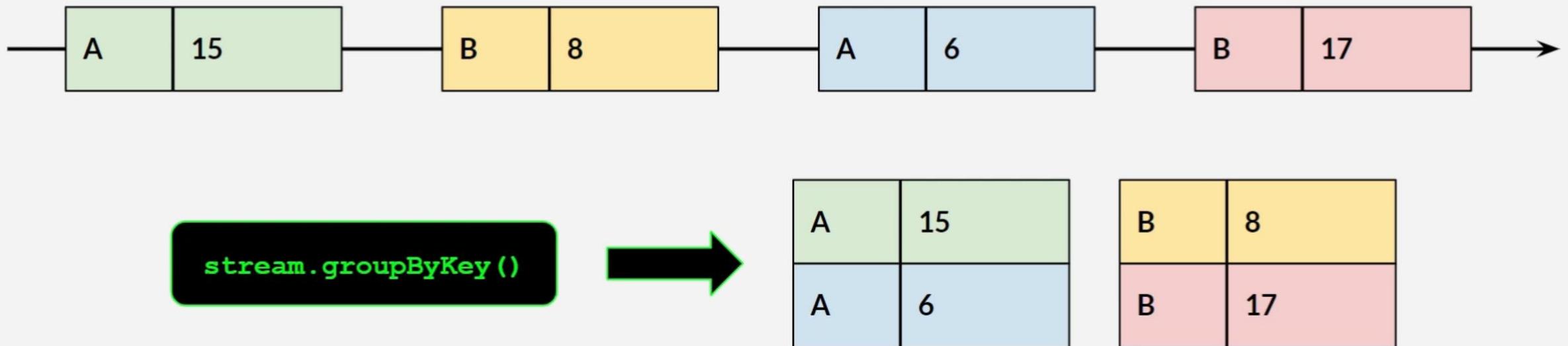


```
stream.flatMap(listPrimeFactorsAndAppendKey())
```



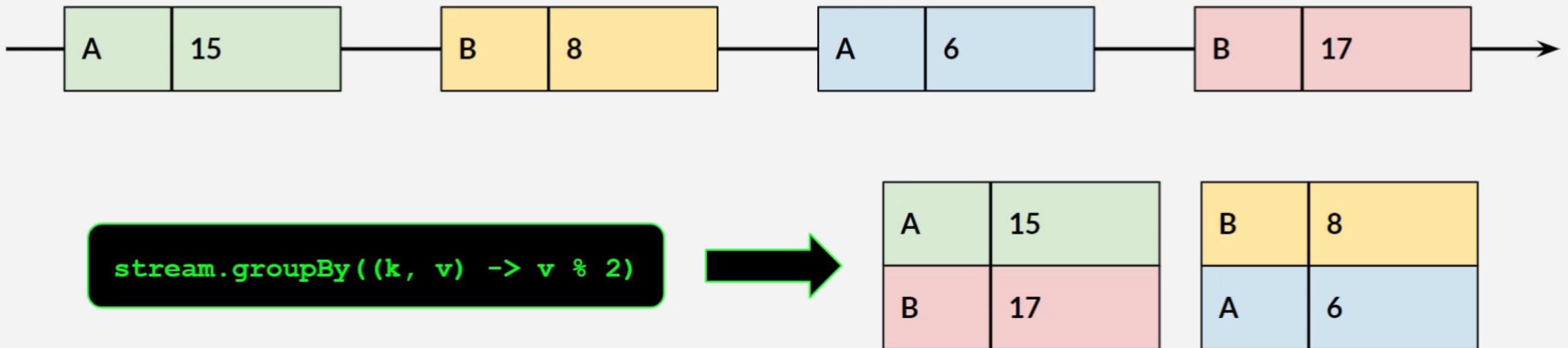
- Takes one record, produces zero or more record
- Change key
- Change value
- Trigger repartition
- Intermediate operation
- KStream

groupByKey



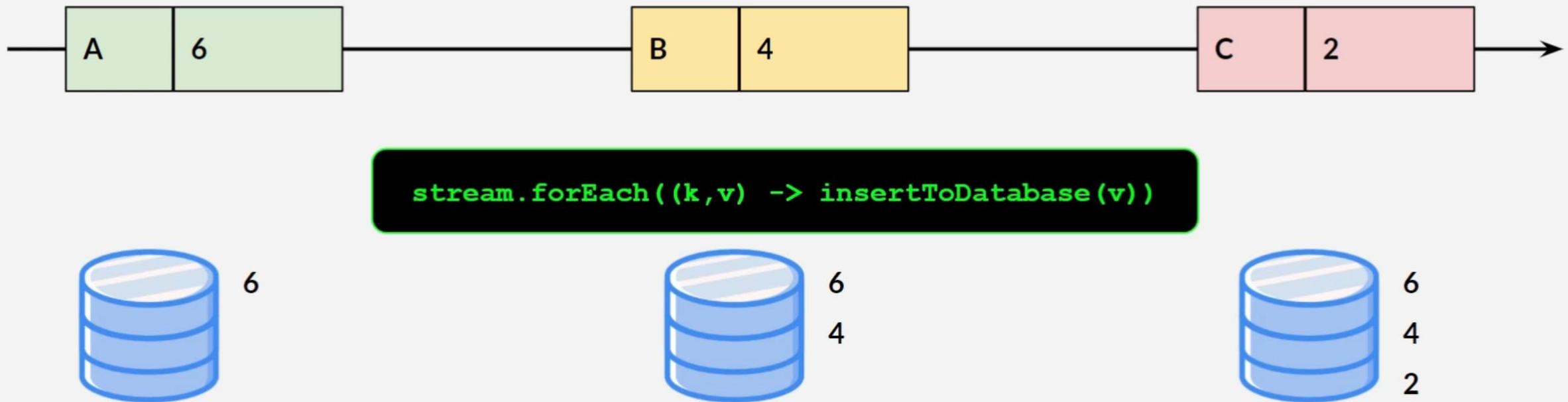
- Intermediate operation
- Group records by existing key
- KStream

groupBy



- Intermediate operation
- Group records by new key
- KStream & KTable

forEach



- Terminal operation
- Takes one record, produces none
- Produces side effect
- Side effect not tracked by kafka
- KStream & KTable

peek



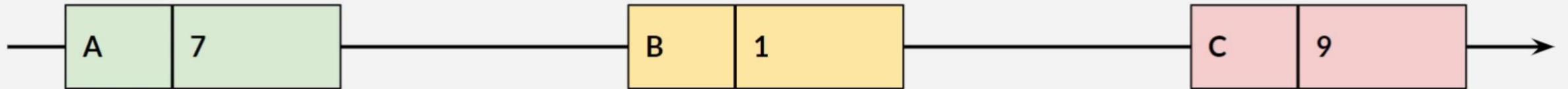
```
stream.peek((k,v) -> insertToDatabase(v)).[nextProcessor]
```

..... Next processor



- Produces unchanged stream
- Produces side effect
- Side effect not tracked by kafka
- Result stream can be processed further
- Intermediate operation
- KStream

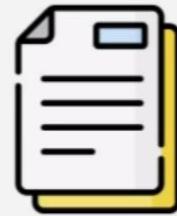
print



```
stream.print(Printed.toSysout())
```

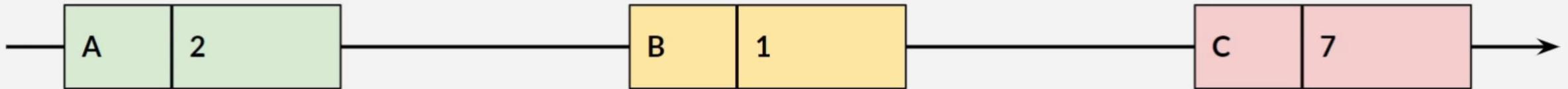


7

7
17
1
9

- Terminal operation
- Print each record
- Something like kafka console consumer
- Print to file or console
- KStream

to

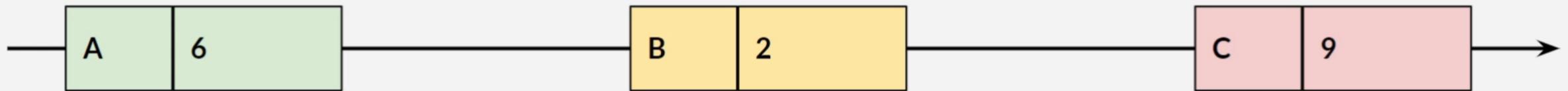


```
stream.to("output-topic")
```



- Terminal operation
- Write stream to destination topic
- KStream

through



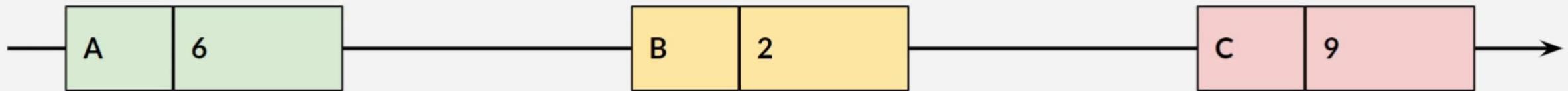
```
stream.through("output-topic").[nextProcessor]
```

..... Next processor



- Intermediate operation
- Write stream to destination topic
- Continue record processing
- KStream

repartition



```
stream.repartition().nextProcessor()  
stream.repartition(Repartitioned.as("output-topic")).nextProcessor()
```

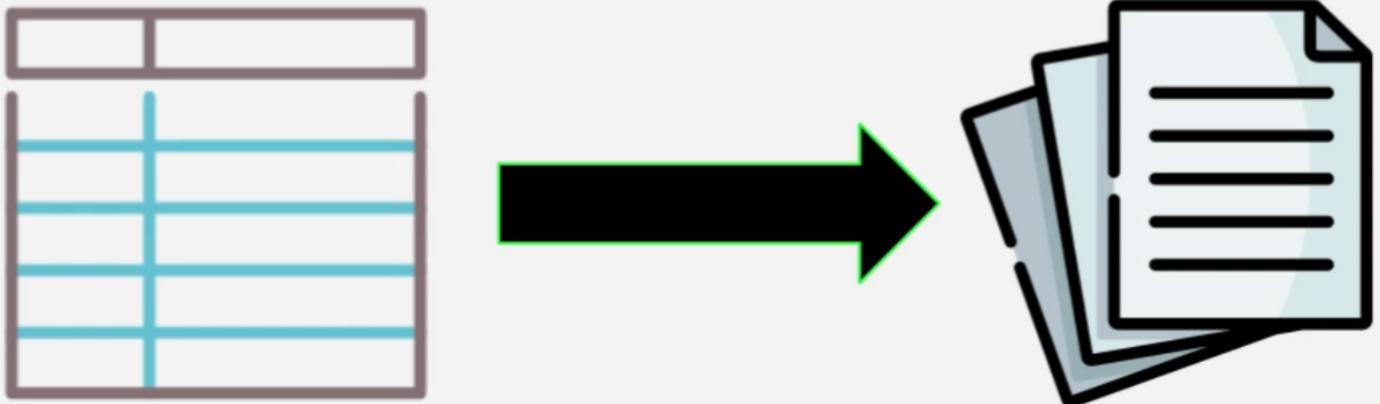
..... Next processor



- Intermediate operation
- Write stream to destination topic
- Continue record processing
- repartition() output-topic name is fixed
- Topic only for kafka internal use

toStream

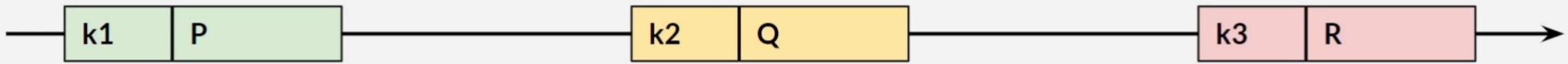
```
table.toStream()
```



- KTable
- Intermediate operation
- Convert KTable to KStream

merge

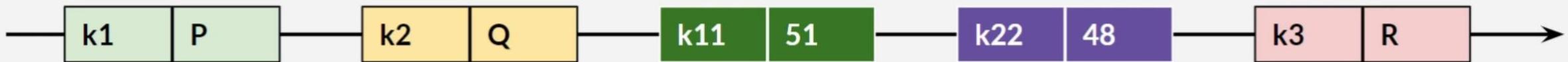
alphabetStream



numericStream



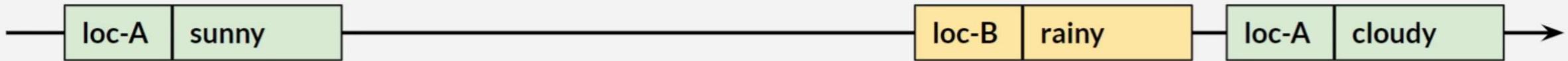
```
var alphaNumericStream = alphabetStream.merge(numericStream)
```



- Merge two streams into one new stream
- No ordering guarantee on resulting stream
- Intermediate operation
- KStream

weather

cogroup



traffic

```
var groupedWeather = weatherStream.groupByKey();  
var groupedTraffic = trafficStream.groupByKey();
```

```
var locationsCogroup = groupedWeather.cogroup(WEATHER_AGGREGATOR)  
    .cogroup(groupedTraffic, TRAFFIC_AGGREGATOR)  
    .aggregate(() -> new Location(), Materialized.with(stringSerde, jsonSerde));
```



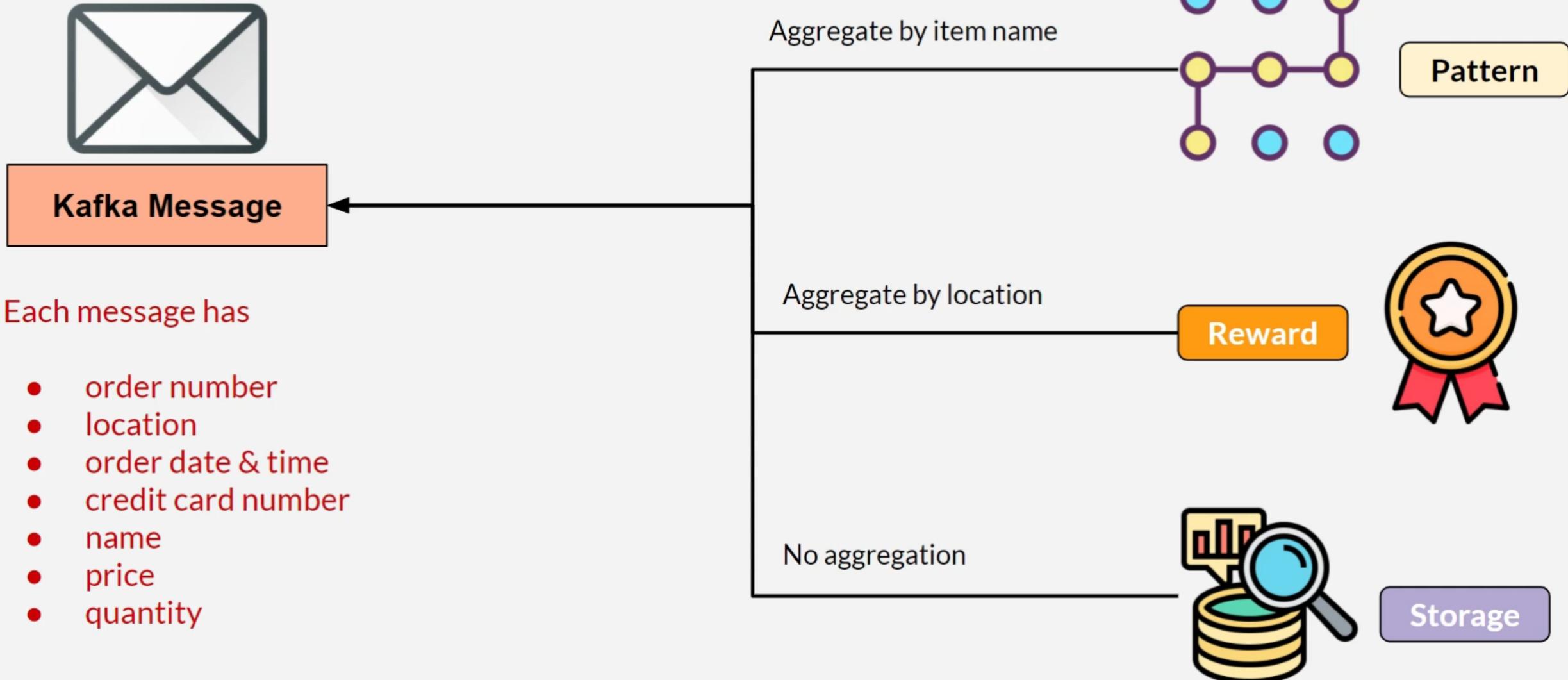
- Intermediate operation
- KStream
- Need aggregator for each cogroup
- This sample : String key, JSON value
- Difference with merge()

cogroup - Aggregator

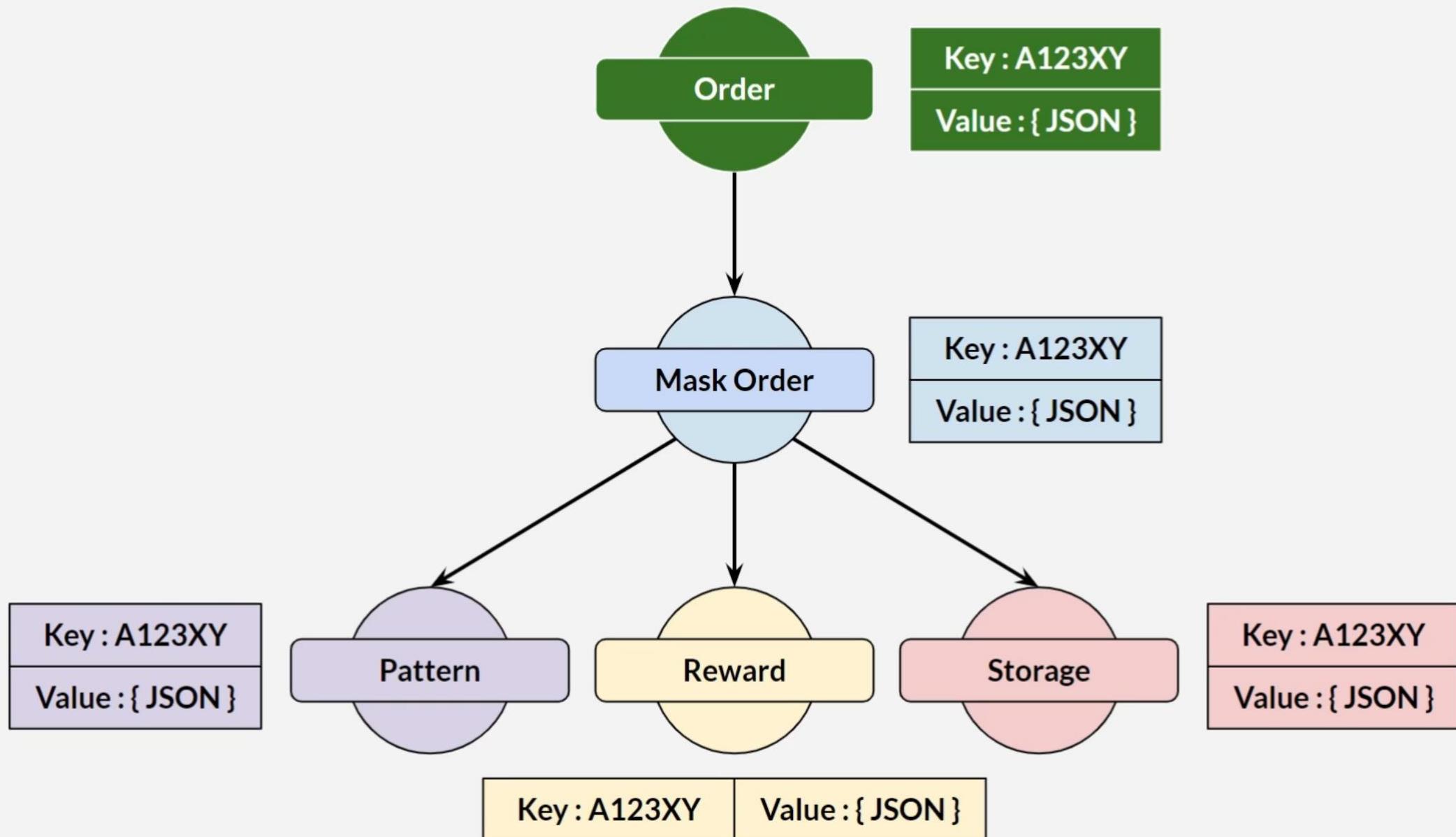
```
Aggregator<String, String, Location> WEATHER_AGGREGATOR = new Aggregator<String, String, Location>() {  
  
    @Override  
    public Location apply(String key, String value, Location aggregate) {  
        aggregate.setWeather(value);  
        return aggregate;  
    }  
};
```

```
Aggregator<String, String, Location> TRAFFIC_AGGREGATOR = new Aggregator<String, String, Location>() {  
  
    @Override  
    public Location apply(String key, String value, Location aggregate) {  
        aggregate.setTraffic(value);  
        return aggregate;  
    }  
};
```

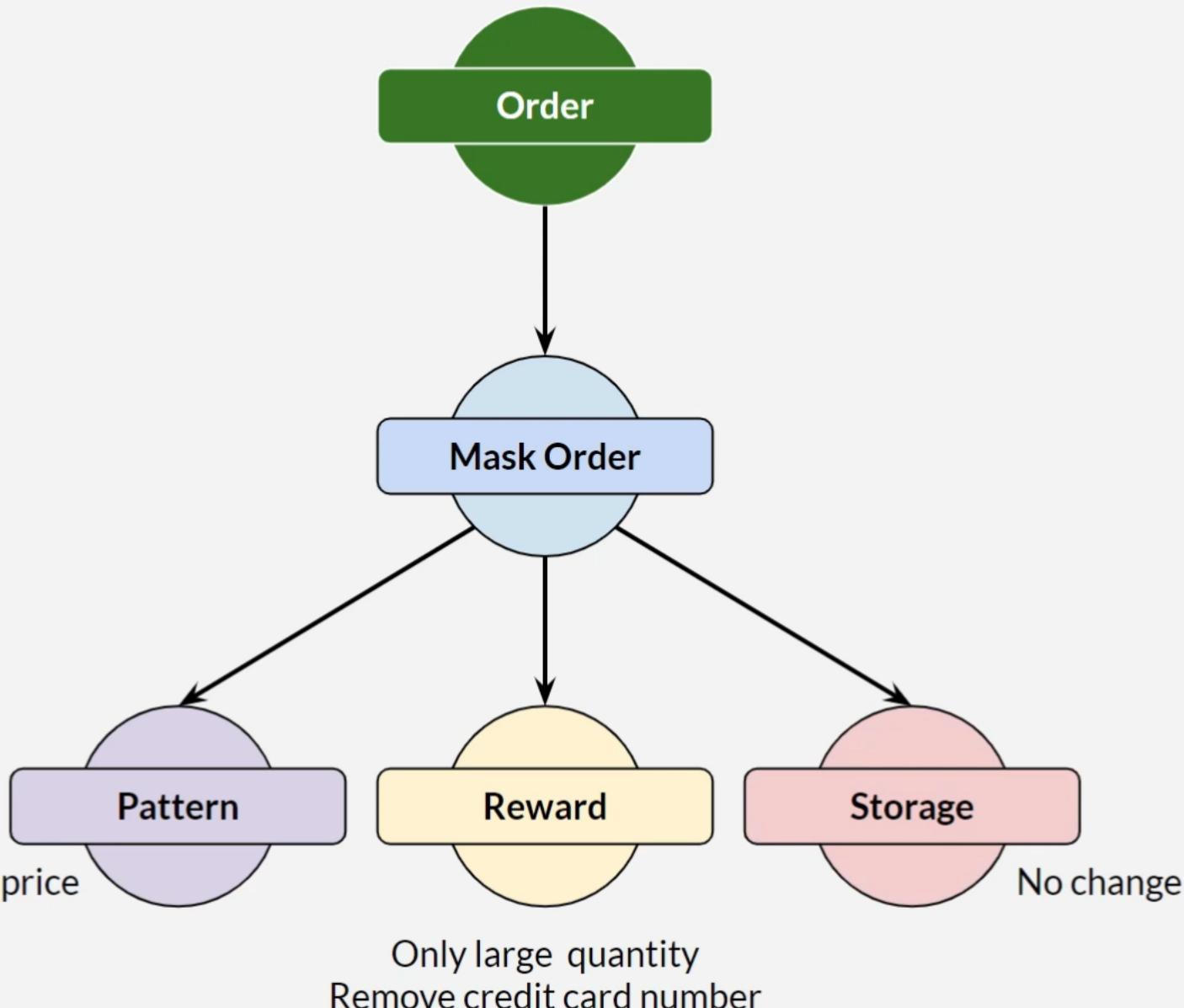
Order - Kafka Message



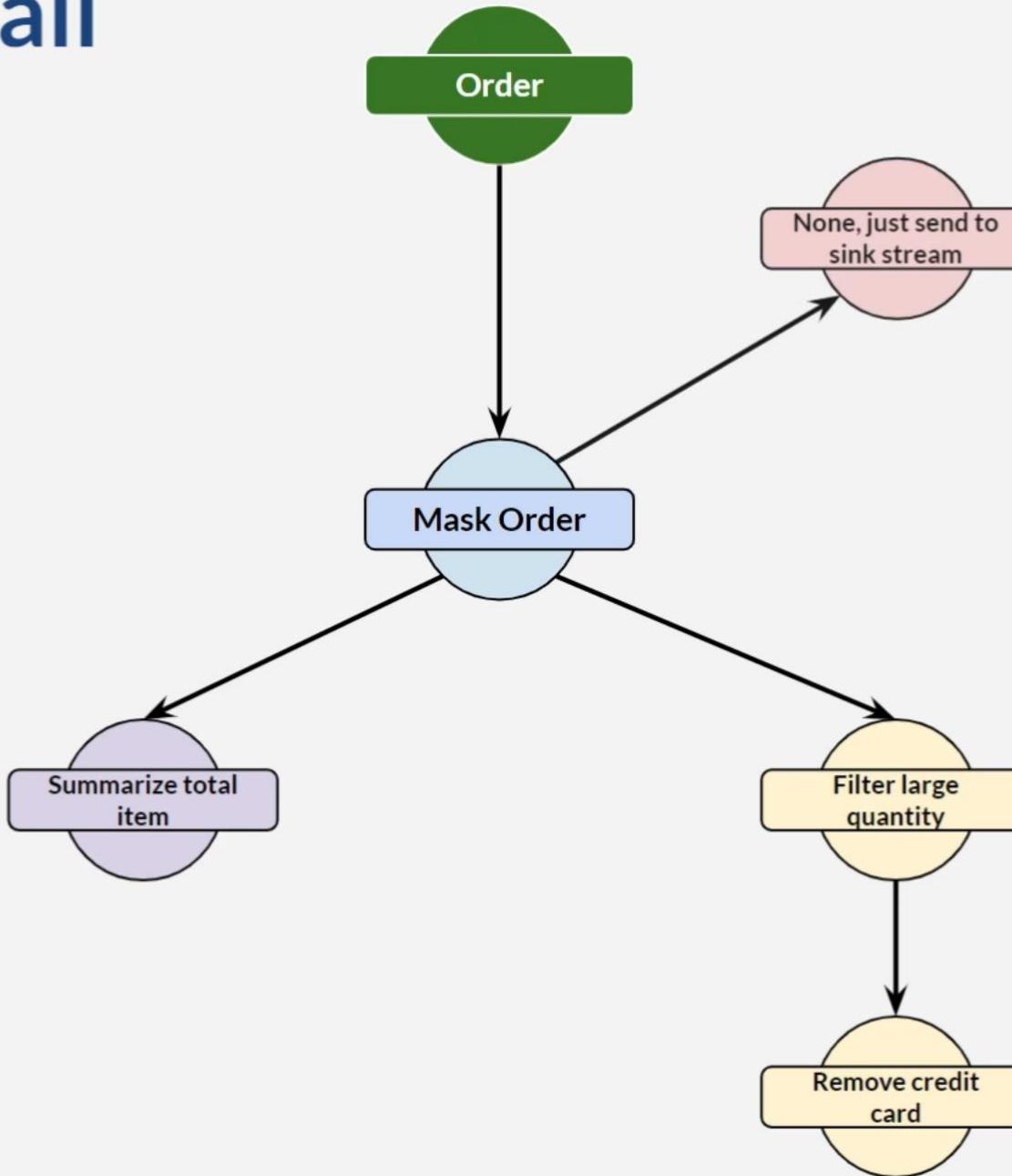
Topology



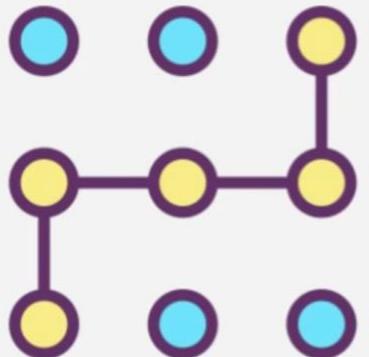
High Level Topology



Topology Detail



Additional Requirement



Pattern

Current : summarize item price * quantity

+ ADDITIONAL: split plastic & non plastic items



Reward

Current : give reward only for item with quantity > xxx

+ ADDITIONAL: give reward only for item that is not cheap

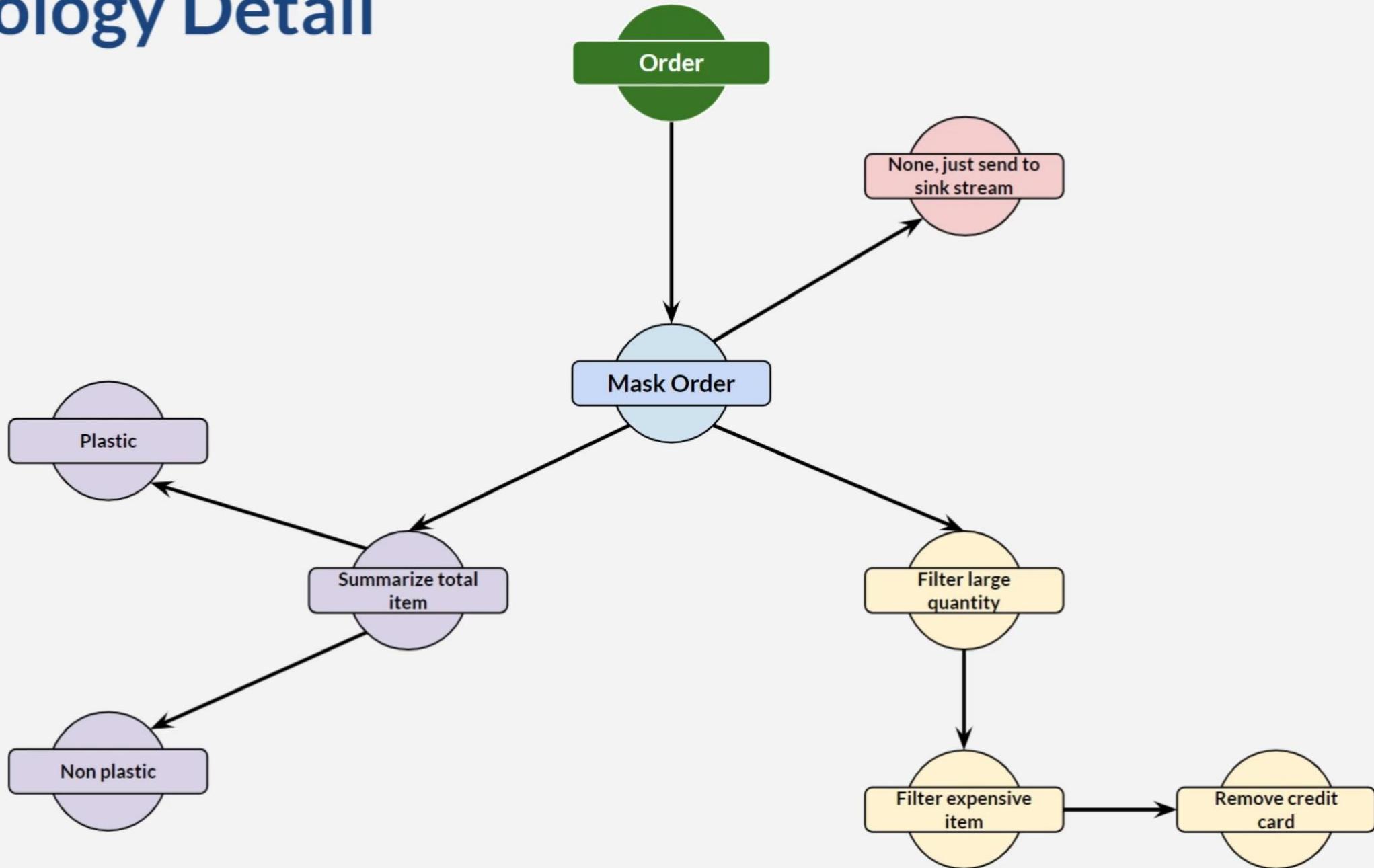


Storage

Current :-

+ ADDITIONAL: key is base64(order number)

