

Distributed Transactions in Spanner



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology

CS 240: Computing Systems and Concurrency Lecture 18

Marco Canini

Credits: Michael Freedman and Kyle Jamieson developed much of the original material.
Contents adapted from Wyatt Lloyd.

Why Google built Spanner

2005 – BigTable [OSDI 2006]

- Eventually consistent across datacenters
- Lesson: “don’t need distributed transactions”

2008? – MegaStore [CIDR 2011]

- Strongly consistent across datacenters
- Option for distributed transactions
 - Performance was not great...

2011 – Spanner [OSDI 2012]

- Strictly Serializable Distributed Transactions
- “We wanted to make it easy for developers to build their applications”

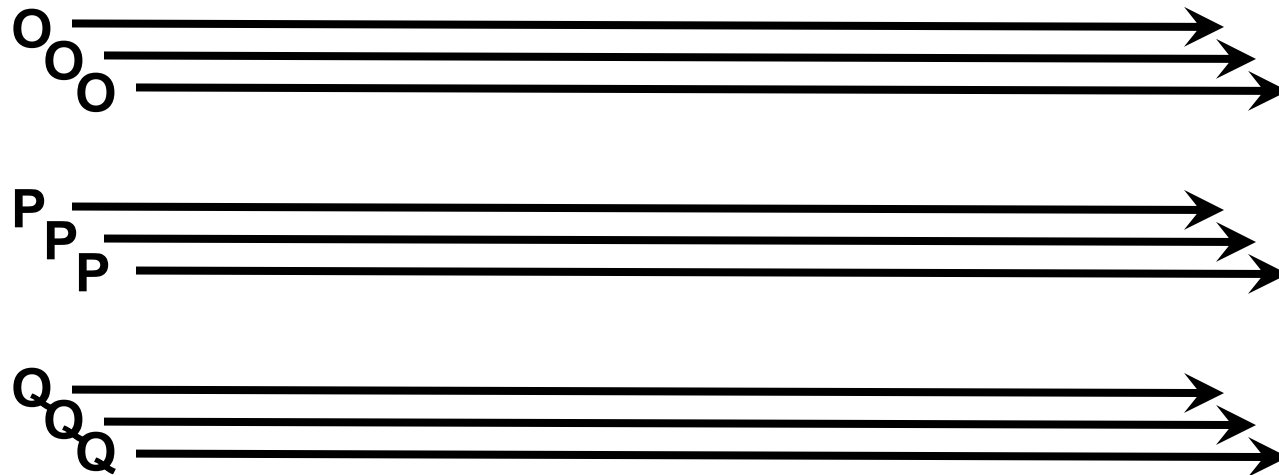
Spanner: Google's Globally-Distributed Database

OSDI 2012

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)

Scale-out vs. fault tolerance



- Every tablet replicated via Paxos (with leader election)
- So every “operation” within transactions across tablets actually is a replicated operation within Paxos RSM
- Paxos groups can stretch across datacenters!

Read-Only Transactions

- Transactions that only read data
 - Predeclared, i.e., developer uses `READ_ONLY` flag / interface
- Reads are dominant operations
 - e.g., FB's TAO had **500 reads** : 1 write [ATC 2013]
 - e.g., Google Ads (F1) on Spanner from 1? DC:
 - 21.5**B** reads in 24h
 - 31.2M single-shard transactions in 24h
 - 32.1M multi-shard transactions in 24h

Make Read-Only Txns Efficient

- Ideal: Read-only transactions that are non-blocking
 - Arrive at shard, read data, send data back
 - Impossible with Strict Serializability
- Goal 1: Lock-free read-only transactions
- Goal 2: Non-blocking stale read-only txns

