Spanner

(Part II)



COS 418: Distributed Systems Lecture 18

Mike Freedman, Jialin Ding

Some slides from the Spanner OSDI talk

1

Recap: Ideas Behind Read-Only Txns

- Tag writes with physical timestamps upon commit
 - Write txns are strictly serializable, e.g., 2PL
- Read-only txns return the writes, whose commit timestamps precede the reads' current time
 - Rotxns are one-round, lock-free, and never abort

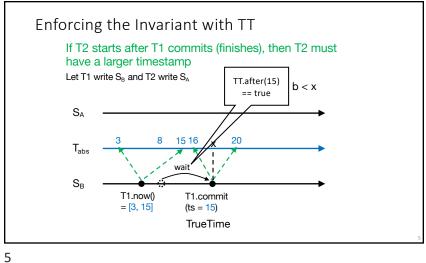
Recap: Spanner is Strictly Serializable

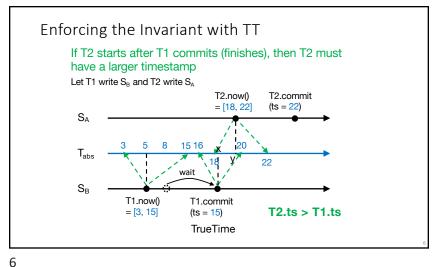
- Efficient read-only transactions in strictly serializable systems
 - Strict serializability is desirable but costly!
 - Reads are prevalent! (340x more than write txns)
 - Efficient rotxns → good system overall performance

2

Recap: TrueTime

- Timestamping writes must enforce the invariant
 - If T2 starts after T1 commits (finishes), then T2 must have a larger timestamp
- TrueTime: partially-synchronized clock abstraction
 - Bounded clock skew (uncertainty)
 - TT.now() → [earliest, latest]; earliest <= T_{abs} <= latest
 - Uncertainty (ε) is kept short
- TrueTime enforces the invariant by
 - Use at least TT.now().latest for timestamps
 - Commit wait





Enforcing the Invariant with TT • What if T1.commit delayed, such that T2 happens after T1.now() but before T1.commit? Tricky as T1.commit.ts = T1.now().latest • Answer: T2 delayed until after T1 commits. Discussed later. T2.commit T2.now() (ts = 22)= [18, 22]15 16 T1.now() T1.commit = [3, 15] (ts = 15)**T2.ts > T1.ts** TrueTime

This Lecture

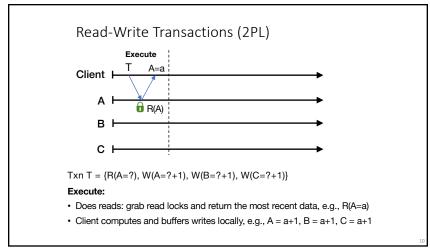
- How write transactions are done
 - 2PL + 2PC (sometimes 2PL for short)
 - How they are timestamped
- How read-only transactions are done
 - How read timestamps are chosen
 - How reads are executed

Read-Write Transactions (2PL)

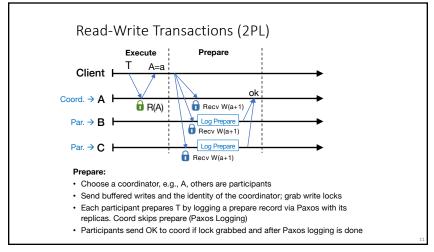
Three phases

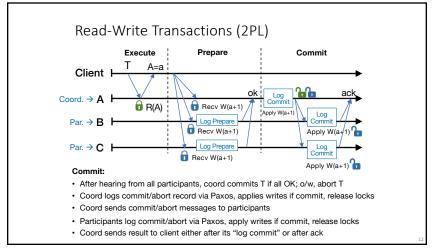


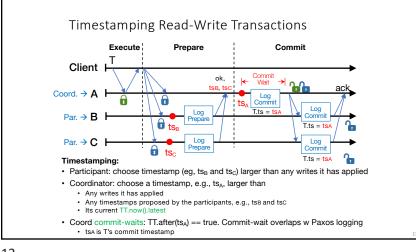
2PC: atomicity



9 10



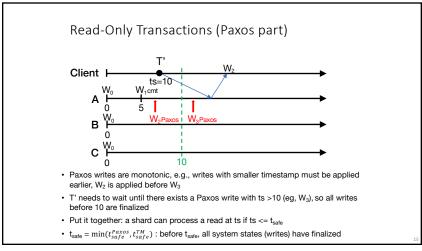


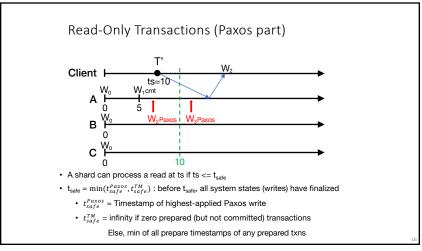


Read-Only Transactions (shards part)

Client T' W_1 W_0 W_0

13 14





Serializable Snapshot Reads

- Client specifies a read timestamp way in the past
 - · E.g., one hour ago
- Read shards at the stale timestamp
- Serializable
 - · Old timestamp cannot ensure real-time order
- Better *performance*
 - · No waiting in any cases
 - E.g., non-blocking, not just lock-free
- Can have performance but still strictly serializable?
 - E.g., one-round, non-blocking, and strictly serializable
 - Coming in next lecture!

Takeaway

- Strictly serializable (externally consistent)
 - Make it easy for developers to build apps!
- Reads dominant, make them efficient
 - One-round, lock-free
- TrueTime exposes clock uncertainty
 - Commit wait and at least TT.now.latest() for timestamps ensure real-time ordering
- Globally-distributed database
 - 2PL w/ 2PC over Paxos!

20