Topology Detail



Sample Data

Cotton Dog

Price : 80
Qty : 250

Plastic Cat

Price : 400
Qty : 500

Wooden Horse

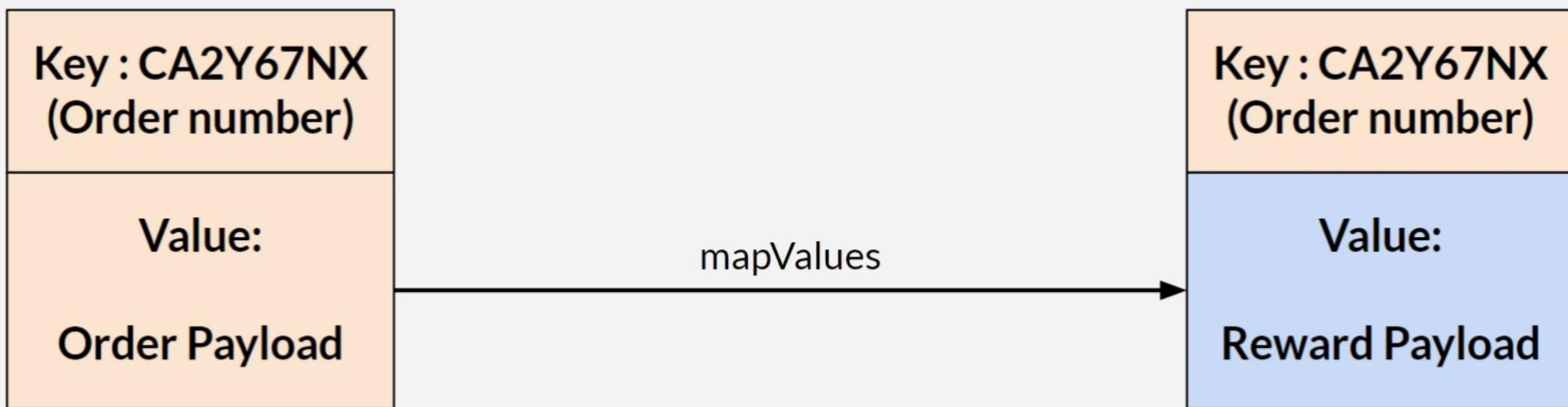
Price : 700
Qty : 90

Steel Pig

Price : 350
Qty : 270

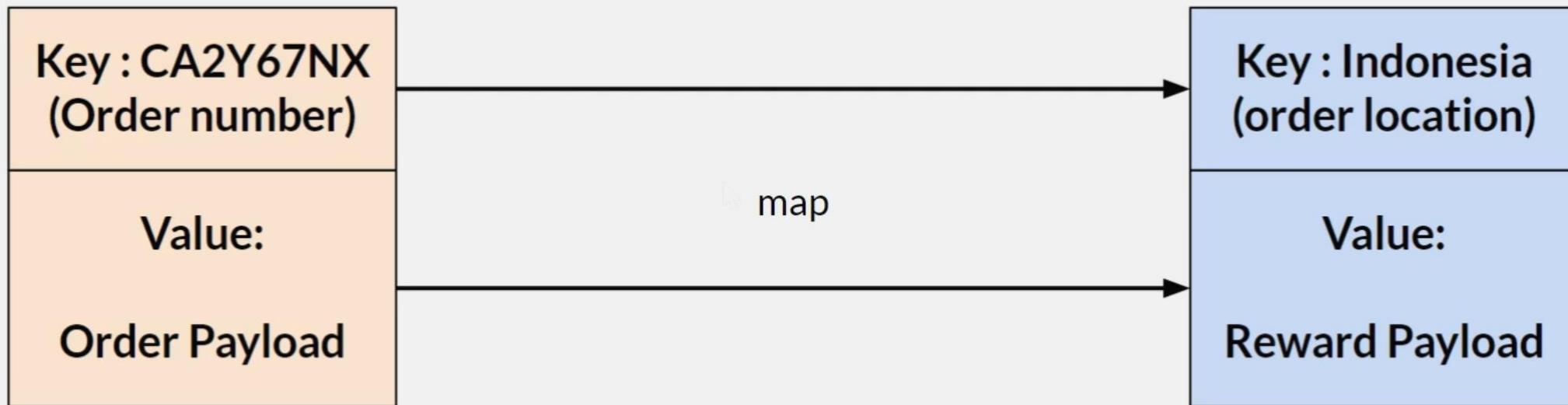
Stream (Kafka Sink Topic)	Data
Pattern - plastic	Plastic Cat
Pattern - not plastic	Cotton Dog, Wooden Horse, Steel Pig
Reward	Plastic Cat, Steel Pig
Storage	Plastic Cat, Cotton Dog, Wooden Horse, Steel Pig

Reward Message



Project Explorer
kafka-stream-order [root] [main]
kafka-stream-pattern [root] [main]
kafka-stream-reward [root] [main]
kafka-stream-sample [root] [destools]
kafka-stream-storage [root] [main]

Reward Message



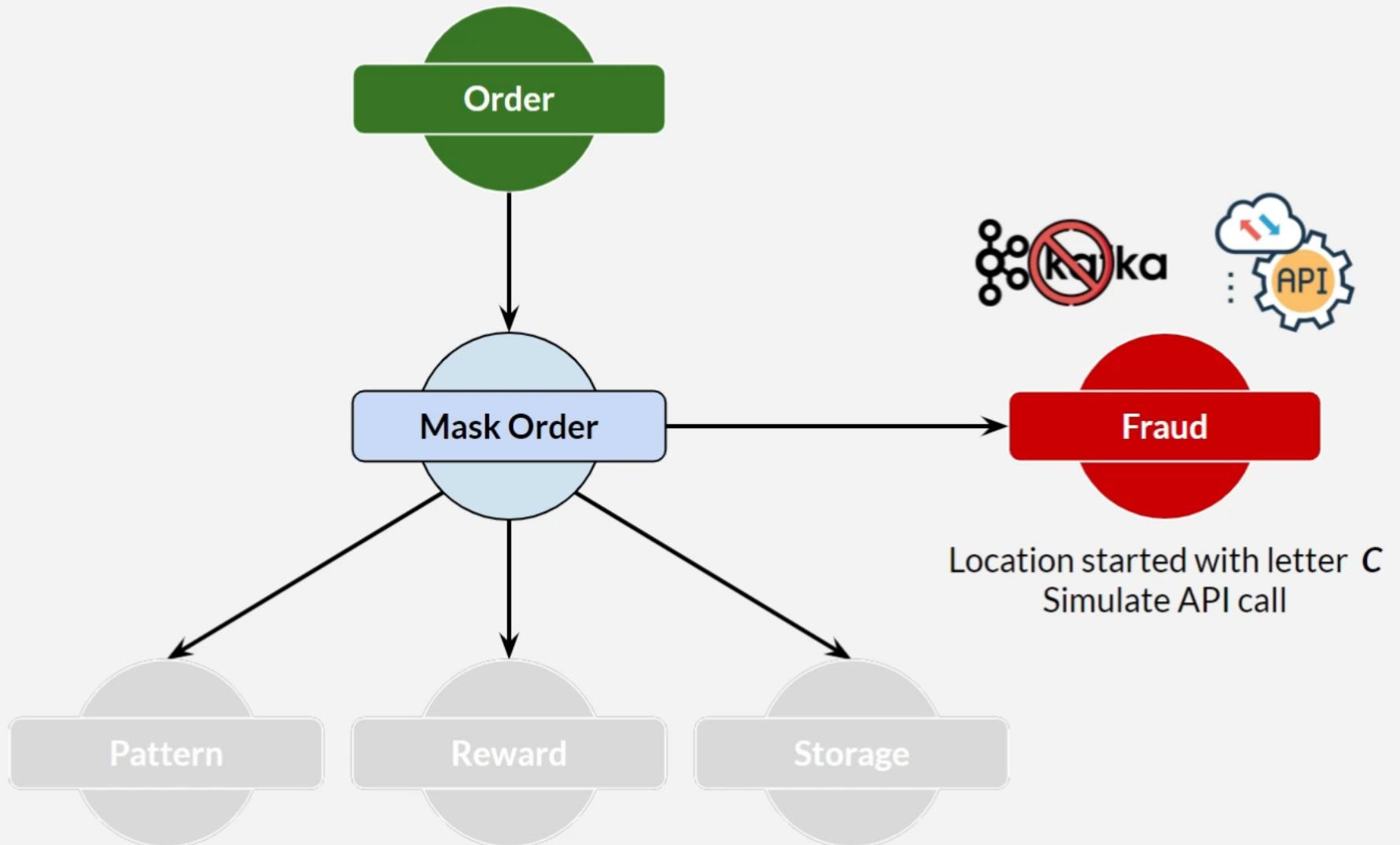
Maven → Console × %Progress × Search → Gradle Tasks → Gradle Executions → JUnit → Call Hierarchy → Scan Pending

No consoles to display at this time.

Something Suspicious



High Level Topology



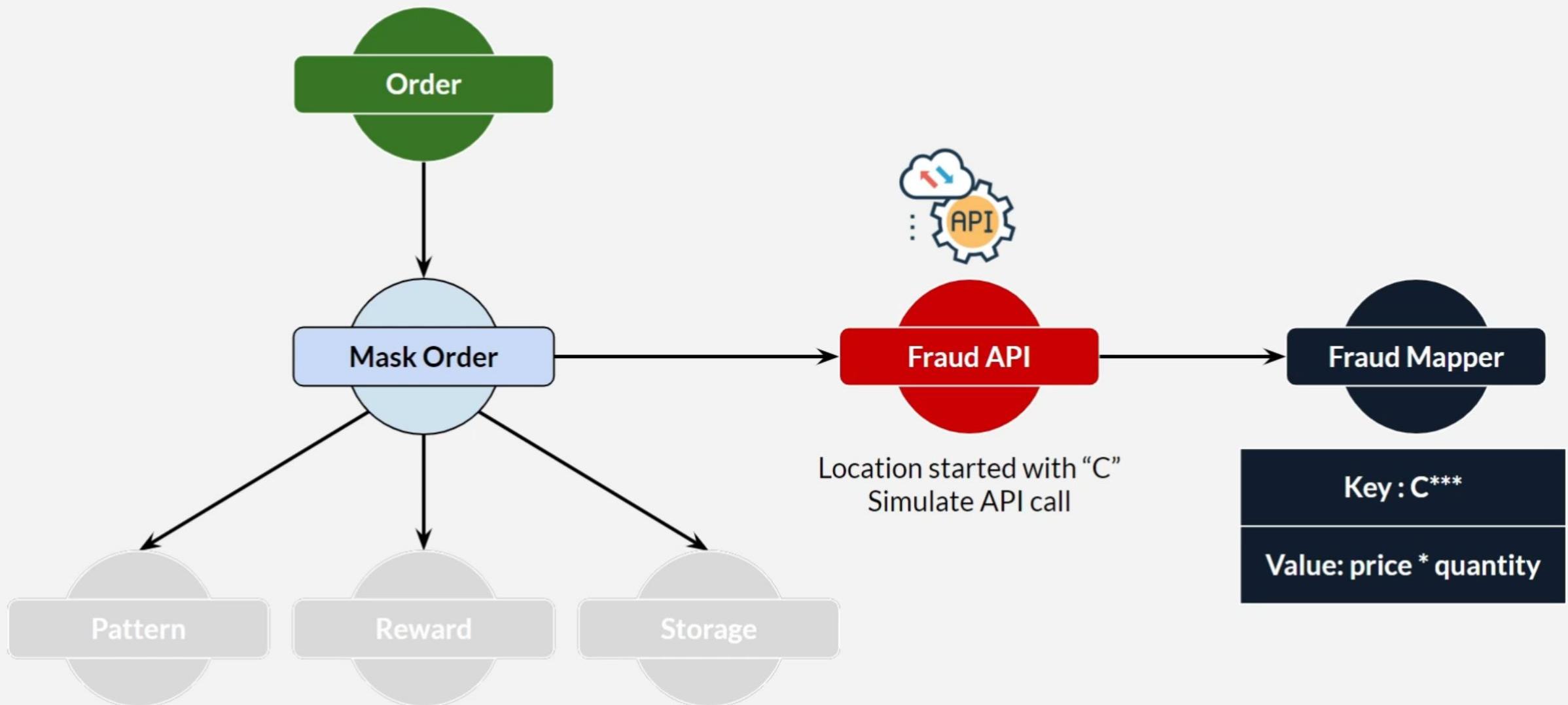
Something Suspicious



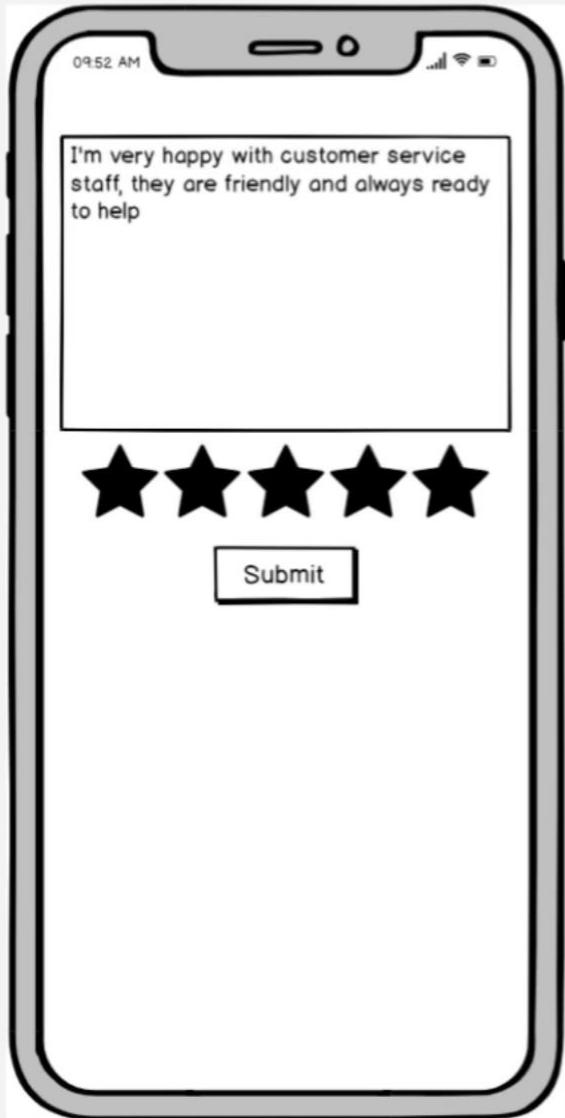
Key: C***

Value: price * quantity

High Level Topology

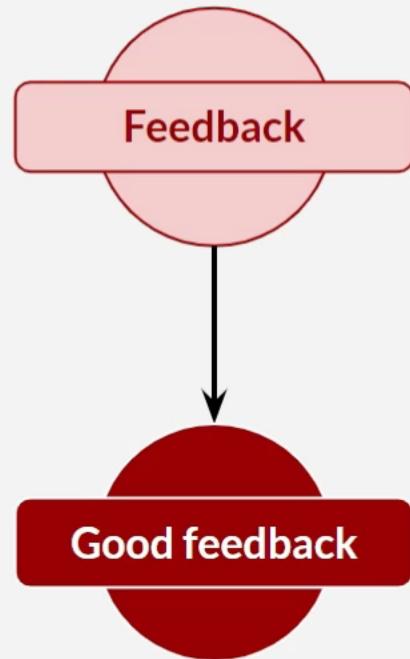


Customer Feedback



happy, good, helpful, etc

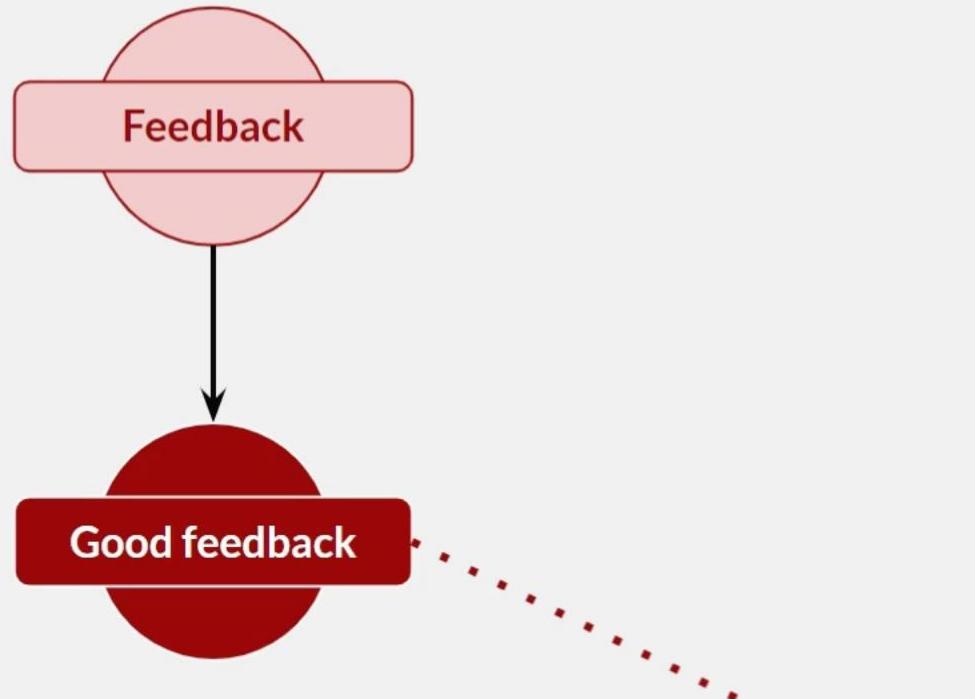
High Level Topology



Order app- project for feedback

- ▶ Feedback*.java
- ▶ Package: **com.virtusa.kafka**
 - ▶ api.request
 - ▶ api.server
 - ▶ broker.message
 - ▶ broker.producer
 - ▶ command.action
 - ▶ command.service

High Level Topology



Who own's the good feedback?

Key : [branch location]

Value: [good word]

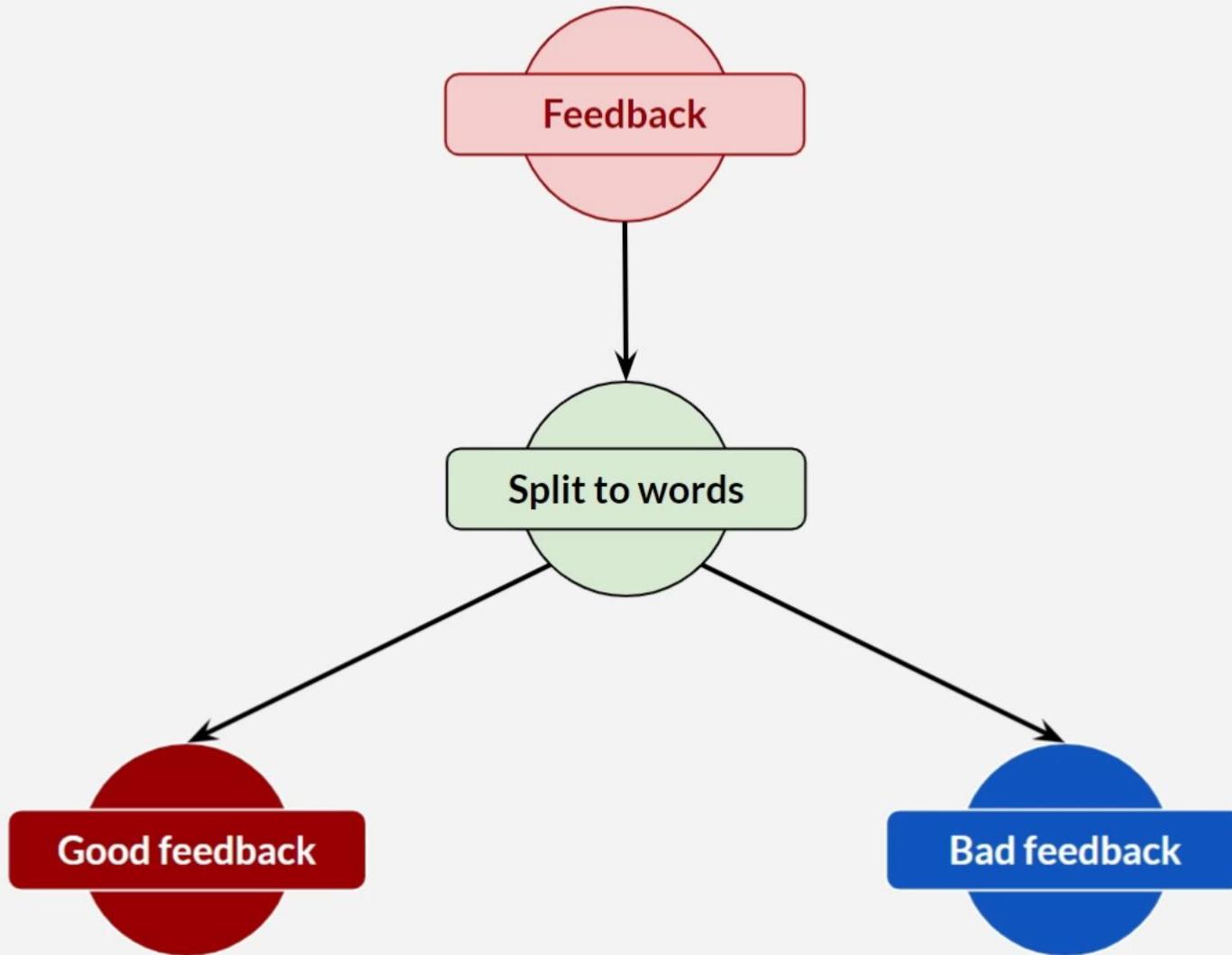
Java Stream API

- ▶ Since Java 8
- ▶ Not kafka stream
- ▶ Some method names:
 - filter
 - filterMap
 - filterEach
 - map
 - peek

Bad Feedback?