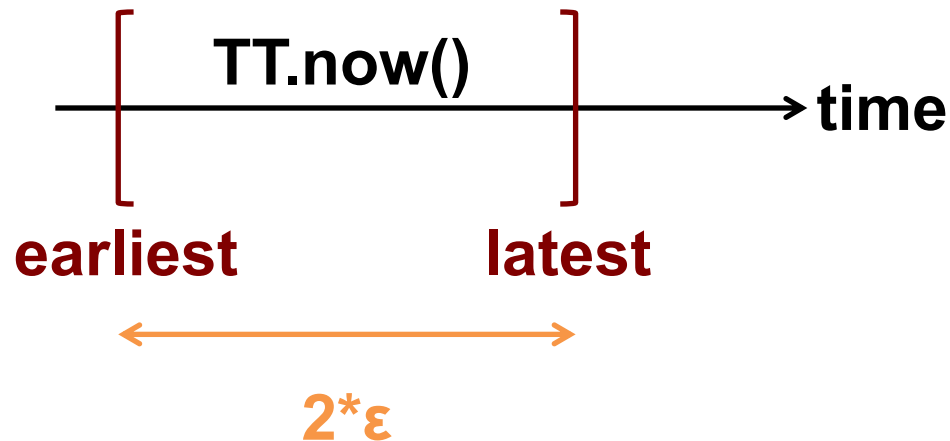
Disruptive idea:

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

Can you engineer some max divergence?

TrueTime

- “Global wall-clock time” with bounded uncertainty
 - ϵ is worst-case clock divergence
 - Timestamps become intervals, not single values



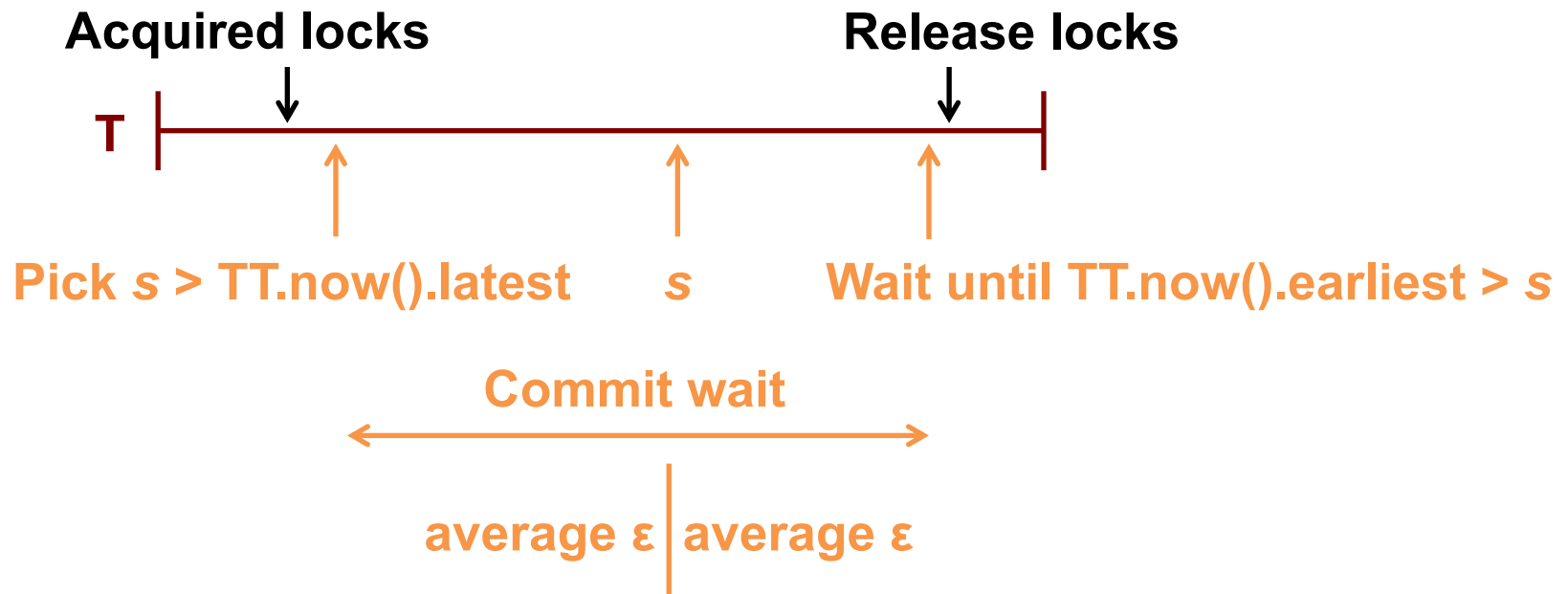
Consider event e_{now} which invoked $tt = \text{TT.now}()$:

