# Spanner: Google's Globally-Distributed Database

## Google, Inc.
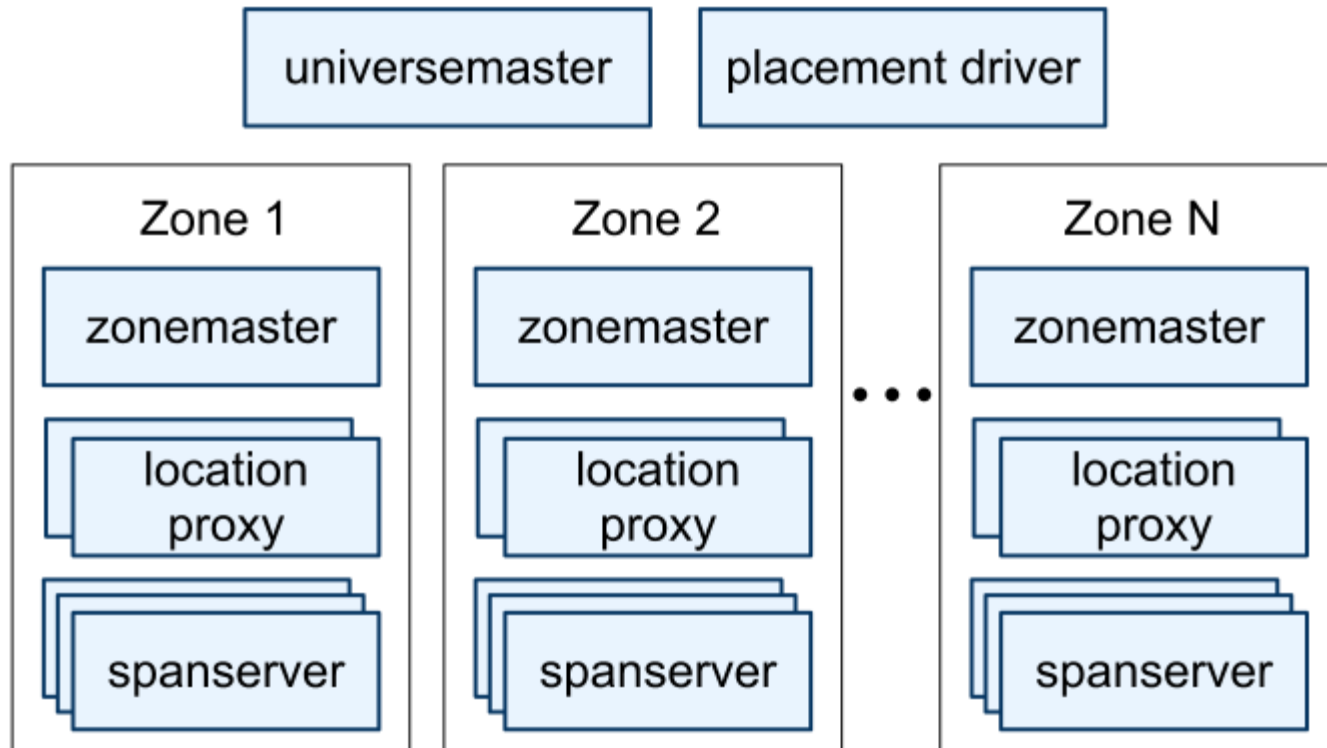
## OSDI 2012

Presented by: Karen Ouyang

# Problem Statement

- Distributed data system with high availability

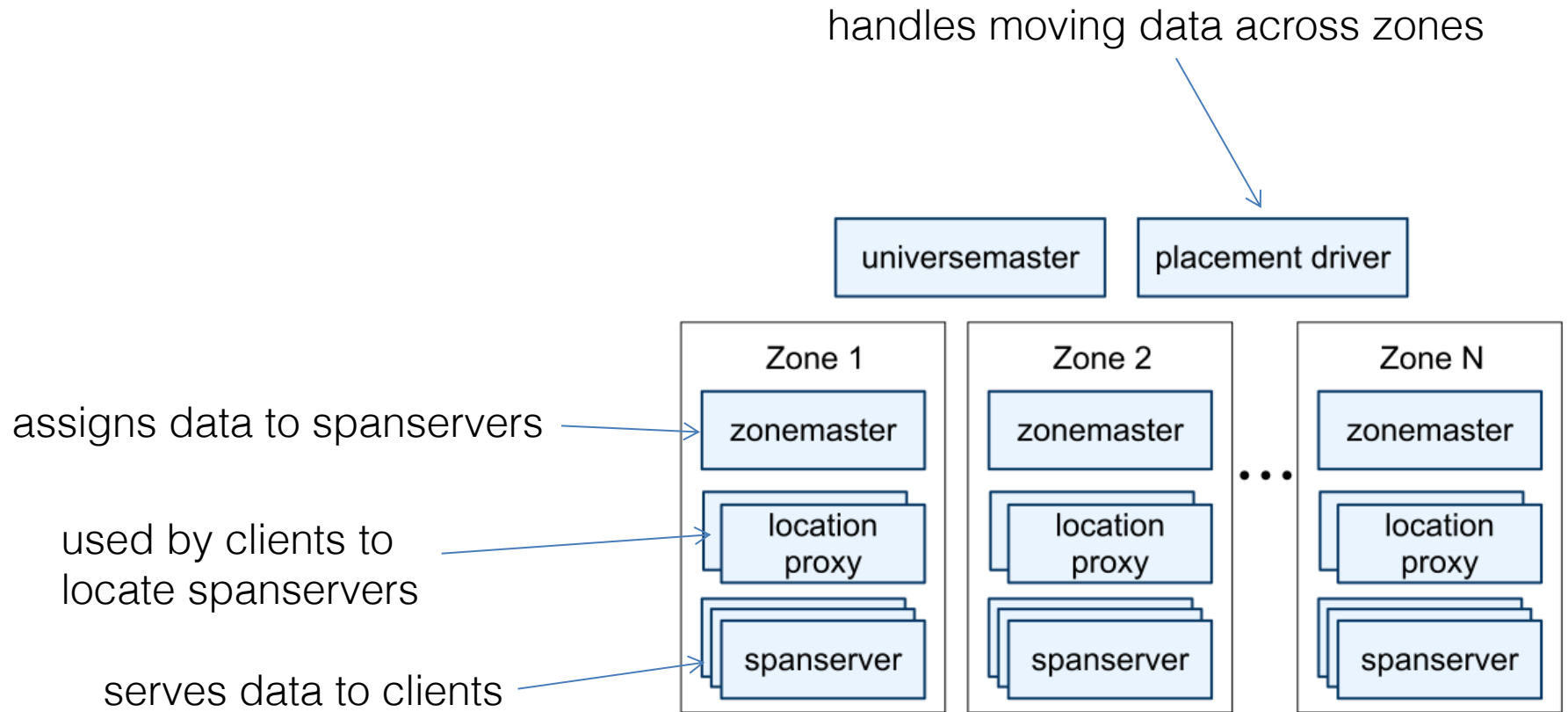- Support external consistency!

# Key Ideas

- Distributed data system with high availability

- Supports external consistency!
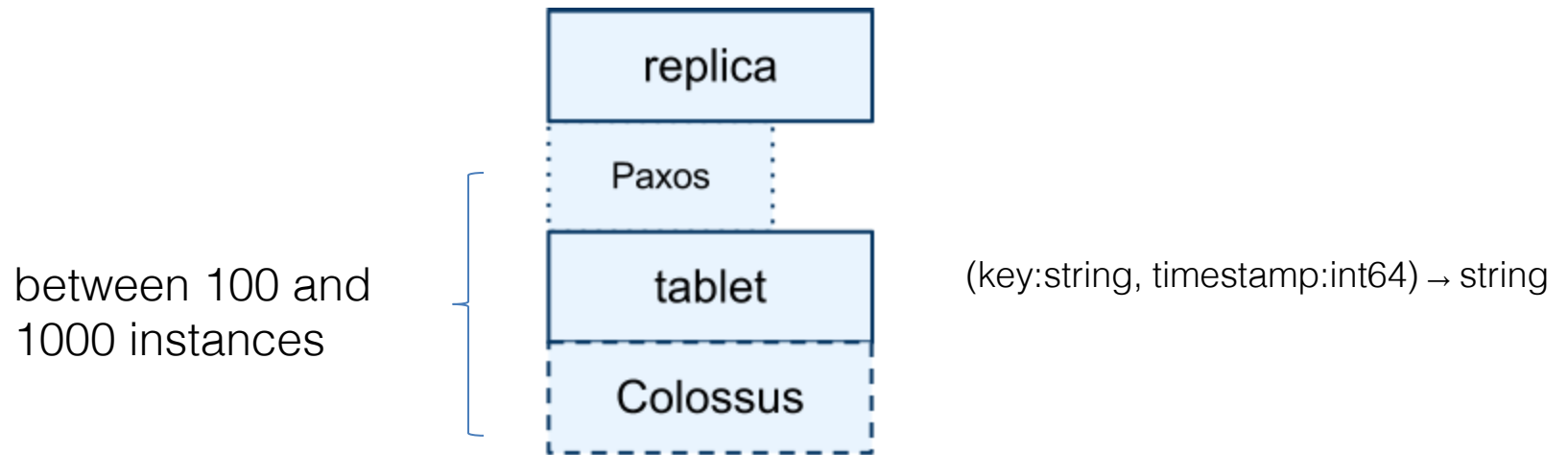
- Enabling technology: TrueTime API

