

High-Performance Transactions

6.5830/6.5831 Lecture 19
Tianyu Li

Recap - Transactions

Single-node

- Disk-based
- 2PL
- Write-ahead Logging + Fuzzy Checkpoints

Recap - Transactions

Multi-node

- 2 PC for multi-node transactions
- Shared-nothing architecture. Use replication for high-availability.

Critique

- “Classical” DBMS matured in the 80s and 90s
- Hardware & workload was much different back then
- DBMS is about dealing with limitation of hardware
- Do the original assumptions still hold?

Critique

Back then: Network is slow

Now: Fiber optics can easily hit 100s of Gbps

Back then: The database machine exists in a basement

Now: The public cloud, global scale

Critique

Back then: Single (or a couple of) processor(s)

Now: AWS offers 96 core machines

Back then: RAM is limited

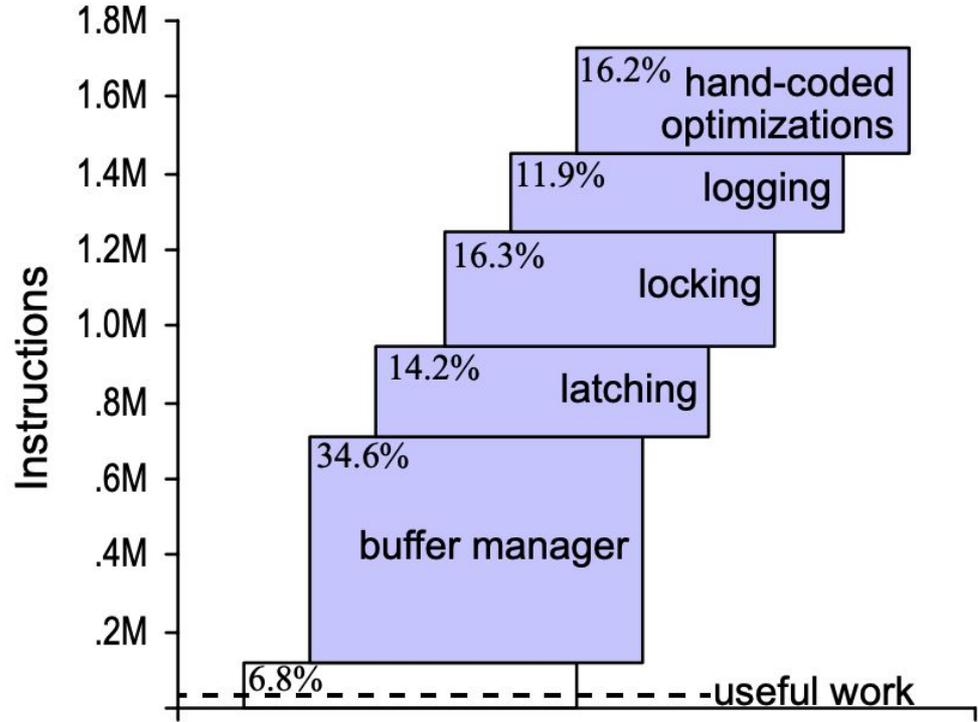
Now: Said 96 core machines has 384 GiB (!) of RAM

What happens now?

Can't we just take the DBMS, run it on faster hardware, and get better performance?

What happens now?

- Running old code on new hardware != speed-up
- New performance bottlenecks
- New architecture required to make use of faster hardware



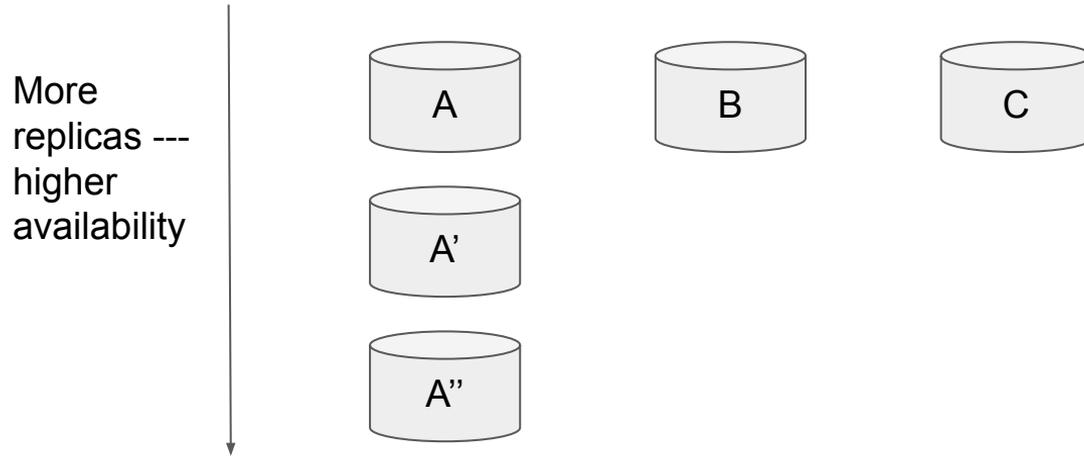
Stavros Harizopoulos, Daniel J. Abadi, Samuel Madden, and Michael Stonebraker. OLTP through the looking glass, and what we found there. SIGMOD 2008

Today -- High Performance Transactions

- Looking Back
- Multi-node
 - Bottleneck: 2-Phase Commit
 - Single-Site Execution
 - Deterministic Transactions
 - Epoch-based Coordination
- Looking Ahead

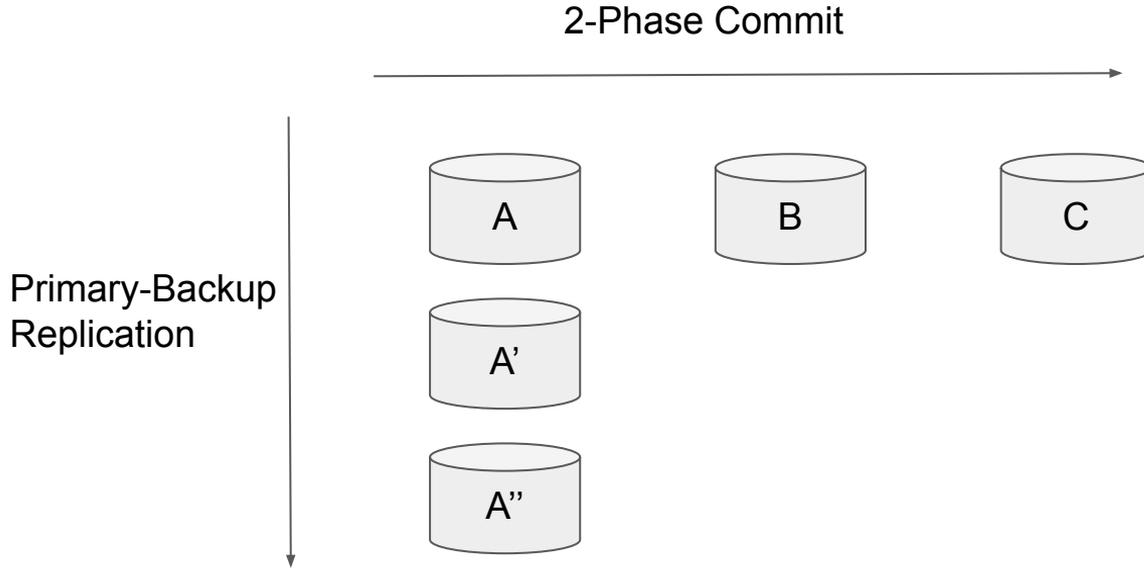
Recap: Scaling a Database

More shard/partitions --- more parallelism and throughput



* Replicas usually also serve read requests

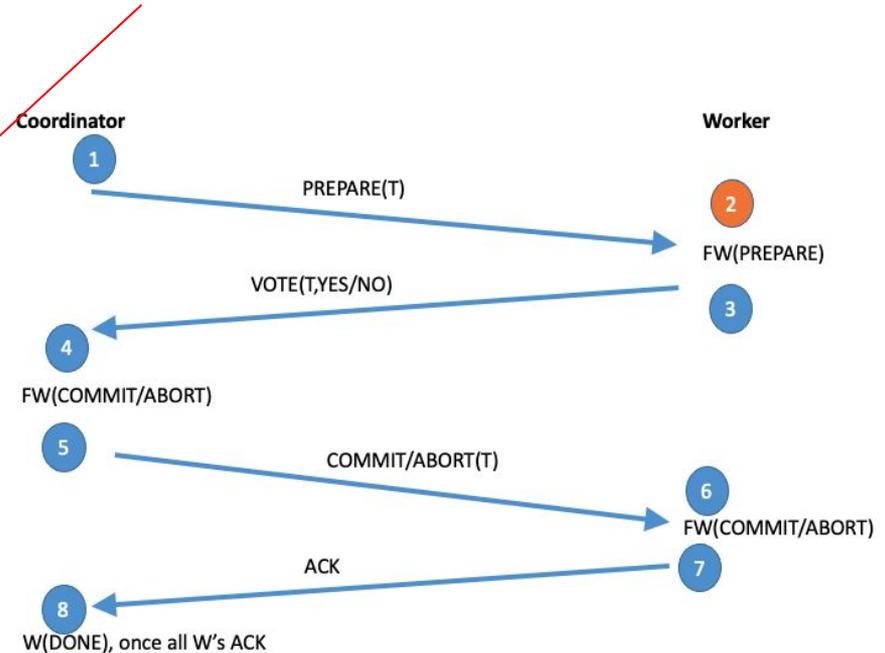
Recap: Scaling a Database



Recap: 2-Phase Commit

1. Log start of transaction
2. Execute transaction on worker nodes
3. PREPARE each worker
4. Log transaction commit if all OK
5. Commit each worker
6. Log Done

Commit Point



Critique: 2-Phase Commit

- 2 network round trips + synchronous logging
 - Worse still — likely need to hold locks throughout process
- 2PC blocks in coordinator fails, until the coordinator can be replaced or recovered
- 2PC basically sacrifices performance for strong guarantees

Example: Google Spanner

- A rare example of geo-distributed strongly consistent transactional system
 - You get the same guarantee as single-node
- Optimized for read-only transactions with TrueTime
- Optimized 2PC (on Paxos)

Spanner: Google's Globally-Distributed Database

James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, JJ Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, Wilson Hsieh, Sebastian Kanthak, Eugene Kogan, Hongyi Li, Alexander Lloyd, Sergey Melnik, David Mwaaura, David Nagle, Sean Quinlan, Rajesh Rao, Lindsay Rolig, Yasushi Saito, Michal Sczymaniak, Christopher Taylor, Ruth Wang, Dale Woodford

Google, Inc.

