

Creating an Apache Kafka® Streaming Data Pipeline



Danica Fine
<https://linktr.ee/thedanicafine>



Lucia Cerchie
<https://luciacerchie.dev/>

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30	COFFEE BREAK	
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00	LUNCH	

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30	COFFEE BREAK	
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00	LUNCH	

What's a data pipeline?

data pipeline

noun \ 'dā-tə 'pīp-līn \

A process (or series of processes) that moves data from a source to a target.

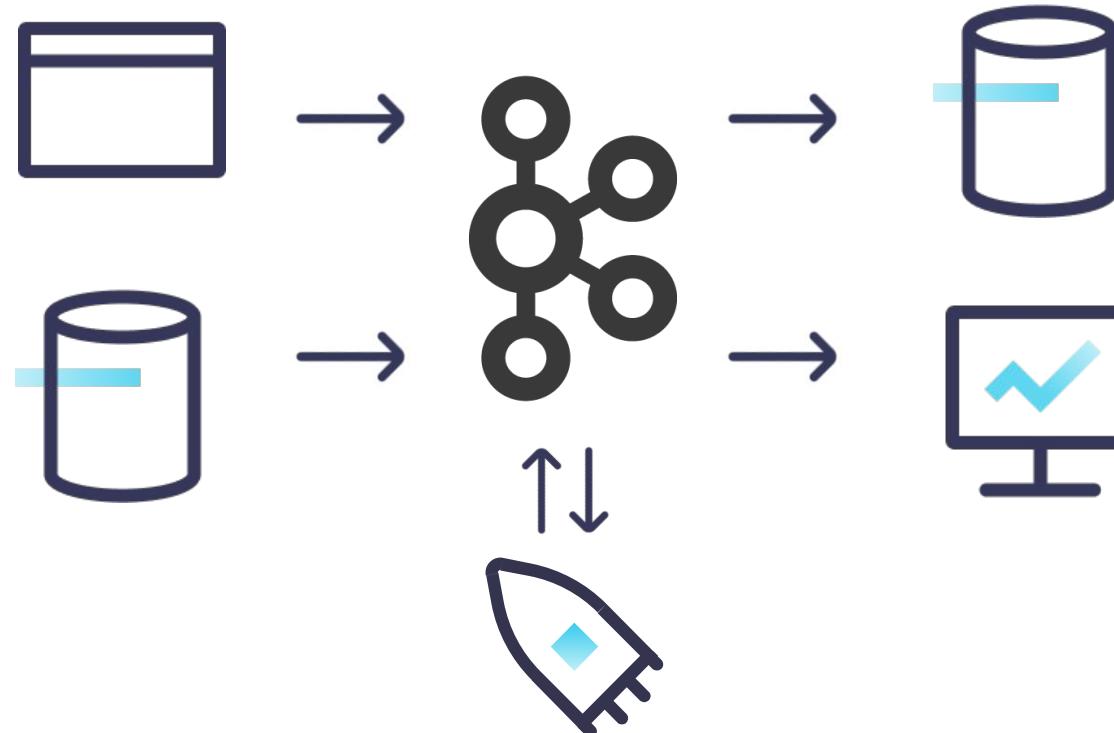
What's a *streaming* data pipeline?

streaming data pipeline

noun \ 'strē-mīn 'dā-tə 'pīp-,līn \

A process (or series of processes) that moves data from a source to a target as the data happens, *in a stream*.

Streaming Data Pipeline



Streaming Data Pipelines

Who cares?

Benefits of Stream Processing

- Event-driven
- Decouples system
 - Agile development
 - Faster time-to-market
- Increase resilience



Stream Processing

What are your use cases?

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30	COFFEE BREAK	
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00	LUNCH	



Kafka? What's *that*?

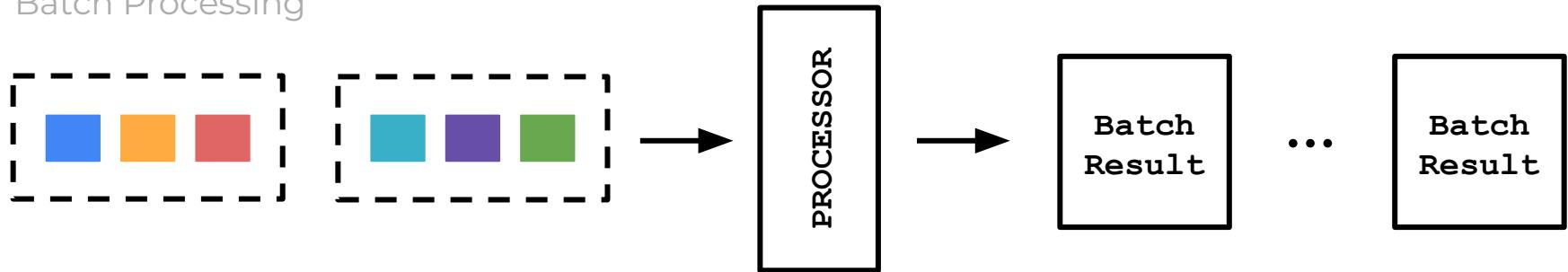
A distributed event streaming platform.



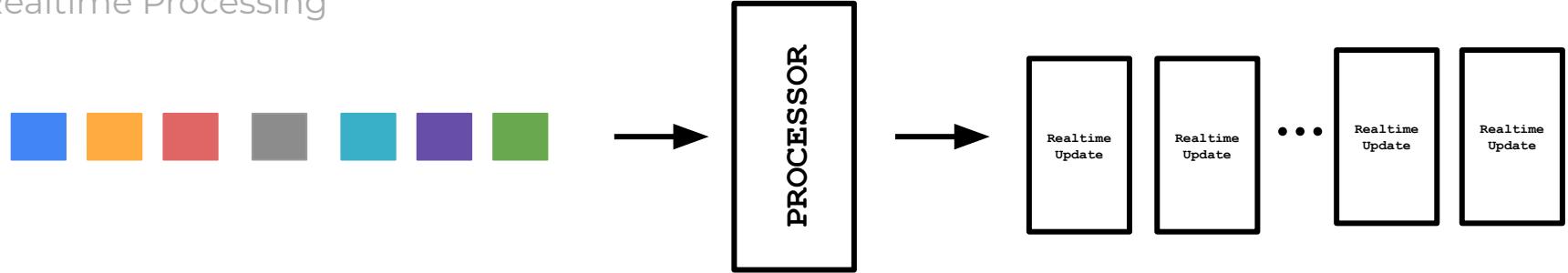
A distributed event ***streaming*** platform.

Paradigm Shift

Batch Processing



Realtime Processing



Some key differences

	Batch Data Processing	Real-Time Data Processing	Streaming Data
Hardware	Most storage and processing resources requirement to process large batches of data.	Less storage required to process the current or recent set of data packets. Less computational requirements.	Less storage required to process current data packets. More processing resources required to "stay awake" in order to meet real-time processing guarantees
Performance	Latency could be minutes, hours, or days	Latency needs to be in seconds or milliseconds	Latency must be guaranteed in milliseconds
Data set	Large batches of data	Current data packet or a few of them	Continuous streams of data
Analysis	Complex computation and analysis of a larger time frame	Simple reporting or computation	Simple reporting or computation



A distributed **event** streaming platform.

Thinking in Events

- Natural way to reason about things
- Indicate that something *has happened*
 - When
 - What/Who
- Immutable pieces of information

Events: what they look like

```
^ Value Header Key  
1 {  
2   "viewtime": 111,  
3   "userid": "User_8",  
4   "pageid": "Page_54"  
5 }
```

Events: what they look like

^ Value Header Key

```
1  [
2    {
3      "key": "task.generation",
4      "stringValue": "0"
5    },
6    {
7      "key": "task.id",
8      "stringValue": "0"
9    },
10   {
11     "key": "current.iteration",
12     "stringValue": "11"
13   }
14 ]
```

Events: what they look like



Events: what they look like

```
▼ {"viewtime":1191,"userid":"User_8","pageid":"Page_56"}  
Partition: 0      Offset: 119      Timestamp: 1666370871970
```

```
▼ {"viewtime":1181,"userid":"User_6","pageid":"Page_28"}  
Partition: 0      Offset: 118      Timestamp: 1666370871421
```

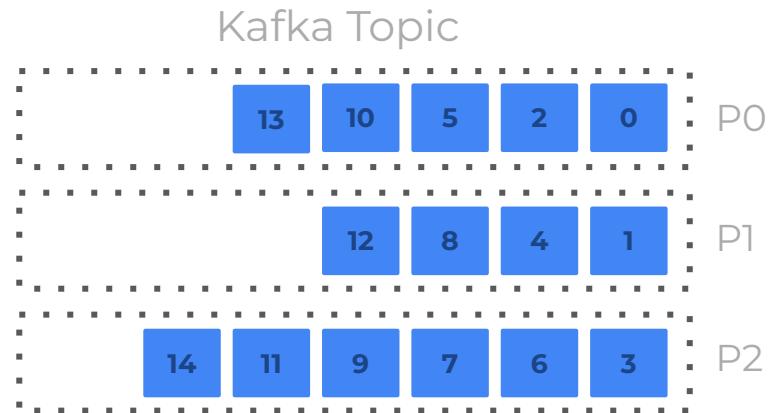
```
▼ {"viewtime":1171,"userid":"User_1","pageid":"Page_15"}  
Partition: 0      Offset: 117      Timestamp: 1666370870791
```



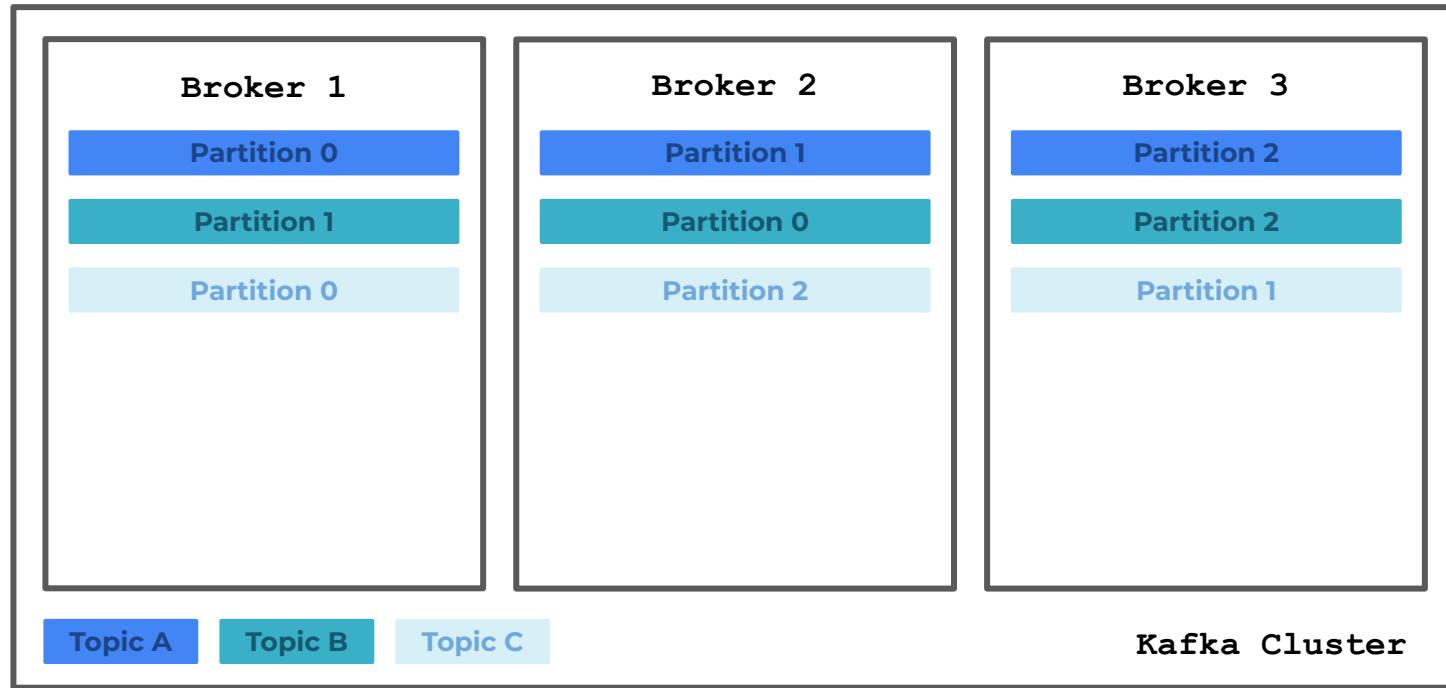
A ***distributed*** event streaming platform.

