

Designing Data Intensive Applications

CHAPTER - 9

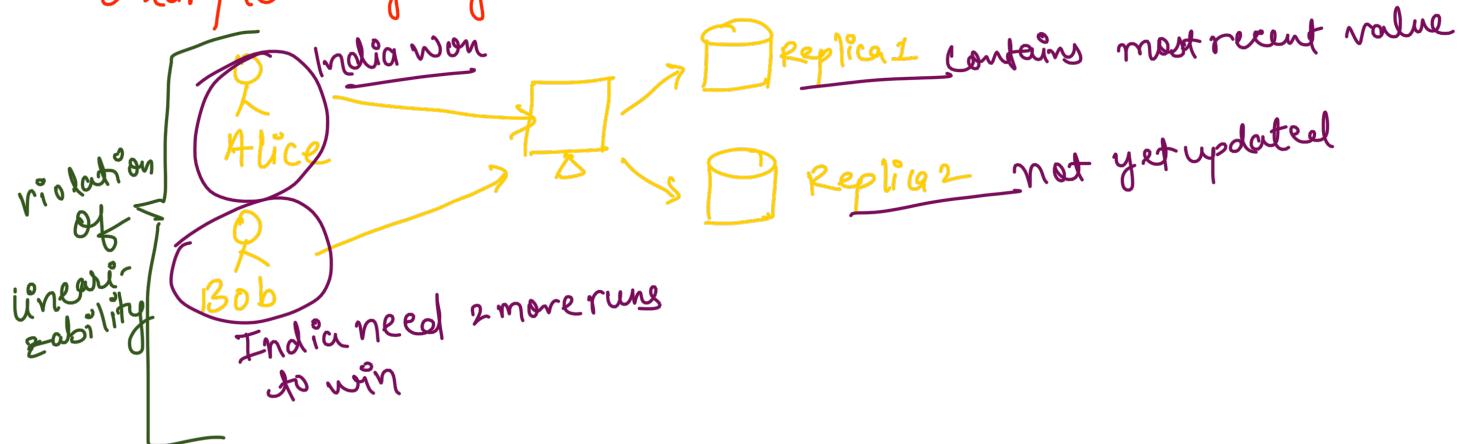
Consistency and Consensus

Subscribe for more such content . . .

YouTube - Ms Deep Singh

Linearizability → There could be multiple replicas in database but to client it appears like a single node. provide strong/atomic consistency.

Example - ongoing cricket match



Making System Linearizable

$x \rightarrow$ current score of India

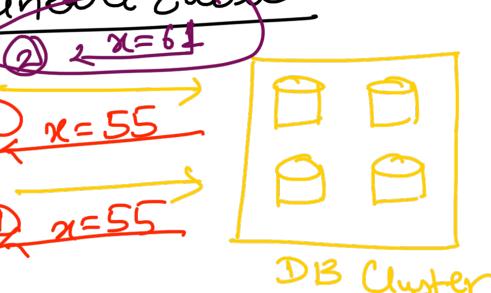
$$x = 55$$

↓

$$x = 61$$

If updated value is sent in response to one client, then Clients all clients should be sent updated score.

Even if write is in progress.



Comparing Linearizability with Serializability

Freshness guarantee on reads

Consistency property of transaction

and writes of an object.

Transaction behaves as same if executed sequentially.

why linearizability?

