

# Designing Data Intensive Applications

## CHAPTER 5

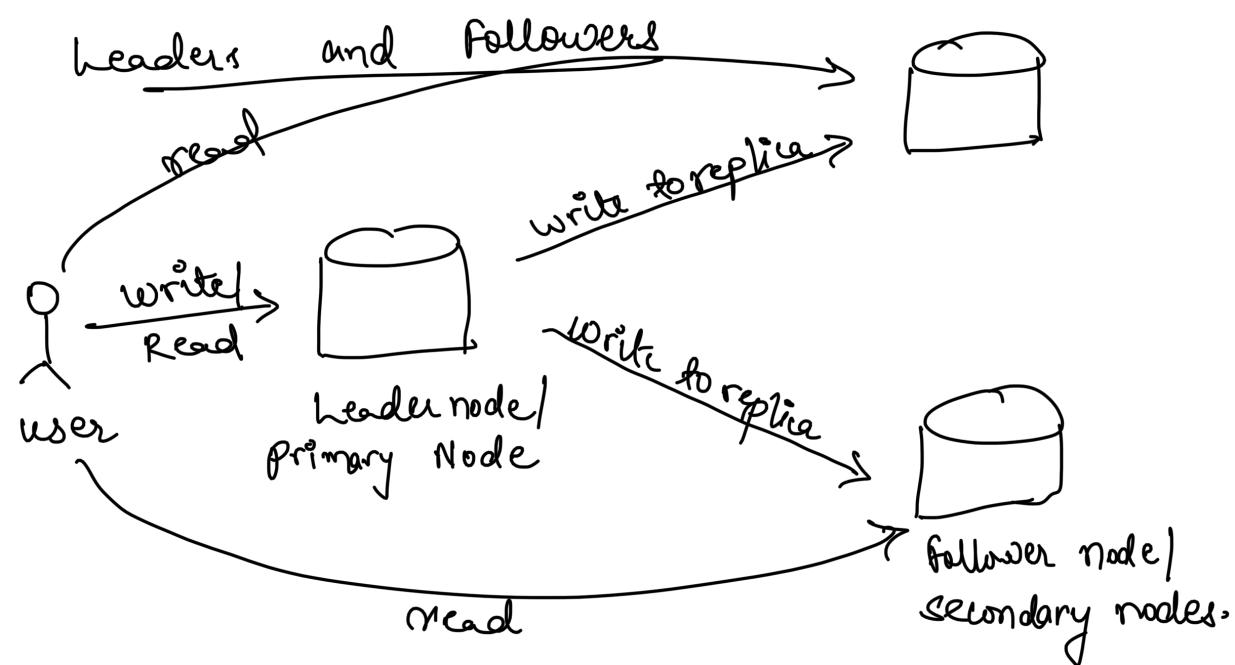
### Replication

Keeping copy of data at multiple places.

why?

- ↳ keep data in some region as user → reduce latency.
- ↳ Availability → System is operational if certain machine goes down
- ↳ can serve more Read TPS.

How to Replicate?



