

# Final Review + Systems Potpourri



**Lecture #26**



**Database Systems**

15-445/15-645

Fall 2017



**Andy Pavlo**

Computer Science Dept.

Carnegie Mellon Univ.

# ADMINISTRIVIA

**Project #4:** TODAY @ 11:59pm



# FINAL EXAM

**Who:** You

**What:** <http://cmudb.io/f17-final>

**When:** Friday Dec 15th @ 5:30pm

**Where:** GHC 4401

**Why:** Because otherwise Joy and I call your family over the break...



# FINAL EXAM

## What to bring:

- CMU ID
- Calculator
- Two pages of handwritten notes (double-sided)

## Optional:

- Spare change of clothes

## What not to bring:

- Your roommate



# COURSE EVALS

Your feedback is strongly needed:

→ <https://cmu.smartevals.com>

Things that we want feedback on:

- Homework Assignments
- Projects
- Reading Materials
- Lectures



# EXTENDED OFFICE HOURS

Andy:

→ Wednesday Dec. 13th @ 12:00pm-1:30pm



# STUFF BEFORE MID-TERM

SQL

Buffer Pool Management

Hash Tables

B+Trees

Storage Models



# PARALLEL EXECUTION

Inter-Query Parallelism

Intra-Query Parallelism

Inter-Operator Parallelism

Intra-Operator Parallelism





# EMBEDDED LOGIC

User-defined Functions

Stored Procedures

Focus on advantages vs. disadvantages



# TRANSACTIONS

## ACID

### Conflict Serializability:

- How to check?
- How to ensure?

### View Serializability

### Recoverable Schedules

### Isolation Levels / Anomalies



# TRANSACTIONS

## Two-Phase Locking

- Strict vs. Non-Strict
- Deadlock Detection & Prevention

## Multiple Granularity Locking

- Intention Locks

## B+Tree Latch Crabbing

## Locks vs. Latches



# TRANSACTIONS

## Timestamp Ordering Concurrency Control

→ Thomas Write Rule

## Optimistic Concurrency Control

→ Read Phase

→ Validation Phase

→ Write Phase

## Multi-Version Concurrency Control

→ Version Storage / Ordering

→ Garbage Collection



# CRASH RECOVERY

Buffer Pool Policies:

→ STEAL vs. NO-STEAL

→ FORCE vs. NO-FORCE

Write-Ahead Logging

Logging Schemes

Checkpoints

ARIES Recovery

→ Log Sequence Numbers

→ CLRs



# DISTRIBUTED DATABASES

System Architectures

Replication

Partitioning Schemes

Two-Phase Commit



## 2015 – Top 10

MongoDB	32
Google Spanner/F1	22
LinkedIn Espresso	16
Apache Cassandra	16
Facebook Scuba	16
Apache Hbase	14
VoltDB	10
Redis	10
Vertica	5
Cloudera Impala	5

## 2016 – Top 10

MongoDB	33
Google Spanner/F1	22
Apache Cassandra	19
Facebook Scuba	17
Redis	16
Apache Hbase	15
CockroachDB	12
LinkedIn Espresso	11
Cloudera Impala	8
Peloton	7

## 2017 – Top 10

Google Spanner/F1	15
MongoDB	14
CockroachDB	10
Apache Hbase	9
Peloton	8
Facebook Scuba	6
Cloudera Impala	6
Apache Hive	6
Apache Cassandra	5
LinkedIn Espresso	5

## Bigtable: A Distributed Storage System for Structured Data

Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach  
Mike Ruber, Tushar Chandra, Andrew Fikes, Robert E. Gruber

{fayjeff,sanjay,wilson,herrn,m3b,nashar,fikes,gruber}@google.com

Google, Inc.

### Abstract

Bigtable is a distributed storage system for managing structured data that is designed to scale to a very large size: petabytes of data across thousands of commodity servers. Many projects at Google store data in Bigtable, including web indexing, Google Earth, and Google Finance. These applications place very different demands on Bigtable, both in terms of data size (from URLs to web pages to satellite imagery) and latency requirements (from backend bulk processing to real-time data serving). Despite these varied demands, Bigtable has successfully provided a flexible, high-performance solution for all of these Google products. In this paper we describe the simple data model provided by Bigtable, which gives clients dynamic control over data layout and format, and we describe the design and implementation of Bigtable.

### 1 Introduction

Over the last two and a half years we have designed, implemented, and deployed a distributed storage system for managing structured data at Google called Bigtable. Bigtable is designed to reliably scale to petabytes of data and thousands of machines. Bigtable has achieved several goals: wide applicability, scalability, high performance, and high availability. Bigtable is used by more than sixty Google products and projects, including Google Analytics, Google Finance, Outlook, Personalized Search, WriteUp, and Google Earth. These products use Bigtable for a variety of demanding workloads, which range from throughput-oriented batch-processing jobs to latency-sensitive serving of data to end users. The Bigtable clusters used by these products span a wide range of configurations, from a handful to thousands of servers, and store up to several hundred terabytes of data.

In many ways, Bigtable resembles a database: it shares many implementation strategies with databases. Parallel databases [14] and main-memory databases [13] have

achieved scalability and high performance, but Bigtable provides a different interface than such systems. Bigtable does not support a full relational data model; instead, it provides clients with a simple data model that supports dynamic control over data layout and format, and allows clients to reason about the locality properties of the data represented in the underlying storage. Data is indexed using row and column names that can be arbitrary strings. Bigtable also treats data as uninterpreted strings, although clients often serialize various forms of structured and semi-structured data into these strings. Clients can control the locality of their data through careful choices in their schemas. Finally, Bigtable schema parameters let clients dynamically control whether to serve data out of memory or from disk.

Section 2 describes the data model in more detail, and Section 3 provides an overview of the client API. Section 4 briefly describes the underlying Google infrastructure on which Bigtable depends. Section 5 describes the fundamentals of the Bigtable implementation, and Section 6 describes some of the refinements that we made to improve Bigtable's performance. Section 7 provides measurements of Bigtable's performance. We describe several examples of how Bigtable is used at Google in Section 8, and discuss some lessons we learned in designing and supporting Bigtable in Section 9. Finally, Section 10 describes related work, and Section 11 presents our conclusions.

### 2 Data Model

A Bigtable is a sparse, distributed, persistent multi-dimensional sorted map. The map is indexed by a row key, column key, and a timestamp; each value in the map is an uninterpreted array of bytes.

`(row:string, column:string, time:int64) → string`

## Megastore: Providing Scalable, Highly Available Storage for Interactive Services

Jason Baker, Chris Bond, James C. Corbett, JJ Furman, Andrey Khorlin, James Larson, Jean-Michel Léon, Yawei Li, Alexander Lloyd, Vadim Yushprakh

{jasonbaker,chrishond,jcorbett,jfurman,akhorlin,jjlarson,jmlewellyn,allroyd,vadim}@google.com

### ABSTRACT

Megastore is a storage system developed to meet the requirements of today's interactive online services. Megastore blends the scalability of a NoSQL database with the convenience of a traditional RDBMS in a novel way, and provides both strong consistency guarantees and high availability. We provide fully serializable ACID semantics within fine-grained partitions of data. This partitioning allows us to synchronously replicate each write across a wide area network with reasonable latency and support seamless failover between datacenters. This paper describes Megastore's semantics and replication algorithm. It also describes our experience supporting a wide range of Google production services built with Megastore.

### Categories and Subject Descriptors

C.2.4 [Distributed Systems]: Distributed databases; H.2.4 [Database Management]: Systems—concurrency, distributed databases

### General Terms

Algorithms, Design, Performance, Reliability

### Keywords

Large databases, Distributed transactions, Bigtable, Paxos

### 1. INTRODUCTION

Interactive online services are forcing the storage community to meet new demands as desktop applications migrate to the cloud. Services like email, collaborative documents, and social networking have been growing exponentially and are testing the limits of existing infrastructure. Meeting these services' storage demands is challenging due to a number of conflicting requirements.

First, the Internet brings a huge audience of potential users, so the applications must be highly scalable. A service

can be built rapidly using MySQL [16] as its datastore, but scaling the service to millions of users requires a complete redesign of its storage infrastructure. Second, services must compete for users. This requires rapid development of features and fast time-to-market. Third, the service must be responsive; hence, the storage system must have low latency. Fourth, the service should provide the user with a consistent view of the data—the result of an update should be visible immediately and durably. Seeing edits to a cloud-hosted spreadsheet vanish, however briefly, is a poor user experience. Finally, users have come to expect Internet services to be up 24/7, so the service must be highly available. The service must be resilient to many kinds of faults ranging from the failure of individual disks, machines, or routers all the way up to large-scale outages affecting entire datacenters.

These requirements are in conflict. Relational databases provide a rich set of features for easily building applications, but they are difficult to scale to hundreds of millions of users. NoSQL datastores such as Google's Bigtable [15], Apache Hadoop's HBase [1], or Facebook's Cassandra [6] are highly scalable, but their limited API and loose consistency models complicate application development. Replicating data across distant datacenters while providing low latency is challenging, as is guaranteeing a consistent view of replicated data, especially during faults.

Megastore is a storage system developed to meet the storage requirements of today's interactive online services. It is novel in that it blends the scalability of a NoSQL database with the convenience of a traditional RDBMS. It uses synchronous replication to achieve high availability and a consistent view of the data. In brief, it provides fully serializable ACID semantics over distant replicas with low enough latencies to support interactive applications.

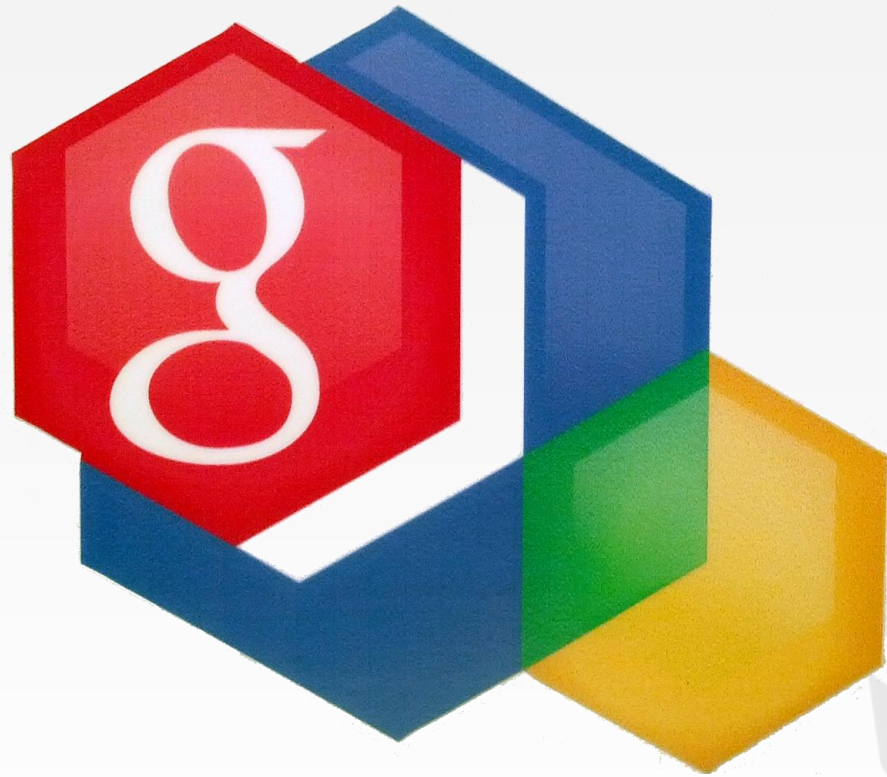
We accomplish this by taking a middle ground in the RDBMS vs. NoSQL design space: we partition the database and replicate each partition separately, providing full ACID semantics within partitions, but only limited consistency guarantees across them. We provide traditional database features, such as secondary indexes, but only those features that can scale within user-tolerable latency limits, and only with the semantics that our partitioning scheme can support. We contend that the data for most Internet services can be suitably partitioned (e.g., by user) to make this approach viable, and that a small, but not spartan, set of features can substantially ease the burden of developing cloud applications.

Contrary to conventional wisdom [24, 28], we were able to use Paxos [27] to build a highly available system that pro-

This article is published under a Creative Commons Attribution License (<http://creativecommons.org/licenses/by/3.0/>), which permits distribution and reproduction in any medium as well allowing derivative works, provided that you attribute the original work to the author(s) and CDRR 2011. <sup>1</sup> 2011 *Network Conference on Innovative Data Systems Research (CDRR '11)*, January 9-12, 2011, Anaheim, California, USA.



# Google Spanner





# GOOGLE SPANNER

Google's geo-replicated DBMS (>2011)

Schematized, semi-relational data model.

Decentralized shared-disk architecture.

Log-structured on-disk storage.

Concurrency Control:

- Strict 2PL + MVCC + Multi-Paxos
- **Externally consistent** global write-transactions with synchronous replication.
- Lock-free read-only transactions.





# PHYSICAL DENORMALIZATION

```
CREATE TABLE users {
  uid INT NOT NULL,
  email VARCHAR,
  PRIMARY KEY (uid)
};
CREATE TABLE albums {
  uid INT NOT NULL,
  aid INT NOT NULL,
  name VARCHAR,
  PRIMARY KEY (uid, aid)
} INTERLEAVE IN PARENT users
ON DELETE CASCADE;
```

## Physical Storage

users(1001)
↳albums(1001, 9990)
↳albums(1001, 9991)
users(1002)
↳albums(1002, 6631)
↳albums(1002, 6634)



# PHYSICAL DENORMALIZATION

```
CREATE TABLE users {  
  uid INT NOT NULL,  
  email VARCHAR,  
  PRIMARY KEY (uid)  
};  
CREATE TABLE albums {  
  uid INT NOT NULL,  
  aid INT NOT NULL,  
  name VARCHAR,  
  PRIMARY KEY (uid, aid)  
} INTERLEAVE IN PARENT users  
ON DELETE CASCADE;
```

## Physical Storage

users(1001)
↳albums(1001, 9990)
↳albums(1001, 9991)
users(1002)
↳albums(1002, 6631)
↳albums(1002, 6634)

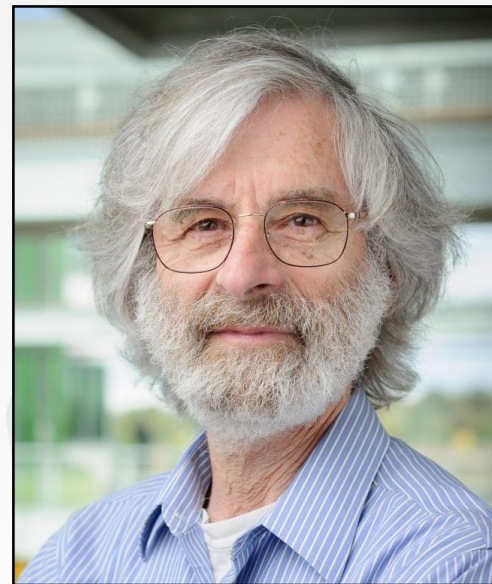


# PAXOS

Consensus protocol where a coordinator proposes an outcome (e.g., commit or abort) and then the participants vote on whether that outcome should succeed.

Does not block if a majority of participants are available and has provably minimal message delays in the best case.

→ First correct protocol that was provably resilient in the face asynchronous networks



Lamport



# PAXOS

Consensus protocol where a coordinator proposes an outcome (e.g., commit or abort) and then the participants vote on whether that outcome should succeed.

Does not block if a majority of participants are available and has provably minimal message delays in the best case.

→ First correct protocol that was provably resilient in the face asynchronous networks

## The Part-Time Parliament

LESLIE LAMPFORT  
Digital Equipment Corporation

Recent archaeological discoveries on the island of Paxos reveal that the parliament functioned despite the peripatetic propensity of its part-time legislators. The legislators maintained consistent copies of the parliamentary record, despite their frequent forays from the chamber and the forgetfulness of their messengers. The Paxos parliament's protocol provides a new way of implementing the state-machine approach to the design of distributed systems.

Categories and Subject Descriptors: C.2.4 [Computer-Communications Networks]: Distributed Systems—Network operating systems; D.4.5 [Operating Systems]: Reliability—Fault-tolerance; J.1 [Administrative Data Processing]: Government

General Terms: Design, Reliability

Additional Key Words and Phrases: State machines, three-phase commit, voting

This submission was recently discovered behind a filing cabinet in the TOCS editorial office. Despite its age, the editor-in-chief felt that it was worth publishing. Because the author is currently doing field work in the Greek isles and cannot be reached, I was asked to prepare it for publication.

The author appears to be an archeologist with only a passing interest in computer science. This is unfortunate; even though the obscure ancient Paxos civilization he describes is of little interest to most computer scientists, its legislative system is an excellent model for how to implement a distributed computer system in an asynchronous environment. Indeed, some of the refinements the Paxos made to their protocol appear to be unknown in the systems literature.

The author does give a brief discussion of the Paxos Parliament's relevance to distributed computing in Section 4. Computer scientists will probably want to read that section first. Even before that, they might want to read the explanation of the algorithm for computer scientists by Lamport [1996]. The algorithm is also described more formally by De Prisco et al. [1997]. I have added further comments on the relation between the ancient protocols and more recent work at the end of Section 4.