- ▶ Analyze bad feedback
- ▶ Stream to analyze feedback : good or bad
- ▶ Topology

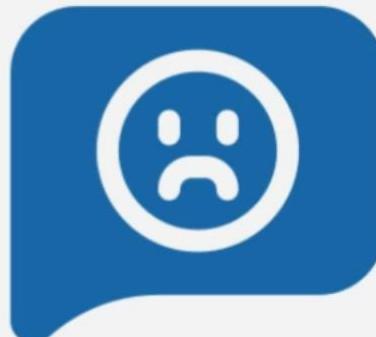
High Level Topology



Customer Feedback

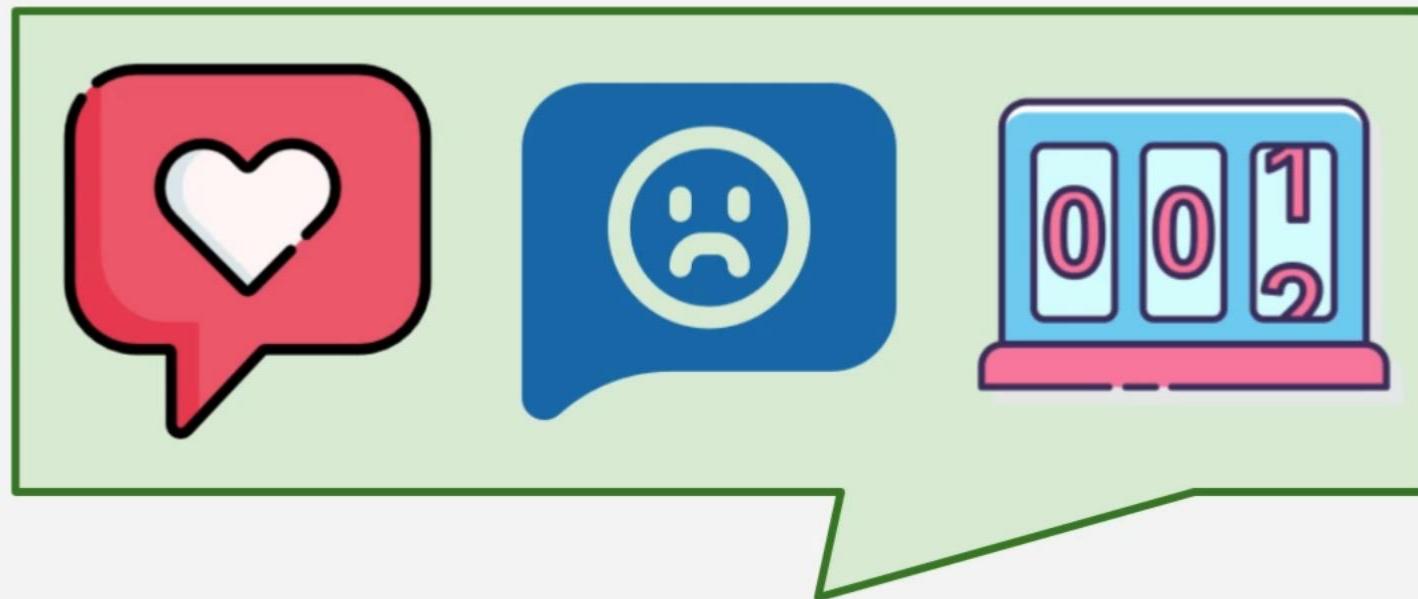


happy, good, helpful, etc

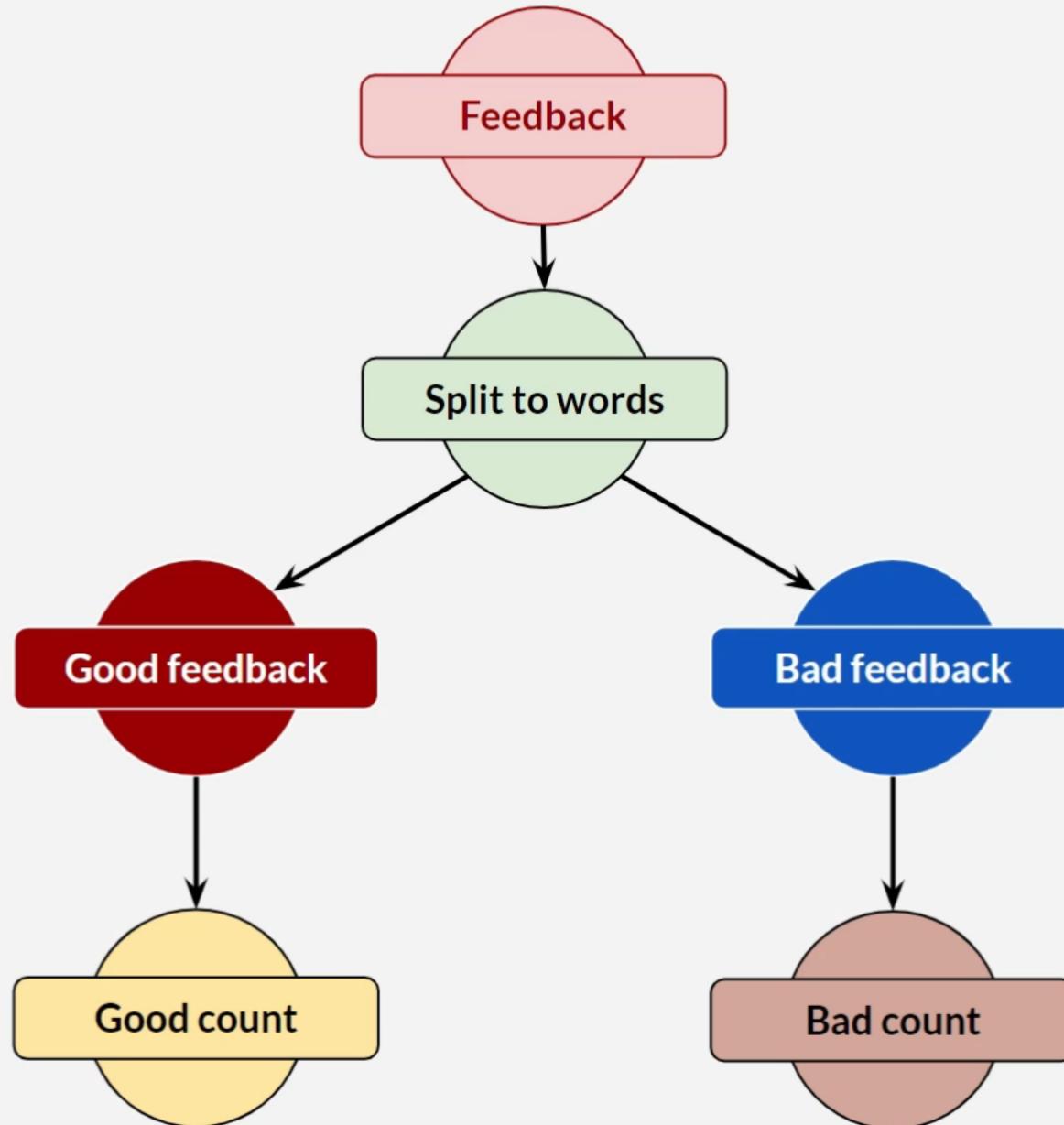


angry, sad, bad, etc

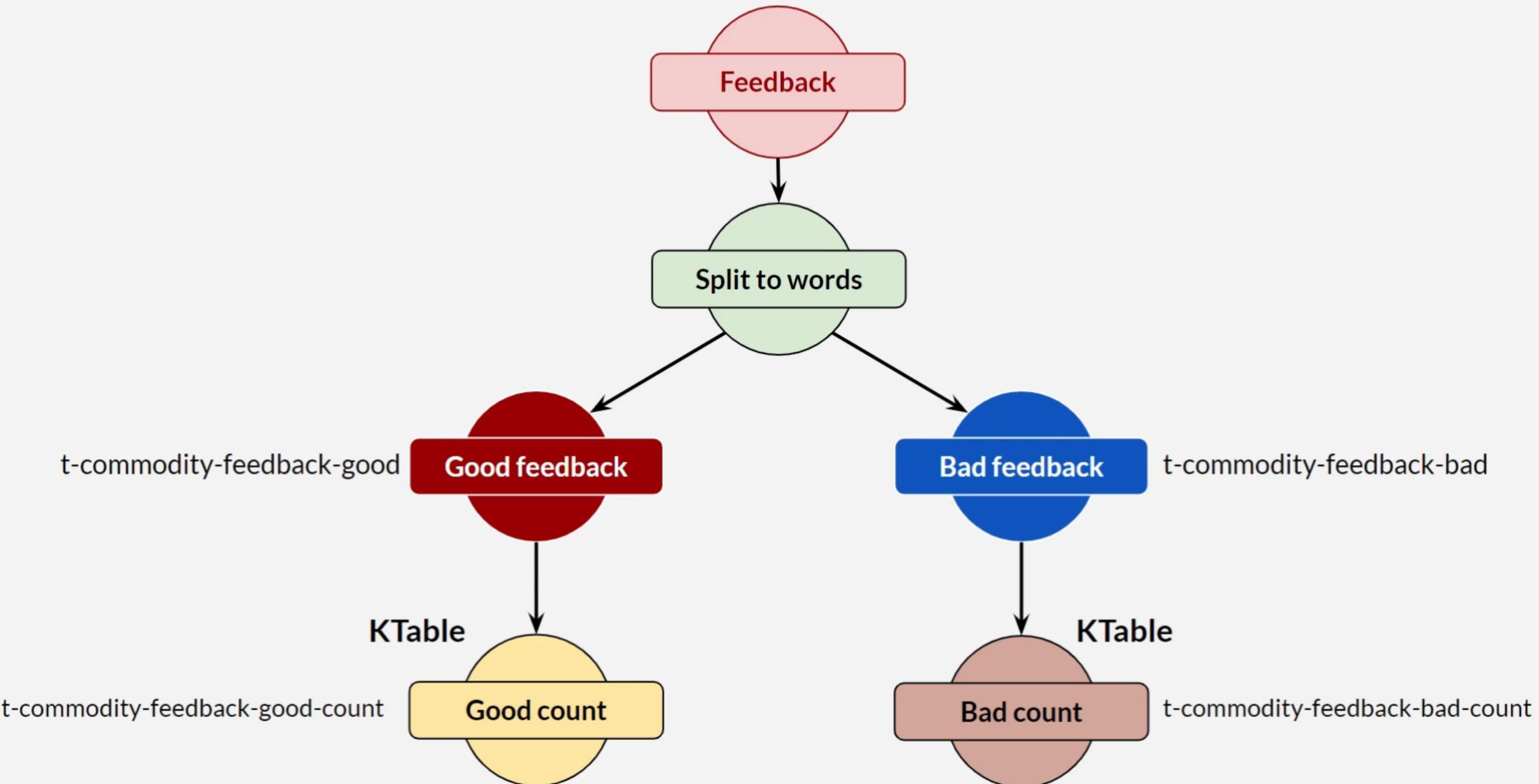
Count The Word



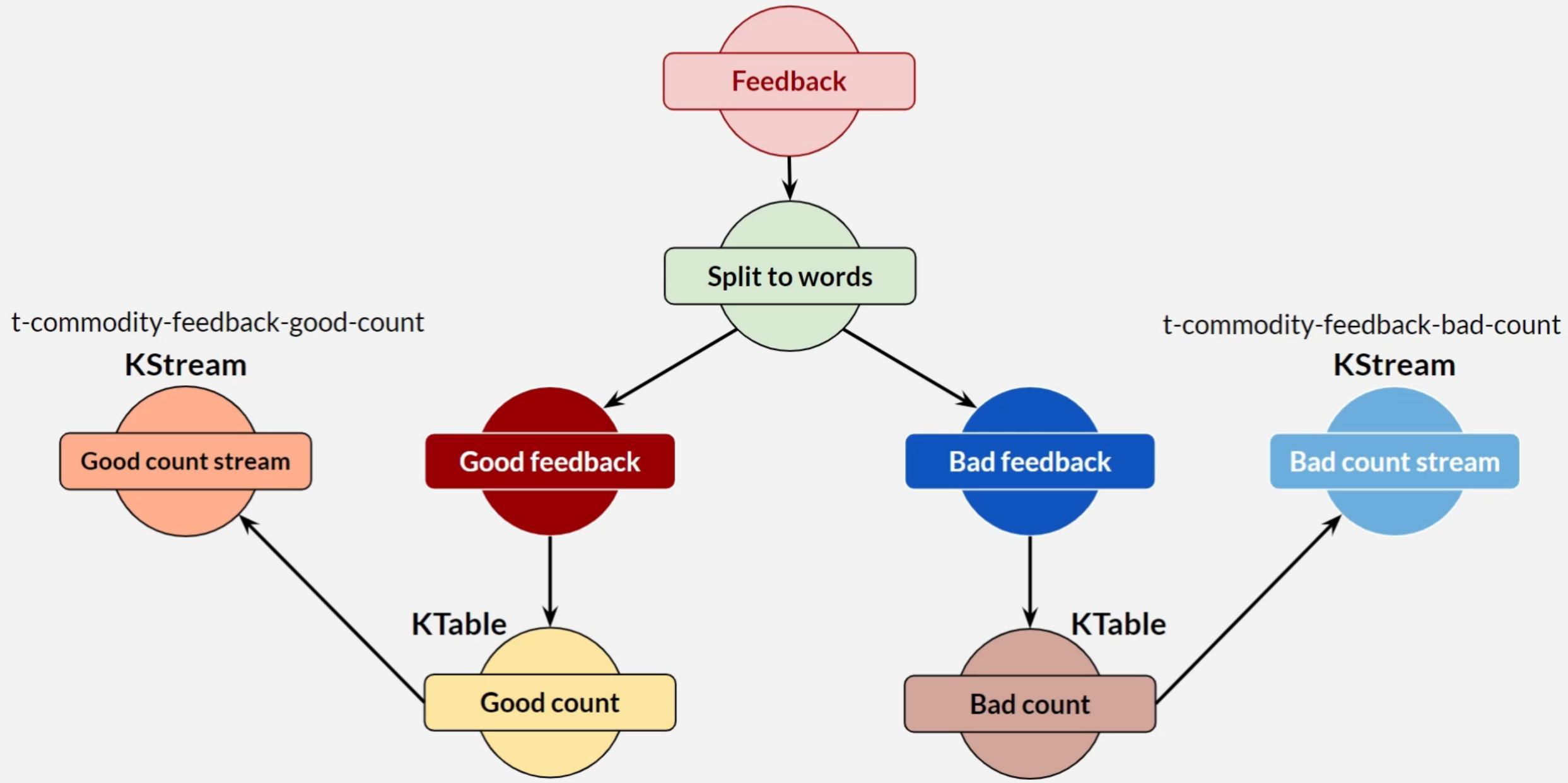
High Level Topology



High Level Topology



High Level Topology



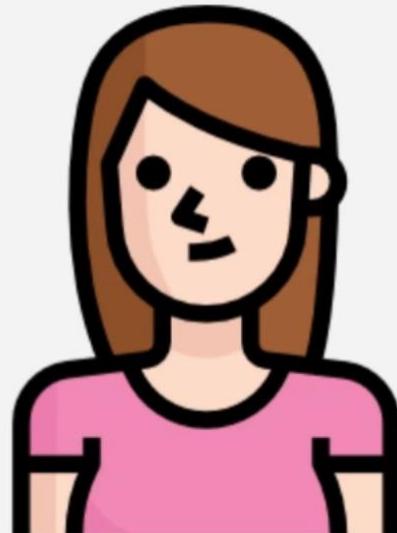
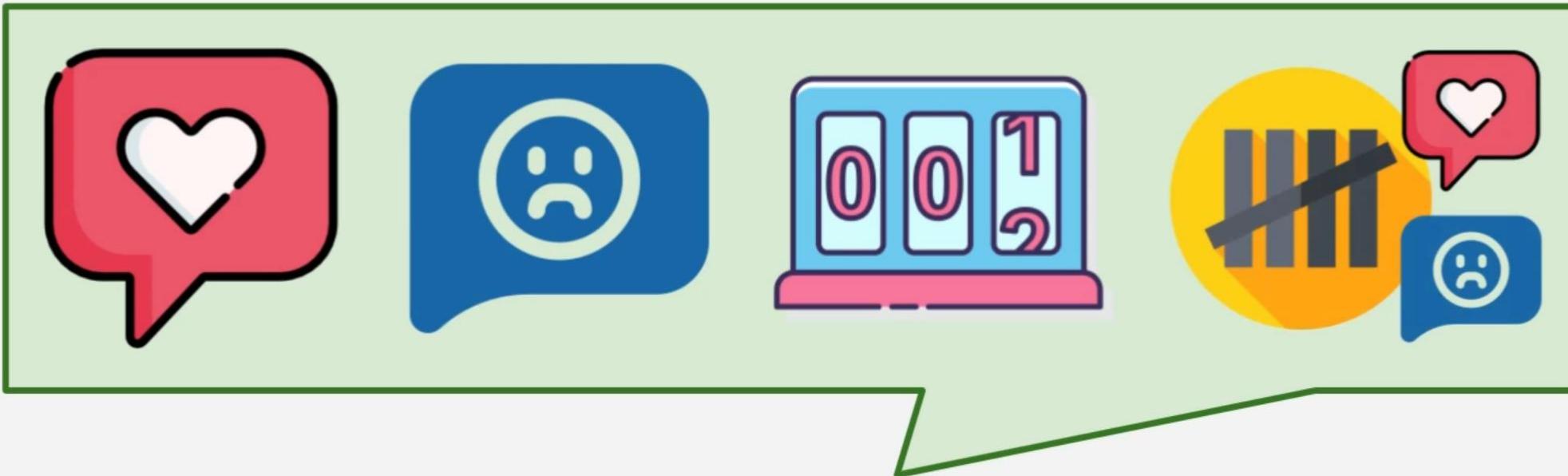
Kafka Stream Configuration

- × Default configuration : 30 seconds
- × Cache and send
- × On commit.interval.ms
- × Adjust configuration

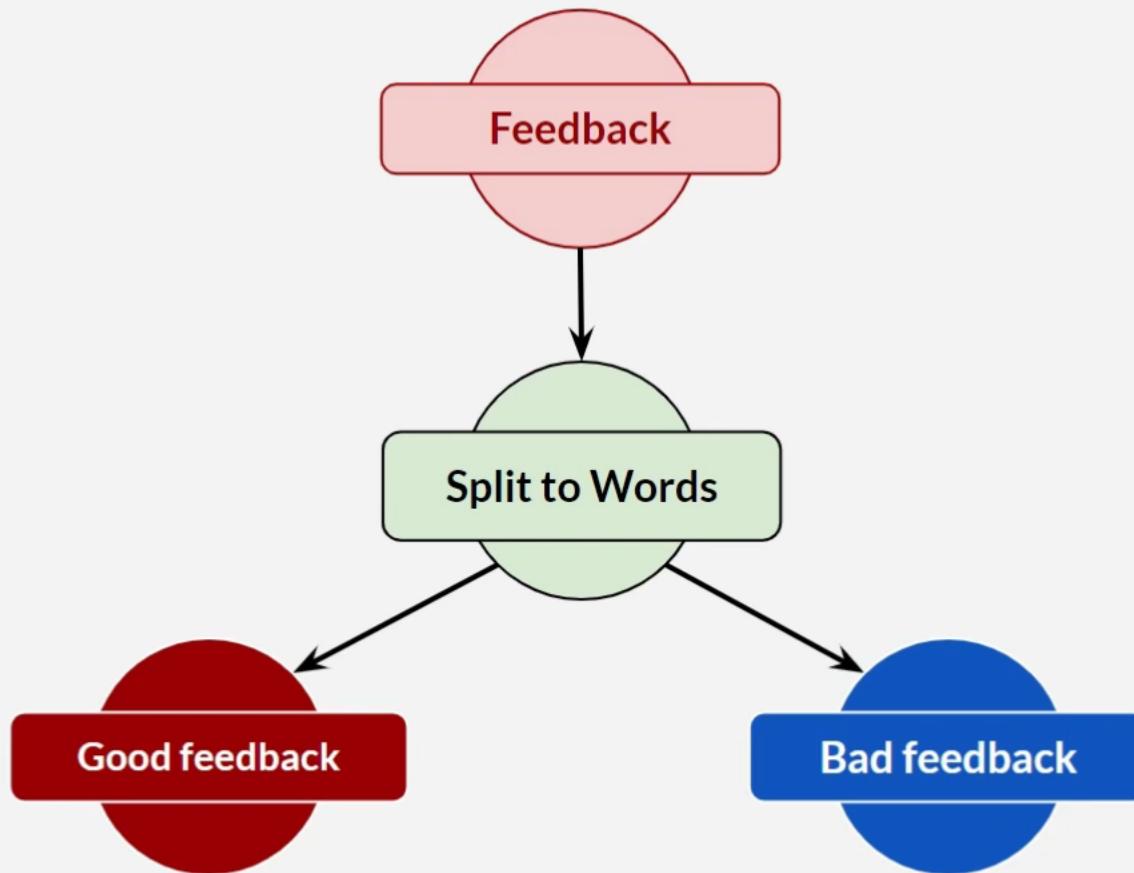
Source code for feedback

- ▶ Kafka-ms-order
- ▶ Kafks-stream-sample

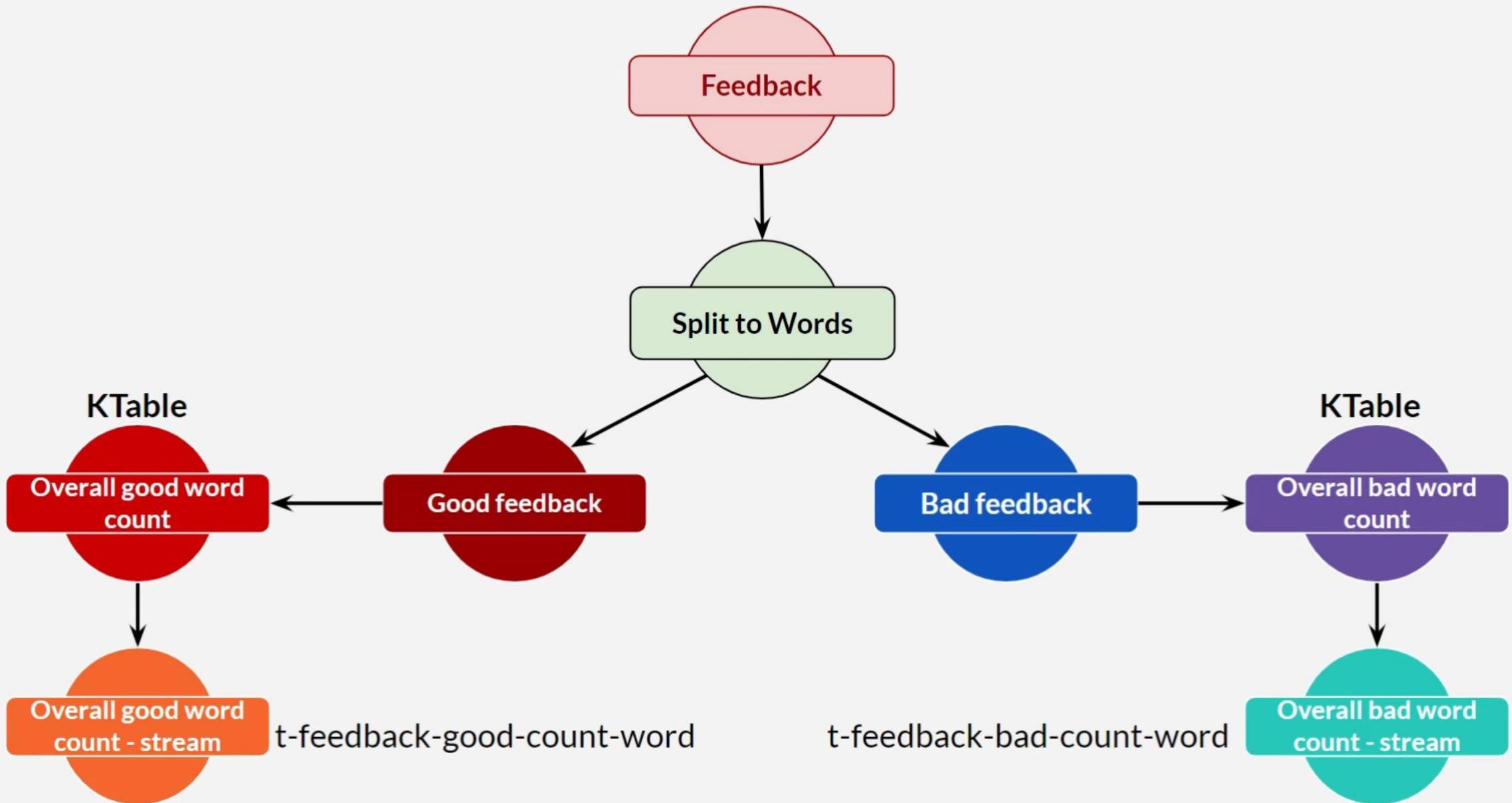
Count The Word



High Level Topology



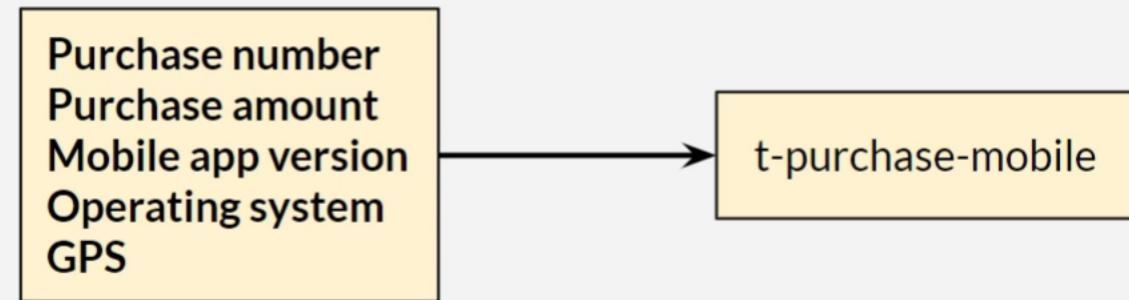
High Level Topology



Customer Stream

- Mobile & Web -

Different Device, Different Data



t-purchase-all

Purchase Web 1

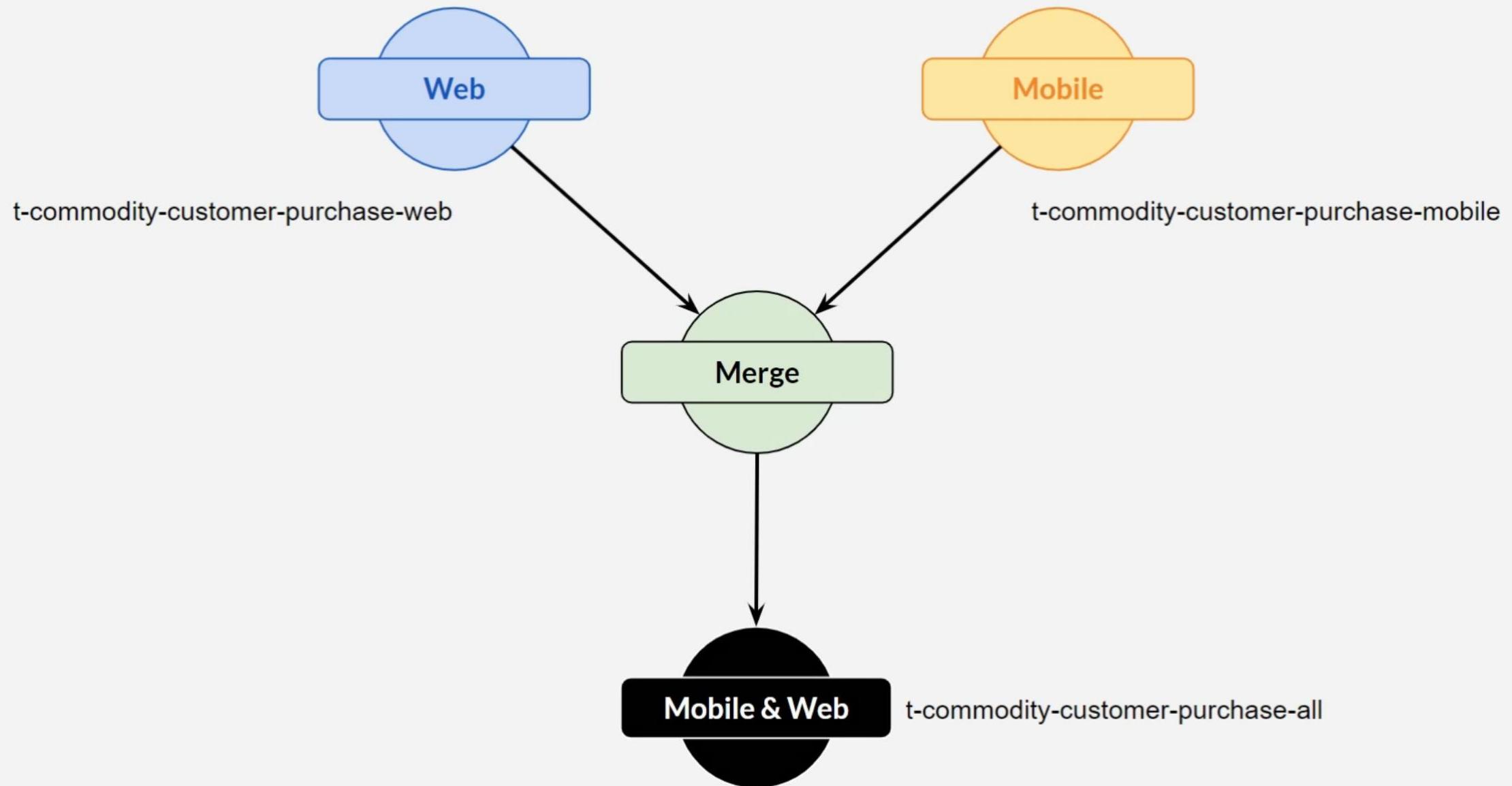
Purchase Mobile 1

Purchase Mobile 2

Purchase Web 2

Purchase Mobile 3

High Level Topology



Create project for Customer Purchase

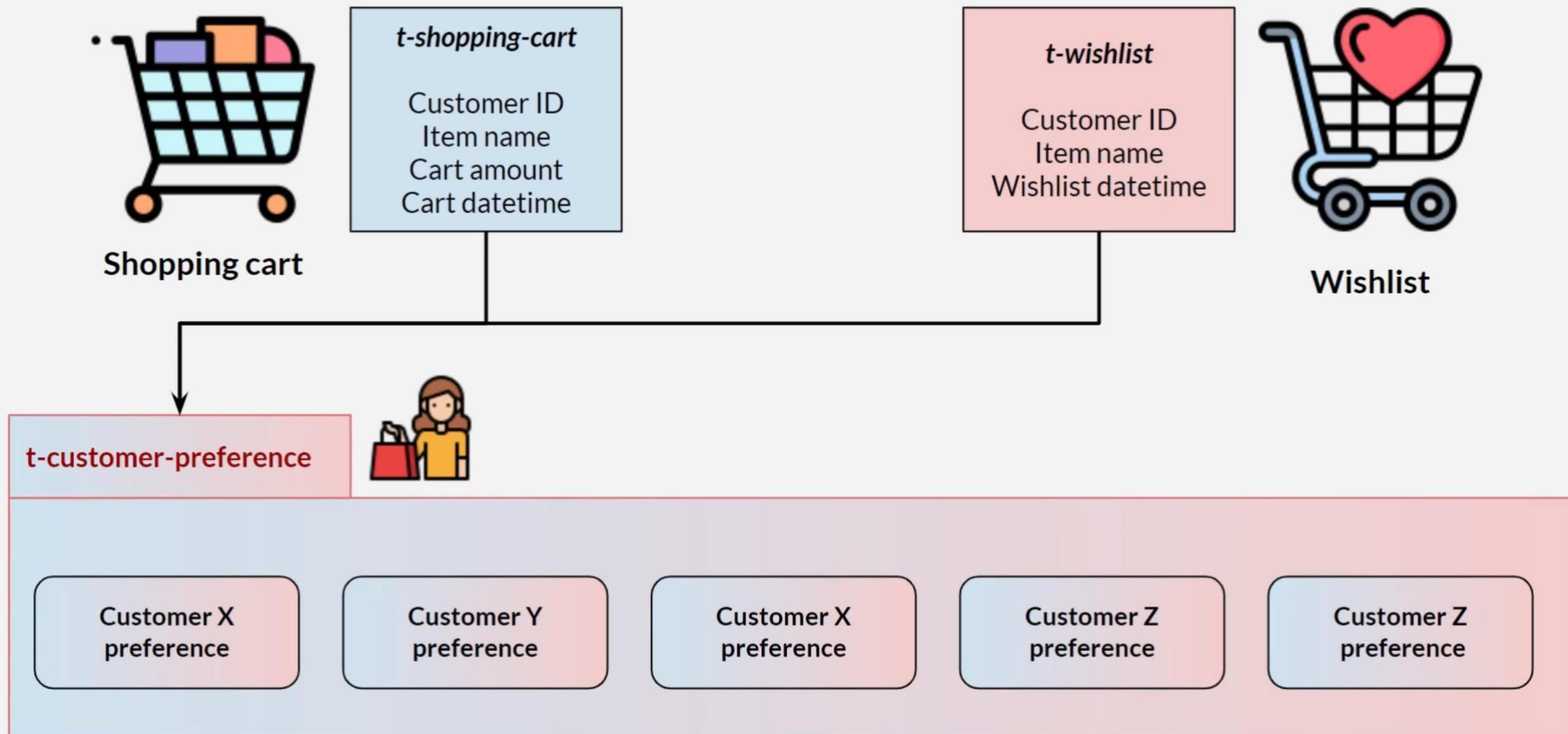
- ▶ CustomerPurchase*.java
- ▶ Package: **com.virtusa.kafka**
 - ▶ api.request
 - ▶ api.server
 - ▶ broker.message
 - ▶ broker.producer
 - ▶ command.action
 - ▶ command.service

