## DISTRIBUTED SYSTEMS (COMP9243)

### Lecture 3a: Replication & Consistency

**Slide 1**

① Replication
② Consistency
- Models
- Protocols
③ Update propagation
④ Replica placement

---

### REPLICATION

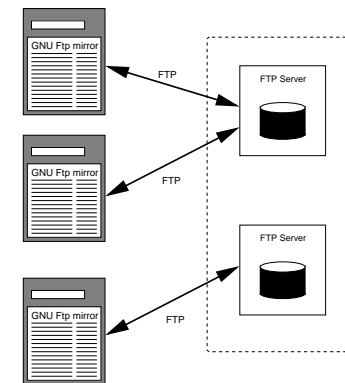Make copies of services on multiple machines.

**Why?:**

➔ Reliability
- Redundancy

➔ Performance
- Increase processing capacity
- Reduce communication

➔ Scalability (prevent centralisation)
- Prevent overloading of single server (*size* scalability)
- Avoid communication latencies (*geographic* scalability)

**Slide 2**

---

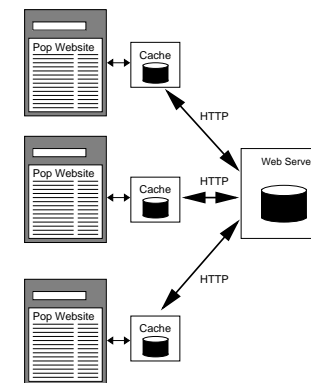### DATA VS CONTROL REPLICATION

Data Replication (Server Replication/Mirroring):

**Slide 3**



---

Data Replication (Caching):

**Slide 4**



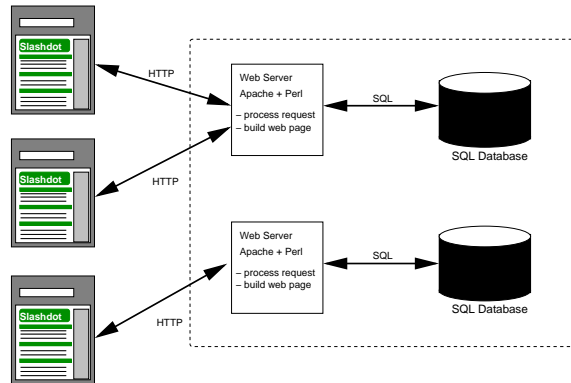What's the difference between mirroring and caching?

## Slide 5

**Slide 5**

Control Replication:



## Slide 6

**Slide 6**

Data and Control Replication:



Will be looking primarily at data replication (including combined data and control replication).

## Slide 7

**Slide 7**

Updates
➜ Consistency (how to deal with updated data)
➜ Update propagation

Replica placement
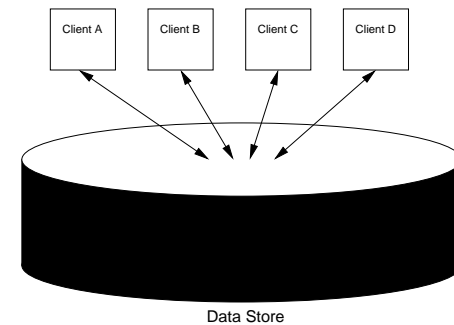➜ How many replicas?
➜ Where to put them?

Redirection/Routing
➜ Which replica should clients use?
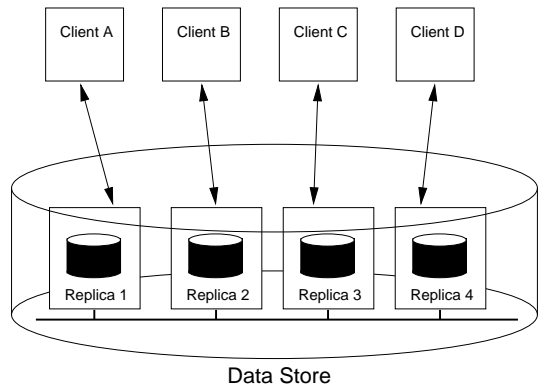
## Slide 8

**Slide 8**

DISTRIBUTED DATA STORE

➜ data-store stores data items

Client's Point of View:



