

CS425 /CSE424/ECE428 – Distributed Systems – Fall 2011

Replica Management

Material derived from slides by I. Gupta, M. Harandi,
J. Hou, S. Mitra, K. Nahrstedt, N. Vaidya

Objective

- Understand replication management
 - Goals
 - Model
 - Group communication & views
 - Consistency
 - Passive vs. active
- Readings
 - §15.1–15.3 (4th ed)
 - §18.1–18.3 (5th ed)

Replication

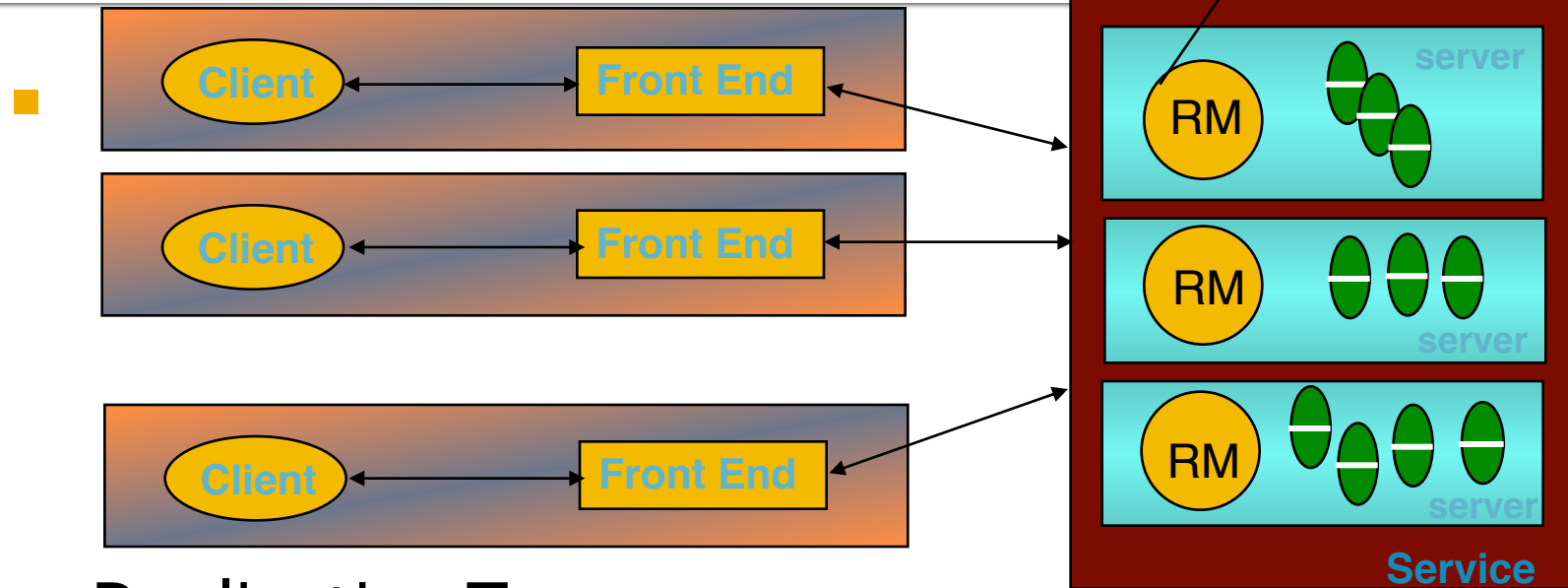
- Enhances a service by replicating data
 - Increased Availability
 - Of service. When servers fail or when the network is partitioned.
 - Fault Tolerance
 - Under the fail-stop model, if up to f of $f+1$ servers crash, at least one is alive.
 - *(later: Byzantine faults, survive f faults with $3f+1$ servers)*
 - Load Balancing
 - One approach: Multiple server IPs can be assigned to the same name in DNS, which returns answers round-robin.

P : probability that one server fails = $1 - P$ = availability of service. e.g. $P = 5\% \Rightarrow$ service is available 95% of the time.

P^n : probability that n servers fail = $1 - P^n$ = availability of service. e.g. $P = 5\%$, $n = 3 \Rightarrow$ service available 99.875% of the time

Goals of Replication

Replica Manager



- Replication Transparency
 - User/client need not know that multiple physical copies of data exist.
- Replication Consistency
 - Data is consistent on all of the replicas (or is converging towards becoming consistent)

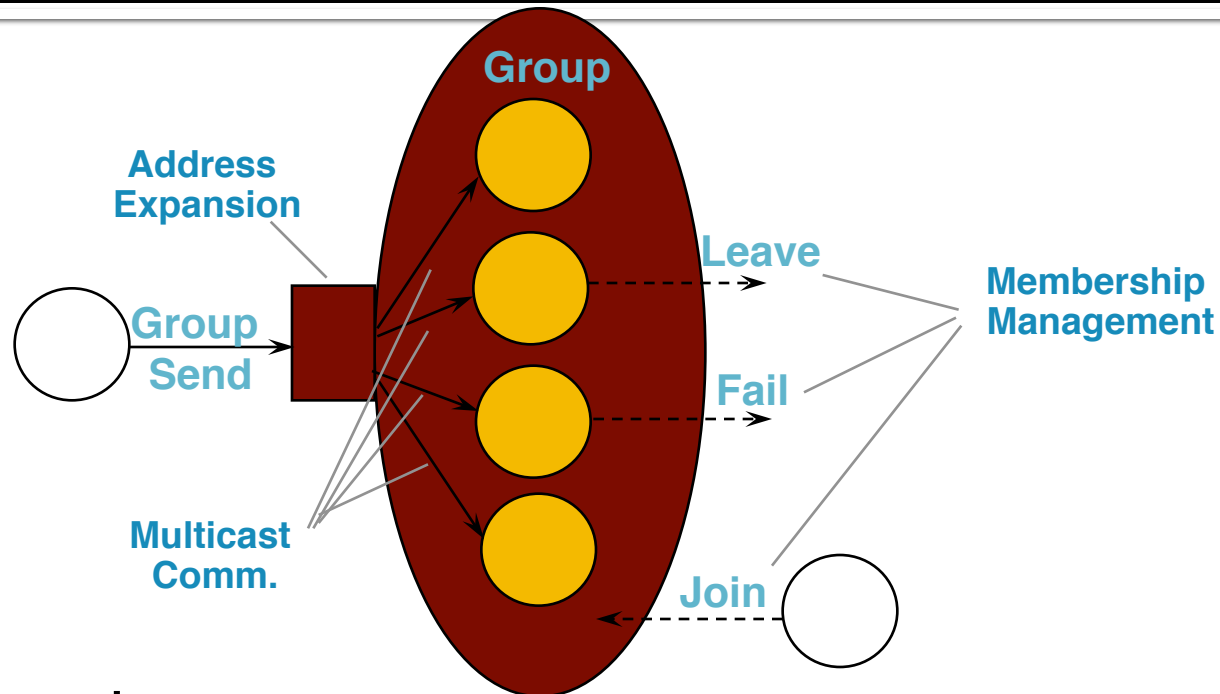
Replication Management

- Request Communication
 - Requests can be made to a single RM or to multiple RMs
- Coordination: The RMs decide
 - whether the request is to be applied
 - the order of requests
 - FIFO ordering: If a FE issues r then r' , then any correct RM handles r and then r' .
 - Causal ordering: If the issue of r "happened before" the issue of r' , then any correct RM handles r and then r' .
 - Total ordering: If a correct RM handles r and then r' , then any correct RM handles r and then r' .
- Execution: The RMs execute the request (often they do this tentatively – why?).

