

Kafka

Apache Kafka is a distributed event streaming platform used to handle high throughput, real-time data streams.

It enables application to publish, store, process and subscribe to stream of records efficiently

eg. Chat solution

Key features of Apache Kafka

Scalability

Fault tolerance

Durability

High throughput

Decoupling



Why do we need Kafka?

Before Kafka, we were using traditional systems

Databases

real time, high volume

message Queue

distributed data streams

ETL batch processing

Traditional MQ

Message retention

msgs removed after consumption

Kafka

Messages can be persisted for a configurable period

Scalability

Limited due to broker-based approach

Scales easily horizontally with partitions

Data throughput

lower

High

Fault tolerance

lose the message

Replication ensures durability

Batch processing

Process data in chunks at intervals

Process data continuously in real time

Batch Processing

Data latency

High

Kafka

Processing type

Process accumulated data over interval

low

Example

Newspaper

Real time as it arrives

live news channel

Kafka Use Cases

E-commerce

Finance

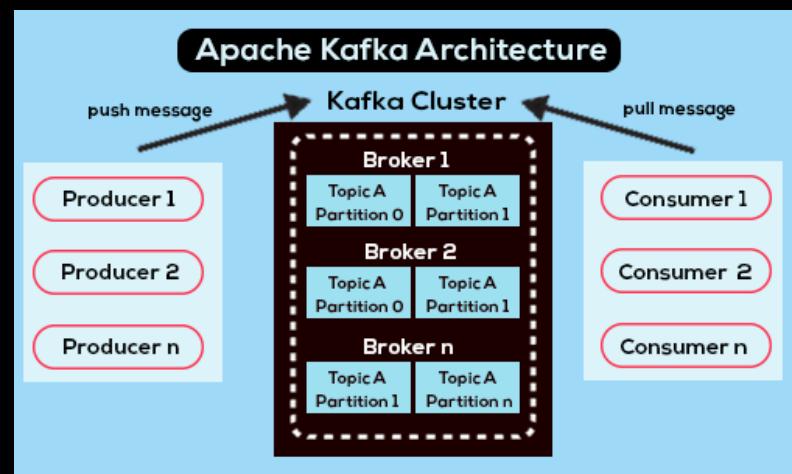
IoT and Telecom

Social Media

Log management

Kafka Architecture

distributed and scalable architecture
for real time data



1. Broker

Server storing and managing the messages.

2. Topic

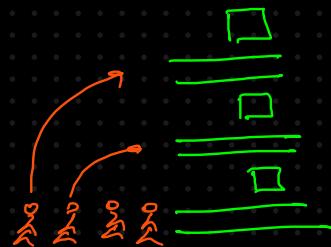
Category, logical grouping or channel where messages are published.

Community → Group 3
Topic 2 → Group 2
Topic 1 → Group 1

Stock prices
↓
trading app
new restaurant
↓
food delivery app

3. Partition

way to divide the data within topic to improve scalability
and parallelism



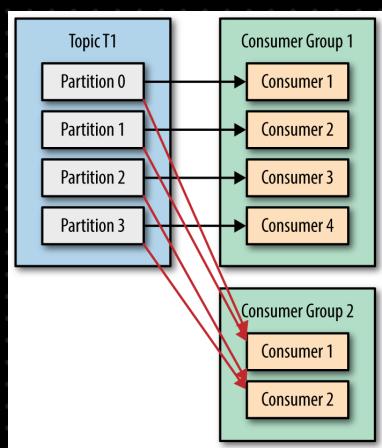
1. Round robin

2. Key-Based Partitioning

4. Consumer Groups

4 people
throwing ball

6 are
catching



Customer support

Agent ⇔ consumer
calls ⇔ messages

⇒ Each consumer from a consumer group reads from a separate Partition

⇒ We cannot have
 $\text{no. of consumers} > \text{no. of groups}$
 no. of partitions

Offset management

Bookmark

Offset is a unique ID
assigned to each message
inside the partition.