① Locking and leader election

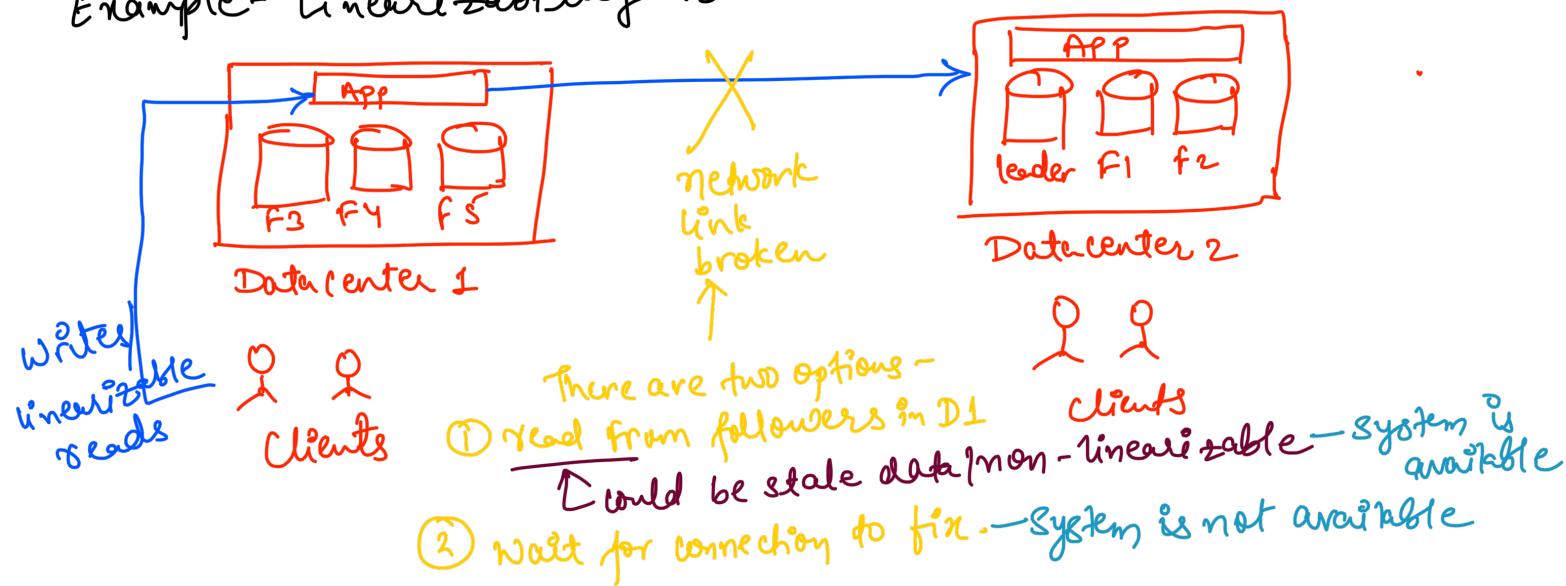
↳ for single leader replication - should only be one leader.
not more → split brain
a node can acquire lock on becoming leader.
linearizable. Example Zookeeper, etcd.

② Uniqueness Guarantees → Example username should be unique for user. should be linearizable operation.

Implementation

- ↳ Keep only single copy → operation will be atomic by default.
System is not fault tolerant.
↳ Consensus Algorithms → can be used to implement linearizable system

Example - Linearizability vs Availability. → Single leader replication



CAP Theorem → Choose either 2 of Consistency, Availability and Partition Tolerance. Impractical, CA system is not possible.

Linearizability and network delays

- ↳ It is generally avoided (if our use-case permits) to improve performance.
It is slow even if there are no network faults.

Ordering Guarantees → operations are executed in particular order

① Ordering and Causality

↳ ordering helps in preserving causality.

↳ Example

② Consistent Prefix Reads

→ question should be asked before so as it can be answered. → effect

Causality impose ordering in events

↑ cause comes first, then effect.

↳ system obeys this ordering
↳ referred as causally consistent.

Linearizability vs Causality

Total order of operations

allows two elements to be compared. e.g. $13 > 5$

we can always say which occurred first.

Database abstracts out that there are no concurrent events.

Partial order

The elements can be incomparable or one could be greater than other.

e.g. $\{a,b\}$ and $\{c,d\}$

↳ two operations are ordered if they're causally related.

↳ one happened before other.

↳ incomparable if concurrently occurred.
Example Git version control system.

↳ Any linearizable system will preserve causality correctly.

Stronger than causal system
hence performance can be slow.

Sequence Number Ordering

↳ maintaining all causal dependencies would be large overhead.

→ for event ordering, sequence no or timestamps can be used.

↳ from logical clock.

assign seq no to each operation → use counters for each operation.
operation & their no can be compared.

