

Tribune: An Externally Consistent Database Common Access API for the Global Data Plane

Matt Weber, Shiyun Huang and Lawrence Supian
Department of Electrical Engineering and Computer Sciences (EECS)
University of California, Berkeley
Email: matt.weber, jane.huang, lsupian@berkeley.edu

Abstract

The Global Data Plane (GDP) is a distributed data storage system with the objective of providing persistent data storage to devices in the internet of things. The GDP has a variety of useful properties including data security, a flat namespace, and location independent routing, that improve the availability and efficiency of universal data storage. As a key component in the Terraswarm project and the SwarmOS, the GDP provides log based storage and routing between sensors and actuators embedded in the physical world. Logs are the fundamental primitive of the GDP, but some CAAPIs (a Common Access Application Programming Interface) require stronger semantics. In this paper we introduce Tribune, a multi-writer log CAAPI that enforces Google Spanner-like external consistency. We compare the trusted Paxos based replication scheme with attack tolerant Byzantine Agreement and evaluate optimistic concurrency against the lock table-based concurrency control scheme used by Google Spanner. In a performance evaluation we found the optimistic algorithm to be generally faster than the comparable lock table approach, particularly when paired with Byzantine agreement. Our long term goal is to provide PTIDES style deterministic execution semantics for swarmlets in the SwarmOS that choose to access time-aware multi-writer logs.

1 Introduction

The Internet of Things (IoT) presents a new paradigm for computer systems, enabling technologies like context-aware apps, extensible cyber-physical systems, large scale sensor data collection for machine learning, and smart cities. But systems of the future will not come for free: engineers will need ubiquitous computing infrastructure in the form of platforms, services, and tools as a foundation for their work. The Terraswarm project [11] seeks to ad-

dress the gap between the development resources of today and those needed to engineer the coming swarm of interconnected devices [15]. There are many dimensions of the problem to be considered, including verification tools, low-power sensors, programming models, and persistent universally available data storage. The last point in particular is of primary concern to swarmlet (an application in the swarm) developers who need data accessible across different usage modes and application contexts [6].

The Global Data Plane (GDP) [7] is a TerraSwarm project that seeks to extend upon the capabilities of the Cloud to meet the needs of decentralized and interoperable swarmlets. Oceanstore [17] provides a starting point, as data must be globally available, durable, private, secure, and efficiently accessible. Additionally, the GDP must support storage and distribution of streaming sensor data, which it achieves through log-based data storage in a flat address space. As such, the log is a principle component of the GDP, and could be used as the fundamental building block for arbitrarily complicated systems for information representation. This design choice does not limit the efficiency or expressiveness of the GDP, because the GDP supports a Common Access Application Programming Interface (CAAPI) for data representations with sophisticated semantics or complex structure. Although possible, the process of reconstructing a full database or key-value store from log data would be prohibitively expensive if every time the sophisticated data structure were needed it had to be built from logs and then immediately thrown away. A CAAPI can maintain the current state of the enhanced structure directly and function as a sort of GDP cache for a particular data representation.

Although CAAPIs can be used in the GDP to store data efficiently, they also have the capacity to enforce semantics on the behavior of operations on data. This work, Tribune, is an example of a CAAPI that enforces external consistency along with standard ACID semantics on the behavior of a multi-writer log. The general concept of this style of CAAPI is illustrated in Figure 1. Tribune's name is a refer-

ence to the role of an ancient Roman Tribune, an elected official with the power to intervene on behalf of the common people by vetoing legislation from the senate. In a similar respect, Tribune can control transactions that attempt to write to a protected log in this multi-writer merge style, and abort transactions that violate its semantics.

This paper is organized as follows: Section II provides background information on the database behavior we seek to emulate in Tribune. Section III gives an overview of the system architecture and design. Section IV goes into the implementation details for our most interesting algorithms. Section V elaborates on our test environment setup and presents benchmark results. In Section VI we address the direction of future work. We conclude with Section VII.

