

# backend

cohort #0 by open camp

# #2's agenda

- 1 ..... admin matters (if any)
- 2 ..... databases
- 3 ..... horizontal scaling
- 4 ..... consensus and gossip protocols
- 5 ..... w2 assignment

2.1 admin matters

### **W1 Assignment**

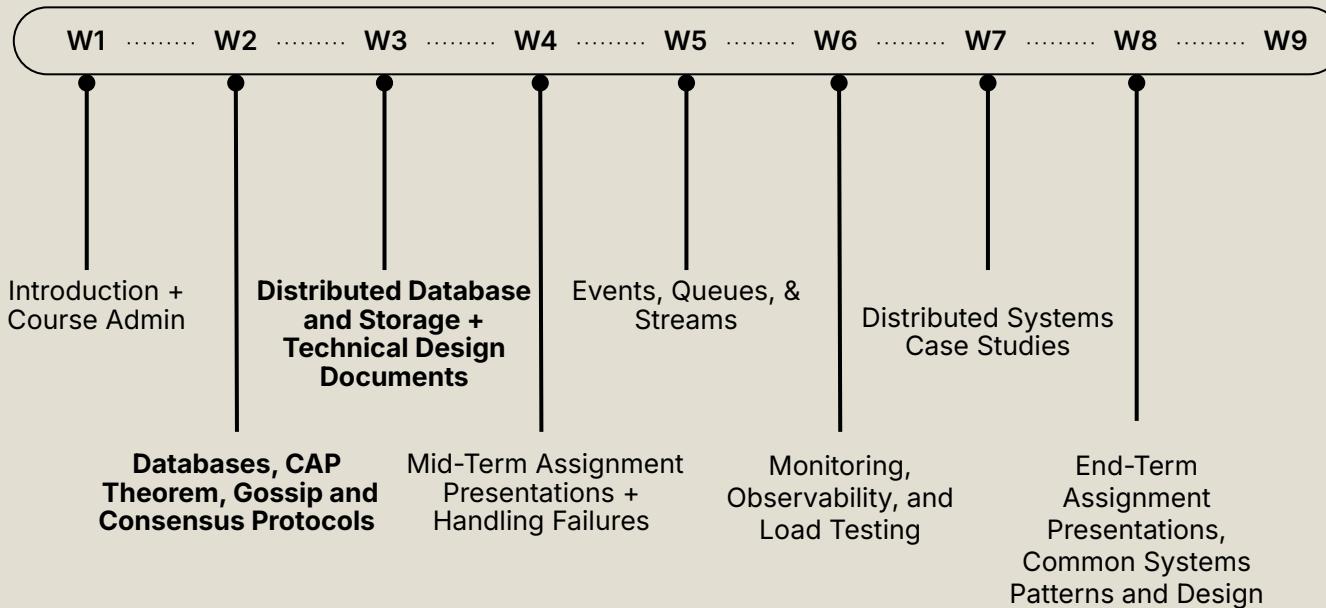
- 2 chose to review AT Protocol
- 2 chose to review ActivityPub
- 2 chose to review and compare both

## 2.2 databases

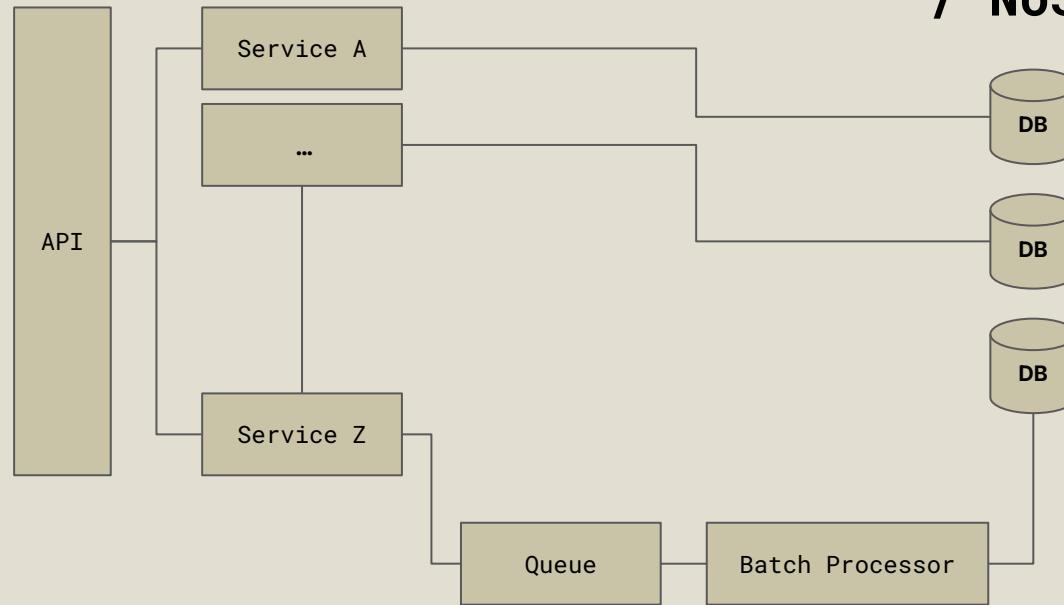
Curriculum: <https://opencamp-cc.github.io/backend-curriculum/>

Start

End



## **Relational / NoSQL**



# **Flow of Data**

# Databases

A traditional single-instance database, by default, runs in a single location on a single machine.

**Assumes 1 machine by default**

# Distributed Databases

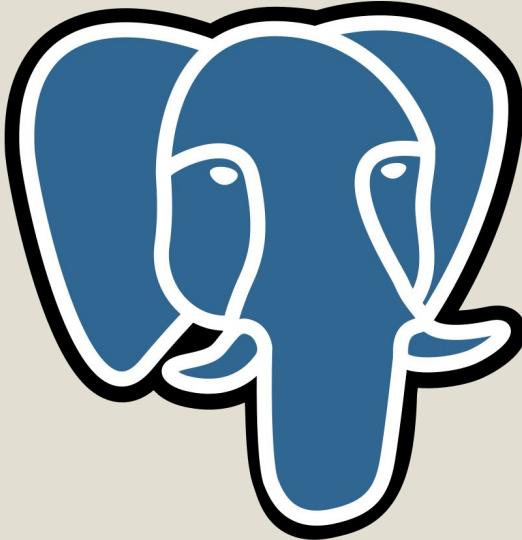
"A distributed database is a database that runs and stores data across multiple computers, as opposed to doing everything on a single machine."

**Assumes 3\* or more machines by default**

**“Which database environments have you done extensive development work in over the past year, and which do you want to work in over the next year?**

> 50%

That's the % of Professional Developers who have used PostgreSQL extensively in the past year and/or want to use it the following year



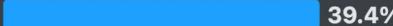
All Respondents

Professional Developers

Learning to Code

Other Coders

PostgreSQL  51.9%

MySQL  39.4%

SQLite  32.1%

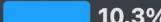
Microsoft SQL Server  27.1%

MongoDB  25.2%

Redis  22.8%

MariaDB  17.1%

Elasticsearch  14.3%

Oracle  10.3%

Dynamodb  9.2%

Firebase Realtime Database  5.6%

Cloud Firestore  5.3%

BigQuery  5%

H2  4.3%

Supabase  3.8%

Cosmos DB  3.8%

Microsoft Access  3.4%

Snowflake  2.8%

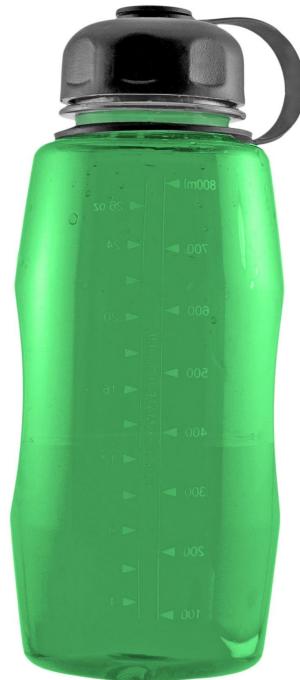
InfluxDB  2.6%

# “Good Database?”



**Insulated**

**Stainless Steel**



**Plastic**

Warm Up

- What are some of the factors you would consider when deciding which database to use or change to?

## People

Existing Skills  
Training / Learning time  
Resistance to Change  
Talent Pool

## Technology

Speed of Development  
Maturity  
Reliability & Consistency  
Documentation  
Security and Updates  
Type of Data (Relational or no)

## Business

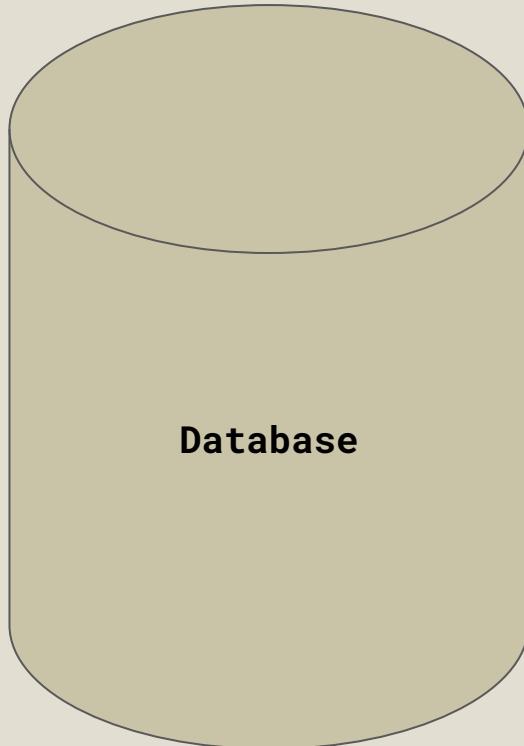
Cost (Buy / Maintenance)  
Legal & Compliance  
Agility

# Considerations

(List is not exhaustive)

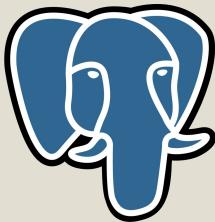
# 2.2 database basics

```
INSERT INTO users  
VALUES ('victor_neo',  
       '...',  
       ...);
```

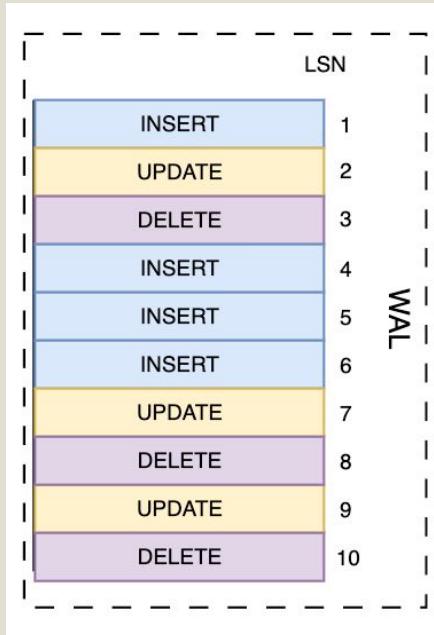


```
INSERT INTO users  
VALUES ('victor_neo', '...', ...);
```

ID	Username	...
1	...	...
2	<b>victor_neo</b>	...



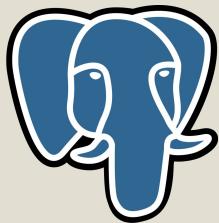
```
INSERT INTO users  
VALUES ('victor_neo', '...', ...);
```



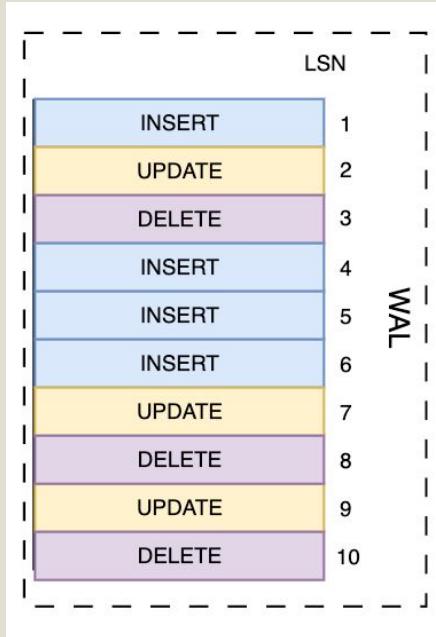
New data does not get written to the table right away!

It's faster to append it to a WAL log, and say that the SQL transaction is complete.

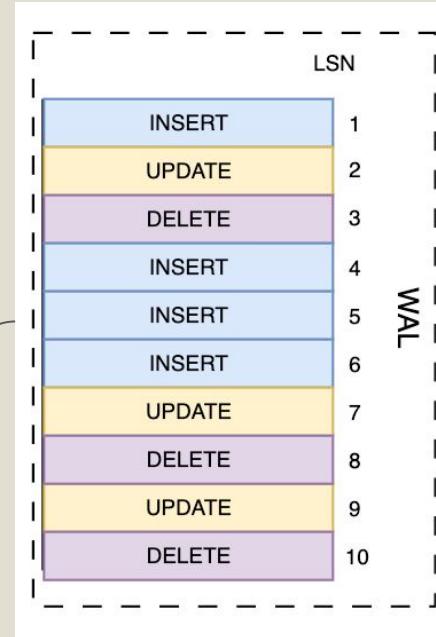
WAL log is append-only.