Guarantee: $tt.\text{earliest} \leq t_{\text{abs}}(e_{\text{now}}) \leq tt.\text{latest}$

TrueTime for Read-Only Txns

- Assign all transactions a wall-clock commit time (s)
 - All replicas of all shards track how up-to-date they are with t_{safe}
 - i.e., all transactions with $s < t_{\text{safe}}$ have committed on this machine
- Goal 1: Lock-free read-only transactions
 - Current time \leq TT.now.latest()
 - $s_{\text{read}} = \text{TT.now.latest}()$
 - wait until $s_{\text{read}} < t_{\text{safe}}$
 - Read data as of s_{read}
- Goal 2: Non-blocking stale read-only txns
 - Similar to above, except explicitly choose time in the past
 - (Trades away consistency for better perf, e.g., lower latency)

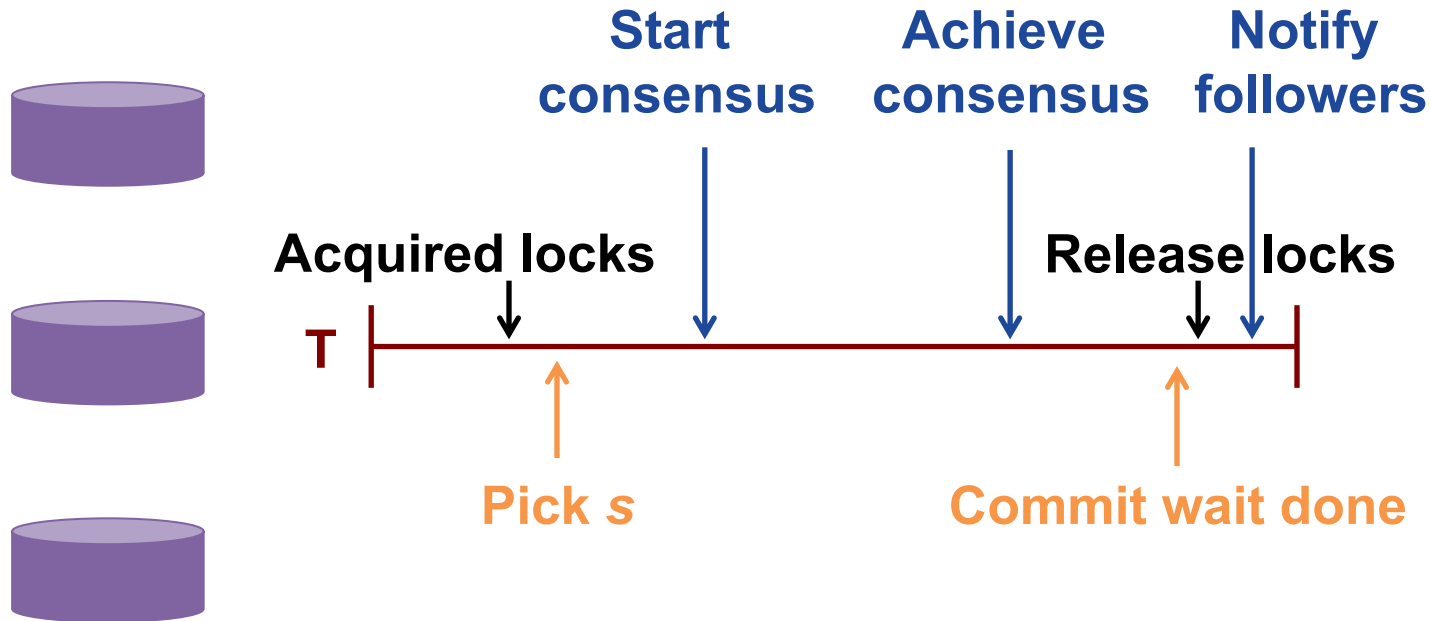
Timestamps and TrueTime



Commit Wait

- Enables efficient read-only transactions
- Cost: 2ε extra latency
- Reduce/eliminate by overlapping with:
 - Replication
 - Two-phase commit

Commit Wait and Replication



Sufficient for single-shard transactions!