Abstract

Spanner is Google's scalable, multi-version, globally-distributed, and synchronously-replicated database. It is the first system to distribute data at global scale and support externally-consistent distributed transactions. This paper describes how Spanner is structured, its feature set, the rationale underlying various design decisions, and a novel time API that exposes clock uncertainty. This API and its implementation are critical to supporting external consistency and a variety of powerful features: non-blocking reads in the past, lock-free read-only transactions, and atomic schema changes, across all of Spanner.

1 Introduction

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 [21] state machines in datacenters spread all over the world. Replication is used for global availability and geographic locality; clients automatically failover between replicas. Spanner automatically reshards data across machines as the amount of data or the number of servers changes, and it automatically migrates data across machines (even across datacenters) to balance load and in response to failures. Spanner is designed to scale up to millions of machines across hundreds of datacenters and trillions of database rows.

Applications can use Spanner for high availability, even in the face of wide-area natural disasters, by replicating their data within or even across continents. Our initial customer was F1 [35], a rewrite of Google's advertising backend. F1 uses five replicas spread across the United States. Most other applications still backhau-

tency over higher availability, as long as they can survive 1 or 2 datacenter failures.

Spanner's main focus is managing cross-datacenter replicated data, but we have also spent a great deal of time in designing and implementing important database features on top of our distributed-systems infrastructure. Even though many projects happily use Bigtable [9], we have also consistently received complaints from users that Bigtable can be difficult to use for some kinds of applications: those that have complex, evolving schemas, or those that want strong consistency in the presence of wide-area replication. (Similar claims have been made by other authors [37].) Many applications at Google have chosen to use Megastore [5] because of its semi-relational data model and support for synchronous replication, despite its relatively poor write throughput. As a consequence, Spanner has evolved from a Bigtable-like versioned key-value store into a temporal multi-version database. Data is stored in schematized semi-relational tables; data is versioned, and each version is automatically timestamped with its commit time; old versions of data are subject to configurable garbage-collection policies; and applications can read data at old timestamps. Spanner supports general-purpose transactions, and provides a SQL-based query language.

As a globally-distributed database, Spanner provides several interesting features. First, the replication configurations for data can be dynamically controlled at a fine grain by applications. Applications can specify constraints to control which datacenters contain which data, how far data is from its users (to control read latency), how far replicas are from each other (to control write latency), and how many replicas are maintained to control durability, availability, and read performance. Data can also be dynamically and transparently moved by

Corbett et. al. Spanner: Google's Globally-Distributed Database. OSDI 2012

Problem

replicas	latency (ms)			throughput (Kops/sec)		
	write	read-only transaction	snapshot read	write	read-only transaction	snapshot read
1D	9.4±.6	—	—	4.0±.3	—	—
1	14.4±1.0	1.4±.1	1.3±.1	4.1±.05	10.9±.4	13.5±.1
3	13.9±.6	1.3±.1	1.2±.1	2.2±.5	13.8±3.2	38.5±.3
5	14.4±.4	1.4±.05	1.3±.04	2.8±.3	25.3±5.2	50.0±1.1

- Read-only transactions scale and perform well
- Read-write transactions not so much

Problem

2PC Scalability

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

2PC end-to-end Latency

operation	latency (ms)		count
	mean	std dev	
all reads	8.7	376.4	21.5B
single-site commit	72.3	112.8	31.2M
multi-site commit	103.0	52.2	32.1M

- 2PC is simply too expensive.

Question: Can we do better?

Aside: Why is this difficult?

Well-known theoretical limitations

- In short, you CANNOT have a “fast and reliable” distributed ACID system.
 - Two Generals Problem [Gray ‘78]
 - CAP Theorem [Brewer ‘00, Gilbert ‘02]
 - Coordination Avoidance in Database Systems [Bailis ‘15]

- We covered this last lecture
 - Many use cases regress to using “NoSQL” systems with more scalability but less guarantees

Why bother with distributed transactions then?

- Really powerful abstraction
- Extremely useful
- Impossibilities are mathematical. We are here to build systems*.

* Often called “NewSQL” systems

Attempt 1: H-Store

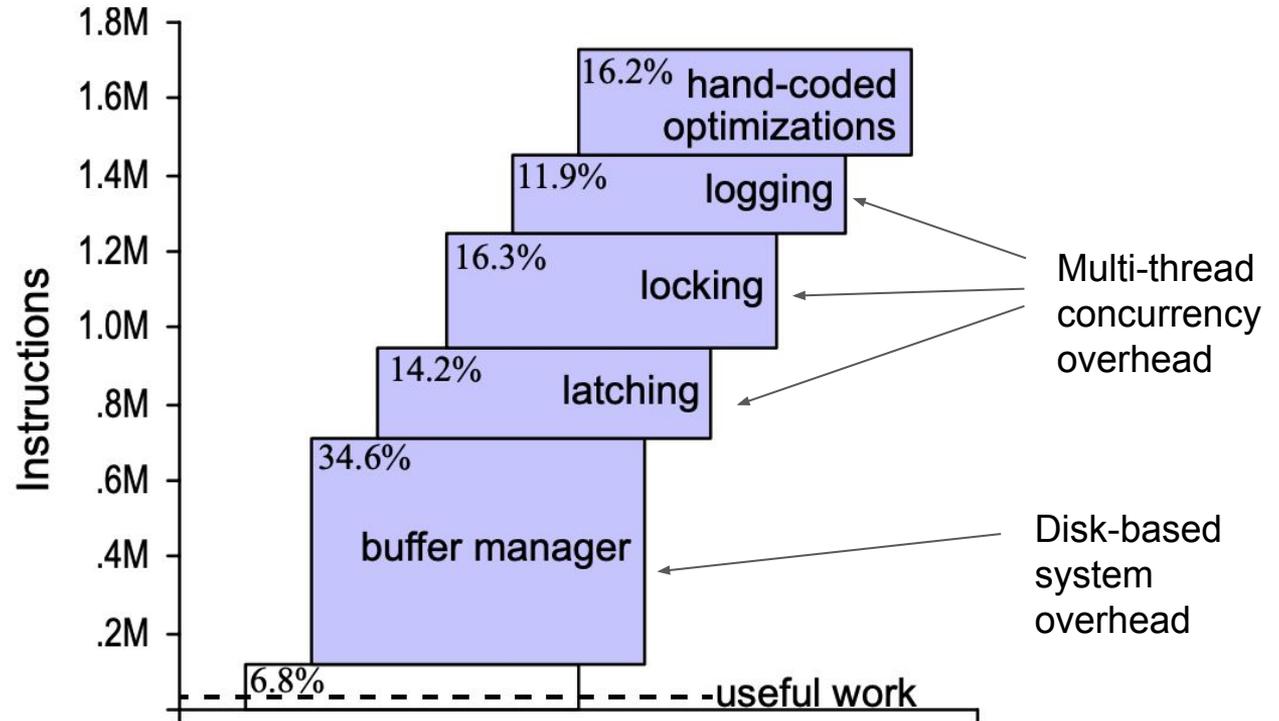
- Large collaborative project @ MIT (among other places) .
- Distributed & main-memory
- Commercialized as VoltDB

The logo for H-Store, featuring a green square with a white 'H' inside, followed by the text '-Store' in a bold, black, sans-serif font.

The logo for VoltDB, featuring a red graphic of three triangles pointing upwards, followed by the text 'VOLTDB' in a bold, red, sans-serif font.

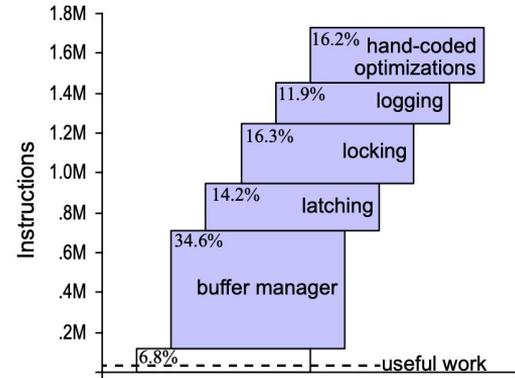
Key Idea

Recall:



Key Idea

- H-Store is a main-memory system
- H-Store partitions data, and executes single-threaded within each partition (site)



M. Stonebraker et. al.
The End of an Architectural Era
(It's Time for a Complete Rewrite). VLDB 2007.

**The End of an Architectural Era
(It's Time for a Complete Rewrite)**

<p style="text-align: center;">M. Stonebraker Samuel Madden Daniel J. Abadi Srinivas Aravamudan MIT CSAIL stonebr@csail.mit.edu</p>	<p style="text-align: center;">Nabil Hachem Microsoft Consulting LLC nhachem@microsoft.com</p>	<p style="text-align: center;">Pat Helariu Microsoft Corporation phelariu@microsoft.com</p>
--	---	--

ABSTRACT

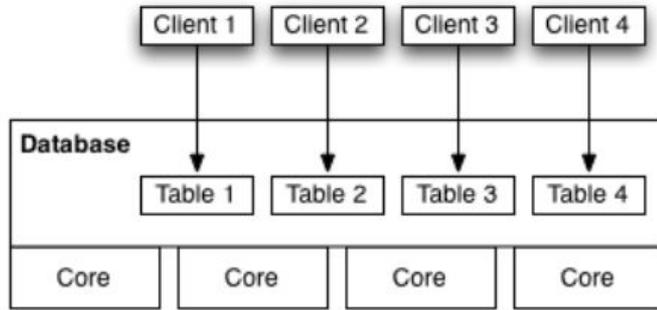
As the architectural paradigm that has dominated the database industry for more than 25 years ages, when hardware characteristics vary much more than when they did, the architectural paradigm of large-scale and general-purpose database systems is being challenged. This paper presents a new architectural paradigm that allows for more efficient use of hardware resources and provides a path to a complete rewrite of the database architecture. The new architecture is based on a set of principles that are designed to be hardware-agnostic and to be able to run on a wide range of hardware architectures. The new architecture is based on a set of principles that are designed to be hardware-agnostic and to be able to run on a wide range of hardware architectures. The new architecture is based on a set of principles that are designed to be hardware-agnostic and to be able to run on a wide range of hardware architectures.

1. INTRODUCTION

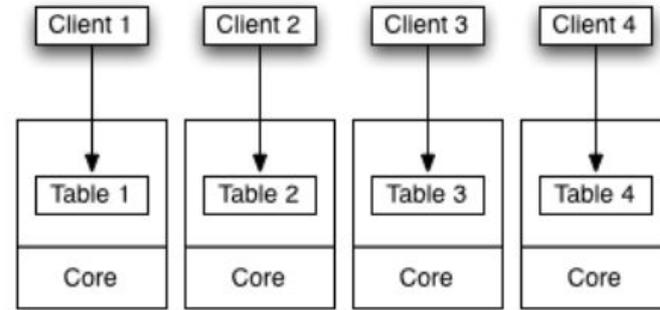
The single-processor database architecture that has dominated the database industry for more than 25 years is being challenged. This paper presents a new architectural paradigm that allows for more efficient use of hardware resources and provides a path to a complete rewrite of the database architecture. The new architecture is based on a set of principles that are designed to be hardware-agnostic and to be able to run on a wide range of hardware architectures.

