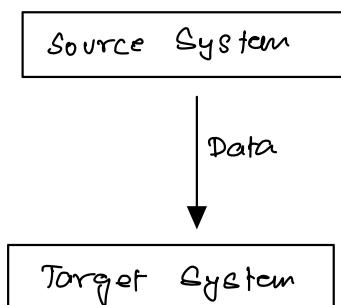
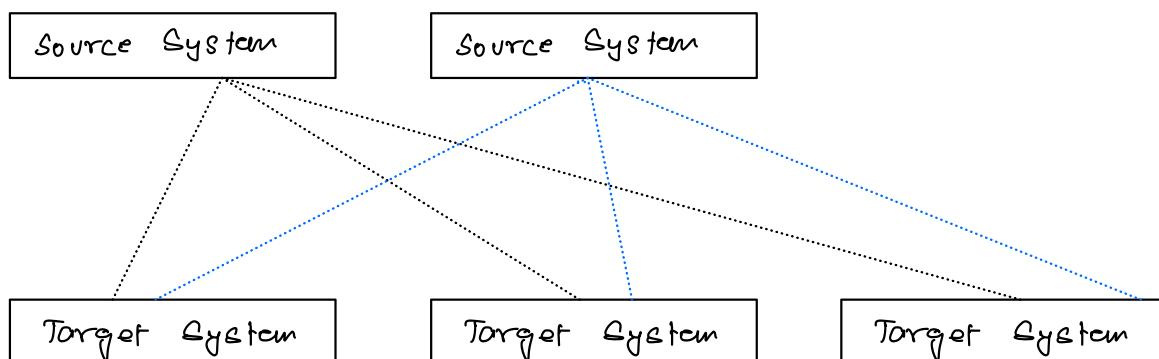


Now maintained by **Kafka** → Created by LinkedIn
Confluent, IBM ←
Cloudflare

At the start of the company we have just
one source system and one target system



As the company grows in size the number
of source systems and the number of
target systems increases.



As the number of source systems and the
target systems increases it becomes difficult
as the source should now communicate with
all the target systems. Also while sending

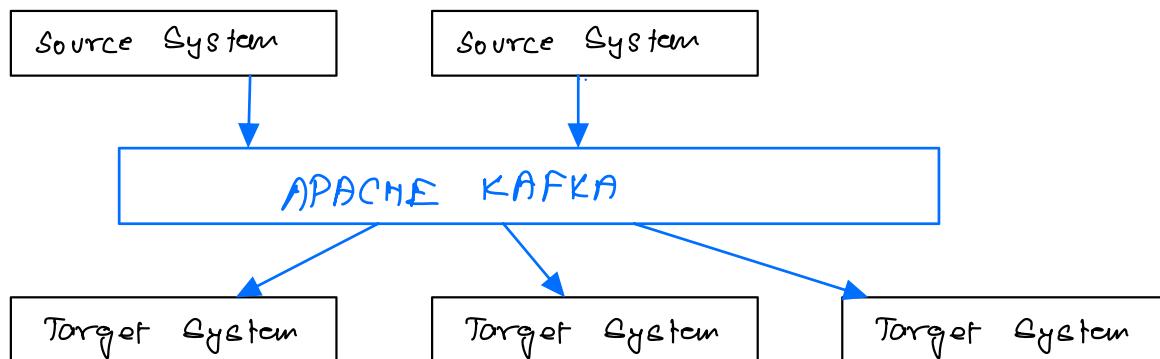
data to the target systems the source systems should take care of the target systems

① Protocol (Eg: TCP, HTTP, REST)

② Data Format (Eg: Binary, CSV, JSON)

③ Data Schema & Evolution

Since each source system will have increased load. This is where Apache Kafka comes.



Now the source systems and the target systems are decoupled

Kafka Use Cases

Use cases: Website Events, Billing Data, Financial Transactions, User Interactions

Targets : Database , Analytics - Email Systems ,
Audit

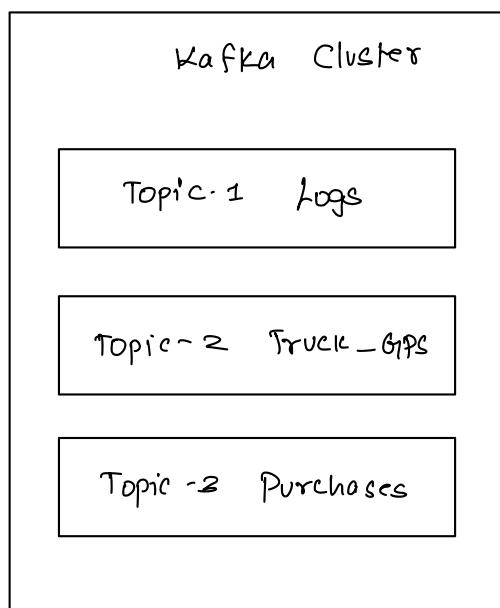
- 1) Messaging System
- 2) Activity Tracking
- 3) Gather metrics from different locations
- 4) Application logs gathering
- 5) De-coupling of system dependencies
- 6) Micro Services Pub/Sub
- 7) Stream Processing → Ingesting a cont. stream
of data to analyze, filter,
transform data in real
time

Advantages of Kafka :

- * Kafka can be horizontally scaled to 100's of nodes
- * can be scaled to millions of messages per second
- * can be scaled to latency less than 10ms
Thus they can be used in real time systems

Kafka Topics

Kafka Topics are a particular stream of data. We can have as many topics as we want in a cluster. A Kafka topic is similar to a table in a database.



