## **Inconsistency in Distributed Systems**

CS 475, Spring 2019 **Concurrent & Distributed Systems** 





# **Recurring Problem: Replication**



OK, we obviously need to actually do something here to replicate the data... but what?

Replication solves some problems, but creates a huge new one: consistency







# Sequentially Consistent DSM







# Ivy Architecture

### Each node keeps a cached copy of each piece of data it reads



### cached data

### cached data





### cached data

GMU CS 475 Spring 2019



4

# IVY VS HW4

- Ivy never copies the actual values until a replica reads them (unlike HW4) Invalidate messages are probably smaller than the actual data! Ivy only sends update (invalidate) messages to replicas who have a copy of
- the data (unlike HW4)
  - Maybe most data is not actively shared
- Ivy requires the lock server to keep track of a few more bits of information (which replica has which data)
- With near certainty Ivy is a lot faster :)



- get?
- Relaxed consistency models
- Reminders:
  - HW3 graded by end of week
  - HW4 is out!



### • Consistency in distributed systems - can we have it all? If not, what can we





# Sequential Consistency



# Availability

۲ will be available!



Our protocol for sequential consistency does NOT guarantee that the system





## **Consistent + Available**





# Still broken...





# **Network Partitions**

- The communication links between nodes may fail arbitrarily But other nodes might still be able to reach that node  $\bullet$



11

# CAP Theorem

- Pick two of three:
  - Consistency: All nodes see the same data at the same time (sequential consistency)
  - Availability: Individual node failures do not prevent survivors from continuing to operate
  - Partition tolerance: The system continues to operate despite message loss (from network and/or node failure)
- You can not have all three, ever





# CAP Theorem vs FLP

- FLP: Can not guarantee both liveness and agreement assuming messages may be delayed but are eventually delivered
- CAP: Can not guarantee consistency, availability, partition-tolerance assuming messages may be dropped
- Nice comparison: <u>http://the-paper-trail.org/blog/flp-and-cap-arent-the-same-thing/</u>



# **CAP** Theorem

- C+A: Provide strong consistency and availability, assuming there are no network partitions
- C+P: Provide strong consistency in the presence of network partitions; minority partition is unavailable
- A+P: Provide availability even in presence of partitions; no sequential consistency guarantee, maybe can guarantee something else



# Still broken...





# **Relaxing Consistency**

- We can relax two design principles:
  - How stale reads can be
  - The ordering of writes across the replicas



<b>P1</b>	W(X) 0		R(X)
<b>P2</b>		W(X) 1	R(X)
<b>P3</b>			R(X)

## Allowing Stale Reads



# Allowing Stale Reads

class MyObj { int x = 0;int y = 0;void thread0() x = 1;if(y==0)System.out.println("OK"); } void thread1() { y = 1;if(x==0)System.out.println("OK"); 7

Java's memory model is "relaxed" in that you can have stale reads

# ntln(<mark>"OK"</mark>);





.....



# **Relaxing Consistency**

 $\bullet$ partition failure

<b>P1</b>	W(X) 0		R(X) [(
<b>P2</b>		W(X) 1	R(X) [(
<b>P3</b>			R(X) [(

Intuition: less constraints means less coordination overhead, less prone to

R(X) [0,1] R(X) [0,1] 0,1] W(X) 0 R(X) [0,1] 0,1] 0,1] R(X) [0,1] R(X) [0,1]



- Assume each machine has a complete copy of memory
- Reads from local memory

```
class Machine1 {
 DSMInt x = 0;
 DSMInt y = 0;
 static void main(String[] args)
  x = 1;
  if(y==0)
      System.out.println("OK");
```

Writes broadcast update to other machines, then immediately continue

```
class Machine2 {
 DSMInt x = 0;
 DSMInt y = 0;
 static void main(String[] args)
  y = 1;
  if(x==0)
      System.out.println("OK");
}
```



- Assume each machine has a complete copy of memory
- Reads from local memory

```
class Machine1 {
 DSMInt x = (1)
 DSMInt y = 0;
 static void main(String[] args)
  ★ = 1;
  if(y==0)
      System.out.println("OK");
}
```

