# Consistency and Replication

- Consistency models
  - Data-centric consistency models
  - Client-centric consistency models

# Why replicate?

- Data replication versus compute replication

- Data replication: common technique in distributed systems
- Reliability
  - If one replica is unavailable or crashes, use another
  - Protect against corrupted data
- Performance
  - Scale with size of the distributed system (replicated web servers)
  - Scale in geographically distributed systems (web proxies)

# Replication Issues

- When to replicate?
- How many replicas to create?
- Where should the replicas located?

- Will return to these issues later (WWW discussion)
- Today: how to maintain *consistency*?
- Key issue: need to maintain *consistency* of replicated data
  - If one copy is modified, others become inconsistent

# CAP Theorem

- Conjecture by Eric Brewer at PODC 2000 conference
  - It is impossible for a web service to provide all three guarantees:
    - **Consistency** (nodes see the same data at the same time)
    - **Availability** (node failures do not the rest of the system)
    - **Partition-tolerance** (system can tolerate message loss)
  - A distributed system can satisfy any two, but not all three, at the same time
- Conjecture was established as a theorem in 2002 (by Lynch and Gilbert)

# CAP Theorem Examples

- Consistency+Availability
  - Single database, cluster database, LDAP, xFS
    - 2 phase commit
- Consistency + partition tolerance
  - distributed database, distributed locking
    - pessimistic locking
- Availability + Partition tolerance
  - Coda, Web caching, DNS
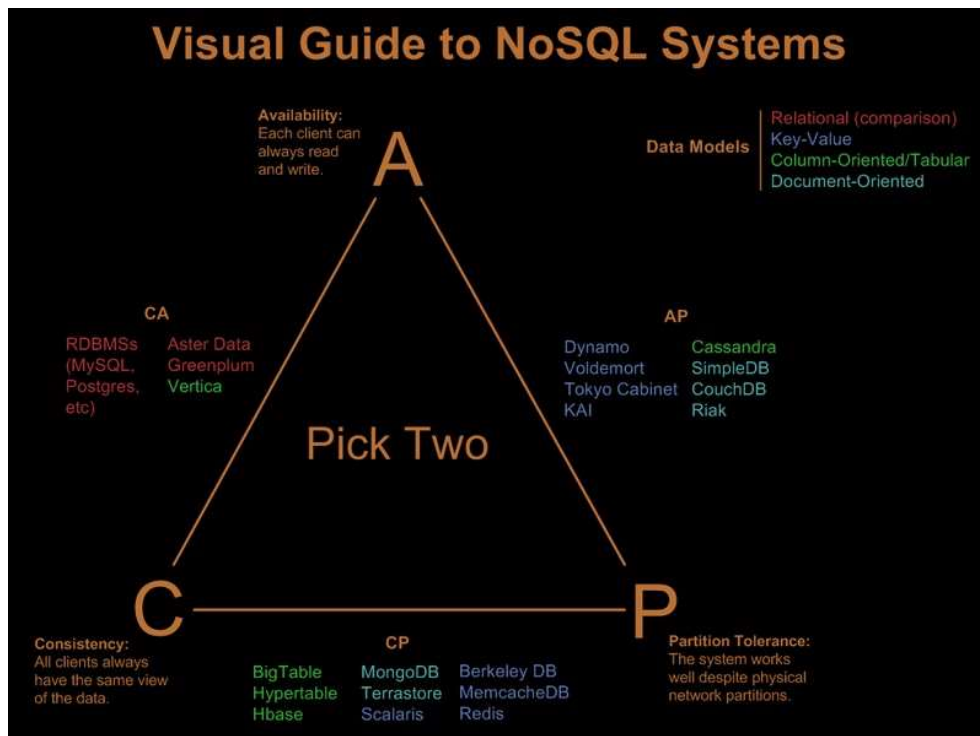    - leases, conflict resolution,
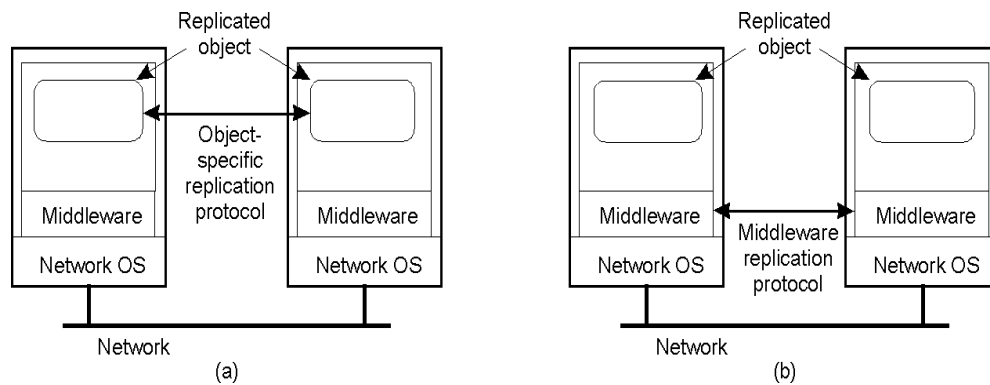
# NoSQL Systems and CAP