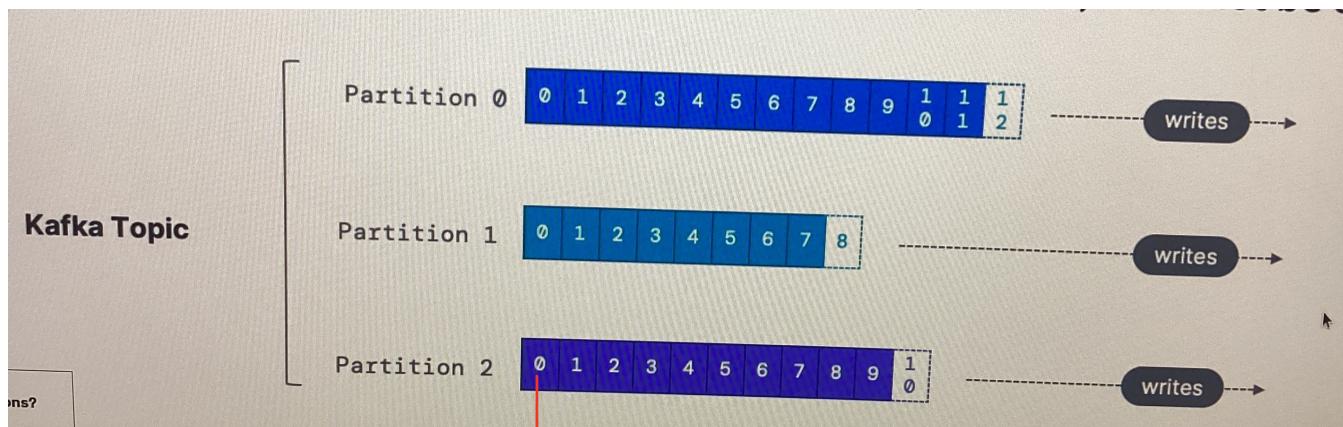
A topic is identified by its name. It can store any kind of data (JSON, CSV, Binary). Sequence of messages is called data stream. Topics cannot be queried and they are also immutable i.e. data added to Kafka cannot be updated / deleted.

Data in a Topic is stored for a limited time (Default: One week - configurable). After that period the data will be deleted.

Kafka Producers add messages to the topics and the Kafka Consumers receive them.

Partitions and Offsets:

Topics are divided into Partitions



As messages are added to the partitions they are assigned with an ID (starting from 0) known as the Kafka Partition Offset.

Even if the messages are deleted the offsets are not deleted ie if I have 100 messages then I delete them instead of the offsets starting

from 0 they start from 101

The messages sent to a topic are stored
properly within these partitions. Order is
guaranteed only within a partition (not
across partitions)

Producers & Message Keys

- The producers decide in which partition the messages are to be stored

Producers:

Producers know to which partition the messages have to be written and which Kafka Broker (Kafka Server) has it.

In case of Kafka Broker failure.

Producers will automatically recover.

Message Keys:

The keys are hashed to the correct positions using the Member 2 algorithm.

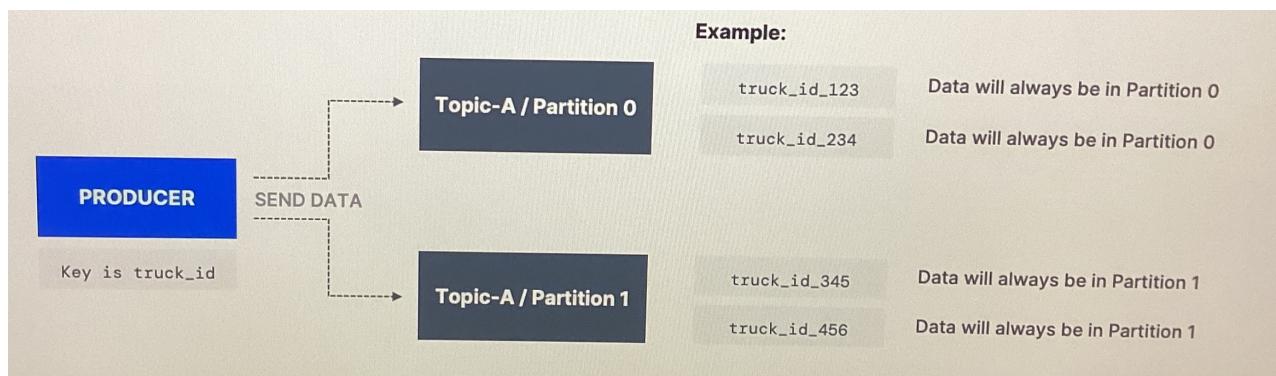
Producers can send a key with the message (optional)

The key can be string, member, binary
When to send a key along with the message?

