Lecture 3a: Replication & Consistency



**Slide 1**

① Replication
② Consistency
  • Models vs Protocols
③ Update propagation

## 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):



What's the difference between mirroring and caching?

**Slide 4**

## Control Replication:



**Slide 5**

What are thethe challenges of doing this?

## Data and Control Replication:



**Slide 6**

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

## REPLICATION ISSUES

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

**Slide 7** Replica placement
➜ How many replicas?
➜ Where to put them?

Redirection/Routing
➜ Which replica should clients use?

## DISTRIBUTED DATA STORE

➜ data-store stores data items

## Client's Point of View:



**Slide 8**

## Slide 9

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



Data Store

## Slide 10

### Data Model:

➔ data item: simple variable

➔ data item values: explicit (`0`, `1`), abstract (`a`,`b`)

➔ data store: collection of data items

### 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:

➔ ClientA/Replica1: `WA(x)1,WA(x)0`

➔ ClientB/Replica2: `RB(x)-,RB(x)1,RB(x)1,RB(x)0`

## CONSISTENCY

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

Consistency Dimensions:
➔ Time and Order

Time:
➔ How old is the data (staleness)?
➔ How old is the data allowed to be?
  • Time, Versions

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

Real world examples of inconsistency?

**Slide 13**

---

## ORDERING

Updates and concurrency result in conflicting operations

Conflicting Operations:
➔ Read-write conflict (only 1 write)
➔ Write-write conflict (multiple concurrent writes)
➔ The order in which conflicting operations are performed affects consistency

Partial vs Total Ordering:
➔ partial order: order of a single client's operations
➔ total order: interleaving of all conflicting operations

**Slide 14**

---

## Example:

**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` |

Can you sanely use a system like this?

**Slide 15**

---

## CONSISTENCY MODEL

*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 16**

## 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

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

## 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.

  ➔ Multiple clients accessing the same data store
  ➔ Described consistency is experienced by all clients
    • 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 replicas 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)a W(x)c} == {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

## LINEARISABLE CONSISTENCY

*All operations are performed in a single sequential order*

➜ Operations ordered according to a global (finite) timestamp.
➜ Program order of each client maintained

**Slide 21**

| | | |
|---|---|---|
| Client A | W(x)a | |
| Client B | W(x)b | |
| Client C | R(x)a | R(x)b |
| Client D | R(x)a | R(x)b |

linearisable

| | | |
|---|---|---|
| Client A | W(x)a | |
| Client B | W(x)b | |
| Client C | R(x)b | R(x)a |
| Client D | R(x)b | R(x)a |

not linearisable

## SEQUENTIAL CONSISTENCY

*All operations are performed in some sequential order*

➜ More than one correct sequential order possible
➜ All clients see the *same* order
➜ Program order of each client maintained
➜ Not ordered according to time Why is this good?

**Slide 22**

| | | |
|---|---|---|
| Client A | W(x)a | |
| Client B | W(x)b | |
| Client C | R(x)b | R(x)a |
| Client D | R(x)b | R(x)a |

sequential

| | | |
|---|---|---|
| Client A | W(x)a | |
| Client B | W(x)b | |
| Client C | R(x)a | R(x)b |
| Client D | R(x)b | R(x)a |

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)

**Slide 23**

| | | |
|---|---|---|
| Client A | W(x)a | W(x)c |
| Client B | W(x)b | |
| Client C | R(x)a R(x)b R(x)c | |
| Client D | R(x)b R(x)a R(x)c | |

causally consistent

| | | |
|---|---|---|
| Client A | W(x)a | W(x)c |
| Client B | R(x)a → W(x)b | |
| Client C | R(x)a R(x)b R(x)c | |
| Client D | R(x)b R(x)a R(x)c | |

not causally consistent

How could we make this valid?

## FIFO (PRAM) CONSISTENCY

*Only partial orderings of writes maintained*

**Slide 24**

| | | |
|---|---|---|
| Client A | W(x)a | W(x)c |
| Client B | R(x)a W(x)b | |
| Client C | R(x)a R(x)b R(x)c | |
| Client D | R(x)b R(x)a R(x)c | |

FIFO consistent

| | | |
|---|---|---|
| Client A | W(x)a | W(x)c |
| Client B | R(x)a W(x)b | |
| Client C | R(x)a R(x)b R(x)c | |
| Client D | R(x)c R(x)a R(x)b | |

not FIFO consistent

How could we make this valid?

## 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

**Slide 25**

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:

➔ 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

**Slide 26**



weak consistent    not weak consistent

How could we make this valid?

---

## 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'

**Slide 27**

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 28**



release consistent

What is an example of an invalid ordering?

