

Data-Intensive Distributed Computing

CS 431/631 451/651 (Winter 2019)

Part 7: Mutable State (2/2) March 19, 2019

> Adam Roegiest Kira Systems

These slides are available at http://roegiest.com/bigdata-2019w/



This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States See http://creativecommons.org/licenses/by-nc-sa/3.0/us/ for details

The Fundamental Problem

We want to keep track of *mutable* state in a *scalable* manner

Assumptions:

State organized in terms of logical records State unlikely to fit on single machine, must be distributed

MapReduce won't do!

Motivating Scenarios

Money shouldn't be created or destroyed: Alice transfers \$100 to Bob and \$50 to Carol The total amount of money after the transfer should be the same

Phantom shopping cart:

Bob removes an item from his shopping cart... Item still remains in the shopping cart Bob refreshes the page a couple of times... item finally gone

Motivating Scenarios

People you don't want seeing your pictures: Alice removes mom from list of people who can view photos Alice posts embarrassing pictures from Spring Break Can mom see Alice's photo?

Why am I still getting messages?

Bob unsubscribes from mailing list and receives confirmation Message sent to mailing list right after unsubscribe Does Bob receive the message?

Three Core Ideas

Why do these scenarios happen?

Partitioning (sharding) To increase scalability and to decrease latency Need distributed transactions!

Replication

To increase robustness (availability) and to increase throughput Need replica coherence protocol!

Caching To reduce latency Need cache coherence protocol!

Source: Wikipedia (Cake)

To: All Graduate Students Due to a recent incident, we would like to remind all Grad Students that refreshments provided in communal areas during an event are for attendees of that event only. Please vacate the communal area and do not consume the refreshments unless you have been specifically invited to participate. To avoid any misunderstanding, you are only invited if you received a specific invitation by e-mail or if it was arranged by your supervisor for you to attend.

Thank you for your cooperation,

The Department Administrator



WWW.PHDCOMICS.COM

Morale of the story: there's no free lunch! (Everything is a tradeoff)

Three Core Ideas

Why do these scenarios happen?

Partitioning (sharding)

To increase scalability and to decrease latency Need distributed transactions!

Replication

To increase robustness (availability) and to increase throughput Need replica coherence protocol!

Caching To reduce latency Need cache coherence protocol!

Relational Databases

... to the rescue!

Source: images.wikia.com/batman/images/b/b1/Bat_Signal.jpg

How do RDBMSes do it?

Transactions on a single machine: (relatively) easy!

Partition tables to keep transactions on a single machine Example: partition by user

What about transactions that require multiple machines? Example: transactions involving multiple users

Solution: Two-Phase Commit







2PC: Assumptions and Limitations

Assumptions:

Persistent storage and write-ahead log at every node WAL is never permanently lost

Limitations:

It's blocking and slow What if the coordinator dies?

Three Core Ideas

Why do these scenarios happen?

Partitioning (sharding) To increase scalability and to decrease latency Need distributed transactions!

Replication

To increase robustness (availability) and to crease throughput

Need replica coherence protocol!

Caching

To reduce latency Need cache coherence protocol!

Replication possibilities

Update sent to a master Replication is synchronous Replication is asynchronous Combination of both

Okay, but if the master fails?

Update sent to an arbitrary replica

Replication is synchronous(?) Replication is asynchronous Combination of both

Distributed Consensus

More general problem: addresses replication and partitioning



Replication possibilities

Update sent to a master Replication is synchronous Replication is asynchronous Combination of both

Okay, but if the master fails?

Update sent to an arbitrary replica

Replication is synchronous(?) Replication is asynchronous Combination of both

Guaranteed consistency with a consensus protocol A buggy mess "Eventual Consistency"

CAP "Theorem" (Brewer, 2000)

Consistency Availability Partition tolerance

... pick two

CAP Tradeoffs

CA = consistency + availability E.g., parallel databases that use 2PC

AP = availability + tolerance to partitions E.g., DNS, web caching

Is this helpful?

CAP not really even a "theorem" because vague definitions More precise formulation came a few years later



Source: Abadi (2012) Consistency Tradeoffs in Modern Distributed Database System Design. IEEE Computer, 45(2):37-42

Abadi Says...

CP makes no sense!

