

Kafka as a Platform: the Ecosystem from the Ground Up

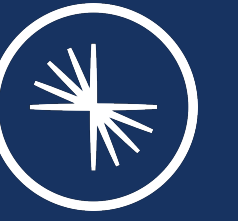
Robin Moffatt | #BudapestData | @rmoff

\$ whoami

- Robin Moffatt (@rmoff)
- Senior Developer Advocate at Confluent
(Apache Kafka, not Wikis 😊)
- Working in data & analytics since 2001
- ♠️ Oracle ACE Director (Alumnus)



<http://rmoff.dev/talks> • <http://rmoff.dev/blog> • <http://rmoff.dev/youtube>



EVENTS

@rmoff



#BudapestData



@confluentinc

EVENTS

EVENTS

- Notification
- State change

Human generated events

A Sale



*A Stock
movement*



Machine generated events

Networking



IoT



Applications



EVENTS

are

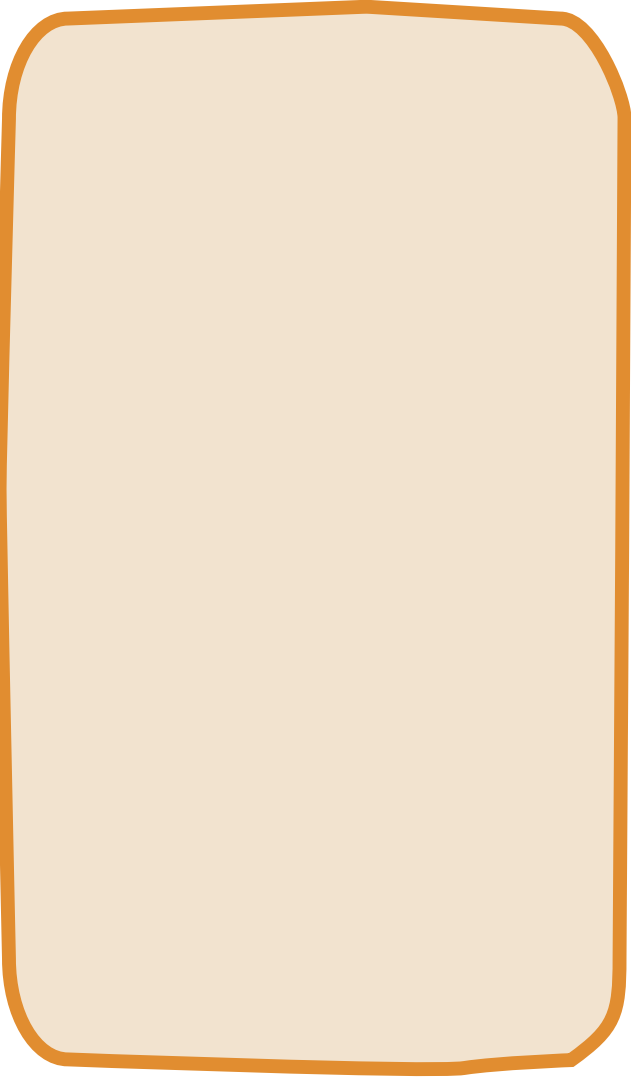
EVERYWHERE

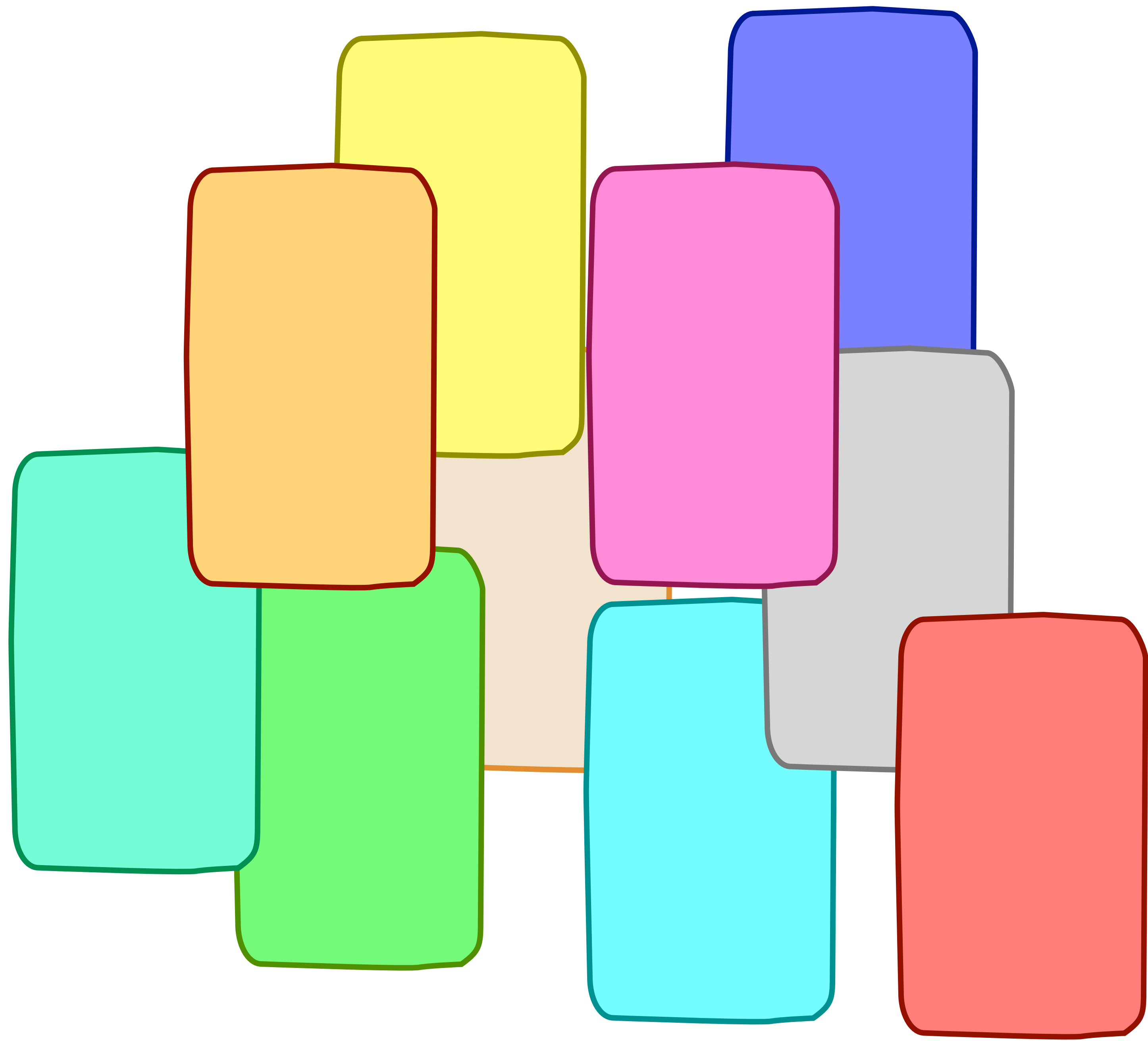
EVENTS

are

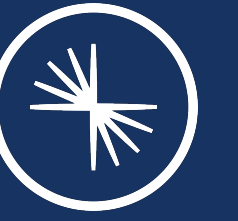
very
^

POWERFUL





K/v



LOGG

@rmoff

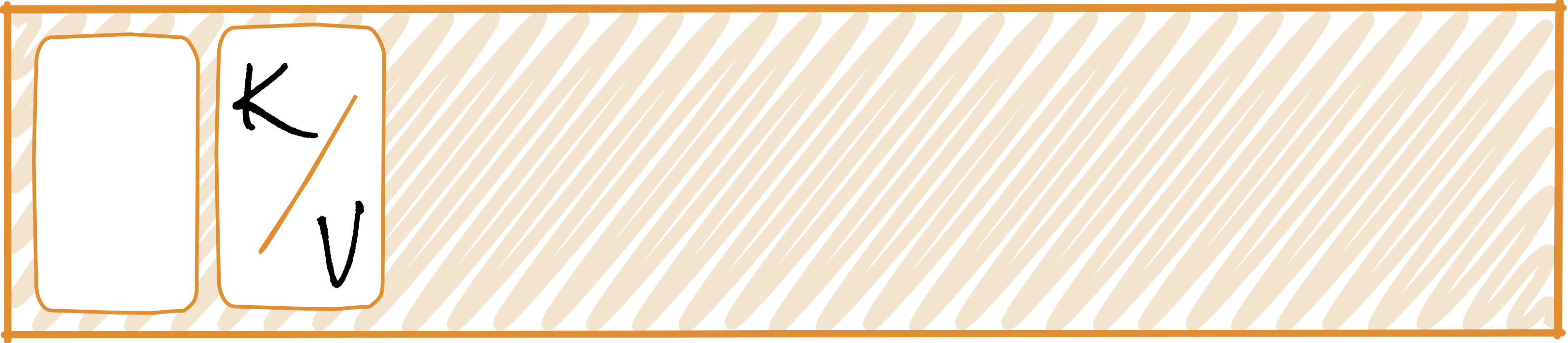


#BudapestData



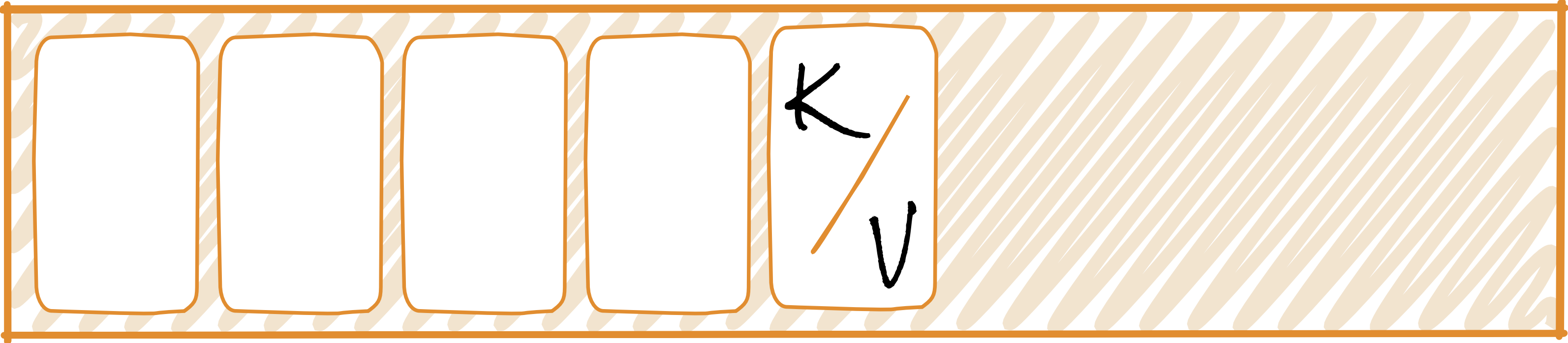
@confluentinc

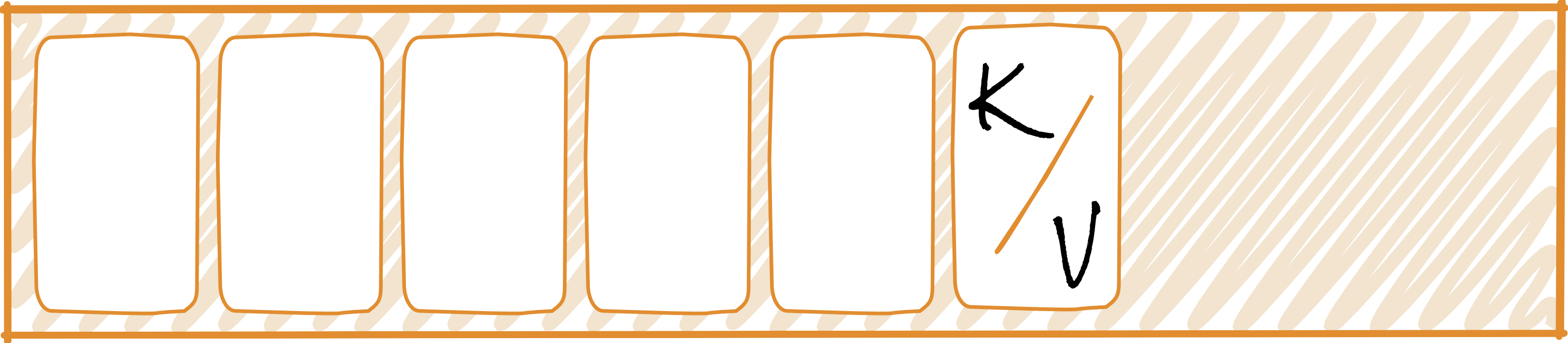


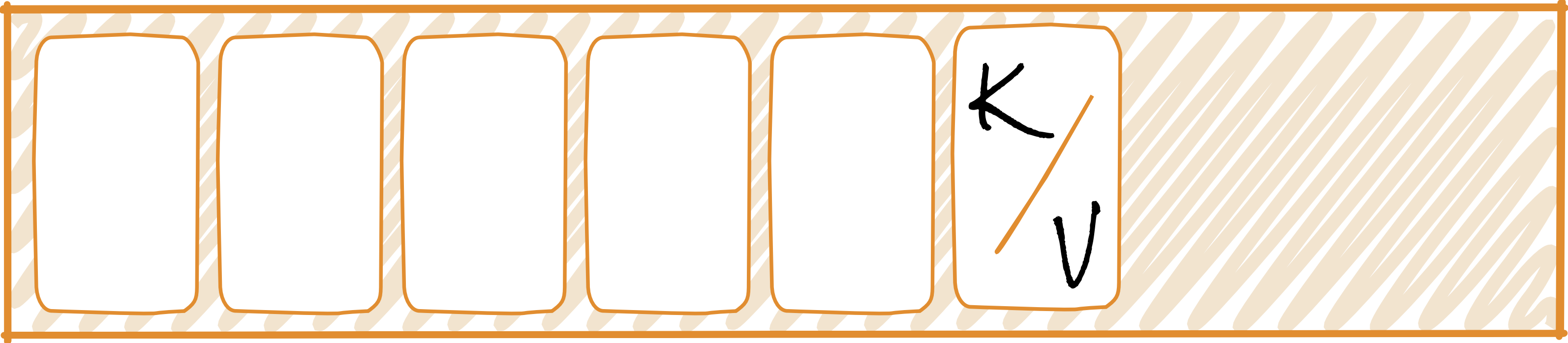








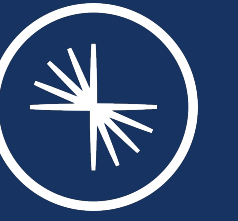




Immutable Event Log



Events are added at the end of the log



TOPICS

@rmoff



#BudapestData



@confluentinc

Topics

Clicks



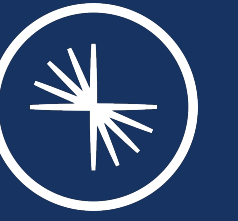
Orders



Customers



Topics are similar in
concept to tables in a
database



PARTITIONS

@rmoff



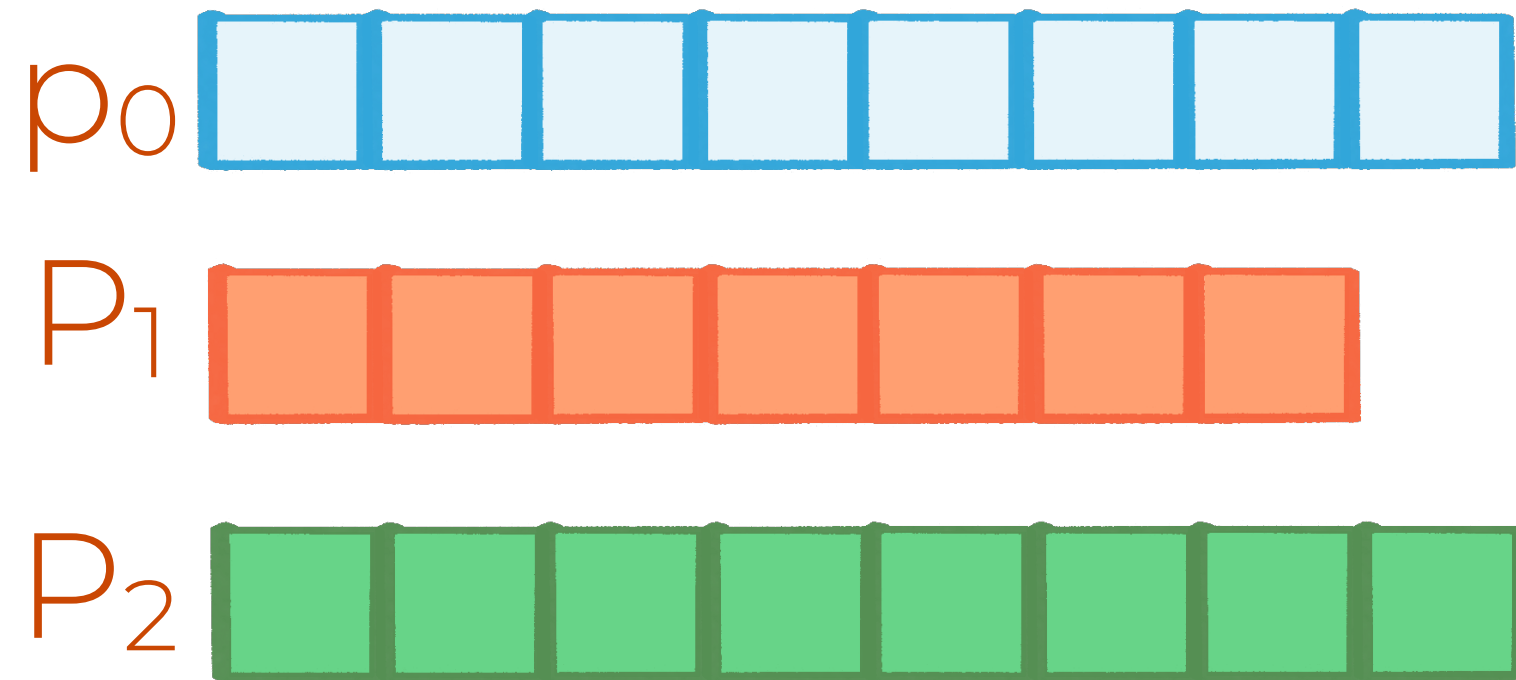
#BudapestData



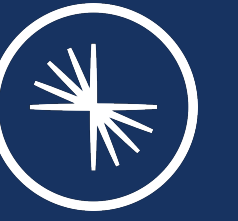
@confluentinc

Partitions

Clicks



Messages are guaranteed to be strictly ordered within a partition



PUB / SUB

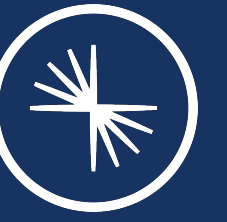
@rmoff



#BudapestData



@confluentinc



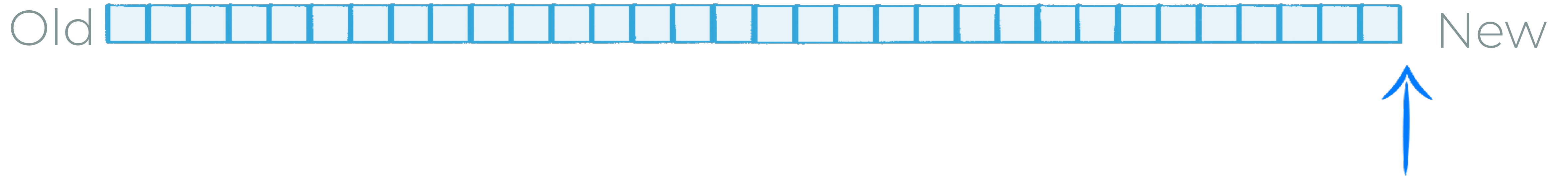
PUB / SUB

@rmoff

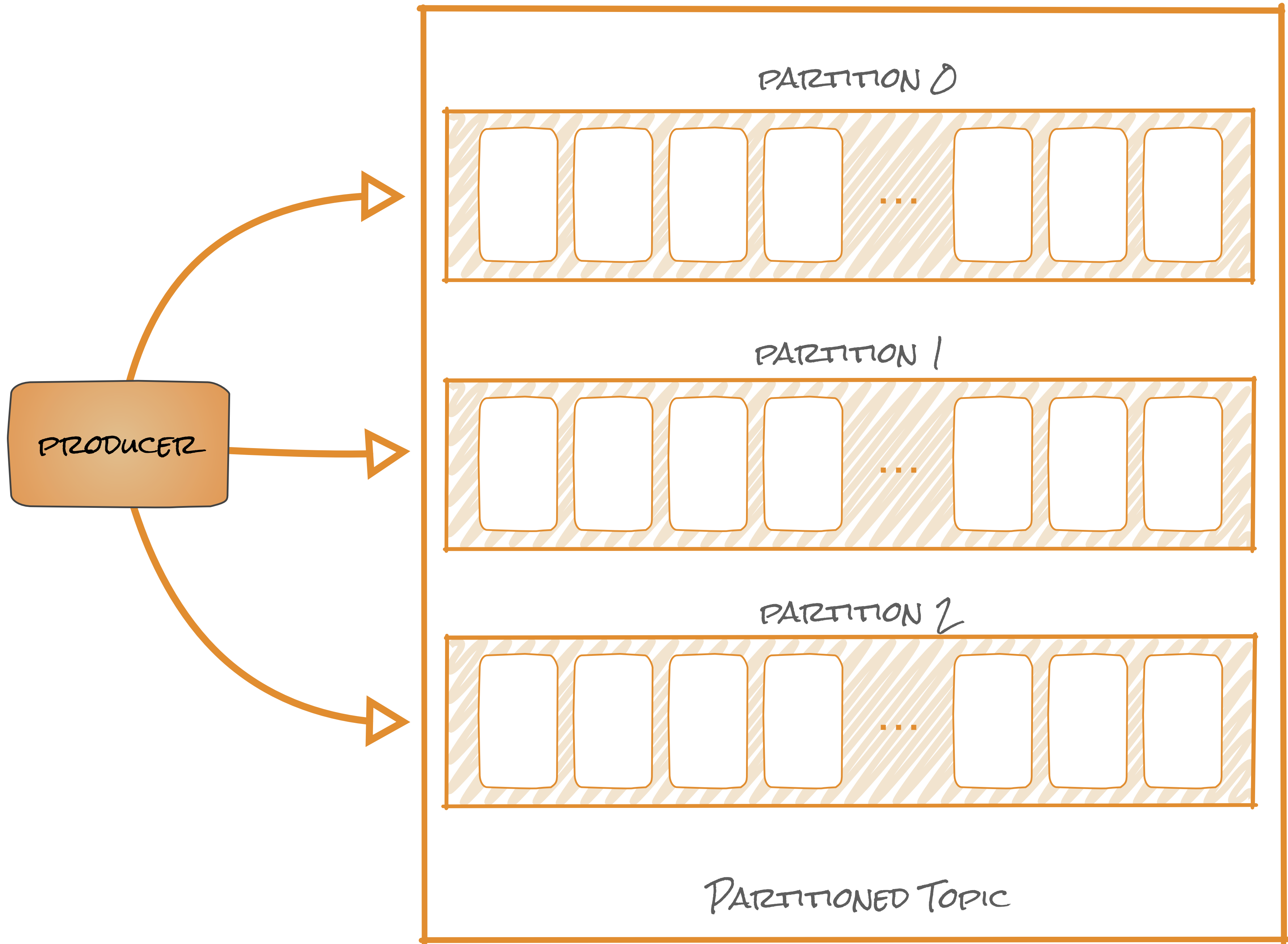
| #BudapestData

| @confluentinc

Producing data



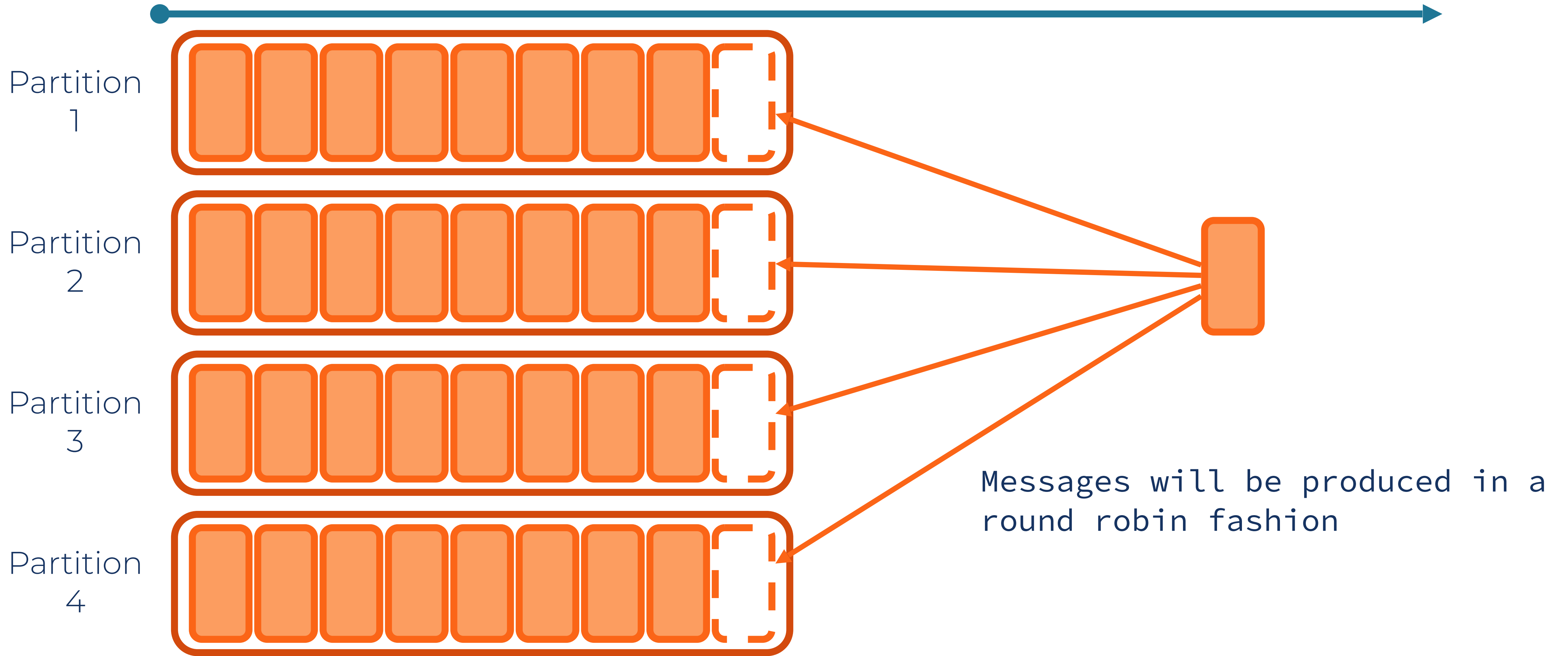
Messages are added at the end of the log



```
try (KafkaProducer<String, Payment> producer = new KafkaProducer<String, Payment>(props)) {  
    for (long i = 0; i < 10; i++) {  
        final String orderId = "id" + Long.toString(i);  
        final Payment payment = new Payment(orderId, 1000.00d);  
        final ProducerRecord<String, Payment> record =  
            new ProducerRecord<String, Payment>("transactions",  
                payment.getId().toString(),  
                payment);  
        producer.send(record);  
    }  
} catch (final InterruptedException e) {  
    e.printStackTrace();  
}
```

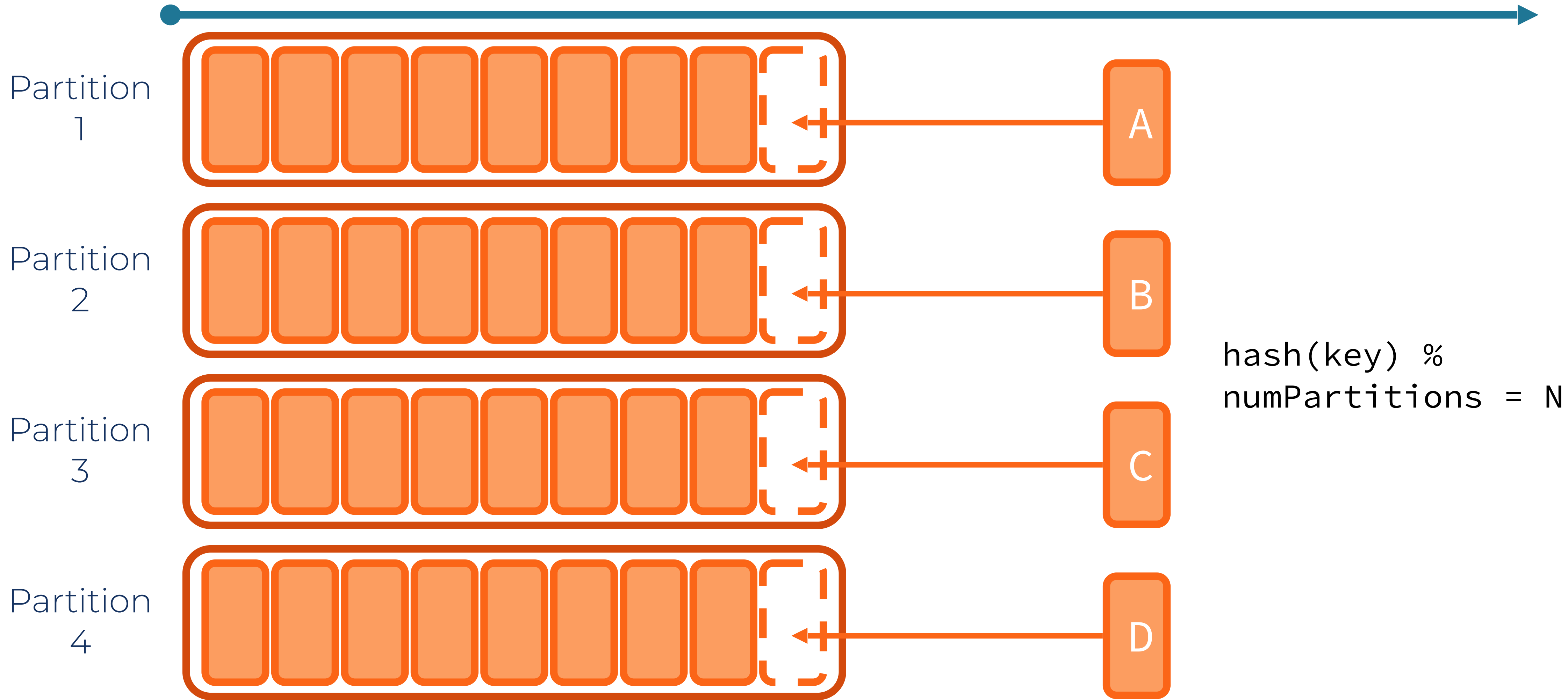
Producing to Kafka - No Key

Time

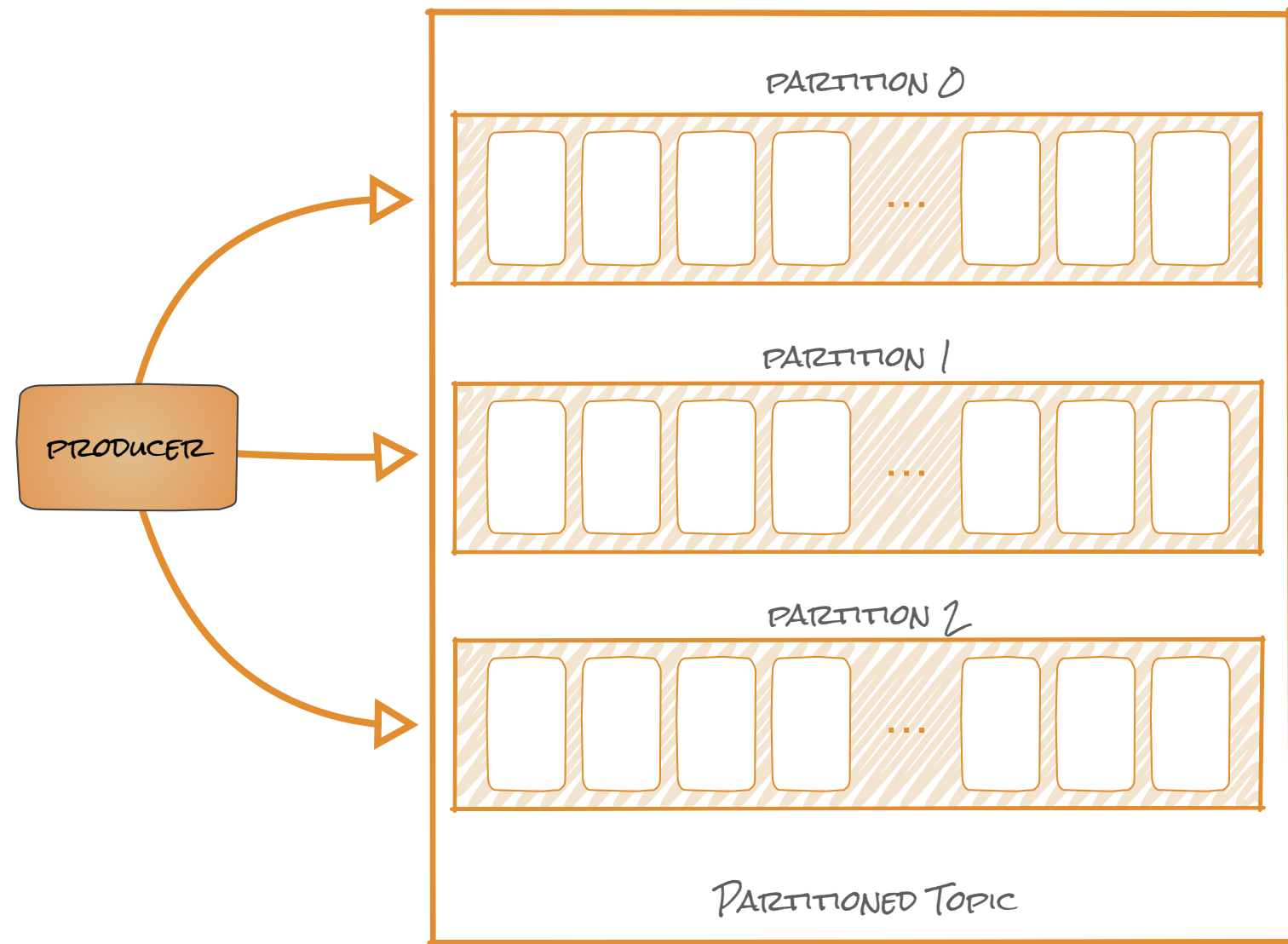


Producing to Kafka - With Key

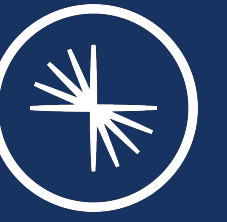
Time



Producers



- **A client application**
- **Puts messages into topics**
- **Handles partitioning, network protocol**
- **Java, Go, .NET, C/C++, Python**
- **Also every other language**