Keith Marzullo  
University of California, San Diego

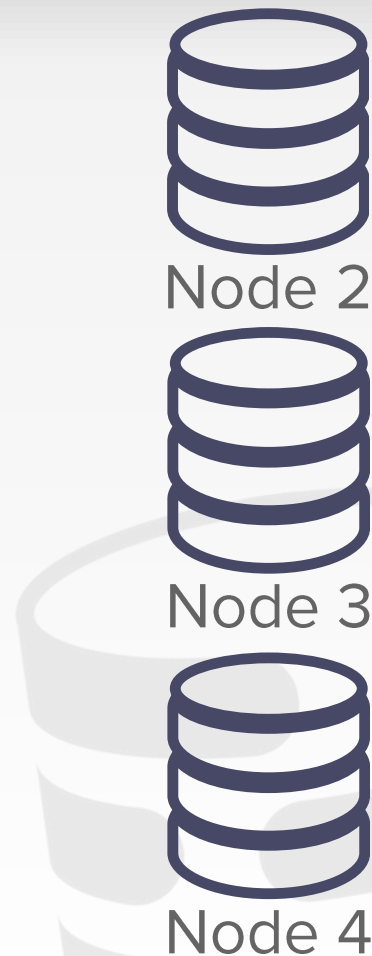
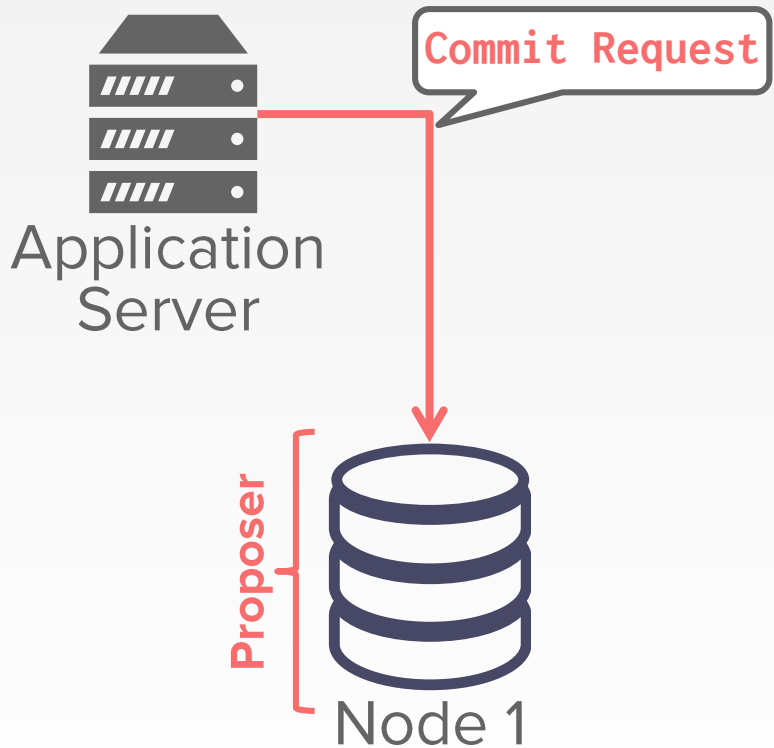
Authors' address: Systems Research Center, Digital Equipment Corporation, 130 Lytton Avenue, Palo Alto, CA 94301.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1998 ACM 0000-0000/98/0000-0000 \$00.00

