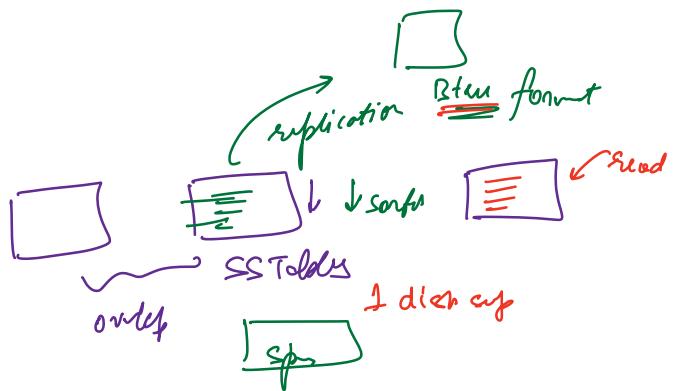


Agenda

- Problem 1
- Zookeeper
- Problem 2
- Kafka

Starting at 9:05

LSM tree



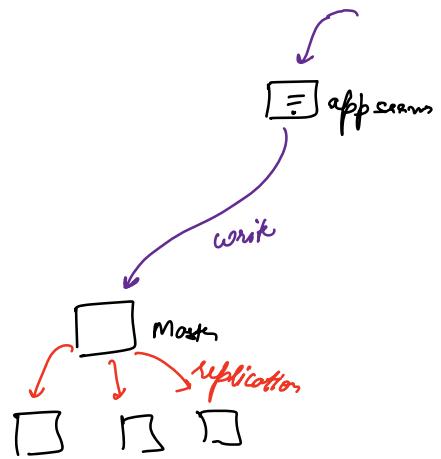
Problem -1

Consistent State Tracking

Master-Slave architecture

↳ all writes go to → master

↳ app server must be aware of who the master is
(all app servers agree on the same master)

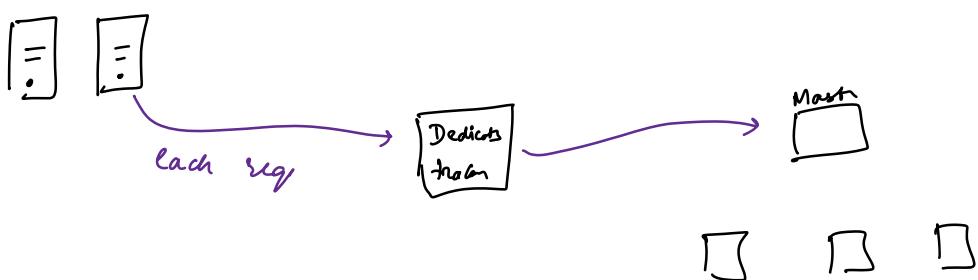


Problem → keep track of the master

↳ when master dies → re-elect the master

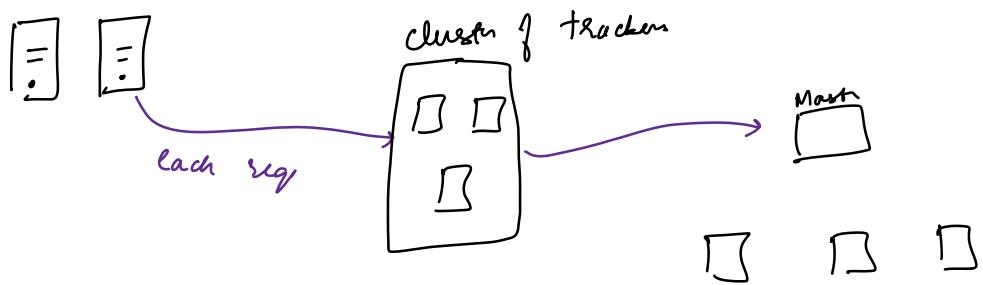
Approach -1

→ dedicated machine → keep track of master



- ① Extra hop for every request
- ② Dedicated Machine → single point of failure

To solve ②



- if master changes, how do they keep track of it
- how to ensure that all machines agree on same master

Zookeeper → Stores configuration in distributed settings

tracks (config) data in a strongly consistent manner.

root folder

/

/ config/

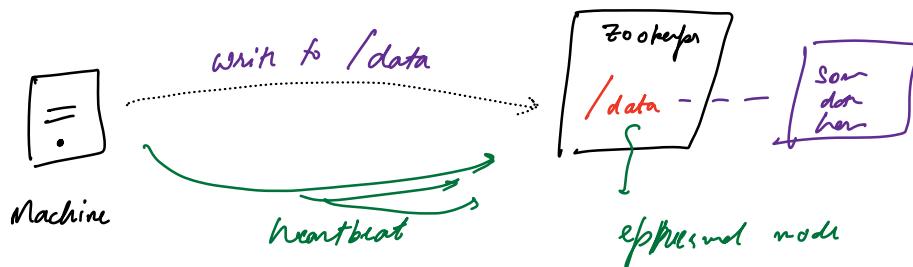
/ config/ aws-key ↗ files are called "nodes"
ZK nodes

/ config / dp-ip in Zookeeper

/ master

Types of ZK Nodes

- ↳ Persistent → write data to it → persists until explicitly deleted
- ↳ **Ephemeral** → data written → **Owner**
this data is only persistent as long as the owner is alive/available