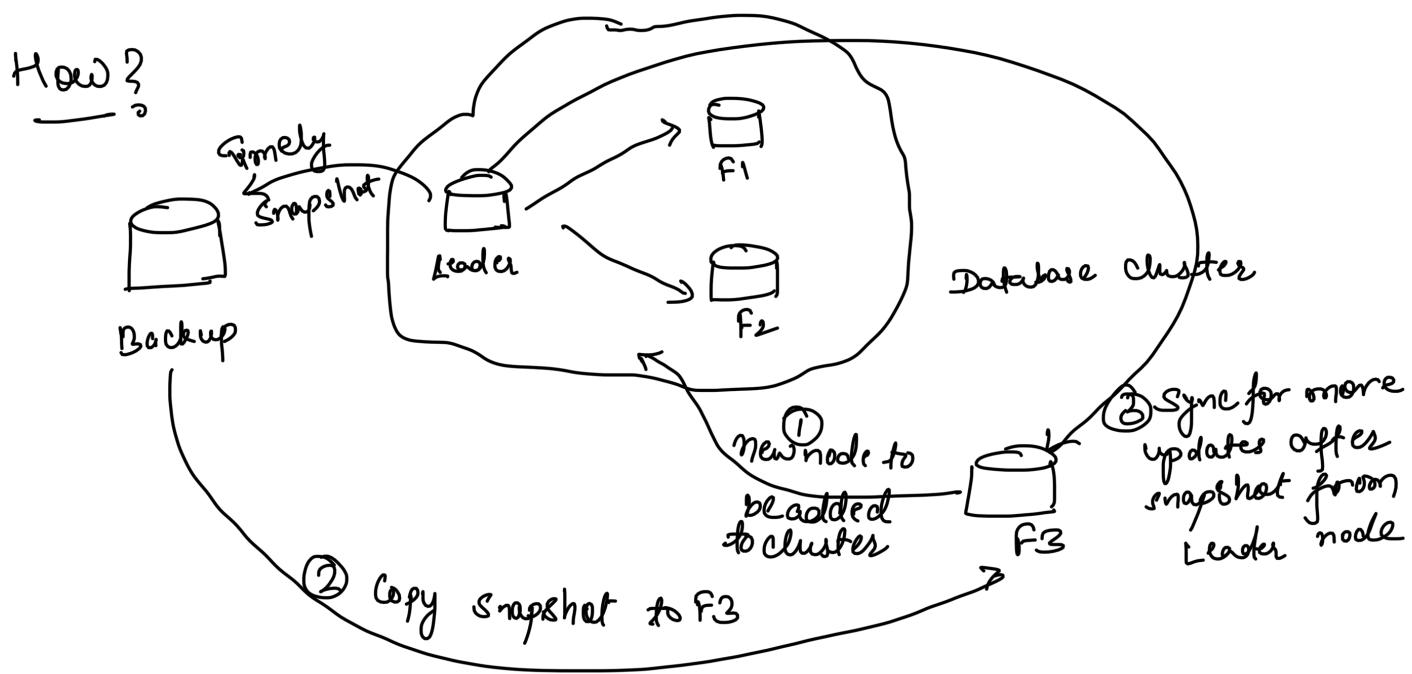
Synchronous vs Asynchronous Replication

- |  |   |
|--|---|
| <p>① User is sent ack once write is completed on leader &amp; followers</p> <p>② Data is always consistent</p> | <p>① User is sent ack once write &amp; complete on leader node - follower nodes are updated in background</p> <p>② eventually consistent.</p> |
|--|---|

③ write can't be processed if any nodes goes down.      ④ For write → leader node should be available.

### Adding new follower

- ↳ in case more replicas are required to handle load.
- ↳ replace failed nodes in cluster.



### Handling Node failures

#### ① Follower Failure

- ↳ keep log of data changes on local disk. Once restarted, recover using log.
- ↳ sync with leader node for updates during downtime.

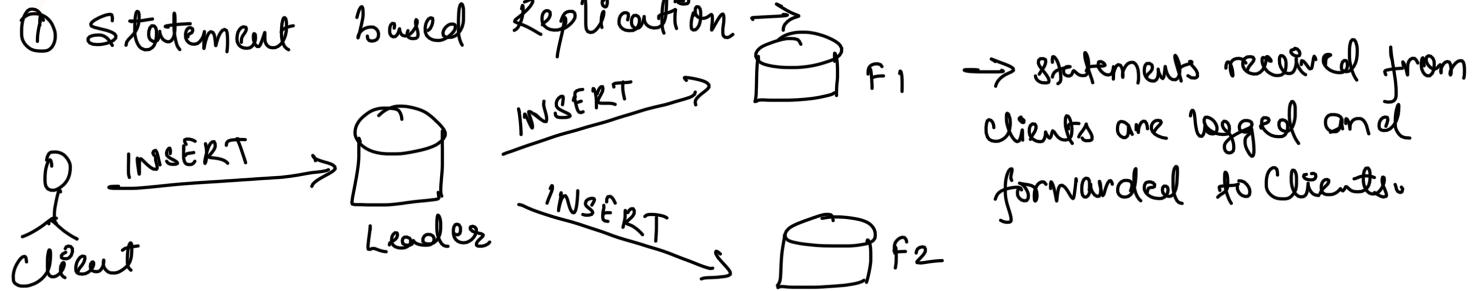
#### ② Leader Failure

- ↳ A follower is promoted as new leader.
- ④ Identify leader has failed → health checks
- ⑤ Choose new leader → best candidate would be node which has most recent replicated data.

