

SPOF  $\rightarrow$  Replication (Introduce Replicas)

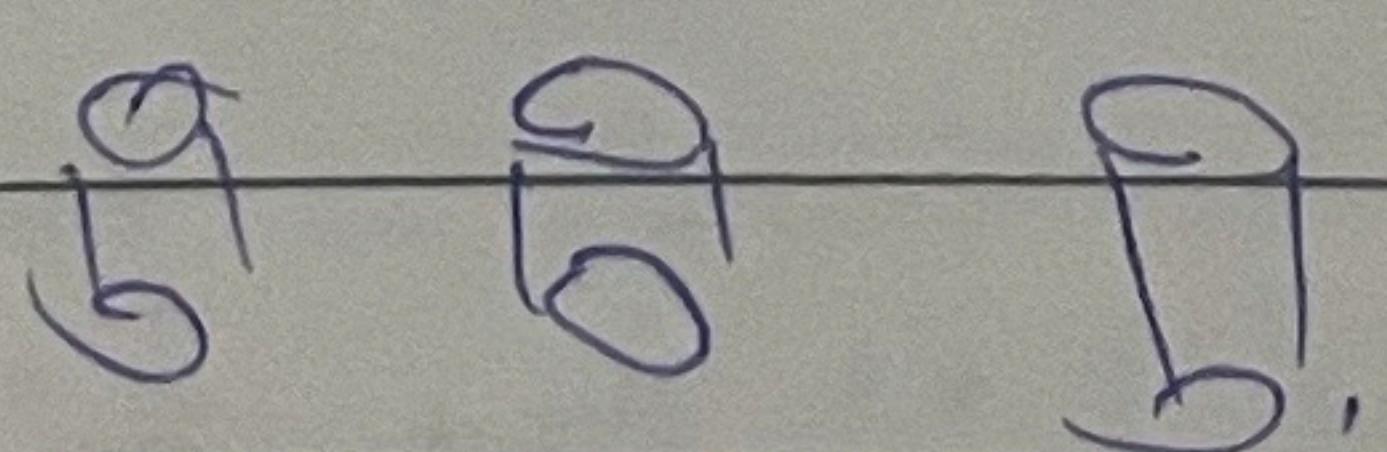
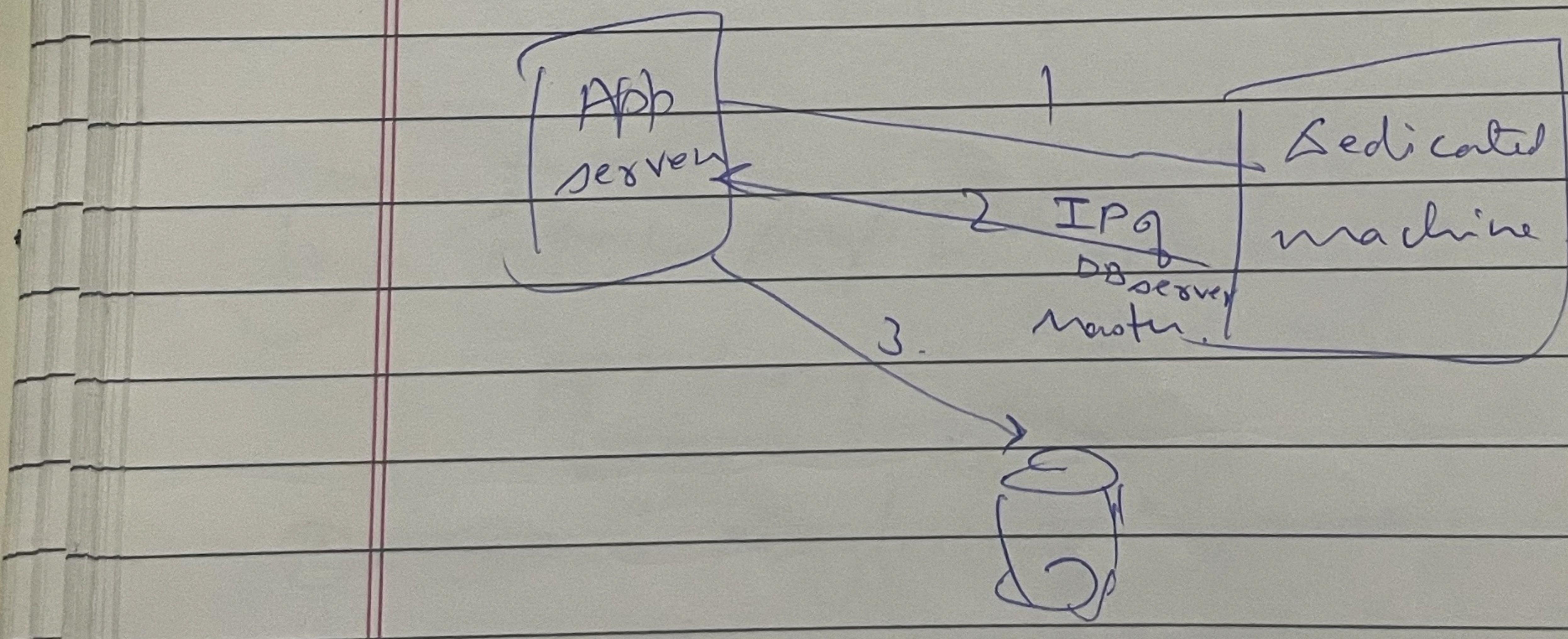
Soln: *classmate*

