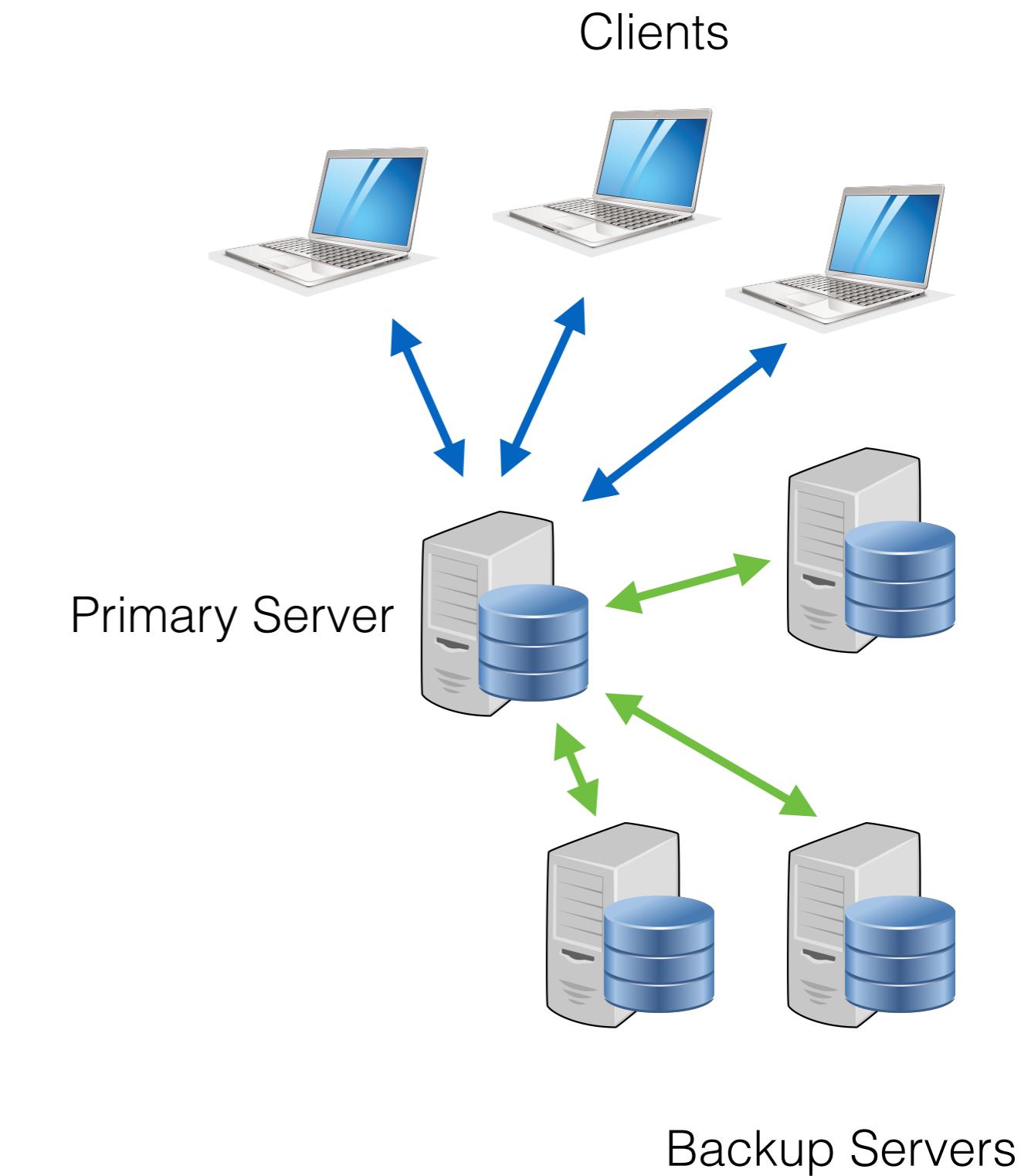