2 Background

Tribune’s primary role is to enforce database semantics on a multi-writer log. We chose the semantics of Google Spanner’s read write transactions [5] because Spanner provides a time-based external consistency guarantee and works at a global scale. External consistency estab-

lishes the invariant that if transaction A commits before transaction B begins as observed from the outside world, timestamp associated with transaction B in the database must be later than than A’s timestamp. Spanner achieves this invariant through a combination of two phase locking within a Paxos group, time-aware two-phase commit between Paxos groups, and precise implementation of the TrueTime API as shown in Figure 2. TrueTime establishes an ordering between timestamps that reflects error in clock synchronization, i.e. it is unknown whether a timestamp proceeds another unless the difference in the timestamps exceeds the known upper bound on the offset between clocks.

Spanner’s high level architecture, shown in Figure 3, shares components and terminology with Chubby [3], Megastore [1], and BigTable [4], all Google projects as well. In Spanner, the tablets that are associated with a particular data model run on a set of replicas in the same Paxos group. The primary responsibility of the Paxos leader is to replicate the state of writes that go through it across the group. The leader also acts as a bottleneck for all read write transactions, which allows it to maintain a lock table that reflects the status of transactions currently in flight. Clients that want to perform read write transactions must acquire locks through two phase commit, and deadlocks are prevented with wound wait [19]. When a read write transaction needs to commit across paxos groups (a “distributed transaction”), one paxos leader steps up for the role of participant leader and runs a transaction manager to organize the two-phase commit. Spanner also supports read only transactions and read transactions (the two are actually distinct types in

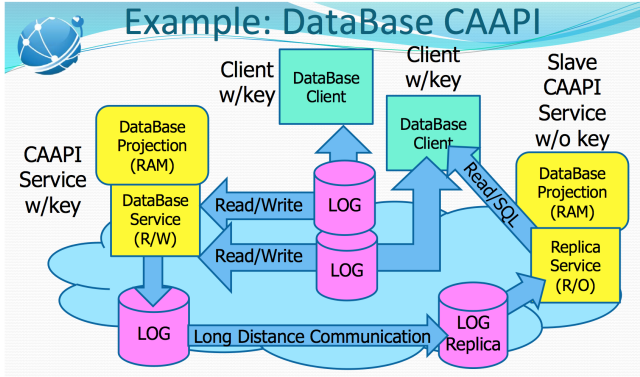


Figure 1. An illustration of a database CAAPI. Image from slide 17 of [7]

Method	Returns
<i>TT.now()</i>	<i>TTinterval</i> : [<i>earliest</i> , <i>latest</i>]
<i>TT.after(t)</i>	true if <i>t</i> has definitely passed
<i>TT.before(t)</i>	true if <i>t</i> has definitely not arrived

Figure 2. Spanner’s TrueTime interface. Image from Table 1 of [5]. Note that *t* is a timestamp

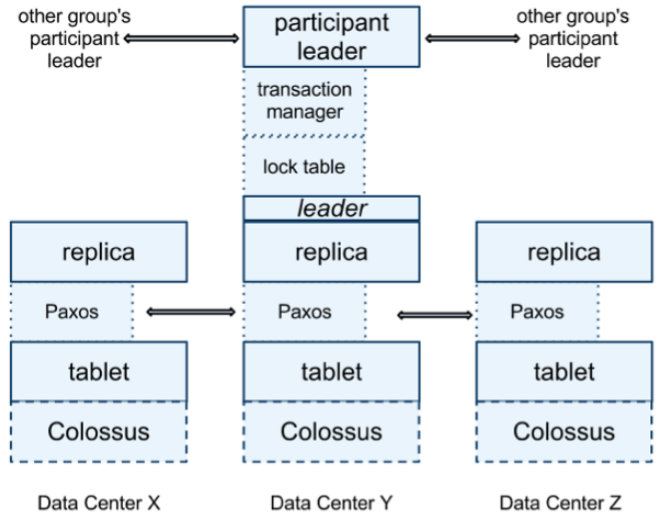


Figure 3. Spanner’s Architecture. Image from Figure 2 of [5]

Spanner) that can be performed at any replica without going through the paxos leader.