### WAL Segment 0001

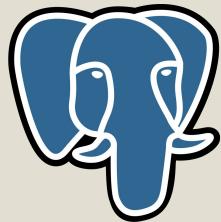


### WAL Segment 0002

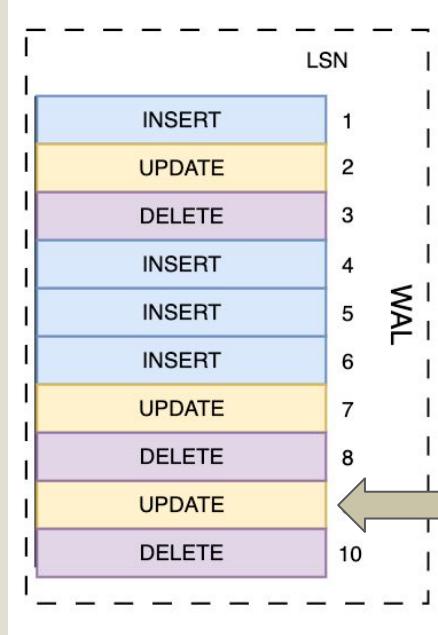


Update all  
changes to  
actual tables  
(commit)

Usual WAL segment size is **16MB**  
Default max\_wal\_size is **1GB**

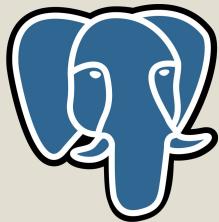


## WAL Segment 0001

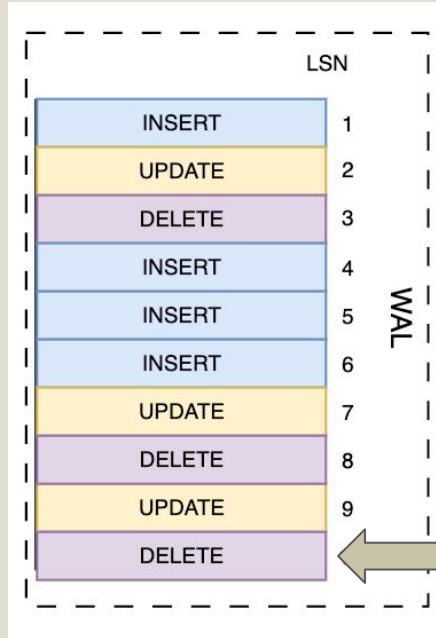


DELETE FROM users...



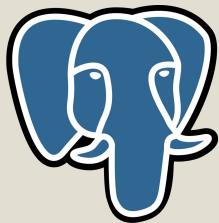


## WAL Segment 0001

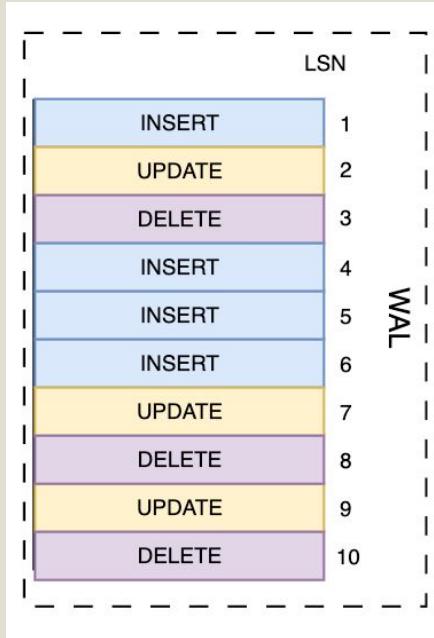


DELETE FROM users...

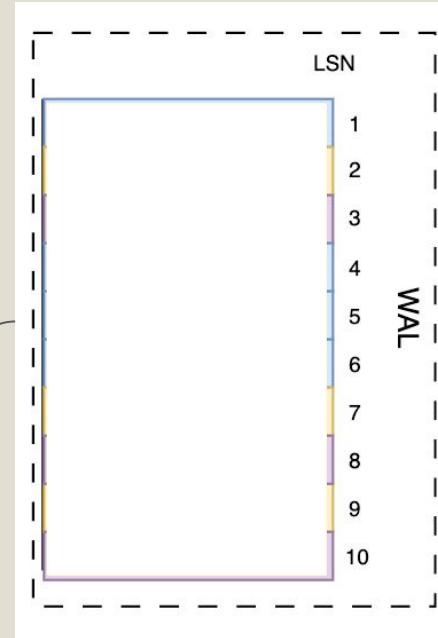
Append DELETE changes  
to the WAL log



### WAL Segment 0001

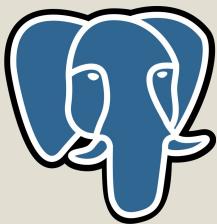


### WAL Segment 0002

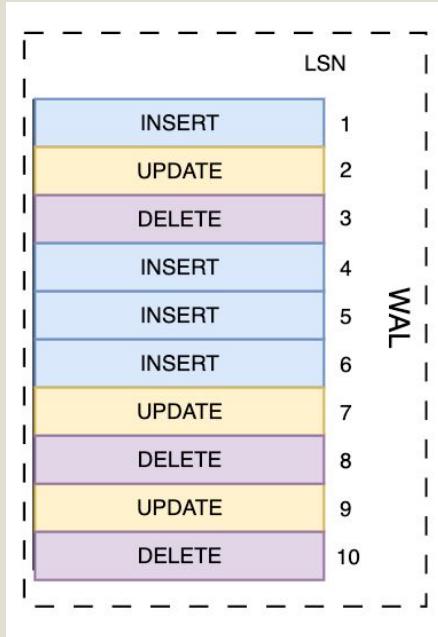


1. Write changes to tables and disk

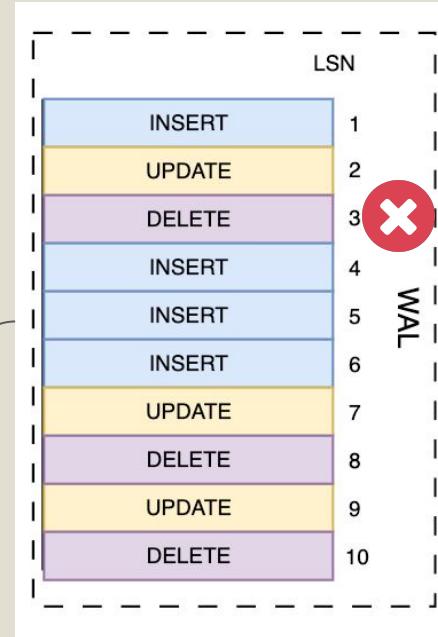
2. New WAL Segment is created after segment 1 is full



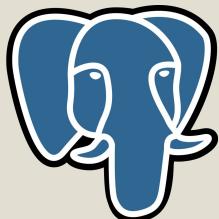
### WAL Segment 0001



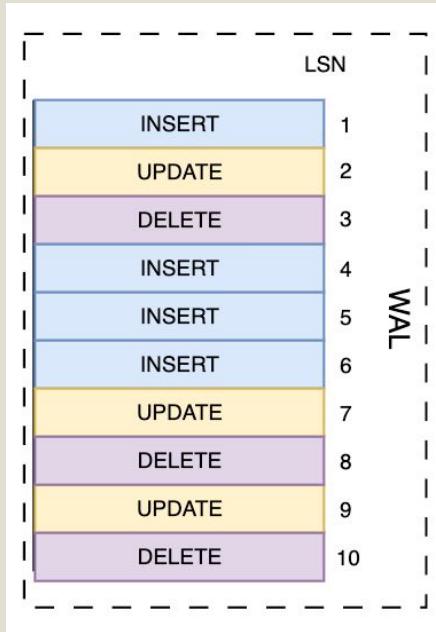
### WAL Segment 0002



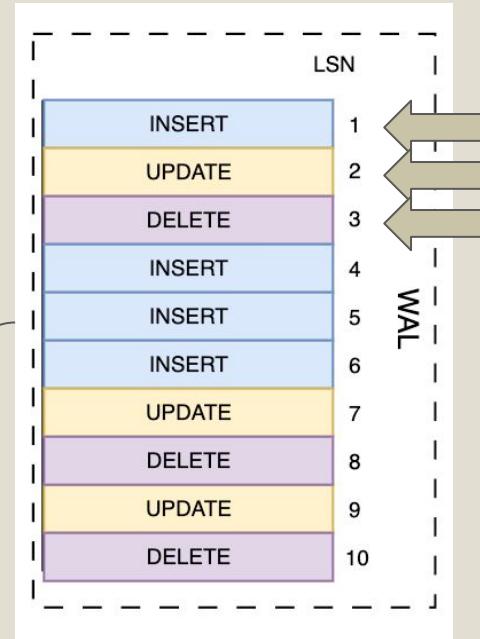
What happens when DB crashes  
after LSN 3 is committed?



WAL Segment 0001



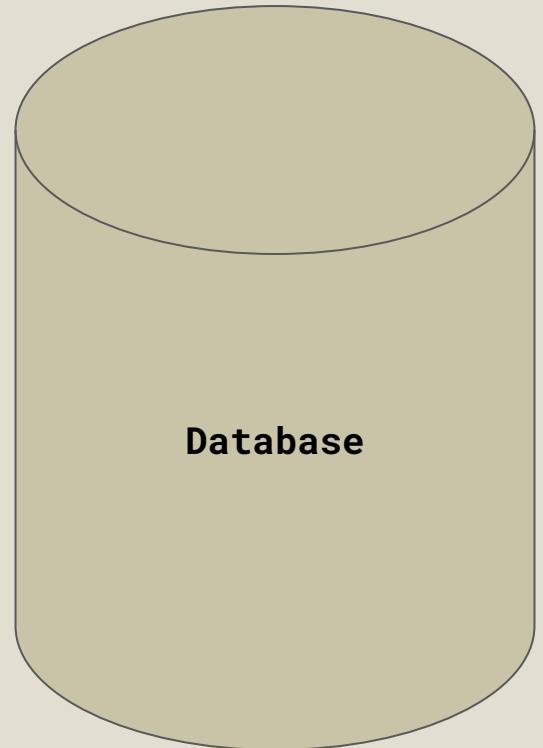
WAL Segment 0002



All changes in 0001 are committed to disk nothing to do here.

"Replay" from LSN 1, 2, and so on to recover data.