Kafka Storage

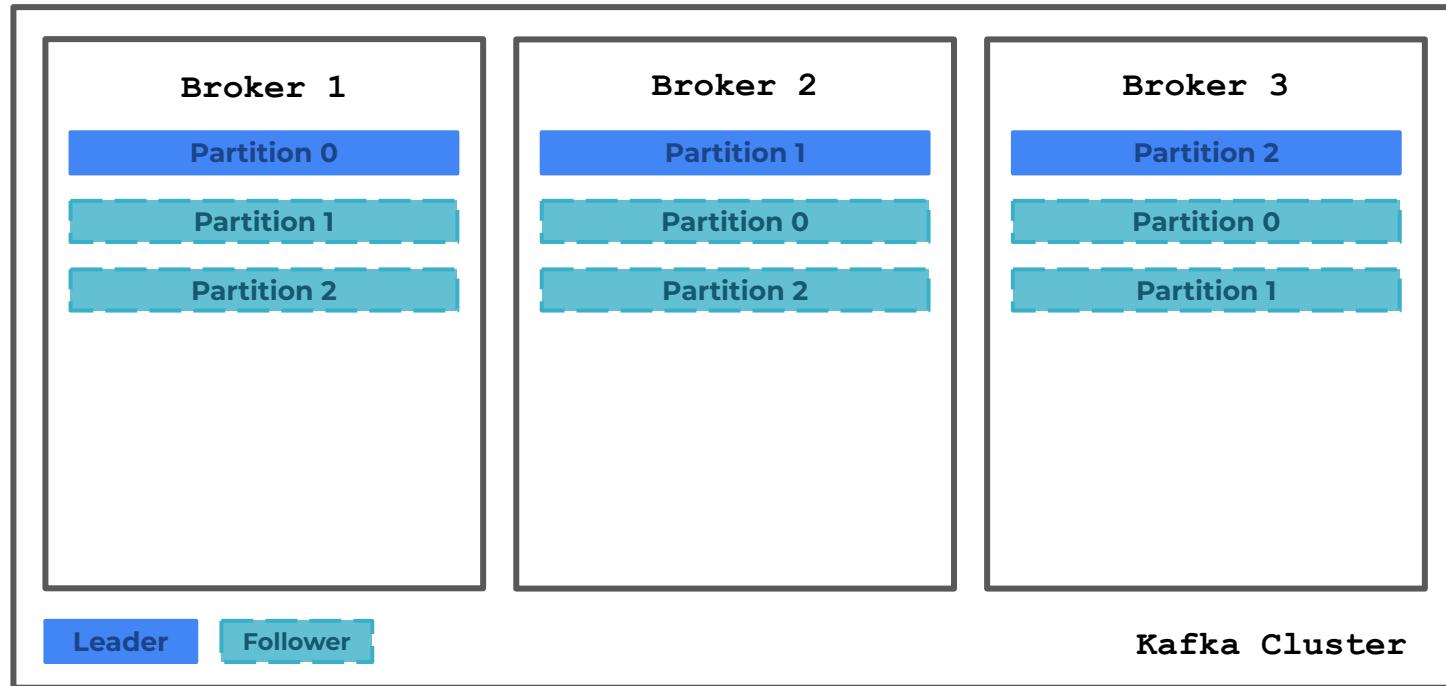
- Topics
 - Basic storage unit
 - Read/WRITE:
 - Producer and consumer clients
 - Completely decoupled
- Partitions
 - Immutable, append-only logs
 - Data is replicated at this level