## Zookeeper & Kafka

↳ Config Management System

State Manager

Approach -



$\rightarrow$  Have a dedicated server

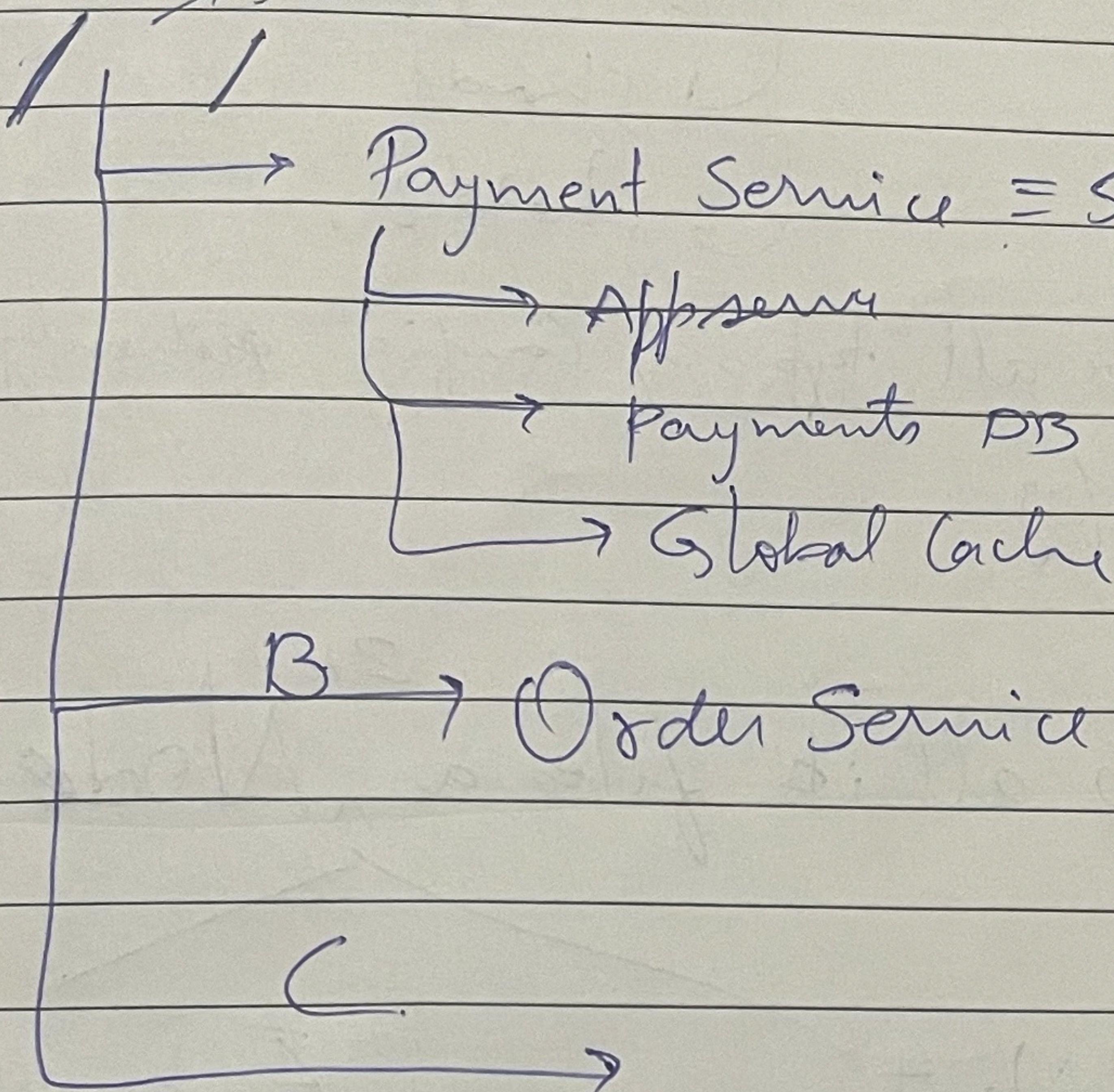
cons 1) SPOF Single Point Of Failure

2) Additional hop (latency ↑↑)

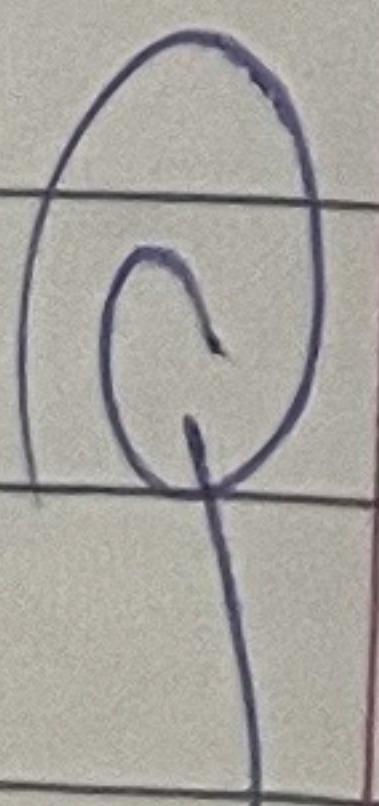
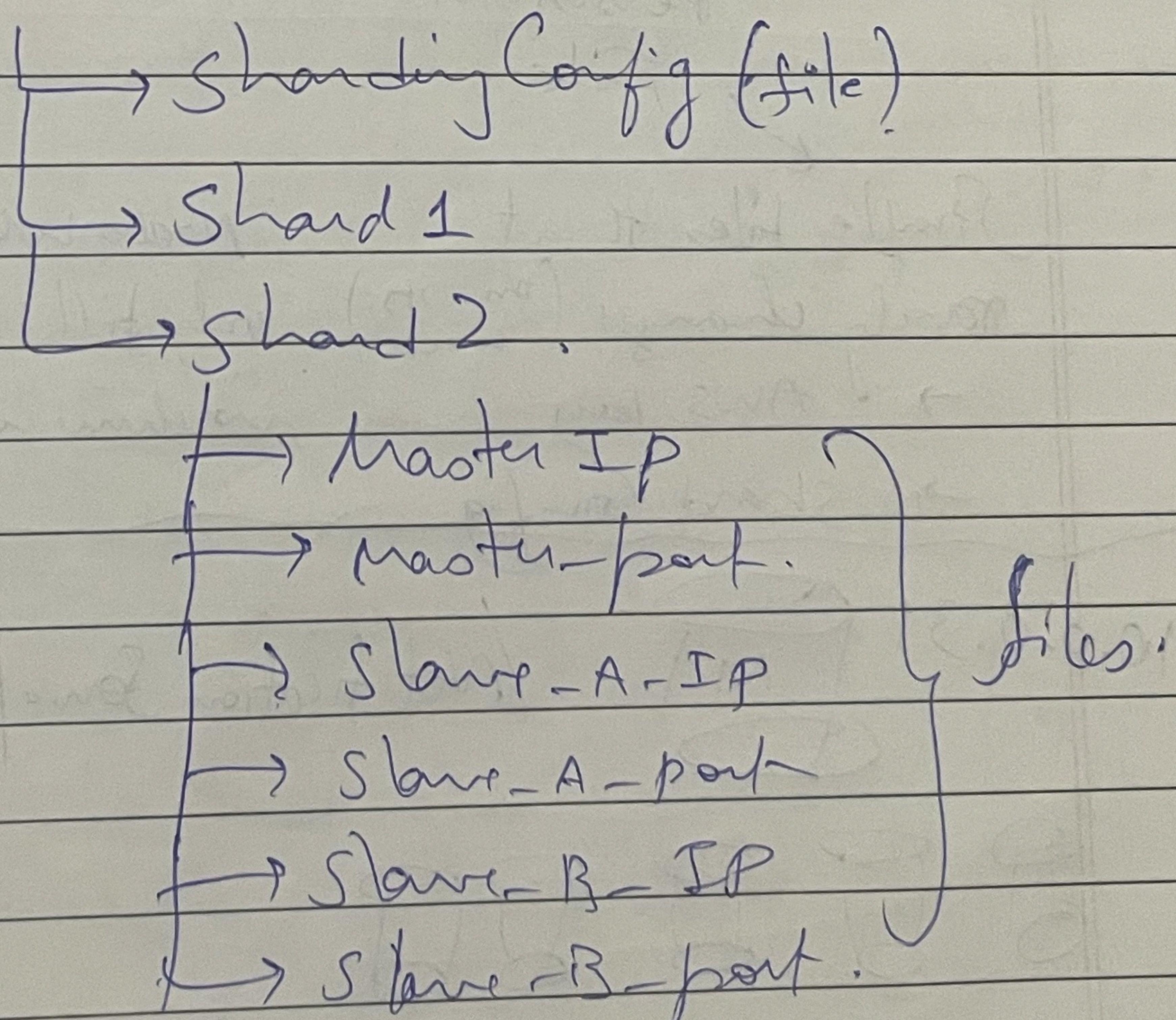
Mostly an large level system to maintain data, usually a proper segregation of data is done, kind of tree like structure.