# Server Organization



datacenters have one or more zones
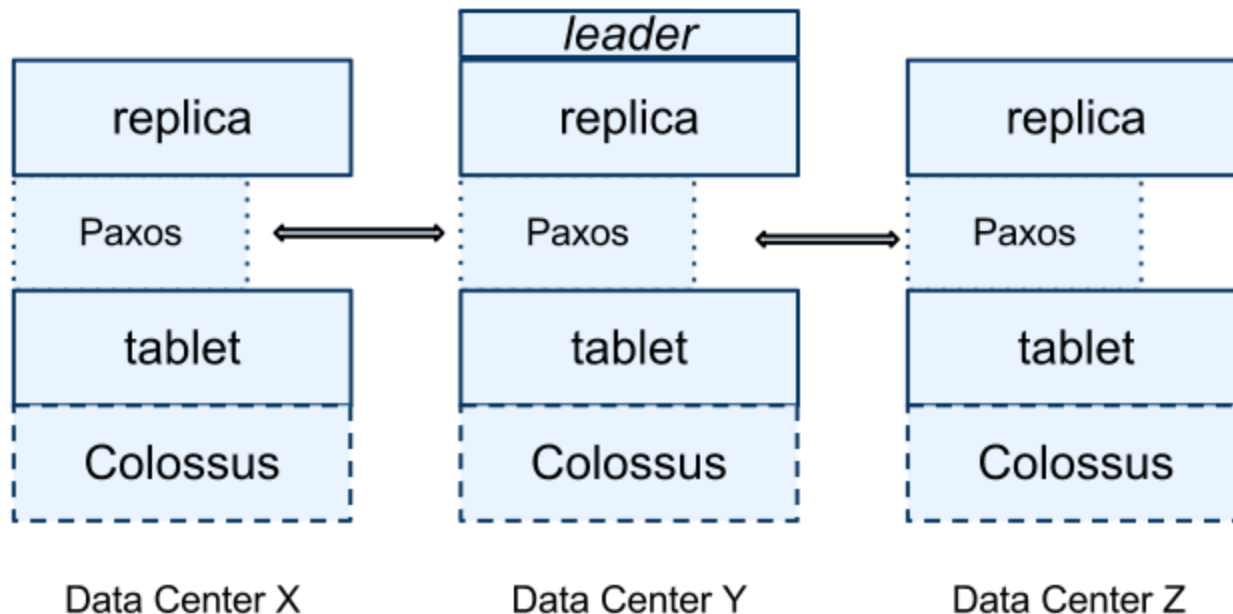
# Server Organization

handles moving data across zones

universemaster    placement driver

Zone 1
zonemaster
location proxy
spanserver

Zone 2
zonemaster
location proxy
spanserver

Zone N
zonemaster
location proxy
spanserver

assigns data to spanservers

used by clients to locate spanservers

serves data to clients

# Spanserver Stack



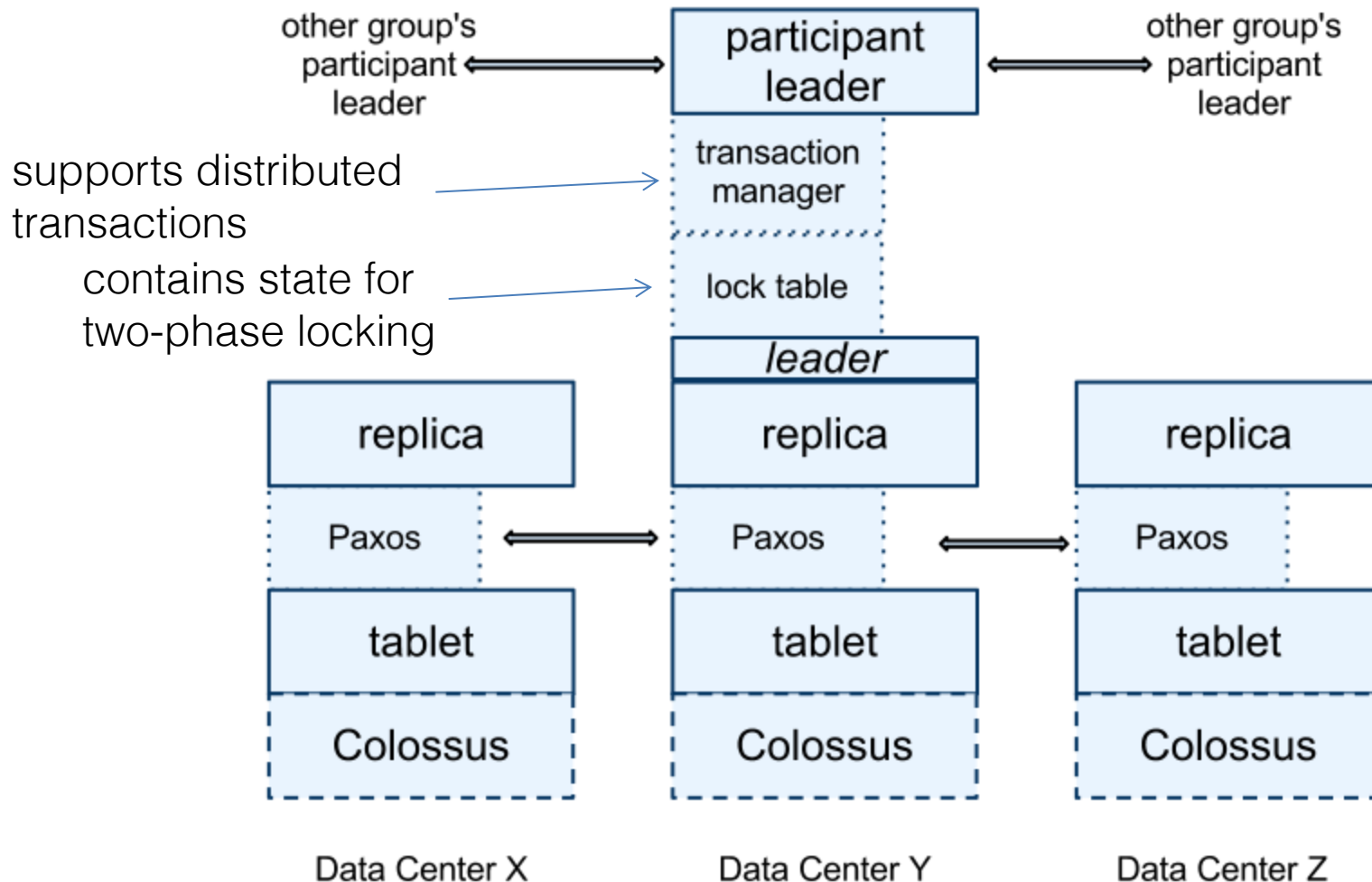between 100 and 1000 instances

(key:string, timestamp:int64) → string

# Spanserver Stack

set of replicas: *Paxos group*

writes initiate Paxos protocol at leader;
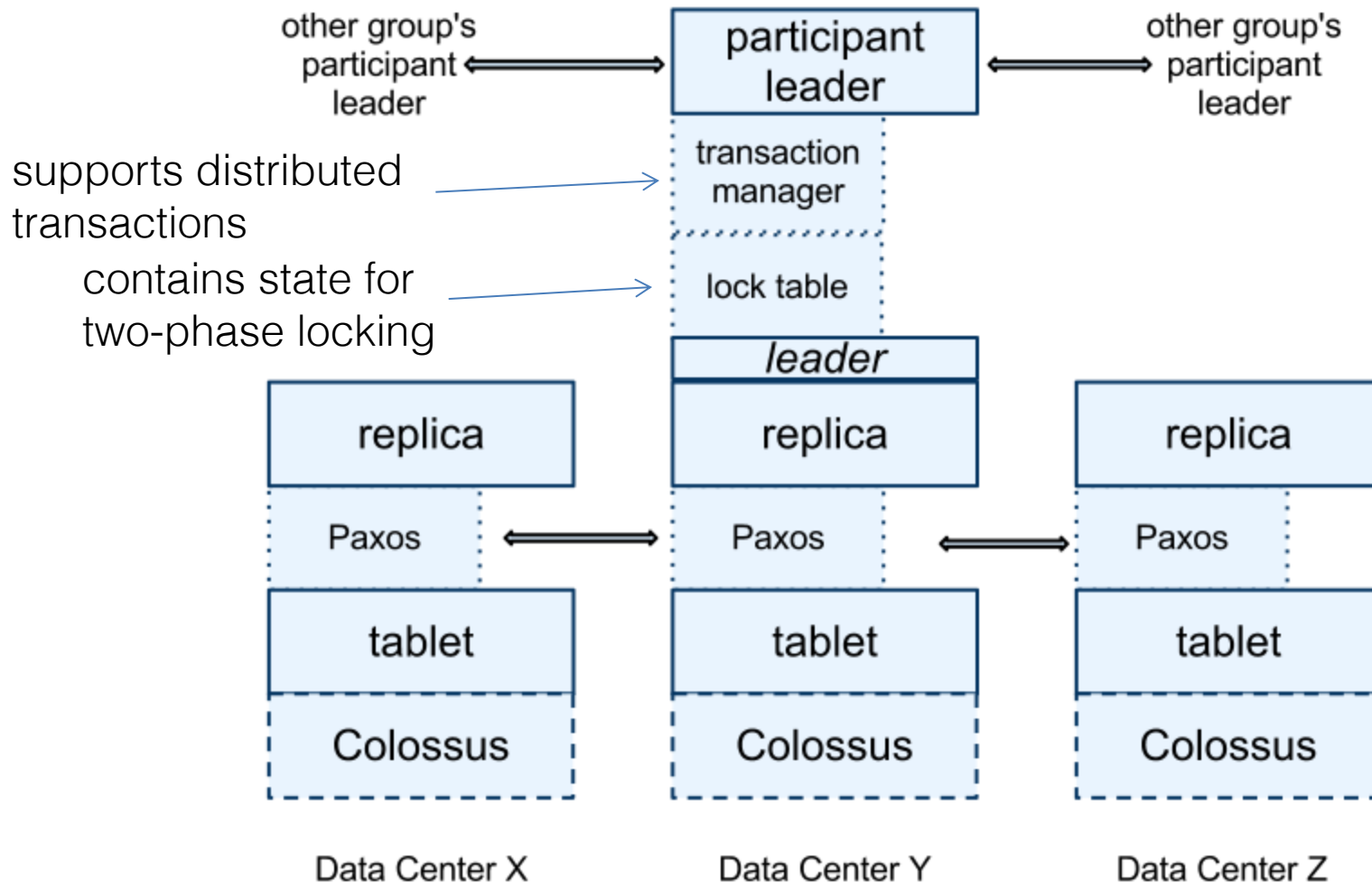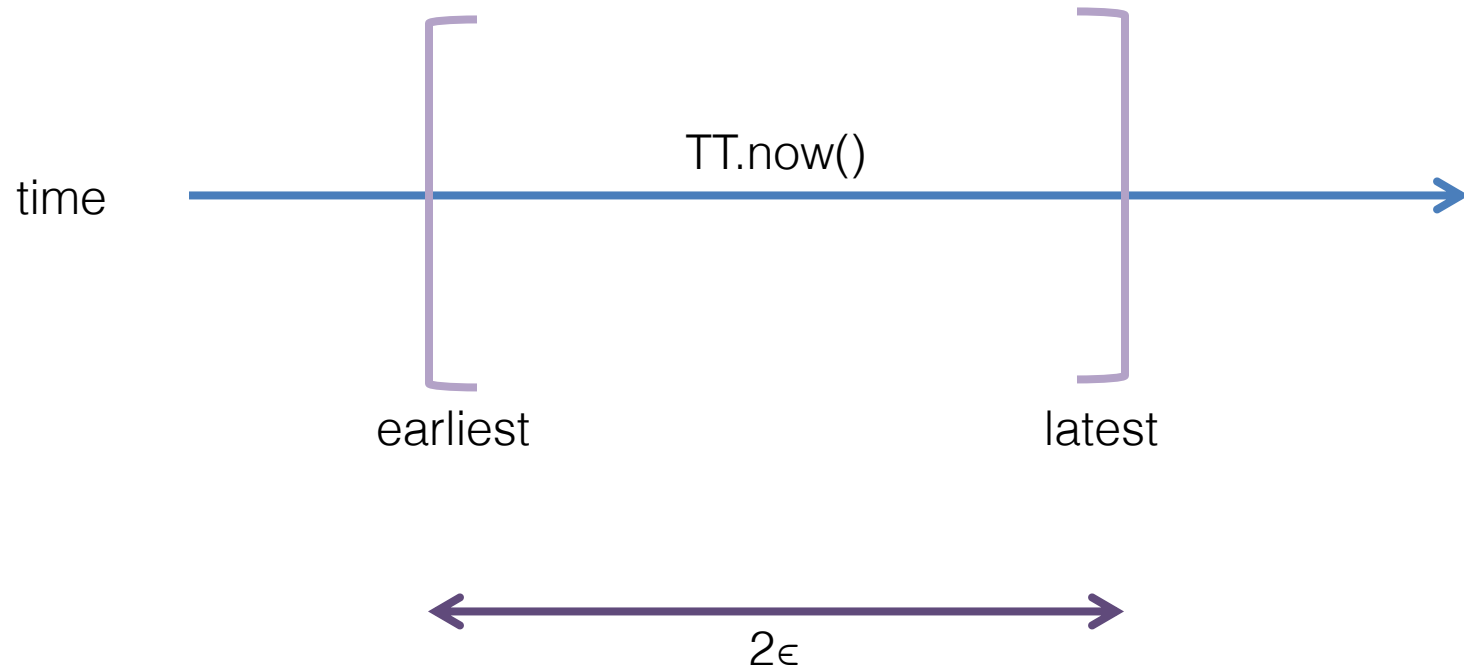reads from any sufficiently up-to-date replica