Writes broadcast update to other machines, then immediately continue

```
class Machine2 {
 DSMInt x = 0:
 DSMInt y = (1)
 static void main(String[] args)
    = 1;
  if(x==0)
      System.out.println("OK");
}
```



- Assume each machine has a complete copy of memory
- Reads from local memory

```
Is this
class Machine1 {
DSMInt x = (1)
 DSMInt y = 0;
 static void main(String[] args)
  x = 1;
  if(y==0)
      System.out.println("OK");
```

Writes broadcast update to other machines, then immediately continue



- It definitely is not sequentially consistent
- Are there any guarantees that it provides though?
  - Reads can be stale  $\bullet$
  - Writes can be re-ordered
  - Not really. •
- Can we come up with something more clever though with SOME guarantee? (Not as is, but with some modifications maybe it's...)  $\bullet$



# Causal Consistency

- An execution is **causally-consistent** if all **causally-related** read/write operations are executed in an order that reflects their causality
- Reads are fresh ONLY for writes that they are dependent on
- Causally-related writes appear in order, but not in order to others
- Concurrent writes can be seen in different orders by different machines
  - Compare to sequential consistency: every machine must see the same order of operations!



# Causal Consistency

<b>P1</b>	W(X)a		
<b>P2</b>		R(X)a	W(X)b
<b>P3</b>		R(X)a	
<b>P4</b>		R(X)a	

**Causally Consistent**. W(X) b and W(X) c are not related, hence could have happened one either order. W(X)a and W(X)B ARE causally related and must occur in this order

W(X)c

R(X)c	R(X)b	
R(X)b	R(x)c	



# Causal Consistency

<b>P1</b>	W(X)a		
<b>P2</b>		R(X)a	W(X)b
<b>P3</b>			
<b>P4</b>			



R(x)b R(x)aR(x)a R(x)b

**NOT Causally Consistent**. X couldn't have been b after it was a

R(x)b R(x)aR(x)a R(x)b

**Causally Consistent**. X can be a or b concurrently



- It is clearly weaker than sequential consistency
  - (Note that anything that is sequentially consistent is also causally  $\bullet$ consistent)
- Many more operations for concurrency  $\bullet$ 
  - Parallel (non-dependent) operations can occur in parallel in different places
    - Sequential would enforce a global ordering
  - E.g. if W(X) and W(Y) occur at the same time, and without dependencies, then they can occur without any locking
- Still requires some perhaps complicated implementation each client must know what is related to what.

# Why Causal Consistency?





# **Eventual Consistency**

- written values
  - Eventually: milliseconds, seconds, minutes, hours, years...
- Writes are NOT ordered as executed  $\bullet$ 
  - Allows for conflicts. Consider: Dropbox
- Git is eventually consistent

• Allow stale reads, but ensure that reads will **eventually** reflect the previously



# **Eventual Consistency**

- More concurrency than strict, sequential or causal
  - These require highly available connections to send messages, and generate lots of chatter
- Far looser requirements on network connections
  - Partitions: OK!
  - Disconnected clients: OK!
  - Always available!
- Possibility for conflicting writes :(



Each node keeps a cached copy of each piece of data it reads

x=1

cached data

read





# Sequential vs Eventual Consistency

- Sequential: "Pessimistic" concurrency control
  - Assume that everything could cause a conflict, decide on an update order as things execute, then enforce it
- Eventual: "Optimistic" concurrency control
  - Just do everything, and if you can't resolve what something should be, sort it out later
  - Can be tough to resolve in general case



31

## **Eventual Consistency: Distributed Filesystem**



When everything can talk, it's easy to synchronize, right? Goal: Everything eventually becomes synchronized. No lost updates (don't replace new version with old)





## **Eventual Consistency: Distributed Filesystem**



When everything can talk, it's easy to synchronize, right? Goal: Everything eventually becomes synchronized. No lost updates (don't replace new version with old)





## **Eventual Consistency: Distributed Filesystem**

- Role of the sync server:
  - Resolve conflicting changes, report conflicts to user
  - Do not allow sync between clients
  - Detect if updates are sequential
  - Enforce ordering constraints





# **Detecting Conflicts**

## t=0write x = a



Do we just use timestamps?

t=1write x = b



# **Detecting Conflicts**

Do we just use timestamps?