Data Store

**Slide 9**

### Distributed Data-Store's Point of View:



**Slide 10**

### Data Model:

➔ data store: collection of data items
➔ data item: simple variable
➔ data item values: explicit (`0`, `1`), abstract (`a`,`b`)

### Operations on a Data Store:

➔ Read. `Ri(x)b` Client i performs a read for data item `x` and it returns `b`
➔ Write. `Wi(x)a` Client i performs write on data item `x` setting it to `a`
➔ Operations not instantaneous
  • Time of issue (when request is sent by client)
  • Time of execution (when request is executed at a replica)
  • Time of completion (when reply is received by client)
➔ Coordination among replicas

**Slide 11**

### Replica Managers:



**Slide 12**

### Timeline:

## INCONSISTENCY

**Slide 13**

Staleness:
➜ How old is the data?
➜ How old is the data allowed to be?
  • Time
  • Versions

Operation order:
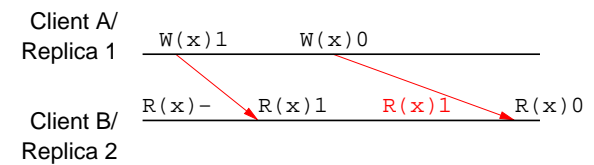➜ Were operations performed in the right order?
➜ What orderings are allowed?

Conflicting Data:
➜ Do replicas have exactly the same data?
➜ What differences are permitted?

## CONSISTENCY

**Slide 14**

Non-distributed data store:
➜ Program order is maintained
➜ Data coherence is respected

Updates and concurrency result in conflicting operations

Conflicting Operations:
➜ Read-write conflict (only 1 write)
➜ Write-write conflict (multiple concurrent writes)

Consistency:
➜ The order in which conflicting operations are performed affects consistency
➜ partial order: order of a single client's operations
➜ total order: interleaving of all conflicting operations

Example:

**Slide 15**

**Client A:** `x = 1; x = 0;`

**Client B:** `print(x);`
`print(x);`

Possible results:
- -, 11, 10, 00
How about 01?

What are the conflicting ops? What are the partial orders?
What are the total orders?

Client A    `W(x)1    W(x)0`

Client B                    `R(x)0   R(x)1`

## CONSISTENCY MODEL

**Slide 16**

*Defines which interleavings of operations are valid (admissible)*

Consistency Model:
➜ Concerned with consistency of a data store.
➜ Specifies characteristics of valid total orderings

A data store that implements a particular model of consistency will provide a total ordering of operations that is valid according to the model.

**Slide 17**

Data Coherence vs Data Consistency:

**Data Coherence** ordering of operations for single data item
➜ e.g. a read of x will return the most recently written value of x

**Data Consistency** ordering of operations for whole data store
➜ implies data coherence
➜ includes ordering of operations on other data items too

**Slide 18**

### DATA-CENTRIC CONSISTENCY MODEL

A contract, between a distributed data store and clients, in which the data store specifies precisely what the results of read and write operations are in the presence of concurrency.

➜ Described consistency is experienced by all clients
➜ Multiple clients accessing the same data store
➜ Client A, Client B, Client C see same kinds of orderings
➜ Non-mobile clients (replica used doesn't change)

**Slide 19**

### STRONG ORDERING VS WEAK ORDERING

Strong Ordering (tight):
➜ All writes must be performed in the order that they are invoked
➜ Example: all clients must see: `W(x)a W(x)b W(x)c`
➜ Strict (Linearisable) Sequential, Causal, FIFO (PRAM)

Weak Ordering (loose):
➜ Ordering of *groups* of writes, rather than individual writes
➜ Series of writes are grouped on a single replica
➜ Only results of grouped writes propagated.
➜ Example: `{W(x)a W(x)b W(x)c} == {W(x)b W(x)a W(x)c}`
➜ Weak, Release, Entry

**Slide 20**

### STRICT CONSISTENCY

*Any read on a data item x returns a value corresponding to the result of the most recent write on x*

Absolute time ordering of all shared accesses



strictly consistent          not strictly consistent

What is *most recent* in a distributed system?
➜ Assumes an absolute global time
➜ Assumes instant communication (atomic operation)
➜ Normal on a uniprocessor
✗ Impossible in a distributed system