- Data-oriented design and storage structures
- Multithreading & data sharing
- Locking based on memory consistency order
- Log-based recovery

Partition vs. Threads



Threads



Partitions

Synchronization
Overhead

No Concurrency
Control Required

Is this reasonable?

Specialized for OLTP (Online Transaction Processing) workload

- Transactions finish quickly
- Transactions are almost always point queries
- Transactions are known beforehand
- Working set fits in memory

Example: TPC-C

- Standardized benchmark used by everyone
- Models a warehouse order processing system
- Several types of transaction issued at random
- E.g. NewOrder Transaction:
 - Check item stock level
 - Create a new order
 - Update item stock level

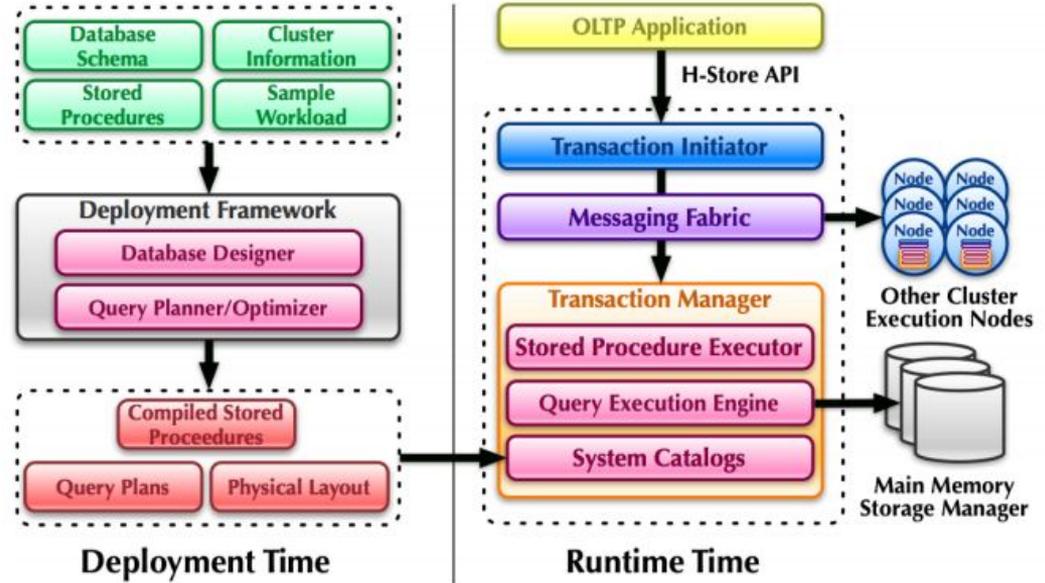
* Technically, TPC-C is strictly specified, but most benchmarks out there only loosely follow the specification.

Partitions

- Turns out, most OLTP workloads mostly partitionable
- TPC-C is about 90% partitionable
- Perform 2PC only for the remaining 10%

H-Store Architecture

- Stored Procedures Only
- Partitioned + Replicated
- Differentiates between:
 - Single-site Transactions
 - Others



Kallman et. al., H-Store: A High-Performance, Distributed Main Memory Transaction Processing System. VLDB 2008.

H-Store: Performance

- Vanilla H-Store can do 70K TPC-C txns compared to a couple thousand from before
- At the time, TPC-C record was about 133 K txn/s on a 128 core server.
 - H-Store can do half of that on low-end desktops.

H-Store: Further Optimizations

- Speculatively execute transactions when blocked
- Predict transaction behavior
- Partition database and schedule work intelligently to minimize percentage of distributed transactions

H-Store: Further Optimizations

- Speculatively execute transactions when blocked
- Predict transaction behavior
- Partition database and schedule work intelligently to minimize percentage of distributed transactions

H-Store: Speculative Execution

- Recall: H-Store single-threaded
- Also recall: 2PC takes > 10 ms to complete
- A partition simply waits out the 10ms instead of doing work

H-Store: Speculative Execution

- Observation: Most transactions succeed
- Idea: Assume transaction succeeds. Do useful work.
- Problem: introduces concurrency, but must not add overhead

Evan P.C. Jones, Daniel J Abadi, Samuel Madden. Low Overhead Concurrency Control for Partitioned Main Memory Databases. SIGMOD 2010

Low Overhead Concurrency Control for Partitioned Main Memory Databases

Evan P. C. Jones
MIT CSAIL
Cambridge, MA, USA
evanj@csail.mit.edu

Daniel J. Abadi
Yale University
New Haven, CT, USA
dab@cs.yale.edu

Samuel Madden
MIT CSAIL
Cambridge, MA, USA
madden@csail.mit.edu

ABSTRACT

Database partitioning is a technique for improving the performance of distributed OLTP databases. Most “single partition” transactions that access data on one partition do not need coordination with other partitions. For workloads that are amenable to partitioning, even when that transaction needs to access multiple data blocks, there is no need for concurrency at all. This strategy makes sense for a main memory database where there are no disk or network delays. Unfortunately, most OLTP applications have some transactions which access multiple partitions. This is the hard work of concurrency control. In this paper, we propose a low overhead concurrency control scheme that allows partitions to work on other transactions while other partitions do not. This is the main idea of a high-weight locking scheme, and the second is an even higher-weight type of speculative concurrency control that avoids the overhead of tracking reads and writes, but maintains performance even when concurrency is not needed. In general, the goal of our work is to use the techniques described here to provide a low overhead concurrency control scheme capable of handling a wide range of workloads. This is a modified TPC-C benchmark, open source concurrency control, and a highly correlated main memory database by up to a factor of five.

Categories and Subject Descriptors

D.1.2 (Database Management): Systems

General Terms

1. INTRODUCTION

Database partitioning is a technique to provide the benefits of replicated structures of transactions, with actually changing multiple transactions independently. However, several previous papers suggest that for some operational workloads, concurrency control is not necessary [1, 2]. In this paper, we show that in some settings and the workload pattern of transactions that can be executed without any delay, there is no need for concurrency control. In fact, we show that in some settings, concurrency control is not necessary to fully utilize a main memory CPU. Instead, we propose a low overhead concurrency control scheme that allows partitions to work on other transactions while other partitions do not. This is the main idea of a high-weight locking scheme, and the second is an even higher-weight type of speculative concurrency control that avoids the overhead of tracking reads and writes, but maintains performance even when concurrency is not needed. In general, the goal of our work is to use the techniques described here to provide a low overhead concurrency control scheme capable of handling a wide range of workloads. This is a modified TPC-C benchmark, open source concurrency control, and a highly correlated main memory database by up to a factor of five.

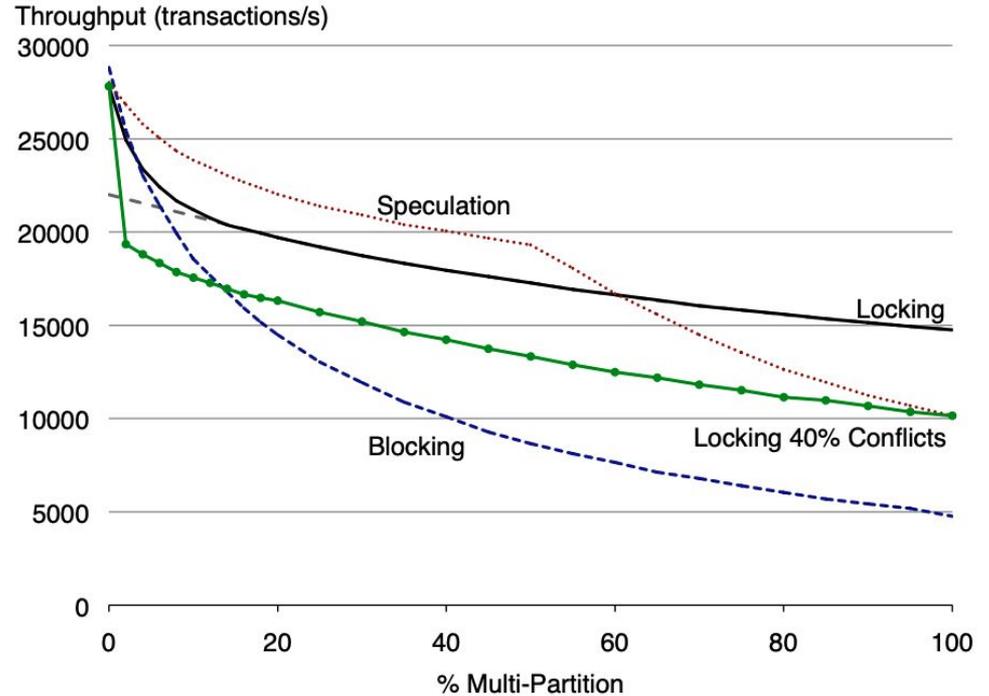
H-Store: Speculative Execution

- Idea 1: Only speculatively execute when waiting for 2PC
 - No locks required
 - Speculative results held back until 2PC finishes
 - Record undo information in-memory

- Idea 2: Speculate whenever stalled
 - E.g., when a multi-partition transaction is reading a remote value
 - Locking required
 - Overhead lower than 2PL, since no latching
 - Still expensive