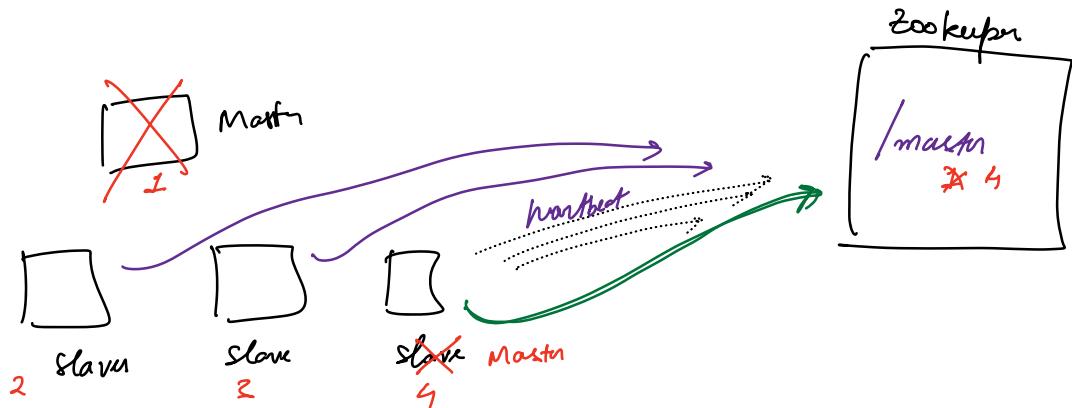


In a master-slave → whoever is the master

↳ can't tell images inside

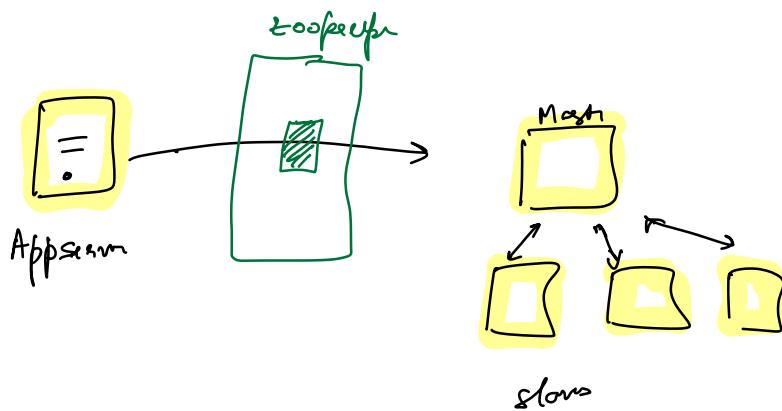
Zookeeper → store their ip in the ephemeral node

Master Election

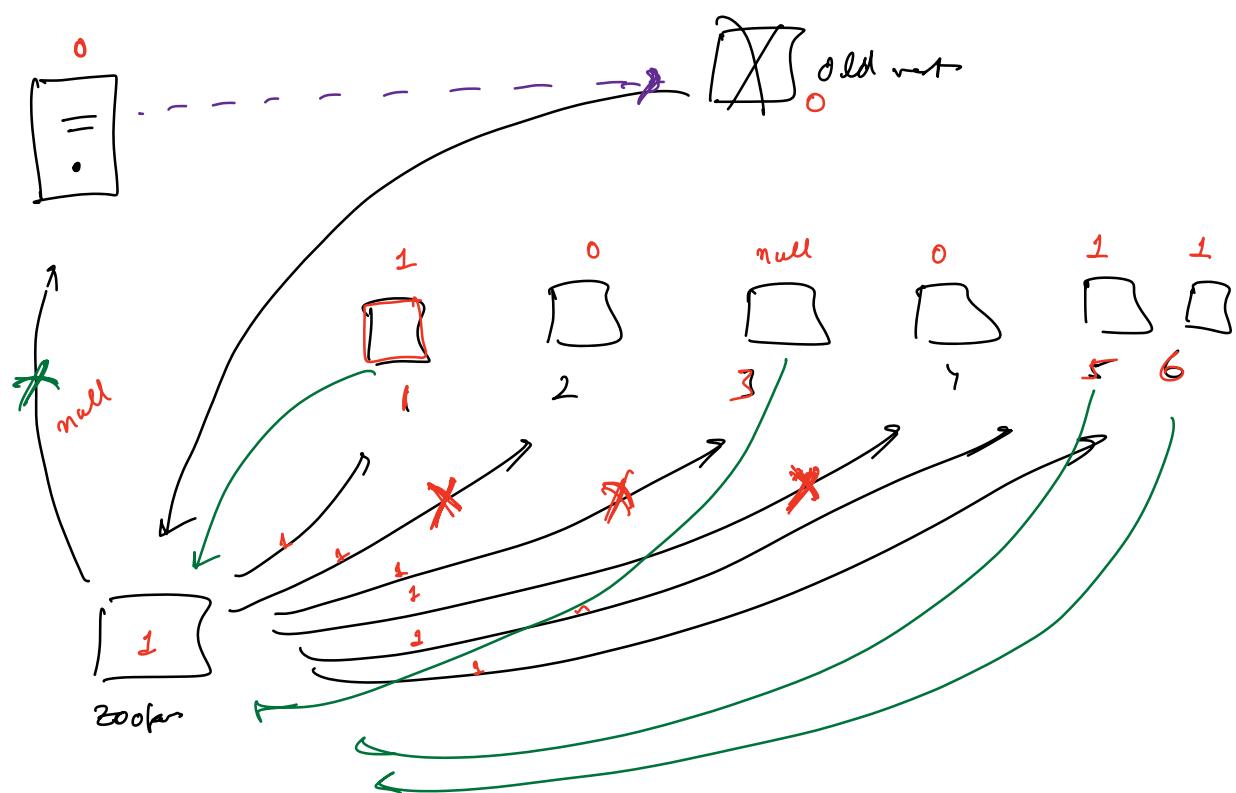
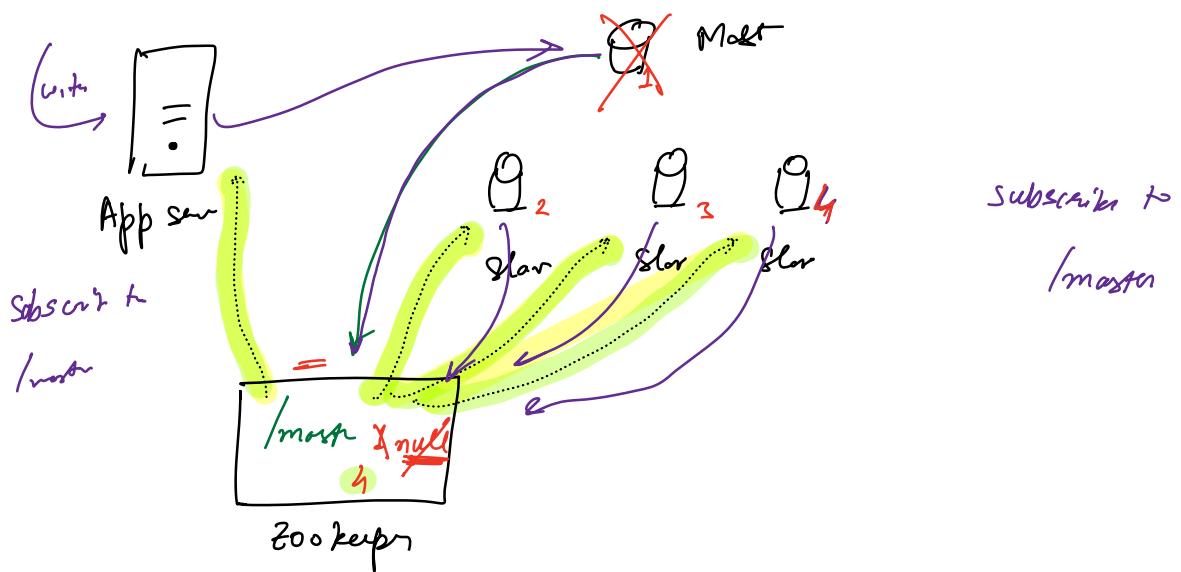