“

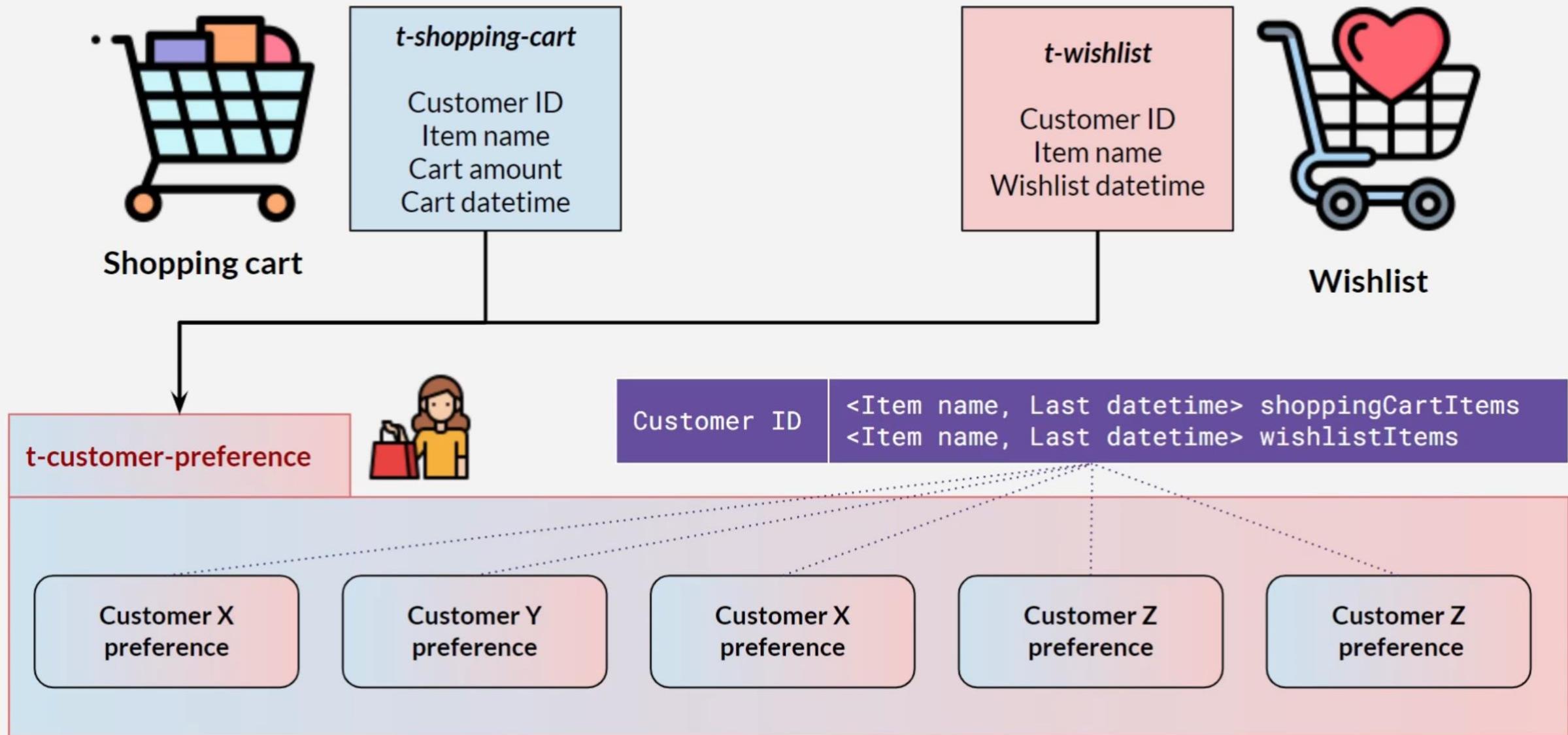
Cart & Wishlist

”

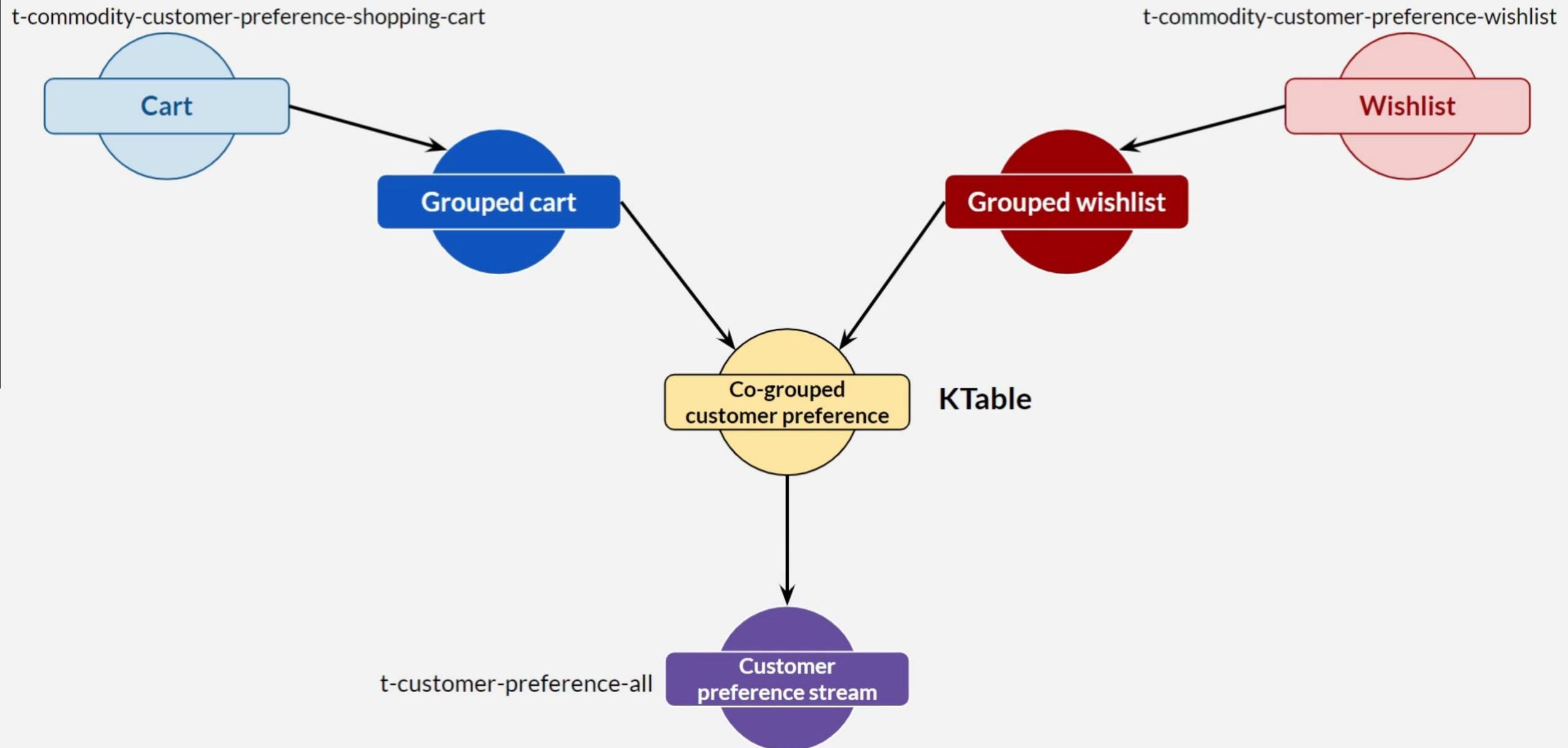
Two Things, One Customer Preference



Two Things, One Customer Preference



High Level Topology



Create project for Customer Preference

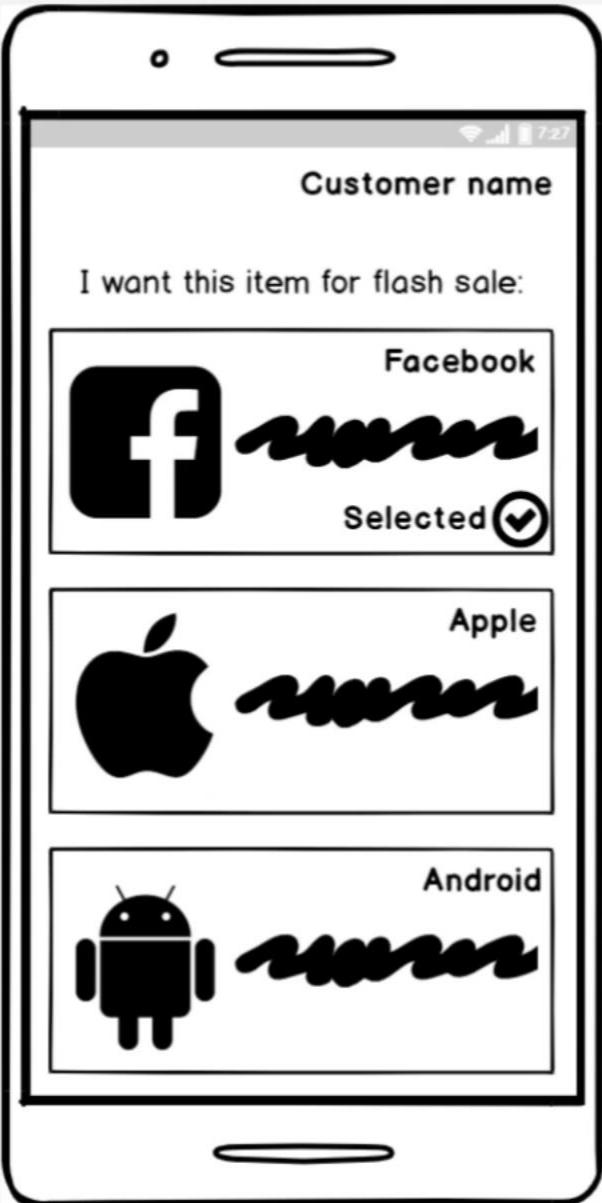
- ▶ CustomerPreference*.java
- ▶ Package: **com.virtusa.kafka**
 - ▶ api.request
 - ▶ api.server
 - ▶ broker.message
 - ▶ broker.producer
 - ▶ command.action
 - ▶ command.service

“

Flash Sale Vote

”

Flash Sale by Customer

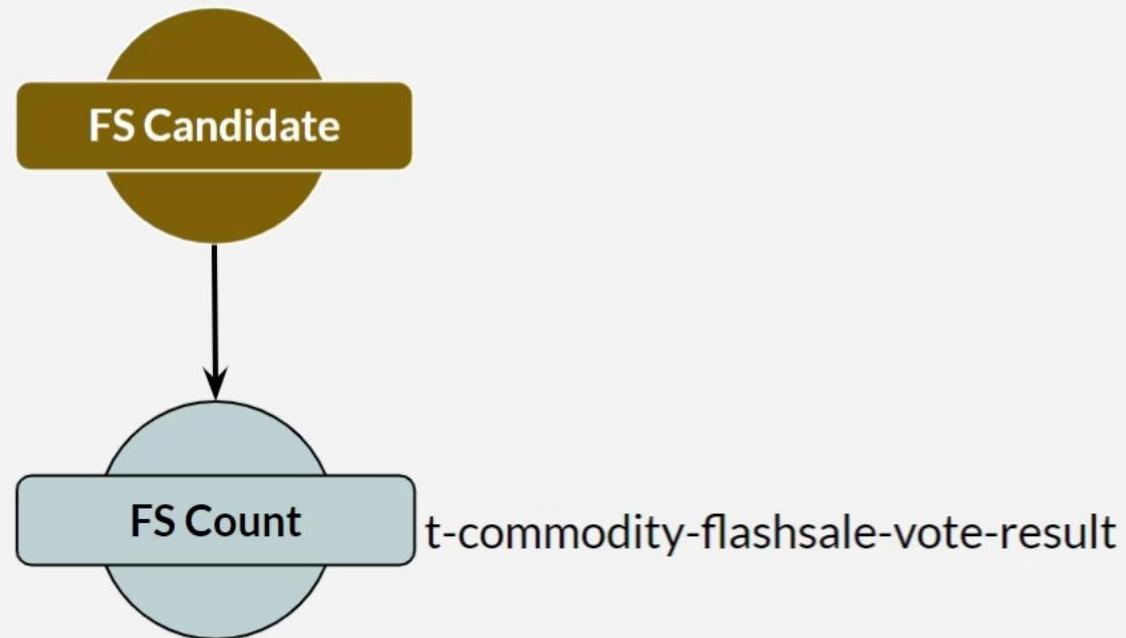


- Customer can vote flash sale candidate
- Next flash sale will be most voted items
- One customer, one candidate
- Can change selection

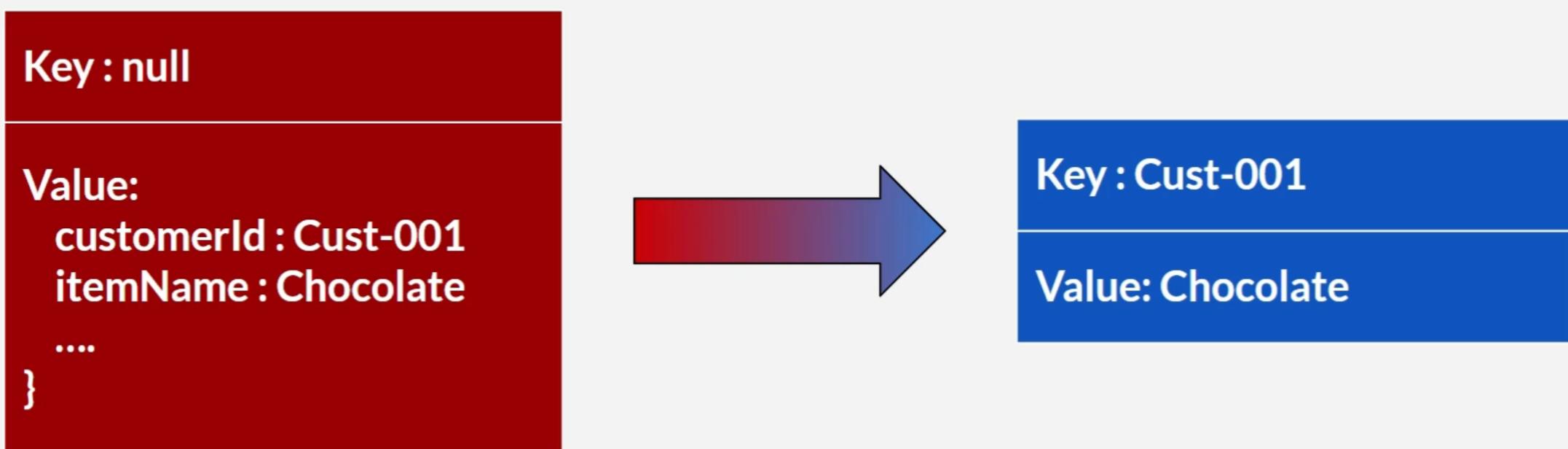
Stream or Table?

- × Track latest selected item per customer
- × Kafka Stream table
- × Record
 - × Key : Customer ID
 - × Value : Flash sale candidate
- × Upserts (update /insert)

High Level Topology



Processing The Message



How Vote Works?

Anna	Cookies	
------	---------	---

Olaf	Cookies	
------	---------	---

Olaf	Cake	
------	------	---

Anna	Cake	
------	------	--

Elsa	Cookies	
------	---------	---

timeline ↓

	Cookies	1
---	---------	---

	Cookies	2
---	---------	---

	Cookies	1
	Cake	1

	Cookies	0
	Cake	2

	Cookies	1
	Cake	2

“

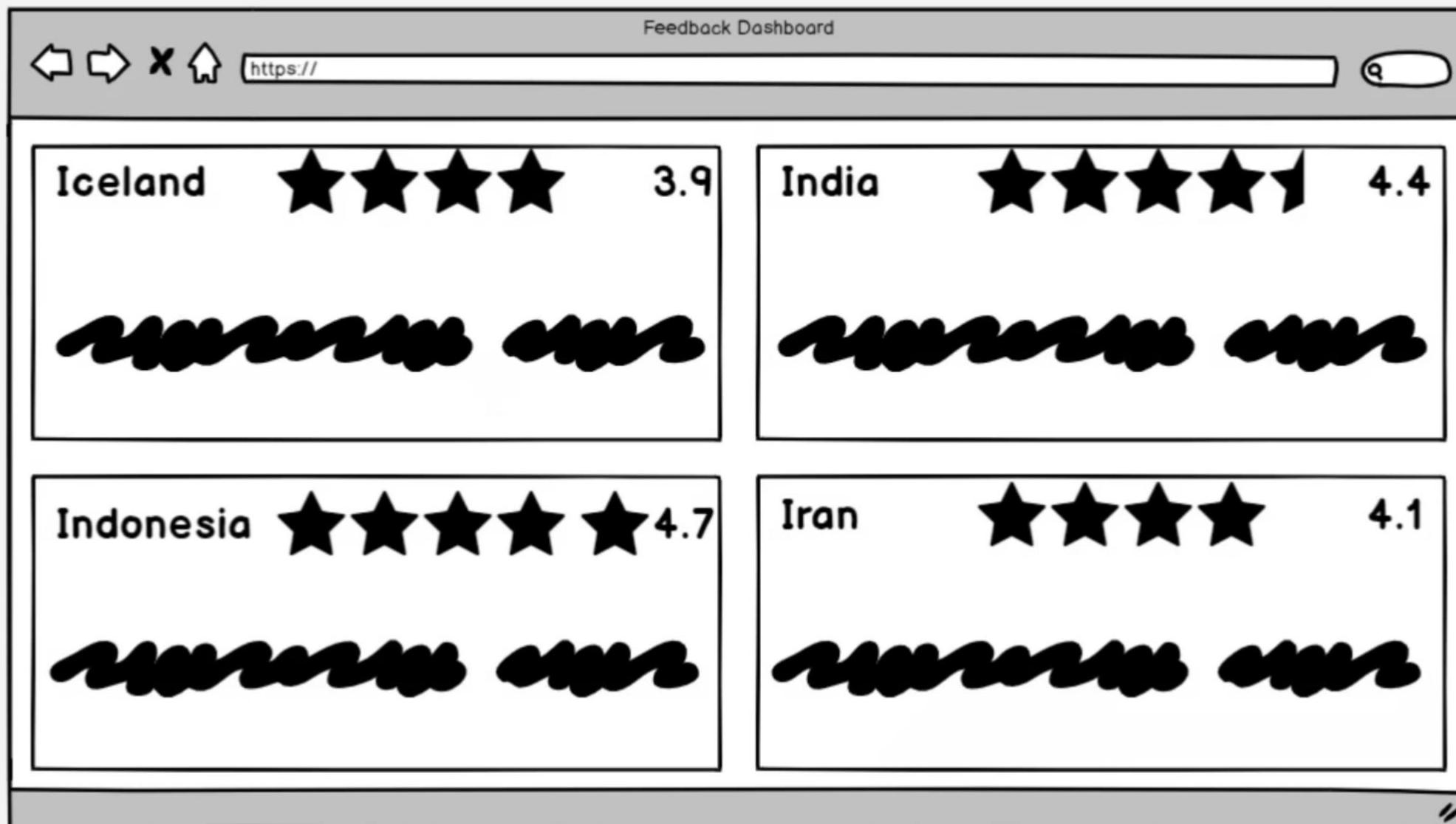
Feedback Rating

”

Feedback Dashboard



Feedback Dashboard



Average Rating

- × $\text{average} = \text{sum(ratings)} / \text{count(ratings)}$
- × Need to know all ratings
- × Use processor API & state store

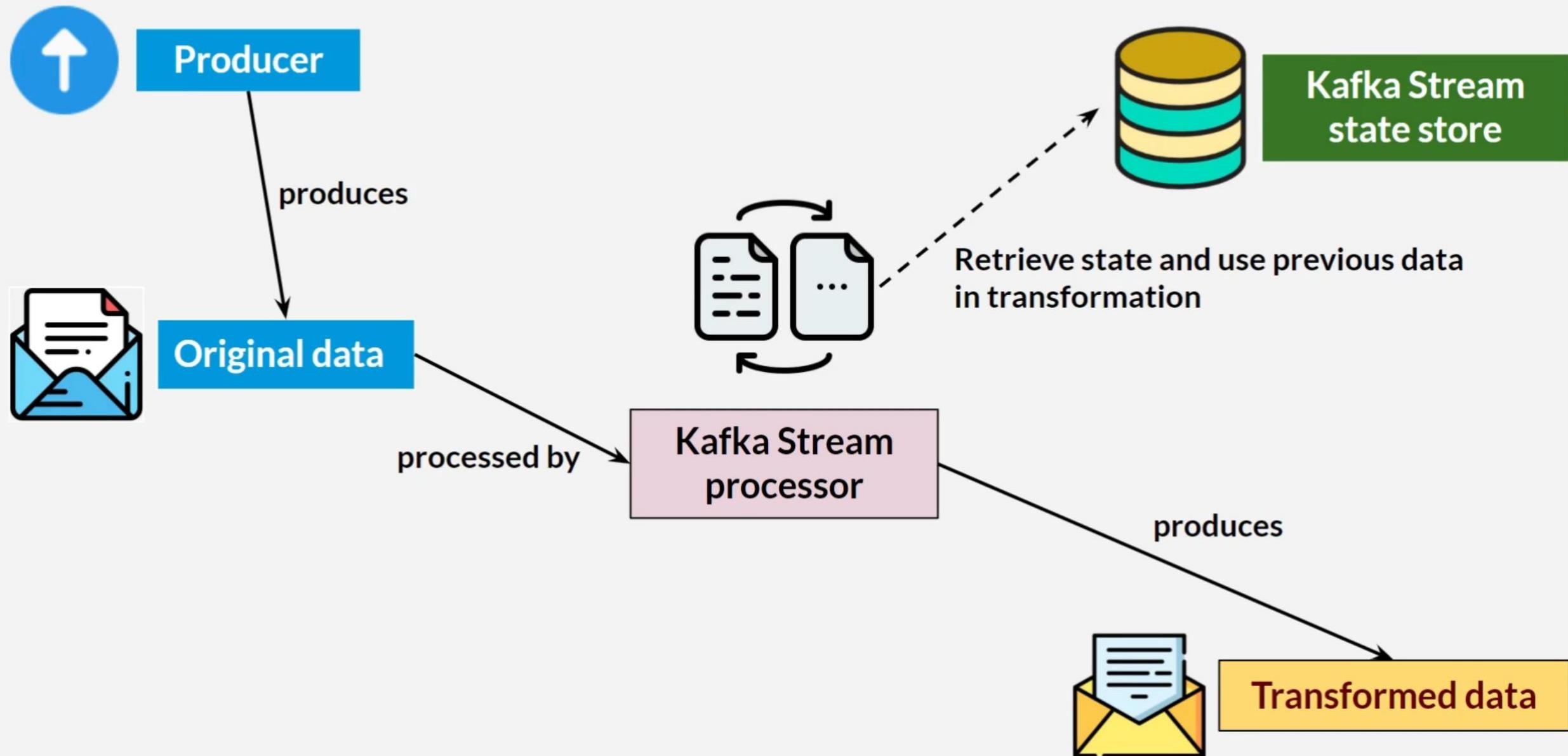
State in Kafka Stream

- × Care about current data
- × No need to know previous data
- × Flash sale example : need to know current / previous choice
- × "User choice" is kafka stream **state**
- × Stateless operations (commodity & feedback)
- × Stateful operations (flash sale vote)

State Store

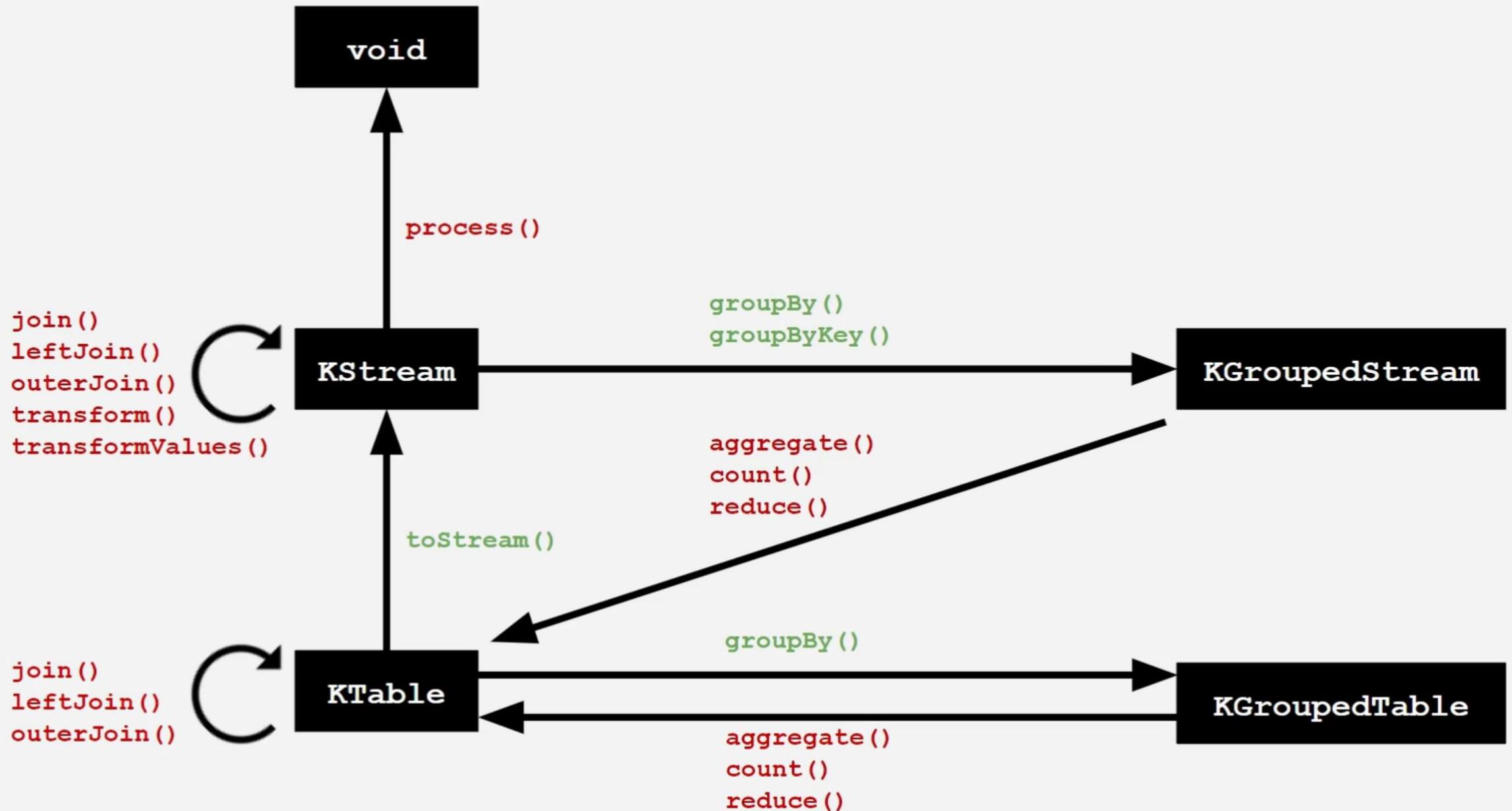
- × See existing information & connect it
- × Kept in state store
- × Key-value data storage
- × Accessed from processor
- × Kafka stream state stores:
 - × In-memory
 - × Persistent (disk based)