2.1 Time Synchronization

Spanner uses a global network of atomic and gps clocks to achieve time synchronization. An advantage of establishing very accurate clocks with very low drift is that very little communication is needed to preserve synchrony. IEEE 1588, the Precision Time Protocol [16] presents an alternative, ip based mechanism to achieve synchronization on the order of microseconds across a wired local area network. Truetime is used in spanner to determine when to pick timestamps for transactions in two-phase commit and how long to hold locks after committing within a paxos group. Note that without any kind of synchronization between clients, relativity makes it impossible to evaluate external consistency. The definition of external consistency requires that transactions outside the system exist on a shared timeline (because they must be comparable) which implies some way to compare times at two locations, whether quantitative or logical [8].

2.2 Consensus Algorithms and Durability

Consensus algorithms like Paxos [9] and Byzantine Agreement [10] can be used to build a state machine with an unambiguous history across distributed, faulty machines. As a consequence, data is unambiguously replicated through the system – a desirable outcome for a global database that must remain consistent and available in the face of datacenter-wide failures.

Paxos and Byzantine Agreement address different failure modes for a distributed state machine. Byzantine agreement tolerates malicious faults while Paxos protocols assume the entire system is trustworthy. There is however, a performance tradeoff. For a system with $2N + 1$ nodes, Paxos consensus allows for N nodes to simultaneously disappear. For a system with $3M + 1$ nodes, Byzantine Agreement can tolerate M malicious nodes but requires expensive pairwise communication to achieve consensus.

Note that availability, as provided by consensus algorithms, is not the same as durability [22]; the former relates to the presence of a backup-copy of a resource accessible in the short-term, where the latter refers to the capacity to protect data from accidental deletion. Reed-Solomon encoding [2] is a good approach for long-term data storage when the top objective is reducing the probability that data will be lost. But for short term data storage, using a consensus algorithm to achieve agreement on the state of a set of replicas is a much faster approach. Perhaps the Reed-Solomon encoding data might be verified by a Byzantine Agreement ring as in [17] but in the short term, other methods become

much more practical. As a long-term objective for the GDP, the system ought to achieve a balance between both durable and responsive storage.

Another dimension of durability concerns the transfer of data from volatile to non-volatile memory such that a particular machine can recover from a crash or power outage. Although we didn't implement this part of a database in Tribune, algorithms that address this concern (such as [12] and [20]) are well known in the literature and are applied in most mature databases.

2.3 Routing in the GDP

Data must be routed to Tribune in the GDP through an overlay network. Oceanstore used Tapestry [24], although a distributed hash table approach such as Chord [21] or Bamboo [18] for routing. A project is currently underway in the GDP to efficiently route to the opaque identifiers that signify the location of a log server or a CAAPL.

3 Background

4 System Architecture

ARCHITECTURE PICTURES GO HERE!!! Our first design of system emulates a simplified version of Google Spanner without any **"distributed transaction"**. Because we opt not to handle **"distributed transactions"**, we do not implement the transaction manager which coordinates two phase commits between different Paxos groups. Our first design of system assumes an one-paxos-group environment so we only implement the lock manager at the leader replica.

The responder receives transactions from client applications. It parses each operation line by line and acquires read locks for all data required for the transaction. The responder buffers writes locally. When all the transaction operation finishes, the responder tries to acquire write locks for

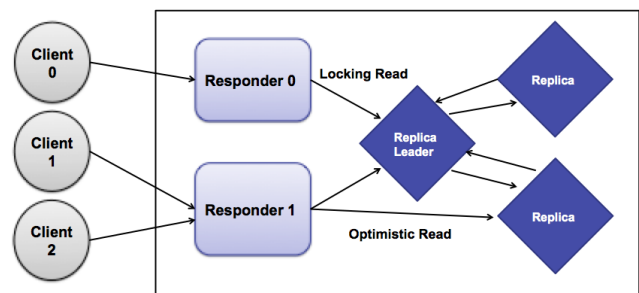


Figure 4. Tribune's Architecture.

the locally modified data. If the attempt turns out successful, the responder will try to commit the transaction at the leader.