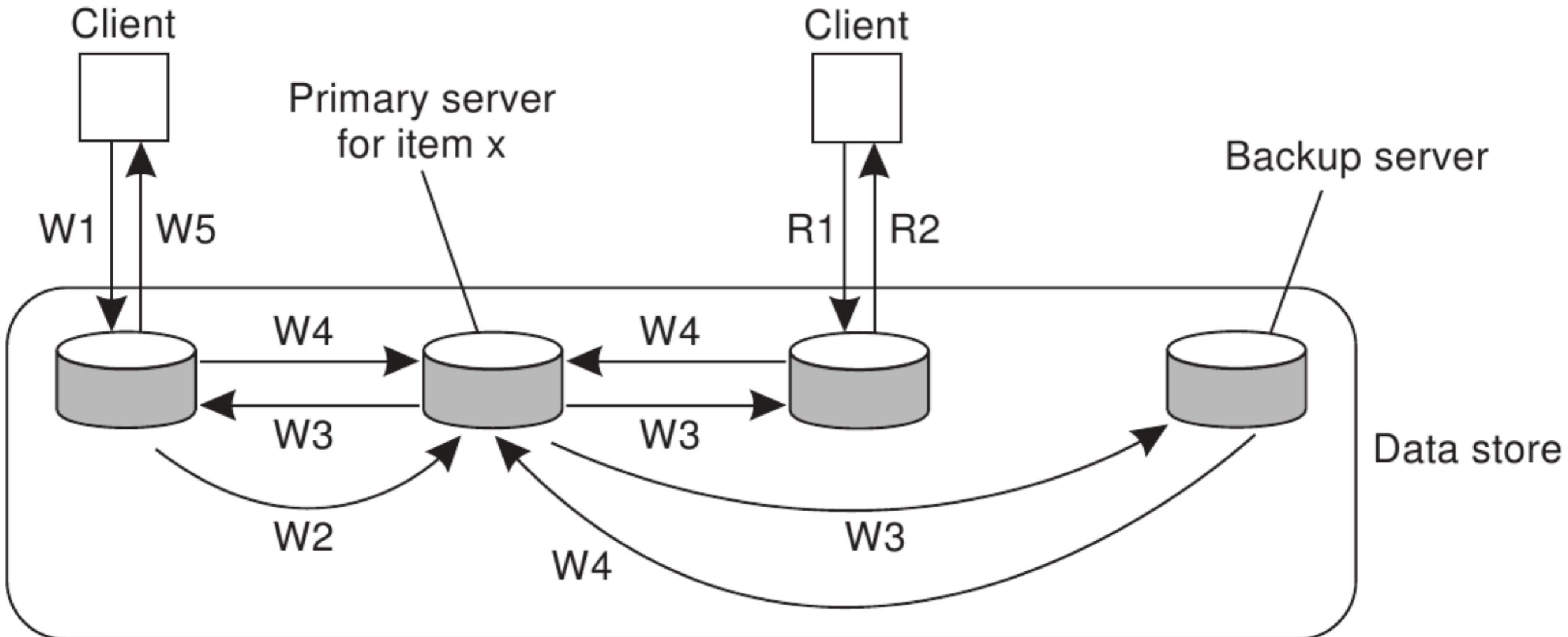