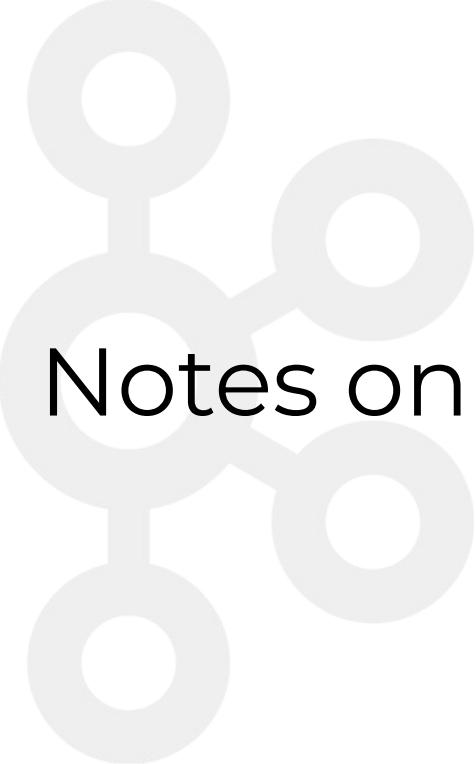


Kafka Cluster

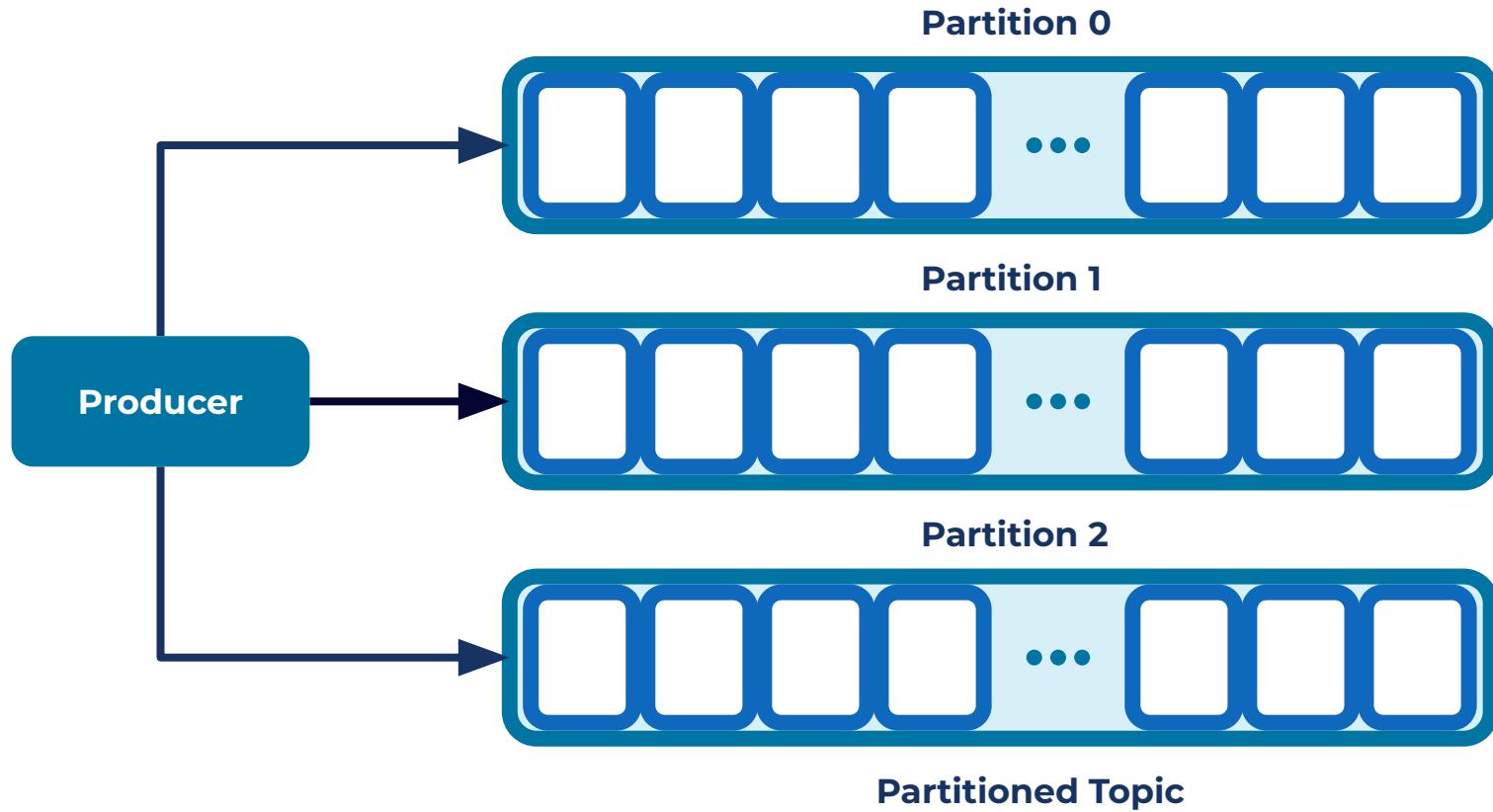


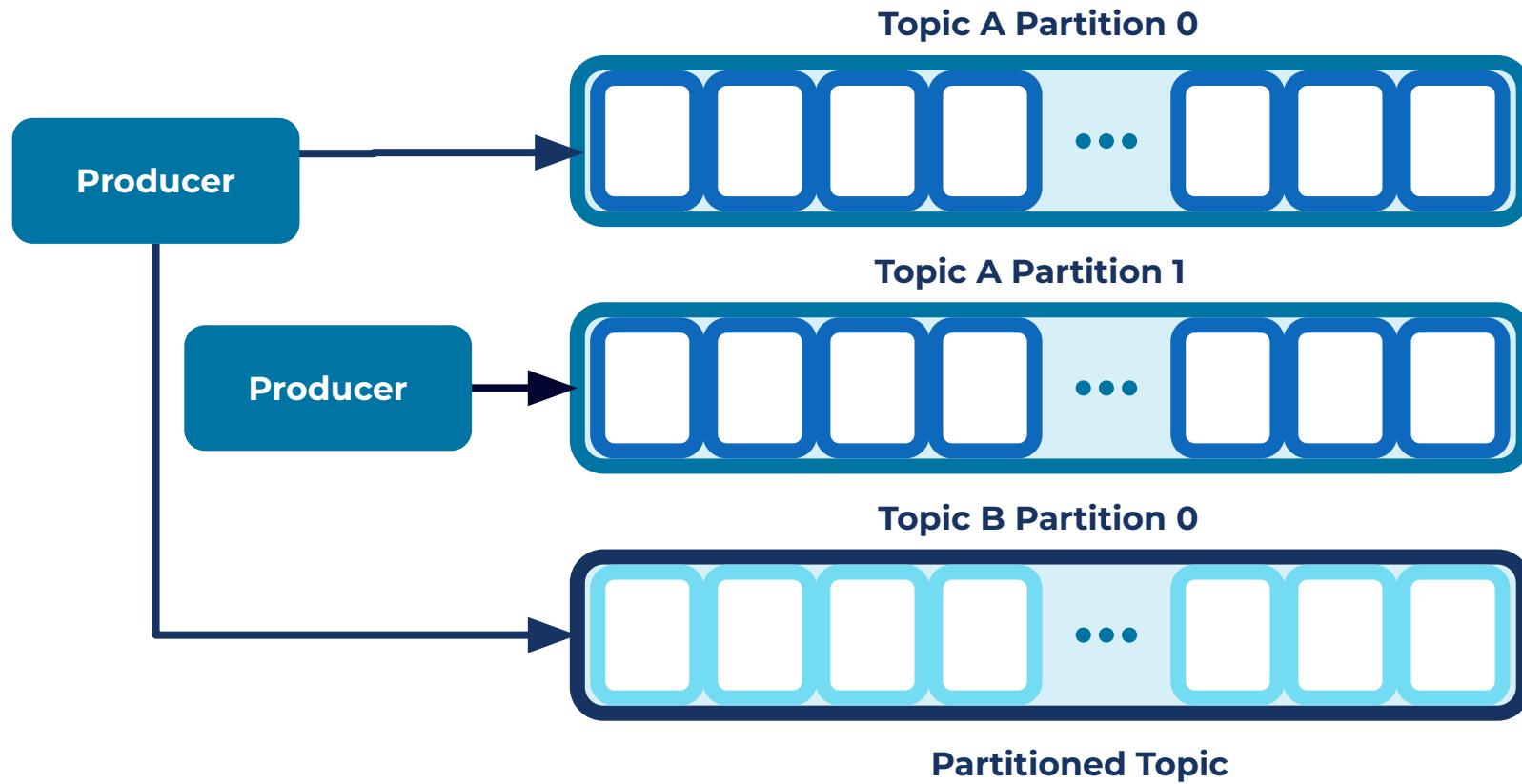
Kafka Cluster



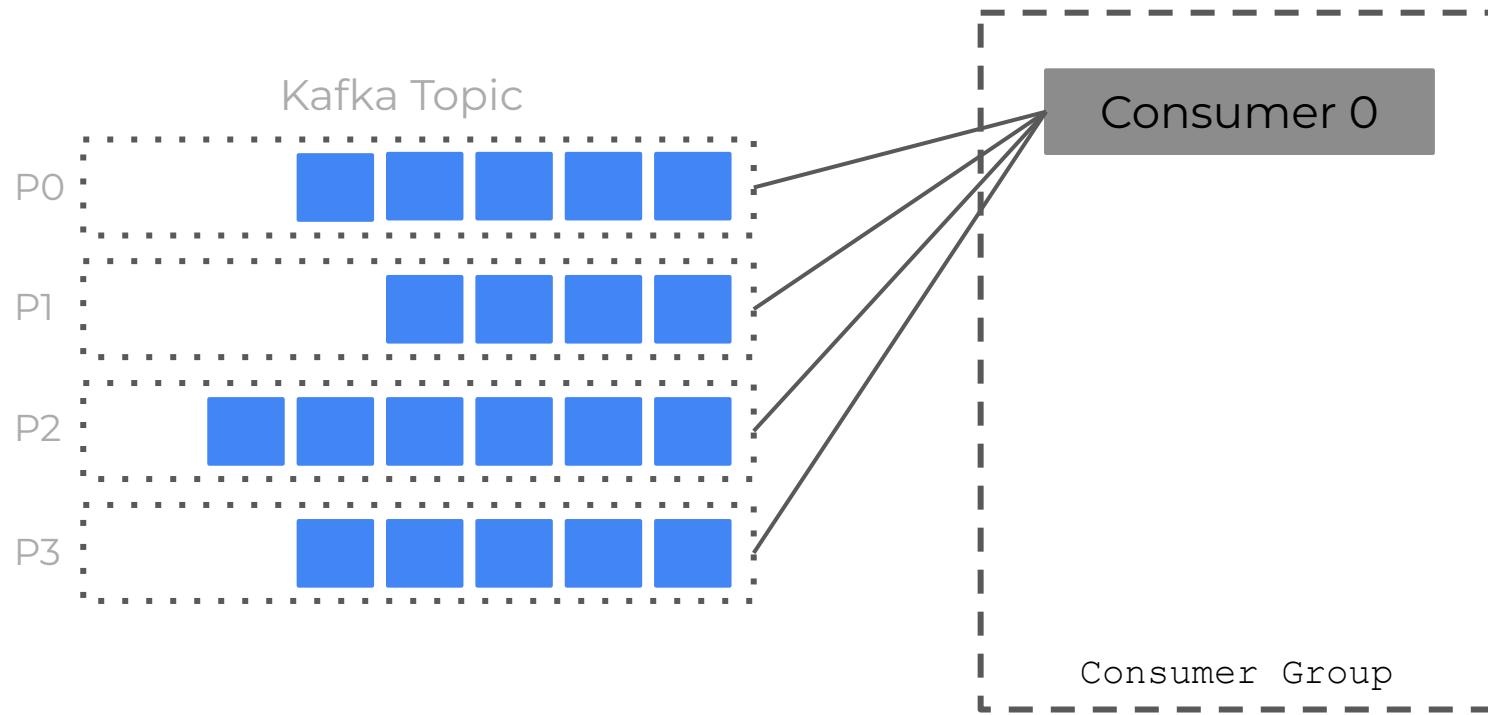


Notes on producers and consumers.