We only implement the necessary components for read-write transactions. So all transactions would go through the leader replica to validate their respective lock leases before committing via paxos protocol. If the validation is successful, the leader would go through all paxos phases and return the final transaction status (abort or commit) back to responder. The responder will return the transaction status back to client apps.

//then talk about modifications on the first version to incorporate optimistic concurrency control and byzantine agreement

5 Algorithms and Implementation

EDIT THIS TO MAKE IT THE RIGHT TENSE In this section, we're going to elaborate on the algorithmic part of the project, which is our choices of concurrency control protocol and commit consensus protocol. We did strict two-phase locking and optimistic concurrency for concurrency control. We also did paxos and byzantine agreement for commit consensus protocol. We would explain in detail how we implement four algorithms in our project and compare the tradeoffs in this section. The benchmark results would be in the next section.

5.1 Programming Tools

We originally choose to implement Tribune in Java because the project team had familiarity with the language, and we felt a strongly typed, garbage collecting, object oriented language would reduce the frequency of programming errors in development. Java is known to perform "stop the world" garbage collection, in which execution periodically halts in all threads for 10s of milliseconds, but considering the low precision of our ultimate testing environment (three laptops) we decided it was an acceptable penalty given Java's other advantages.

As shown in Figure 4, client processes need to establish a connection with Responders, and Responders must communicate with Replicas, likely all across a network. Early on we experimented with basic java socket IO, but we realized we would have to repetitively implement the same programming idiom of a server that handles the connection and a parser thread that is responsible for reading a custom defined language of intended commands to be called on the receiving object. A little investigation found that this functionality was already available in the Java Remote Method Invocation framework [14], so we chose to use it instead of building a similar system from scratch.

However, RMI was not a perfect solution. The RMI server that is started for each registered object does not use a thread pool; it simply forks off a new thread whenever it gets a connection. The Sun implementation of RMI exposes a parameter for the number of TCP connections that can be ongoing at once, but this doesn't limit the number of worker threads that can be alive in the system at a given time. Therefore with a sufficiently large number of clients, spawning an arbitrarily large number of threads the RMI server may be ill-conditioned. Switching to SEDA [23] for future work, would alleviate this problem, but significantly change our programming model.

RMI also presents a naming problem. Surprisingly, RMI does not support a remote rmiregistry where objects can be found directly by registration name. A remote object may only be registered on a local rmiregistry server with a name of the form "<ip address of host machine>/<registration name>". Since our performance evaluation was performed on machines without static ip addresses, we potentially faced an addressing nightmare where the ip address of remote objects would change sporadically. To address this issue, we implemented a naming service over RMI we called Remote Registry, that maps a "<registration name>" to a full network-ready string of the form "///<ip address of host machine>/<registration name>". With Remote Registry we could publish not only the location of an object when its ip address changed, but also the availability of resources during initialization.

5.2 Initialization

5.3 Strict Two-Phase Locking

```
// Hello.java
import javax.swing.JApplet;
import java.awt.Graphics;

public class Hello extends JApplet {
    public void paintComponent(Graphics g) {
        g.drawString("Hello, world!", 65,
            95);
    }
}
```

5.4 Optimistic Concurrency Control

5.5 Initialization

5.6 Pseudo Paxos vs Pseudo Byzantine Agreement

We had to deal first-hand with a non-obvious implementation detail that arises in the practical implementation

of Paxos: out-of-date replicas. We took inspiration from Chubby’s solution to the problem [3] by using sequence numbers to catch a replica back up to the current state of the system when possible. It would be impractical to maintain a permanent map between paxos sequence numbers and changes to the database, so instead we truncate the log at a fixed distance in the past. If a replica is further behind than the end of the log, it asks for a snapshot of the current state of the system from the leader, and achieves up-to-date status.

