



<b>Title</b>	<b>APUS: Fast and Scalable PAXOS on RDMA</b>
<b>Author(s)</b>	<b>Wang, C; Jiang, J; Chen, X; YI, N; Cui, H</b>
<b>Citation</b>	<b>ACM Symposium on Cloud Computing 2017 (SoCC '17), Santa Clara, CA, 24-27 September 2017. In Proceedings of SoCC '17, 2017, p. 94-107</b>
<b>Issued Date</b>	<b>2017</b>
<b>URL</b>	<b><a href="http://hdl.handle.net/10722/245447">http://hdl.handle.net/10722/245447</a></b>
<b>Rights</b>	<b>This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.</b>

# APUS: Fast and Scalable Paxos on RDMA

Cheng Wang  
The University of Hong Kong  
Hong Kong  
cwang2@cs.hku.hk

Jianyu Jiang  
The University of Hong Kong  
Hong Kong  
jyjiang@cs.hku.hk

Xusheng Chen  
The University of Hong Kong  
Hong Kong  
chenxus@hku.hk

Ning Yi  
The University of Hong Kong  
Hong Kong  
ezreal10@hku.hk

Heming Cui  
The University of Hong Kong  
Hong Kong  
heming@cs.hku.hk

## ABSTRACT

State machine replication (SMR) uses Paxos to enforce the same inputs for a program (e.g., Redis) replicated on a number of hosts, tolerating various types of failures. Unfortunately, traditional Paxos protocols incur prohibitive performance overhead on server programs due to their high consensus latency on TCP/IP. Worse, the consensus latency of extant Paxos protocols increases drastically when more concurrent client connections or hosts are added. This paper presents APUS, the first RDMA-based Paxos protocol that aims to be fast and scalable to client connections and hosts. APUS intercepts inbound socket calls of an unmodified server program, assigns a total order for all input requests, and uses fast RDMA primitives to replicate these requests concurrently.

We evaluated APUS on nine widely-used server programs (e.g., Redis and MySQL). APUS incurred a mean overhead of 4.3% in response time and 4.2% in throughput. We integrated APUS with an SMR system Calvin. Our Calvin-APUS integration was 8.2X faster than the extant Calvin-ZooKeeper integration. The consensus latency of APUS outperformed an RDMA-based consensus protocol by 4.9X. APUS source code and raw results are released on github.com/hku-systems/apus.

## CCS CONCEPTS

• **Computer systems organization** → **Reliability; Availability;**

## KEYWORDS

State Machine Replication, Fault Tolerance, Remote Direct Memory Access, Software Reliability

### ACM Reference Format:

Cheng Wang, Jianyu Jiang, Xusheng Chen, Ning Yi, and Heming Cui. 2017. **APUS: Fast and Scalable Paxos on RDMA**. In *Proceedings of SoCC '17, Santa Clara, CA, USA, September 24–27, 2017*, 14 pages. <https://doi.org/10.1145/3127479.3128609>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

SoCC '17, September 24–27, 2017, Santa Clara, CA, USA

© 2017 Association for Computing Machinery.

ACM ISBN 978-1-4503-5028-0/17/0...\$15.00

<https://doi.org/10.1145/3127479.3128609>

## 1 INTRODUCTION

State machine replication (SMR) runs the same program on replicas of hosts and invokes a distributed consensus protocol (typically, Paxos [53]) to enforce the same total order of inputs among replicas. Since the consensus on an input can be reached as long as a quorum (typically, majority) of replicas agree, SMR tolerates various errors, including hardware failures of minor replicas. SMR is deployed on clouds to make the metadata (e.g., leadership) of a distributed system highly available.

The strong fault-tolerance of SMR makes it an ideal high-availability service for general server programs. Recent SMR systems [29, 40, 49] use Paxos to enforce the same inputs for a server program, and they use advanced techniques (e.g., deterministic inter-thread synchronization [29, 78]) to make the program transit the same execution states across replicas. These SMR systems tolerate hardware failures for server programs.

Unfortunately, despite much effort, state-of-the-art still lacks a fast, scalable Paxos protocol for general server programs. A main reason is that traditional Paxos protocols [29, 66, 75] go through software network layers in OS kernels [72], which incurs high consensus latency. For efficiency, Paxos protocols typically take the Multi-Paxos approach [54]: it assigns one replica as the “leader” to invoke consensus requests, and the other replicas as “backups” to agree on requests. To agree on an input, at least one round-trip time (RTT) is required between the leader and a backup. Given that a ping RTT in LAN typically takes hundreds of  $\mu$ s, and that the request processing time of key-value store servers (e.g., Redis) is at most hundreds of  $\mu$ s, Paxos incurs high overhead in the response time of server programs.

Worse, the consensus latency of extant consensus protocols is often *scale-limited*: it increases drastically when the number of concurrent requests or replicas increases. For instance, the consensus latency of ZooKeeper [42] increases by 2.6X when the number of concurrent proposing requests increases from 1 to 20 (on 3 replicas). Scatter [37] shows that the consensus latency of its Paxos protocol increases by 1.6X when the number of replicas increases from 3 to 9.

Our evaluation found that the scalability problem in traditional consensus protocols mainly stem from OS kernels. We ran 4 popular consensus protocols [10, 21, 29, 75] on 24-core hosts with 40Gbps network (i.e., network bandwidth was not a bottleneck), we then ran 24 concurrent request connections. When the number of replicas increased from 3 to 9, the consensus latency of 3 protocols increased

by 105.4% to 168.3%, and 36.5% to 63.7% of the increase was in OS kernels.

As modern server programs tend to support more concurrent client connections, and advanced SMR systems tend to deploy more replicas (e.g., Azure [52] deploys seven or nine replicas) to support both replica failures and upgrades, the limited scalability in extant consensus protocols becomes even more pronounced.

Recent hardware-accelerated consensus protocols [33, 43, 44, 74] are effective on reducing consensus latency, but they are either unsuitable for general server programs or are not designed to be scalable on concurrent client connections. For instance, DARE [73], a novel consensus protocol, achieves the lowest consensus latency on a small number of client connections, but both its evaluation and ours show that its consensus latency increases quickly when more connections are added. Other recent works [32, 33, 55, 74] leverage the synchronous network ordering in a datacenter to safely skip consensus if packets arrive at replicas in the same order. These works require rewriting a server program to use their new libraries for checking the order of packets, so they are not designed to run legacy server programs.

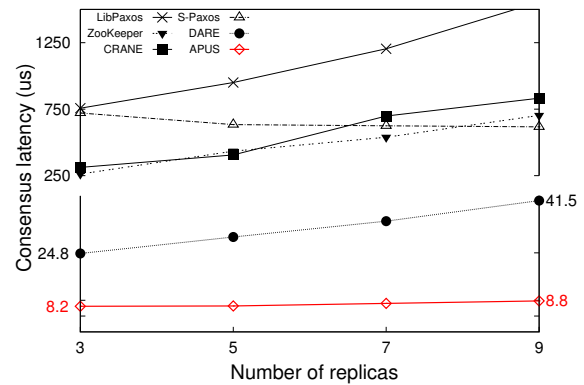
We argue that the problem of high, scale-limited consensus latency is not fundamental in Paxos. OS kernels, a major source of this problem, can be bypassed with advanced network features such as Remote Direct Memory Access (RDMA) within the same datacenter.

We present APUS,<sup>1</sup> the first RDMA-based Paxos protocol and runtime system. APUS intercepts an unmodified server program’s inbound socket calls (e.g., `recv()`), assigns a total order for all received requests in all connections, and uses fast RDMA primitives to invoke consensus on these requests concurrently. To ensure the same robustness as regular Paxos, APUS’s runtime system efficiently tackles several reliability challenges, including atomic delivery of messages (§4.2), transparent replication (§5.1), and failure recovery (§5.2).

A fast and scalable Paxos protocol, APUS has many practical applications, and we elaborate two below. First, it can be integrated into existing SMR systems (e.g., Calvin [78]), making the response time of a server program running in these systems almost as fast as the program’s unreplicated execution.

Second, it can support many server programs that are already well-tested or deterministic, including single-threaded ones such as Redis [76] and multi-processed ones such as Nginx [67] and MediaTomb [12]. Even if a program is pre-mature and undergoing debugging, enforcing the same order of inputs by APUS can still help debugging tools (e.g., PRES [71]) easily reproduce bugs. §3.2 further illustrates APUS’s broad applications.

We implemented APUS in Linux and compared it with five open source consensus protocols, including four traditional ones (libPaxos [75], ZooKeeper [10], CRANE [29] and S-Paxos [21]), and an RDMA-based one (DARE [73]). We evaluated APUS on nine widely used or studied programs, including 4 key-value stores (Redis [76], Memcached [62], SSDB [77], and MongoDB [65]), a SQL server MySQL [13], an anti-virus server ClamAV [26], a multimedia server MediaTomb [12], an LDAP server OpenLDAP [70], and



**Figure 1: Comparing APUS to five existing consensus protocols. All six protocols ran a client with 24 concurrent connections. The Y axis is broken to fit in all protocols.**

Calvin [78], a SMR-like database built on top of ZooKeeper [10]. Evaluation shows that

- (1) APUS is fast and scalable. Figure 1 shows that APUS’s consensus latency outperformed four traditional consensus protocols by at least 32.3X. Its consensus latency stayed almost constant to the number of concurrent requests and replicas. Its consensus latency was faster than DARE by 4.9X in average.
- (2) APUS is easy to work with SMR. The Calvin-APUS integration took only 39 lines of code. Calvin-APUS’s response time was 8.2X faster than the extant Calvin-ZooKeeper integration, and it incurred only 10.6% overhead in response time and 4.1% in throughput over Calvin’s unreplicated execution.
- (3) APUS achieves low overhead on real-world server programs. Compared to all nine server programs’ unreplicated executions, APUS incurred 4.3% overhead in response time and 4.2% in throughput.
- (4) It is robust on replicas failures and packet losses.

Our major contribution is the first Paxos protocol that achieves low performance overhead on diverse, widely-used server programs. A fast, scalable, and deployable Paxos protocol, APUS can widely promote the adoption of SMR and improve the fault-tolerance of various systems [20, 21, 29, 40, 48, 52] within a datacenter.

The remaining of this paper is organized as follows. §2 introduces Paxos and RDMA background. §3 gives an overview of APUS. §4 presents APUS’s consensus protocol with its runtime system. §5 presents implementation details. §6 compares APUS with DARE. §7 does evaluation, §8 discusses related work, and §9 concludes.