Figure Courtesy of Nathan Hurst
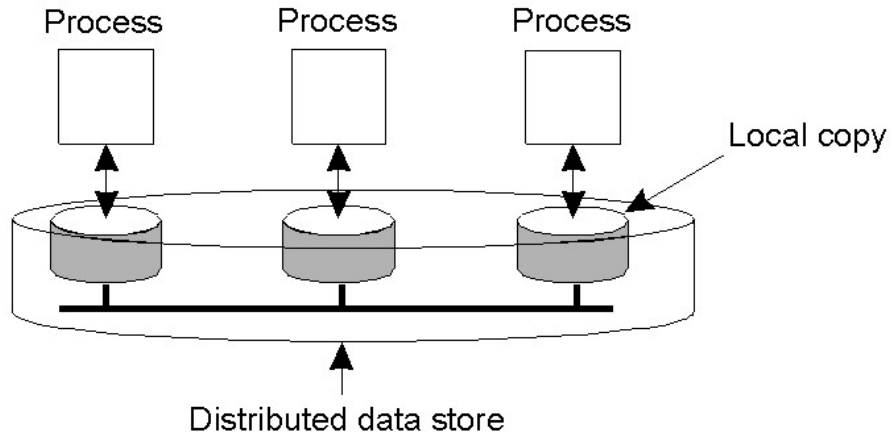
# Object Replication



- •Approach 1: application is responsible for replication
  - – Application needs to handle consistency issues
- •Approach 2: system (middleware) handles replication
  - – Consistency issues are handled by the middleware
  - – Simplifies application development but makes object-specific solutions harder

# Replication and Scaling

- • Replication and caching used for system scalability
- • Multiple copies:
  - – Improves performance by reducing access latency
  - – But higher network overheads of maintaining consistency
  - – Example: object is replicated $N$ times
    - • Read frequency $R$, write frequency $W$
    - • If $R<<W$, high consistency overhead and wasted messages
    - • Consistency maintenance is itself an issue
      - – What semantics to provide?
      - – Tight consistency requires globally synchronized clocks!
- • Solution: loosen consistency requirements
  - – Variety of consistency semantics possible

# Data-Centric Consistency Models



Distributed data store

- Consistency model (aka *consistency semantics*)
  - Contract between processes and the data store
    - If processes obey certain rules, data store will work correctly
  - All models attempt to return the results of the last write for a read operation
    - Differ in how "last" write is determined/defined

# Strict Consistency

- Any read always returns the result of the most recent write
  - Implicitly assumes the presence of a global clock
  - A write is immediately visible to all processes
    - Difficult to achieve in real systems (network delays can be variable)