Client-driven transactions for multi-shard transactions

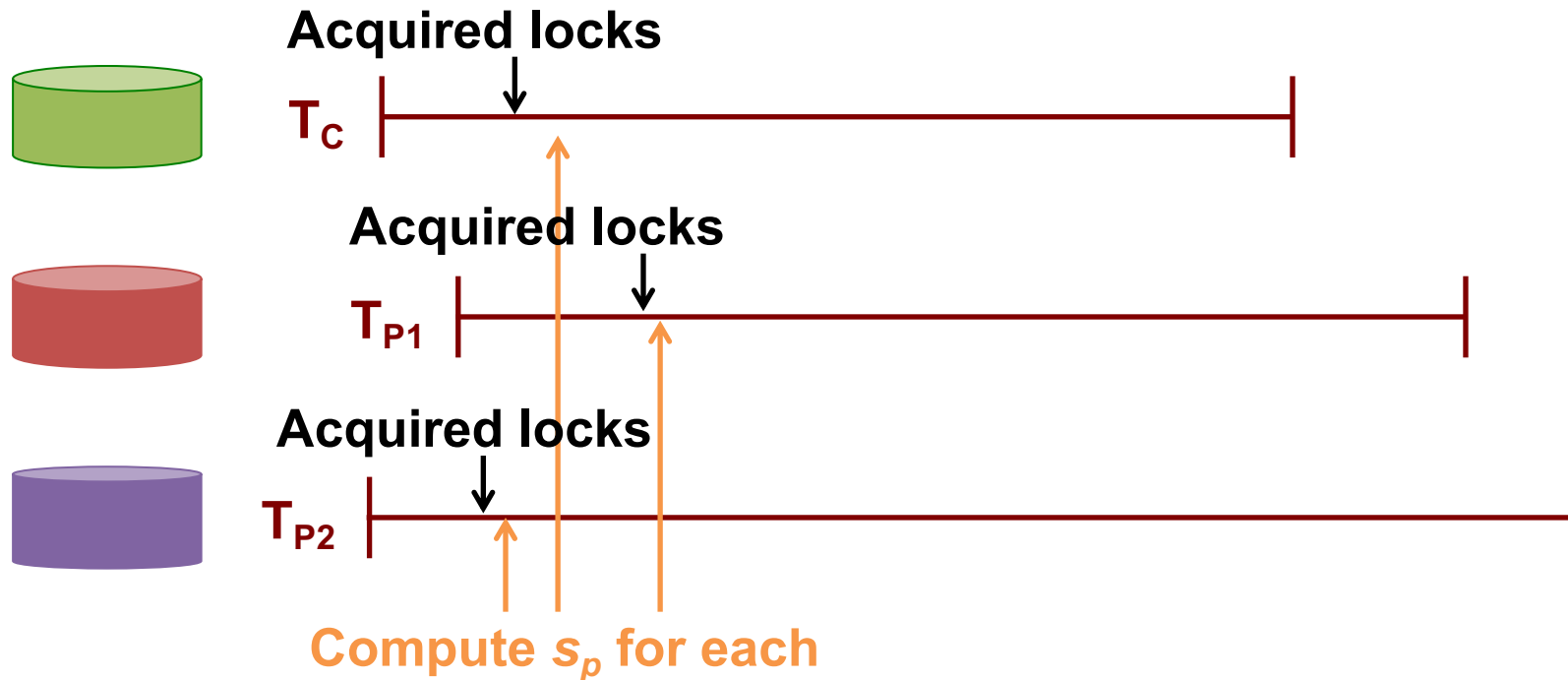
Client: 2PL w/ 2PC

1. Issues reads to leader of each shard 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 identity of coordinator and buffered writes
5. Waits for commit from coordinator