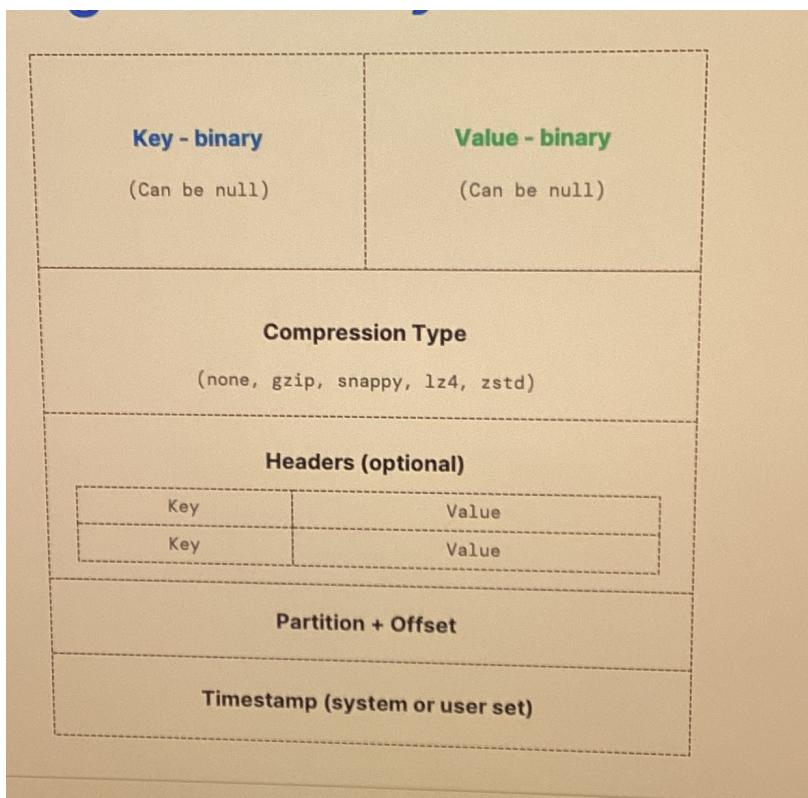
When a key is not sent by the producer then the messages sent will be distributed among the partitions in round robin.

But if a key is sent by the producers then it is guaranteed that the messages will end up in the same partition.

Eg:



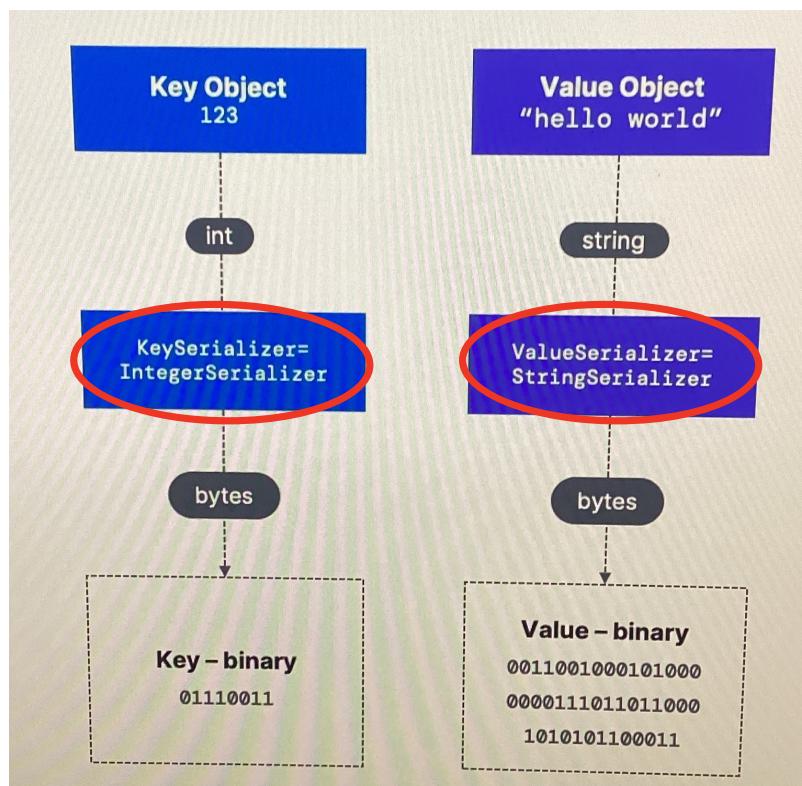
Message format:



Kafka only accepts bytes as an input from producers and sends bytes as output to consumers

Message Serialization

The process of transformation of objects / data into bytes



Common Serializers:

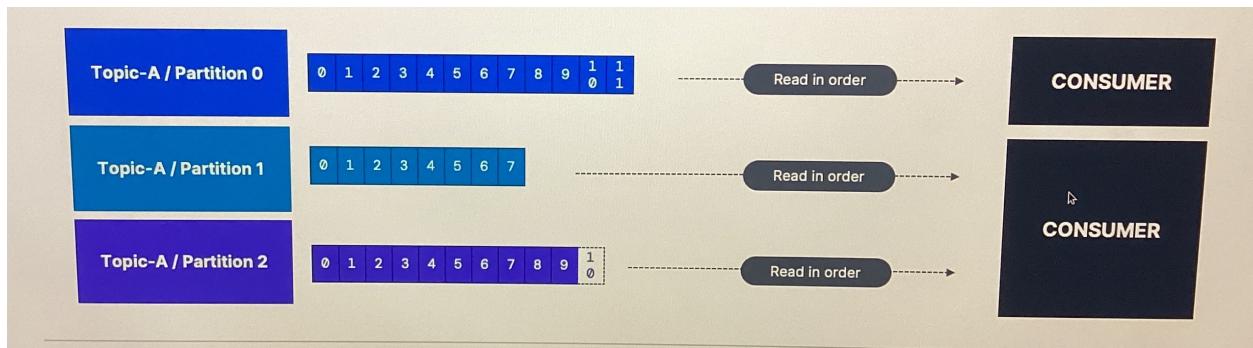
- ① Strings (inc. JSON)
- ② Int, Float
- ③ Avro
- ④ Protobuf

Consumers & Deserialization

Consumer should pull (request) data from the kafka topic. Consumers know to which broker and which partition the request should be sent.

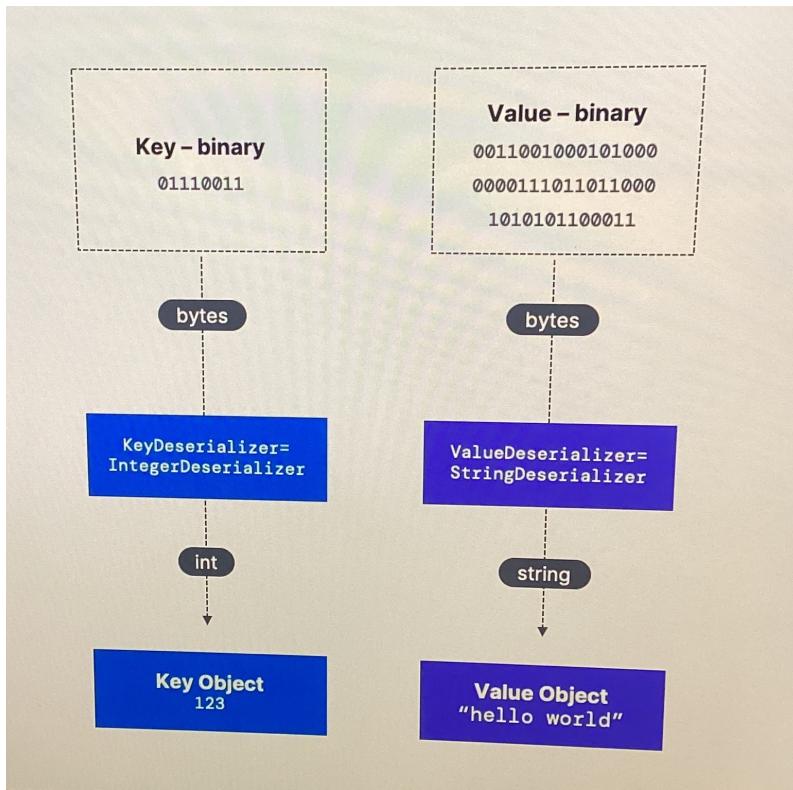
In case of broker failure the consumer how to recover.

Data is read from low to high offset within each partitions



Consumer Deserialization

This is a process of converting bytes into objects / data. The consumer should know in advance what is the format of the message. Eg $\text{key} \rightarrow \text{int}$ $\text{value} \rightarrow \text{string}$