Data Center X      Data Center Y      Data Center Z

# Spanserver Stack

transactions with 1+ group: two-phase commit
select *coordinator leader* from participant leaders

other group's
participant
leader

participant
leader

other group's
participant
leader

supports distributed
transactions

transaction
manager

contains state for
two-phase locking

lock table

*leader*

| replica | replica | replica |
|---------|---------|---------|
| Paxos | Paxos | Paxos |
| tablet | tablet | tablet |
| Colossus | Colossus | Colossus |

Data Center X          Data Center Y          Data Center Z

# TrueTime API

- Exposes clock uncertainty by expressing time as an interval
- Uses GPS and atomic clocks
- *Time master* machines per datacenter
- Client polls multiple masters to compute time interval

# TrueTime API



time

TT.now()

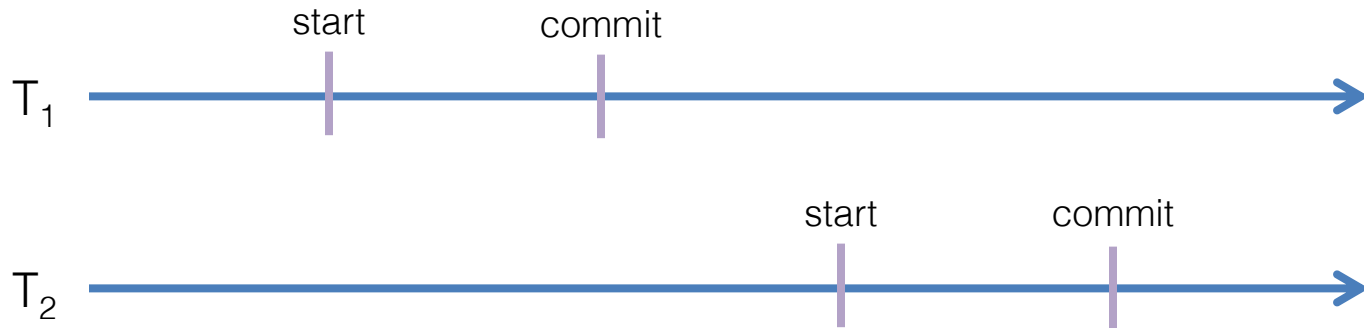earliest                    latest

$2\epsilon$

# Consistency

- Ensure external consistency by ensuring timestamp order

- All transactions are assigned timestamp

- Data written by $T$ is timestamped with $s$

# Read-Write Transactions

- Two-phase locking: assign timestamps at any time that locks are held

- Assign timestamps to Paxos writes in increasing order across leaders
  - A leader only assigns timestamps within its leader lease; leader leases are disjoint

# Read-Write Transactions

- Transactions: two-phase commit
- Two transactions



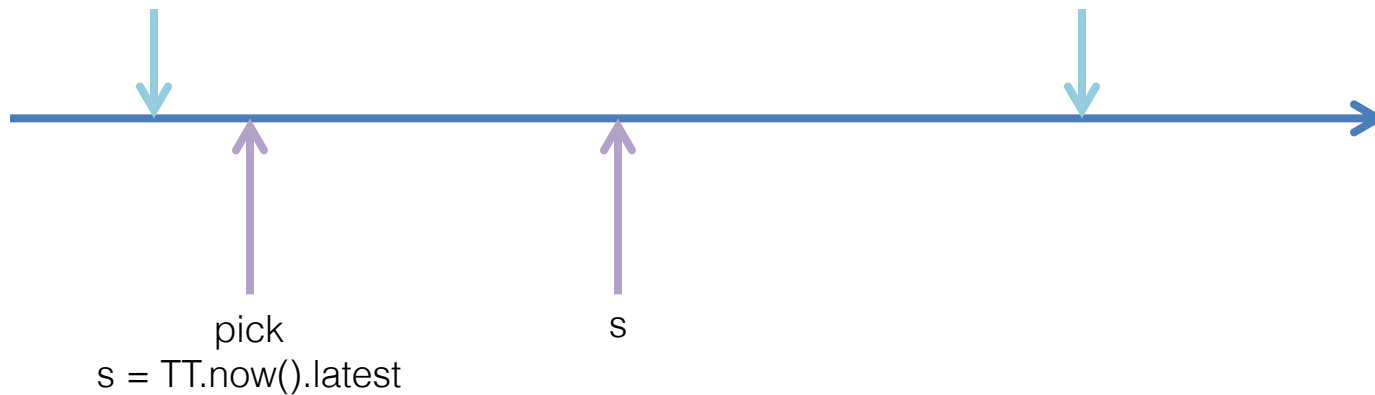- Assign commit timestamps with $s_1 < s_2$
- How?

# Read-Write Transactions

Start: commit timestamp is *after* time of commit request at server

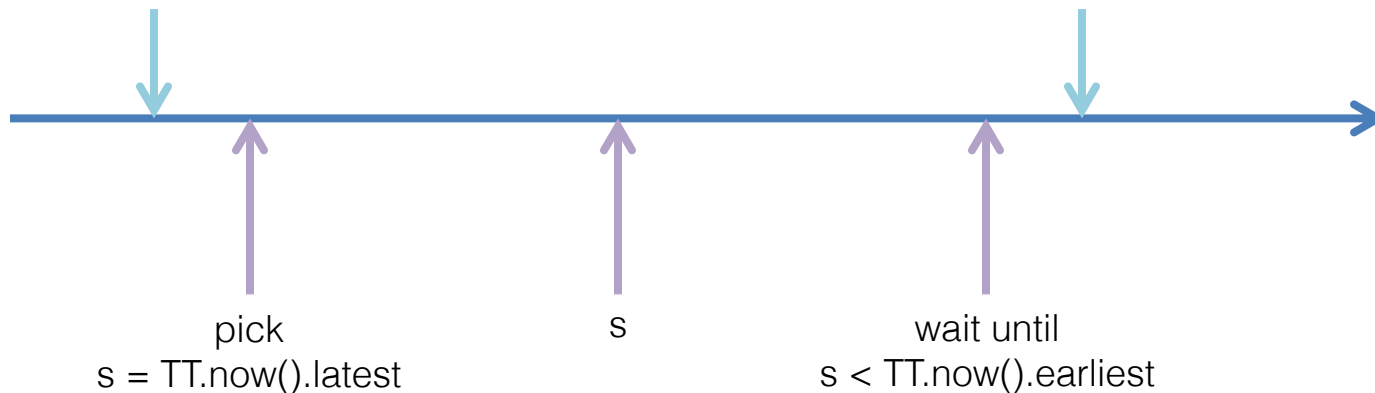- or: $t_{abs}(e_2^{server}) \leq s$

# Read-Write Transactions

Commit wait: cannot see data committed by T until s (assigned timestamp) has passed



