

# Spanner

## Storage insights



COS 518: *Advanced Computer Systems*  
Lecture 6

Michael Freedman

## 2PL & OCC = strict serialization

- Provides semantics as if only one transaction was running on DB at time, in serial order
- + Real-time guarantees
- 2PL: Pessimistically get all the locks first
- OCC: Optimistically create copies, but then recheck all read + written items before commit

2

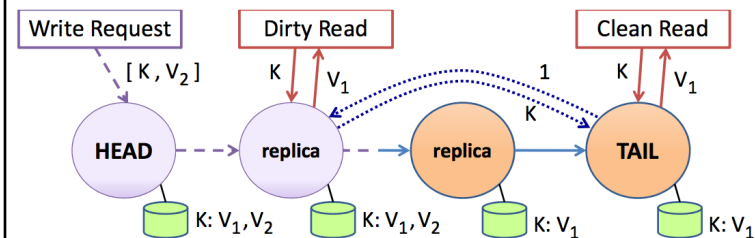
## Multi-version concurrency control

Generalize use of multiple versions of objects

3

## Multi-version concurrency control

- Maintain multiple versions of objects, each with own timestamp. Allocate correct version to reads.
- Prior example of MVCC:



4

## Multi-version concurrency control

---

- Maintain multiple versions of objects, each with own timestamp. Allocate correct version to reads.
- Unlike 2PL/OCC, reads never rejected
- Occasionally run garbage collection to clean up

5

## MVCC Intuition

---

- Split transaction into read set and write set
  - All reads execute as if one “snapshot”
  - All writes execute as if one later “snapshot”
- Yields snapshot isolation < serializability

6

## Serializability vs. Snapshot isolation

---

- Intuition: Bag of marbles:  $\frac{1}{2}$  white,  $\frac{1}{2}$  black
- Transactions:
  - T1: Change all white marbles to black marbles
  - T2: Change all black marbles to white marbles
- Serializability (2PL, OCC)
  - T1  $\rightarrow$  T2 or T2  $\rightarrow$  T1
  - In either case, bag is either ALL white or ALL black
- Snapshot isolation (MVCC)
  - T1  $\rightarrow$  T2 or T2  $\rightarrow$  T1 or T1 || T2
  - Bag is ALL white, ALL black, or  $\frac{1}{2}$  white  $\frac{1}{2}$  black

7

## Timestamps in MVCC

---

- Transactions are assigned timestamps, which may get assigned to objects those txns read/write
- Every object version  $O_v$  has both read and write TS
  - ReadTS: Largest timestamp of txn that reads  $O_v$
  - WriteTS: Timestamp of txn that wrote  $O_v$

8

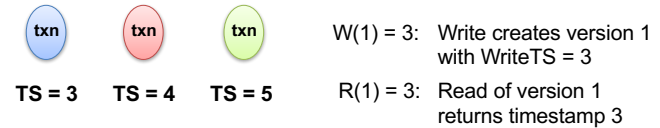
## Executing transaction T in MVCC

- Find version of object O to read:
  - # Determine the last version written before read snapshot time
  - Find  $O_v$  s.t.  $\max \{ \text{WriteTS}(O_v) \mid \text{WriteTS}(O_v) \leq \text{TS}(T) \}$
  - $\text{ReadTS}(O_v) = \max(\text{TS}(T), \text{ReadTS}(O_v))$
  - Return  $O_v$  to T
- Perform write of object O or abort if conflicting:
  - Find  $O_v$  s.t.  $\max \{ \text{WriteTS}(O_v) \mid \text{WriteTS}(O_v) \leq \text{TS}(T) \}$
  - # Abort if another T' exists and has read O after T
  - If  $\text{ReadTS}(O_v) > \text{TS}(T)$ 
    - Abort and roll-back T
  - Else
    - Create new version  $O_w$
    - Set  $\text{ReadTS}(O_w) = \text{WriteTS}(O_w) = \text{TS}(T)$

