

# Apache KAFKA (theory + imagination)

## # what is Apache Kafka?

Imagine : Your system is not one app. Its many services talking to each other - : orders, payments, emails, analytics, notifications etc

If they call each other directly, the system becomes slow, tightly coupled and fragile.

Kafka Solves this :- It is a high speed data pipeline that lets system send and receive data in real time without depending on each other. Imagine it as - :

App A → Kafka → App B

Apps dont talk directly, they drop a message in Kafka, others pick them up when ready.

Imagine Kafka like a post office for Data.

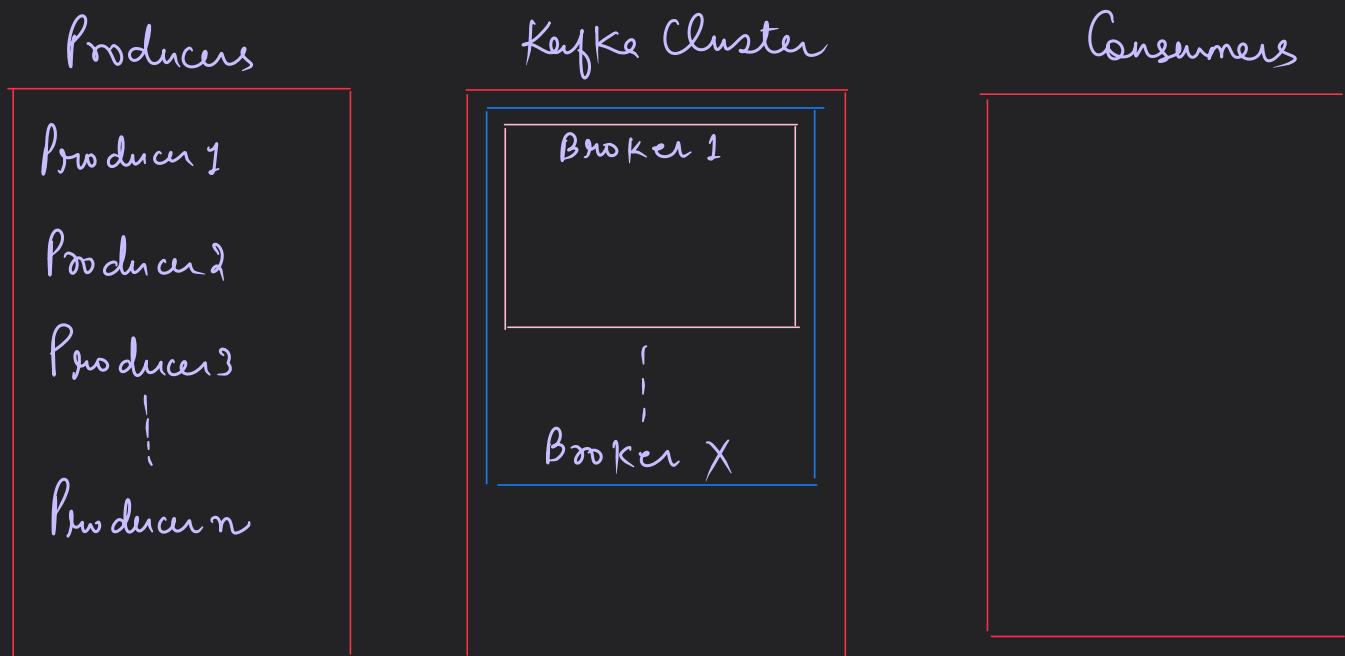
Producer → Person sending letter

Kafka → Post Office

Topic → Mailbox category

Consumer → Person Receiving letter

# X — Kafka Architecture — X



# think of Kafka as a huge warehouse of boxes. Apps don't talk to each other. They drop boxes in warehouse. Other apps come & pick boxes.

# Kafka Cluster :- Imagine it as a group of machines (servers)

working together. Each machine is called a Broker.

Kafka Cluster = many Brokers (servers)

# Topic :- Inside a warehouse we don't mix everything. We create separate sections Eg.

Section name what goes inside?

Orders Order messages

Payments Payment messages

Emails Email events

Topic ≈ Category of message

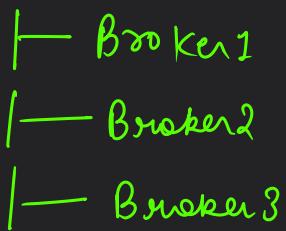
# Partitions Imagine the order section is huge and just one shelf is not enough. So we split it into multiple shelves. these shelves  $\approx$  Partitions

Topic : Orders

- shelf 0 (Partition 0)
- shelf 1 (Partition 1)
- shelf 2 (Partition 2)

$\Rightarrow$  topics are big category and partitions are smaller storage unit inside it.