pick
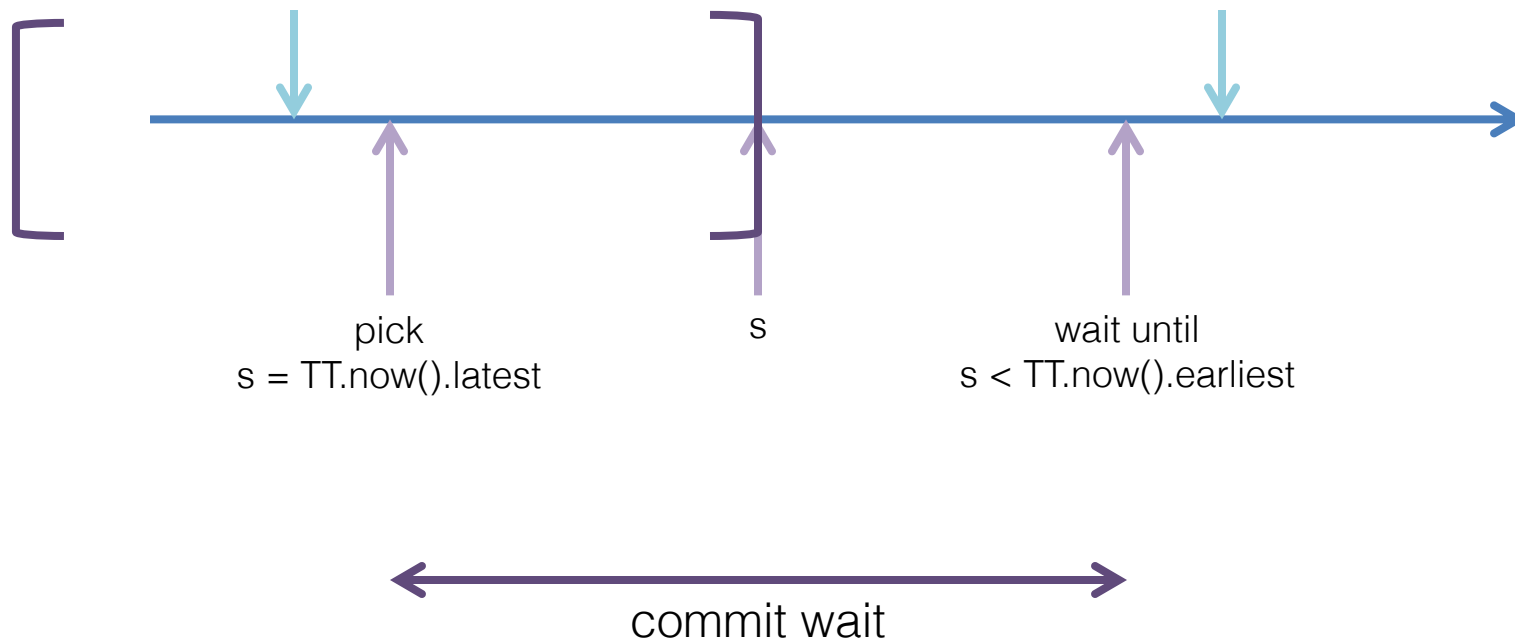s = TT.now().latest

s

# Read-Write Transactions

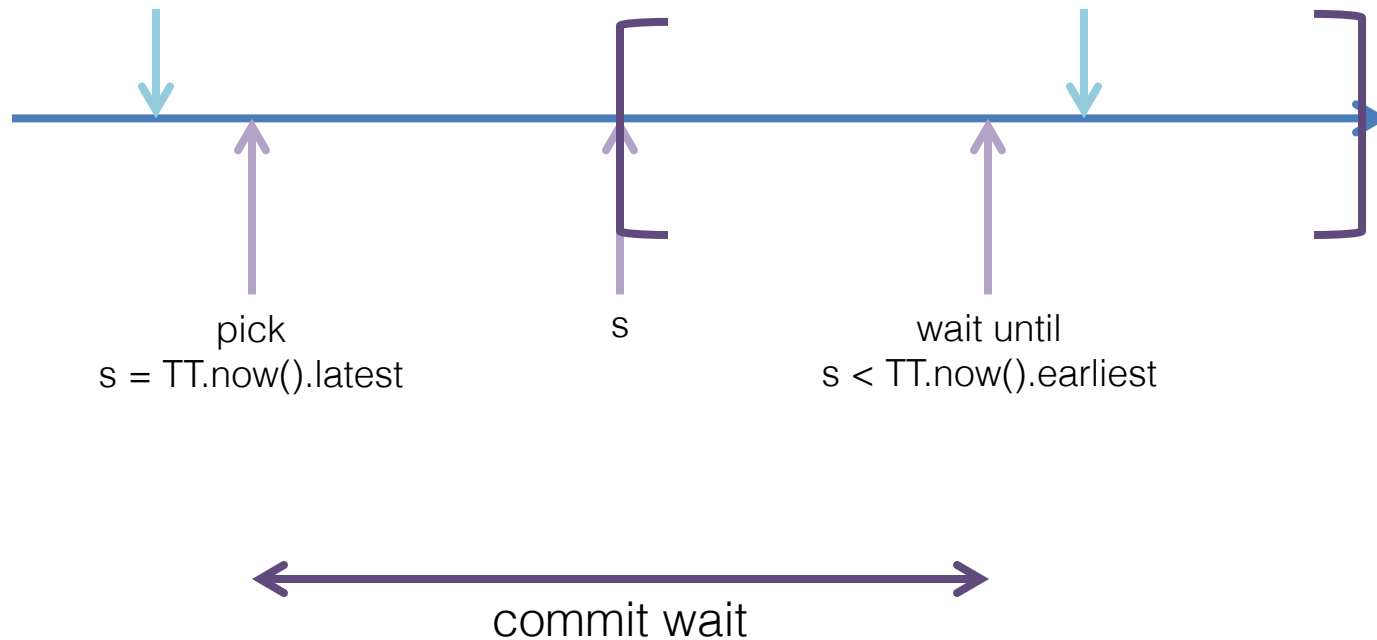Commit wait: cannot see data committed by T until s (assigned timestamp) has passed



pick
s = TT.now().latest

s

wait until
s < TT.now().earliest

# Read-Write Transactions

Commit wait: cannot see data committed by T until s (assigned timestamp) has passed



pick
s = TT.now().latest

s

wait until
s < TT.now().earliest

commit wait

# Read-Write Transactions

Commit wait: cannot see data committed by T until s (assigned timestamp) has passed



