

## CSE 486/586 Distributed Systems Consistency --- 1

Steve Ko  
Computer Sciences and Engineering  
University at Buffalo

CSE 486/586

## Recap: Concurrency (Transactions)

- Question: How to support transactions (with locks)?
  - Multiple transactions share data.
- First strategy: Complete serialization
  - One transaction at a time with one big lock
  - Correct, but at the cost of performance
- How to improve performance?
  - Let's see if we can concurrently execute transactions.

CSE 486/586

2

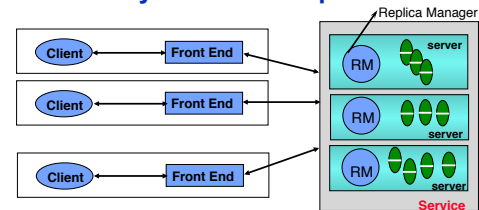
## Recap: Concurrency (Transactions)

- Problem: Not all current executions produce a correct outcome
  - Serial equivalence & strict execution must be met.
- How do we meet the requirements using locks?
  - Overall strategy: using more and more fine-grained locking
  - No silver bullet. Fine-grained locks have their own implications.
  - Exclusive locks (per-object locks)
  - Non-Exclusive locks (read/write locks)
  - Other finer-grained locks (e.g., two-version locking)
- Atomic commit problem
  - Commit or abort (consensus)
  - 2PC

CSE 486/586

3

## Consistency with Data Replicas



- Consider that this is a distributed storage system that serves read/write requests.
- Multiple copies of a same object stored at different servers
- Question: How to maintain consistency across different data replicas?

CSE 486/586

4

## Consistency

- Why replicate?
- Increased availability of service. When servers fail or when the network is partitioned.
  - $P$ : 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
- Fault tolerance
  - Under the fail-stop model, if up to  $f$  of  $f+1$  servers crash, at least one is alive.
- Load balancing
  - One approach: Multiple server IPs can be assigned to the same name in DNS, which returns answers round-robin.

CSE 486/586

5

## This Week

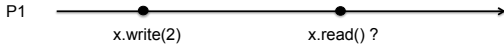
- We will look at different consistency guarantees (models).
- We'll start from the strongest guarantee, and gradually relax the guarantees.
  - Linearizability (or sometimes called strong consistency)
  - Sequential consistency
  - Causal consistency
  - Eventual consistency
- Different applications need different consistency guarantees.
- This is all about client-side perception.
  - When a read occurs, what do you return?
- First
  - Linearizability: we'll look at the concept first, then how to implement it later.

CSE 486/586

6

## Our Expectation with Data

- Consider a single process using a filesystem
- What do you expect to read?



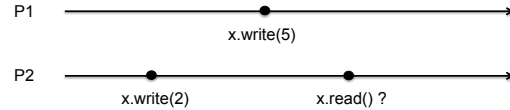
- Our expectation (as a user or a developer)
  - A read operation returns the most recent write.
  - This forms our basic expectation from any file or storage system.
- **Linearizability** meets this basic expectation.
  - But it extends the expectation to handle **multiple processes**...
  - ...and **multiple replicas**.
  - The strongest consistency model

CSE 486/586

7

## Expectation with Multiple Processes

- What do you expect to read?
  - A single filesystem with multiple processes



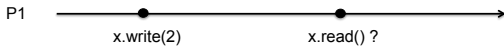
- Our expectation (as a user or a developer)
  - A read operation returns the most recent write, **regardless of the clients**.
  - We expect that a read operation returns the most recent write **according to the single actual-time order**.
  - In other words, read/write should behave **as if there were a single (combined) client making all the requests**.

CSE 486/586

8

## Expectation with Multiple Copies

- What do you expect to read?
  - A single process with multiple servers with copies



- Our expectation (as a user or a developer)
  - A read operation returns the most recent write, **regardless of how many copies there are**.
  - Read/write should behave **as if there were a single copy**.

CSE 486/586

9

## Linearizability

- **Three aspects**
  - A read operation returns the most recent write,
    - ...regardless of the clients,
    - ...according to the single actual-time ordering of requests.
- Or, put it differently, read/write should behave as if there were,
  - ...a single client making all the (combined) requests in their original actual-time order,
  - ...over a single copy.
- You can say that **your storage system guarantees linearizability when it provides single-client, single-copy semantics where a read returns the most recent write**.

