

COLUMBIA UNIVERSITY

MASTER THESIS

**Reliable Synchronization in
Multithreaded Servers**

Author:
Rui GU

Supervisor:
Dr. Junfeng YANG

*A thesis submitted in fulfillment of the requirements
for the degree of Master of Science*

in the

System Software Laboratory
Department of Computer Science

Committee:

Associate Professor Junfeng Yang, Chair Professor

Professor Jason Nieh

Assistant Professor Yinzhi Cao, Lehigh University

May 15, 2017

“It’s good to learn from your mistakes. It’s better to learn from other people’s mistakes.”

Warren Buffett

Columbia University

Abstract

Junfeng Yang
Department of Computer Science

Master of Science

Reliable Synchronization in Multithreaded Servers

by Rui GU

State machine replication (SMR) leverages distributed consensus protocols such as PAXOS to keep multiple replicas of a program consistent in face of replica failures or network partitions. This fault tolerance is enticing on implementing a principled SMR system that replicates general programs, especially server programs that demand high availability. Unfortunately, SMR assumes deterministic execution, but most server programs are multithreaded and thus nondeterministic. Moreover, existing SMR systems provide narrow state machine interfaces to suit specific programs, and it can be quite strenuous and error-prone to orchestrate a general program into these interfaces. This paper presents CRANE, an SMR system that transparently replicates general server programs. CRANE achieves distributed consensus on the socket API, a common interface to almost all server programs. It leverages deterministic multithreading (specifically, our prior system PARROT) to make multithreaded replicas deterministic. It uses a new technique we call time bubbling to efficiently tackle a difficult challenge of nondeterministic network input timing. Evaluation on five widely used server programs (e.g., Apache, ClamAV, and MySQL) shows that CRANE is easy to use, has moderate overhead, and is robust.

Acknowledgements

My first and biggest thanks go to my adviser, Professor Junfeng Yang, for both his excellent guidance and motivation. I would also like to thank the other members of my dissertation committee - Professor Jason Nieh, Professor Yinzhi Cao - for their invaluable comments and constructive criticisms that greatly improved my thesis. It is an honor to have them on my committee.

I want to express my thanks to Professor Heming Cui from University of Hongkong. Heming served as my “secondary” advisor at Columbia and he taught me how to be focus and hard working. I also want to have my special thanks to Professor Shan Lu, University of Chicago. Shan was once my research advisor at University of Wisconsin, Madison. She guided me through my first research project and I’m still encouraged by her research ideology.

I also want to express my biggest thanks to all my previous lab mates, Guoliang Jin, Linhai Song, Po-chun Chang, Dongdong Deng, Yinzhi Cao, Yang Tang, Gang Hu, Xinhao Yuan, David Williams-King, Lingmei Weng, Linjie Zhu and Kexin Pei. My research life will never be that colorful without them.

My last thanks to my Mom, Dad and my fiancée Xi Liu. It is my Mom who always stand and encourage me in every sleepless night.

Contents

Abstract	ii
Acknowledgements	iii
1 Introduction	1
1.1 Motivation	1
1.2 Background	2
1.3 Preview	2
2 Transparent State Machine Replication	4
2.1 CRANE's Overview	4
2.2 Architecture	4
2.3 CRANE's Synchronization Wrapper	6
2.4 The Time Bubbling Technique	6
3 Discussion	8
3.1 Limitation	8
3.2 Application	8
4 Evaluation	9
4.1 Ease of Use	10
4.2 Consistency of Network Outputs	10
4.3 Performance Overhead in Normal Case	11

Chapter 1

Introduction

1.1 Motivation

State machine replication (SMR) models a program as a deterministic state machine, where the states are important program data and the transitions are deterministic executions of program code under input requests. SMR runs replicas of the program and invokes a distributed consensus protocol (typically PAXOS [22, 21, 36]) to ensure the same sequence of input requests for replicas, as long as a quorum (typically a majority) of the replicas agrees on the input request sequence. Under the deterministic execution assumption, this quorum of replicas must reach the same exact state despite replica failures or network partitions. SMR is proven safe in theory and provides high availability in practice [10, 27, 8, 33, 7, 26, 16, 14].

The fault-tolerant benefit of SMR makes it particularly attractive on implementing a principled replication system for general programs, especially server programs that require high availability. Unfortunately, doing so remains quite challenging; the core difficulty lies in the deterministic state machine abstraction required by SMR, elaborated below. First, the deterministic execution assumption breaks down in today's server programs because they are almost universally multithreaded. Even on the same exact sequence of input requests, different executions of the same exact multithreaded program may run into different thread interleaving, or schedules, depending on such factors as OS scheduling and physical arrival times of requests. Thus, they can easily exercise different schedules and reach divergent execution states - a difficult problem well recognized by the community [6, 16, 15, 14]. To tackle this problem, one prior approach, execute-verify [16], detects divergence of execution states and retries, but it relies on developers to manually annotate states, a strenuous and error-prone process.