Replication Management

- **Agreement:** The RMs attempt to reach consensus on the effect of the request.
 - E.g., Two phase commit through a coordinator
 - If this succeeds, effect of request is made permanent
- **Response**
 - One or more RMs responds to the front end.
 - The first response to arrive is good enough because all the RMs will return the same answer.
 - Thus each RM is a **replicated state machine**
 - "Multiple copies of the same State Machine begun in the Start state, and receiving the same Inputs in the same order will arrive at the same State having generated the same Outputs." [Wikipedia, Schneider 90]

Group Communication: A building block



- "Member"= process (e.g., an RM)
- Static Groups: group membership is pre-defined
- Dynamic Groups: Members may join and leave, as necessary

Views

- A group membership service maintains group views, which are lists of current group members.
 - This is NOT a list maintained by a one member, but...
 - Each member maintains its own local view
- A view $V_p(g)$ is process p 's understanding of its group (list of members)
 - Example: $V_{p.0}(g) = \{p\}$, $V_{p.1}(g) = \{p, q\}$, $V_{p.2}(g) = \{p, q, r\}$, $V_{p.3}(g) = \{p, r\}$
 - The second subscript indicates the "view number" received at p
- A new group view is disseminated, throughout the group, whenever a member joins or leaves.
 - Member detecting failure of another member reliable multicasts a "view change" message (requires causal-total ordering for multicasts)
 - The goal: the compositions of views and the order in which the views are received at different members is the same. (i.e., view deliveries are "virtually synchronous")

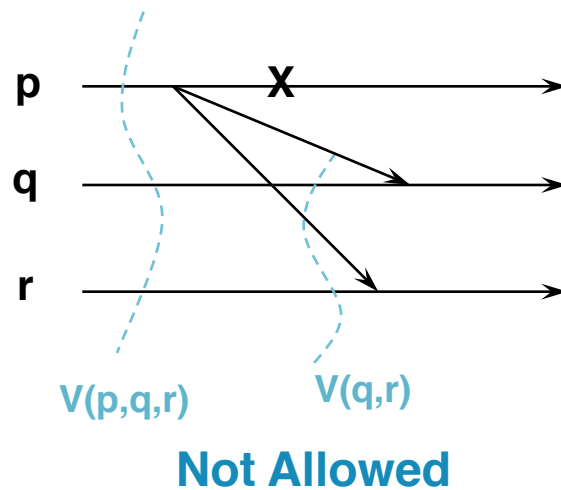
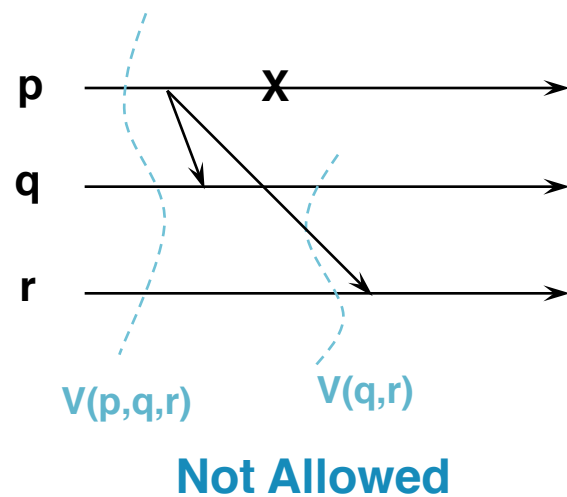
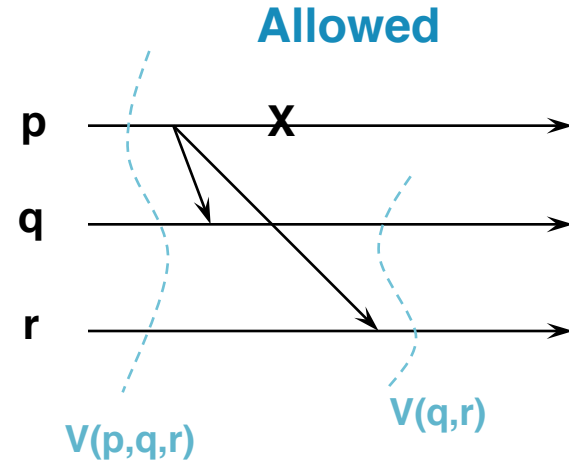
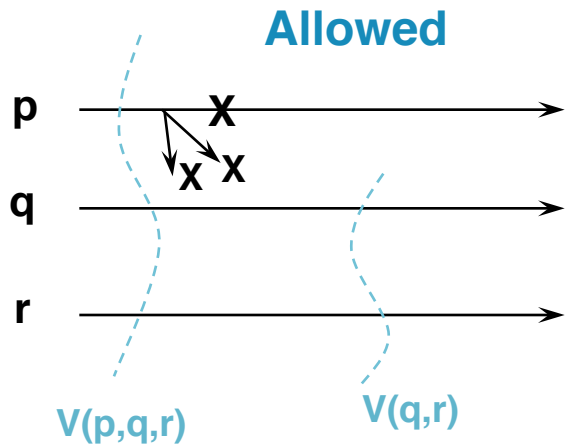
Views

- An event is said to occur in a view $v_{p,i}(g)$ if the event occurs at p , and at the time of event occurrence, p has delivered $v_{p,i}(g)$ but has not yet delivered $v_{p,i+1}(g)$.
- Messages sent out in a view i need to be delivered in that view at all members in the group ("What happens in the View, stays in the View")
- Requirements for view delivery
 - **Order**: If p delivers $v_i(g)$ and then $v_{i+1}(g)$, then no other process q delivers $v_{i+1}(g)$ before $v_i(g)$.
 - **Integrity**: If p delivers $v_i(g)$, then p is in all $v_{*,i}(g)$.
 - **Non-triviality**: if process q joins a group and becomes reachable from process p , then eventually, q will always be present in the views that delivered at p .
 - Exception: partitioning of group
 - We'll discuss partitions next lecture. Ignore for now.