pick
s = TT.now().latest

s

wait until
s < TT.now().earliest

commit wait

# Read-Write Transactions

$$s_1 < t_{abs}(e_1^{commit})$$

$$t_{abs}(e_1^{commit}) < t_{abs}(e_2^{start})$$

$$t_{abs}(e_2^{start}) < t_{abs}(e_2^{server})$$

$$t_{abs}(e_2^{server}) \leq s_2$$

$$s_1 < s_2$$

# Read-Write Transactions

Two-phase commit

coordinator
leader →

participant →

participant →

# Read-Write Transactions
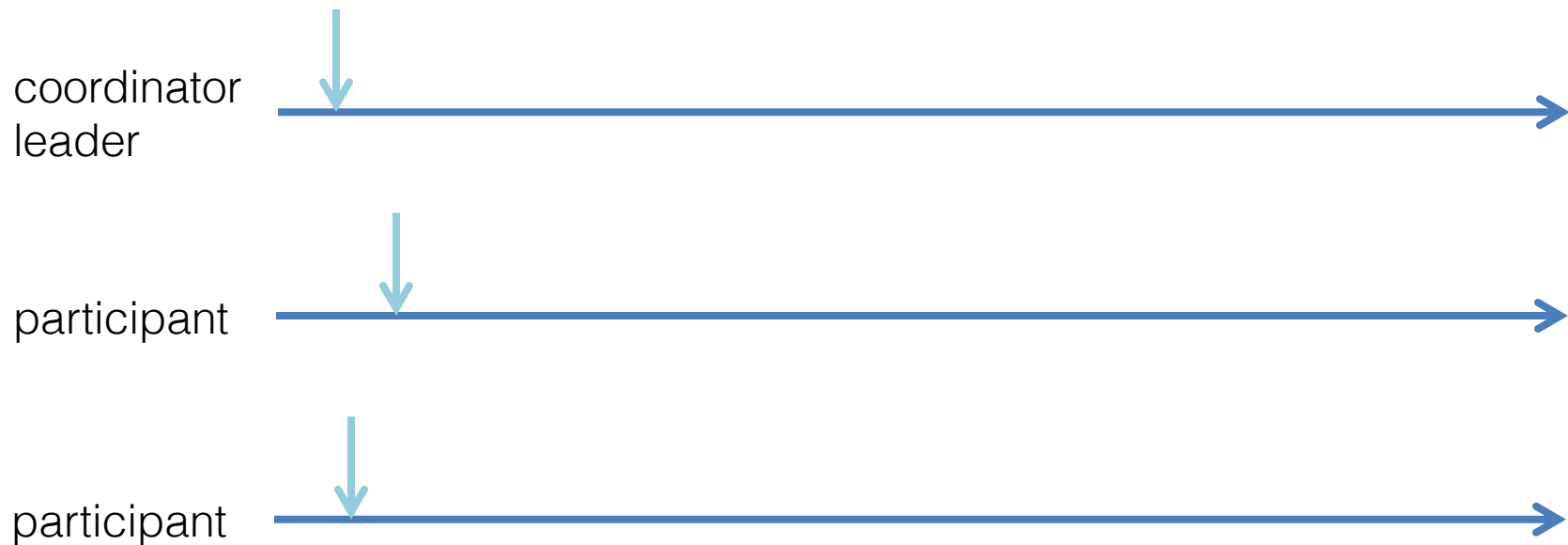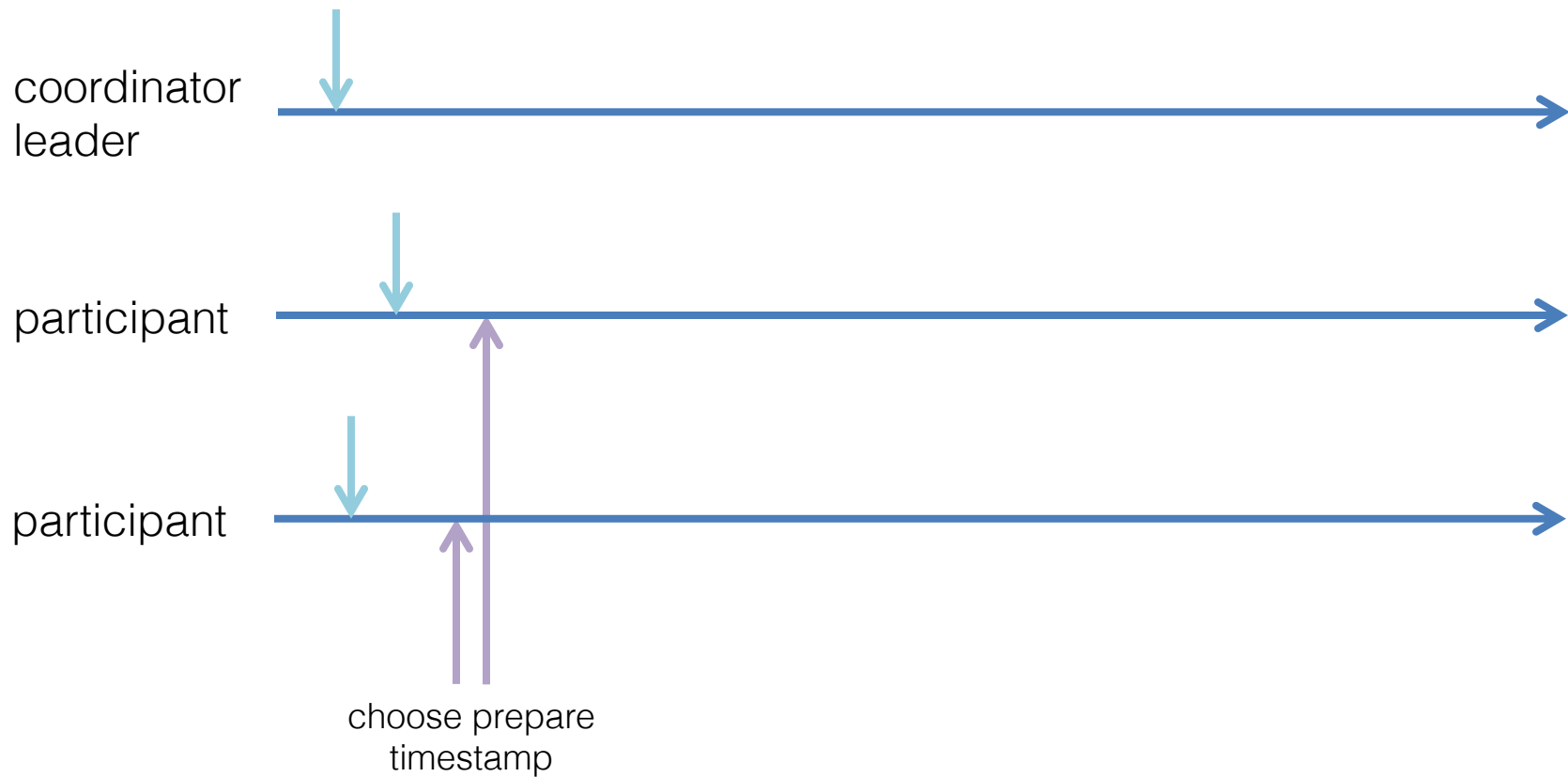
Two-phase commit: client begins