Commit Wait and 2-Phase Commit

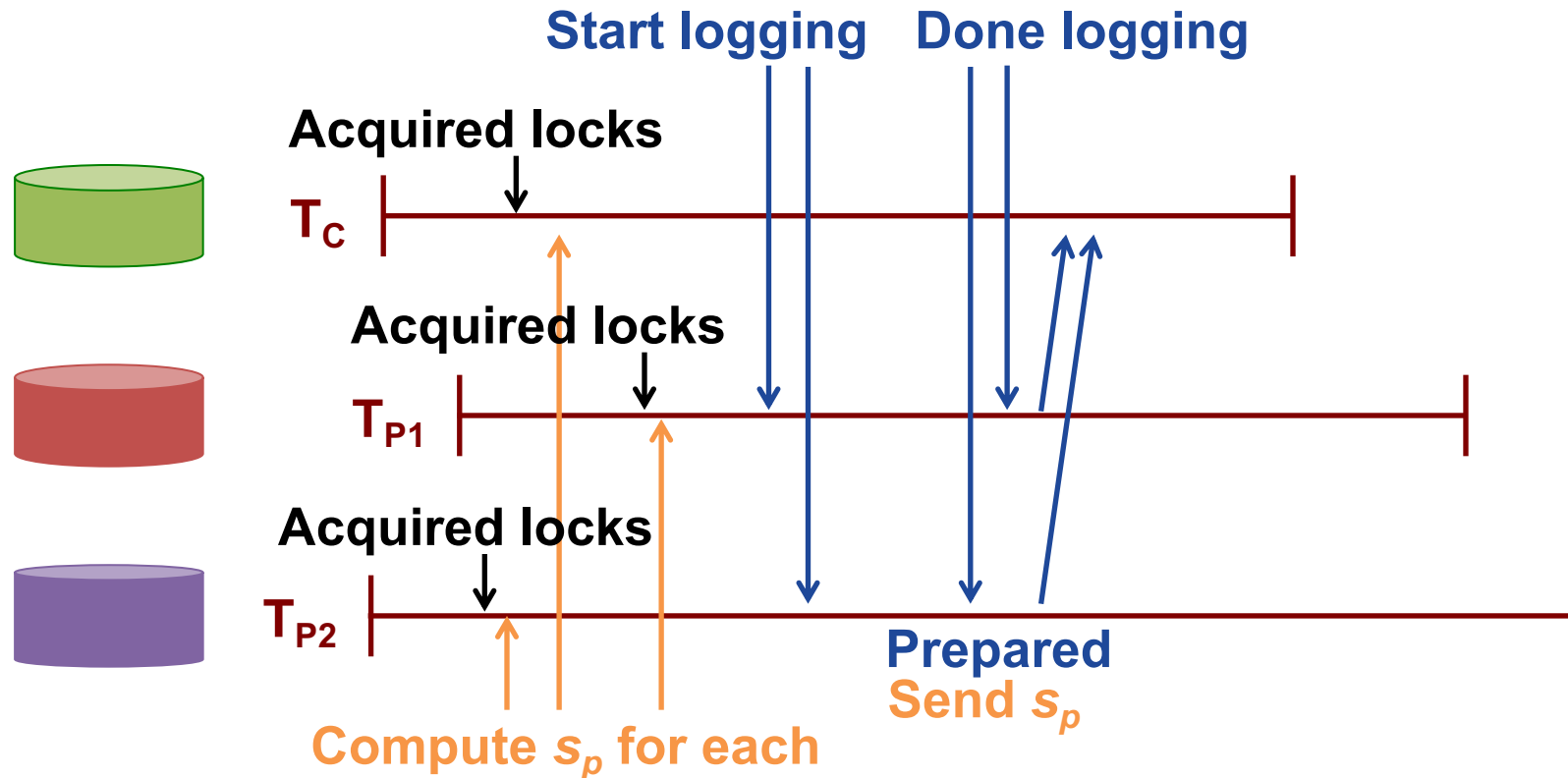
- On commit msg from client, leaders acquire local write locks
 - If non-coordinator:
 - Choose prepare ts $>$ previous local timestamps
 - Log prepare record through Paxos
 - Notify coordinator of prepare timestamp
 - If coordinator:
 - Wait until hear from other participants
 - Choose commit timestamp \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