9

## Digging deeper

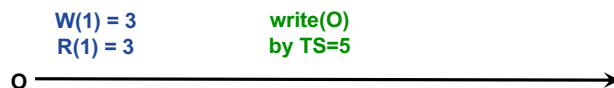
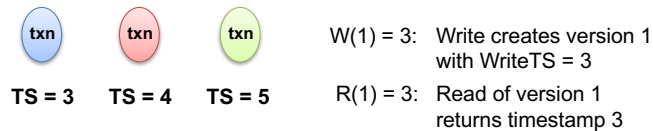
### Notation



10

## Digging deeper

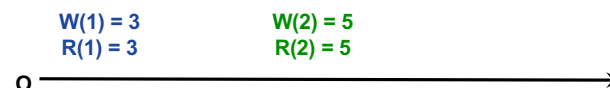
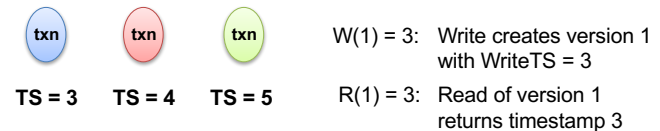
### Notation



11

## Digging deeper

### Notation

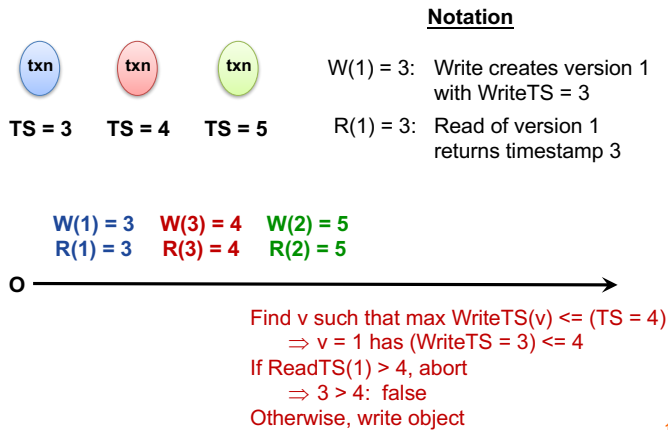


$\text{write}(O)$   
 by  $\text{TS} = 4$

Find  $v$  such that  $\max \text{WriteTS}(v) \leq (\text{TS} = 4)$   
 $\Rightarrow v = 1$  has  $(\text{WriteTS} = 3) \leq 4$   
 If  $\text{ReadTS}(1) > 4$ , abort  
 $\Rightarrow 3 > 4$ : false  
 Otherwise, write object