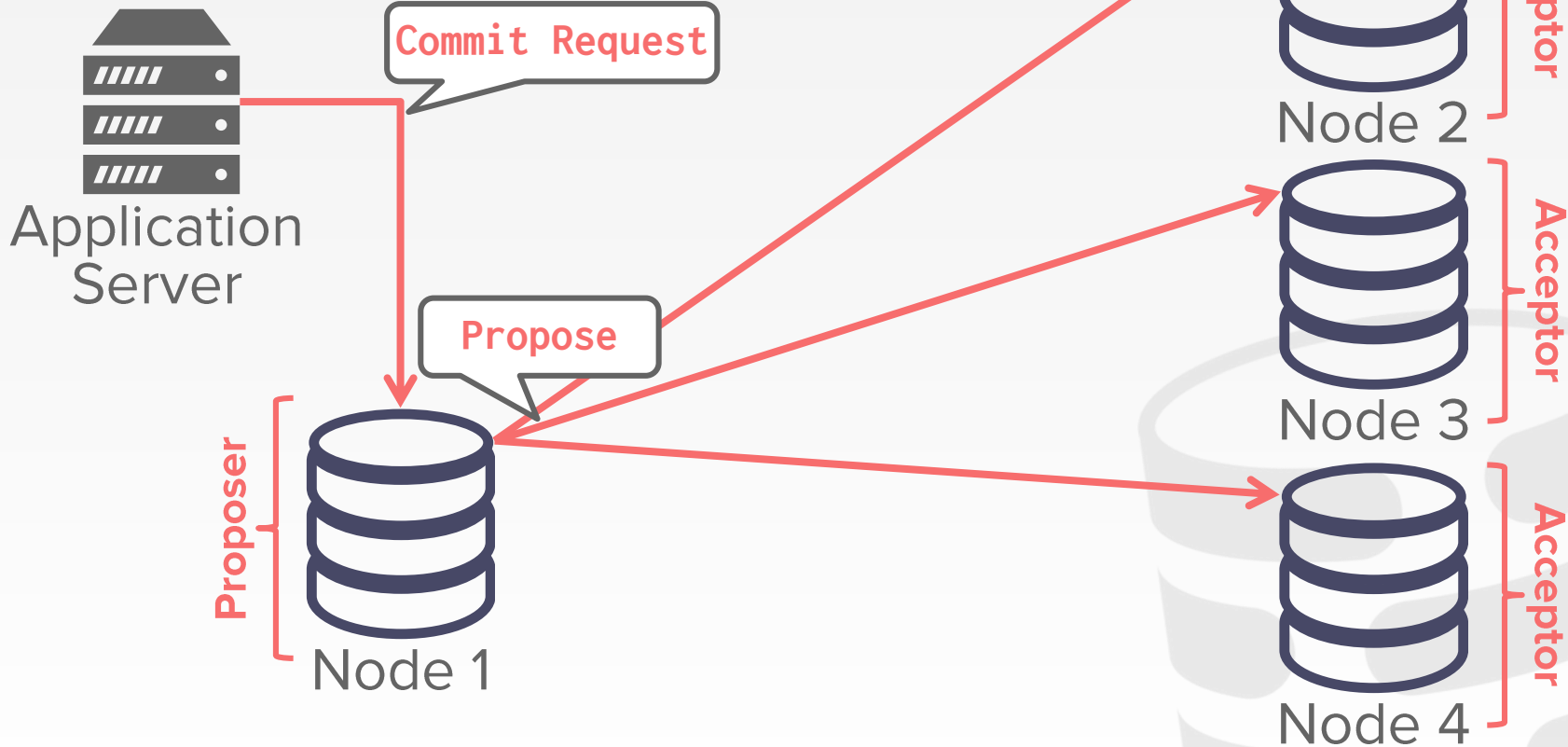


# PAXOS

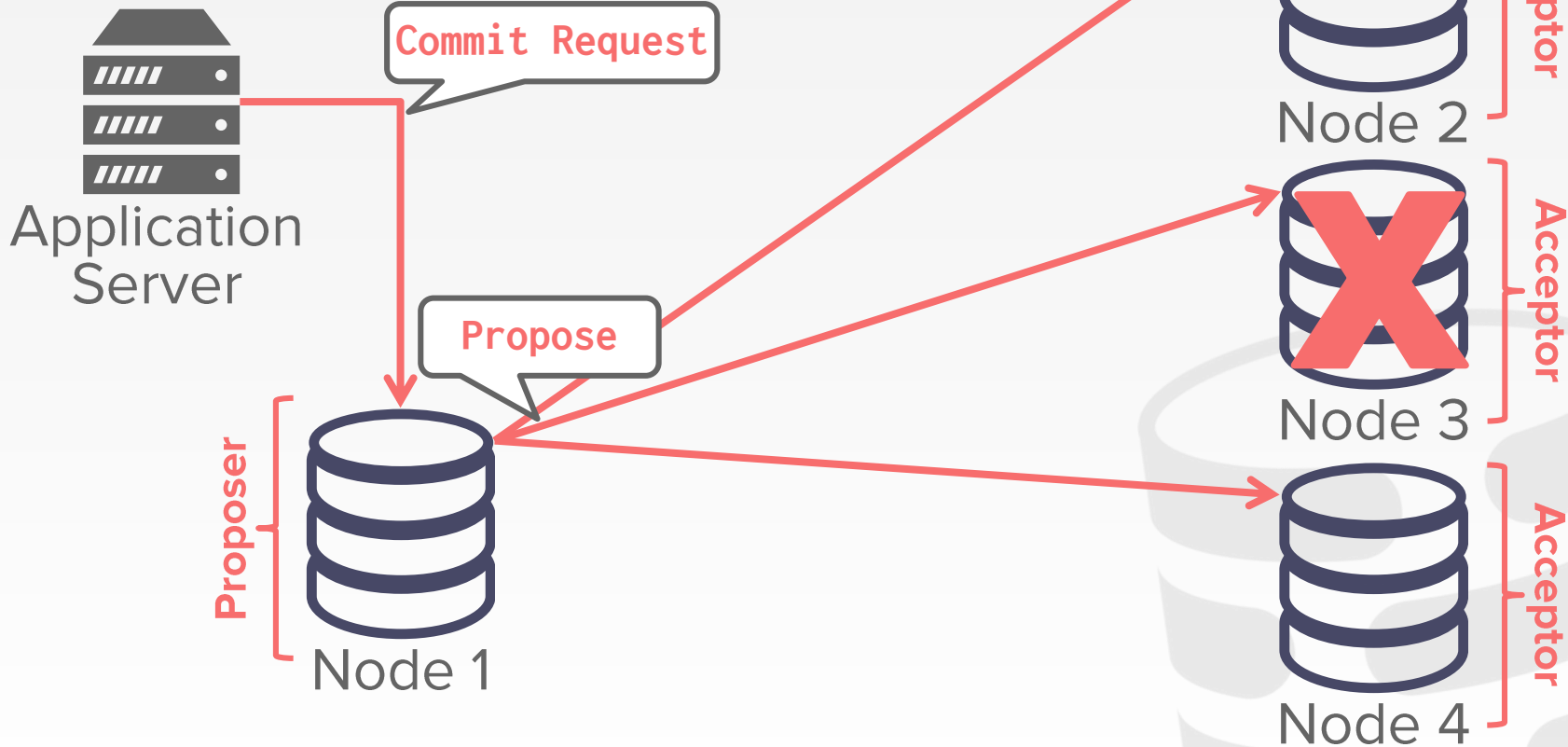




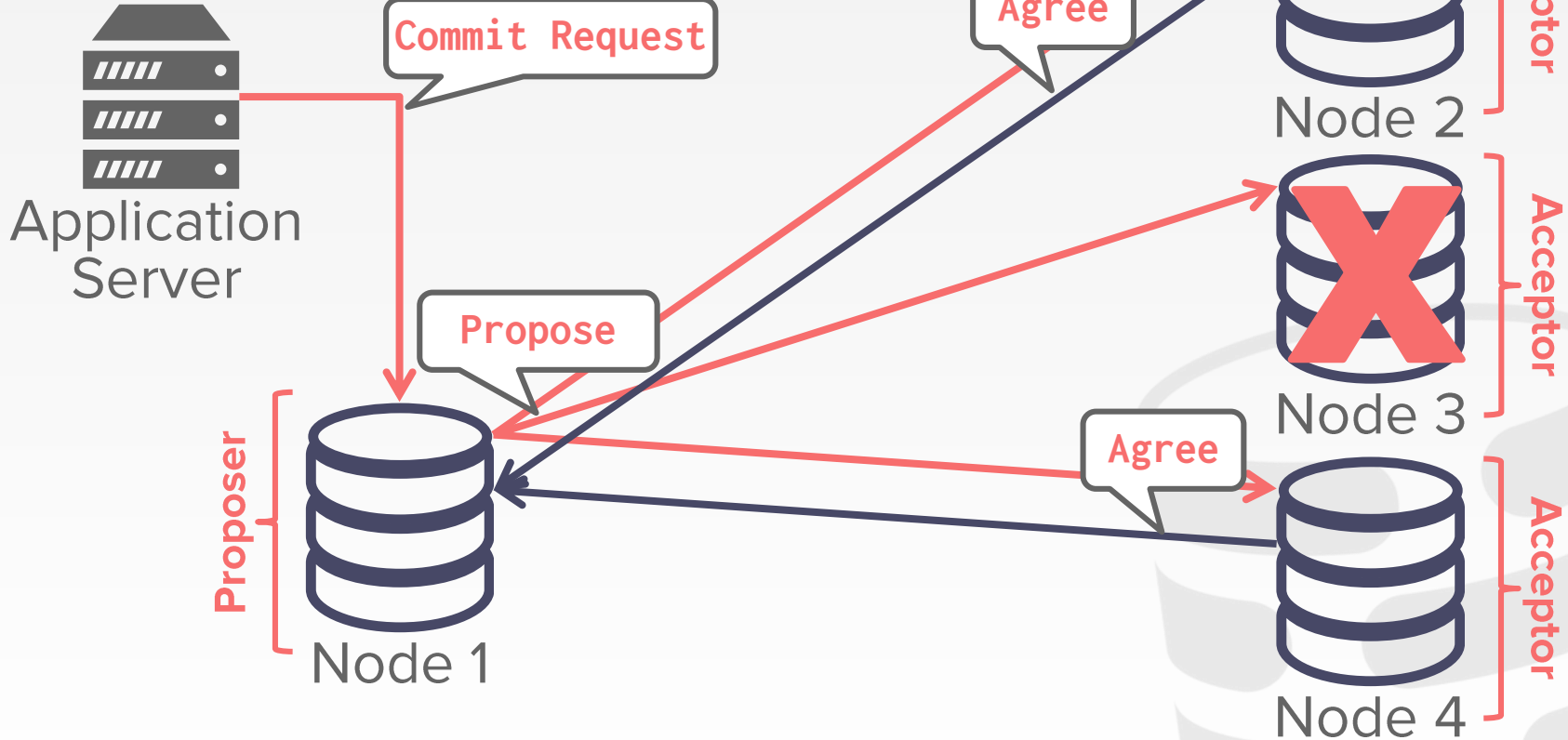
# PAXOS



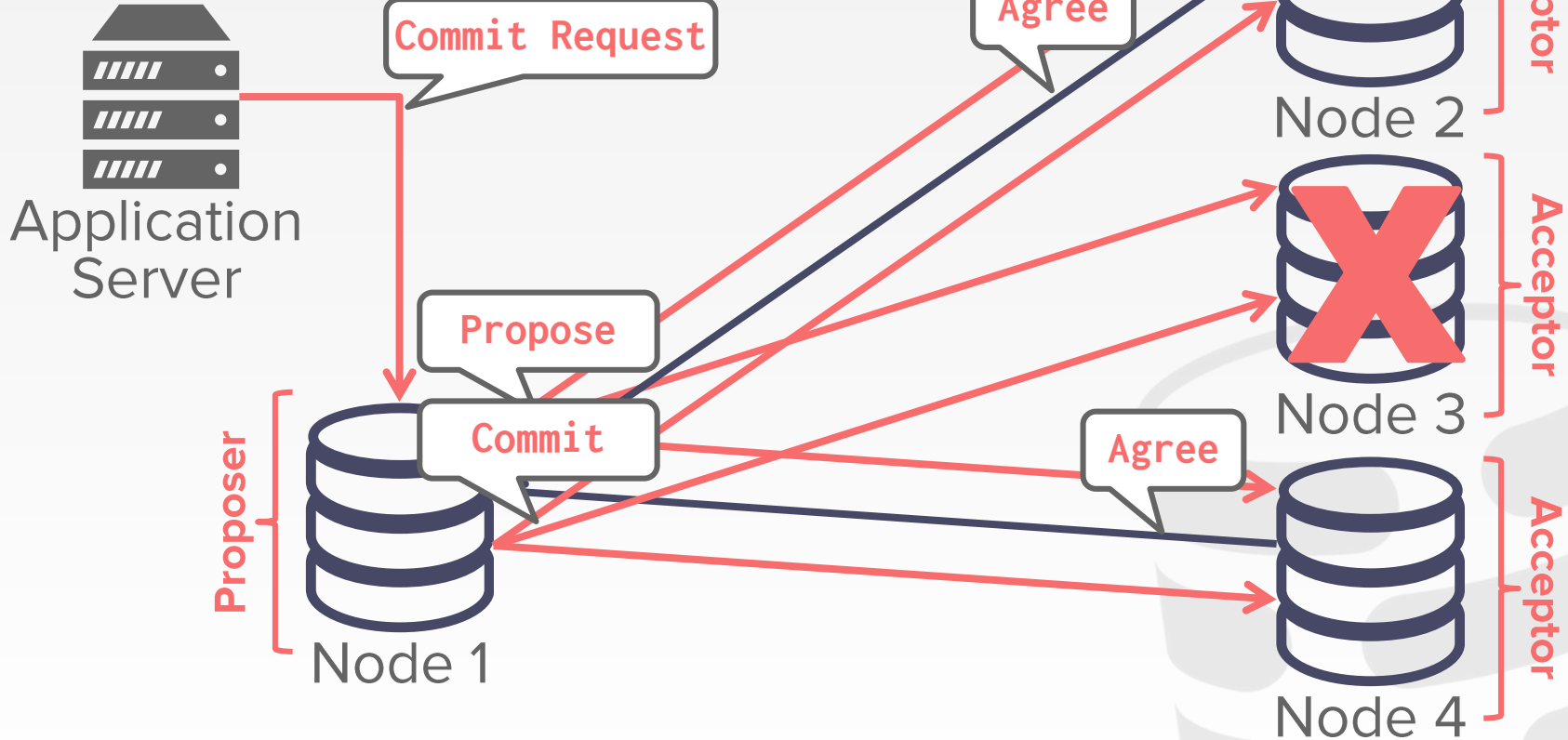
# PAXOS



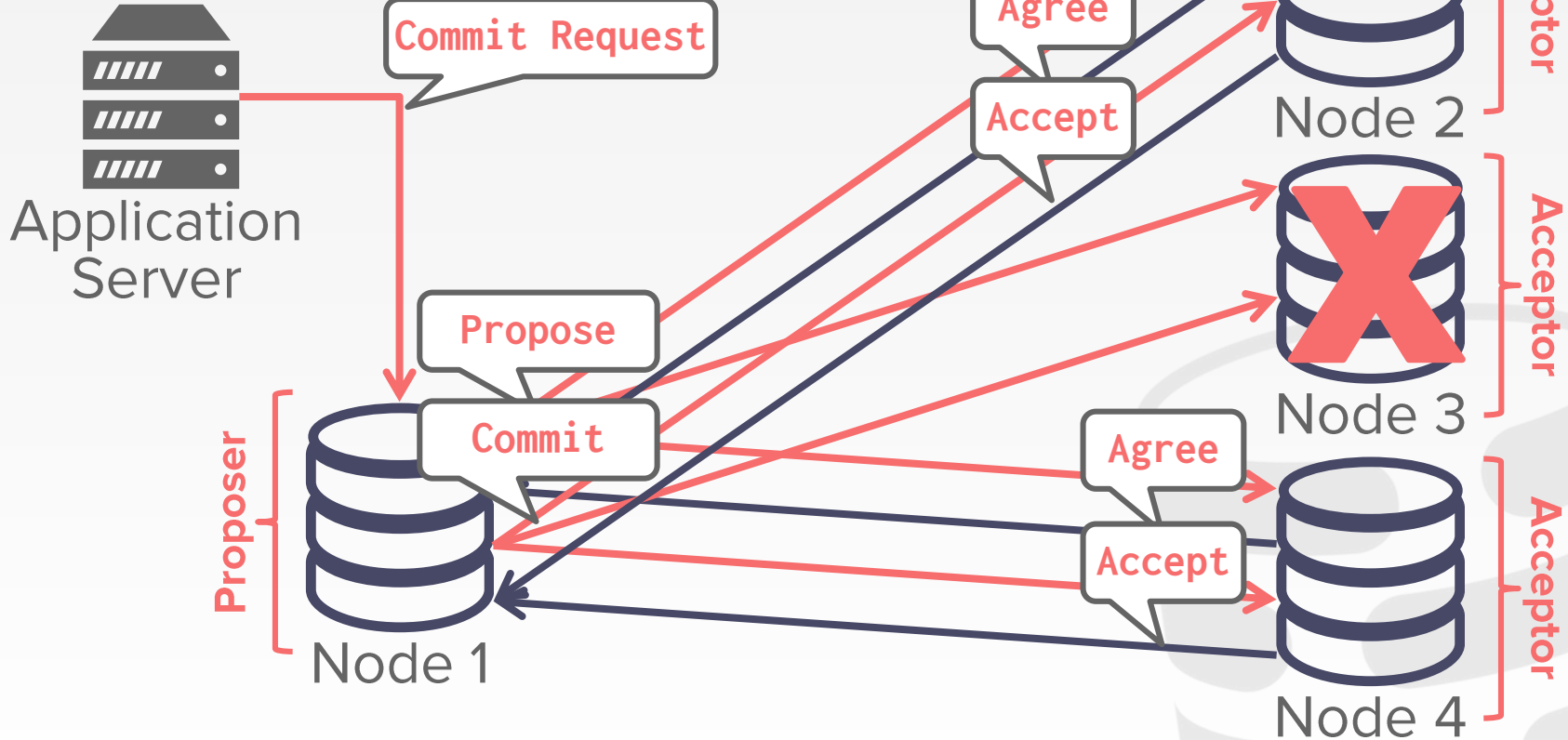
# PAXOS



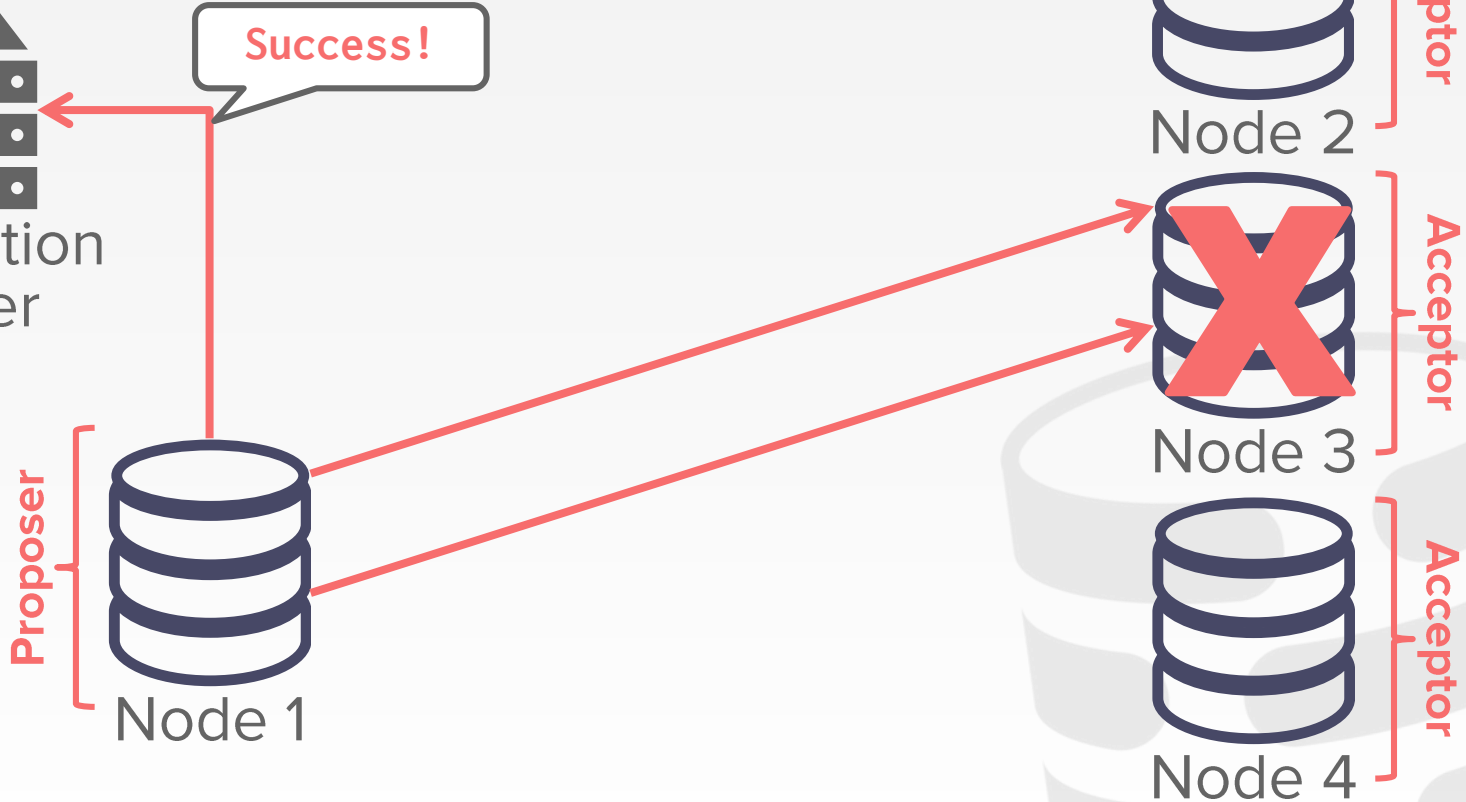
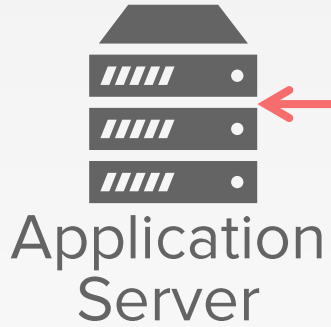
# PAXOS