## 2 BACKGROUND

### 2.1 Paxos

Paxos [53, 54] enforces a total order of inputs for a program running across replicas. Because a consensus can be reached as long as a majority of replicas agree, Paxos is known for tolerating various

<sup>1</sup>We name our system after apus, one of the fastest birds.

faults, including hardware failures of minor replicas and packet losses.

SMR systems [29, 40] often use Paxos to replicate important online services. An typical SMR system contains two orthogonal parts: (1) a Paxos protocol that enforces a total order of inputs for the same program running across replicas; and (2) a technique (e.g., deterministic mutex locks [29, 78]) that makes the program transit same execution states on the same inputs.

The consensus latency of Paxos protocols is notoriously high and unscalable [10, 37]. As datacenters incorporate faster networking hardware and more CPU cores, traditional consensus protocols [10, 21, 29, 40, 75] are having fewer performance bottlenecks on network bandwidth and CPU resources.

However, software TCP/IP layers in OS kernels remain performance bottlenecks [72]. To quantify this bottleneck, we evaluated four traditional consensus protocols [10, 21, 29, 75] on 24-core hosts with 40Gbps network, and we spawned 24 concurrent consensus connections. When changing the replica group size from 3 to 9, although network and CPUs were not saturated, the consensus latency of 3 protocols drastically increased by 105.4% to 168.3% (Figure 1), and 36.5% to 63.7% of this increase was in OS kernel. When only one consensus connection was spawned, the latency increase on the number of replicas was more gentle (Table 2 in §7.1).

This evaluation shows that both the number of concurrent requests and replicas make consensus latency increase drastically. This problem becomes worse as server programs tend to support more concurrent requests and advanced SMR systems (e.g., Azure [52]) deploy seven to nine replicas to in case replica failures and upgrades.

## 2.2 RDMA

RDMA architectures (e.g., Infiniband [1] and RoCE [8]) become common within a datacenter due to its ultra low latency, high throughput, and its decreasing prices. The ultra low latency of RDMA not only comes from bypassing the OS kernel, but also its dedicated network stack implemented in hardware. Therefore, RDMA is considered the fastest kernel bypassing technique [46, 64, 73]; it is several times faster than software-only kernel bypassing techniques (e.g., DPDK [7] and Arrakis [72]).

RDMA has three operation types, from fast to slow: one-sided read/write operations, two sided send/recv operations, and IPoIB (IP over Infiniband). IPoIB runs unmodified socket programs, but it is a few times slower than the other two types. A one-sided RDMA write can directly write from one replica’s memory to a remote replica’s memory without involving the remote OS kernel or CPU. Prior work [64] shows that one-sided operations are up to 2X faster than two-sided operations [47], so APUS uses one-sided operations (or “WRITE” in this paper). On a WRITE success, the remote NIC (network interface card) sends an RDMA ACK to local NIC.

A one-sided RDMA communication between a local and a remote NIC has a Queue Pair (QP), including a send queue and a receive queue. Such a QP is a global data structure between every two replicas, but pushing a message into a local QP takes at most 0.2 μs in our evaluation. Different QPs between different replicas work in parallel (leveraged by APUS in §4.1). Each QP has a Completion Queue (CQ) to store ACKs. A QP belongs to a type of “XY”: X can

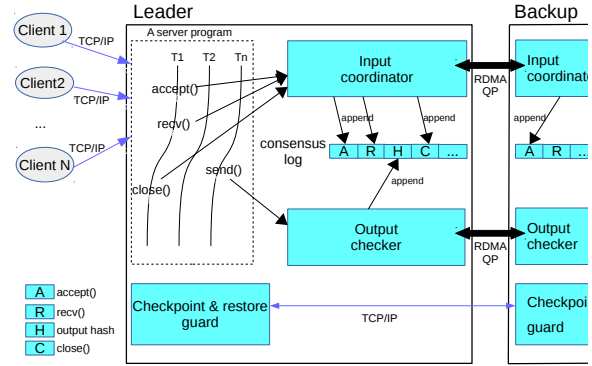


Figure 2: APUS Architecture (key components are in blue).

be R (reliable) or U (unreliable), and Y can be C (connected) or U (unconnected). HERD [46] shows that WRITES on RC and UC OPs incur almost the same latency, so APUS uses RC QPs.

Normally, to ensure a WRITE resides in remote memory, the local replica busily polls an ACK from the CQ before it proceeds (or *signaling*). Polling ACK is time consuming as it involves synchronization between the NICs on both sides of a CQ. We looked into the ACK pollings in a recent RDMA-based consensus protocol DARE [73]. We found that, although it is highly optimized (its leader maintains one global CQ to receive all backups’ ACKs in batches), busily polling ACKs slowed DARE down (§7.3): when the CQ was empty, each poll took 0.039~0.12 μs; when the CQ has one or more ACKs, each poll took 0.051~0.42 μs.

Fortunately, depending on protocol logic, one can do *selective signaling* [46]: it only checks for an ACK after pushing a number of WRITES. Because APUS’s protocol logic does not rely on RDMA ACKs, it just occasionally invokes selective signaling to clean up ACKs.

## 3 OVERVIEW

### 3.1 APUS Architecture

APUS deployment is similar to a typical SMR’s: it runs a program on replicas within a datacenter. Replicas connect with each other using RDMA QPs. Client programs located in LAN or WAN. The APUS leader handles client requests and runs its RDMA-based protocol to enforce the same total order for all requests across replicas.

Figure 2 shows APUS’s architecture. APUS intercepts a server program’s inbound socket calls (e.g., `recv()`) using a Linux technique called `LD_PRELOAD`. APUS involves four key components: a Paxos consensus protocol for input coordination (in short, the *coordinator*), a circular in-memory consensus log (the *log*), a guard process that handles checkpointing and recovering a server’s process and file system state (the *guard*), and an optional output checking tool (the *checker*).

The coordinator is involved when a thread of a program running on the APUS leader calls an inbound socket call (e.g., `recv()`). The thread executes the Libc call, gets the received data, appends a log

entry on the leader’s local consensus log, and replicates this entry to backups’ consensus logs using our Paxos protocol (§4).

In this protocol, all threads in the server program running on the leader replica can concurrently invoke consensus on their log entries (requests), but APUS enforces a total order for all entries in the leader’s local consensus log. As a consensus request, each thread does an RDMA WRITE to replicate its log entry to the corresponding log entry position on all APUS backups. Each APUS backup polls from the latest unagreed entry on its local consensus log; if it agrees with the proposed log entry, it does an RDMA WRITE to write a consensus reply on the leader’s corresponding entry.

To ensure Paxos safety [60], all APUS backups agree on the entries proposed from the leader in a total order without allowing any entry gap. When a majority of replicas (including the leader) has written a consensus reply on the leader’s local entry, this entry has reached a consensus. By doing so, APUS consistently enforces the same consensus log for both the leader and backups. §4.5 presents a proof sketch on the correctness of the protocol, and §4.6 analyzes why it is fast and scalable.

The output checker is periodically invoked as a program replicated in APUS executes outbound socket calls (e.g., `send()`). For every 1.5KB (MTU size) of accumulated outputs per connection, the checker unions the previous hash with current outputs and computes a new CRC64 hash. For simplicity, the output checker uses APUS’s input consensus protocol (§4) to compare hashes across replicas.

Our evaluation found that the output checker had negligible performance impact and all output divergence were due to physical times (§7.4). This suggests that many server programs are well-tested, and the output checker can be turned on only in program debug phase. If APUS is integrated into an SMR system, the output checker is not needed because SMR already has techniques to enforce the same program executions.

A guard runs on each APUS replica to cope with replica management, including checkpointing program states and adding/recovering replicas (§5.2).

## 3.2 Motivating Applications of APUS

**Building fast SMR systems.** Extant SMR systems (e.g., CRANE [29], Rex [40], and Calvin [78]) use TCP/IP-based consensus protocols, thus they incur high overhead in server programs’ response time. APUS can greatly alleviate this overhead. Evaluation (§7.2) shows that the response time of our Calvin-APUS integration on realistic SQL workloads was 8.2X than its extant Calvin-ZooKeeper integration. Compared to Calvin’s unreplicated execution, APUS incurred only 10.6% overhead in response time and 4.1% in throughput.

**Improving the availability of server programs.** Many real-world server programs handle online requests and store important data, so they naturally demand high availability against hardware failures. Many programs are suitable to run with APUS because they are already well-tested or deterministic (e.g., single-threaded ones such as Redis and multi-processed ones such as Nginx and MediaTomb). Other orthogonal techniques such as deterministic multithreading [16, 30, 56, 78] can be combined with APUS to make a replicated server program behave the same on the same inputs.

Our evaluation (§7.4) shows that, compared to all nine evaluated programs’ unreplicated executions at peak performance, APUS incurs 4.2% overhead in throughput and 4.3% in response time.

**Improving debugging efficiency.** Even if a server program is under development and may contain nondeterministic concurrency bugs, APUS can still benefit extant debug tools [14, 50, 71] because these tools often require extra mechanisms to frequently replay the same total order of inputs. APUS logs program inputs persistently, and it can efficiently replay these inputs in the same order when integrated into debug tools (e.g., PRES [71]).

## 4 THE RDMA-BASED PAXOS PROTOCOL

### 4.1 Normal Case

APUS’s consensus protocol has three main elements. First, a Paxos consensus log. Second, threads of a server program running on the leader host (or *leader threads*). APUS hooks the inbound socket calls (e.g., `recv()`) of these leader threads and invoke consensus requests on these calls. We denote the data received from each of these calls as a consensus request (i.e., an entry in the consensus log). Third, an APUS internal thread running on every backup (or *backup threads*), which agrees on consensus requests. The APUS leader enables the first and second elements, and backups enable the first and third elements.

```
struct log_entry_t {
    consensus_ack reply[MAX]; // Per replica consensus reply.
    viewstamp_t vs;
    viewstamp_t last_committed;
    int node_id;
    viewstamp_t conn_vs; // client connection ID.
    int call_type; // socket call type.
    size_t data_sz; // data size in the call.
    char data[0]; // data, with a canary value in the last byte.
} log_entry;
```

Figure 3: APUS’s log entry for each socket call.

Figure 3 depicts the format of a log entry in APUS’s consensus log. Most fields are the same as those in a typical Paxos protocol [60] except three: the `reply` array, `conn_vs`, and `call_type`. The `reply` array is a piece of memory on the leader side, preserved for backups to do RDMA WRITES for their consensus replies. The `conn_vs` is for identifying which TCP connection this socket call belongs to (see §4.3). The `call_type` identifies different types of socket calls (e.g., the `accept()` type and the `recv()` type) for the entry.