# Sequential Consistency

- Sequential consistency: weaker than strict consistency
  - Assumes all operations are executed in some sequential order and each process issues operations in program order
    - Any valid interleaving is allowed
    - All agree on the same interleaving
    - Each process preserves its program order
    - Nothing is said about "most recent write"

```
P1:  W(x)a                               P1:  W(x)a
P2:        W(x)b                         P2:        W(x)b
P3:              R(x)b      R(x)a        P3:              R(x)b      R(x)a
P4:                  R(x)b  R(x)a        P4:                  R(x)a  R(x)b

            (a)                                      (b)
```

# Linearizability

- Assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
  - Difference between linearizability and serializability?
    - Granularity: reads/writes versus transactions
- Example:

| Process P1 | Process P2 | Process P3 |
|------------|------------|------------|
| x = 1;     | y = 1;     | z = 1;     |
| print ( y, z); | print (x, z); | print (x, y); |

# Linearizability Example

- Four valid execution sequences for the processes of the previous slide. The vertical axis is time.

| | | | |
|---|---|---|---|
| x = 1; | x = 1; | y = 1; | y = 1; |
| print ((y, z); | y = 1; | z = 1; | x = 1; |
| y = 1; | print (x,z); | print (x, y); | z = 1; |
| print (x, z); | print(y, z); | print (x, z); | print (x, z); |
| z = 1; | z = 1; | x = 1; | print (y, z); |
| print (x, y); | print (x, y); | print (y, z); | print (x, y); |
| | | | |
| Prints:  001011 | Prints: 101011 | Prints: 010111 | Prints: 111111 |
| | | | |
| Signature: | Signature: | Signature: | Signature: |
| 001011 | 101011 | 110101 | 111111 |
| (a) | (b) | (c) | (d) |

# Causal consistency

- Causally related writes must be seen by all processes in the same order.
  – Concurrent writes may be seen in different orders on different machines

| | | | | |
|---|---|---|---|---|
| P1: W(x)a | | | | |
| P2: | R(x)a | W(x)b | | |
| P3: | | | R(x)b | R(x)a |
| P4: | | | R(x)a | R(x)b |
| | | (a) | | |

| | | | | |
|---|---|---|---|---|
| P1: W(x)a | | | | |
| P2: | | W(x)b | | |
| P3: | | | R(x)b | R(x)a |
| P4: | | | R(x)a | R(x)b |
| | | (b) | | |

Not permitted                              Permitted

# Other models

- FIFO consistency: writes from a process are seen by others in the same order. Writes from different processes may be seen in different order (even if causally related)
  - Relaxes causal consistency
  - Simple implementation: tag each write by (Proc ID, seq #)
- Even FIFO consistency may be too strong!
  - Requires all writes from a process be seen in order
- Assume use of critical sections for updates
  - Send final result of critical section everywhere
  - Do not worry about propagating intermediate results
    - Assume presence of synchronization primitives to define semantics

# Other Models

**Use granularity of critical sections, instead of individual read/write**

- Weak consistency
  - Accesses to synchronization variables associated with a data store are sequentially consistent
  - No operation on a synchronization variable is allowed to be performed until all previous writes have been completed everywhere
  - No read or write operation on data items are allowed to be performed until all previous operations to synchronization variables have been performed.
- Entry and release consistency
  - Assume shared data are made consistent at entry or exit points of critical sections

# Summary of Data-centric Consistency Models

| Consistency | Description |
| --- | --- |
| Strict | Absolute time ordering of all shared accesses matters. |
| Linearizability | All processes must see all shared accesses in the same order. Accesses are furthermore ordered according to a (nonunique) global timestamp |
| Sequential | All processes see all shared accesses in the same order. Accesses are not ordered in time |
| Causal | All processes see causally-related shared accesses in the same order. |
| FIFO | All processes see writes from each other in the order they were used. Writes from different processes may not always be seen in that order |

(a)