COMMENT HOW WE MAKE AN ASSUMPTION ABOUT THE LEADER NOT CHANGING

Also, we considered using Raft [13] as an alternative to Paxos We implement the

5.7 Remote Method Invocation

6 Experiments and Results

In this section we present Tribune’s performance when given a variety of transaction workloads. All transaction workloads are randomly generated and execute concurrently in Tribune. We first evaluate a suite of long high conflict transactions, next we present the performance of a sample bank application that runs a high volume of short transactions between accounts, and finish with a high volume test of reads to a small set of memory locations to simulate a memory hotspot.

6.1 Experimental Setup

A graphic representation of our setup can be seen in Figure 5. We elected to run Tribune on three Macbook

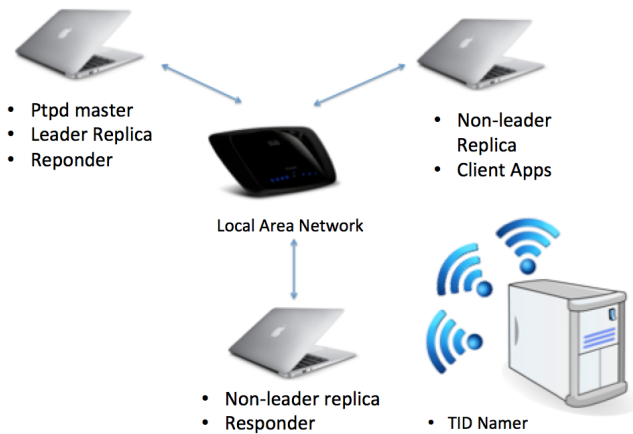


Figure 5. Experimental Setup for Evaluation

Pros as opposed to a cluster because we need easy physical access to the machines and admin status to run a Precision Time Protocol Daemon. Specifically we ran ptpd2, an open source implementation of IEEE1588. Using a Linksys router to establish a local area network over ethernet, we established a 10 μ s average offset from the master clock to the slaves. Going six standard deviations from the mean, we determined a 700 μ s value for the clock error in our implementation of TrueTime would upper bound the actual clock error in almost every case.

As for hardware and operating system specifications, the IEEE 1588 master machine in the upper left of Figure /reftest-setup has OS X Version 10.9.3, 2.4 GHz Intel Core i7, and 8 GB 1600 MHz DDR3. It runs the ptpd master and the consensus leader because it is our strongest machine The laptop at the bottom of the diagram has OS X Version 10.10.1, 2.5 GHz Intel Core i5, and 4 GB 1600 MHz DDR3, and the one in the upper right has OS X Version 10.9.3, 2.3GHz Intel Core i7, and 16 GB 1600 MHz DDR3. Note that all Client programs run on the machine in the upper right. It evenly distributes client program requests between each of the two responders. Transaction running time is measured on this laptop using a call to Java 8’s `Time.Instant.Now()` function before sending the transaction to a responder and after receiving a “commit” or “abort” message from Tribune.

The final component of our test setup is a transaction GUID namer running on a server in the same building on UC Berkeley’s EECS department network. Laptops connect to it through WiFi routers in our lab area. Ping time to the server is only slightly longer than to the LAN laptops so we expect this to have had a negligible impact on performance.

6.2 High Conflict Transactions

Our first benchmark consists of very long transactions with a high probability of accessing the same memory address for reads and writes. Therefore for both optimistic concurrency and locking control we expect to have a high percentage of aborts when many transactions are running on Responders. The specific format of the experiment shown in Figures 6 and 7 is to run a variable number of clients running from 1 to 10 concurrently in Tribune. This means we begin the trial with one client running and when it finishes we start two clients running, and so on. The name “Locking Paxos” refers to lock table based concurrency control with Paxos replication, “Optimistic Byzantine” refers to optimistic concurrency control with Byzantine replication, and the other runs are configurations of those features that correspond to the names in the legend. We use the same naming convention for the other experiments. Each client runs a list of three transactions, each consisting of one thousand commands to be executed. Because many transactions abort

shortly after starting, we only report the running time in Figure 6 for transactions that commit.