State Store

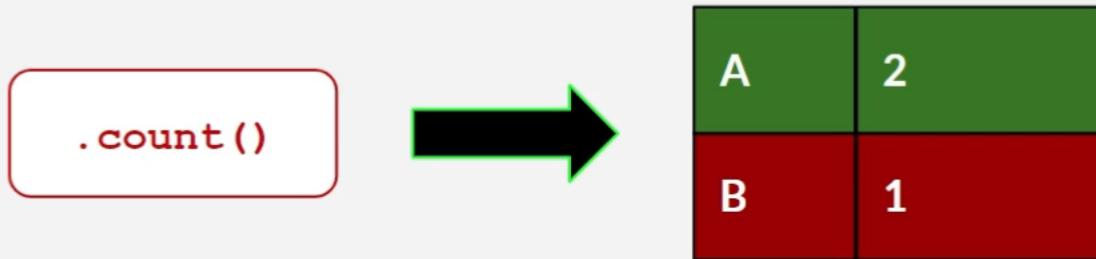
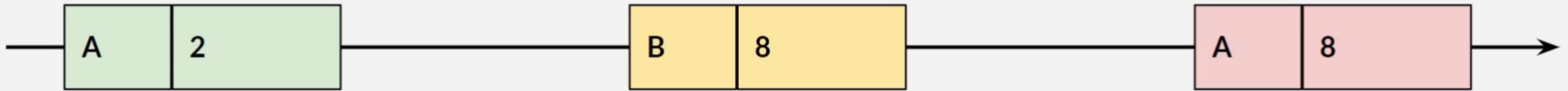


Important Aspects

- × Data locality
 - × Same machine with processing node
 - × No network overhead
 - × No sharing store
- × Fault tolerance
 - × Recover quickly in case application failure
 - × Use changelog topic

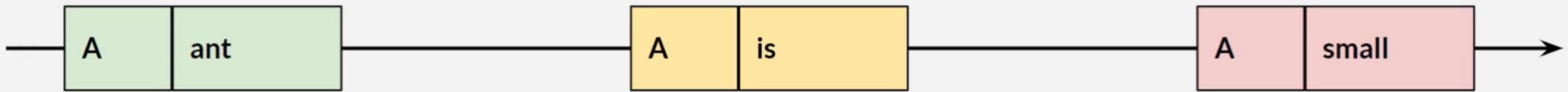


count



- Counts number of record based on key
- Ignore null keys or values

aggregate stream



```
mapValues( (k, v) -> v.length() )
    .groupByKey()
    .aggregate(() -> 0, (aggKey, newValue, aggValue) -> aggValue + newValue)
```



- Aggregate record based on key
- Need initializer and adder
 - Initializer on example : () -> 0
 - Adder on example : (aggKey, newValue, aggValue) -> aggValue + newValue
- Aggregation result can be different type with input
- Ignore null keys
- On first each non-null key : call initializer, then call adder
- On non-null value : call adder

reduce stream

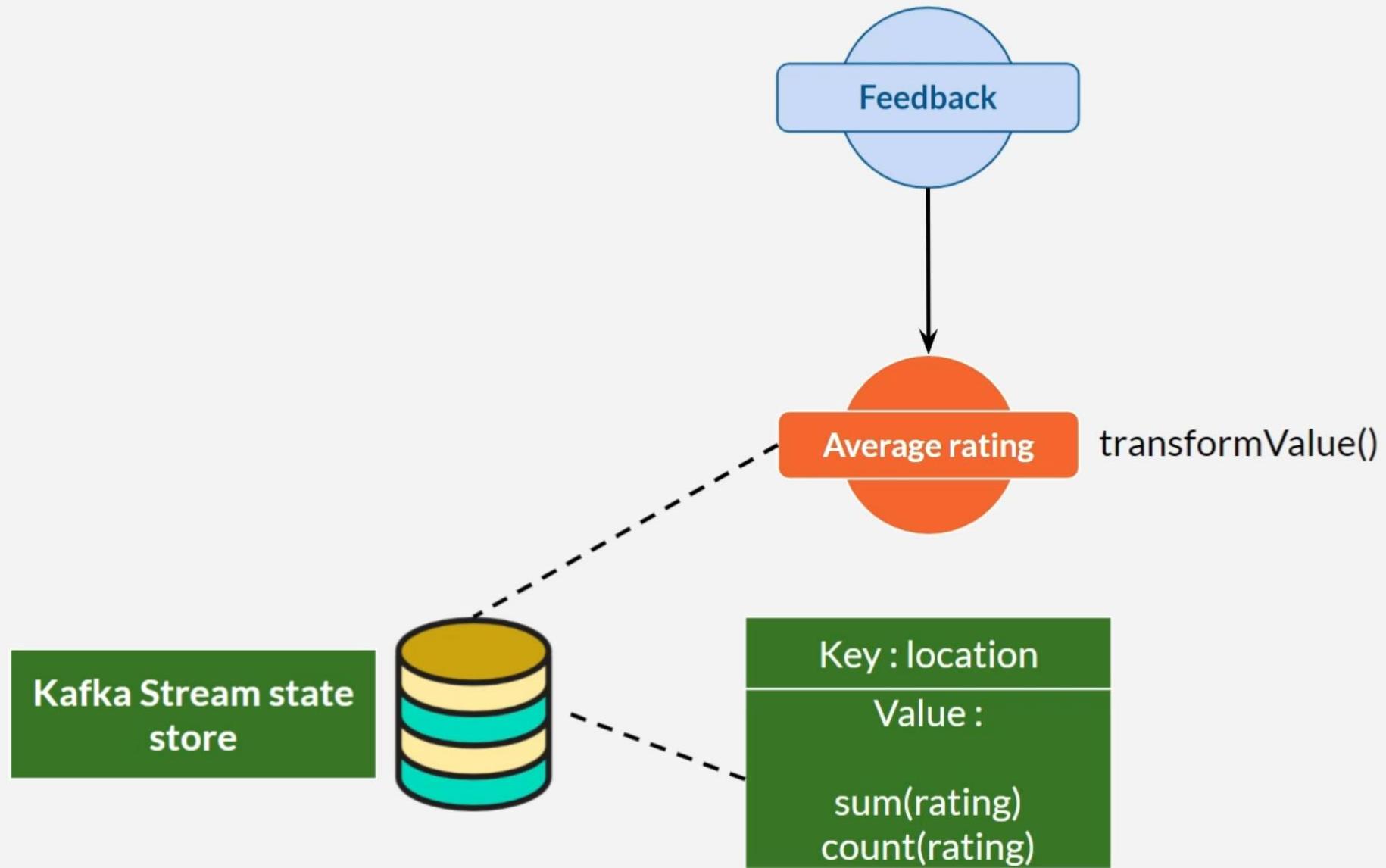


```
mapValues((k, v) -> v.length())
    .groupByKey()
    .reduce(String::concat)
```

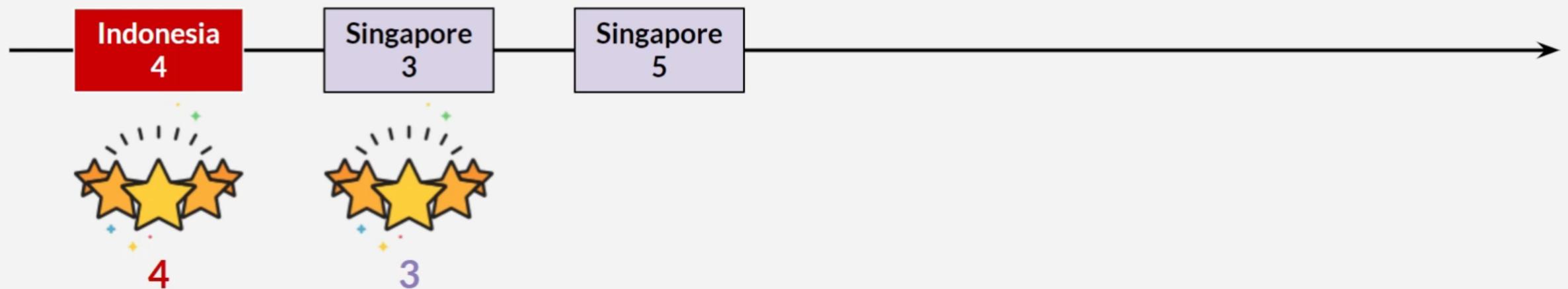


- Specialized form / shorter syntax for aggregate
- Result and input type not change
- Need reducer
 - In example : `String.concat()`

High Level Topology



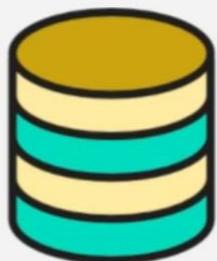
State Store



Indonesia
Sum : 4
Count : 1

Singapore
Sum : 3
Count : 1

State Store



Indonesia
Sum : 12
Count : 4

Singapore
Sum : 8
Count : 2

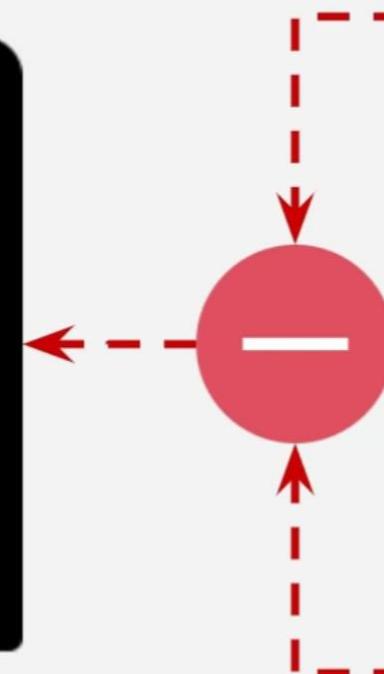
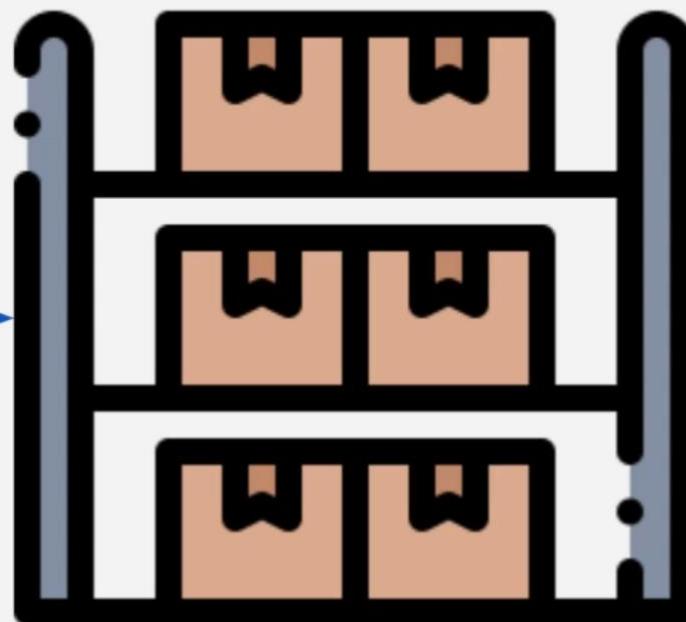
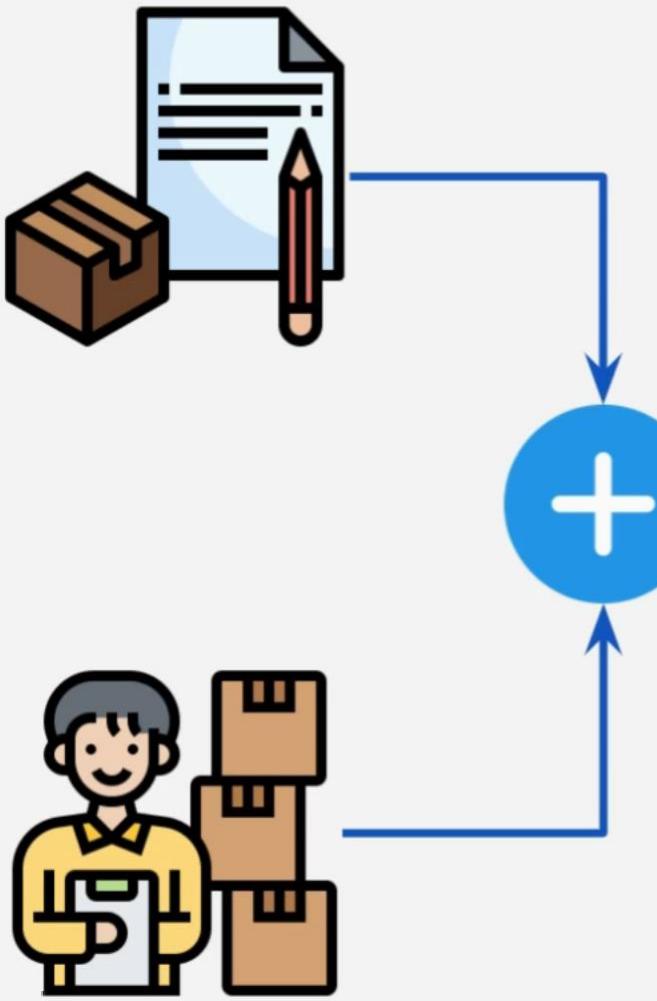
“

Inventory

”

Inventory

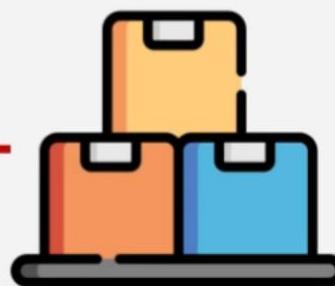
Supplier delivery



Item sold

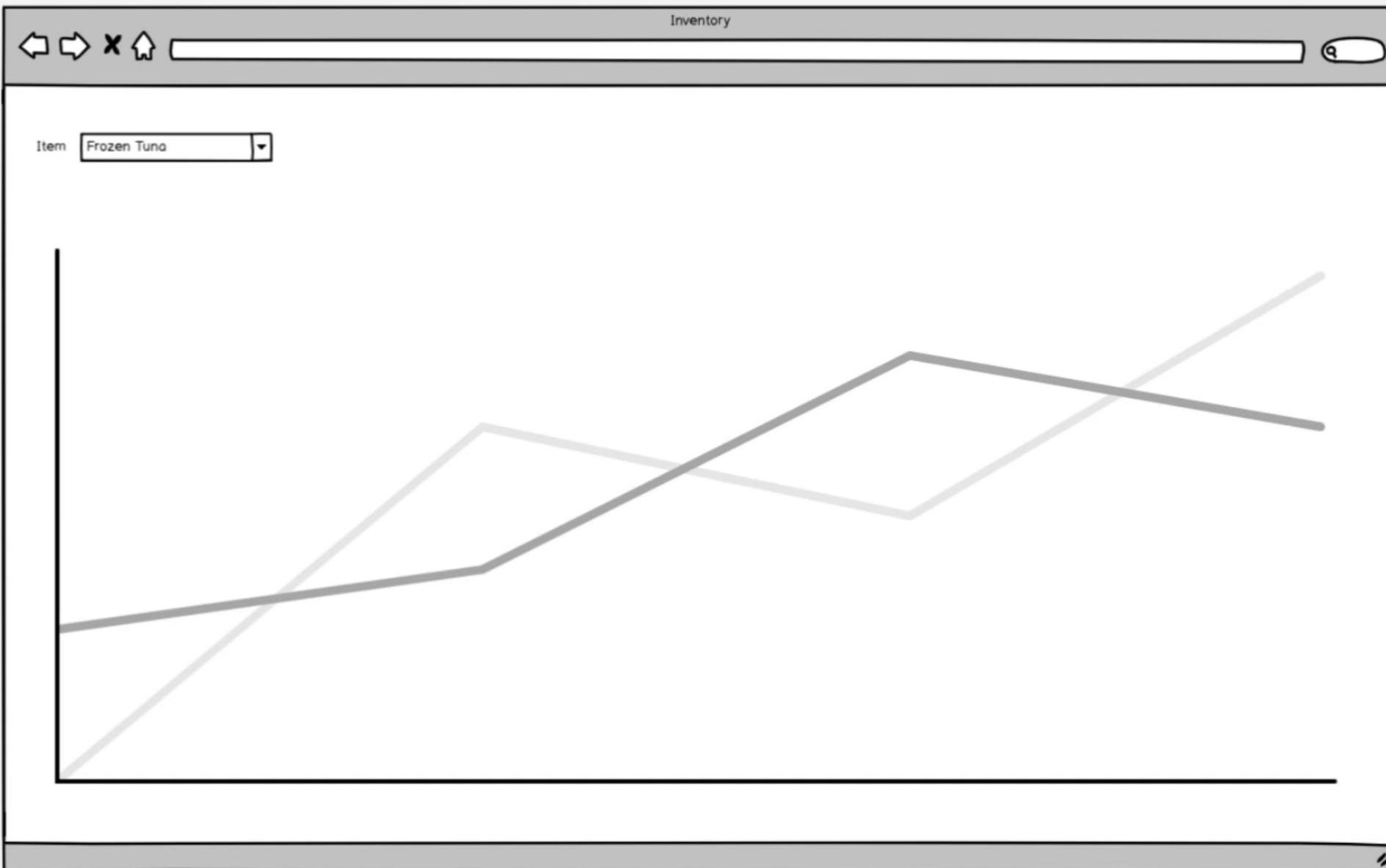


Item refund, etc



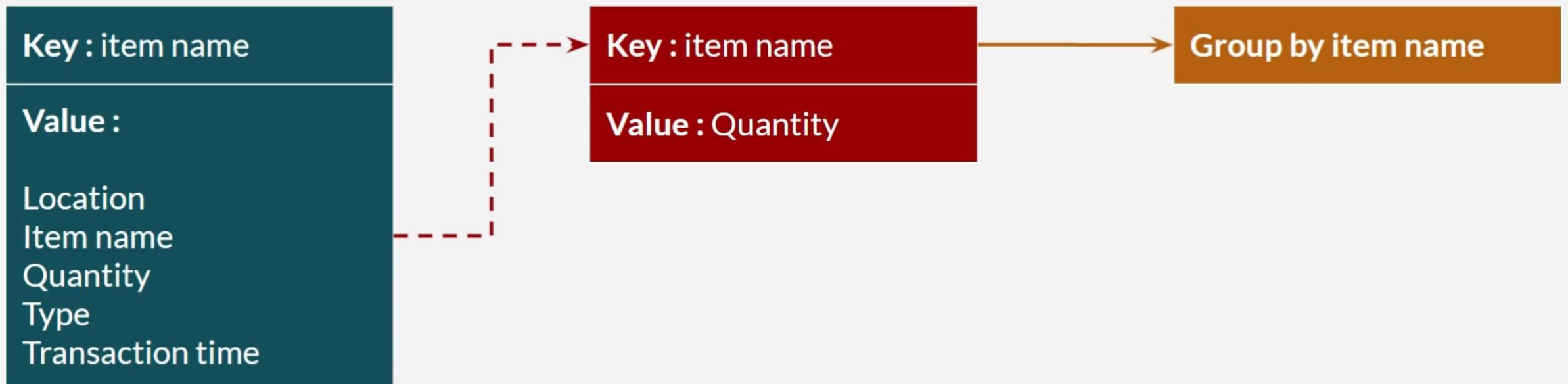
Item moved to store display,

Inventory Chart



- × Kafka stream aggregation (sum)
- × Inventory message sent
 - × Key : item name
 - × Value : item name, location, quantity, type
 - × Value (optional) : transaction time
- × No validation, e.g:
 - × Negative amount
 - × Inventory quantity still enough for removal

Grouping & Aggregate



Grouping & Aggregate

Record #	Record after map values	Group by key	KGroupedStream aggregate		Aggregated
			Initializer	Adder	
1	(candy, 4)	(candy, 4)	0 (for candy)	(candy, 0 + 4)	(candy, 4)
2	(apple, 2)	(apple, 2)	0 (for apple)	(apple, 0 + 2)	(apple, 2)
3	(candy, 3)	(candy, 3)		(candy, 4 + 3)	(candy, 7)
4	(candy, 1)	(candy, 1)		(candy, 7 + 1)	(candy, 8)
5	(apple, 4)	(apple, 4)		(apple, 2 + 4)	(apple, 6)
6	(apple, 5)	(apple, 5)		(apple, 6 + 5)	(apple, 11)

“

Joining Order and Payment Stream

”

Kafka Stream Join

- × Taking a two stream / table and make new stream / table out of them
- × Join based on record key
- × Join on:
 - × KStream & KStream
 - × KTable & KTable
 - × KStream & KTable
 - × KStream & GlobalKTable
- × Join types : Inner, left, outer

Join Combinations

Left Type	Right Type	Co-Partition	Windowed	Join Types		
				Inner Join	Left Join	Outer Join
KStream	KStream	Required	Yes	Supported	Supported	Supported
KTable	KTable	Required	No	Supported	Supported	Supported
KStream	KTable	Required	No	Supported	Supported	Not Available
KStream	GlobalKTable	Not required	No	Supported	Supported	Not Available

Co-Partition

- × Combinations:
 - × KStream / KStream
 - × KTable / KTable
 - × KStream / KTable
- × Must be co-partitioned
- × Otherwise, **runtime error**

Co-partition

- × Left type must have same number of partition with right type
 - × Left : x & Right : x = co-partitioned
 - × Left : x & Right : y = Not co-partitioned
- × Same partitioning strategy
 - × Key with same value must go same partition
 - × If use custom logic to send to partition n , both must be the same

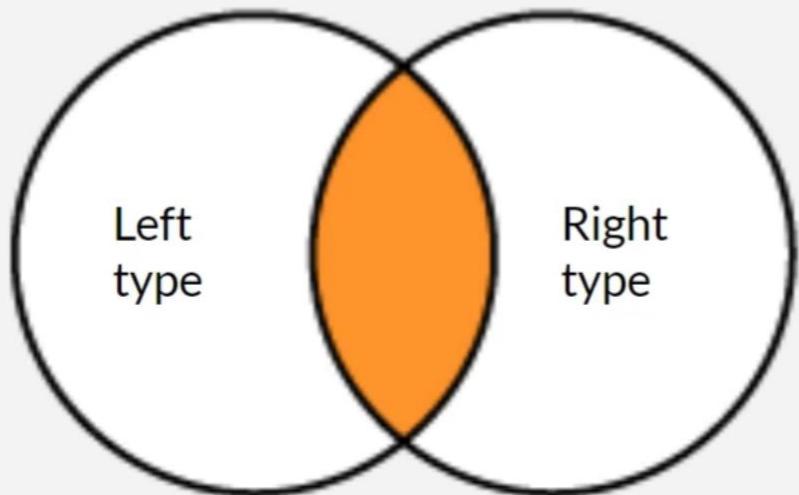
GlobalKTable

- × No need co-partition
- × GlobalKTable:
 - × Fetch / copy all data from all partitions of certain topic
 - × Live in kafka stream instance
 - × Consumes memory / disk on kafka stream instance
 - × Relatively small data & almost static
- × Can join stream to GlobalKTable, although different partition number

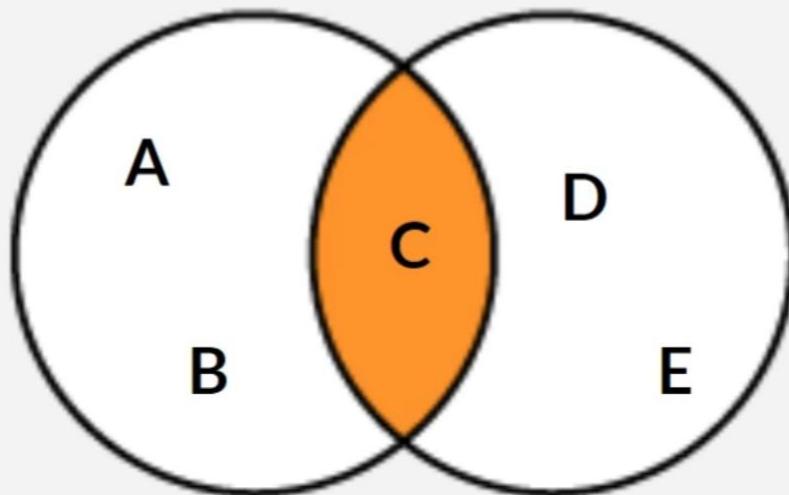
Windowed KStream-KStream

- ✗ KStream : unbounded, unlimited data
- ✗ KStream / KStream join need constraints
- ✗ Otherwise both KStream will be scanned when new data arrives -> **huge performance cost**
- ✗ Use window for constraint

Inner Join



example →

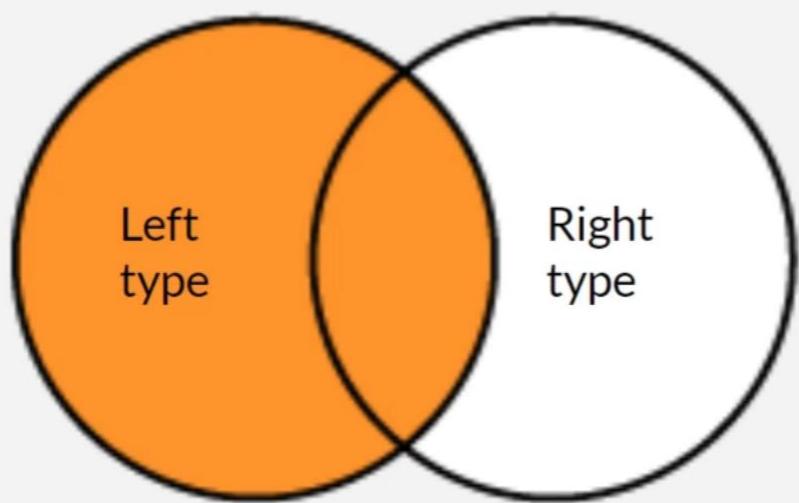


(left value, right value)

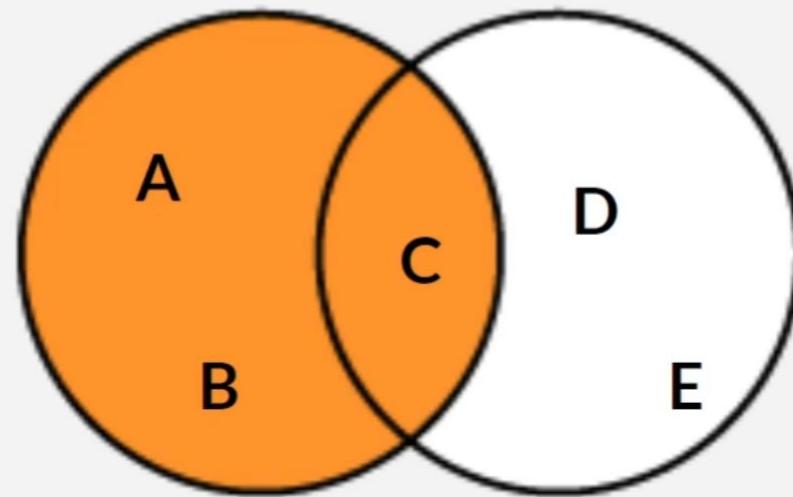
Result is:
(C, C)

Left (Primary)	Right (Secondary)
KStream	KStream
KTable	KTable
KStream	KStream
KStream	GlobalKTable

Left Join



example →

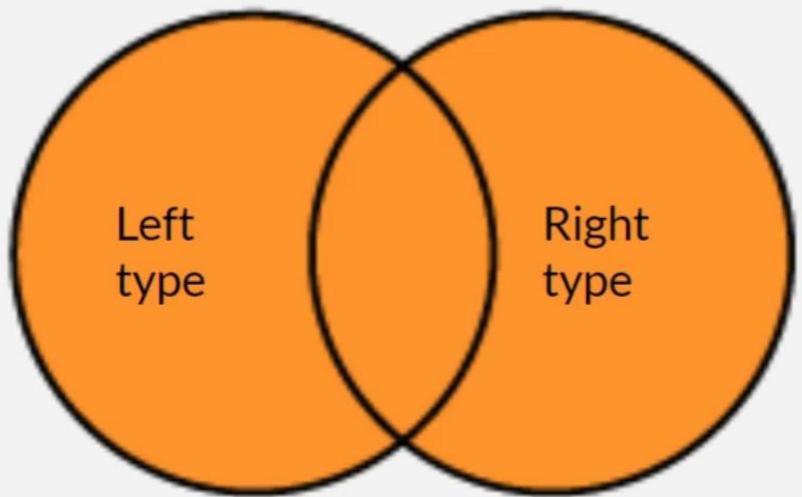


Left (Primary)	Right (Secondary)
KStream	KStream
KTable	KTable
KStream	KStream
KStream	GlobalKTable