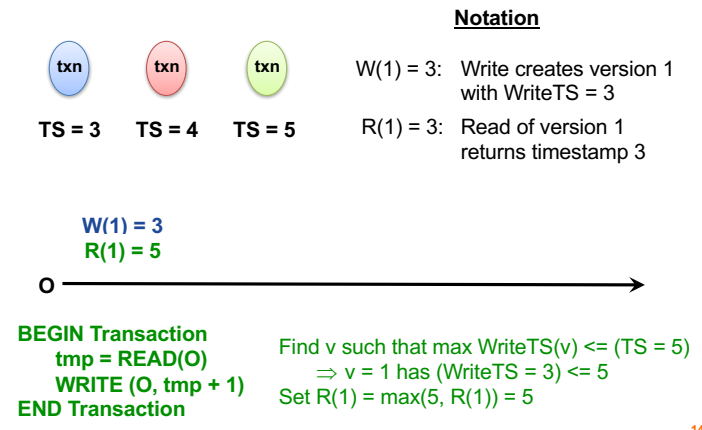
12

## Digging deeper



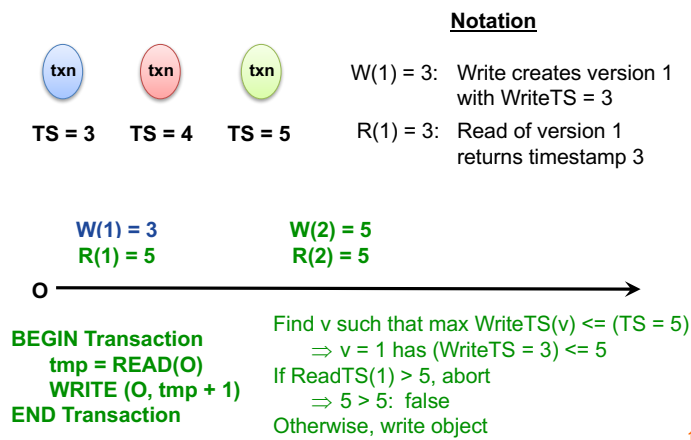
13

## Digging deeper



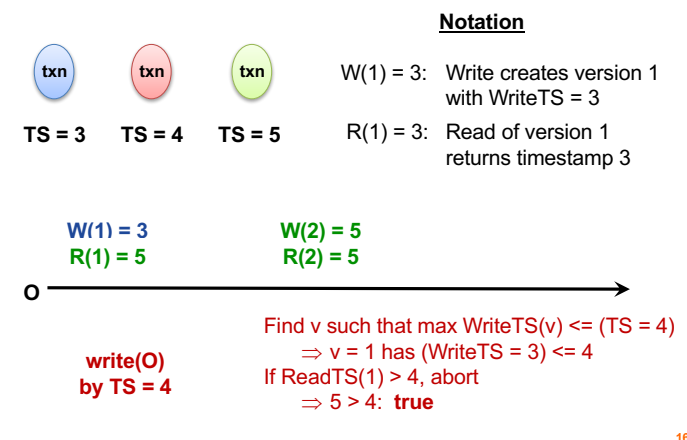
14

## Digging deeper



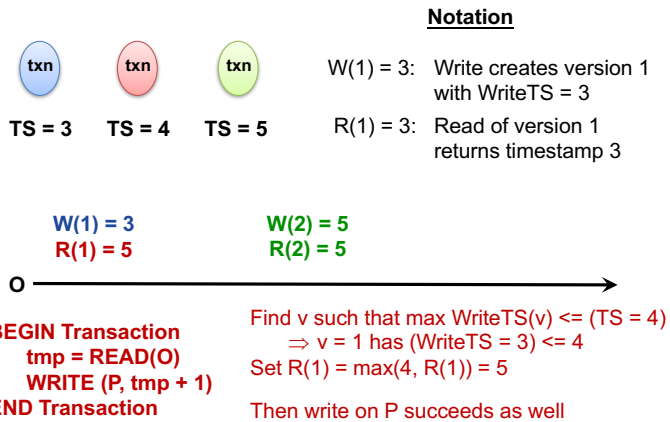
15

## Digging deeper



16

## Digging deeper

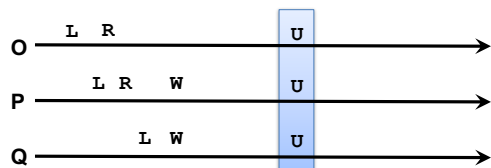


17

## Distributed Transactions

18

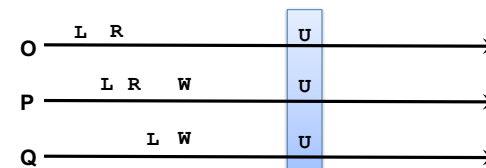
## Consider partitioned data over servers



- Why not just use 2PL?
  - Grab locks over entire read and write set
  - Perform writes
  - Release locks (at commit time)

19

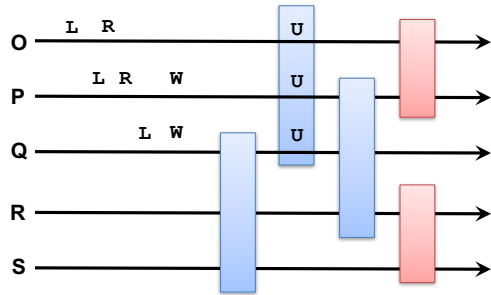
## Consider partitioned data over servers