PUB / SUB

@rmoff



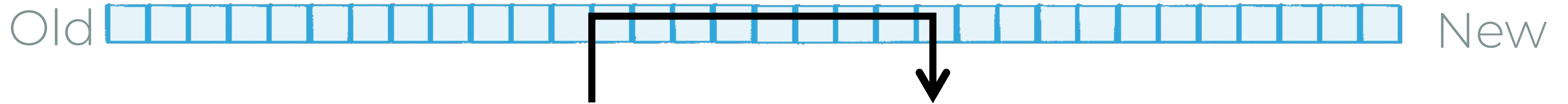
#BudapestData



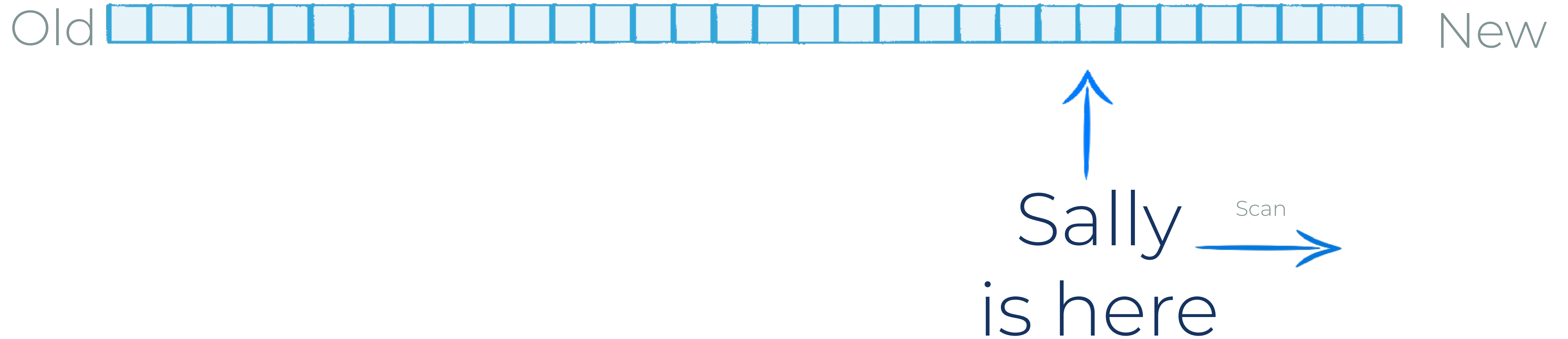
@confluentinc

Consuming data - access is only sequential

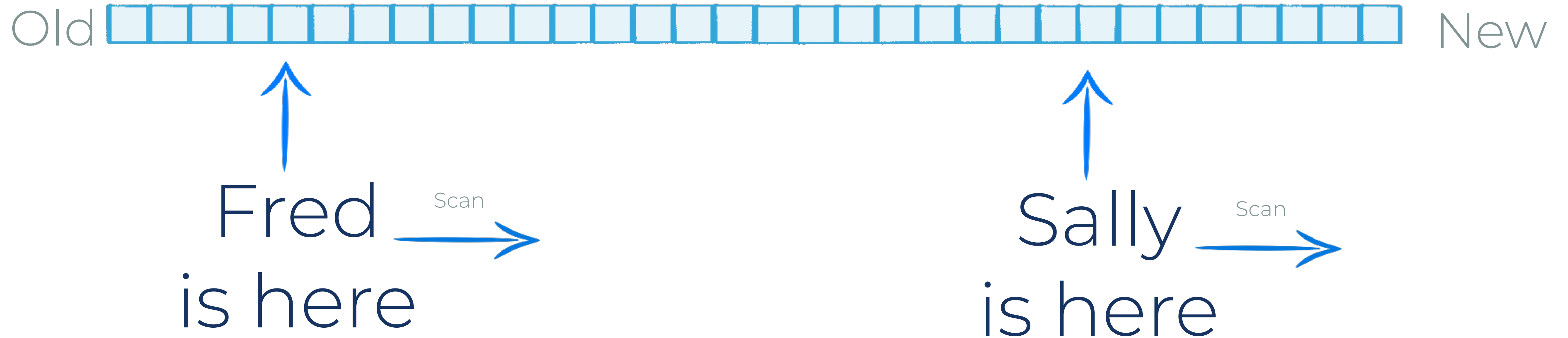
Read to offset & scan



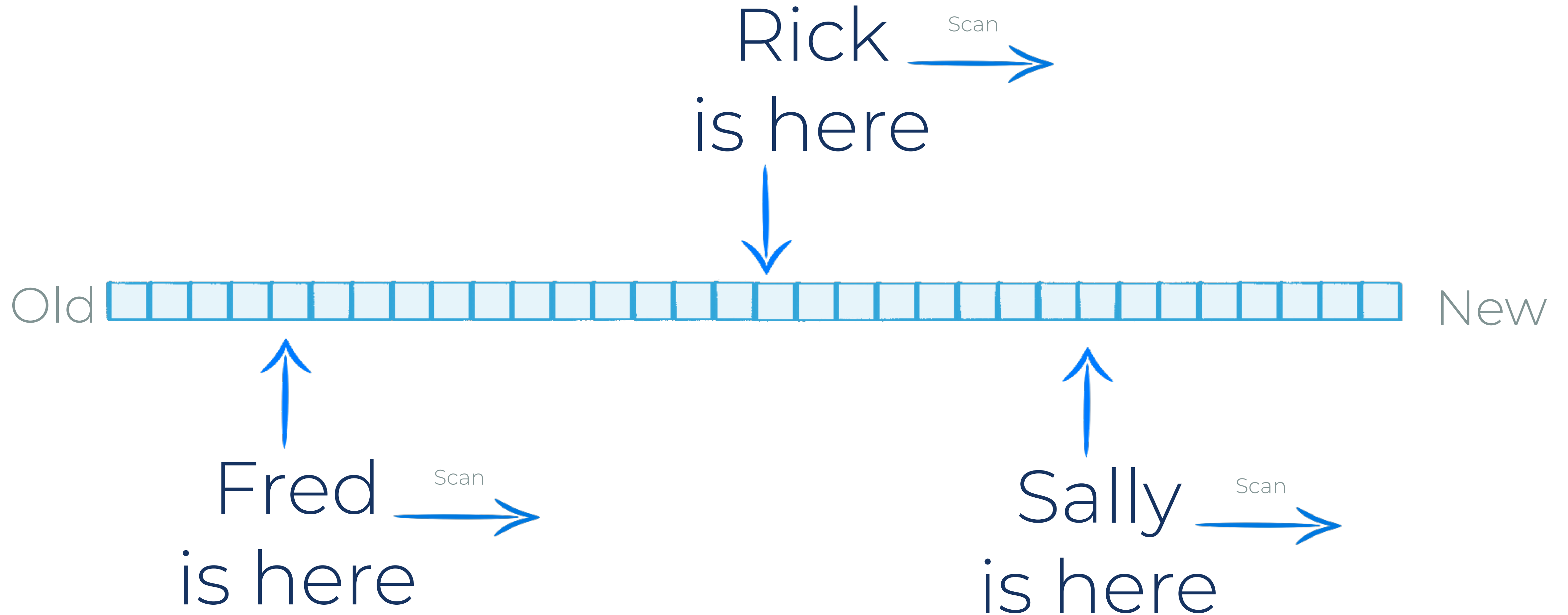
Consumers have a position of their own



Consumers have a position of their own

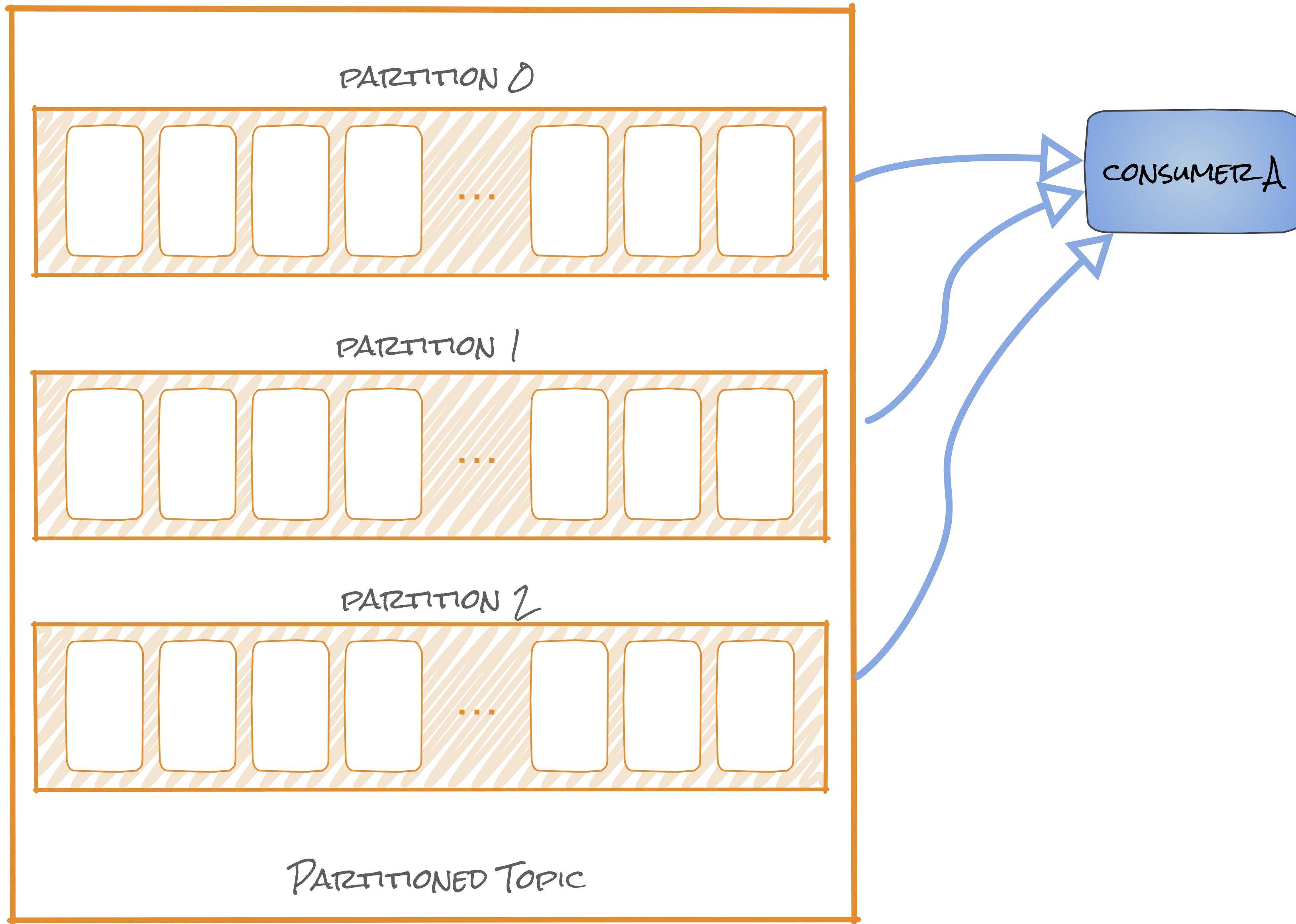


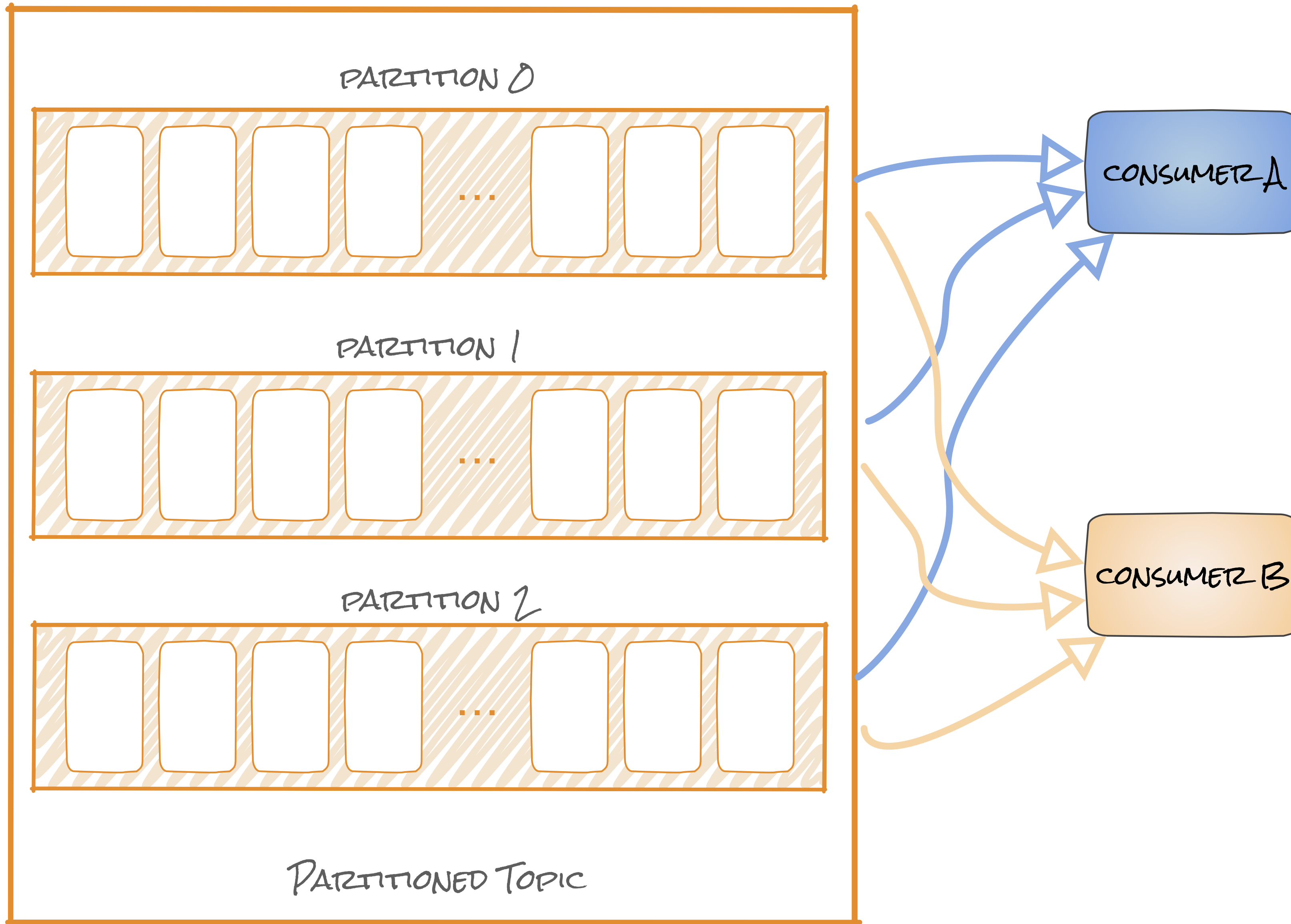
Consumers have a position of their own



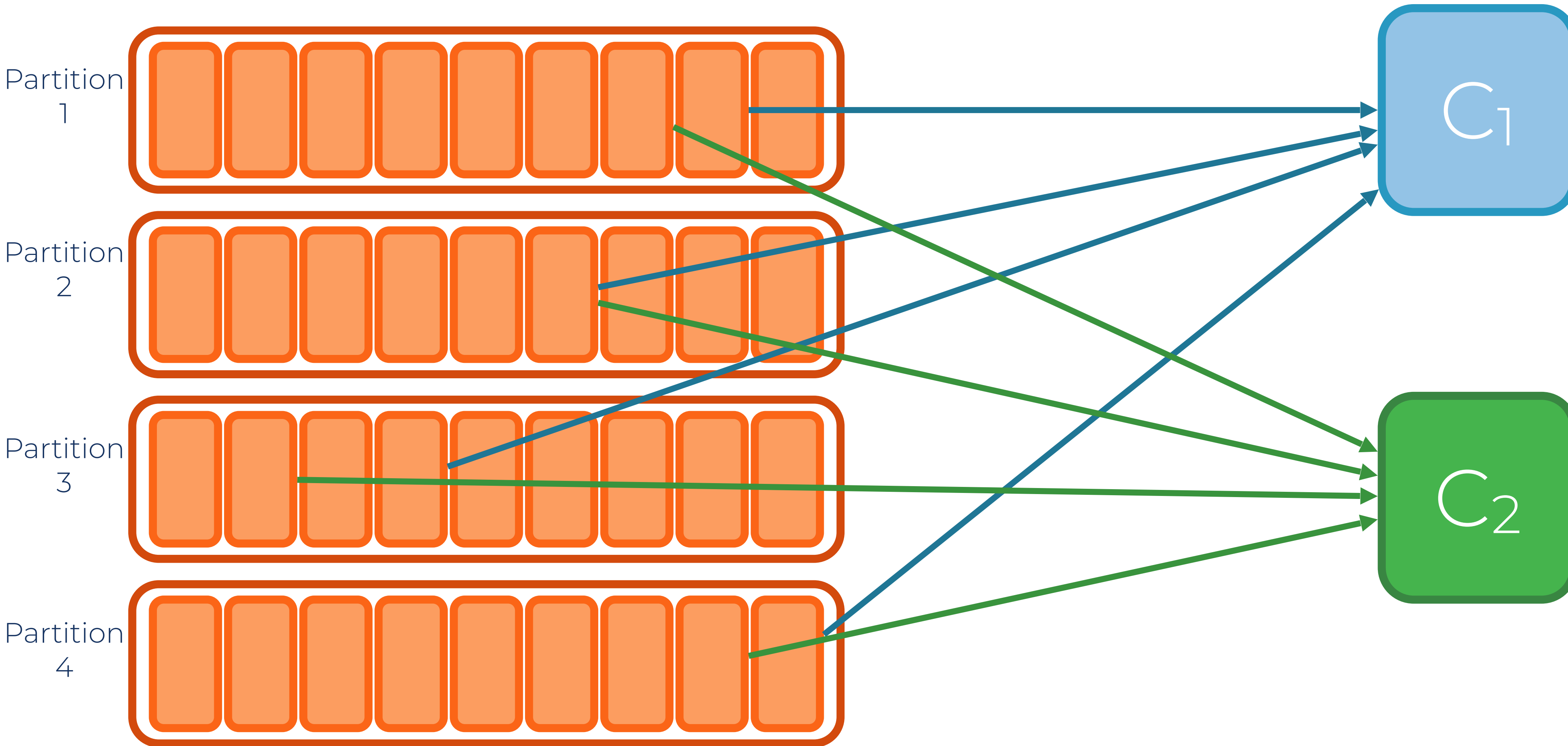
```
try (final KafkaConsumer<String, Payment> consumer = new KafkaConsumer<>(props)) {
    consumer.subscribe(Collections.singletonList(TOPIC));

    while (true) {
        ConsumerRecords<String, Payment> records = consumer.poll(100);
        for (ConsumerRecord<String, Payment> record : records) {
            String key = record.key();
            Payment value = record.value();
            System.out.printf("key = %s, value = %s%n", key, value);
        }
    }
}
```

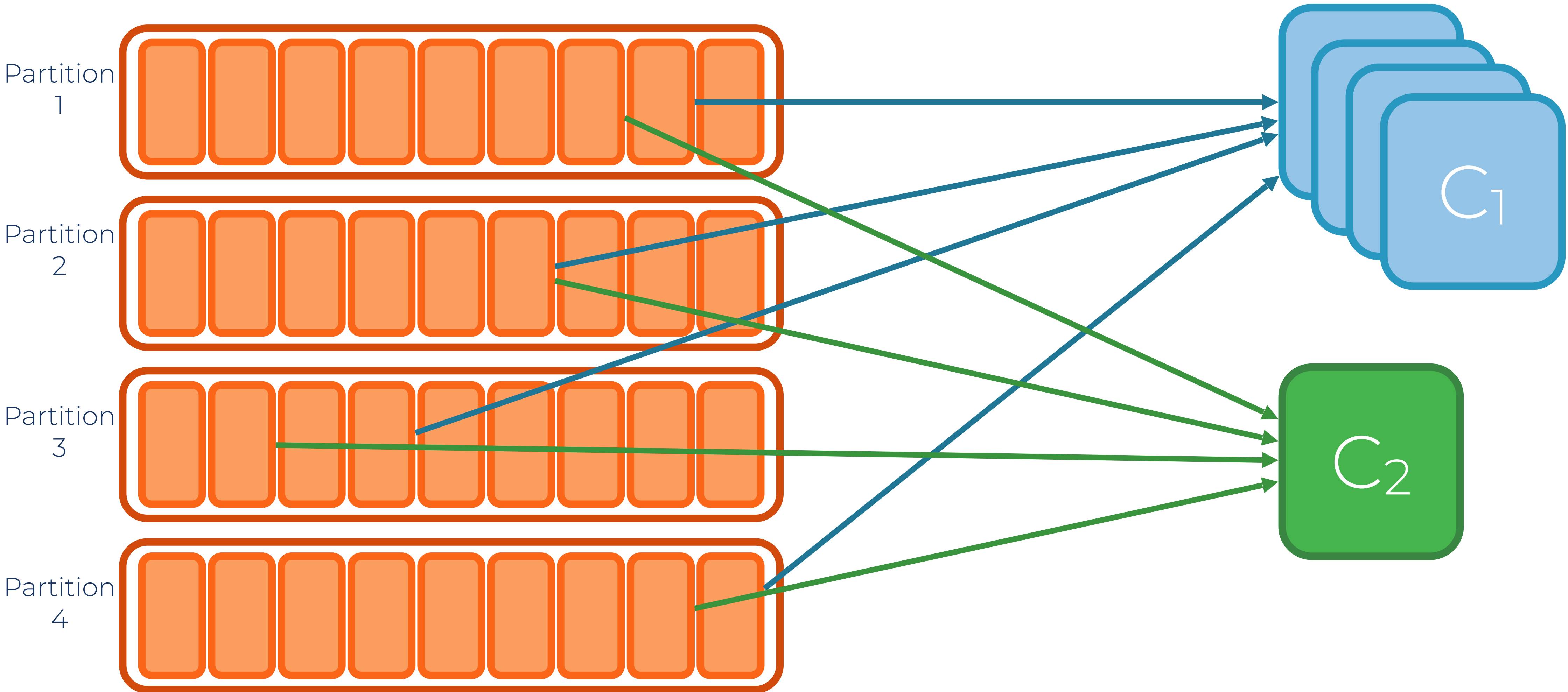




Consuming From Kafka - Multiple Consumers



Consuming From Kafka - Grouped Consumers



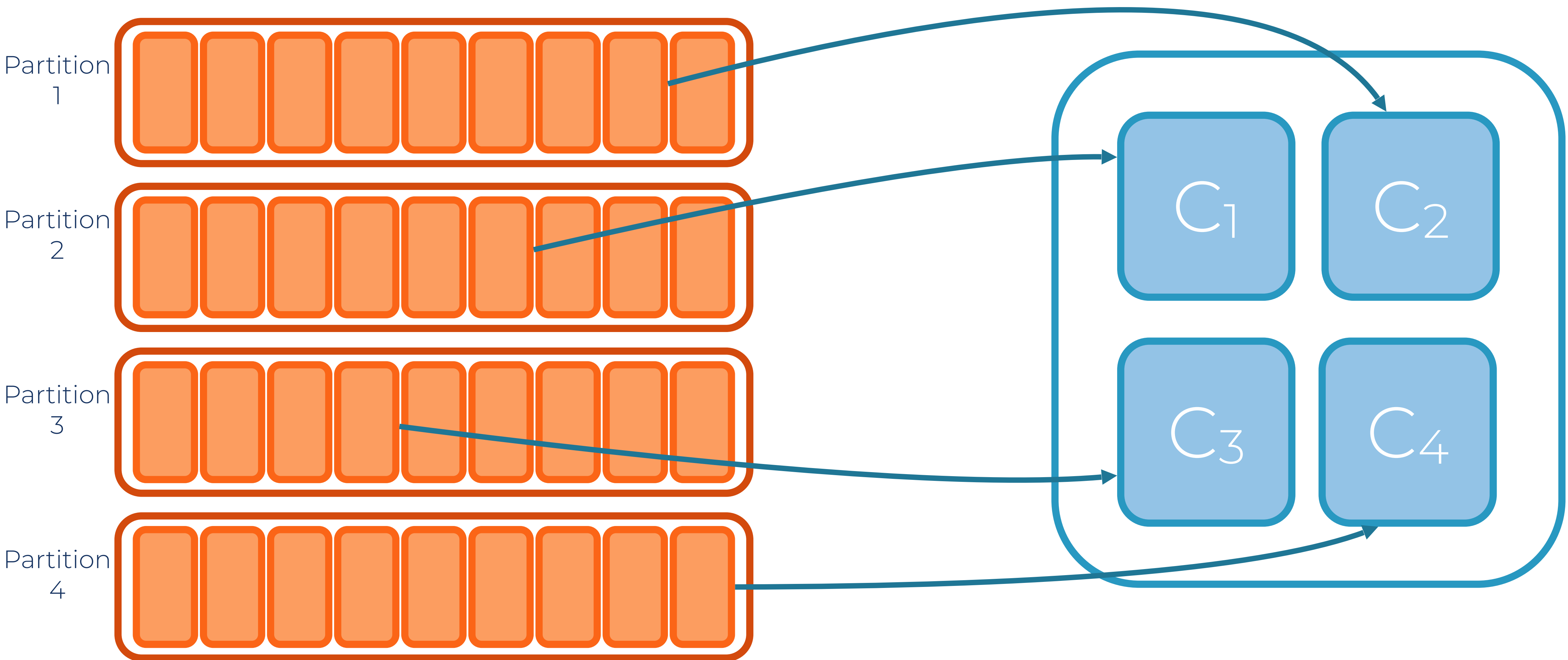
**CONSUMER GROUP
COORDINATOR**

CONSUMERS

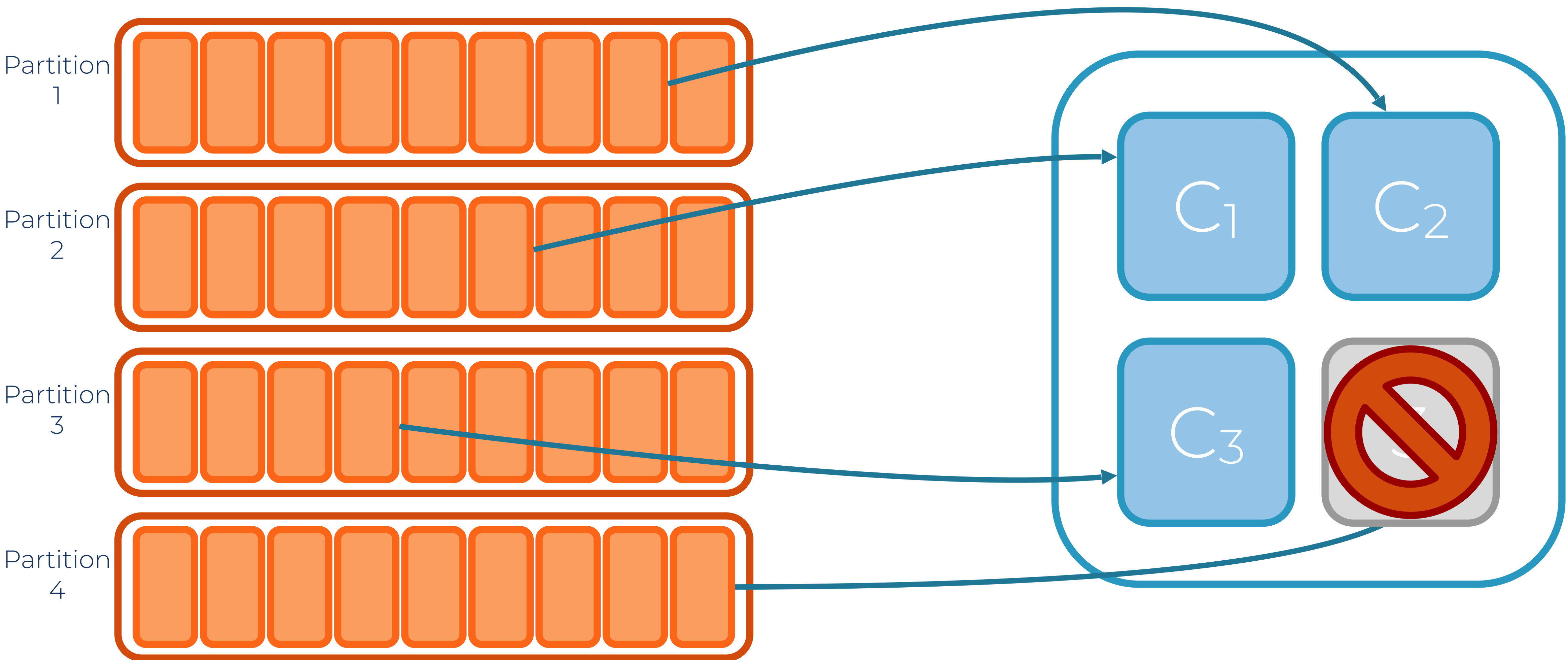
CONSUMER GROUP



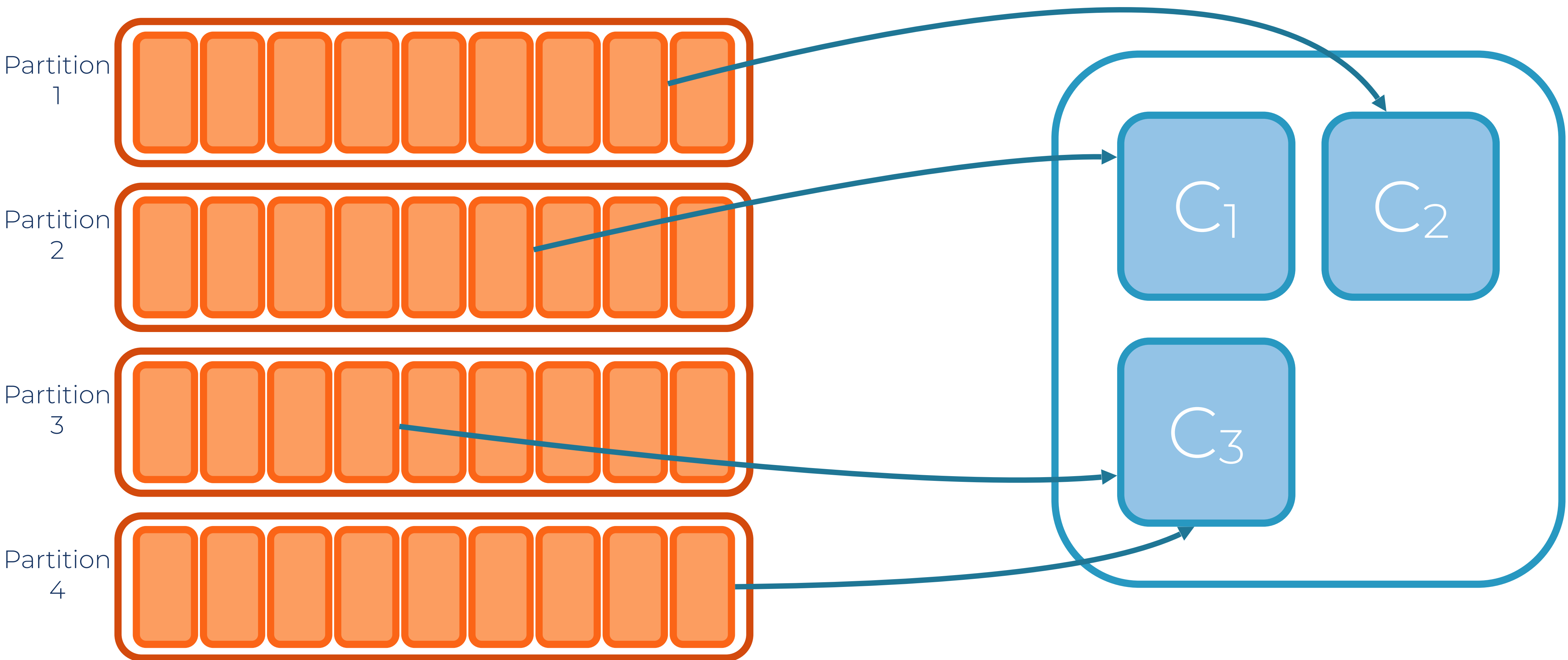
Consuming From Kafka - Grouped Consumers



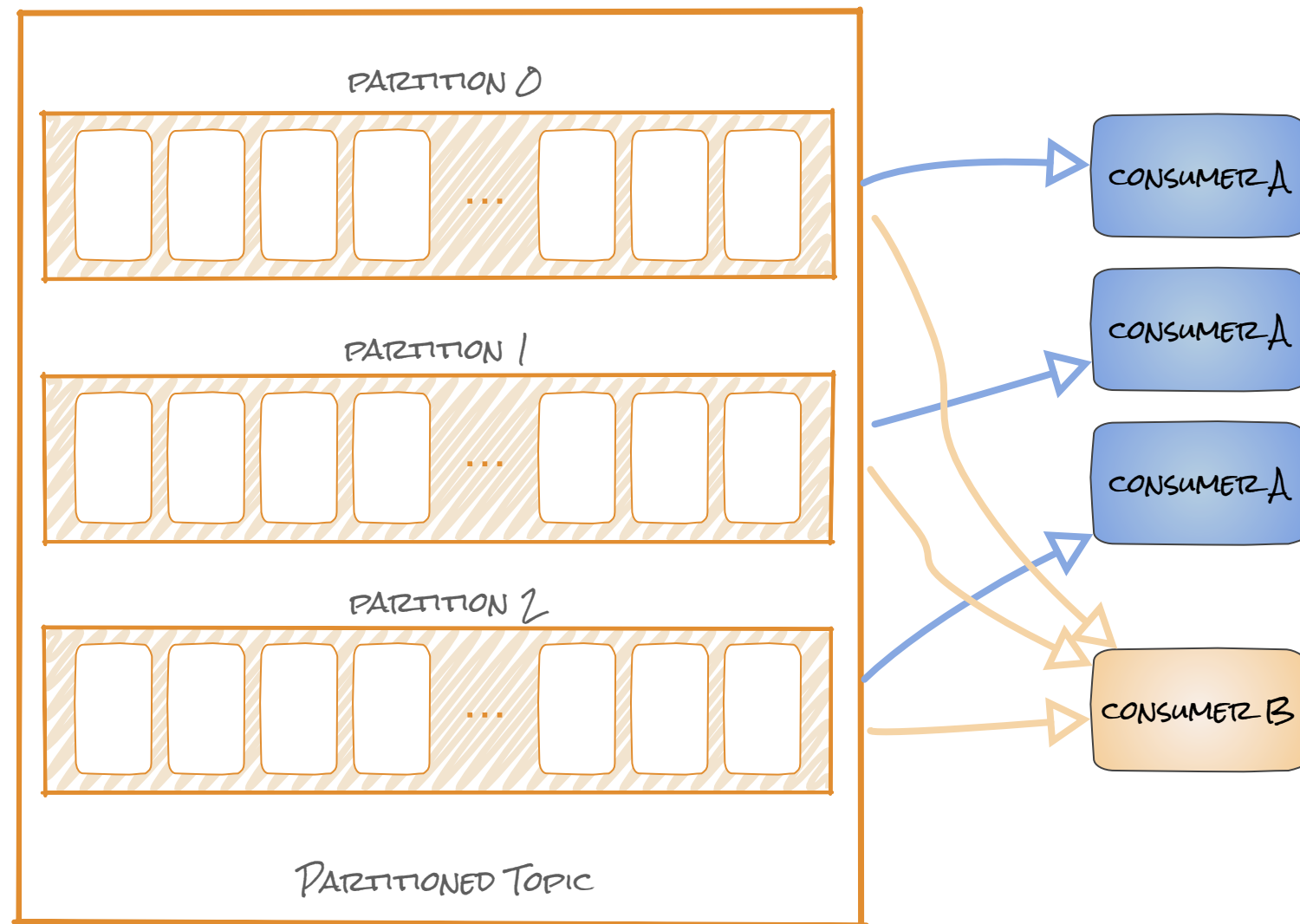
Consuming From Kafka - Grouped Consumers



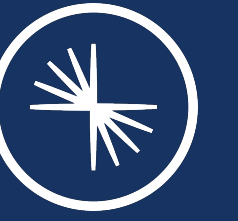
Consuming From Kafka - Grouped Consumers



Consumers



- **A client application**
- **Reads messages from topics**
- **Horizontally, elastically scalable (if stateless)**
- **Java, Go, .NET, C/C++, Python, everything else**