# Passive Replication

- Clients communicate with primary server
- WRITES are atomically forwarded from primary server to backup servers
- READS are replied by the primary server
- Also known as Primary Copy (or Backup) Replication
- Specifications:
  - At most one replica can be the primary server at any time.
  - Each client maintains a variable  $L$  (leader) that specifies the replica to which it will send requests. Requests are queued at the primary server.
  - Backup servers ignore client requests.



# Passive Replication Protocol



W1. Write request

W2. Forward request to primary

W3. Tell backups to update

W4. Acknowledge update

W5. Acknowledge write completed

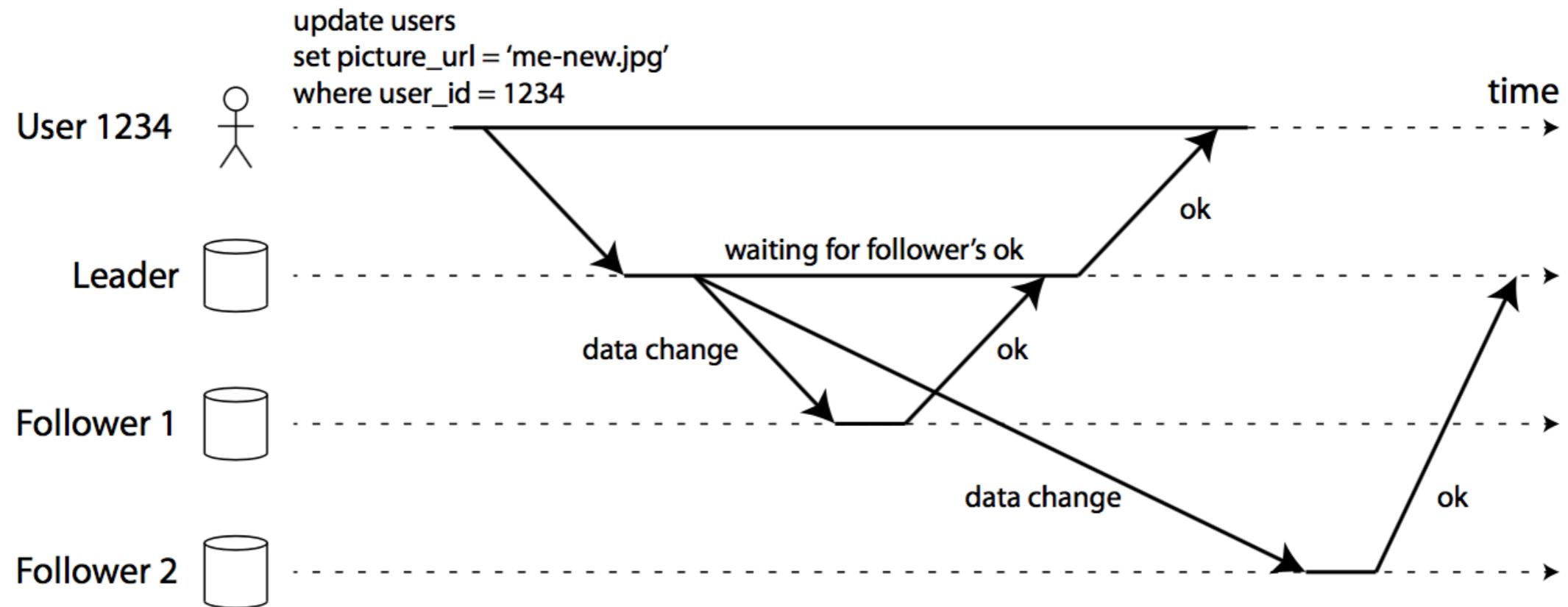
R1. Read request

R2. Response to read

# Request Phases

- **Request:** The front end issues the request, containing a unique identifier, to the primary replica manager.
- **Coordination:** The primary takes each request atomically, in the order in which it receives it. It checks the unique identifier, in case it has already executed the request, and if so it simply resends the response.
- **Execution:** The primary executes the request and stores the response.
- **Agreement:** If the request is an update, then the primary sends the updated state, the response and the unique identifier to all the backups. The backups send an acknowledgement.
- **Response:** The primary responds to the front end, which hands the response back to the client.

# Synchronous vs Asynchronous



- The **advantage** of synchronous replication is that the follower is guaranteed to have an up-to-date copy of the data that is consistent with the leader.
  - If the leader suddenly fails, we can be sure that the data is still available on the follower.
- The **disadvantage** is that if the synchronous follower doesn't respond (because it has crashed, or there is a network fault, or for any other reason), the write cannot be processed.
  - The leader must block all writes and wait until the synchronous replica is available again.

# Primary Failure

- When the primary replica fails, a failover procedure is required (can be manual or automatic)
  1. Determining that the leader has failed.
  2. Choosing a new leader.
  3. Reconfiguring the system to use the new leader.
- Issues:
  - If asynchronous replication is used, the new leader may not have received all writes from the old leader before it failed.
  - Discarding writes is especially dangerous if other storage systems outside of the database need to be coordinated with the database contents.
  - It could happen that two nodes both believe that they are the leader

# Replication Log

- The primary replica stores all changes locally, in a **replication log**
- Applying the replication allows us to perform the correct update on any object (given the correct **log sequence number**)
- This is used so set up **new backup replicas**, given a snapshot of the objects, the corresponding log sequence number, and the primary's replication log.
- The same holds for catch-up recovery of failing backup replicas.

# Implementing Replication Log

- **Statement Log**: the primary logs every write request (*statement*) that it executes, and sends that statement log to its followers.
  - *Potential issues*: non-deterministic values (rand()), concurrency issues, side effects on other components
- **Write-Ahead Log**: similarly to LSM trees, the primary append every write requests in the log, and sends the whole sequence of write requests
  - *Potential issues*: the log describes the data on a very low level: a WAL contains details of which bytes were changed in which disk block. What if we update something?

# Simple Protocol

- System model:
  - point-to-point communication
  - no communication failures → no network partitions
  - upper bound on message delivery time → synchronous communications
  - FIFO channels
  - at most one server crashes
- Two servers:
  - The primary  $p_1$
  - The backup  $p_2$
- Variables:
  - At server  $p_i$ ,  $\text{primary} = \text{true}$  if  $p_i$  acts as the current primary
  - At clients,  $\text{primary}$  is equal to the identifier of the current primary