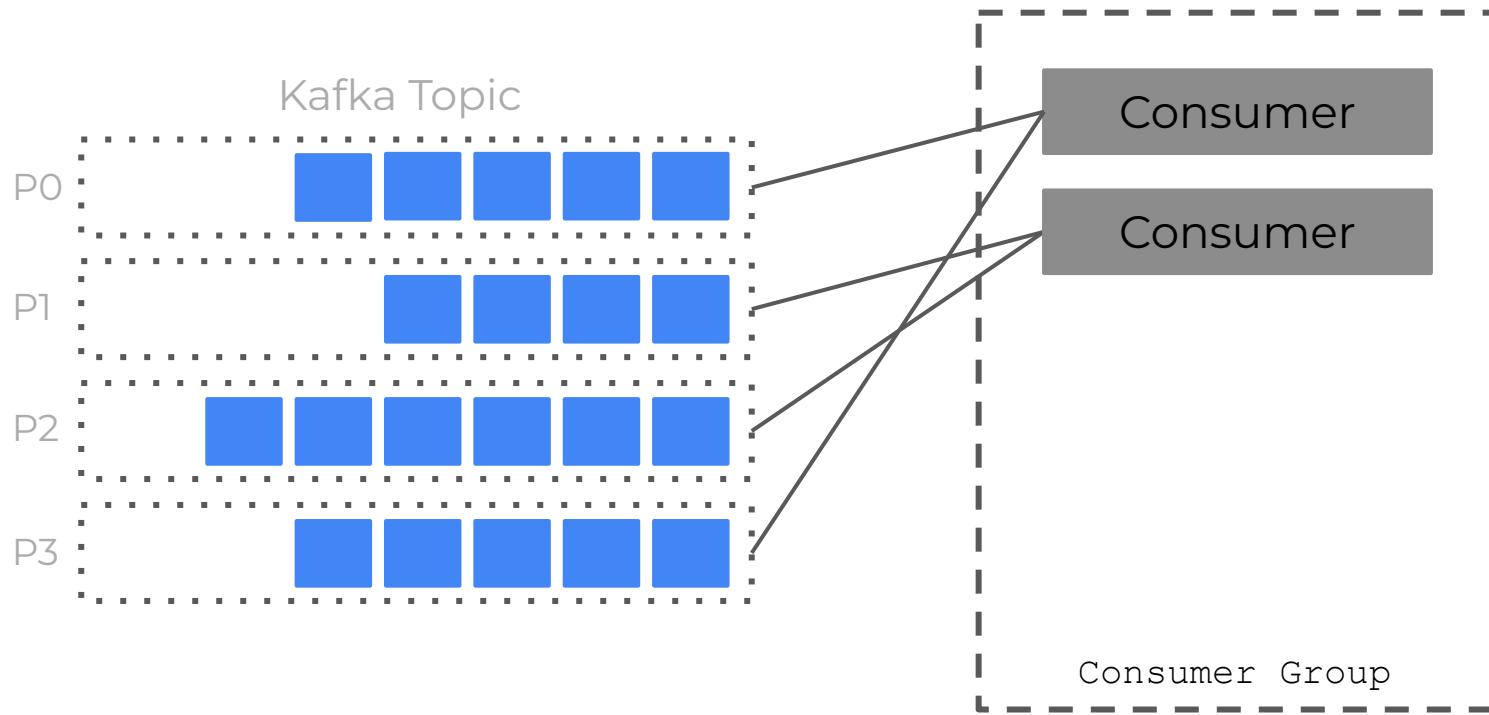




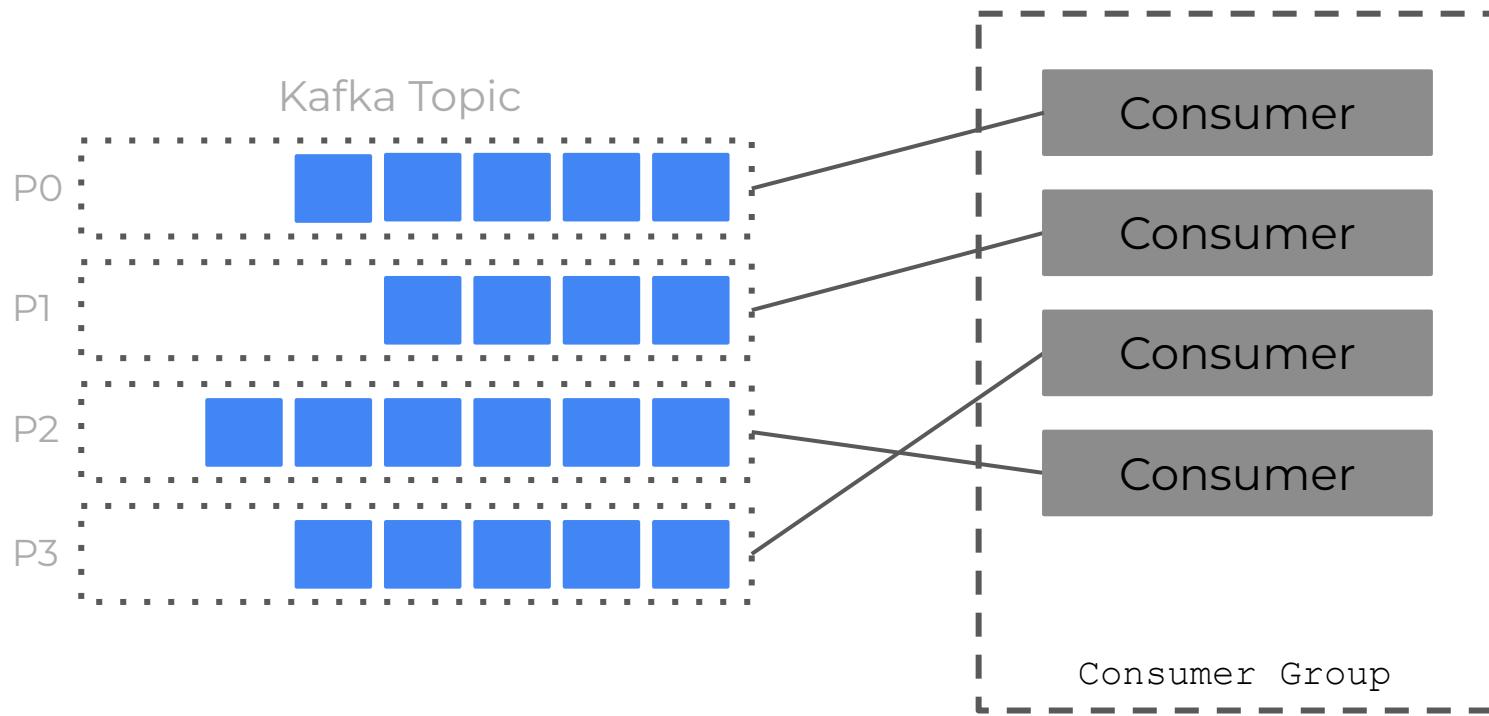
Consumer Groups



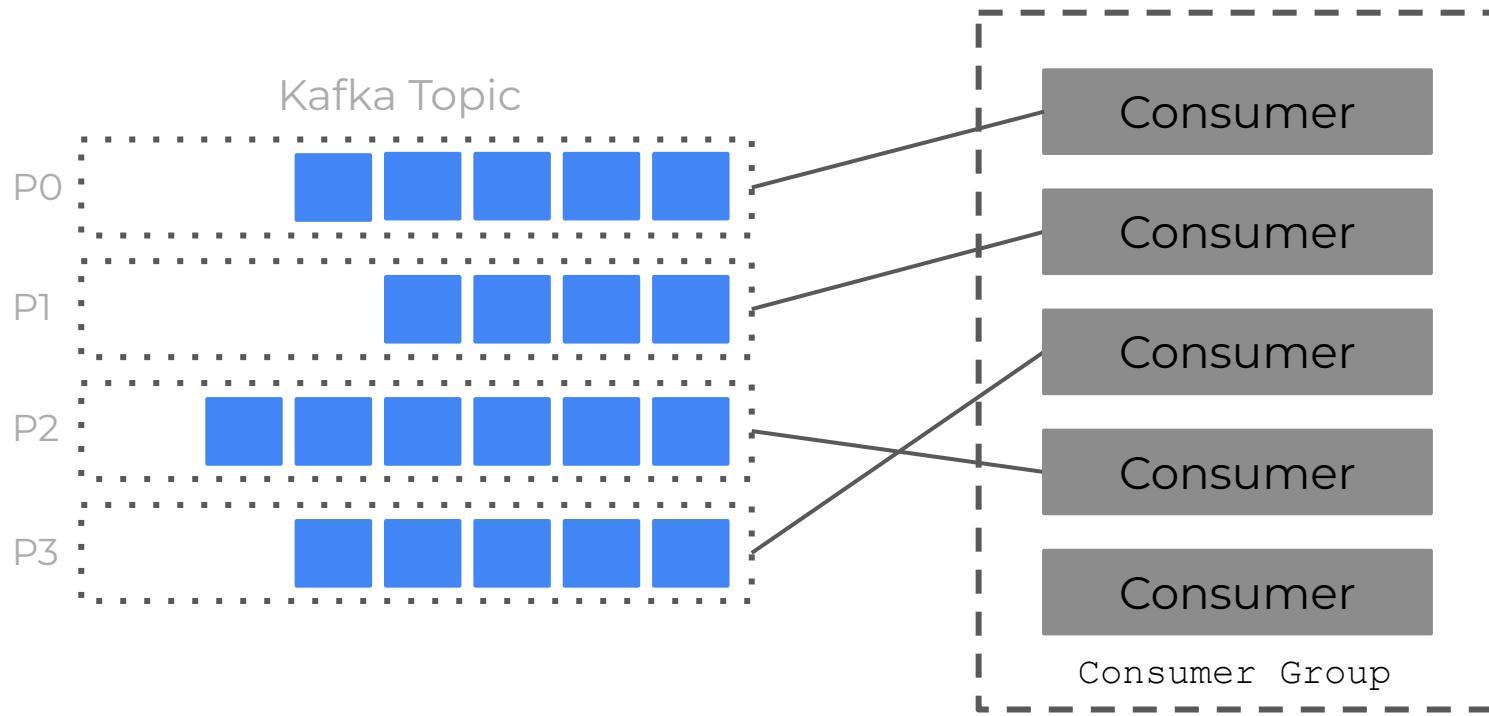
Consumer Groups



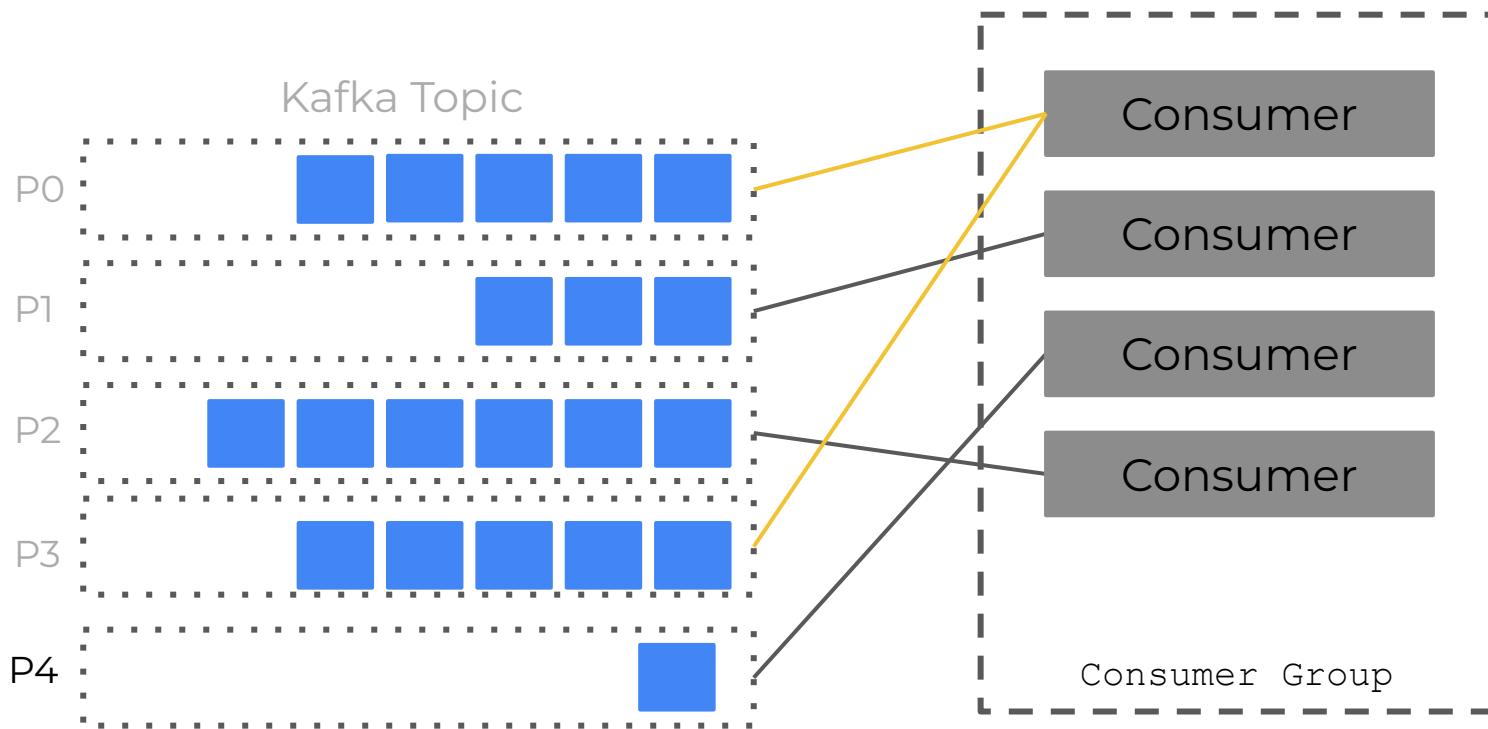
Consumer Groups



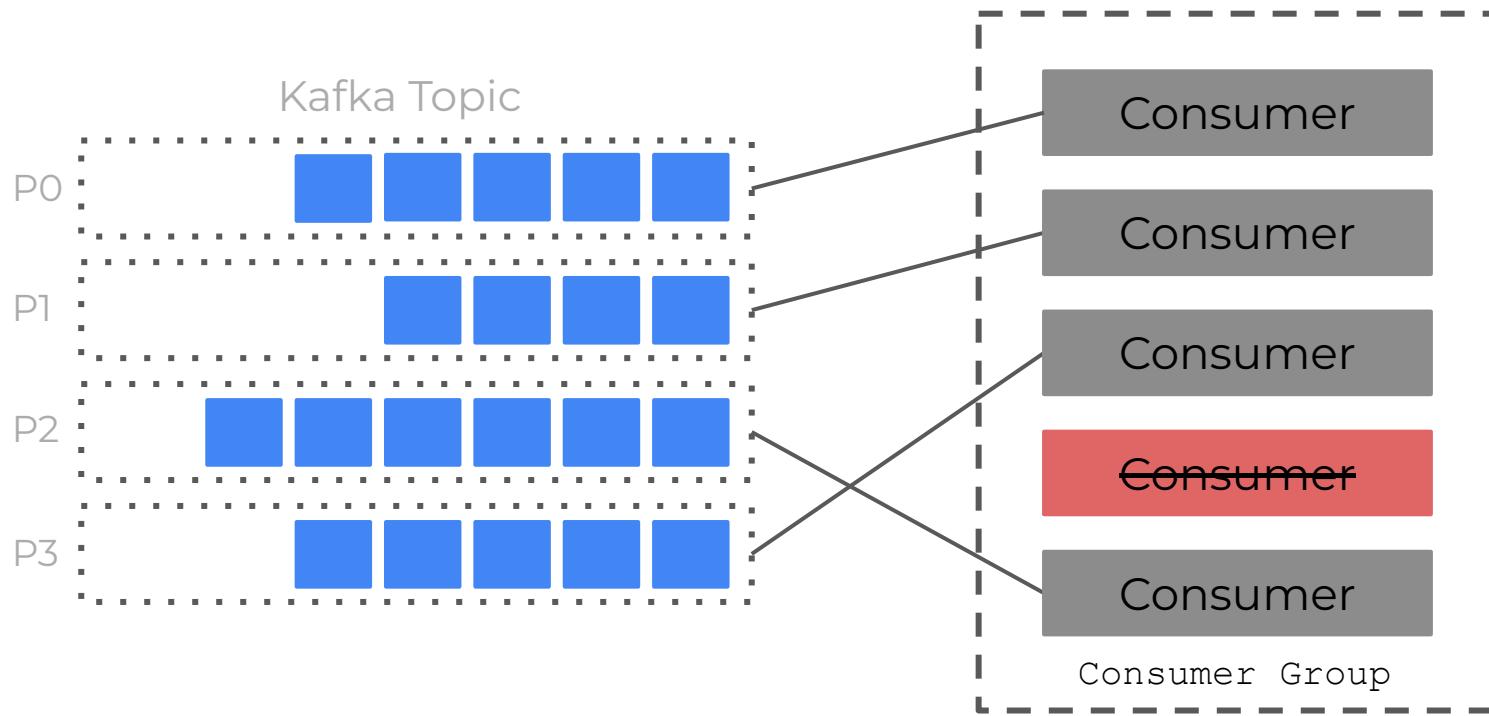
Consumer Groups



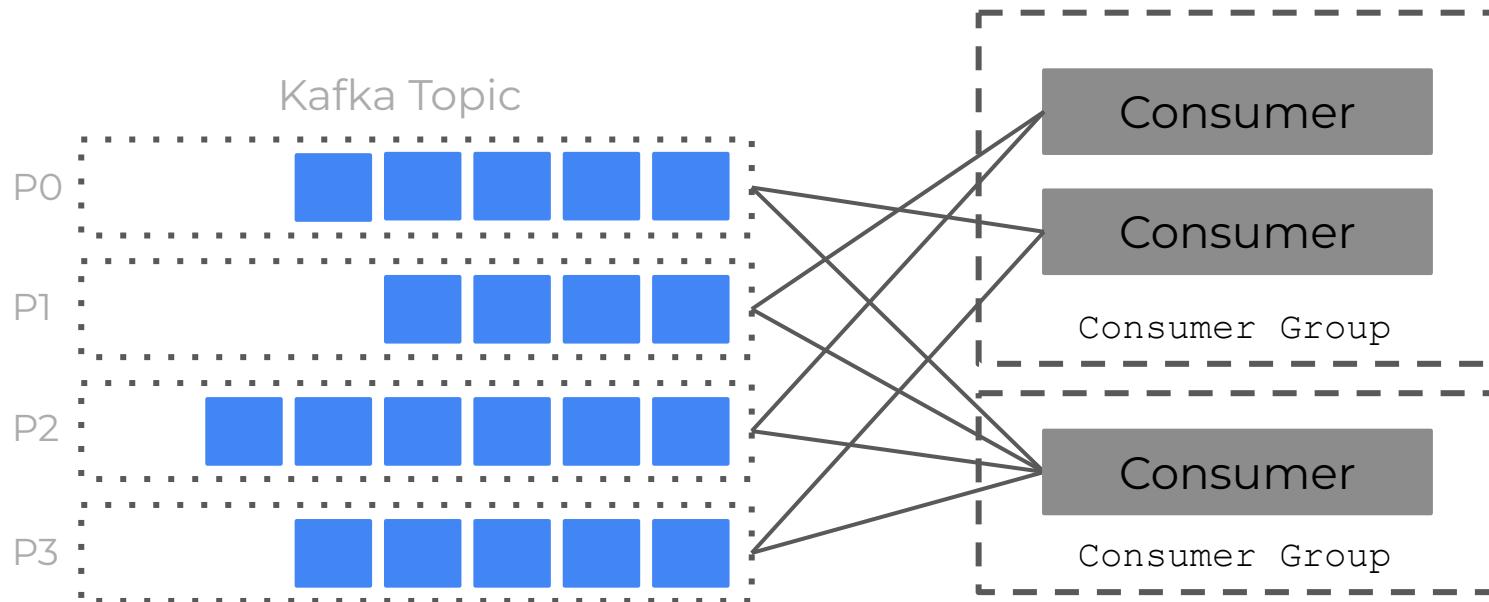
Consumer Groups



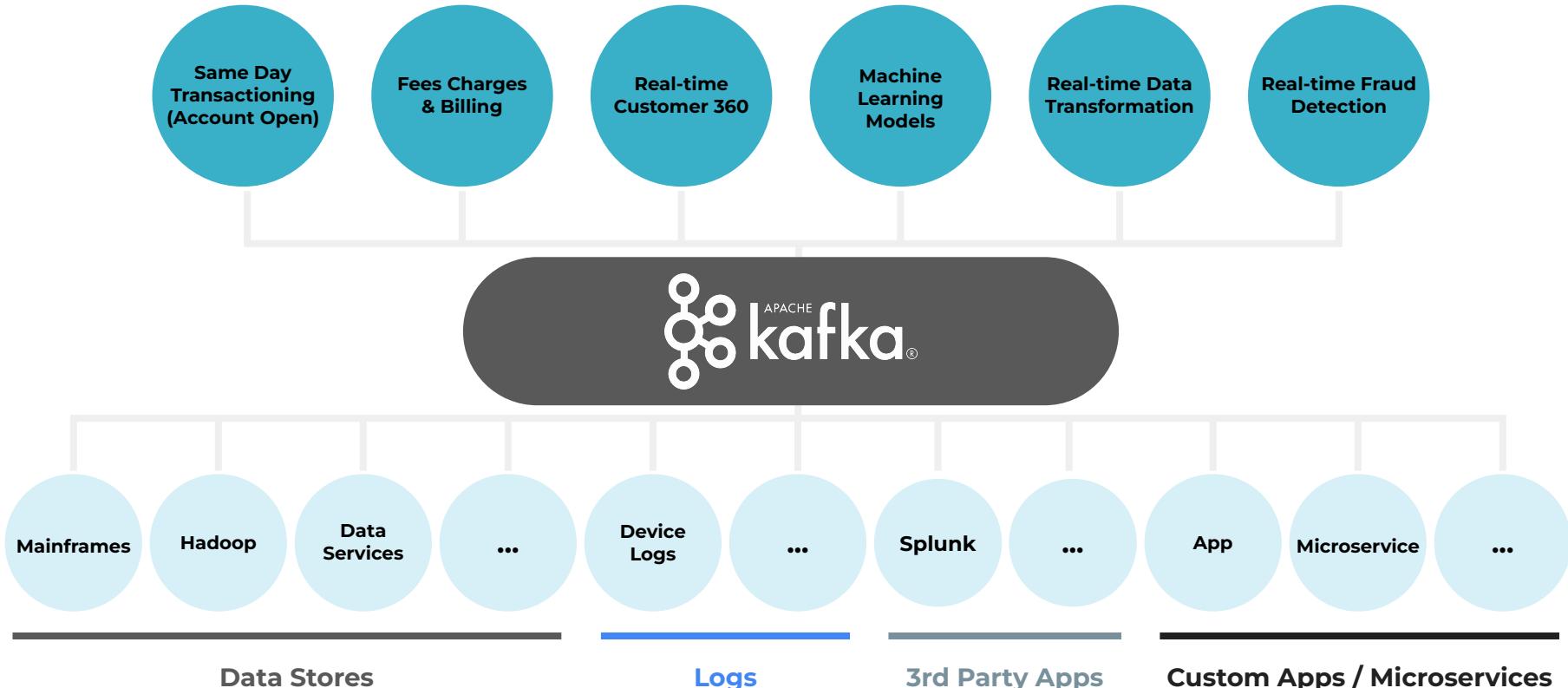
Consumer Groups



Consumer Groups



Unified Fabric for App & Data Modernization





Confluent? What's *that*?