For re-election every slave will try to write its own ip in zookeeper's /master node.

Inside zookeeper → lock is acquired & only 1 writer is allowed to pass through



Set a watch / subscribe to any node

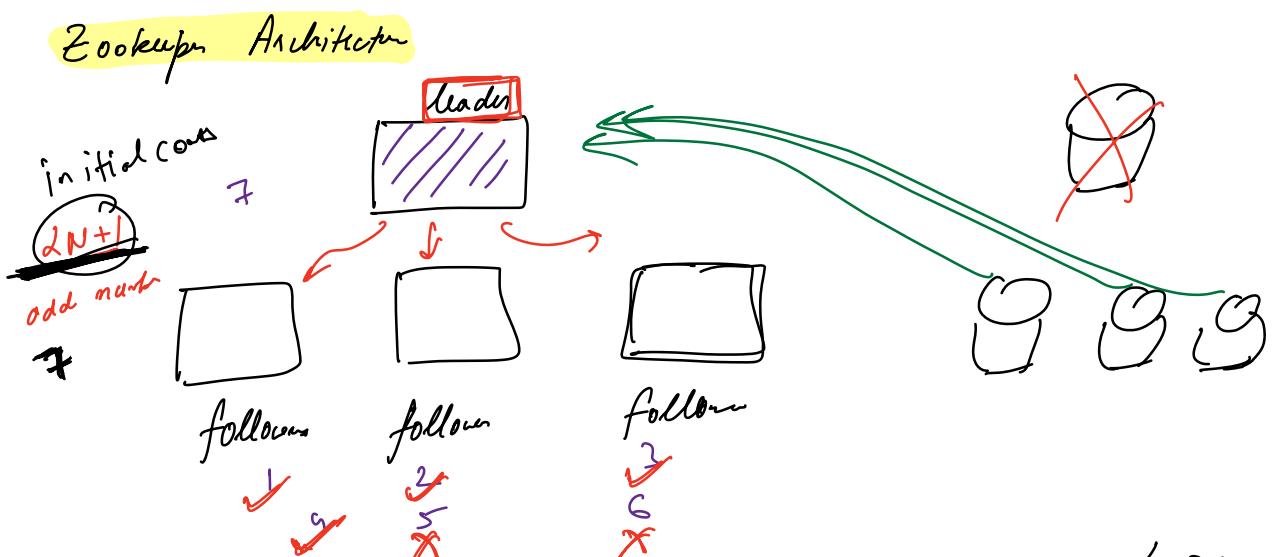


zookeeper must notify all subscribers → ~~1 way channel.~~

2 way connection.

the moment any app slave detects that it is no longer connected to Zookeeper or if it gets a notification from Zookeeper saying that master is null → Stop serving write requests.

Extra hop → solved
 Subscriptions to ZK nodes
 ↳ 2 way connection



When a client comes to ZK's leader,

→ it writes to self

→ it also waits for a majority of slaves. of $(2N+1)$

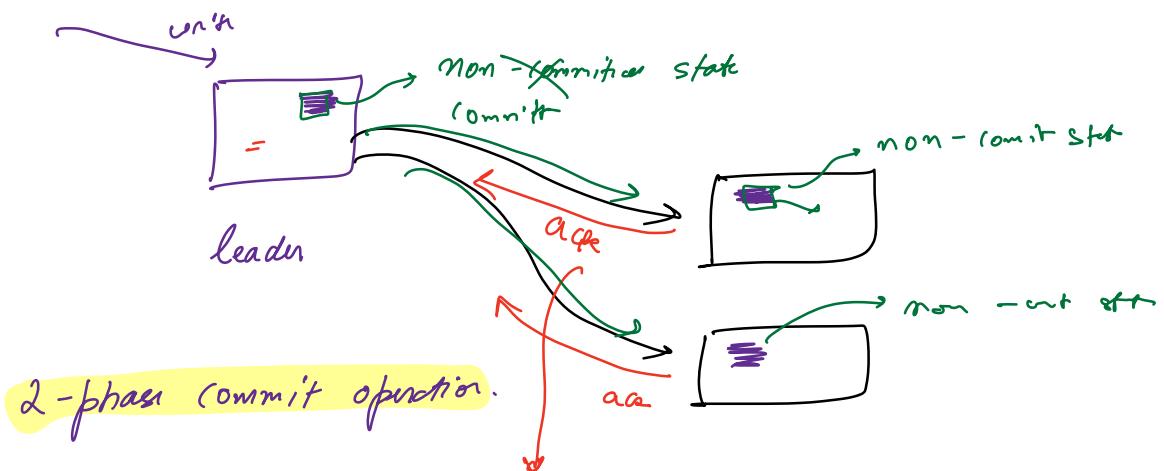
totally have been
 committed to.

read → from multiple nodes.

if we establish a quorum.

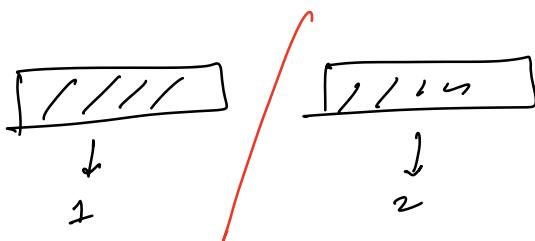
collection
 consensus

members 50%



guarantees that if in future you fail
you to commit data → I will be able
to do so.