scaling  
databases

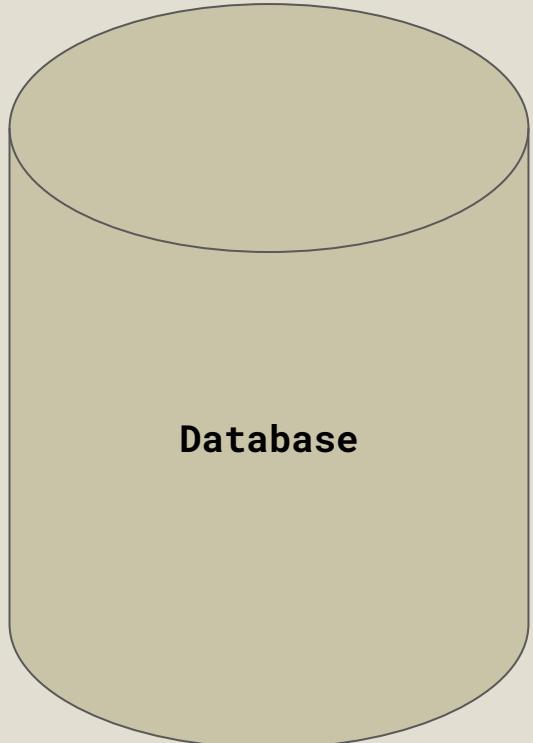


Data Size

Connections

CPU Load

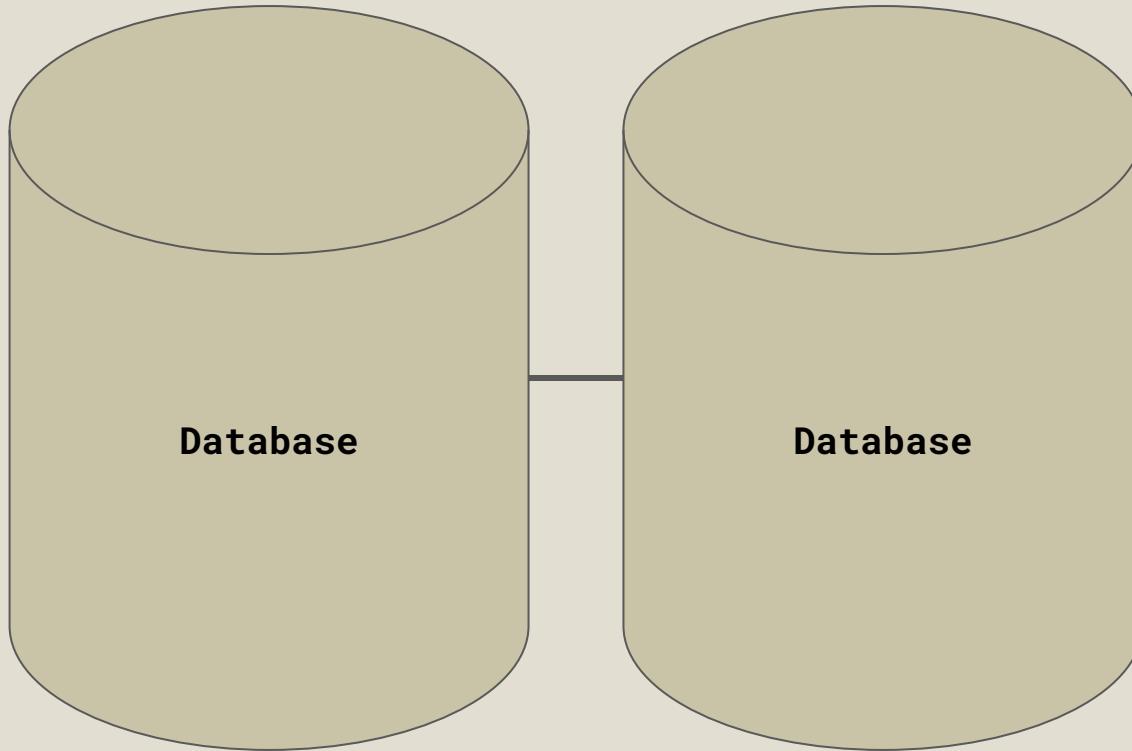
Memory Load



Writes Heavy

Reads Heavy

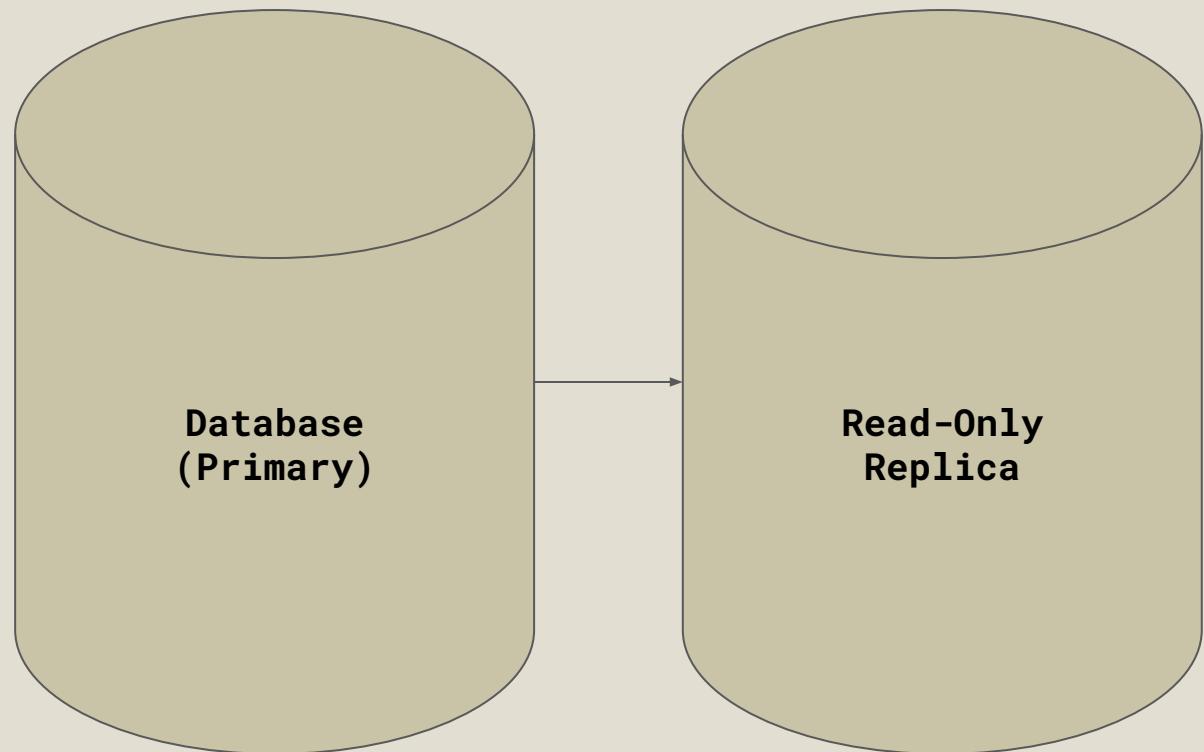
Transactional vs  
Warehousing



**Horizontal Scaling**

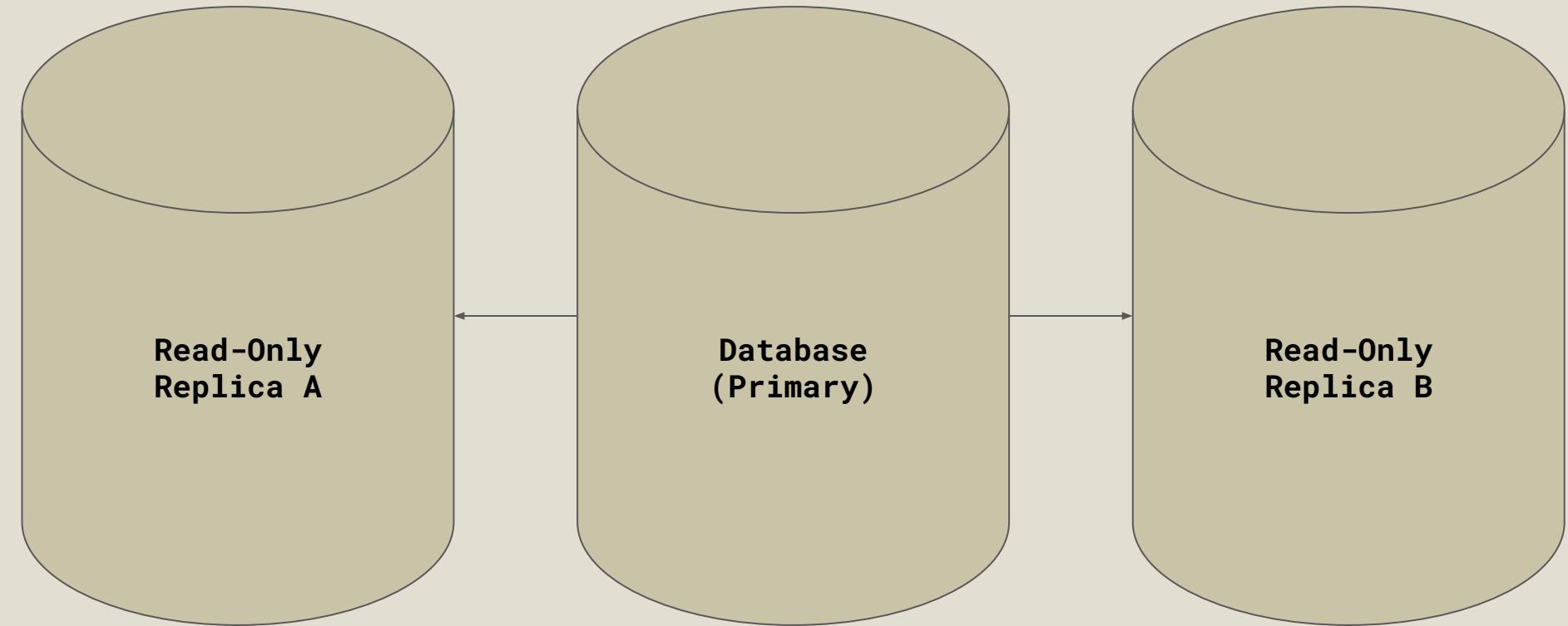
## 2.3 horizontal scaling

Node

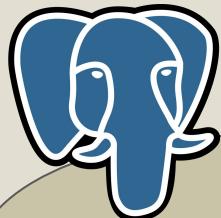


**Replication**

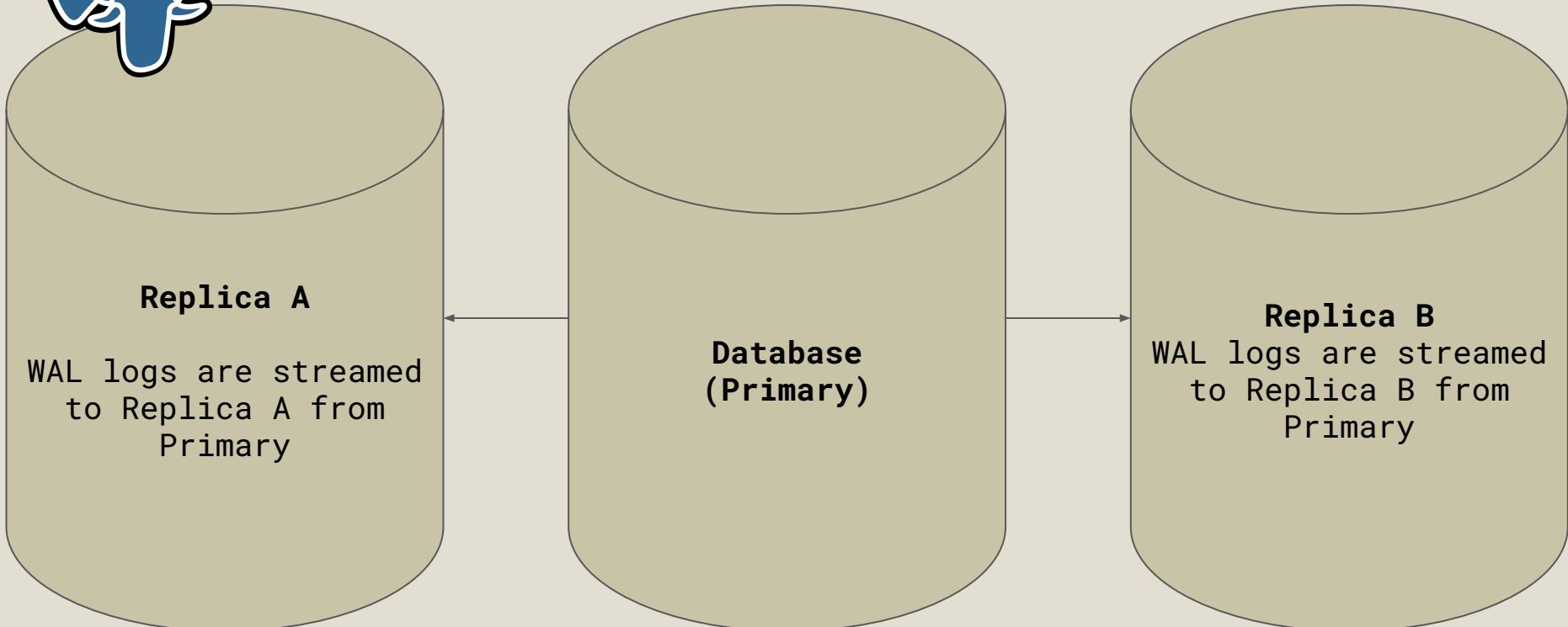
**Writes go to  
Primary only**



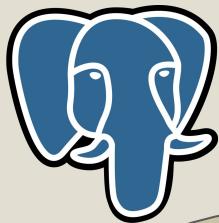
**Replication**



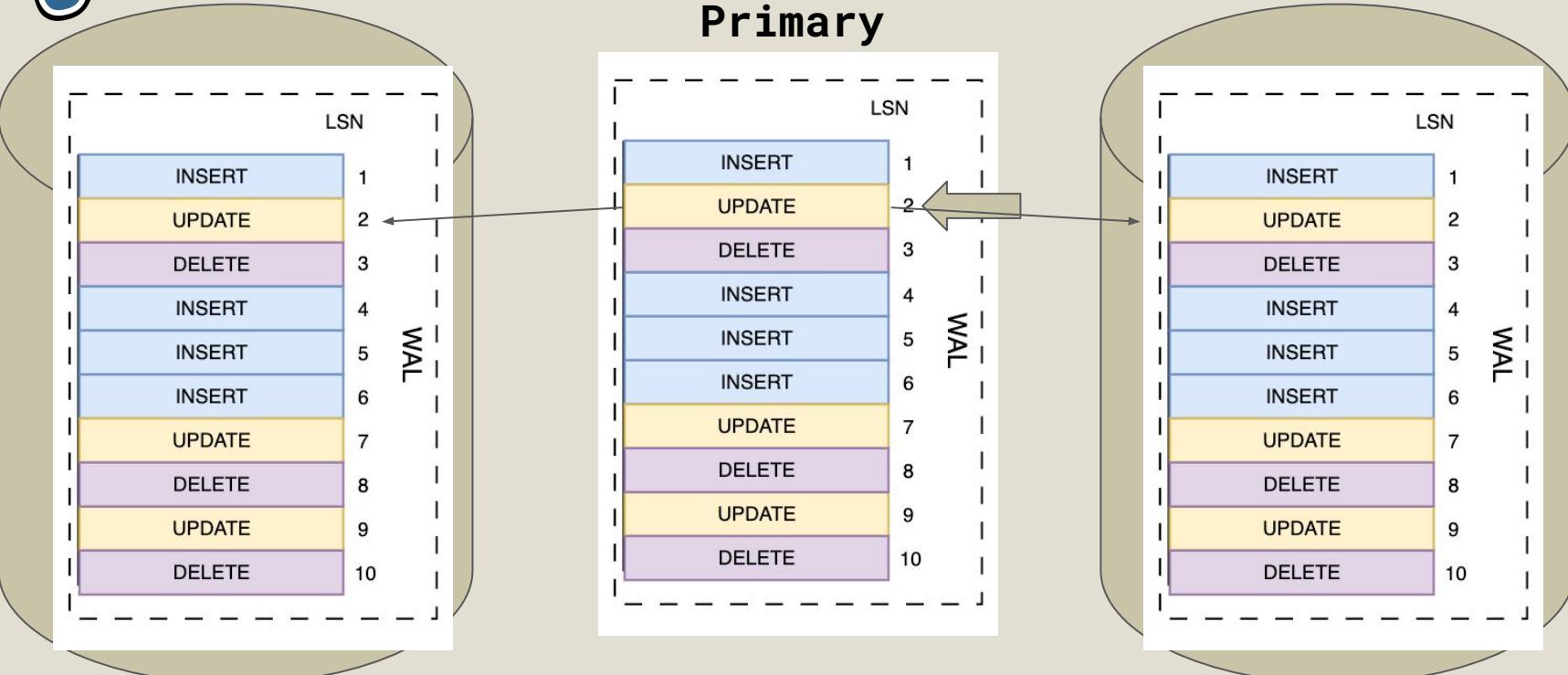
Writes go to  
Primary only



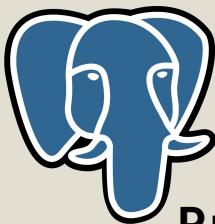
# Streaming Replication



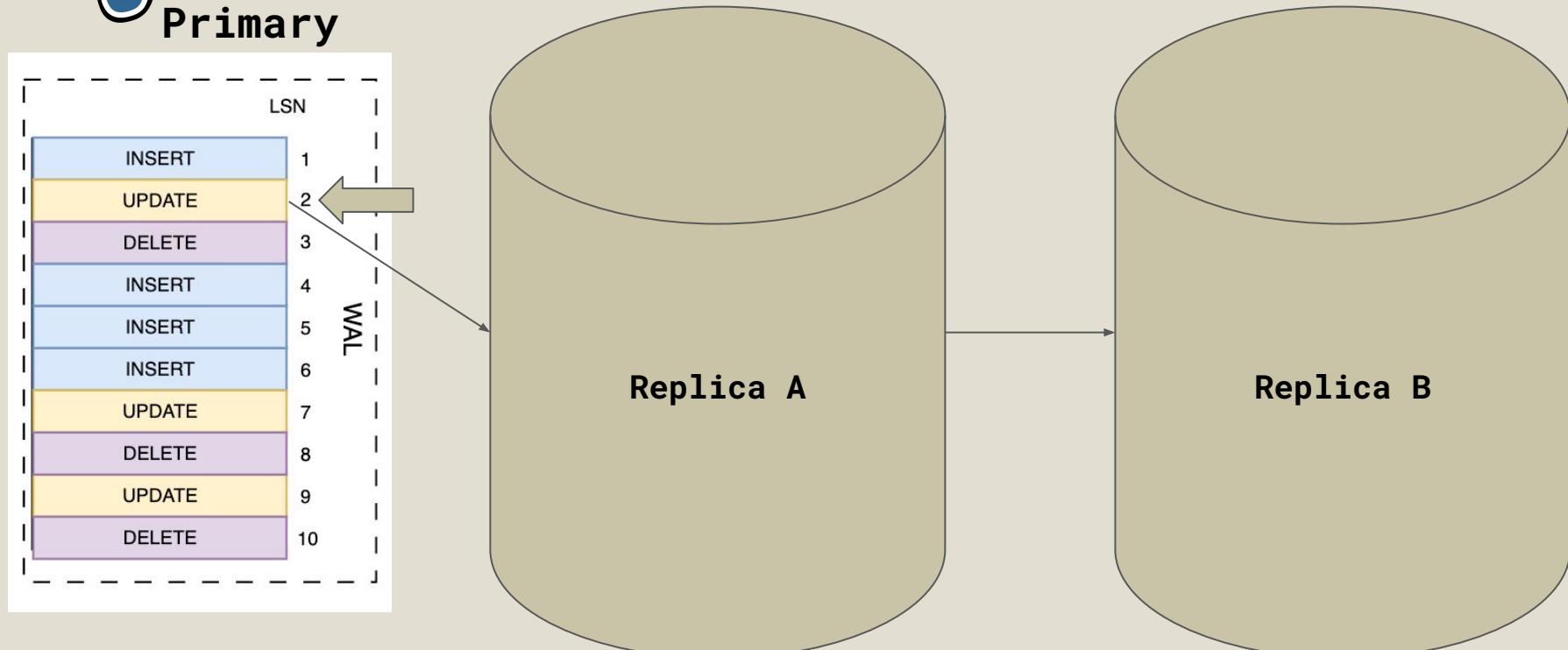
Burden is on the Primary node to keep in sync

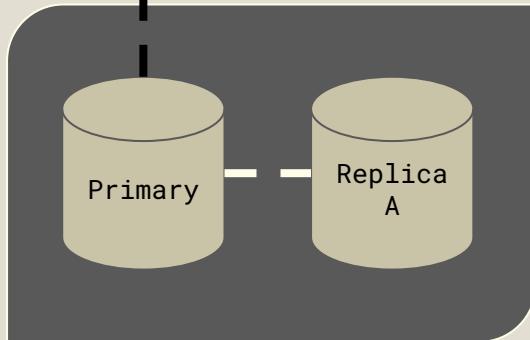


**Streaming Replication**



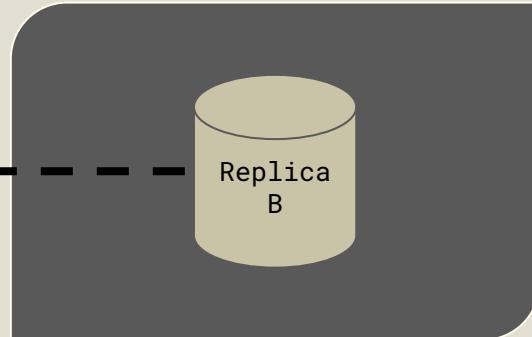
Lesser burden on Primary node, at the cost of more latency





**Region A**

Handles read and  
**write requests from  
everywhere**

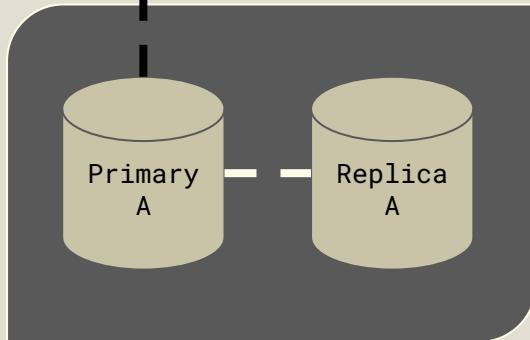


**Region B**

Handles read requests coming  
from Region B

All nodes have  
access to all the data

## Multi-Region Reads (R)



**Region A**

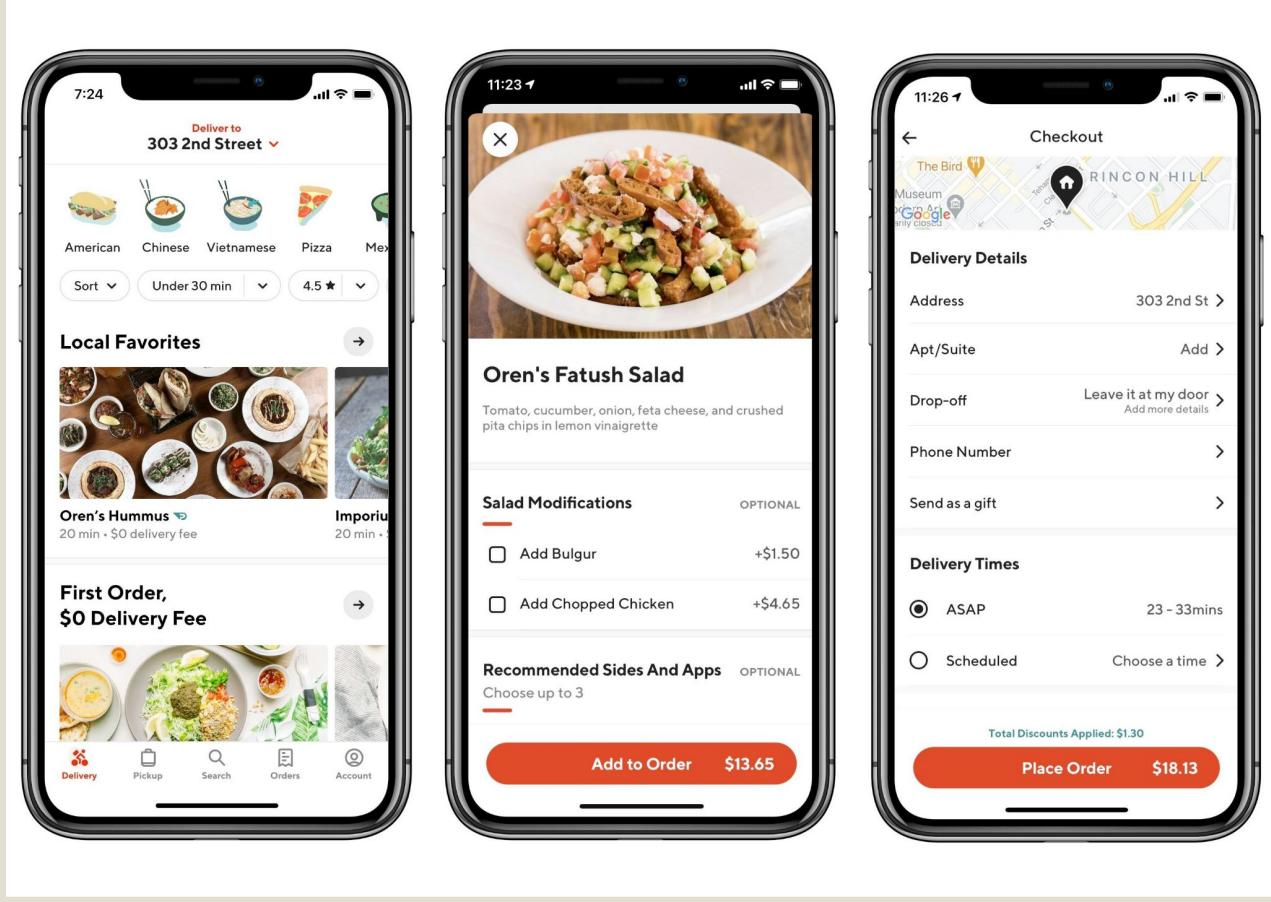
## Multi-Region Reads + Writes (RW)

Handles read and  
**write requests from**  
**Region A++**

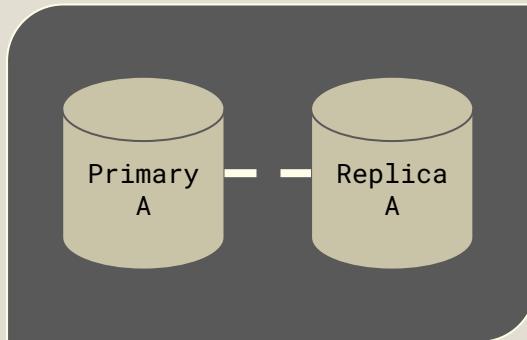