Second, to leverage existing SMR systems such as ZooKeeper [38], developers often have to shoehorn their programs into the narrowly defined state machine interfaces provided by these SMR systems. Ideally, experts - those with intimate knowledge of the arcane (think how many papers [22, 21, 36, 10, 27] are needed to explain PAXOS), under-specified [27] SMR protocols and subtle failure scenarios in distributed systems - should build a solid SMR system, which all other developers then leverage. However, an SMR system often has to settle for a specific state and transitional interface because it cannot anticipate all possibilities in which developers structure their programs. For example, Chubby [8] defines a lock server interface, and ZooKeeper a pseudo file system interface. Orchestrating a server program into such a narrow interface not only requires intrusive and error-prone modifications to the program's structure and code, but also disrupts the SMR system itself at times. For instance, developers abused Chubby for storage [8], causing the Chubby developers to add quota support. This paper presents CRANE, an SMR system that

transparently replicates server programs for high availability. With CRANE, a developer focuses on implementing her program's intended functionality, not replication. When she is ready to replicate her program for availability, she simply runs CRANE with her program on multiple replicas. Within each replica, CRANE interposes on the socket and the thread synchronization interfaces to keep replicas in sync. Specifically, it considers each incoming socket call (e.g., `accept()` a client's connection or `recv()` a client's data) an input request, and runs a PAXOS consensus protocol [27] to ensure that a quorum of the replicas sees the same exact sequence of the incoming socket calls.

1.2 Background

Two prior approaches attempted to tackle this challenge. Execute-agree-follow [14] records a partially ordered schedule of Pthreads synchronizations on one replica and replays it on the other replicas, which may incur high network bandwidth consumption and performance overhead. dOS [6] also leverages DMT for replication, but it determines the logical admission time for each request using two-phase commit. Aside from two-phase commit's known intolerance of primary failures, per-request commit is also costly.

1.3 Preview

One may consider solving this challenge by leveraging the underlying distributed consensus protocol to determine the logical admission time for each request. Specifically, when running the consensus protocol to decide each request's position in the request sequence, a predicted logical admission time can be carried as part of the decision as well. Unfortunately, predicting a logical admission time for each request accurately is quite challenging because typical server programs have background threads which may frequently tick logical clocks. A too-small predication leads to replica divergence if another replica has already run past the predicted logical time. A too-large predication blocks the system unnecessarily because replicas cannot admit the request before reaching the predicted time.

Our key insight is that many requests need no admission time consensus because their admission times are already deterministic. Hypothetically, if the requests arrive faster than they are admitted at each replica, each request's admission time is fully deterministic because each replica simply admits requests as fast as it can. In practice, requests do not arrive this fast. However, there are still frequent bursts of requests that arrive together. Among replicas, as long as the first request of a burst is admitted at a deterministic logical time, all the other requests in the burst are admitted at deterministic logical times without requiring consensus.

Leveraging this insight, we created a technique called time bubbling to enforce deterministic logical times efficiently. It ensures that the first request in a burst is admitted at each replica deterministically by inserting a deterministic wait after the previous burst of requests are all admitted. During this wait, each replica only processes already admitted requests, and does not admit new requests. CRANE negotiates a consistent duration of the wait via the underlying distributed consensus protocol, and enforces this wait at each replica via DMT. These waits are like deterministic time bubbles between bursts of requests (hence the name of the technique), creating the illusion that the requests arrive faster than they are admitted.

In short, by converting per-request admission time consensus to per-burst, time bubbling efficiently combines the input determinism of PAXOS and the execution determinism of DMT. For busy servers, requests in bursts greatly outnumber the other requests. (We observed that 66.65% to 93.88% of requests are in bursts) They rarely need to invoke time consensus, enjoying good performance. For idle servers, time consensus overhead does not matter much because the servers are idle anyway.

We implemented CRANE by interposing on the POSIX socket and the Pthreads synchronization interfaces. It intercepts operations along these interfaces by hijacking dynamically linked library calls for transparency. It implements the PAXOS protocol atop libevent [24] for distributed consensus, and leverages our PARROT system for deterministic multithreading. Unlike prior SMR systems with narrow interfaces, CRANE's checkpoint and recovery must work with general programs. To this end, it leverages CRIU [11] to checkpoint and restore process states, and LXC [25] for file system states. An additional benefit of using the LXC container is that CRANE isolates the replicated server program from the environment, avoiding nondeterministic systems resource contentions.