Kafka Cluster



Topics : Exists at cluster level

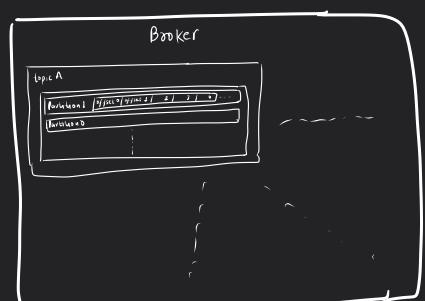
- |— Orders
- |— Payments
- |— Emails

Each topic  $\rightarrow$  split into partitions

Partitions  $\rightarrow$  stored across brokers.

Broker1  $\rightarrow$  orders partition 0  
Broker2  $\rightarrow$  orders partition 1  
Broker3  $\rightarrow$  orders partition 2  
⋮

Brokers do not own any topic. they store some partitions of many topics.



Remember A topic is split into partitions and those partitions are stored across different brokers in a kafka cluster.

# A partition number is unique inside a topic. but copies(replica) of that partition can live on multiple brokers.

If you see

Broker 1 → Topic A Partition 0

Broker 2 → Topic A → Partition 1

this is correct, but one is leader and other is Replica

# Replication Kafka doesn't keep only 1 copy of data (dangerous)

It keeps Backups. This is called Replication factor.

If Topic A

Partitions = 2 (0 and 1)

then replication factor = 2.

means Each partition has 2 copies on different brokers.

Topic A (2 partitions, RF=2)

Partition	Leader Broker	Replica Broker
0	Broker 1	Broker 2
1	Broker 2	Broker 1

So Physically,

Broker 1 stores -:

Topic A → Partition 0 (Leader)

Topic A → Partition 1 (Replica)

Leader :- handles reads & writes

Replica :- just keeps copy of data

Broker 2 stores -:

Topic A → Partition 1 (Leader)

Topic A → Partition 0 (Replica)

# # When and how is Replica Actually get used?

- ①  $\Rightarrow$  Suppose Broker 1 (leader) Crashes
- $\rightarrow$  Kafka automatically
- $\Rightarrow$  detects leader is dead
- $\Rightarrow$  chooses one replica (say broker 2)
- $\Rightarrow$  makes it a new leader.
- $\Rightarrow$  This is called leader election
- $\Rightarrow$  Replica becomes active only during failure
- fault Tolerance  
in  
Kafka
- ②  $\Rightarrow$  Case 2: Continuous Data Backup
- $\Rightarrow$  Even when leader is alive,
- $\rightarrow$  Replicas continuously copy new message.
- $\rightarrow$  They stay in sync.
- $\rightarrow$  this is called ISR (In-Sync-Replicas)
- $\rightarrow$  these replicas are fully up to date.

- ③  $\Rightarrow$  Case 3 : Prevention of Data loss

$\rightarrow$  If Replication factor is 3,

You can loose 2 brokers and still not loose data (if atleast 1

ISR Survives)

think it of  
leader  $\rightarrow$  main driver  
replica  $\rightarrow$  backup driver in case  
the backup don't drive until leader finds

- $\star\star$  Kafka replicas are passive followers that keep copying the data from the leader and become active only if the leader fails.
- $\Rightarrow$  Kafka only picks in-sync replicas to become leader. If replica was slow/outdated  $\Rightarrow$  not chosen

## # How does Kafka know a particular replica is outdated?

Kafka compares message positions (offset) between leader & replica.

## # offset's — Every message in a partition gets a number.

message 1 → offset 0  
message 2 → offset 1  
message 3 → offset 2 } offset = message position number

→ leader partition is like a main book. whenever producer sends data, leader writes message in book & tells replicas, to copy the same message! -

→ Replica keeps copying messages from leader in order. It also keeps track: my last copied offset is 105.

→ Kafka checks if Replica is healthy by checking offset lag, or heartbeat / fetch requests.