**Region B**

Handles read and **write requests**  
coming from **Region B++**

All nodes have access  
to all the data

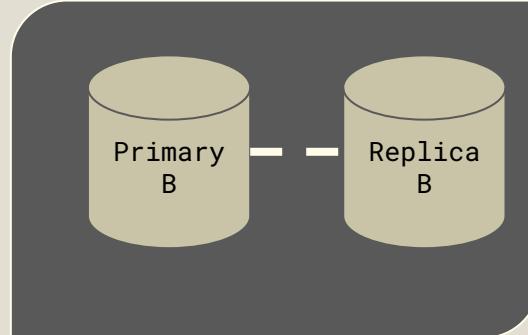


# Food Delivery: Multi-Region R or RW?



**Region A**

Handles read and  
**write requests** from  
Region A++



**Region B**

Handles read and **write requests**  
coming from Region B++

All nodes only have data in  
their own region

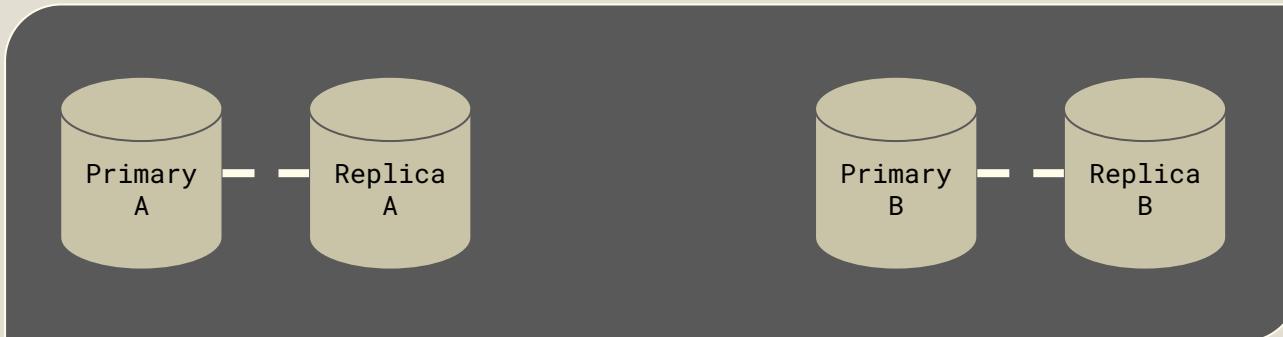
## Partitioning

Nodes only keep data they are responsible for

## 2.3.1 partitioning

```
SELECT * FROM users WHERE ID = '...'
```

Handles read and **write requests**  
based on some **partition or routing key**



Eg. Primary = Primary ID % 2

# Partitioning

## Pros

- Scale each primary and replicas independently, eg. storage
- Data can be local and situated nearer to users

## Cons

- Partition key selection is crucial to avoid overloading of a set of nodes
- Disaster Recovery needs to be tested for each set of nodes
- Queries that need to span across different set of nodes are expensive

## Most-loved picks



Crock-Pot 7 Quart Oval Manual Slo...

4.6 ★★★★★ 35,134

\$34<sup>99</sup> List: \$49.99

Delivery Fri, Mar 14



UGG Women's Classic Ultra Mini Boot

4.5 ★★★★★ 8,411

\$149<sup>95</sup>

Delivery Mar 17 - 25



LILLUSORY Women's Oversized...

4.4 ★★★★★ 4,381

\$39<sup>99</sup>

Delivery Sat, Mar 15



Streamlight 88811 Wedge 300...