$2N$ nodes



Split Brain

~~$2N+1$ nodes~~



$2N+1$ machine initiation → read / write will succeed
only if it is successful
at at least $(N+1)$ machines

CAP / AP CP

ZK's leader is down

↳ leader re-election is happening.

↳ any writes to ZK are forbidden.

↳ db writes/reads → still happening.

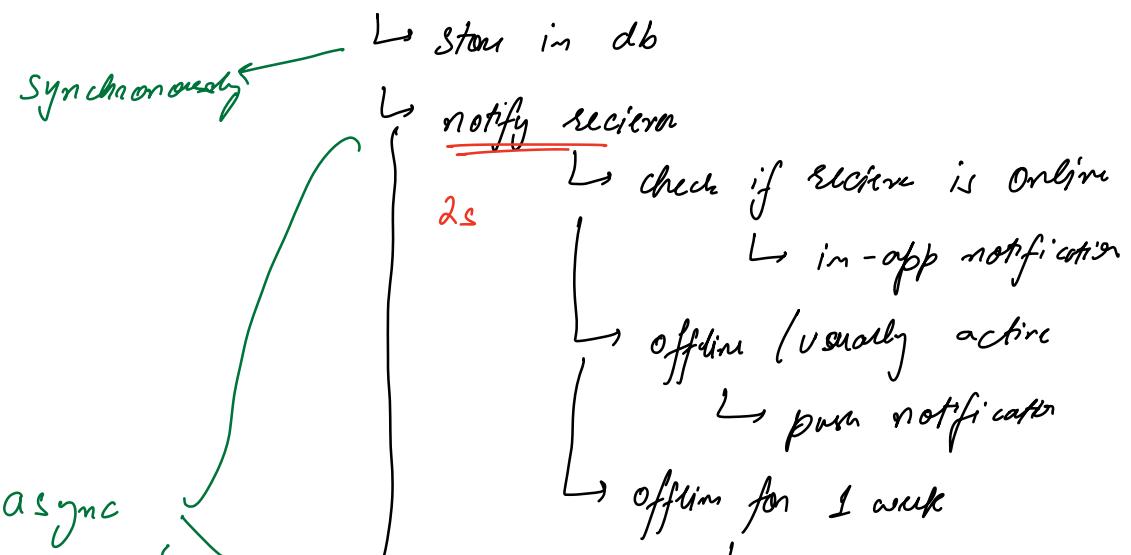
↳ any app can still read from ZK

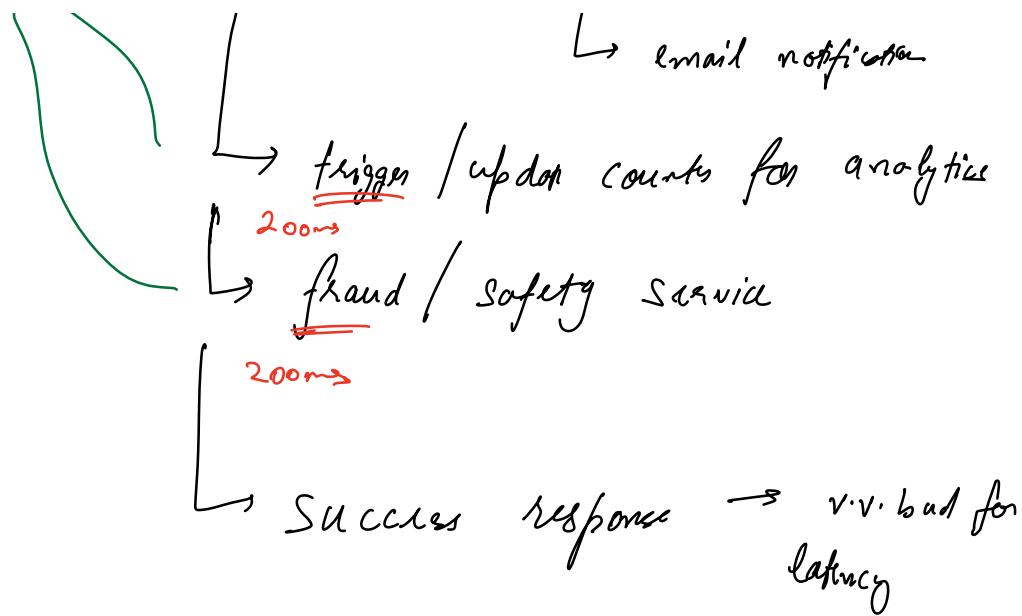
10:30 → 10:40

Problem - 2

Async Tasks / Low Latency

Whenever a message arrives

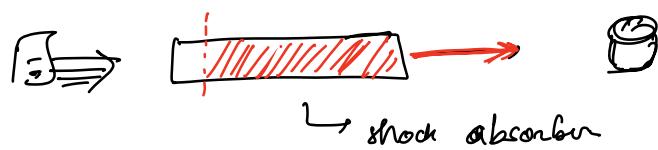




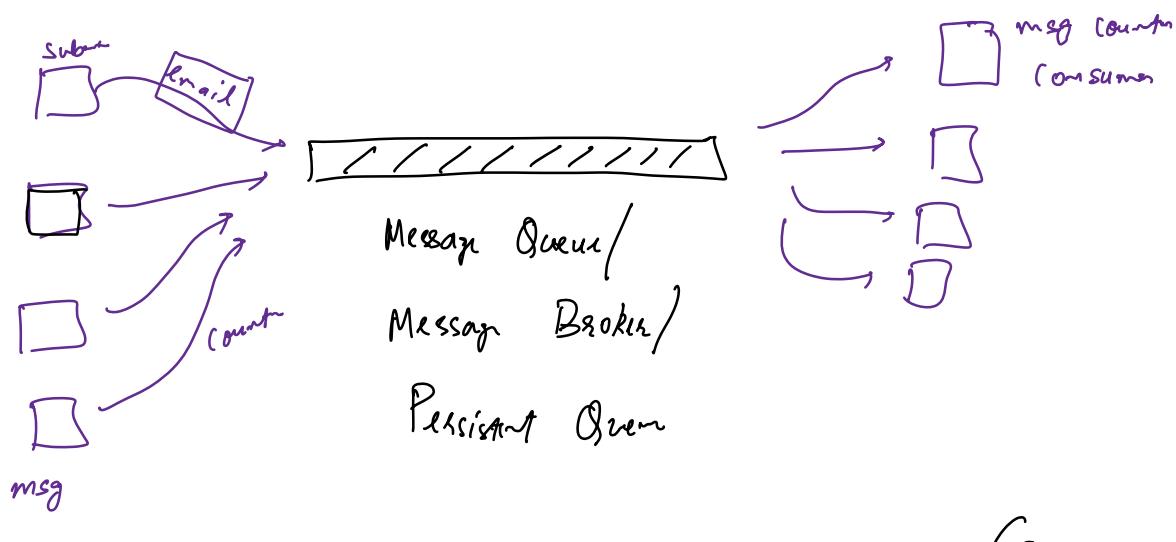
- ① we want to keep latency low
 - ↳ return success as soon as sync tasks are done
- ② async tasks are still important.
 - ↳ even though we're permanently returned success
 - ↳ these tasks are actually performed without failure