Figure 4 shows APUS’s consensus protocol. Suppose a leader thread invokes a consensus request when it calls a socket call `recv()`. This thread’s consensus request has four steps. The first step (L1, not shown in Figure 4) is executing the actual socket call, because the thread needs to get the received data and returned value, to allocate a distinct log entry, and to replicate the entry in backups’ consensus logs.

The second step (L2) is local preparation, including assigning a `viewstamp` (a totally-ordered Paxos consensus request ID [60]) for this entry in the consensus log, allocating a distinct entry in the log, and storing the entry to a local storage. We denote the time taken on storing an entry as  $t_{SSD}$ .

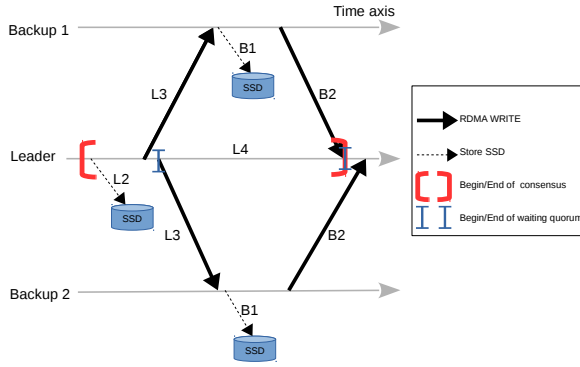


Figure 4: APUS consensus algorithm in normal case.

Third, each leader thread concurrently invokes a consensus via the third step (L3): WRITE the log entry to remote backups. This step is thread-safe because each leader thread works on its own distinct entry and remote backups’ corresponding entries. An L3 WRITE returns quickly after pushing the entry to its local QP connecting the leader and each backup. We denote the time taken for this push as  $tp_{USH}$ , which took at most  $0.2 \mu s$  in our evaluation.  $tp_{USH}$  is serial for concurrently arriving requests on each QP, but the WRITES (all L3 arrows in Figure 3) to different QPs run in parallel.

The fourth step (L4) is that the leader thread polls on its reply field in its local log entry to wait for backups’ consensus replies. It breaks the poll if a number of heartbeats fail (§4.4). If a majority of replicas agrees on the entry, an input consensus is reached, the leader thread leaves this `recv()` call and proceeds with its program logic.

On each backup, a backup thread polls from the latest unagreed log entry. It breaks the poll if a number of heartbeats fail (§4.4). If no heartbeat fails, the backup thread then agrees on entries in the same total order as those on the leader’s consensus log, using three steps. First (B1), it does a regular Paxos view ID check [60] to see whether the leader’s view ID matches its own one, it then stores the log entry in its local SSD. To scale to concurrently arriving requests, the backup thread scans multiple entries it agrees with at once. It then stores them in APUS’s parallel storage.

Second (B2), on each entry the backup agrees, the backup thread does an RDMA WRITE to send back a consensus reply to the reply array element in the leader’s corresponding entry. Third (B3, not shown in Figure 4), the backup thread does a regular Paxos check [60] on `last_committed` and to know the latest entry that has reached consensus. It then “executes” the committed entries by forwarding the data in these entries to the server program on its local replica. Carrying latest committed entries in next consensus requests is a common, efficient Paxos implementation method [60].

To ensure Paxos safety, the backup thread agrees on log entries in order without allowing any gap [60]. If the backup suspects it misses some log entries (e.g., because of packet loss), it invokes

a learning request to the leader asking for the missing entries. We found one backup thread per backup suffices to achieve low overhead on concurrent connections (§7.3).

## 4.2 Atomic Message Delivery

On a backup side, one tricky challenge is that atomicity must be ensured on the leader’s RDMA WRITES on all entries and backups’ polls. For instance, while a leader thread is doing a WRITE on vs to a remote backup, the backup’s thread may be reading vs concurrently, causing a corrupted read value.

To address this challenge, one prior approach [34, 46] leverages the left-to-right ordering of RDMA WRITES and puts a special non-zero variable at the end of a fix-sized log entry because they mainly handle key-value stores with fixed value length. As long as this variable is non-zero, the RDMA WRITE ordering guarantees that the log entry WRITE is complete. However, because APUS aims to support general server programs with largely variant received data lengths, this approach cannot be applied in APUS.

Another approach is using atomic primitives provided by RDMA hardware, but a prior evaluation [80] has shown that RDMA atomic primitives are much slower than normal RDMA WRITES and local memory reads.

APUS tackles this challenge by using the leader to add a canary value after the data array. A backup thread always first checks the canary value according to `data_size` and then starts a standard Paxos consensus reply decision [60]. This synchronization-free approach ensures that a APUS backup thread always reads a complete entry efficiently.

## 4.3 Handling Concurrent Connections

Unlike traditional Paxos protocols which mainly handle single-threaded programs due to the deterministic execution assumption in SMR, APUS aims to support both single-threaded as well as multi-threaded or -processed programs running on multi-core machines. Therefore, a strongly consistent mechanism is needed to map each connection on the leader and its corresponding connection on backups. A naive approach is matching a leader connection’s socket descriptor to the same one on a backup, but programs on backups may return nondeterministic descriptors due to systems resource contention.

Fortunately, Paxos already makes viewstamps [60] of requests (log entries) strongly consistent across replicas. For TCP connections, APUS adds the `conn_vs` field, the viewstamp of the the first socket call in each connection (i.e., `accept()`) as the connection ID for log entries.

## 4.4 Leader Election

Leader election on RDMA raises a main challenge: because backups do not communicate with each other in normal case, a backup proposing itself as the new leader does not know the remote memory locations where the other backups are polling. Writing to a wrong remote memory location may cause the other backups to miss all leader election messages. A recent system [73] establishes an extra control QP to handle leader election, complicating deployments.

APUS addresses this challenge with a simple, clean design. It runs leader election on the normal-case consensus log and QP. In normal case, the leader does WRITES to remote logs as heartbeats with a period of  $T$ . Each consensus log maintains an `elect` [MAX] array, one array element for each replica. This `elect` array is only used in leader election. Once backups miss heartbeats from the leader for  $3^*T$ , they suspect the leader to fail, close the leader's QPs, and start to work on the `elect` array to elect a new leader.

Backups use a standard Paxos leader election algorithm [60] with three steps. Each backup writes to its own `elect` element indexed by its replica ID on other replicas' `elect`. First, each backup waits for a random time (similar to random election timeouts in Raft [68]), and it proposes a new view with a standard two-round Paxos consensus [54] by including both its view and the index of its latest log entry. The other backups also propose their views and poll on this `elect` array in order to agree on an earlier proposal or confirm itself as the winner. The backup with a more up-to-date log will win the proposal. A log is more up-to-date if its latest entry has either a higher view or the same view but a higher index.

Second, the winner proposes itself as a leader candidate using this `elect` array. Third, after the second step reaches a quorum, the new leader notifies remote replicas itself as the new leader and it starts to WRITE periodic heartbeats. Overall, APUS safely avoids multiple "leaders" to corrupt consensus logs, because only one leader is elected in each view, and backups always close an outdated leader's QPs before electing a new leader. For robustness, the above three steps are inherited from a practical Paxos election algorithm [60], but APUS makes the election efficient and simple in an RDMA domain.

#### 4.5 Correctness

APUS's protocol derives from Paxos Made Practical (PMP) [60], a practical viewstamp-based Paxos protocol. We made this design choice because Paxos is notoriously difficult to understand [53, 54, 79], implement [25, 60], and verify [39, 82]. Deriving from a practical protocol [60] helps us incorporate these readily mature understanding and theoretically verified safety rules into APUS.

We made two major modifications from PMP. *Modification 1*: APUS replicas use the faster and more scalable one-sided RDMA WRITE to replicate log entries (§4.1) and do leader elections (§4.4). *Modification 2*: to prevent outdated leaders from messing up log entries, APUS's backups conservatively close the QP with the outdated leader right after suspecting it has failed (§4.4).

These two modifications empower APUS's protocol to comply with Paxos *safety* guarantee: all replicas see the same total order of request entries in their local consensus logs. We will show that APUS satisfies three properties: (1) *leader completeness*: all agreed entries should present in the logs of subsequent leaders; (2) *log matching*: two replicas' logs cannot have different *agreed entries* (entries agreed by a majority) on the same log position; (3) *data integrity*: all replicas cannot read corrupted data from log entries. The first two properties are widely considered sufficient to ensure safety [68] in traditional TCP/IP based Paxos protocols, while we add data integrity because our protocol is based on RDMA.

The leader completeness property of APUS inherits from PMP. APUS follows PMP's view change protocol to ensure at most one

leader in each view and the newly elected leader is most up-to-date (§4.4). Therefore, we admit this property and omit the proof here. Below, we are going to prove the log matching and the data integrity properties.

**Log matching.** We prove the log matching property by induction on the view number. The base case holds because the initial log is empty. For the inductive case, we hypothesize that the property holds for the views up to  $v$  and prove the property still holds in the next view. We can safely assume there is a leader (*the current leader*) in the new view, because if no leader presents, no new entries can get agreed and the property will not be broken.

The inductive case can be proved in two steps. First, since the current leader is elected, a majority of replicas must have closed QPs with old leaders to prevent them writing to their logs (modification 2). Therefore, only the current leader can replicate its log entries to a majority of replicas and get the entries agreed. Second, when one leader thread replicates data messages in one log entry, it WRITES to the same position on remote backups' logs (§4.1). Therefore, the newly replicated entries to be agreed are all identical to those of the current leader. Combining the two steps, the newly agreed entries in the new view are all identical. Therefore, the inductive case is proved and thus the log matching property holds.

**Data integrity.** All replicas will not read corrupted entries because of three reasons. First, each leader thread replicates log entries to disjoint memory addresses (§4.1). Therefore, the replication mechanism is thread-safe and will not cause contentions. Second, RDMA provides error detection mechanisms to prevent data corruption during network transmission, which is a basic requirement of Paxos deployment. Third, APUS has an atomic log entry read/write mechanism (§4.2) to prevent replicas from reading incomplete log entries. These three factors work together to ensure all replicas see correctly replicated, non-corrupted log entries and hence the property holds.