## Solution: Zookeeper [Z.K.]

→ ZK is exactly File system. (Same like Linux)



/A | payments DB



Why can't we maintain all the above config in just a ~~single~~ (xml), as it is also structured? This being a multi-server system, we need data to be consistent. So if the data is in diff files, we can obtain LOCK on the files to keep it consistent. we can't have race conditions.

/A (payment DB/)

↳ aws-key  
 ↳ shard config  
 ↳ shard 1  
 ↳ shard 2.

It can have all types of config not only  
 Shard config.

~~ZK calls all its files as Node~~

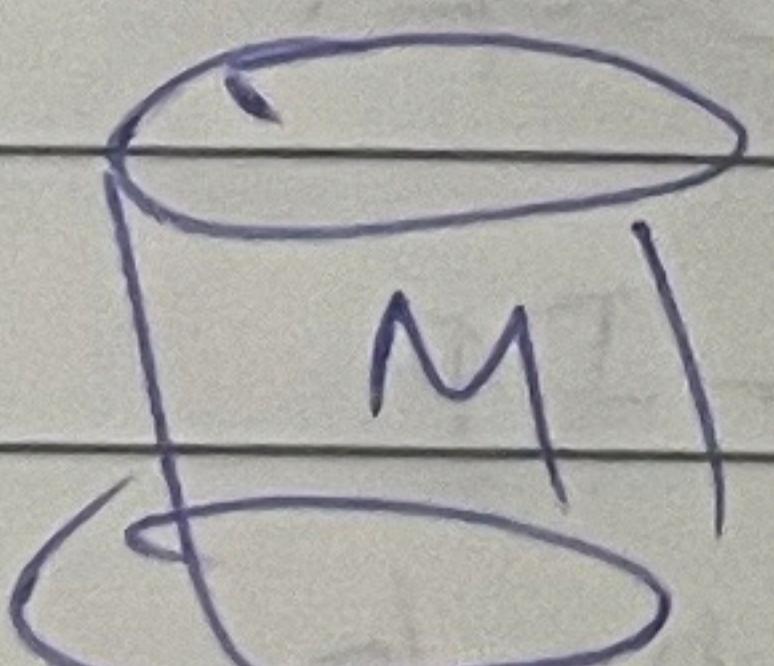
Persistent Node  
 (OK)  
 Persistent File

Ephemeral Node  
 (OR)

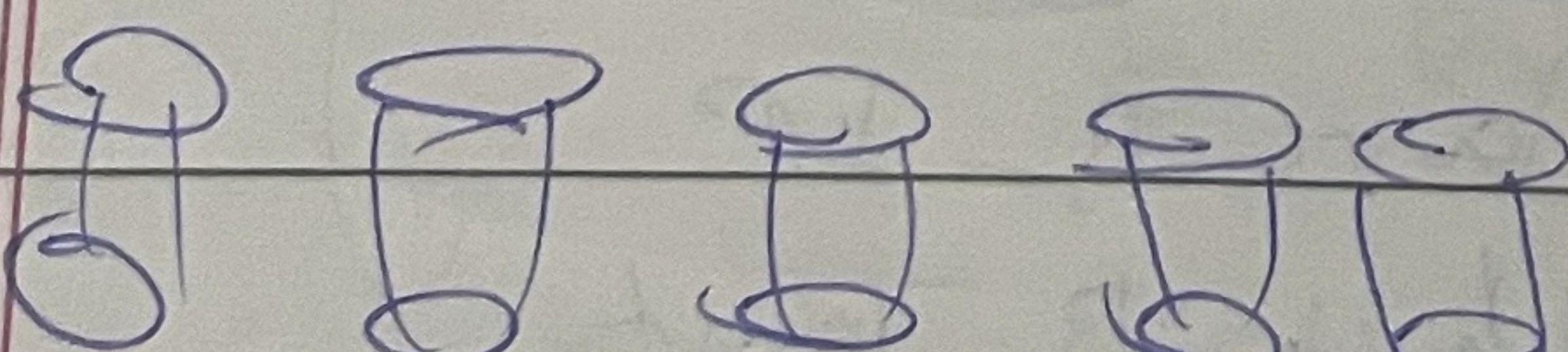
Temp files  
 ↳ Master IP

Simple files that rarely changes (once in a week) Nodes which has to live only till its underlying machine is ACTIVE  
 → AWS-key  
 → shard config

10.0.0.31



(Notification Service | Database | shard 37)



