**Inspire…Educate…Transform.**

# Kafka Activity

**Ref: https://kafka.apache.org/intro**

# Why Kafka

- In the real world data exists on many systems in parallel, all of which need to interact with Hadoop and with each other.
- The situation quickly becomes more complex, ending with a system where multiple data systems are talking to one another over many channels.
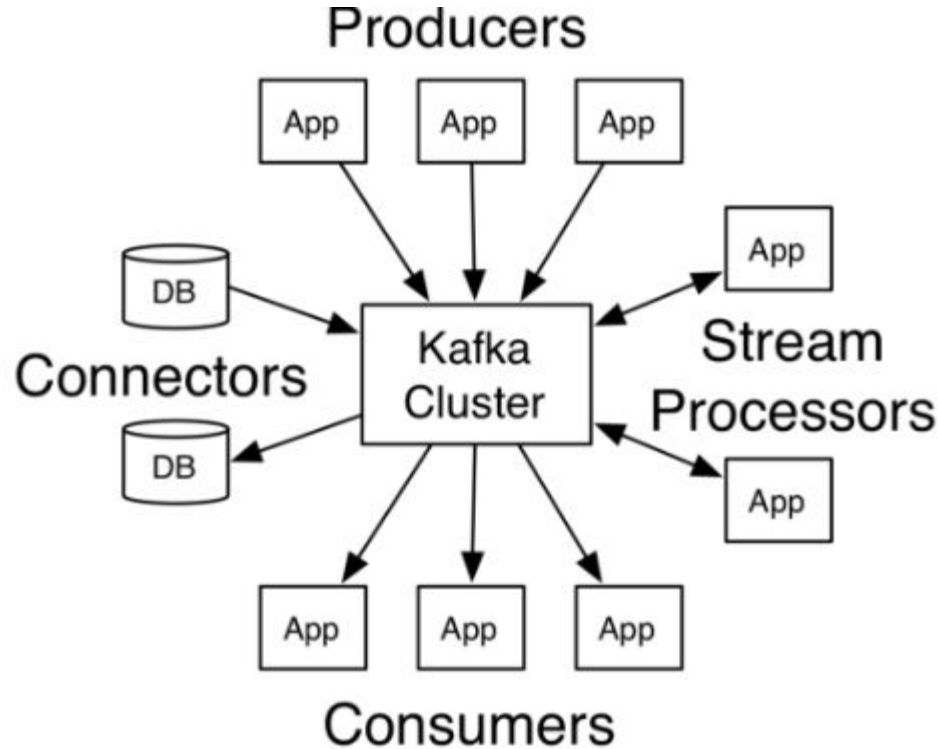
# Contd..

- Each of these channels requires separate way and communication methods and moving data between these systems becomes a full-time job for a team of developers.
- Solution for this is a single component to act like a broker which will serves the purpose of taking the data from different sources and sending them to the different sources.

# Think of it

# Kafka

A streaming platform has three key capabilities:
- Publish and subscribe to streams of records, similar to a message queue or enterprise messaging system.
- Store streams of records in a fault-tolerant durable way.
- Process streams of records as they occur.

# Contd..

- It is a distributed publish-subscribe messaging system and a robust queue that can handle a high volume of data and enables you to pass messages from one end-point to another.

- Kafka is suitable for both offline and online message consumption.

- Kafka messages are persisted on the disk and replicated within the cluster to prevent data loss

## Contd..

- Kafka is built on top of the ZooKeeper synchronization service.

- It integrates very well with Apache Storm and Spark for real-time streaming data analysis.

# Applications

- Kafka is generally used for two broad classes of applications:
  - Building real-time streaming data pipelines that reliably get data between systems or applications
  - Building real-time streaming applications that transform or react to the streams of data

# Contd..

- Kafka is run as a cluster on one or more servers that can span multiple datacenters.
- The Kafka cluster stores streams of *records* in categories called *topics*.
- Each record consists of a key, a value, and a timestamp.

# Terminologies in Kafka

## Topics

- A stream of messages belonging to a particular category is called a topic. Data is stored in topics.
- Topics are split into partitions. For each topic, Kafka keeps a minimum of one partition.
- Each such partition contains messages in an immutable ordered sequence that is continuously append to .

# Contd..

## Partition

- Topics may have many partitions, so it can handle an arbitrary amount of data.

- Partition is like a Single log.

- Messages are written to it in an append-only fashion, and are read in order from beginning to end.
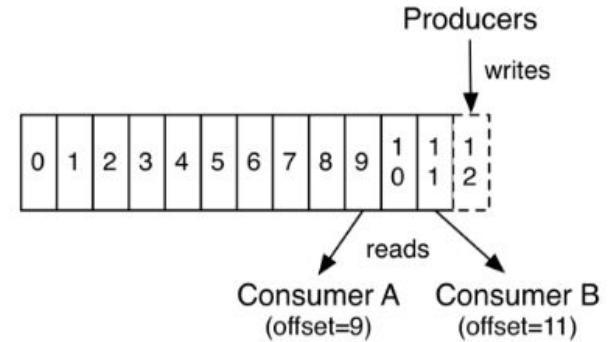
## Partition offset

- The records in the partitions are each assigned a sequential id number called the *offset* that uniquely identifies each record within the partition.

# Contd..

- Metadata retained on a per-consumer basis is the offset or position of that consumer in the log.
- This offset is controlled by the consumer
- Normally a consumer will advance its offset linearly as it reads records, but, in fact,
- Since the position is controlled by the consumer it can consume records in any order it likes.