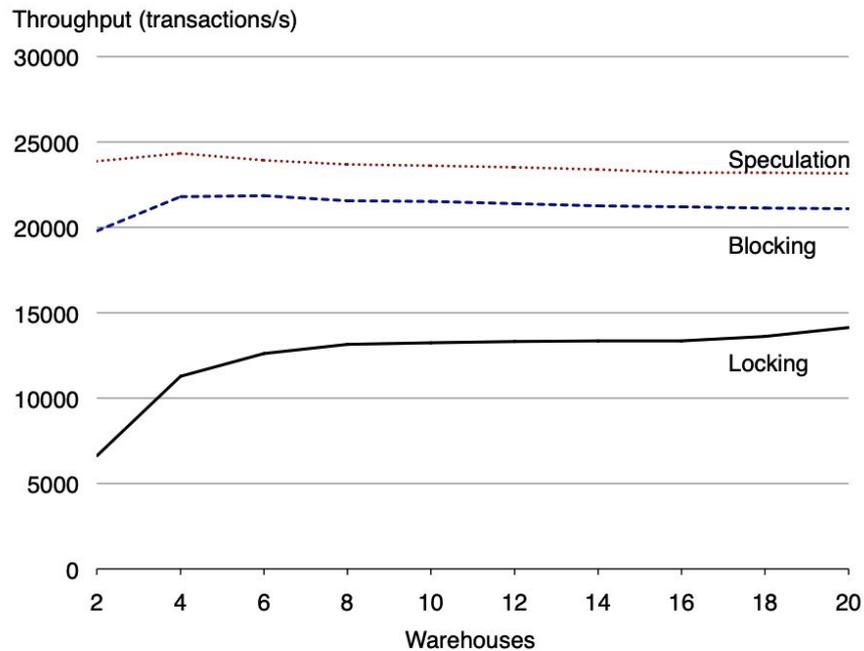
H-Store: Speculative Execution

- Synthetic benchmark ---
single operation
transactions
- Baseline no conflict



H-Store: Speculative Execution

- TPC-C
- Locking overhead increases with complex workload
- Speculation better

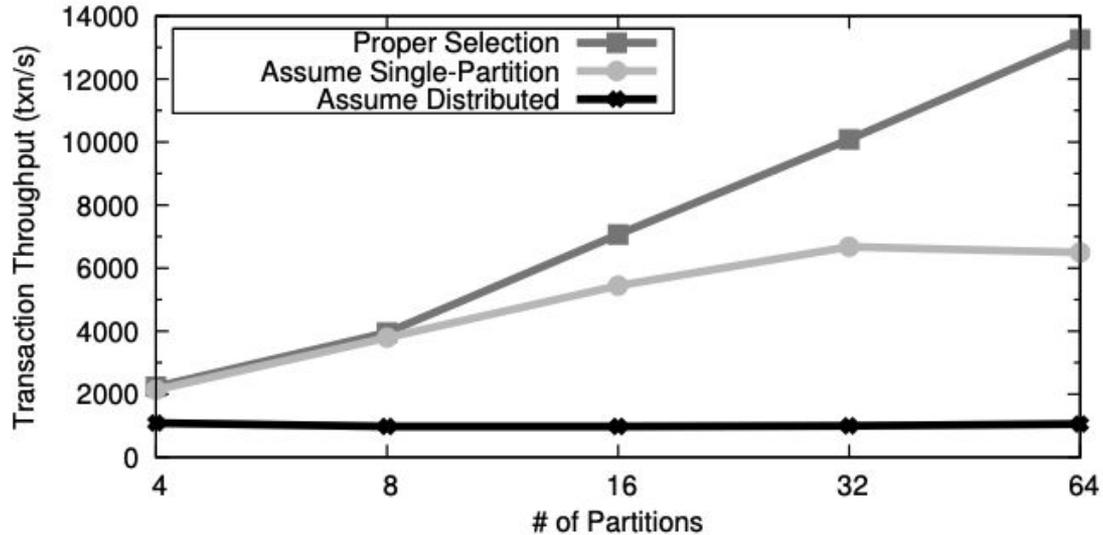


H-Store: Further Optimizations

- Speculatively execute transactions when blocked
- Predict transaction behavior
- Partition database and schedule work intelligently to minimize percentage of distributed transactions

H-Store: Predictive Modeling

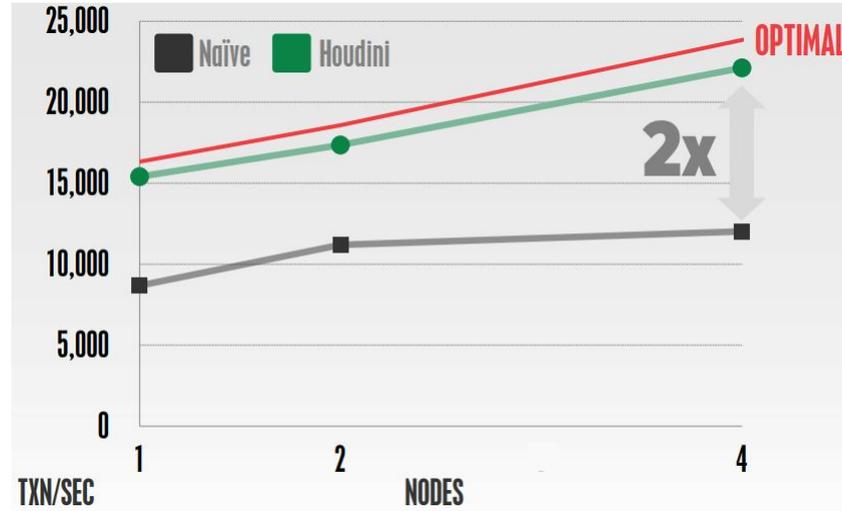
- Observe: many transaction optimizations available. None applies to all situations.



H-Store: Predictive Modeling

- Guessing game to apply the right optimization
- Question: Can we make educated guesses instead?

H-Store: Predictive Modeling



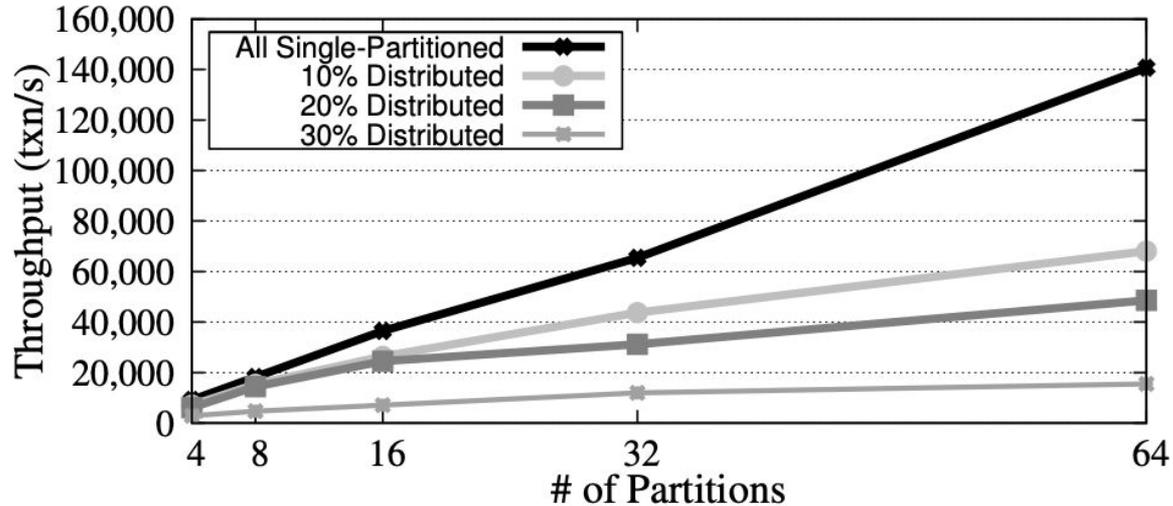
Transaction Behavior Prediction

H-Store: Further Optimizations

- Speculatively execute transactions when blocked
- Predict transaction behavior
- Partition database and schedule work intelligently to minimize percentage of distributed transactions

H-Store: Partitioning

- H-Store performance hinges on percentage of one-site txns
- Huge win if we can maximize one-site probability
- Intelligent partitioning required



H-Store: Partitioning

- Hiring someone to do partitioning is expensive and unreliable
- Can automate in H-Store with Large Neighborhood Search
- Details omitted

Skew-Aware Automatic Database Partitioning in Shared-Nothing, Parallel OLTP Systems

Andrew Pavlo
Brown University
pavlo@brown.edu

Carlo Curino
Google Research
krc@google.com

Stan Zdonik
Brown University
szdonik@brown.edu

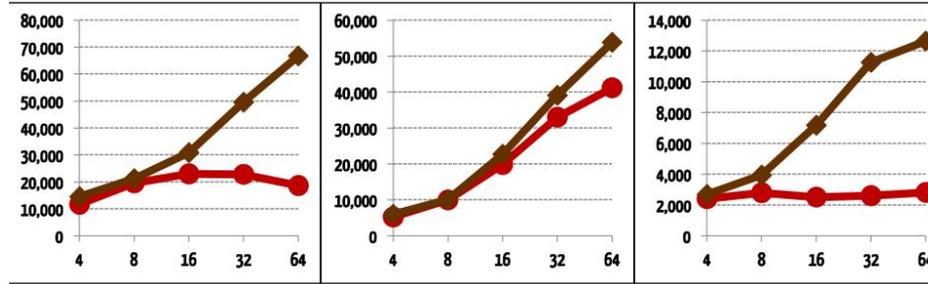
Andrew Pavlo, Carlo Curino, Stan Zdonik.
Skew-Aware Automatic Database
Partitioning in Shared-Nothing, Parallel
OLTP Systems. SIGMOD 2012.

ABSTRACT
The advent of commodity-based multi-terabyte computing systems has led to the development of parallel database management systems (DBMSs) for online transaction processing (OLTP) applications that scale when executing multi-tenant workloads. In the partitioning of these DBMSs, a predominant concern is the maintenance of an optimal database design that minimizes the number of partitions and the size of each partition. However, the cost of most partitioning algorithms is exponential in the number of partitions, which makes them infeasible for large numbers of partitions. In this paper, we present a novel approach to automatic database partitioning that addresses this problem. Our approach is based on the idea of the large neighborhood search (LNS) algorithm for combinatorial optimization problems. We describe the design of our LNS-based partitioning algorithm and its implementation in the context of a shared-nothing, parallel OLTP system. We evaluate our approach using a real-world OLTP workload and compare our results to those of a state-of-the-art partitioning algorithm. Our results show that our approach is able to find partitions that are significantly smaller than those of the state-of-the-art algorithm, and that our approach is able to find partitions that are significantly smaller than those of the state-of-the-art algorithm. Our approach is able to find partitions that are significantly smaller than those of the state-of-the-art algorithm, and that our approach is able to find partitions that are significantly smaller than those of the state-of-the-art algorithm.

Categories and Subject Descriptors
D.1.2 Database Management: Physical Design

Keywords
OLTP, Parallel, Shared Nothing, H-Store, KDD, Stored Procedures

H-Store: Partitioning



TATP
+88%

TPC-C
+16%

TPC-C Skewed
+183%

Optimized Partition

Takeaway

- The Good
 - Specialization is good
 - Partitioning is really powerful
 - Seemingly simple optimizations can lead to huge speed-ups
- The Not-so-good
 - If you really need distributed transactions, you are out of luck*

* People are still working on it!

Xinjing Zhou, Xiangyao Yu, Goetz Graefe,
Michael Stonebraker. Lotus: Scalable
Multi-Partition Transactions on
Single-Threaded Partitioned Databases.
VLDB 2022.



Attempt 2: Calvin / Aria

- Why is H-Store faster without concurrency?
- No non-determinism from threading
 - Limits cross thread/node coordination need
 - Coordination often a bottleneck
- Can the same idea be applied to truly distributed transactions?

Key Idea: Calvin

- Have a global *deterministic* ordering of transaction execution.
- Take the input and execute anywhere. Get the same result.