4.7 ★★★★★ 4,271

\$91<sup>94</sup>

Delivery Fri, Mar 14



SaphiRose Women's Long Hood...

4.6 ★★★★★ 10,233

\$49<sup>98</sup> List: \$59.99

Delivery Fri, Mar 14



# Partition Key for an E-Commerce Website?

## Partition Key = Primary ID % 2

Show list of most-loved products

### Most-loved picks



Crock-Pot 7 Quart Oval Manual Slo...

4.6 ★★★★★ 35,134

\$34<sup>99</sup> List: \$49.99

Delivery Fri, Mar 14



UGG Women's Classic Ultra Mini Boot

4.5 ★★★★★ 8,411

\$149<sup>95</sup>

Delivery Mar 17 - 25



LILLUSORY Women's Oversized...

4.4 ★★★★★ 4,381

\$39<sup>99</sup>

Delivery Sat, Mar 15



Streamlight 88811 Wedge 300...

4.7 ★★★★★ 4,271

\$91<sup>94</sup>

Delivery Fri, Mar 14



SaphiRose Women's Long Hood...

4.6 ★★★★★ 10,233

\$49<sup>98</sup> List: \$59.99

Delivery Fri, Mar 14



Poor key choice for an E-Commerce Website

SELECT \* WHERE ID IN ...



### Most-loved picks

A stainless steel Crock-Pot slow cooker with a black lid and a handle, containing a meal of meat and vegetables.	A tan-colored UGG Women's Classic Ultra Mini Boot.	A woman wearing a black, long-sleeved, oversized sweatshirt.	A black Streamlight 88811 Wedge 300 LED flashlight.	A black SaphiRose Women's Long Hooded Raincoat with a plaid lining visible at the hem.
Crock-Pot 7 Quart Oval Manual Slo... 4.6 ★★★★★ 35,134 \$34 <sup>99</sup> List: \$49.99 Delivery Fri, Mar 14	UGG Women's Classic Ultra Mini Boot 4.5 ★★★★★ 8,411 \$149 <sup>95</sup> Delivery Mar 17 - 25	LILLUSORY Women's Oversized... 4.4 ★★★★★ 4,381 \$39 <sup>99</sup> Delivery Sat, Mar 15	Streamlight 88811 Wedge 300-... 4.7 ★★★★★ 4,271 \$91 <sup>94</sup> Delivery Fri, Mar 14	SaphiRose Women's Long Ho... 4.6 ★★★★★ 10,233 \$49 <sup>98</sup> List: \$59.99 Delivery Fri, Mar 14

Data split between partitions

**Peer Discussion**

- 1. What other kinds of products / services can benefit from partitioning? (Name at least 2)**
  
- 2. What sort of keys would make sense to be partition keys for these products / services?**

**~10 mins, discuss in pairs and pick a person to share.**

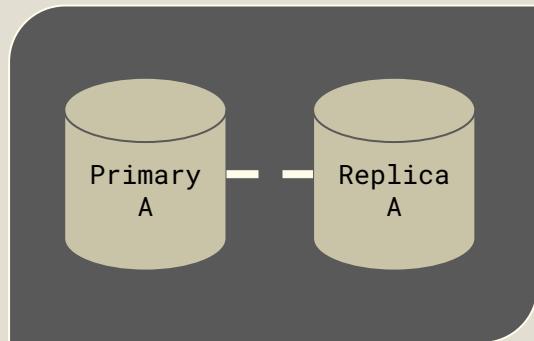
**General Guideline: Select a key that best allow  
50 - 80% of Reads to be within a single partition**



**What is data partitioning,  
and how to do it right**

<https://www.cockroachlabs.com/blog/what-is-data-partitioning-and-how-to-do-it-right/>

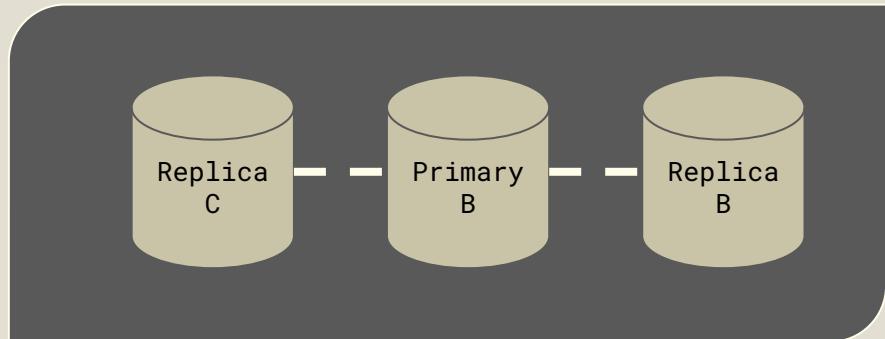
Handles read and  
**write requests** from  
Region A++



Region A

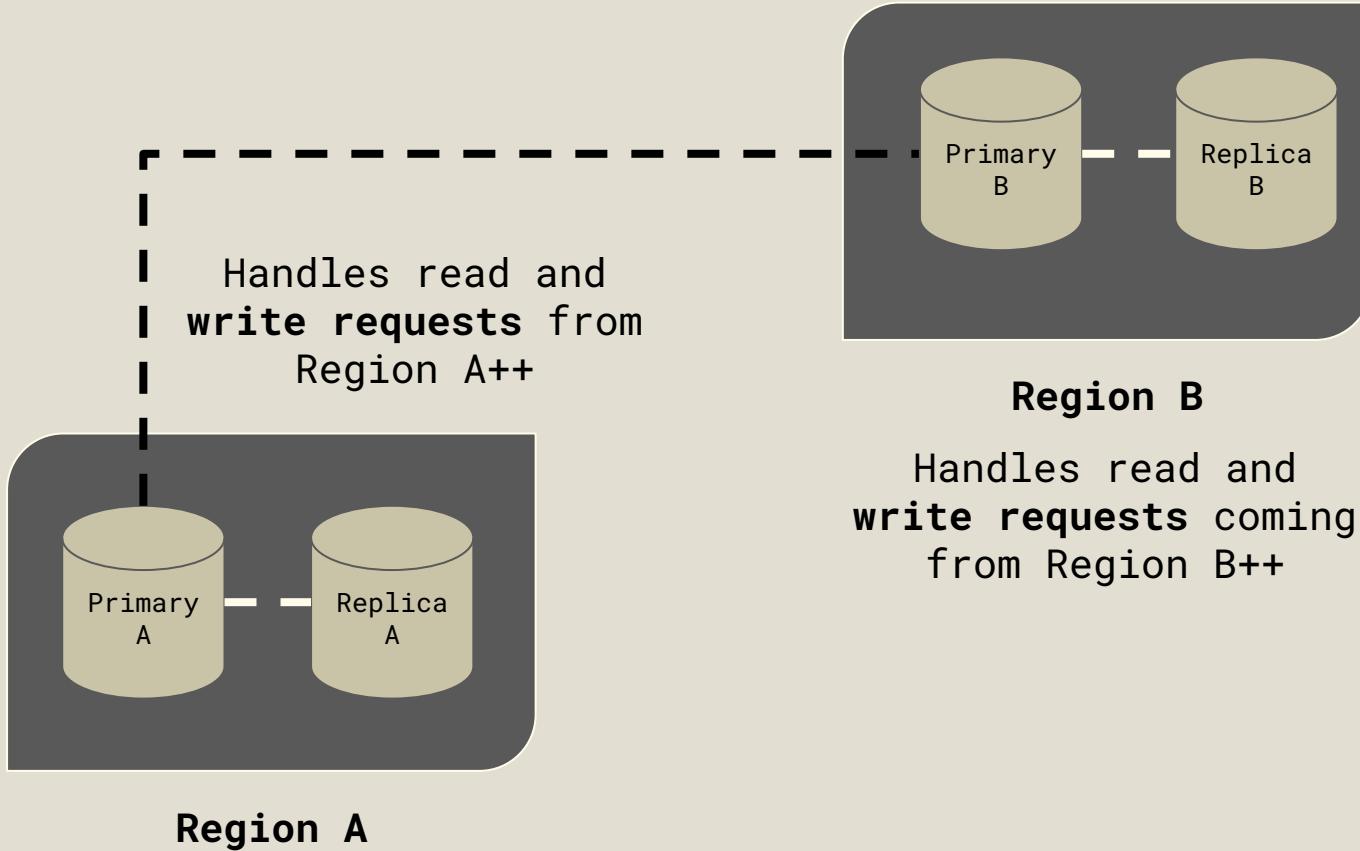
Region B

Handles read and  
**write requests** coming  
from Region B++



# Regional Scaling

## 2.3.2 multi-primary



## Multi-Primary

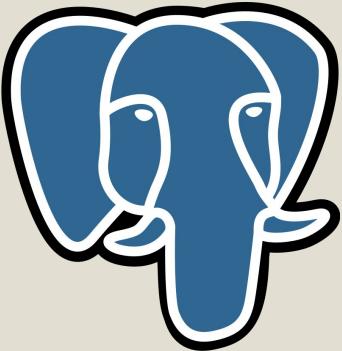
All nodes keep all the data

## Pros

- All nodes have full copies of all data\*  
(easier time with disaster recovery)
- Each DB has a complete view of the  
entire Database

## Cons

- Storage space multiplied by N nodes
- Usually need at least 3 or more nodes
- Primaries need to achieve **consensus**  
for a write to be confirmed

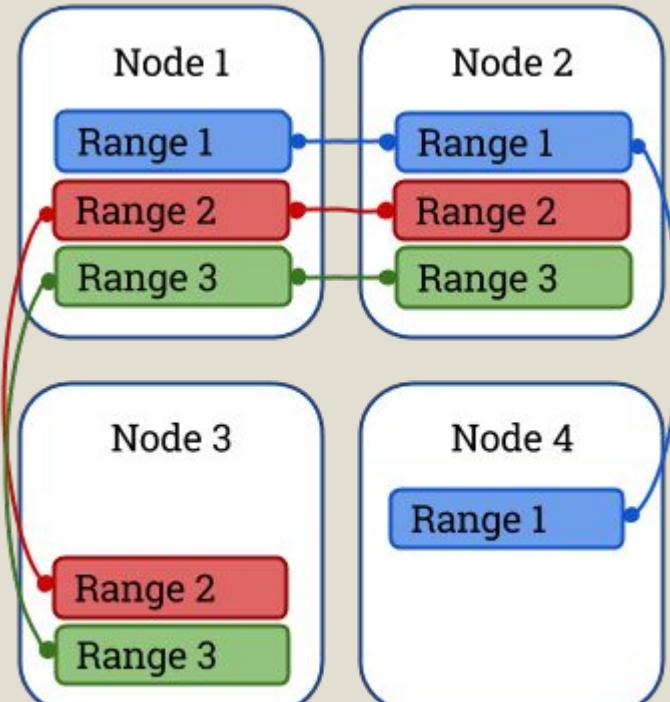


Does not natively Multi-Primary replication!

(3rd party tools exist to do so)



Supports Multi-Primary replication



Range Replication in CockroachDB

**Number of Nodes and Number of Replicas (Replication Factor)** available determines the number of node failures we can tolerate.

Number of replicas = 3 by default in CockroachDB

## More on Multi-Primary in W3

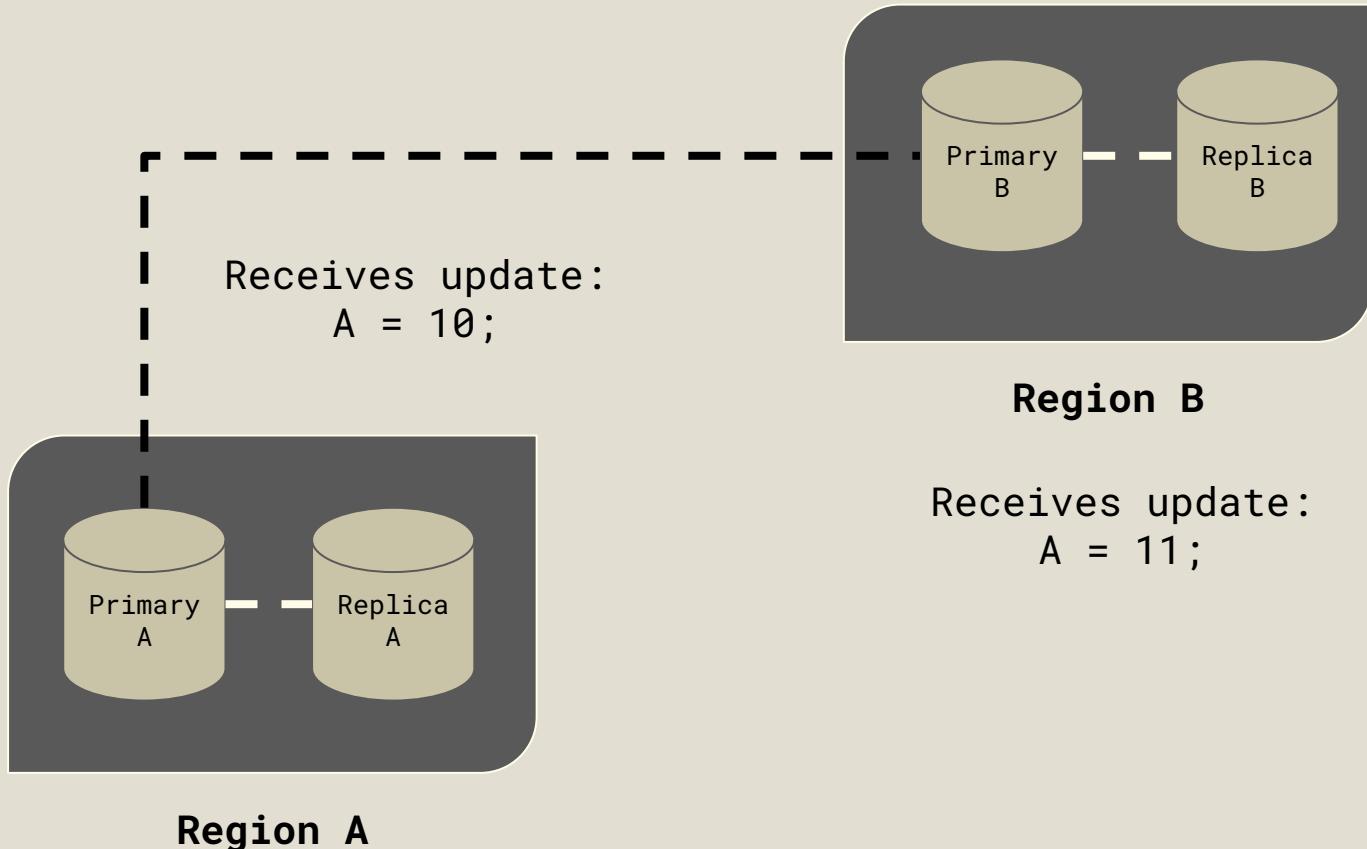
<https://www.cockroachlabs.com/blog/the-new-stack-meet-cockroachdb-the-resilient-sql-database/>