Kafka, but make it easy.

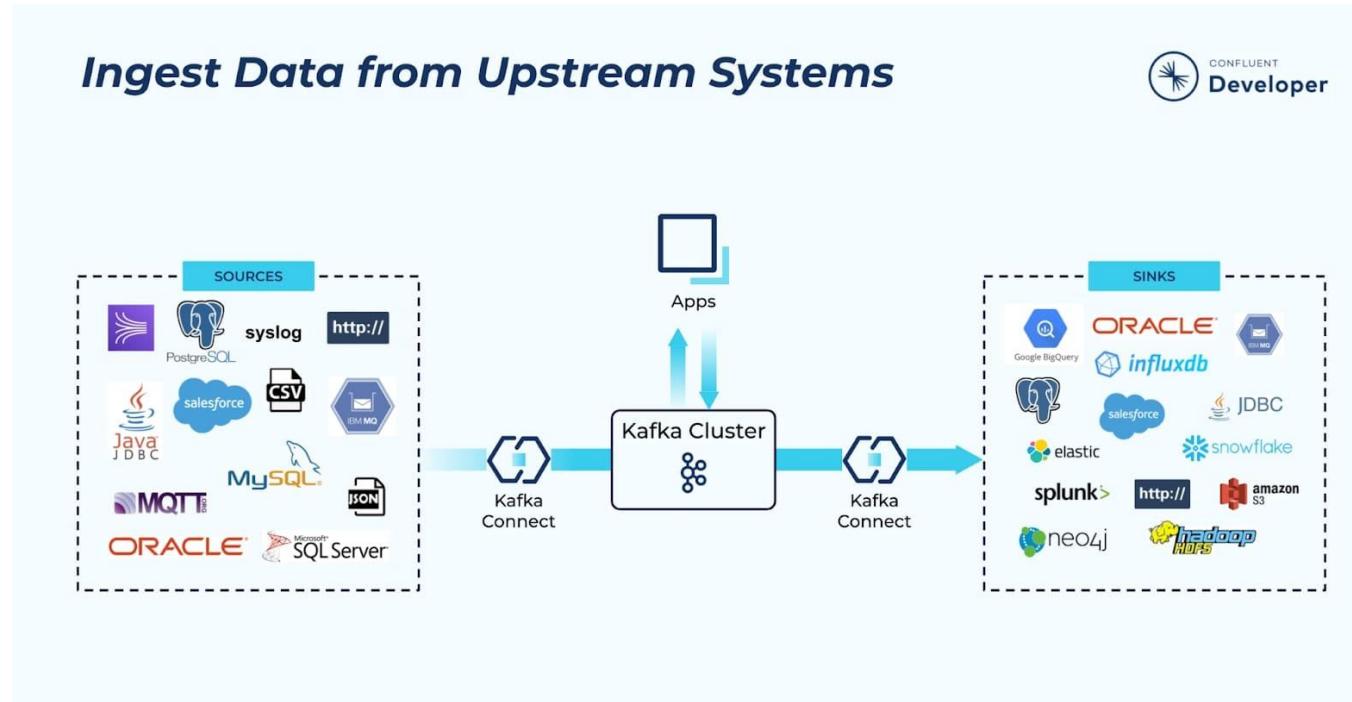
Confluent

- Fully-managed, cloud-based Kafka
- Auxiliary tools:
 - Kafka Connect
 - ksqlDB
 - Schema management

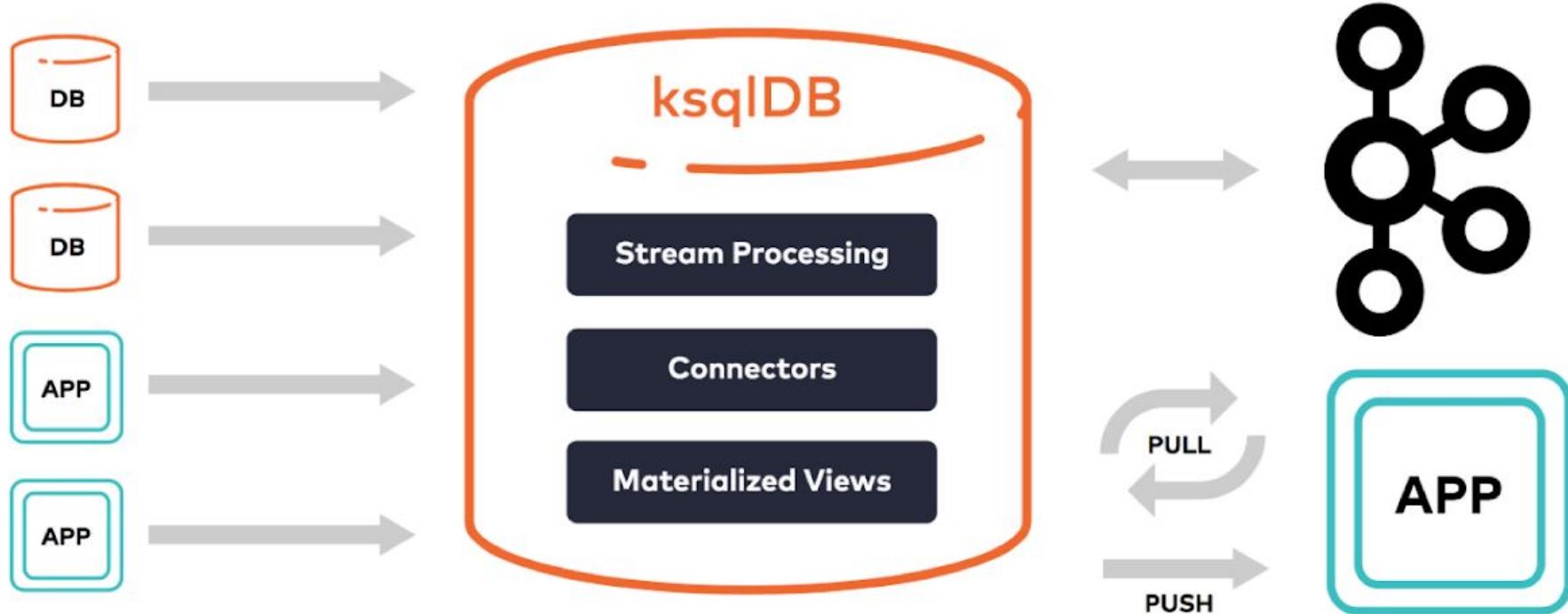


CONFLUENT

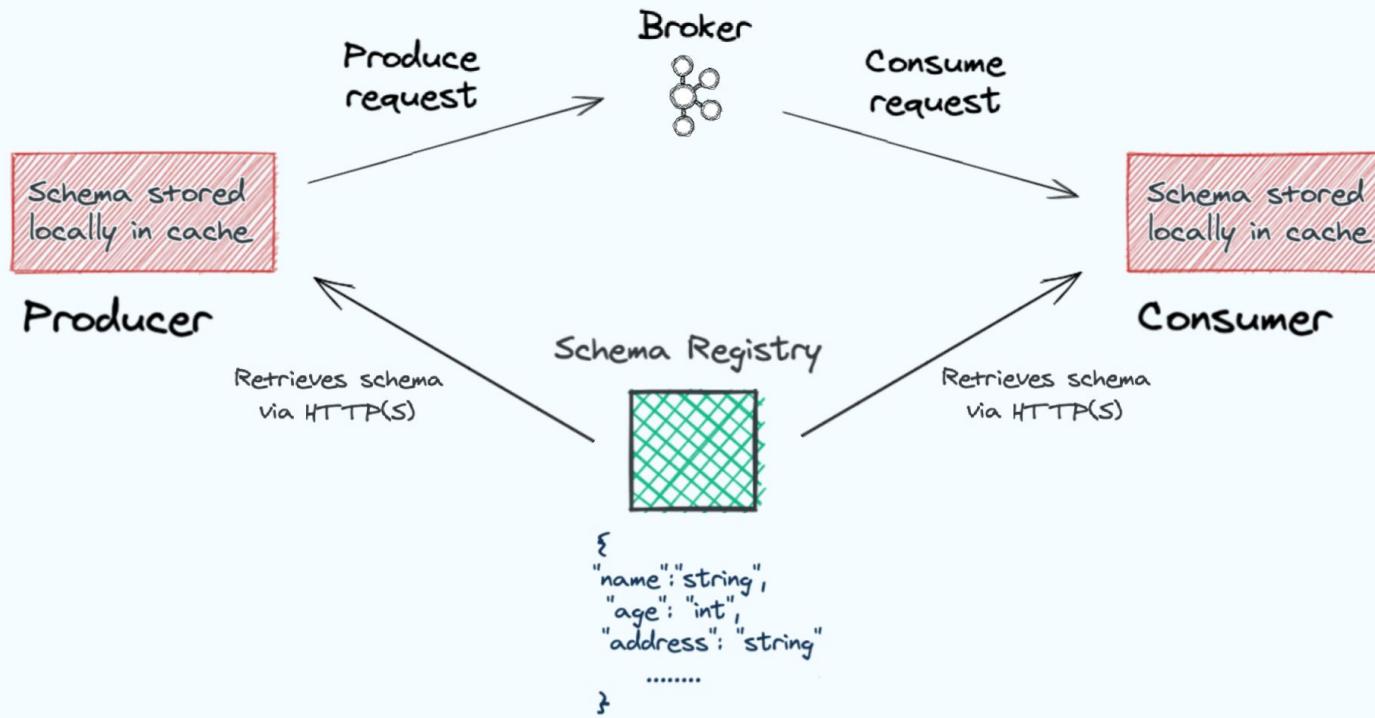
Connectors



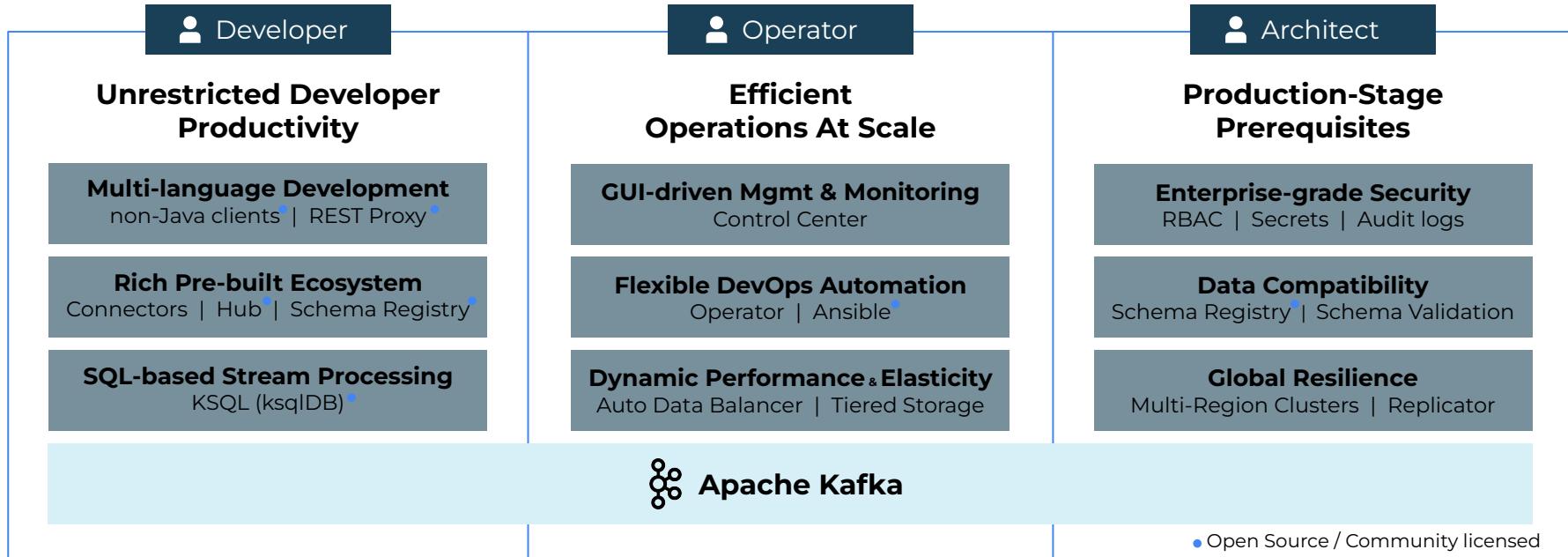
<https://developer.confluent.io/learn-kafka/ksqldb/intro/>