As expected, the Byzantine configurations were slower on average than the Paxos setups. The mean transaction execution times were also consistent

Note that the abort rate in Figure 7 is quite low for larger numbers of transactions dropping below 20% in some cases. As a general trend within a particular system configuration, the commit rate decreases as the contention increases. However, it's very surprising that the optimistic approaches resulted in fewer aborts than the lock based ones - as one would typically expect locking free reads and writes in high conflict transactions to consistently step on each other. We give two possible explanations for why this happened. First, the failure mode of optimistic concurrency in high conflict scenarios is mitigated in our experiment because there are only two replicas where a read might make place, and given the assumptions of our experiment, neither one is likely to be considerably out-of-date. Second, would wait deadlock resolution in the lock table may overzealously abort transactions that ought to have been able to commit in an optimistic scenario.

6.3 Bank Application

The second benchmark focuses on a more realistic application where data values might be important: a bank scenario similar to an ATM or teller for bank. Like the high conflict experiment, we generate random workloads, but for each client, we reduce the number of commands for each transaction inside the client from 1000 to 20. To compensate we start 10 times as many clients on each iteration of the experiment and run 10 transactions per client. We start with ten client apps and increment by ten for every run up to one hundred. The objective is to present a realistic use

case where conflicts over data appear, but we aren't forced to abort transactions as often as in the high conflict scenario.

It is very odd that the mean execution time of Optimistic Paxos is slower than Optimistic Byzantine in Figure 8, considering that Byzantine agreement is simulated in Tribune at this stage by doing strictly more work than Optimistic Paxos in the form of pairwise communication between replicas.

Spanner-style semantics concerning time are most certainly not the only option for a CAAPI in the GDP. PTIDES [25] provides both an execution model and a simulation environment for time-aware computation. Rather than dealing with transactions, PTIDES determines the execution of atomic events and enforces that sensor-to-actuator deadlines are met. As an advantage of tight time synchronization and atomic events, PTIDES can define a deterministic time-ordered merge between two input streams by simply waiting out the clock uncertainty before writing. This policy guarantees no future event could be timestamped further in the