coordinator
leader ──────────────────────────────────────────►

participant ──────────────────────────────────────────►

participant ──────────────────────────────────────────►

# Read-Write Transactions

Two-phase commit

coordinator
leader

participant

participant

# Read-Write Transactions

Two-phase commit



coordinator
leader

participant

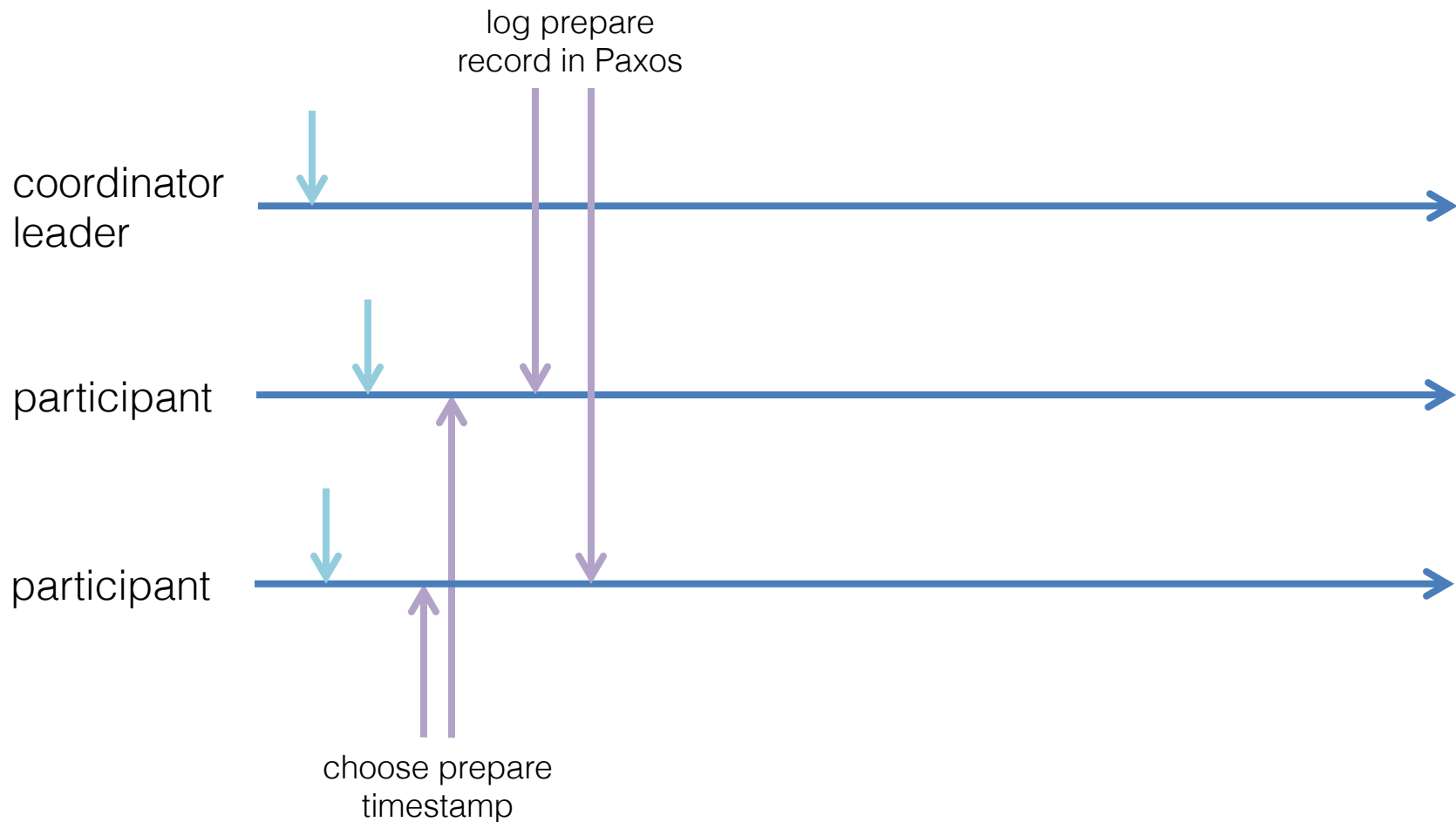participant

choose prepare
timestamp

# Read-Write Transactions

Two-phase commit



log prepare
record in Paxos

coordinator
leader

participant

participant

choose prepare
timestamp

# Read-Write Transactions

## Two-phase commit



log prepare
record in Paxos

coordinator
leader

participant

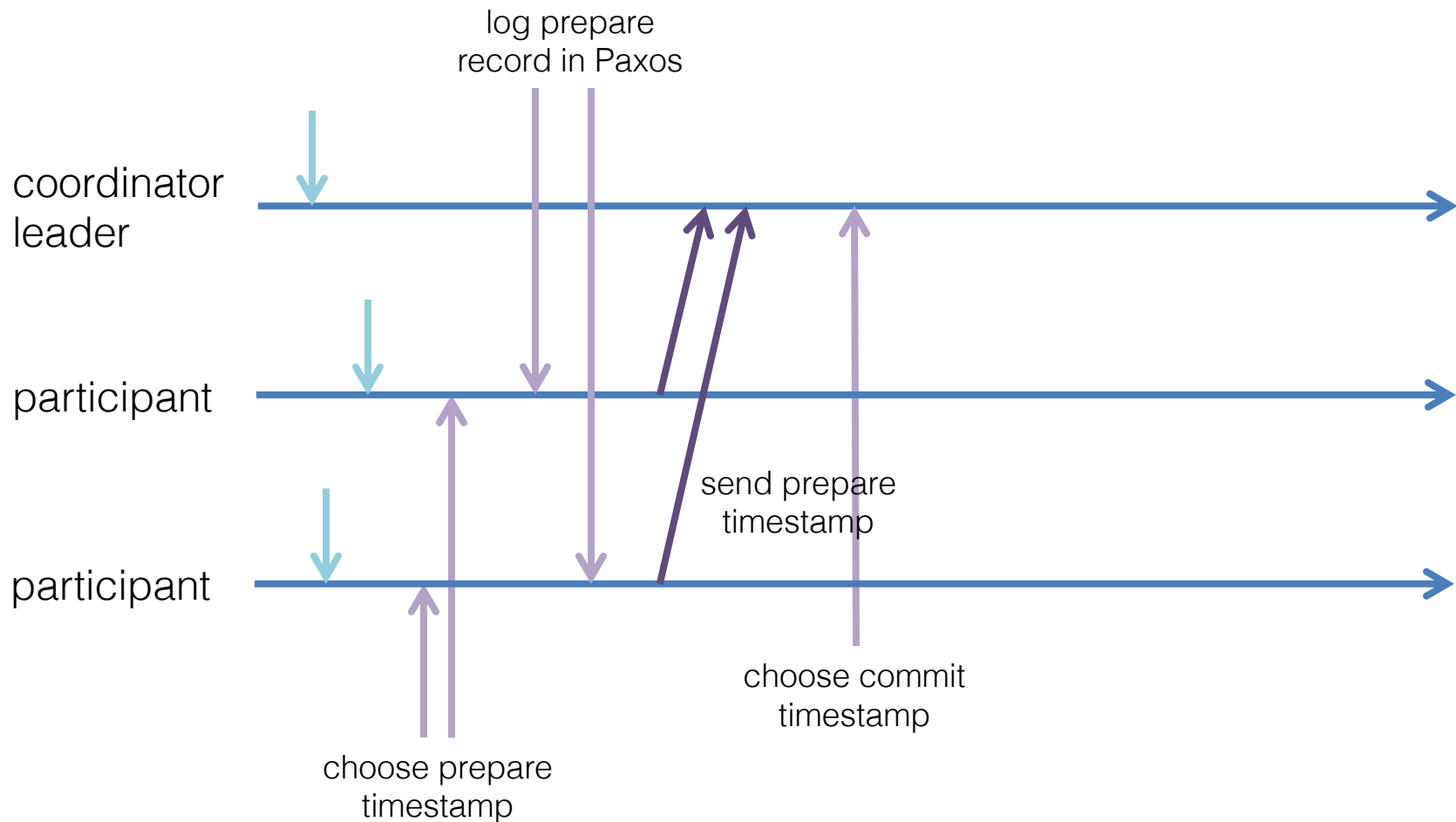participant

send prepare
timestamp

choose prepare
timestamp
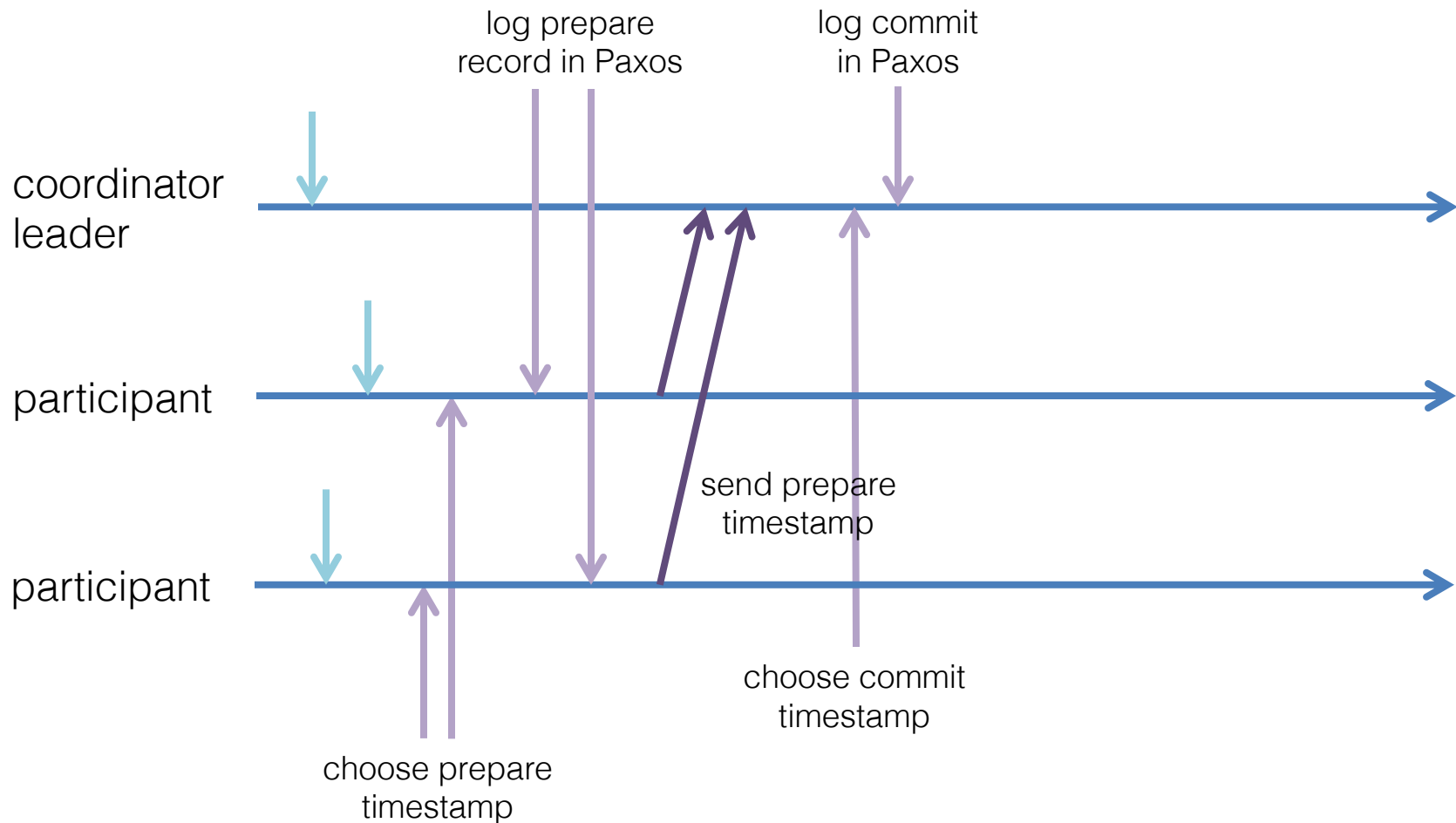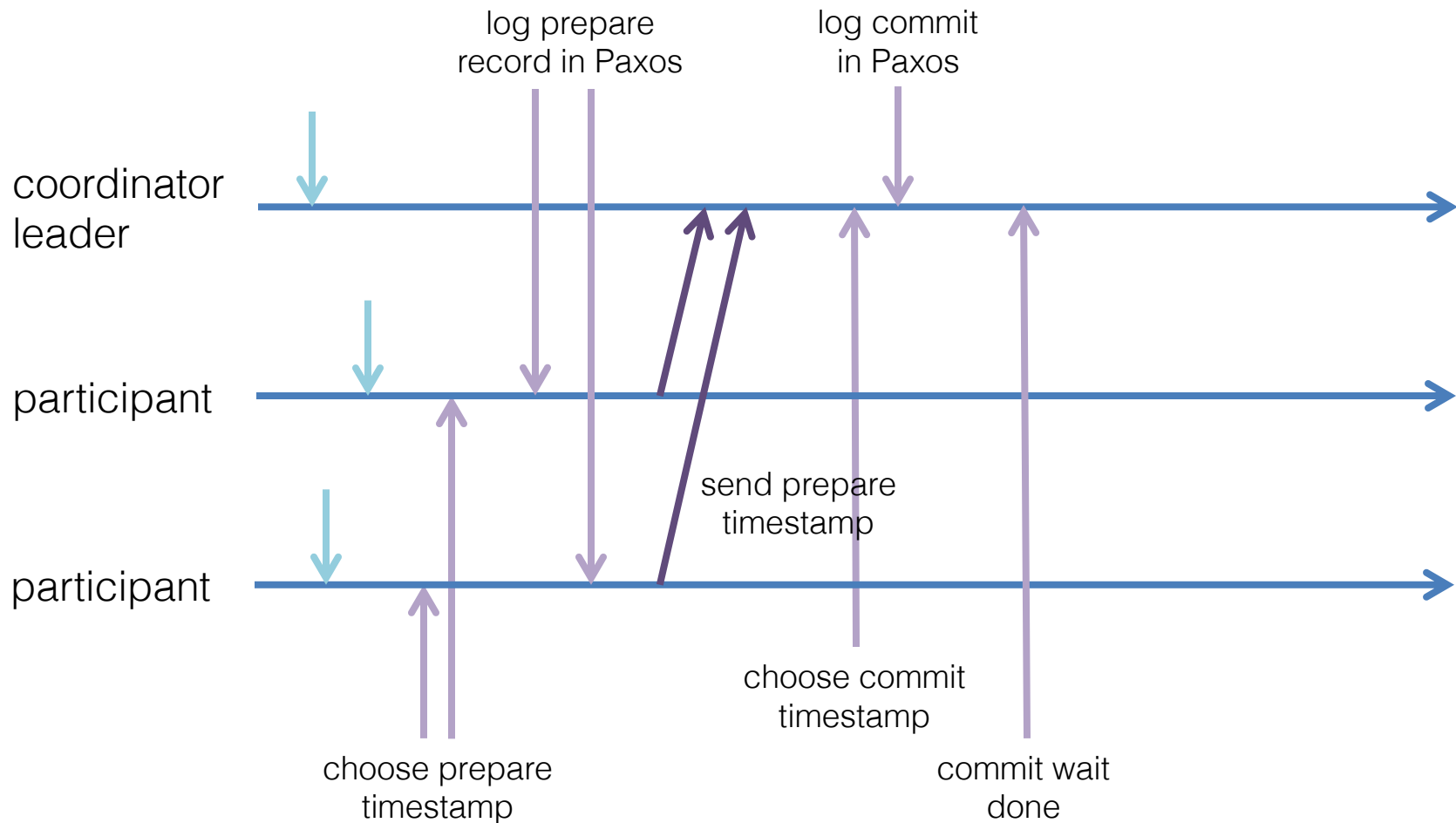
# Read-Write Transactions