# Simple Protocol

Protocol executed by the primary  $p_1$

upon initialization do

```
└ primary ← true
```

upon receive  $\langle \text{REQ}, r \rangle$  from  $c$  do

$state \leftarrow \text{update}(state, r)$

% Update local state

**send**  $\langle \text{STATE}, state \rangle$  to  $p_2$

% Send update to backup

send  $\langle \text{REP}, \text{reply}(r) \rangle$  to  $c$

% Reply to client

**repeat** every  $\tau$  seconds

| send ⟨HB⟩ to  $p_2$

% Heartbeat message

**upon recovery after a failure do**

| { start behaving like a backup }

# Simple Protocol

---

Protocol executed by the backup  $p_2$

---

**upon** initialization **do**

$\lfloor primary \leftarrow \text{false}$

**upon receive**  $\langle \text{STATE}, s \rangle$  **do**

$\lfloor state \leftarrow s$  % Update local state

**upon** not receiving a heartbeat for  $\tau + \delta$  seconds **do**

$\lfloor primary \leftarrow \text{true}$  % Becomes new primary

**send**  $\langle \text{NEWP} \rangle$  **to**  $c$  % Inform the client of new primary

$\lfloor \{ \text{start behaving like a primary} \}$

---

# Simple Protocol

---

Protocol executed by client  $c$

---

**upon** initialization **do**

$primary \leftarrow p_1$  % Initial primary

**upon receive**  $\langle \text{NEWP} \rangle$  **from**  $p_2$  **do**

$primary \leftarrow p_2$  % Backup

**upon** operation( $r$ ) **do**

**while not received a reply do**

**send**  $\langle \text{REQ}, r \rangle$  **to**  $primary$

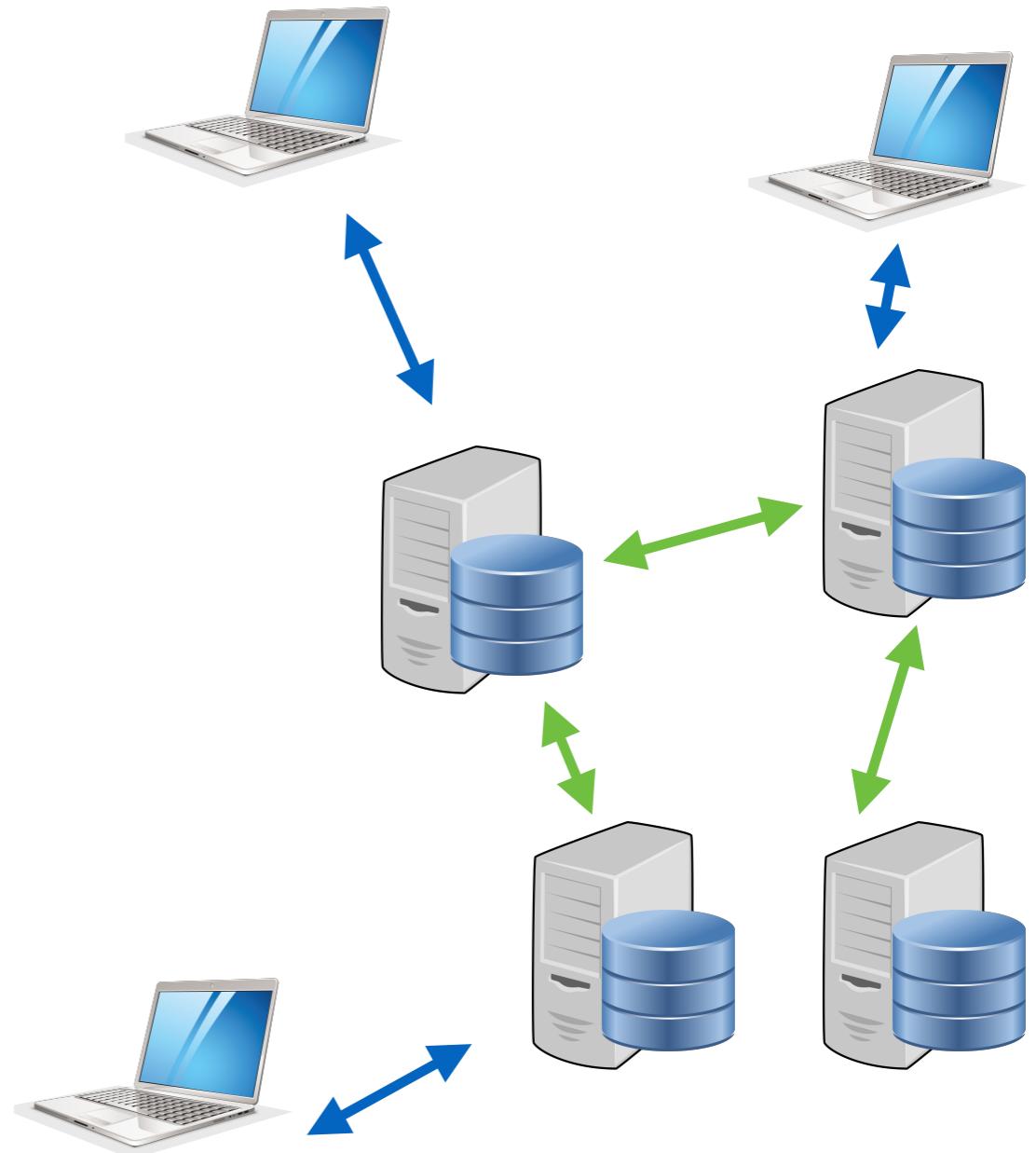
**wait receive**  $\langle \text{REP}, v \rangle$  **or receive**  $\langle \text{NEWP} \rangle$

**return**  $v$

---

# Active Replication

- a.k.a. multi-leader replication
- Clients communicate with several/all servers
- Every server handles any operation and sends the response
- WRITES must be applied in the same order (**total order broadcast**)
- One way to implement totally-ordered multicast is to use logical clocks



# Use Cases

- When does it make sense to use active replication?
  - multi-datacenter operations (increased performance and fault tolerance w.r.t. passive replication)
- Clients with offline operations
  - Each client device is a 'datacenter'
- Collaborative editing
  - Each copy is a 'datacenter'

# Request Phases

- **Request:** The front end attaches a unique identifier to the request and multicasts it to the group of replica managers, using a totally ordered, reliable multicast primitive.
- **Coordination:** The group communication system delivers the request to every correct replica manager in the same (total) order.
- **Execution:** Every replica manager executes the request. Since they are state machines and since requests are delivered in the same total order, correct replica managers all process the request identically. The response contains the client's unique request identifier.
- **Agreement:** No agreement phase is needed, because of the multicast delivery semantics.
- **Response:** Each replica manager sends its response to the front end. The number of replies that the front end collects depends upon the failure assumptions and the multicast algorithm.