### NO, what if clocks are out of sync? NO does not actually detect conflicts

$$t=1$$
  
vrite x = b



# **Detecting Conflicts**





### Still doesn't tell us what to do with a conflict

Solution: Track version history on clients

V=0write x = b



# **Client-Centric Consistency**

- What can we guarantee in disconnected operation?
- Monotinic-reads: any future reads will return the same or newer value (never older)
- Monotonic-writes: A processes' writes are always processed in order
- Read-you-writes
- Writes follow reads



## Eventually Consistent + Available + Partition Tolerant



# Choosing a consistency model

- Sequential consistency
  - All over it's the most intuitive
- Causal consistency
  - "Increasingly useful" but not really widely used still pay coordination cost, unclear what the performance benefits are
- Eventual consistency
  - Very popular in industry and academia
  - File synchronizers, Amazon's Bayou and more

GMU CS 475 Spring 2019



40

# Example: Facebook

- Problem: >1 billion active users
- Solutions: Thousands of servers across the world
- What kind of consistency guarantees are reasonable? Need 100% availability!
- If I post a story on my news feed, is it OK if it doesn't immediately show up on yours?
  - Two users might not see the same data at the same time
  - Now this is "solved" anyway because there is no "sort by most recent first" option anyway





# **Example: Airline Reservations**

- GDS needs to sell as many seats as possible within given constraints
- If I have 100 seats for sale on a flight, does it matter if reservations for flights are reconciled immediately?
- If I have 5 seats for sale on a flight, does it matter if reservations are reconciled immediately?

 Reservations and flight inventory are managed by a GDS (Global Distribution) System), who acts as a middle broker between airlines, ticket agencies and consumers [Except for Southwest and Air New Zealand and other oddballs]



42

# **Example: Airline Reservations**

- Result: Reservations can be made using either a strong consistency model or a weak, eventual one
- Most reservations are made under the normal strong model (reservation is confirmed immediately)
- GDS also supports "Long Sell" issue a reservation without confirmed availability, need to eventually reconcile it
- Long sells require the seller to make clear to the customer that even though there's a confirmation number it's not confirmed!





# Filesystem consistency

- What consistency guarantees do a filesystem provide?
- read, write, sync, close
- On sync, guarantee writes are persisted to disk
- Readers see most recent
- What does a network file system do?



# Network Filesystem Consistency

- How do you maintain these same semantics?
- (Cheat answer): Very, very expensive
  - EVERY write needs to propagate out
  - EVERY read needs to make sure it sees the most recent write
  - Oof. Just like Ivy.



# **Consistency Takeaways**

- availability
- Weaker consistency also has a tradeoff (weaker consistency)
- But: applications can make these design choices clear to end-users  $\bullet$ 
  - Facebook
  - Dropbox
- Next week: examples of two systems that involve replication and handle consistency differently: DNS, NFS

• Strong consistency (sequential or strict) comes at a tradeoff: performance,



## This work is licensed under a Creative Commons Attribution-ShareAlike license

- This work is licensed under the Creative Commons Attribution-ShareAlike 4.0 International
- You are free to:
  - Share copy and redistribute the material in any medium or format
  - Adapt remix, transform, and build upon the material
  - for any purpose, even commercially.
- Under the following terms:
  - suggests the licensor endorses you or your use.
  - contributions under the same license as the original.
  - legally restrict others from doing anything the license permits.

License. To view a copy of this license, visit <u>http://creativecommons.org/licenses/by-sa/4.0/</u>

• Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that

• ShareAlike — If you remix, transform, or build upon the material, you must distribute your

No additional restrictions — You may not apply legal terms or technological measures that