## Sequential Consistency

*All operations are performed in some sequential order*

➔ More than one correct sequential order
➔ All clients see the *same* order
➔ Program order of each client maintained
➔ Not ordered according to time

```
Client A  ─────────W(x)a──────────      Client A  ─────────W(x)a──────────
Client B  ──W(x)b────────────────      Client B  ──W(x)b────────────────
Client C  ──────R(x)b──────R(x)a──      Client C  ──────R(x)a──────R(x)b──
Client D  ────────────R(x)b──R(x)a─      Client D  ────────────R(x)b──R(x)a─
              sequential                        not sequential
```
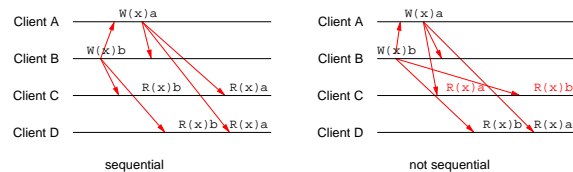
### Performance:

read time + write time >= minimal packet transfer time

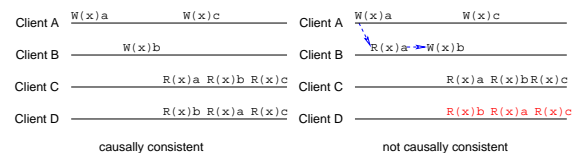## Causal Consistency

*Potentially causally related writes are executed in the same order everywhere*
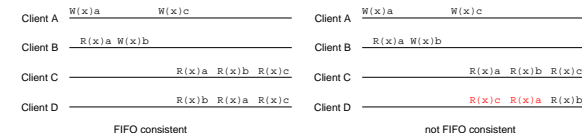
### Causally Related Operations:

➔ Read followed by a write (in same client)
➔ `W(x)` followed by `R(x)` (in same or different clients)

```
Client A  W(x)a          W(x)c       Client A  W(x)a          W(x)c
Client B  ──────W(x)b───────────     Client B  ───R(x)a──W(x)b──────
Client C  ────R(x)a R(x)b R(x)c──    Client C  ────R(x)a R(x)b R(x)c──
Client D  ────R(x)b R(x)a R(x)c──    Client D  ────R(x)b R(x)a R(x)c──
          causally consistent              not causally consistent
```

## FIFO (PRAM) Consistency

*Only partial orderings of writes maintained*

```
Client A  W(x)a        W(x)c       Client A  W(x)a        W(x)c
Client B  ──R(x)a W(x)b──────      Client B  ──R(x)a W(x)b──────
Client C  ────R(x)a R(x)b R(x)c─   Client C  ────R(x)a R(x)b R(x)c─
Client D  ────R(x)b R(x)a R(x)c─   Client D  ────R(x)c R(x)a R(x)b─
          FIFO consistent                   not FIFO consistent
```

## Weak Consistency

*Shared data can be counted on to be consistent only after a synchronisation is done*

Enforces consistency on a *group of operations*, rather than single operations

➔ Synchronisation variable (`S`)
➔ Synchronise operation (`synchronise(S)`)
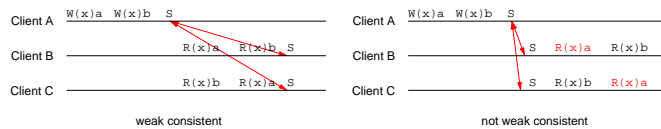➔ Define `critical section' with synchronise operations

### Properties:

➔ Order of synchronise operations sequentially consistent
➔ Synchronise operation cannot be performed until all previous writes have completed everywhere
➔ Read or Write operations cannot be performed until all previous synchronise operations have completed

## Example:

**Slide 25**

➜ `synchronise(S) W(x)a W(y)b W(x)c synchronise(S)`
➜ Writes performed locally
➜ Updates propagated only upon synchronisation
➜ Only `W(y)b` and `W(x)c` have to be propagated



weak consistent                    not weak consistent

---

**Slide 26**

### RELEASE CONSISTENCY

Explicit separation of synchronisation tasks

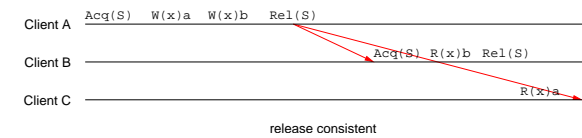➜ `acquire(S)` - bring local state up to date
➜ `release(S)` - propagate all local updates
➜ acquire-release pair defines 'critical region'