Persistent Queues

↳ any data stored is durable → stored on disk / replicated



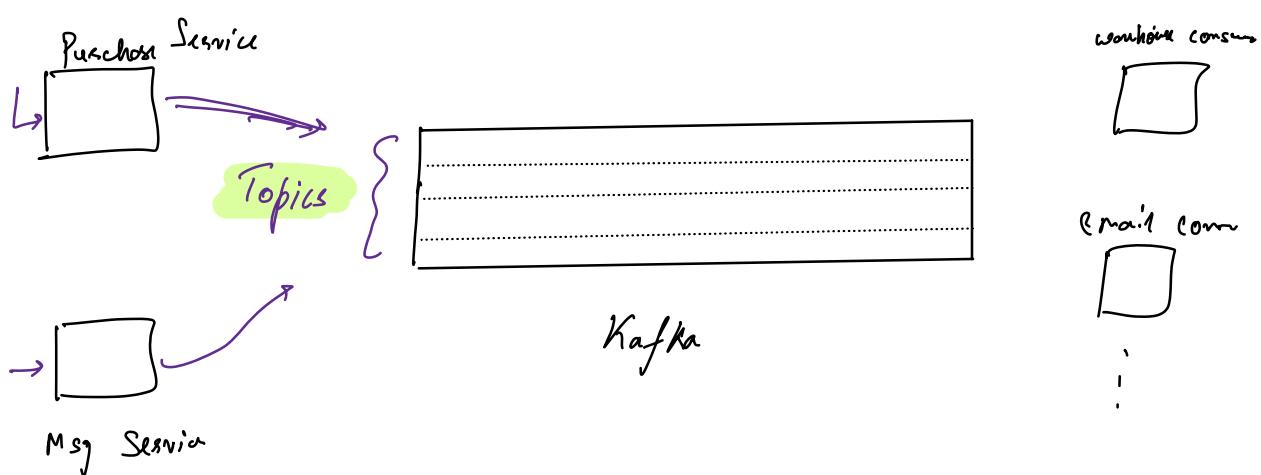
Publisher / Subscriber (Pub - Sub)

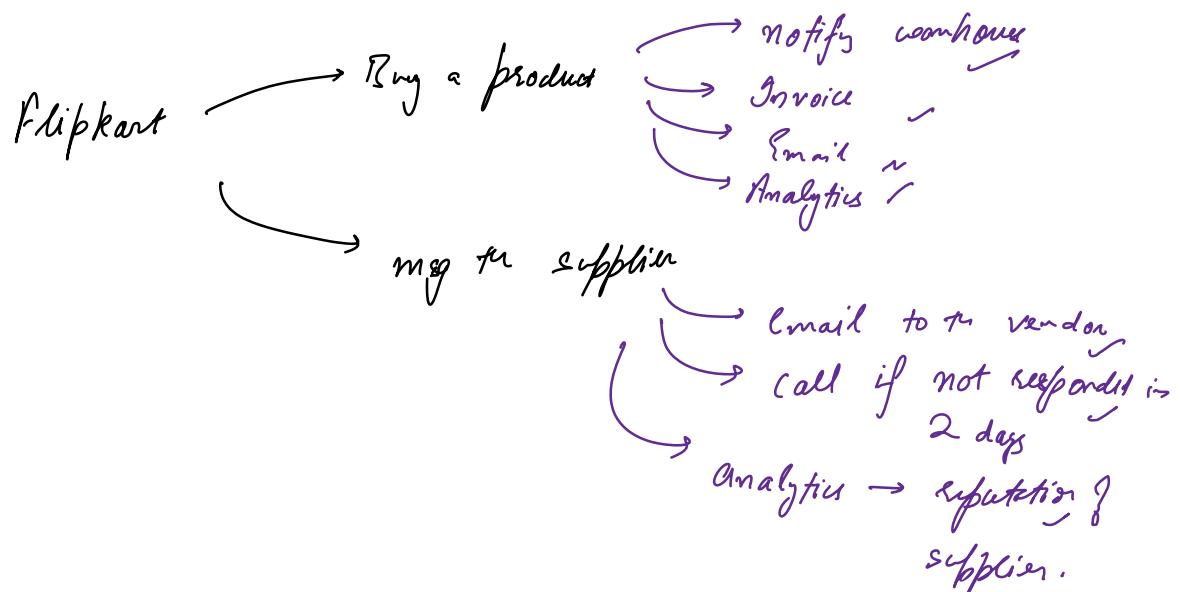


Producers

(On screen)

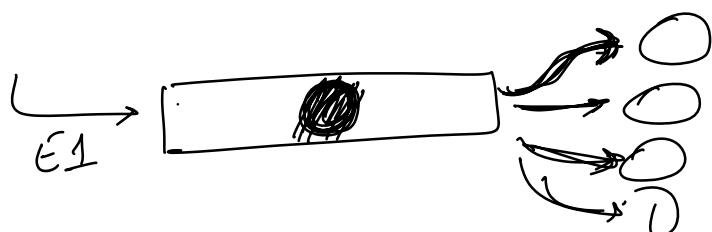
Kafka → internally uses Zookeeper.



