# Distributed Systems

## CS425/ECE428

*Instructor: Radhika Mittal*

# Logistics

- MP2 due soon (Monday).

# Distributed Transactions and Replication

- Transaction processing can be *distributed* across multiple servers.

  - Different objects can be stored on different servers.

  - An object may be replicated across multiple servers.

- **Case study: Google's Spanner System**
  - Note for exams:
    - no detailed questions from Spanner paper.
    - only some high-level objective questions from materials in slides.

# Spanner: Google's Globally-Distributed Database

- First three lines from the paper:

  - Spanner is a scalable, globally-distributed database designed, built, and deployed at Google.

  - At the highest level of abstraction, it is a database that shards data across many sets of Paxos state machines in datacenters spread all over the world.

  - Replication is used for global availability and geographic locality; clients automatically failover between replicas.

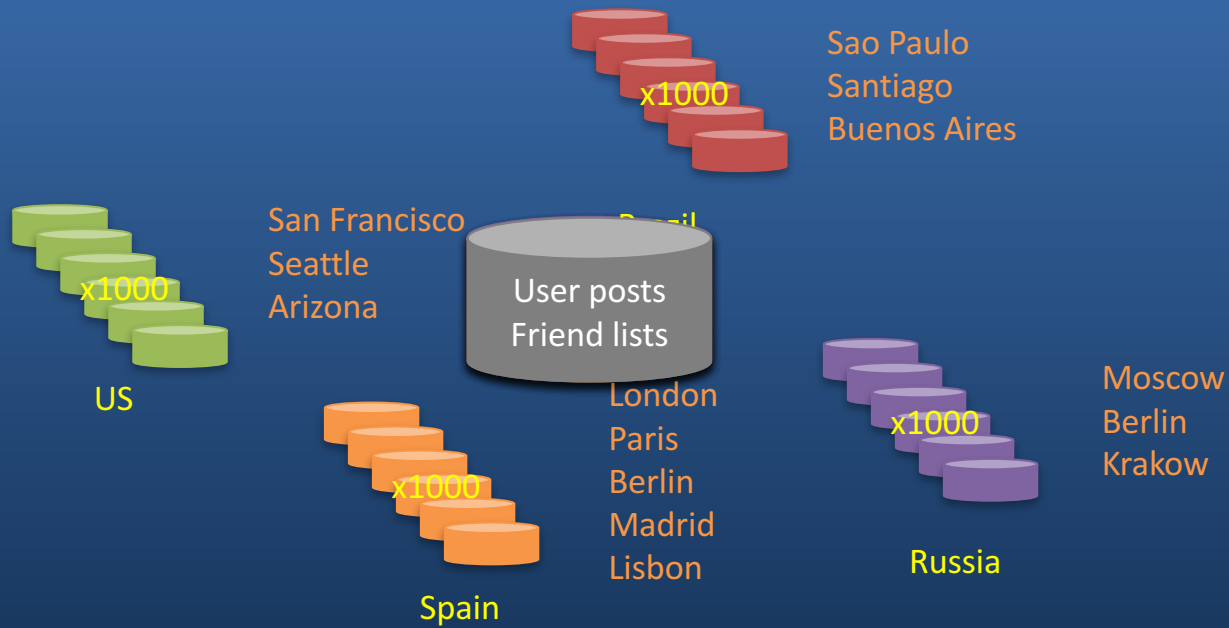# Spanner: Google's Globally-Distributed Database

Wilson Hsieh

representing a host of authors

OSDI 2012

Google™

# What is Spanner?

- Distributed multiversion database
  - General-purpose transactions (ACID)
  - SQL query language
  - Schematized tables
  - Semi-relational data model

- Running in production
  - Storage for Google's ad data
  - Replaced a sharded MySQL database