#### 4.6 Analytical Analysis on Performance

APUS is designed to be scalable to the number of concurrent client connections for general server programs. In contrast, a recent RDMA-based protocol DARE [73] is designed to achieve the lowest latency on a small number of connections for its own key-value store server. Below is an analytical analysis on APUS's consensus latency, and we compare APUS and DARE in §6.1.

Suppose the APUS leader has  $N$  client connections, and  $N$  requests arrive at the same time. APUS invokes consensus on all requests in the same way without distinguishing them as "read only" or "write". Suppose there are only three replicas.

According to the leader's four steps **L1~L4**, to reach consensus for all these  $N$  requests, the time taken on the leader's  $i_{th}$  request includes five parts: (1) an SSD storage time  $t_{SSD}$  in **L2** (each leader thread does a SSD store in parallel); (2) because an RDMA QP is a global data structure between every two replicas, pushing a message to a QP is serialized, which costs  $i \times t_{PUSH}$  for  $i_{th}$  request; (3) a  $\frac{1}{2}t_{RTT}$  in **L3**; (4) an SSD storage time  $t_{SSD}$  in **B1** for each backup (done by backups in parallel); and (5) a  $\frac{1}{2}t_{RTT}$  in **B2**. On APUS's leader, the average consensus latency for all  $N$  requests sums up as the equation below:

$$\begin{aligned}
 APUS &= \left( \sum_{i=1}^N (2t_{SSD} + i \times t_{PUSH} + t_{RTT}) \right) / N \\
 &= 2t_{SSD} + \frac{(N+1)}{2} t_{PUSH} + t_{RTT}
 \end{aligned}
 \tag{1}$$

This equation shows that APUS’s consensus latency is scalable to  $N$  because  $t_{PUSH}$  is often below  $0.2 \mu s$  (§2.2).

## 5 IMPLEMENTATION DETAILS

### 5.1 Replicating an unmodified server program

To replicate an unmodified server program, APUS leverages a Linux technique called `LD_PRELOAD`. This technique enables interception of `libc` function calls and customized code injection. With this mechanism, APUS intercepts all the `libc` inbound socket calls and invokes the RDMA-based consensus protocol (§4) to replicate the inputs from the leader to backups. By doing so, any server program using the POSIX socket can run in APUS without being modified.

### 5.2 Checkpoint and Restore

To handle replica failures, a Paxos protocol must provide a persistent input logging storage. We explicitly designed the input logging storage mechanism in APUS to be thread-safe and scalable to the number of concurrent client connections. Specifically, the input logging operation bypasses the kernel cache, and is directly applied onto the disk.

In APUS, the guard process on a backup replica checkpoints the local server program’s process state and file system state (including the input logging storage) of current working directory within a one-minute duration.

Such a checkpoint operation and its duration is not sensitive to normal case performance because it is invoked on one backup replica, and hence the other backups can still reach quorum rapidly. Each checkpoint is associated with a last committed socket call viewstamp of the server program. After each checkpoint, the backup dispatches the checkpoint zip file to the other replicas.

Specifically, APUS leverages CRIU [28], a popular, open source tool, to checkpoint a server program’s process state (e.g., CPU registers and memory). Since CRIU does not support checkpointing RDMA connections, APUS’s guard first sends a “close RDMA QP” request to an APUS internal thread, lets this thread closes all remote RDMA QPs, and then invokes CRIU.

### 5.3 Network Output Checking Tool

Server programs often send replies with non-blocking IO. To align outputs across replicas, APUS uses a bucket-based hash computation mechanism. When a server calls a `send()` call, APUS puts the sent bytes into a local, per-connection bucket with 1.5KB (MTU size). Whenever a bucket is full, APUS computes a new CRC64 hash on a union of the current hash and this bucket. To compare a hash across replicas, the output checker uses APUS’ input consensus protocol (§4.1). Because this protocol is invoked rarely, we did not observe its performance impact. The output checker is mainly for server programs’ development purpose (§3.1).

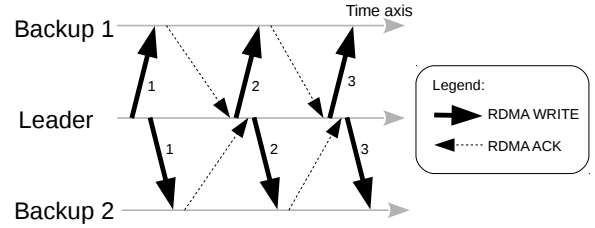


Figure 5: *DARE’s RDMA-based protocol*. It is a sole-leader, two-round protocol with three steps: (1) the leader WRITES a consensus request to all backups’ consensus logs and waits for ACKs to check if they succeed; (2) for the successful backups in (1), the leader does WRITES to update tail pointer of their consensus logs; and (3) on receiving a majority of ACKs in (2), a consensus is reached, the leader does WRITES to notify backups.

## 6 DISCUSSIONS

### 6.1 Comparing APUS with DARE

DARE [73] deviates from Paxos due to its centralized, sole-leader protocol: in normal case, the leader does all consensus work via RDMA, and the other replicas are silent and do not consume CPU. Figure 5 shows DARE’s protocol with two-rounds: first, leader does RDMA WRITES of consensus requests on each replica; second, leader does RDMA WRITES on each replica to update a global variable that points to the latest request (tail of consensus log) in each backup. DARE backups are silent in both rounds, and only their RDMA NICs send back RDMA ACKs to the leader’s NIC. Because the second round updates a global variable on every backup, which *serializes* all consensus requests, DARE is not designed to be scalable to concurrent connections.

DARE is mainly designed to achieve the lowest consensus latency on a small number of concurrent key-value connections. To this end, it has two clever features. First, on an input consensus, DARE needs to store the input only once on the leader, because its backups are silent. In current DARE implementation, leader does not store inputs and works purely in-memory. Second, it batches SET and GET requests separately. For GET requests, leader does only one-round RDMA READs to check view IDs from backups. Both DARE’s evaluation and ours (§7.3) show that, when there were at most six concurrent connections, DARE achieved the lowest consensus latency in extant evaluation [33, 43, 44, 55, 74].

Despite the two features, the serialization problem in DARE still affects its scalability, especially when many SET and GET requests arrive concurrently. DARE’s evaluation [73] confirmed this problem: on three replicas and nine concurrent connections, DARE’s



throughput on the 50% SET and 50% GET randomly arrival workload was 43.5% lower than that on the 100% SET workload. Our evaluation (§7.3) reproduces this problem when increasing the number of concurrent connections from 1 to 24: DARE’s consensus latency increased approximately linearly to the number of connections; APUS’s consensus latency was faster than DARE by 4.9X in average.

Overall, APUS differs from DARE in three aspects. First, APUS is a Paxos protocol for general server programs; DARE is a novel, sole-leader consensus protocol for its own key-value store. Second, APUS is designed to be scalable on many concurrent client connections; DARE is mainly designed to achieve lowest consensus latency on a smaller number of connections. Third, APUS is a persistent protocol; DARE currently works purely in-memory. These differences show that APUS is more suitable for general server programs, and DARE more suitable for maintaining metadata.

## 6.2 APUS Limitations

APUS currently does not replicate physical times such as `time()` because these physical results are often explicit and easy to examine from network outputs (e.g., a timestamp in the header of a reply). Existing Paxos approaches [49, 60] can be leveraged to intercept these functions and make programs produce the same results among replicas.

To replicate general client requests [29, 73], APUS totally orders all types of requests and it has not incorporated read-only optimization [49], because its performance overhead is already low (§7.4).

## 7 EVALUATION

Evaluation was done on nine RDMA-enabled Dell R430 and five Supermicro SuperServer 1019P hosts. Each host has Linux 3.16.0 and 2.6 GHz Intel Xeon CPU. The Dell R430 hosts are equipped with 24 hyperthreading cores, 64 GB memory, and 1 TB SSD. The SuperServer 1019P hosts have 28 hyperthreading cores, 32 GB memory, and 375GB SSD. All NICs are Mellanox ConnectX-3 (40Gbps) connected with RoCE [8]. All programs’ unreplicated executions run on IPoIB (§2.2). Workloads run on idle replicas.

We compared APUS with five open source consensus protocols, including four traditional ones (libPaxos [75], ZooKeeper [10], CRANE [29] and S-Paxos [21]) and an RDMA-based one (DARE [73]). S-Paxos is designed to achieve scalable throughput on more replicas.

We evaluated APUS on nine widely used or studied programs, including 4 key-value stores Redis, Memcached, SSDB, MongoDB; MySQL, a SQL server; ClamAV, an anti-virus server that scans files and delete malicious ones; MediaTomb, a multimedia storage server that stores and transcodes video and audio files; OpenLDAP, an LDAP server; Calvin [78], a popular SMR system for databases. We picked Calvin because: (1) it replicates inputs with a highly-engineered consensus protocol ZooKeeper [10], a good comparison target for APUS; and (2) it implements deterministic synchronization, which can make a program run deterministically. Table 1 shows workloads. The rest of this section focuses on five questions:

- §7.1: Is APUS much faster than traditional consensus?
- §7.2: How easy is APUS to integrate into SMR systems?
- §7.3: How faster is APUS compared to DARE?

**Table 1: Benchmarks and workloads. “Self” in the Benchmark column means we used a program’s own benchmark.**

Program	Benchmark	Workload/Input
ClamAV	clamscan [3]	Files in /lib from a replica
MediaTomb	ApacheBench [11]	Transcoding videos
Memcached	mcperf [2]	50% set, 50% get operations
MongoDB	YCSB [5]	Insert operations
MySQL	Sysbench [4]	SQL transactions
OpenLDAP	Self	LDAP queries
Redis	Self	50% set, 50% get operations
SSDB	Self	Eleven operation types
Calvin	Self	SQL transactions

§7.4: What is the performance overhead of running APUS with server programs? How well does it scale?

§7.5: How stable is APUS’s performance in a congested network?

§7.6: How well does APUS handle replica failures?

## 7.1 Comparing w/ Traditional Consensus

We ran APUS and four traditional consensus protocols using their own client programs or popular client programs with 100K requests of similar sizes. For each protocol, we ran a client with 24 concurrent connections on a 24-core machine located in LAN, and we used up to nine replicas. Both the number of concurrent connections and replicas are common high values [10, 29, 40, 73].

All four traditional protocols were run on IPoIB (§2.2). Figure 1 shows that the consensus latency of three traditional protocols increased almost linearly to the number of replicas (except S-Paxos). S-Paxos batches requests from replicas and invokes consensus when the batch is full. More replicas can take shorter time to form a batch, so S-Paxos incurred a slightly better consensus latency with more replicas. Nevertheless, its latency was always over 600  $\mu$ s. APUS’s consensus latency outperforms these four protocols by at least 32.3X.