Commit Wait and 2-Phase Commit



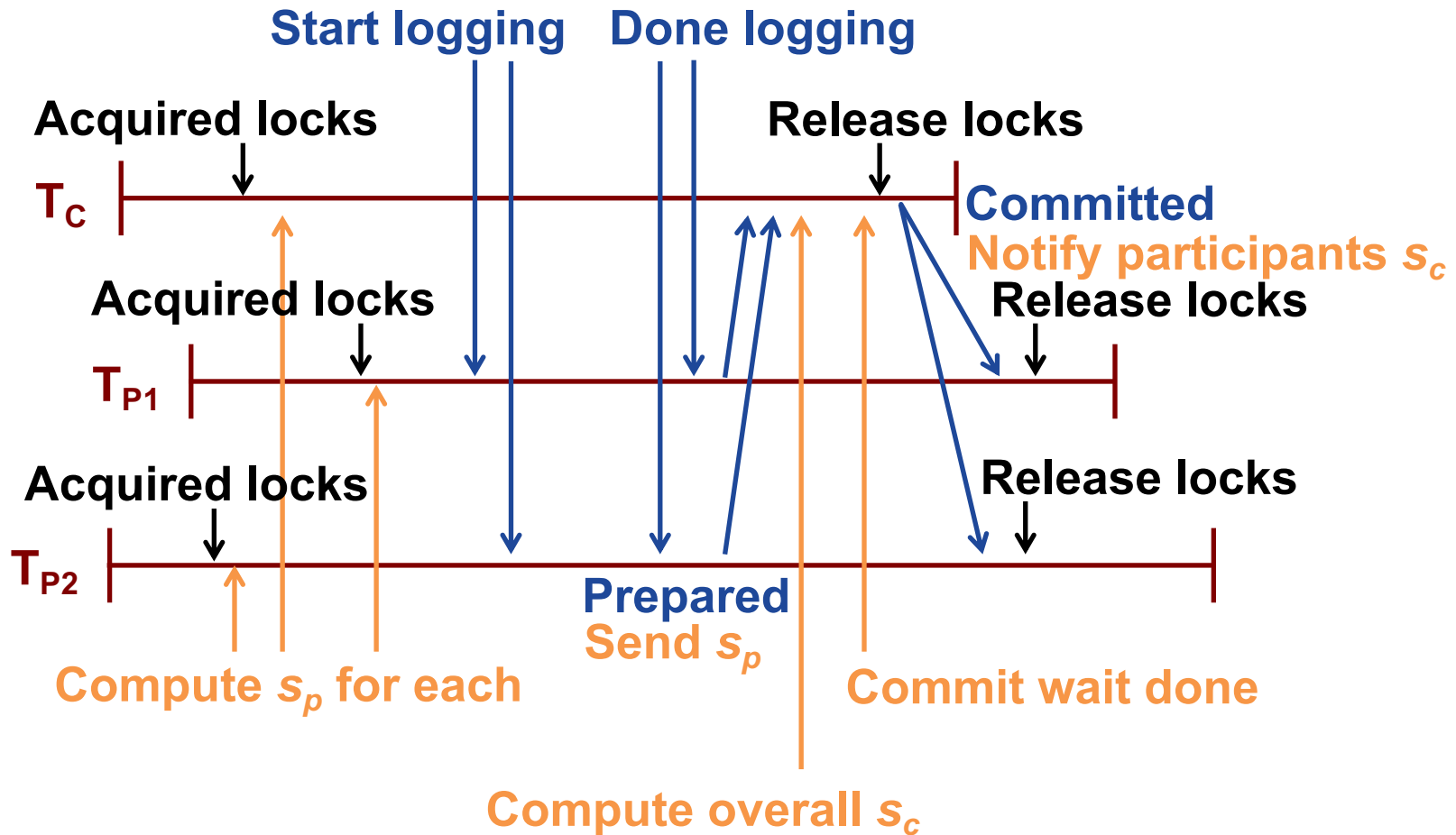
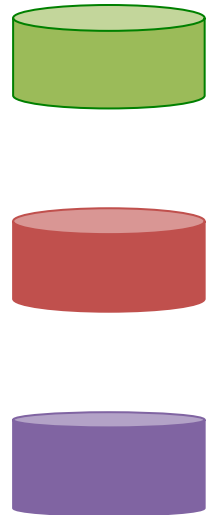
1. Client issues reads to leader of each shard group, which acquires read locks and returns most recent data