CAP says, in the presence of P, choose A or C But you'd want to make this tradeoff even when there is no P

Fundamental tradeoff is between consistency and latency Not available = (very) long latency

Move over, CAP

PACELC ("pass-elk")

PAC

If there's a partition, do we choose A or C?

ELC

Otherwise, do we choose Latency or Consistency?

At the end of the day...

Guaranteed consistency with a consensus protocol A buggy mess "Eventual Consistency"

Sounds reasonable in theory... What about in practice? To: All Graduate Students Due to a recent incident, we would like to remind all Grad Students that refreshments provided in communal areas during an event are for attendees of that event only. Please vacate the communal area and do not consume the refreshments unless you have been specifically invited to participate. To avoid any misunderstanding, you are only invited if you received a specific invitation by e-mail or if it was arranged by your supervisor for you to attend.

Thank you for your cooperation,

The Department Administrator



WWW.PHDCOMICS.COM

Morale of the story: there's no free lunch! (Everything is a tradeoff)



Machine fails: What happens?

HBase



Three Core Ideas

Why do these scenarios happen?

Partitioning (sharding) To increase scalability and to decrease latency Need distributed transactions!

Replication

To increase robustness (availability) and to increase throughput

Need replica coherence protocol!

Caching

To reduce latency Need cache coherence protocol!

Facebook Architecture



Read path:

Look in memcached Look in MySQL Populate in memcached

Write path:

Write in MySQL Remove in memcached

Subsequent read:

Look in MySQL Populate in memcached

Facebook Architecture: Multi-DC



- 1. User updates first name from "Jason" to "Monkey".
- 2. Write "Monkey" in master DB in CA, delete memcached entry in CA and VA.
- 3. Someone goes to profile in Virginia, read VA replica DB, get "Jason".
- 4. Update VA memcache with first name as "Jason".
- 5. Replication catches up. "Jason" stuck in memcached until another write!

Facebook Architecture: Multi-DC



Solution: Piggyback on replication stream, tweak SQL

REPLACE INTO profile (`first_name`) VALUES ('Monkey') WHERE `user_id`='jsobel' MEMCACHE_DIRTY 'jsobel:first_name'

Source: www.facebook.com/note.php?note_id=23844338919

Three Core Ideas

Why do these scenarios happen?

Partitioning (sharding)

To increase scalability and to decrease latency Need distributed transactions!

Replication

To increase robustness (availability) and to increase throughput

Need replica coherence protocol!

Caching

To reduce latency Need cache coherence protocol!

Now imagine multiple datacenters...

What's different?

Source: Google

tl;dr -

Single row transactions Easy to implement, obvious limitations

Implement a global consensus protocol for every transaction Guarantee consistency, but slow

Eventual consistency Who knows?

tl;dr -

Single row transactions Easy to implement, obvious limitations

Implement a global consensus protocol for every transaction Guarantee consistency, but slow

Eventual consistency

Who knows?

Can we cheat a bit?

tl;dr -



Entity groups Groups of entities that share affinity Example: user + user's photos + user's posts etc.



But what if that's not enough?

Preserving commit order: example schema





Preserving commit order







Initial state	
T1@ts1	INSERT INTO ads VALUES (2, "elkhound puppies")
T2@ts2	INSERT INTO impressions VALUES (US, 2PM, 2)

Google's Spanner

Features:

Full ACID translations across multiple datacenters, across continents! External consistency (= linearizability):

system preserves happens-before relationship among transactions

How?

Given write transactions A and B, if A *happens-before* B, then timestamp(A) < timestamp(B)

TrueTime \rightarrow write timestamps





Why this works





TrueTime







What's the catch?

Source: The Matrix

Three Core Ideas

Partitioning (sharding) To increase scalability and to decrease latency Need distributed transactions!

Replication

To increase robustness (availability) and to increase throughput Need replica coherence protocol!

Caching To reduce latency Need cache coherence protocol!

Source: Wikipedia (Cake)

To: All Graduate Students Due to a recent incident, we would like to remind all Grad Students that refreshments provided in communal areas during an event are for attendees of that event only. Please vacate the communal area and do not consume the refreshments unless you have been specifically invited to participate. To avoid any misunderstanding, you are only invited if you received a specific invitation by e-mail or if it was arranged by your supervisor for you to attend.

Thank you for your cooperation,

The Department Administrator



WWW.PHDCOMICS.COM

Morale of the story: there's no free lunch! (Everything is a tradeoff)