To find scalability bottlenecks in traditional protocols, we used only one client connection and broke down their consensus latency on leader (Table 2). From 3 to 9 replicas, the consensus latency (the “Latency” column) of these protocols increased more gently than that on 24 concurrent connections. For instance, when the number of replicas increased from three to nine, ZooKeeper latency increased by 30.3% with one connection; this latency increased by 168.3% with 24 connections (Figure 1). This indicates that concurrent consensus requests are the major scalability bottleneck for these protocols.

Specifically, three protocols had scalable latency on the arrival of their first consensus reply (the “First” column), which implies that network is not saturated. libPaxos is an exception because its two-round protocol consumed much bandwidth. However, on the leader, there is a big gap between the arrival of the first consensus reply and the “majority” reply (the “Major” column). Given that the replies’ CPU processing time was small (the “Process” column), we can see that various systems layers, including OS kernels, network libraries, and language runtimes (e.g., JVM), are another major

**Table 2: Performance breakdown of traditional protocols on leader with only one connection.** The “Proto-#Rep” column is the protocol name and replica group size; “Latency” is the consensus latency; “First” is the latency of leader’s first received consensus reply; “Major” is the latency of leader’s consensus; “Process” is leader’s time spent in processing all replies; and “Sys” is leader’s time spent in systems (OS kernel, network stacks, and JVM) between the “First” and “Major” reply. Times are in  $\mu$ s.

Proto-#Rep	Latency	First	Major	Process	Sys
libPaxos-3	81.6	74.0	81.6	2.5	5.1
libPaxos-9	208.3	145.0	208.3	12.0	51.3
ZooKeeper-3	99.0	67.0	99.0	0.84	31.2
ZooKeeper-9	129.0	76.0	128.0	3.6	49.4
CRANE-3	78.0	69.0	69.0	13.0	0
CRANE-9	148.0	83.0	142.0	30.0	35.0
S-Paxos-3	865.1	846.0	846.0	20.0	0
S-Paxos-9	739.1	545.0	731.0	35.0	159.1

scalable bottleneck (the “Sys” column). This indicates that RDMA is useful on bypassing systems layers.

Both CRANE and S-Paxos’s leader handles consensus replies rapidly, so they two had same “First” and “Major” arrival times (i.e., “Sys” times were 0 on three replicas).

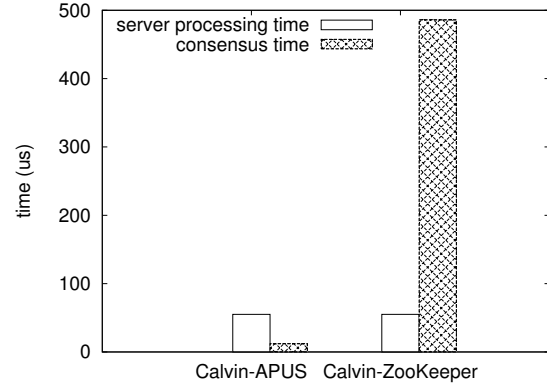
### 7.2 Integrating APUS into Calvin

Calvin [78] is a SMR-like distributed database which leverages ZooKeeper [10] for replicating client requests to achieve high availability. We replaced ZooKeeper with APUS in Calvin to replicate inputs and compared the performance of Calvin-ZooKeeper and Calvin-APUS.

The Calvin-APUS integration took 39 lines of code. Calvin currently uses ZooKeeper to batch inputs and then replicate them. To reduce response time, Calvin-APUS replicates each request immediately on its arrival. Figure 6 shows that the consensus latency of ZooKeeper was 7.6X higher than Calvin’s own request processing time, which indicates that ZooKeeper added a high overhead in Calvin’s response time. Calvin-APUS’s response time was 8.2X faster than Calvin-ZooKeeper’s because APUS’s consensus latency was 45.7X faster than ZooKeeper’s. Calvin’s unreplicated execution throughput is 19825 requests/s, and Calvin-ZooKeeper was 16241 requests/s. Calvin-APUS was 19039 requests/s, a 4.1% overhead over Calvin’s unreplicated execution.

### 7.3 Comparing with DARE

Because DARE only supported a key-value server written by the authors, we ran APUS with Redis, a popular key-value server for comparison. Figure 7 shows APUS and DARE’s consensus latency on variant concurrent connections. Both APUS and DARE ran seven replicas with randomly arriving, update-heavy (50% SET and 50% GET) and read-heavy (10% SET and 90% GET) workloads. DARE performance on two workloads were different because it handles GETs with only one consensus round [73]. APUS handles all requests with the same protocol. When there was only one

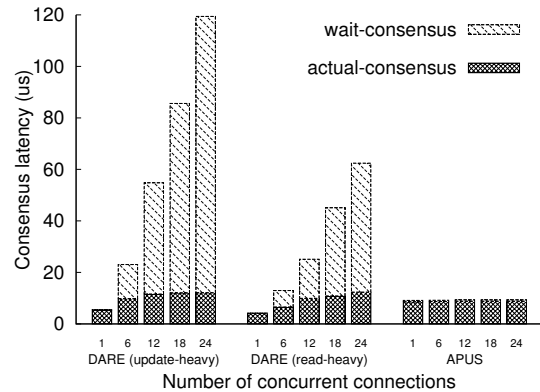


**Figure 6: Comparing Calvin-APUS and Calvin-ZooKeeper.**

connection, DARE achieved the lowest consensus latency we have seen in prior work because it is a sole-leader protocol (§6.1). On variant connections, APUS’s average consensus latency was faster than DARE by 4.9X for two main reasons.

First, APUS is a one-round protocol and DARE is a two-round protocol (for SETs), so DARE’s “actual-consensus” time was 53.2% higher than APUS. Even using read-heavy workloads (DARE uses one-round for GETs) with APUS, APUS’s actual consensus time was still slightly faster than DARE’s on over six connections, because APUS avoids expensive ACK pollings (§2.2).

Second, DARE’s second consensus round updates a global variable for each backup and serializes consensus requests (§6.1). Although DARE mitigates this limitation by batching same SET or GET types, randomly arriving requests often break batches, causing a large “wait-consensus” time (a new batch can not start consensus until prior batches reach consensus). DARE evaluation [73] confirmed such a high wait duration: with three replicas and nine concurrent connections, DARE’s throughput on real-world inspired workloads (50% SET and 50% GET arriving randomly) was 43.5%



**Figure 7: APUS and DARE consensus latency (divided into two parts) on variant connections.** “Wait-consensus” is the time an input request spent on waiting consensus to start. “Actual-consensus” is the time spent on running consensus.

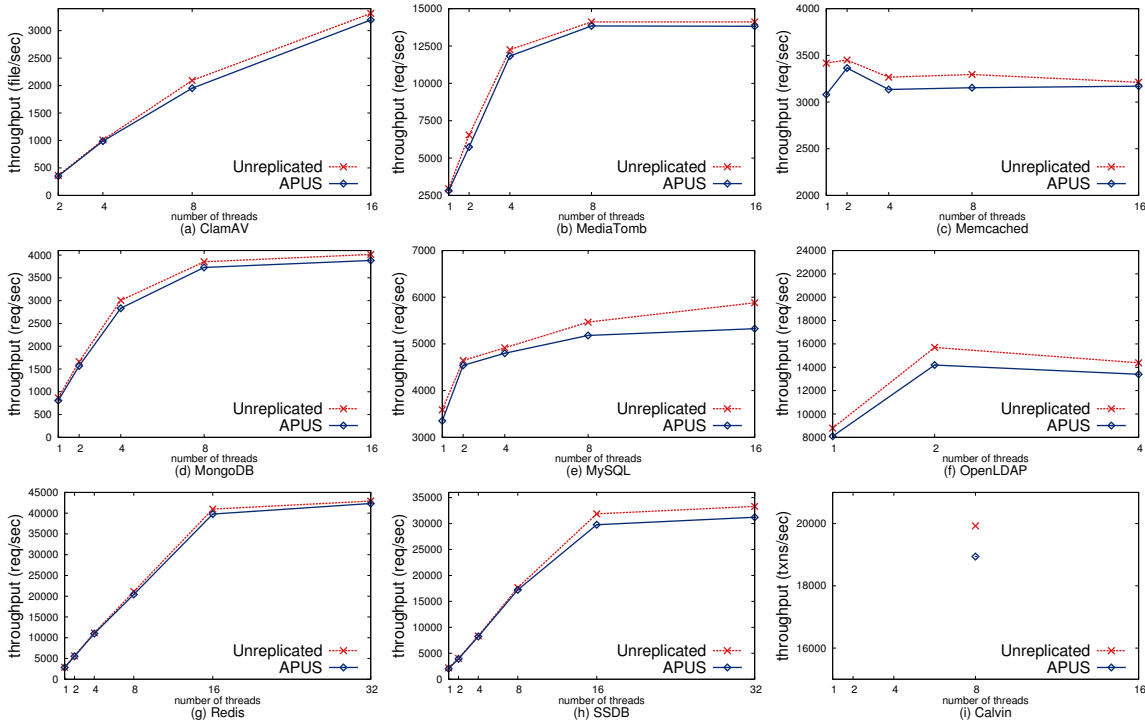


Figure 8: APUS throughput compared to server programs' unreplicated executions.

lower than that on 100% SET workloads. APUS's "wait-consensus" was almost 0 as it enables concurrent consensus requests (§4.1).

DARE evaluation also showed that, with 100% SET workloads, its throughput decreased by 30.1% when the number of replicas increased from three to seven. We reproduced a similar result: we used the same workloads and 24 concurrent connections, and we varied the number of replicas from three to nine. We found that APUS consensus latency increased merely by 7.3% and DARE increased by 67.3% (shown in Figure 1).

Overall, we found DARE better on smaller number of concurrent connections and replicas (e.g., metadata [10, 23]), and APUS better on larger number of connections or replicas (e.g., replicating server programs [29, 40]).

## 7.4 Performance Overhead

To stress APUS, we used nine replicas to run all nine server programs without modifying them. We used up to 32 concurrent client connections (most evaluated programs reached peak throughput at 16), and then we measured mean response time and throughput in 50 runs.

We turned on output checking (§5.3) and didn't observe a performance impact. Only two programs (MySQL and OpenLDAP) have different output hashes caused by physical times (an approach [60] can be leveraged to enforce same physical times across replicas).