Commit Wait and 2-Phase Commit

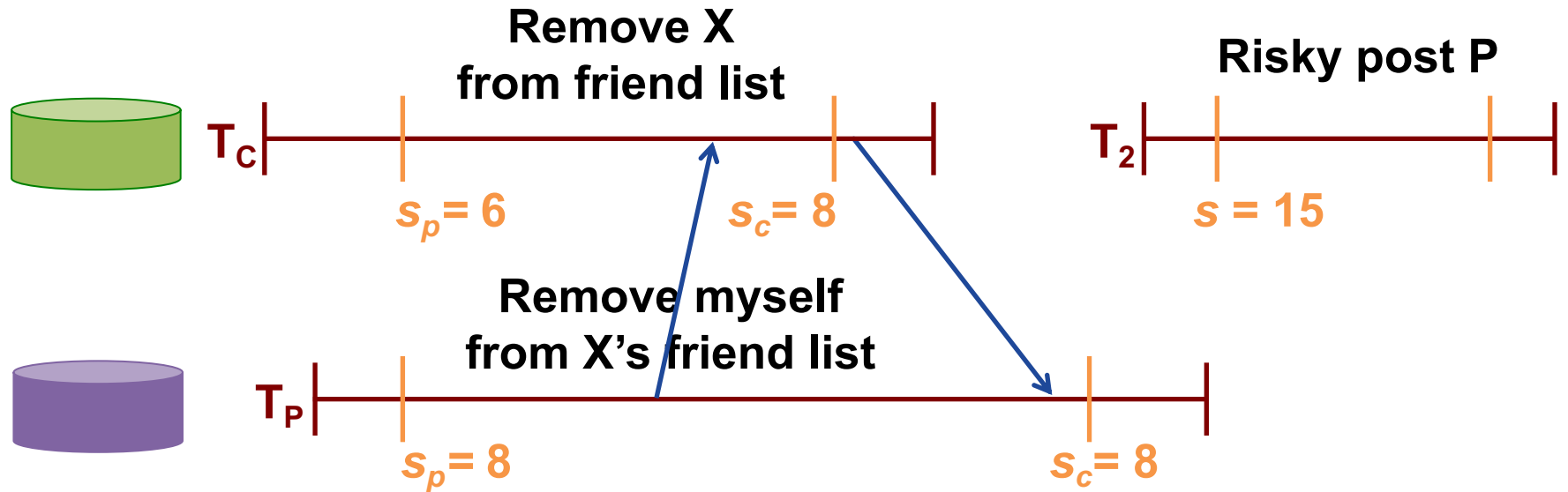





2. Locally performs writes
3. Chooses coordinator from set of leaders, initiates commit
4. Sends commit msg to each leader, incl. identity of coordinator

Commit Wait and 2-Phase Commit



Example



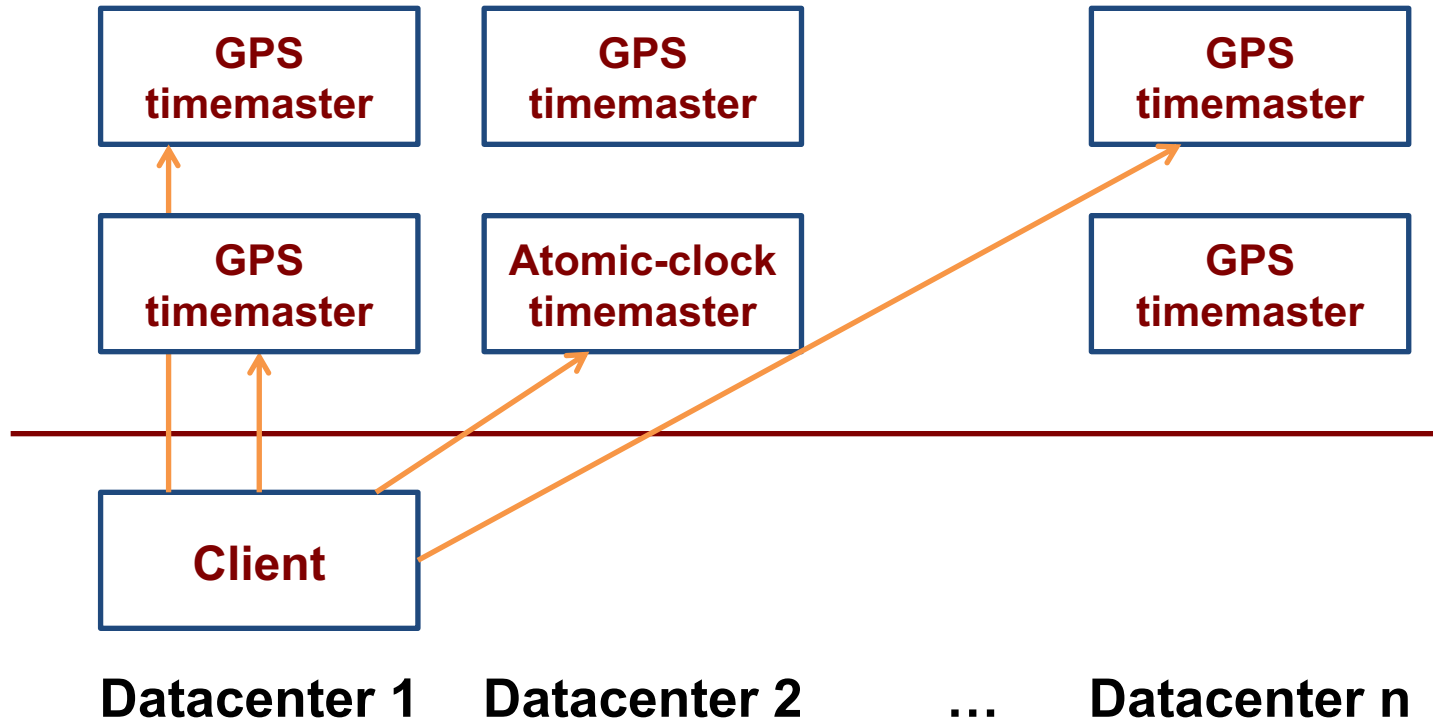
	Time	<8	8	15
 My friends		[X]	[]	
 My posts				[P]
 X's friends		[me]	[]	