Common Deserializers

- ① Strings (inc. JSON)
- ② Int, Float
- ③ Avro
- ④ ProtocolBuf

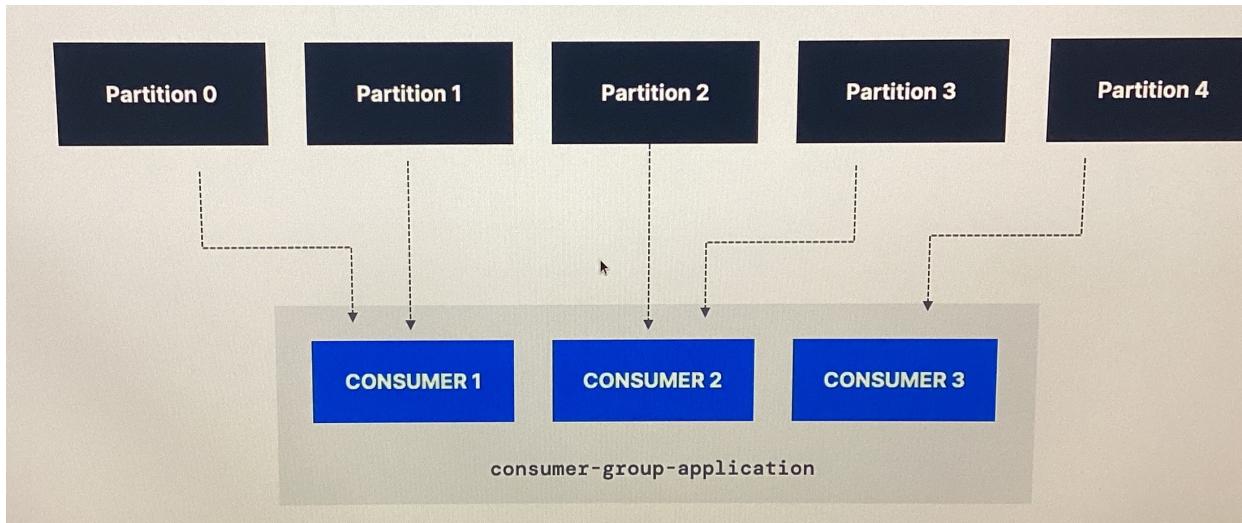
NOTE



The serialization / deserialization type should not change during a topic life cycle (create a new topic instead)

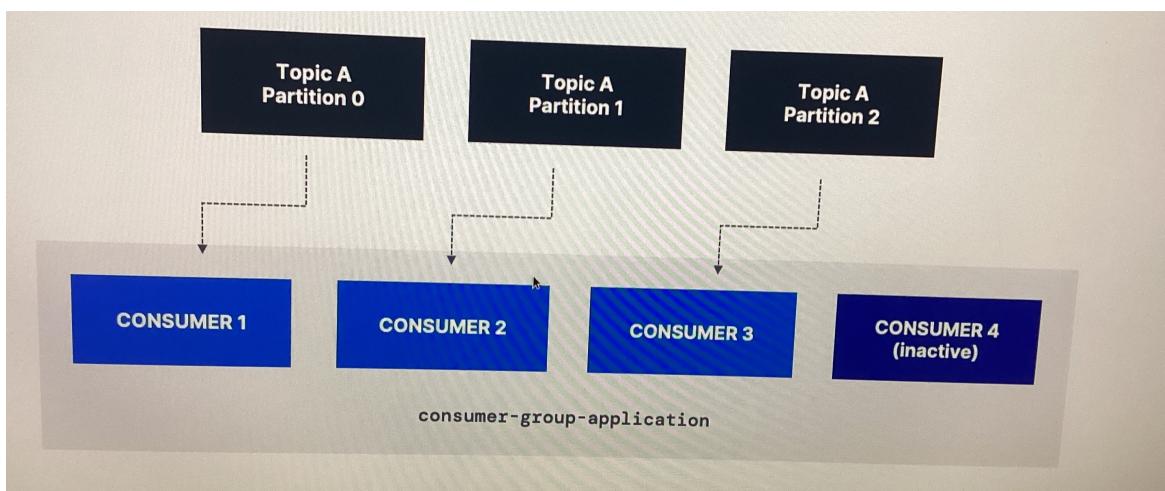
Consumer Group 2 Consumer Offsets

All the consumers in an application read data as consumer groups

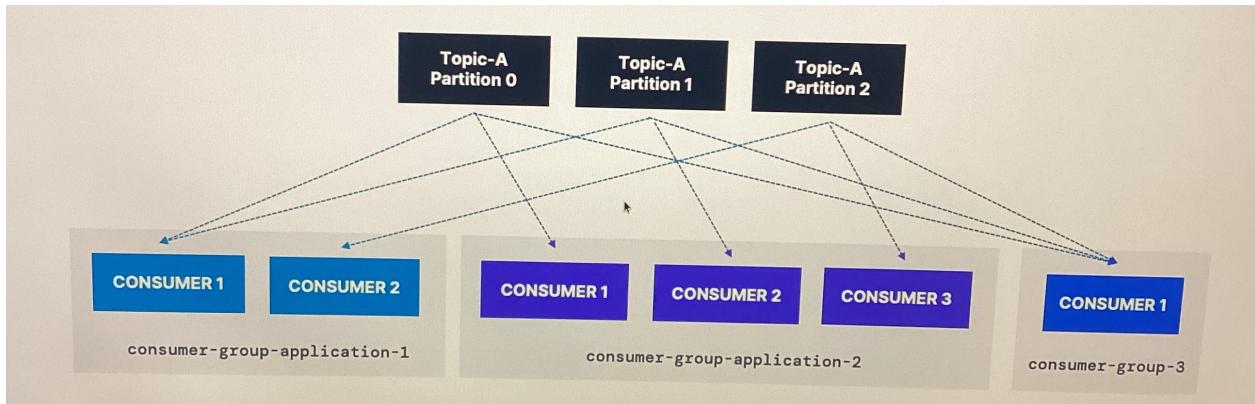


Each consumer within a group reads data from **distinct partitions**

Inactive Consumers



Multiple Consumer Groups



To create distinct consumer groups the consumer group property `group.id` is used.

Consumer Offsets

Kafka stores the offsets at which the consumer group has been reading.

Eg: If a consumer has read a message from 23rd offset in Partition Topic-1 then 23 is stored.

These offsets are stored in a Kafka topic called `-- consumer-offsets`.