Google™

# Example: Social Network

x1000

Sao Paulo
Santiago
Buenos Aires

San Francisco
Seattle
Arizona

Brazil

User posts
Friend lists

x1000

US

x1000

London
Paris
Berlin
Madrid
Lisbon

Moscow
Berlin
Krakow

x1000

Russia

x1000

Spain

Google™

# Overview

- Feature: Lock-free distributed read transactions
- Property: External consistency of distributed transactions
    - First system at global scale
- Implementation: Integration of concurrency control, replication, and 2PC
    - Correctness and performance
- Enabling technology: TrueTime
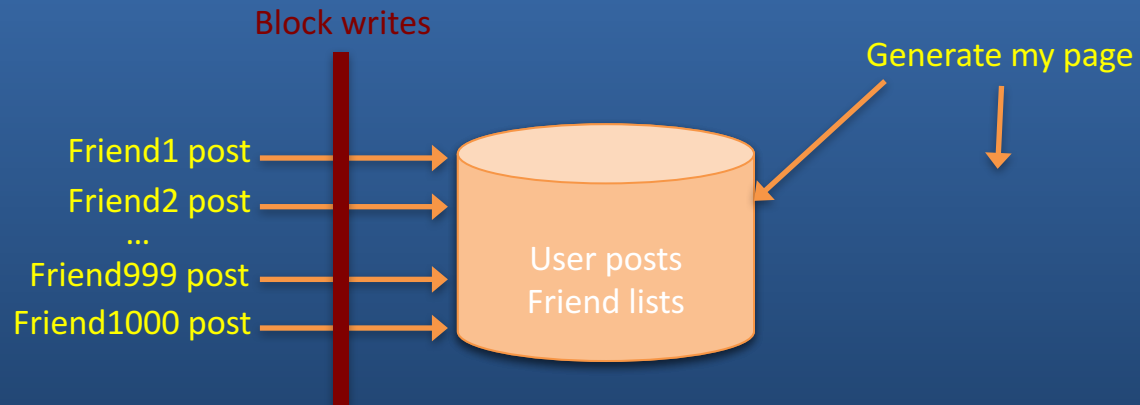    - Interval-based global time

Google™

# Read Transactions

- Generate a page of friends' recent posts
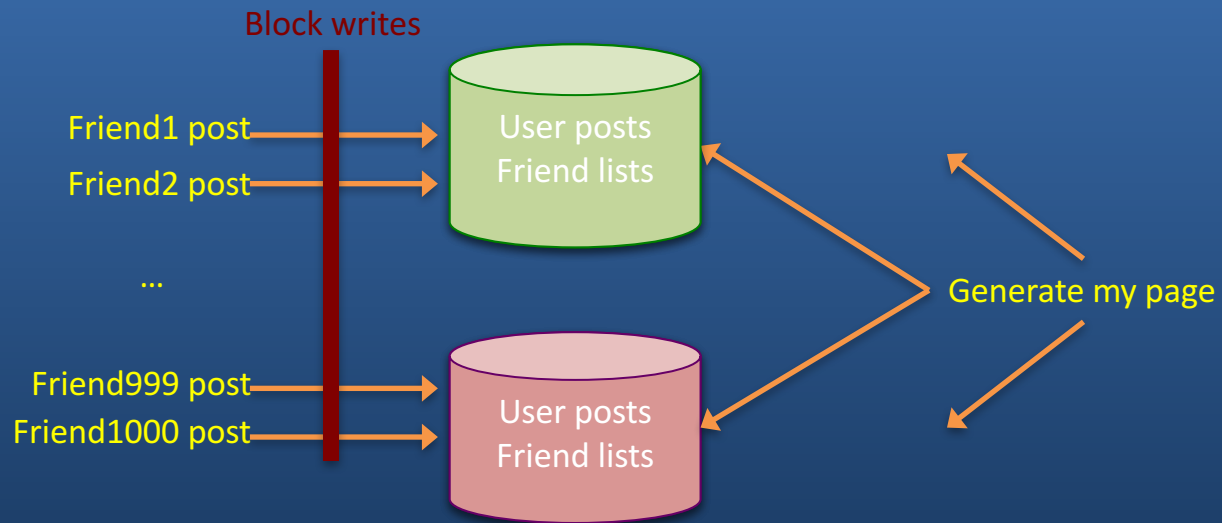  - Consistent view of friend list and their posts