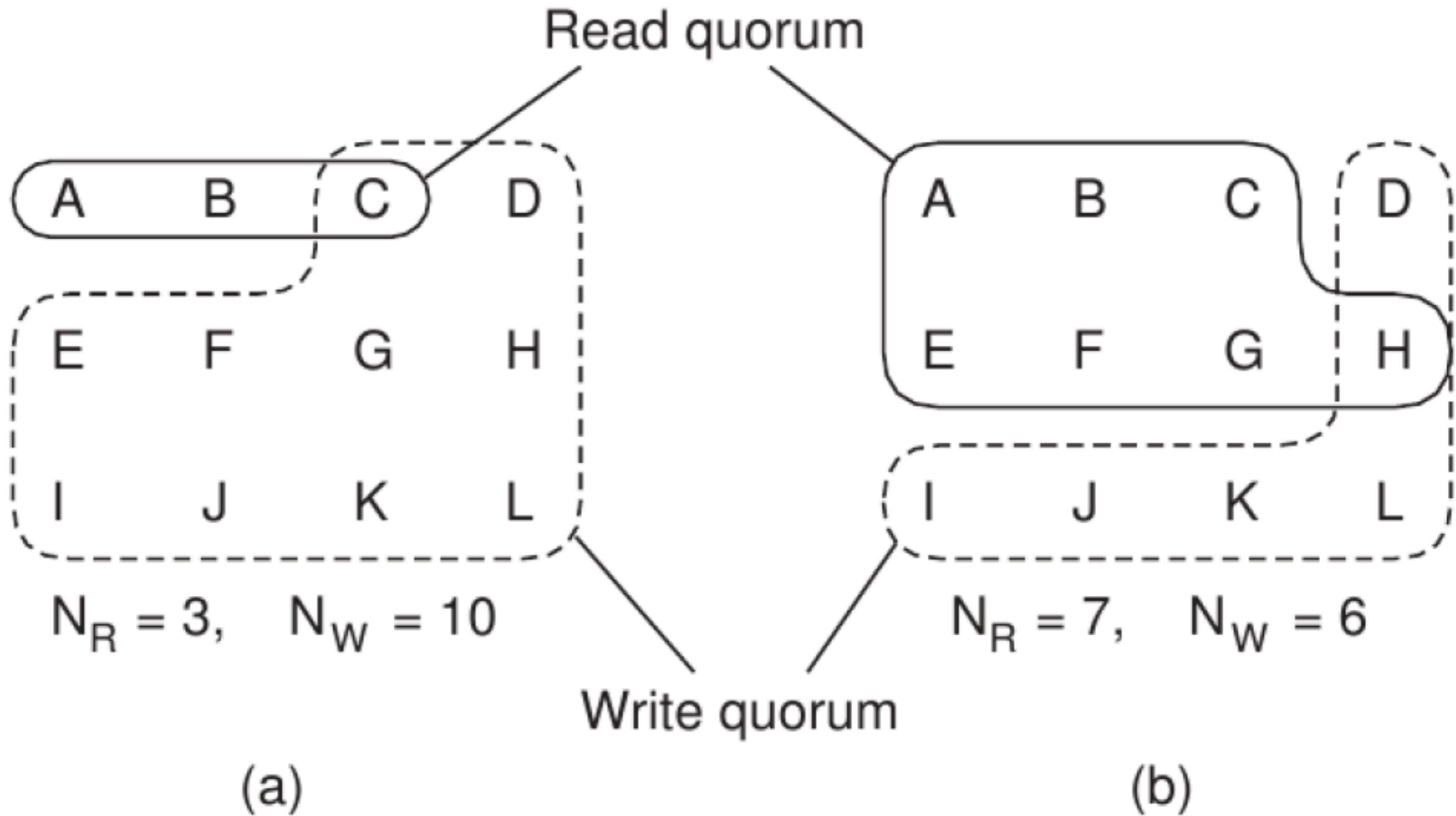
# Quorum Protocols

- Proposed by Gifford in 1979
- **Quorum-based protocols** guarantee that each operation is carried out in such a way that a *majority vote* (a quorum) is established.
  - *Write quorum W*: the number of replicas that need to acknowledge the receipt of the update to complete the update
  - *Read quorum R*: the number of replicas that are contacted when a data object is accessed through a read operation

# Quorum Systems

- Formally, a **quorum system**  $S = \{S_1, \dots, S_N\}$  is a collection of **quorum sets**  $S_i \subseteq U$  such that two quorum sets have at least an element in common
- For replication, we consider two quorum sets, a **read quorum**  $R$  and a **write quorum**  $W$
- **Rules:**
  1. Any read quorum must overlap with any write quorum
  2. Any two write quorums must overlap
- $U$  is the set of replicas, i.e.,  $|U| = N$