### Properties:

➜ Order of synchronisation operations are FIFO consistent
➜ Release cannot be performed until all previous reads and writes done by the client have completed
➜ Read or Write operations cannot be performed until all previous acquires done by the client have completed

---

**Slide 27**



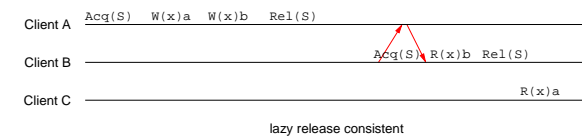release consistent

---

## Lazy Release Consistency:

**Slide 28**

➜ Don't send updates on release
➜ Acquire causes client to get newest state
➜ Added efficiency if acquire-release performed by same client (e.g., in a loop)



lazy release consistent

## ENTRY CONSISTENCY

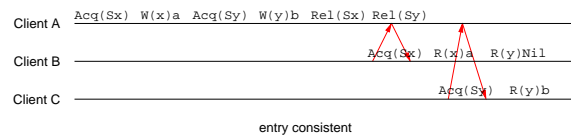*Synchronisation variable associated with specific shared data item (guarded data item)*

➜ Each shared data item has own synchronisation variable
➜ `acquire()`
  - Provides ownership of synchronisation variable
  - Exclusive and nonexclusive access modes
  - Synchronises data
  - Requires communication with current owner
➜ `release()`
  - Relinquishes exclusive access (but not ownership)
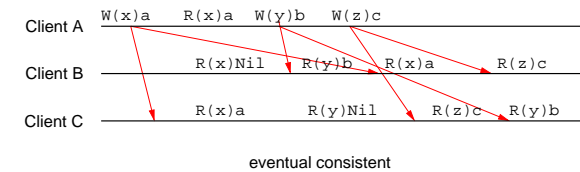
**Slide 29**

---

Properties:

➜ Acquire does not complete until all guarded data is brought up to date locally
➜ If a client has exclusive access to a synchronisation variable, no other client can have any kind of access to it
➜ When acquiring nonexclusive access, a client must first get the updated values from the synchronisation variable's current owner

**Slide 30**



entry consistent

---

## EVENTUAL CONSISTENCY

*If no updates take place for a long time, all replicas will gradually become consistent*



eventual consistent

**Slide 31**

Requirements:

➜ Few read-write conflicts (R » W)
➜ Few write-write conflicts
➜ Clients accept inconsistency (i.e., old data)

---

Examples:

➜ DNS:
  - no write-write conflicts
  - updates slowly (1-2 days) propagate to all caches
➜ WWW:
  - few write-write conflicts
  - mirrors eventually updated
  - cached copies (browser or proxy) eventually replaced

**Slide 32**

## CAP Theory

C: Consistency: Linearisability
A: Availability: Timely response
P: Partition-Tolerance: Functions in the face of a partition
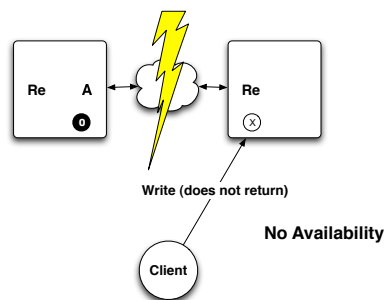
**Slide 33**



Consistency

Partition Tolerance    Availability

---

CAP Impossibility Proof:

**Slide 34**



Re    A
    0

Re
⊗

Write (does not return)

No Availability

Client

---

## CAP Consequences

For wide-area systems:

➜ must choose: Consistency or Availability
➜ choosing Availability
   • Eventual consistency
➜ choosing Consistency
   • delayed (and potentially failing) operations