Why consistency matters

1. Remove untrustworthy person X as friend
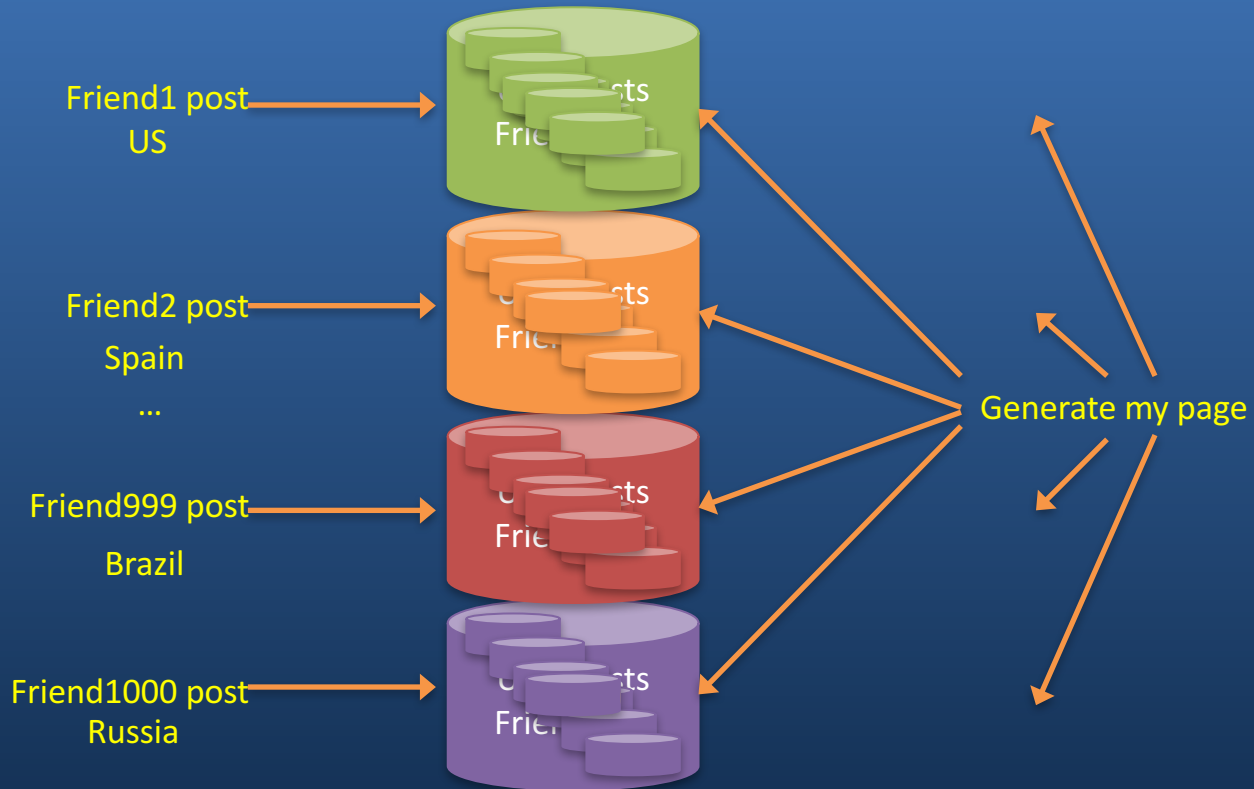2. Post P: "My government is repressive…"

Google™

# Single Machine



Block writes

Generate my page

Friend1 post

Friend2 post

…

Friend999 post

Friend1000 post

User posts
Friend lists

# Multiple Machines



Block writes

Friend1 post

Friend2 post

…

Friend999 post

Friend1000 post

User posts
Friend lists

User posts
Friend lists

Generate my page

# Multiple Datacenters



Friend1 post
US

Friend2 post
Spain
…

Friend999 post
Brazil

Friend1000 post
Russia

Generate my page

# Version Management

- Transactions that write use strict 2PL
  - Each transaction $T$ is assigned a timestamp $s$
  - Data written by $T$ is timestamped with $s$