Calvin: Fast Distributed Transactions for Partitioned Database Systems

Alexander Thomson Yale University thomson@cs.yale.edu	Thaddeus Diamond Yale University diamond@cs.yale.edu	Shu-Chun Weng Yale University scweng@cs.yale.edu
Kun Ren Yale University kun@cs.yale.edu	Philip Shao Yale University shao-philip@cs.yale.edu	Daniel J. Abadi Yale University dna@cs.yale.edu

ABSTRACT

Many distributed storage systems achieve high data access throughput via partitioning and replication, each system with its own advantages and tradeoffs. In order to achieve high scalability, however, today's systems generally reduce transactional support, disallowing single transactions from spanning multiple partitions. Calvin is a practical transaction scheduling and data replication layer that uses a deterministic ordering guarantee to significantly reduce the normally prohibitive contention costs associated with distributed transactions. Unlike previous deterministic database system prototypes, Calvin supports disk-based storage, scales near-linearly on a cluster of commodity machines, and has no single point of failure. By replicating transaction inputs rather than effects, Calvin is also able to support multiple consistency levels—including Posix-based strong consistency across geographically distant replicas—at no cost to transactional throughput.

Categories and Subject Descriptors

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

General Terms

Algorithms, Design, Performance, Reliability

1. BACKGROUND AND INTRODUCTION

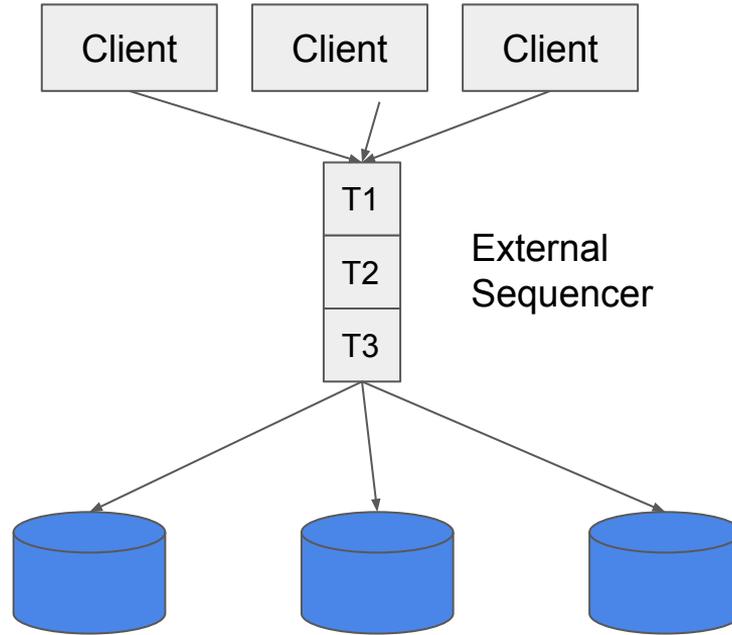
One of several current trends in distributed database system design is a move away from supporting traditional ACID database transactions. Some systems, such as Amazon's Dynamo [13], MongoDB [24], CouchDB [6], and Cassandra [17] provide no transactional support whatsoever. Others provide only limited transactionality, such as single-row transactional updates (e.g. Bigtable [11]) or transactions whose accesses are limited to small subsets of a database (e.g. Acute [9], Magnetite [7], and the Oracle NoSQL Database [26]). The primary reason that each of these systems does not support fully ACID transactions is to provide linear outward scalability. Other systems (e.g. VoltDB [27, 16]) support full ACID, but cease (or limit) concurrent transaction execution when processing a transaction that accesses data spanning multiple partitions.

Reducing transactional support greatly simplifies the task of building linearly scalable distributed storage solutions that are designed to serve “un embarrassingly partitionable” applications. For applications that are not easily partitionable, however, the burden of ensuring atomicity and isolation is generally left to the application programmer, resulting in increased code complexity, slower application development, and low performance client-side transaction scheduling.

Calvin is designed to run alongside a non-transactional storage system, transforming it into a shared-nothing (near-linearly scalable database system that provides high-availability) and full ACID transactions. These transactions can potentially span multiple parti-

Alexander Thomson et. al. Calvin: Fast Distributed Transactions for Partitioned Database Systems. SIGMOD 2012.

Deterministic Transactions



Deterministic Transactions

- Observe: this is not so different from serializable, where execution is equivalent to a serial schedule
- However: Calvin fixes the schedule **before** execution
- Therefore: coordination also largely done **before** execution

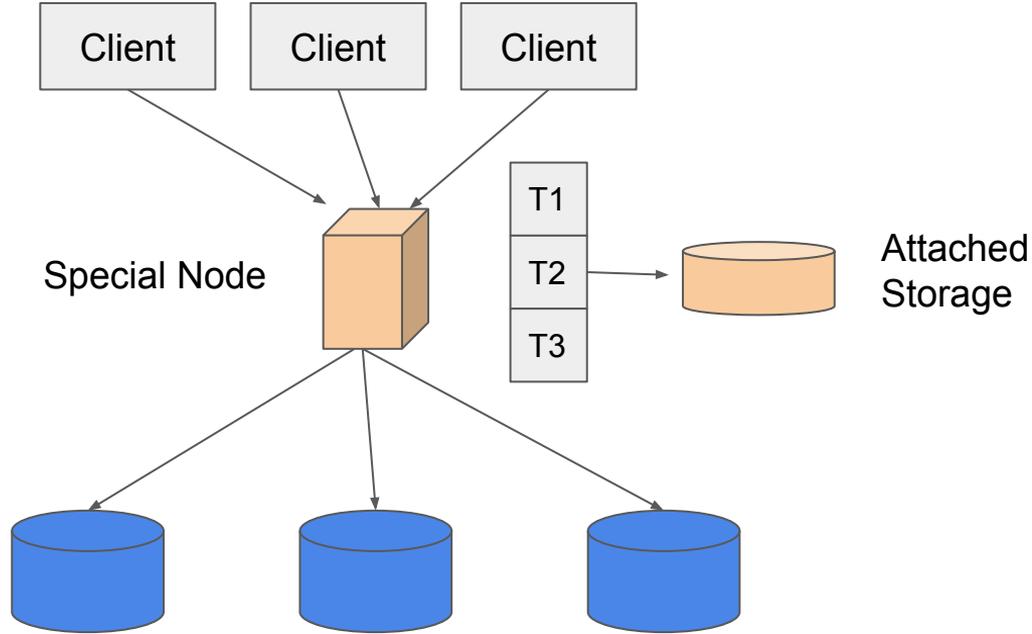
Practical Considerations

- Sequencer is a bottleneck of the system and single-point of failure
- We still want concurrency for performance on a single node
- We still need to track progress of replicas to give guarantees

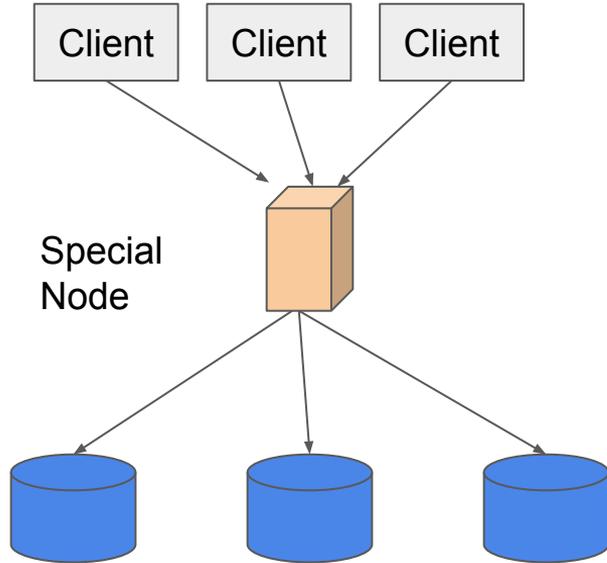
Practical Considerations

- Sequencer is a bottleneck of the system and single-point of failure
- We still want concurrency for performance on a single node
- We still need to track progress of replicas to give guarantees

Sequencer: Initial Attempt

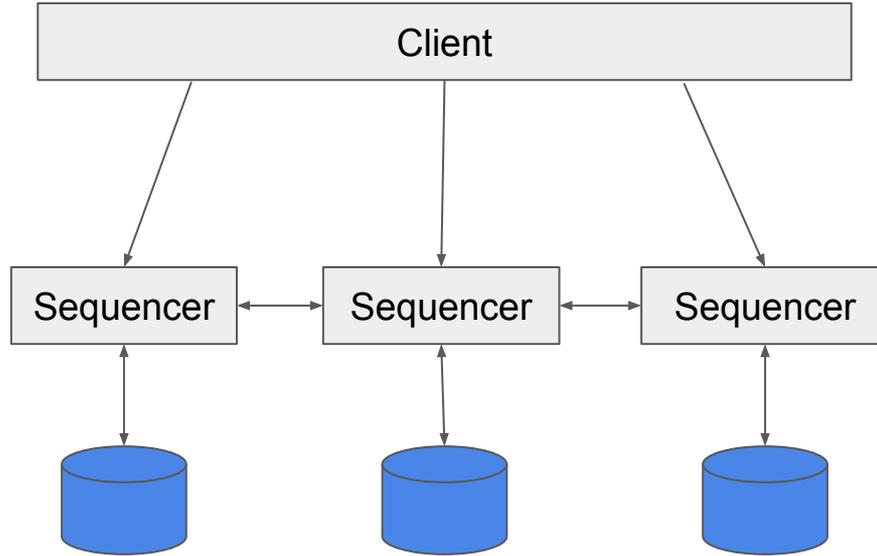


Sequencer: Initial Attempt



- Special node failure difficult to handle
- Txn throughput bottlenecked by special node throughput

Distributed Sequencer



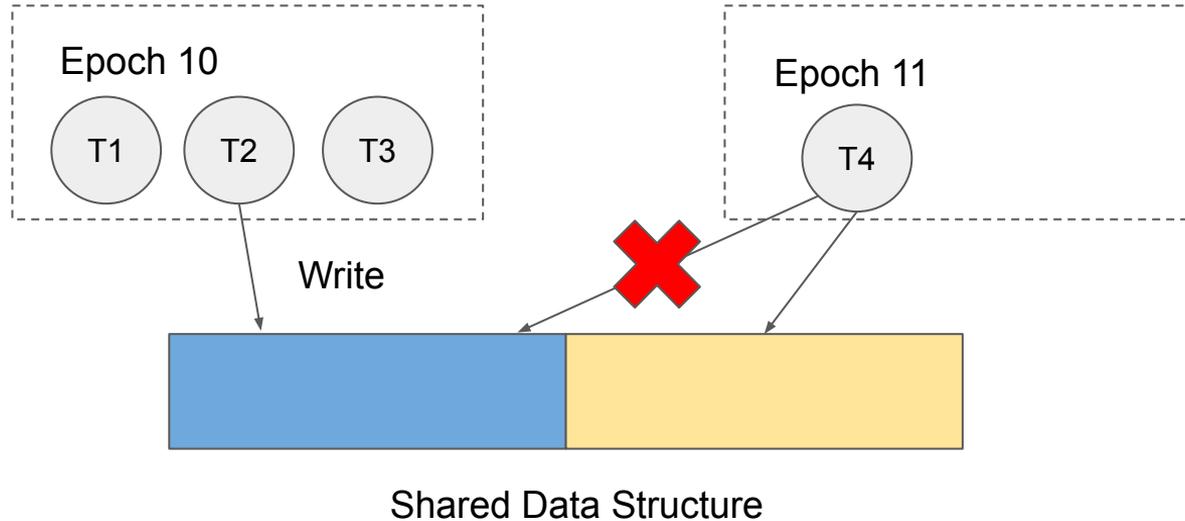
- Don't synchronize for every request
- Each sequencer collects a batch of requests
- Periodically replicate / persist and exchange batches