→ Master IP  
 10.0.0.31

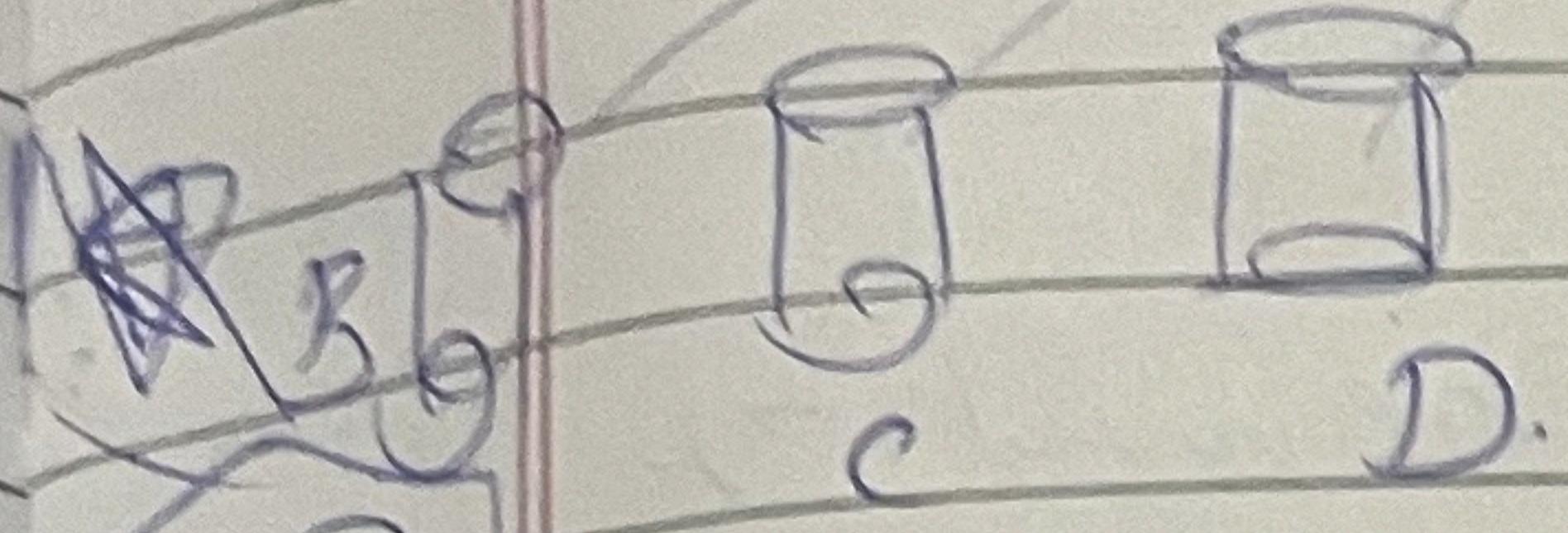
If M dies its file in the above folder should also be deleted as its no more connected

ZK

## Notification Service / Database / Shard

IP: A

Heartbeat.

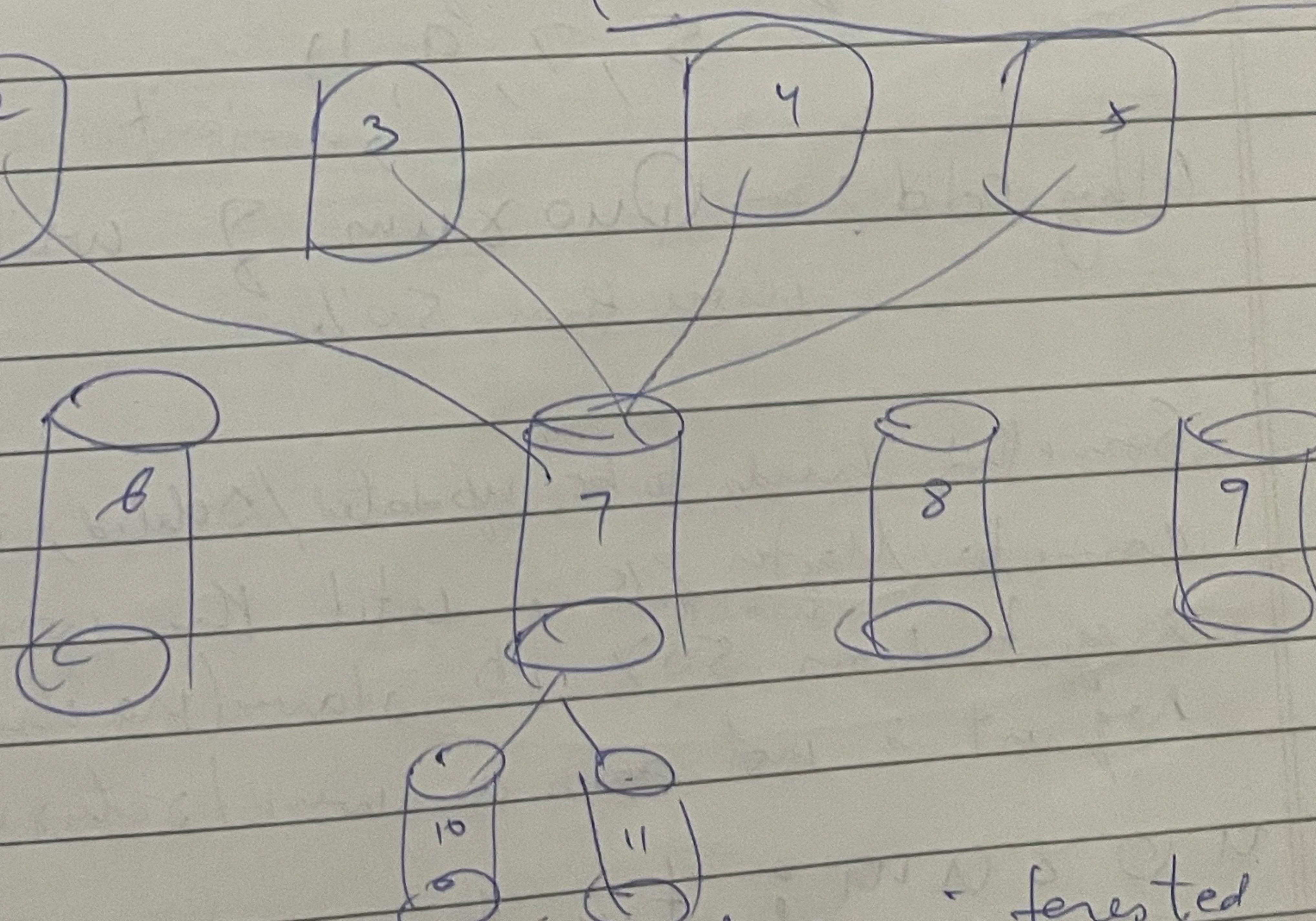


→ Master: A  
→ Slave 1: B  
→ Slave 2: C  
→ Slave 3: D

Until the heartbeat is available, files remains basically the owner (n) is the owner of Master-IP file. If it is alive it remains in heart beat missed, files get deleted.

## ZK Watch (Subscribe to a file to know its status change)

When 7 dies, all subs are notified, 10 and 11 will start election and winner goes to ZK as new ownership on new master file with its IP



