

# Part 1: Replication Basics



- 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)

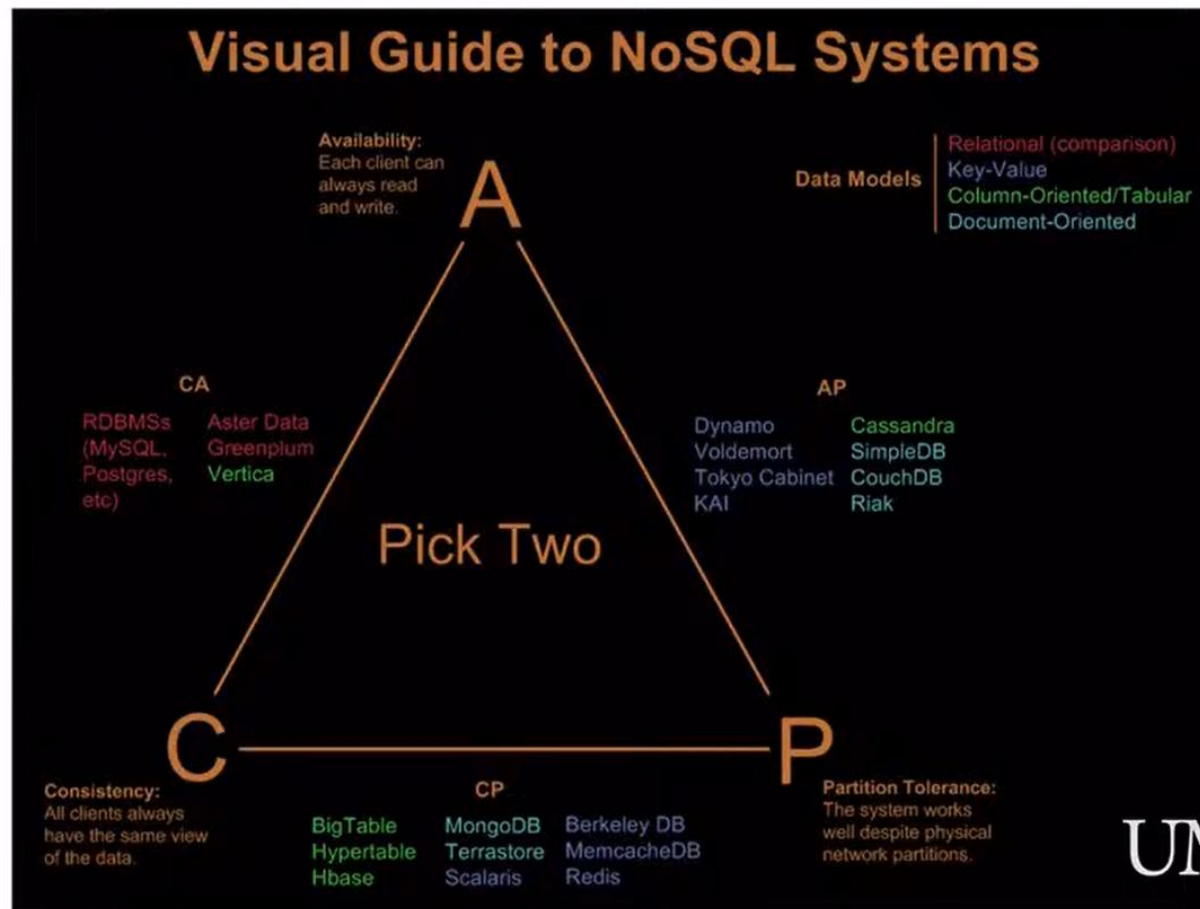
# CAP Theorem Examples

- 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



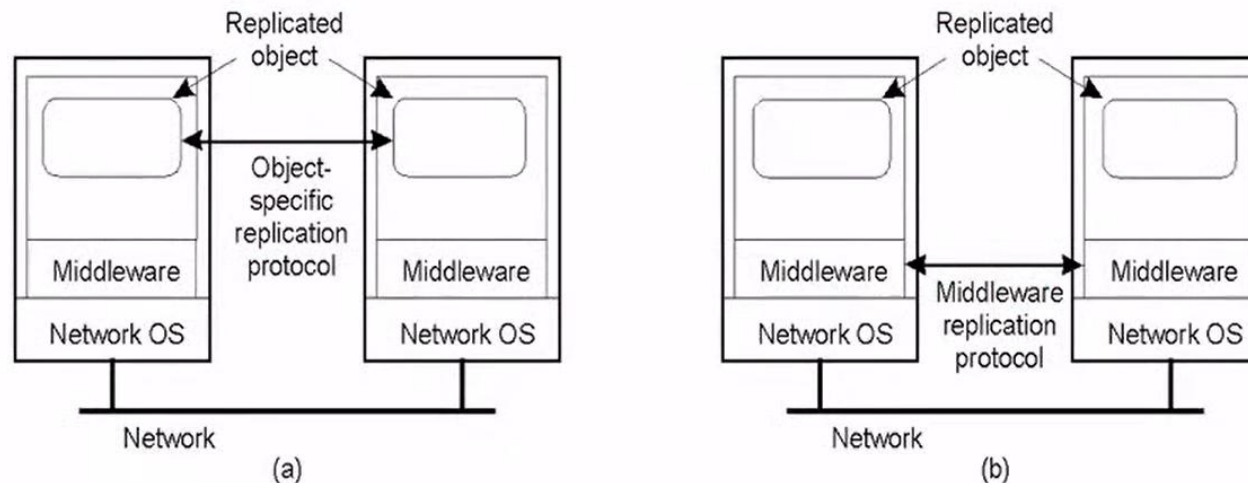
Prashant Shen...

UMassAmherst

Figure Courtesy of Nathan Hurst

Powered by Zoom

# 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

Amherst

Powered by Zoom

# Replication and Scaling

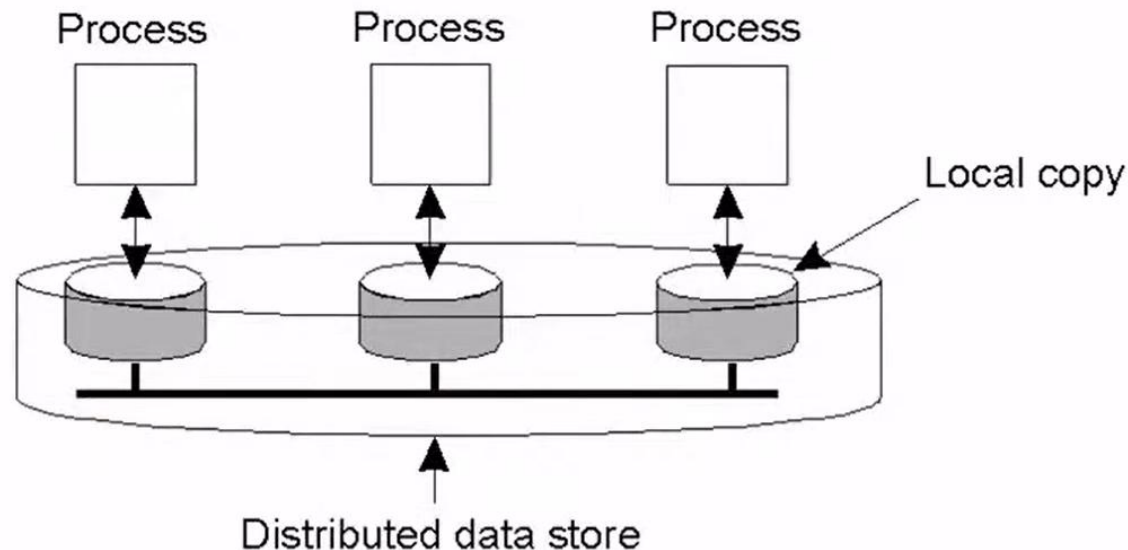


- Replication and caching used for system scalability
- Multiple copies:
  - Improves performance by reducing access latency
  - But higher network overheads of maintaining consistency
  - Example: object is replicated  $N$  times
    - Read frequency  $R$ , write frequency  $W$
    - If  $R \ll W$ , 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

Amherst

Powered by Zoom

## Part 2: Data-Centric Consistency Model



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

(a)

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

(b)

# Linearizability



Prashant Shen...