## Two-phase commit

# Read-Write Transactions

## Two-phase commit



log prepare
record in Paxos

log commit
in Paxos

coordinator
leader

participant

participant

send prepare
timestamp

choose commit
timestamp

choose prepare
timestamp

# Read-Write Transactions

## Two-phase commit



log prepare
record in Paxos

log commit
in Paxos

coordinator
leader

participant

participant

send prepare
timestamp

choose prepare
timestamp
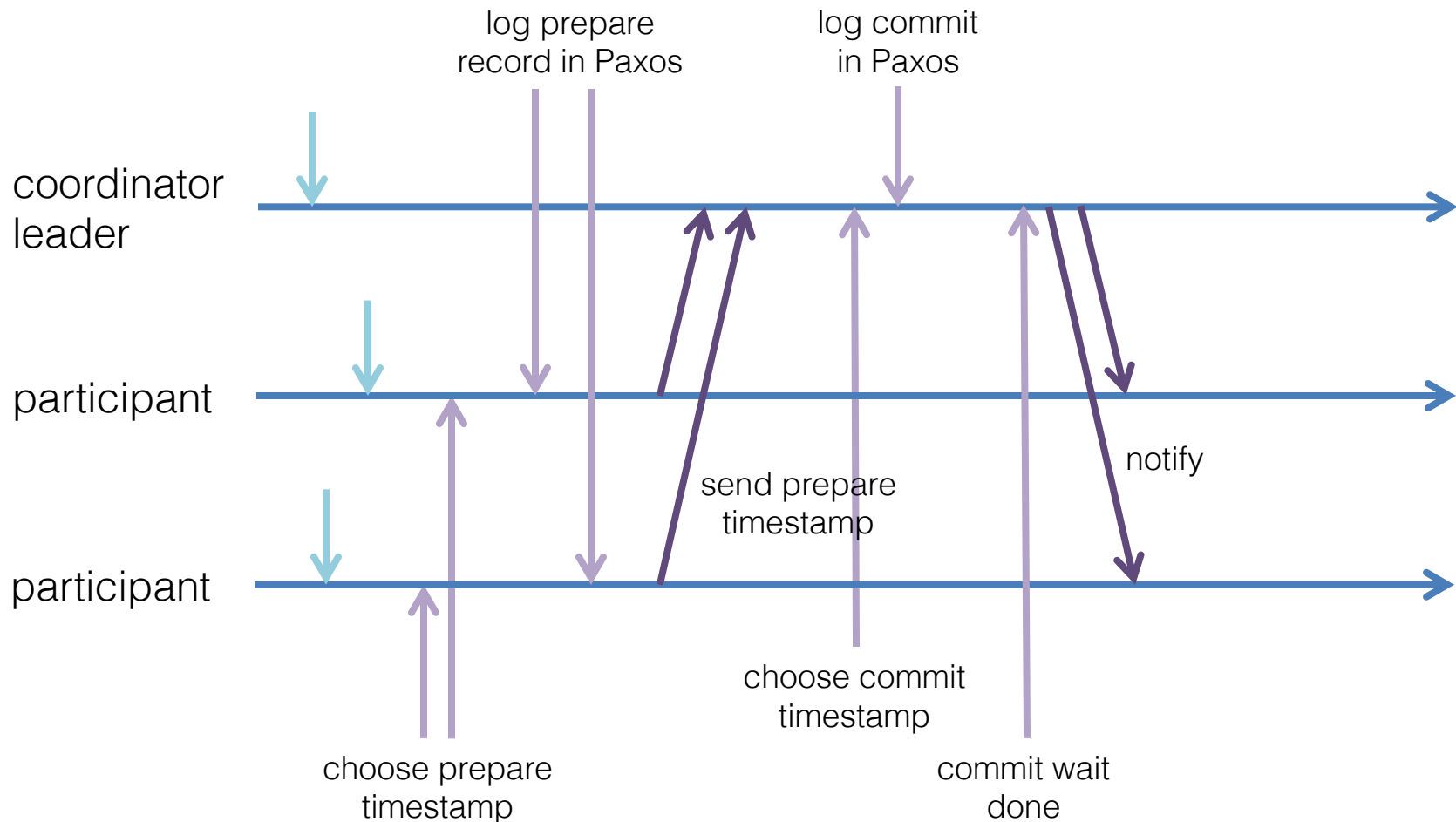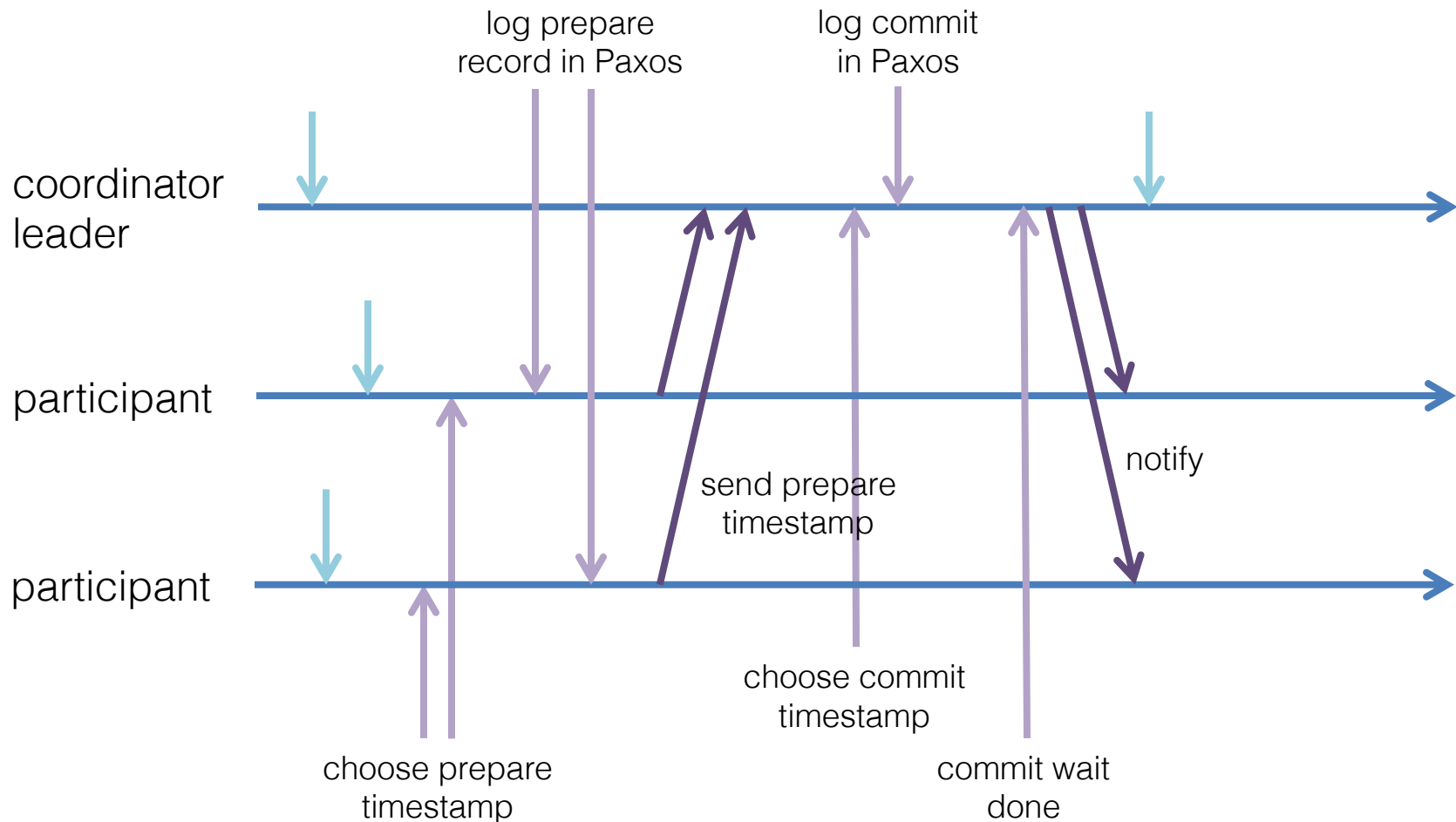
choose commit
timestamp