We evaluated CRANE on five widely used server programs, including HTTP servers Apache and Mongoose, an anti-virus server ClamAV, a uPnP multimedia server MediaTomb, and a database server MySQL. Our results on popular performance benchmarks show that CRANE works with all the servers easily (three servers require no modification, and the other two servers each require only two lines of PARROT hints [12] to improve performance); that CRANE's performance overhead is moderate (an average of 34.19% overhead at the servers' peak performance setups on our 24-core machines); and that CRANE is robust on replica failures.

Our key conceptual contribution is the idea of transparent SMR for general programs, which has the potential to expand SMR's adoption and improve availability of many systems. This idea also applies to other replication concepts (e.g., byzantine fault tolerance [9, 17]). This idea has other broad applications as well. Our engineering contributions include the CRANE system and its evaluation on widely used server programs. All CRANE's source code (including a standalone, libevent based PAXOS implementation), benchmarks, and evaluation results are available at github.com/columbia/crane.

Chapter 2

Transparent State Machine Replication

2.1 CRANE's Overview

CRANE is deployed as a typical SMR system. A set of three or five replicas is set up within a LAN, and each replica runs an CRANE instance containing the same server program. On the CRANE system starts, one replica becomes the *primary* (or leader) replica which proposes the order of requests to execute, and the others become backup replicas which follow the primary's proposals. A number of clients in LAN or WAN send network requests to the primary and get responses. If the primary machine fails, the other replicas run a leader election to elect a new primary.

2.2 Architecture

To support general server programs transparently, CRANE chooses the POSIX socket API as its consensus interface. CRANE enforces two kinds of orders for socket calls. First, for client programs' outgoing socket calls (e.g., `connect()` and `send()`), CRANE enforces that all replicas see the same sequence of client socket calls with PAXOS. CRANE does not need to order the clients' blocking socket calls because CRANE is not designed to replicate clients. Second, for a server program's blocking socket calls (e.g., `poll()`, `accept()`, and `recv()`), CRANE enforces that these calls are scheduled and returned in the same sequence of logical times across replicas. CRANE responses to the clients only using the server program on the primary, and it drops the responses of the server programs on backups. For a server program's outgoing socket calls (e.g., `send()`), CRANE simply schedules them using DMT and does not invoke consensus. The reason is that these calls readily have consistent contents via enforcing the same logical admission times of input requests and the same thread schedules for server programs across replicas.

Figure 2.1 shows a CRANE instance running on the primary. The instance contains five main components, the proxy, the PAXOS consensus, the DMT scheduler, the time bubbling component that enforces the same logical clocks for servers' blocking socket calls across replicas via inserting time bubbles, and the checkpoint component that periodically checkpoints the server program. A server program runs transparently in a CRANE instance without being aware of CRANE's components. A backup replica runs the same CRANE instance except that its proxy does not accept connections from clients and does not invoke consensus.

The proxy component is a CRANE instance's gateway. It accepts socket requests from clients and forwards the requests to the server program on its own replica. It accepts responses from the server program and forwards the responses to the clients.

logical clocks in a time bubble, it either admits new client socket call (if any) or inserts another time bubble. If the scheduler does not exhaust the logical clocks after serving current requests, PARROT has a mechanism to exhaust them rapidly.

To recover from replica failures or add new replicas, the checkpoint component is invoked every minute on a backup replica. It checkpoints the server process running with DMT. While one can always start a server replica from scratch and replay the entire sequence of socket calls, this replay can be extremely time-consuming for long-running servers. Prior SMR systems rely on narrow state machine interfaces for checkpoint and recovery, which does not work for general server programs. Instead, CRANE leverages two popular open source tools: CRIU, to checkpoint process state such as CPU registers and memory; and LXC, to checkpoint the file system state of a server program's current working directory and installation directory.

Each checkpoint in CRANE is associated with a global index in PAXOS's consensus order, so if one replica needs recovery, CRANE ships the latest checkpoint from a backup replica, restores the process running DMT and the server program, and re-executes socket calls starting from this index. The proxy and consensus components do not require checkpoints because we explicitly designed their execution states independent to the server's process.

2.3 CRANE's Synchronization Wrapper

This section describes how CRANE handles a server program's synchronizations, including Pthreads synchronizations and blocking socket calls. Because how to handle these synchronizations is tightly relevant to the PARROT DMT scheduler we leverage, in this section, we first introduce some background on the PARROT scheduler, including its primitives and wrappers. And then we describe how CRANE leverages PARROT's primitives and wrappers to implement its own synchronization wrappers.