The offsets are committed to the `-- consumer-offsets` only once in a while. This ensures that if a consumer dies and comes back, the consumer will be able to read from where it left.

By default your consumer will automatically commit offsets (at least once) if we choose to manually commit then there are 3 semantics

① At Least Once

Offsets are committed after the message is processed. If the processing goes wrong then the message will be read again. This results in message duplication. This should be used when processing of message again won't impact the system)

② At most Once:

Offsets are processed as soon as the messages are received. If the processing goes wrong then some messages will be lost.

③ Exactly Once

Kafka Workflows \Rightarrow Use Transactional API (in the Kafka Streams API)

External System \Rightarrow use an idempotent workflow consumer.

Brokers & Topics

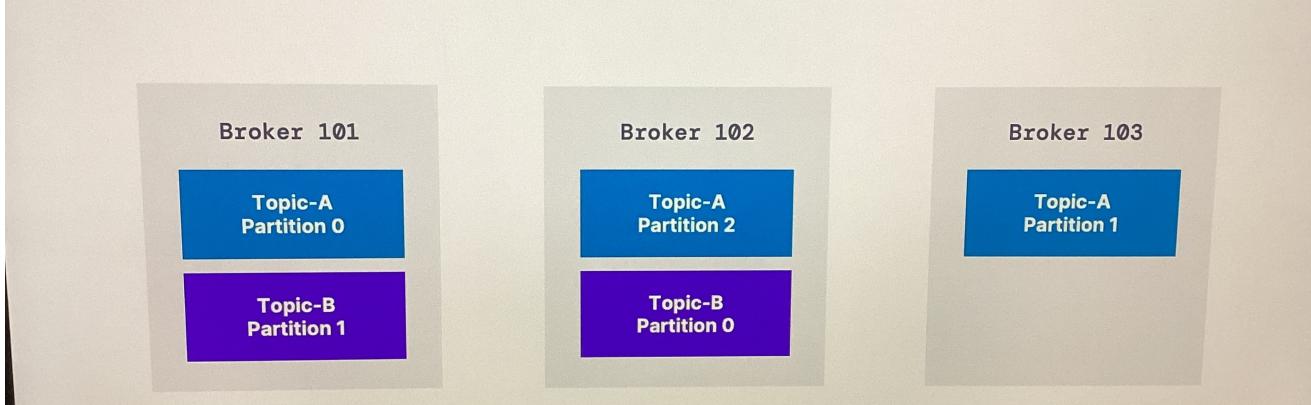
Kafka Cluster:

A kafka cluster is composed of multiple brokers (servers) & good number of brokers to start with is Brokers → 3 A kafka broker is just a server but since they receive and send data they are known as brokers.

A broker is identified with an ID which is an integer.

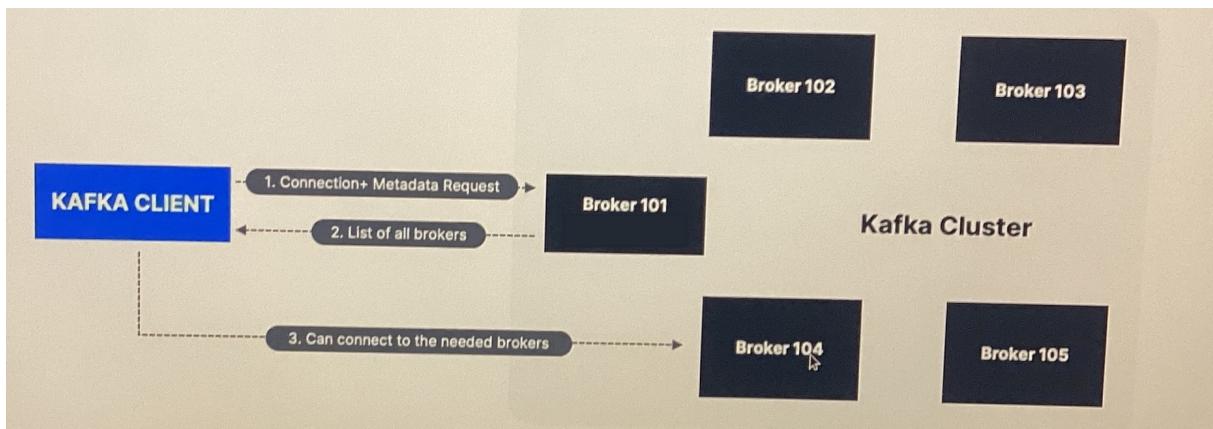
Each broker contains certain topic partitions i.e. the data is distributed among the brokers.