commit wait
done

# Read-Write Transactions
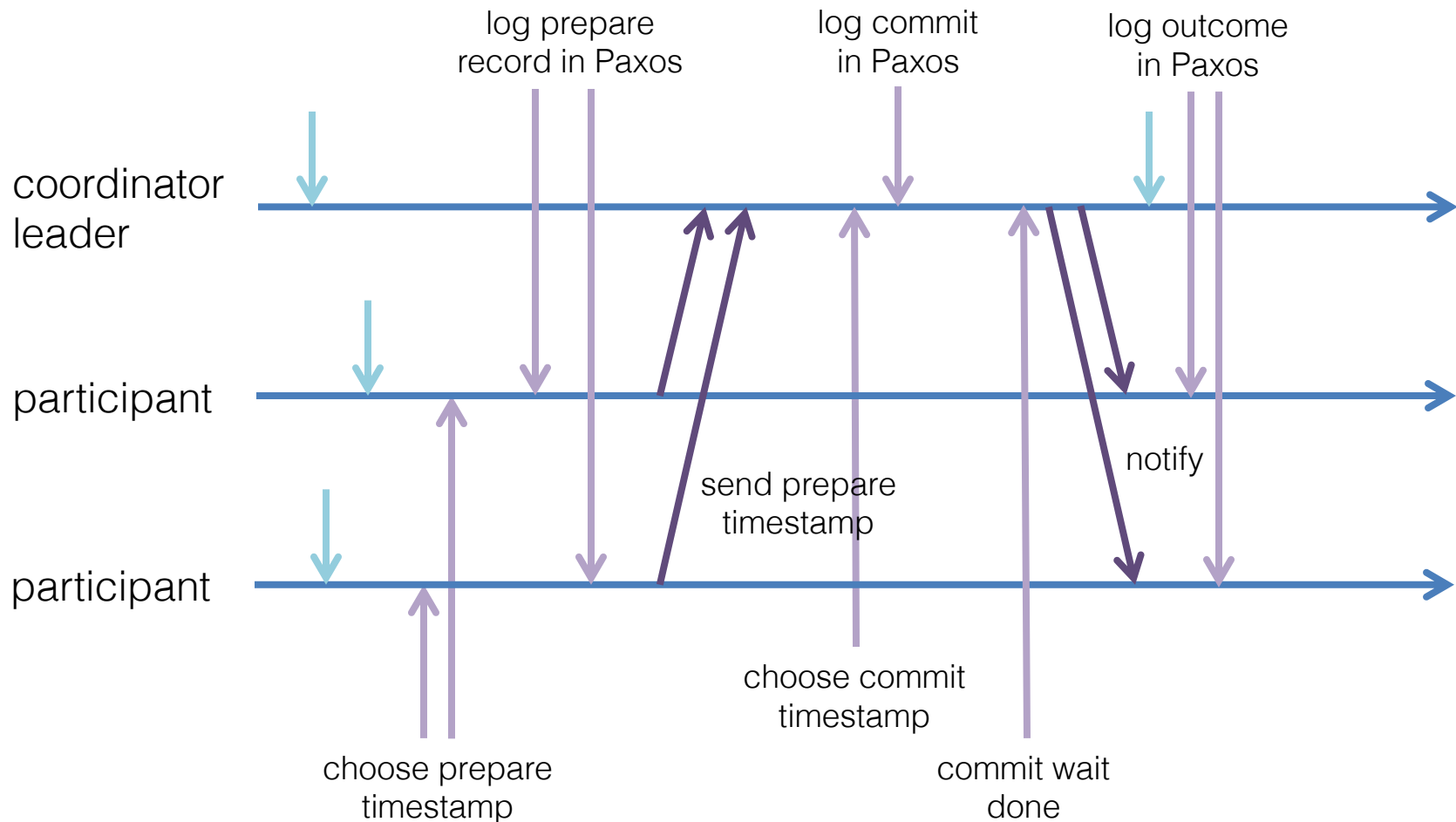
## Two-phase commit

# Read-Write Transactions

## Two-phase commit

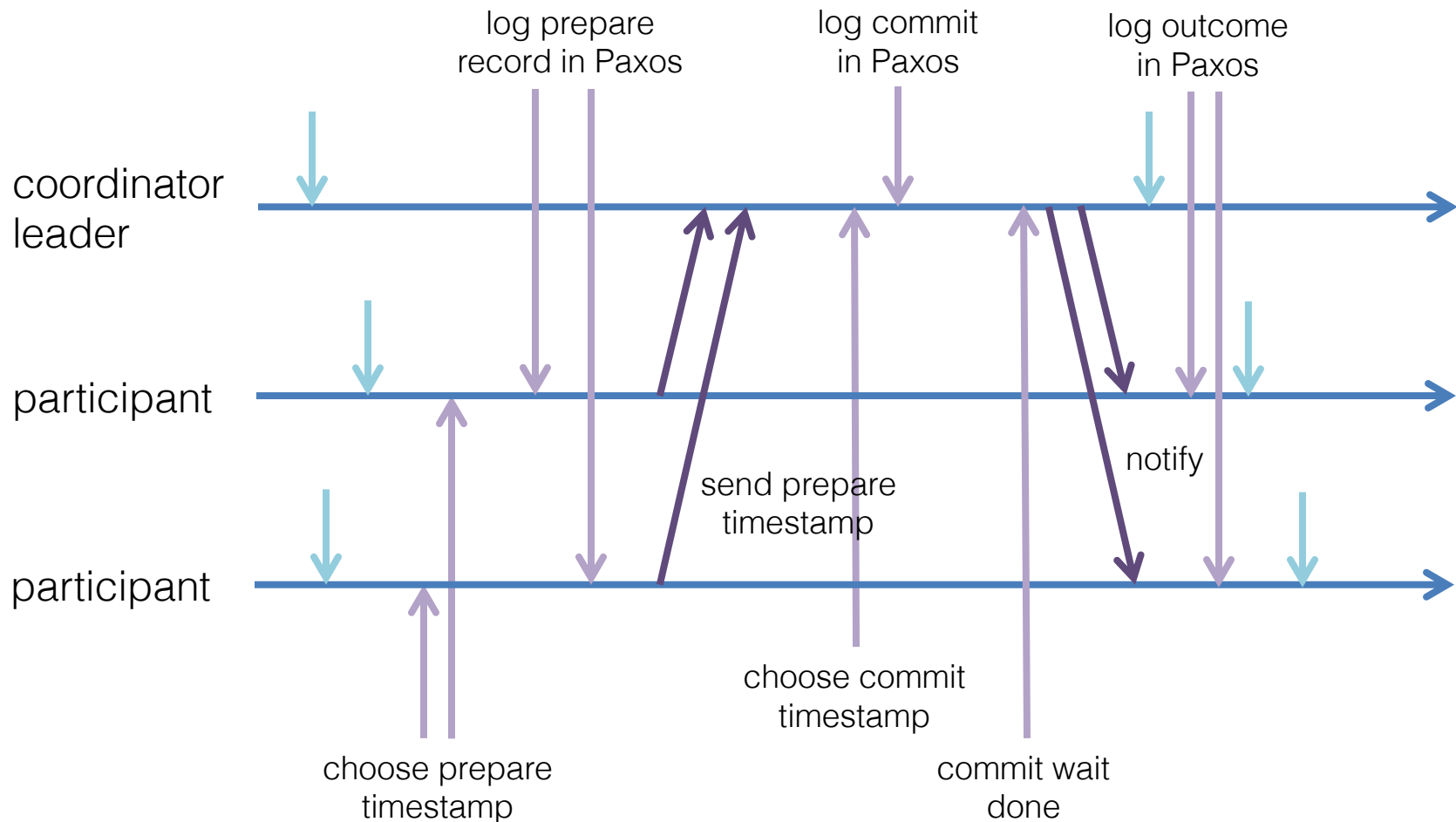# Read-Write Transactions

## Two-phase commit

# Read-Write Transactions

## Two-phase commit
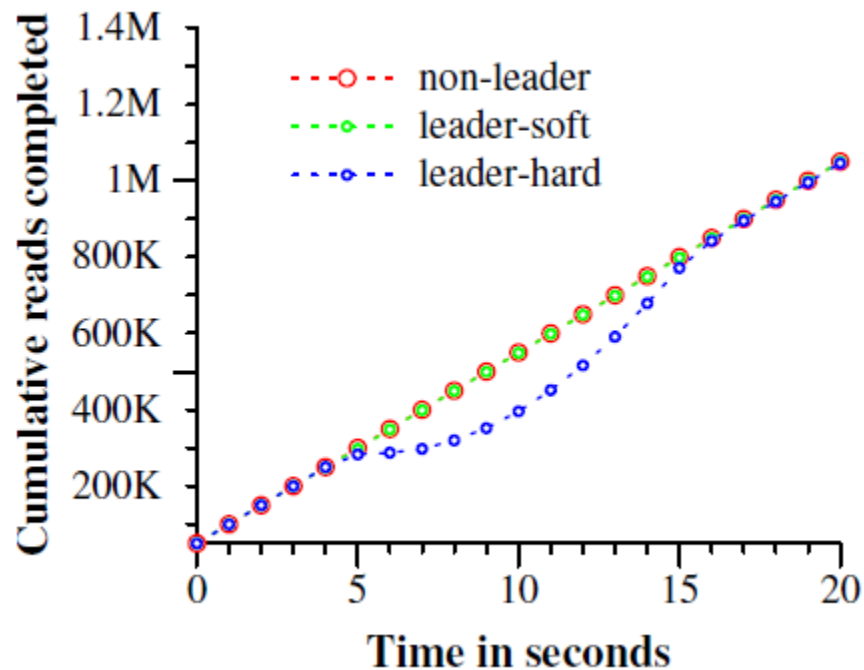
# Read-Only Transactions

- Serving reads at a timestamp
  - Replica tracks safe time $t_{safe}$: can read $t \leq t_{safe}$
  - Define $t_{safe} = min(t^{Paxos}, t^{TM})$
- Assigning timestamps to RO transactions
  - Simplest: assign $s_{read} = $ TT.now().latest
  - May block; should assign oldest timestamp that preserves external consistency

# Microbenchmarks

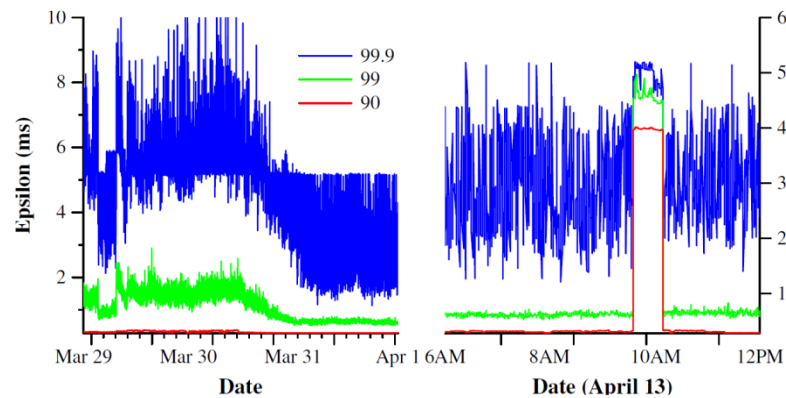| participants | latency (ms) | |
| --- | --- | --- |
| | mean | 99th percentile |
| 1 | 17.0 ±1.4 | 75.0 ±34.9 |
| 2 | 24.5 ±2.5 | 87.6 ±35.9 |
| 5 | 31.5 ±6.2 | 104.5 ±52.2 |
| 10 | 30.0 ±3.7 | 95.6 ±25.4 |
| 25 | 35.5 ±5.6 | 100.4 ±42.7 |
| 50 | 42.7 ±4.1 | 93.7 ±22.9 |
| 100 | 71.4 ±7.6 | 131.2 ±17.6 |
| 200 | 150.5 ±11.0 | 320.3 ±35.1 |

Two-phase commit scalability

# Microbenchmarks



Effect of killing servers on throughput

# Performance

- TrueTime



- F1, Google's advertising backend
  - Automatic failover ☺
  - High standard deviation for latency?

# Final Thoughts

- Implemented at a large scale (F1)!
- Commit wait is pretty clever
- Very dependent on clocks
- Security?

# References

- Corbett et al. "Spanner: Google's Globally-Distributed Database." *Proc. Of OSDI.* 2012.

- http://research.google.com/archive/spanner.html