BROKERS

and REPLICATION

@rmoff



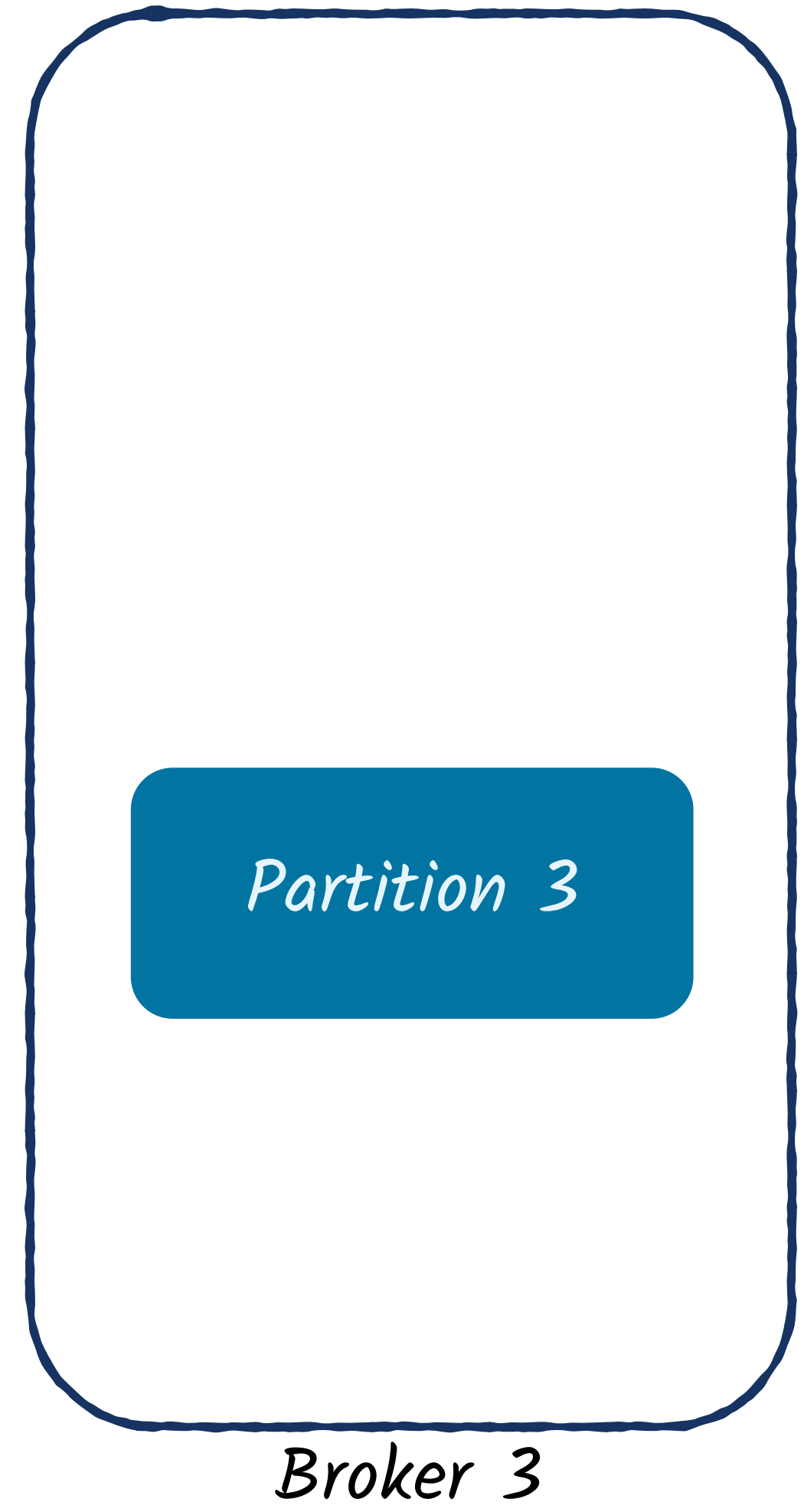
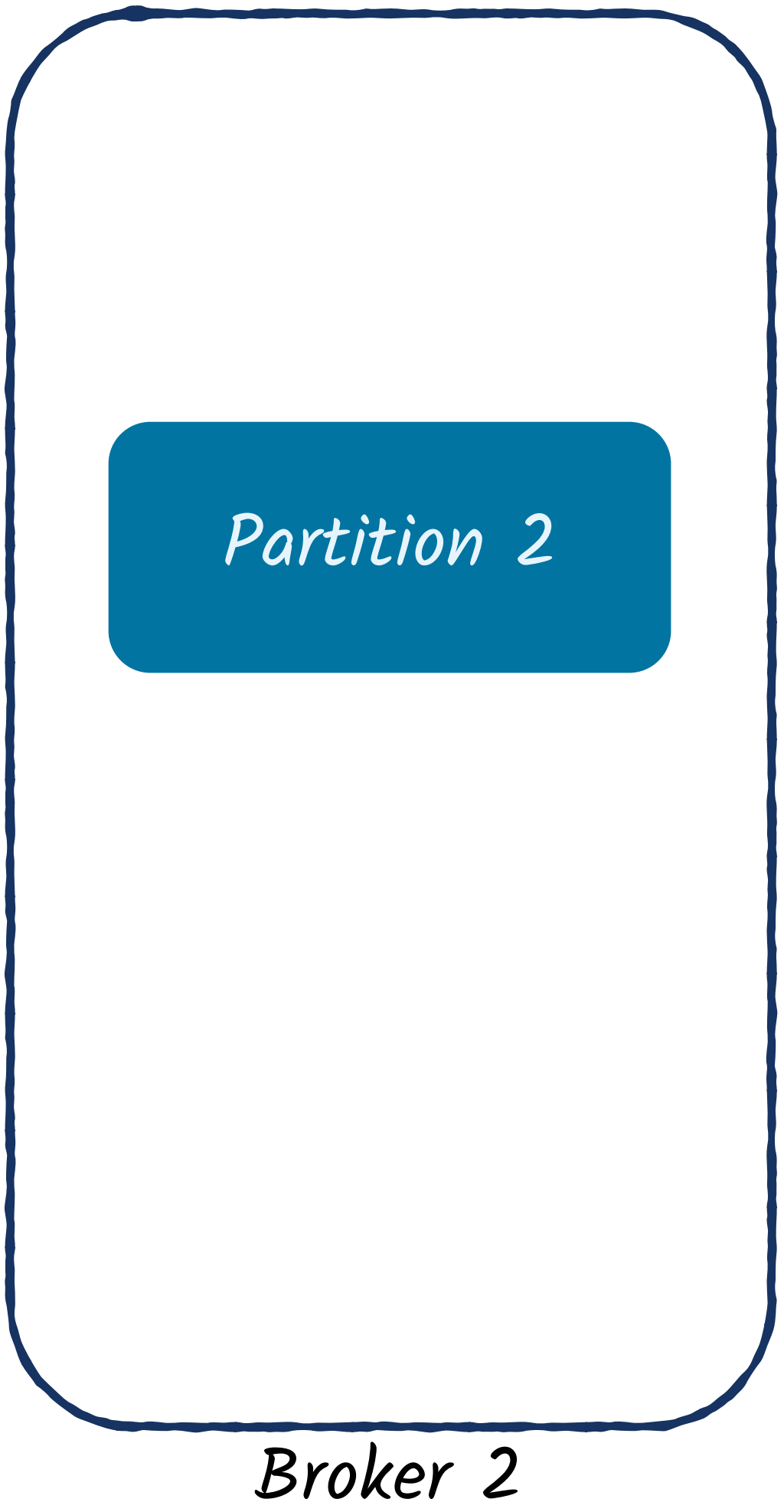
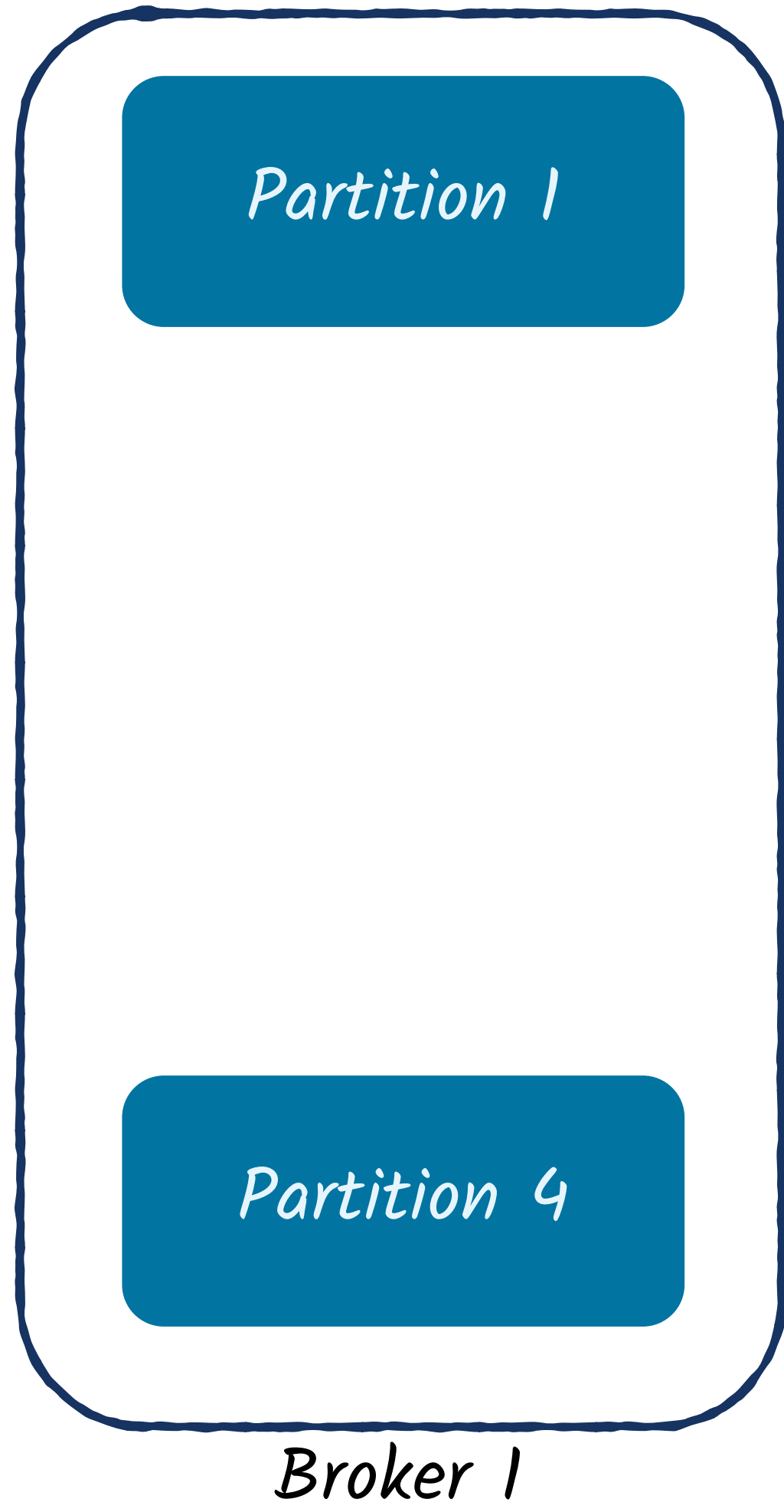
#BudapestData



@confluentinc

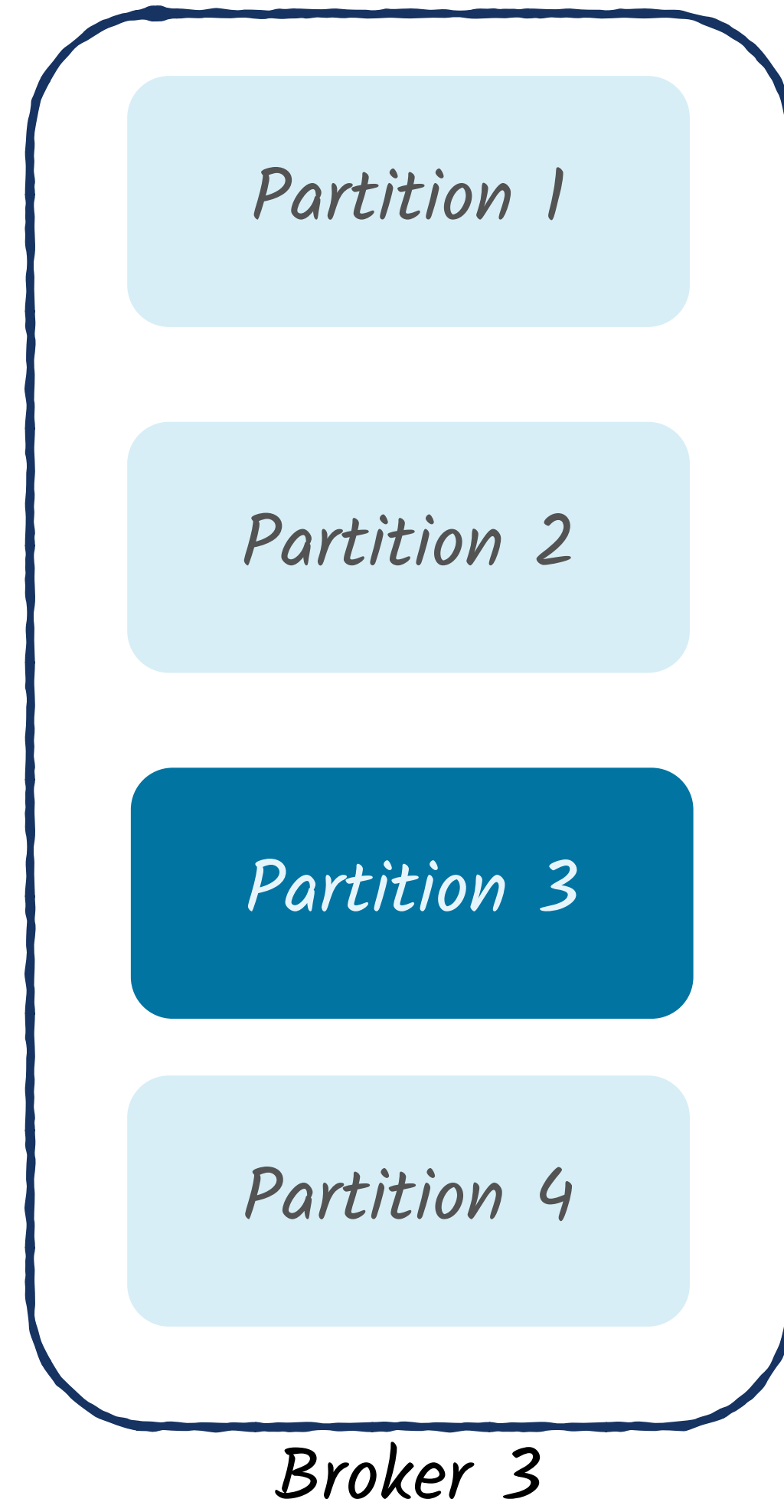
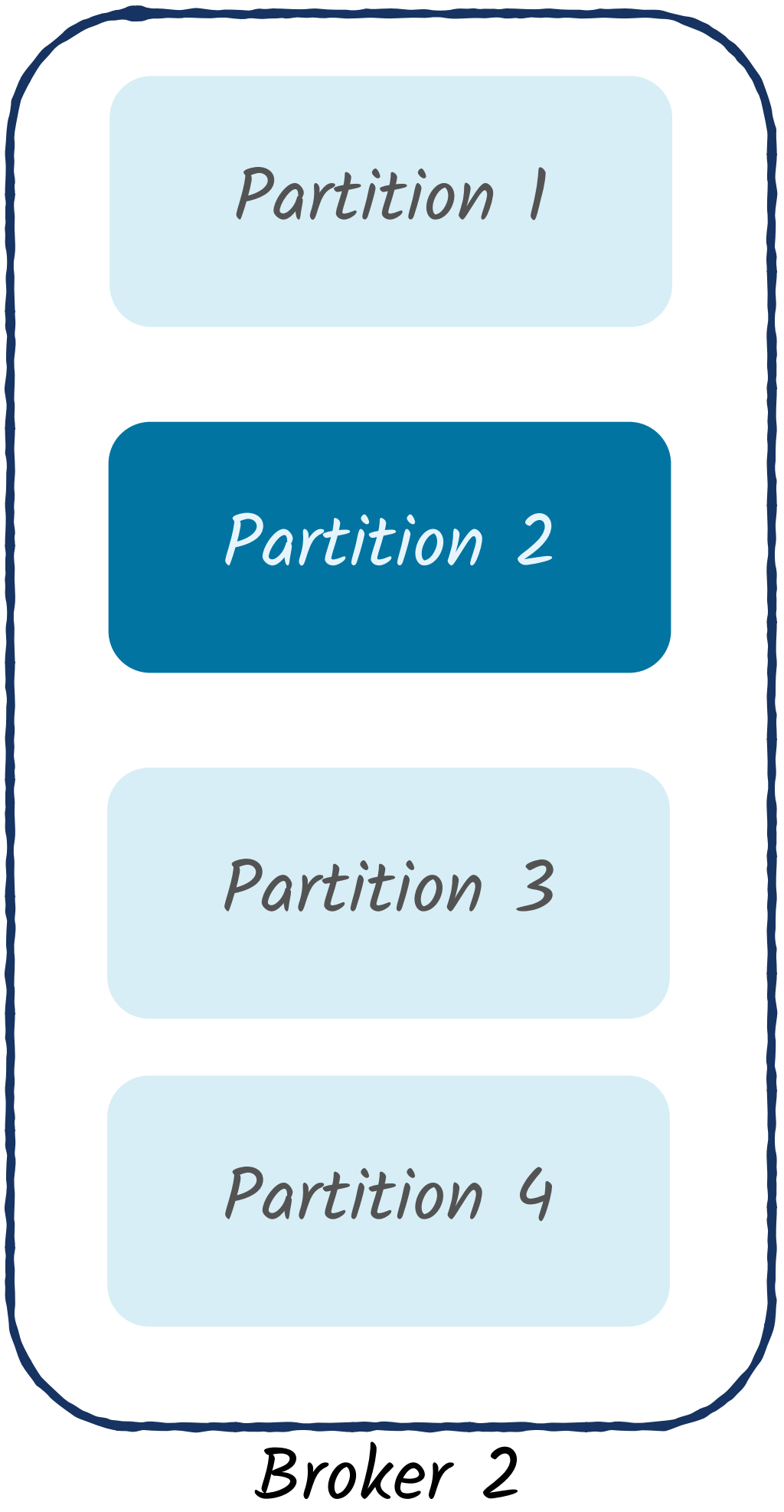
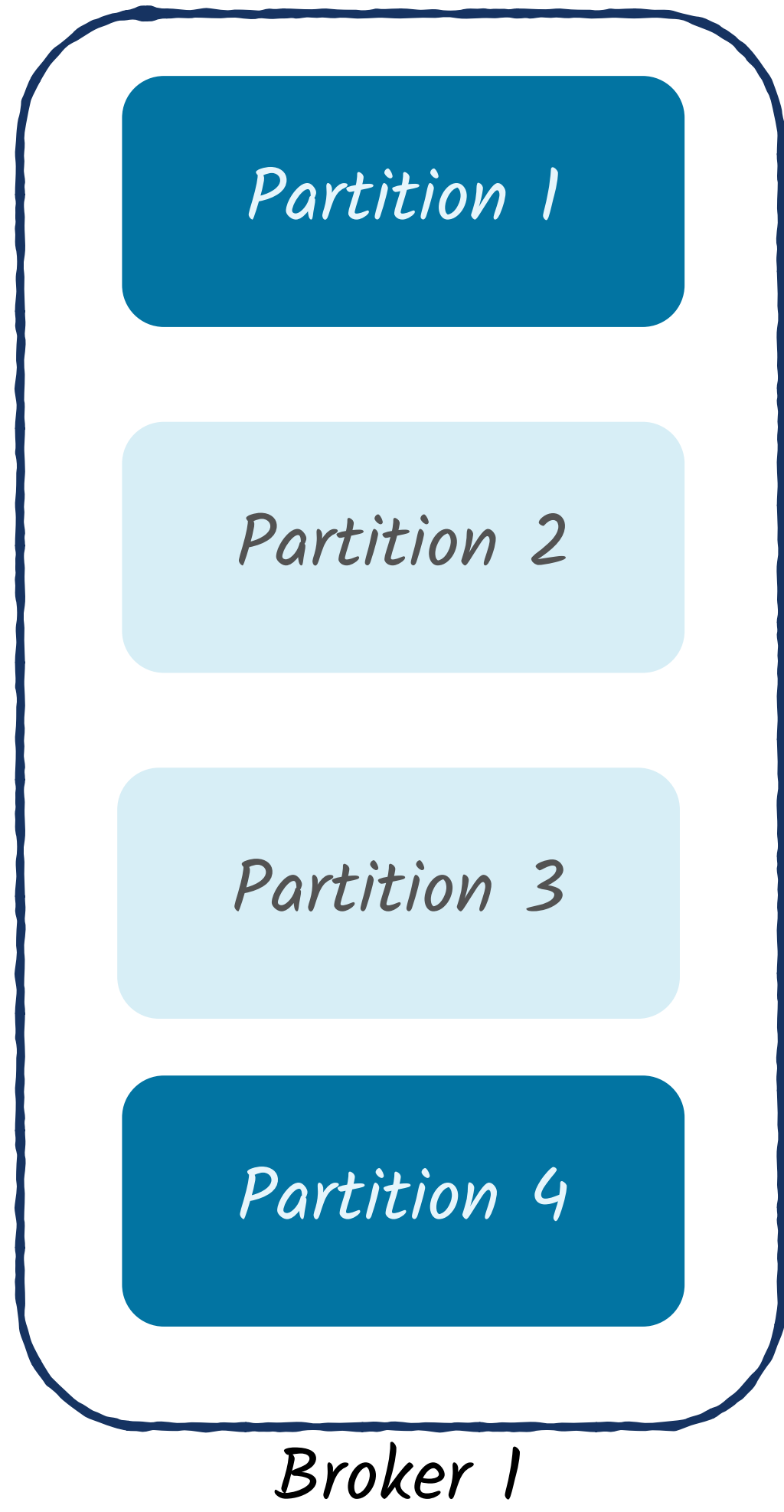
Partition Leadership and Replication

Leader
Follower



Partition Leadership and Replication

Leader
Follower



Partition Leadership and Replication

Leader
Follower

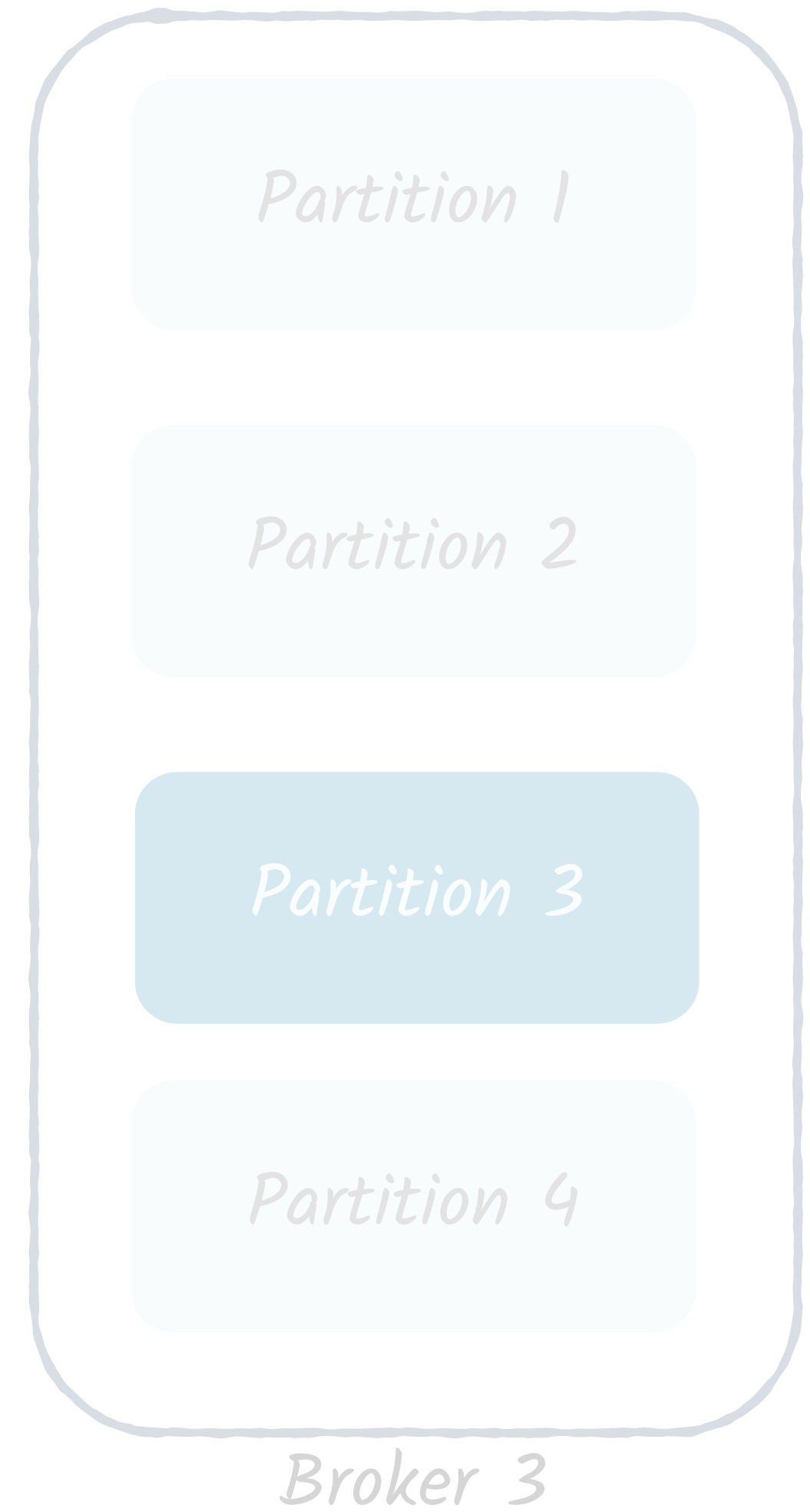
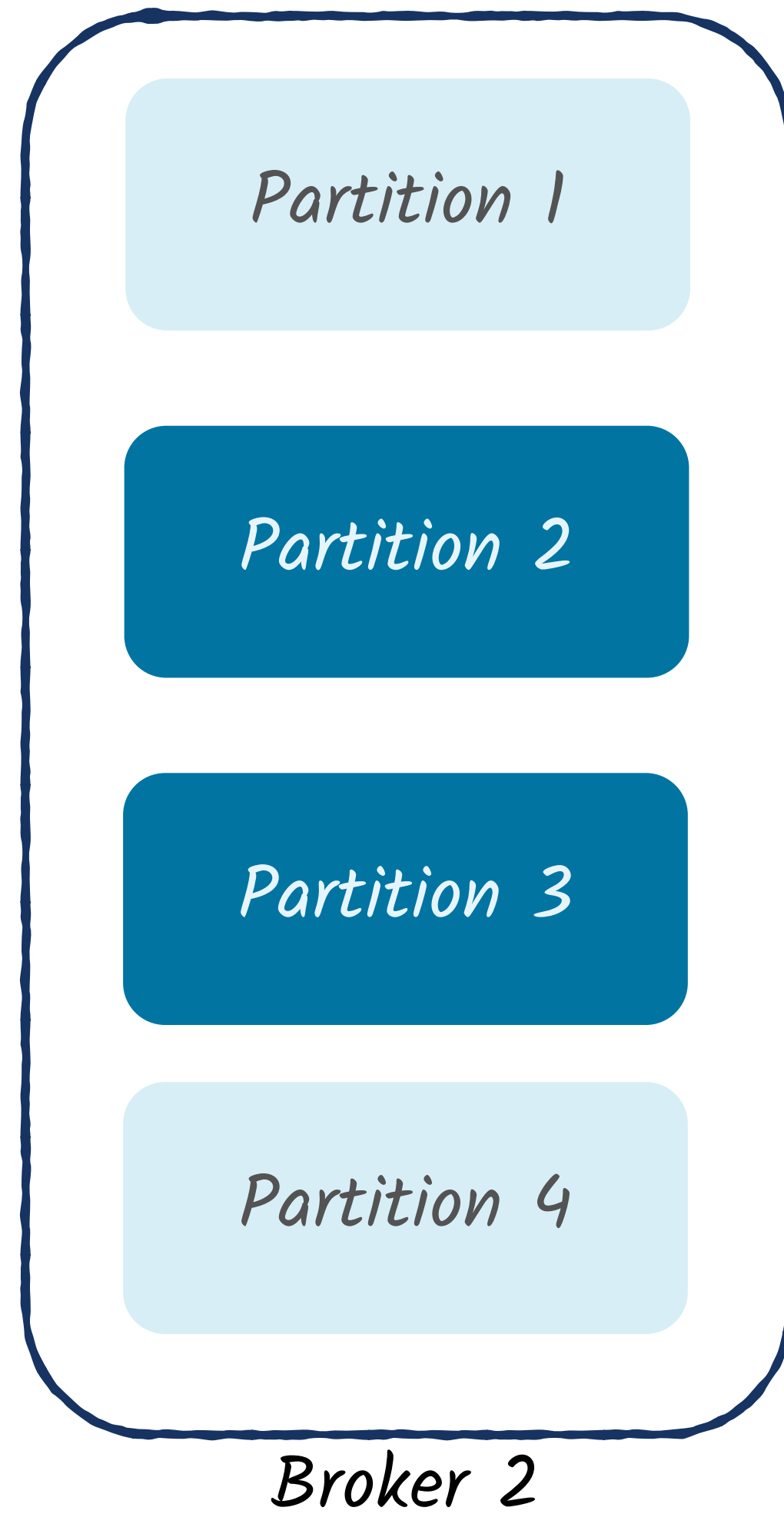
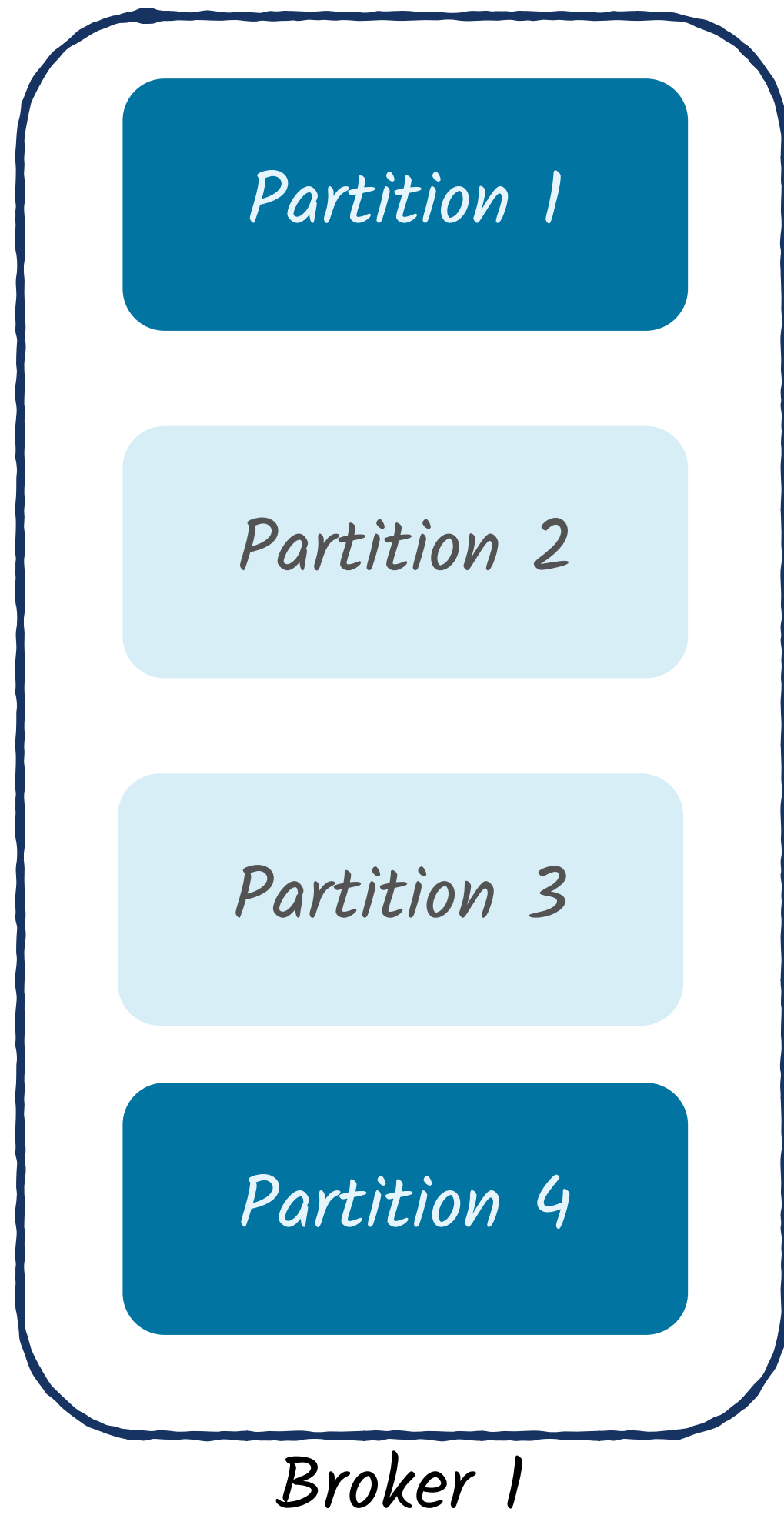
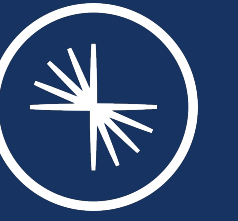




Photo by Raoul Droog on Unsplash

DEMO



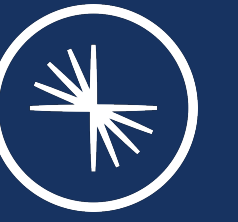
So far, this is

Pretty good

@rmoff

| #BudapestData

| @confluentinc

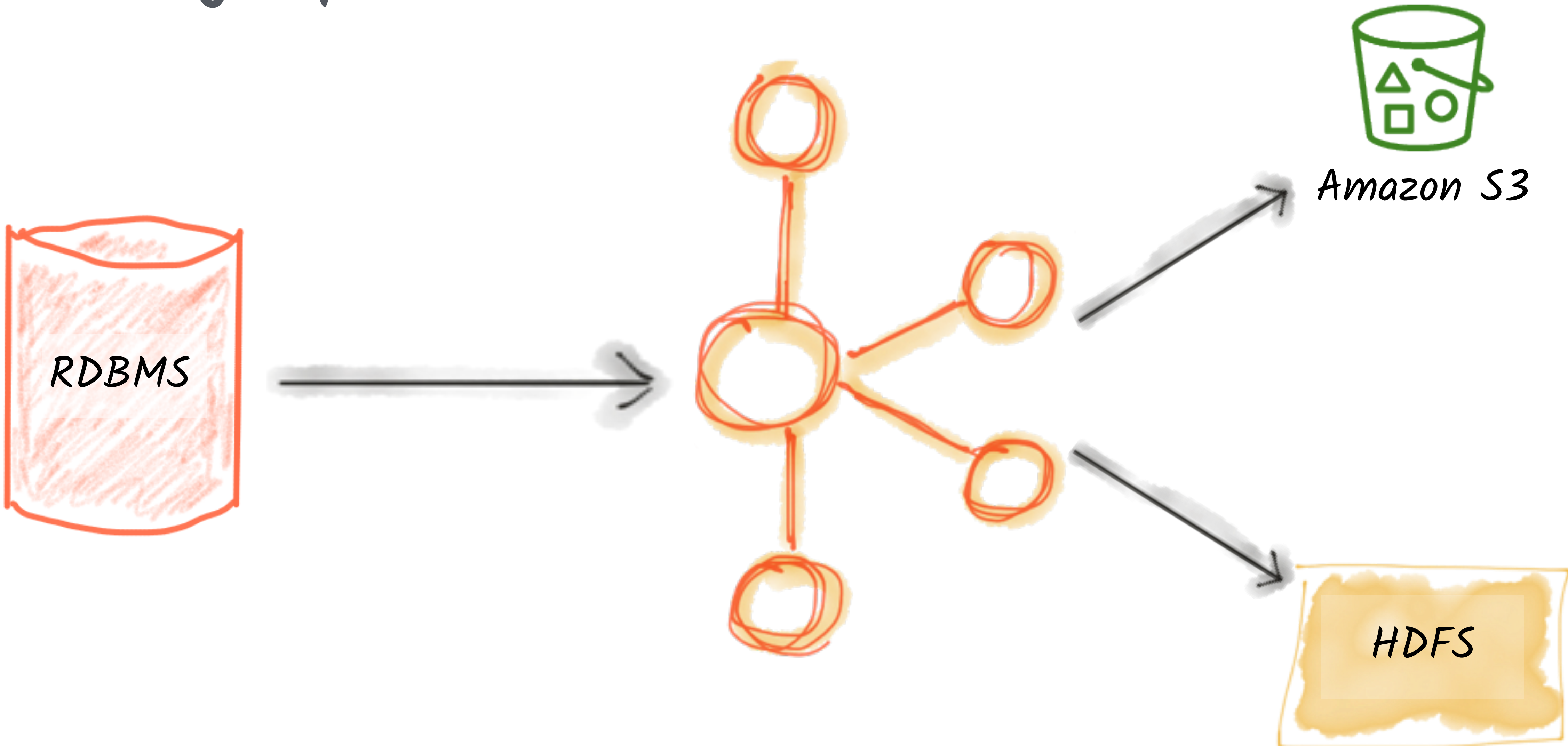


So far, this is

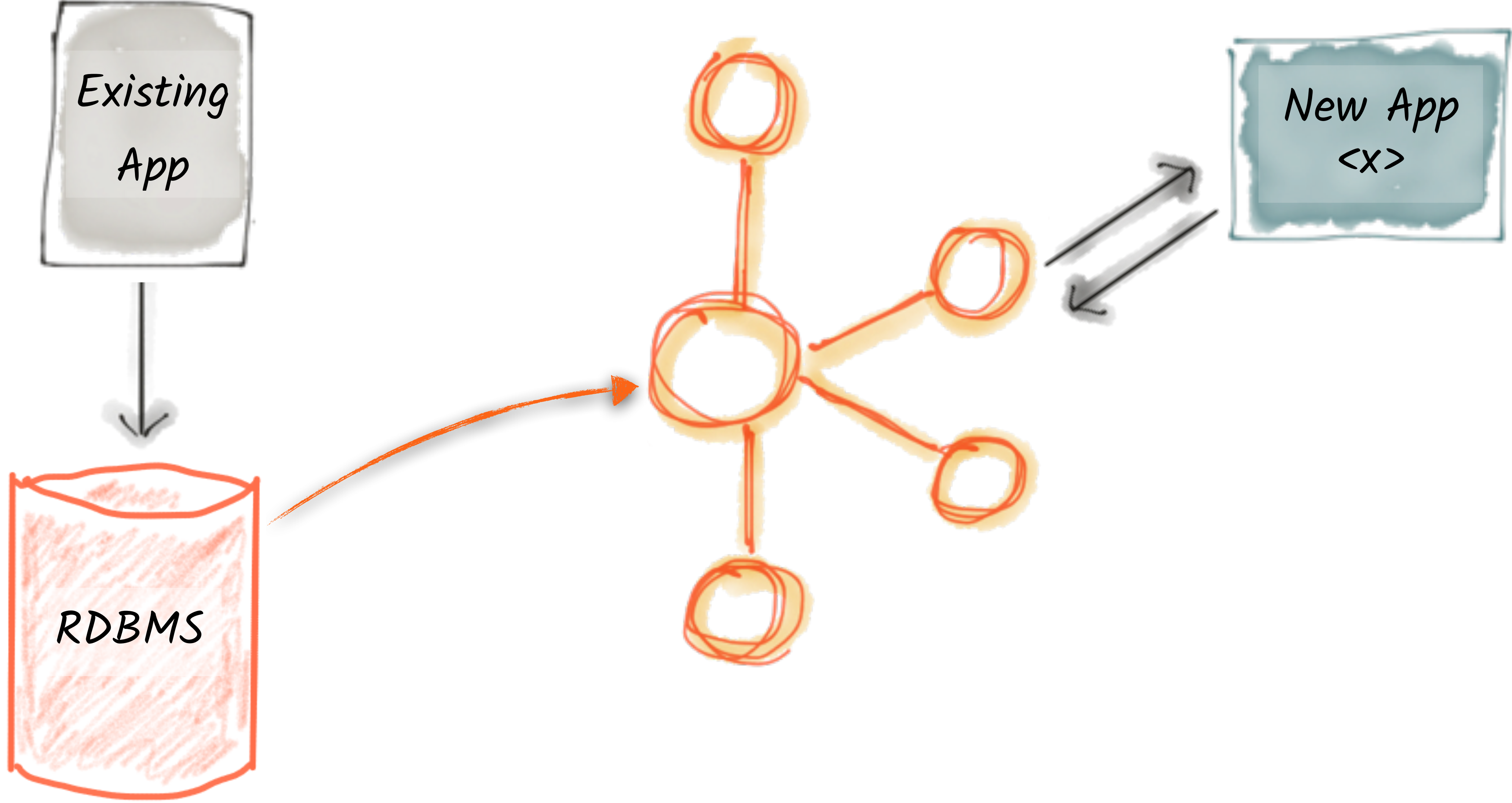
Pretty good

but I've not finished yet...