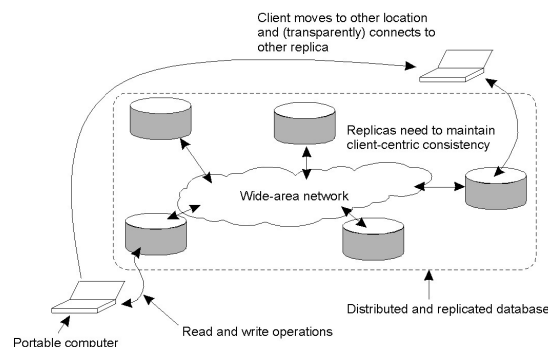| Consistency | Description |
| --- | --- |
| Weak | Shared data can be counted on to be consistent only after a synchronization is done |
| Release | Shared data are made consistent when a critical region is exited |
| Entry | Shared data pertaining to a critical region are made consistent when a critical region is entered. |

(b)

# Client-centric Consistency Models

- Assume read operations by a single process *P* at two *different* local copies of the same data store
  - Four different consistency semantics
- *Monotonic reads*
  - Once read, subsequent reads on that data items return same or more recent values
- *Monotonic writes*
  - A write must be propagated to all replicas before a successive write by the *same process*
  - Resembles FIFO consistency (writes from same process are processed in same order)
- *Read your writes*: read(x) always returns write(x) by that process
- *Writes follow reads*: write(x) following read(x) will take place on same or more recent version of x

# Eventual Consistency

- Many systems: one or few processes perform updates
  - How frequently should these updates be made available to other read-only processes?
- Examples:
  - DNS: single naming authority per domain
  - Only naming authority allowed updates (no write-write conflicts)
  - How should read-write conflicts (consistency) be addressed?
  - NIS: user information database in Unix systems
    - Only sys-admins update database, users only read data
    - Only user updates are changes to password

# Eventual Consistency

- Assume a replicated database with few updaters and many readers
- Eventual consistency: in absence of updates, all replicas converge towards identical copies
  - Only requirement: an update should eventually propagate to all replicas
  - Cheap to implement: no or infrequent write-write conflicts
  - Things work fine so long as user accesses same replica
  - What if they don't:

Client moves to other location and (transparently) connects to other replica

Replicas need to maintain client-centric consistency

Wide-area network

Distributed and replicated database

Read and write operations

Portable computer

# Epidemic Protocols

- Used in Bayou system from Xerox PARC
- Bayou: weakly connected replicas
  - Useful in mobile computing (mobile laptops)
  - Useful in wide area distributed databases (weak connectivity)
- Based on theory of epidemics *(spreading infectious diseases)*
  - Upon an update, try to "infect" other replicas as quickly as possible
  - Pair-wise exchange of updates (*like pair-wise spreading of a disease)*
  - Terminology:
    - Infective store: store with an update it is willing to spread
    - Susceptible store: store that is not yet updated
- Many algorithms possible to spread updates