Figure 8 shows APUS's throughput. For Calvin, we only collected the 8-thread result because Calvin uses this constant thread count in their code to serve requests. Compared to these server programs' unreplicated executions, APUS merely incurred a mean

throughput overhead of 4.2% (note that in Figure 8, the Y-axes of most programs start from a large number). As the number of threads increases, all programs' unreplicated executions got a performance improvement except Memcached. Prior work [40] also showed that Memcached itself scaled poorly. Overall, APUS scaled as well as unreplicated executions on concurrent requests.

Table 3: Leader's input consensus events per 10K requests, 8 threads. The "# Calls" column means the number of socket calls that went through APUS input consensus; "Input" means average bytes of a server's inputs received in these calls; "First" is the latency of leader's first received consensus reply; and "Quorum" means the average time leader has spent on waiting for quorum replies.

Program	# Calls	Input	First	Quorum
ClamAV	30,000	37.0	10.4 $\mu$ s	10.9 $\mu$ s
MediaTomb	30,000	140.0	16.9 $\mu$ s	17.4 $\mu$ s
Memcached	10,016	38.0	6.5 $\mu$ s	7.0 $\mu$ s
MongoDB	10,376	490.6	8.3 $\mu$ s	9.2 $\mu$ s
MySQL	10,009	28.8	7.1 $\mu$ s	7.8 $\mu$ s
OpenLDAP	10,016	27.3	5.8 $\mu$ s	6.4 $\mu$ s
Redis	10,016	40.5	5.2 $\mu$ s	6.0 $\mu$ s
SSDB	10,016	47.0	5.7 $\mu$ s	6.2 $\mu$ s
Calvin	10,002	128.0	10.1 $\mu$ s	10.8 $\mu$ s

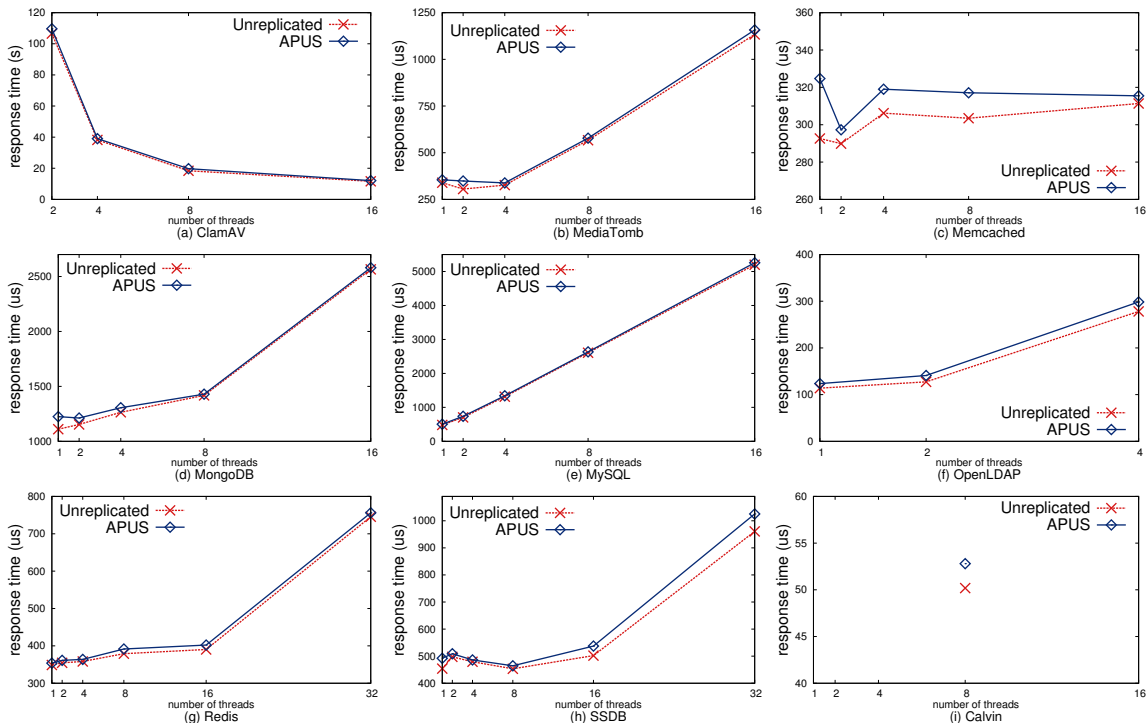


Figure 9: APUS response time compared to server programs' unreplicated executions.

To understand APUS's performance overhead, we broke down its consensus latency on the leader replica. Table 3 shows these statistics per 10K requests, 8 or max (if less than 8) threads. According to the consensus algorithm in Figure 4, for each socket call, APUS's leader does an "L2": SSD write, and an "L4": quorum waiting phase (the "quorum time" column). L4 implies backups' performance because each backup stores consensus requests in local SSD and then WRITES a reply to the leader.

The small consensus latency shown in Table 3 makes APUS achieve a low overhead of 4.3% on response time in Figure 9. Figure 7 and Table 3 also indicate that APUS had a low overhead of on programs' response time.

### 7.5 APUS Performance in a Congested Network

In a production datacenter, network bandwidth is often shared by many services. This raises a performance concern of APUS because RoCE-based network may perform poorly in a congested network due to packet retransmission or congestion control [38, 84].

To evaluate the performance of APUS in a congested network, we generated traffic to consume the network bandwidth between Paxos leader and a randomly chosen backup replica and varied this consumption from 0 to 30Gbps. We ran APUS with SSDB with seven replicas and measured the performance. Figure 10 shows the throughput of APUS and the unreplicated executions with an increasing bandwidth consumption. Even when the traffic consumes 30 Gbps bandwidth, APUS's throughput does not drop significantly.

To understand APUS's performance overhead in the congested network, we collected the consensus latency on the leader replica.

As shown in Figure 11, the consensus latency has increased only  $4.5\mu s$  when the traffic bandwidth consumption increases from 0 to 30Gbps. Because the request processing time of the server program and the TCP network latency between the client and leader replica are much longer than the increased consensus latency, the incurred overhead is negligible.

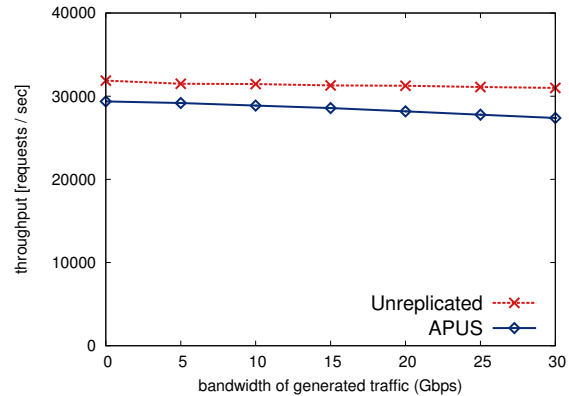
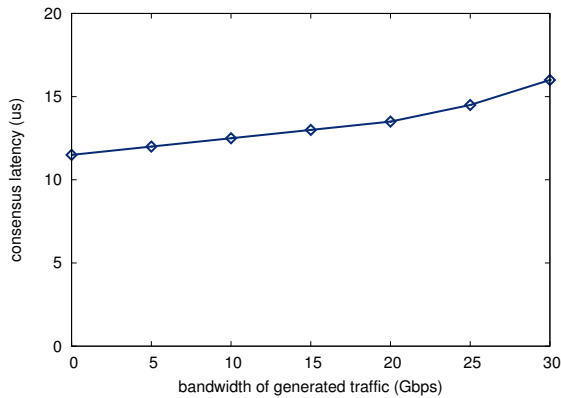


Figure 10: SSDB throughput on APUS with network traffic congestion.

### 7.6 Checkpoint and Recovery

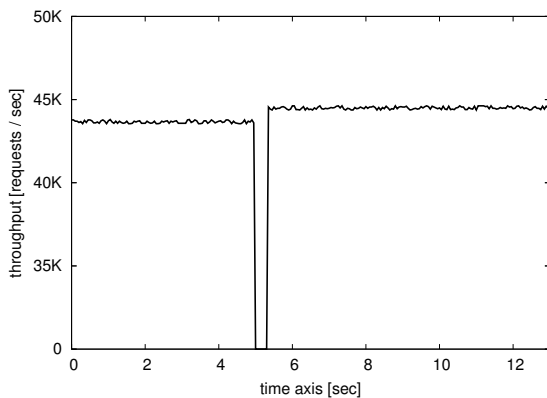
We ran same performance benchmark as in §7.4 and measured programs' checkpoint timecost. Each program checkpoint operation



**Figure 11: APUS consensus latency running SSDB with network traffic congestion.**

(§5.2) costs 0.12s to 11.6s depending on the amount of modified memory and files since a program’s last checkpoint. ClamAV incurred the largest checkpoint time (11.6s) because it loaded and scanned files in the `/lib` directory. Checkpoints did not affect APUS performance in normal case because they were done on only one backup. Leader and other backups still formed majority and reached consensus rapidly.

To evaluate APUS’s Paxos robustness, we ran APUS with Redis with three replicas. We manually killed one backup and then modified another backup’s code to drop all its consensus reply messages. We did not observe a performance change, as other seven replicas still reach consensus. We then manually killed the APUS leader and measured Redis throughput on the leader election approach (§4.4). APUS’s default heartbeat period was 100 ms, and its three-round leader election took only 10.7  $\mu$ s. Redis throughput is shown in Figure 12. After a new leader was elected, Redis throughput went up slightly because there were only two replicas left.



**Figure 12: Redis throughput on APUS leader election.**

## 8 RELATED WORK

**Software-based consensus.** Various Paxos algorithms [53, 54, 60, 66, 79] and implementations [23, 25, 29, 60] exist. Paxos is notoriously difficult to be fast and scalable [37, 48, 63], so server programs carry a weaker asynchronous replication approach (e.g., Redis [76]). Consensus is essential in datacenters [6, 41, 83] and worldwide distributed systems [27, 57], so much work is done to improve Paxos’s input commutativity [58, 66], understandability [53, 68], and verification [39, 82]. Paxos is extended to tolerate byzantine faults [15, 17, 22, 24, 51, 59, 61] and hardware faults [18].

Three SMR systems, Eve [49], Rex [40], and CRANE [29], use traditional Paxos protocols to improve the availability of programs. None of these systems has evaluated their response time overhead on key-value servers, which are extremely sensitive on latency. APUS is the first SMR system that achieves low overhead on both response time and throughput on real-world key-value servers.