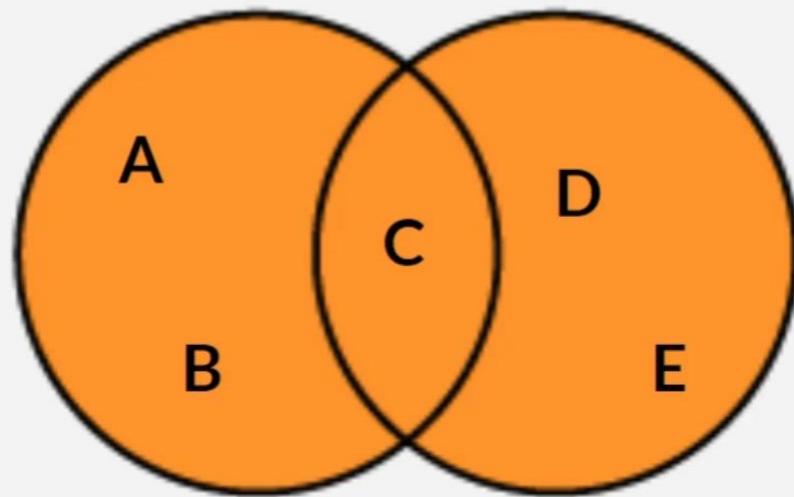
(left value, right value)

Result are:
(A, null)
(B, null)
(C, C)

Outer Join



example →



Left (Primary)	Right (Secondary)
KStream	KStream
KTable	KTable

(right value, left value)

Result are:

(A, null)

(B, null)

(C, C)

(null, D)

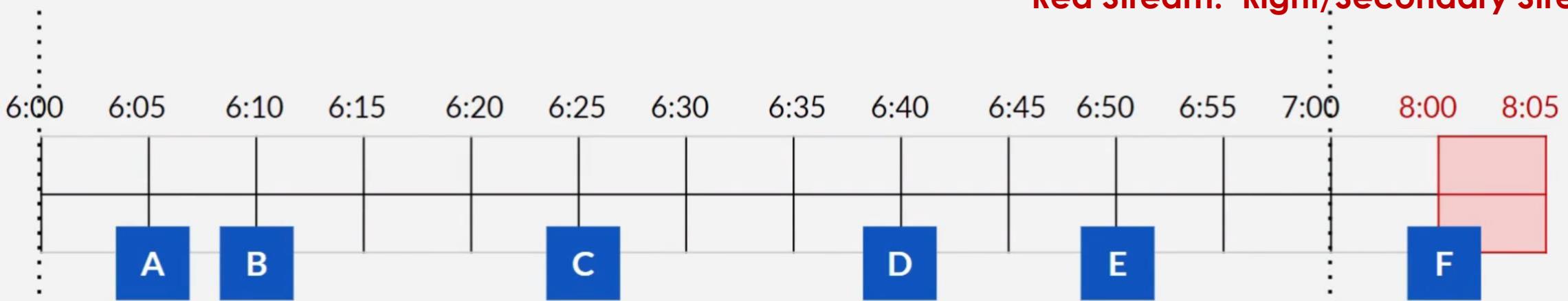
(null, E)

Join Window

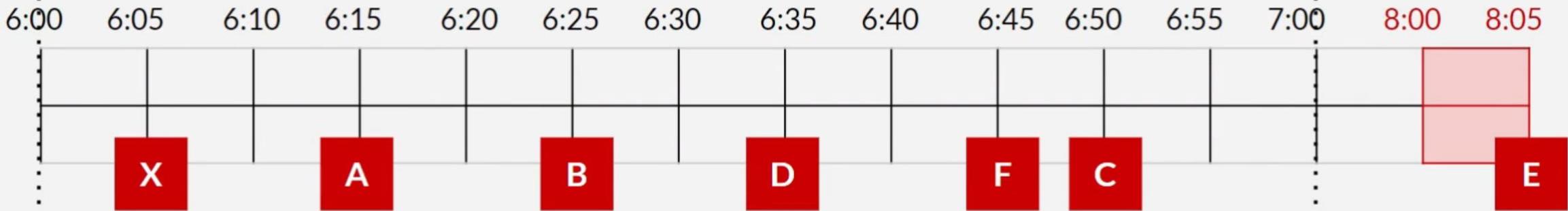
- × Class `JoinWindows`
- × Define maximum time difference
- × Record key K arrives on left stream at 7:15
 - × Matching record on right stream between 7:15 - 8:15
 - × 7:16 : match
 - × 7:50 : match
 - × 8:12 : match
 - × 8:17 : not match

Input Streams

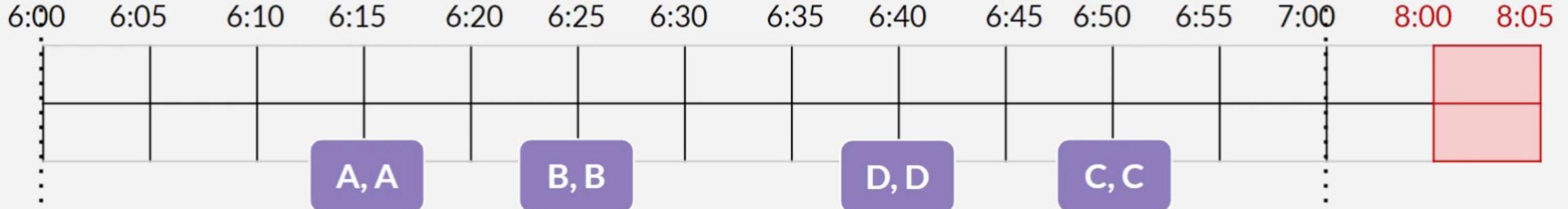
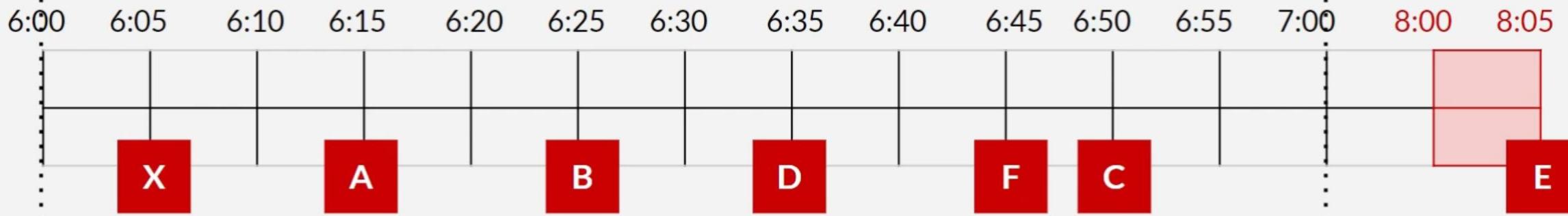
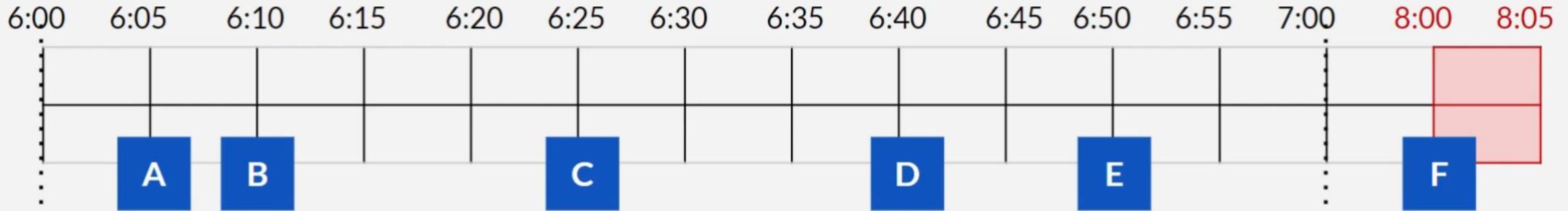
Blue Stream: Left/Primary Stream
Red Stream: Right/Secondary Stream



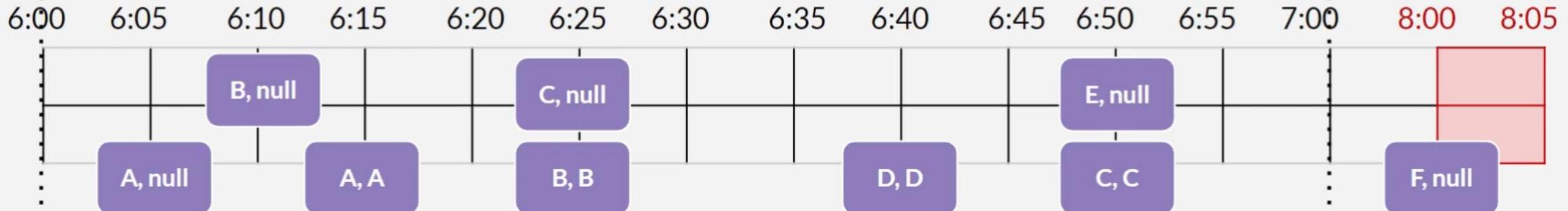
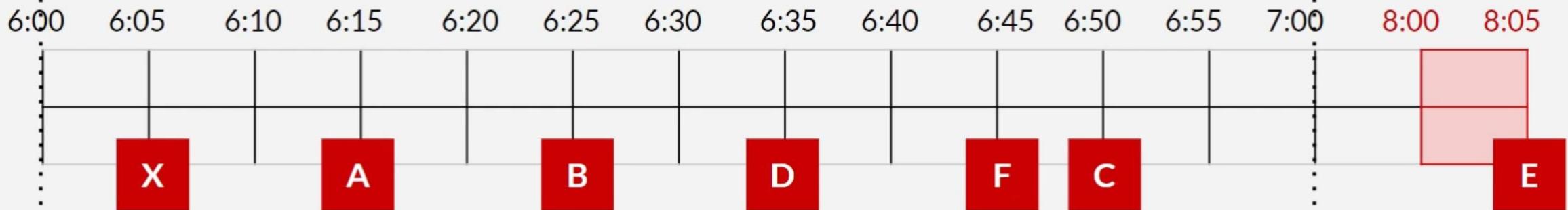
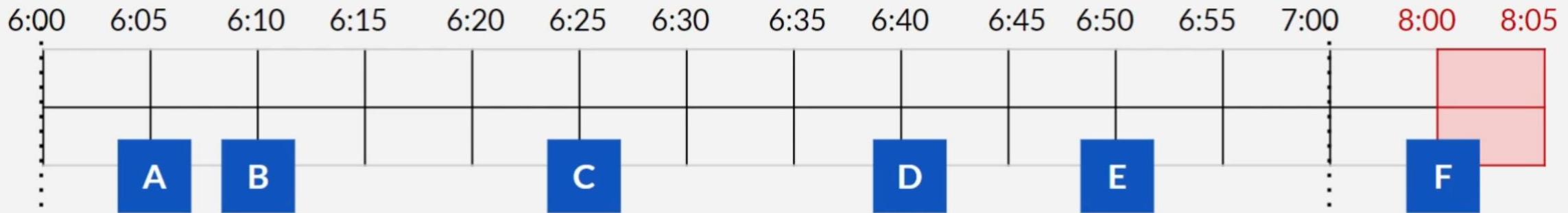
Time window : 1 hour



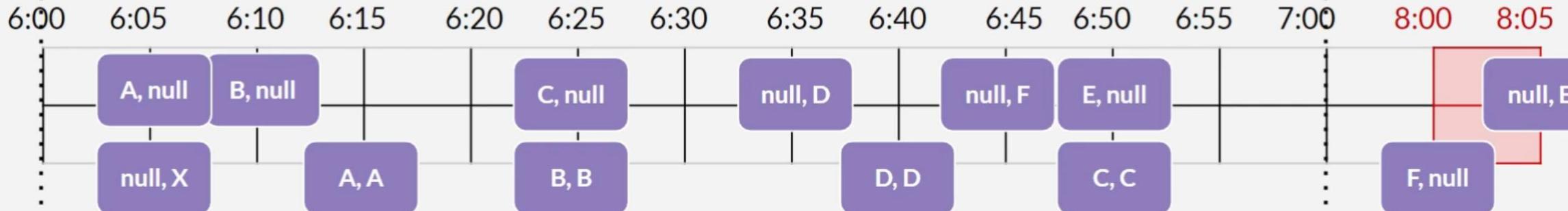
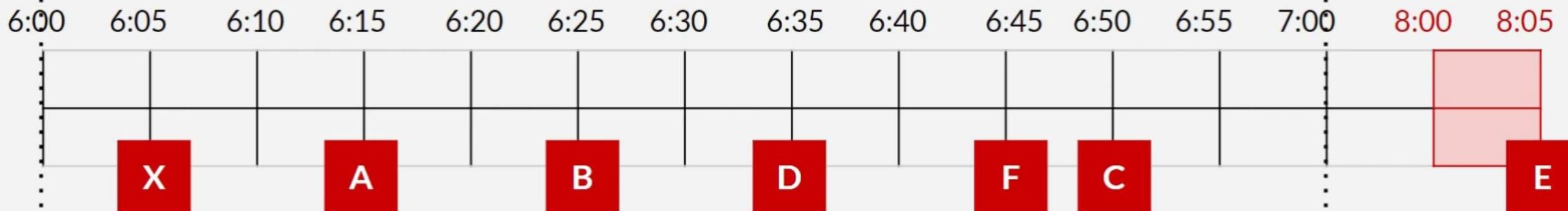
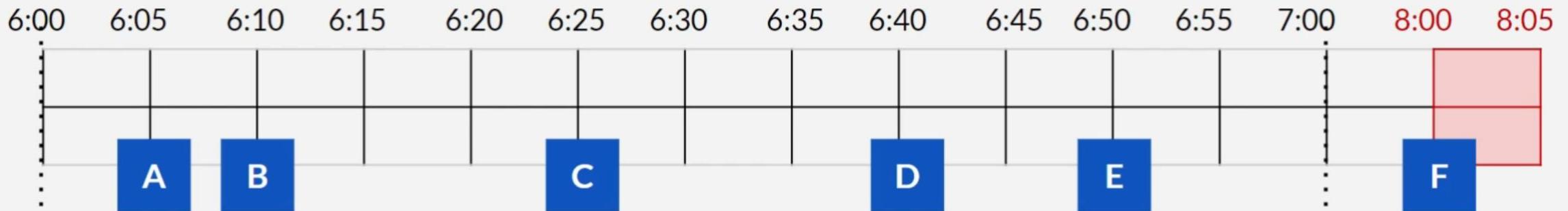
Inner Join Stream-Stream



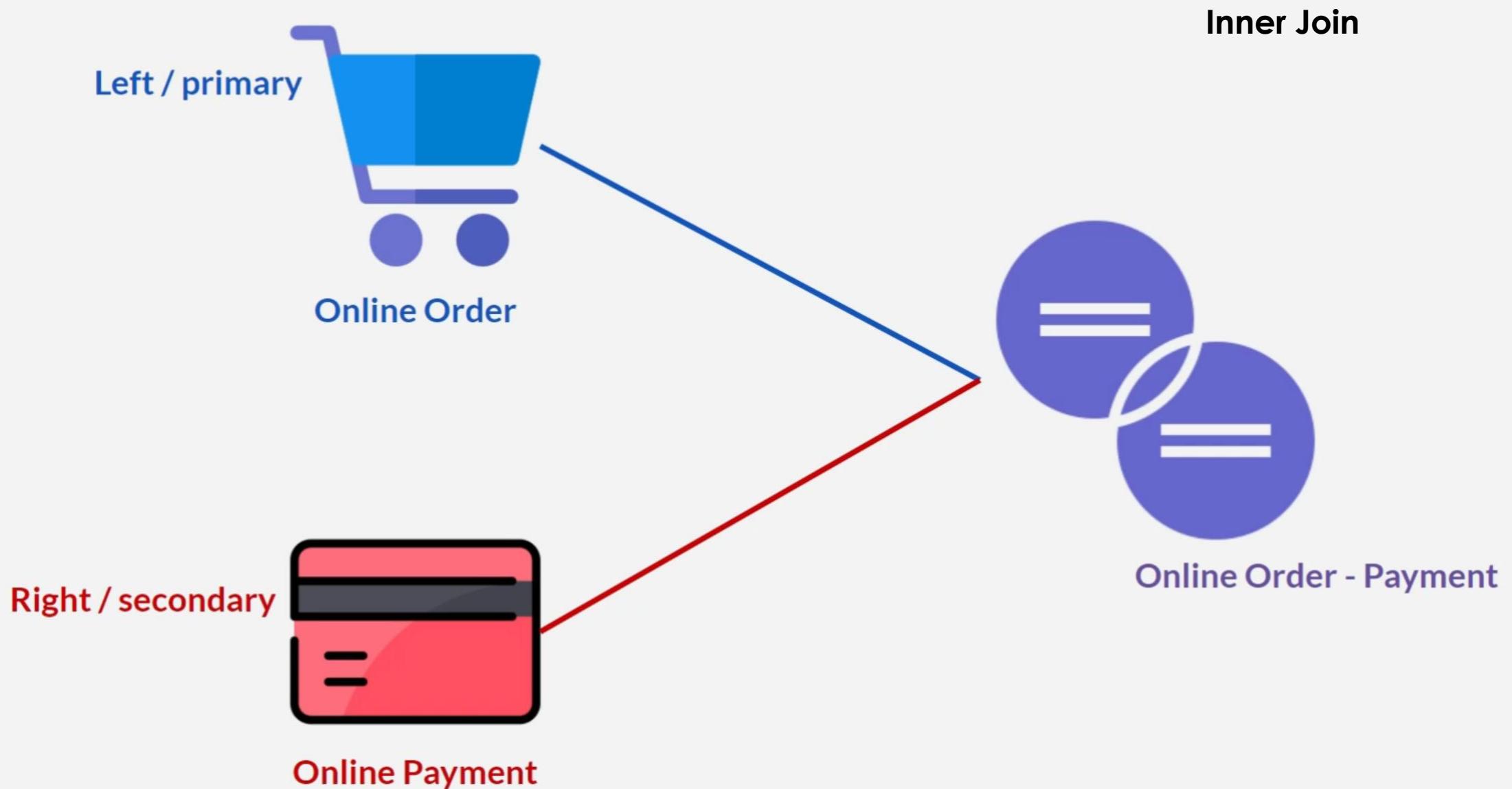
Left Join Stream-Stream



Outer Join Stream-Stream



Online Order & Payment



Join Class

OnlineOrderMessage.java

1st class : left source stream

```
int a  
Int b
```



OnlinePaymentMessage.java

2nd class : right source stream

```
double x
```

OnlineOrderPaymentMessage.java

3rd class : join class

```
int a  
int b  
double x
```

join class example 2

```
int a  
double x
```

join class example 3

```
int a  
String someField  
double aPlusX()
```

Join Syntax

```
leftStream.join(rightStream, joiner, windows, joined)
leftStream.leftJoin(rightStream, joiner, windows, joined)
leftStream.outerJoin(rightStream, joiner, windows, joined)
```

```
// joiner

private JoinClass joiner(LeftClass left, RightClass right) { ... }
```

```
// windows

JoinWindows.of(...)
```

```
// joined

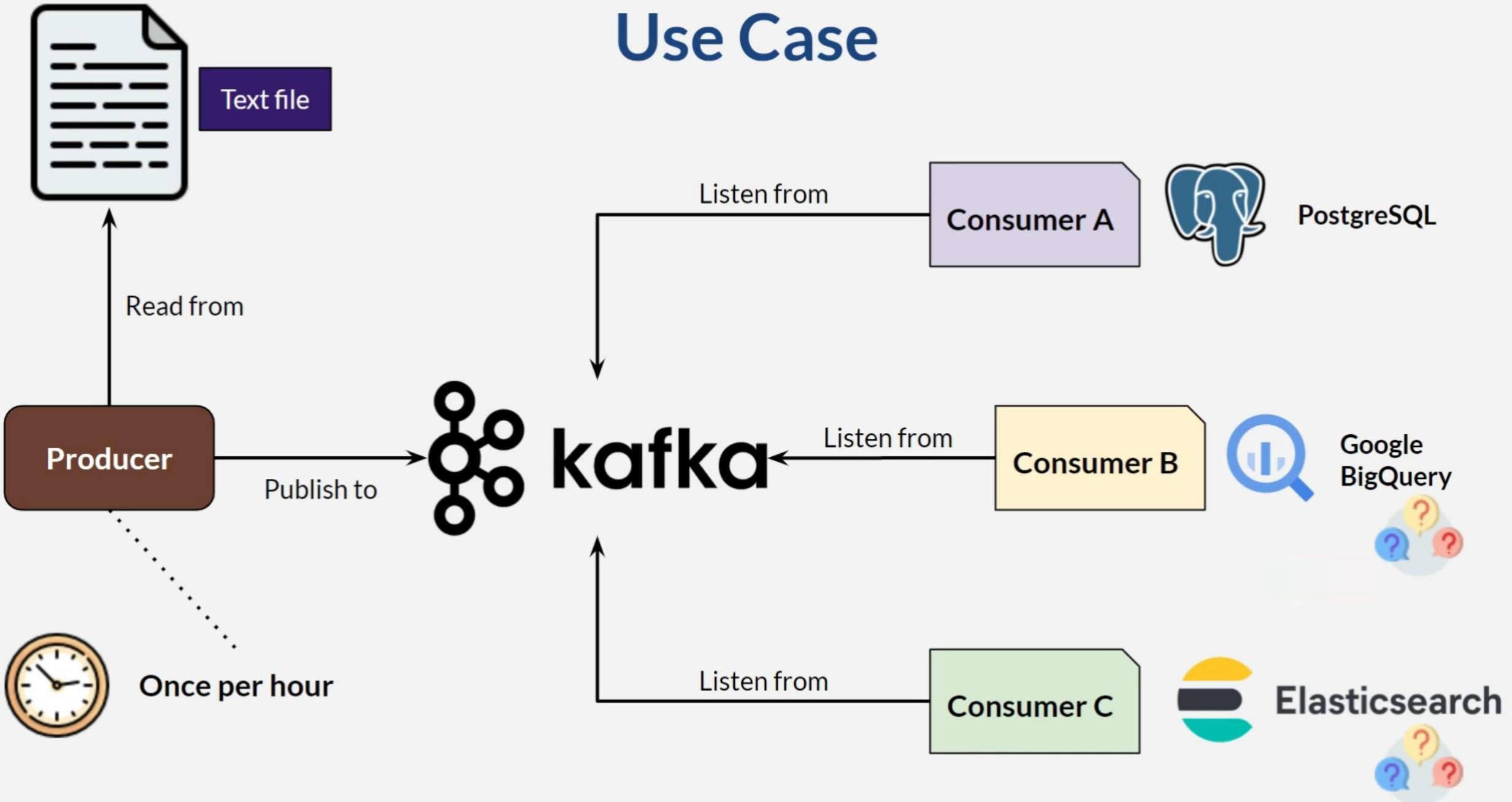
Joined.with(keySerde, leftSerde, rightSerde)
```

“

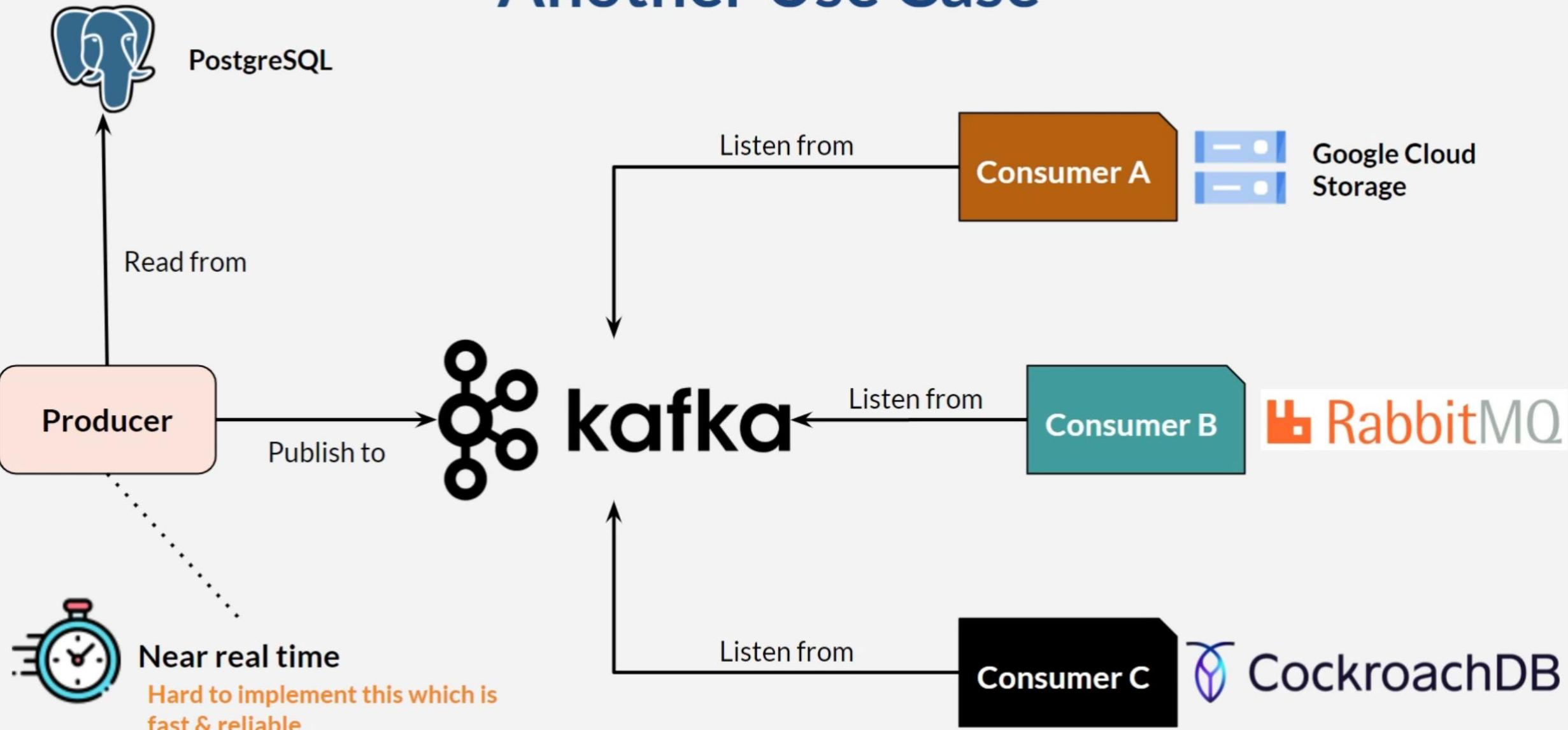
Kafka Connect

”