For example: a consumer can reset to an older offset to reprocess data from the past or skip ahead to the most recent record and start consuming from "now".

# Contd..

**Replicas of partition**

- Each partition is replicated across a configurable number of servers for fault tolerance.

- The Kafka cluster durably persists all published records—whether or not they have been consumed—using a configurable retention period.

**Leader**

- Each partition has one server which acts as the "leader" and zero or more servers which act as "followers". The leader handles all read and write requests for the partition while the followers passively replicate the leader.If the leader fails one of the partition will be selected as a leader by zookeeper.

# Producers

- Producers publish data to the topics of their choice.
- The producer is responsible for choosing which record to assign to which partition within the topic.
- This can be done in a round-robin fashion simply to balance load or it can be done according to some semantic partition function
- Every time a producer publishes a message to a broker, the broker simply appends the message to the last segment file.
- Actually, the message will be appended to a partition. Producer can also send messages to a partition of their choice.

# Contd..

**Brokers**

- Brokers are simple system responsible for maintaining the published data. Each broker may have zero or more partitions per topic.
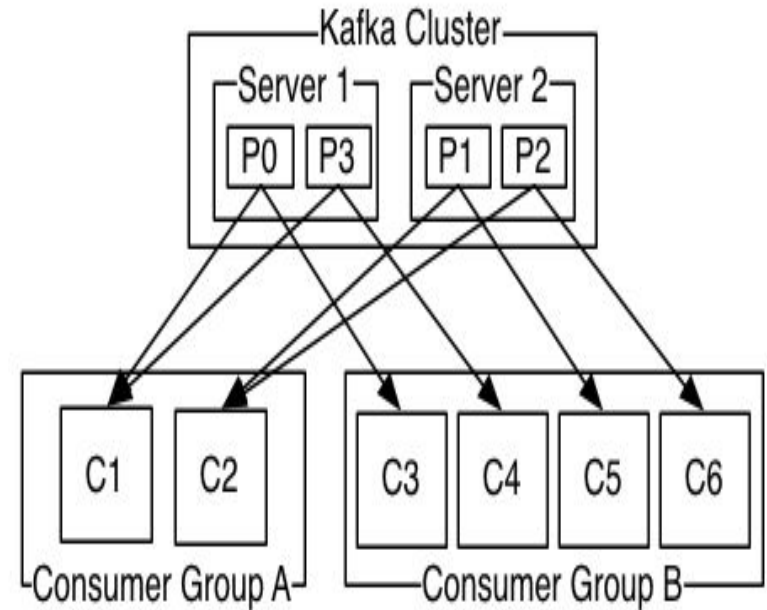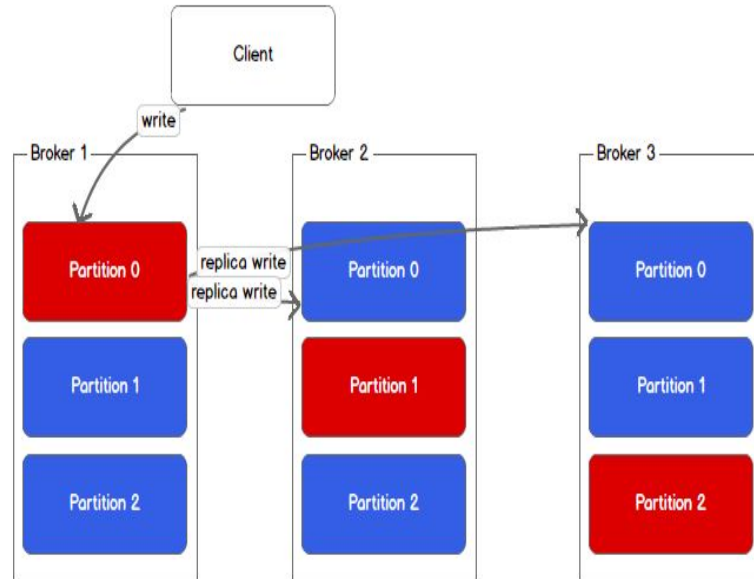
**Kafka Cluster**

- Kafka's having more than one broker are called as Kafka cluster. These clusters are used to manage the persistence and replication of message data.

# Read and Write



Leader (red) and replicas (blue)

# Zookeeper

- For the purpose of managing and coordinating, Kafka broker uses **ZooKeeper**.
- Also, uses it to notify producer and consumer about the presence of any new broker in the Kafka system or failure of the broker in the Kafka system.
- As soon as Zookeeper send the notification regarding presence or failure of the broker then producer and consumer, take the decision and starts coordinating their task with some other broker

# Zookeeper and Broker nodes

Zookeeper Nodes:

a.insofe.edu.in:2181

b.insofe.edu.in:2181

c.insofe.edu.in:2181

e.insofe.edu.in:2181

Broker ID:

c.insofe.edu.in:9092

# Kafka Topic Creation

*export PATH=$PATH:/usr/hdp/current/kafka-broker/bin*

<span style="color:red">*## list the topics on the zookeper*</span>

*kafka-topics.sh --list --zookeeper <zookeeper-ip>*

<span style="color:red">*## Create a topic using the below command*</span>

<span style="color:red">*## Topic Name is insofe_<batch >_<topic name>*</span>

*kafka-topics.sh --create --zookeeper <zookeper-ip>*
*--replication-factor 1 --partitions 1 --topic <topic_name>*

# Producer

*kafka-console-producer.sh --broker-list <broker-id> --topic <topic_name>*

# Consumer

kafka-console-consumer.sh --zookeeper <zookeeper-ip> --topic <topic_name> --from-beginning