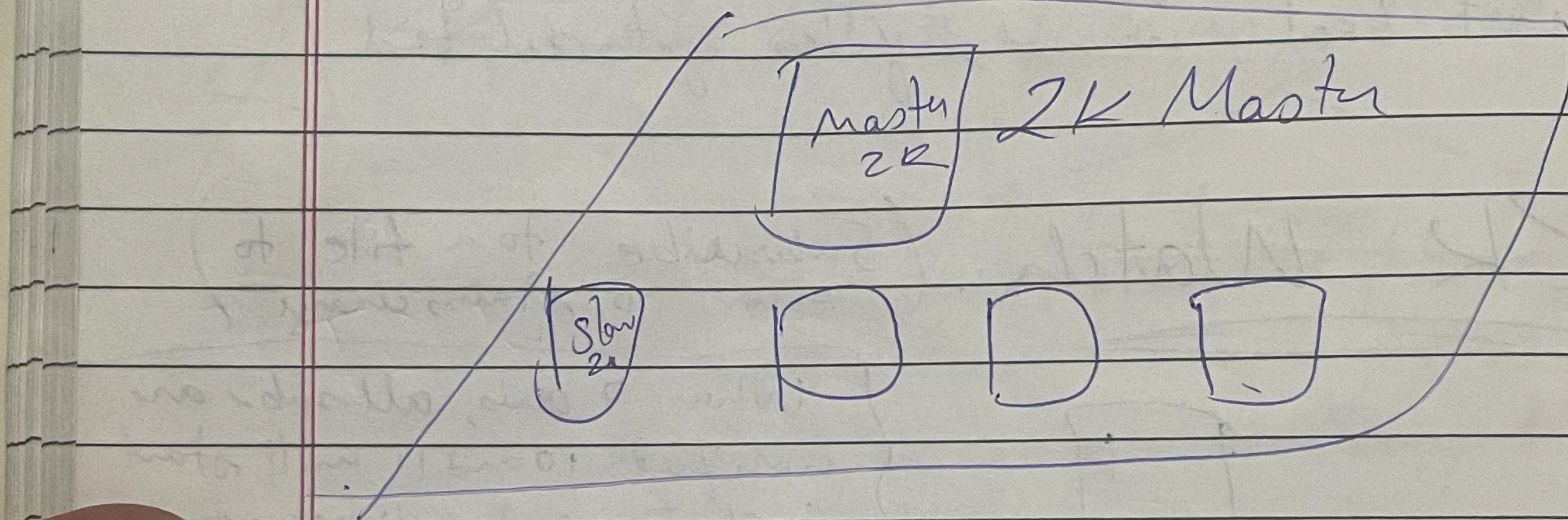
10 and 11 (slaves) would be interested to know when the Master(7) dies, so 10 and 11 subscribe to 7(file), App servers 2, 3, 4, 5 will also subscribe to 7.

ZK = File System

ZK files = ZK nodes.

(1) Persistent ZK Node      (2) Ephemeral ZK node

ZK is internally itself is a cluster  
so it's not a SPoF.



Cluster of odd no. of machines

5, 7, 9, 11

majority

Why odd? Quorum of writes is more than 50%.

Something needs to be updated/deleted, request comes to Master ZK, until the write is registered on  $\lceil \frac{n}{2} + 1 \rceil$  slaves (machines). Request is not answered/returned.

Quorum: try to write on all slaves, the request on Master is answered back only when  $\lceil \frac{n}{2} + 1 \rceil$  servers are slaves.

Even when any server dies majority of servers will have correct data

# Gossip Protocol.

~~Under Master election~~

For Master election.

## Async Tasks

can handle.

→ 100 req/sec.

System  
A

can handle.

→ 20 req/sec.

System  
B

①

② Rate of Production ≠ Rate of Consumption.

③ If we want to do aggregations in b/w Queues can be done. [Digests].  
outlook digest

ASYNC Tasks → Persistent Queues.

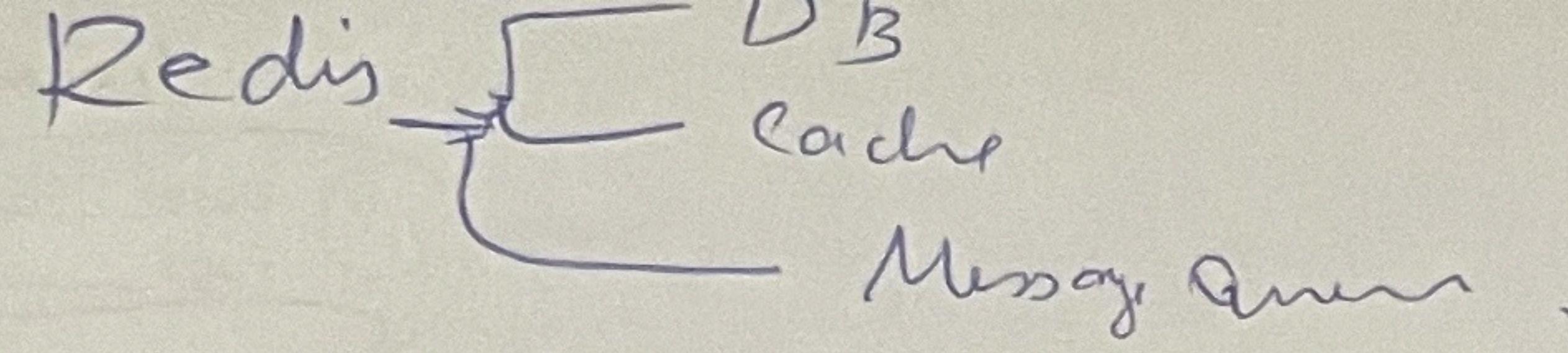
Publish Subscribe

Pub-sub

Message Queues.

Ex: Kafka.

RabbitMQ.



classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

# Kafka

## Principle

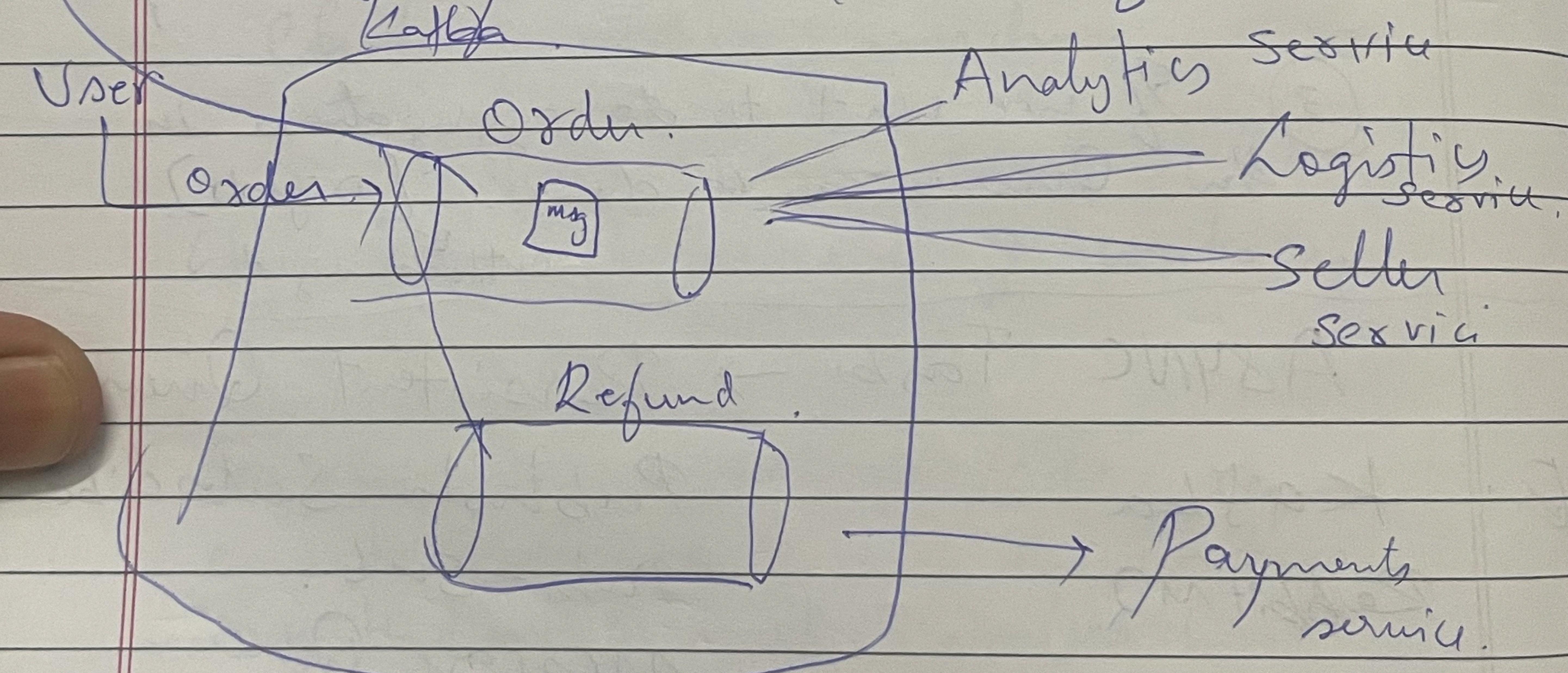
- 1) Distributed. ~ Horizontally scaled easily  
Ex: 10 Million msg/sec.
- 2) Persisted Queue.