View Synchronous Communication

- View Synchronous Communication = Group Membership Service + Reliable multicast
- The following guarantees are provided for multicast messages:
 - **Integrity**: If p delivered message m , p will not deliver m again. Also $p \in \text{group}(m)$, i.e., p is in the latest view.
 - **Validity**: Correct processes always deliver all messages. That is, if p delivers message m in view $v(g)$, and some process $q \in v(g)$ does not deliver m in view $v(g)$, then the next view $v'(g)$ delivered at p will not include q .
 - **Agreement**: Correct processes deliver the same sequence of views, and the same set of messages in any view.
 - if p delivers m in V , and then delivers V' , then all processes in $V \cap V'$ deliver m in view V
 - All View Delivery conditions (Order, Integrity and Non-triviality conditions, from last slide) are satisfied
- "What happens in the View, stays in the View"
- View and message deliveries are allowed to occur at different physical times at different members!

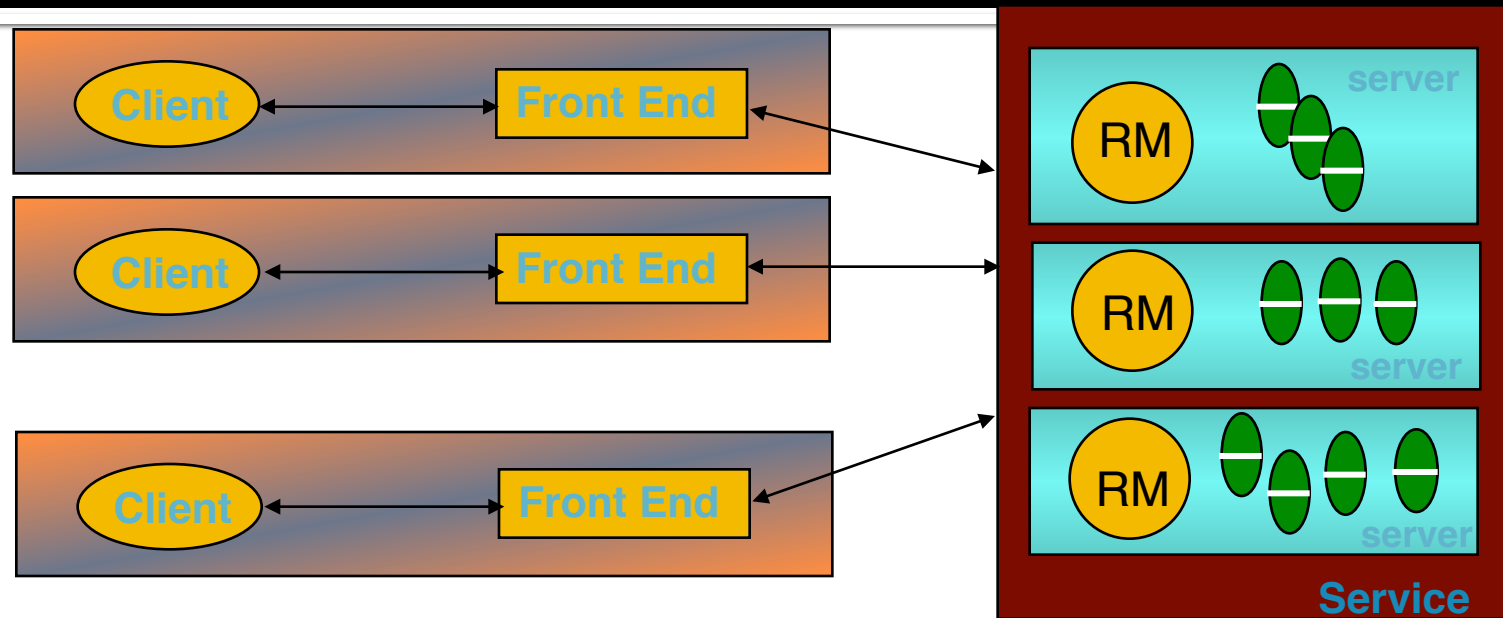
Example: View Synchronous Communication



State Transfer

- When a new process joins the group, state transfer may be needed (at view delivery point) to bring it up to date
 - "state" may be list of all messages delivered so far (wasteful)
 - "state" could be list of current server object values (e.g., a bank database) – could be large
 - Important to optimize this state transfer
- View Synchrony = "Virtual Synchrony"
 - Provides an abstraction of a synchronous network that hides the asynchrony of the underlying network from distributed applications
 - But does not violate FLP impossibility (since can partition)
- Used in ISIS toolkit (NY Stock Exchange)

Back to Replication



- Need consistent updates to all copies of object
 - Linearizability
 - Sequential consistency

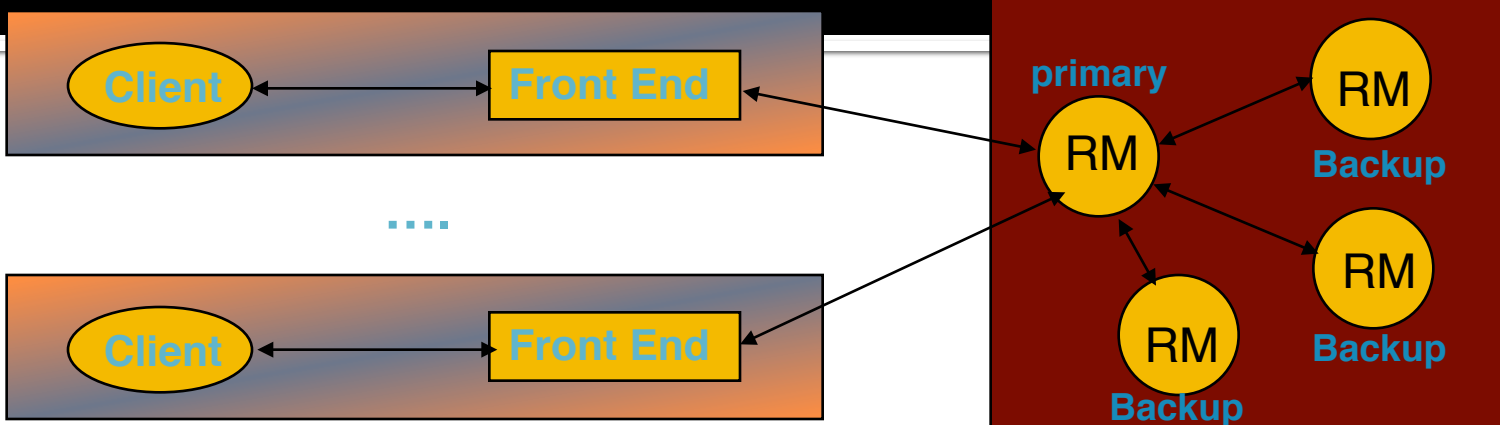
Linearizability

- Let the sequence of read and update operations that client i performs in some execution be oi_1, oi_2, \dots
 - "Program order" for the client
- A replicated shared object service is **linearizable** if for any execution (real), there is some interleaving of operations (virtual) issued by all clients that:
 - meets the specification of a single correct copy of objects
 - is consistent with the real times at which each operation occurred during the execution
- Main goal: any client will see (at any point of time) a copy of the object that is correct and consistent