- 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                         |
|------------------------------------|------------------------------------|------------------------------------|
| $x = 1;$<br>$\text{print} (y, z);$ | $y = 1;$<br>$\text{print} (x, z);$ | $z = 1;$<br>$\text{print} (x, y);$ |

# Linearizability Example



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

```
x = 1;  
print ((y, z);  
y = 1;  
print (x, z);  
z = 1;  
print (x, y);
```

Prints: 001011

Signature:  
001011  
(a)

```
x = 1;  
y = 1;  
print (x,z);  
print(y, z);  
z = 1;  
print (x, y);
```

Prints: 101011

Signature:  
101011  
(b)

```
y = 1;  
z = 1;  
print (x, y);  
print (x, z);  
x = 1;  
print (y, z);
```

Prints: 010111

Signature:  
110101  
(c)

```
y = 1;  
x = 1;  
z = 1;  
print (x, z);  
print (y, z);  
print (x, y);
```

Prints: 111111

Signature:  
111111  
(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)

Not permitted

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

(b)

Permitted

# 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)

Not permitted

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

(b)

Permitted



Live chat replay was turned off for this video.

## UMass CS677 Distributed Systems - ...

UMass OS - 16 / 27



- 15 UMass CS 677 - Spring 22 - Lecture 15 - Distributed Transactions UMass OS 1:05:56
- ▶ 16 UMass CS677 - Spring'22 - Lecture 16 - Consistency Models in... UMass OS 1:18:43
- 17 UMass CS677 (Spring 22) Lecture 17: Replication UMass OS 1:10:29
- 18 UMass CS677 (Spring'22) Lecture 18: Byzantine Generals Problem,... UMass OS 1:11:54
- 19 UMass CS677 (Spring'22) Lecture 19: Distributed Consensus UMass OS 1:11:49
- 20 UMass CS677 (Spring '22) Lecture 20: Distributed Web Applications... UMass OS 1:13:04

48:39 / 1:18:42 • Causal Consistency > UMassAmherst CS677: Distributed and Operating Systems

CC lectur, pag Powered by Zoom

UMass CS677 Distributed Systems - Spring 22

## UMass CS677 - Spring'22 - Lecture 16 - Consistency Models in Distributed Systems

U UMass OS 5.81K subscribers

Subscribe

👍 5

👎

➦ Share

🔼 Save

⋮

# 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 Model



- 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*:  $\text{read}(x)$  always returns  $\text{write}(x)$  by that process
- *Writes follow reads*:  $\text{write}(x)$  following  $\text{read}(x)$  will take place on same or more recent version of  $x$

# Part 3: Eventual Consistency



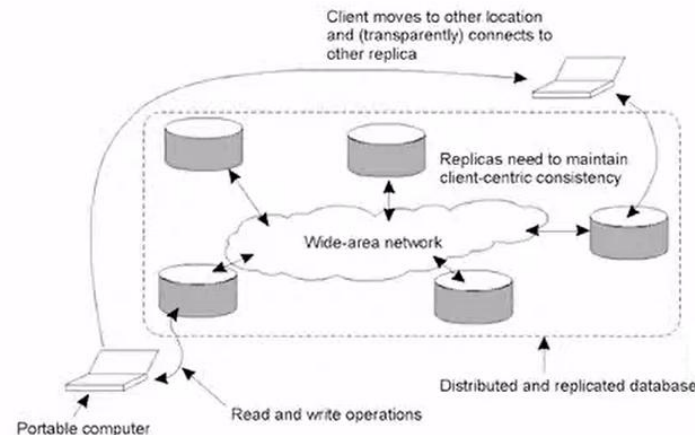
- Many systems: one or few processes perform updates
  - How frequently should these updates be made available to other read-only processes?
- Examples:
  - DNS: single naming authority per domain
  - Only naming authority allowed updates (no write-write conflicts)
  - How should read-write conflicts (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
  - Cloud storage: dropbox, OneDrive, iCloud all use eventual consistency



# 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:





# 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  $\Rightarrow$  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