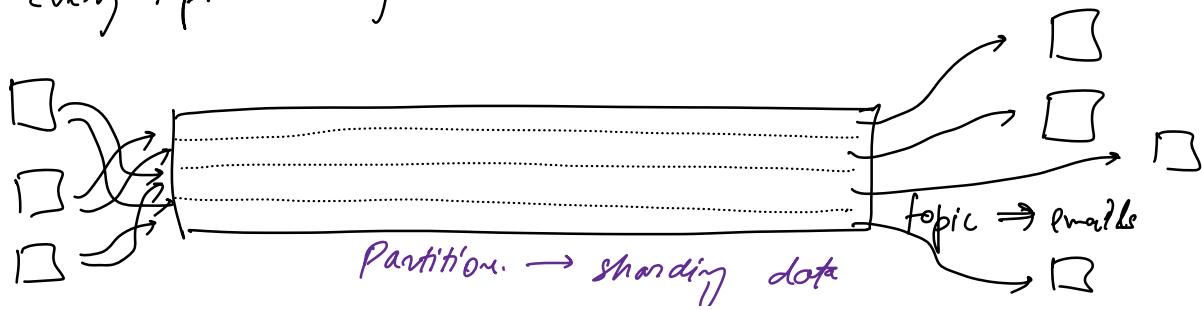


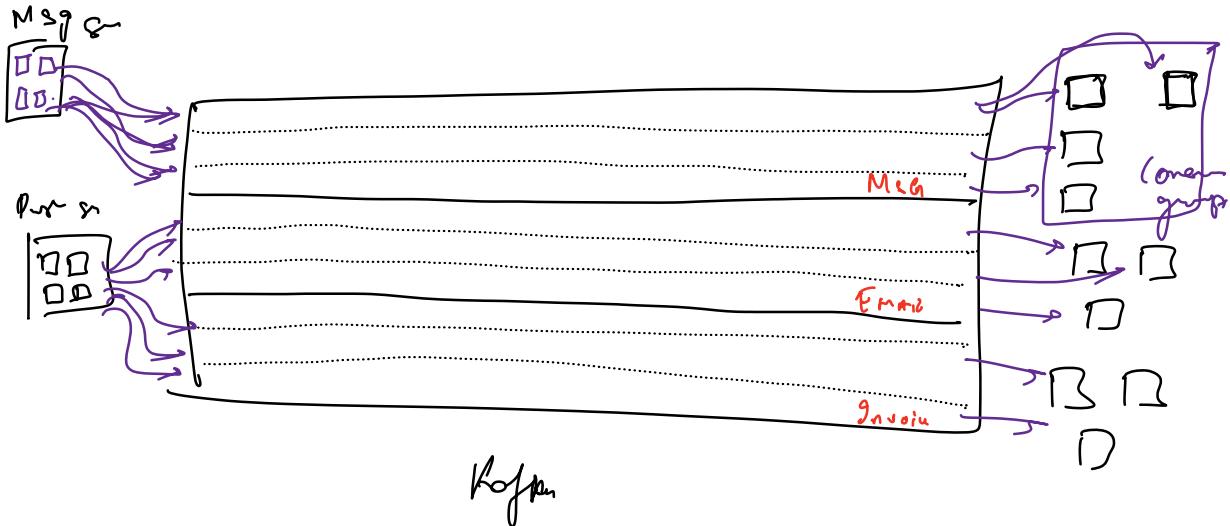
Producers can produce at different rates
 Consumers can consumer at different rates.

any data produced is persisted for **1 week** (retention period)

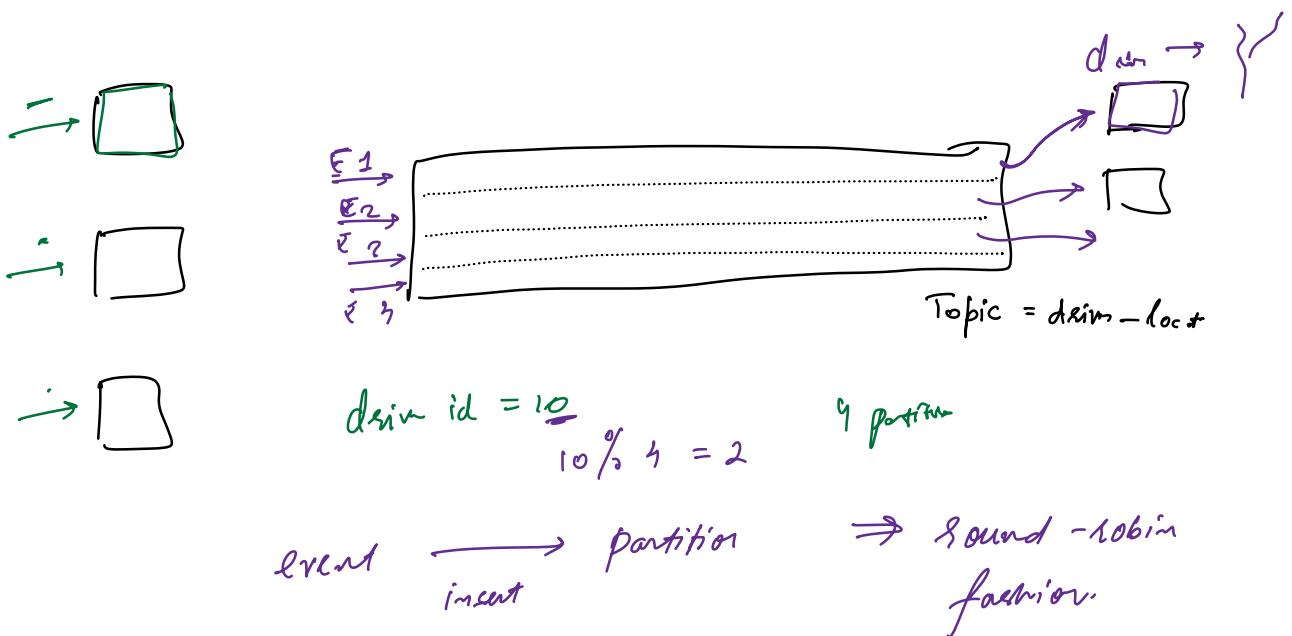


Every topic is further subdivided into partitions





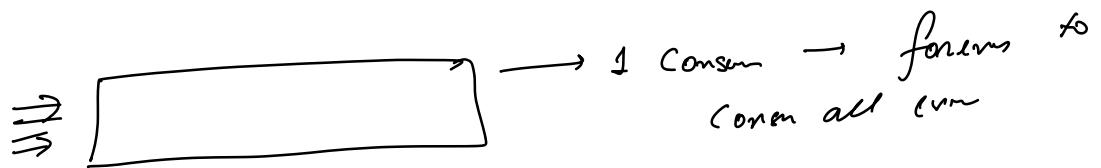
partition in a topic \leq # consumer for the topic
 Uber