Sequential Consistency

- The real-time requirement of linearizability is hard, if not impossible, to achieve in real systems
- A less strict criterion is **sequential consistency**: A replicated shared object service is sequentially consistent if for any execution (real), there is some interleaving of clients' operations (virtual) that:
 - meets the specification of a single correct copy of objects
 - is consistent with the program order in which each individual client executes those operations.
- This approach does not require absolute time or total order. Only that for each client the order in the sequence be consistent with that client's program order (~ FIFO).
- Linearizability implies sequential consistency. Not vice-versa!
- Challenge with guaranteeing seq. cons.?
 - Ensuring that all replicas of an object are consistent.

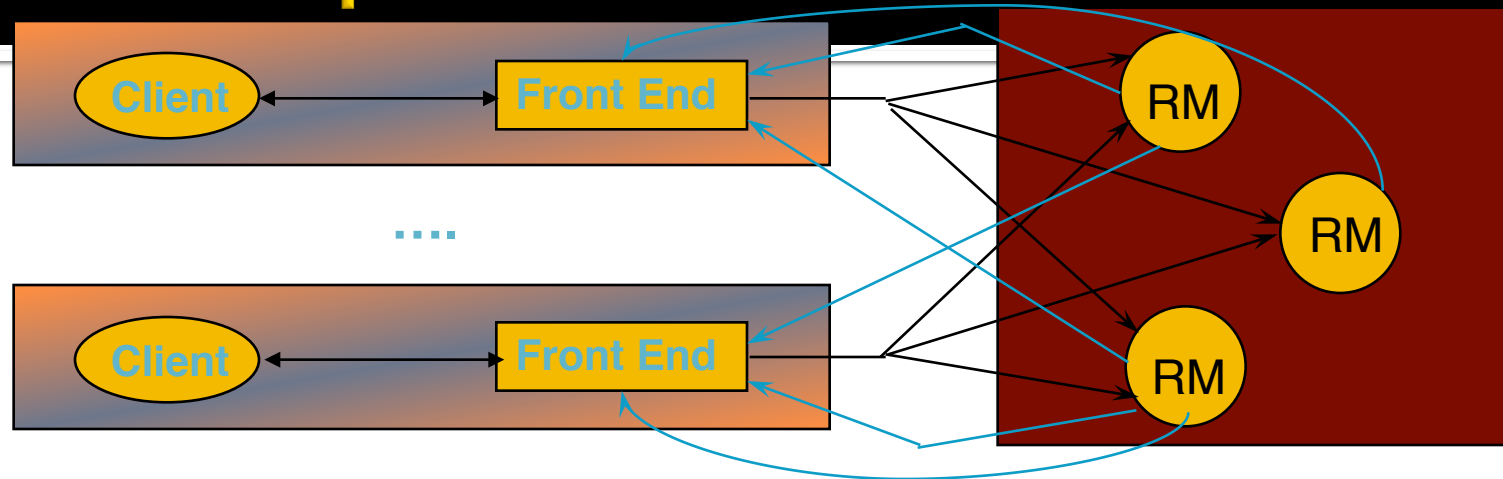
Passive (Primary-Backup) Replication

- 
- The diagram illustrates the architecture of passive (primary-backup) replication. On the left, two client nodes are shown, each consisting of a 'Client' (yellow oval) and a 'Front End' (yellow rectangle) connected by a bidirectional arrow. These client nodes are connected to a central 'primary' node (yellow circle) within a red-shaded replication group. This primary node is also connected to three 'Backup' nodes (yellow circles), each labeled 'RM Backup'. The primary node is labeled 'primary' and 'RM'. The backup nodes are labeled 'RM Backup'. Arrows indicate the flow of data and control between the primary and backup nodes, and between the client front ends and the primary node.
- **Request Communication:** the request is issued to the primary RM and carries a unique request id.
 - **Coordination:** Primary takes requests atomically, in order, checks id (resends response if not new id.)
 - **Execution:** Primary executes & stores the response
 - **Agreement:** If update, primary sends updated state/ result, req-id and response to all backup RMs (1-phase commit enough).
 - **Response:** primary sends result to the front end

Fault Tolerance in Passive Replication

- If the primary fails, a backup becomes primary by leader election, and the replica managers that survive agree on which operations had been performed at the point when the new primary takes over.
 - The above requirement can be met if the replica managers (primary and backups) are organized as a group and if the primary uses view-synchronous group communication to send updates to backups.
- Thus the system remains sequentially consistent in spite of crashes