- How do you get serializability?
  - On single machine, single COMMIT op in the WAL
  - In distributed setting, assign global timestamp to txn (at sometime after lock acquisition and before commit)
    - Centralized txn manager
    - Distributed consensus on timestamp (not all ops)

20

## Strawman: Consensus per txn group?



- Single Lamport clock, consensus per group?
  - Linearizability composes!
  - But doesn't solve concurrent, non-overlapping txn problem

21

## Spanner: Google's Globally-Distributed Database

OSDI 2012

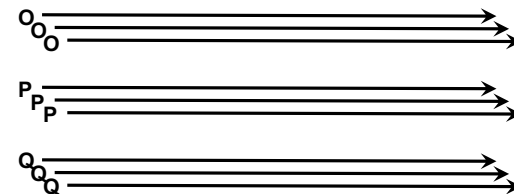
22

## Google's Setting

- Dozens of zones (datacenters)
- Per zone, 100-1000s of servers
- Per server, 100-1000 partitions (tablets)
- Every tablet replicated for fault-tolerance (e.g., 5x)

23

## Scale-out vs. fault tolerance



- Every tablet replicated via Paxos (with leader election)
- So every "operation" within transactions across tablets actually a replicated operation within Paxos RSM
- Paxos groups can stretch across datacenters!
  - (COPS took same approach *within* datacenter)

24

## Disruptive idea:

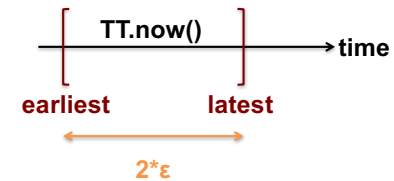
Do clocks **really** need to be arbitrarily unsynchronized?

Can you engineer some max divergence?

25

## TrueTime

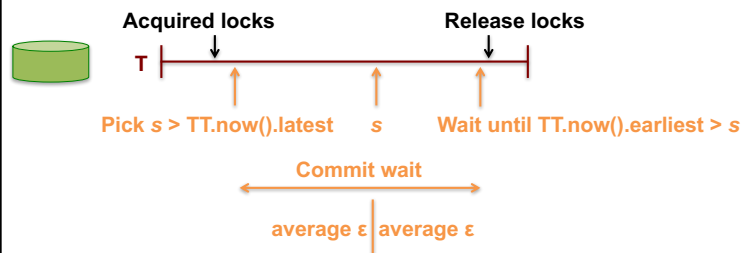
- “Global wall-clock time” with bounded uncertainty



Consider event  $e_{\text{now}}$  which invoked  $tt = \text{TT.now}()$ :  
Guarantee:  $tt.\text{earliest} \leq t_{\text{abs}}(e_{\text{now}}) \leq tt.\text{latest}$

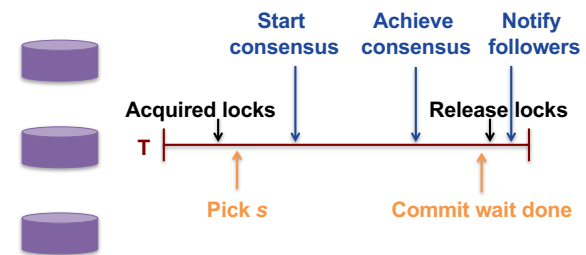
26

## Timestamps and TrueTime



27

## Commit Wait and Replication



28

## Client-driven transactions

Client:

1. Issues reads to leader of each tablet group, which acquires read locks and returns most recent data
2. Locally performs writes
3. Chooses coordinator from set of leaders, initiates commit
4. Sends commit message to each leader, include identify of coordinator and buffered writes
5. Waits for commit from coordinator