CRANE wraps a rich set of common blocking socket operations, including `select()`, `poll()`, `epoll wait()`, `accept()`, and `recv()`. CRANE also modifies the wrappers of Pthreads synchronizations. These wrappers are sufficient for the server programs in our evaluation. CRANE needs to modify the `pthread_mutexlock()` wrapper to do three things. First, if the PAXOS request sequence has been empty for a physical duration `Wtimeout`, CRANE requests a time bubble with `Nclock` logical clocks. Second, if the head of the PAXOS sequence is a time bubble, CRANE decreases the logical clock in the time bubble by one, or it removes this bubble if zero clock is left. Third, CRANE signals a thread that blocks on a socket operation (e.g., `recv()`) if there is a matching client socket call (e.g., `send()`) at the head of the PAXOS sequence.

CRANE also needs to modify PARROT's idle thread mechanism because sometimes this thread is the only thread in the run queue, and CRANE needs to frequently check whether a new client socket call comes or a time bubble insertion is needed. To do so, CRANE replaces PARROT's `get turn()` and `put turn()` primitives within the idle thread to be mutex lock and unlock operations, then the idle thread also runs the function to check and insert time bubbles.

2.4 The Time Bubbling Technique

Figure 2.2 shows the time bubbles inserted by the time bubbling technique. The technique groups clients' socket operations as bursts. A request burst can be a group of real socket requests (rectangles), or can be a time bubble with a fixed number of

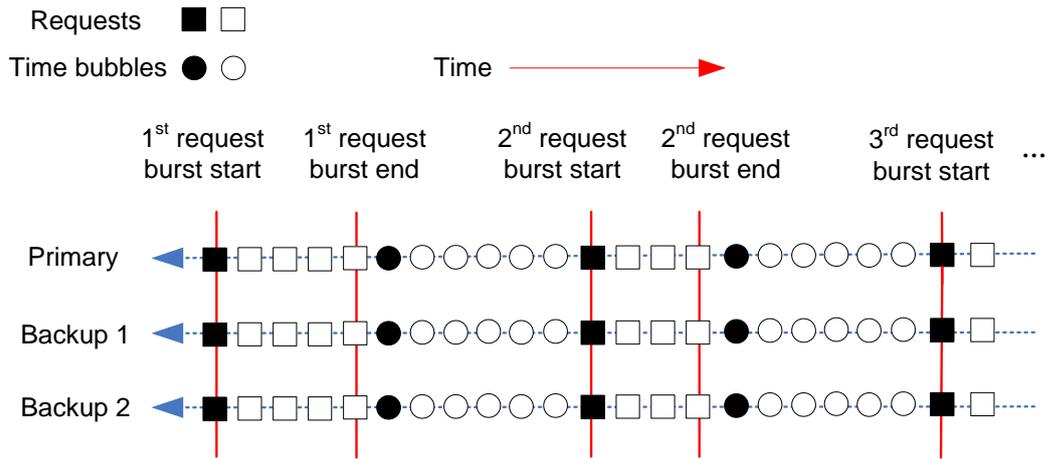


FIGURE 2.2: The request and time bubble flow.

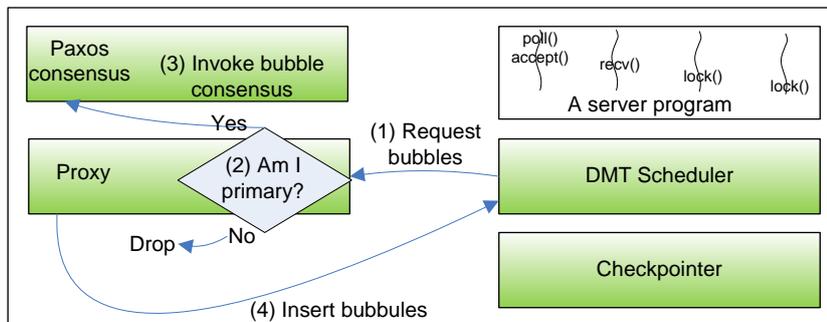


FIGURE 2.3: The work flow of inserting a time bubble.

logical clocks (circles). In this figure, black requests are the first operation for each burst. In a conceptual level, CRANE uses three rules to enforce the same sequence of logical times for socket requests (rectangles) and thus the same schedules across different replicas. First, CRANE uses PAXOS to ensure the same sequence of client socket calls as well as inserted time bubbles as a “PAXOS request sequence” for each replica, as shown in each horizontal arrow. Second, CRANE uses DMT to guarantee that it only ticks logical clocks (i.e., schedules Pthreads synchronizations or socket operations) when this sequence is not empty. Third, the time bubbling technique ensures that this sequence is not empty, otherwise it inserts a time bubble.

Figure 2.3 shows the work flow of our time bubbling technique with four steps. Each replica’s DMT just waits for a physical duration $W_{timeout}$, if no further requests come, (1) the DMT requests its own proxy to insert a time bubble. (2) The proxy then checks whether it sees itself as the primary in the PAXOS protocol. If so, it asks (3) the consensus component to invoke consensus on whether inserting this bubble; otherwise it drops this request. After a consensus on this bubble insertion is reached, (4) each machine’s proxy simply inserts the bubble into the PAXOS sequence, granting N_{clock} logical clocks to the DMT scheduler. If a server has not exhausted the logical clocks in a time bubble after serving current requests, PARROT’s idle thread mechanism exhausts these clocks rapidly. Then, the server can continue to process further requests in time.