Decoupling with producers and consumers



Confluent



Self Managed Software

Freedom Of Choice



Fully Managed Cloud Service



Enterprise Support



Professional Services

Committer-Led Expertise



Training



Partners

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30		COFFEE BREAK
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00		LUNCH

Hands On

Confluent Cloud and Apache Kafka Installation

Kafka Installation Process



- Python Client Installation
`pip install confluent-kafka`
- Bundled with librdkafka
- Note: M1 MacBooks ***must*** install from source.

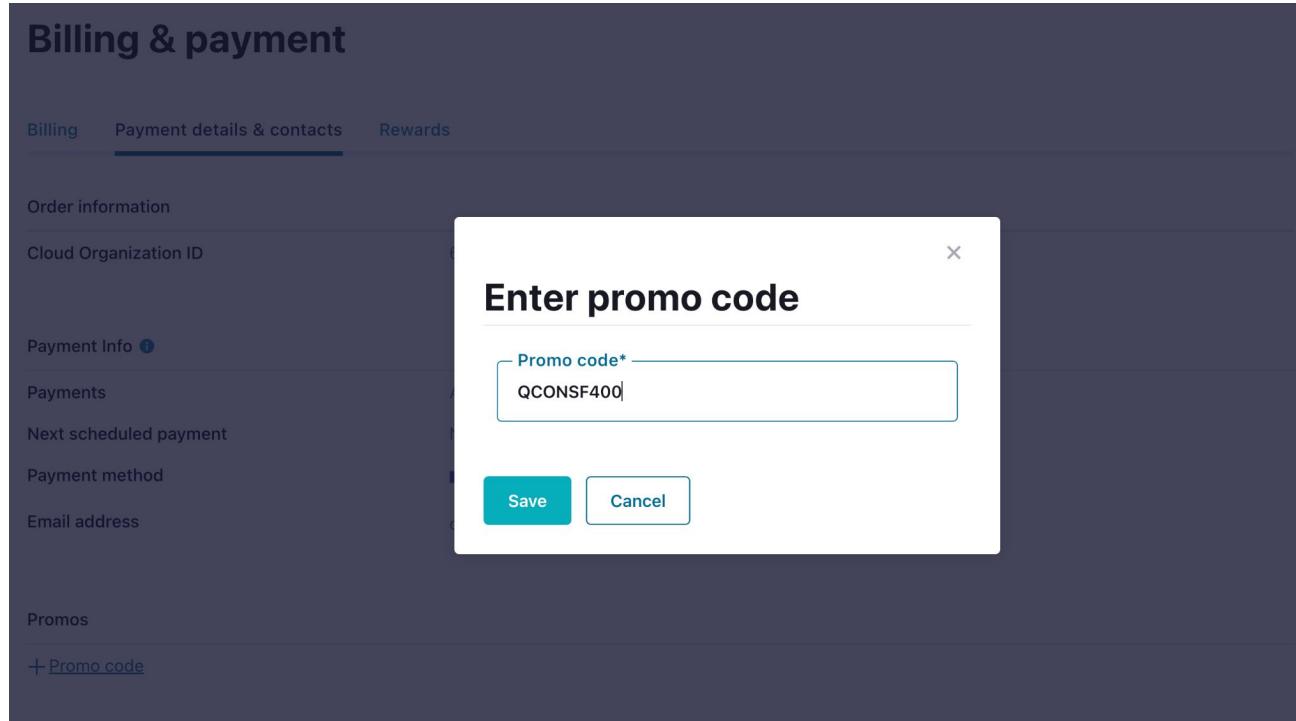
Confluent Cloud Account



Use Promo Code
QCONSF400

Cloud Account Set Up - Apply Promo Code

Menu → Billing & payment → Payment details & contacts → + Promo code



Cloud Account Set Up - Create Cluster

Default Env → + Add cluster → Basic cluster

Create cluster

1. Select cluster type —— 2. Region/zones —— 3. Review and launch

Cluster name ⓘ

Base cost	\$0 /hr
Write	\$0.1265 /GB
Read	\$0.1265 /GB
Storage	\$0.00015972 /GB-hour
Partitions	\$0.0046 /Partition-hour

[Go back](#)

[Review payment method](#)

[Launch cluster](#)

Cloud Account Set Up - Create API Key

Cluster Overview → + API Keys → + Add key → Global Access

Create key

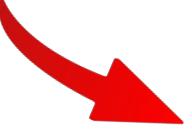
1. Access control 2. Get your API key

Use this API key to connect with the cluster. Store the API key and secret below somewhere safe. This is the only time you'll see the secret.

These credentials can take up to one minute to propagate.

	E657UJ24FWYN5VVW	
	TyURIZf8SsKWcSPrnVR9k9Wg7xRr3dwk Zl8XXKZZ/HyXuufPvP03c/JdMC033LC9	

Description
qconf

Download and continue

Confluent Cloud - Create SR API Key

Cluster Overview → Schema Registry → Change URL → + API Key

 confluent.cloud/environments/env-mozdq/schema-registry/schemas

 confluent.cloud/environments/env-mozdq/schema-registry/api-keys

Morning Agenda

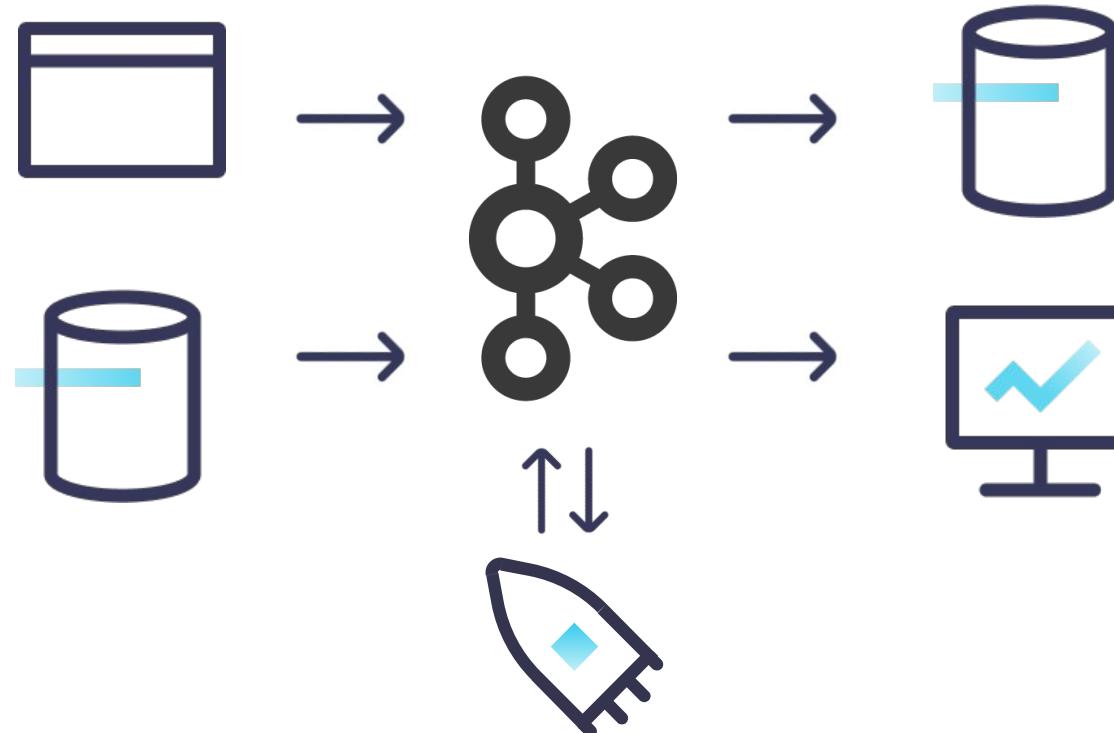
Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30	COFFEE BREAK	
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00	LUNCH	

Coffee Break

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30		COFFEE BREAK
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00		LUNCH

Streaming Data Pipeline



Producing Data Into Kafka

Producer API vs Kafka Connect

Writing Data to Kafka

- Kafka Producer API
 - All of your favorite languages
 - Great when you own the application producing the data
- Kafka Connect
 - Great for integrating with data *at rest*
 - Low- to no-code option
 - Hundreds of data sources (and sinks)

Kafka Connect

- Handles ingest and egress
- Many data sources/sinks:
 - Cloud Warehouses
 - Relational Databases
 - NoSQL Stores
 - SaaS Platforms
- Configuration-driven (no-code, ftw!)
- In the cloud or self-managed



Checkpoint

Think of a dataset you use regularly; would you use the Producer API or Kafka Connect to capture and produce that data into Kafka?

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30	COFFEE BREAK	
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00	LUNCH	

Hands On

DataGen Source Connector

Use Case: Website Analytics

- E-commerce website
- Goal: Assess popularity of products
 - How long do users spend on a product page?
 - Are users from certain regions spending more time on product pages?

Users

Key: User_9

```
{  
    "registertime": 1511778995516,  
    "userid": "User_9",  
    "regionid": "Region_6",  
    "gender": "FEMALE"  
}
```

Pageviews

Key: 49601

```
{  
  "viewtime": 49601,  
  "userid": "User_6",  
  "pageid": "Page_33"  
}
```

DataGen Datasets

- Quickstart options – explore them!
- `schema.string`
 - Optional configuration parameter
 - Provide any JSON string outlining an AVRO schema
 - Must adhere to avro-random-generator standards

Example – Base Schema

```
{  
  "type": "record",  
  "namespace": "survey.bot",  
  "name": "survey",  
  "doc": "Single-question survey.",  
  "fields": [  
    {  
      "doc": "Unique survey ID.",  
      "name": "survey_id",  
      "type": "int"  
    },  
    {  
      "doc": "Survey question.",  
      "name": "question",  
      "type": "string"  
    },  
    {  
      "doc": "Survey summary.",  
      "name": "summary",  
      "type": "string"  
    },  
    {  
      "doc": "Response options.",  
      "name": "options",  
      "type": {  
        "type": "array",  
        "items": "string"  
      }  
    },  
    {  
      "doc": "Survey enabled flag.",  
      "name": "enabled",  
      "type": "boolean"  
    }  
  ]  
}
```

Example – String Options

```
{  
    "doc": "Response options.",  
    "name": "options",  
    "type": {  
        "type": "array",  
        "items": "string"  
        "arg.properties": {  
            "options": [  
                ["Always", "Sometimes", "Never"],  
                ["18-30", "31-50", "50+"],  
                ["Red", "Blue", "Yellow", "Purple"]  
            ]  
        }  
    }  
}
```