Use Case



Another Use Case

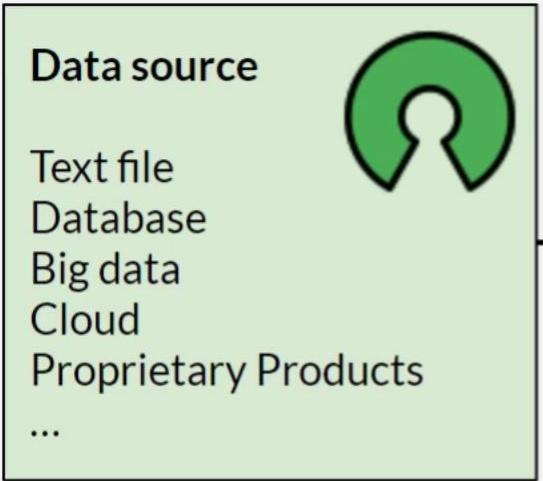


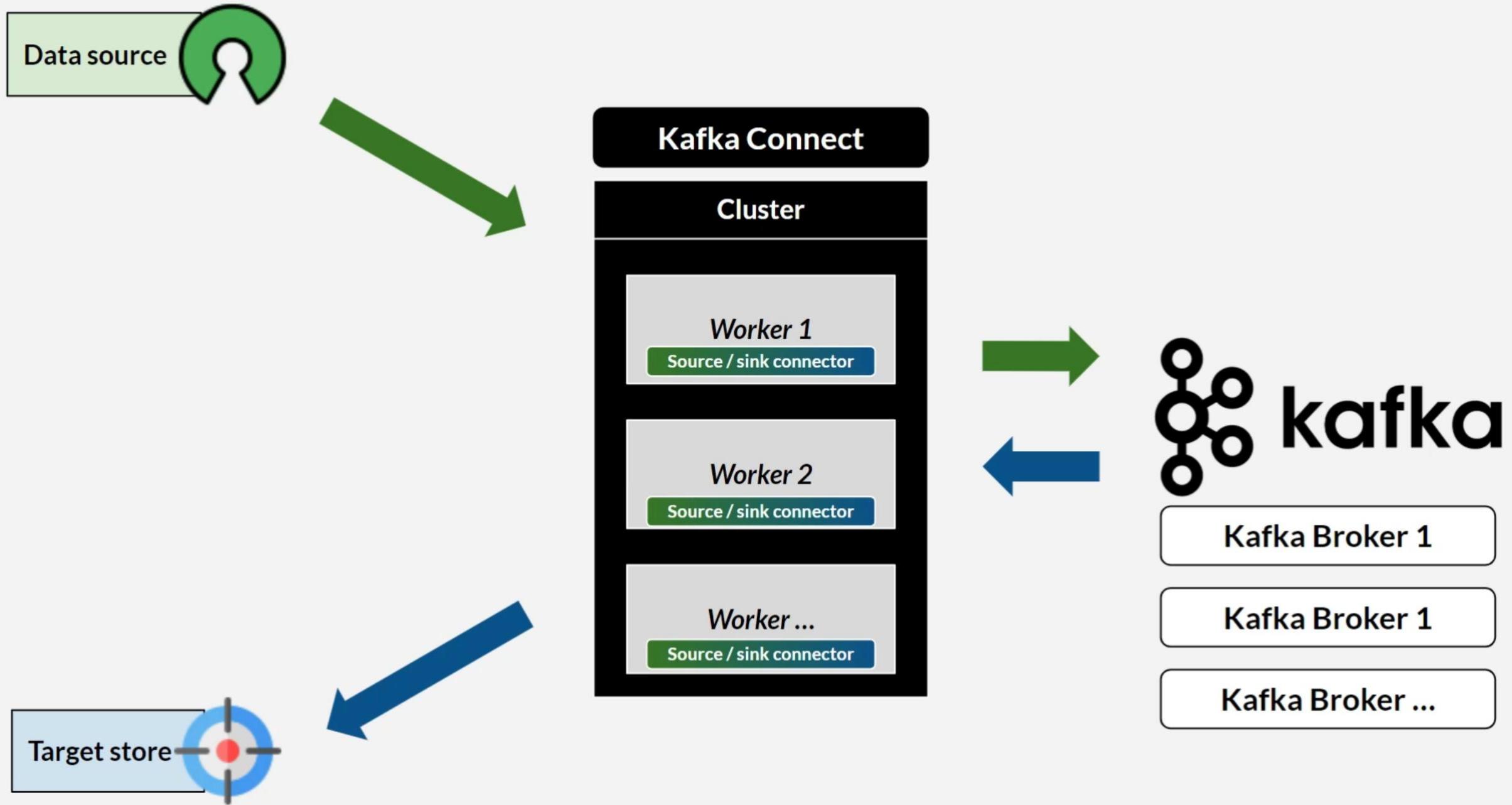
Message In, Message Out

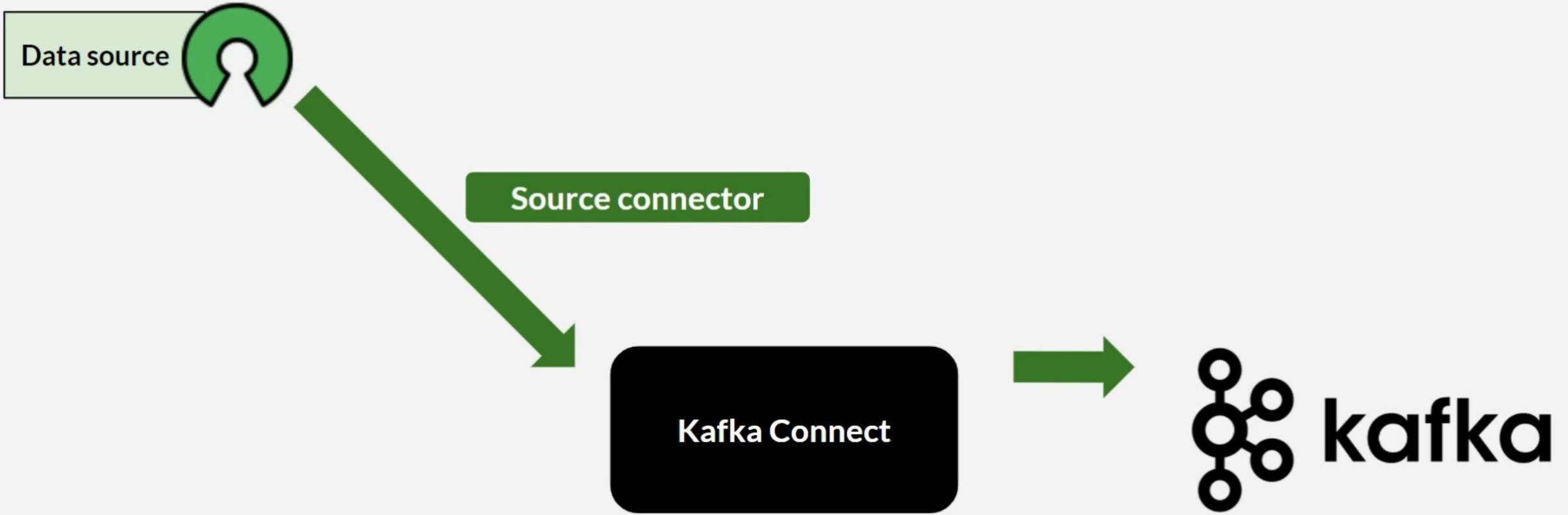
- × Happens in real life
- × Today we have multiple data sources (**in**)
 - × text file, relational databases
 - × Non relational database, email, cloud, social media, specific software, etc
 - × Can be more than 100
- × Multiple target data store (**out**)
 - × local text file, FTP, relational / non-relational databases, big data, cloud, etc

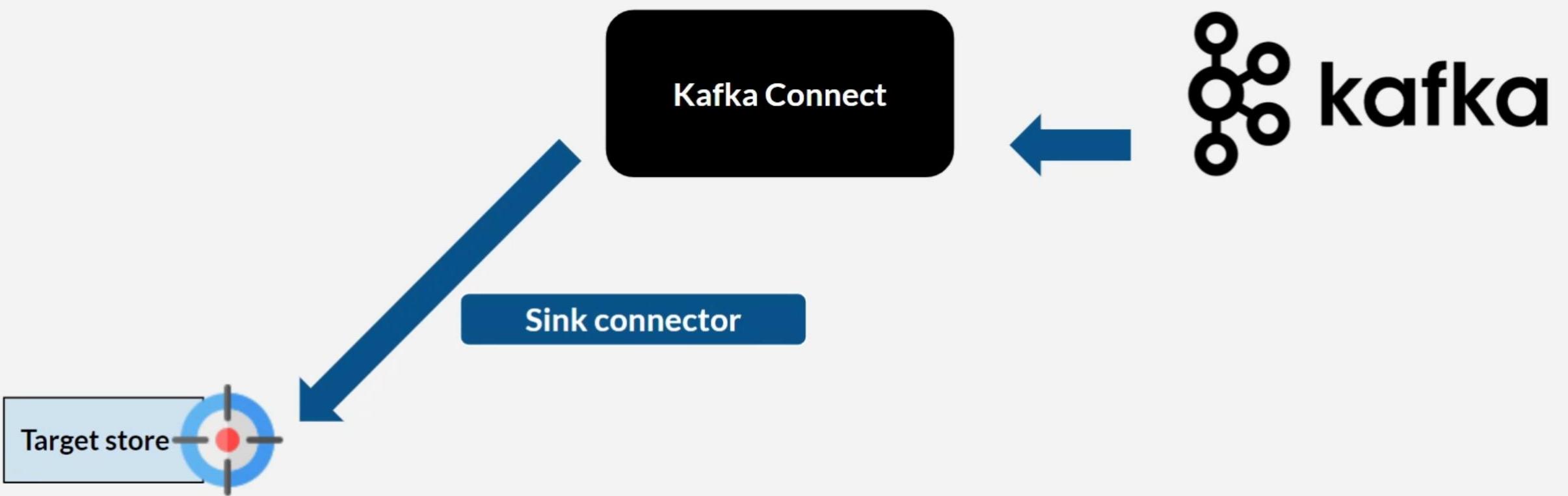
On Kafka

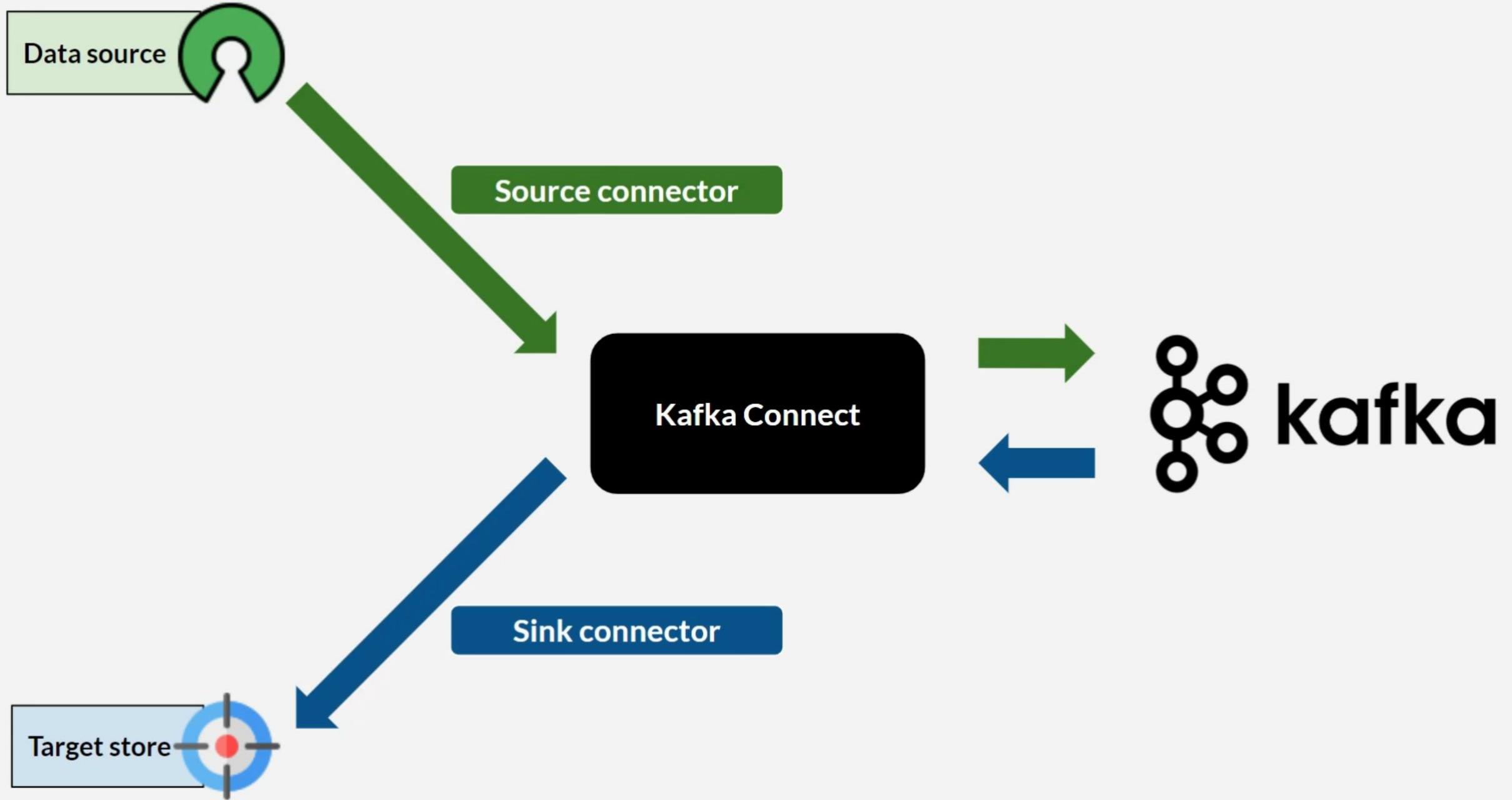
- × Write a lot of producers and consumers
- × Extra time & effort for performance & reliability
- × Good news : ***you don't have to write your own!***
- × People / companies already wrote them
- × Read plugin : from non-kafka into kafka
- × Write plugin : from kafka into non-kafka
- × **Kafka Connect**











Kafka Connect

- × Additional platform for kafka
- × Data integration
- × Transfer data between Kafka - non kafka
- × Horizontally scalable & fault tolerant
- × Uses connectors for interact with kafka server

Connectors

- × Java jar file
- × Plugin for kafka connect
- × Interface between kafka and non-kafka
- × **Source** connector : read (ingest) into kafka (**producer**)
- × **Sink** connector : write from kafka to non-kafka (**consumer**)
- × Install connector for specific need
- × Write configuration (json)
- × Declarative configuration
- × Fast & reliable

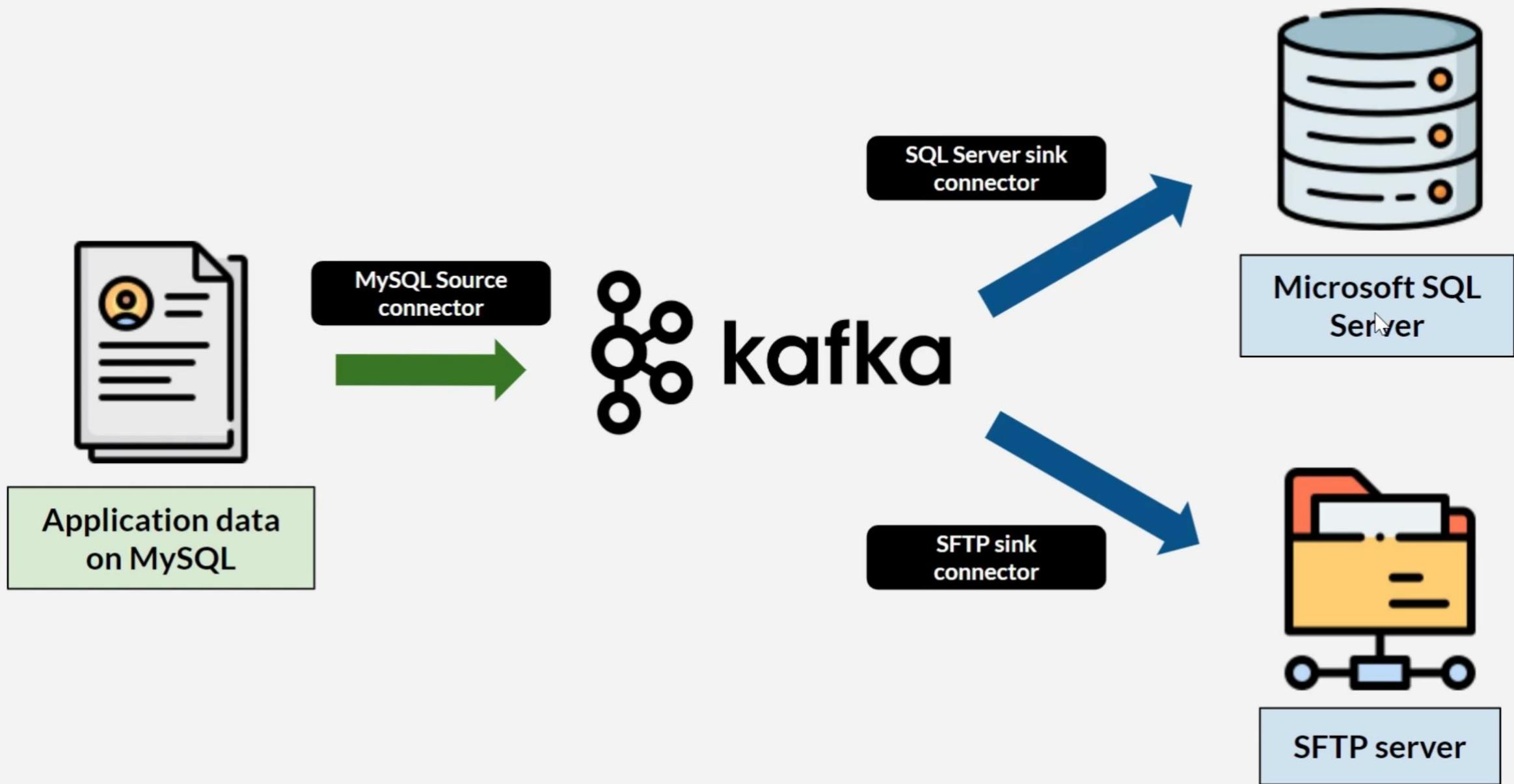
Connectors

- ✗ A lot of source / sink connectors
- ✗ Shorten time and effort

How to Get Connectors?

- × Curated list on [confluent.io/hub](https://www.confluent.io/hub)
- × Google : *kafka source / sink connector for xxx*
- × Build your own

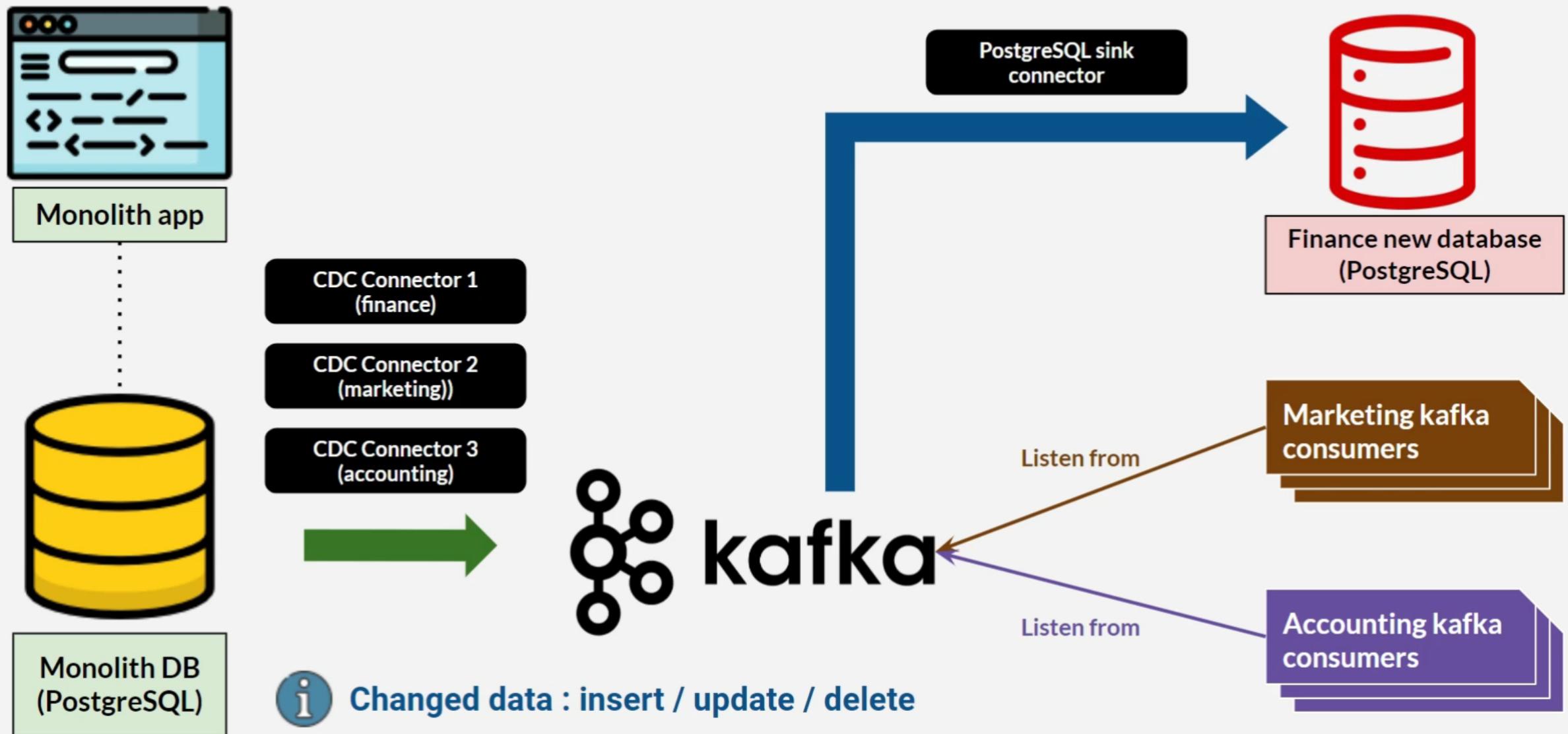
Use Case : Write to Data Stores



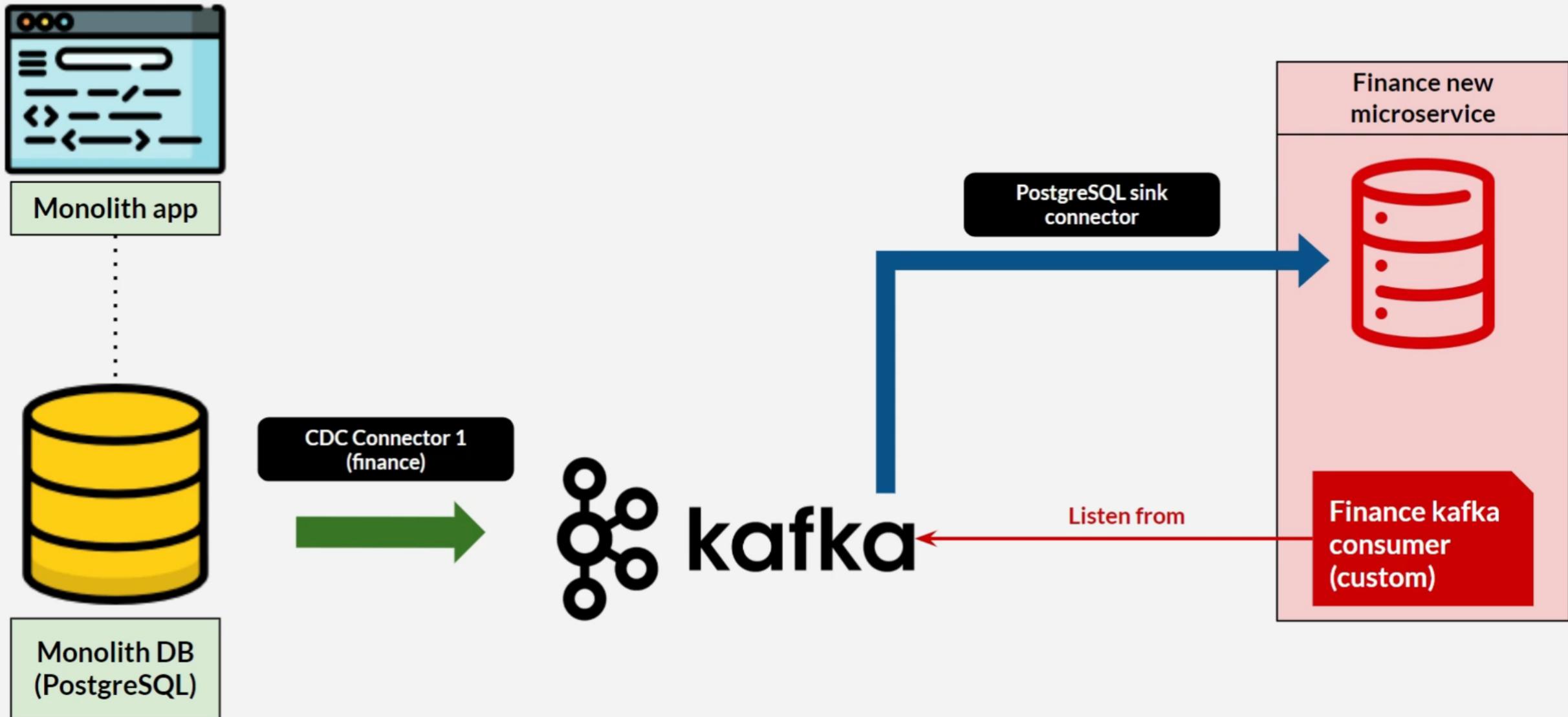
Use Case : Modernize Legacy System

- × Modernize legacy monolith into microservices
- × Modernization is hard and long
- × Legacy system still needs to run during modernization
- × Modernize functionalities part by part
- × Use kafka connect CDC (Change Data Capture) connectors

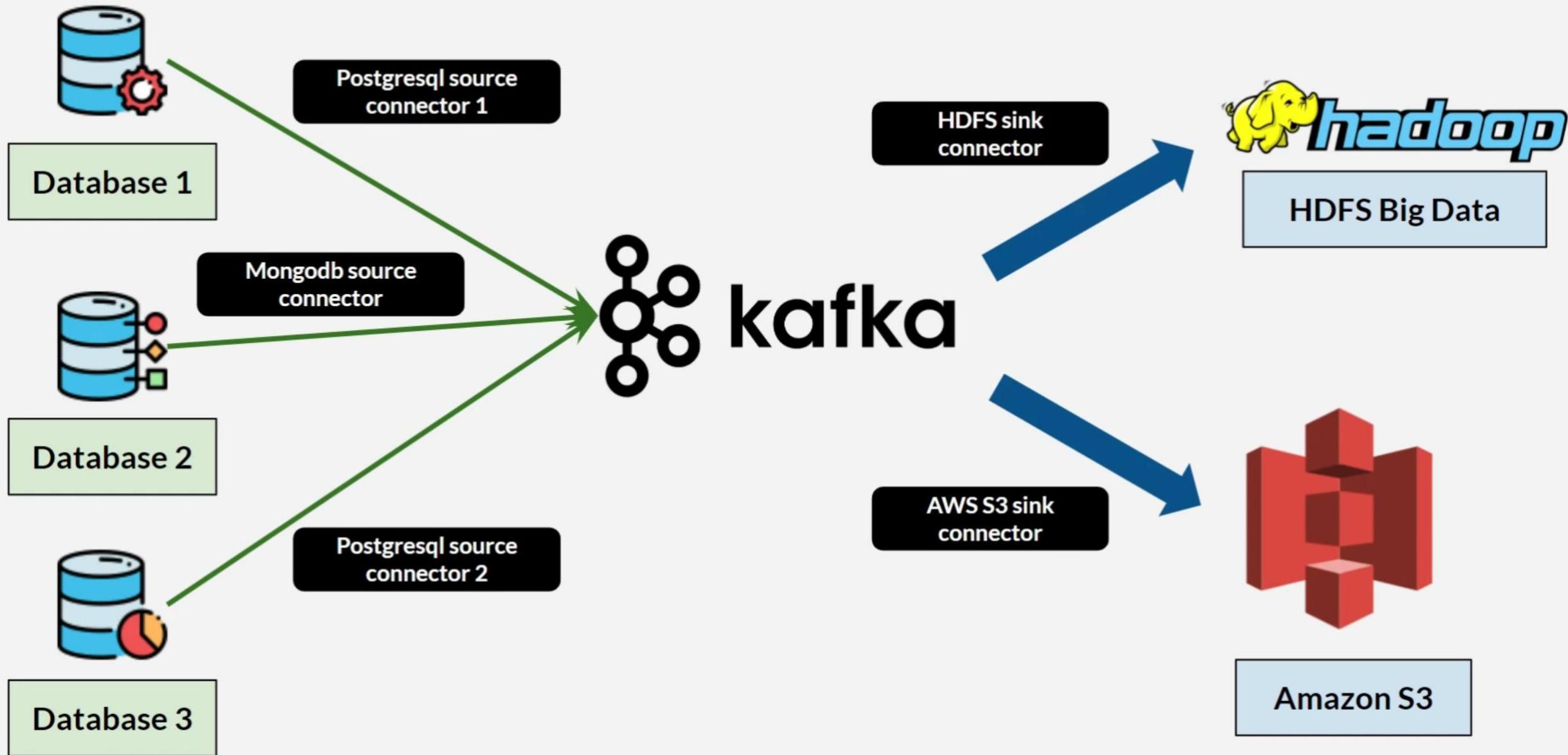
Use Case : Modernize Legacy System



Use Case : Modernize Legacy System



Use Case : Data Engineering ETL Pipeline



```
#> docker-compose -f [script-file] -p [project] down
```

Part	What to run (in sequence)
1 - core kafka	#> docker-compose -f docker-compose-core.yml -p core up -d
2 - kafka connect	#> docker-compose -f docker-compose-core.yml -p core down #> docker-compose -f docker-compose-connect.yml -p connect up -d #> docker-compose -f docker-compose-connect-sample.yml -p connect-sample up -d



“

CONDUCTOR

Kafka UI

”

Use docker-compose-full.yml

[Download Conductor:](https://www.conduktor.io/download/) <https://www.conduktor.io/download/>

[Official Documentation:](https://www.conduktor.io/kafka) <https://www.conduktor.io/kafka>



Hello, Hamsika B!
On trial until 2022-07-10

My Account

Are you connecting to an existing cluster?



+ New Kafka Cluster

Learn more →

Don't have a Kafka cluster? Let's start one for you!

+ Start local Kafka cluster

Learn more →

Welcome!

Conduktor helps you to
manage your Kafka clusters.
Start by adding one and let's rock on!

Discover Kafkademy



Hello, Hamsika B!
On trial until 2022-07-10

My Account

Cluster Configuration

Add new Cluster

[Documentation](#)

[Ask us Anything](#)

Kafka Cluster

Schema Registry

Kafka Connect

ksqlDB (beta)

Metrics

SSH

Plugins

Conduktor can connect to your Kafka cluster *only* if your computer can also access the Kafka servers.