Sequence number for non-causal example multi leader replication

① Each node generate its own numbers.

node should have some identifier so as numbers are unique across cluster
eg one node generate odd and other even.

② Attach Timestamp → used in Last write Wins algorithm.

③ Assign range of numbers to each node.

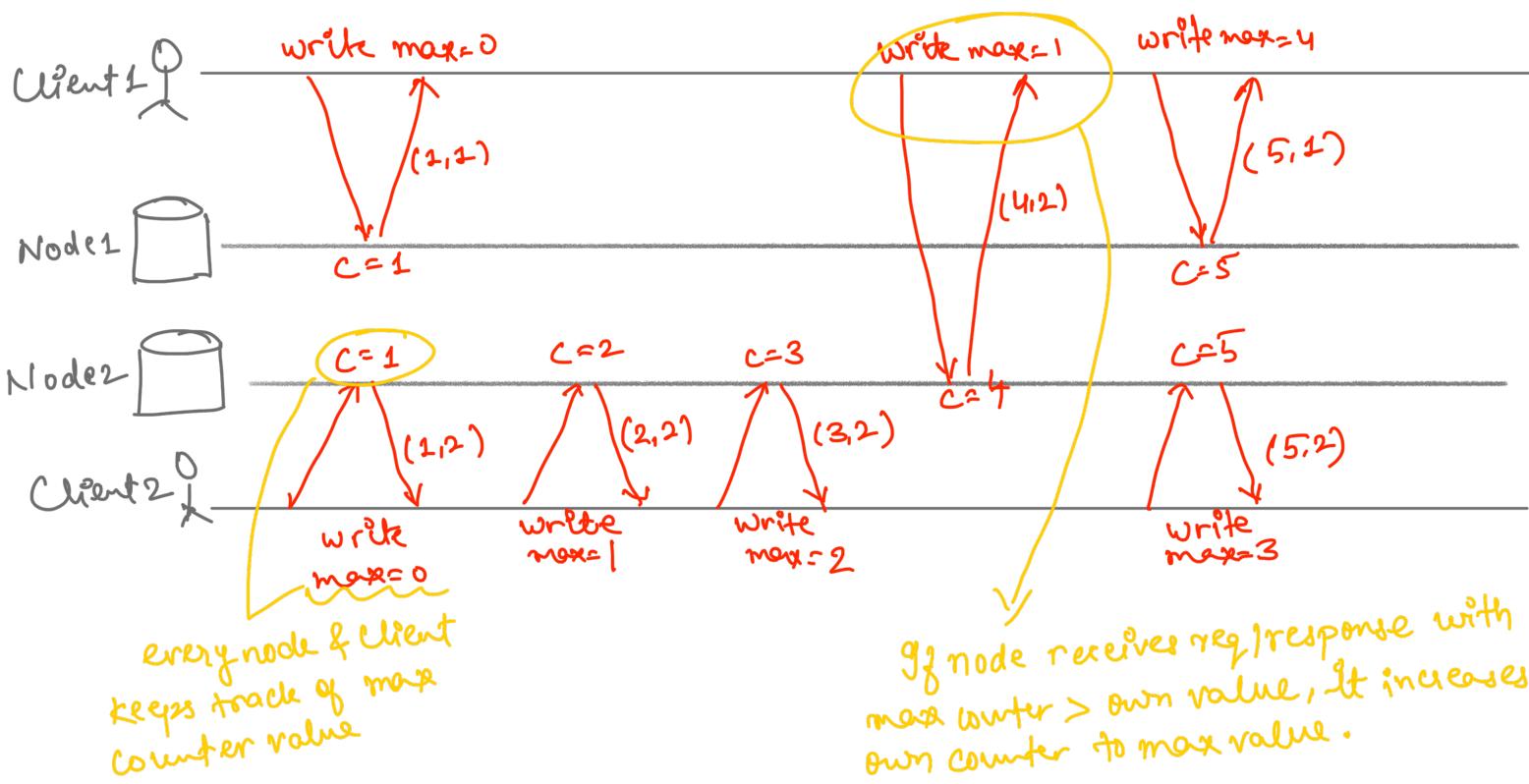
Generated seq numbers are not causally consistent:
they are not capturing order of events
across the nodes.