# PAXOS



# PAXOS





# PAXOS

Proposer



Acceptors

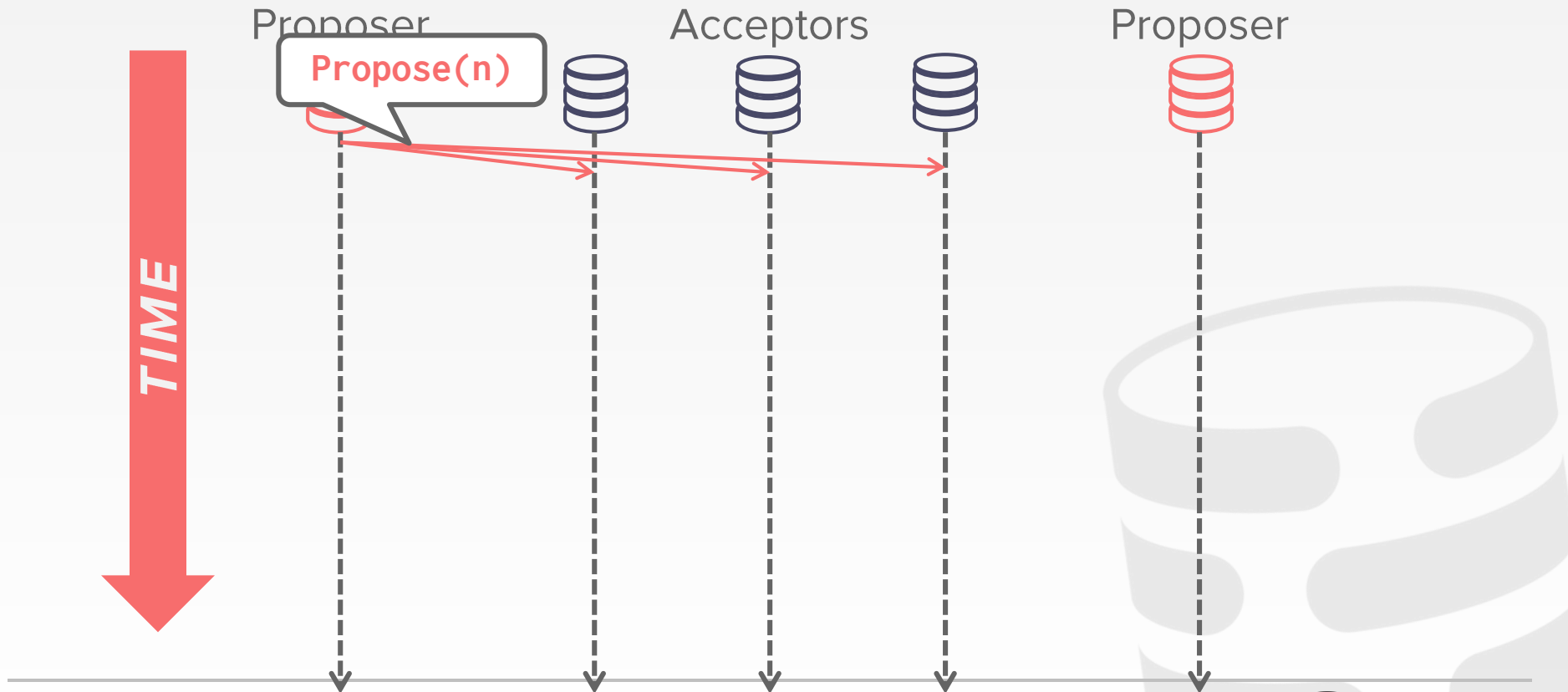


Proposer





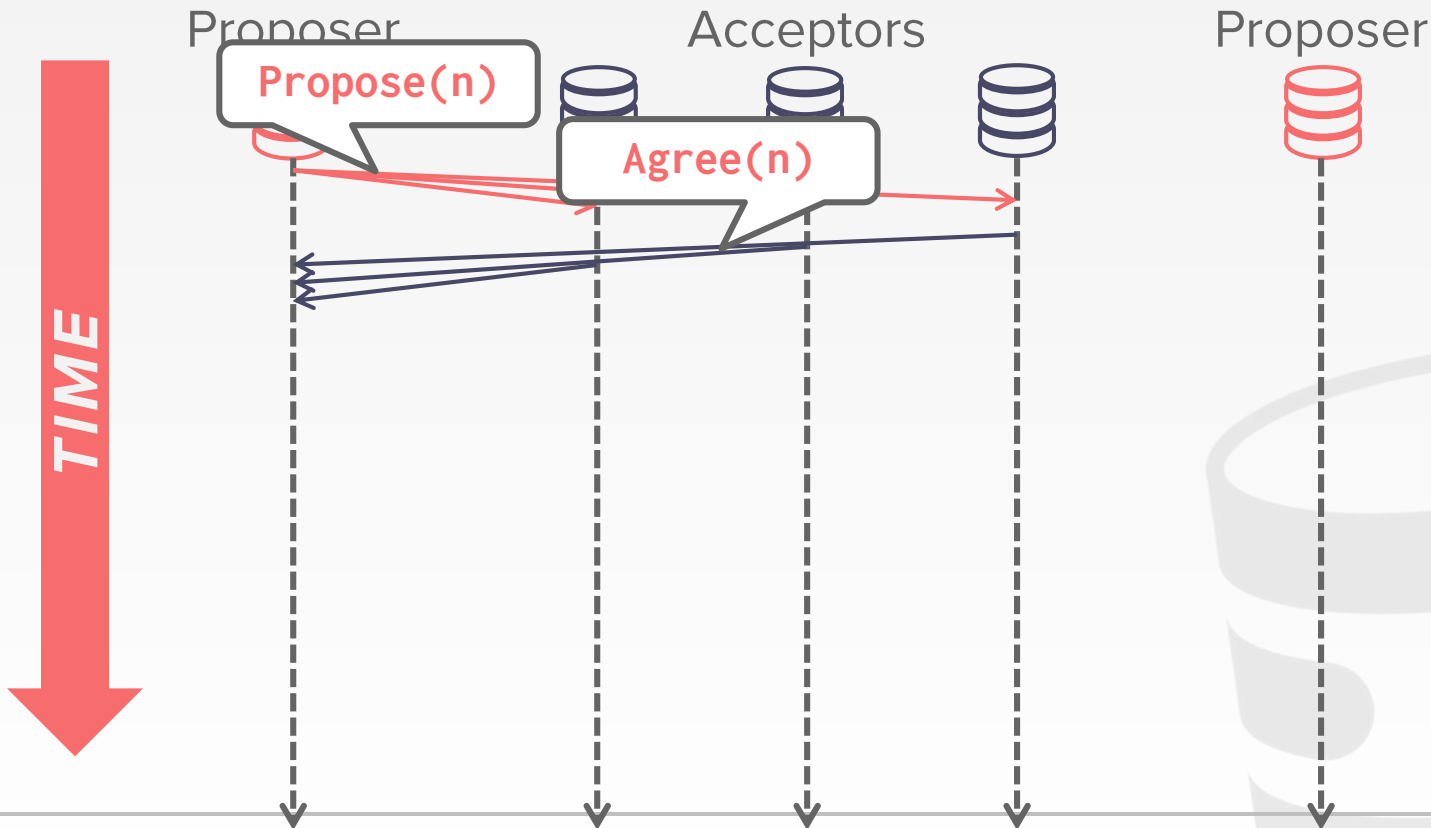
# PAXOS





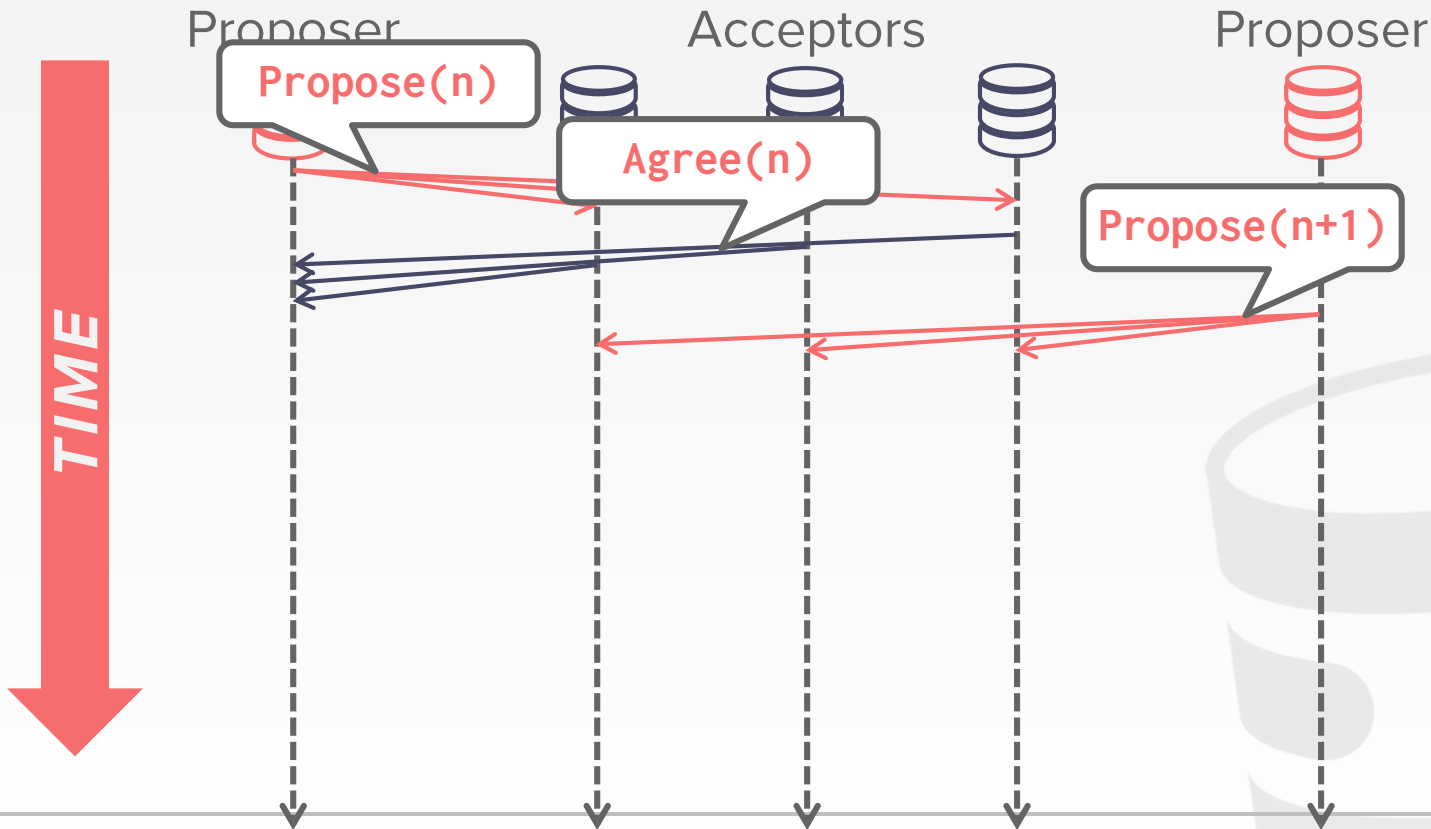


# PAXOS



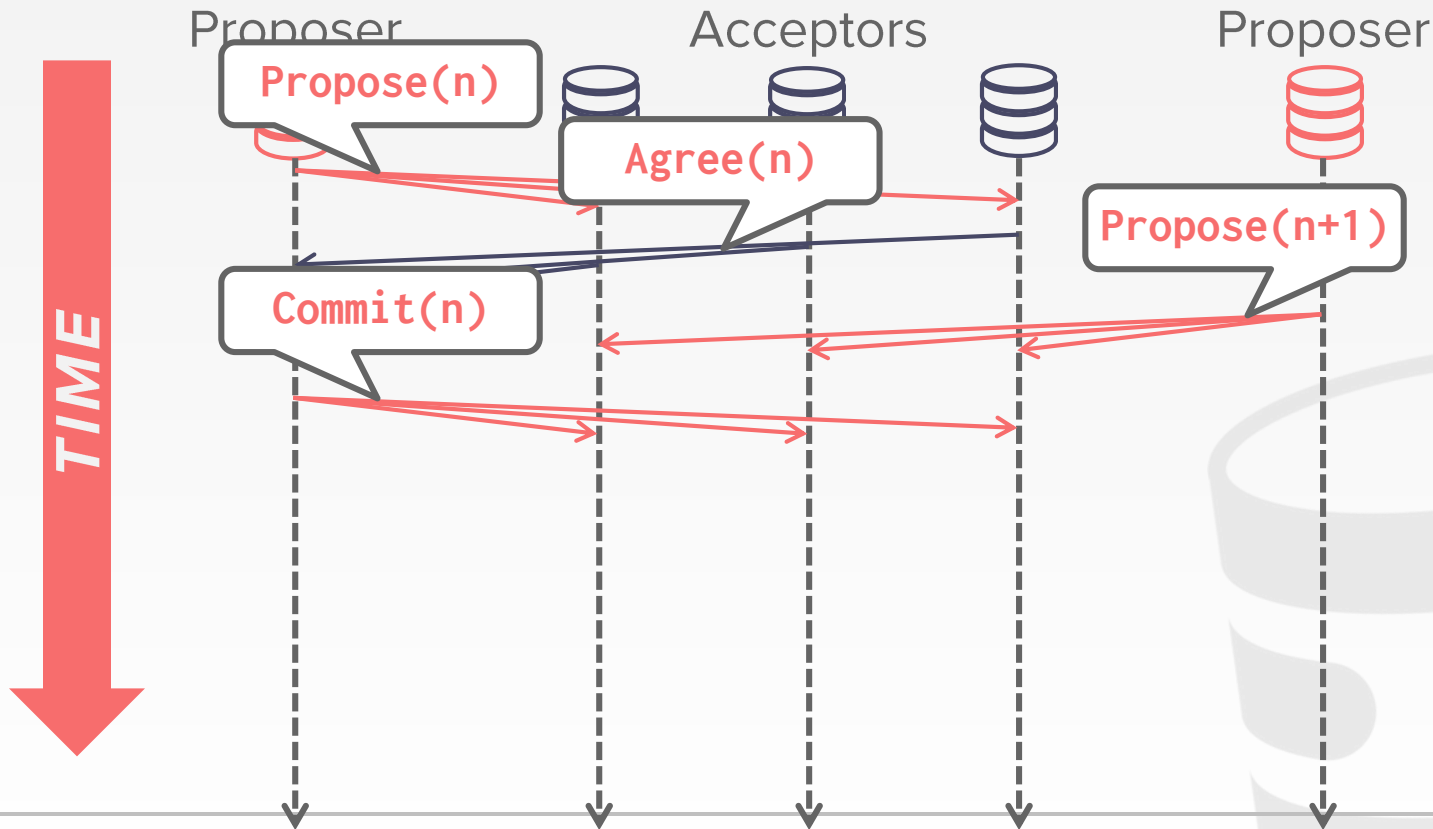


# PAXOS



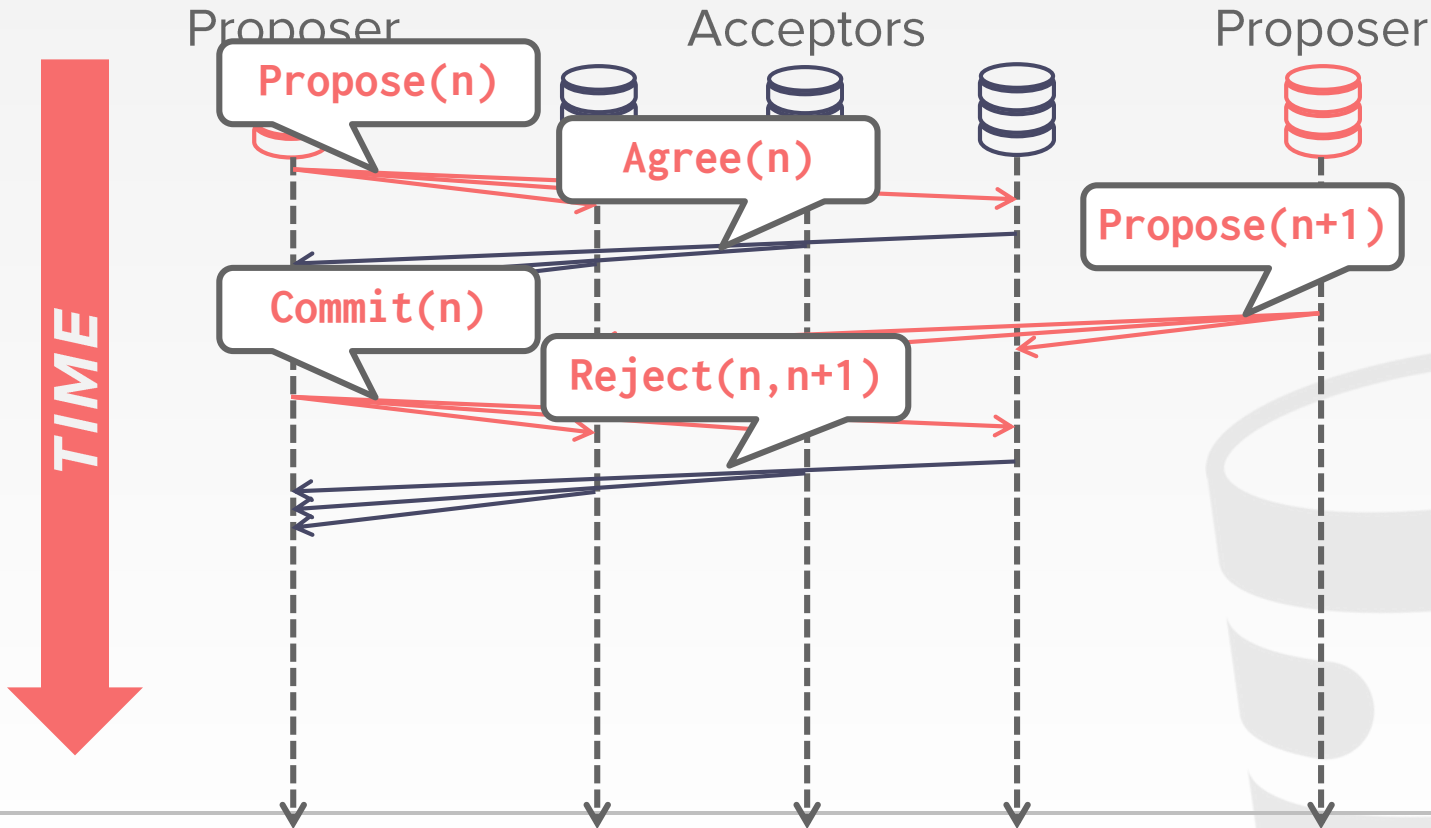


# PAXOS



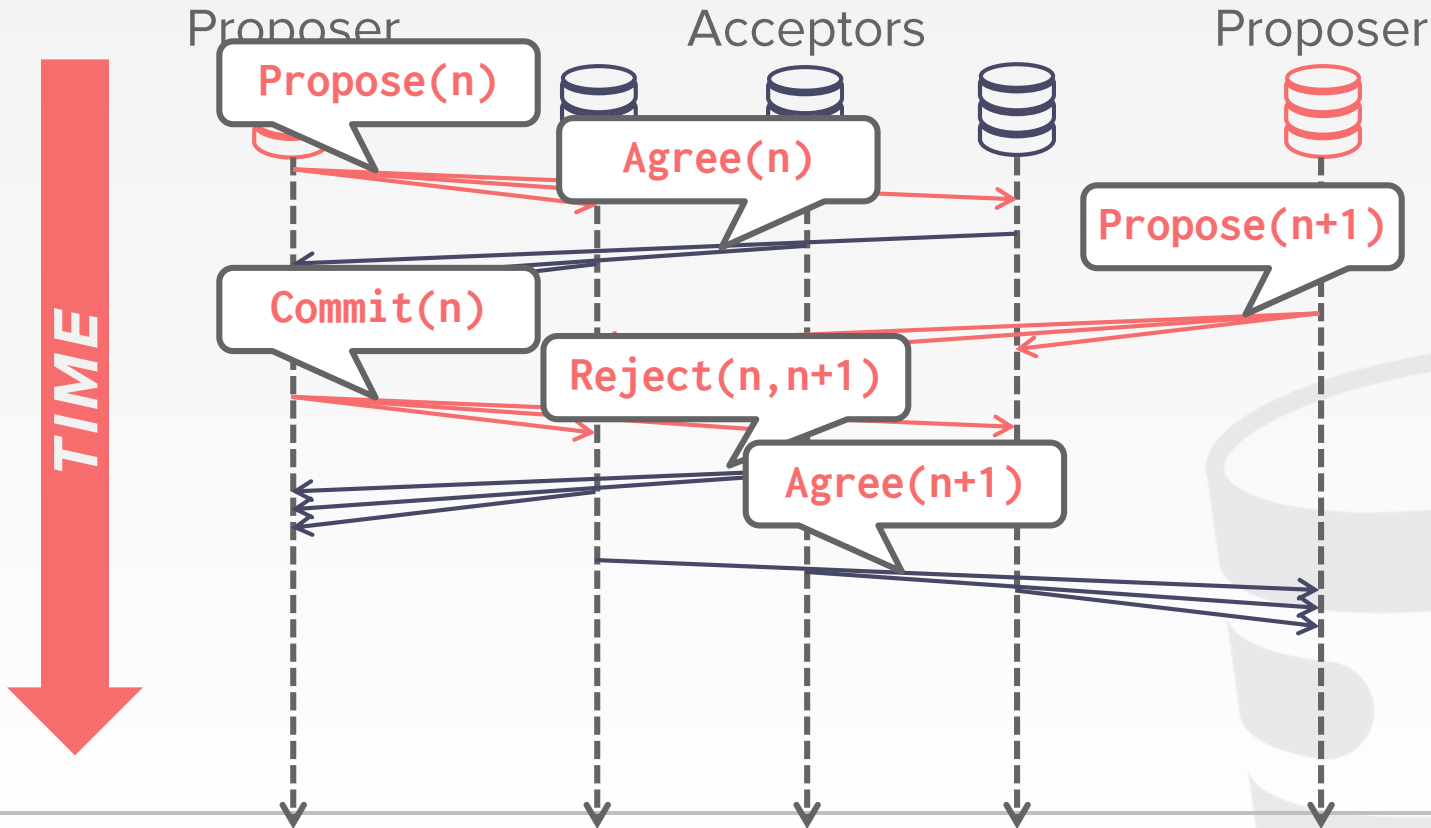


# PAXOS



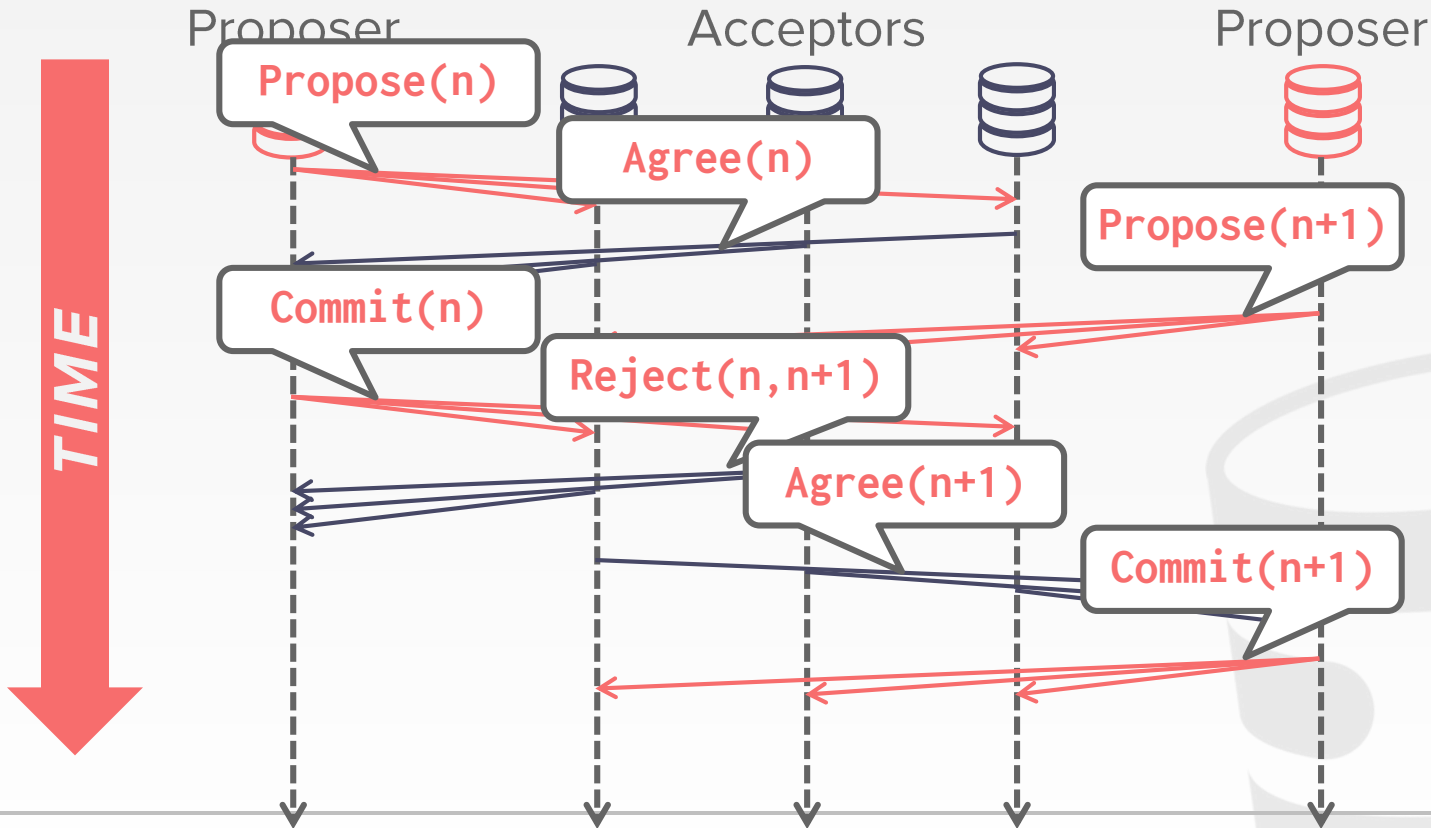


# PAXOS



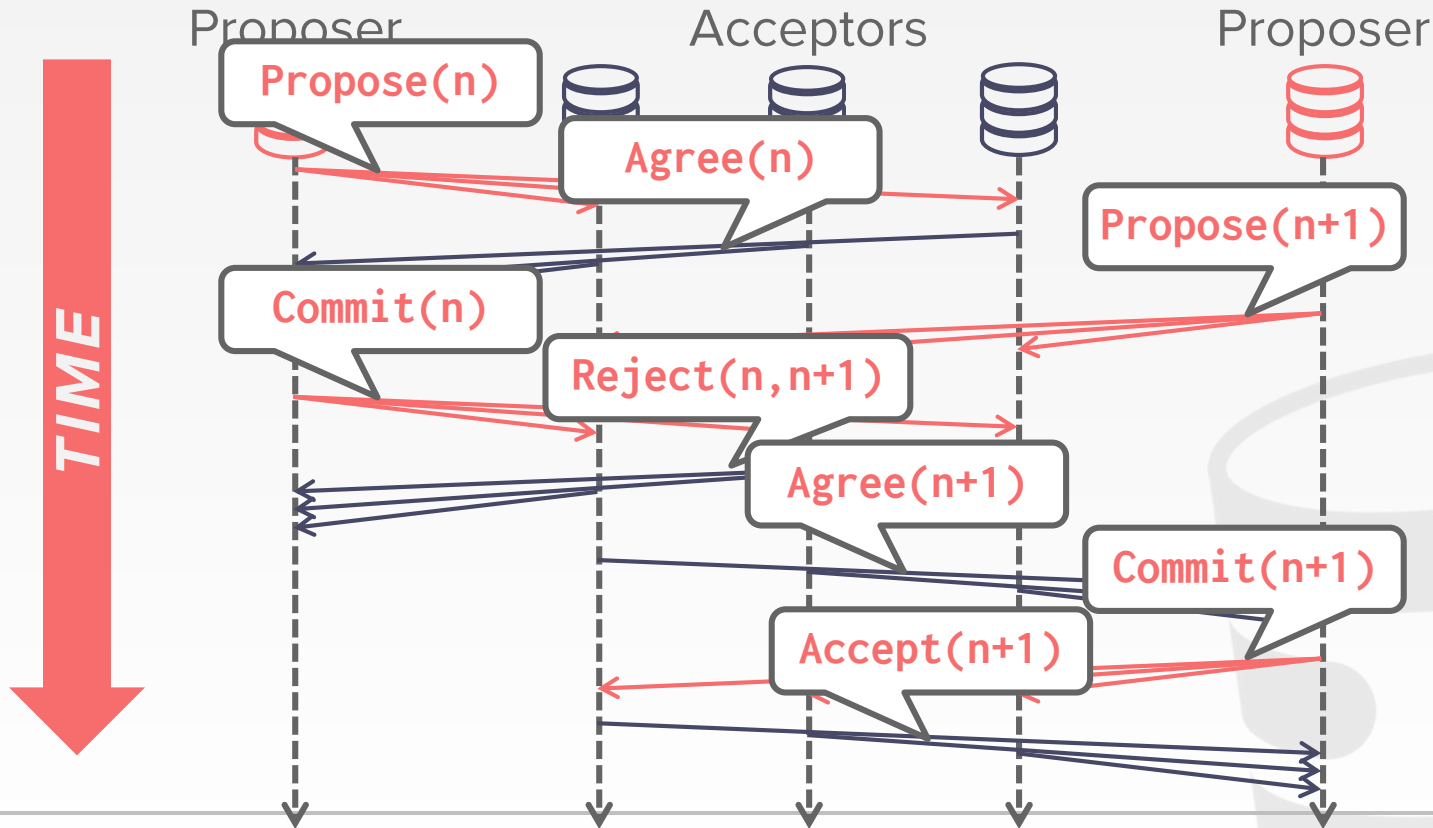