# Spreading an Epidemic

- Anti-entropy
  - Server *P* picks a server *Q* at random and exchanges updates
  - Three possibilities: only push, only pull, both push and pull
  - Claim: A pure push-based approach does not help spread updates quickly (Why?)
    - Pull or initial push with pull work better
- Rumor mongering (aka *gossiping*)
  - Upon receiving an update, *P* tries to push to *Q*
  - If *Q* already received the update, stop spreading with prob *1/k*
  - Analogous to "hot" gossip items => stop spreading if "cold"
  - Does not guarantee that all replicas receive updates
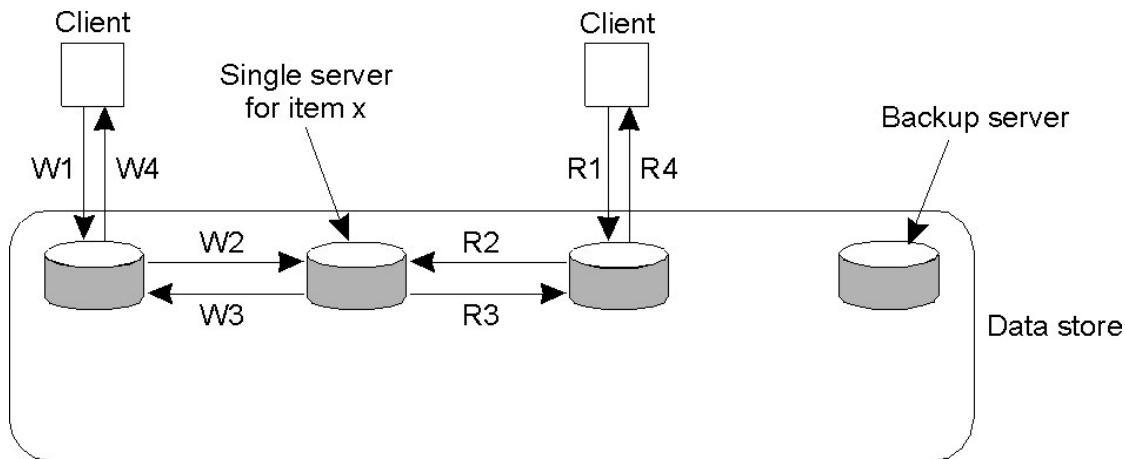    - Chances of staying susceptible: $s = e^{-(k+1)(1-s)}$

# Removing Data

- Deletion of data items is hard in epidemic protocols
- Example: server deletes data item *x*
  - No state information is preserved
    - Can't distinguish between a deleted copy and no copy!
- Solution: death certificates
  - Treat deletes as updates and spread a death certificate
    - Mark copy as deleted but don't delete
    - Need an eventual clean up
      - Clean up dormant death certificates

# Implementation Issues

- Two techniques to implement consistency models
  - Primary-based protocols
    - Assume a primary replica for each data item
    - Primary responsible for coordinating all writes
  - Replicated write protocols
    - No primary is assumed for a data item
    - Writes can take place at any replica

# Remote-Write Protocols

Client

Client

Single server
for item x

Backup server

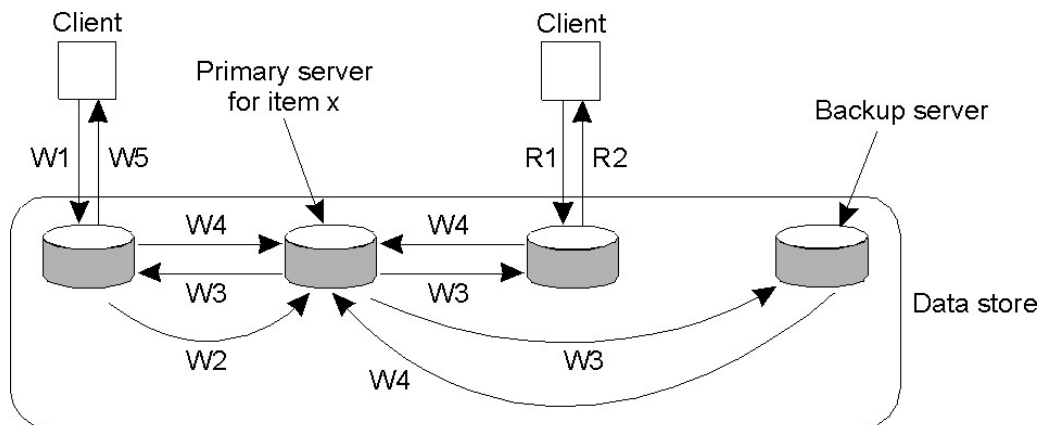W1  W4

R1  R4

W2

R2

W3

R3

Data store

W1. Write request
W2. Forward request to server for x
W3. Acknowledge write completed
W4. Acknowledge write completed

R1. Read request
R2. Forward request to server for x
R3. Return response
R4. Return response