Streaming Pipelines



Evolve processing from old systems to new

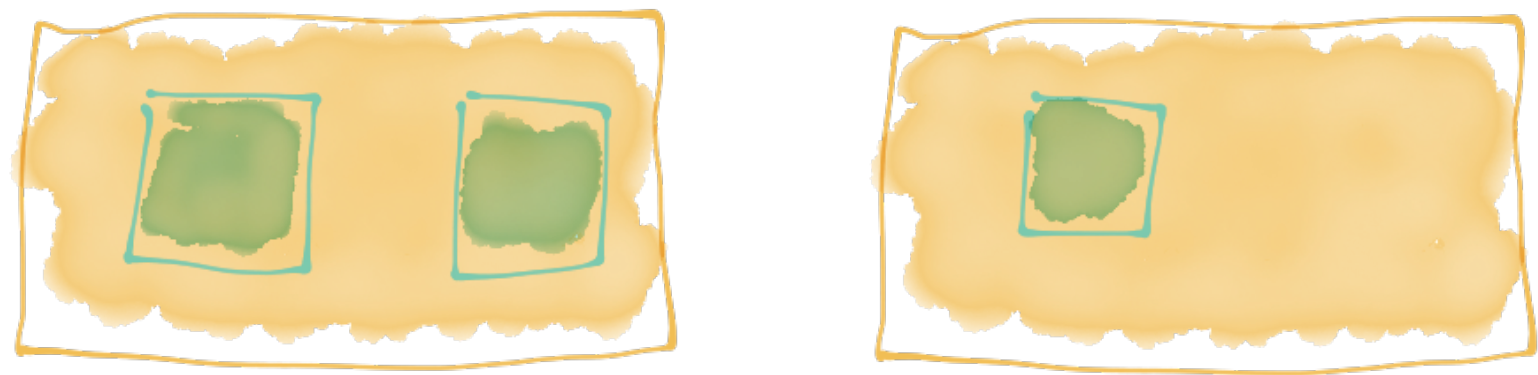




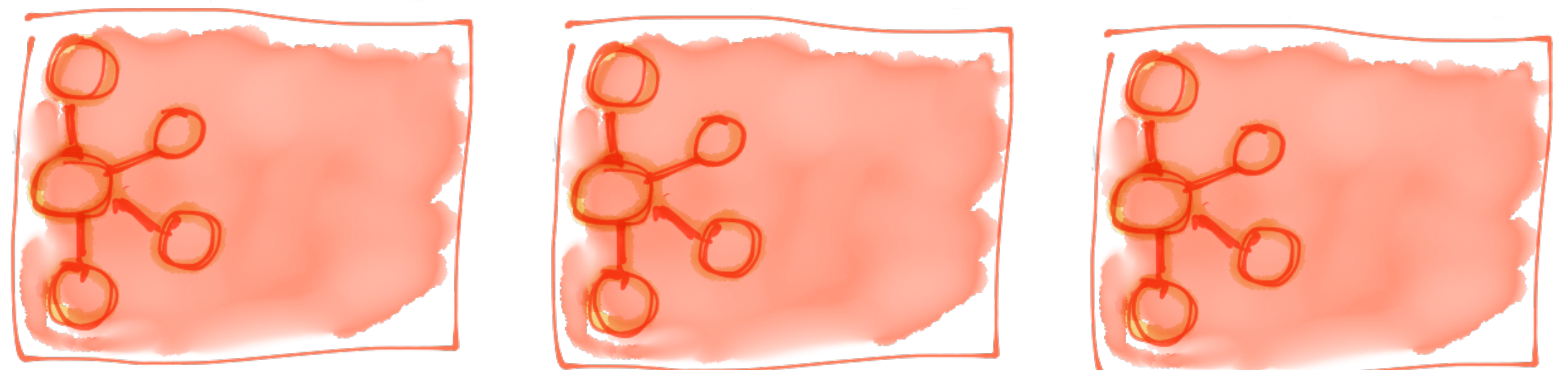
Streaming Integration with Kafka Connect



Sources



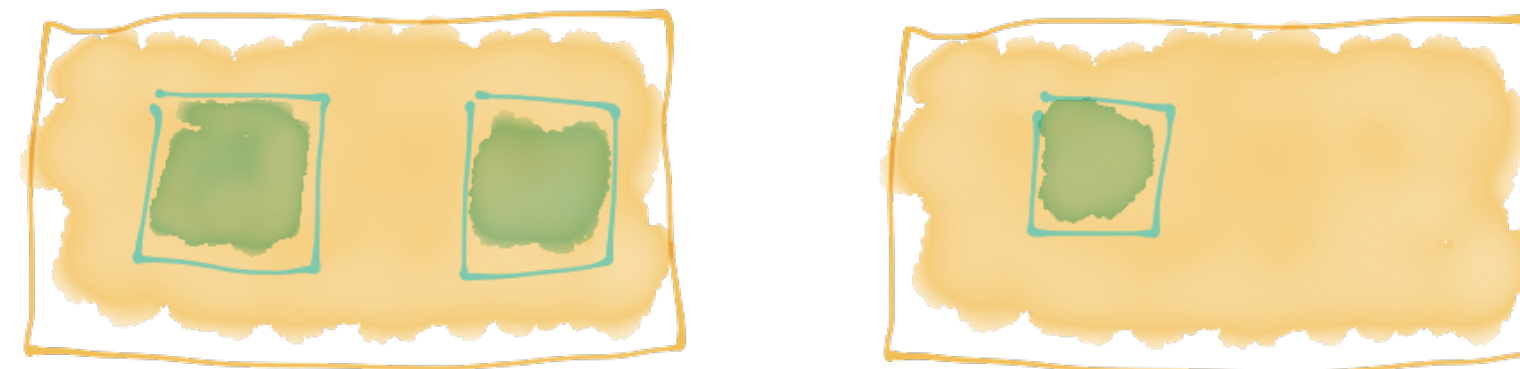
Kafka Connect



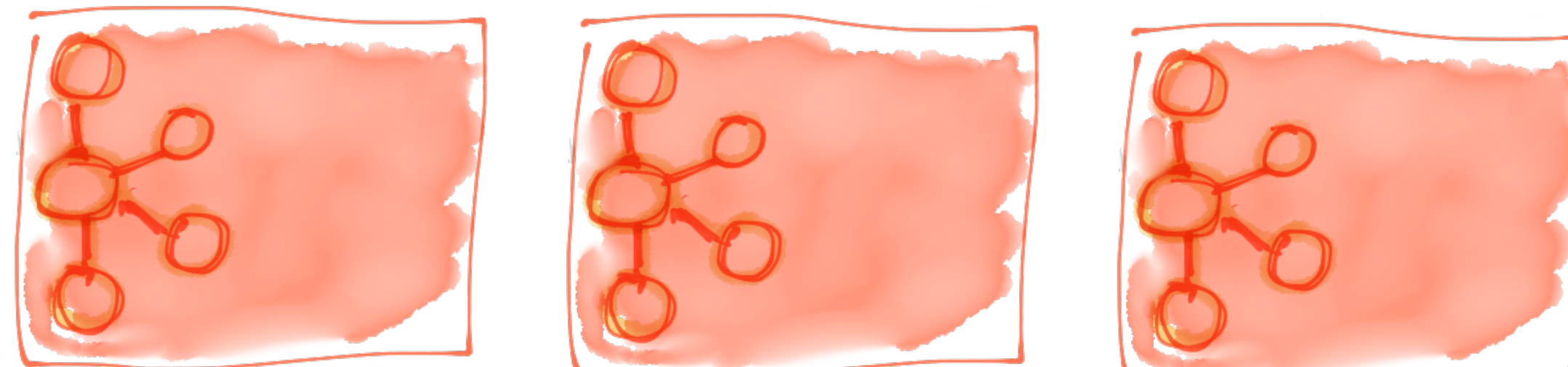
Kafka Brokers

Streaming Integration with Kafka Connect

Sinks

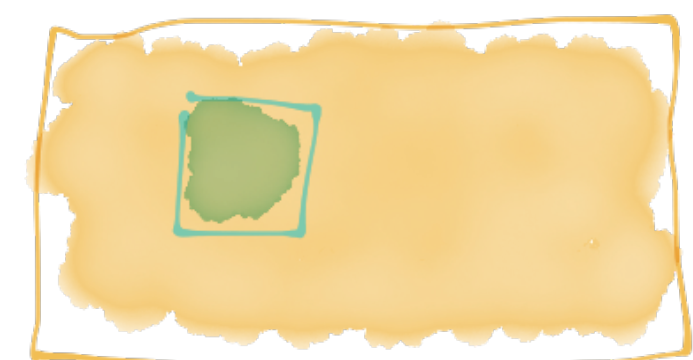
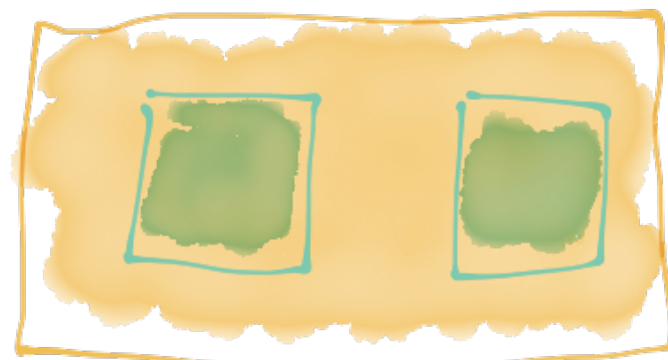
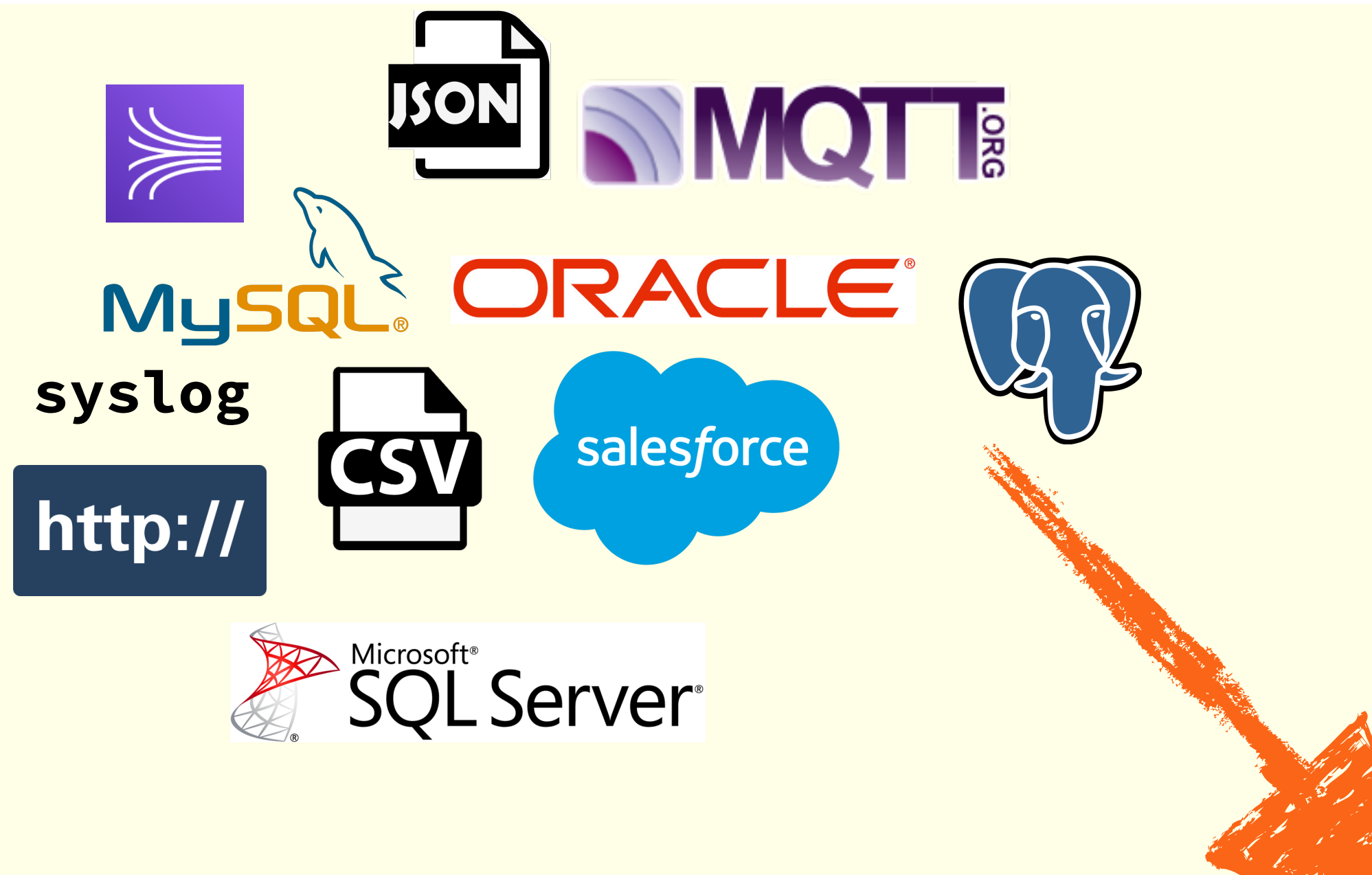


Kafka Connect

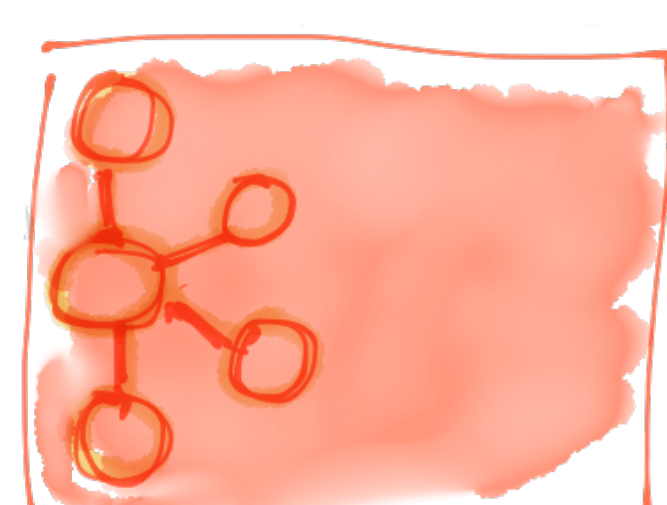
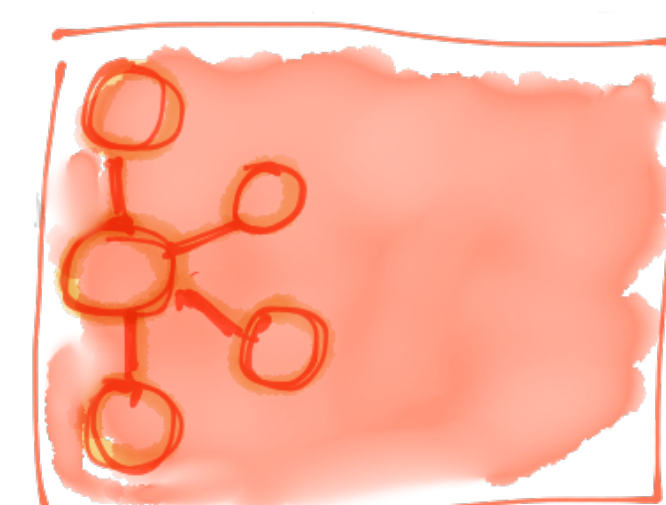
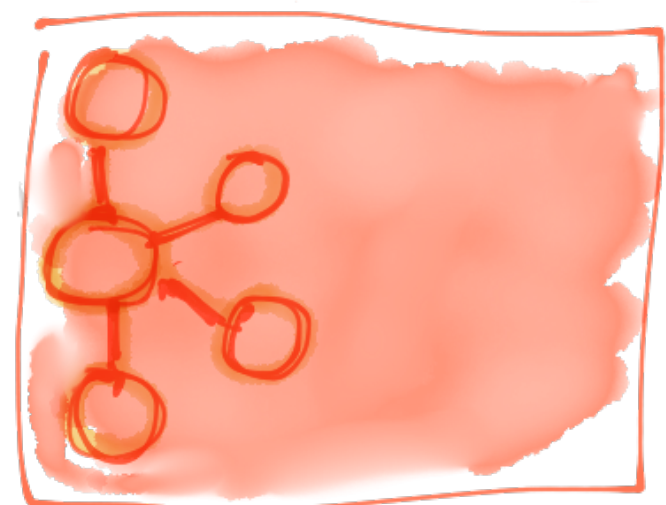


Kafka Brokers

Streaming Integration with Kafka Connect



Kafka Connect



Kafka Brokers

Look Ma, No Code!

```
{
```

```
  "connector.class":
```

```
    "io.confluent.connect.jdbc.JdbcSourceConnector",
```

```
  "connection.url":
```

```
    "jdbc:mysql://asgard:3306/demo",
```

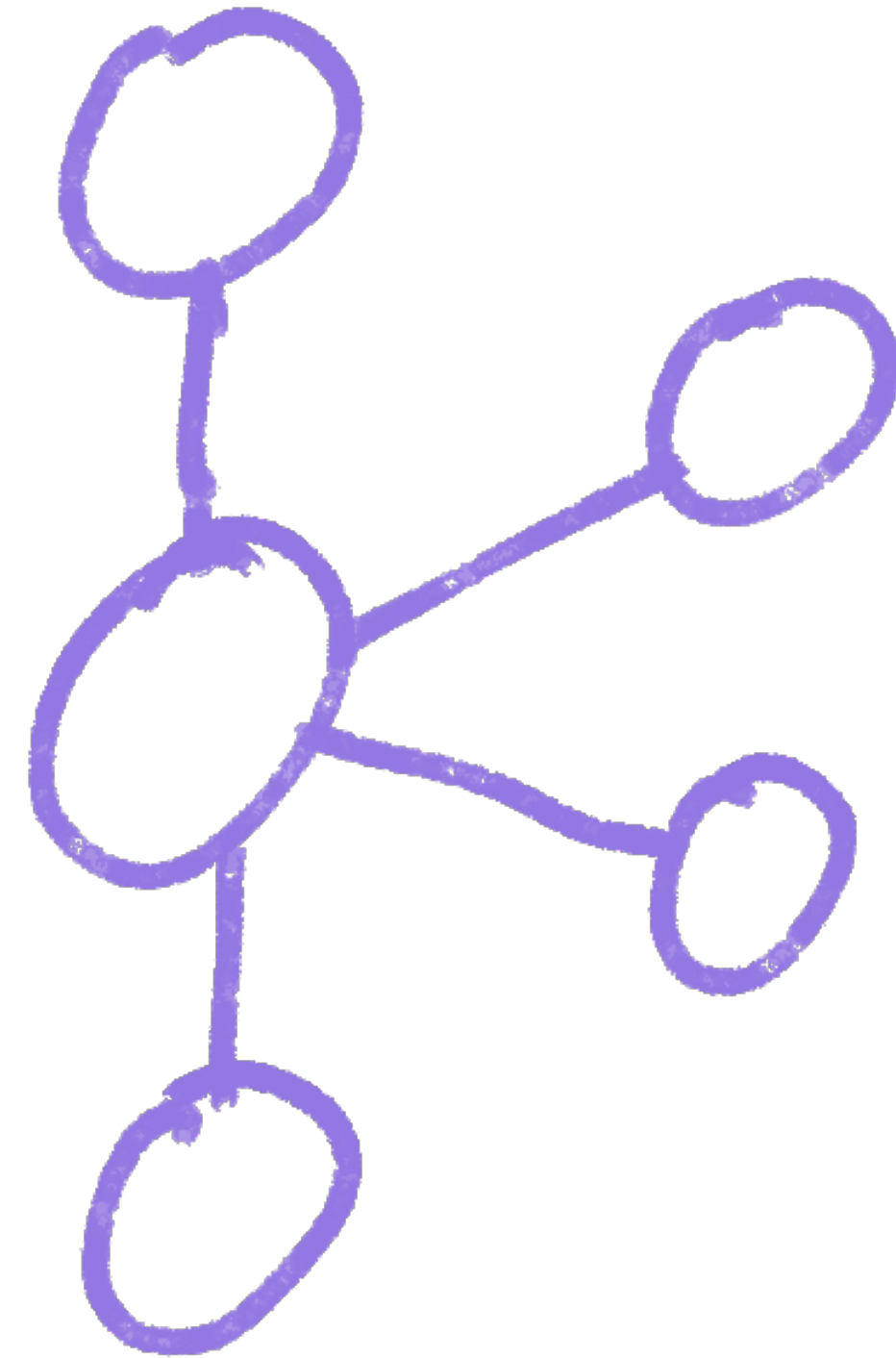
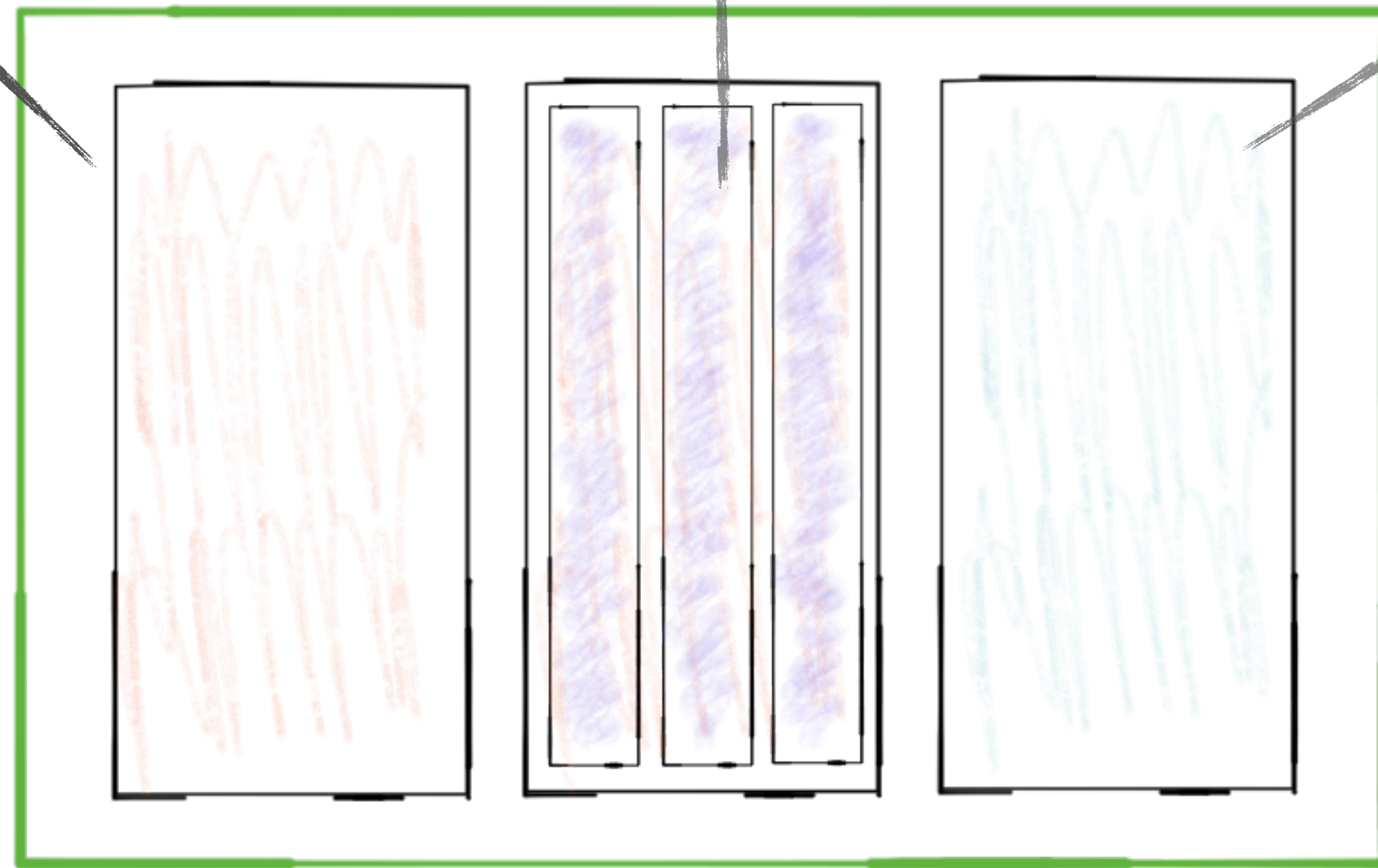
```
  "table.whitelist":
```

```
    "sales,orders,customers"
```

```
}
```

Extensible

Connector *Transform(s)* *Converter*



hub.confluent.io

Confluent Hub

Discover Kafka® connectors
and more

Filters

Plugin type ⊕

- Sink
- Source
- Transform
- Converter

Enterprise support ⊕

- Confluent supported
- Partner supported
- None

Results (1)


[+ Submit a plugin](#)

Kafka Connect JDBC

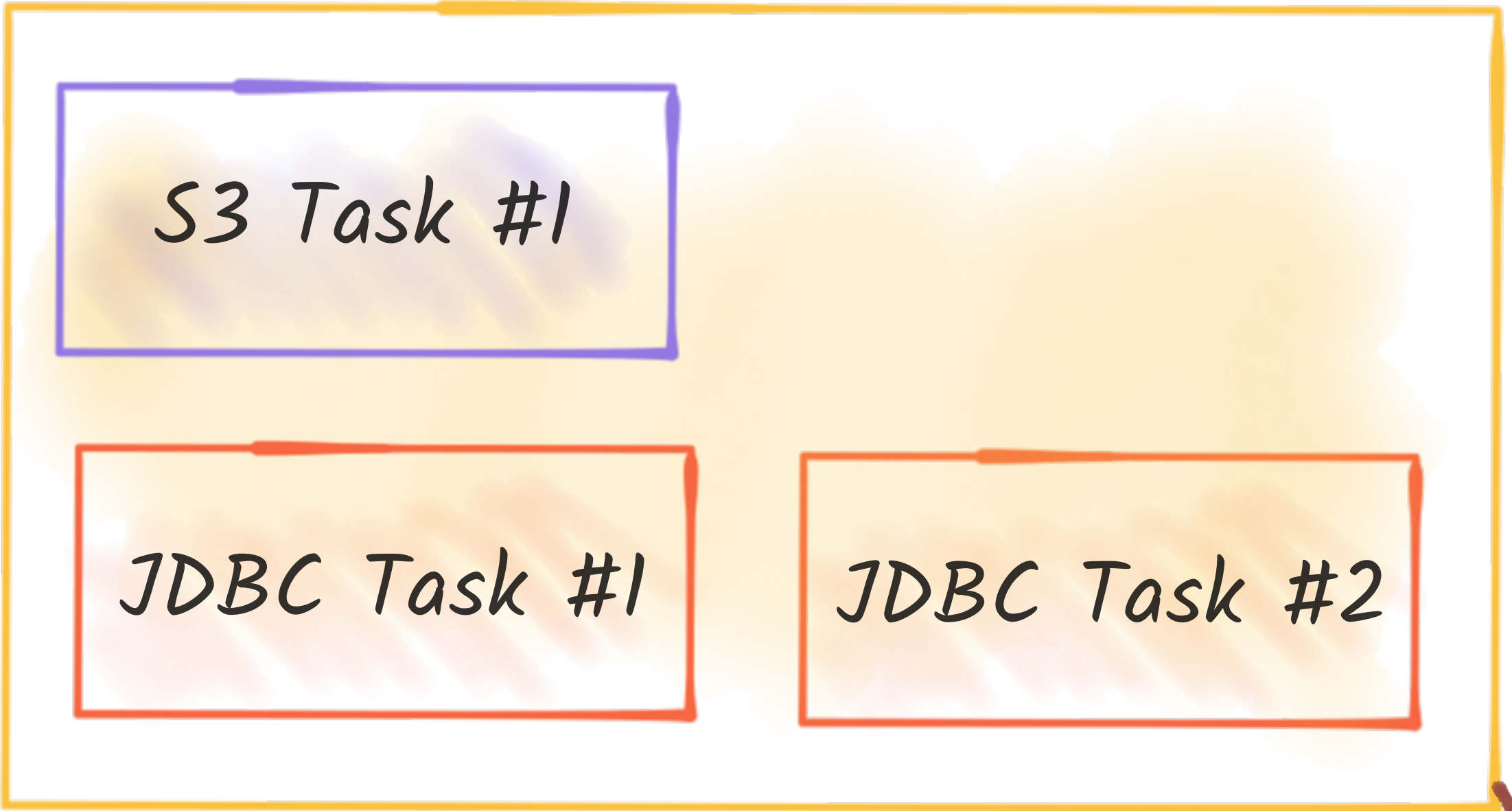
SINK, SOURCE CONNECTOR

The JDBC source and sink connectors allow you to exchange data between relational databases and Kafka. The JDBC source connector allows you to import data from any relational database with a JDBC driver into Kafka topics

Enterprise support: Confluent supported	Verification: Confluent built	License: Free
Installation: Confluent Hub CLI, Download	Author: Confluent, Inc.	Version: 5.4.1

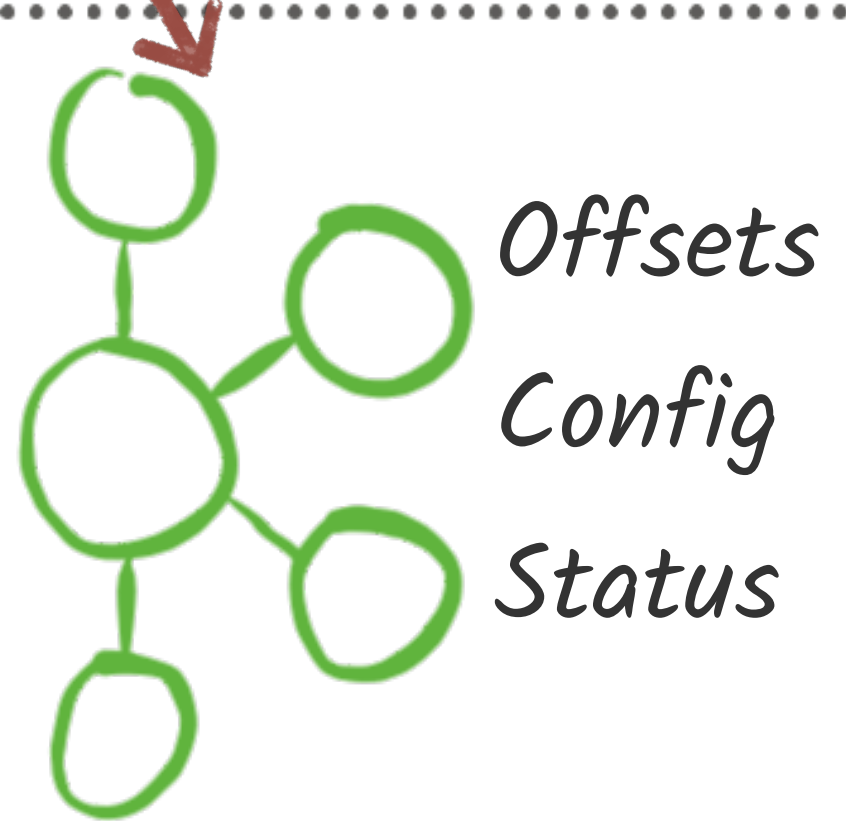


Single Kafka Connect node

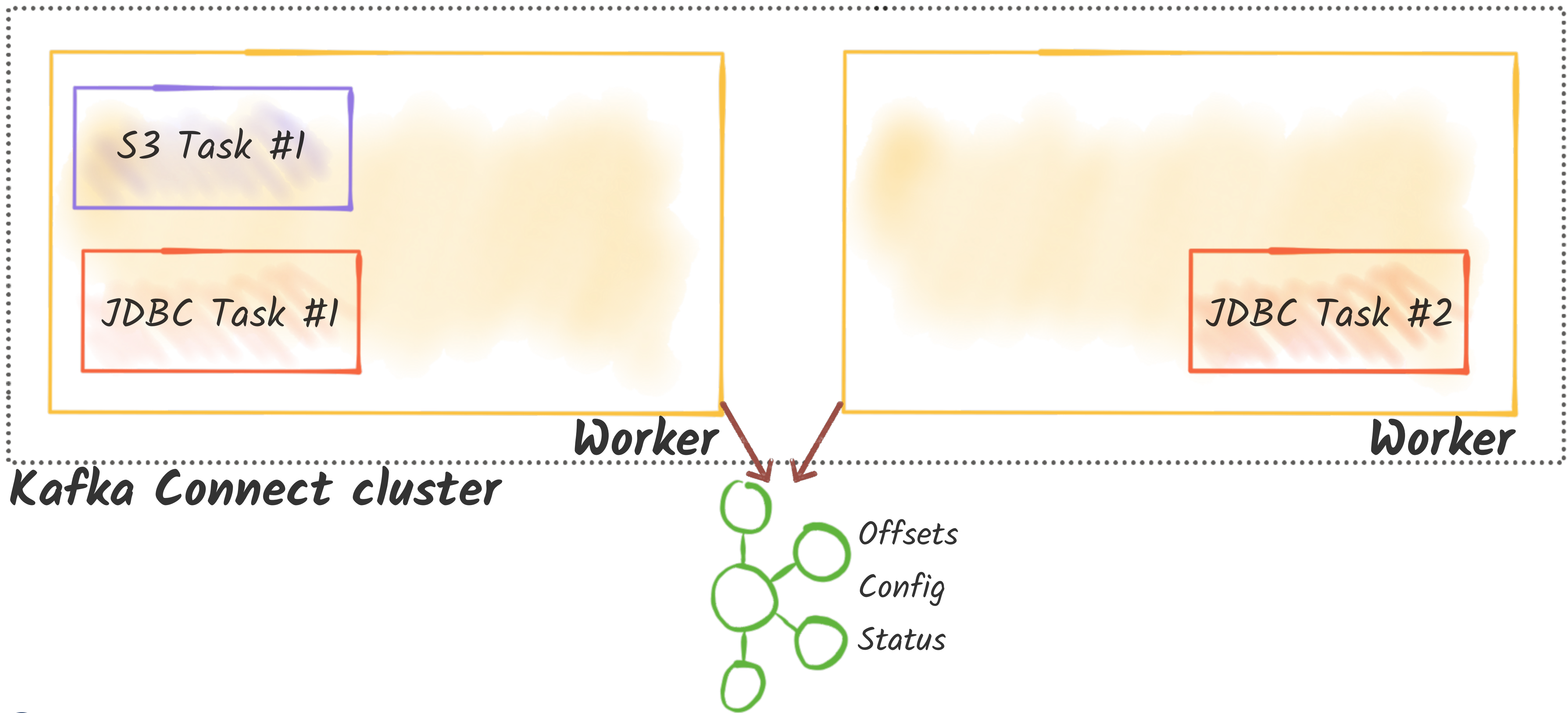


Worker

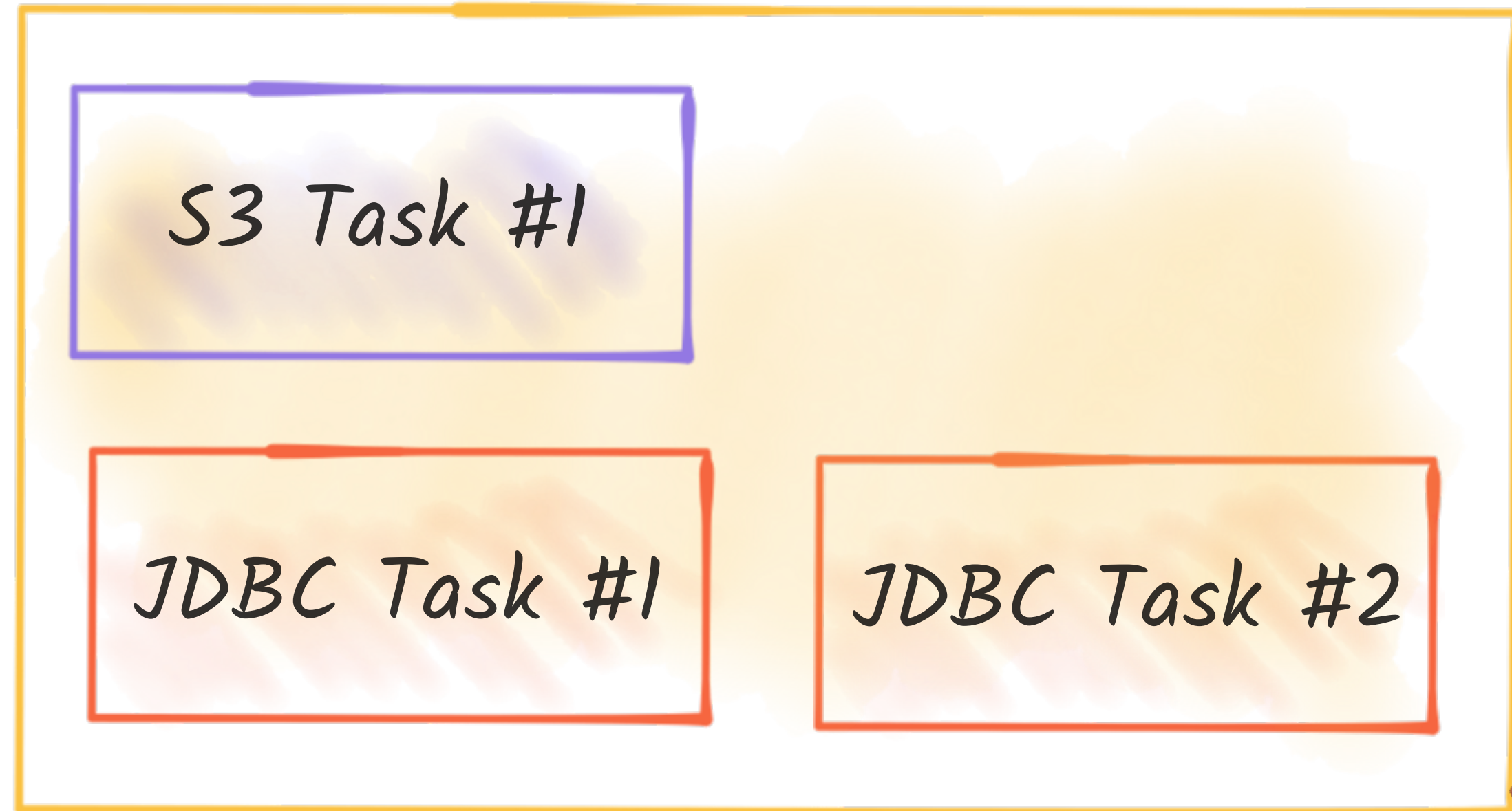
Kafka Connect cluster



Kafka Connect - scalable and fault-tolerant

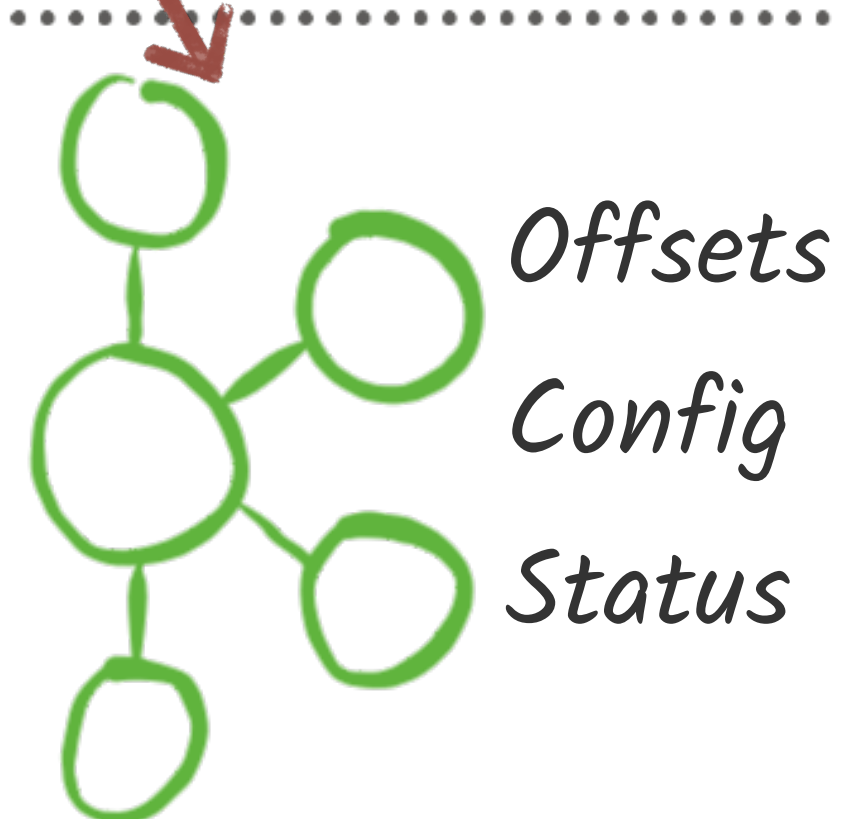


Automatic fault tolerance



Worker

Kafka Connect cluster

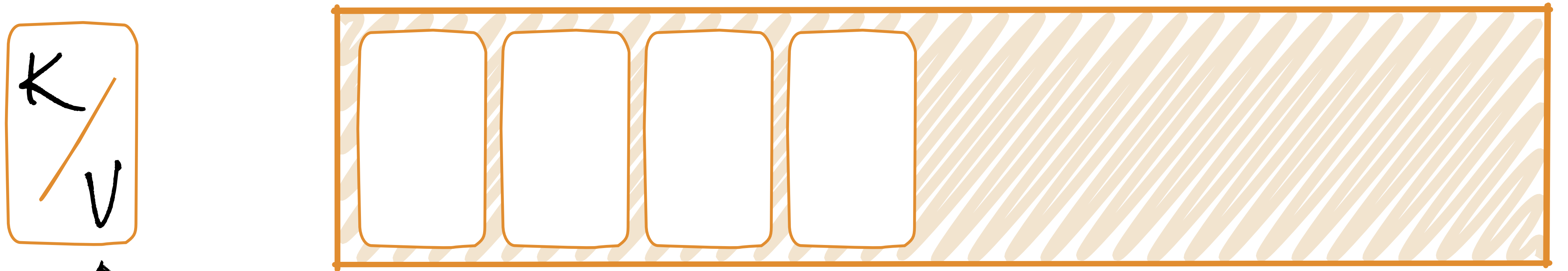










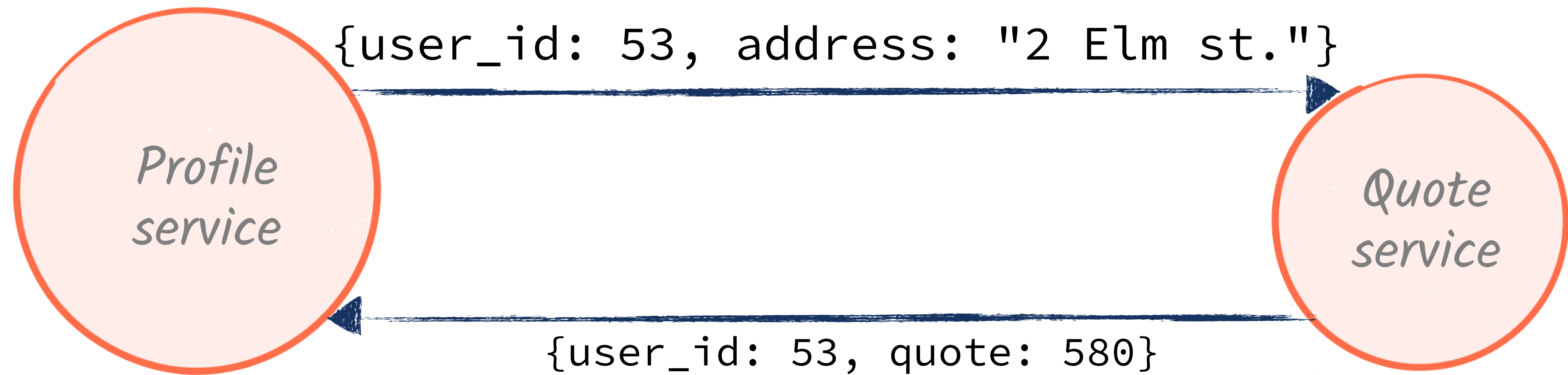


Wait...what's this?