## Kafka Queue / Message Queue

Kafka  
TOPIC

6 distributed.

Collection → Machine / cluster of machines



Msg should be read by all 3 services and only once.

## Web hooks.

### Working:

FIFO

msg-id.

101 102 103 104 105

106

→ Maintains last read msg as its FIFO seq. can be maintained.

Registered Consumers.

Analytics

Offset

Seller.

- Msg doesn't get deleted after reading
- Msg id will be set with TTL
- Msg will be deleted only when TTL expires

~~Kafka internally uses Zookeeper to save configs such as OFFSET~~

Msg could/should be kept like, have only ordered or any required info only any more info required, DB call should be made to get more info.

### Problems

- As Queue is FIFO (ordered) all msg will have to be in one server (One topic in one server)
- Can't be read parallelly as it FIFO outputted.

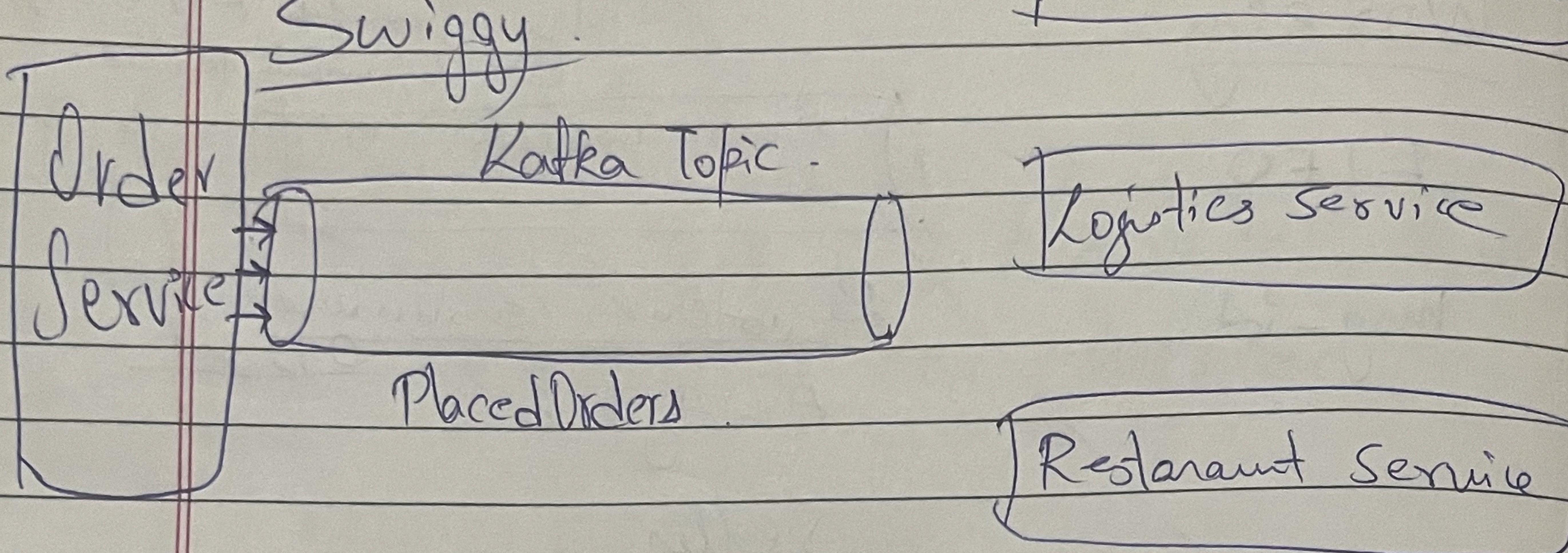
18/12/2023

## Kafka Contd...

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

Notification Service



- If there are many Topics in a server
- If we encounter some issue, separate Topics into separate servers.
- There may be a situation where 1 entire machine may not be sufficient to store all msgs  
1 Queue(Topic) also