Chapter 3

Discussion

This section first discusses CRANE’s limitations and then introduces its applications.

3.1 Limitation

CRANE leverages PARROT to make synchronizations deterministic. PARROT is explicitly designed not to handle data races. However, in the context of CRANE, data races are less harmful because, if they cause backups to crash, CRANE can still operate and recover as long as a quorum of the replicas is still alive. Moreover, leveraging CRANE’s replication architecture, one can deploy a race detector on a backup replica [13], achieving both good CRANE performance and full determinism. There are other sources of nondeterminism besides thread scheduling and request timing. These other sources of non-determinism may cause backups to diverge, too. For example, backups may do different things based on their IP addresses, data read from `/dev/random`, addresses returned by `malloc`, physical time observed via `gettimeofday`, or delivery time of signals. Prior work has shown how to eliminate these sources of nondeterminism using record-replay [19, 23] or OS-level techniques [6], which CRANE can leverage. Another solution is to treat all these sources as inputs and leverage distributed consensus to let all replicas observe the same input. We leave these ideas for future work. We inspected server programs’ network outputs among replicas, and we found that these outputs were consistent in CRANE except physical times. For a server program that spawns multiple processes which communicate via IPC, CRANE currently does not make these IPC operations deterministic. We expect that it should be easy to support deterministic IPC in CRANE because it already makes socket API deterministic. In addition, `dos` [6] and `DDOS` [15] have many effective techniques for tackling this problem, which CRANE can leverage.

3.2 Application

We envision three applications for CRANE. First, CRANE can be leveraged by other replication concepts (e.g., byzantine fault tolerance [9, 17]) and record-replay [18, 20, 23] because they also suffer from nondeterminism. Second, promising results in [13] have shown that CRANE’s transparent replication architecture can enable multiple types of program analysis tools within one execution, making a server program enjoy benefits of multiple analyses. Third, CRANE’s determinism as well as its time bubbling technique alone can be applied to mitigate timing channels [4, 37, 5].

Chapter 4

Evaluation

Our evaluation was done on a set of three replica machines, with each having Linux 3.13.0, 1Gbps bandwidth LAN, 2.80 GHz dual-socket hex-core Intel Xeon with 24 hyperthreading cores, 64GB memory, and 1TB SSD. We evaluated CRANE on five widely used server programs, including HTTP servers Apache [3] and Mongoose [30]; ClamAV [1], an anti-virus scanning server that scans files in parallel and deletes malicious ones; MediaTomb [28], a uPnP multimedia server that uploads, shares, and transcodes pictures and videos in parallel; and MySQL [31], an SQL database. Although MySQL has a replication feature [32], this feature is mainly for improving read performance, not for providing SMR fault tolerance. SMR's high availability and fault-tolerance are attractive to these servers programs, because these programs provide online service and contain important in-memory execution states and storage (e.g., ClamAV's security database, MediaTomb's SQLite [34] database, and MySQL).

For Apache and Mongoose, we used Apache's own concurrency stress testing benchmark ApacheBench to invoke concurrent HTTP requests for a PHP page, which takes about 70 ms for a PHP interpreter to generate the page contents. For ClamAV, we used its own client utility clamscan to request the server to scan ClamAV's own source code and installation directories in parallel. For MediaTomb, because it has a web interface, we used ApacheBench to invoke concurrent requests which use mencoder [29] to transcode a 15MB video from AVI to MP4. For MySQL, we used SysBench [35] to generate random select queries. These workloads triggered 8-12 threads in each server program to process requests concurrently at peak performance on our machines. These popular benchmarks and workloads cover CPU, network, and file-IO bounded operations.

CRANE has two parameters for the time bubbling technique. The first parameter, *Wtimeout*, is the physical duration that the primary's DMT scheduler waits before it requests consensus on a time bubble insertion. To prevent this parameter significantly deferring responses, CRANE sets its default value 100us, two orders of magnitudes smaller than the workloads' response times and wide-area network latencies. The second parameter, *Nclock*, is the number of logical clocks within each time bubble. CRANE sets its default value 1000, because we observed that the amounts of executed Pthreads synchronizations to process each request in most of the evaluated servers are closed to this scale. We used these default values in all evaluations unless explicitly specified. A sensitivity evaluation on these two parameters showed that their default values were reasonable choices. To mitigate network latency, benchmark clients were ran within the replicas' LAN. Larger latency will mask CRANE's overhead. We measured each workload's response time as it has direct impact on users. For each data point, we ran 1K requests for 20 times and then picked the median value. The rest of this section focuses on these questions:

1. Is CRANE easy to use?
2. Compared to nondeterministic executions, does CRANE consistently enforce the same sequence of network outputs among replicas?
3. What is CRANE's performance overhead compared to nondeterministic executions?