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

Google™

# Synchronizing Snapshots

Global wall-clock time

==

External Consistency:
Commit order respects global wall-time order

==

Timestamp order respects global wall-time order
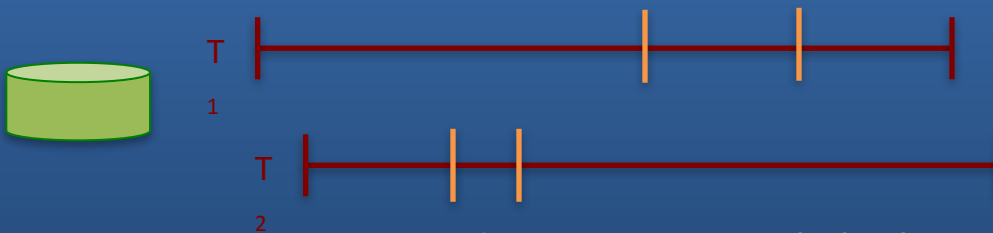
given

timestamp order == commit order

Google™

# Timestamps, Global Clock

- Strict two-phase locking for write transactions
- Assign timestamp while locks are held

Acquired locks                    Release locks

T |————————————————————————|
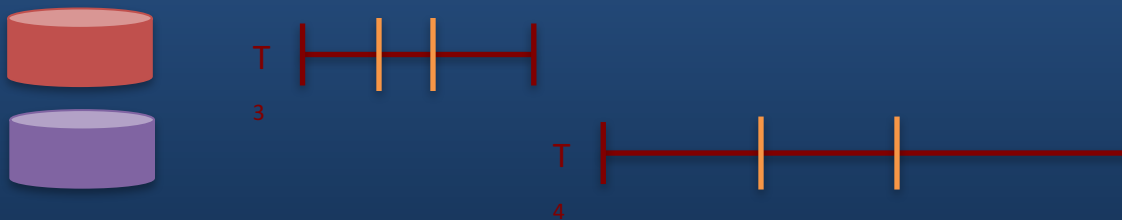
Pick $s$ = now()