# PAXOS





# PAXOS





# MULTI-PAXOS

If the system elects a single leader that is in charge of proposing changes for some period of time, then it can skip the **PREPARE** phase.

→ Fall back to full Paxos whenever there is a failure.

The system has to periodically renew who the leader is.







# CONCURRENCY CONTROL

MVCC + Strict 2PL with Wound-Wait  
Deadlock Prevention

Ensures ordering through globally unique timestamps generated from atomic clocks and GPS devices.

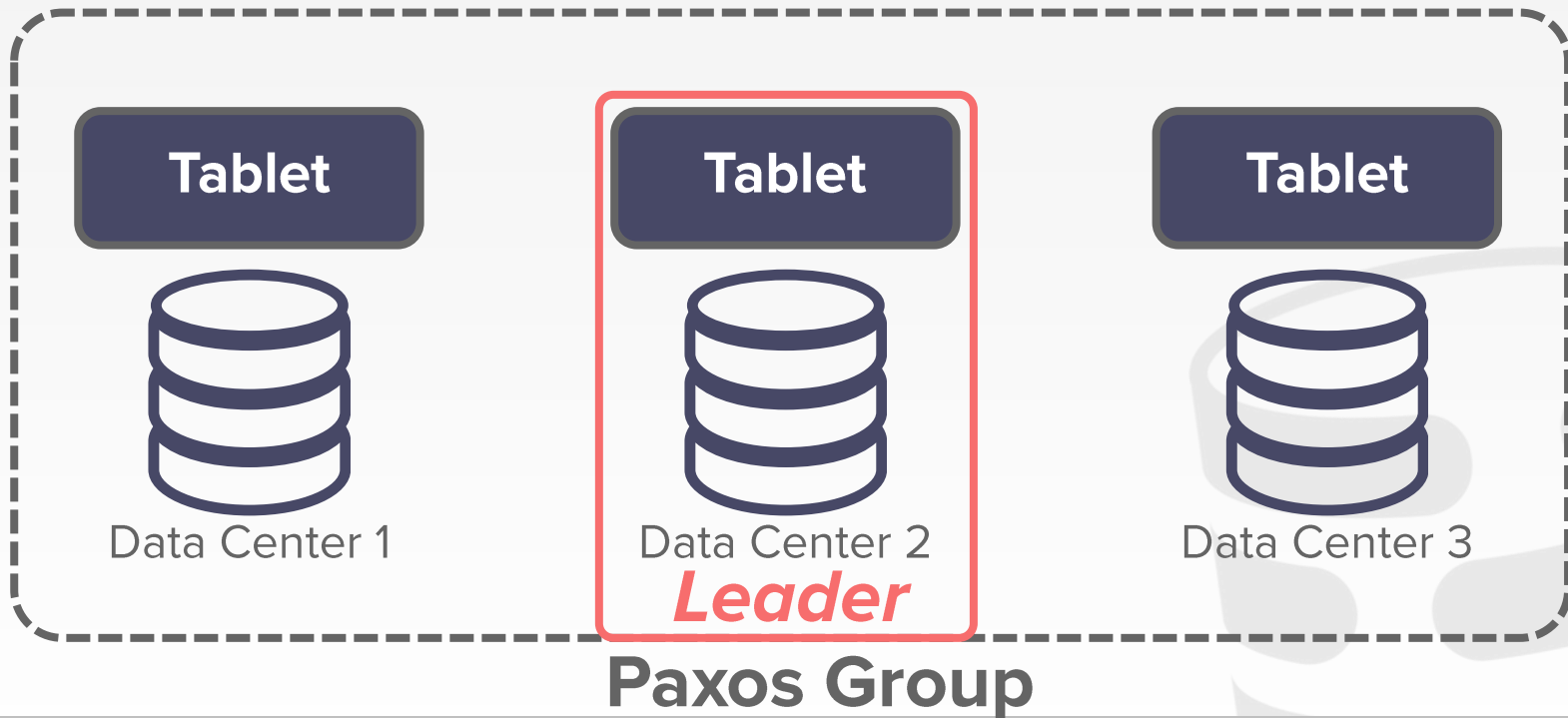
Database is broken up into tablets:

- Use Paxos to elect leader in tablet group.
- Use 2PC for txns that span tablets.