↓

is replica actively syncing?

show for the  
replica is

### Case 1 Replica is healthy

leader offset = 200  
Replica offset = 200 } Replica is in-sync  $\Rightarrow$  It will be part of ISR (in Sync Replica)

### Case 2 Replica is slow

leader offset = 200  
Replica offset = 180 } Replica is 20 msg behind. Kafka waits for a configured time (replica.lag.time.max.ms)  
If it doesn't catch up  $\rightarrow$  Remove from ISR

Case 3

## Replica Steps talking

### of Replicar

- stops sending fetch request
- or network breaks
- Leader says: - this guy is dead or too slow to be a leader
- Removed from ISR.

# What happens when a producer sends a message - inside KAFKA?

# Lets follow 1 message from Producer → Kafka Cluster

Inside Kafka Cluster we have -

- \* multiple brokers (servers)
- \* Metadata manager → Zookeeper for older, KRAFT for new Kafka version
- \* Topics
- \* Partitions

Step 0 - Topic Creation (happens before messages)

(we / system has to create topics)

→ topic: orders

partitions 3

Replication factor: 2

→ 1 leader, 1 replica.

What happens internally?

→ metadata system chooses brokers

→ assigns partitions

these are created during topic creation.

Eg.

Partition	Leader	Replica
0	Broker 1	Broker 2
1	Broker 2	Broker 3
2	Broker 3	Broker 1

⇒ this mapping is stored in cluster metadata.

→ Brokers don't decide at runtime where message goes.  
Partition-to-broker mapping is already decided when topic is created.

### Now producer Sends Message

Step 1 :- Producer asks cluster for metadata.

Producer says, where is topic Orders?

Cluster replies →

Partition 0 → Broker 1 (leader)  
1 → Broker 2 (leader)  
2 → Broker 3 (leader)

Producer cache this.

Step 2 :- Producer chooses Partition based on (a) key provided

(b) No Key.

→ round robin.

E.g. Key = Customer ID 55

hash → Partition 2

Step 3 Producer Sends Message to Leader Broker

Producer directly contacts Broker 3 (leader of partition 2)

Kafka Cluster is not a load balancer. Producer knows exactly where to send.

Step 4 Leader broker stores messages.

→ partition is just a log file on disk

Leader → opens partition log files → appends message at end → assigns offset

Before write : last offset = 100

After write : last offset = 101

offset : simply next number in file.

Step 5 Replication Happens

Leader tells replicas - fetch new data from me.

Replicas copy message & reach offset (101) → now they are ISR  
(in Sync Replica)

Step 6 Acknowledgement to Producer

depending on acks :

Setting 1 when producer gets success

1 leader wrote message.

all leader + all ISR wrote

## # How are offsets created?

Offset is not random.

for each partition, offset = position of message in append-only-log

Each partition has its own offset sequence.

Partition 0 offset 10 ≠ Partition 1 offset 10

## # How consumers read using offsets and consumer groups?

Consumer reads message from a topic. But Kafka doesn't push messages. Consumers pull messages.

## # What is an offset for Consumers?

As we already know producers create offset when writing  
Eg. (Partition 0)

Offset	Message
0	Order A
1	Order B
2	Order C
3	Order D

Consumes keep track :-

→ I have read until offset 1, so next read starts from offset 2.

→ Kafka stores consumer progress in an internal topic :-  
- Consumer\_offsets

So, Kafka remembers reading position even if app starts.

Part 2 Consumer Group :- team of consumers working together

topic :- Orders (3 partitions)

Confluent Cloud  
Book's Family Child

Case 1 → 1 Consumer

Consumer A → Reads P<sub>0</sub> → Reads P<sub>1</sub> → Reads P<sub>2</sub>

Case 2 → 3 Consumers (same group)

Consumer A → Partition 0

Consumer B → Partition 1

Consumer C → Partition 2

Kafka divides work :- One partition → One consumer at a time (within a group)

Case 3 → More Consumers than partitions.

topic partitions = 3

Consumers = 5

⇒ 2 consumers will sit ideal.

~~Part 3~~ How consumers read step by step.

(a) Consumer joins a group

↓

Kafka assigns partitions

↓

Consumer sends fetch requests

↓

Broker returns messages starting from stored offset

↓

Consumer processes message

↓

Consumer commits new offset

→

Auto Commit

Manual Commit

~~Part 4~~ Rebalancing

when Consumer Joins / Consumer Leaves / Broker fails, Kafka re-assigns partitions → this is automatic

Producer tracks where to write (partitions)

Consumer tracks where to start read (offset)