	<b>Dataset on each node</b>	<b>Risks</b>
<b>Partitioning</b>	Partial	Routing key needs to be well selected
<b>Multi-Primary</b>	Partial	Consensus* needs to be achieved, latency
<b>Streaming Replication</b>	Full	Primary node might be overwhelmed
<b>Chained Replication</b>	Full	Latency for data to reach nodes that are far in the chain

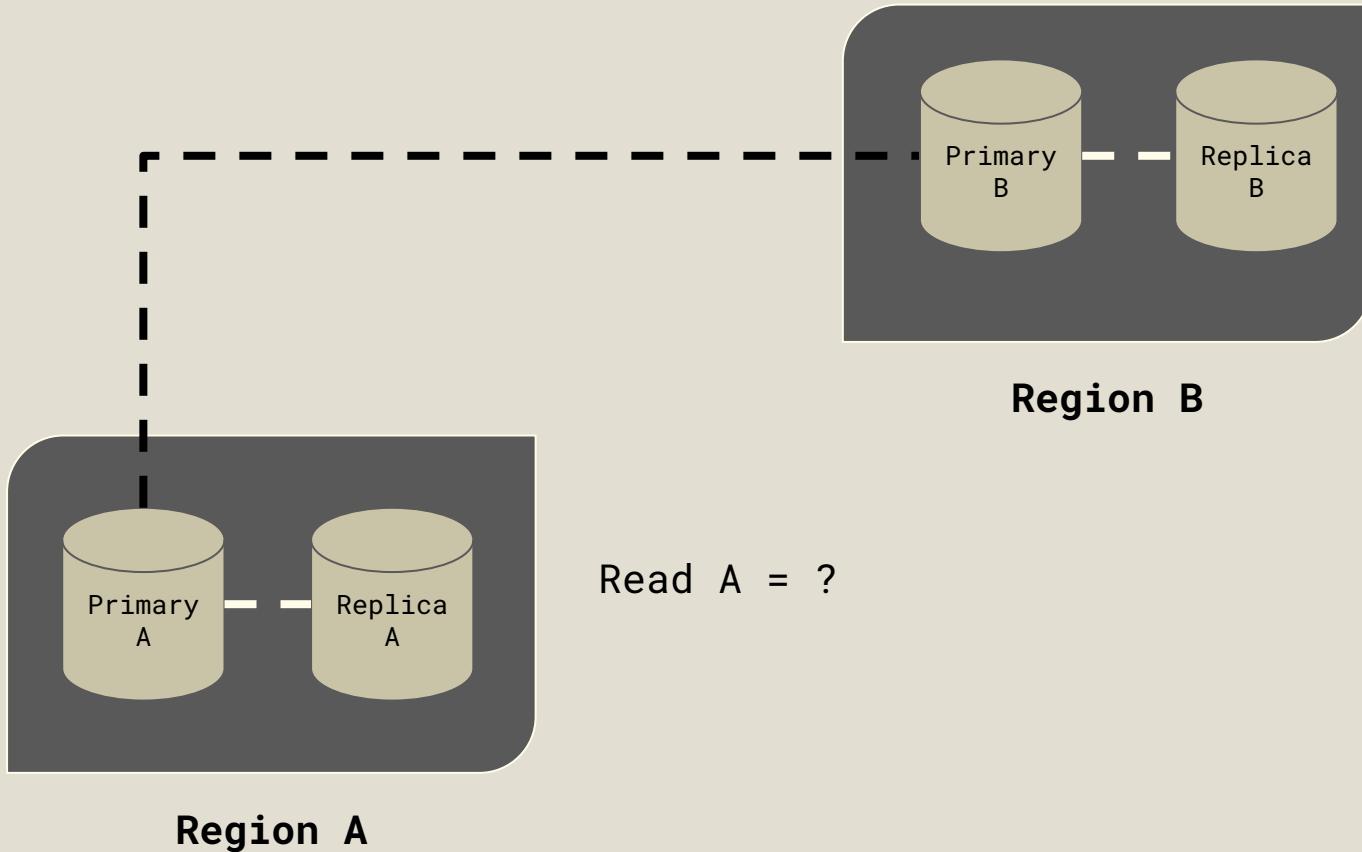


# Horizontal Scaling Methods

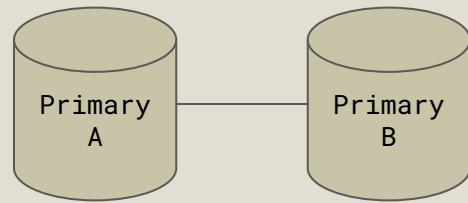
## 2.4: consensus and gossip



# Consensus Problem

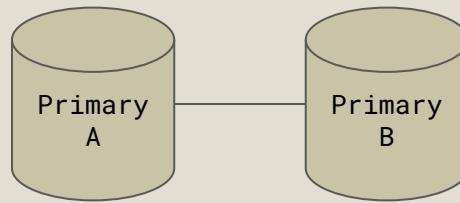


Receives update: A = 10;



Receives update: A = 11;

Receives update: A = 10;

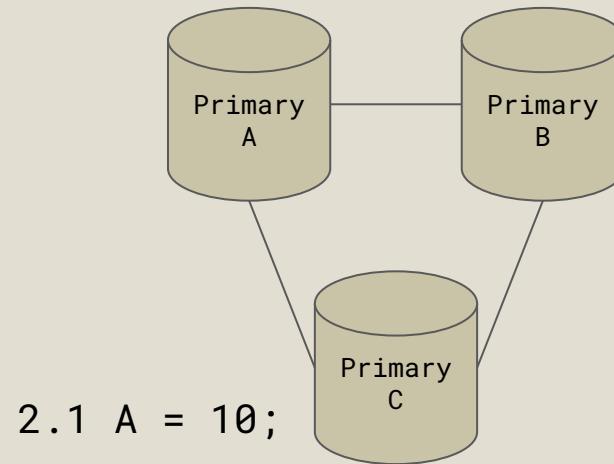


Receives update: A = 11;

**Unlikely to achieve consensus  
without intervention**

1.1 Receives update:  $A = 10$ ;

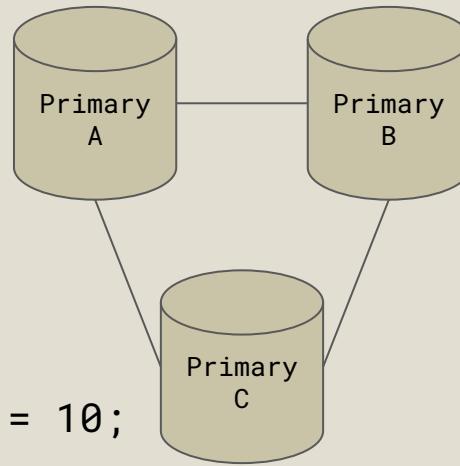
1.2 Receives update:  $A = 11$ ;



**Consensus: A and C agrees**

1.1 Receives update: A = 10;

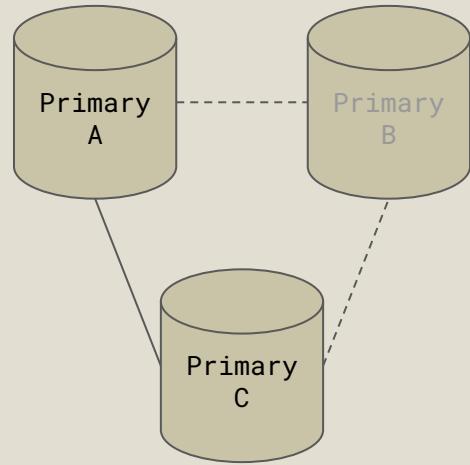
3.1 Receives update: A = 10;



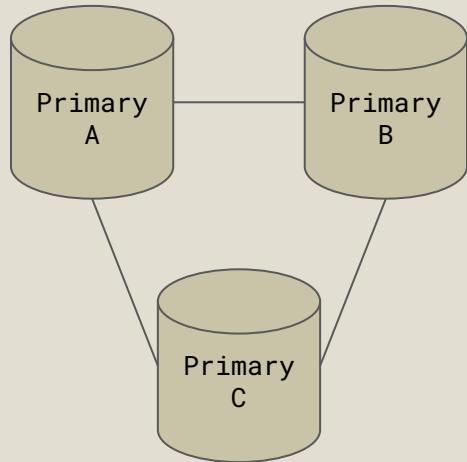
2.1 Receives update: A = 10;

**Consensus achieved: % agrees**

# Even More Questions



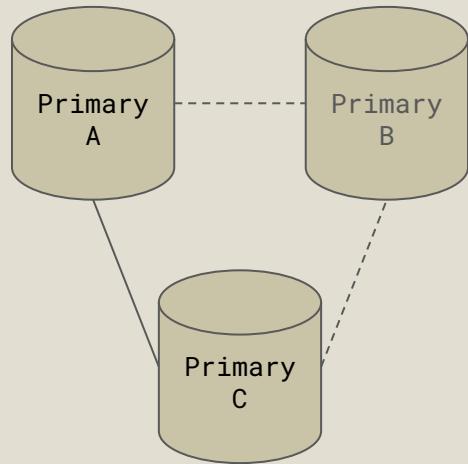
**Q. What happens when a node disappears?**



**Q. How do nodes talk amongst each other?**

**Break Time**

# CAP Theorem

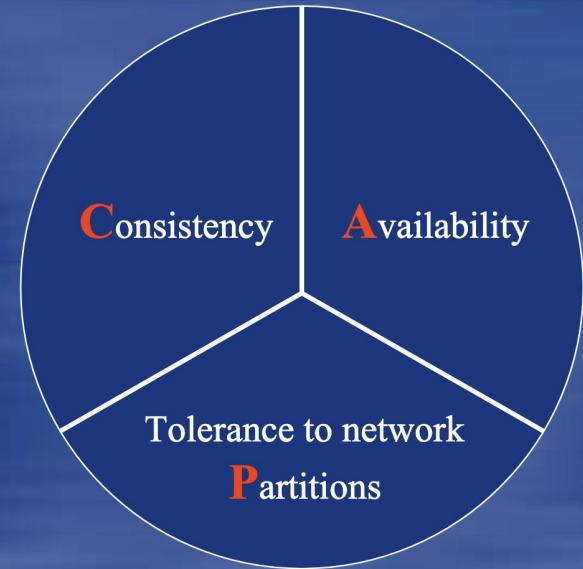


**Q. What happens when a node disappears  
in a multi-primary setup?**



Inktomi

# The CAP Theorem



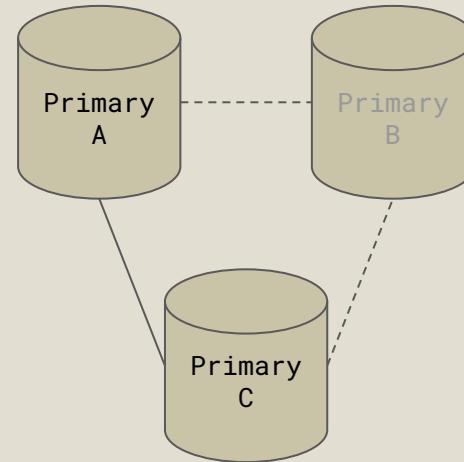
Theorem: You can have **at most two** of these properties for any shared-data system

PODC Keynote, July 19, 2000

# CAP Theorem

[https://sites.cs.ucsb.edu/~rich/class/cs293b-cloud/papers/Brewer\\_podc\\_keynote\\_2000.pdf](https://sites.cs.ucsb.edu/~rich/class/cs293b-cloud/papers/Brewer_podc_keynote_2000.pdf)

When network partitions happen, you need to decide if you want to prioritize **Availability** or **Consistency**.



**CAP Theorem**

ATM network encounters a partition, and the ATM may be unable to get the current actual account balance.

Will you:

- a. Deny all transactions
- b. Allow all transactions
- c. Allow some types of transactions? (what types?)



**ATM: Consistency or Availability?**

Types of transactions:

- Deposit
- **Withdrawal**
- Check balance

Can withdrawals happen without knowing the perfect account balance?