**Slide 35**

---

## Client-Centric Consistency Models

*Provides guarantees about ordering of operations for a single client*
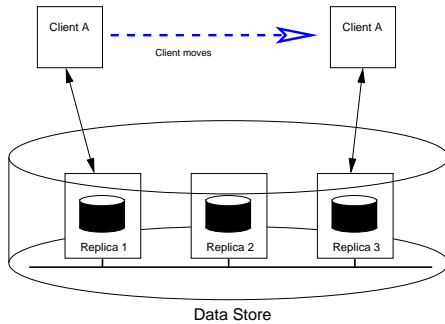
➜ Single client accessing data store
➜ Client accesses different replicas (modified data store model)
➜ Data isn't shared by clients
➜ Client A, Client B, Client C may see different kinds of orderings

**Slide 36**

In other words:

➜ The effect of an operation depends on the client performing it
➜ Effect also depends on the history of operations that client has performed.

---

## Slide 37

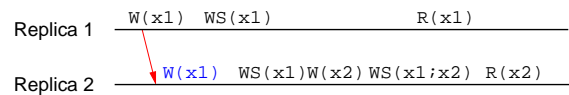**Data-Store Model for Client-Centric Consistency:**



Data Store

- Data-items have an owner

- No write-write conflicts

## Slide 38

**Notation and Timeline for Client-Centric Consistency:**

- ➔ `xi[t]`: version of x at replica i at time t
- ➔ Write Set: `WS(xi[t])`: set of writes at replica i that led to xi(t)
- ➔ `WS(xi[t1];xj[t2])`: WS(xj(t2)) contains same operations as WS(xi(t1))
- ➔ `WS(!xi[t1];xj[t2])`: WS(xj(t2)) does not contain the same operations as WS(xi(t1))
- ➔ `R(xi[t])`: a read of x returns xi(t)

## Slide 39

### MONOTONIC READS

*If a client has seen a value of x at a time t, it will never see an older version of x at a later time*



monotonic–read consistent     not monotonic–read consistent

**When is Monotonic Reads sufficient?**

## Slide 40

### MONOTONIC WRITES

*A write operation on data item x is completed before any successive write on x by the same client*

All writes by a single client are sequentially ordered.



monotonic–write consistent     not monotonic–write consistent

**How is this different from FIFO consistency?**

- ➔ Only applies to write operations of single client.
- ➔ Writes from clients not requiring monotonic writes may appear in different orders.

## READ YOUR WRITES

*The effect of a write on x will always be seen by a successive read of x by the same client*

```
Replica 1 ──W(x1)──────────────        Replica 1 ──W(x1)──────────────
              ╲                                    
Replica 2 ────╲─WS(x1;x2)──────R(x2)    Replica 2 ────WS(!x1;x2)──────R(x2)

       read─your─writes consistent            not read─your─writes consistent
```

## WRITE FOLLOWS READS

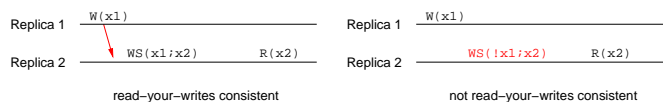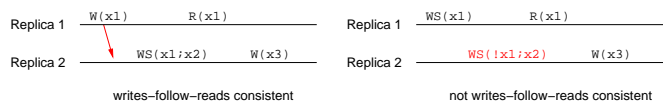*A write operation on x will be performed on a copy of x that is up to date with the value most recently read by the same client*

```
Replica 1 ──W(x1)────────R(x1)──────        Replica 1 ──WS(x1)────────R(x1)──────
              ╲                                    
Replica 2 ────╲─WS(x1;x2)──────W(x3)──      Replica 2 ────WS(!x1;x2)──────W(x3)──

       writes─follow─reads consistent            not writes─follow─reads consistent
```

Trade-offs

### Consistency and Redundancy:
- ➜ All copies must be strongly consistent
- ➜ All copies must contain full state
- ➜ Reduced consistency → reduced reliability

### Consistency and Performance:
- ➜ Consistency requires extra work
- ➜ Consistency requires extra communication
- ✗ Can result in loss of overall performance

### Consistency and Scalability:
- ➜ Implementation of consistency must be scalable
  - • don't take a centralised approach
  - • avoid too much extra communication