Replication and Fault tolerance

2 15 ↴
6 5 8 3 7 8 9
. . .

Kafka ensures data
availability and fault
tolerance through
replication



How to run and use Kafka

1. Local

2. Confluent → Kafka → LinkedIn (2010)

Open sourced (2011)

founders of Kafka created Confluent → 2014

easy to set up
Cloud native
managed service

Amazon (AWS) → MSK

GCP → Pub Sub + Kafka on GKE

Azure → Event Hubs + Confluent Cloud

Kafka on Confluent

Environment



Cluster



Topic



Partitions

← send a message

use coupon code

CONFLUENTDEV1

to avoid entering
credit card

valid as of
Jan 25

if not valid
try

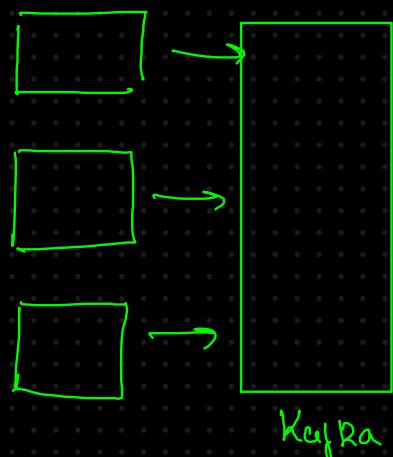
KAFKA101

else enter card details

to earn \$400 credits

Producing messages to a Kafka

1. Kafka Confluent UI
2. Produce data via Python
 - (a) Create a client on Confluent for details
 - (b) Copy the secret key and save it
 - (c) get the code changed to json



Callback, poll and Flush

Function	Purpose	Why is it Needed?	How it Works?
<code>callback in produce()</code>	Handles message delivery status	Kafka produces messages asynchronously, meaning we don't know immediately if the message was successfully delivered or failed. The callback function is used to track message success or failure.	When a message is sent, Kafka does not confirm its delivery immediately. The callback function is triggered when the message is acknowledged (or failed).
<code>poll()</code>	Executes callback functions	Kafka does not call the <code>callback</code> function automatically. Instead, it queues them until <code>poll()</code> is called.	When <code>poll()</code> is executed, Kafka processes all pending callbacks (e.g., checking if messages were successfully sent).
<code>flush()</code>	Ensures all messages are sent before shutting down the producer	Since Kafka sends messages asynchronously, messages might still be in the queue when the script ends. If we don't flush, unsent messages could be lost.	It waits for all in-flight messages to be delivered before shutting down.

Step	Function Used	What It Does?	What Happens If We Don't Use It?
1. Send Message	<code>producer.produce(topic, key, value, callback=delivery_report)</code>	Sends message asynchronously to Kafka	The message will be sent, but we won't know if it was successfully delivered.
2. Trigger Callbacks	<code>producer.poll(1)</code>	Ensures the callback function is executed to check delivery status	The callback function will never execute, and we won't know if the message was delivered or failed.
3. Finalize Sending	<code>producer.flush()</code>	Ensures all messages are sent before shutting down the producer	Unsent messages may be lost when the script exits.

Consuming messages from Kafka cluster

Python

Do make sure to install confluent before running

Command	Purpose
confluent login	Log into Confluent Cloud
confluent environment create <name>	Create a new environment
confluent environment use <name>	Set active environment
confluent kafka cluster create <name> --cloud aws --region us-east-1	Create a Kafka cluster
confluent kafka cluster use <cluster-id>	Select active Kafka cluster
confluent kafka topic create <topic-name> --partitions <num>	Create a topic with partitions
confluent kafka topic list	List all topics
confluent kafka topic describe <topic-name>	Describe topic details
confluent kafka topic produce <topic-name>	Send messages to Kafka topic
confluent kafka topic consume <topic-name> --from-beginning	Read messages from Kafka topic
confluent kafka topic delete <topic-name>	Delete a Kafka topic

Set the key will be required for the first time