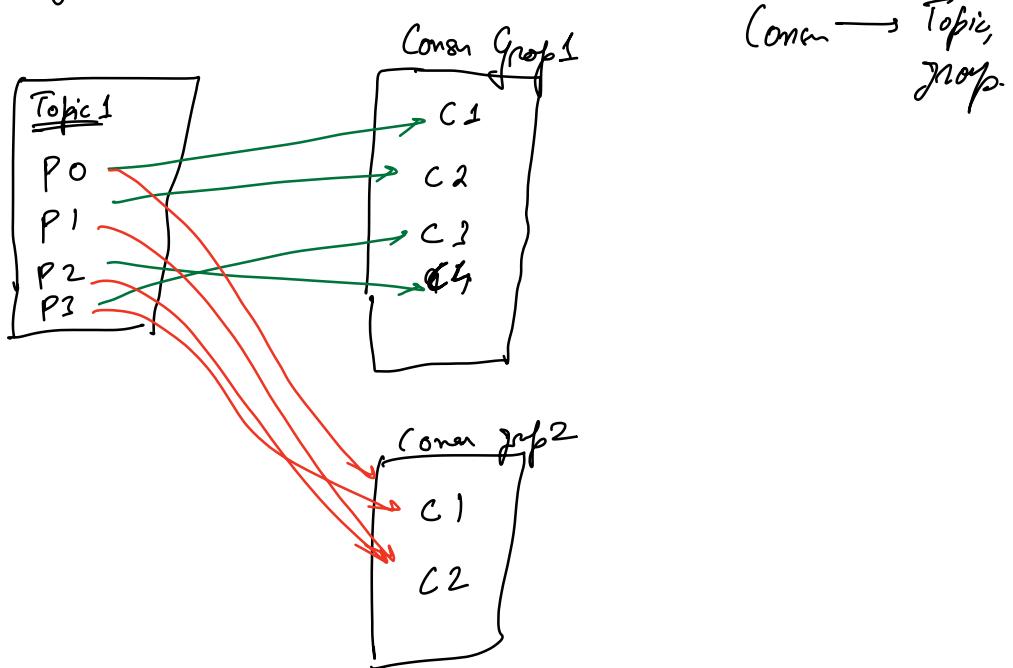
allows you to supply a key

$$(\text{hash}(\text{key}) \% \# \text{partition})$$

If topic is huge

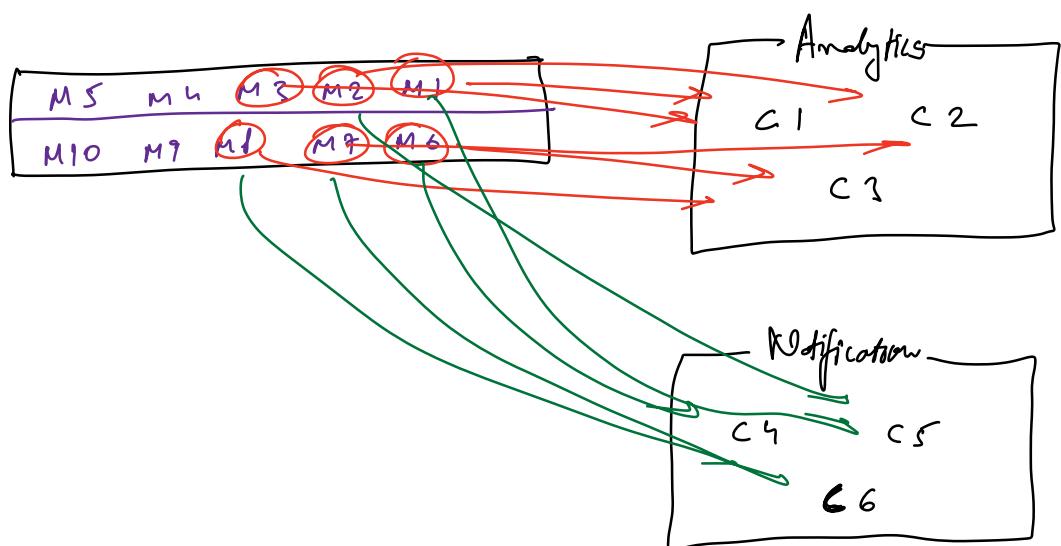
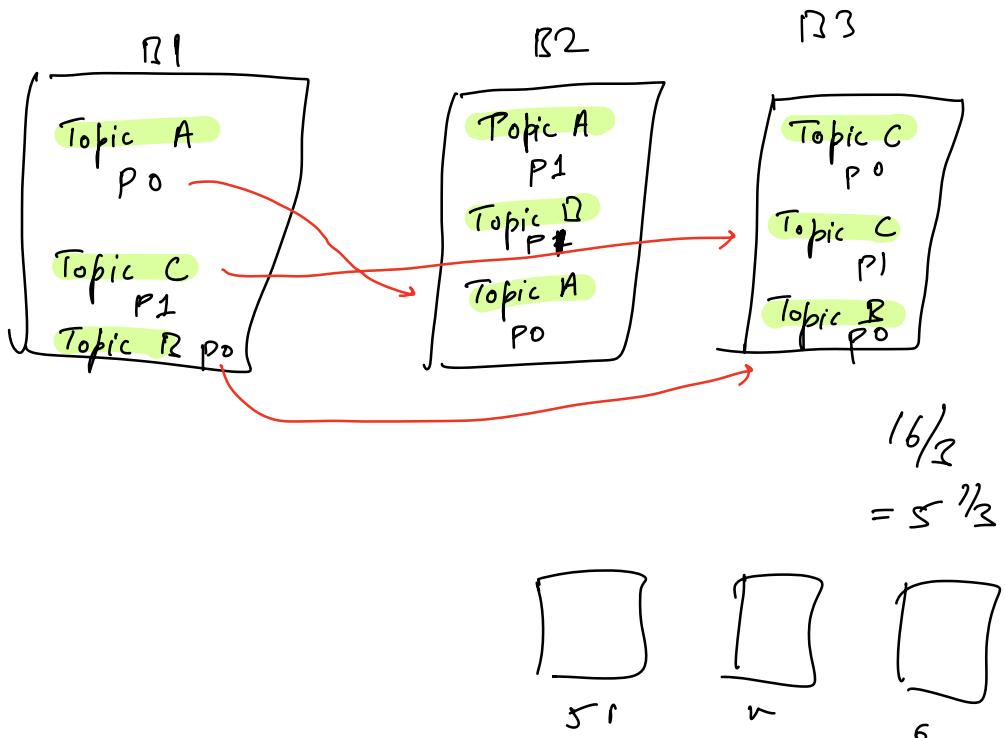