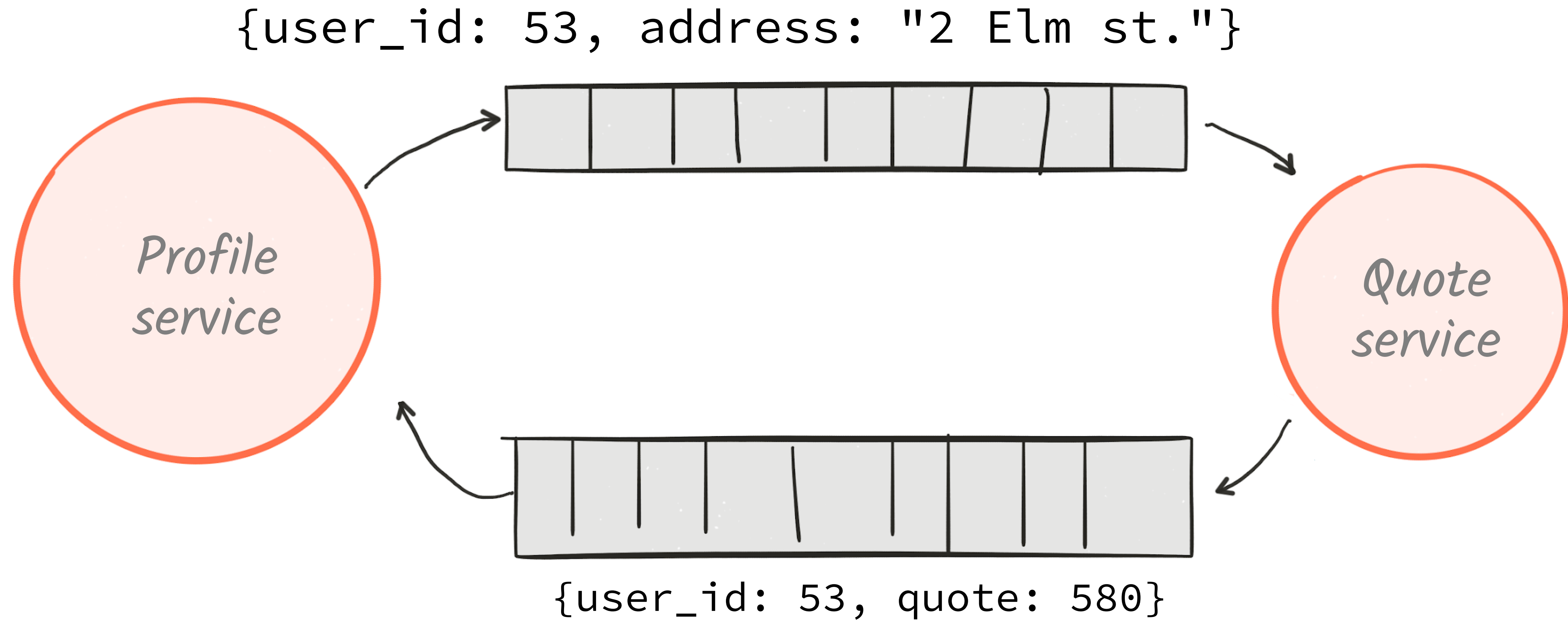
How do you serialise your data?

Avro Protobuf JSON Schema JSON CSV

APIs are contracts between services

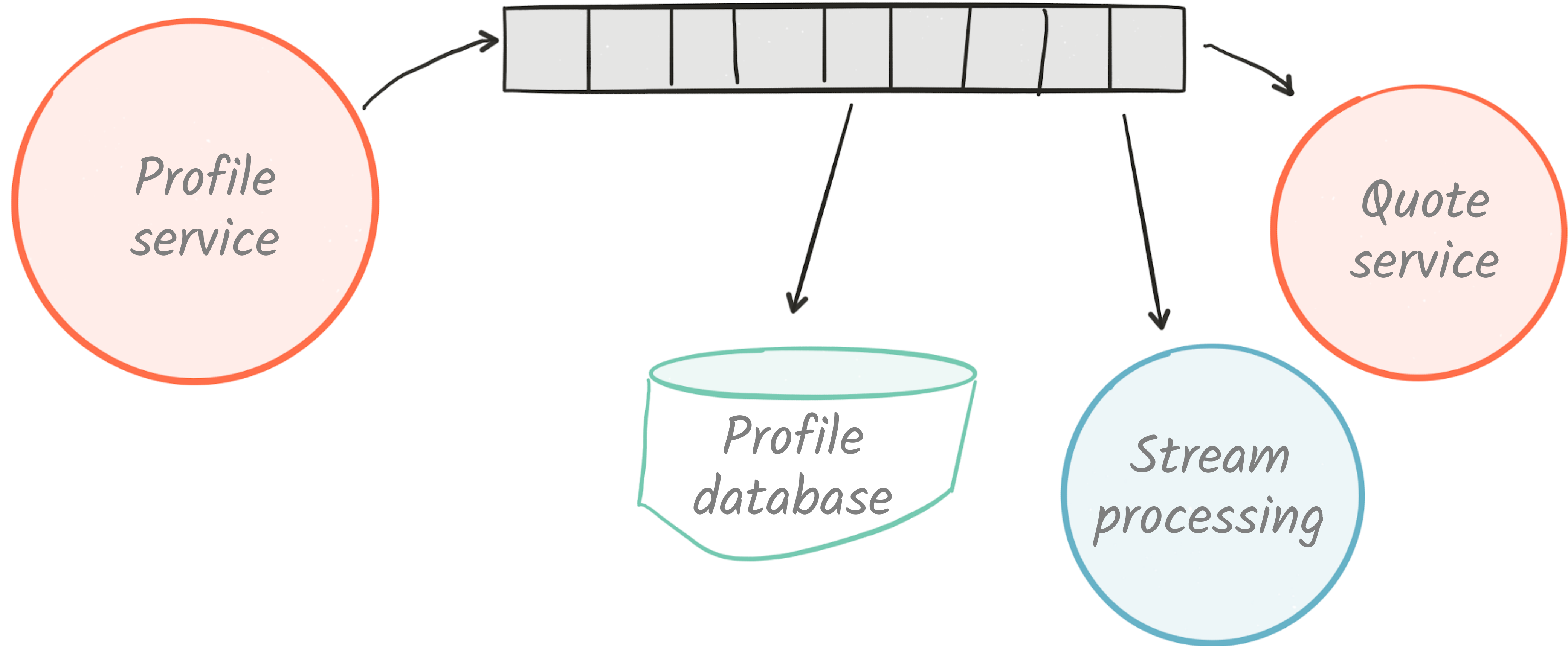


But not all services



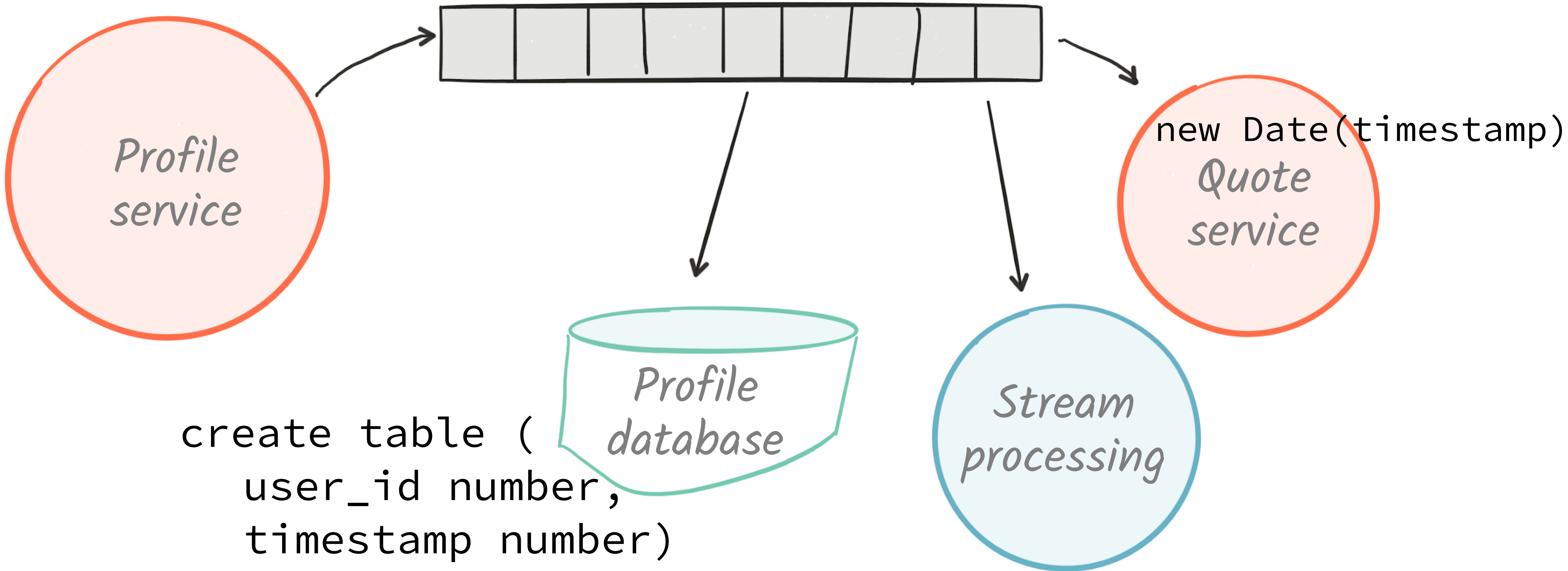
And naturally...

```
{user_id: 53, address: "2 Elm st."}
```



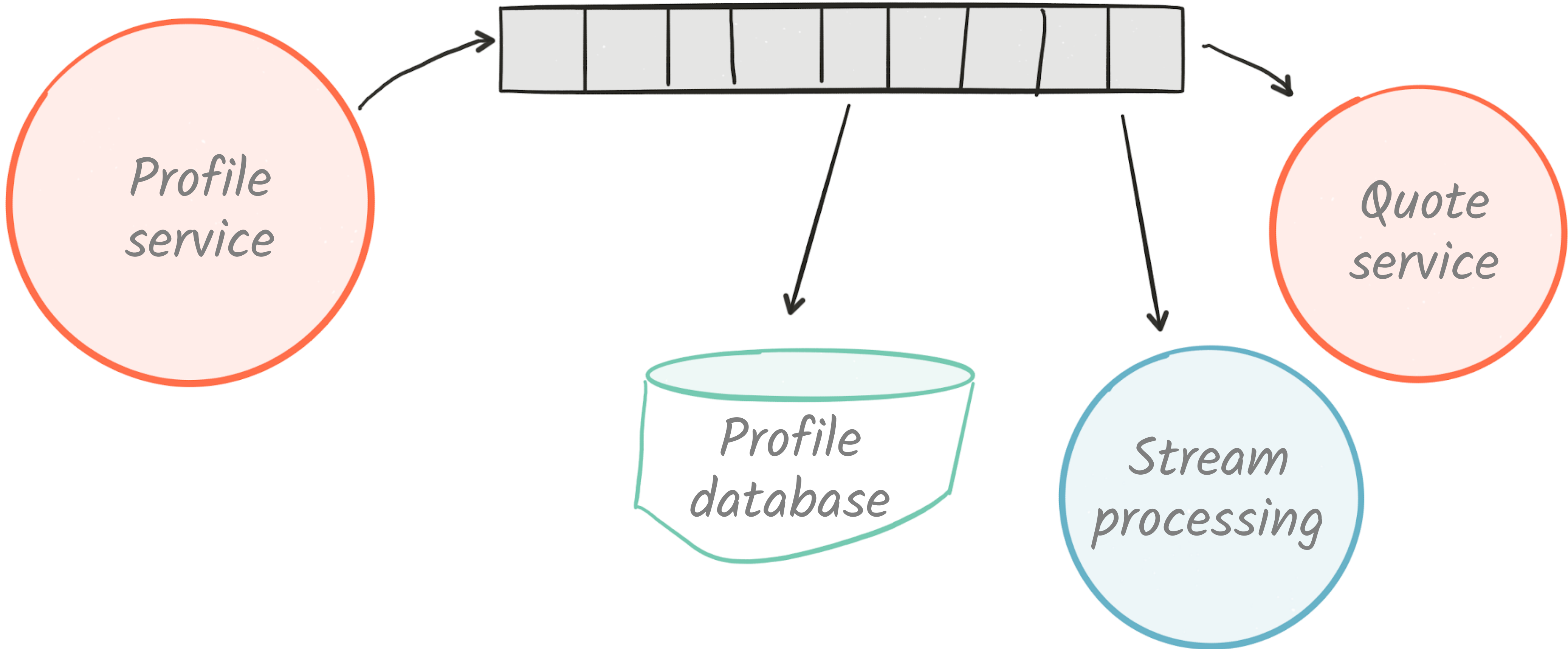
Schemas are about how teams work together

{user_id: 53, timestamp: 1497842472}



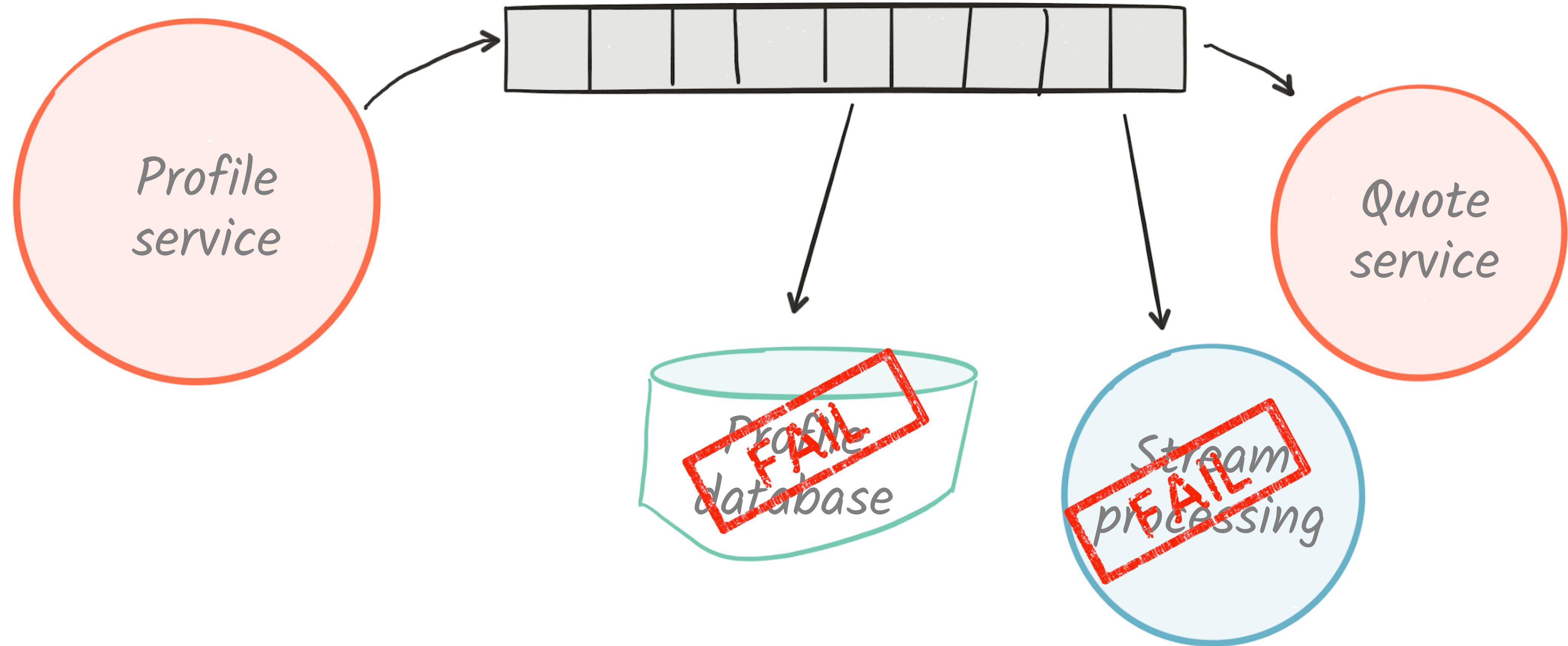
Things change...

{user_id: 53, timestamp: "June 28, 2017 4:00pm"}



Moving fast and breaking things

{user_id: 53, timestamp: "June 28, 2017 4:00pm"}



Lack of schemas - Coupling teams and services

```
2001 2001 Citrus Heights-Sunrise Blvd  
Citrus_Hghts 60670001 3400293 34 SAC  
Sacramento SV Sacramento Valley SAC  
Sacramento County APCD SMA8 Sacramento  
Metropolitan Area CA 6920 Sacramento 28 6920  
13588 7400 Sunrise Blvd 95610 38 41 56  
38.6988889 121 16 15.98999977  
-121.271111 10 4284781 650345 52
```


Serialisation & Schemas

Avro Protobuf JSON Schema JSON CSV

Serialisation & Schemas

Avro



Protobuf



JSON Schema




JSON



CSV

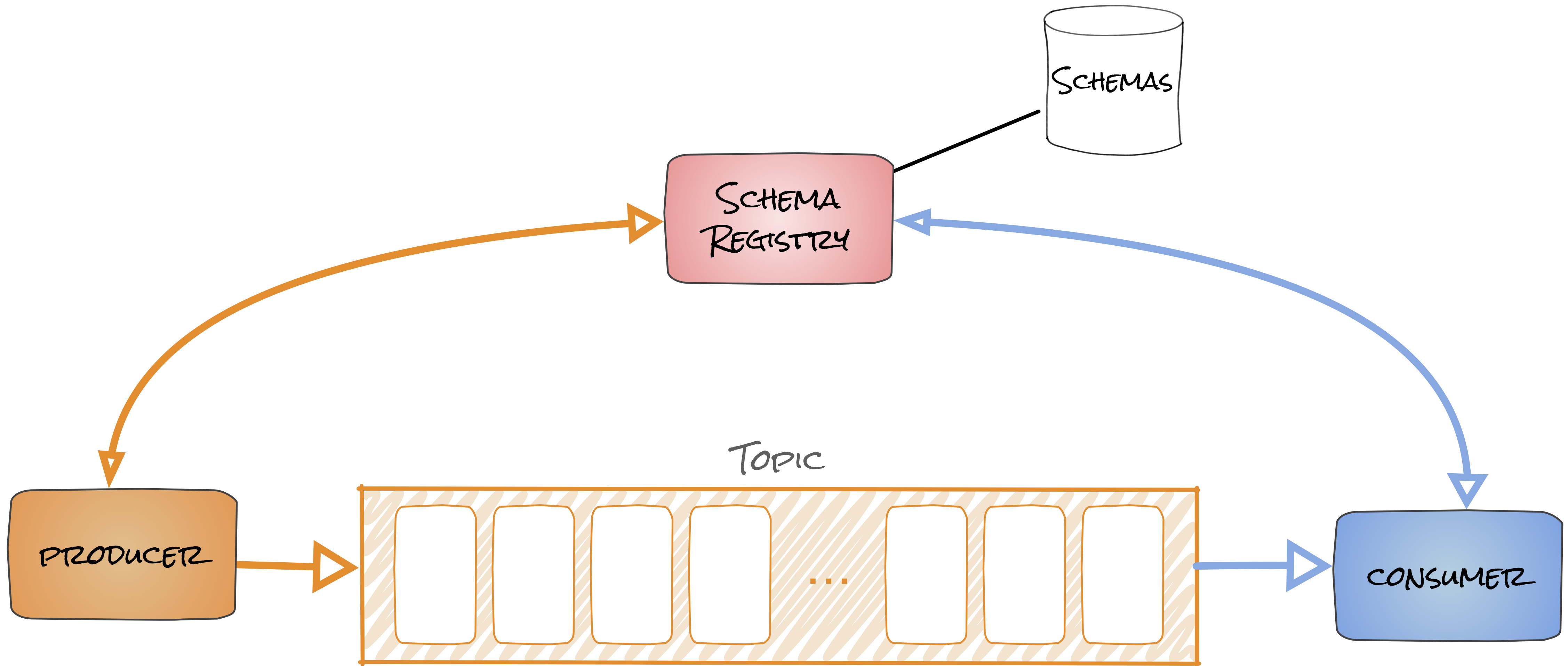


 **Gwen (Chen) Shapira**
@gwenshap

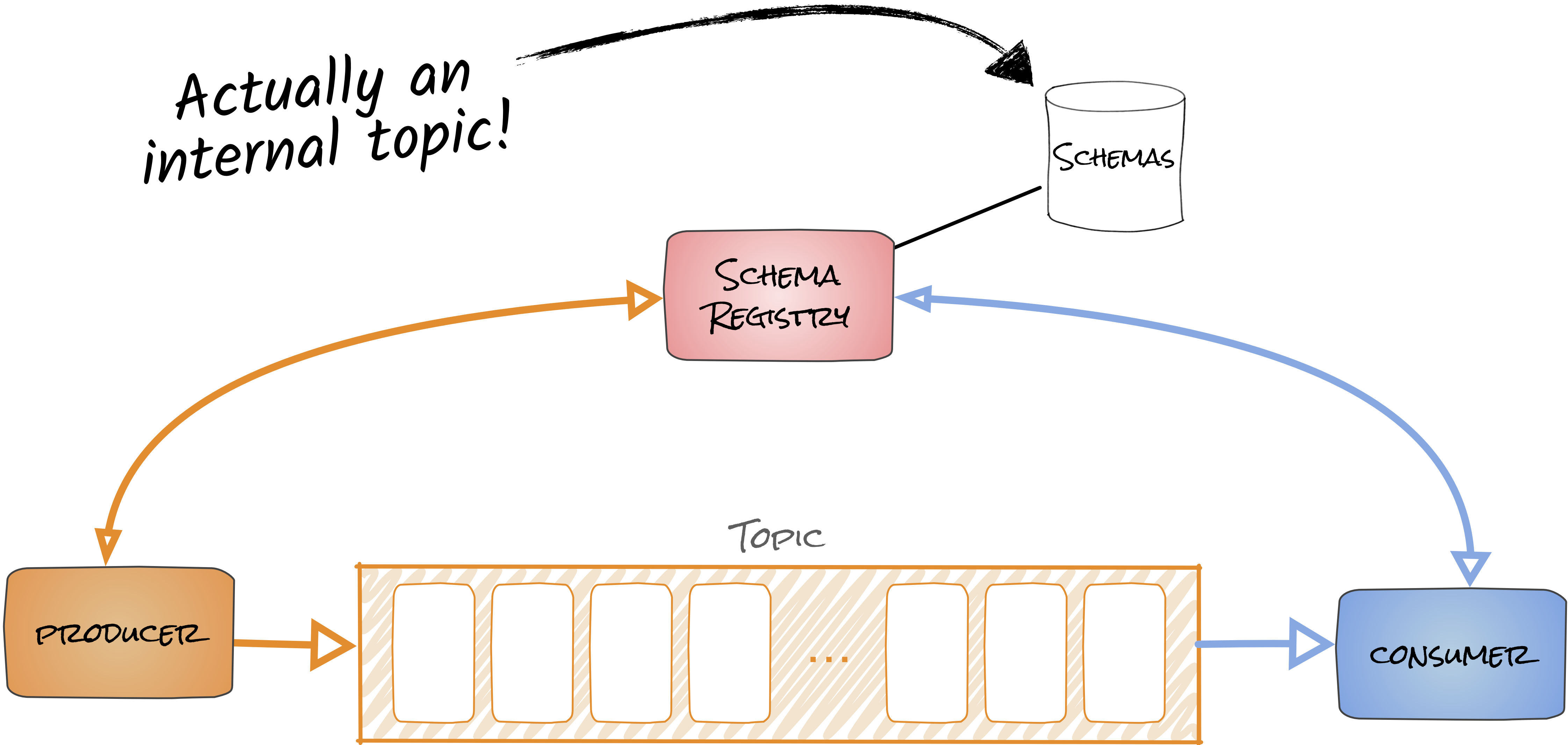
If your dev process doesn't validate schema compatibility somewhere between your IDE and production - you are screwed and don't know it.

5:50 AM - 5 Apr 2017

https://qconnewyork.com/system/files/presentation-slides/qcon_17_-_schemas_and_apis.pdf