Key Idea: Epochs

- Recall: Coordination off critical path = win
- Idea: Synchronize loosely at set intervals (i.e., epochs)
- Wait for next epoch if uncertain
- Common idea in concurrent programming. Also used in single-node systems (e.g., Silo, Bw-Tree)

Example: Epoch

- Concurrently updating a buffer in-place
- Periodically “freeze” a log chunk to flush to disk



Example: Epoch-based Sequencer



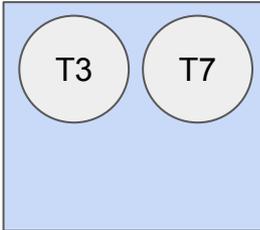
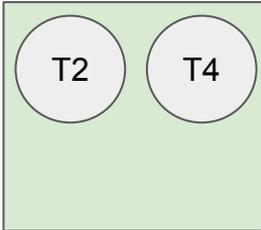
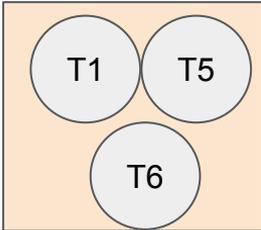
Epoch: 12

Global Txn Ordering

Sequencer 1

Sequencer 2

Sequencer 3



Seq 1 , Epoch 11
1, 5, 6

Seq 2, Epoch 11
2, 4

Seq 3, Epoch 11
3, 7

Example: Epoch-based Sequencer



Global Txn Ordering

Sequencer 1

Sequencer 2

Sequencer 3

Seq 1 , Epoch 11
1, 5, 6

Seq 2, Epoch 11
2, 4

Seq 3, Epoch 11
3, 7

- Now: Global ordering can be obtained through round-robin

What's the price?

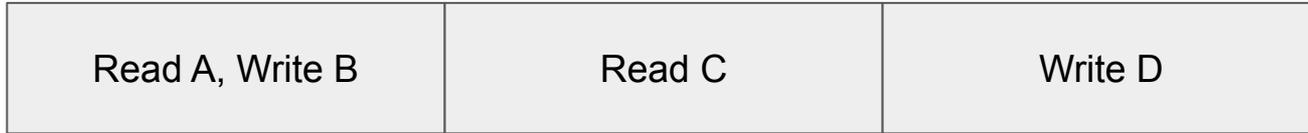
- Transactions are not sequenced until epoch end
- Recurring theme for epoch-based schemes: throughput vs. latency trade-off
- More on this later

Practical Considerations

- Sequencer is a bottleneck of the system and single-point of failure
- We still want concurrency for performance on a single node
- We still need to track progress of replicas to give guarantees

Scheduler: Deterministic Concurrency Control

Consider Schedule:



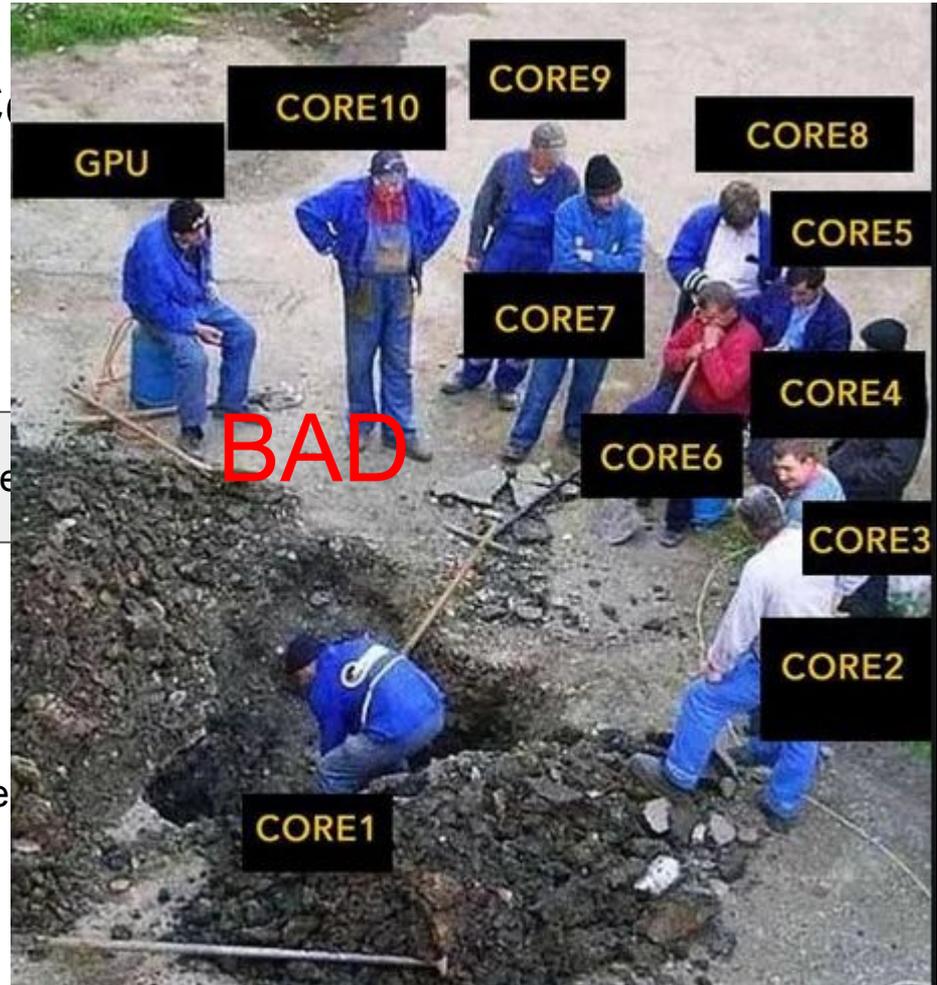
- No actual conflict
- No reason to execute in-order
- Challenge: concurrent execution that preserves deterministic schedule

Scheduler: Deterministic C

Consider Schedule:

Read A, Write B	Re
-----------------	----

- No actual conflict
- No reason to execute in-order
- Challenge: concurrent execution that prese



Scheduler: Deterministic Concurrency Control

- Need to allow for concurrent execution

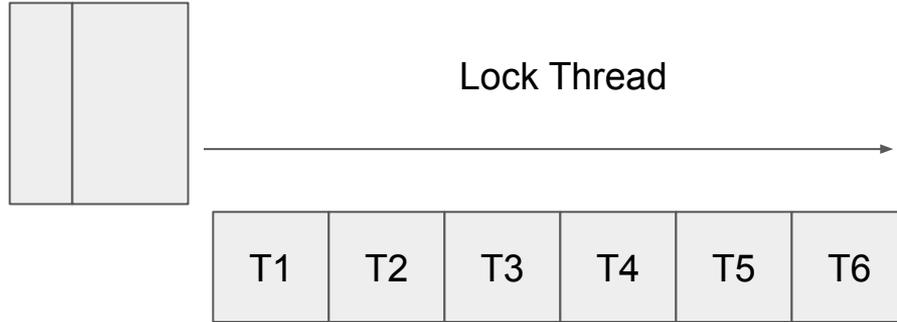
- However, concurrent execution has to follow predetermined schedule

Scheduler: Deterministic Concurrency Control

- Similar to 2 PL
- Allow arbitrary concurrent execution permitted by lock manager
- However, control how locks are granted

Scheduler: Deterministic Concurrency Control

Lock Table



Deterministic Schedule

- Don't **request** locks, **grant** locks.
- Dedicated lock thread assigns locks strictly in predetermined order
- Transaction executes when all locks are granted
- Assumption: read/write set known / can be determined before execution

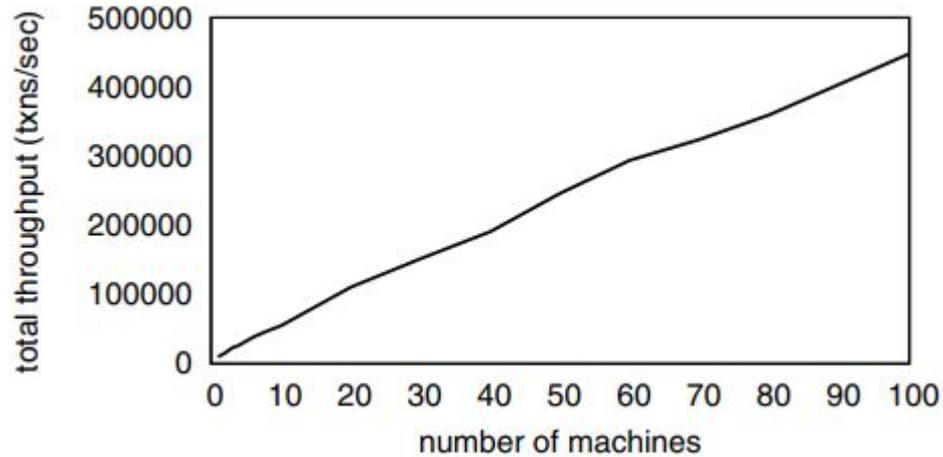
Practical Considerations

- Sequencer is a bottleneck of the system and single-point of failure
- We still want concurrency for performance on a single node
- We still need to track progress of replicas to give guarantees