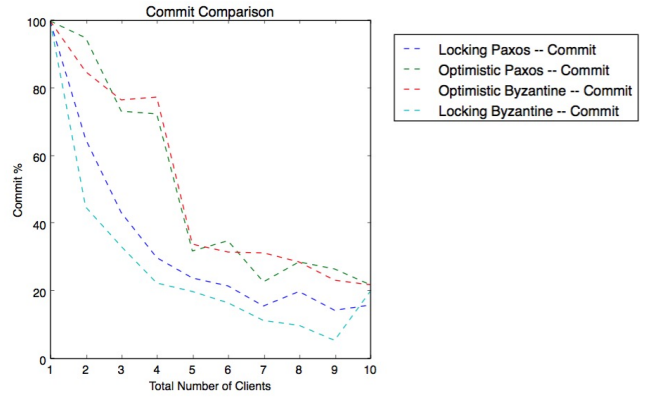


Figure 7. High Conflict Rate: Commit Percentage

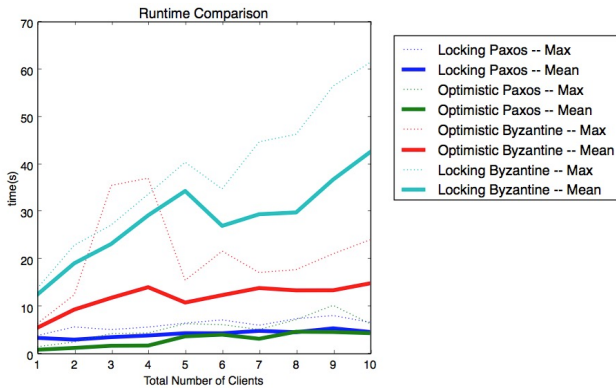


Figure 6. High Conflict Rate: Committing Transaction Execution Time

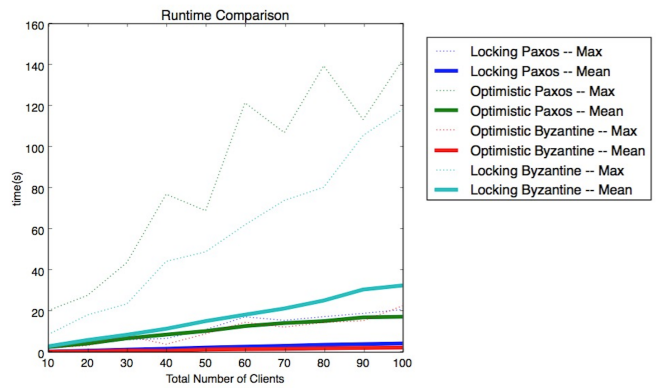


Figure 8. Bank Application: Committing Transaction Execution Time

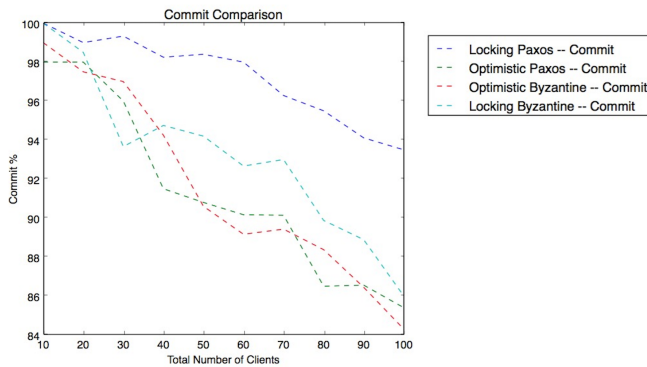
past than the written value.

Reconsidering our implementation choices now the project is complete, **THREADS ARE EVIL Finish this section!**

8 Conclusion

9 Acknowledgements

Thanks to Dr. Patricia Derler for her input on Spanner and time synchronization. Also thanks to Prof. John Kubiawicz for numerous discussions about the GDP and



6.4 Hot Spot Reads

Mention the log scale on the last time graph. Also mention that the commit rate for the third test is 100 percent

7 Future Work

Figure 9. Bank Application: Commit Percentage

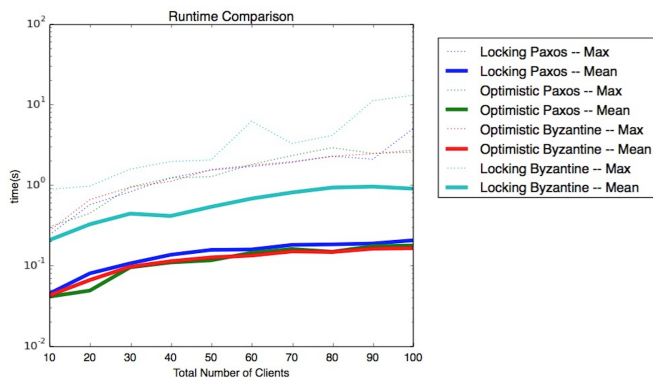


Figure 10. Hot Spot Reads: Committing Transaction Execution Time

CAAPIs.

References

- [1] J. Baker, C. Bond, J. C. Corbett, J. J. Furman, A. Khorlin, J. Larson, J.-M. Lon, Y. Li, A. Lloyd, and V. Yushprakh. Megastore: Providing scalable, highly available storage for interactive services. In *CIDR*, volume 11, pages 223–234, 2011.
- [2] J. Blmer, M. Kalfane, R. Karp, M. Karpinski, M. Luby, and D. Zuckerman. An xor-based erasure-resilient coding scheme, 1995.
- [3] T. D. Chandra, R. Griesemer, and J. Redstone. Paxos made live: an engineering perspective. In *Proceedings of the twenty-sixth annual ACM symposium on Principles of distributed computing*, pages 398–407. ACM, 2007.
- [4] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber. Bigtable: A distributed storage system for structured data. *ACM Transactions on Computer Systems (TOCS)*, 26(2):4, 2008.
- [5] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, and P. Hochschild. Spanner: Google's globally-distributed database. In *Proceedings of OSDI*, volume 1, 2012.
- [6] P. Dabbelt and J. D. Kubiawicz. A case for the universal dataplane, September 2013. Presented at theFirst International Workshop on the Swarm at the Edge of the Cloud (SEC'13 @ ESWeek).
- [7] J. D. Kubiawicz. Enabling the swarm through the global data plane, November 2014.
- [8] L. Lamport. Time, clocks, and the ordering of events in a distributed system. *Communications of the ACM*, 21(7):558–565, 1978.
- [9] L. Lamport. The part-time parliament. *ACM Trans. Comput. Syst.*, 16(2):133–169, May 1998.
- [10] L. Lamport, R. Shostak, and M. Pease. The byzantine generals problem. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 4(3):382–401, 1982.
- [11] E. A. Lee, J. D. Kubiawicz, J. M. Rabaey, A. L. Sangiovanni-Vincentelli, S. A. Seshia, J. Wawrzynek, D. Blaauw, P. Dutta, K. Fu, and C. Guestrin. The TerraSwarm research center (TSRC)(a white paper). Technical report, Technical report UCB/EECS-2012-207, 2012.