Solution - Lamport Timestamps

{counter, nodeId}

one with greater timestamp is greater

if same time, greater nodeId is greater



The Same username problem

Two users trying to create account with same username at same time.
 ↳ Username creation operation is not collected yet, Lamport timestamp will not help.

↳ A node should check with all other nodes if any request with same username.

↑ What if network fault between node communication.

To uniquely identify usernames, total ordering of operations
is not sufficient

↳ you should when total order is finalized.

Total order Broadcast.

How can all nodes in cluster agree on same total ordering of operation

for single leader database? - one node as leader handles this.

what if leader fails?

Total order Broadcast

→ also used in consensus algorithms.

↳ protocol for exchanging msgs between nodes.

→ Reliable delivery

→ Totally ordered delivery → delivered to every node in same order

→ If msg is already delivered, node is not allowed to insert it back in earlier position.

Distributed Transactions and

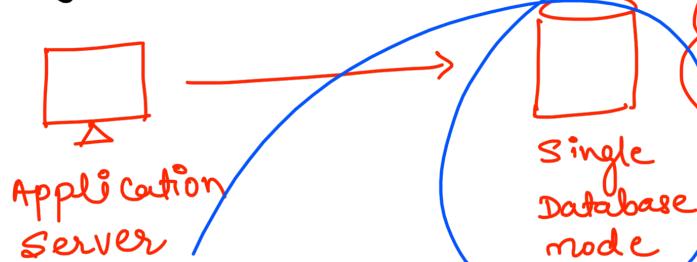
Consensus

→ Several nodes to agree on something

e.g. leader selection

Atomic Commit and Two Phase Commit (2PC)

① Single node Commit



On transaction Commit
① Update WAL on disk
② append commit record to log on disk

what if DB crashes before #1 or after #2?
↳ rollback changes

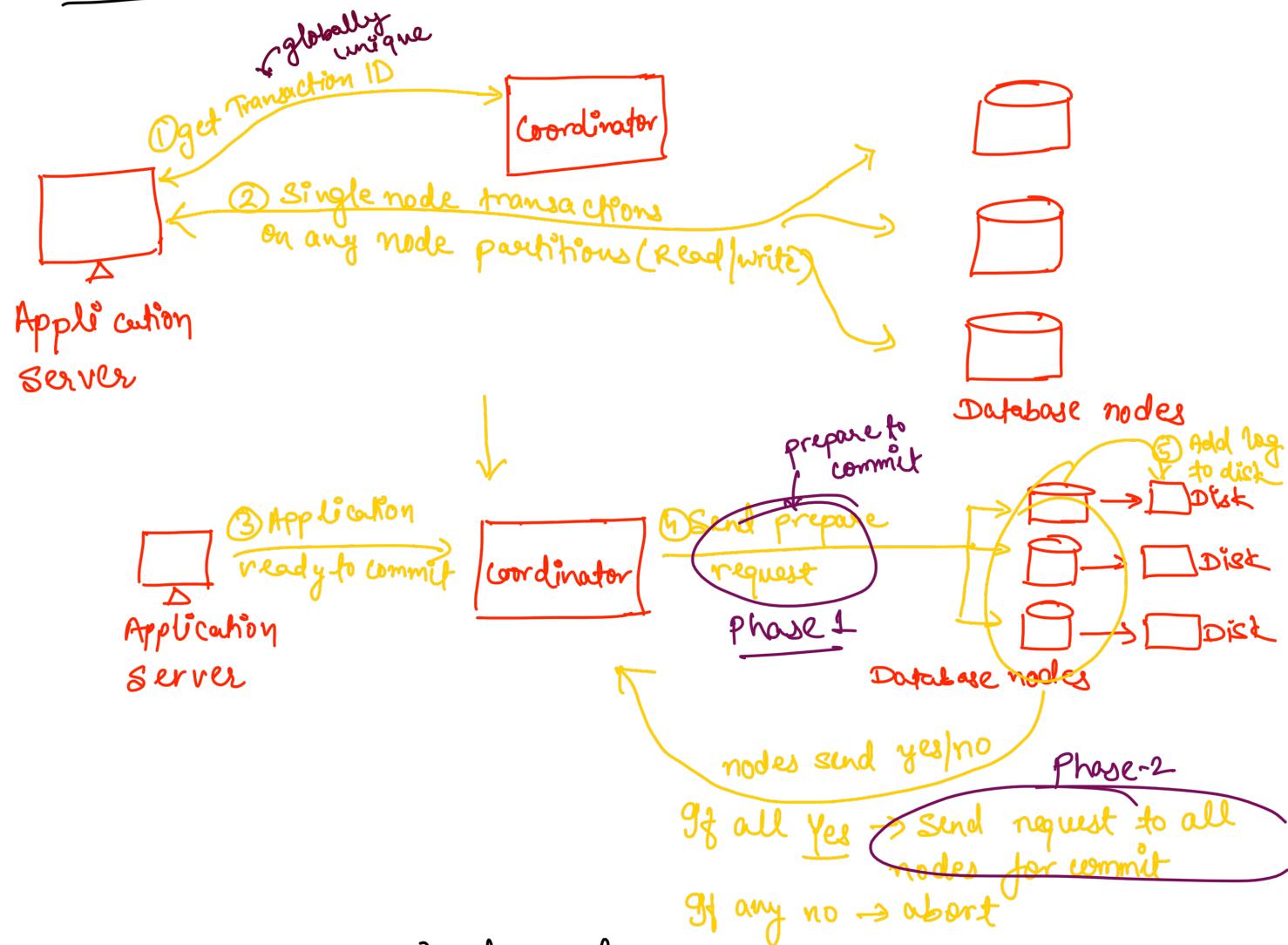
what if crashed after #2?
↳ recover from disk once node restarts.