Multi-node Transactions: Idea

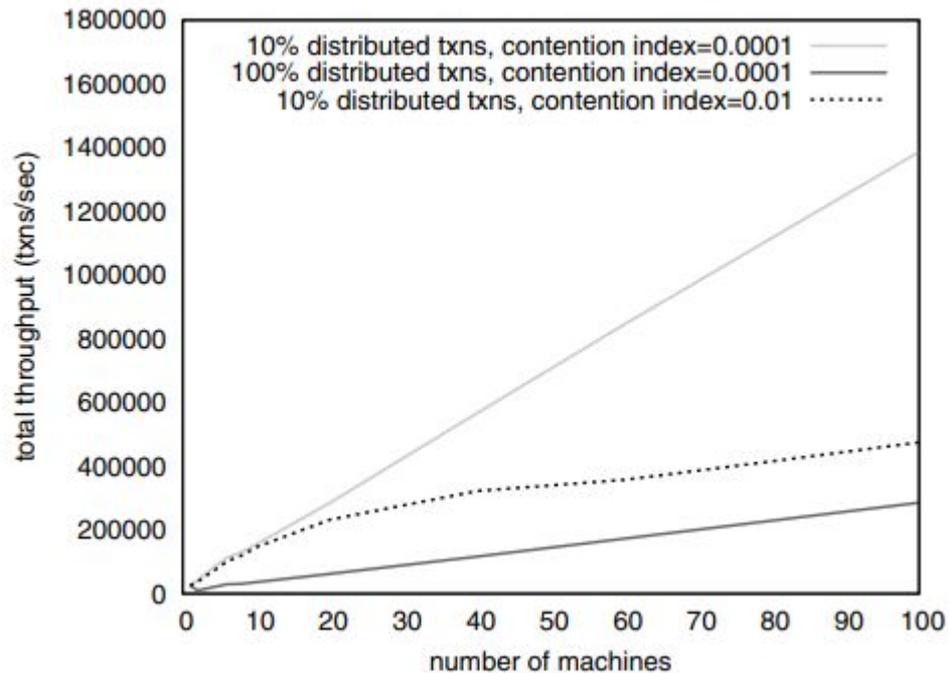
- Read remotely if needed
 - No overhead
- Write locally
- Multi-node transaction is done when every participant done with local writes
 - Use locks/node progress along the serial order to ensure correct reads

Calvin: Results

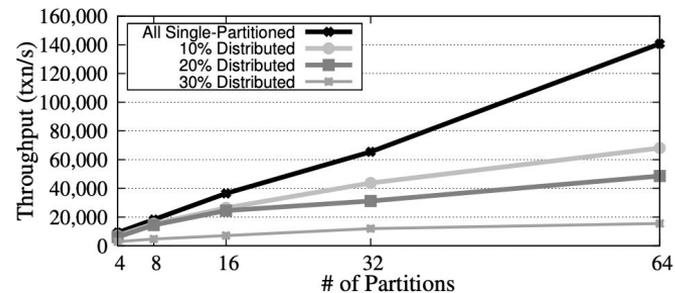


- TPC-C (100% New Order)

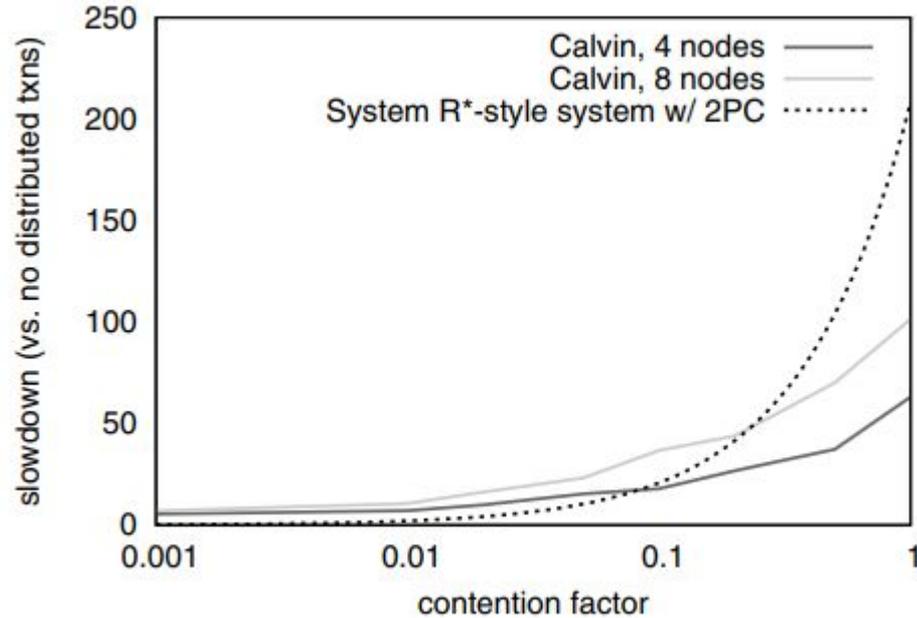
Calvin: Results



- Synthetic Microbenchmark



Calvin: Results



- 100 % Multi-partition
- Y-axis is slow down factor

Calvin: Criticism

- Transaction read/write sets must be known beforehand
- Not always practical

Aria: Practical Deterministic OLTP

- Relaxes the requirement to know r/w sets beforehand
- Speculatively execute first, repair later
- Details omitted

Yi Lu, Xiangyao Yu, Lei Cao, Samuel Madden
Madden. Aria: A Fast and Practical
Deterministic OLTP Database. VLDB 2020.

Aria: A Fast and Practical Deterministic OLTP Database

Yi Lu¹, Xiangyao Yu², Lei Cao¹, Samuel Madden¹
¹Massachusetts Institute of Technology, Cambridge, MA, USA
²University of Wisconsin-Madison, Madison, WI, USA
{yli,leo,madden}@csail.mit.edu, yyyc@facstaff.wisc.edu

ABSTRACT

Deterministic databases are able to efficiently run transactions across different replicas without coordination. However, existing state-of-the-art deterministic databases require that transaction read/write sets are known before execution, making such systems impractical in many OLTP applications. In this paper, we present Aria, a new distributed and deterministic OLTP database that does not have this limitation. The key idea behind Aria is that it first executes a batch of transactions against the same database snapshot in an execution phase, and then deterministically (without communication between replicas) chooses those that should commit to ensure serializability in a commit phase. We also propose a novel deterministic scheduling mechanism that allows Aria to order transactions in a way that reduces the number of conflicts. Our experiments on a cluster of eight nodes show that Aria outperforms systems with conventional nondeterministic concurrency control algorithms and the state-of-the-art deterministic databases by up to a factor of two on two popular benchmarks (YCSB and TPC-C).

VLDB Reference Format:

Yi Lu, Xiangyao Yu, Lei Cao and Samuel Madden. Aria: A Fast and Practical Deterministic OLTP Database. VLDB, 2020. 2057–2060, 2020.
DOI: <https://doi.org/10.1145/3479763.3487808>

1. INTRODUCTION

Modern database systems employ replication for high availability and data partitioning for scale-out. Replication allows systems to provide high availability, i.e., tolerant to machine failures, but also incurs additional network round trips to ensure writes are synchronized to replicas. Partitioning across several nodes allows systems to scale to larger databases. However, most implementations require the use of two-phase commit (2PC) [7] to address the issues caused by nondeterministic events such as system failures and race conditions in concurrency control. This introduces additional

latency to distributed transactions and implies availability and scalability (e.g., due to coordinator failures). Deterministic concurrency control algorithms [8, 15, 31, 52] provide a new way of building distributed and highly available database systems. They avoid the use of expensive commit and replication protocols by ensuring different replicas always independently produce the same results as long as the same input transactions are given. Therefore, rather than replicating and synchronizing the updates of distributed transactions, deterministic databases only have to replicate the input transactions across different replicas, which can be done asynchronously and often with much less communication. In addition, deterministic databases avoid the use of two-phase commit, since they naturally eliminate nondeterministic race conditions in concurrency control and are able to recover from system failures by re-executing the same original input transactions.

The state-of-the-art deterministic databases, DORM [16], PAVV [18], and Calyx [52], achieve determinism through dependency graphs or ordered locks. The key idea in DORM and PAVV is that a dependency graph is built from a batch of input transactions based on the read/write sets. In this way, the database can produce deterministic results as long as the transactions are run following the dependency graph. The key idea in Calyx is that read/write locks are acquired prior to executing the transactions, and according to the ordering of input transactions. A transaction is assigned to a worker thread for execution once all needed locks are granted. As shown in the left side of Figure 1, existing deterministic databases perform dependency analysis before transaction execution, which requires that the read/write set of a transaction be known a priori. For very simple transactions, e.g., that only access records through simple lookups on a primary key, this can be done easily. However, in reality, many transactions access records through complex predicates over many attributes; for such queries, these systems must execute the query at least twice: once to determine the read/write set, once to execute the query, and possibly more times if the pre-determined read/write set changes between these two executions. In addition, Calyx requires the use of a single-threaded lock manager per database partition, which

Takeaways

- Determinism can be a good thing
- Distributed coordination off the critical path = win

Attempt 3: COCO

- H-Store leveraged workload patterns
- Calvin/Aria introduces sequencer and increases latency
- Can we attack the problem of 2PC head-on?
 - Surprisingly, yes. Well, sometimes.
 - Area of (very) recent work

COCO

- Multiple transactions perform 2PC together at epoch granularity
- Failures result in all transactions within an epoch to fail together (does not include conflict aborts or user aborts)
- Replicas kept up-to-date at epoch boundary

Yi Lu, Xiangyao Yu, Lei Cao, Samuel Madden.
Epoch-based Commit and Replication in Distributed
OLTP Databases. VLDB 2021

Epoch-based Commit and Replication in Distributed OLTP Databases

Yi Lu
Massachusetts Institute of Technology
ylu@csail.mit.edu

Lei Cao
Massachusetts Institute of Technology
lcao@csail.mit.edu

Xiangyao Yu
University of Wisconsin-Madison
xyy@cs.wisc.edu

Samuel Madden
Massachusetts Institute of Technology
madden@csail.mit.edu

ABSTRACT

Many modern data-oriented applications are built on top of distributed OLTP databases for both scalability and high availability. Such distributed databases enforce atomicity, durability, and consistency through two-phase commit (2PC) and synchronous replication at the granularity of every single transaction. In this paper, we present COCO, a new distributed OLTP database that supports epoch-based commit and replication. The key idea behind COCO is that it separates transactions into epochs and treats a whole epoch of transactions as the commit unit. In this way, the overhead of 2PC and synchronous replication is significantly reduced. We support two variants of optimistic concurrency control (OCC) using physical time and logical time with various optimizations, which are enabled by the epoch-based execution. Our evaluation on two popular benchmarks (YCSB and TPC-C) show that COCO outperforms systems with fine-grained 2PC and synchronous replication by up to a factor of four.