## CONSISTENCY PROTOCOLS

Consistency Protocol: implementation of a consistency model

### Primary-Based Protocols:
- ➜ Remote-write protocols
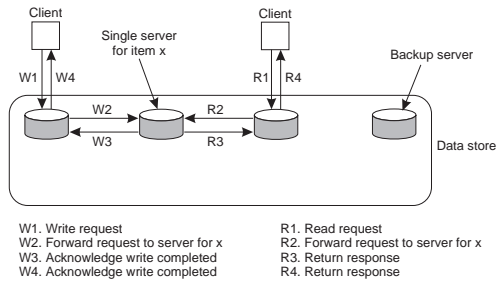- ➜ Local-write protocols

### Replicated-Write Protocols:
- ➜ Active Replication
- ➜ Quorum-Based Protocols

## Single Server:

➜ All writes and reads executed at single server
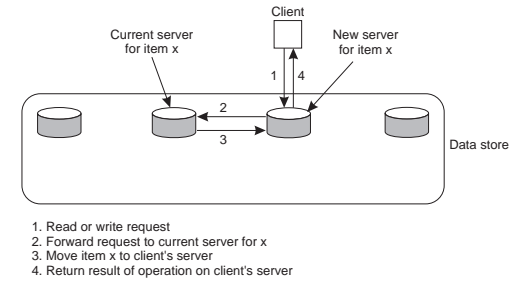➜ No replication of data

**Slide 45**



```
Client              Client
              Single server                   Backup server
              for item x
W1   W4                    R1   R4
        W2         R2
        W3         R3                          Data store
```

W1. Write request          R1. Read request
W2. Forward request to server for x   R2. Forward request to server for x
W3. Acknowledge write completed    R3. Return response
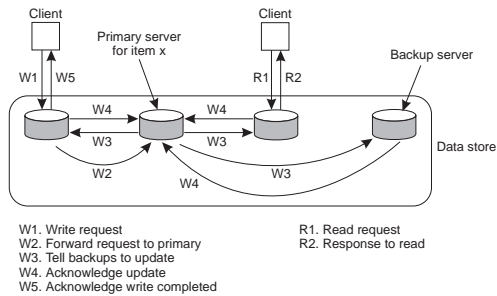W4. Acknowledge write completed    R4. Return response

## Primary-Backup:
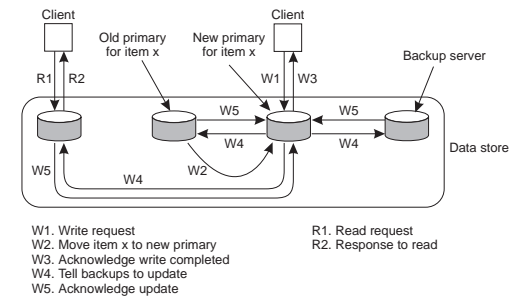
➜ All writes executed at single server, Reads are local
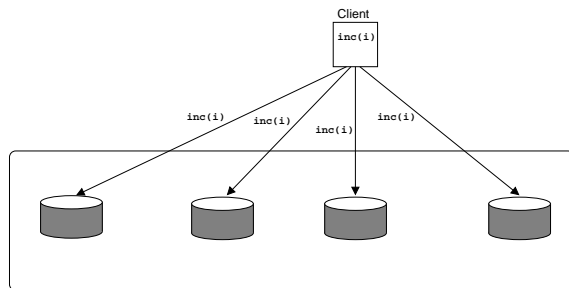➜ Updates block until executed on all backups
✗ Performance

**Slide 46**



```
Client              Client
              Primary server                  Backup server
              for item x
W1   W5                    R1   R2
        W4         W4
        W3         W3                          Data store
        W2         W4       W3
```

W1. Write request          R1. Read request
W2. Forward request to primary  R2. Response to read
W3. Tell backups to update
W4. Acknowledge update
W5. Acknowledge write completed

## Migration:

➜ Data item migrated to local server on access
➜ Distributed, non-replicated, data store

**Slide 47**



```
                         Client
         Current server        New server
         for item x            for item x
                           1   4
                2
                3                            Data store
```

1. Read or write request
2. Forward request to current server for x
3. Move item x to client's server
4. Return result of operation on client's server