# Timestamp Invariants

- Timestamp order == commit order



- Timestamp order respects global wall-time order

# TrueTime

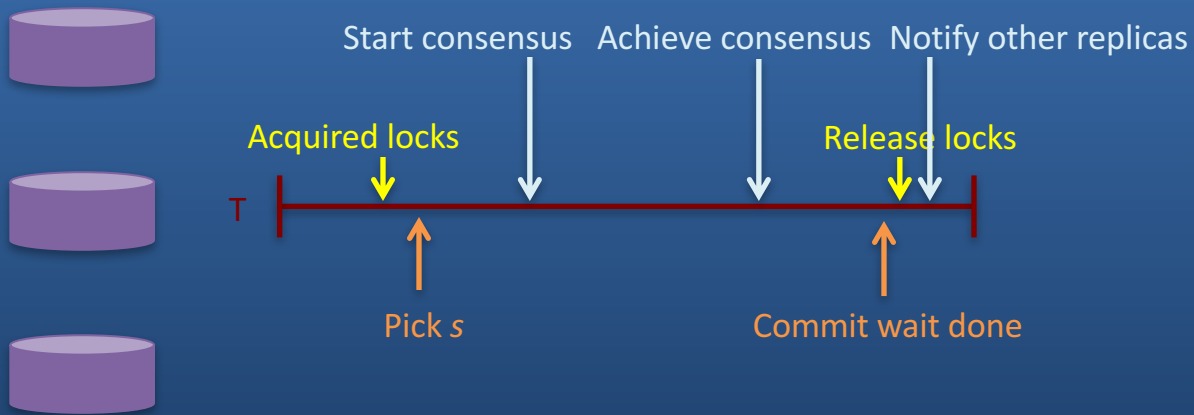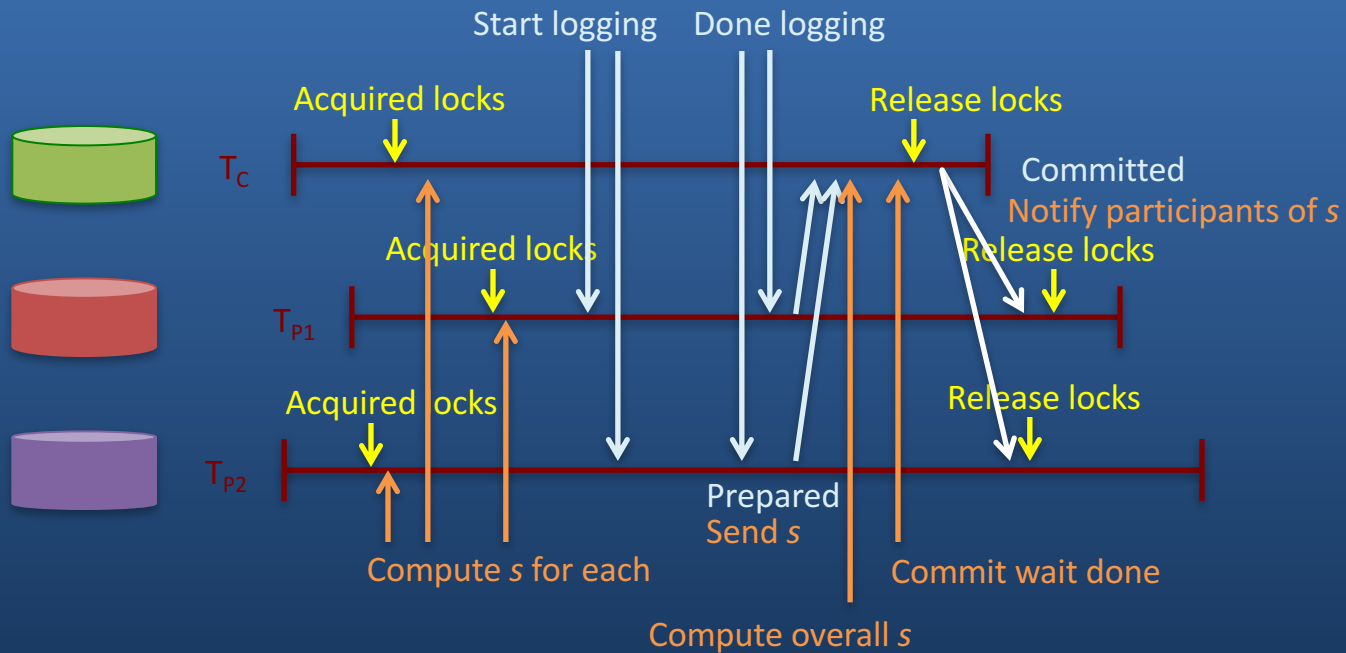- "Global wall-clock time" with bounded uncertainty



TT.now()

earliest          latest

$2*\varepsilon$

time

# Timestamps and TrueTime

Acquired locks                                    Release locks

T

Pick $s$ = TT.now().latest          $s$     Wait until TT.now().earliest > $s$

Commit wait

average ε | average ε

# Commit Wait and Replication

Google™

# Commit Wait and 2-Phase Commit



Start logging    Done logging

Acquired locks                                    Release locks

$T_C$                                              Committed
                                                   Notify participants of $s$