in notes

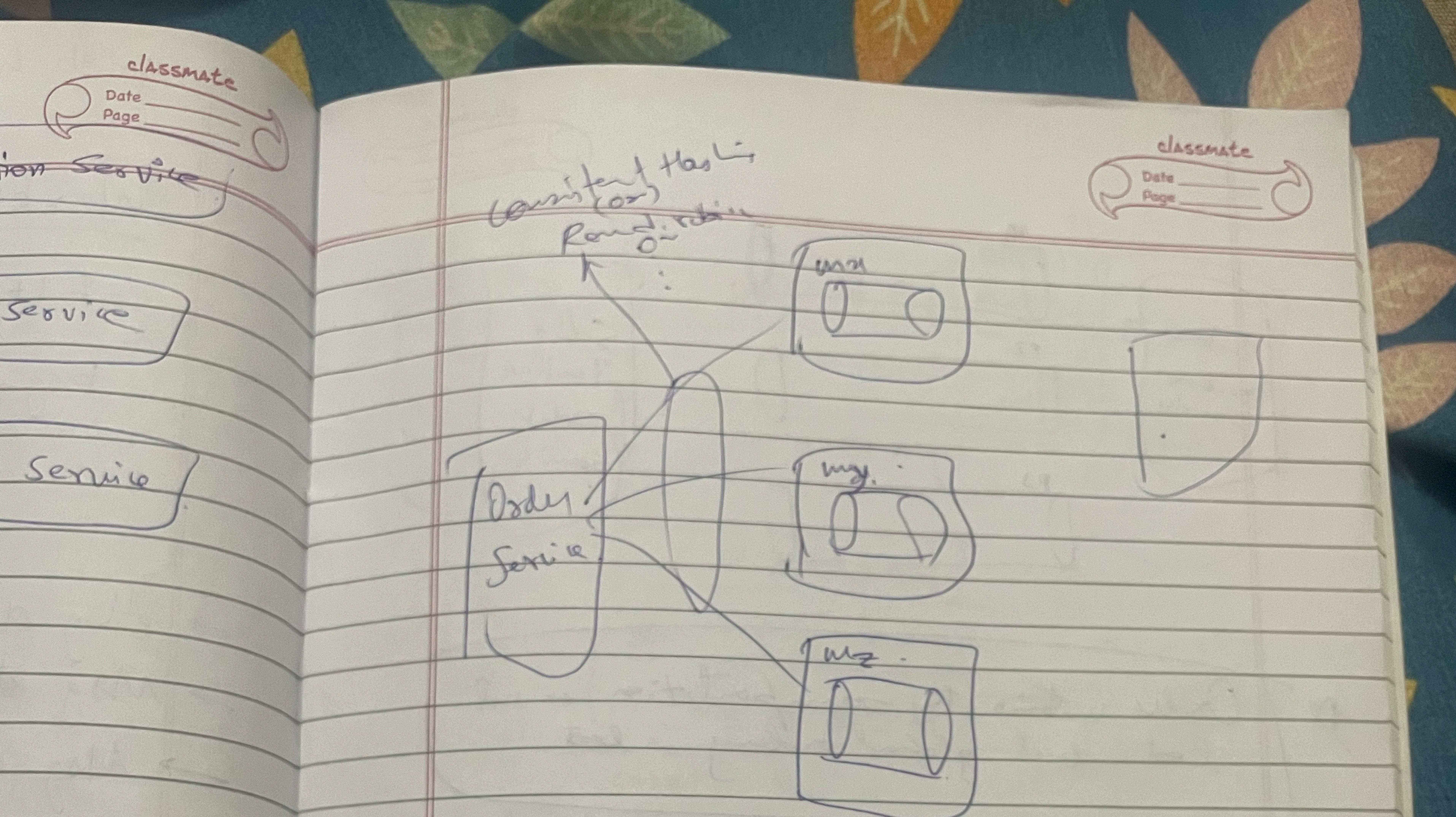
The Usual Way to solve above problem → SHARDING

But we can't use sharding directly before understanding it...

↳ Queue is FIFO [order should be maintained]

Problem: When we divide a Queue across multiple machines, it no longer remains One Queue. It becomes multiple pieces of One Queue.

→ If we put up msg in any fashion (CHAI RR) it will be difficult to maintain order (FIFO) to read the msg.



→ One possible soln. get Timestamp. from all partitions ( $m_1, m_2, m_3$ ) and pick the latest (Very very unoptimized)

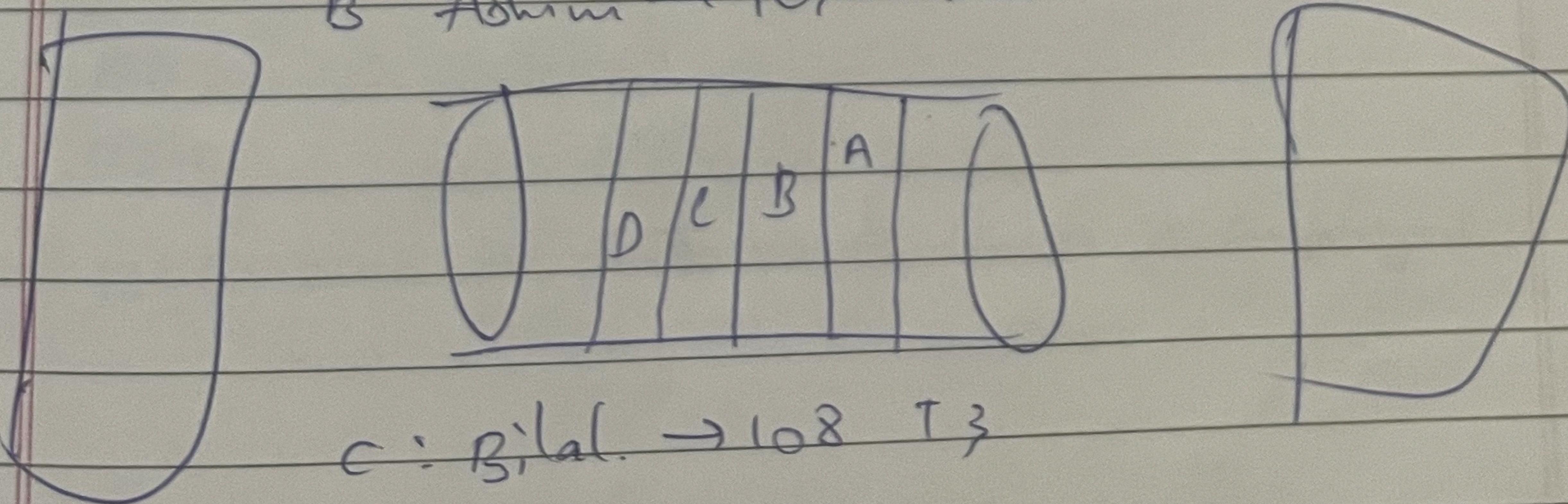
in notes

The Kafka Topic when partitions, Kafka no longer gives FIFO guarantees across all partitions