**Hardware- or Network- assisted consensus.** Recent systems [33, 43, 44, 55, 74] leverage augmented network hardware or topology to improve Paxos consensus latency. Three systems [33, 43, 44] implement consensus protocols in hardware devices (e.g., switches). “Consensus in a Box” [44] implemented ZooKeeper’s protocol in FPGA. These systems reported similar performance as DARE and they are suitable to maintain compact metadata (e.g., leader election). Prior work [55] pointed out that these systems’ programmable hardware are not suitable to store large amount of replicated states (e.g., server programs’ continuously arriving inputs).

Speculative Paxos [74] and NOPaxos [55] use the datacenter topology to order requests, so they can eliminate consensus rounds if packets are not reordered or lost. These systems require rewriting a server program to use their new libraries for checking the order of packets, so they are not designed to run legacy server programs. Moreover, these two systems’ consensus modules are TCP/UDP-based and incur high consensus latency, which APUS can help.

**RDMA-based systems.** RDMA techniques have been implemented in various architectures, including Infiniband [1], RoCE [8], and iWRAP [9]. RDMA is used to speed up high performance computing [36], key-value stores [34, 45, 46, 64], transactional systems [35, 47, 80], distributed programming languages [19], and file systems [81]. For instance, FaRM [34] runs on RDMA and it provides in a primary-backup replication [31, 69]. Paxos provides better availability than primary-backup. These systems use RDMA to speed up different aspects, so they are complementary to APUS.

## 9 CONCLUSION

We have presented APUS, the first RDMA-based Paxos protocol and its runtime system. Evaluation on five consensus protocols and nine widely used programs shows that APUS is fast, scalable, and deployable. It has the potential to greatly promote the deployments of SMR and improve the reliability of many real-world programs.

## ACKNOWLEDGMENTS

We thank Dan R. K. Ports (our shepherd) and anonymous reviewers for their many helpful comments. This paper is funded in part by a research grant from the Huawei Innovation Research Program (HIRP) 2017, HK RGC ECS (No. 27200916), HK RGC GRF (No. 17207117), and a Croucher innovation award.

## REFERENCES

- [1] 2001. An Introduction to the InfiniBand Architecture. <http://buyya.com/superstorage/chap42.pdf>. (2001).
- [2] 2004. A tool for measuring memcached server performance. <https://github.com/twitter/twemperf>. (2004).
- [3] 2004. clamscan - scan files and directories for viruses. <http://linux.die.net/man/1/clamscan>. (2004).
- [4] 2004. SysBench: a system performance benchmark. <http://sysbench.sourceforge.net>. (2004).
- [5] 2004. Yahoo! Cloud Serving Benchmark. <https://github.com/brianfrankcooper/YCSB>. (2004).
- [6] 2011. Why the data center needs an operating system. [https://cs.stanford.edu/~matei/papers/2011/hotcloud\\_datacenter\\_os.pdf](https://cs.stanford.edu/~matei/papers/2011/hotcloud_datacenter_os.pdf). (2011).
- [7] 2012. Data Plane Development Kit (DPDK). <http://dpdk.org/>. (2012).
- [8] 2012. Mellanox Products: RDMA over Converged Ethernet (RoCE). [http://www.mellanox.com/page/products\\_dyn?product\\_family=79](http://www.mellanox.com/page/products_dyn?product_family=79). (2012).
- [9] 2012. RDMA - iWARP. <http://www.chelsio.com/nic/rdma-iwarp/>. (2012).
- [10] 2012. ZooKeeper. <https://zookeeper.apache.org/>. (2012).
- [11] 2014. ab - Apache HTTP server benchmarking tool. <http://httpd.apache.org/docs/2.2/programs/ab.html>. (2014).
- [12] 2017. MediaTomb - Free UPnP MediaServer. <http://mediatomb.cc/>. (2017).
- [13] 2017. MySQL Database. <http://www.mysql.com/>. (2017).
- [14] Gautam Altekar and Ion Stoica. 2009. ODR: output-deterministic replay for multicore debugging. In *Proceedings of the 22nd ACM Symposium on Operating Systems Principles (SOSP '09)*, 193–206.
- [15] Yair Amir, Claudiu Danilov, Danny Dolev, Jonathan Kirsch, John Lane, Cristina Nita-Rotaru, Josh Olsen, and David Zage. 2010. Steward: Scaling Byzantine fault-tolerant replication to wide area networks. *IEEE Transactions on Dependable and Secure Computing* 7, 1 (2010), 80–93.
- [16] Amitai Aviram, Shu-Chun Weng, Sen Hu, and Bryan Ford. 2010. Efficient System-Enforced Deterministic Parallelism. In *Proceedings of the Ninth Symposium on Operating Systems Design and Implementation (OSDI '10)*.
- [17] Bharath Balasubramanian and Vijay K. Garg. 2014. Fault Tolerance in Distributed Systems Using Fused State Machines. *Distrib. Comput.* (2014).
- [18] Diogo Behrens, Dmitrii Kuvaiskii, and Christof Fetzer. 2014. HardPaxos: Replication Hardened against Hardware Errors. In *Reliable Distributed Systems (SRDS), 2014 IEEE 33rd International Symposium on*.
- [19] Jonathan Behrens, Ken Birman, Sagar Jha, Matthew Milano, Edward Tremel, Eugene Bagdasaryan, Theo Gkountouvas, Weijia Song, and Robbert van Renesse. 2016. Derecho: Group Communication at the Speed of Light. (2016).
- [20] Carlos Eduardo Bezerra, Fernando Pedone, and Robbert Van Renesse. 2014. Scalable State-Machine Replication. In *Proceedings of the 2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN '14)*.
- [21] Martin Biely, Zarko Milosevic, Nuno Santos, and Andre Schiper. 2012. S-Paxos: Offloading the Leader for High Throughput State Machine Replication. In *Proceedings of the 2012 IEEE 31st Symposium on Reliable Distributed Systems (SRDS '12)*.
- [22] Yuriy Brun, George Edwards, Jae Young Bang, and Nenad Medvidovic. 2011. Smart Redundancy for Distributed Computation. In *Proceedings of the 2011 31st International Conference on Distributed Computing Systems (ICDCS '11)*.
- [23] Mike Burrows. 2006. The Chubby lock service for loosely-coupled distributed systems. In *Proceedings of the Seventh Symposium on Operating Systems Design and Implementation (OSDI '06)*, 335–350.
- [24] Miguel Castro and Barbara Liskov. 1999. Practical Byzantine Fault Tolerance. In *Proceedings of the Third Symposium on Operating Systems Design and Implementation (OSDI '99)*.
- [25] Tushar D. Chandra, Robert Griesemer, and Joshua Redstone. 2007. Paxos Made Live: An Engineering Perspective. In *Proceedings of the Twenty-sixth Annual ACM Symposium on Principles of Distributed Computing (PODC '07)*.
- [26] Clam AntiVirus 2017. <http://www.clamav.net/>. (2017).
- [27] James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, J. J. Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, Wilson Hsieh, Sebastian Kanthak, Eugene Kogan, Hongyi Li, Alexander Lloyd, Sergey Melnik, David Mwaura, David Nagle, Sean Quinlan, Rajesh Rao, Lindsay Rolig, Yasushi Saito, Michal Szymaniak, Christopher Taylor, Ruth Wang, and Dale Woodford. 2012. Spanner: Google's Globally-distributed Database. In *Proceedings of the 12th Symposium on Operating Systems Design and Implementation (OSDI '16)*.
- [28] criu 2015. CRIU. <http://criu.org>. (2015).
- [29] Heming Cui, Rui Gu, Cheng Liu, and Junfeng Yang. 2015. Paxos Made Transparent. In *Proceedings of the 25th ACM Symposium on Operating Systems Principles (SOSP '15)*.
- [30] Heming Cui, Jiri Simsa, Yi-Hong Lin, Hao Li, Ben Blum, Xinan Xu, Junfeng Yang, Garth A. Gibson, and Randal E. Bryant. 2013. Parrot: a Practical Runtime for Deterministic, Stable, and Reliable Threads. In *Proceedings of the 24th ACM Symposium on Operating Systems Principles (SOSP '13)*.
- [31] Brendan Cully, Geoffrey Lefebvre, Dutch Meyer, Mike Feeley, Norm Hutchinson, and Andrew Warfield. 2008. Remus: High availability via asynchronous virtual machine replication. In *Proceedings of the 5th USENIX Symposium on Networked Systems Design and Implementation*. San Francisco, 161–174.
- [32] Huynh Tu Dang, Marco Canini, Fernando Pedone, and Robert Soulé. 2016. Paxos made switch-y. *ACM SIGCOMM Computer Communication Review* 46, 1 (2016), 18–24.
- [33] Huynh Tu Dang, Daniele Sciascia, Marco Canini, Fernando Pedone, and Robert Soulé. 2015. NetPaxos: Consensus at Network Speed. In *Proceedings of the 1st ACM SIGCOMM Symposium on Software Defined Networking Research (SOSR '15)*.
- [34] Aleksandar Dragojević, Dushyanth Narayanan, Orion Hodson, and Miguel Castro. 2014. FaRM: Fast Remote Memory. In *Proceedings of the 11th USENIX Conference on Networked Systems Design and Implementation (NSDI '14)*.
- [35] Aleksandar Dragojević, Dushyanth Narayanan, Edmund B. Nightingale, Matthew Renzelmann, Alex Shamis, Anirudh Badam, and Miguel Castro. 2015. No Compromises: Distributed Transactions with Consistency, Availability, and Performance. In *Proceedings of the 25th ACM Symposium on Operating Systems Principles (SOSP '15)*.
- [36] Message Passing Interface Forum. 2009. Open MPI: Open Source High Performance Computing. (Sept. 2009).
- [37] Lisa Glendenning, Ivan Beschastnikh, Arvind Krishnamurthy, and Thomas Anderson. 2011. Scalable Consistency in Scatter. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP '11)*.
- [38] Chuanxiong Guo, Haitao Wu, Zhong Deng, Gaurav Soni, Jianxi Ye, Jitu Padhye, and Marina Lipshteyn. 2016. RDMA over commodity ethernet at scale. In *Proceedings of the 2016 conference on ACM SIGCOMM 2016 Conference*. ACM, 202–215.
- [39] Huayang Guo, Ming Wu, Lidong Zhou, Gang Hu, Junfeng Yang, and Lintao Zhang. 2011. Practical Software Model Checking via Dynamic Interface Reduction. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP '11)*, 265–278.
- [40] Zhenyu Guo, Chuntao Hong, Mao Yang, Dong Zhou, Lidong Zhou, and Li Zhuang. 2014. Rex: Replication at the Speed of Multi-core. In *Proceedings of the 2014 ACM European Conference on Computer Systems (EUROSYS '14)*. ACM, 11.
- [41] Benjamin Hindman, Andy Konwinski, Matei Zaharia, Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, and Ion Stoica. 2011. Mesos: A Platform for Fine-grained Resource Sharing in the Data Center. In *Proceedings of the 8th USENIX conference on Networked Systems Design and Implementation (NSDI '11)*. USENIX Association, Berkeley, CA, USA.
- [42] Patrick Hunt, Mahadev Konar, Flavio P. Junqueira, and Benjamin Reed. 2010. ZooKeeper: Wait-free Coordination for Internet-scale Systems. In *Proceedings of the 2010 USENIX Conference on USENIX Annual Technical Conference (USENIX-ATC '10)*.
- [43] Dang Huynh Tu, Bressana Pietro, Wang Han, Lee Ki Shu, Weatherspoon Hakim, Canini Marco, Pedone Fernando, and Soule Robert. 2016. *Network Hardware-Accelerated Consensus*. Technical Report. USI Technical Report Series in Informatics.
- [44] Zsolt István, David Sidler, Gustavo Alonso, and Marko Vukolic. 2016. Consensus in a Box: Inexpensive Coordination in Hardware. In *Proceedings of the 13th Usenix Conference on Networked Systems Design and Implementation (NSDI '16)*.
- [45] Jithin Jose, Hari Subramoni, Krishna Kandalla, Md. Wasi-ur Rahman, Hao Wang, Sundeep Narravula, and Dhableswar K. Panda. 2012. Scalable Memcached Design for InfiniBand Clusters Using Hybrid Transports. In *Proceedings of the 2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (Ccgri 2012) (CCGRID '12)*.
- [46] Anuj Kalia, Michael Kaminsky, and David G. Andersen. 2014. Using RDMA Efficiently for Key-value Services. (Aug. 2014).
- [47] Anuj Kalia, Michael Kaminsky, and David G. Andersen. 2016. FaSST: Fast, Scalable and Simple Distributed Transactions with Two-Sided (RDMA) Datagram RPCs. In *Proceedings of the 12th Symposium on Operating Systems Design and Implementation (OSDI '16)*.
- [48] Manos Kapritsos and Flavio P. Junqueira. 2010. Scalable Agreement: Toward Ordering As a Service. In *Proceedings of the Sixth International Conference on Hot Topics in System Dependability (HotDep'10)*.
- [49] Manos Kapritsos, Yang Wang, Vivien Queama, Allen Clement, Lorenzo Alvisi, Mike Dahlin, et al. 2012. 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, 237–250.
- [50] Baris Kasikci, Benjamin Schubert, Cristiano Pereira, Gilles Pokam, and George Candea. 2015. Failure Sketching: A Technique for Automated Root Cause Diagnosis of In-production Failures. In *Proceedings of the 25th ACM Symposium on Operating Systems Principles (SOSP '15)*.
- [51] Ramakrishna Kotla, Lorenzo Alvisi, Mike Dahlin, Allen Clement, and Edmund Wong. 2007. Zyzzyva: Speculative Byzantine Fault Tolerance. In *Proceedings of the 21st ACM Symposium on Operating Systems Principles (SOSP '07)*.
- [52] Sriram Krishnan. 2010. *Programming Windows Azure: Programming the Microsoft Cloud*.