**ATM: Consistency or Availability?**

Limit withdrawals to some amount,  
eg. \$100 or less → **prioritize Availability**

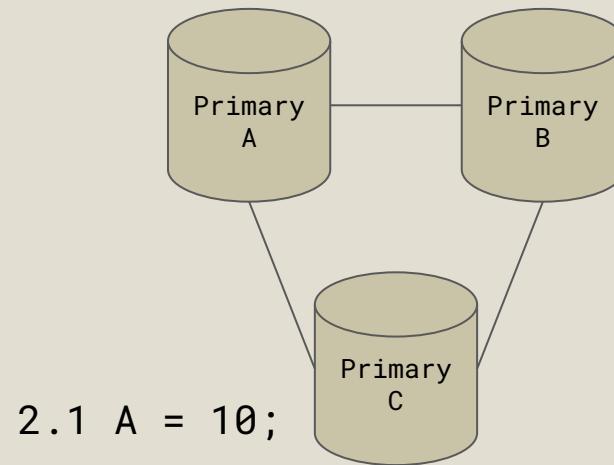
Deny all withdrawals and deposits  
→ **prioritize Consistency**



**ATM: Consistency or Availability?**

1.1 Receives update: A = 10;

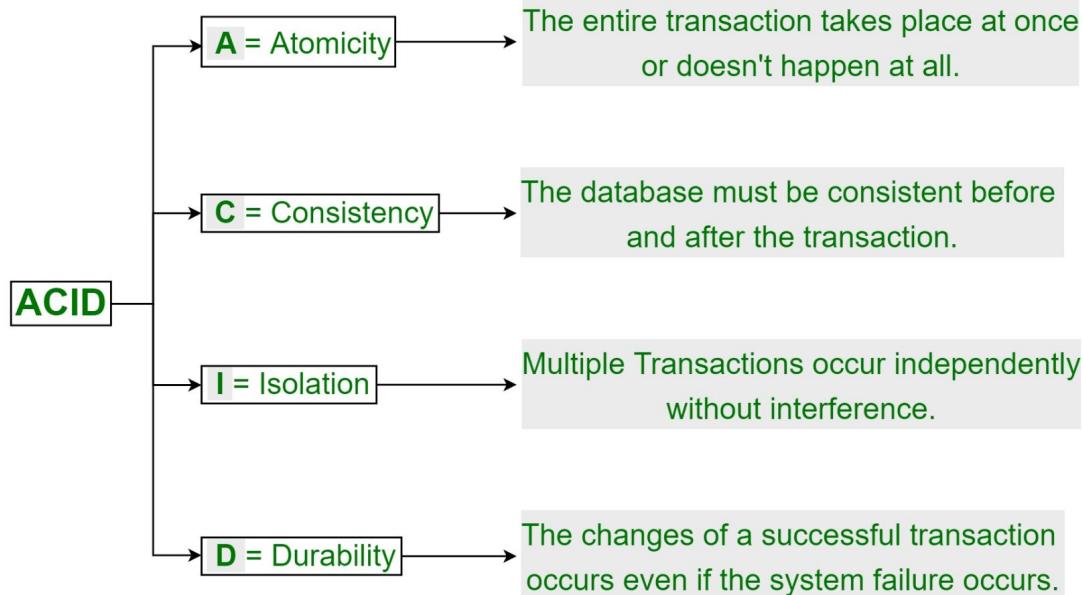
3.1 Receives update: A = 10;



2.1 A = 10;

**Note about Consistency in CAP**

# ACID Properties in DBMS



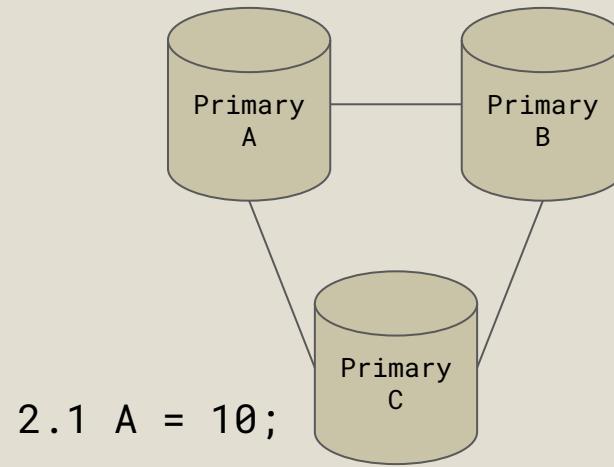
ĐG

**CAP Consistency != ACID Consistency**

<https://www.geeksforgeeks.org/acid-properties-in-dbms/>

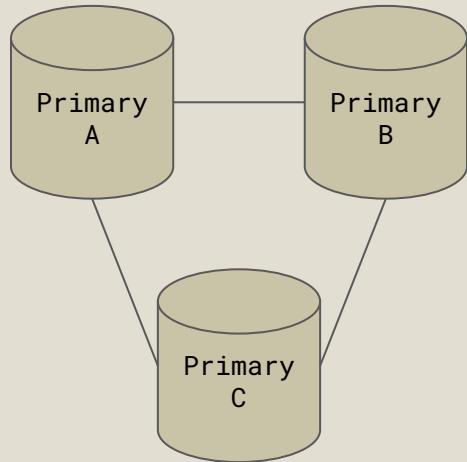
1.1 Receives update: A = 10;

3.1 Receives update: A = 10;



**Eventual Consistency:**  
All nodes will be consistent, eventually

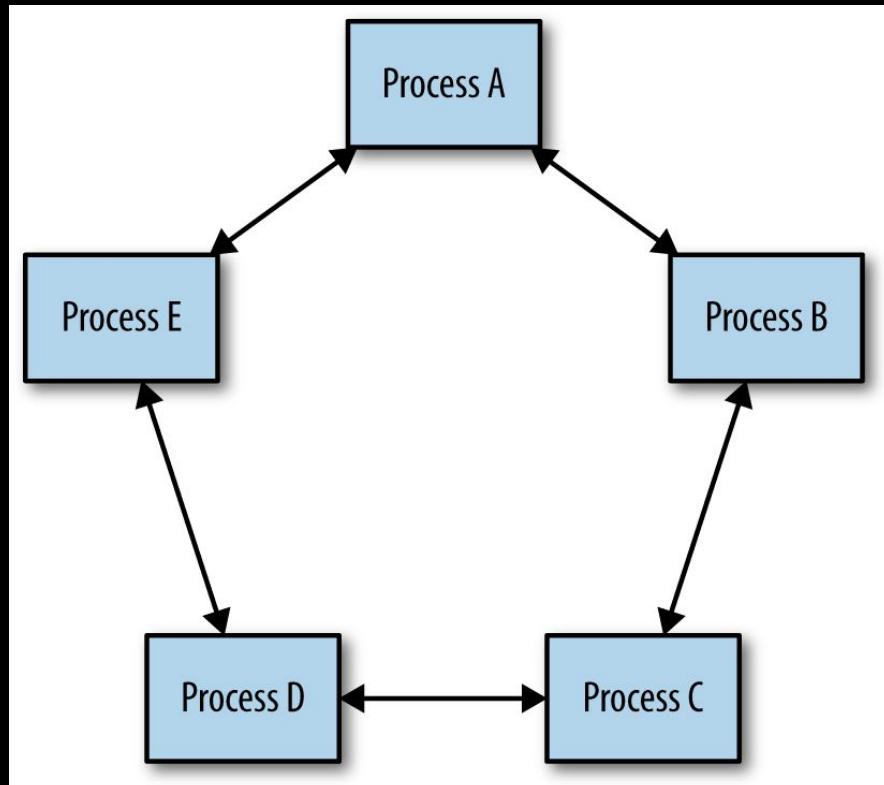
# Consensus and Gossip Protocols



**Q. How do nodes talk amongst each other?**

"We've found distributed consensus to be effective in building reliable and highly available systems that require a consistent view of some system state.

The distributed consensus problem deals with reaching agreement among a group of processes connected by an unreliable communications network."



"We've found distributed consensus to be effective in building reliable and highly available systems that require a consistent view of some system state.

The distributed consensus problem deals with reaching agreement among a group of processes connected by an unreliable communications network."

**System State = Data in Database**

**Group of Processes = Database Nodes**



**Raft**

**Paxos**  
(no mascot)

**Popular Consensus Algorithms**

"A Raft cluster of 3 nodes can tolerate a single node failure while a cluster of 5 can tolerate 2 node failures.

The recommended configuration is to either run 3 or 5 Consul servers per datacenter.

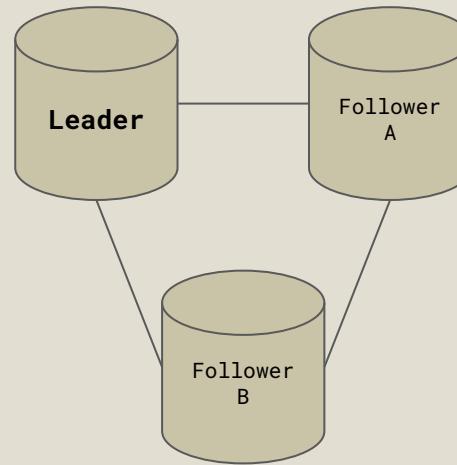
This maximizes availability without greatly sacrificing performance. The deployment table below summarizes the potential cluster size options and the fault tolerance of each."

Consensus Protocol

<https://developer.hashicorp.com/consul/docs/architecture/consensus>



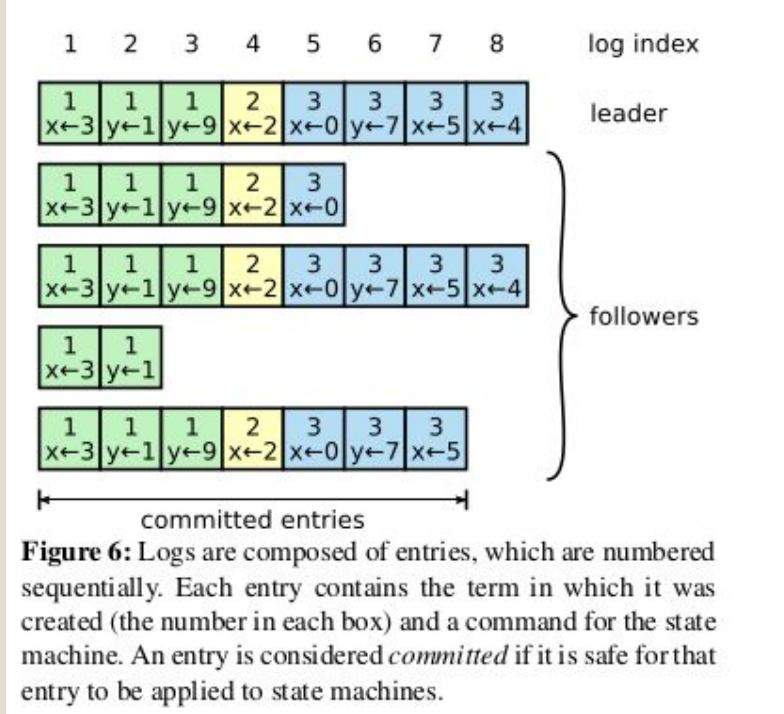
- There is always a **Leader**
- **Leader decides what data will be committed**
  - When Leader goes missing, Followers will put themselves up for voting to choose next Leader (**Leader Election**)
  - A quorum is a majority of members: for a cluster of size N, quorum requires at least  $(N/2)+1$  members, eg. cluster of 3 nodes can tolerate 1 failure only



## Raft Algorithm: Quick Summary

<https://raft.github.io/raft.pdf>

## Append-only log in Raft



**Figure 6:** Logs are composed of entries, which are numbered sequentially. Each entry contains the term in which it was created (the number in each box) and a command for the state machine. An entry is considered *committed* if it is safe for that entry to be applied to state machines.