HARDING

A: Ashim → 107 T1

B: Ashim → 107 T2

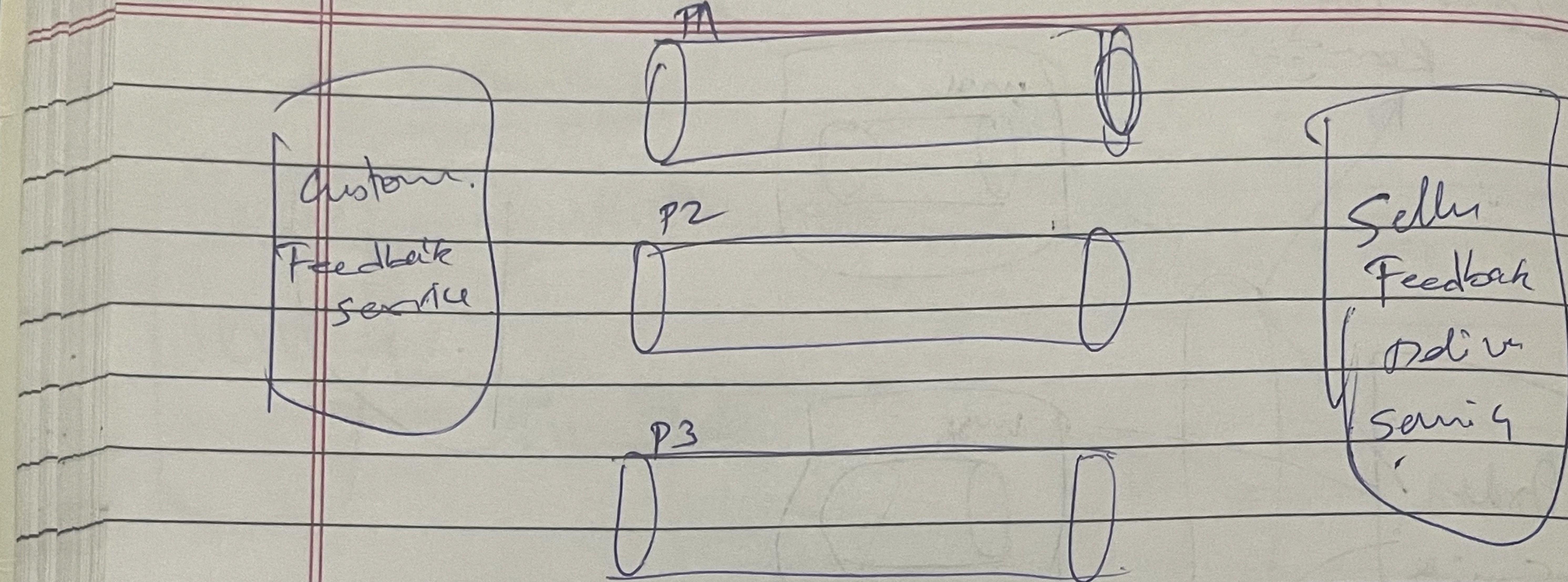


C: Bilal → 108 T3

D: Bilal → 107 T4

→ the order is imp. w/ respect to Review as well - decrease  $T_2$

$T_2$

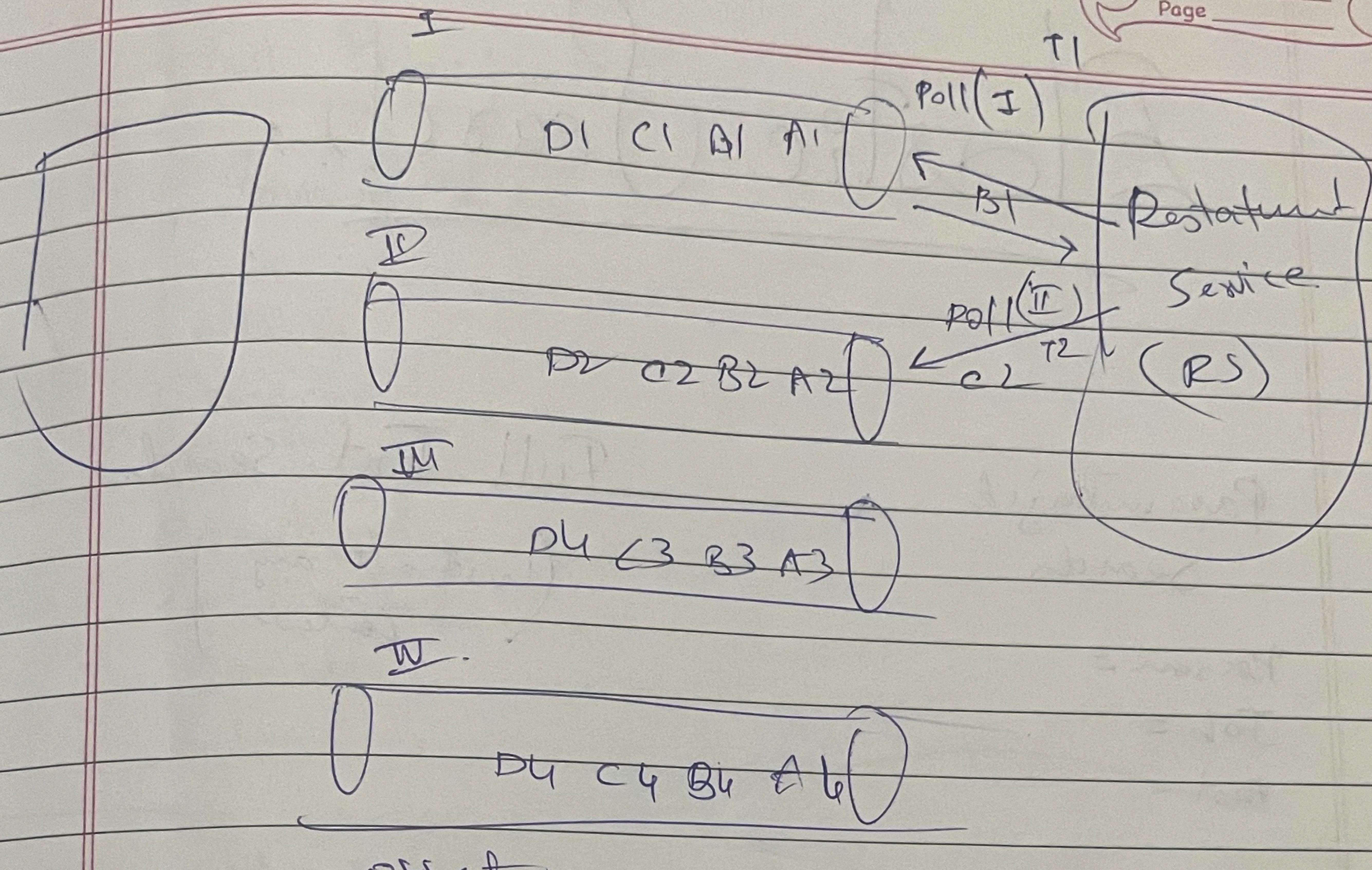


→ Order within the same partition is not lost but global ordering is lost

If seller wants to read in sequence (or) it is imp. to seller is read in order maintain (H or seller ID)

→ basically depends on consumer service

→ If we have multiple consumers, we have to decide on one CF key and every other consumer will have to leave with it (Trade-offs)



RS  
 I → AT B1  
 II → B2 C2 + 2  
 III → A3.  
 IV → C4.

If RS is reading one by one from partitions  
 we can use this idea to increase performance  
 and employ different RS web servers to  
 read from diff partitions separately

RS1 read from partition I.

RS2	"	"	"	II
RS3	"	"	"	III
RS4	"	"	"	IV

RS5 will be ~~loss~~ waste even if employed