- Example of Topic-A with 3 partitions and Topic-B with 2 partitions



Kafka Broker Discovery

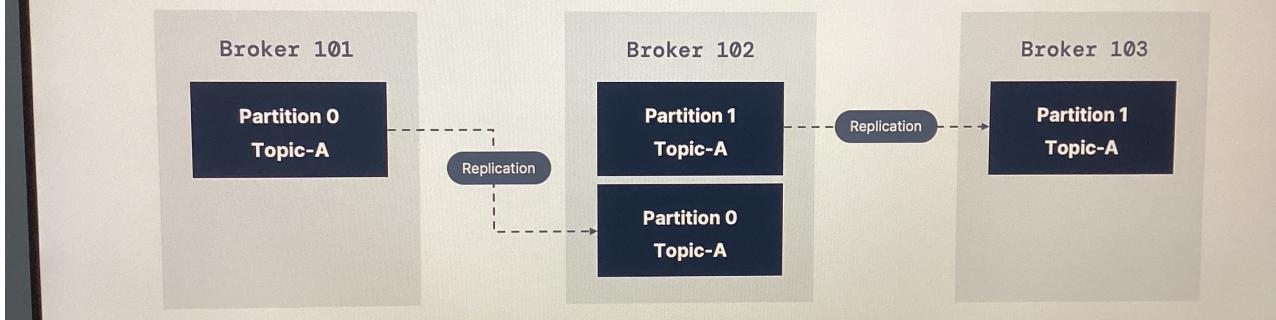
Every Kafka Broker is called a **Bootstrap Server** because if a consumer connects to just **one broker** then the consumer knows how to get connected to all the **brokers** in the cluster.



Topic Replication

Topics should have a replication factor > 1 (For production it should be 3). Thus if a broker is down another broker can serve the data.

- Example: Topic-A with 2 partitions and replication factor of 2

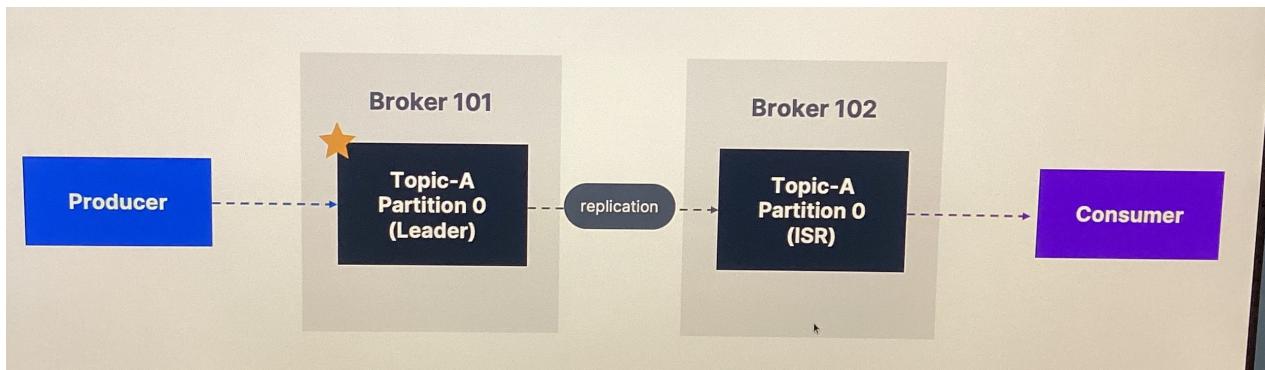


Leader for a Partition

At any time only one broker can be a leader for a given partition and the producers send data ^{only} to the broker that is the leader of the partition and the other brokers replicate the data from the leader.

Similarly the consumer read messages only from the leader till in the latest version (Kafka v2.4+). Of Kafka there is a feature known as the

Kafka Consumer Replica Fetching which allows the consumer to read data from replicas



This might help to reduce latency (the consumer is close to the replica broker)

In Sync Replica (ISR)

The replica brokers send fetch requests to the leader to get the latest messages. If the leader doesn't receive a fetch request within the configured `replica.log.time.max.ms` or the follower has not consumed the message within this period then the replica broker is removed from the ISR list

Producer Acknowledgement

The producer can choose to receive acknowledgement that the data has been stored successfully by the broker. There are 3 settings

$acks = 0 \Rightarrow$ Producer won't wait for data acknowledgement (possible data loss)

$acks = 1 \Rightarrow$ Producer will wait for the broker acknowledgement (limited data loss)

$acks = all \Rightarrow$ Leader + Replicas acknowledgement (no data loss)

Zookeeper

Versions

4.x => Doesn't need Zookeeper

3.x => can work without Zookeeper (KIP500) using
Kafka Raft instead (RAFT)

2.x => Cannot work without Zookeeper.

Zookeeper is slowly disappearing and is going
to be replaced because Zookeeper is less
secure than Kafka and we need to configure
Zookeeper in such a way that it
accepts connection only from Kafka. With
Kafka > v0.10 it doesn't store any consumer offsets
or any other consumer data.

Zookeeper is a software that has a list
of all the Kafka brokers and whenever a
broker goes down Zookeeper performs the
leader election.