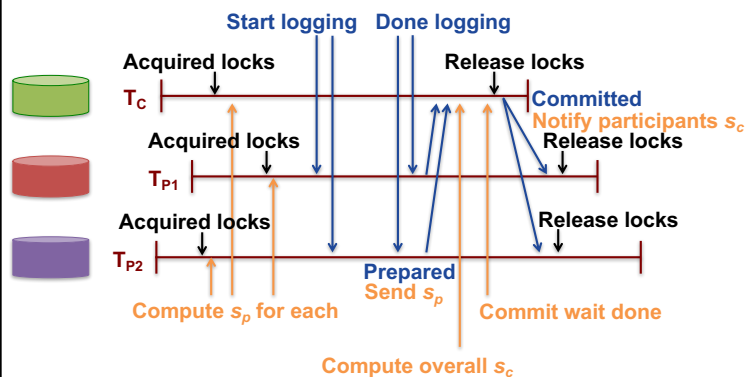
29

## Commit Wait and 2-Phase Commit

- On commit msg from client, leaders acquire local write locks
  - If non-coordinator:
    - Choose prepare ts > previous local timestamps
    - Log prepare record through Paxos
    - Notify coordinator of prepare timestamp
  - If coordinator:
    - Wait until hear from other participants
    - Choose commit timestamp  $\geq$  prepare ts, > local ts
    - Logs commit record through Paxos
    - Wait commit-wait period
    - Sends commit timestamp to replicas, other leaders, client
- All apply at commit timestamp and release locks

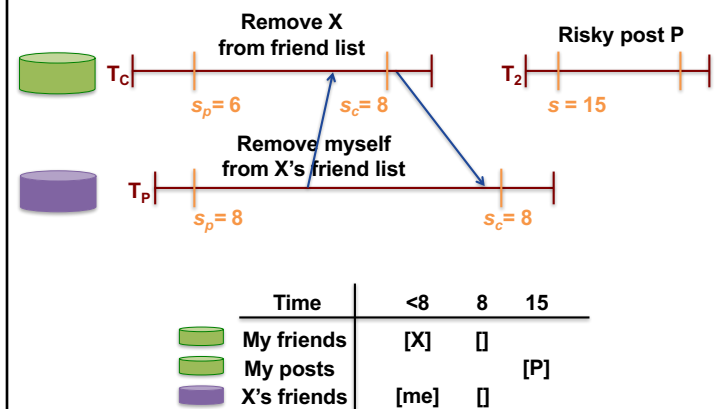
30

## Commit Wait and 2-Phase Commit



31

## Example



32



## Read-only optimizations

- Given global timestamp, can implement read-only transactions lock-free (snapshot isolation)
- Step 1: Choose timestamp  $s_{read} = TT.now.latest()$
- Step 2: Snapshot read (at  $s_{read}$ ) to each tablet
  - Can be served by any up-to-date replica

33

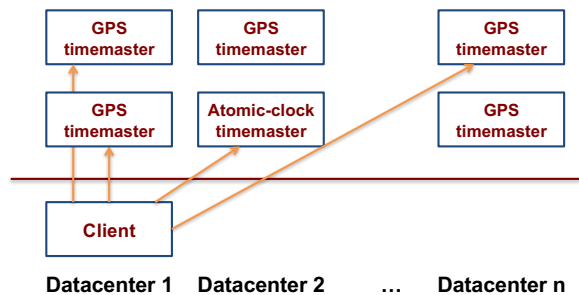
## Disruptive idea:

Do clocks **really** need to be arbitrarily unsynchronized?

**Can you engineer some max divergence?**

34

## TrueTime Architecture



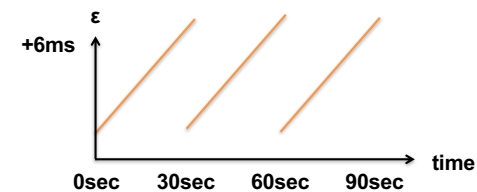
Compute reference [earliest, latest] =  $now \pm \epsilon$

35

## TrueTime implementation

$now = reference\ now + local-clock\ offset$

$\epsilon = reference\ \epsilon + worst-case\ local-clock\ drift$   
 $= 1ms + 200\ \mu s/sec$



- What about faulty clocks?
  - Bad CPUs 6x more likely in 1 year of empirical data