Acquired locks                                     Release locks

$T_{P1}$

Acquired locks                                     Release locks

$T_{P2}$

                          Prepared
                          Send $s$

Compute $s$ for each           Commit wait done

Compute overall $s$

# Example

Remove X
from my friend
list

Risky post P

$T_C$ |————|————————|————|

$s_C=6$      $s=8$

$T_2$ |————|————————|————|

$s=15$

Remove myself
from X's friend
list

$T_P$ |————————————————|————|

$s_P=8$      $s=8$

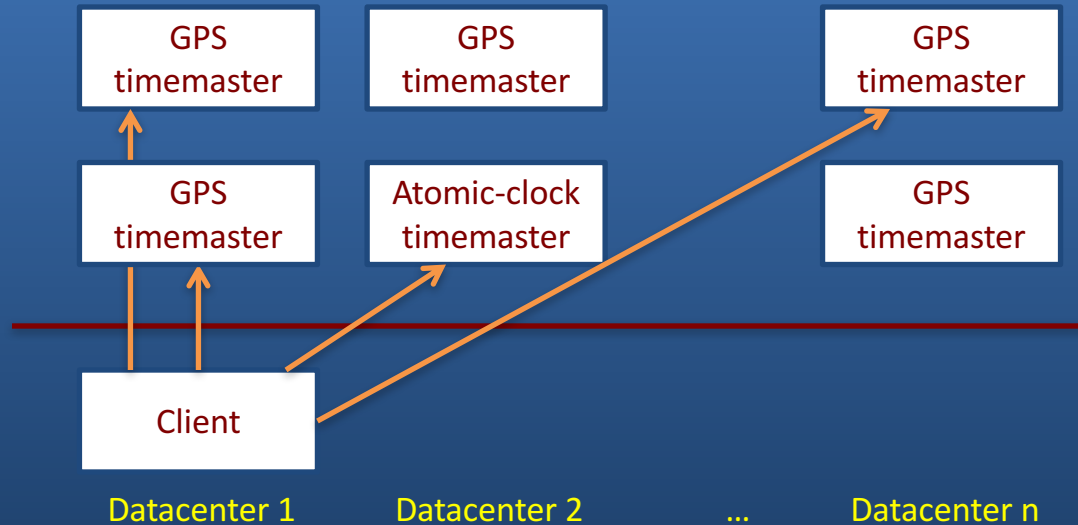| Time | <8 | 8 | 15 |
|---|---|---|---|
| My friends | [X] | [] | |
| My posts | | | [P] |
| X's friends | [me] | [] | |

Google™

# What Have We Covered?

- Lock-free read transactions across datacenters
- External consistency
- Timestamp assignment
- TrueTime
  - Uncertainty in time can be waited out

Google™

# What Haven't We Covered?

- How to read at the present time
- Atomic schema changes
  - Mostly non-blocking
  - Commit in the future
- Non-blocking reads in the past
  - At any sufficiently up-to-date replica
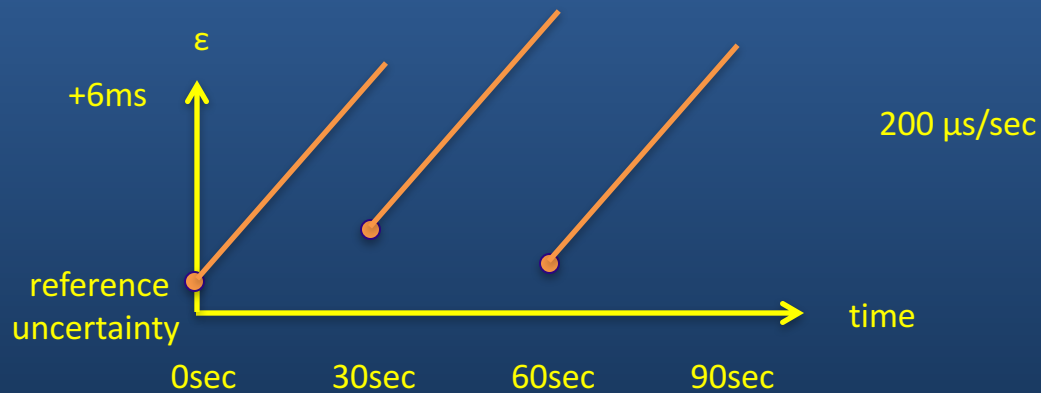
Google™