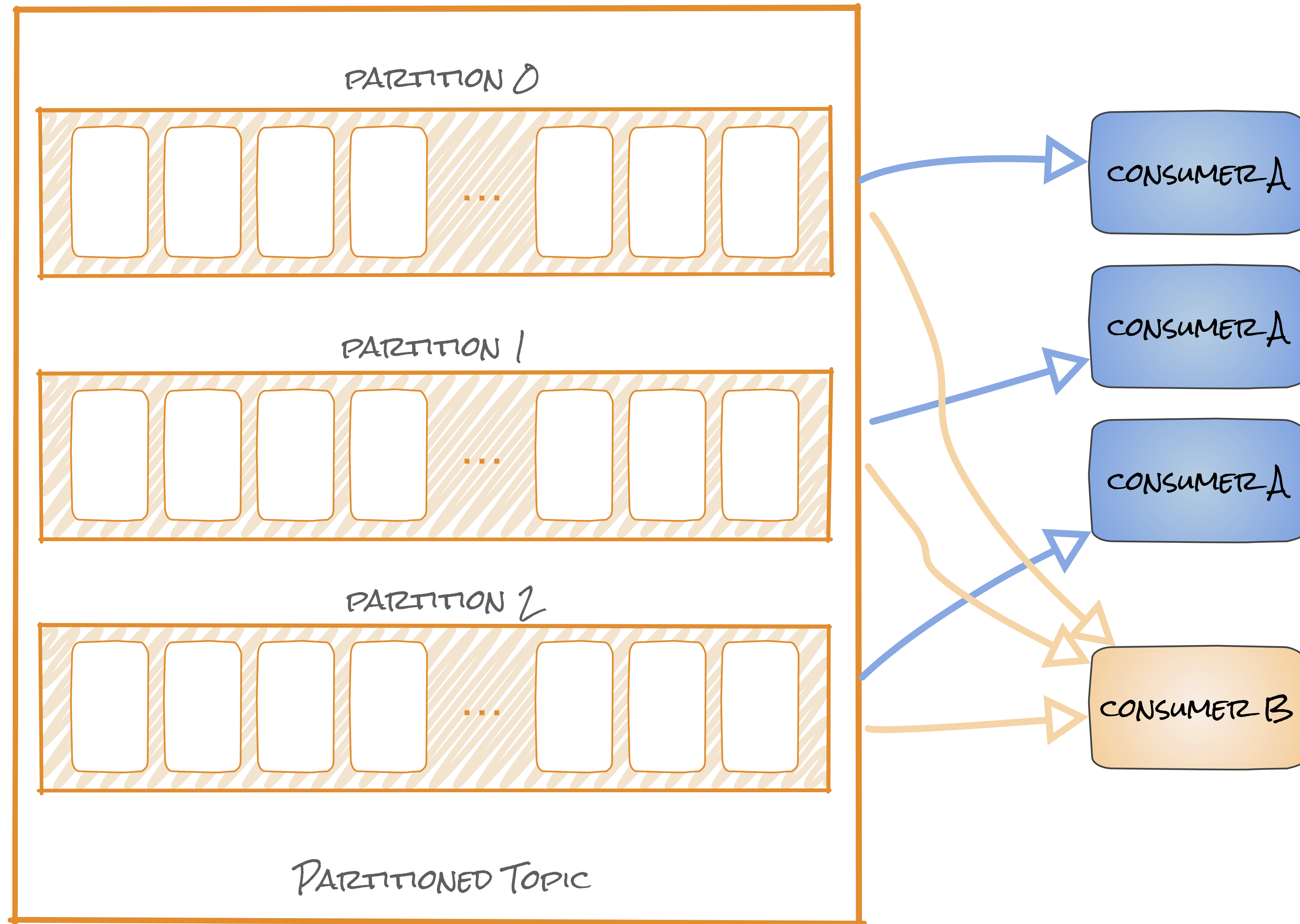


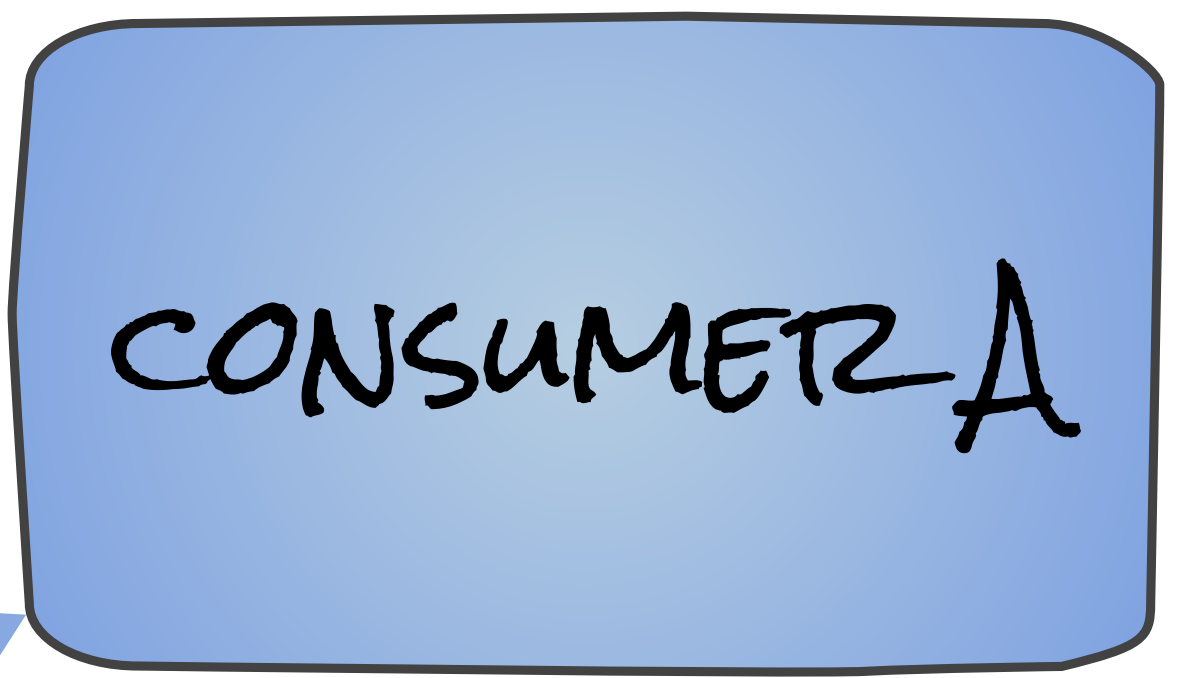
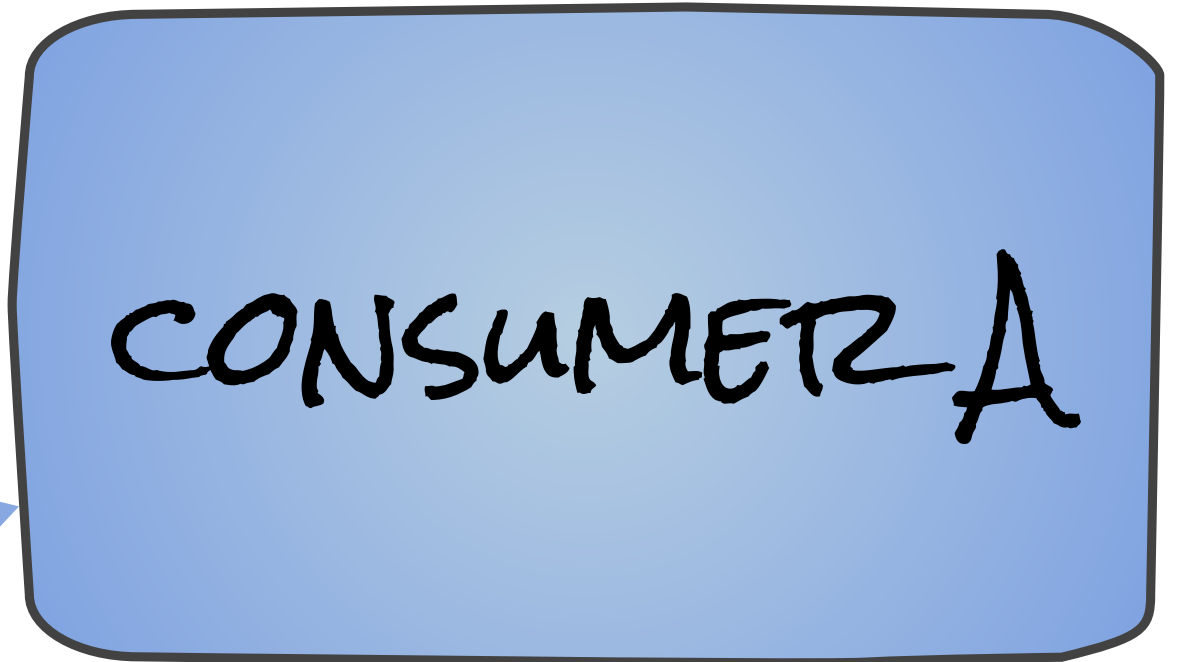
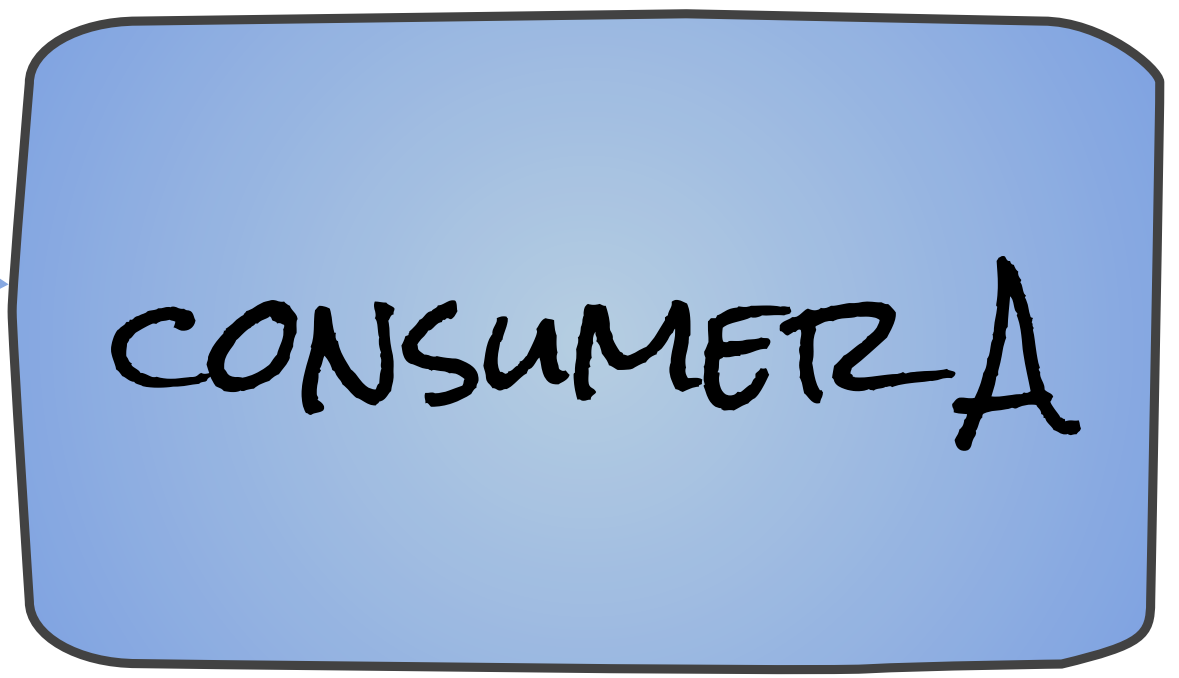
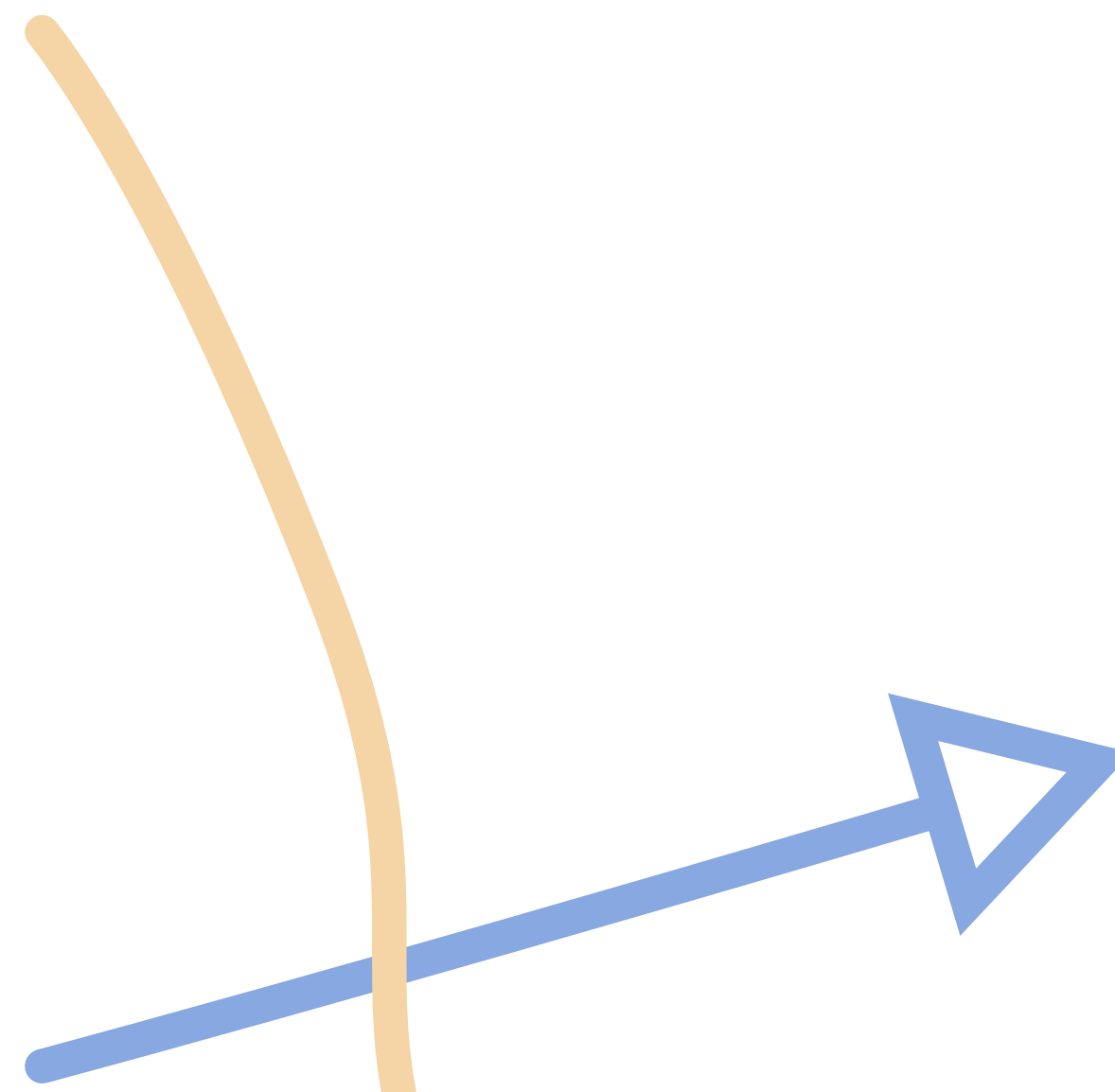
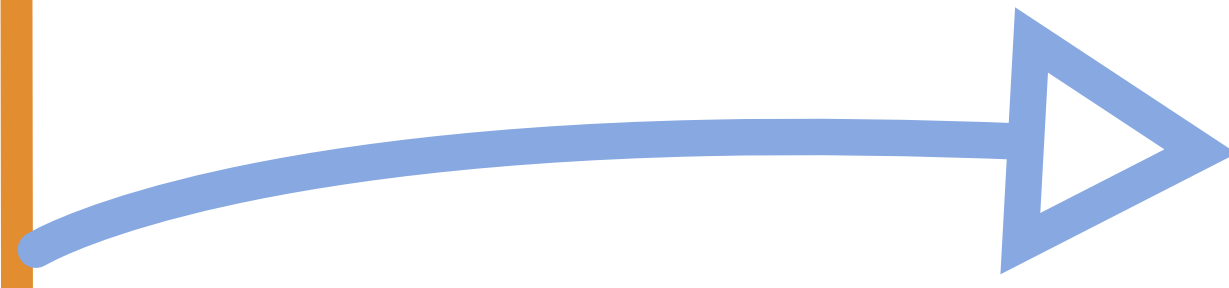
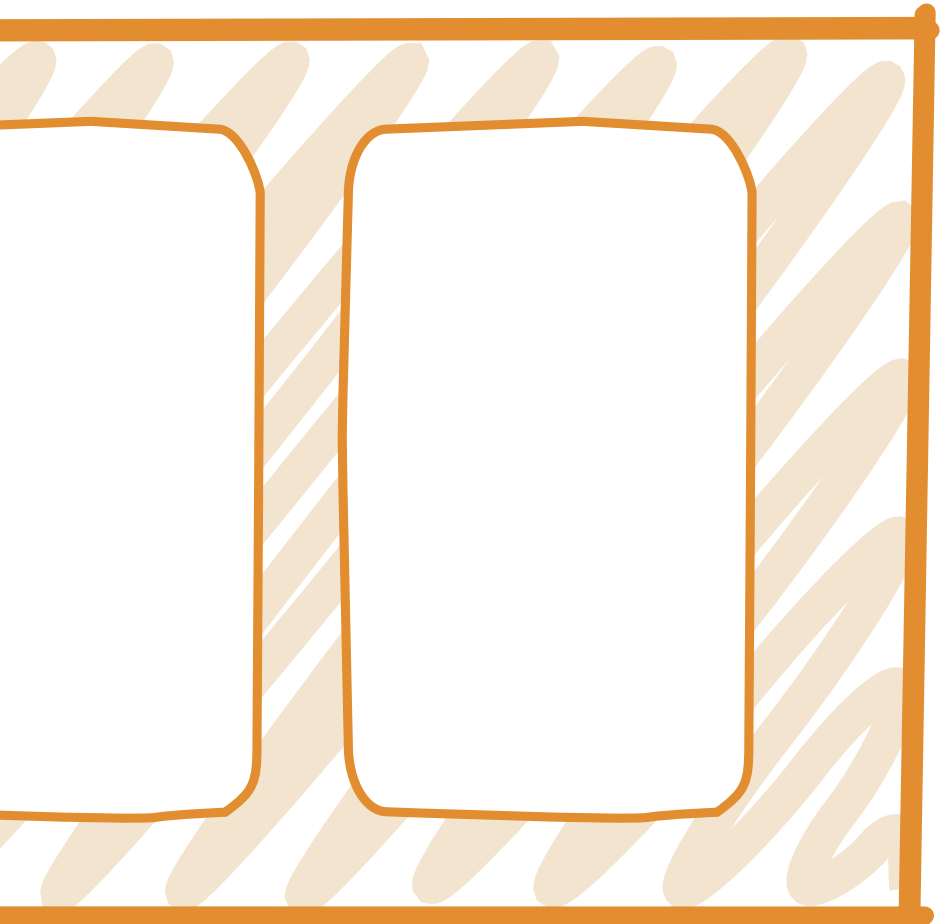
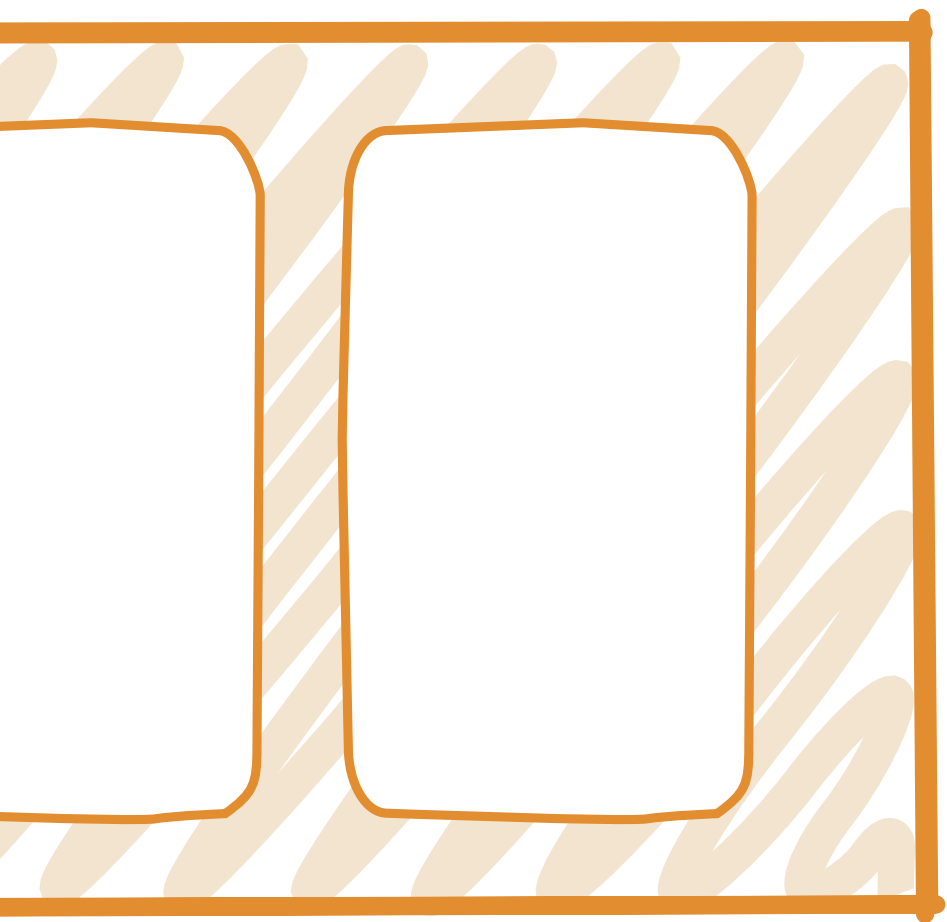
Actually an internal topic!



Producers contain serializers

```
props.put("key.serializer",  
    "org.apache.kafka.serializers.StringSerializer");  
props.put("value.serializer",  
    "io.confluent.kafka.serializers.KafkaAvroSerializer");  
props.put("schema.registry.url",  
    "http://schema-registry:8081");  
...  
producer<String, LogLine> producer = new  
    KafkaProducer<String, LogLine>(props);
```

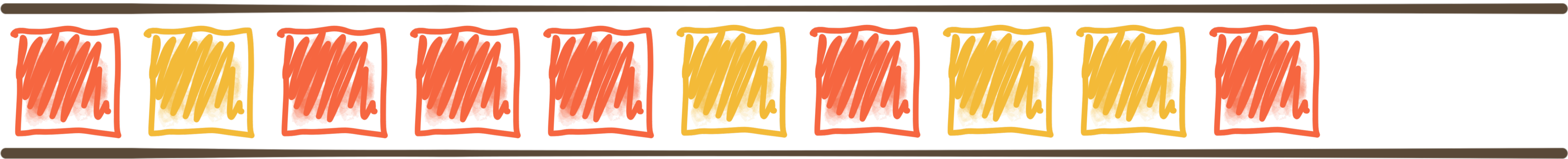




```
{  
  "reading_ts": "2020-02-14T12:19:27Z",  
  "sensor_id": "aa-101",  
  "production_line": "w01",  
  "widget_type": "acme94",  
  "temp_celcius": 23,  
  "widget_weight_g": 100  
}
```



Streams of events



Time →

Stream Processing



Stream: widgets



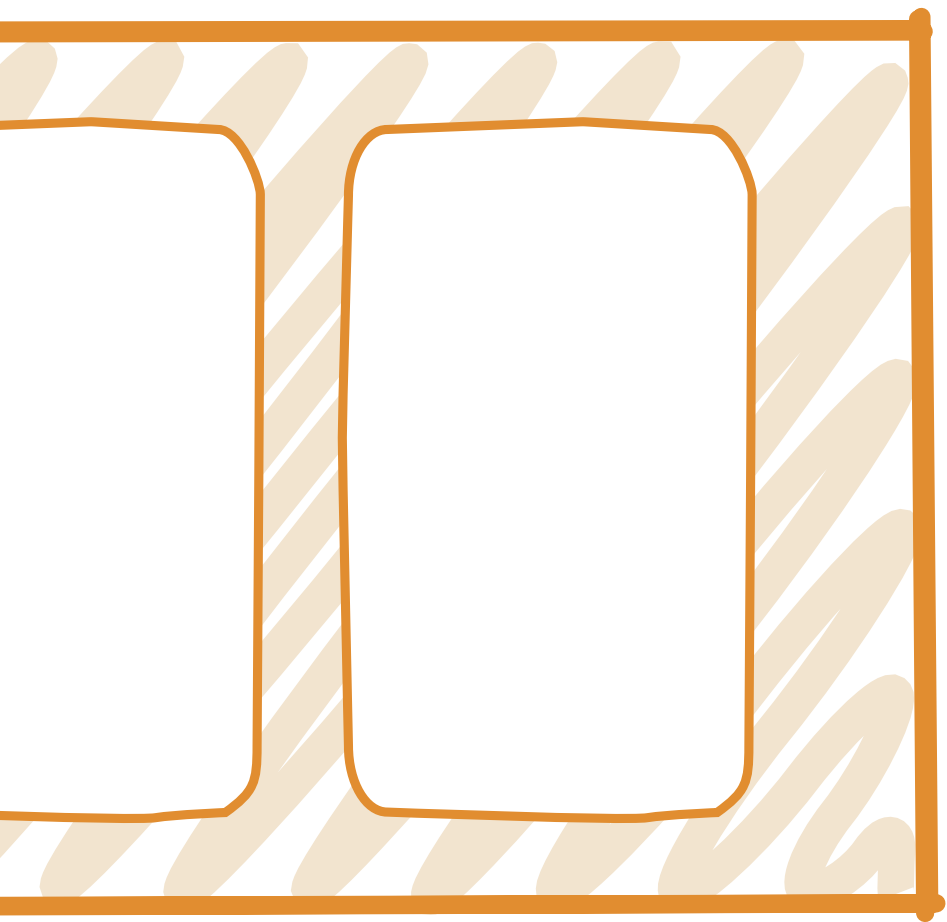
Stream: widgets_red

Stream Processing with Kafka Streams



```
final StreamsBuilder builder = new StreamsBuilder()
    .stream("widgets", Consumed.with(stringSerde, widgetsSerde))
    .filter( (key, widget) -> widget.getColour().equals("RED") )
    .to("widgets_red", Produced.with(stringSerde, widgetsSerde));
```





STREAMS
APPLICATION

STREAMS
APPLICATION

STREAMS
APPLICATION

Stream Processing with ksqlDB

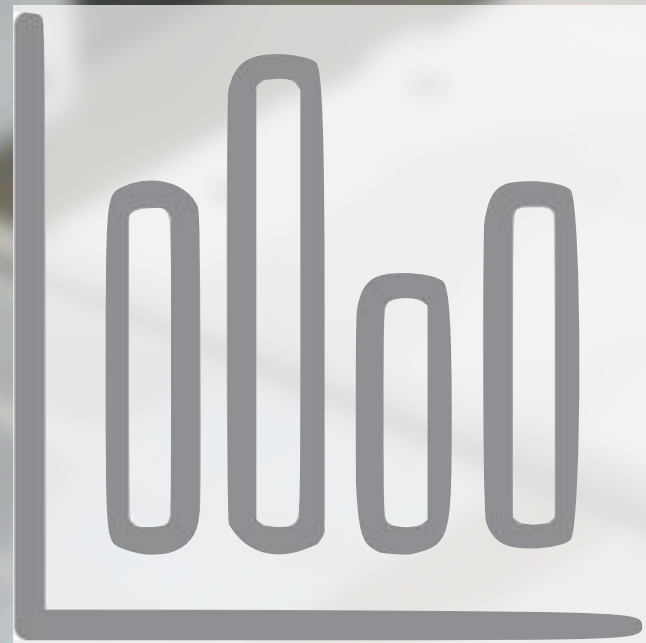


ksqlDB

```
CREATE STREAM widgets_red AS  
SELECT * FROM widgets  
WHERE colour='RED';
```



```
{  
  "reading_ts": "2020-02-14T12:19:27Z",  
  "sensor_id": "aa-101",  
  "production_line": "w01",  
  "widget_type": "acme94",  
  "temp_celcius": 23,  
  "widget_weight_g": 100  
}
```





```
SELECT *  
FROM WIDGETS  
WHERE WEIGHT_G > 120
```



```
SELECT COUNT(*)  
FROM WIDGETS  
GROUP BY PRODUCTION_LINE
```



```
SELECT AVG(TEMP_CELCIUS) AS TEMP  
FROM WIDGETS  
GROUP BY SENSOR_ID  
HAVING TEMP > 20
```

```
{  
  "reading_ts": "2020-02-14T12:19:27Z",  
  "sensor_id": "aa-101",  
  "production_line": "w01",  
  "widget_type": "acme94",  
  "temp_celcius": 23,  
  "widget_weight_g": 100  
}
```



*Object store,
data warehouse,
RDBMS*

```
CREATE SINK CONNECTOR dw WITH (  
  'connector.class' = 'S3Connector',  
  'topics' = 'widgets'  
...);
```

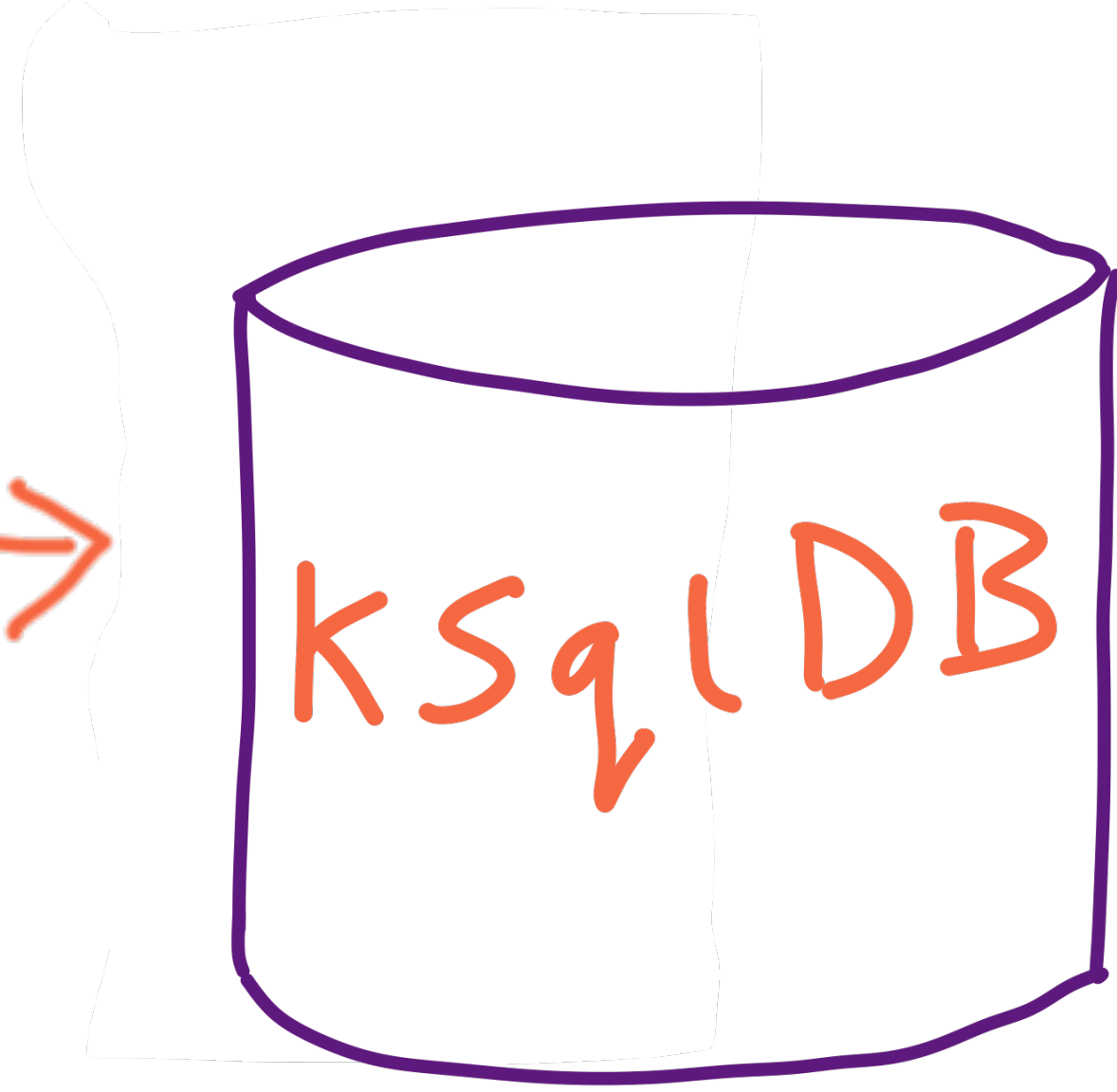
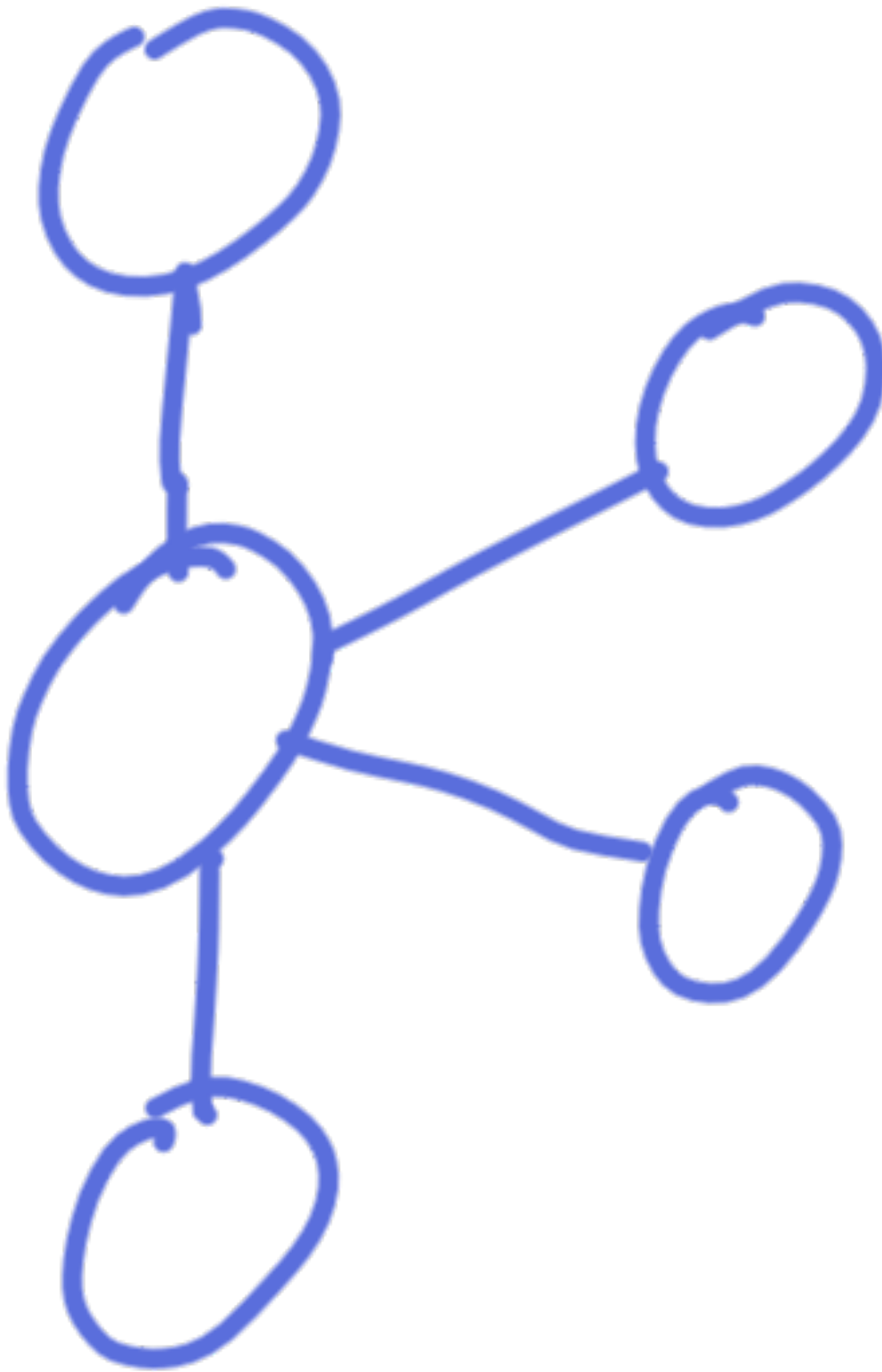
@rmoff

#BudapestData

@confluentinc

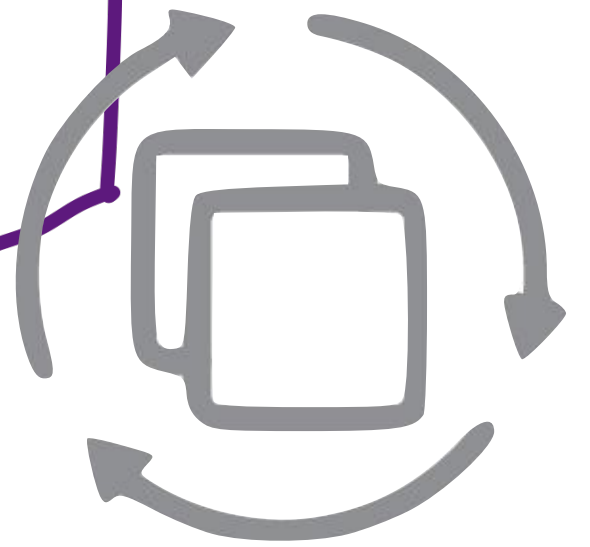
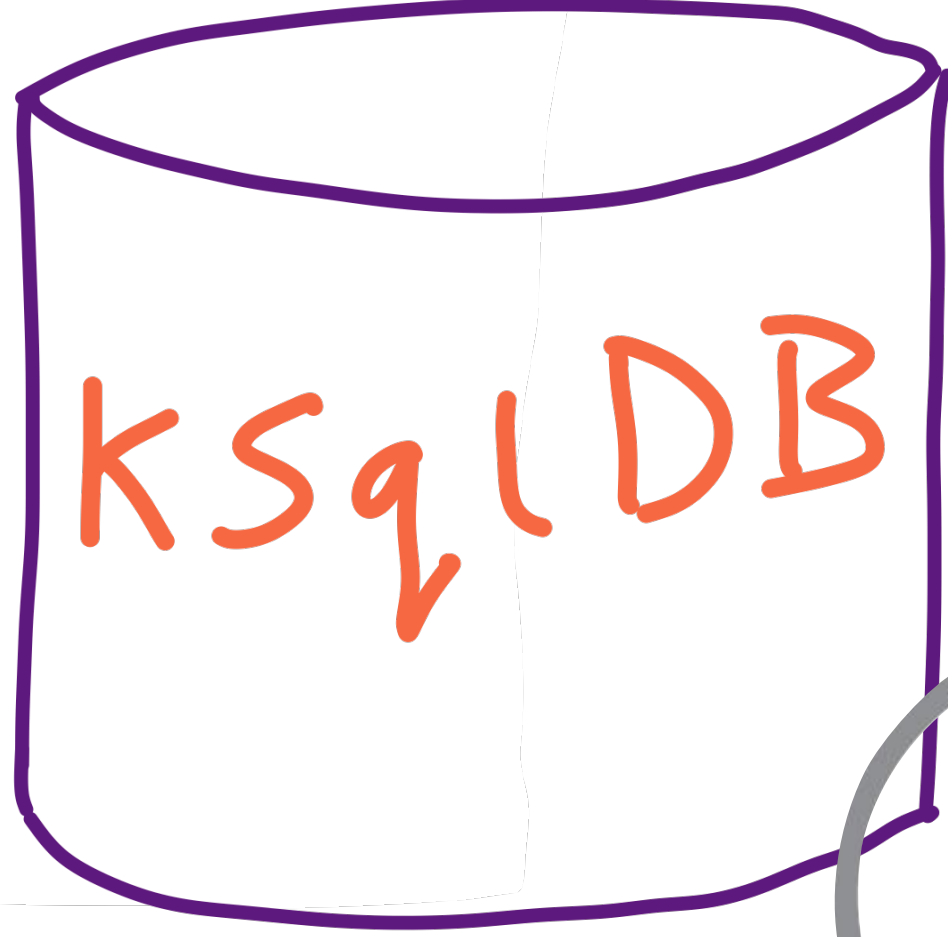
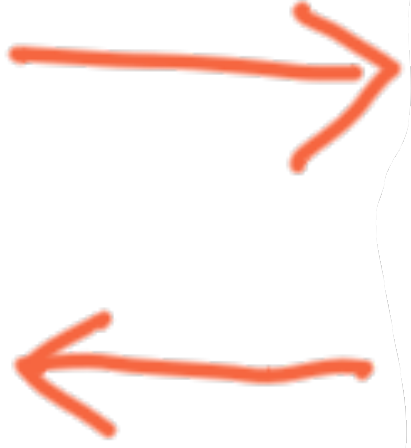
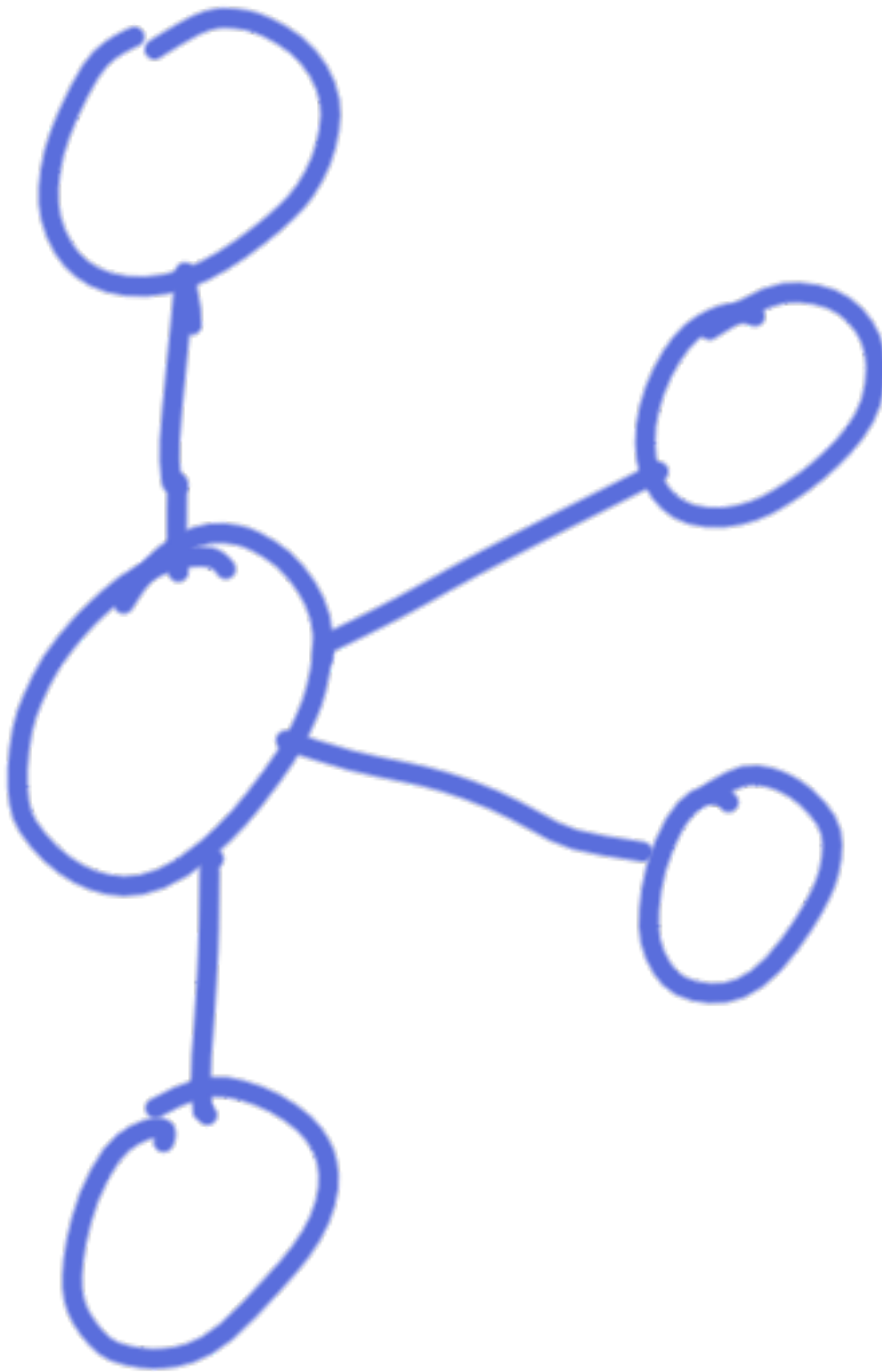
Stream Processing with ksqlDB