Example – Int Range

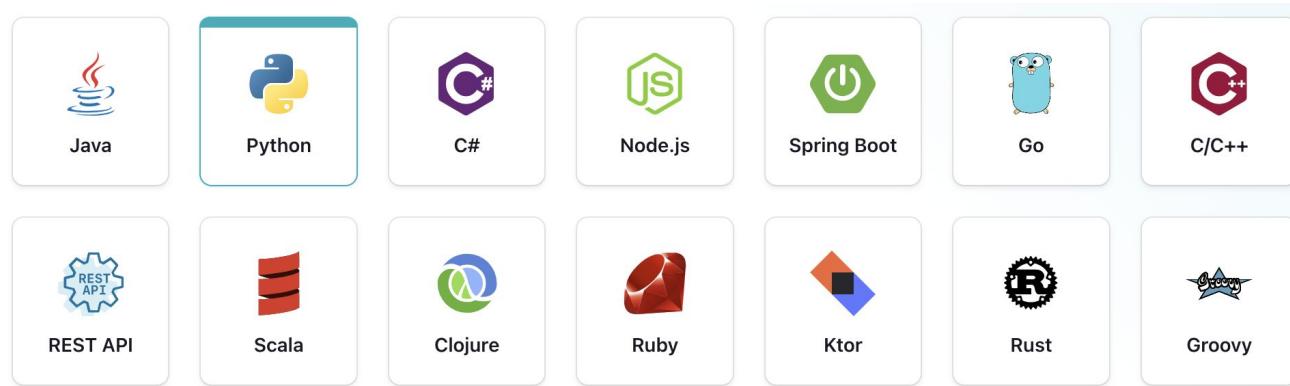
```
{  
    "doc": "Unique survey ID.",  
    "name": "survey_id",  
    "type": {  
        "type": int,  
        "arg.properties": {  
            "range": {  
                "min": 100,  
                "max": 1000  
            }  
        }  
    }  
}
```

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30	COFFEE BREAK	
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00	LUNCH	

Kafka Clients

- Low-level access to Kafka clusters
- High-level stream processing
- Available in most programming languages



Confluent Cloud - Python Client Configuration

```
# Required connection configs for Kafka producer, consumer, and admin
bootstrap.servers=pkc-lzvrd.us-west4.gcp.confluent.cloud:9092
security.protocol=SASL_SSL
sasl.mechanisms=PLAIN
sasl.username={{ CLUSTER_API_KEY }}
sasl.password={{ CLUSTER_API_SECRET }}

# Best practice for higher availability in librdkafka clients prior to 1.7
session.timeout.ms=45000

# Required connection configs for Confluent Cloud Schema Registry
schema.registry.url=https://psrc-22250.us-central1.gcp.confluent.cloud
basic.auth.credentials.source=USER_INFO
basic.auth.user.info={{ SR_API_KEY }}:{{ SR_API_SECRET }}
```

Producer Configurations

- key.serializer and value.serializer
- acks
- partitioner.class
- batch.size
- retries
- delivery.timeout.ms
request.timeout.ms, and linger.ms



Confluent Producer
Documentation

Consumer Configurations

- key.deserializer and value.deserializer
- group.id and partition.assignment.strategy
- session.timeout.ms
- auto.offset.reset and enable.auto.commit



Confluent Consumer
Documentation

Checkpoint Quiz

Which configuration options apply to the following scenarios?

Durable Writes

- Scenario:
 - Highly-available system
 - Minimal impact in a disaster recovery situation
- Goal: Ensure that all brokers containing an in-sync replica will receive every message.
- Producer Configuration:
 - acks=all

Non-Committing Consumer

- Scenario:
 - Compacted topic
 - Consumer is building up state
- Goal: every time the consumer comes up, it should be able to read from the beginning of the Kafka topic to build up state.
- Consumer Configuration:
 - `auto.offset.reset=earliest`
 - `enable.auto.commit=false`

Exactly Once Consumer

- Scenario: A mission-critical application that needs to receive every message once.
- Goal: eliminate data loss or duplicate reads from consumers
- Consumer Configuration:
 - enable.auto.commit=false
 - Manually commit after processing each message.
- Or:
 - enable.auto.commit=true
 - auto.commit.interval.ms=1000

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30	COFFEE BREAK	
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00	LUNCH	

Hands On

Create a Python Client

Confluent Cloud - Create Python Client

Cluster Overview → Clients → + New client → Python

<p>Cluster API Key</p> <p>An API key is required to connect to your cluster. You can either use an existing key or create a new one here.</p>	<p> Copy</p>
<p></p> <p>Schema Registry API Key</p> <p>Schema Registry allows you to enforce a strict schema for your data. Learn more </p>	<pre>1 # Required connection configs for Kafka producer, consumer, and admin 2 bootstrap.servers=pkc-lzvrd.us-west4.gcp.confluent.cloud:9092 3 security.protocol=SASL_SSL 4 sasl.mechanisms=PLAIN 5 sasl.username={{ CLUSTER_API_KEY }} 6 sasl.password={{ CLUSTER_API_SECRET }} 7 8 # Best practice for higher availability in librdkafka clients prior 9 # to 1.7 10 session.timeout.ms=45000 11 12 # Required connection configs for Confluent Cloud Schema Registry 13 schema.registry.url=https://psrc-22250.us- 14 central1.gcp.confluent.cloud 15 basic.auth.credentials.source=USER_INFO 16 basic.auth.user.info={{ SR_API_KEY }}:{{ SR_API_SECRET }}</pre>

Python Client Sandbox

- Update client configurations to point to your cluster using appropriate API Key
- Execute producer and consumer
- Play in the sandbox:
 - Alter configuration parameters
 - Add new datasets/schemas and produce them to Confluent Cloud
 - Conduct some analysis
 - Make the users and pageviews data more interesting



<https://github.com/danicafine/qconsf-python>

Morning Agenda

Time	Module Type	Description
09:00	Instruction	Introductions Streaming Data Pipelines – Benefits and Use Cases
09:15	Instruction	Apache Kafka – Architecture, Ecosystem, and Confluent
10:00	Hands On	Apache Kafka Installation Confluent Cloud Account Set Up
10:30	COFFEE BREAK	
10:45	Instruction	Introduction to Kafka Connect
10:55	Hands On	Kafka Connect Data Generator on Confluent Cloud
11:15	Instruction	Producer and Consumer Configurations Python Code Samples and Walk Through
11:45	Hands On	Setting Up Producers and Consumers
12:00	LUNCH	

Lunch

Afternoon Agenda

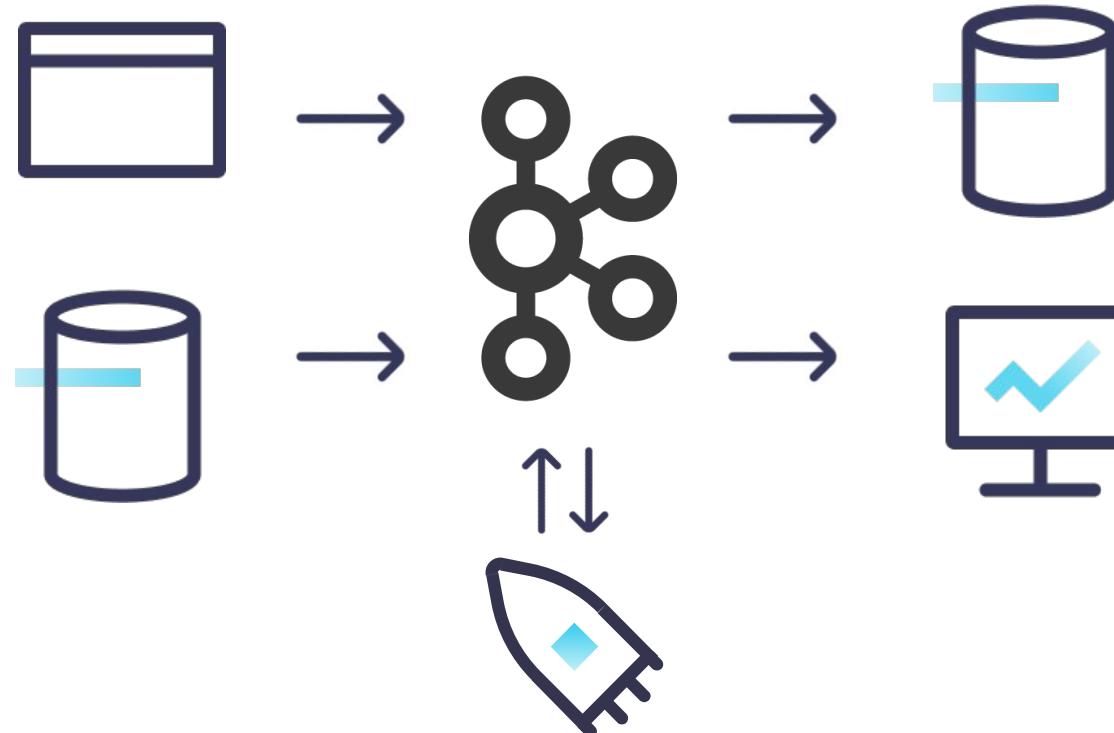
Time	Module Type	Description
12:00		LUNCH
13:00	Instruction	Stream Processing Techniques in the Apache Kafka Ecosystem ksqlDB in Confluent Cloud
13:45	Hands On	Build a Stream Processing Application with ksqlDB
14:30		COFFEE BREAK
14:45	Instruction	Advanced Stream Processing Pipelines Use Cases and Inspiration
15:00	Hands On	Outline an End-to-End Pipeline Independent Working Time to Complete Stream Processing Pipelines

Afternoon Agenda

Time	Module Type	Description
12:00		LUNCH
13:00	Instruction	Stream Processing Techniques in the Apache Kafka Ecosystem ksqlDB in Confluent Cloud
13:45	Hands On	Build a Stream Processing Application with ksqlDB
14:30		COFFEE BREAK
14:45	Instruction	Advanced Stream Processing Pipelines Use Cases and Inspiration
15:00	Hands On	Outline an End-to-End Pipeline Independent Working Time to Complete Stream Processing Pipelines

Stream Processing Techniques

Streaming Data Pipeline



Consumer/Producer API

- Lowest level
- Consume – Process – Produce
- Handle state and fault-tolerance

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

Example: Consumer/Producer API

```
# set up Kafka Consumer for Users
consumer = clients.consumer(clients.user_deserializer(),
    'consumer-group-usermask', [config['topics']['users']])

# set up Kafka Producer for UserMask
producer = clients.producer(clients.usermask_serializer())

# start consumption loop
while True:
    msg = consumer.poll(1.0)

    if msg is not None:
        # received a message
        user = msg.value()

        # transform the data
        usermask = UserMask(user.registertime, user.userid, user.regionid)

        # send data to Kafka
        producer.produce(config['topics']['usermasks'], key=str(user.userid), value=usermask)
```

Kafka Streams

- Java library for stream processing
- Rich API of transformations
- Built-in state handling and failover

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

Example: Kafka Streams

```
static Topology buildTopology(String inputTopic, String outputTopic) {  
    Serde<String> stringSerde = Serdes.String();  
    StreamsBuilder builder = new StreamsBuilder();  
    builder  
        .stream(inputUserTopic, Consumed.with(stringSerde, stringSerde))  
        .peek((k,v) -> logger.info("Observed event: {}", v))  
        .mapValues(u -> u.maskUser())  
        .peek((k,v) -> logger.info("Transformed event: {}", v))  
        .to(outputUserTopic, Produced.with(stringSerde, stringSerde));  
    return builder.build();  
}
```

ksqldb

- SQL syntax
- Kafka Streams under the hood
- Available as self-managed and cloud-based offerings

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

Example: ksqlDB

```
CREATE STREAM usermasks  
  WITH (kafka_topic='usermasks', value_format='avro', partitions=1) AS  
    SELECT  
      REGISTERTIME,  
      USERID,  
      REGIONID,  
      MASK(GENDER)  
    FROM USERS;
```

ksqlDB

Use Case: Website Analytics

- Register data: Stream or Table?
- Pre-processing:
 - Filtering out unnecessary data
 - Enrichment
- Aggregation

Creating a Table

```
CREATE TABLE users (
    id STRING PRIMARY KEY
) WITH (
    kafka_topic='users',
    value_format='AVRO'
);
```

```
1 CREATE TABLE users (
2     id STRING PRIMARY KEY
3 ) WITH (
4     kafka_topic='users',
5     value_format='AVRO'
6 );
7
```

● [Add query properties](#)

auto.offset.reset

= Earliest



+Add another field

Stop

Run query

Creating a Stream

```
CREATE STREAM pageviews WITH (
    kafka_topic='pageviews',
    value_format='AVRO'
);
```

Rekeying a Stream

```
CREATE STREAM pageviews_rekeyed WITH (
    kafka_topic='pageviews-rekeyed',
    value_format='AVRO'
) AS
SELECT
    *
FROM PAGEVIEWS
PARTITION BY USERID;
```

Filtering a Stream

```
CREATE STREAM pageviews_filtered WITH (
    kafka_topic='pageviews-filtered',
    value_format='AVRO'
) AS
SELECT
    ROWTIME AS pageviewtime,
    *
FROM PAGEVIEWS_REKEYED
WHERE USERID != 'User_1'
EMIT CHANGES;
```