Active Replication

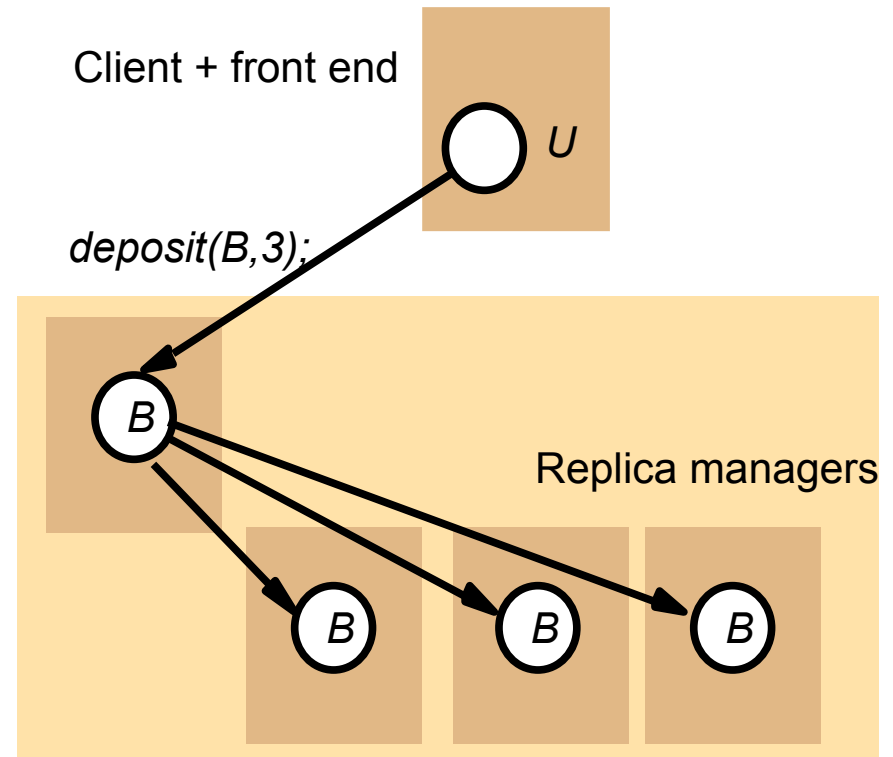
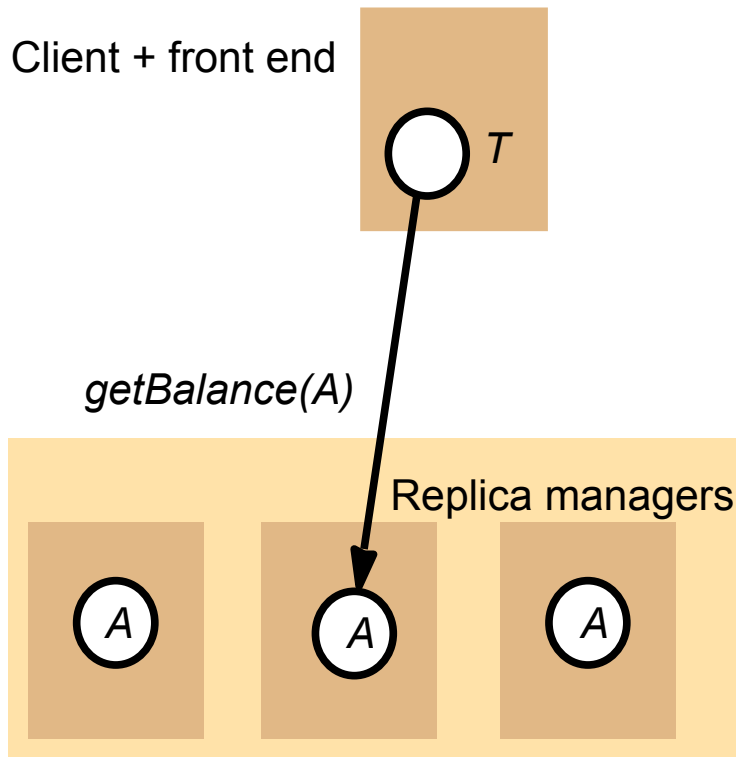


- **Request Communication:** The request contains a unique identifier and is multicast to all by a reliable totally-ordered multicast.
- **Coordination:** Group communication ensures that requests are delivered to each RM in the same order (but may be at different physical times!).
- **Execution:** Each replica executes the request. (Correct replicas return same result since they are running the same program, i.e., they are replicated protocols or replicated state machines)
- **Agreement:** No agreement phase is needed, because of multicast delivery semantics of requests
- **Response:** Each replica sends response directly to FE

Fault Tolerance in Active Replication

- RMs work as replicated state machines, playing equivalent roles. That is, each responds to a given series of requests in the same way. One way of achieving this is by running the same program code at all RMs (but only one way – why?).
- If any RM crashes, state is maintained by other correct RMs.
- This system implements sequential consistency
 - The total order ensures that all correct replica managers process the same set of requests in the same order.
 - Each front end's requests are served in FIFO order (because the front end awaits a response before making the next request).
- So, requests are FIFO-total ordered.
- Caveat (Out of band): If clients are multi-threaded and communicate with one another while waiting for responses from the service, we may need to incorporate causal-total ordering.

Transactions on Replicated Data



One Copy Serialization

- In a non-replicated system, transactions appear to be performed one at a time in some order. This is achieved by ensuring a serially equivalent interleaving of transaction operations.
- **One-copy serializability**: The effect of transactions performed by clients on replicated objects should be the same as if they had been performed one at a time on a single set of objects (i.e., 1 replica per object).
 - Equivalent to combining serial equivalence + replication transparency/consistency

Two Phase Commit Protocol For Transactions on Replicated Objects

- Two level nested 2PC
- In the first phase, the coordinator sends the `canCommit?` command to the participants, each of which then passes it onto the other RMs involved (e.g., by using view synchronous communication) and collects their replies before replying to the coordinator.
- In the second phase, the coordinator sends the `doCommit` or `doAbort` request, which is passed onto the members of the groups of RMs.

Primary Copy Replication

- For now, assume no crashes/failures
- All the client requests are directed to a single primary RM.
- Concurrency control is applied at the primary.
- To commit a transaction, the primary communicates with the backup RMs and replies to the client.
- View synchronous comm. gives → one-copy serializability
- Disadvantage? Performance is low since primary RM is bottleneck.