VLDB Reference Format:

Yi Lu, Xiangyao Yu, Lei Cao, and Samuel Madden. Epoch-based Commit and Replication in Distributed OLTP Databases. VLDB, 14(3): 743–756, 2021.
doi:10.14778/3440095.3440098

1 INTRODUCTION

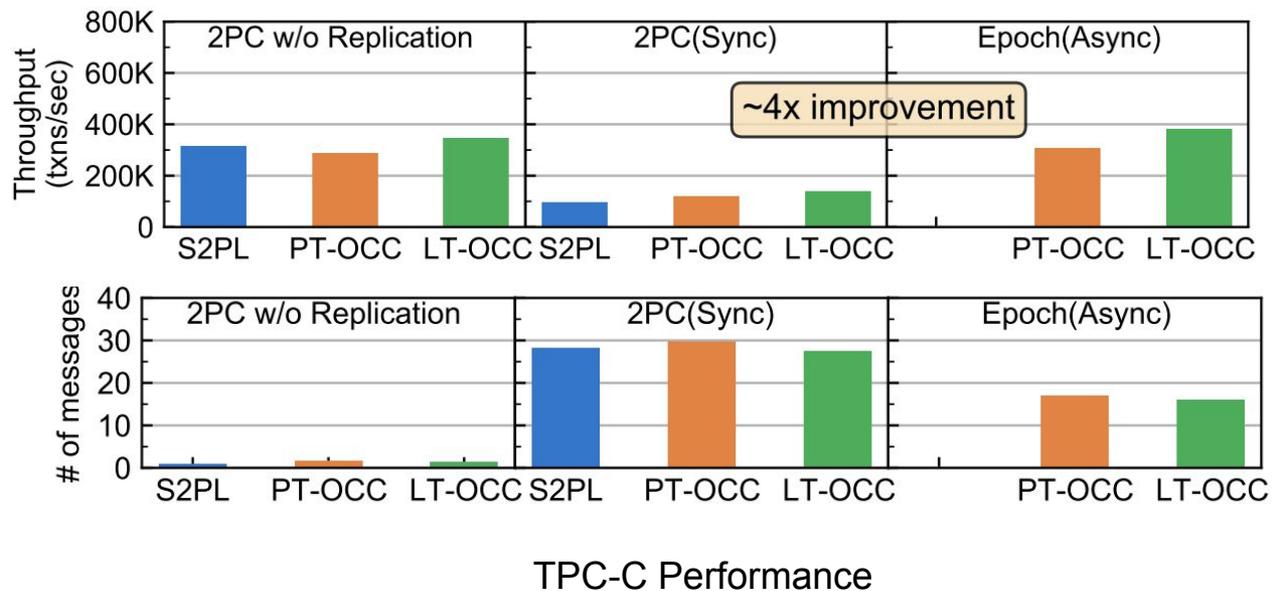
Many modern distributed OLTP databases use a shared-nothing

effect of committed transactions recorded to persistent storage and survives server failures.

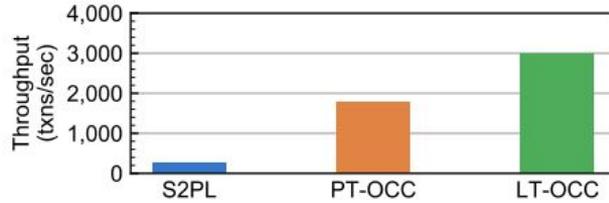
It is well known that 2PC causes significant performance degradation in distributed databases [3, 12, 46, 61], because a transaction is not allowed to release locks until the second phase of the protocol, blocking other transactions and reducing the level of concurrency [21]. In addition, 2PC requires two network round-trip delays and two sequential durable writes for every distributed transaction, making it a major bottleneck in many distributed transaction processing systems [21]. Although there have been some efforts to eliminate distributed transactions or 2PC, unfortunately, existing solutions either introduce impractical assumptions (e.g., the read/write set of each transaction has to be known a priori in deterministic databases) [5, 64, 62] or significant runtime overhead (e.g., dynamic data partitioning [12, 31]).

In addition, a desirable property of any distributed database is high availability, i.e., when a server fails, the system can mask the failure from end users by replacing the failed server with a standby machine. High availability is typically implemented using data replication, where all writes are handled at the primary replica and are shipped to the backup replicas. Conventional high availability protocols must make a tradeoff between performance and consistency. On one hand, asynchronous replication allows a transaction to commit once its writes arrive at the primary replica; propagation to backup replicas happens in the background asynchronously [14]. Transactions can achieve high performance but

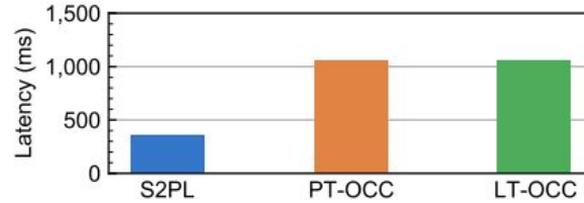
COCO: Results



COCO: Results over WAN



TPC-C Throughput



TPC-C P99 Latency

- Recall: Throughput vs. Latency trade-off

Epochs: Takeaways

- Surprisingly effective in alleviating bottleneck
- Throughput vs. Latency Trade-off
- General pattern in high-performance concurrent/distributed programming

Tianyu Li, Badrish Chandramouli, Samuel Madden.
Performant Almost-Latch-Free Data Structures
Using Epoch Protection. DaMoN 2022



What have we achieved?

- A class of new transactional systems (aka. NewSQL) that retains the strong guarantees of traditional relational DBMS, while being much more scalable and performant like NoSQL systems
 - These systems are largely main-memory systems
 - These system optimize around partitioning and sharding for performance
 - These system feature new, interesting concurrency control / distributed commit schemes
- And we have the performance numbers to show off
 - Exciting.*
 - Txn throughput went from a couple of thousands to millions per second
 - Current record holder for TPC-C does 707 M TpmC
 - OceanBase from Alibaba's Ant Financial

* Yes, call us nerds.

Criticism

We Are Boring

Sam Madden
madden@csail.mit.edu

AI is enjoying a renaissance, with popular press and major corporations building a variety of smart, AI-based applications, from self-driving cars to household robots to household gadgets that learn our behaviors and

Despite all of these applications revolving around data, the database community has content to cede these domains to our AI colleagues. This is absurdly the world-wide web, and (nearly) big data, we risk being an also-ran in computer science in the coming decade. These smart systems will work, and play, and the database community ought to be thinking

What Are We Doing With Our Lives? Nobody Cares About Our Concurrency Control Research

Andrew Pavlo
Carnegie Mellon University
pavlo@cs.cmu.edu

ABSTRACT

Most of the academic papers on concurrency control published in the last five years have assumed the following two design decisions: (1) applications execute transactions with serializable isolation and (2) applications execute most (if not all) of their transactions using stored procedures. I know this because I am guilty of writing these papers too. But results from a recent survey of database administrators indicates that these assumptions are not realistic. This survey includes both legacy deployments where the cost of changing the application to use either serializable isolation or stored procedures

1. ACKNOWLEDGEMENTS

This work was supported (in part) by the Intel Science and Technology Center for Big Data and the U.S. National Science Foundation (CCF-1438955).

2. BIOGRAPHIES

Andrew Pavlo is an Assistant Professor of Databaseology in the Computer Science Department at Carnegie Mellon University. At CMU, he is a member of the Database Group and the Parallel Data Laboratory. His research is in collaboration with the Intel Science

Criticism

- Do we really need many more transactions per second?
 - Probably not
 - The most you will see is around 750 M req/s on China's 11/11 Single's Day
 - Most of this workload embarrassingly parallel
 - Ultimately OLTP demand driven by economic growth, not database speed

- Are these new algorithms practical?
 - Maybe. But less than you'd think.
 - TPC-C != real-world (many dirty tricks to squeeze numbers out of TPC-C)
 - People don't need transactions all the time
 - Most people don't write highly optimized stored procedures

Should we just call it done?

Yes and No.

- (Almost) nobody cares if you improve TPC-C by 10%
- Guarantees are nice though
 - Many “NoSQL” systems ended up retrofitting to add transactions and/or strong consistency.
- New hardware
 - E.g. NVM and RDMA can change the equation (e.g., Microsoft’s FaRM)
- Cloud-native?
 - Traditional DBMSs fundamentally designed for shared-nothing
 - Cloud services change that
 - Open question on how to build performant, cheap, elastic cloud-native DBMS

Transactions in the Cloud

- Several different assumptions:
 - Disks are not persistent – must replicate across availability zones and even data centers
 - Strong primitives such as highly-available shared object storage exist
 - Economics dominate – people care about paying for just what they need

- Recent years saw a number of “cloud-native” database systems
 - E.g., AWS Aurora, Microsoft Socrates, SingleStore, FoundationDB, etc.
 - Most of these systems build on top of existing cloud storage services in “shared-disk” fashion
 - Most of these systems separate compute and storage for flexibility and service-ify important components for performance (e.g., logging)
 - “Serverless” is the hot word where these systems hide away infrastructure details and charge users on a per request basis

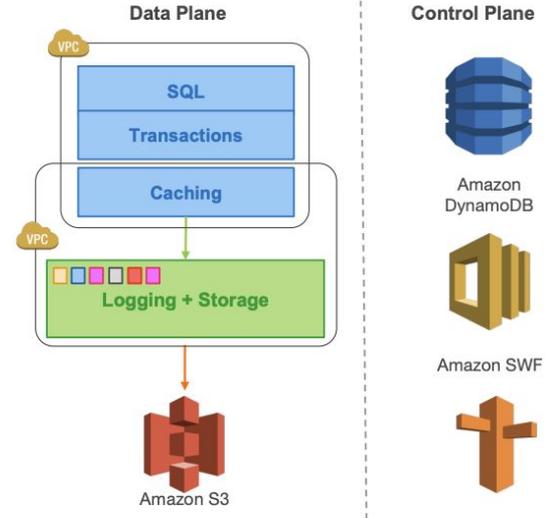
Cloud-native OLTP

- Key Idea: Storage & Compute Separation
 - Use cloud object storage (e.g., S3) for persistent storage layer
 - Attach ephemeral machines to storage when needed
 - Allows for separate scaling of resources

- Key Challenge: Performance
 - Object storage is often slow & over the network (upwards of 10ms instead of hundreds of microseconds of fast SSDs, and often rate-limited to tens of MBs per second)

Example: Amazon Aurora

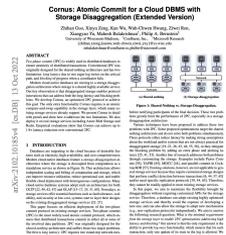
- Idea: take existing DBMS (e.g., PostgreSQL), and replace the storage layer
- Buffer pool loads data from S3, evicts data by erasing it
 - I.e., used exclusively as a cache
- Optimized logging layer to reduce commit latency and materialize pages in S3 by replaying



Cloud-native OLTP

- Area of active research to investigate what is fundamentally new about cloud-native databases
 - Example: can simplify protocol as underlying cloud storage guarantees consensus, fault-tolerance, etc.
 - Other potential directions: leveraging autoscaling serverless functions, converting existing components into multi-tenant services, etc.

Zhihan Guo, Xinyu Zeng, Kan Wu, Wuh-Chwen Hwang, Ziwei Ren, Xiangyao Yu, Mahesh Balakrishnan, Philip A. Bernstein. Cornus: Atomic Commit for a Cloud DBMS with Storage Disaggregation. (to appear) VLDB 2022



Takeaways

- Transactions have come a long way since the classical 2PL + ARIES + 2PC
- A host of new systems leveraging workload specialization and other clever insights to boost transactional performance by many orders of magnitude
 - Whether all of this speed-up is real is debatable
 - Regardless, many of the innovations run in production today
- Transactions research is alive and well in new settings such as the autoscaling cloud