- [53] Leslie Lamport. 1998. The part-time parliament. *ACM Trans. Comput. Syst.* 16, 2 (1998), 133–169.
- [54] Leslie Lamport. 2001. Paxos made simple. <http://research.microsoft.com/en-us/um/people/lamport/pubs/paxos-simple.pdf>. (2001).
- [55] Jialin Li, Ellis Michael, Naveen Kr. Sharma, Adriana Szekeres, and Dan R. K. Ports. 2016. Fast Replication with NOPaxos: Replacing Consensus with Network Ordering. In *Proceedings of the 12th Symposium on Operating Systems Design and Implementation (OSDI '16)*.
- [56] Tongping Liu, Charlie Curtsinger, and Emery D. Berger. 2011. DTHREADS: efficient deterministic multithreading. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP '11)*. 327–336.
- [57] Yanhua Mao, Flavio Paiva Junqueira, and Keith Marzullo. 2008. Mencius: building efficient replicated state machines for WANs. In *Proceedings of the 8th USENIX conference on Operating systems design and implementation*, Vol. 8. 369–384.
- [58] Parisa Jalili Marandi, Carlos Eduardo Bezerra, and Fernando Pedone. 2014. Re-thinking State-Machine Replication for Parallelism. In *Proceedings of the 2014 IEEE 34th International Conference on Distributed Computing Systems (ICDCS '14)*.
- [59] Rolando Martins, Rajeev Gandhi, Priya Narasimhan, Soila Pertet, António Casimiro, Diego Kreutz, and Paulo Verissimo. 2013. Experiences with fault-injection in a Byzantine fault-tolerant protocol. In *ACM/IFIP/USENIX International Conference on Distributed Systems Platforms and Open Distributed Processing*. Springer, 41–61.
- [60] David Mazieres. 2007. *Paxos made practical*. Technical Report. Technical report, 2007. <http://www.scs.stanford.edu/dm/home/papers>.
- [61] Hein Meling, Keith Marzullo, and Alessandro Mei. 2012. When You Don't Trust Clients: Byzantine Proposer Fast Paxos. In *Proceedings of the 2012 IEEE 32Nd International Conference on Distributed Computing Systems (ICDCS '12)*.
- [62] Memcached 2017. <https://memcached.org/>. (2017).
- [63] Ellis Michael. 2015. *Scaling Leader-Based Protocols for State Machine Replication*. Ph.D. Dissertation. University of Texas at Austin.
- [64] Christopher Mitchell, Yifeng Geng, and Jinyang Li. 2013. Using One-sided RDMA Reads to Build a Fast, CPU-efficient Key-value Store. In *Proceedings of the USENIX Annual Technical Conference (USENIX '14)*.
- [65] mongodb 2017. MongoDB. <http://www.mongodb.org>. (2017).
- [66] Iulian Moraru, David G. Andersen, and Michael Kaminsky. 2013. There is More Consensus in Egalitarian Parliaments. In *Proceedings of the 13th ACM Symposium on Operating Systems Principles (SOSP '11)*.
- [67] Nginx 2012. Nginx Web Server. <https://nginx.org/>. (2012).
- [68] Diego Ongaro and John Ousterhout. 2014. In Search of an Understandable Consensus Algorithm. In *Proceedings of the USENIX Annual Technical Conference (USENIX '14)*.
- [69] Diego Ongaro, Stephen M. Rumble, Ryan Stutsman, John Ousterhout, and Mendel Rosenblum. 2011. Fast Crash Recovery in RAMCloud. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP '11)*.
- [70] OpenLDAP 2017. OpenLDAP. (2017). <https://www.openldap.org/>
- [71] Soyeon Park, Yuanyuan Zhou, Weiwei Xiong, Zuoning Yin, Rini Kaushik, Kyu H. Lee, and Shan Lu. 2009. PRES: probabilistic replay with execution sketching on multiprocessors. In *Proceedings of the 22nd ACM Symposium on Operating Systems Principles (SOSP '09)*. 177–192.
- [72] Simon Peter, Jialin Li, Irene Zhang, Dan R. K. Ports, Doug Woos, Arvind Krishnamurthy, Thomas Anderson, and Timothy Roscoe. 2014. Arrakis: The Operating System is the Control Plane. In *Proceedings of the Eleventh Symposium on Operating Systems Design and Implementation (OSDI '14)*.
- [73] Marius Poke and Torsten Hoefler. 2015. DARE: High-Performance State Machine Replication on RDMA Networks. In *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing (HPDC '15)*.
- [74] Dan R. K. Ports, Jialin Li, Vincent Liu, Naveen Kr. Sharma, and Arvind Krishnamurthy. 2015. Designing Distributed Systems Using Approximate Synchrony in Data Center Networks. In *Proceedings of the 12th USENIX Conference on Networked Systems Design and Implementation (NSDI'15)*.
- [75] Marco Primi. 2016. LibPaxos. <http://libpaxos.sourceforge.net/>. (2016).
- [76] Redis 2017. <http://redis.io/>. (2017).
- [77] SSDB 2017. [ssdb.io/](http://ssdb.io/). (2017).
- [78] Alexander Thomson, Thaddeus Diamond, Shu-Chun Weng, Kun Ren, Philip Shao, and Daniel J. Abadi. 2014. Fast Distributed Transactions and Strongly Consistent Replication for OLTP Database Systems. (May 2014).
- [79] Robbert Van Renesse and Deniz Altinbuken. 2015. Paxos Made Moderately Complex. *ACM Computing Surveys (CSUR)* 47, 3 (2015), 42:1–42:36.
- [80] Xingda Wei, Jiaxin Shi, Yanzhe Chen, Rong Chen, and Haibo Chen. 2015. Fast In-memory Transaction Processing Using RDMA and HTM. In *Proceedings of the 25th ACM Symposium on Operating Systems Principles (SOSP '15) (SOSP '15)*.
- [81] Garth Gibson Wittawat Tantisiroj. 2008. *Network File System (NFS) in High Performance Networks*. Technical Report CMU-PDLSVD08-02. Carnegie Mellon University.
- [82] Junfeng Yang, Tisheng Chen, Ming Wu, Zhilei Xu, Xuezheng Liu, Haoxiang Lin, Mao Yang, Fan Long, Lintao Zhang, and Lidong Zhou. 2009. MODIST: Transparent Model Checking of Unmodified Distributed Systems. In *Proceedings of the Sixth Symposium on Networked Systems Design and Implementation (NSDI '09)*. 213–228.
- [83] Matei Zaharia, Benjamin Hindman, Andy Konwinski, Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, and Ion Stoica. 2011. The Datacenter Needs an Operating System. In *Proceedings of the 3rd USENIX Conference on Hot Topics in Cloud Computing*.
- [84] Yibo Zhu, Haggai Eran, Daniel Firestone, Chuanxiong Guo, Marina Lipshteyn, Yehonatan Liron, Jitendra Padhye, Shachar Raindel, Mohamad Haj Yahia, and Ming Zhang. 2015. Congestion control for large-scale RDMA deployments. In *ACM SIGCOMM Computer Communication Review*, Vol. 45. ACM, 523–536.