# TrueTime Architecture



Compute reference [earliest, latest] = now ± ε

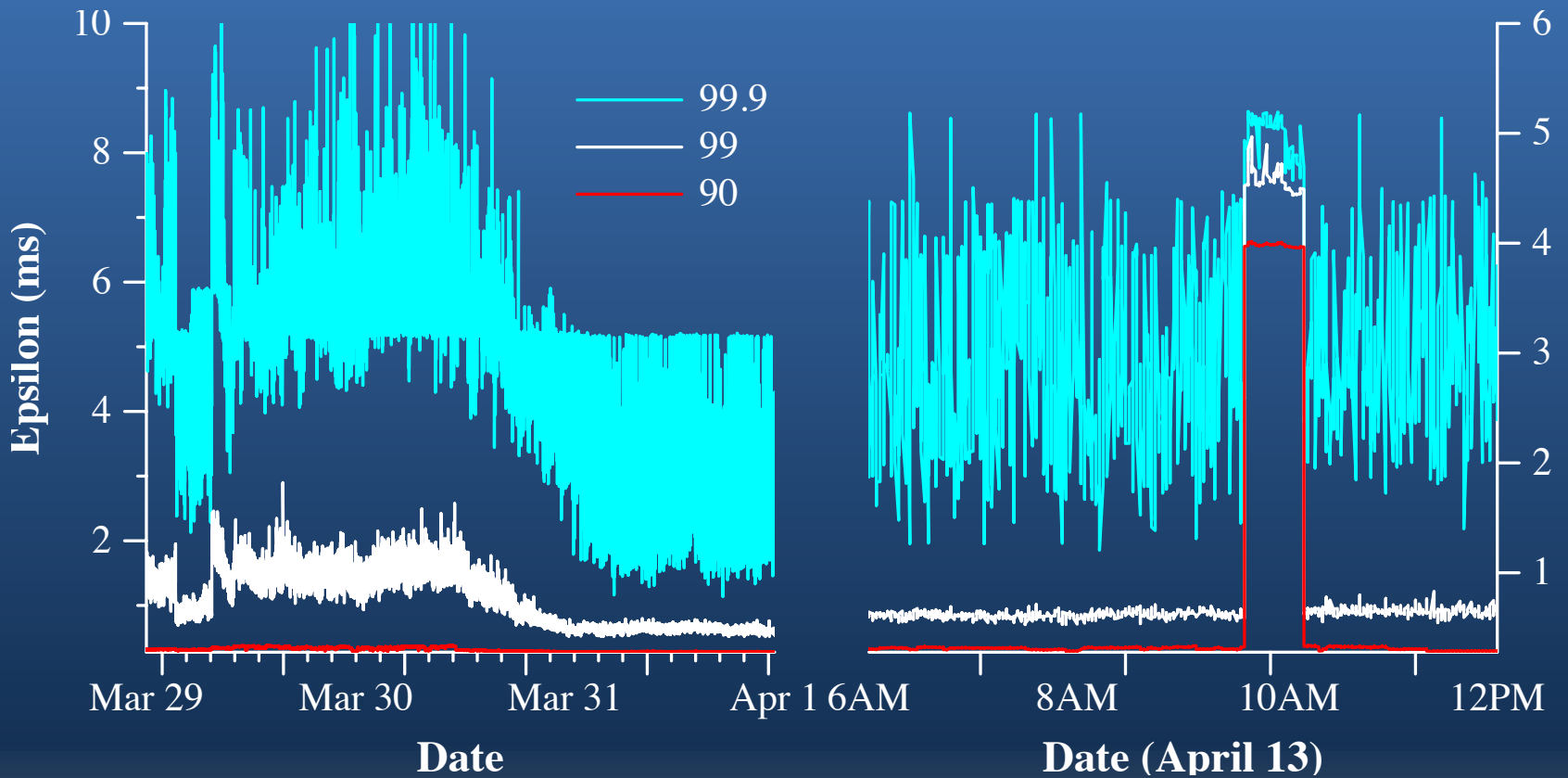# TrueTime implementation

now = reference now + local-clock offset

ε = reference ε + worst-case local-clock drift

# What If a Clock Goes Rogue?

- Timestamp assignment would violate external consistency
- Empirically unlikely based on 1 year of data
  - Bad CPUs 6 times more likely than bad clocks

# Network-Induced Uncertainty

# What's in the Literature

- External consistency/linearizability
- Distributed databases
- Concurrency control
- Replication
- Time (NTP, Marzullo)

# Future Work

- Improving TrueTime
  - Lower $\varepsilon < 1$ ms
- Building out database features
  - Finish implementing basic features
  - Efficiently support rich query patterns

Google™

# Conclusions

- Reify clock uncertainty in time APIs
  - Known unknowns are better than unknown unknowns
  - Rethink algorithms to make use of uncertainty
- Stronger semantics are achievable
  - Greater scale != weaker semantics

Google™

# Thanks

- To the Spanner team and customers
- To our shepherd and reviewers
- To lots of Googlers for feedback
- To you for listening!

- Questions?

Google™