Source stream



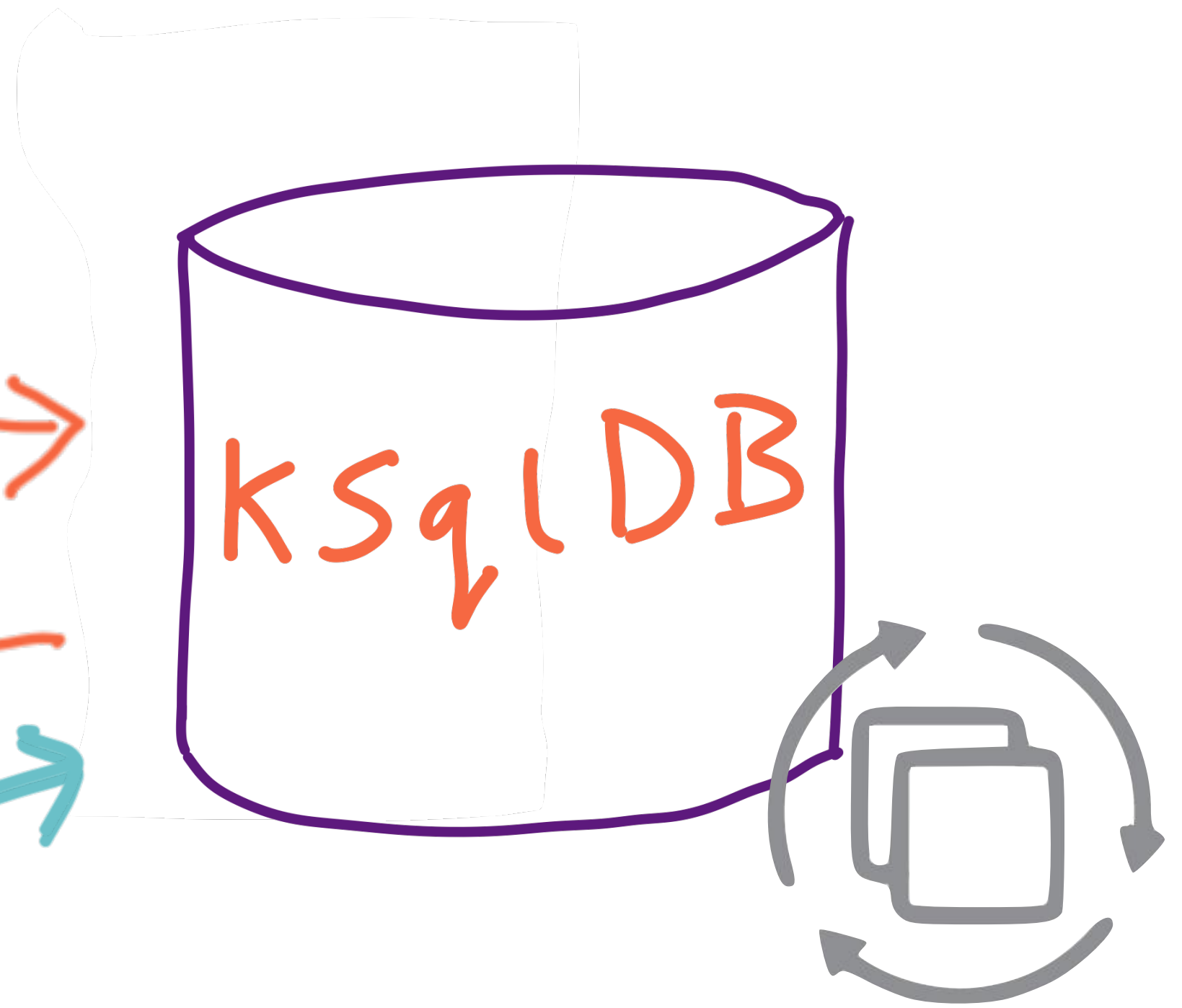
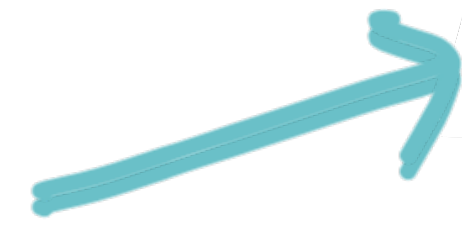
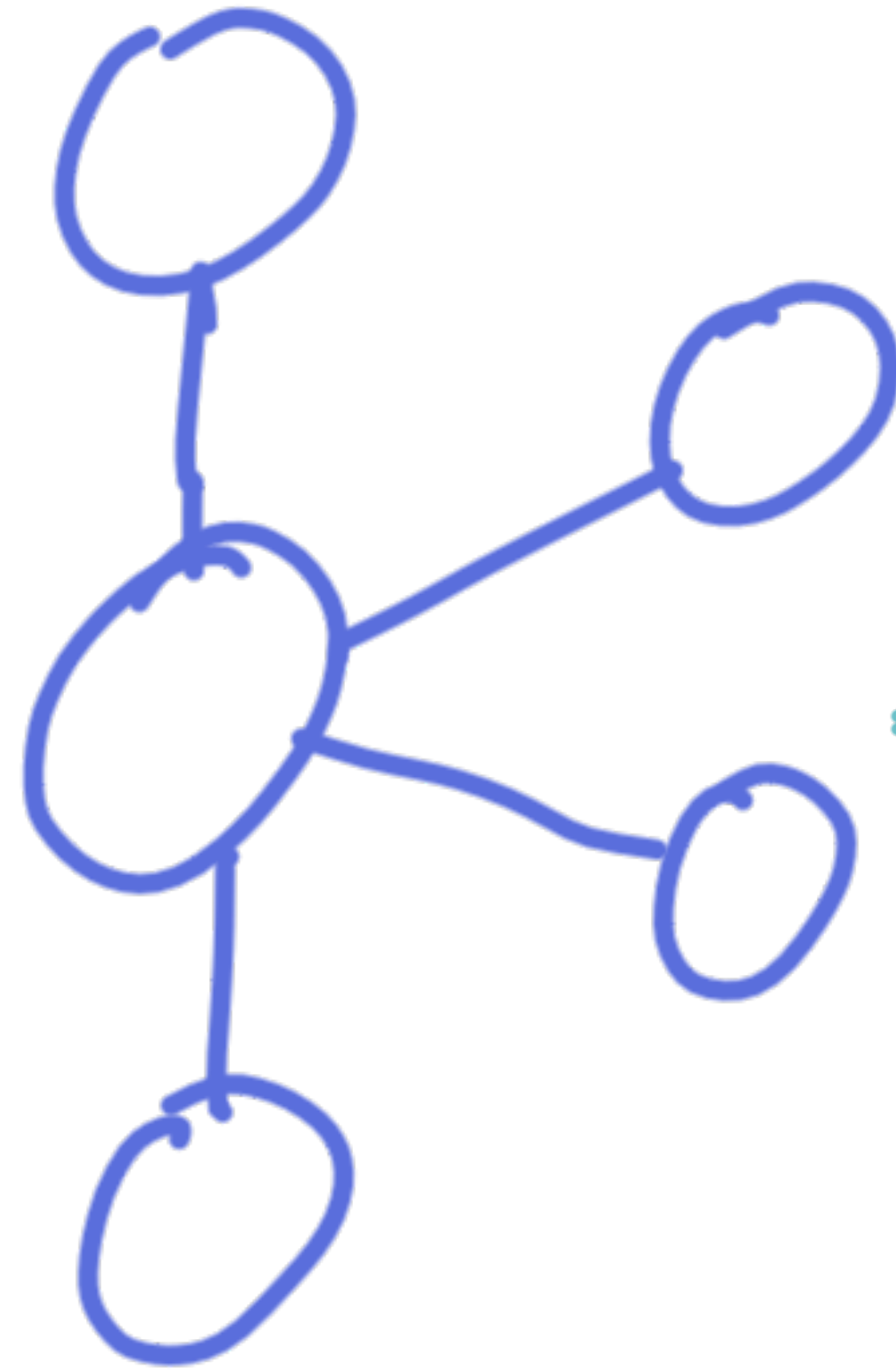
Stream Processing with ksqlDB

Source stream



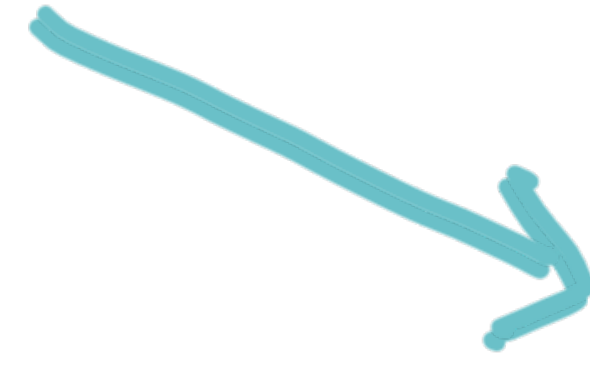
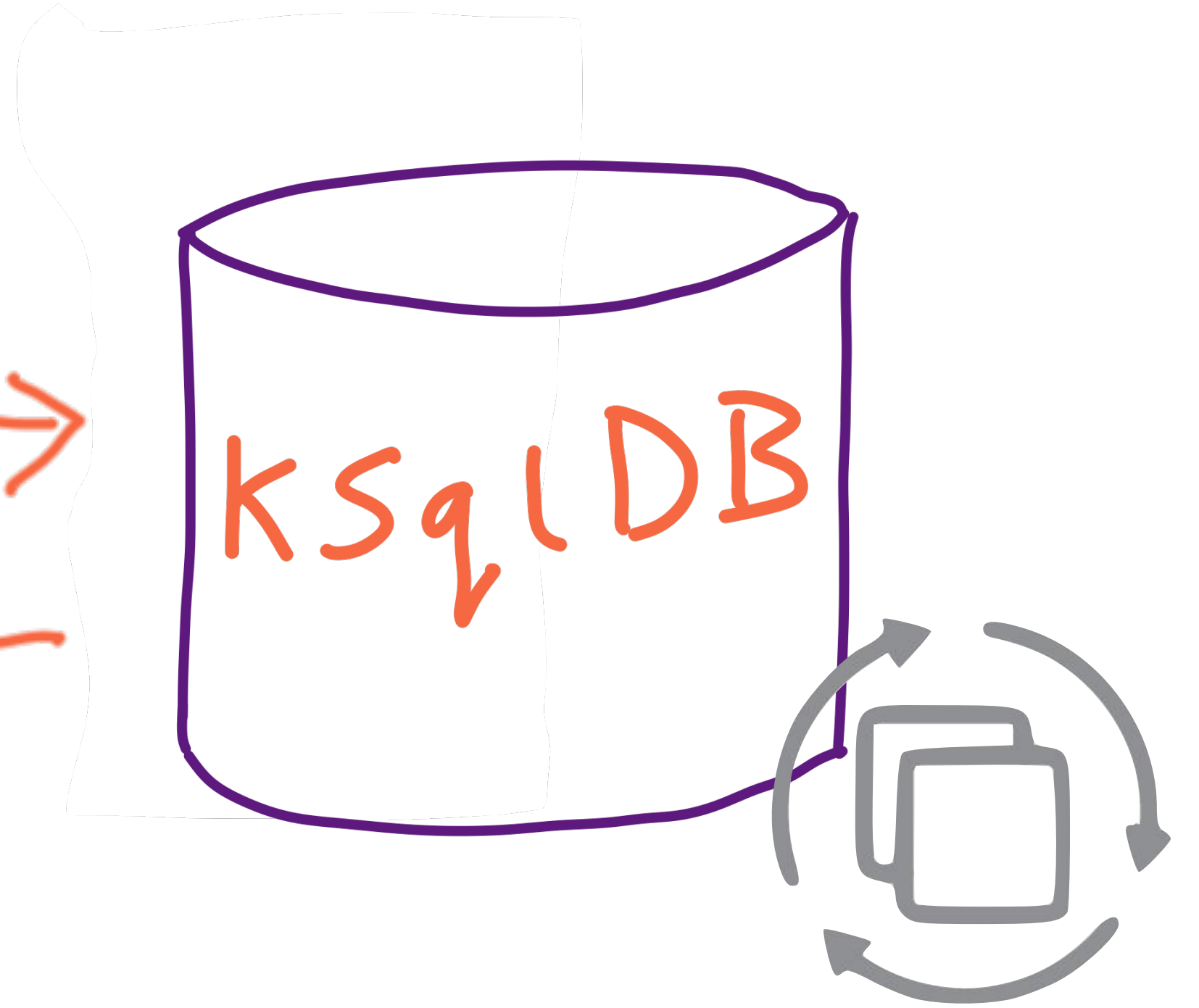
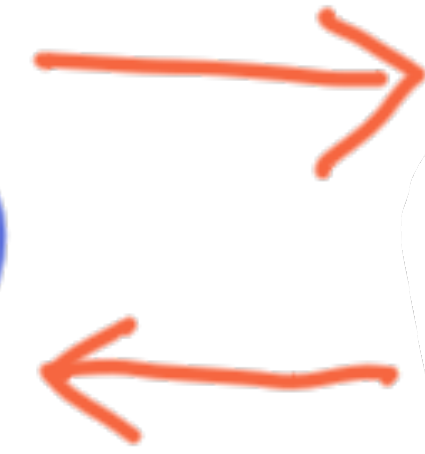
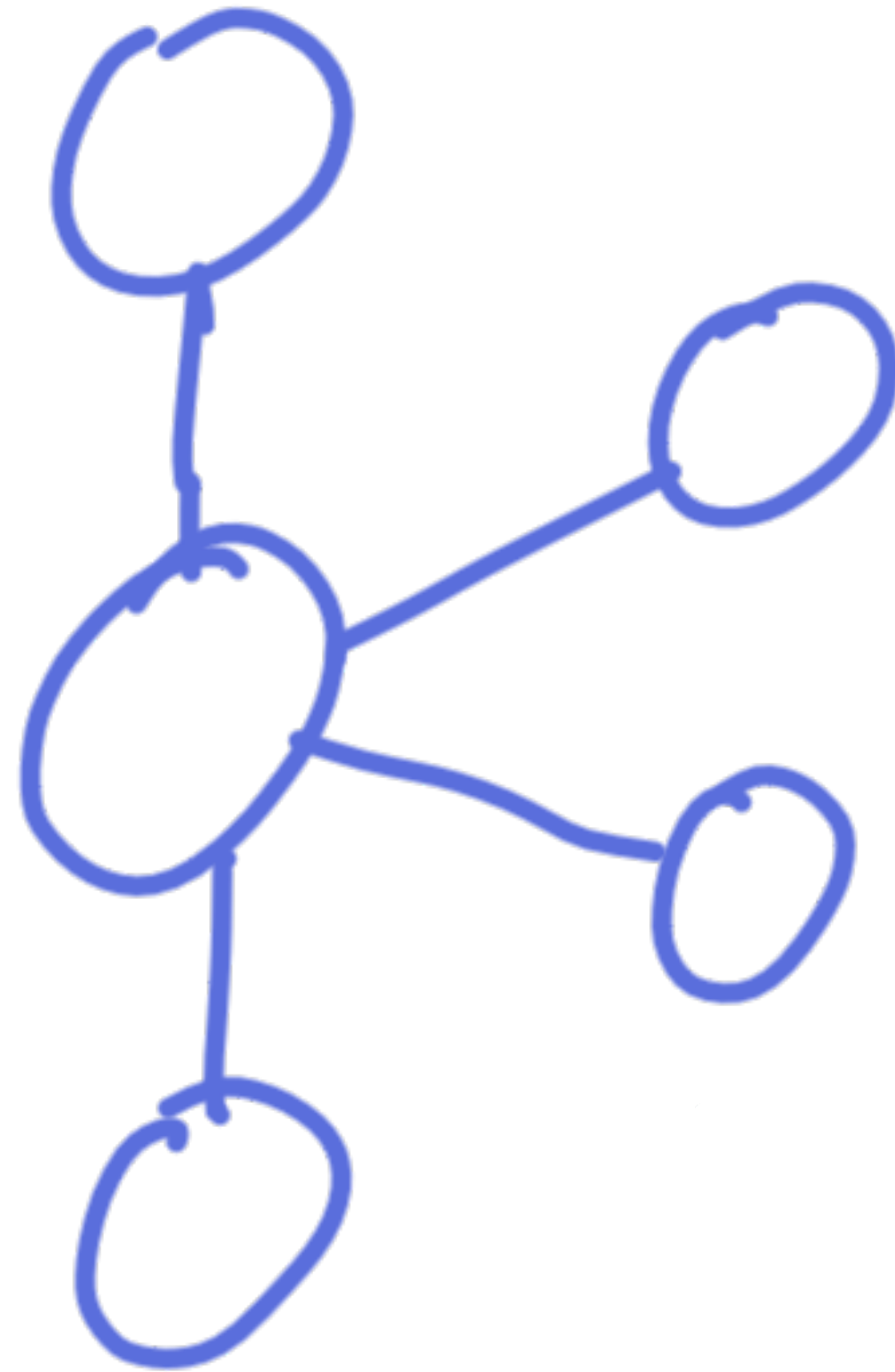
Stream Processing with ksqlDB

Source stream



Stream Processing with ksqlDB

Source stream

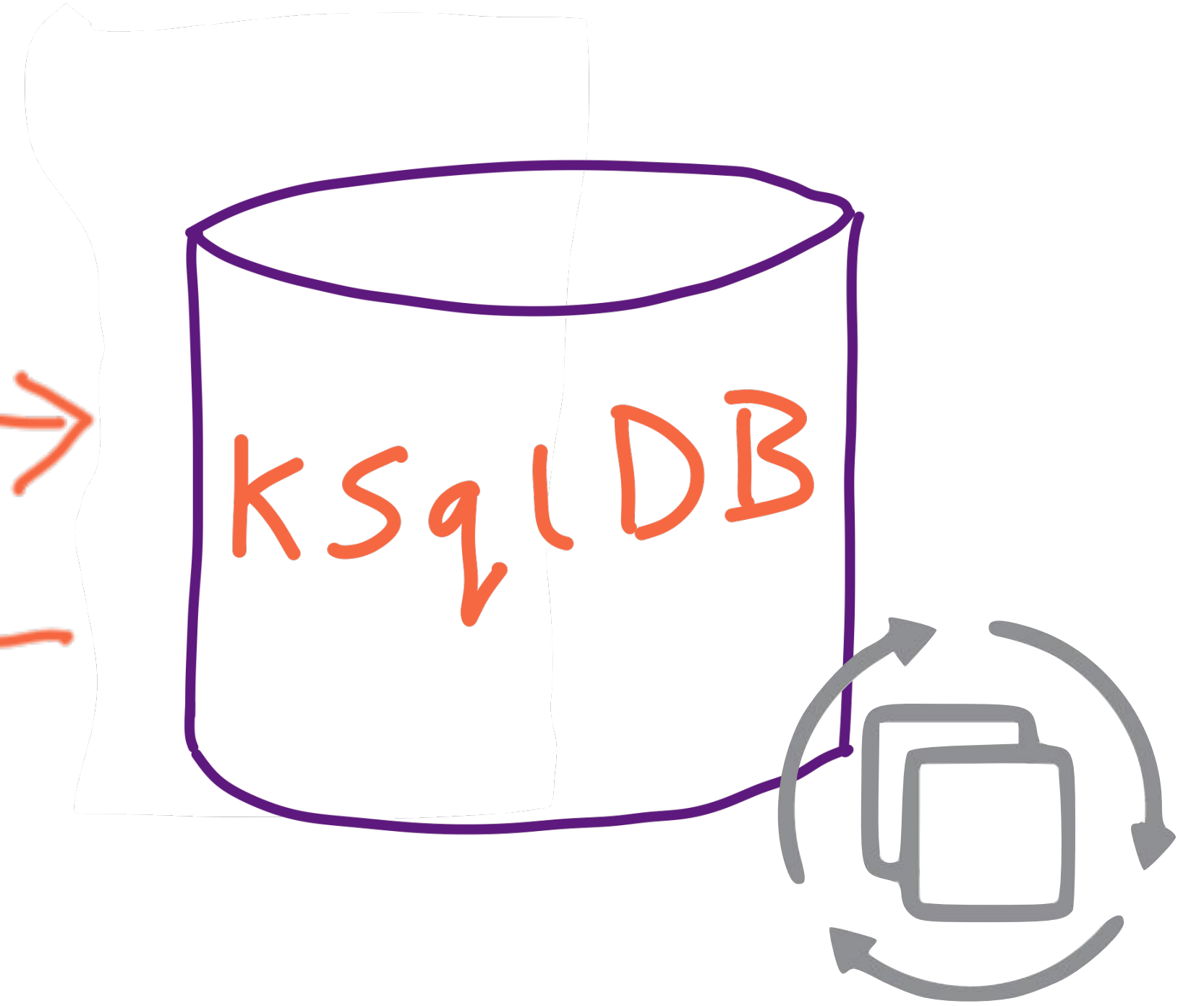
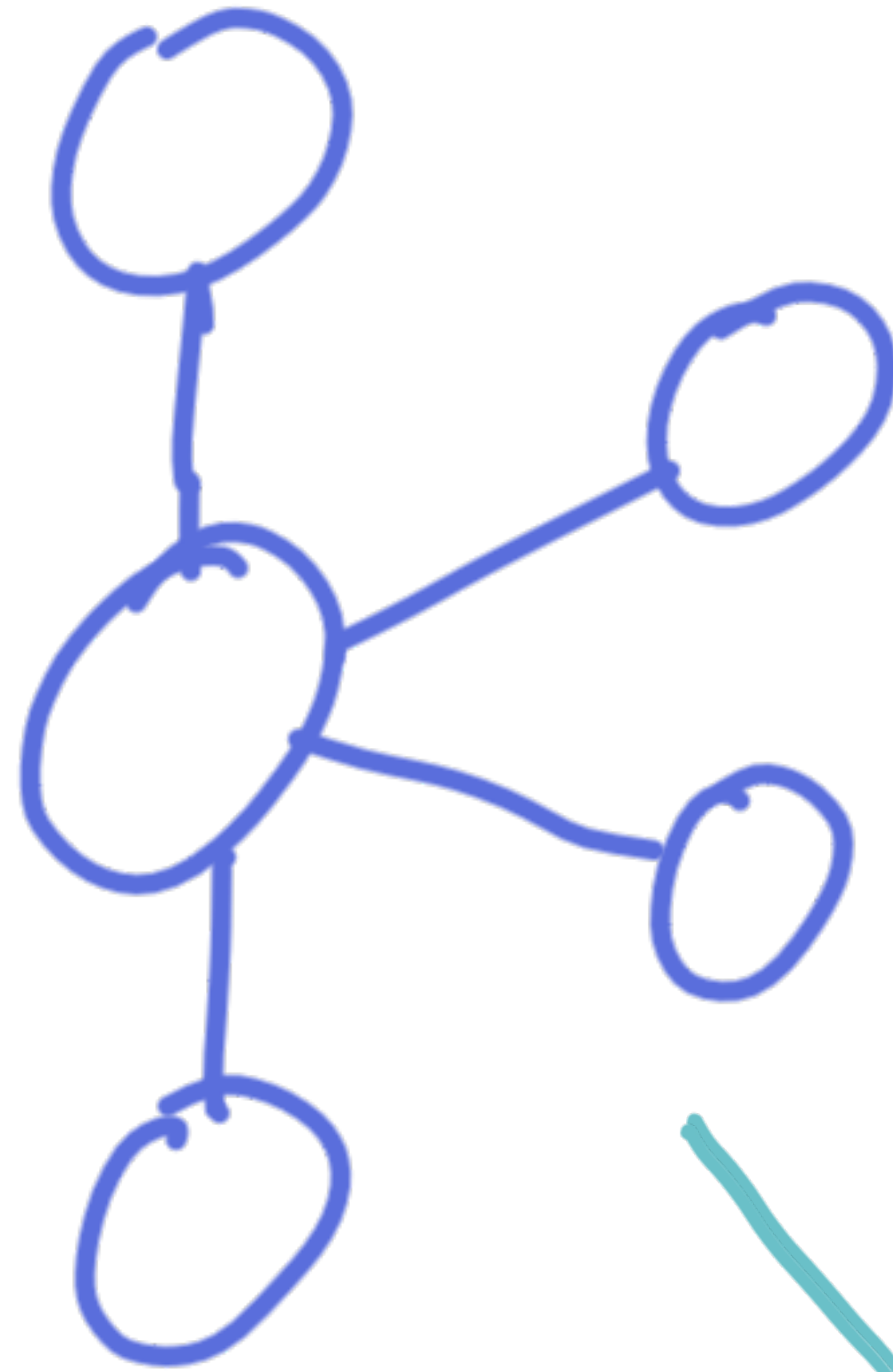


Analytics

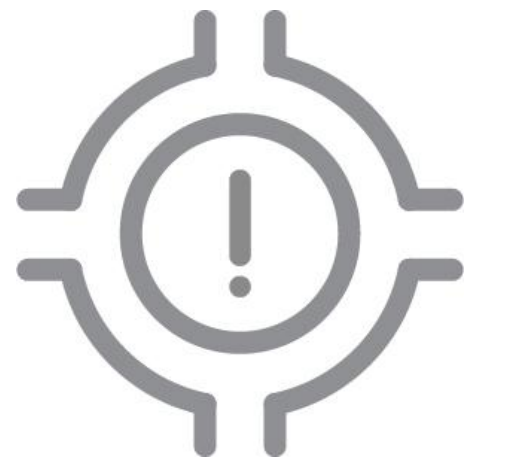


Stream Processing with ksqlDB

Source stream



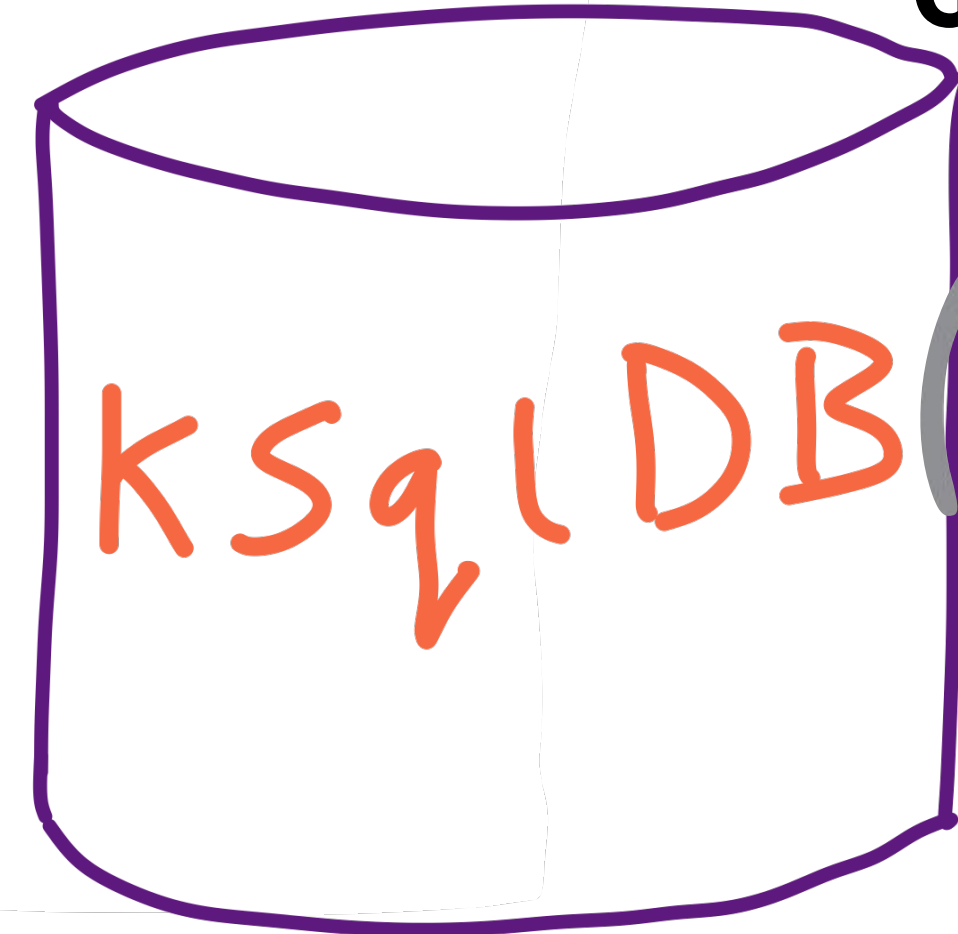
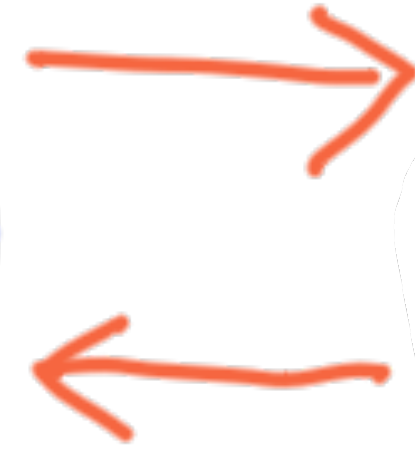
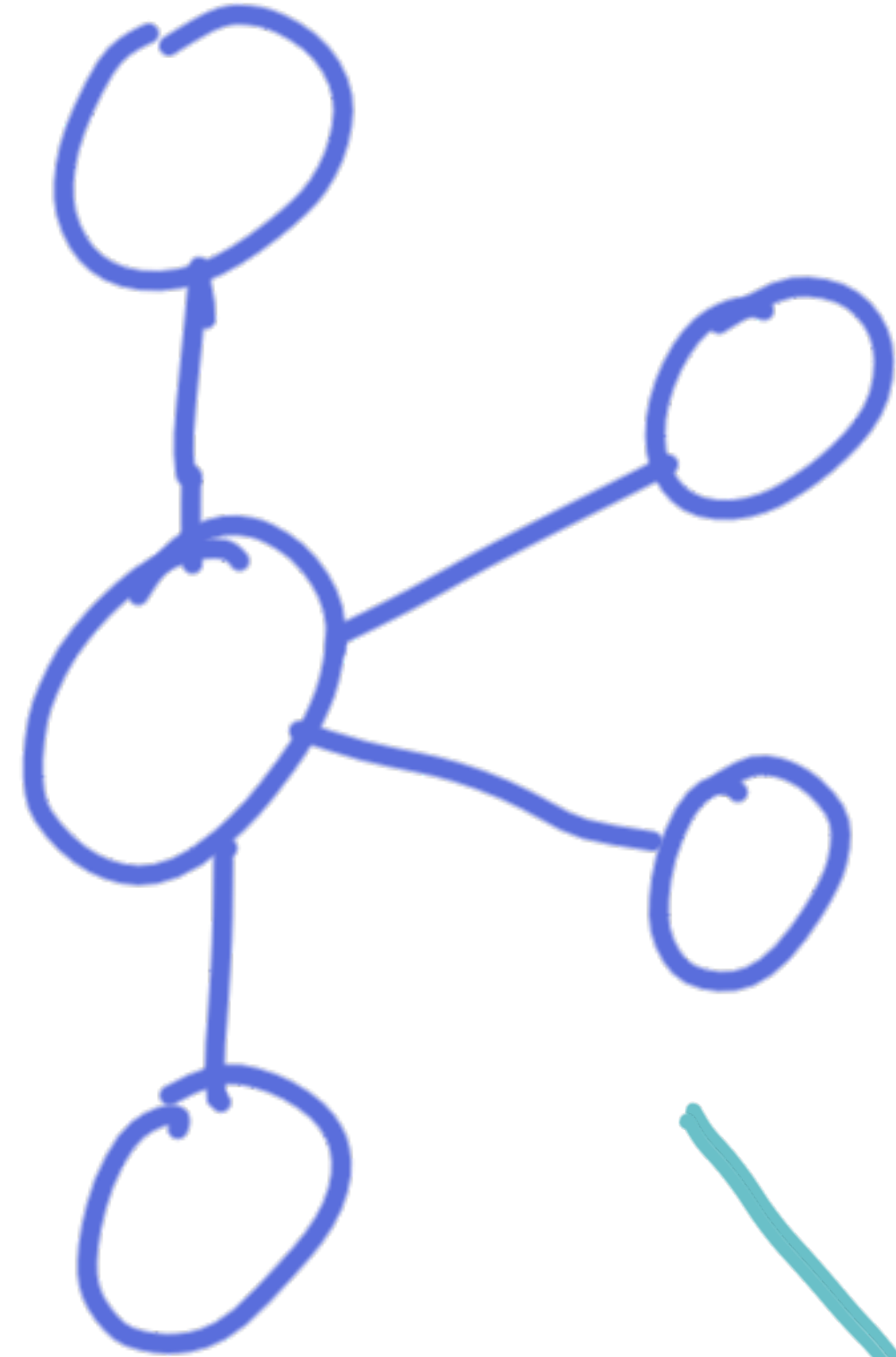
Applications /
Microservices



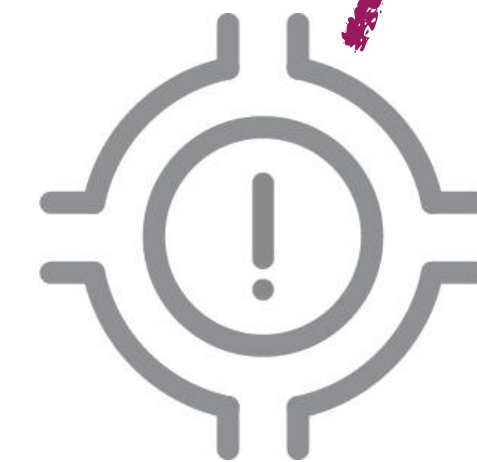
Stream Processing with ksqlDB

...SUM(TXN_AMT)
GROUP BY AC_ID

Source stream →



Applications /
Microservices



AC_ID=42

BALANCE=94.00
AC_ID=42

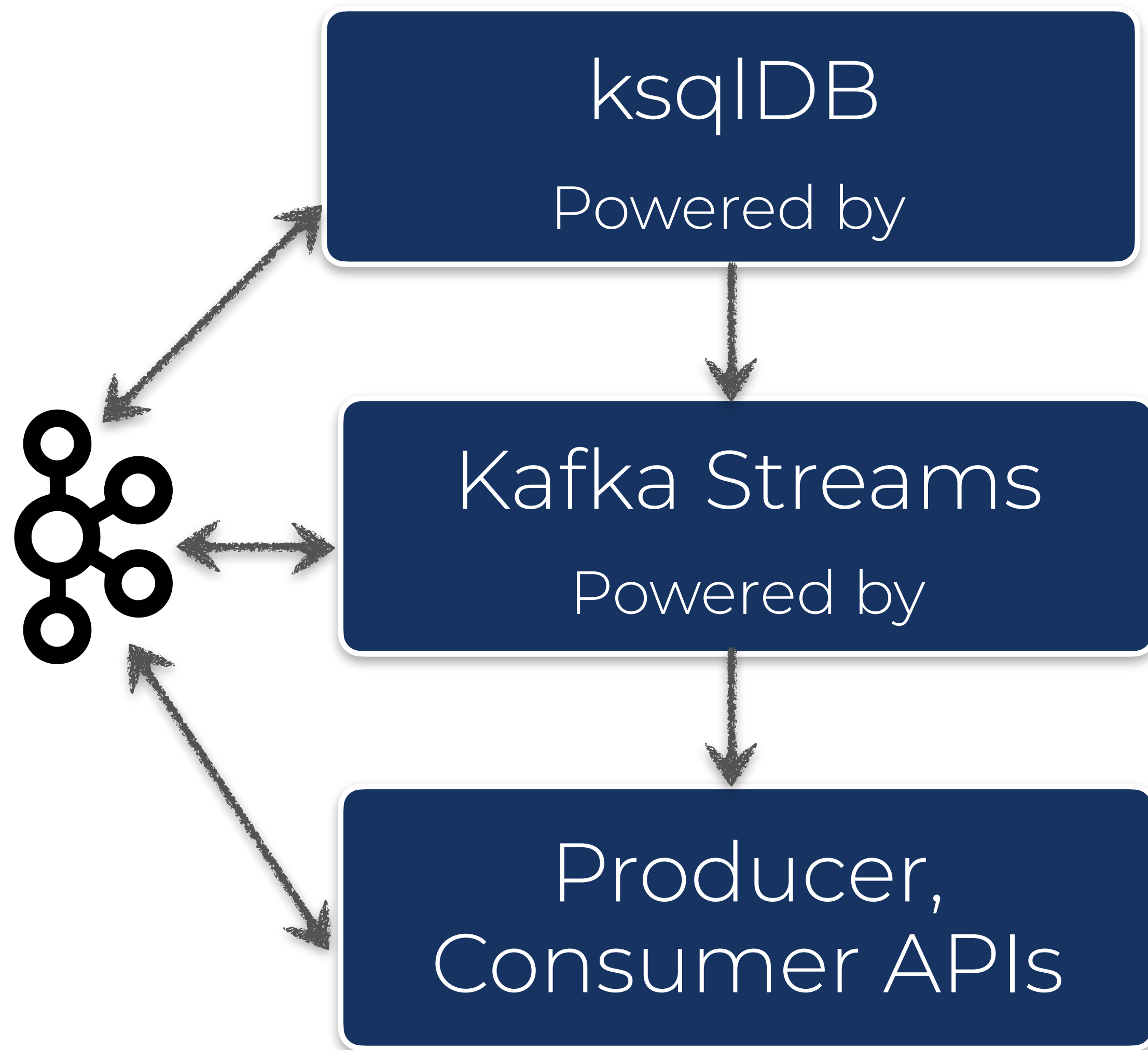
ksqlDB

or

Kafka Streams?