- ⑥ System use new leader → new leader accepts writes from clients.
  - ↳ If crashed node comes back, it is promoted to follower.
- {more details on above in coming discussions}

## Replication logs

### ① Statement based Replication →



→ statements received from clients are logged and forwarded to clients.

#### Problems

- ↳ functions like NOW() or RAND() can store different value on nodes.
- ↳ if query depend on existing data, they must be executed in same order on each replica. → else can cause different effect.

### ② WAL - Write ahead log

Append only file to maintain the operations on DB → discussed in Ch-3

### ③ Logical (Row based) Replication

Use different log format for replication & storage engine.

↳ issue in #2 where WAL is used by replication & internal BTree.

① Row Insert → new values for all columns

② Row delete → primary key of row

③ Row Update → new values for columns that changed.

### ④ Trigger Based Replication

write custom code that is triggered for any operation in DB.

## Replication Log

↳ application can see outdated data in case of async replication.  
 ↳ eventual consistency.

### ① Read After Write consistency → "read your own writes"

The user who made an update, will see recent data on read.

This might not be the case for other users.

↳ a Decide when to read from leader and when from Followers.

⑥ for 1 minute after write, read only from leader.

what if same user reads info on multiple devices?

↳ can user requests routed to same datacenter to handle this.

## ② Monotonic Reads

Users gets different data on subsequent reads.

↳ can same user always read from same replica.

## ③ Consistent Prefix Reads

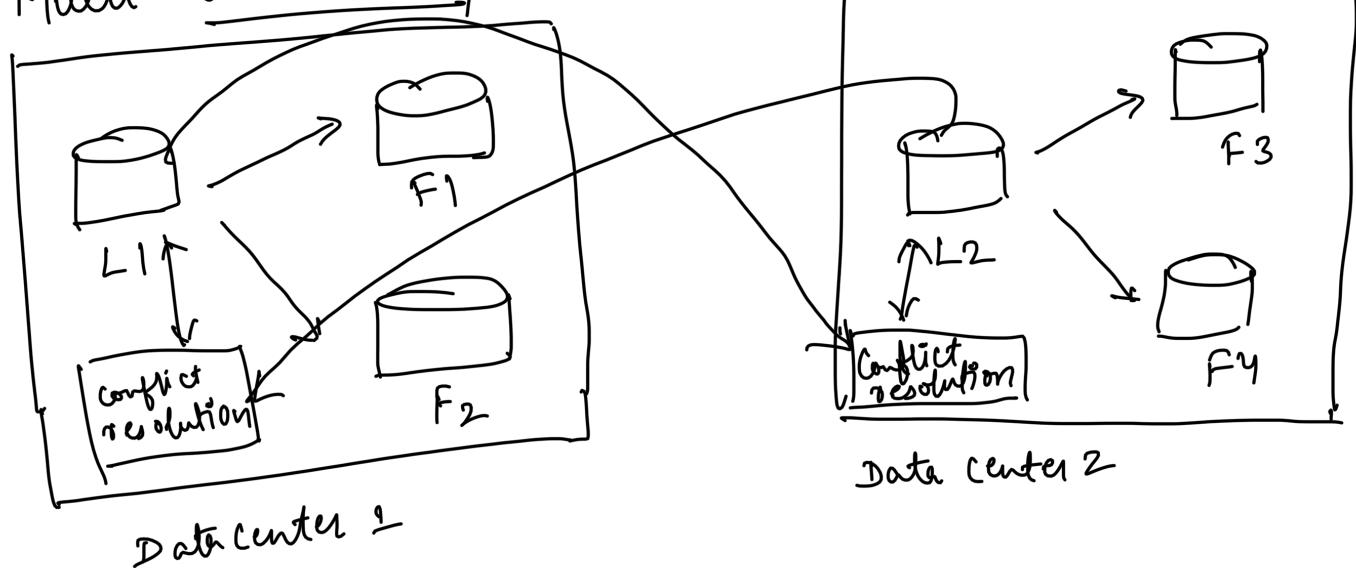
if writes are performed in certain order, reads should happen in same order. → violation of causality.  
→ generally happens in partitioned databases.