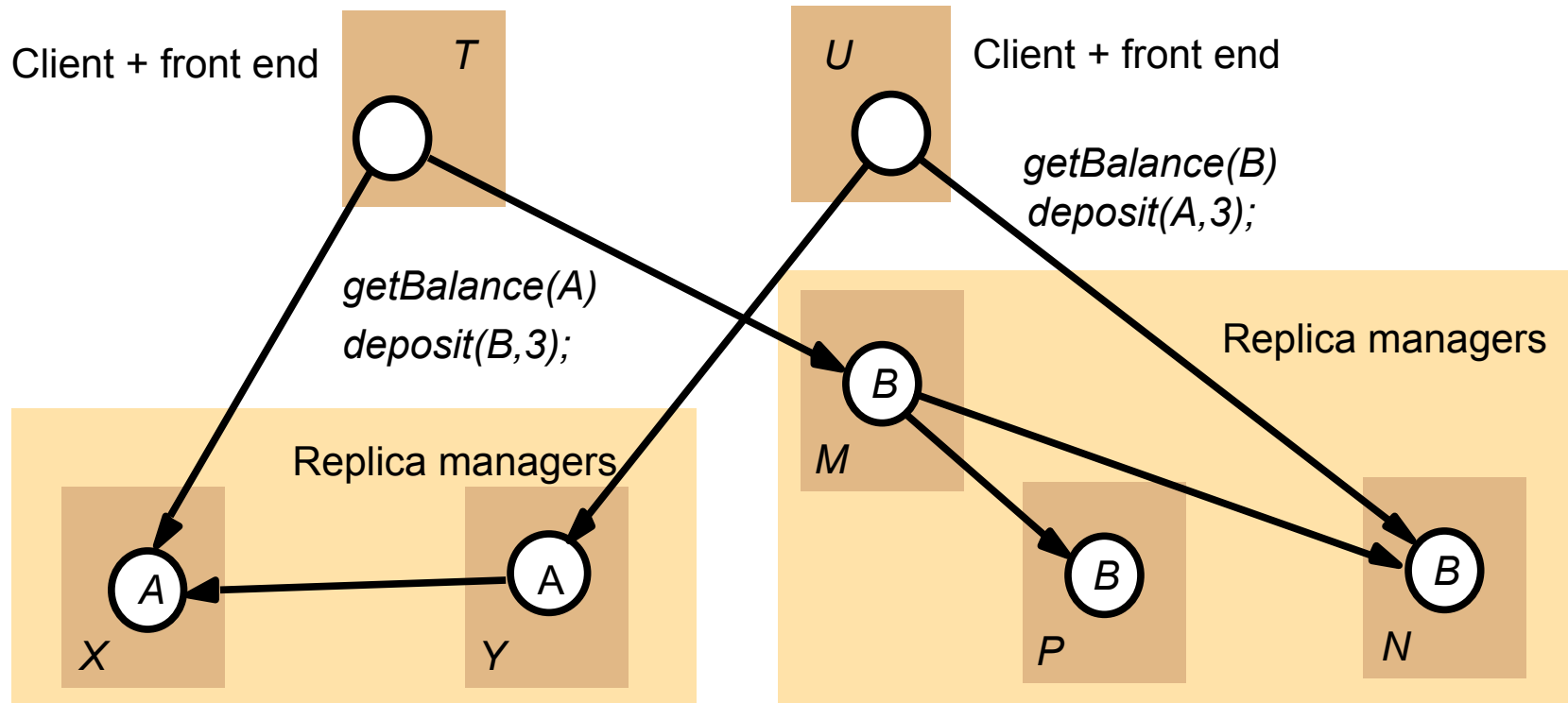
Read One/Write All Replication

- An FE (client front end) may communicate with any RM.
- Every write operation must be performed at all of the RMs
 - Each contacted RM sets a write lock on the object.
- A read operation can be performed at any single RM
 - A contacted RM sets a read lock on the object.
- Consider pairs of conflicting operations of different transactions on the same object.
 - Any pair of write operations will require locks at all of the RMs
→ not allowed
 - A read operation and a write operation will require conflicting locks at some RM → not allowed
 - One-copy serializability is achieved.
- Disadvantage? Failures block the system (esp. writes).

Available Copies Replication

- A client's read request on an object can be performed by any RM, but a client's update request must be performed across all available (i.e., non-faulty) RMs in the group.
- As long as the set of available RMs does not change, local concurrency control achieves one-copy serializability in the same way as in read-one/write-all replication.
- May not be true if RMs fail and recover during conflicting transactions.

Available Copies Approach



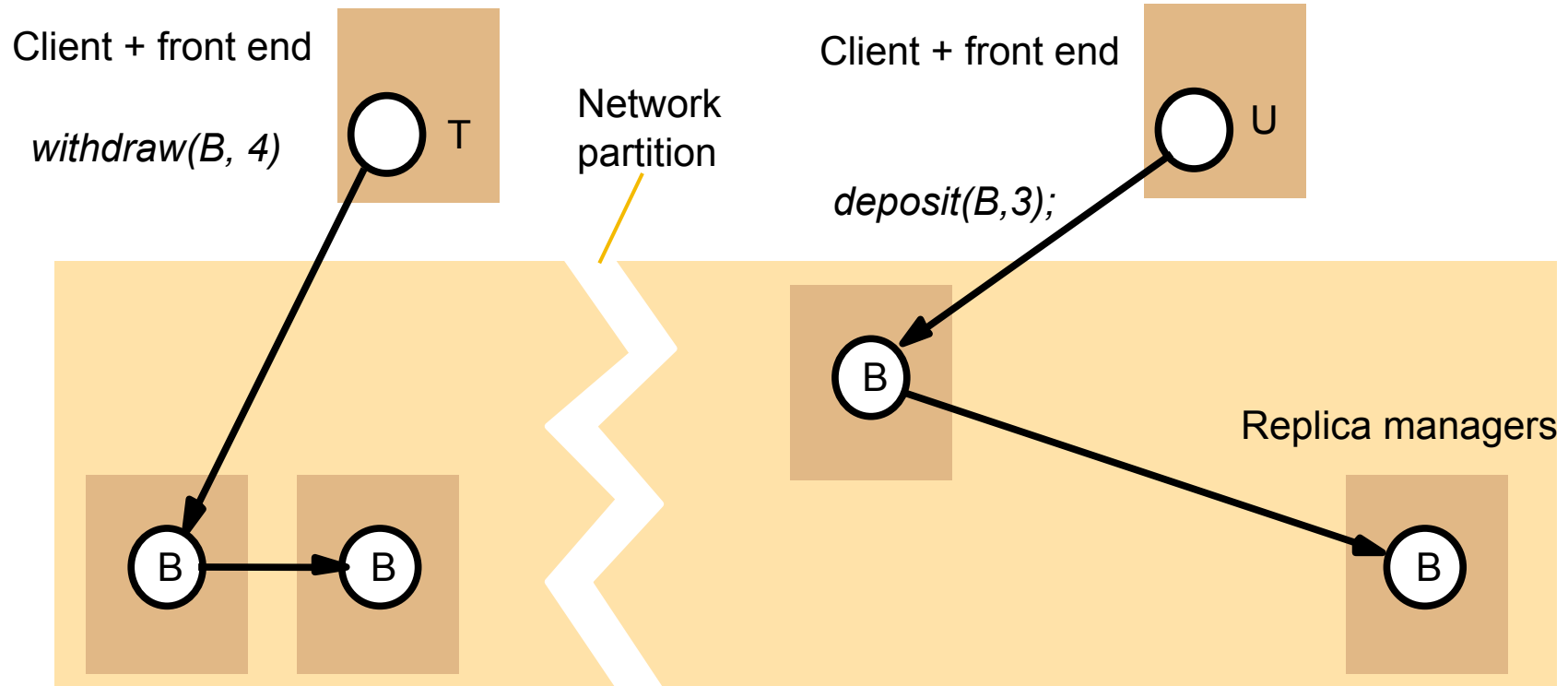
The Impact of RM Failure

- Assume that (i) RM X fails just after T has performed getBalance; and (ii) RM N fails just after U has performed getBalance. Both failures occur before any of the deposit()'s.
- Subsequently, T's deposit will be performed at RMs M and P, and U's deposit will be performed at RMY.
- The concurrency control on A at RM X does not prevent transaction U from updating A at RMY.
- Solution: Must also serialize RM crashes and recoveries with respect to entire transactions.