Disruptive idea:

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

Can you engineer some max divergence?

TrueTime Architecture

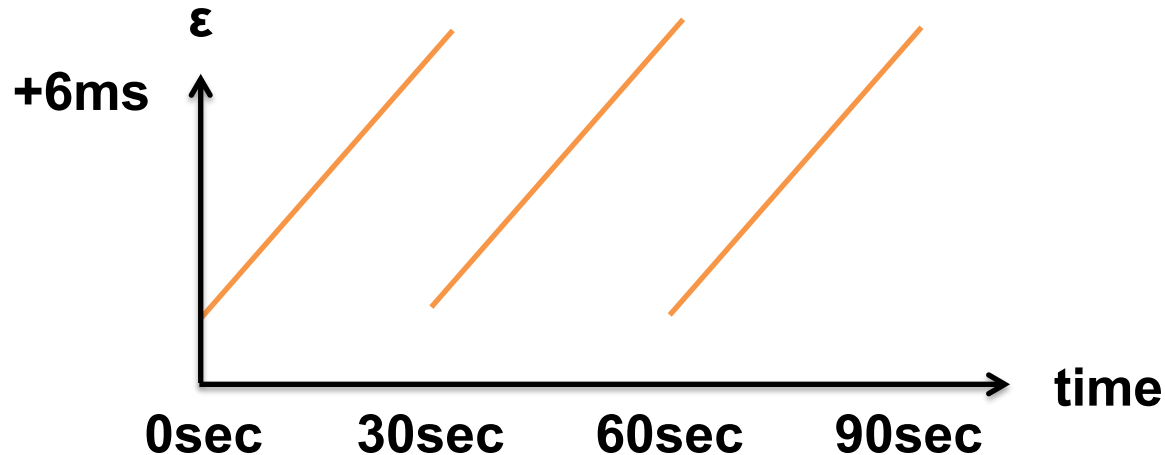


Compute reference [earliest, latest] = now $\pm \epsilon$

TrueTime implementation

now = reference now + local-clock offset

ϵ = reference ϵ + worst-case local-clock drift
= 1ms + 200 μ s/sec



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

Spanner

- Make it easy for developers to build apps!
- Reads dominant, make them lock-free
- TrueTime exposes clock uncertainty
 - Commit wait ensures transactions end after their commit time
 - Read at `TT.now.latest()`
- Globally-distributed database
 - 2PL w/ 2PC over Paxos!

Known unknowns > unknown unknowns

**Rethink algorithms to reason about
uncertainty**