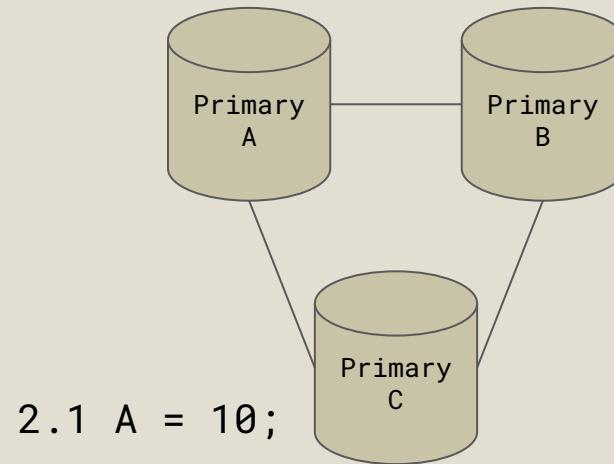


## Raft Algorithm: Quick Summary

<https://raft.github.io/raft.pdf>

1.1 Receives update: A = 10;

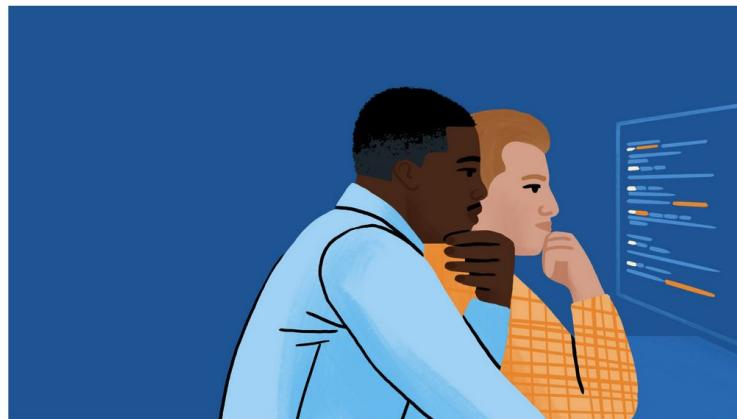
3.1 Receives update: A = 10;



**Raft used for ensuring  
eventual consistency**

POSTED ON MAY 16, 2023 TO DATA INFRASTRUCTURE, PRODUCTION ENGINEERING

## Building and deploying MySQL Raft at Meta



By Anirban Rahut, Abhinav Sharma, Yichen Shen, Ahsanul Haque



**Before: MySQL semisynchronous databases (which uses Paxos internally)**

**After: MySQL Raft**

# Building and deploying MySQL Raft at Meta

<https://engineering.fb.com/2023/05/16/data-infrastructure/mysql-raft-meta/>

# Summary

- Understand how the underlying mechanics work for different databases, eg. WAL logs for PostgreSQL
- Prefer simple horizontal scaling options (eg. replication) over complicated ones such as Multi-Primary
- Partitioning is another option if data can be split cleanly between regions

- CAP Theorem: Managing risks of network failures (partitions) by prioritizing either Consistency or Availability
- Consensus: How distributed nodes talk to each other and agree, eg. Raft algorithm
- We will explore a little more on CAP when we talk about Distributed Databases next week

p. s.  
Current Trends

searchcode.com itself is now officially a ship of Theseus project, as every part I started with has now been replaced. A brief history,

- First version released using PHP, CodeIgniter, MySQL, Memcached, Apache2 and Sphinx search
- Rewritten using Python, Django, MySQL, Memcached, Sphinx search, Nginx and RabbitMQ
- Never publicly released version using Java, MySQL, Memcached, Nginx and Sphinx search
- Start of Covid 19 rewritten using Go, MySQL, Redis, Caddy and Manticore search
- Replaced Manticore search with [custom index](#) now the stack consists Go, MySQL, Redis and Caddy
- As of a few days ago Go, SQLite, Caddy

The constant between everything till now was the use of MySQL as the storage layer. The reasons for using it initially was it was there, I knew how to work it and it would scale along with my needs fairly well. So what changed? If you look at my previous choices you will see there is in general a move to reducing the number of dependencies. The older and more crusty I get the more I appreciate having a single binary I can just deploy. Single binary deploys are very simple to reason about.

**searchcode.com's SQLite database is probably 6 terabytes bigger than yours**

<https://boyter.org/posts/searchcode-bigger-sqlite-than-you/>



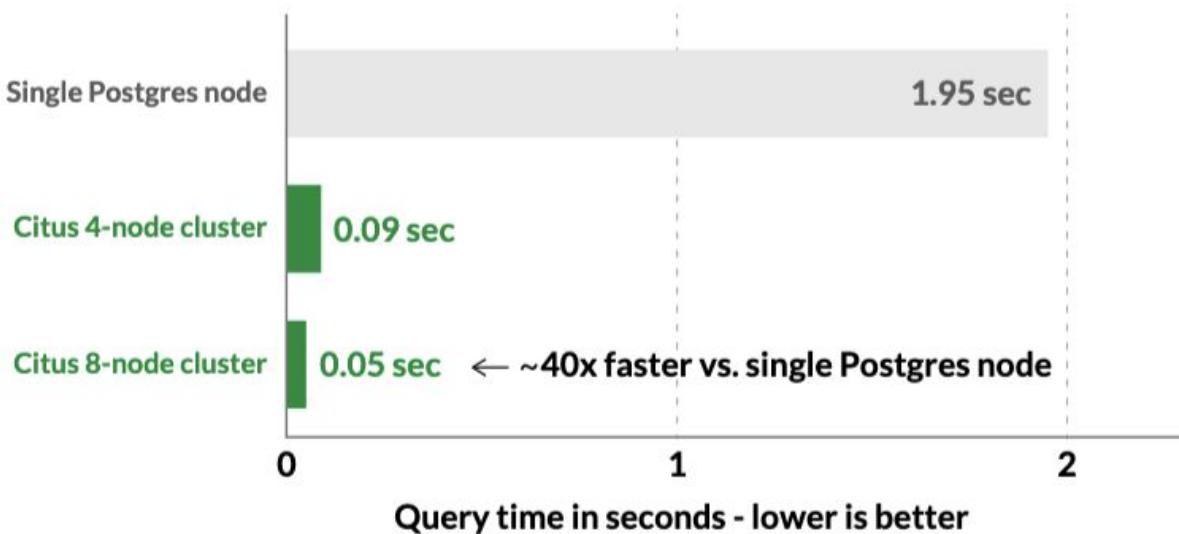
Database Category	Database Name	Wire-Protocol Compatibility	Syntax Compatibility	Feature Compatibility	Runtime Compatibility
Sharded Postgres	CitusData	Full	Full	High	High
	Azure CosmosDB for PostgreSQL	Full	Full	High	High

Database Category	Database Name	Wire-Protocol Compatibility	Syntax Compatibility	Feature Compatibility	Runtime Compatibility
Distributed SQL Databases	Google Spanner	Full	High	Low	Low
	CockroachDB	Full	High	Moderate	Low
	YugabyteDB	Full	Full	High	High

## Distributed Databases offering PostgreSQL Compatibility

<https://www.yugabyte.com/postgresql/compare-postgresql-compatibility/>

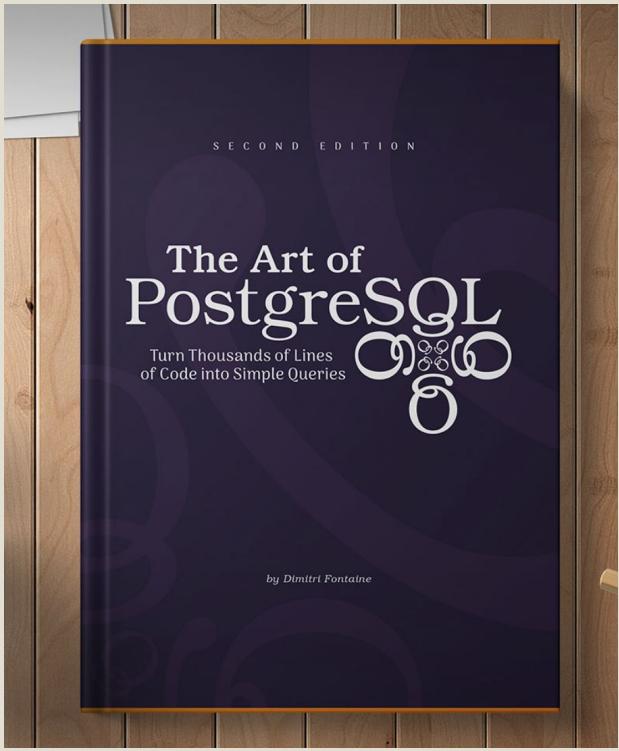


Time to complete a SQL query, using a benchmark that measures analytical query performance. Run on ~100 GB of GitHub archive data in JSON format. All servers are Azure VMs with 16 vCPUs, 64 GB of memory, and network-attached disks with 7500 IOPS –running Postgres 13 and Citus 9.5, with default settings.

## Citus's Claims for Performance for Distributed PostgreSQL

<https://www.citusdata.com/getting-started/>

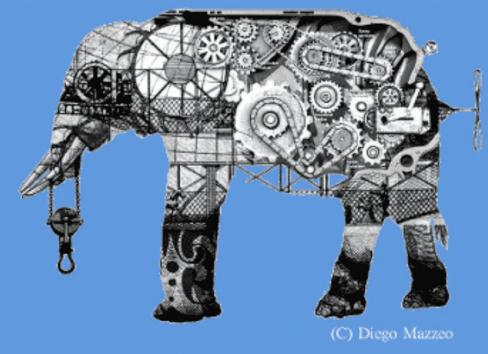
optional  
readings



**Art of PostgreSQL**  
<https://theartofpostgresql.com/>

## The Internals of PostgreSQL

for database administrators and system developers



**The Internals of PostgreSQL**  
<https://www.interdb.jp/pg/index.html>

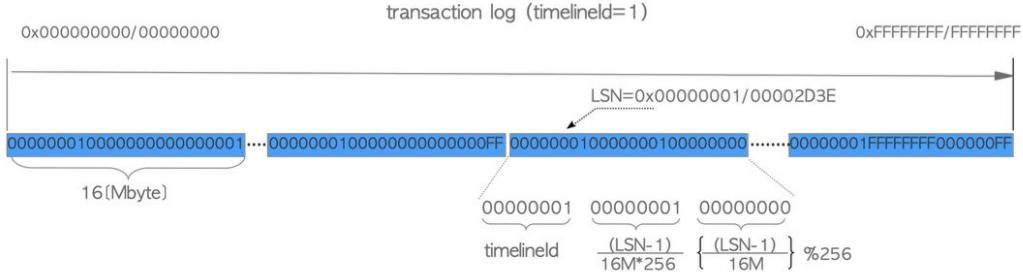


Figure 9.6. Transaction log and WAL segment files

The WAL segment filename is in hexadecimal 24-digit number and the naming rule is as follows:

$$\text{WAL segment file name} = \text{timelineld} + (\text{uint32}) \frac{\text{LSN} - 1}{16M * 256} + (\text{uint32}) \left( \frac{\text{LSN} - 1}{16M} \right)$$

#### **i timelineld**

PostgreSQL's WAL contains the concept of **timelineld** (4-byte unsigned integer), which is for Point-in-Time Recovery (PITR) described in [Chapter 10](#). However, the timelineld is fixed to 0x00000001 in this chapter because this concept is not required in the following descriptions.

O'REILLY®

# Database Internals

A Deep-Dive into How Distributed Data  
Systems Work



Alex Petrov

**Database Internals**  
<https://www.databass.dev/>

# assignment

# Assignment

## Extend a Distributed Social Media Platform



### Assignment #1

- You are a Senior Engineer in a team responsible for building a new product that will extend upon one of these platforms / protocols.
- Product: Let user post and upload their location check-ins, similar to Foursquare / Swarm) with GPS coordinates and Photos, with a global newsfeed.
- **Choose one of the two protocols, analyze it, and share in a 3-page technical document about why it would be a good or bad choice for this project.**
- Audience: Tech Leads, Software Engineers, Product Managers, Senior Business Leaders
- **Due date: Sunday, 23:59 (UTC+8)**

**W1 Assignment**

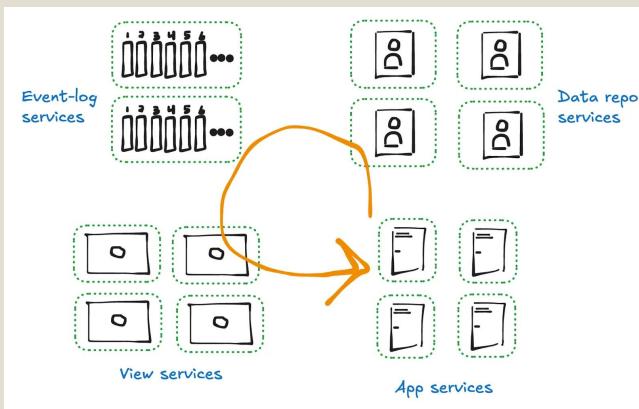
- 2 chose to review AT Protocol
- 2 chose to review ActivityPub
- 2 chose to review and compare both

Assignment

- W1 → Distributed Social Media Research
- W2 → Review W1 + Build PoC
- W3 → Design 1st Draft of System
- W4 → Implementation + Sharing
- W5 → Implementation (Queues)
- W6 → Implementation (Load Testing)
- W7 → Technical Retrospective
- W8 → Complete System Implementation

## Bluesky: AT Protocol

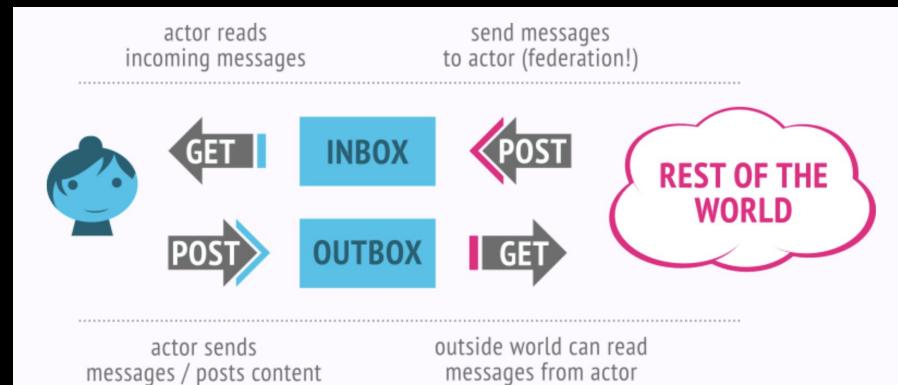
"The AT Protocol is an open, decentralized network for building social applications."



AT Protocol: <https://atproto.com/>

## Mastodon: ActivityPub

"ActivityPub is a decentralized social networking protocol based on the ActivityStreams 2.0 data format."



ActivityPub: <https://activitypub.rocks/>

## Assignment #2

1. Review 1 other technical design document written by others
  - Choose one that has decided on a different protocol from yours
  - Focus on understanding their whys, and insights from them
  - Where and how did their analysis differ from yours?
  - Are there potential Pros and Cons that they have missed out?
2. Build a quick and hacky PoC to create a Post on the platform you have chosen (share link to code repository)

**Due date: Sunday, 23:59 (Anywhere on Earth, UTC-12)**

- Which protocol will you use? Decide after reviewing the documents.
- For the PoC, just focus on creating a Post first; avoid doing any form of systems design until you have created a Post