Standing on the Shoulders of Streaming Giants



Ease of use

Flexibility

```
CREATE STREAM,  
CREATE TABLE, SELECT, JOIN,  
GROUP BY, SUM, ...
```

ksqlDB UDFs

```
KStream<>, KTable<>, filter(), map(),  
flatMap(), join(), aggregate(),  
transform(), ...
```

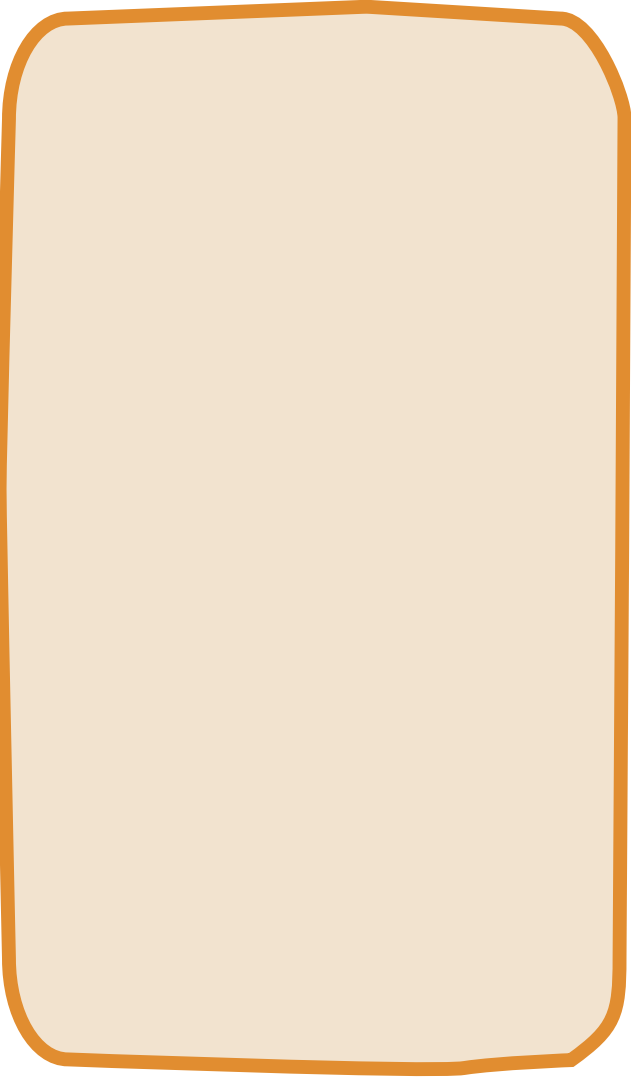
```
subscribe(), poll(), send(),  
flush(), beginTransaction(), ...
```



Photo by Raoul Droog on Unsplash

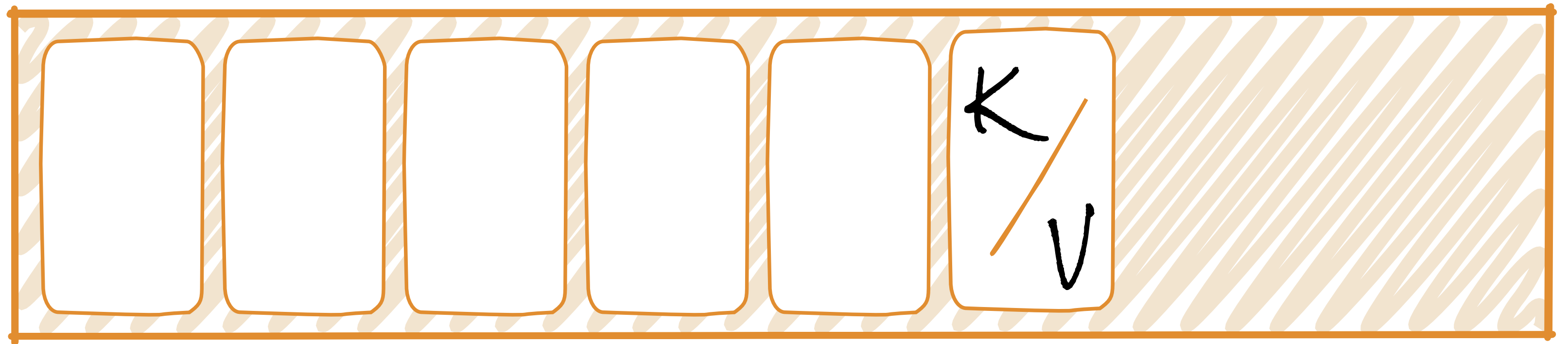
DEMO

Summary

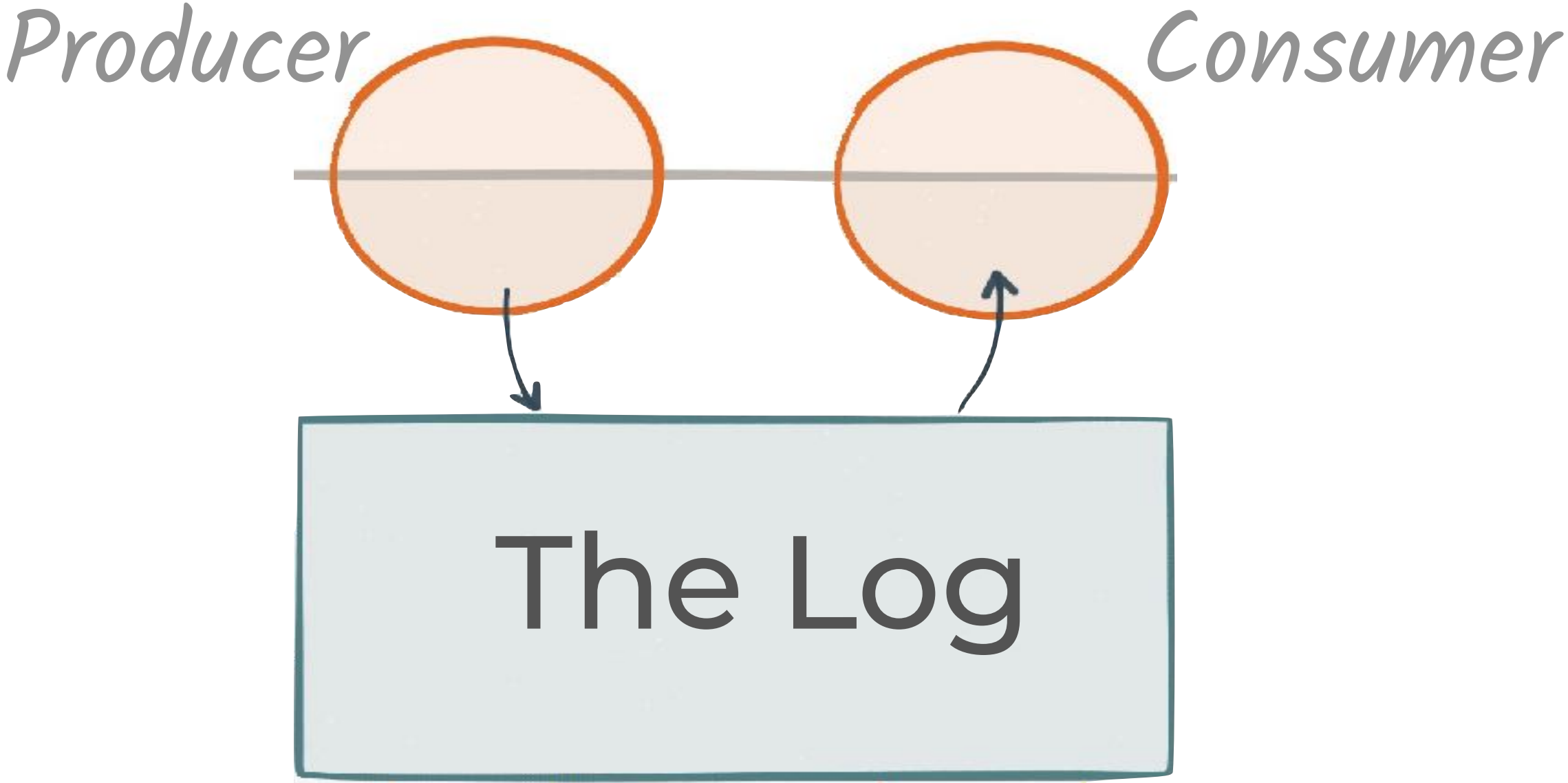


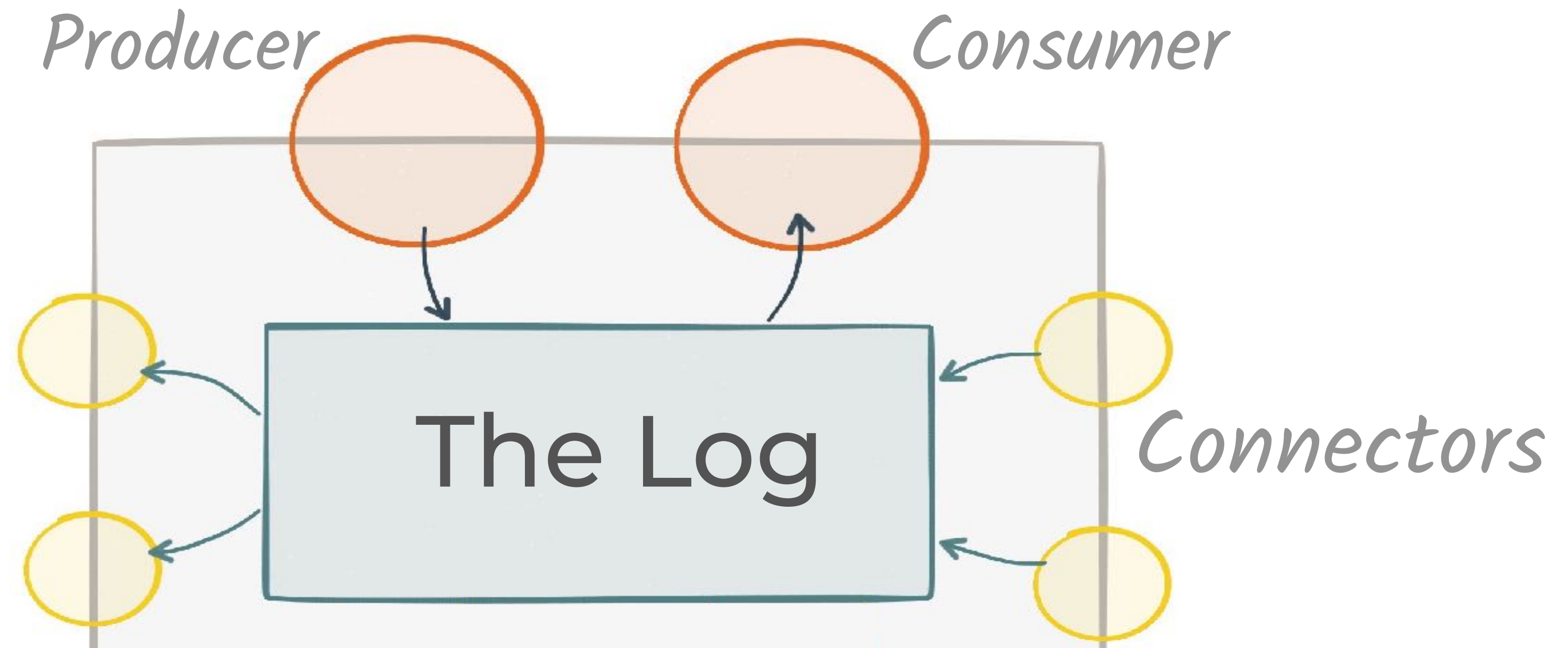
K
/

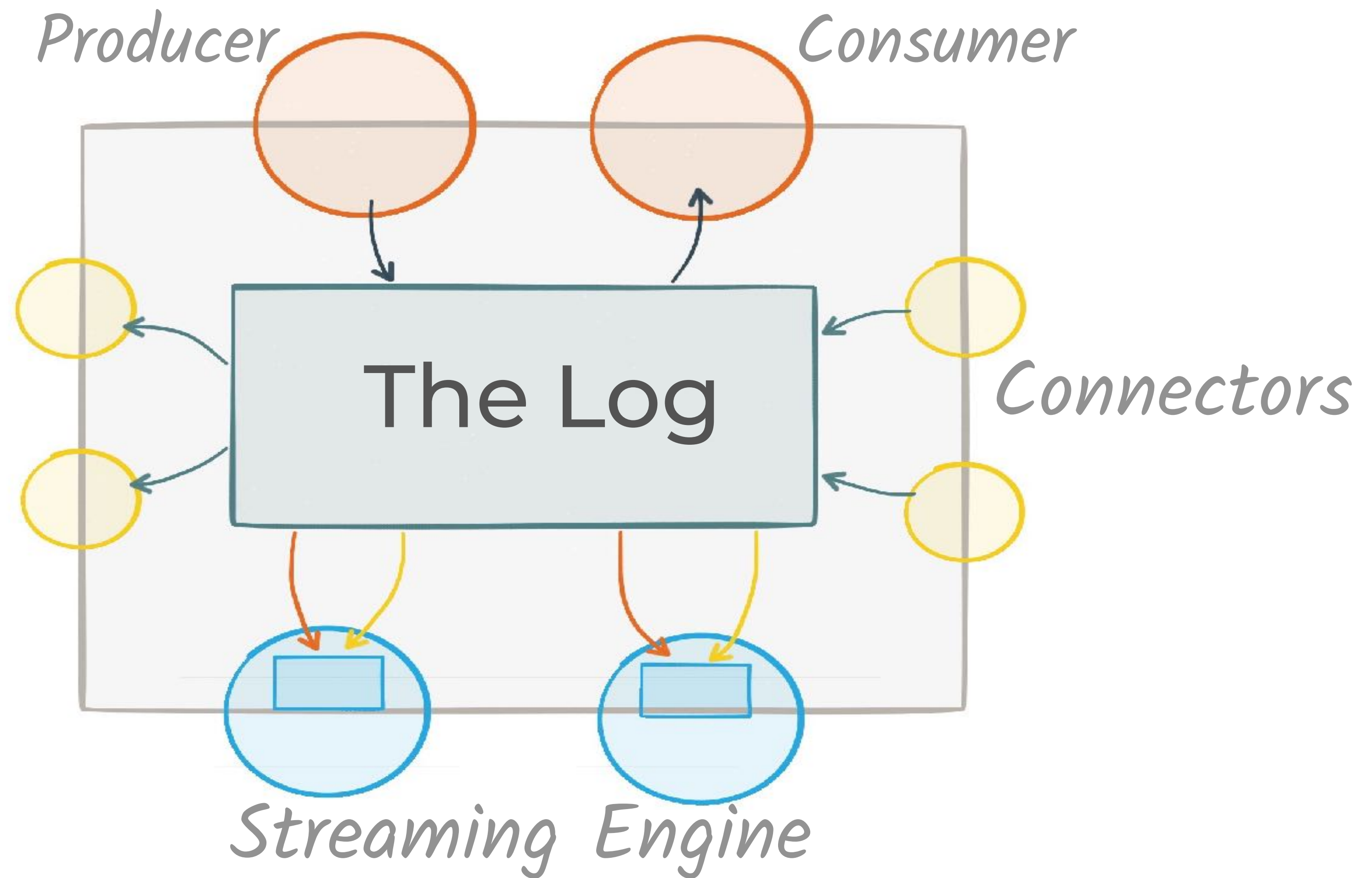
V



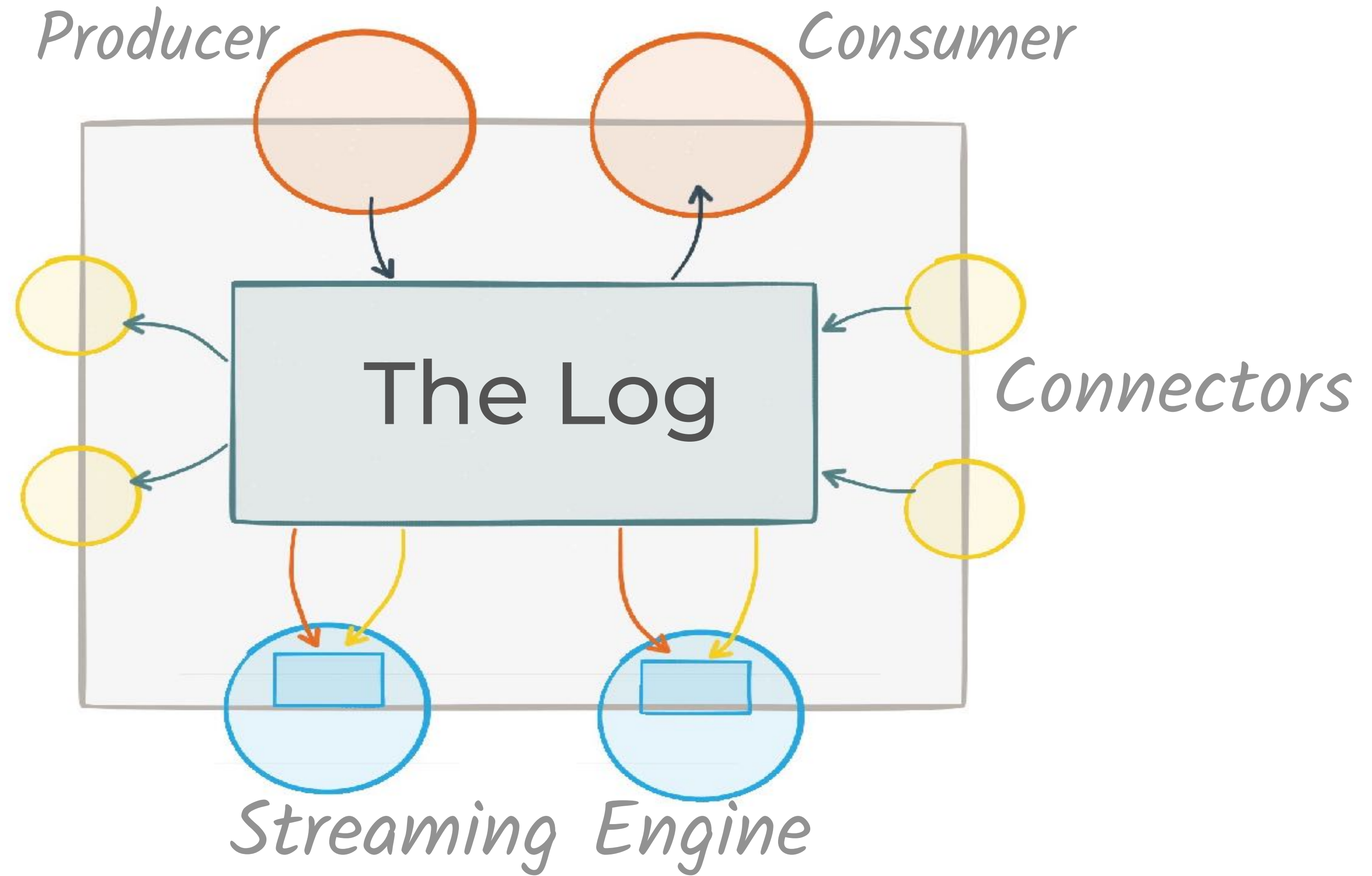
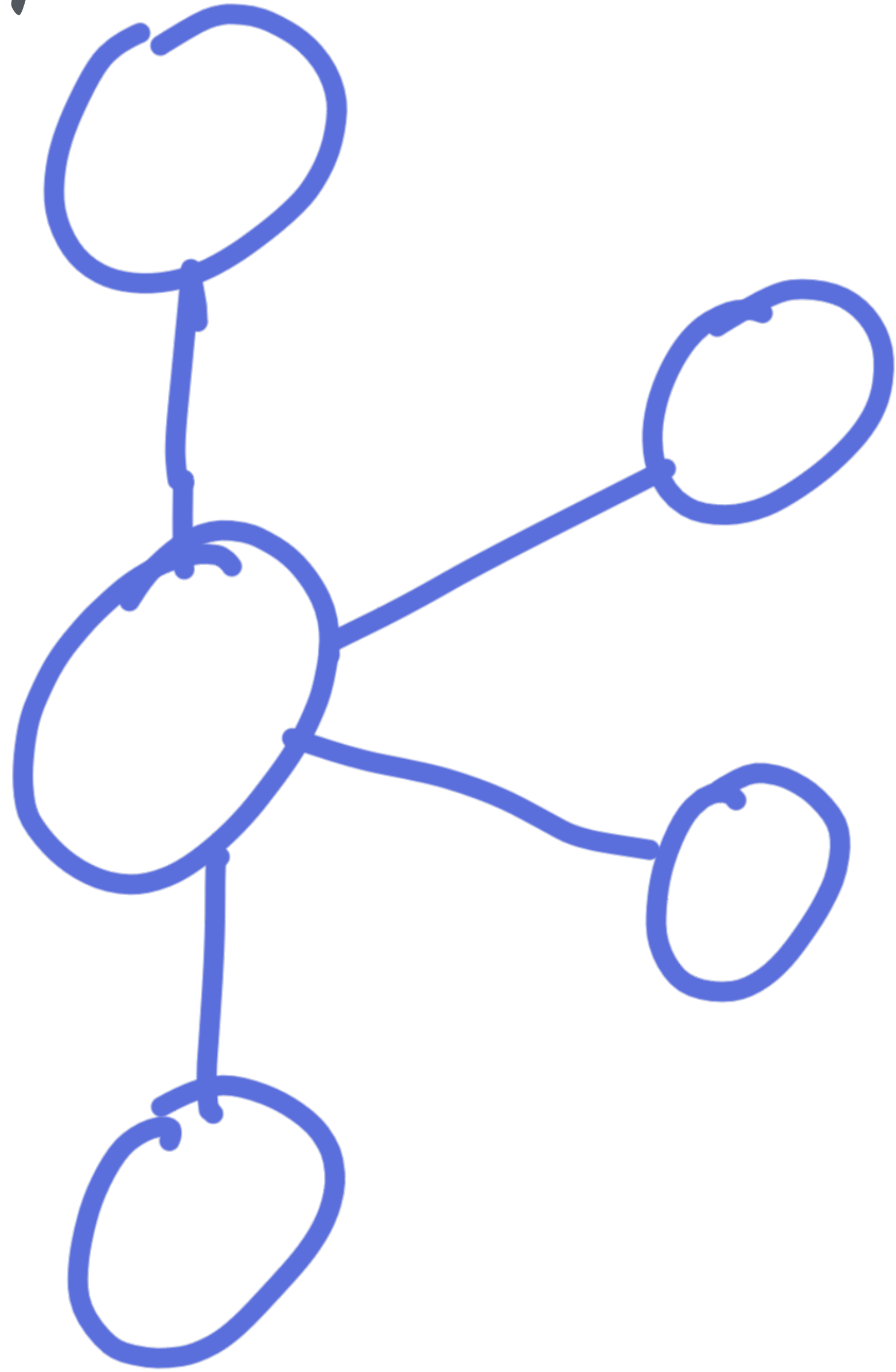
The Log







Apache Kafka



EVENTS

are

EVERYWHERE

EVENTS

are

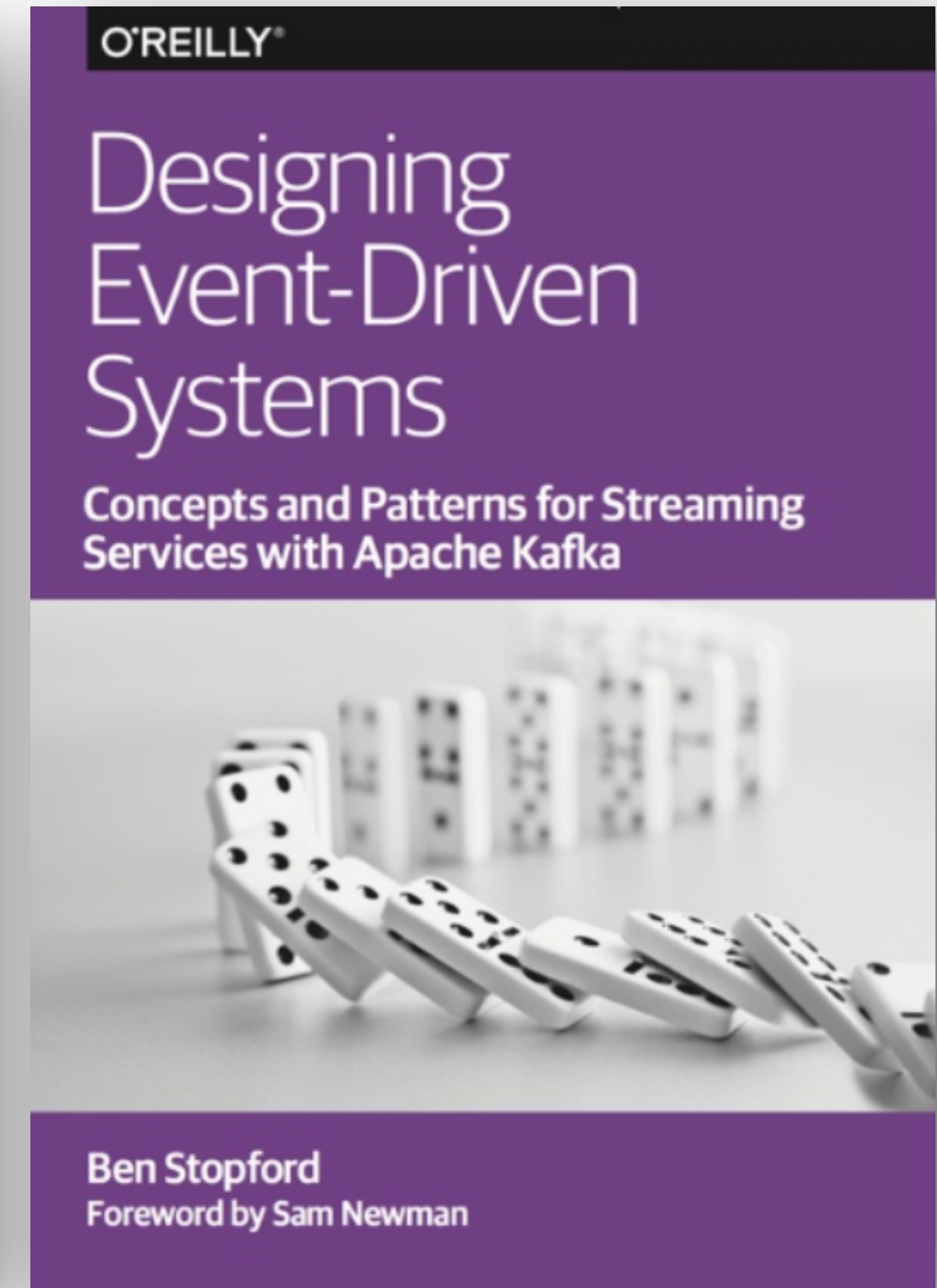
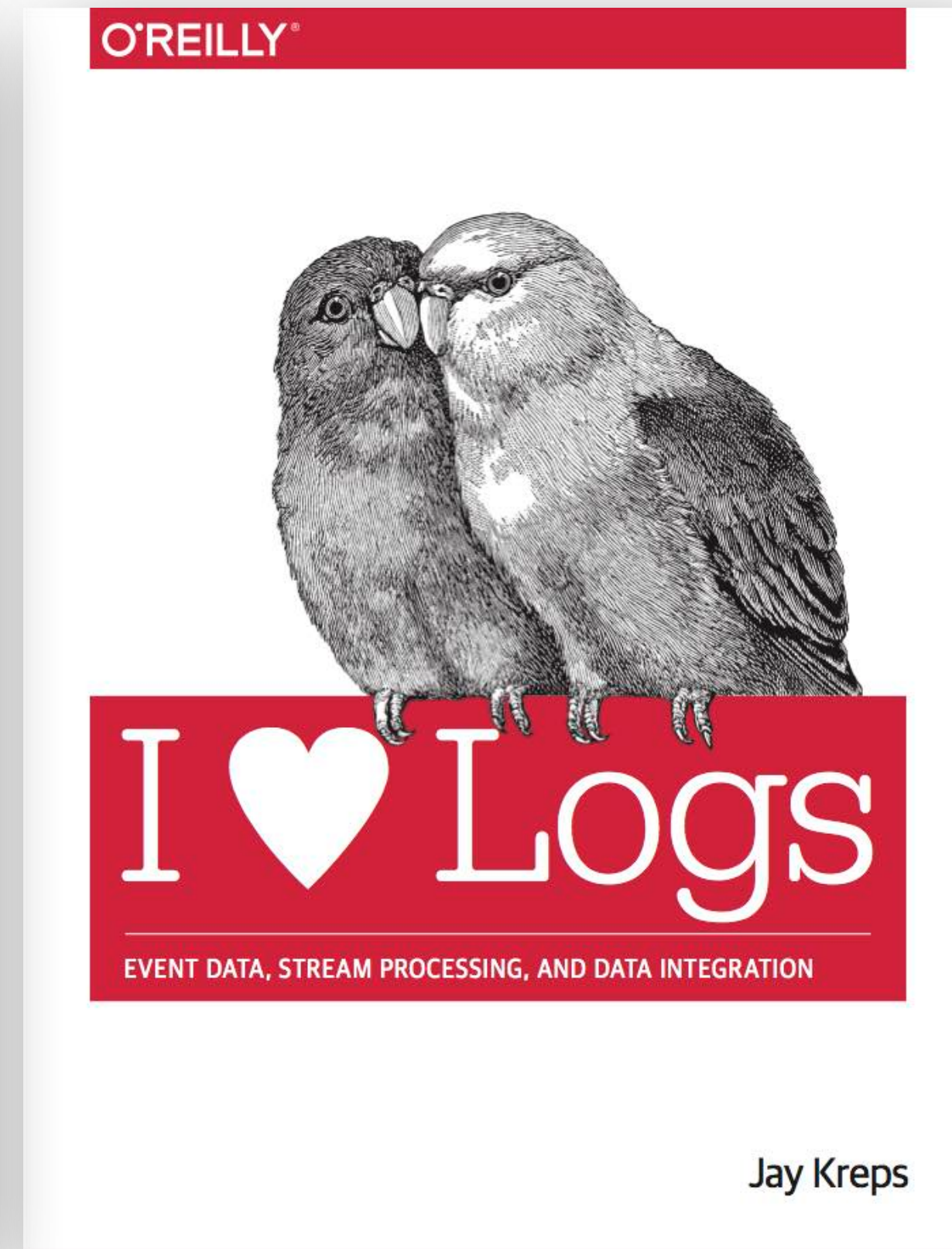
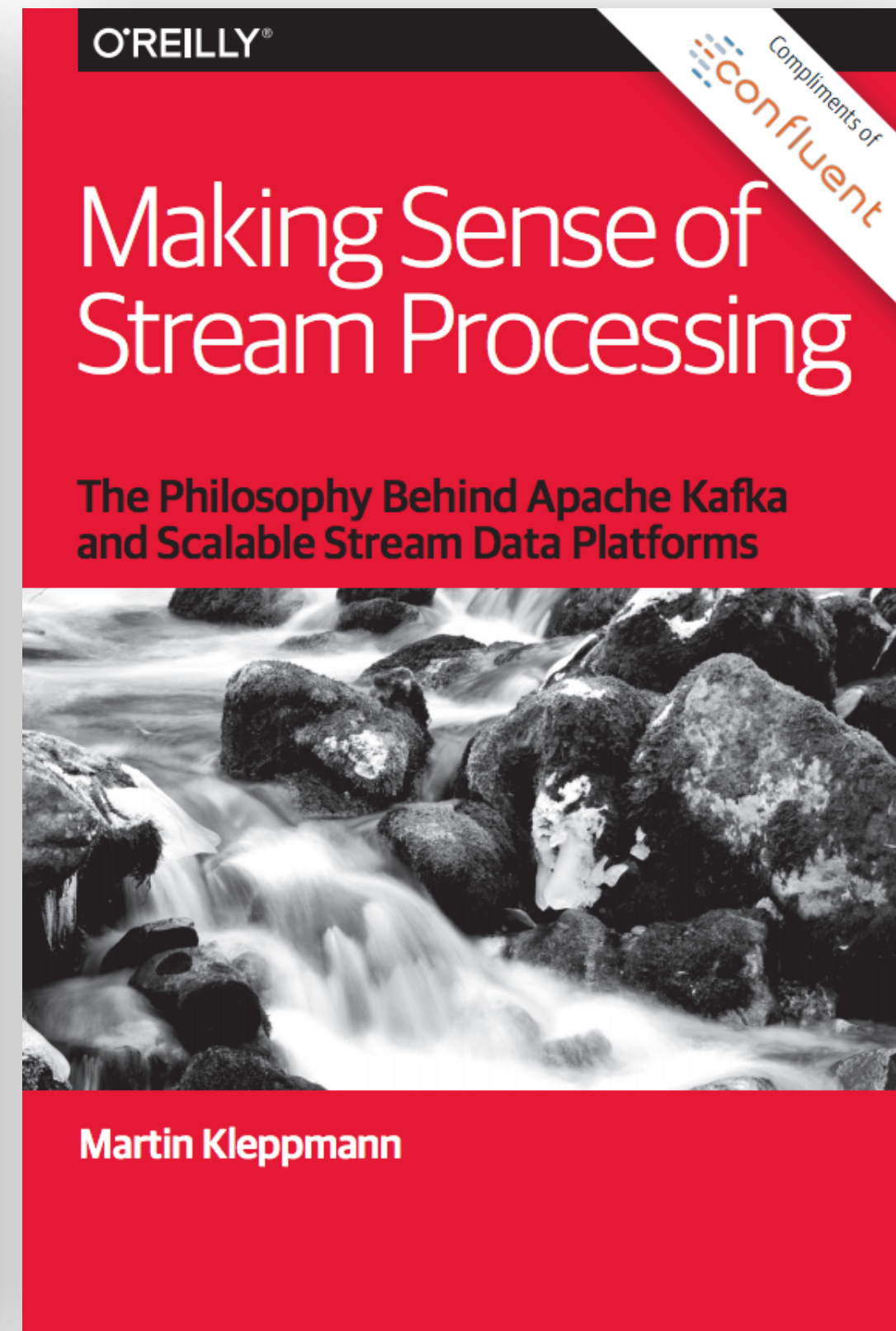
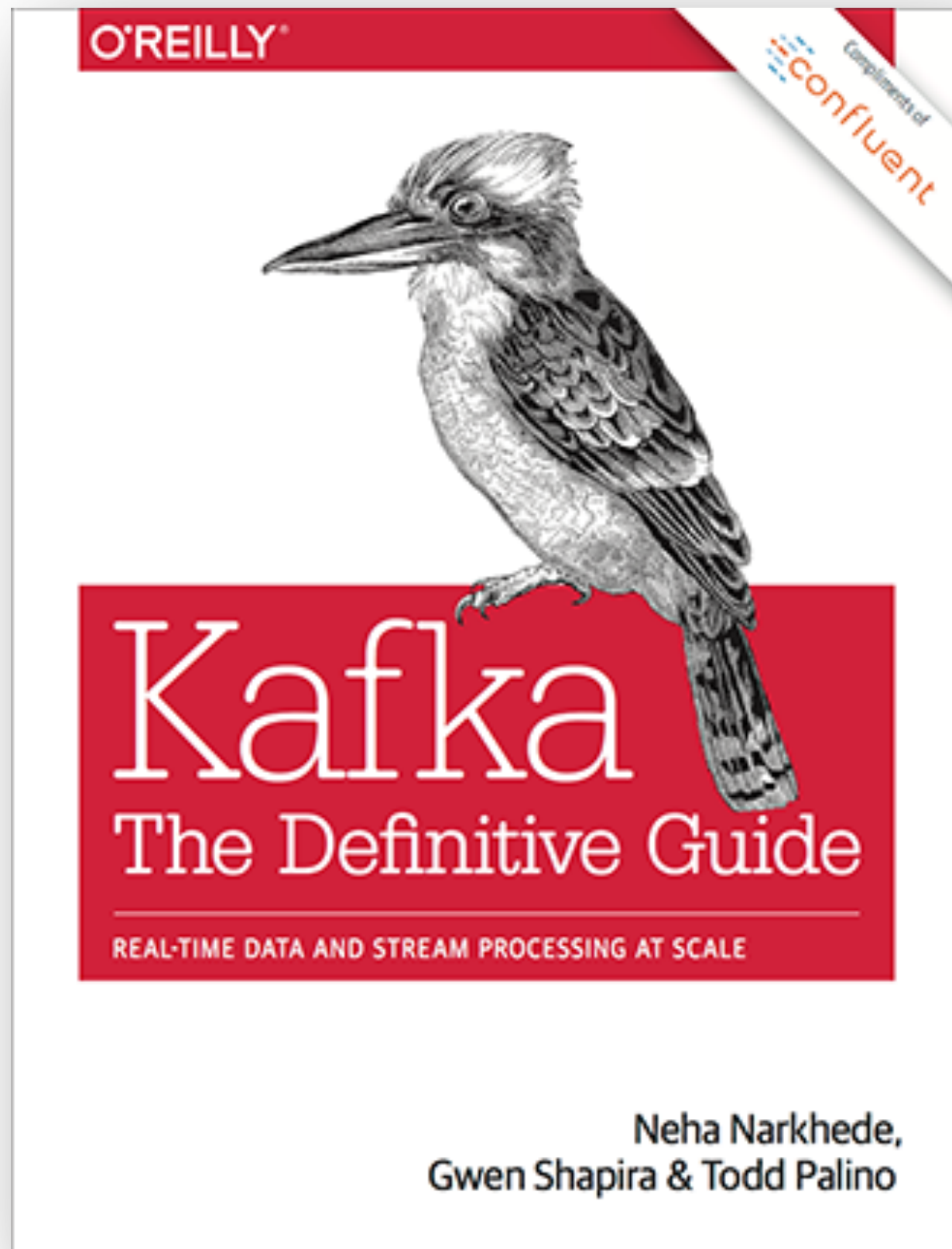
very
^

POWERFUL

Free Books!



<https://rmoff.dev/budapestdata>



60DEVADV

\$50 USD off your bill each calendar month
for the first three months when you sign up
<https://rmoff.dev/690>

Free money!
(additional \$60 towards
your bill 😊)



confluent cloud

Fully Managed Kafka as a Service

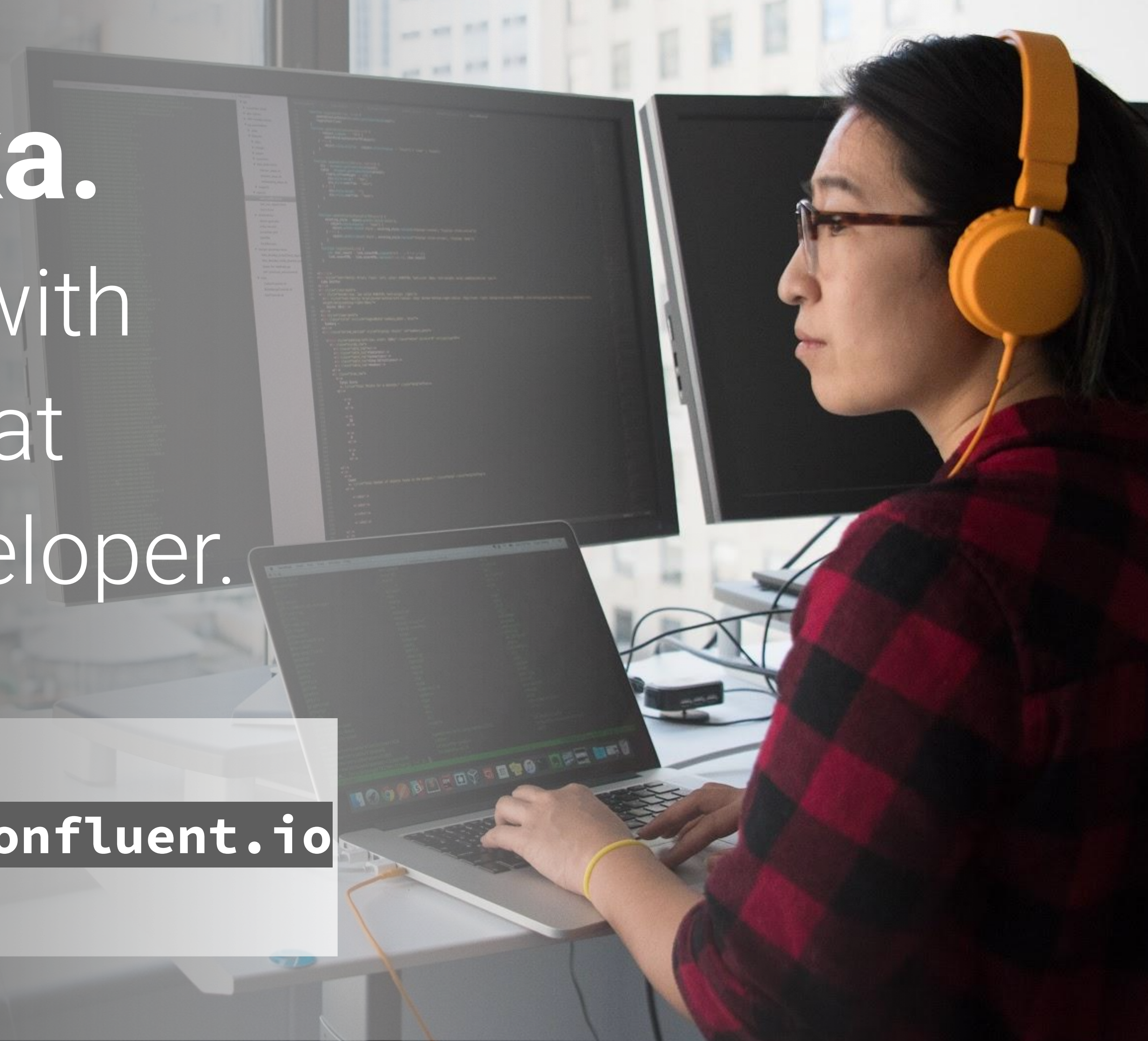


Learn Kafka.

Start building with
Apache Kafka at
Confluent Developer.



developer.confluent.io



Confluent Community Slack group



cnfl.io/slack

@rmoff

#BudapestData

#EEOF

rmoff.dev/talks