what if multiple Nodes?

↳ we can't send commit msgs to all nodes. → what if it fails on few nodes? → network delay, node crash, etc.

Difficult to handle.

Two Phase Commit (2PC) → the commit/abort is split into two phases.



what if coordinator is down?

↳ Coordinator is down after sending commit request to one node → other nodes will not know what to do with transaction but wait.

① Either database nodes interact with each to decide what to do with specific transaction.

② Coordinator dumps all actions to disk
So it can recover once back online.

↳ what if disk is also corrupt?
↳ manual intervention by admin.

Distributed Transaction Examples

① Database Internal → internal transactions among the nodes in cluster.

② Heterogeneous distributed transactions

↳ different technologies
↳ two databases from different vendors, message broker & database

extended Architecture

XA Transactions → 2PC Implementation across heterogeneous tech

SC API for interfacing with transaction coordinator.

In Java - Java Transaction API → network driver or client library

The driver can send required info to coordinator (prepare | commit | abort)

Fault Tolerant Consensus → safety properties

- ① Uniform agreement among nodes → no two nodes decide differently → consensus algo properties
- ② Integrity → node decides only once
- ③ Validity → nodes decides value v , then v was proposed by some node.
- ④ Termination → every node decides a value that doesn't crash → system should reach a decision even if node crashes.
↳ Liveness property

Single Leader Replication and consensus

How is leader chosen?
in case of failover → manually by operation team
or automatic selection

Can split brain problem occur?

Consensus algo are based on total order broadcast.
like a single leader replication
and it requires a leader

We need all nodes agree for leader.

↳ liveness

Spiral loop

Solution → Epoch numbering for quorum

Protocols define epoch number → within each epoch, leader is unique.
If leader dead → election is given an incremented epoch number → increasing
monotonic
helpful in resolving conflicts.

Coordination Services eg. zookeeper, etc

- ① Linearizable atomic operations → for concurrent operation on nodes, only one succeed using atomic compare and set operation.

consensus

② Total ordering of messages

↳ fencing tokens are used to avoid conflicts on lock lease.

↳ Zookeeper uses monotonically increasing transaction ID & version number

③ Failure Detection → via heartbeats.

④ Change notifications → new nodes added / removed

See you in next video —

Happy learning ☺

Don't forget to subscribe . . ~