Local Validation (using Our Example)

- From T' s perspective,
 - T has read from an object at X → X must have failed after T' s operation.
 - T observes the failure of N when it attempts to update the object B → N' s failure must be before T.
 - Thus: N fails → T reads object A at X; T writes objects B at M and P → T commits → X fails.
- From U' s perspective,
 - Thus: X fails → U reads object B at N; U writes object A at Y → U commits → N fails.
- At the time T tries to commit,
 - it first checks if N is still not available and if X, M and P are still available. Only then can T commit.
 - It then checks if the failure order is consistent with that of other transactions (T cannot commit if U has committed)
 - If T commits, U' s validation will fail because N has already failed.
- Can be combined with 2PC.
- Caveat: Local validation may not work if partitions occur in the network

Network Partition



Dealing with Network Partitions

- During a partition, pairs of conflicting transactions may have been allowed to execute in different partitions. The only choice is to take corrective action after the network has recovered
 - Assumption: Partitions heal eventually
- Abort one of the transactions after the partition has healed
- Basic idea: allow operations to continue in partitions, but finalize and commit trans. only after partitions have healed
- But to optimize performance, better to avoid executing operations that will eventually lead to aborts...how?

Quorum Approaches

- **Quorum** approaches used to decide whether reads and writes are allowed
- There are two types: **pessimistic quorums** and **optimistic quorums**
- In the pessimistic quorum philosophy, updates are allowed only in a partition that has the majority of RMs
 - Updates are then propagated to the other RMs when the partition is repaired.

Static Quorums

- The decision about how many RMs should be involved in an operation on replicated data is called Quorum selection
- Quorum rules state that:
 - At least r replicas must be accessed for read
 - At least w replicas must be accessed for write
 - $r + w > N$, where N is the number of replicas
 - $w > N/2$
 - Each object has a version number or a consistent timestamp
- Static Quorum predefines r and w , & is a pessimistic approach: if partition occurs, update will be possible in at most one partition

Voting with Static Quorums

- A version of quorum selection where each replica has a number of votes. Quorum is reached by majority of votes (N is the total number of votes)
- e.g., a cache replica may be given a 0 vote
- with $r = w = 2$, Access time for write is 750 ms (parallel writes). Access time for read without cache is 750 ms. Access time for read with cache can be in the range 175ms to 825ms – why?.

Replica	votes	access time	version chk	P(failure)
Cache	0	100ms	0ms	0%
Rep1	1	750ms	75ms	1%
Rep2	1	750ms	75ms	1%
Rep3	1	750ms	75ms	1%

Optimistic Quorum Approaches

- An Optimistic Quorum selection allows writes to proceed in any partition.
- This might lead to write-write conflicts. Such conflicts will be detected when the partition heals
 - Any writes that violate one-copy serializability will then result in the transaction (that contained the write) to abort
 - Still improves performance because partition repair not needed until commit time (and it's likely the partition may have healed by then)
- Optimistic Quorum is practical when:
 - Conflicting updates are rare
 - Conflicts are always detectable
 - Damage from conflicts can be easily confined
 - Repair of damaged data is possible or an update can be discarded without consequences
 - Partitions are relatively short-lived

View-based Quorum

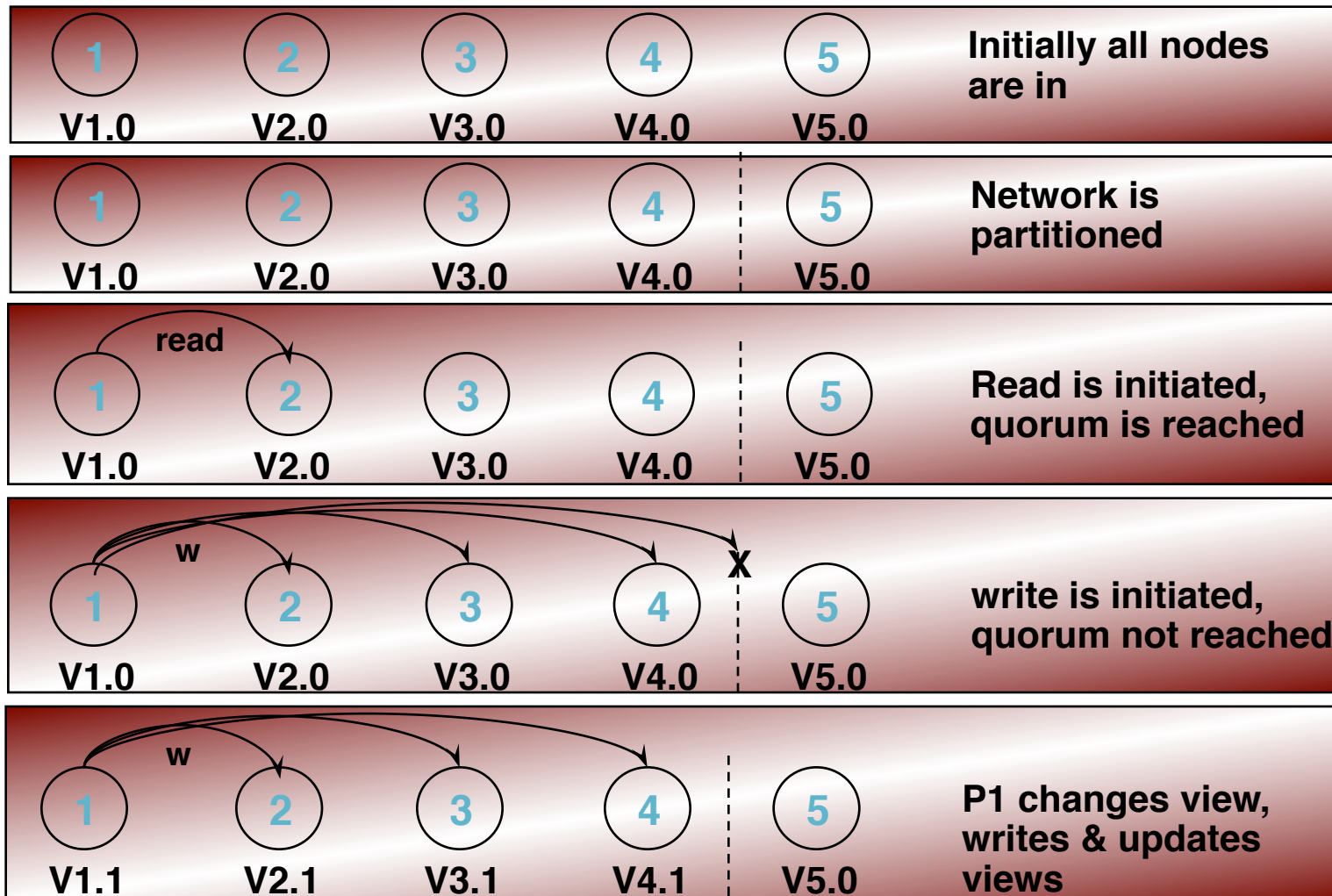
- An optimistic approach
- Quorum is based on views at any time
 - Uses group communication as a building block (see previous lecture)
- In a partition, inaccessible nodes are considered in the quorum as ghost participants that reply “Yes” to all requests.
 - Allows operations to proceed if the partition is large enough (need not be majority)
- Once the partition is repaired, participants in the smaller partition know whom to contact for updates.