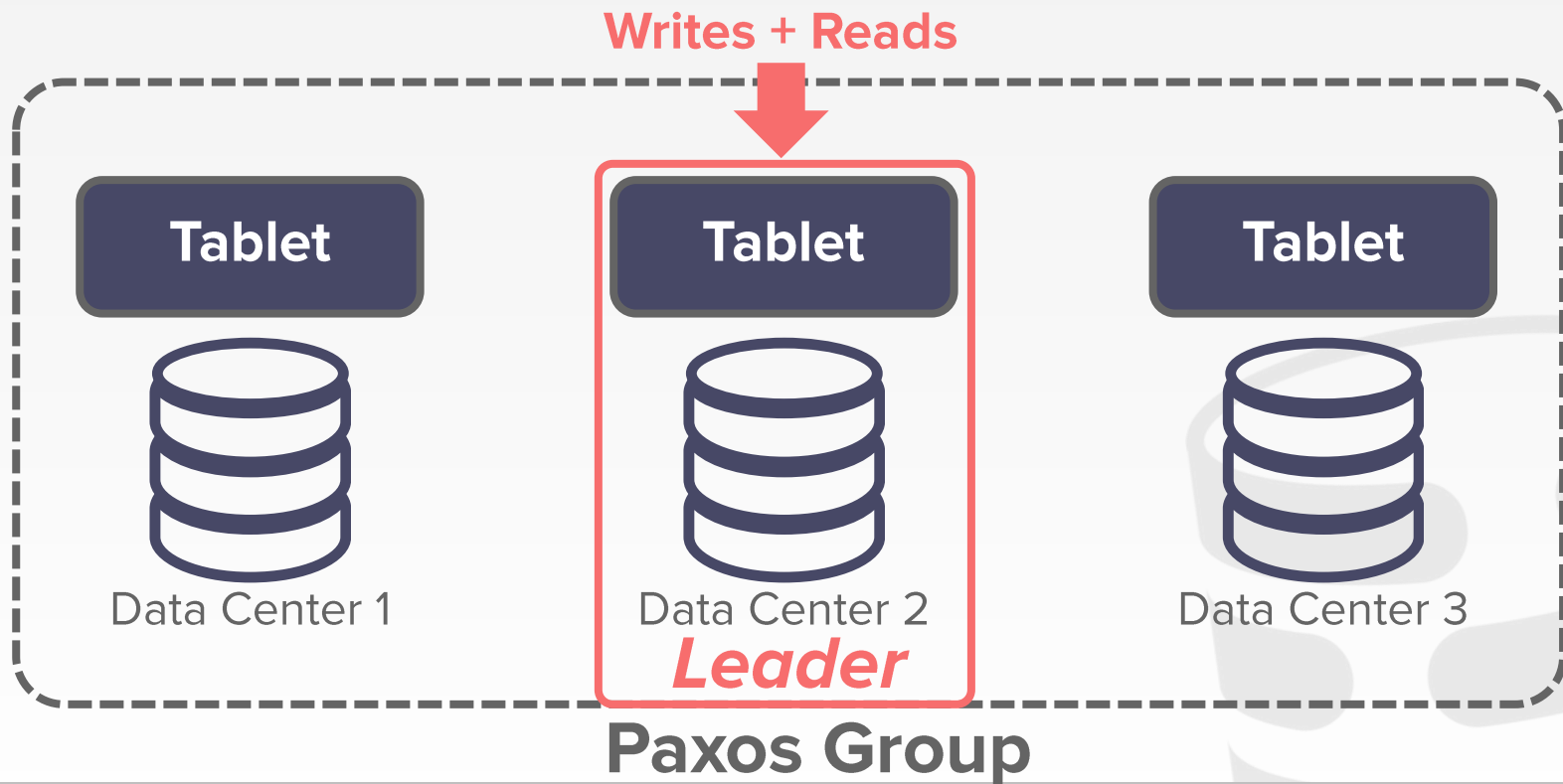


# SPANNER TABLETS



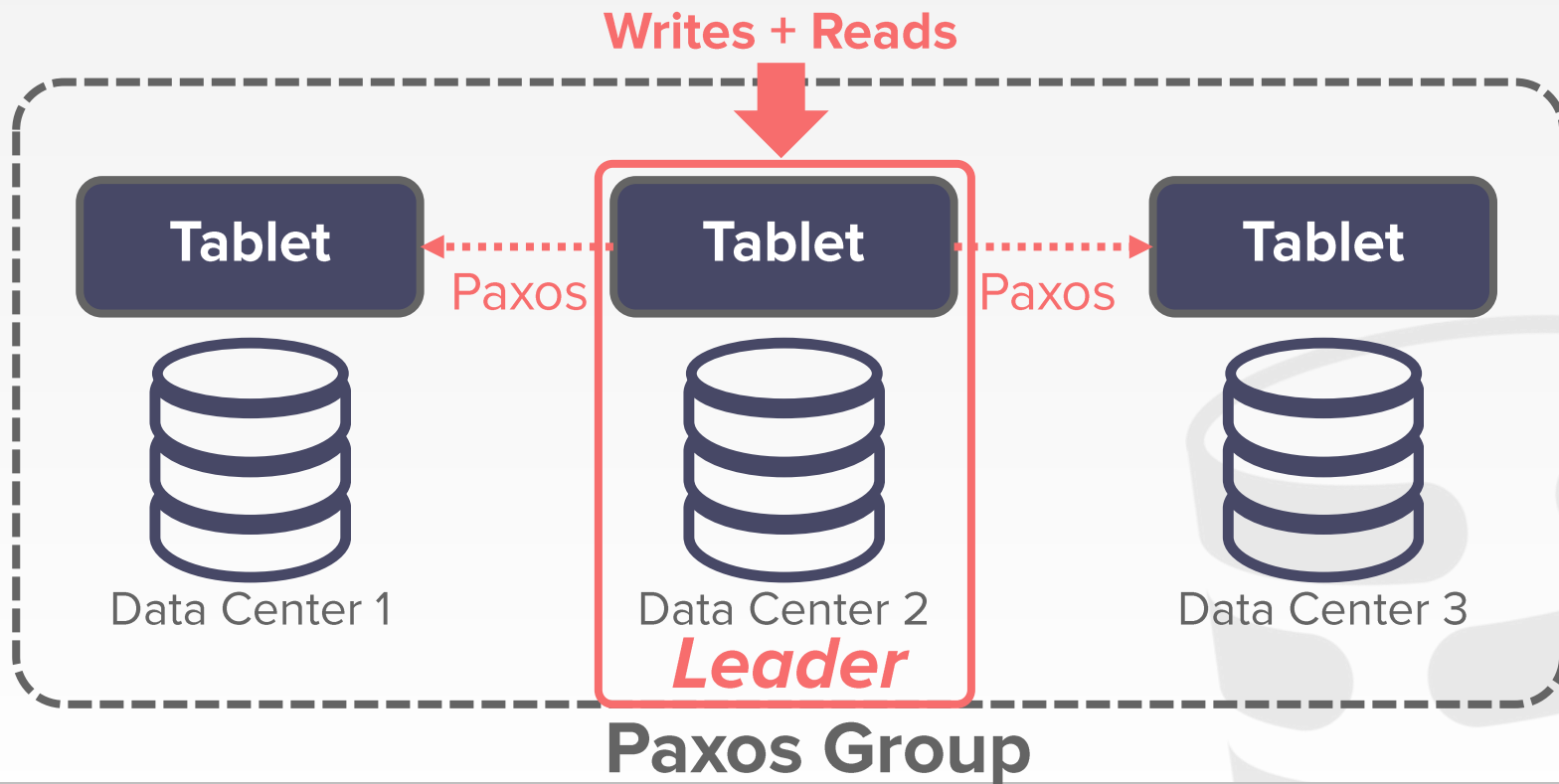


# SPANNER TABLETS



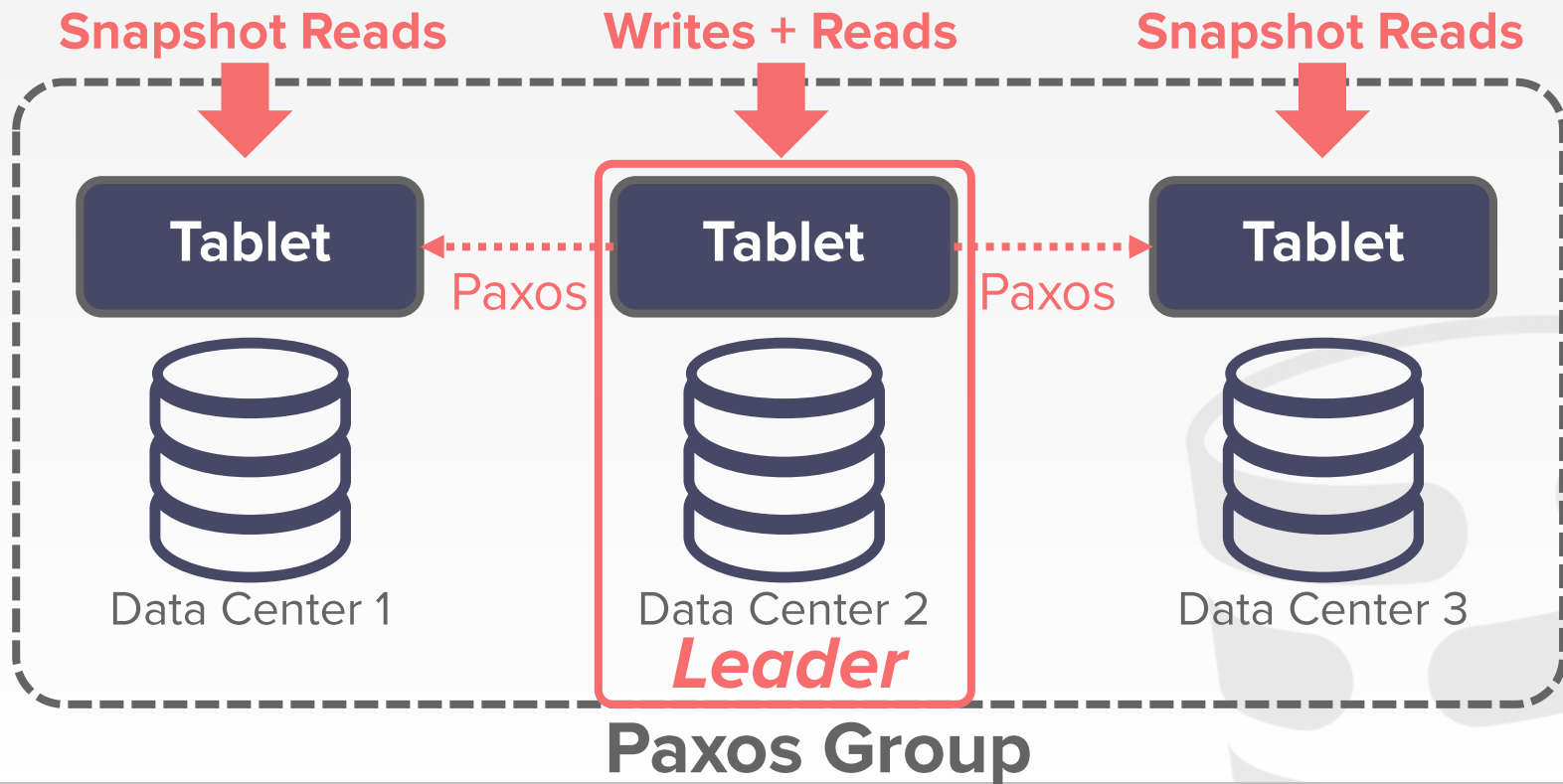


# SPANNER TABLETS



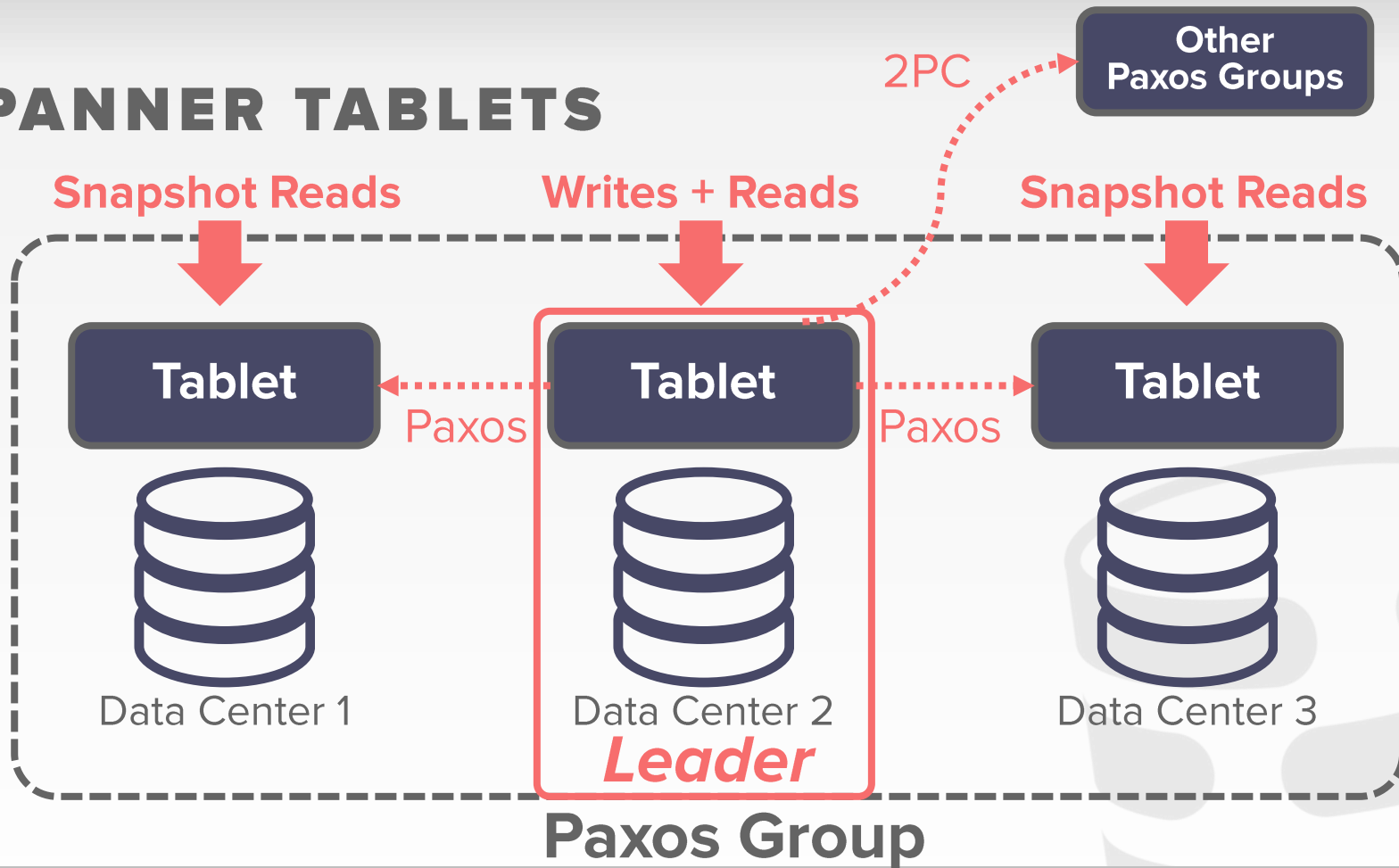


# SPANNER TABLETS





# SPANNER TABLETS



# TRANSACTION ORDERING

Spanner orders transactions based on physical "wall-clock" time.

- This is necessary to guarantee linearizability.
- If  $T_1$  finishes before  $T_2$ , then  $T_2$  should see the result of  $T_1$ .

Each Paxos group decides in what order transactions should be committed according to the timestamps.

- If  $T_1$  commits at  $time_1$  and  $T_2$  starts at  $time_2$  where  $time_1 < time_2$ , then  $T_1$ 's timestamp should be less than  $T_2$ 's.



# SPANNER TRUETIME

The DBMS maintains a global wall-clock time across all data centers with bounded uncertainty.

Timestamps are intervals, not single values





# SPANNER TRUETIME

The DBMS maintains a global wall-clock time across all data centers with bounded uncertainty.

Timestamps are intervals, not single values



# SPANNER TRUETIME

Each data center has GPS and atomic clocks

- These two provide fine-grained clock synchronization down to a few milliseconds.
- Every 30 seconds, there's maximum 7 ms difference.

Multiple sync daemons per data center

- GPS and atomic clocks can fail in various conditions.
- Sync daemons talk to each other within a data center as well as across data centers.



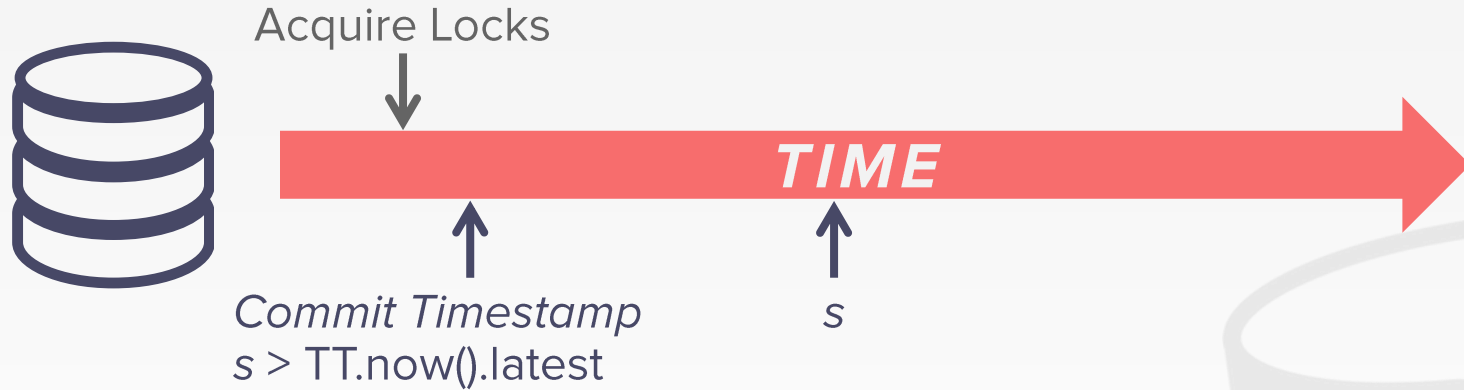
# SPANNER TRUETIME



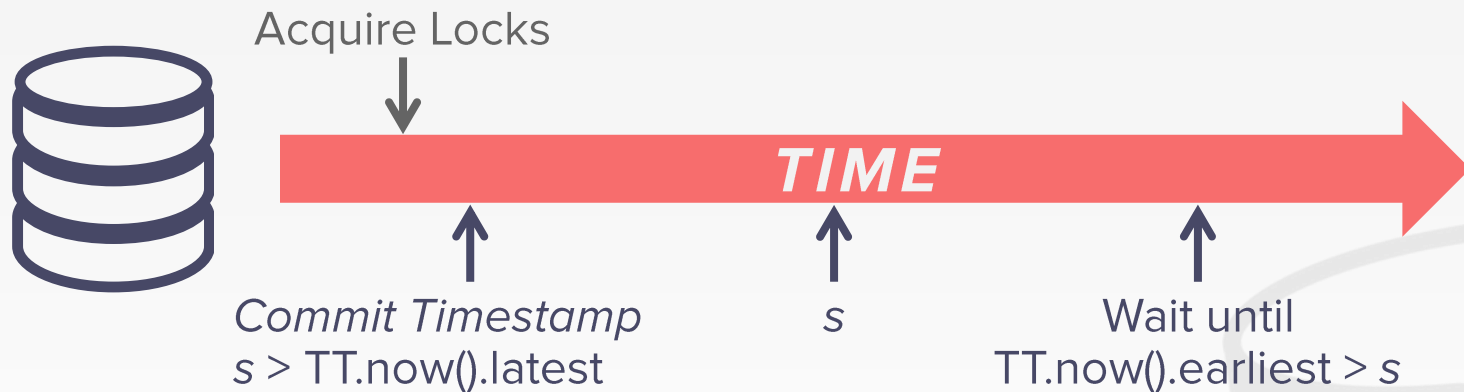
Acquire Locks



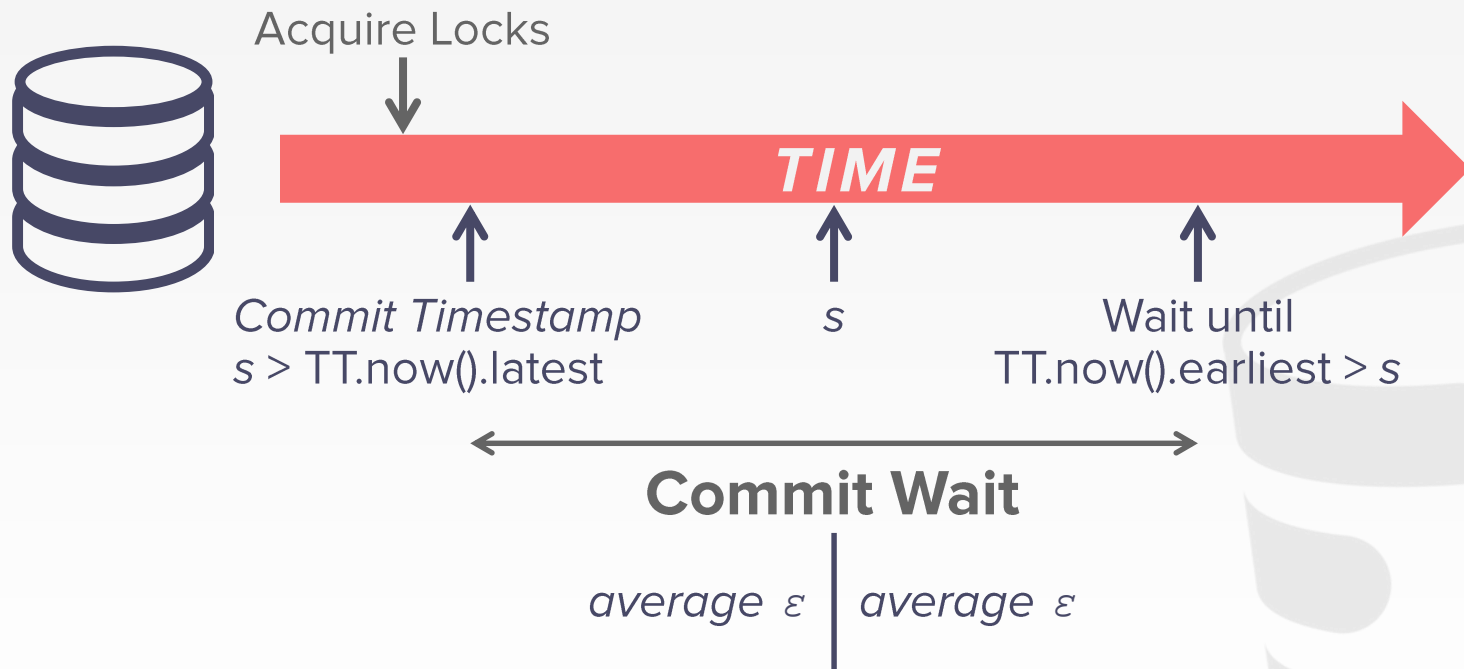
# SPANNER TRUETIME



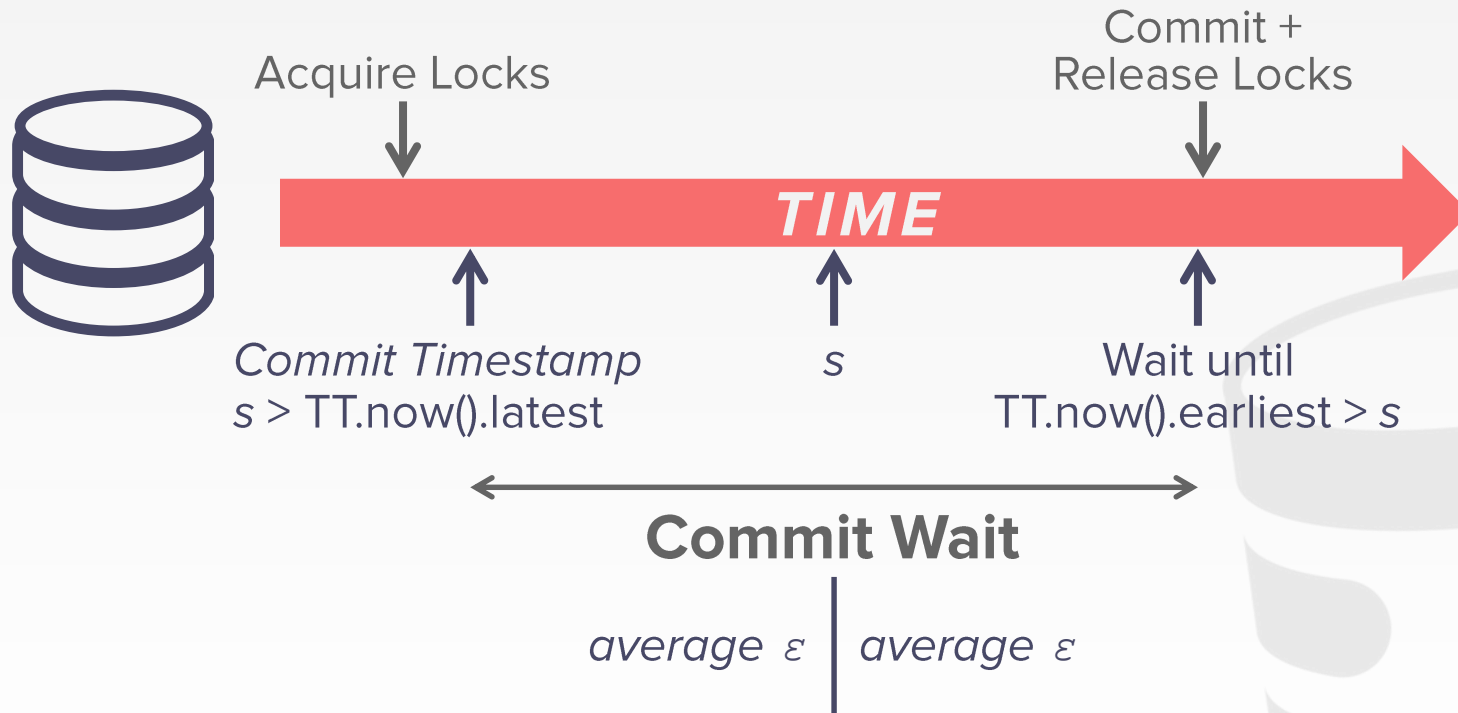
# SPANNER TRUETIME



# SPANNER TRUETIME



# SPANNER TRUETIME



# GOOGLE F1 (2013)

OCC engine built on top of Spanner.