- **Traditionally used in client-server systems (no replication)**

# Remote-Write Protocols (2)

Client

Client

Primary server
for item x

Backup server

W1  W5
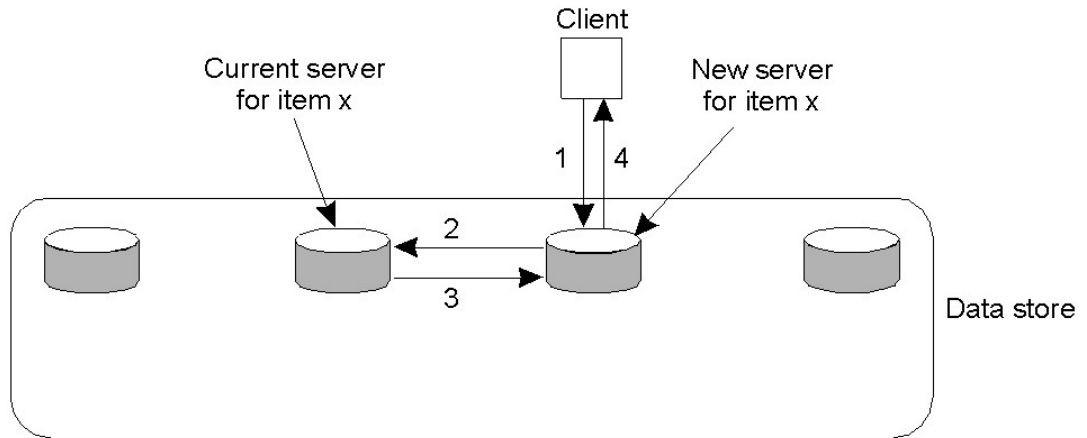
R1  R2

W4

W4

W3

W3

W2

W3

W4

Data store

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

- Primary-backup protocol
  - Allow local reads, sent writes to primary
  - Block on write until all replicas are notified
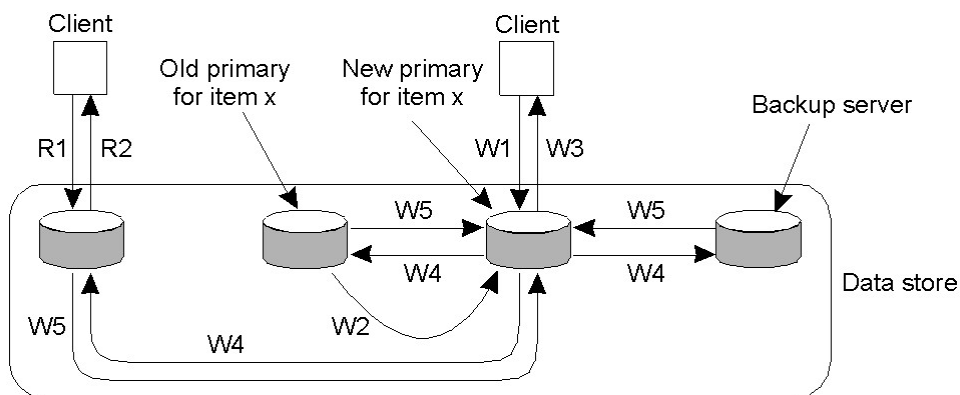  - Implements sequential consistency