CSE 486/586

10

## Linearizability Exercise

- Assume that the following happened with object x over a linearizable storage.
  - C1: x.write(A)
  - C2: x.write(B)
  - C3: x.read() → B, x.read() → A
  - C4: x.read() → B, x.read() → A
- What would be an actual-time ordering of the events?
  - One possibility: C2 (write B) → C3 (read B) → C4 (read B) → C1 (write A) → C3 (read A) → C4 (read A)
- How about the following?
  - C1: x.write(A)
  - C2: x.write(B)
  - C3: x.read() → B, x.read() → A
  - C4: x.read() → A, x.read() → B

CSE 486/586

11

## CSE 486/586 Administrivia

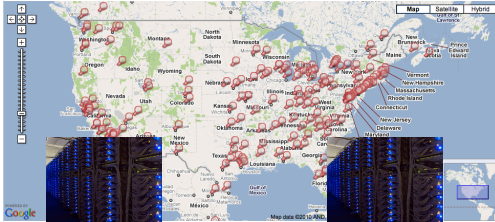
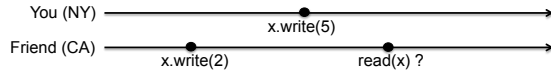
- PA3 deadline: 4/3 (Friday)
- Grading is going on with PA2B and midterm.

CSE 486/586

12

## Linearizability Subtleties

- Notice any problem?



California CSE 486/586

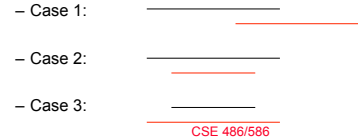
North Carolina

13

## Linearizability Subtleties

- A read/write operation is never a dot!
  - It takes time. Many things are involved, e.g., network, multiple disks, etc.
  - Read/write latency: the time measured right before the call and right after the call from the client making the call.
- Clear-cut (e.g., black---write & red---read)

- Not-so-clear-cut (parallel)



CSE 486/586

14

## Linearizability Subtleties

- Let's go back to the single-client, single-copy semantics.
- With a single process and a single copy, can overlaps happen?
  - No, these are cases that do not arise with a single process and a single copy.
- Thus, we (as a system designer) have **freedom to impose an order**.
  - Linearizability does not mandate any particular order for overlapping operations.
  - You can implement a particular ordering strategy.
  - As long as there is a single, interleaving ordering for overlapping operations, it's fine.
  - This ordering should still provide the single-client, single-copy semantics.

CSE 486/586

15

## Linearizability Subtleties

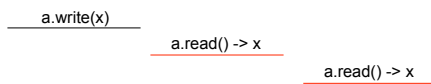
- Definite guarantee
- Relaxed guarantee when overlap
- Case 1  $\text{---}$   $\text{---}$
- Case 2  $\text{---}$   $\text{---}$
- Case 3  $\text{---}$   $\text{---}$

CSE 486/586

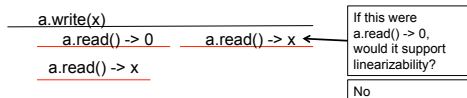
16

## Linearizability Examples

- Example 1



- Example 2

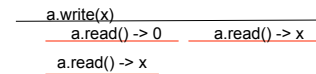


CSE 486/586

17

## Linearizability Examples

- In example 2, what are the constraints?



- Constraints
  - $\text{a.read() -> 0}$  happens before  $\text{a.read() -> x}$  (cannot change the order).
  - $\text{a.read() -> x}$  happens before  $\text{a.read() -> x}$  (cannot change the order).
  - The rest are up for grabs.

CSE 486/586

18

## Linearizability Examples

- In example 2, why would a.read() return 0 and x when they're overlapping?

```
_____
a.write(x)
_____
a.read() -> 0   a.read() -> x
_____
a.read() -> x
```

- This assumes that there's a particular storage system that shows this behavior.
- At some point between a read/write request sent and returned, the result becomes visible.
  - E.g., you read a value from physical storage, *prepare it for return (e.g., putting it in a return packet, i.e., making it visible)*, and actually return it.
  - Or you *actually write a value to a physical disk, making it visible* (out of multiple disks, which might actually write at different points).

CSE 486/586

19

## Linearizability Examples

- Example 3

```
_____
a.write(x)
_____
a.read() -> x   a.read() -> x
_____
a.read() -> y
_____
a.write(y)
```

- Constraints
  - a.read() → x and a.read() → x: we cannot reorder these.
  - a.read() → y and a.read() → x: we cannot reorder these.
  - The rest is up for grabs.

CSE 486/586

20

## Linearizability (Textbook Definition)

- 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 actual 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
- The strongest form of consistency

CSE 486/586

21

## Summary

- Linearizability
  - Single-client, Single-copy semantics
- A read operation returns **the most recent write, regardless of the clients, according to their actual-time ordering.**

CSE 486/586

22

## Acknowledgements

- These slides contain material developed and copyrighted by Indranil Gupta (UIUC).

CSE 486/586

23