## Multi Leader Replication

↳ write can happen on any of leader nodes.

↳ Each leader simultaneously acts as follower to other leaders.

### ① Multi data center operation



### Benefits

- ① latency improvement - write can happen in specific data center.
- ② fault tolerant → data center outage.  
↳ network issues

### Downsides

- ① same write can happen in both data centers concurrently.  
↳ requires conflict resolution.

## Examples

- ↳ Clients with offline operation - Calendar application on multi-devices.
- ↳ Collaborative editing - Google docs - conflict resolution for concurrent users.

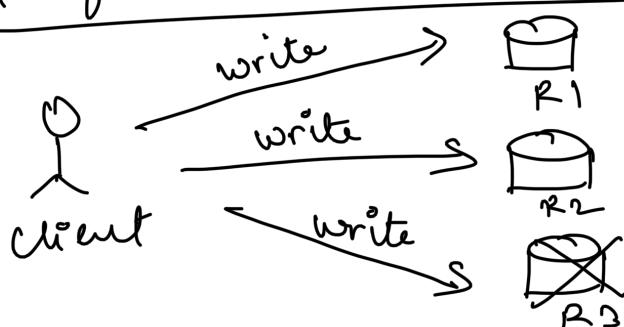
## Handling of write conflicts

- ① Conflict Avoidance → avoid conflicts to occur.
  - ↳ requests from particular user are always routed to same leader
    - ↳ single reader in user's view
- ② Converging to Consistent State
  - ↳ assign each write a unique ID. example Timestamp. Conflicts can be resolved basis "Last Write Wins".
  - ↳ assign each replica unique ID. writes from higher number replica takes precedence over other.
- ③ Custom Application Code
  - ↳ write your custom conflict handler.

## Leaderless Replication

- ↳ any replica can accept writes from clients.  
ex - Cassandra, Voldemort

### What if a node is down on write?



- ↳ only R1 & R2 are available
- ↳ write is successful basis quorum
- ↳ on Read as well, request is sent to all nodes, and then recent value can be picked up basis version no.

### How will be R3 in sync once it comes online?

- ① Read Repair → on read, we identified which Replica has stale data, so write back to R3.
- ② Anti-entropy process → background process to check differences between replicas.

## Quorums

n-replicates

w - no of nodes to confirm on write for it to be success  
r - no of nodes to be queried for read.

$$w+r > n$$

We'll get up-to-date reads.

## General recommendation

$n \rightarrow$  odd number

$$w.r \rightarrow (n+1)/2$$

- ① If you set lower values of  $r$  and  $w$  ( $r \neq w \leq n$ )  $\rightarrow$  possibility of stale reads.

- ② even with ( $w+r > n$ ), there is possibility of stale reads.

- ② Sloppy quorum  $\rightarrow$

↳ In case of network disruptions

↳ will you return errors?

↳ or return response (maybe stale data) without reaching quorum.

 Sloppy grammar

greeds and writes still require  $r$  and  $w$  successful responses but the response might not be from designated node.

→ return temporary response and once node is back, restore it to correct state.

- ⑥ Two concurrent writes →

Two concurrent writes multiple nodes will have confusion around which is most recent write.

- ④ Last write wins → each replica stores most recent value and overwrites older values.

- ↳ as we say writes are concurrent, so write order is undefined. Attach timestamps to each write and choose most recent one.
- ↳ This might lead to data lose.
- ↳ Approach for handling this could be to assign version numbers to each write operation and subsequently merge basis that on re-writes.
- ↳ The old values are deleted directly, kept in database with help of delete markers - called tombstones.

Hope you enjoyed the chapter Summary.

Subscribe for more such content.

YT channel - Ms Deep Singh

Happy Learning ☺