Concuren Groups



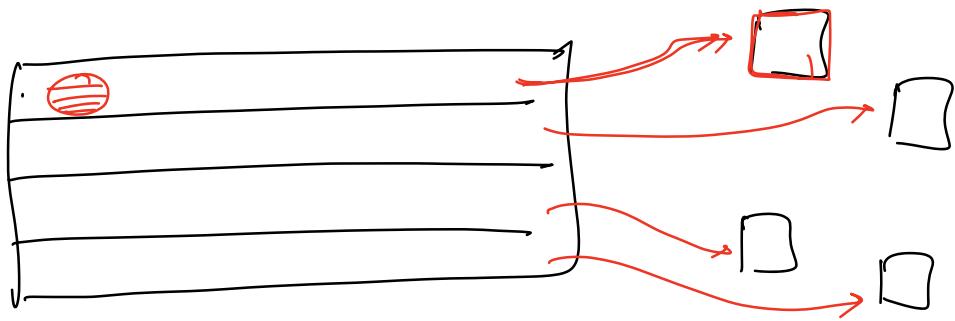
What happens if 1 machine (broken) dies → prevent data loss?
replication.

$\frac{4 \text{ topics}}{3 \text{ brokers (machines)}}$ → $\frac{2 \text{ partition each}}{\text{is 2}}$ → $\frac{\text{replication factor}}{6 \text{ per}}$



offset is added to even more
|

↳ incremental



N topics, each topic has P partition, replicate factor X

R brokers

$$N \cdot P \cdot X \rightarrow R \text{ brok}$$

$$N \cdot P \cdot X / R \text{ ifr}$$

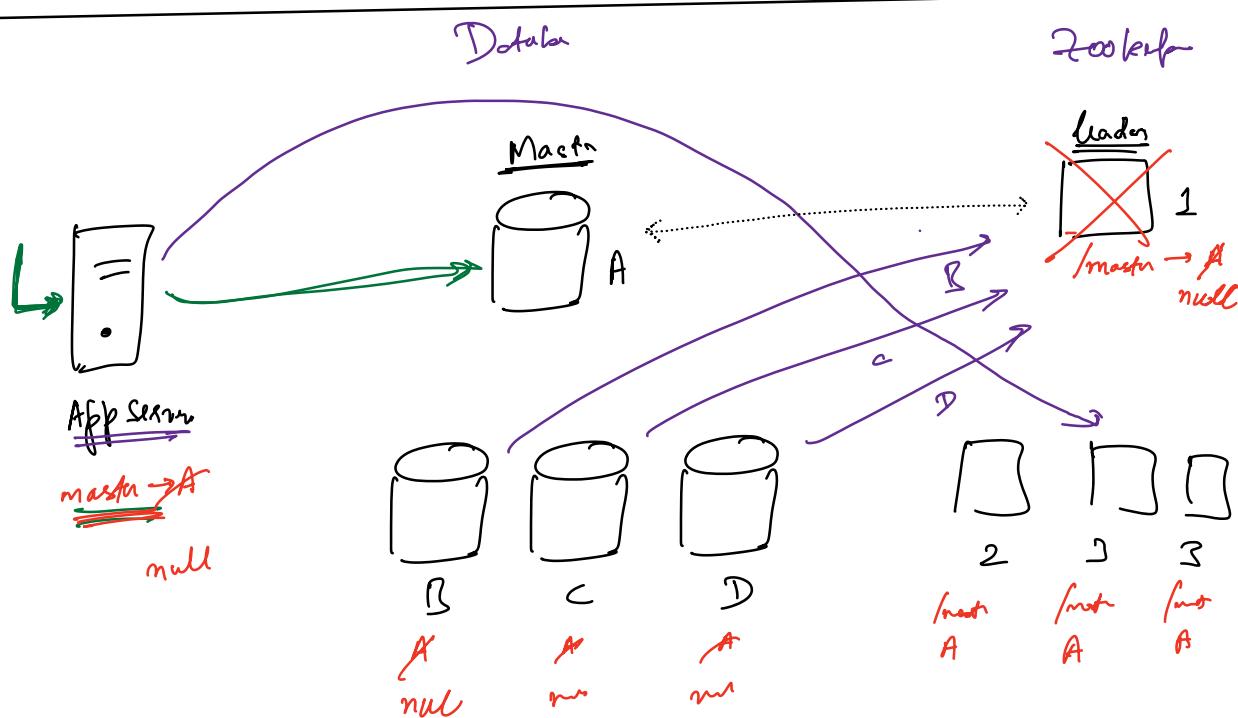
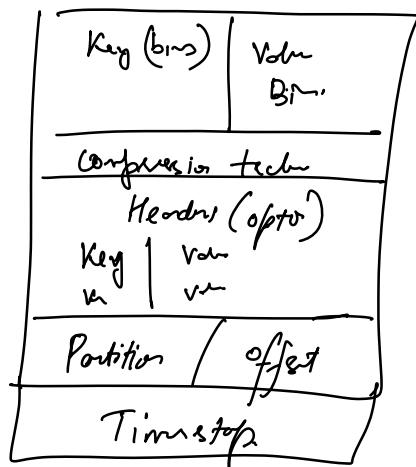
1000 events / sec

each consumer takes 100 ms to process an event

↳ 10 events / sec

$$\frac{1000}{10} = \underline{100} \text{ consumers}$$

125 Consuming



Once a month → master goes down → unavailability for 5 seconds

$$\left(1 - \frac{5}{30 \times 86400}\right)^{100\%} \approx 99.999807\%$$

SQL databases → not sharded

ACID guarantees