4.1 Ease of Use

All five servers we evaluated were able to be transparently plugged and played in CRANE without modification. For ClamAV, MediaTomb, and MySQL, we did not need to modify any line of code and they already have moderate performance overhead compared to the unreplicated nondeterministic executions. For Apache and Mongoose, the default schedules serialized parallel computations. For each of the two servers, we added two lines of soft barrier performance hints invented by PARROT [12] to line up parallel computations as much as possible and compute efficient DMT schedules.

4.2 Consistency of Network Outputs

To verify whether the server programs running in different replicas maintain the same execution states, we compared each server program's network outputs logged in three replicas. Network outputs imply a server's execution states, including the outcomes of ad-hoc synchronizations and data races, which synchronization schedules can not capture. We ran the performance workloads and logged the order and contents of server programs' outgoing socket calls, including `send()`, `sendto()`, `sendmsg()`, `write()`, and `pwrite()`. These calls are sufficient to capture all network outputs of the evaluated programs. We then used `diff` to compare the logs across replicas.

We designed two experiment plans. In plan I, we ran CRANE with the programs. In plan II, we disabled only the time bubbling component in CRANE for three reasons: (1) we wanted to know whether time bubbling is needed to keep replicas in sync, (2) enabling PAXOS made us easy to ship the same workload to replicas, and (3) enabling PARROT made us easy to intercepted and logged network outputs. Among the five programs, three server programs, Apache, MediaTomb, and Mongoose, used ApacheBench to spawn workloads. In plan I, CRANE's logs from all three replicas had the same order and contents of outputs except physical times in the responded HTTP headers. In plan II, despite that we disabled only the time bubbling component, the logs' order of responded HTTP headers and contents across replicas were different. Two server programs, ClamAV and MySQL, used specific benchmarks to spawn workloads. In plan I, the logs showed that CRANE enforced the same network outputs. In plan II, the orders of the outputs across replicas were different. These experiments suggest that simply combining PAXOS and DMT is not sufficient to keep replicas in sync, and the time bubbling technique is needed.

To diagnose consistency of network outputs more concisely, we wrote a micro-benchmark for Apache. We used the `curl` utility to spawn two concurrent HTTP requests: a PUT request of a PHP page and a GET request on this page, and then we inspected the outcome of the GET request. We ran Apache in CRANE with this micro benchmark for 100 times and found that three replicas consistently reported the same GET result in each run, either "200 OK" or "404 Not Found", depending

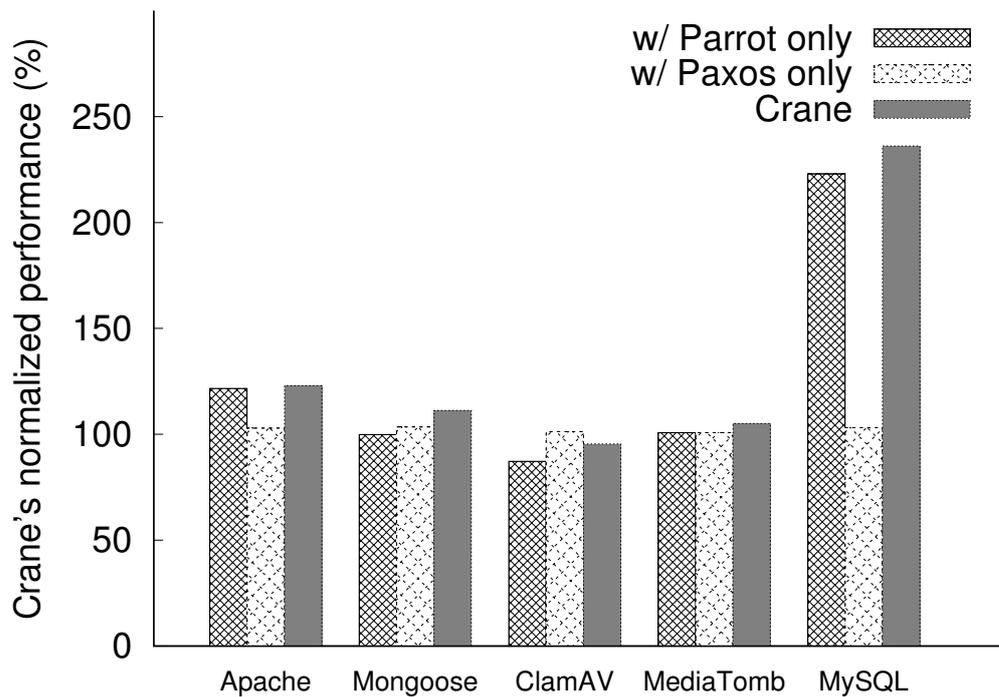


FIGURE 4.1: CRANE's performance normalized to un-replicated nondeterministic execution.

on the order of the PUT and GET request arriving at the primary's proxy. And then we ran Apache's un-replicated execution for 100 times on each replica, and three replicas reported "404 Not Found" for 6, 8, 11 times respectively.