Using SASL or SSL? AWS or Docker? Check out our documentation if you have issues.

Integration



aiven



confluent



Red Hat



Github

Welcome!

Conduktor helps you to
manage your Kafka clusters.
Start by adding one and let's rock on!

Cluster Name *

Bootstrap Servers *

Zookeeper

Additional Properties

Color

None

Logs

Test Kafka Connectivity

Test the connection

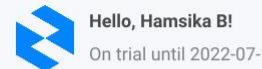
DUPLICATE

RESET

CANCEL

SAVE

Discover Kafkademy



My Account

Cluster Configuration

Add new Cluster

[Documentation](#)[Ask us Anything](#)

Kafka Cluster

Schema Registry

Kafka Connect

ksqlDB (beta)

Metrics

SSH

Plugins

URL

http://localhost:8081

Security

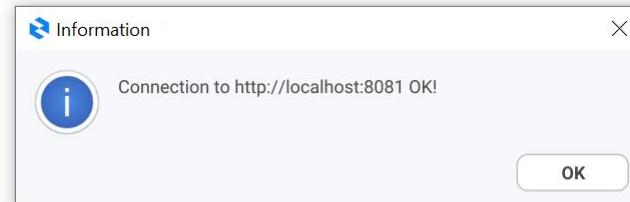
None

Basic Auth

Bearer Token

Additional Properties

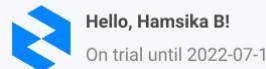
foo=bar



Welcome!

Conduktor helps you to
manage your Kafka clusters.
Start by adding one and let's rock on!

[Test Schema Registry Connectivity](#)[DUPLICATE](#)[RESET](#)[CANCEL](#)[SAVE](#)[Discover Kafkademy](#)



Cluster Configuration

Add new Cluster

[Documentation](#)[Ask us Anything](#)

Kafka Cluster

Schema Registry

Kafka Connect

ksqlDB (beta)

Metrics

SSH

Plugins

My Kafka Connect

Name

My Kafka Connect

URL

http://localhost:8083

HTTP Headers ?

eg: X-App-Source=Conduktor, X-API-Token=abc

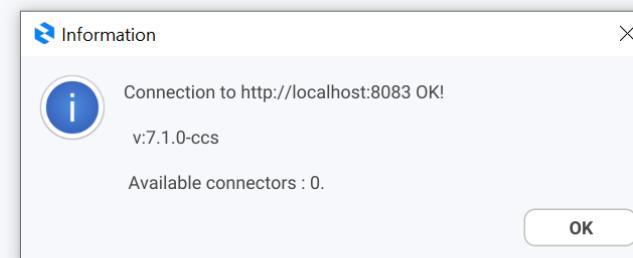
Security

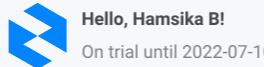
None

Basic Auth

Bearer Token

▶ HTTPS Configuration

[Test Kafka Connect Connectivity](#)[+ ADD](#)[DUPLICATE](#)[RESET](#)[CANCEL](#)[SAVE](#)[Discover Kafkademy](#)



Cluster Configuration

Add new Cluster

[Documentation](#)[Ask us Anything](#)

Kafka Cluster

Schema Registry

Kafka Connect

ksqlDB (beta)

Metrics

SSH

Plugins

My ksqlDB

Name

My ksqlDB

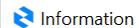
URL

http://localhost:8088

HTTP Headers ?

eg: X-App-Source=Conduktor, X-API-Token=abc

Security



Information

Connection to http://localhost:8088 OK! (Healthy: true)



OK

[Test ksqlDB Connectivity](#)

▶ HTTPS Co

We're looking for ksqlDB feedback (0.10+).

We still have many features and improvements planned, but don't hesitate to contact us. Thanks!

[+ ADD](#)[DUPLICATE](#)[RESET](#)[CANCEL](#)[SAVE](#)[Discover Kafkademy](#)



Cluster Configuration

Add new Cluster

[Documentation](#)[Ask us Anything](#)[Kafka Cluster](#)[Schema Registry](#)[Kafka Connect](#)[ksqlDB \(beta\)](#)[Metrics](#)[SSH](#)[Plugins](#)

My ksqlDB

Name

My ksqlDB

URL

http://localhost:8088

HTTP Headers ?

eg: X-App-Source=Conduktor, X-API-Token=abc

Security

None

Basic Auth

[Using Confluent Cloud?](#)[Test ksqlDB Connectivity](#)

► HTTPS Configuration

We're looking for ksqlDB feedback (0.10+).

We still have many features and improvements planned, but don't hesitate to contact us. Thanks!

[+ ADD](#)[DUPLICATE](#)[RESET](#)[CANCEL](#)[SAVE](#)[Discover Kafkademy](#)

My Local Cluster
localhost:9092
PLAINTEXT



Hello, Hamsika B!
On trial until 2022-07-10

My Account



Are you connecting to an existing cluster?

+ New Kafka Cluster

Learn more →

Don't have a Kafka cluster? Let's start one for you!

+ Start local Kafka cluster

Learn more →

Discover Kafkademy



 My Local Cluster[Overview](#)[Brokers](#)[Topics](#)[Consumers](#)[Schema Registry](#)[Kafka Connect](#)[Kafka Streams](#)[ksqlDB](#)[Security](#)

TOPICS

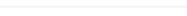
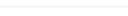
8
Topics134
Partitions0
URP0
No Leader0
< Min ISR[+ CREATE](#)[Check Reassignments](#)[Eye](#) [Download](#) [🔍](#)

Topic Name	Partitions	Count	Size	Consumers	Activity
------------	------------	-------	------	-----------	----------

kafka-ksqldbksql_processing_log	1	x1 100%	0 B	0	...
---	---	------------	-----	---	---------------------

TOPICS

8

Topic Name	Cleanup Policy	Retention (time or size)	Compaction (key-based)	
kafka-ksqldb				
▼ Advanced Configuration				
Property	Kafka Default	Broker Override	Topic Override	
min.insync.replicas	1	1		 Edit
retention.bytes	-1	-1		 Edit
retention.ms	604800000	-		 Edit
compression.type	producer	producer		 Edit
delete.retention.ms	86400000	86400000		 Edit
file.delete.delay.ms	60000	60000		 Edit
flush.messages	92233720368547...	92233720368547...		 Edit
 	0000000000000000			 
 				

+ CREATE

	Count	Size	Consumers	Activity	
	0	0 B	-		

“

Schema Registry

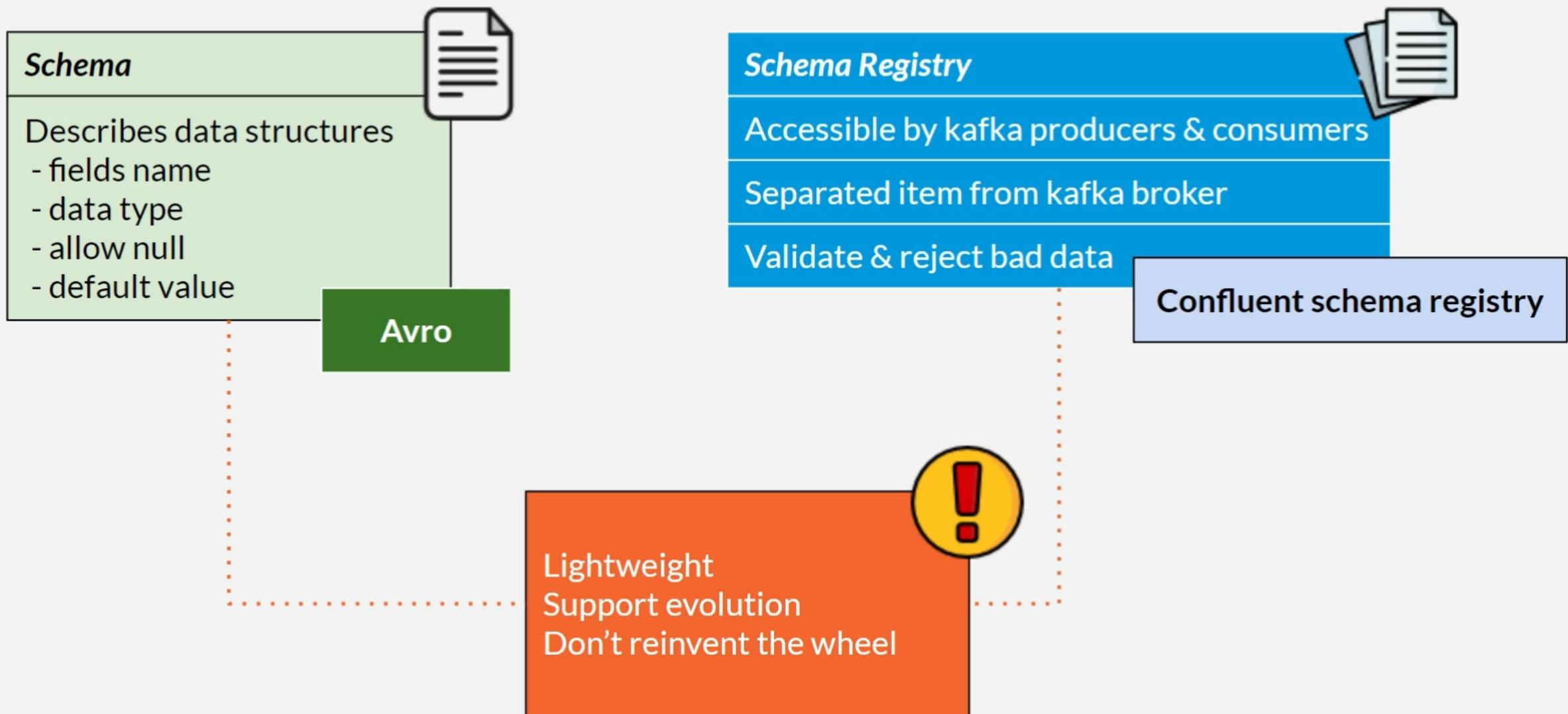
”

Use docker-compose-full.yml

[Download Conductor:](https://www.conduktor.io/download/) <https://www.conduktor.io/download/>

[Official Documentation:](https://www.conduktor.io/kafka) <https://www.conduktor.io/kafka>

What We Need?



Confluent Schema Registry

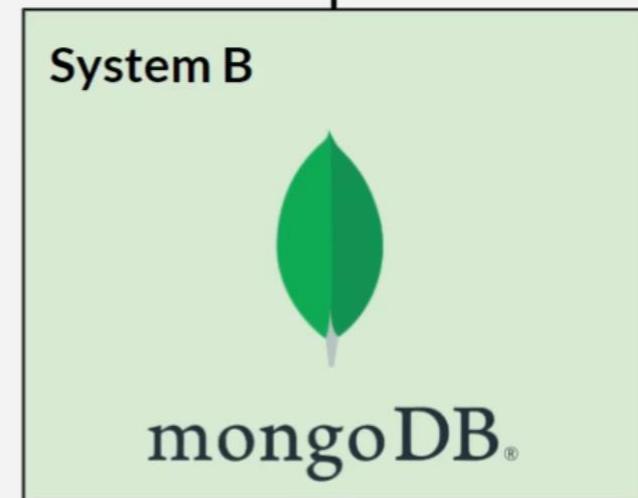
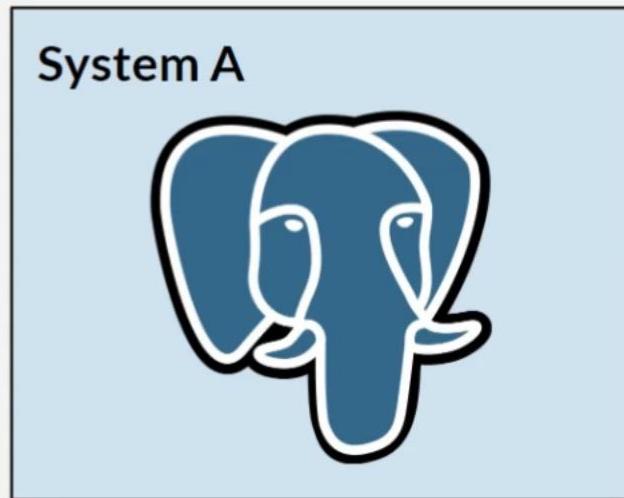
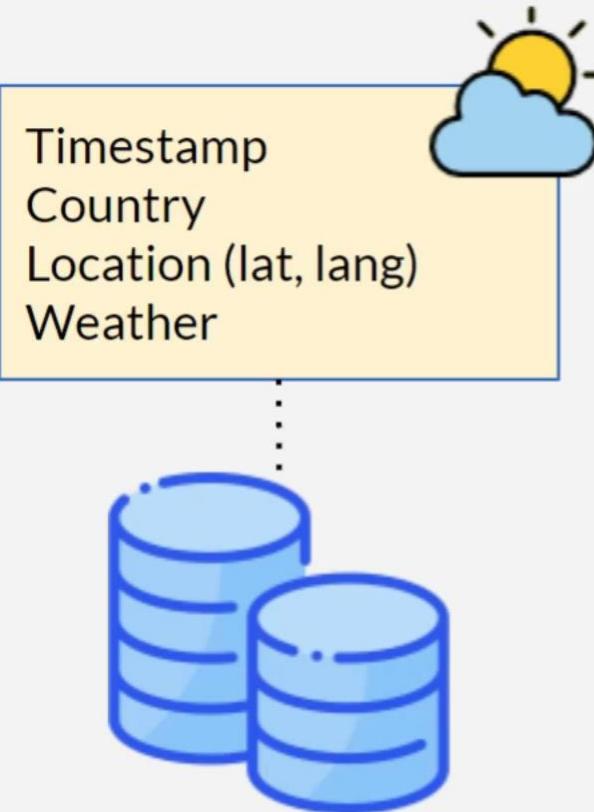
- × Free from confluent.io
- × Save & restore schema
- × Create, read, update, delete schema
- × Validates vs schema on topic
- × Works for key and / or value

“

AVRO

”

Data Transfer



timestamp	country	location	weather
<i>Long, epoch time</i>	<i>String, ISO 3166-1 country code</i>	<i>Point, (lat, long)</i>	<i>String, custom weather code</i>
1640945914	ID	-6.3040, 106.6435	STORM
1641010362	ID	-6.3040, 106.6435	LIGHT_RAIN

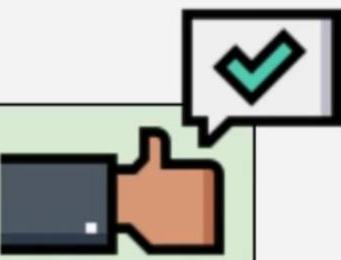
Data Transfer

- × From postgre to mongodb?
- × Database-to-database is not practical
- × Kafka
- × Data format on kafka?
- × Recognized by both system

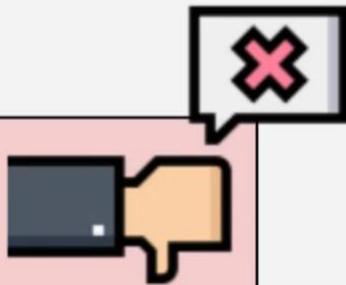
CSV

1640945914, ID, -6.3040, 106.6435, STORM

1641010362, ID, -6.3040, 106.6435, LIGHT_RAIN



- Good programming language support



- Sensitive column order
- No data type validation : change timestamp from epoch second to string 2021-12-31T10:40:25
- No data structure (lat & long must be on separated column)

JSON

```
{  
  "timestamp":1640945914,  
  "country":"ID",  
  "geolocation":{  
    "lat":-6.3040,  
    "long":106.6435  
  },  
  "weather":"STORM"  
}
```

```
{  
  "weather":"LIGHT_RAIN",  
  "country":"ID",  
  "timestamp":1641010362,  
  "geolocation":{  
    "lat":-6.3040,  
    "long":106.6435  
  }  
}
```

- Good programming language support
- Column order is not matter
- Good data structure (JSON elements, array)

- No data type validation : change timestamp from epoch second to string 2021-12-31T10:40:25

AVRO

Avro Schema

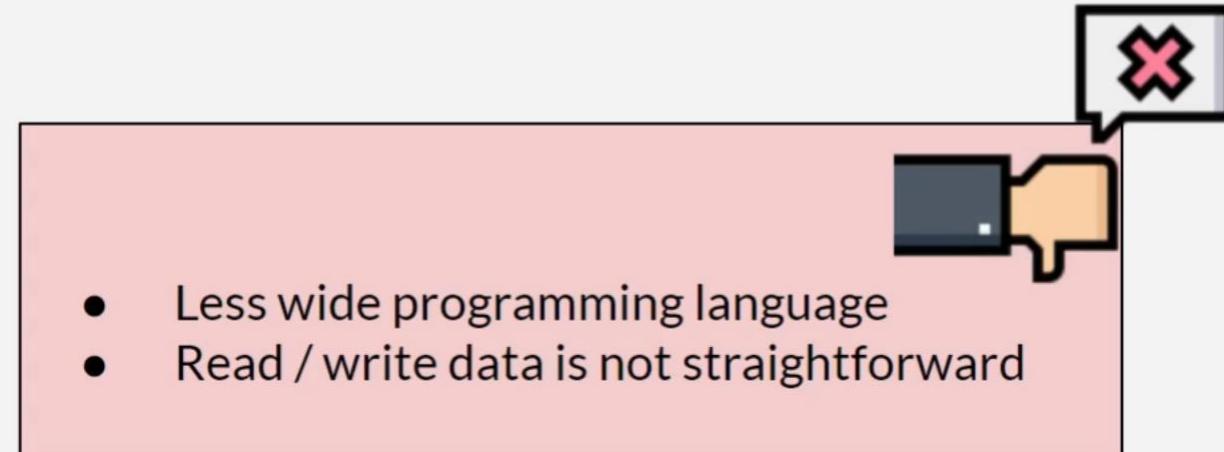
```
Timestamp : long, not null  
Country : string, not null  
Weather : string, not null  
Geolocation : geolocation  
(another avro), not null
```

Avro Body (binary)

```
{  
    "weather": "LIGHT_RAIN",  
    "country": "ID",  
    "timestamp": 1641010362,  
    "geolocation": {  
        "lat": -6.3040,  
        "long": 106.6435  
    }  
}
```



- Column order is not matter
- Good data structure (AVRO elements, array)
- Enforce data validation with schema
- Binary (less size)
- Integrates well with hadoop (big data)
- Schema evolution



- Less wide programming language
- Read / write data is not straightforward

Why Avro?

- × Current schema registry supports avro, protobuf, & json schema
- × Why avro?
 - × Big data ecosystem
- × Kafka supports avro better

Avro Schema Definition

```
"type": schema type as defined on avro specification,  
"namespace": like package in Java,  
"name": schema name,  
"aliases": (optional) array of name alias for this schema,  
"doc": (optional) documentation,  
"fields": array of fields specification  
    "name": field name,  
    "type": field type, sometimes combined with "logicalType",  
    "default": (optional) default value for field.  
    "doc": (optional) field documentation,  
    "order": (optional) field sort order  
    "aliases": (optional) array of field alias
```

Avro Schema Definition

```
{  
  "type": "record",  
  "namespace": "com.course",  
  "name": "AvroSample",  
  "aliases": ["AvroExample", "SampleAvro"],  
  "doc": "This is just a sample doc",  
  "fields": [  
    {"name": "name", "type": "string"},  
    {"name": "birthDate", "type": "long", "logicalType": "timestamp-millis"},  
    {"name": "email", "type": "string"},  
    {"name": "maritalStatus", "type": "string", "default": "UNKNOWN"}  
  ]  
}
```

Avro Schema

- × **type**
 - × Primitive : int, string, boolean, etc
 - × Complex :
 - × record (mostly will use this)
 - × enum
 - × array
 - × map
 - × union
 - × fixed
- × **logicalType**
 - × Give more meaning to primitive

Primitive Types

- × **null**: no value
- × **boolean**: a binary value (true/false)
- × **int**: 32-bit signed integer
- × **long**: 64-bit signed integer
- × **float**: 32-bit floating-point number
- × **double**: 64-bit floating-point number
- × **bytes**: sequence of 8-bit unsigned bytes
- × **string**: unicode character sequence

Avro Schema Definition

```
{  
  "type": "record",  
  "namespace": "com.course",  
  "name": "AvroSample",  
  "aliases": ["AvroExample", "SampleAvro"],  
  "doc": "This is just a sample doc",  
  "fields": [  
    {"name": "name", "type": "string"},  
    {"name": "email", "type": "string"},  
    {"name": "customerRating", "type": "float"},  
    {"name": "acceptPromotionEmail", "type": "boolean"}  
  ]  
}
```

Logical Types

- × More than primitive
- × Data type is common
- × Logical type gives more meaning to existing primitive
- × Example:
 - × **decimal** - bytes
 - × **uuid** - string
 - × **date** - int
 - × **time-millis** - int
 - × **timestamp-millis** - long
 - × ... (complete at avro documentation)

Avro Complex Types

- ✗ Record
- ✗ Union
- ✗ Enum
- ✗ Array
- ✗ Map

<https://avro.apache.org/docs/current/spec.html>

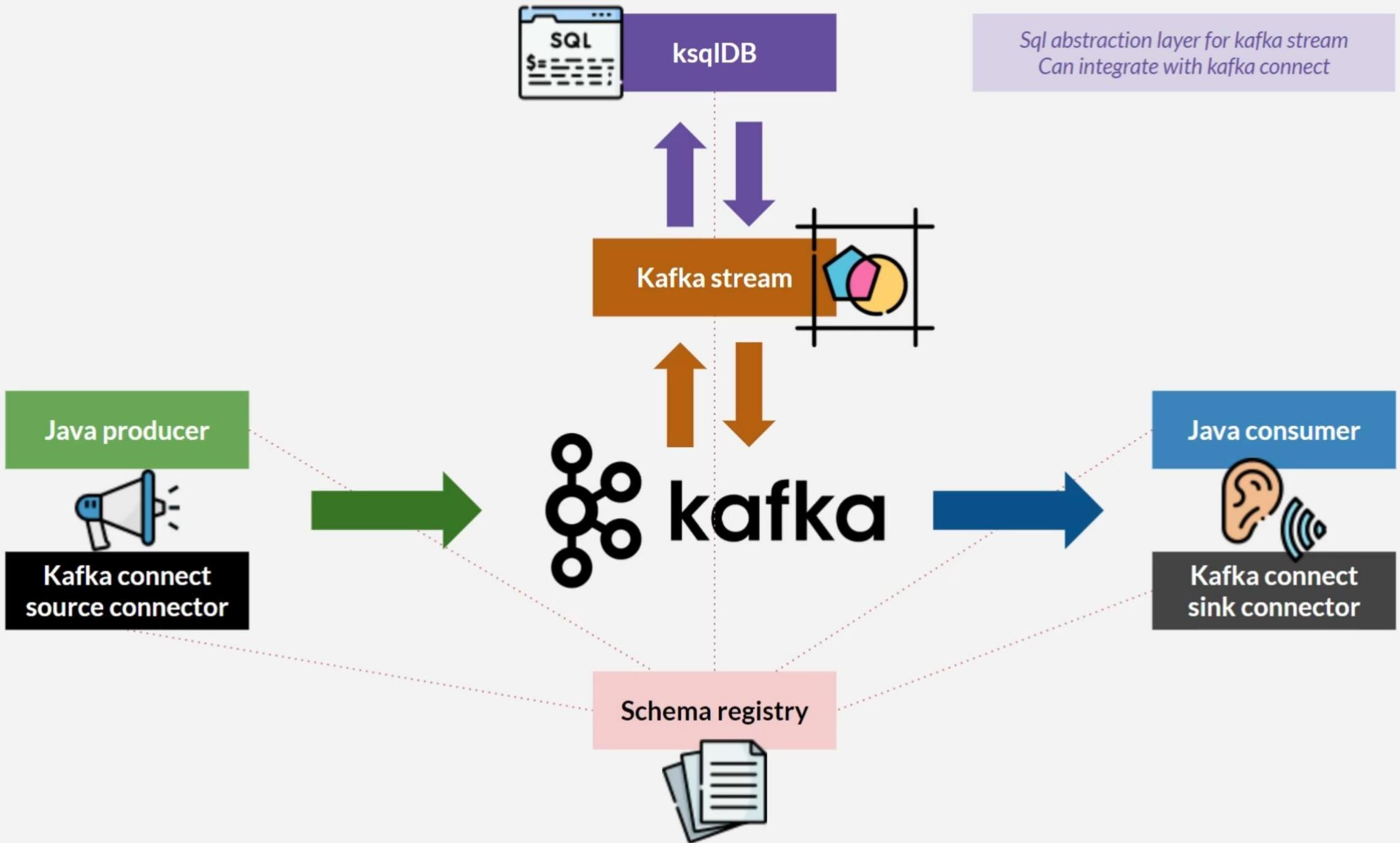
Spring initializr/IDE

- ▶ Dependency :
 - ▶ Spring for apache kafka
- ▶ Settings:
 - ▶ Java Project with Grade
 - ▶ Spring Boot 2.X
 - ▶ Group: com.virtusa
 - ▶ Artifact: kafka-avro-producer & kafka-avro-consumer
 - ▶ Adjust package name as necessary
 - ▶ Java Version: 11 or later

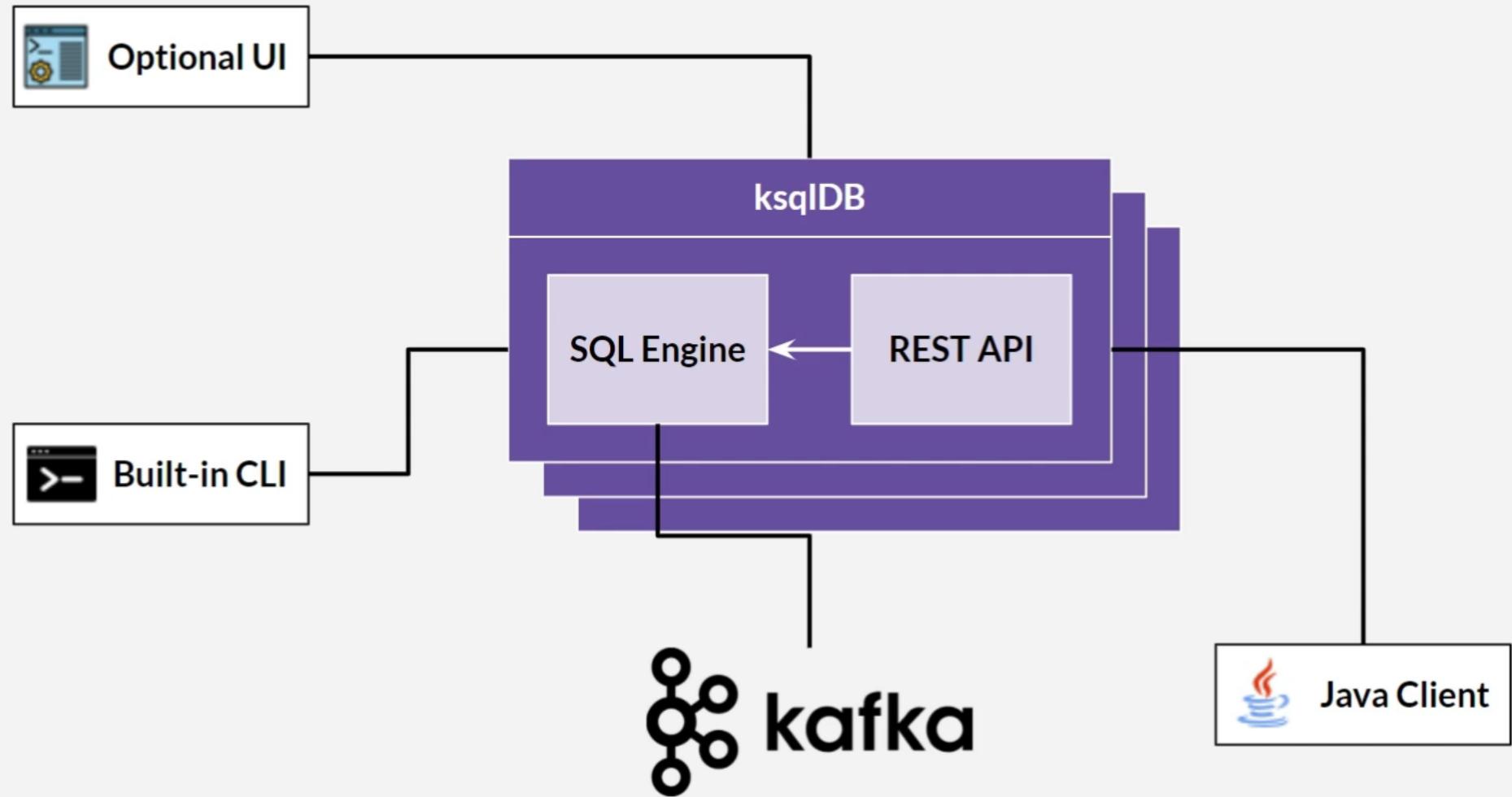
“

ksqlDB

”



ksqlDB Architecture



Venkat
Corporate Trainer & Motivational Speaker