## Migrating Primary (multiple reader/single writer):

**Slide 48**



```
      Client              Client
        Old primary    New primary
        for item x     for item x           Backup server
R1   R2                    W1   W3
                    W5              W5
                    W4              W4        Data store
W5          W4         W2
```

W1. Write request          R1. Read request
W2. Move item x to new primary  R2. Response to read
W3. Acknowledge write completed
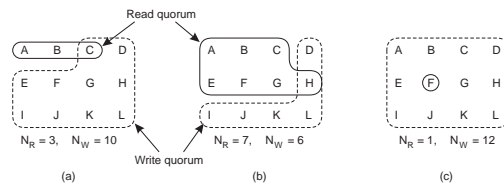W4. Tell backups to update
W5. Acknowledge update

## ACTIVE REPLICATION

➜ Updates (write operation) sent to all replicas
➜ Need totally-ordered multicast
➜ e.g. sequencer/coordinator to add sequence numbers

## QUORUM-BASED PROTOCOLS

➜ Voting
➜ Versioned data
➜ Read Quorum: Nr
➜ Write Quorum: Nw
➜ $N_r + N_w > N$ Why?
➜ $N_w > N/2$ Why?
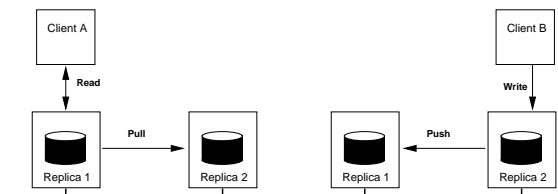
## UPDATE PROPAGATION

What to propagate?

➜ Data
  • R/W high
➜ Update operation
  • low bandwidth costs
➜ Notification/Invalidation
  • R/W low

## PUSH VS PULL



Pull:
➜ Updates propagated only on request
➜ Also called *client-based*
➜ R/W low
➜ Polling delay

Push:
➜ Push updates to replicas
➜ Also called *server-based*
➜ When low staleness required
➜ R » W
✗ Have to keep track of all replicas

**Slide 53**

Compromise: Leases:

Server promises to push updates until lease expires

Lease length depends on:

**age:** Last time item was modified

**renewal-frequency:** How often replica needs to be updated

**state-space overhead:** lower expiration time to reduce bookkeeping when many clients

---

REPLICA PLACEMENT

**Slide 54**



---

Permanent Replicas:
➜ Initial set of replicas
➜ Created and maintained by data-store owner(s)
➜ Allow writes

Server-Initiated Replicas:
➜ Enhance performance
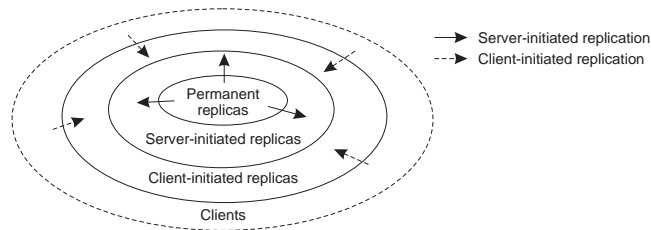➜ Not maintained by owner
➜ Placed close to groups of clients
  ● Manually
  ● Dynamically

**Slide 55**

Client-Initiated Replicas:
➜ Client caches
➜ Temporary
➜ Owner not aware of replica
➜ Placed close to client
➜ Maintained by host (often client)

---

DYNAMIC REPLICATION

Situation changes over time
➜ Number of users, Amount of data
➜ Flash crowds
➜ R/W ratio

**Slide 56** Dynamic Replica Placement:
➜ Network of replica servers
➜ Keep track of data item requests at each replica
➜ Deletion threshold
➜ Replication threshold
➜ Migration threshold
➜ Clients always send requests to nearest server

## MISCELLANEOUS IMPLEMENTATION AND DESIGN ISSUES

End-to-End argument:

➔ Where to implement replication mechanisms?

➔ Application? Middleware? OS?

Policy vs Mechanism:

➔ Consistency models built into middleware?

➔ One-size-fits-all?

Determining Policy:

➔ Who determines the consistency model used?

- Application
- Middleware
- Client
- Server

## READING LIST

**Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services** An overview of the CAP theorem and its proof.

**Eventual Consistency** An overview of eventual consistency and client-centric consistency models.