4.3 Performance Overhead in Normal Case

To understand the performance impact of CRANE's components, we divided CRANE's components into two major parts: the DMT part ran by PARROT; and the proxy (with PAXOS) part which enforces the same sequence of client socket calls across replicas. Each part ran independently without the other part. The proxy part represents the performance overhead of invoking PAXOS consensus for client socket calls, and the DMT part represents the PARROT DMT scheduler's overhead.

Figure 4.1 shows the servers' performance running in CRANE normalized by their un-replicated nondeterministic executions. The mean overhead of CRANE for the five evaluated programs is 34.19% due to two main reasons. First, except for MySQL, which does fine-grained, per-table mutex and read-write locks frequently, the DMT schedules were efficient on the other four servers. The reason is that PARROT's scheduling primitives are already highly optimized for multi-core. The proxy-only part incurred 0.82%-3.46% overhead, which is not surprising, because the number of socket calls is much smaller than the number of Pthreads synchronizations in these programs. In short, CRANE's performance mainly depends on the DMT schedules' performance. MediaTomb incurred modest speedup because its transcoder mencoder had significant speedup with PARROT. We inspected MediaTomb's micro performance counters with the Intel VTune [2] profiling tool. When running in CRANE, MediaTomb only made 6.6K synchronization context switches,

while in the Pthreads runtime it made 0.9M synchronization context switches. This saving caused MediaTomb running with PARROT a 12.76% speedup compared to its nondeterministic execution. The PARROT evaluation [12] also observed a 49% speedup on the mencoder program.

The time bubbling technique saves most of needs on invoking consensus for the logical times of clients' socket operations, confirmed by the low frequency of inserted time bubbles in Table 1. Apache, MediaTomb, and Mongoose uses ApacheBench as its benchmark, and each request contained a connect(), send(), and close() call. ClamAV uses its own clamscan benchmark, and each request contained 18 socket calls. MySQL's benchmark contained 6-7 socket calls for each query. The ratio of inserted bubbles is merely 6.12%-33.35%. MediaTomb had the highest ratio of time bubbles because it took the longest time (9,703ms) to process each request. Note that the number of inserted time bubbles across replicas is the same within the same run of CRANE. Within different runs of CRANE, this number can be different because Wtimeout is a physical duration.

Bibliography

- [1] <http://www.clamav.net/>.
- [2] <http://software.intel.com/en-us/intel-vtune-amplifier-xe/>.
- [3] *Apache Web Server*. <http://www.apache.org>. 2012.
- [4] Aslan Askarov, Danfeng Zhang, and Andrew C. Myers. “Predictive Black-box Mitigation of Timing Channels”. In: *Proceedings of the 17th ACM conference on Computer and communications security (CCS '10)*. Oct. 2010.
- [5] Amittai Aviram et al. “Determinating Timing Channels in Compute Clouds”. In: *Proceedings of the 2010 ACM Workshop on Cloud Computing Security Workshop (CCSW '10)*. Oct. 2010.
- [6] Tom Bergan et al. “Deterministic process groups in dOS”. In: *Proceedings of the Ninth Symposium on Operating Systems Design and Implementation (OSDI '10)*. Oct. 2010.
- [7] William J. Bolosky et al. “Paxos Replicated State Machines As the Basis of a High-performance Data Store”. In: *Proceedings of the 8th USENIX Conference on Networked Systems Design and Implementation*. NSDI'11. Boston, MA: USENIX Association, 2011.
- [8] Mike Burrows. “The Chubby lock service for loosely-coupled distributed systems”. In: *Proceedings of the Seventh Symposium on Operating Systems Design and Implementation (OSDI '06)*. 2006, pp. 335–350.
- [9] Miguel Castro and Barbara Liskov. “Practical Byzantine Fault Tolerance”. In: *Proceedings of the Third Symposium on Operating Systems Design and Implementation (OSDI '99)*. Oct. 1999.
- [10] Tushar D. Chandra, Robert Griesemer, and Joshua Redstone. “Paxos Made Live: An Engineering Perspective”. In: *Proceedings of the Twenty-sixth Annual ACM Symposium on Principles of Distributed Computing (PODC '07)*. Aug. 2007.
- [11] CRIU. <http://criu.org>. 2015.
- [12] Heming Cui et al. “Parrot: a Practical Runtime for Deterministic, Stable, and Reliable Threads”. In: *Proceedings of the 24th ACM Symposium on Operating Systems Principles (SOSP '13)*. Nov. 2013.
- [13] Heming Cui et al. “REPFRAME: An Efficient and Transparent Framework for Dynamic Program Analysis”. In: *Proceedings of 6th Asia-Pacific Workshop on Systems (APSys '15)*. July 2015.
- [14] Zhenyu Guo et al. “Rex: Replication at the Speed of Multi-core”. In: *Proceedings of the 2014 ACM European Conference on Computer Systems (EUROSYS '14)*. ACM. 2014, p. 11.
- [15] Nicholas Hunt et al. “DDOS: Taming Nondeterminism in Distributed Systems”. In: *Eighteenth International Conference on Architecture Support for Programming Languages and Operating Systems (ASPLOS '13)*. 2013, pp. 499–508.