# Quorum Examples



# Quorum Examples

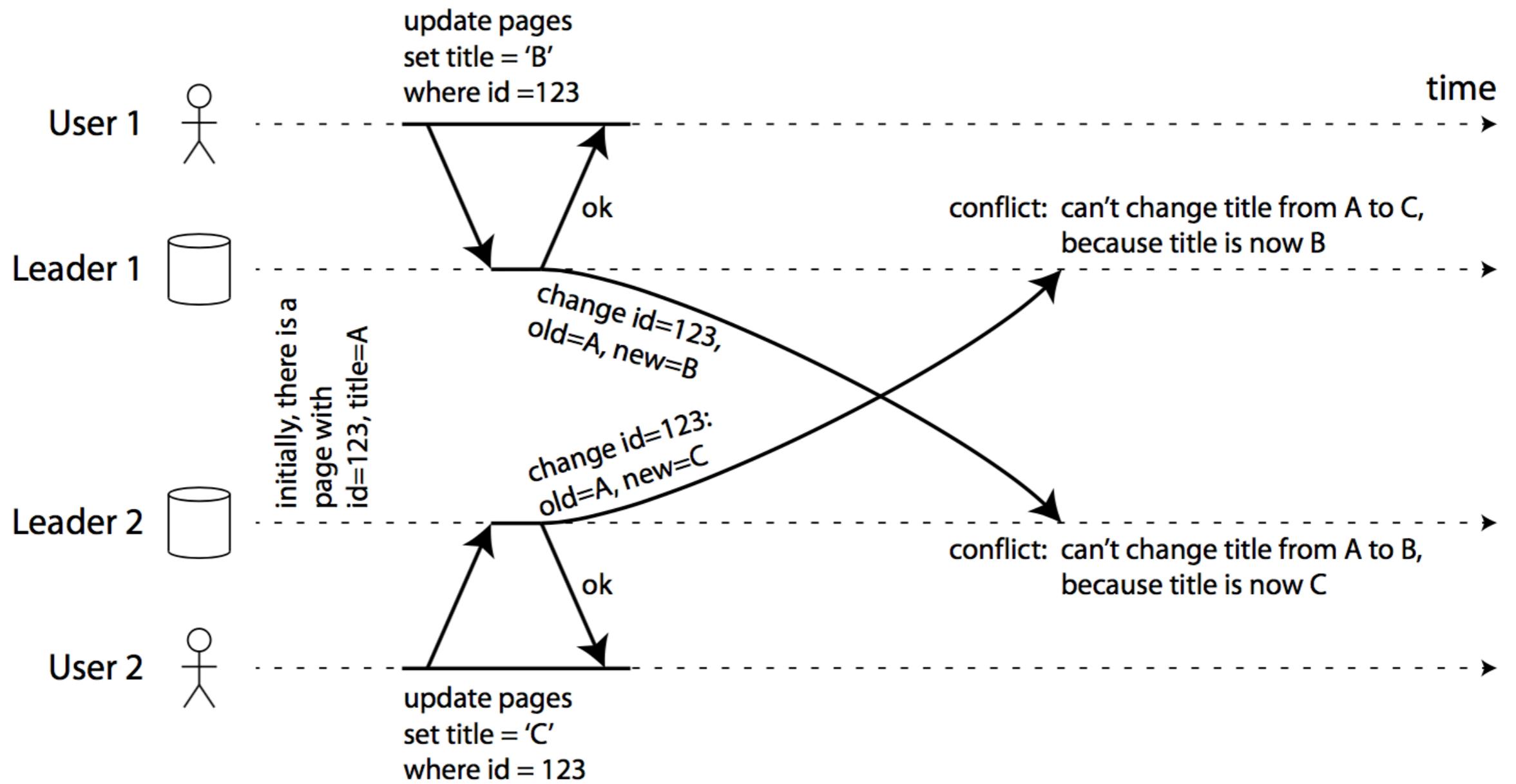
- Read rule:  $|R| + |W| > N \Rightarrow$  read and write quorums overlap
- Write rule:  $2|W| > N \Rightarrow$  two write quorums overlap
- The quorum sizes determine the costs for read and write operations
- Minimum quorum sizes for are

$$\min |W| = \left\lfloor \frac{N}{2} \right\rfloor + 1 \quad \min |R| = \left\lceil \frac{N}{2} \right\rceil$$

- Write quorums requires majority
- Read quorum requires at least half of the nodes
- ROWA ( $R, W, N$ ) = ( $N = N$ ,  $R = 1$ ,  $W = N$ )
- Amazon's Dynamo ( $N = 3$ ,  $R = 2$ ,  $W = 2$ )
- LinkedIn's Voldemort ( $N = 2$  or  $3$ ,  $R = 1$ ,  $W = 1$  default)
- Apache's Cassandra ( $N = 3$ ,  $R = 1$ ,  $W = 1$  default)

# Write Conflicts

The biggest problem with active replication is that write conflicts can occur, which means that conflict resolution is required.