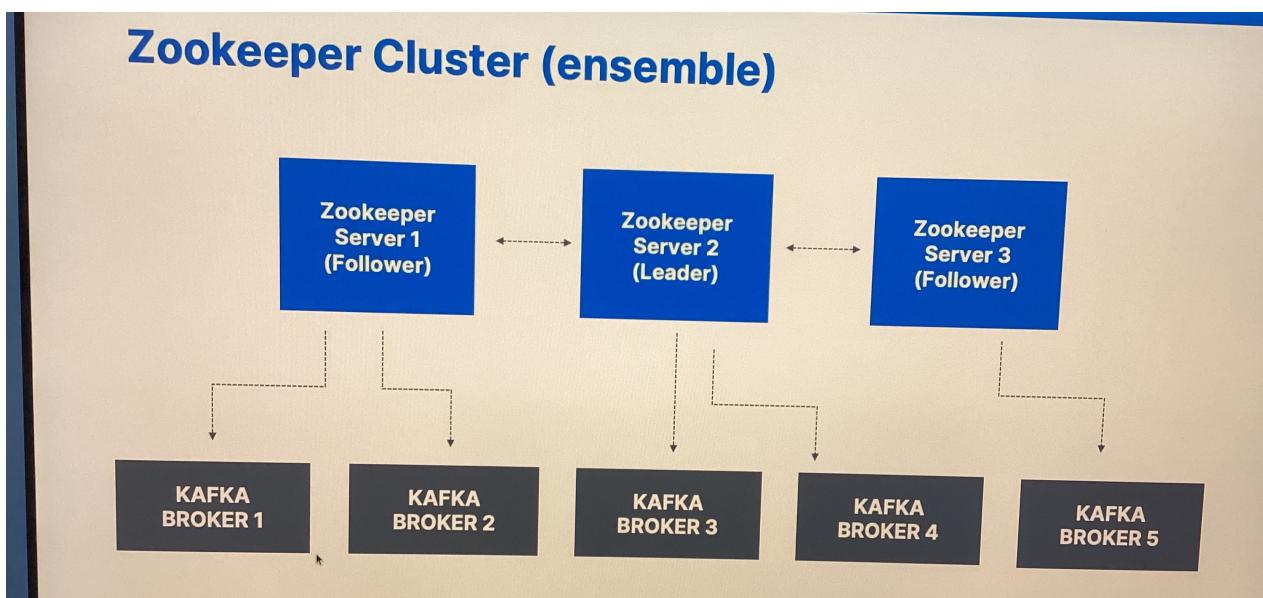
Zookeeper also sends notifications to the
brokers when

- ① A topic is created / deleted
- ② A broker goes down / comes up.

Zookeeper

Zookeeper Cluster

Zookeeper by design works with an odd number of (Zookeeper) servers i.e. 1, 3, 5, 7(max)
 Zookeeper has a leader which events to the rest of the follower servers



Should Zookeeper be used ??

with Kafka brokers?

Yes, until Kafka ID is out

With **Kafka clients**?

Overcome Kafka clients to connect to Kafka brokers directly instead of using the Zookeeper's endpoint. So for Kafka clients **NO**

KAFKA RAFT (or) KRAFT

Zookeeper had issues when Kafka clusters had $> 1,00,000$ partitions.

By removing Zookeeper, Kafka can

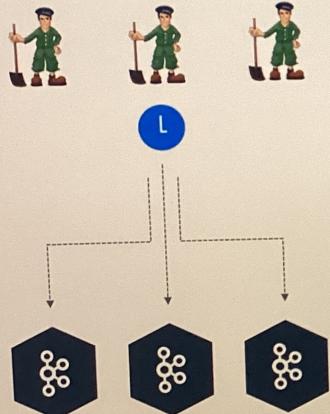
- ① Scale to millions of partitions
- ② Single security model for the whole system (Before we had to take care of zookeeper)
- ③ Single process to start and faster controller shutdown and recovery time

So Kafka 3.x implemented the KRAFT protocol to replace Zookeeper (production ready since 3.3.1) KIP 833

In Zookeeper we have the Zookeeper cluster in which the leader will control all the workers but in KRAFT we have Quorum controller where we have just the Kafka workers among which one will be the Quorum leader

Kafka KRaft Architecture

With Zookeeper



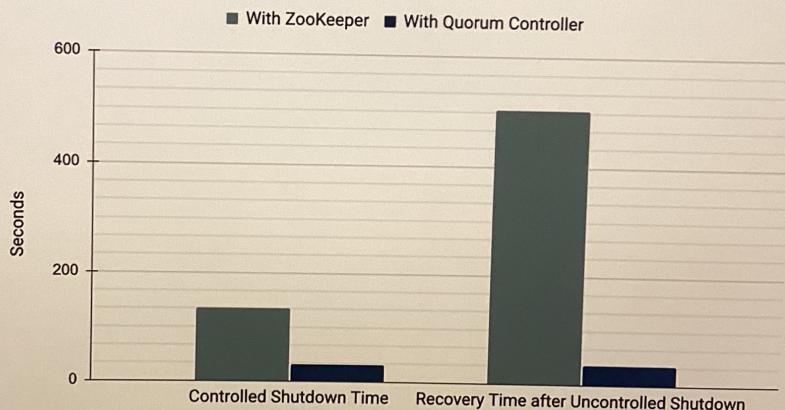
With Quorum Controller



L Quorum Leader

KRaft Performance Improvements

Timed Shutdown Operations In Apache Kafka with 2 Million Partitions
Faster is better



<https://www.confluent.io/blog/kafka-without-zookeeper-a-sneak-peek/>