# Local-Write Protocols (1)



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

- Primary-based local-write protocol in which a single copy is migrated between processes.
  - Limitation: need to track the primary for each data item
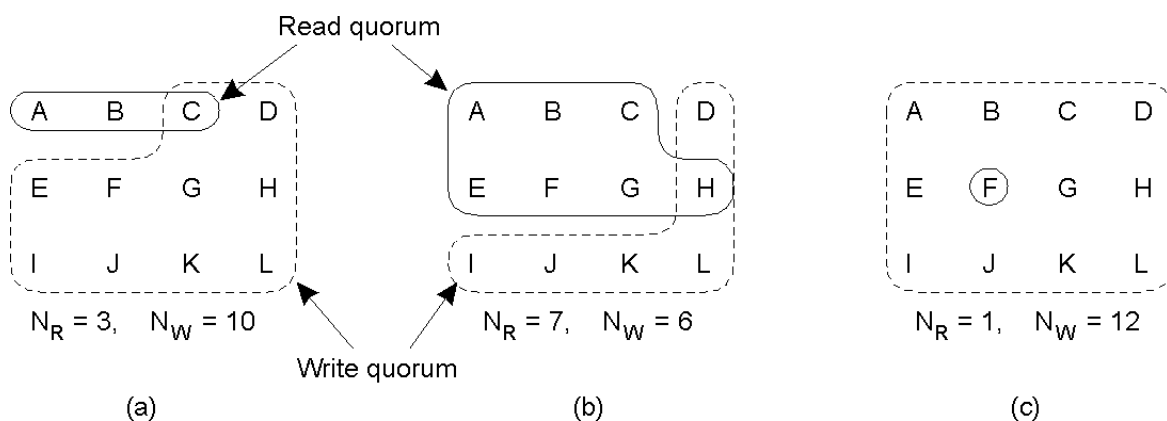
# Local-Write Protocols (2)



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

R1. Read request
R2. Response to read

- Primary-backup protocol in which the primary migrates to the process wanting to perform an update

# Replicated-write Protocols

- Relax the assumption of one primary
  - No primary, any replica is allowed to update
  - Consistency is more complex to achieve
- Quorum-based protocols
  - Use voting to request/acquire permissions from replicas
  - Consider a file replicated on $N$ servers
    - $N_R + N_W > N$      $N_W > N/2$
  - Update: contact $N_W$ servers and get them to agree to do update (associate version number with file)
  - Read: contact $N_R$ and obtain version number
    - If all servers agree on a version number, read

# Gifford's Quorum-Based Protocol

Read quorum

| A | B | C | D |
|---|---|---|---|
| E | F | G | H |
| I | J | K | L |

$N_R = 3, \quad N_W = 10$

| A | B | C | D |
|---|---|---|---|
| E | F | G | H |
| I | J | K | L |

$N_R = 7, \quad N_W = 6$

Write quorum

(a)                                (b)

| A | B | C | D |
|---|---|---|---|
| E | F | G | H |
| I | J | K | L |

$N_R = 1, \quad N_W = 12$

(c)

- Three examples of the voting algorithm:
- a)  A correct choice of read and write set
- b)  A choice that may lead to write-write conflicts
- c)  A correct choice, known as ROWA (read one, write all)

# Replica Management

- Replica server placement
  - Web: geophically skewed request patterns
  - Where to place a proxy?
    - K-clusters algorithm
- Permanent replicas versus temporary
  - Mirroring: all replicas mirror the same content
  - Proxy server: on demand replication

- Server-initiated versus client-initiated

# Content Distribution

- Will come back to this in Chap 12

- CDN: network of proxy servers
- Caching:
  - update versus invalidate
  - Push versus pull-based approaches
  - Stateful versus stateless
- Web caching: what semantics to provide?

# Final Thoughts

- Replication and caching improve performance in distributed systems
- Consistency of replicated data is crucial
- Many consistency semantics (models) possible
  - Need to pick appropriate model depending on the application
  - Example: web caching: weak consistency is OK since humans are tolerant to stale information (can reload browser)
  - Implementation overheads and complexity grows if stronger guarantees are desired

# Fault Tolerance

- Single machine systems
  - Failures are all or nothing
    - OS crash, disk failures
- Distributed systems: multiple independent nodes
  - Partial failures are also possible (some nodes fail)
- *Question:* Can we automatically recover from partial failures?
  - Important issue since probability of failure grows with number of independent components (nodes) in the systems
  - Prob(failure) = Prob(Any one component fails)=1-P(no failure)