36

Known unknowns > unknown unknowns

Rethink algorithms to reason about uncertainty

37

The case for log storage:  
Hardware tech affecting software design

38

## Latency Numbers Every Programmer Should Know

latency.txt June 7, 2012

1	Latency Comparison Numbers			
2	-----			
3	L1 cache reference	0.5 ns		
4	Branch mispredict	5 ns		
5	L2 cache reference	7 ns		
6	Mutex lock/unlock	25 ns		
7	Main memory reference	100 ns		
8	Compress 1K bytes with Zip	3,000 ns	3 us	
9	Send 1K bytes over 1 Gbps network	10,000 ns	10 us	
10	Read 4K randomly from SSD*	150,000 ns	150 us	
11	Read 1 MB sequentially from memory	250,000 ns	250 us	
12	Round trip within same datacenter	500,000 ns	500 us	
13	Read 1 MB sequentially from SSD*	1,000,000 ns	1,000 us	1 ms
14	Disk seek	10,000,000 ns	10,000 us	10 ms
15	Read 1 MB sequentially from disk	20,000,000 ns	20,000 us	20 ms
16	Send packet CA->Netherlands->CA	150,000,000 ns	150,000 us	150 ms

From <https://gist.github.com/jboner/2841832>

See also [https://people.eecs.berkeley.edu/~rscs/research/interactive\\_latency.html](https://people.eecs.berkeley.edu/~rscs/research/interactive_latency.html)

39

~2016



**Seagate (\$50)**  
1TB HDD 7200RPM  
Model: STD1000DM003-1SB10C

Operation	HDD Performance
Sequential Read	176 MB/s
Sequential Write	190 MB/s
Random Read 4KiB	0.495 MB/s 121 IOPS
Random Write 4KiB	0.919 MB/s 224 IOPS
DQ Random Read 4KiB	1.198 MB/s 292 IOPS
DQ Random Write 4KiB	0.929 MB/s 227 IOPS

<http://www.tomshardware.com/answers/id-3201572/good-normal-read-write-speed-hdd.html>

40

~2016



**Seagate (\$50)**  
1TB HDD 7200RPM  
Model: STD1000DM003-1SB10C



**Samsung (\$330)**  
512 GB 960 Pro NVMe PCIe M.2  
Model: MZ-V6P512BW

Operation	HDD Performance	SSD Performance
Sequential Read	176 MB/s	2268 MB/s
Sequential Write	190 MB/s	1696 MB/s
Random Read 4KiB	0.495 MB/s 121 IOPS	44.9 MB/s 10,962 IOPS
Random Write 4KiB	0.919 MB/s 224 IOPS	151 MB/s 36,865 IOPS
DQ Random Read 4KiB	1.198 MB/s 292 IOPS	348 MB/s 84961 IOPS
DQ Random Write 4KiB	0.929 MB/s 227 IOPS	399 MB/s 97,412 IOPS

<http://www.tomshardware.com/answers/id-3201572/good-normal-read-write-speed-hdd.html>

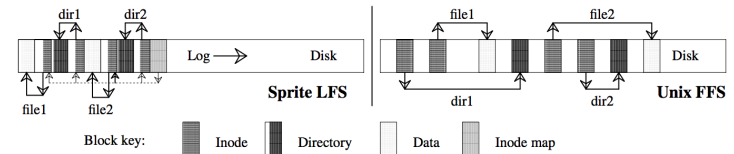
<http://ssd.userbenchmark.com/SpeedTest/182182/Samsung-SSD-960-PRO-512GB>

41

## The Design and Implementation of a Log-Structured File System

Mendel Rosenblum and John K. Ousterhout

- **Idea:** Traditionally disks laid out with spatial locality due to cost of seeks
- **Observation:** main memory getting bigger → most reads from memory
- **Implication:** Disk workloads now write-heavy → avoid seeks → write log
- **New problem:** Many seeks to read, need to occasionally defragment
- **New tech solution:** SSDs → seeks cheap, erase blocks change defrag



This paper will appear in the *Proceedings of the 13th ACM Symposium on Operating Systems Principles* and the February 1992 *ACM Transactions on Computer Systems*.

42