## Lazy Release Consistency:

➔ 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)

```
            Acq(S)  W(x)a  W(x)b   Rel(S)
Client A _____
                              Acq(S) R(x)b Rel(S)
Client B _____
                                       R(x)a
Client C _____
```

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)

## 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

```
            Acq(Sx)  W(x)a  Acq(Sy)  W(y)b  Rel(Sx) Rel(Sy)
Client A _____
                              Acq(Sx)  R(x)a    R(y)Nil
Client B _____
                                          Acq(Sx)   R(y)b
Client C _____
```

entry consistent

## CAP THEORY

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



You can only choose two of
C A or P

## CAP Theory

**Slide 33**

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

Consistency

Partition Tolerance

Availability

You can only choose two of
C A or P

## CAP Theory

**Slide 34**

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

Consistency

Partition Tolerance

Availability

You can only choose two of
C A or P

## CAP Theory

**Slide 35**

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

Consistency

Partition Tolerance

Availability

You can only choose two of
C A or P

**Slide 36**

CAP Impossibility Proof:

Replica A    0

Replica B    0

Client

## CAP Impossibility Proof:

**Slide 37**



## CAP Impossibility Proof:

**Slide 38**



## CAP Impossibility Proof:

**Slide 39**



No Consistency

## CAP Impossibility Proof:

**Slide 40**



No Availability

CAP Impossibility Proof:



Write (fails)

**No Partition Tolerance**

Client

---

## CAP CONSEQUENCES

For wide-area systems:

➜ Must choose: Consistency or Availability
➜ Choosing Availability
  • Give up on consistency?
  • Eventual consistency
➜ Choosing Consistency
  • No availability
  • delayed (and potentially failing) operations

Why can't we choose C and A and forget about P?

---

## EVENTUAL CONSISTENCY

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



eventual consistent

Requirements:

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

---

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
  • manual merging for write-write conflicts

---

## 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

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.

---

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

• Data-items have an owner

• No write-write conflicts

---

### 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)

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

---

## 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*

```
           WS(x1)        R(x1)                  WS(x1)      R(x1)
Replica 1 ─────────────────────       Replica 1 ────────────────────

              WS(x1;x2)      R(x2)                WS(!x1;x2)  R(x2)  WS(x1;x2)
Replica 2 ─────────────────────       Replica 2 ────────────────────

        monotonic-read consistent              not monotonic-read consistent
```

When is Monotonic Reads sufficient?

## 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.

**Slide 49**

| | |
|---|---|
| Replica 1  W(x1) | Replica 1  W(x1) |
| Replica 2  W(x1) WS(x1)  W(x2) | Replica 2  WS(!x1;x0)  W(x2) |
| 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*

**Slide 50**

| | |
|---|---|
| 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 |

When is Read Your Writes sufficient?

## 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*

**Slide 51**

| | |
|---|---|
| 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 |

When is Write Follows Reads sufficient?

## CHOOSING THE RIGHT MODEL

Trade-offs

Consistency and Redundancy:

➜ All copies must be strongly consistent

➜ All copies must contain full state

➜ Reduced consistency → reduced reliability

**Slide 52**

Consistency and Performance:

➜ Consistency requires extra work and communication

✗ Can result in loss of overall performance

☑ Weaker consistency possible

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

---

# REMOTE-WRITE PROTOCOLS

## Single Server:

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



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

---

## Primary-Backup:

➜ All writes executed at single server, Reads are local
➜ Updates block until executed on all backups
✗ Performance

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

---

# LOCAL-WRITE PROTOCOLS

## Migration:

➜ Data item migrated to local server on access
☑ Performance (when not sharing data)

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):

- ☑ Performance for concurrent reads
- ✗ Performance for concurrent writes



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

---

## QUORUM-BASED PROTOCOLS

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

$N_R = 3, \quad N_W = 10$ (a)

$N_R = 7, \quad N_W = 6$ (b)

$N_R = 1, \quad N_W = 12$ (c)

---

## ACTIVE REPLICATION

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

---

## 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

## Push Update Propagation:

What to propagate?
→ Data
  • R/W high
→ Update operation
  • low bandwidth costs
→ Notification/Invalidation
  • R/W low

## 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

→ Server-initiated replication
---▶ Client-initiated replication

Permanent replicas
Server-initiated replicas
Client-initiated replicas
Clients

## 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

## 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

Dynamic Replica Placement:

**Slide 65**
➜ Network of replica servers
➜ Keep track of data item requests at each replica
➜ Thresholds:
  • 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?
**Slide 66**
➜ One-size-fits-all?

Determining Policy:

➜ Who determines the consistency model used?
  • Application, Middleware
  • Client, Server

Keep It Simple, Stupid:

➜ Will the programmer understand the consistency model?

## READING LIST

**Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services** An overview of
**Slide 67** the CAP theorem and its proof.

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

## HOMEWORK

Consistency Models:

➜ Research consistency models used in existing Distributed Systems
➜ Why are those models being used?
➜ In the systems you looked at, could other models have been used? Would that have made the system better?

**Slide 68** Hacker's Edition:

➜ Find a system that provides Eventual Consistency
  ➜ (alternatively, implement (possibly in Erlang) a system that provides Eventual Consistency)
➜ Replicate some data and perform queries. How often do you get inconsistent results?
➜ If you can tweak replication parameters, how do they affect the consistency of results?