# Handling Write Conflicts

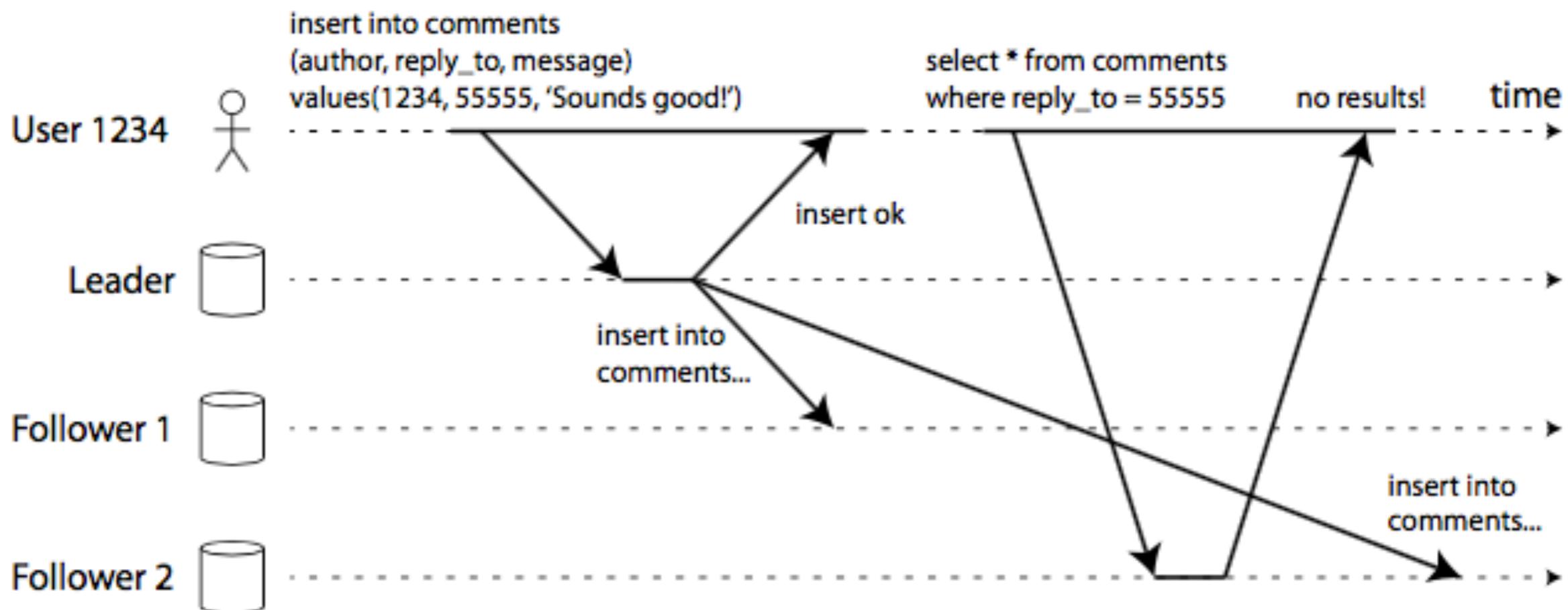
- Avoid them by '**normally**' using a **single leader**, and change leader for exceptional conditions only.
- **Converge** towards a consistent state
  - Give each write a unique ID, pick the write with the highest ID as the winner, and throw away the other writes. If a timestamp is used, this technique is known as **last write wins** (LWW). Although this technique is popular, it is dangerously prone to data loss.
  - Give each replica a unique ID, and let writes that originated at a **higher-numbered replica** always take precedence over writes that originated at a lower-numbered replica. This also implies data loss.
  - Record the conflict in an explicit data structure that preserves all information, and write application code which resolves the conflict at **some later time** (perhaps by prompting the user).
- Use custom logic
- Use automatic logic (e.g., conflict-free replicated data types, CRDTs)

# Replication Lag

- If an application reads from **asynchronous followers**, it may see **outdated information** if the follower has fallen behind.
- This leads to **apparent inconsistencies** in the database
  - if you run the same query on the leader and a follower at the same time, you may get different results, because not all writes have been reflected in the follower.
  - This inconsistency is just a **temporary state**
    - if you stop writing to the database and wait a while, the followers will eventually catch up and become consistent with the leader.
  - For that reason, this effect is known as **eventual consistency**.

# Read Your Writes Consistency

if the user views the data shortly after making a write, the new data may have not yet reached the replica.



# Client-centric Consistency Models

- **Each WRITE operation** is assigned a unique identifier
  - Done by the replica manager where the operation is requested
- For each replica manager, we keep track of:
  - Write set  $WS$  : contains the write operations executed so far
- For each client  $c$ , we keep track of:
  - Read set  $WS_R$  : contains write operations relevant to the read operations performed by  $c$
  - Write set  $WS_W$  : contains write operations relevant to the write operations performed by  $c$