- [16] Manos Kapritsos et al. "All about Eve: Execute-Verify Replication for Multi-Core Servers." In: *Proceedings of the Tenth Symposium on Operating Systems Design and Implementation (OSDI '12)*. Vol. 12. 2012, pp. 237–250.
- [17] Ramakrishna Kotla et al. "Zyzyva: Speculative Byzantine Fault Tolerance". In: *Proceedings of the 21st ACM Symposium on Operating Systems Principles (SOSP '07)*. Oct. 2007.
- [18] Oren Laadan, Nicolas Viennot, and Jason Nieh. "Transparent, Lightweight Application Execution Replay on Commodity Multiprocessor Operating Systems". In: *Proceedings of the ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems (SIGMETRICS '10)*. June 2010, pp. 155–166.
- [19] Oren Laadan, Nicolas Viennot, and Jason Nieh. "Transparent, lightweight application execution replay on commodity multiprocessor operating systems". In: *ACM SIGMETRICS Performance Evaluation Review*. Vol. 38. 1. 2010, pp. 155–166.
- [20] Oren Laadan et al. "Pervasive Detection of Process Races in Deployed Systems". In: *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP '11)*. Oct. 2011.
- [21] Leslie Lamport. *Paxos made simple*. <http://research.microsoft.com/en-us/um/people/lamport/pubs/paxos-simple.pdf>.
- [22] Leslie Lamport. "The part-time parliament". In: *ACM Trans. Comput. Syst.* 16.2 (1998), pp. 133–169.
- [23] Dongyoon Lee et al. "Respec: efficient online multiprocessor replay via speculation and external determinism". In: *Fifteenth International Conference on Architecture Support for Programming Languages and Operating Systems (ASPLOS '10)*. Mar. 2010, pp. 77–90.
- [24] *libevent*. libevent.org/. 2015.
- [25] *LXC*. <https://linuxcontainers.org/>.
- [26] Yanhua Mao, Flavio Paiva Junqueira, and Keith Marzullo. "Mencius: building efficient replicated state machines for WANs". In: *Proceedings of the 8th USENIX conference on Operating systems design and implementation*. Vol. 8. 2008, pp. 369–384.
- [27] David Mazieres. *Paxos made practical*. Tech. rep. Technical report, 2007. <http://www.scs.stanford.edu/dm/home/papers,2007>.
- [28] *MediaTomb - Free UPnP MediaServer*. <http://mediatomb.cc/>. 2014.
- [29] *MEncoder*. <https://www.mplayerhq.hu/>. 2015.
- [30] *Mongoose*. <https://code.google.com/p/mongoose/>.
- [31] *MySQL*. <http://www.mysql.com/>.
- [32] *MySQL Replication*. <https://dev.mysql.com/doc/refman/5.0/en/replication.html>.
- [33] Jun Rao, Eugene J. Shekita, and Sandeep Tata. "Using Paxos to Build a Scalable, Consistent, and Highly Available Datastore". In: *Proc. VLDB Endow.* (Jan. 2011).
- [34] *SQLite*. <https://www.sqlite.org/>.
- [35] *SysBench: a system performance benchmark*. <http://sysbench.sourceforge.net>. 2004.

-
- [36] Robbert Van Renesse and Deniz Altinbuken. “Paxos Made Moderately Complex”. In: *ACM Computing Surveys (CSUR)* 47.3 (2015), 42:1–42:36.
 - [37] Danfeng Zhang, Aslan Askarov, and Andrew C. Myers. “Predictive Mitigation of Timing Channels in Interactive Systems”. In: *Proceedings of the 18th ACM conference on Computer and communications security (CCS '11)*. Oct. 2011.
 - [38] ZooKeeper. <https://zookeeper.apache.org/>.