View-based Quorum - details

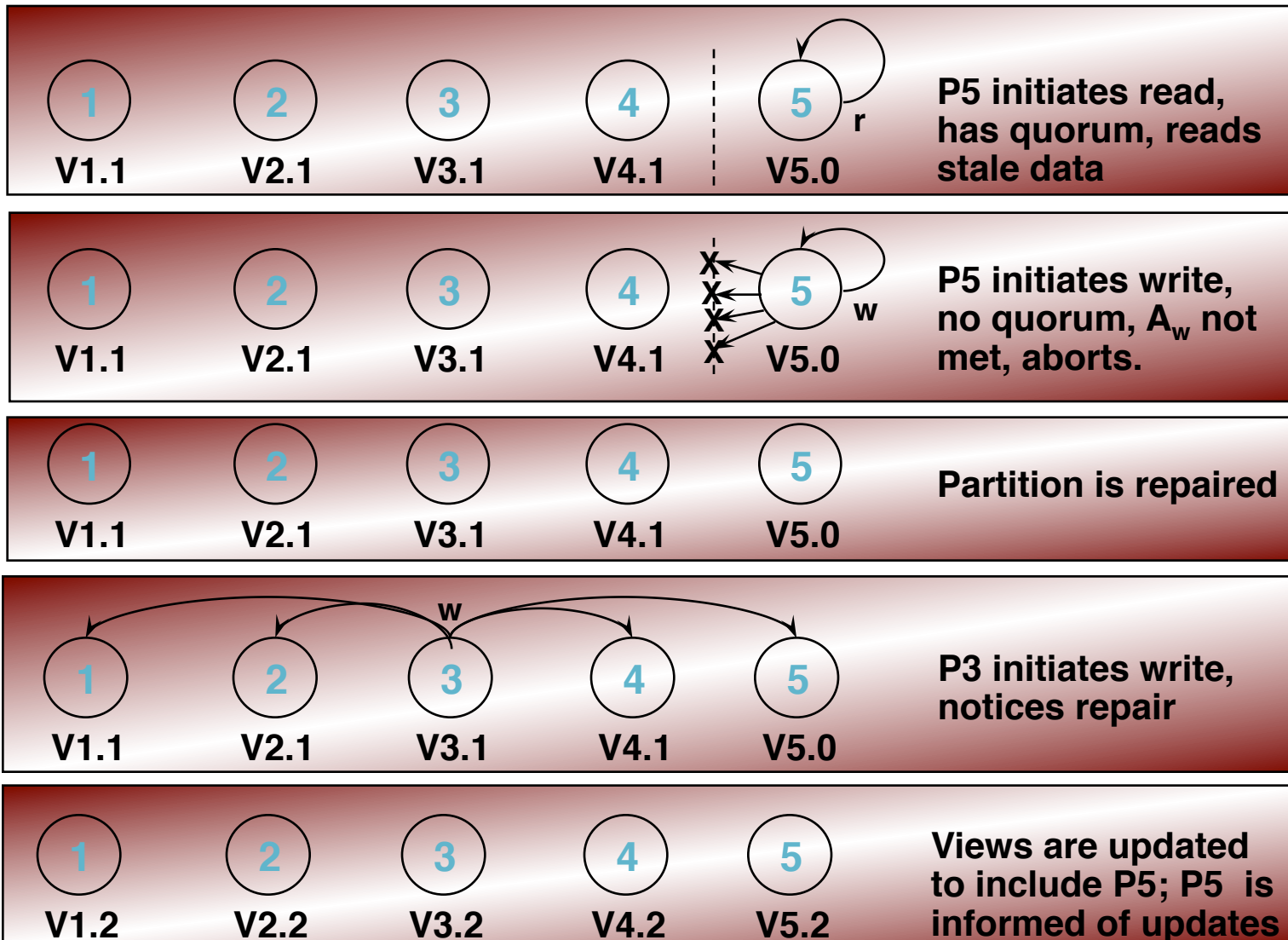
- Uses view-synchronous communication as a building block (see previous lecture)
- Views are per object, numbered sequentially and only updated if necessary
- We define thresholds for each of read and write :
 - A_w : minimum nodes in a view for write, e.g., $A_w > N/2$
 - A_r : minimum nodes in a view for read
 - E.g., $A_w + A_r > N$
- If ordinary quorum cannot be reached for an operation, then we take a straw poll, i.e., we update views
- In a large enough partition for read, $\text{Viewsize} \geq A_r$
- In a large enough partition for write, $\text{Viewsize} \geq A_w$ (inaccessible nodes are considered as ghosts that reply Yes to all requests.)
- The first update after partition repair forces restoration for nodes in the smaller partition

Example: View-based Quorum

- Consider: $N = 5$, $w = 5$, $r = 1$, $A_w = 3$, $A_r = 1$



Example: View-based Quorum (cont'd)



Summary

- Transactions
- Concurrency Control
- Replicated Data
- Replication with Transactions
- Different types of replication
- Quorums
- Reading for this lecture was: Section 15.5
- Reading for next week: Section 15.4, Handout
- Topics: Gossiping, Self-stabilization
- MP2 due Oct 31. Demo Signup Sheet on class newsgroup – please sign up! (Demos on Wednesday Nov 3).

Optional Slides

Quorum Consensus Examples

[Gifford]'s examples
for a replicated file system

Ex1:
High R to W ratio
Single RM on Replica 1

Ex2:
Moderate R to W ratio
Accessed from local
LAN of RM 1

Ex3:
V. High R to W ratio
All RM's equidistant

		Example 1	Example 2	Example 3
<i>Latency</i> (milliseconds)	Replica 1	75	75	75
	Replica 2	65	100	750
	Replica 3	65	750	750
<i>Voting configuration</i>	Replica 1	1	2	1
	Replica 2	0	1	1
	Replica 3	0	1	1
<i>Quorum sizes</i>	R	1	2	1
	W	1	3	3
Derived performance of file suite:				
<i>Read</i>	Latency	75	75	75
	Blocking probability	0.01	0.0002	0.000001
<i>Write</i>	Latency	75	100	750
	Blocking probability	0.01	0.0101	0.03

0.01 failure prob.
per RM

Summary

- Replicating objects across servers improves performance, fault-tolerance, availability
- Raises problem of Replica Management
- Group communication an important building block
- View Synchronous communication service provides totally ordered delivery of views +multicasts
- RMs can be built over this service
- Passive and Active Replication