→ In the read phase, F1 returns the last modified timestamp with each row. No locks.

→ The timestamp for a row is stored in a hidden lock column. The client library returns these timestamps to the F1 server.

→ If the timestamps differ from the current timestamps at the time of commit the transaction is aborted.

## F1: A Distributed SQL Database That Scales

Jeff Shute  
Chad Whitley  
David Menestrina

Radek Vingralek  
Eric Rolins  
Stephan Elmer  
Traian Stancescu

Bart Samwel  
Miroslav Cancian  
John Cieslewicz  
Himani Apte

Ben Handy  
Kyle Littenfeld  
Ian Rae

Google, Inc.  
\*University of Wisconsin-Madison

### ABSTRACT

F1 is a distributed relational database system built at Google to support the AdWords business. F1 is a hybrid database that combines high availability, the scalability of NoSQL systems like Bigtable, and the consistency and usability of traditional SQL databases. F1 is built on Spanner, which provides synchronous cross-dataserver replication and strong consistency. Synchronous replication implies higher commit latency, but we mitigate that latency by using a hierarchical schema model with structured data types and through smart application design. F1 also includes a fully functional distributed SQL query engine and automatic change tracking and publishing.

### 1. INTRODUCTION

F1<sup>1</sup> is a fault-tolerant globally-distributed OLTP and OLAP database built at Google as the new storage system for Google's AdWords system. It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements.

The key goals of F1's design are:

- Scalability:** The system must be able to scale up, trivially and transparently, just by adding resources. Our sharded database based on MySQL was hard to scale up, and even more difficult to re-allocate. Our users needed complex queries and joins, which meant they had to carefully shard their data, and re-sharding data without breaking applications was challenging.
- Availability:** The system must never go down for any reason - dataserver outages, planned maintenance, schema changes, etc. The system stores data for Google's core business. Any downtime has a significant revenue impact.
- Consistency:** The system must provide ACID transactions, and must always prevent applications with

<sup>1</sup>Previously described briefly in [2].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were printed to present the results of The 20th International Conference on Very Large Databases, August 26th - 30th 2013, Rio de Janeiro, Brazil.  
Proceedings of the VLDB Endowment, Vol. 6, No. 11  
Copyright 2012 VLDB Endowment 2150-5805/13/06... \$ 10.00.

consistent and correct data.

Designing applications to cope with concurrency anomalies in their data is very error-prone, time-consuming, and ultimately not worth the performance gains.

- Usability:** The system must provide full SQL query support and other functionality users expect from a SQL database. Features like indexes and ad hoc query are not just nice to have, but absolute requirements for our business.

Recent publications have suggested that these design goals are mutually exclusive [5, 11, 23]. A key contribution of this paper is to show how we achieved all of these goals in F1's design, and where we made trade-offs and sacrifices. The name F1 comes from genetics, where a *F1 hybrid* is the first generation offspring resulting from a cross mating of distinctly different parental types. The F1 database system is indeed such a hybrid, combining the best aspects of traditional relational databases and scalable NoSQL systems like Bigtable [6].

F1 is built on top of Spanner [7], which provides extremely scalable data storage, synchronous replication, and strong consistency and ordering properties. F1 inherits those features from Spanner and adds several more:

- Distributed SQL queries, including joining data from external data sources
- Transactionally consistent secondary indexes
- Asynchronous schema changes including database re-organizations
- Optimistic transactions
- Automatic change history recording and publishing

Our design choices in F1 result in higher latency for typical reads and writes. We have developed techniques to hide that increased latency, and we found that user-facing transactions can be made to perform as well as in our previous MySQL system.

- An F1 schema makes data clustering explicit, using tables with hierarchical relationships and columns with structured data types. This clustering improves data locality and reduces the number and cost of RPCs required to read remote data.



# GOOGLE CLOUD SPANNER (2017)

## Spanner Database-as-a-Service.

## AFAIK, it is based on Spanner SQL not F1.

### Spanner: Becoming a SQL System

David F. Bacon    Nathan Bales    Nico Bruno    Brian F. Cooper    Adam Dickinson  
Andrew Flakes    Campbell Fraser    Andrey Gubarev    Milind Joshi    Eugene Kogan  
Alexander Lloyd    Sergey Melnik    Rajesh Rao    David Shue    Christopher Taylor  
                                 Marcel van der    Holst    Dale Woodford

Google, Inc.

#### ABSTRACT

Spanner is a globally-distributed data management system that backs hundreds of mission-critical services at Google. Spanner is built on ideas from both the systems and database communities. The first Spanner paper published at OSDI 12 focused on the system aspects such as scalability, automatic sharding, fault tolerance, consistent replication, external consistency, and wide-area distribution. This paper highlights the database DNA of Spanner. We describe distributed query execution in the presence of rehashing, query restarts upon transient failures, range extraction that drives query routing and index seeks, and the improved blockwise-columnar storage format. We touch upon migrating Spanner to the common SQL dialect shared with other systems at Google.

#### 1. INTRODUCTION

Google's Spanner [1] started out as a key-value store offering multi-row transactions, external consistency, and transparent failover across datacenters. Over the past 7 years it has evolved into a relational database system. In that time we have added a strongly-typed schema system and a SQL query processor, among other features. Initially, some of these database features were "bolted on" – the first version of our query system used high-level APIs almost like an external application, and its design did not leverage many of the unique features of the Spanner storage architecture. However, as we have developed the system, the desire to make it behave more like a traditional database has forced the system to evolve. In particular,

- The architecture of the distributed storage stack has driven fundamental changes in our query compilation and execution, and
- The demands of the query processor have driven fundamental changes in the way we store and manage data.

These changes have allowed us to preserve the massive scalability of Spanner, while offering customers a powerful platform for database applications. We have previously described the distributed architecture and data and concurrency model of Spanner [1]. In

Permission is granted to reproduce or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and full citation on the first page. Copyrights for third party components of this work must be honored. For all other uses, contact the owner/author(s).

SIGMOD 17, May 14–19, 2017, Chicago, IL, USA

© 2017 Copyright held by the owner/authors.

ACM ISBN 978-1-4503-4917-0/17/06

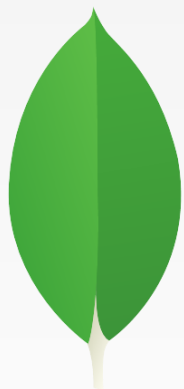
DOI: <http://dx.doi.org/10.1145/3103918.3105610>

this paper, we focus on the "database system" aspects of Spanner, in particular how query execution has evolved and forced the rest of Spanner to evolve. Most of these changes have occurred since [1] was written, and in many ways today's Spanner is very different from what was described there.

A prime motivation for this evolution towards a more "database-like" system was driven by the experiences of Google developers trying to build on previous "key-value" storage systems. The prototypical example of such a key-value system is Bigtable [1], which continues to see massive usage at Google for a variety of applications. However, developers of many OLTP applications found it difficult to build these applications without a strong schema system, cross-row transactions, consistent replication and a powerful query language. The initial response to these difficulties was to build transaction processing systems on top of Bigtable; an example is Megastore [2]. While these systems provided some of the benefits of a database system, they lacked many traditional database features that application developers often rely on. A key example is a robust query language, meaning that developers had to write complex code to process and aggregate the data in their applications. As a result, we decided to turn Spanner into a full featured SQL system, with query execution tightly integrated with the other architectural features of Spanner (such as strong consistency and global replication). Spanner's SQL interface borrows ideas from the F1 [3] system, which was built to manage Google's AdWords data, and included a federated query processor that could access Spanner and other data sources.

Today, Spanner is widely used as an OLTP database management system for structured data at Google, and is publicly available in beta as Cloud Spanner<sup>1</sup> on the Google Cloud Platform (GCP). Currently, over 5,000 databases run in our production instances, and are used by teams across many parts of Google and its parent company Alphabet. This data is the "source of truth" for a variety of mission-critical Google databases, incl. AdWords. One of our large users is the Google Play platform, which executes SQL queries to manage customer purchases and accounts. Spanner serves tens of millions of QPS across all of its databases, managing hundreds of petabytes of data. Replicas of the data are served from datacenters around the world to provide low latency to scattered clients. Despite this wide replication, the system provides transactional consistency and strongly consistent replicas, as well as high availability. The database features of Spanner, operating at this massive scale, make it an attractive platform for new development as well as migration of applications from existing data stores, especially for "big" customers with lots of data and large workloads. Even "small" customers benefit from the robust database features, strong

<sup>1</sup><http://cloud.google.com/spanner>



# mongoDB

# MONGODB

## Document Data Model

- Think JSON, XML, Python dicts
- Not Microsoft Word documents

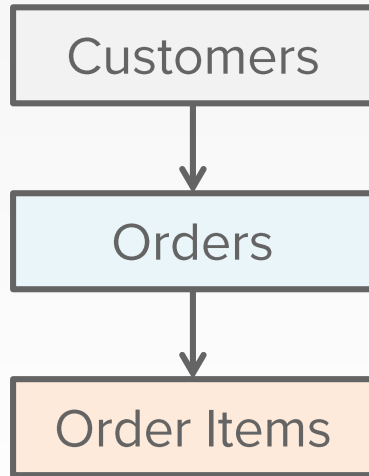
## Different terminology:

- Document → Tuple
- Collection → Table/Relation



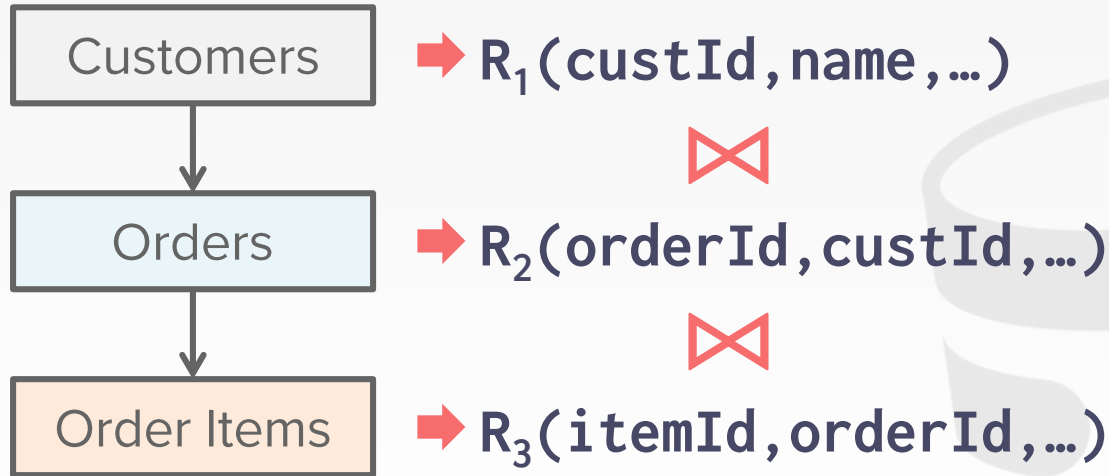
# BCNF EXAMPLE

A customer has orders and each order has order items.



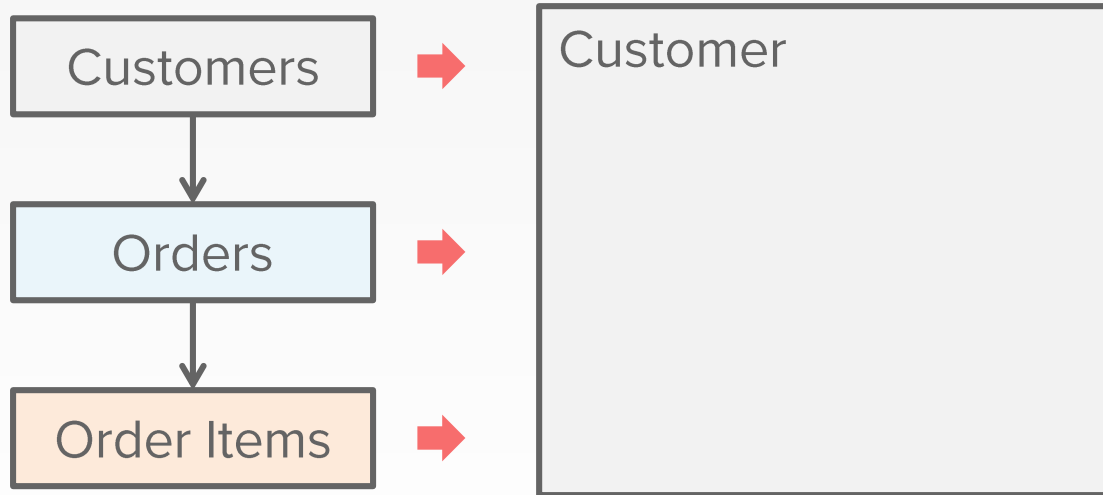
# BCNF EXAMPLE

A customer has orders and each order has order items.



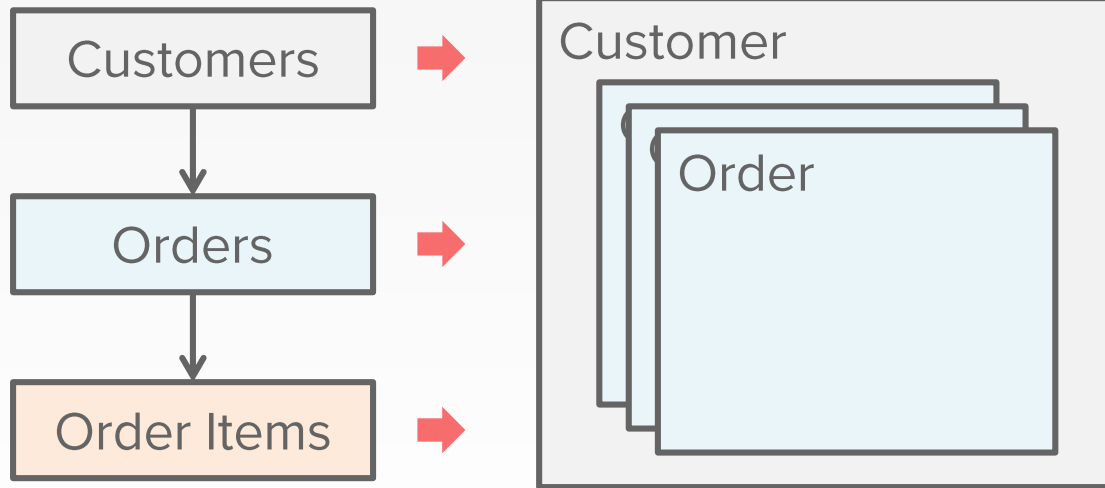
# BCNF EXAMPLE

A customer has orders and each order has order items.



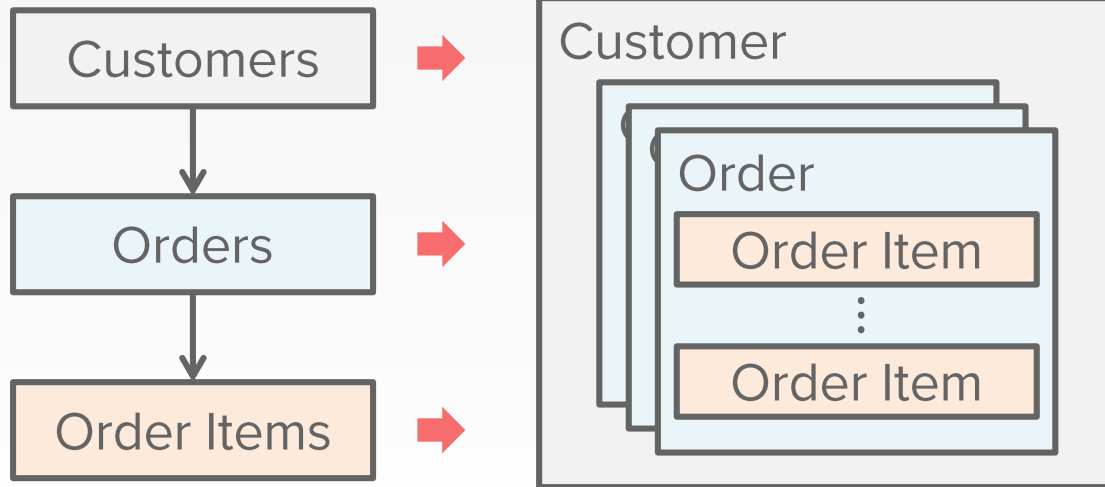
# BCNF EXAMPLE

A customer has orders and each order has order items.