Stream-to-Table Join

```
CREATE STREAM pageviews_enriched WITH (
    kafka_topic='pageviews-enriched',
    value_format='AVRO'
) AS
SELECT
    pageviews.viewtime          AS viewtime,
    FROM_UNIXTIME(pageviews.pageviewtime) AS pageviewtime,
    pageviews.pageid            AS pageid,
    pageviews.userid            AS userid,
    users.regionid              AS regionid,
    users.gender                AS gender,
    FROM_UNIXTIME(users.registertime) AS registertime
FROM pageviews_filtered AS pageviews
INNER JOIN users
ON pageviews.userid = users.id
EMIT CHANGES;
```

Pull Query

SELECT

*

FROM pageviews

WHERE userid != User_0' ;



Filter by keyword

▼ {"VIEWTIME":7311,"USERID":"User_3","PAGEID":"Page_68"}

▼ {"VIEWTIME":7301,"USERID":"User_7","PAGEID":"Page_99"}

▼ {"VIEWTIME":7291,"USERID":"User_6","PAGEID":"Page_30"}

Windowing

```
CREATE TABLE pageviewRegionalAnalysis WITH (
    kafka_topic='pageview-regional-analysis',
    format='AVRO'
) AS
SELECT
    regionid,
    COUNT(*) AS views,
    AVG(viewtime) AS viewtime_avg
FROM pageviews_enriched
WINDOW TUMBLING (SIZE 30 MINUTES)
GROUP BY regionid
EMIT FINAL;
```

Afternoon Agenda

Time	Module Type	Description
12:00		LUNCH
13:00	Instruction	Stream Processing Techniques in the Apache Kafka Ecosystem ksqldb in Confluent Cloud
13:45	Hands On	Build a Stream Processing Application with ksqldb
14:30		COFFEE BREAK
14:45	Instruction	Advanced Stream Processing Pipelines Use Cases and Inspiration
15:00	Hands On	Outline an End-to-End Pipeline Independent Working Time to Complete Stream Processing Pipelines

Hands On

Stream Processing With ksqlDB

ksqldb - Create a Cluster

Cluster Overview → ksqldb → Create with Tutorial → Follow Prompts

New cluster

1. Access control —— 2. Name and sizing

Cluster name

Cluster size

1	\$0.255553 USD /hr	▼
---	--------------------	---

i Confluent streaming units (CSUs) determine the amount of data your cluster can process. High Availability is enabled by default on clusters with 8 CSUs or more.

[Learn more ↗](#)

Configuration options

Default

Advanced

ksqldb Stream Processing

- Use Case: Website Analytics
 - Register data: Stream or Table?
 - Pre-processing:
 - Filtering out unnecessary data
 - Enrichment
 - Aggregation
- Another use case?



ksqldb Reference Materials

Afternoon Agenda

Time	Module Type	Description
12:00		LUNCH
13:00	Instruction	Stream Processing Techniques in the Apache Kafka Ecosystem ksqlDB in Confluent Cloud
13:45	Hands On	Build a Stream Processing Application with ksqlDB
14:30		COFFEE BREAK
14:45	Instruction	Advanced Stream Processing Pipelines Use Cases and Inspiration
15:00	Hands On	Outline an End-to-End Pipeline Independent Working Time to Complete Stream Processing Pipelines

Coffee Break

Afternoon Agenda

Time	Module Type	Description
12:00		LUNCH
13:00	Instruction	Stream Processing Techniques in the Apache Kafka Ecosystem ksqlDB in Confluent Cloud
13:45	Hands On	Build a Stream Processing Application with ksqlDB
14:30		COFFEE BREAK
14:45	Instruction	Advanced Stream Processing Pipelines Use Cases and Inspiration
15:00	Hands On	Outline an End-to-End Pipeline Independent Working Time to Complete Stream Processing Pipelines

Use Cases and Inspiration

Online Multiplayer System

- Kafka forms the backbone of the system
 - All user-input is treated as an event
 - Allows for easy conflict-resolution
- Wide variety of applications
 - Gaming
 - Collaborative work tools (think Google Docs)

Houseplant Alerting System

- Raspberry Pi-based, IoT Project
 - Collect moisture readings from plants using sensors
 - Use Python script to send readings to Kafka as events
- Stream Processing
 - Enrich readings with houseplant metadata
 - Determine if plants need to be watered and send a text alert

Kafka Connect - Elastic Sink Connector

Create a free account

The screenshot shows the Elastic homepage with the following content:

- Welcome home**: A header section featuring four colored cards: yellow (Enterprise Search), pink (Observability), teal (Security), and blue (Analytics). Each card has an icon and a brief description.
- Enterprise Search**: Create search experiences with a refined set of APIs and tools.
- Observability**: Consolidate your logs, metrics, application traces, and system availability with purpose-built UIs.
- Security**: Prevent, collect, detect, and respond to threats for unified protection across your infrastructure.
- Analytics**: Explore, visualize, and analyze your data using a powerful suite of analytical tools and applications.
- Get started by adding integrations**: A section with a callout graphic showing various data types (logs, metrics, files) being processed. It includes a description and three buttons: **Add integrations**, **Try sample data**, and **Upload a file**.
- Description**: To start working with your data, use one of our many ingest options. Collect data from an app or service, or upload a file. If you're not ready to use your own data, add a sample data set.

Kafka Connect - Elastic Sink Connector

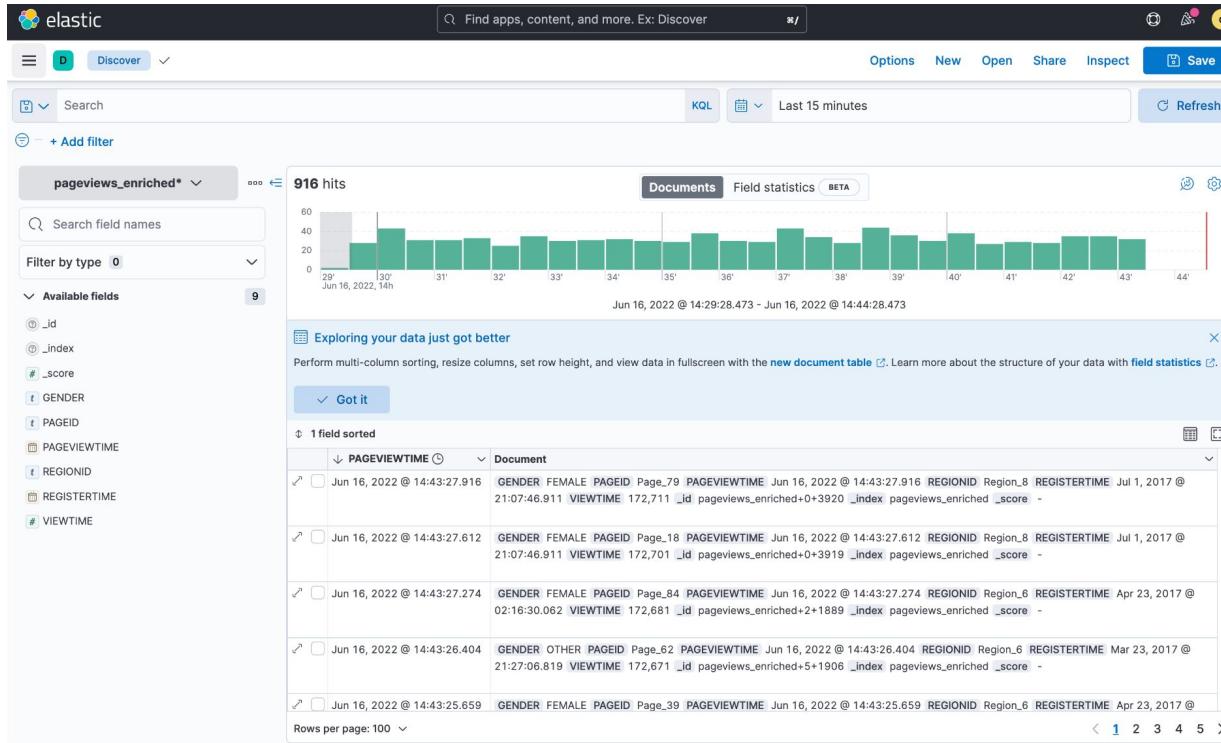
Configure Elastic Sink Connector

The screenshot shows the Confluent Cloud interface with the following details:

- Header:** CONFLUENT logo, Stream Catalog, LEARN, Notifications, and a three-dot menu.
- Breadcrumbs:** HOME > ENVIRONMENTS > DEFAULT > CLUSTER_0 >
- Left sidebar:** Cluster overview, Topics, Data integration (Clients, Connectors selected), Stream lineage, and ksqlDB.
- Search bar:** A search bar with the text "elasticsearch".
- Results:** A search result card for the "Elasticsearch Service Sink" connector. It includes:
 - Icon:** A yellow, black, and teal circular icon.
 - Name:** Elasticsearch Service Sink
 - Type:** Sink
 - Description:** The Kafka Connect Elasticsearch Service sink connector for Confluent Cloud moves data from Apache Kafka® to Elasticsearch.
 - Documentation:** A link to the connector's documentation.
 - Enterprise support:** Confluent supported.
 - Performance:** Available fully managed, Push data in less than 10 minutes.
 - Call-to-action:** A "Get started" button.
- Bottom navigation:** Schema Registry, CLI and tools.

Kafka Connect - Elastic Sink Connector

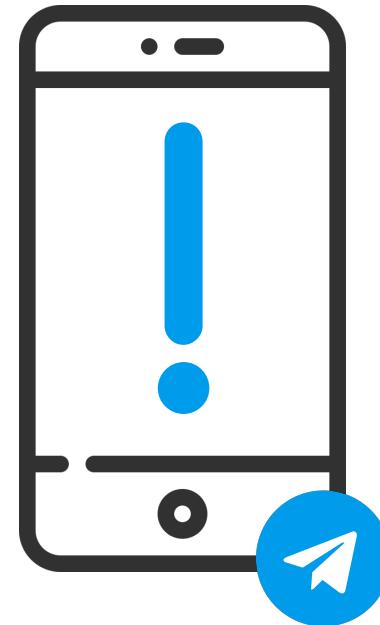
Conduct analysis



Kafka Connect - HTTP Sink Connector

Send alerts to your phone

- Telegram App
- Create a bot
- Send data to bot endpoint



Sample Data Sets

- US Census Bureau -
<https://www.census.gov/quickfacts/fact/table/US/PST045221>
- Northern Data Hub -
https://datahub.bradford.gov.uk/ebase/datahubext.eb?ebd=0&ebp=10&ebz=1_1666813862928
- Network Rail -
<https://datafeeds.networkrail.co.uk/ntrod/login;jsessionid=1C499ED60E92B2BCD2332267555EE21B>

DataGen Schema String – Base Schema

```
{  
  "type": "record",  
  "namespace": "survey.bot",  
  "name": "survey",  
  "doc": "Single-question survey.",  
  "fields": [  
    {  
      "doc": "Unique survey ID.",  
      "name": "survey_id",  
      "type": "int"  
    },  
    {  
      "doc": "Survey question.",  
      "name": "question",  
      "type": "string"  
    },  
    {  
      "doc": "Survey summary.",  
      "name": "summary",  
      "type": "string"  
    },  
    {  
      "doc": "Response options.",  
      "name": "options",  
      "type": {  
        "type": "array",  
        "items": "string"  
      }  
    },  
    {  
      "doc": "Survey enabled flag.",  
      "name": "enabled",  
      "type": "boolean"  
    }  
  ]  
}
```

DataGen Schema String – String Options

```
{  
    "doc": "Response options.",  
    "name": "options",  
    "type": {  
        "type": "array",  
        "items": "string"  
        "arg.properties": {  
            "options": [  
                ["Always", "Sometimes", "Never"],  
                ["18-30", "31-50", "50+"],  
                ["Red", "Blue", "Yellow", "Purple"]  
            ]  
        }  
    }  
}
```

DataGen Schema String – Int Range

```
{  
    "doc": "Unique survey ID.",  
    "name": "survey_id",  
    "type": {  
        "type": int,  
        "arg.properties": {  
            "range": {  
                "min": 100,  
                "max": 1000  
            }  
        }  
    }  
}
```

What are your use cases?

How can we implement them now?

Afternoon Agenda

Time	Module Type	Description
12:00		LUNCH
13:00	Instruction	Stream Processing Techniques in the Apache Kafka Ecosystem ksqlDB in Confluent Cloud
13:45	Hands On	Build a Stream Processing Application with ksqlDB
14:30		COFFEE BREAK
14:45	Instruction	Advanced Stream Processing Pipelines Use Cases and Inspiration
15:00	Hands On	Outline an End-to-End Pipeline Independent Working Time to Complete Stream Processing Pipelines

Hands On

End-to-End Pipelines