# Read-Your-Writes Implementation

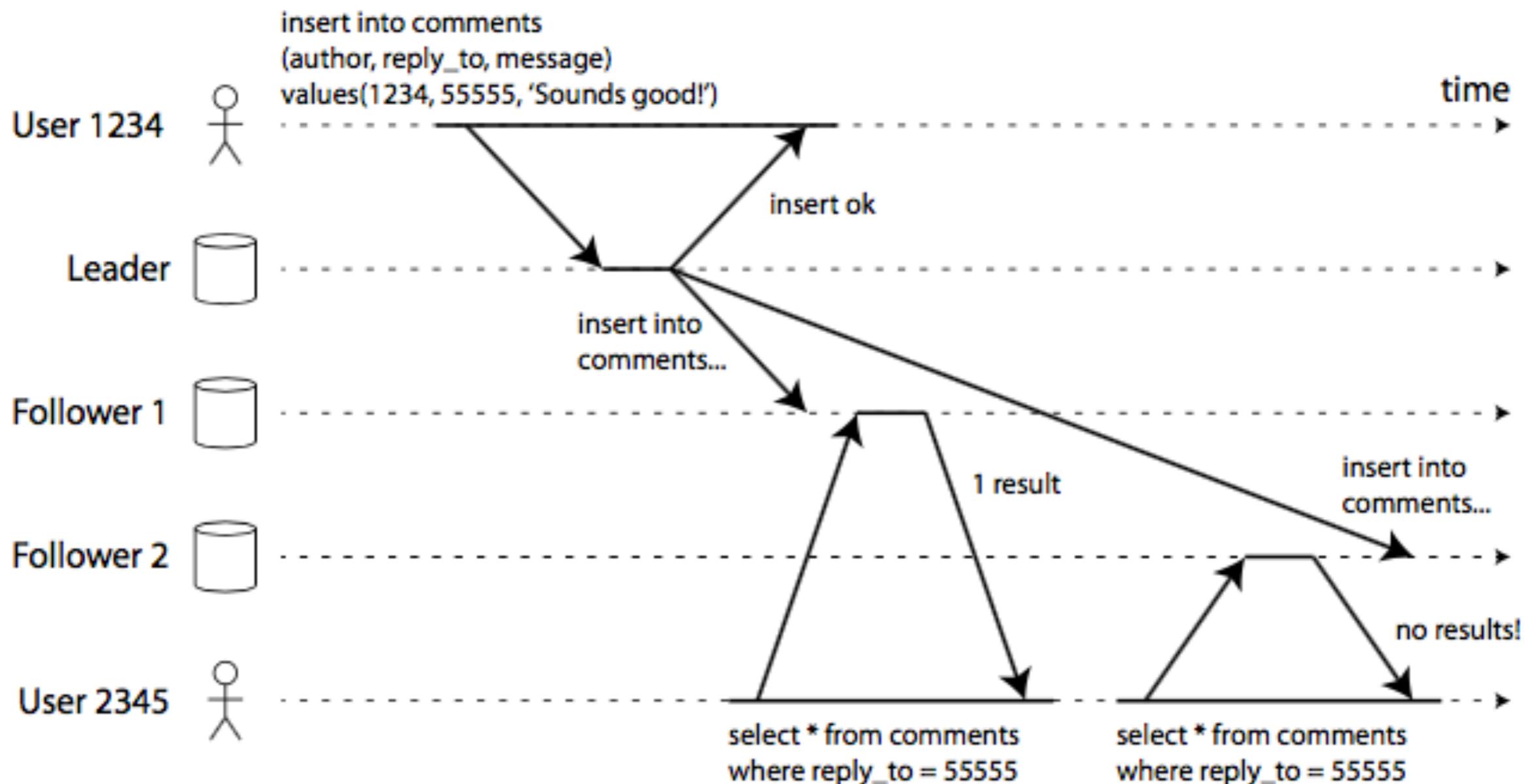
- To perform a READ:
  - A client
    - sends READ and its  $WS_w$  to a replica manager S.
  - The replica manager S:
    - Checks if the  $WS_w \subseteq WS$ , i.e., all the WRITES seen from the client have been applied by the replica manager
    - If not, asks the other replica managers the missing WRITES
    - Applies the missing WRITES locally and update its WS
    - Return the requested value to the client

# Read-Your-Write Implementation

- To perform a WRITE:
  - A client
    - sends WRITE and adds it to its  $WS_w$
  - The replica manager S:
    - Perform the WRITE
    - adds it to its WS

# Monotonic Reads Consistency

if a user makes several reads from different replicas, it's possible for a user to see things moving backwards in time.



# Monotonic-Read Implementation

- To perform a READ:
  - A client
    - sends READ and its  $WS_R$  to a replica manager S.
  - The replica manager S:
    - Checks if the  $WS_R \subseteq WS$ , i.e., all the WRITES seen from the client have been applied by the replica manager
    - If not, asks the other replica managers the missing WRITES
    - Applies the missing WRITES locally and update its WS
    - Return the requested value and WS to the client
  - The client
    - adds WS to its  $WS_R$

# Monotonic-Read Implementation

- To perform a WRITE:
  - A client
    - sends WRITE
  - The replica manager S:
    - Perform the WRITE
    - adds it to its WS

# Extra: Monotonic-Write

- In a monotonic-write consistent data storage system, a write operation by a client on a data item is completed before any successive write operation on the same object by the same client

# Extra: Writes Follow Reads

- In a writes-follows-reads consistent data storage system, a write operation by a client on a data item following a previous read operation on the same item by the same client is guaranteed to take place on the same or a more recent value of the item that was read

# Additional References

- D. Terry et al., *Session Guarantees for Weakly Consistent Replicated Data*,  
<https://www.cis.upenn.edu/~bcpierce/courses/dd/papers/SessionGuaranteesPDIS.ps>
- D. Terry, *Replicated Data Consistency Explained Through Baseball*, <http://research.microsoft.com/pubs/157411/ConsistencyAndBaseballReport.pdf>