# BCNF EXAMPLE

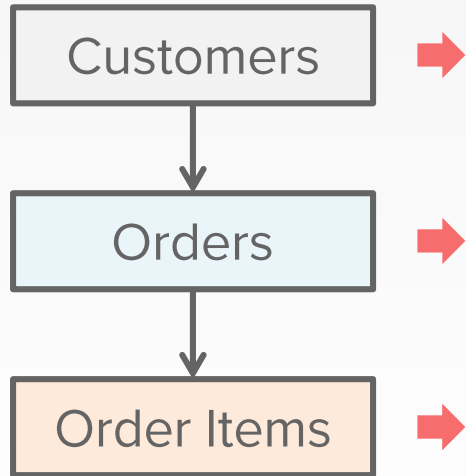
A customer has orders and each order has order items.





# BCNF EXAMPLE

A customer has orders and each order has order items.



```

{
  "custId": 1234,
  "custName": "Andy",
  "orders": [
    { "orderId": 9999,
      "orderItems": [
        { "itemId": "XXXX",
          "price": 19.99 },
        { "itemId": "YYYY",
          "price": 29.99 },
      ] }
  ]
}
  
```

# QUERY EXECUTION

## JSON-only query API

Single-document atomicity.

→ **OLD**: No server-side joins. Had to "pre-join" collections by embedding related documents inside of each other.

→ **NEW**: Server-side joins (only left-outer equi)

No cost-based query planner / optimizer.



# Distributed Architecture

Heterogeneous distributed components.

- Shared nothing architecture
- Centralized query router.

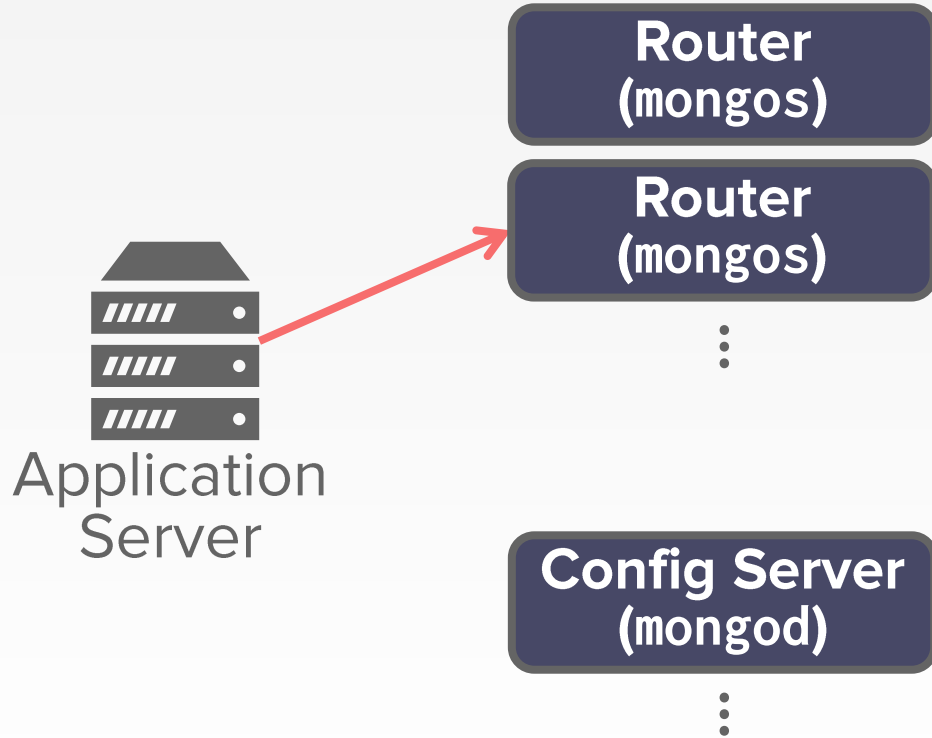
Master-slave replication.

Auto-sharding:

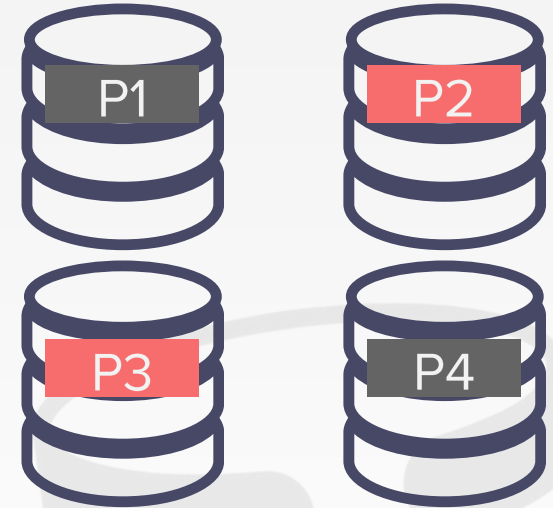
- Define 'partitioning' attributes for each collection (hash or range).
- When a shard gets too big, the DBMS automatically splits the shard and rebalances.



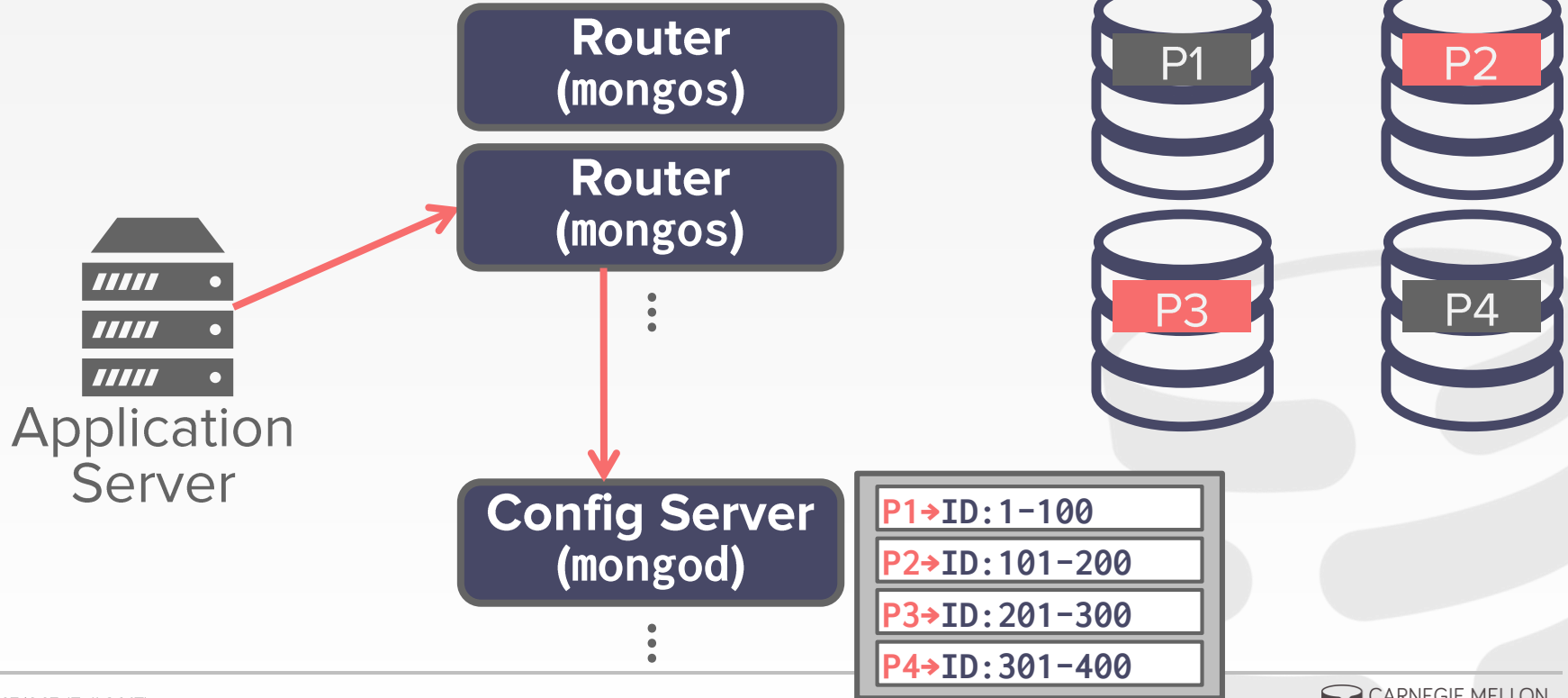
# MONGODB CLUSTER ARCHITECTURE



## Shards (mongod)

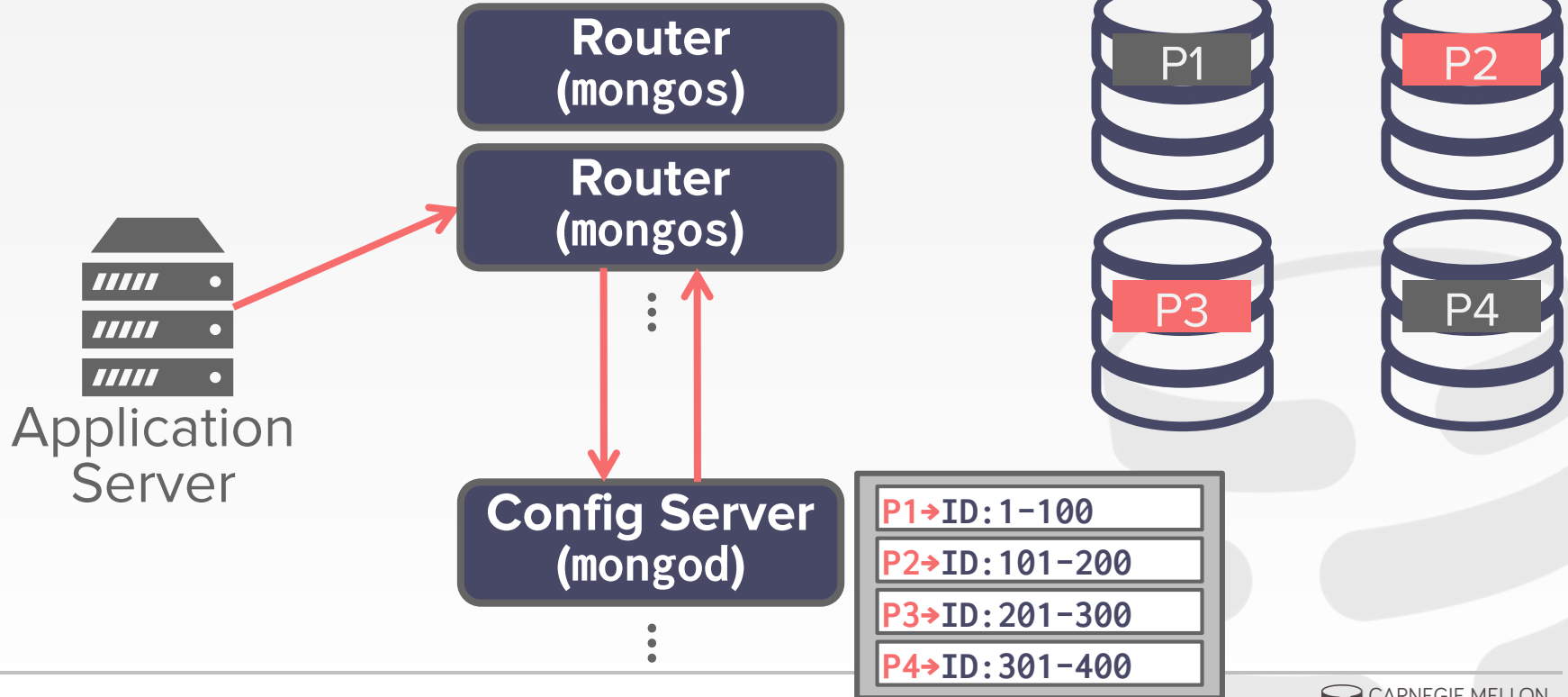


# MONGODB CLUSTER ARCHITECTURE

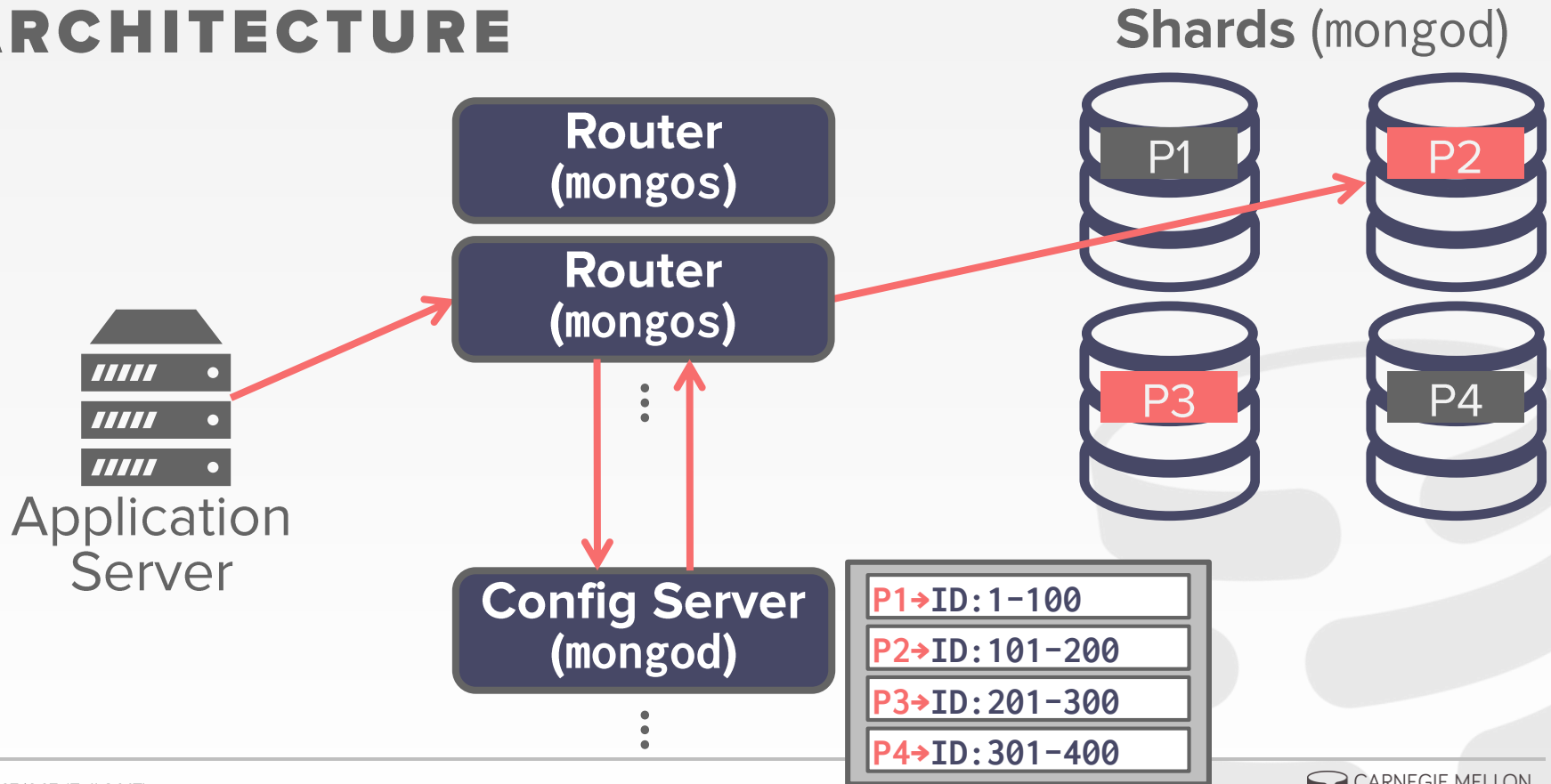


# MONGODB CLUSTER ARCHITECTURE

Shards (mongod)



# MONGODB CLUSTER ARCHITECTURE



# STORAGE ARCHITECTURE

Originally used **mmap** storage manager

- No buffer pool.
- Let the OS decide when to flush pages.
- Single lock per database.

Version 3 (2015) now supports pluggable storage managers.

- **WiredTiger** from BerkeleyDB alumni.  
<http://cmudb.io/lectures2015-wiredtiger>
- **RocksDB** from Facebook (“MongoRocks”)  
<http://cmudb.io/lectures2015-rocksdB>





# ANDY'S CONCLUDING REMARKS

Databases are awesome.

- They cover all facets of computer science.
- We have barely scratched the surface...

Going forth, you should now have a good understanding how these systems work.

This will allow you to make informed decisions throughout your career.

- Both MySQL and Postgres are getting really good...
- Avoid premature optimizations.