- [12] C. Mohan, D. Haderle, B. Lindsay, H. Pirahesh, and P. Schwarz. ARIES: a transaction recovery method supporting fine-granularity locking and partial roll-backs using write-ahead logging. *ACM Transactions on Database Systems (TODS)*, 17(1):94–162, 1992.
- [13] D. Ongaro and J. Ousterhout. *In search of an understandable consensus algorithm (extended version)*. 2014.
- [14] E. Pitt and K. McNiff. *Java.Rmi: The Remote Method Invocation Guide*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2001.
- [15] J. M. Rabaey. The swarm at the edge of the cloud—a new perspective on wireless. In *VLSI Circuits (VLSIC), 2011 Symposium on*, pages 6–8, 2011.
- [16] R. Ratzel and R. Greenstreet. Toward higher precision. *Communications of the ACM*, 55(10):38–47, 2012.
- [17] S. C. Rhea, P. R. Eaton, D. Geels, H. Weatherspoon, B. Y. Zhao, and J. Kubiatawicz. Pond: The OceanStore prototype. In *FAST*, volume 3, pages 1–14, 2003.
- [18] S. C. Rhea, D. Geels, T. Roscoe, and J. Kubiatawicz. *Handling churn in a DHT*. Computer Science Division, University of California, 2003.
- [19] D. J. Rosenkrantz, R. E. Stearns, and P. M. Lewis, II. System level concurrency control for distributed database systems. *ACM Trans. Database Syst.*, 3(2):178–198, June 1978.
- [20] R. Sears and E. Brewer. Segment-based recovery: write-ahead logging revisited. *Proceedings of the VLDB Endowment*, 2(1):490–501, 2009.
- [21] I. Stoica, R. Morris, D. Liben-Nowell, D. R. Karger, M. F. Kaashoek, F. Dabek, and H. Balakrishnan. Chord: a scalable peer-to-peer lookup protocol for internet applications. *Networking, IEEE/ACM Transactions on*, 11(1):17–32, 2003.
- [22] H. Weatherspoon, P. Eaton, B.-G. Chun, and J. Kubiatawicz. Antiquity: exploiting a secure log for wide-area distributed storage. *ACM SIGOPS Operating Systems Review*, 41(3):371–384, 2007.
- [23] M. Welsh, D. Culler, and E. Brewer. SEDA: an architecture for well-conditioned, scalable internet services. In *ACM SIGOPS Operating Systems Review*, volume 35, pages 230–243. ACM, 2001.
- [24] B. Zhao, L. Huang, J. Stribling, S. Rhea, A. Joseph, and J. Kubiatawicz. Tapestry: A resilient global-scale overlay for service deployment. *IEEE Journal on Selected Areas in Communications*, 22(1):41–53, Jan. 2004.
- [25] J. Zou, S. Matic, E. A. Lee, T. H. Feng, and P. Dierler. Execution strategies for PTIDES, a programming model for distributed embedded systems. In *Real-Time and Embedded Technology and Applications Symposium, 2009. RTAS 2009. 15th IEEE*, pages 77–86. IEEE, 2009.