

Prime: Byzantine Replication under Attack

Yair Amir, *Member, IEEE Computer Society*, Brian Coan,
Jonathan Kirsch, *Member, IEEE*, and John Lane, *Member, IEEE*

Abstract—Existing Byzantine-resilient replication protocols satisfy two standard correctness criteria, safety and liveness, even in the presence of Byzantine faults. The runtime performance of these protocols is most commonly assessed in the absence of processor faults and is usually good in that case. However, faulty processors can significantly degrade the performance of some protocols, limiting their practical utility in adversarial environments. This paper demonstrates the extent of performance degradation possible in some existing protocols that do satisfy liveness and that do perform well absent Byzantine faults. We propose a new performance-oriented correctness criterion that requires a consistent level of performance, even with Byzantine faults. We present a new Byzantine fault-tolerant replication protocol that meets the new correctness criterion and evaluate its performance in fault-free executions and when under attack.

Index Terms—Performance under attack, Byzantine fault tolerance, replicated state machines, distributed systems.



1 INTRODUCTION

EXISTING Byzantine fault-tolerant state machine replication protocols are evaluated against two standard correctness criteria: *safety* and *liveness*. Safety means that correct servers do not make inconsistent ordering decisions, while liveness means that each update to the replicated state is eventually executed. Most Byzantine replication protocols are designed to maintain safety in all executions, even when the network delivers messages with arbitrary delay. However, the well-known FLP impossibility result [2] implies that no asynchronous Byzantine agreement protocol can always be both safe and live, and thus these systems ensure liveness only during periods of sufficient synchrony and connectivity [3] or in a probabilistic sense [4], [5].

When the network is sufficiently stable and there are no Byzantine faults, Byzantine fault-tolerant replication systems can satisfy much stronger performance guarantees than liveness. The literature has many examples of systems that have been evaluated in such benign executions and that achieve throughputs of thousands of update operations per second (e.g., [6], [7]). It has been a less common practice to assess the performance of Byzantine fault-tolerant replication systems when some of the processors actually exhibit Byzantine faults. In this paper, we point out that in many systems, a small number of Byzantine processors can degrade performance to a level far below what would be achievable with only correct processors. Specifically, the Byzantine processors can cause the system to make progress at an extremely slow rate, even when the network is stable and could support much higher throughput. While “correct” in the traditional sense (both safety and liveness are met),

systems vulnerable to such performance degradation are of limited practical use in adversarial environments.

We experienced this problem firsthand in 2005, when the US Defense Advanced Research Projects Agency (DARPA) conducted a red team experiment on our Steward system [8]. Steward survived all of the tests according to the metrics of safety and liveness, and most attacks did not impact performance. However, in one experiment, we observed that the system was slowed down to 20 percent of its potential performance. After analyzing the attack, we found that we could slow the system down to roughly one percent of its potential performance. This experience led us to a new way of thinking about Byzantine fault-tolerant replication systems. We concluded that liveness is a necessary but insufficient correctness criterion for achieving high performance when the system actually has Byzantine faults. This paper argues that new, *performance-oriented* correctness criteria are needed to achieve a practical solution for Byzantine replication.

Preventing the type of performance degradation experienced by Steward requires addressing what we call a *Byzantine performance failure*. Previous work on Byzantine fault tolerance has focused on mitigating Byzantine failures in the value domain (where a faulty processor tries to subvert the protocol by sending incorrect or conflicting messages) and the time domain (where messages from a faulty processor do not arrive within protocol time-outs, if at all). Processors exhibiting performance failures operate arbitrarily but correctly enough to avoid being suspected as faulty. They can send valid messages slowly but without triggering protocol time-outs; reorder or drop certain messages, both of which could be caused by a faulty network; or, with malicious intent, take one of a number of possible actions that a correct processor in similar circumstances might legitimately take. Thus, processors exhibiting performance failures are correct in the value and time domains yet have the potential to significantly degrade performance. The problem is magnified in wide-area networks, where time-outs tend to be large and it may be difficult to determine what type of performance should be expected. A performance failure is not a new failure mode but rather is a strategy taken by an adversary that controls Byzantine processors.

• Y. Amir, J. Kirsch, and J. Lane are with the Department of Computer Science, Johns Hopkins University, 3400 North Charles Street, Baltimore, MD 21218. E-mail: {yairamir, jak, johllane}@cs.jhu.edu.

• B. Coan is with Telcordia Technologies, One Telcordia Drive, 1P397 Piscataway, NJ 08854. E-mail: coan@research.telcordia.com.

Manuscript received 3 Feb. 2009; revised 27 Feb. 2010; accepted 2 Aug. 2010; published online 23 Nov 2010.

For information on obtaining reprints of this article, please send e-mail to: tdsc@computer.org, and reference IEEECS Log Number TDSC-2009-02-0018. Digital Object Identifier no. 10.1109/TDSC.2010.70.

In order to better understand the challenges associated with building Byzantine fault-tolerant replication protocols that can resist performance failures, we analyzed existing protocols to assess their vulnerability to performance degradation by malicious servers. We observed that most of the protocols (e.g., [6], [7], [8], [9], [10], [11], [12]) share a common feature: they rely on an elected leader to coordinate the agreement protocol. We call such protocols *leader based*. We found that leader-based protocols are vulnerable to performance degradation caused by a malicious leader.

Based on the understanding gained from our analysis, we developed Prime, the first Byzantine fault-tolerant state machine replication protocol capable of making a meaningful performance guarantee even when some of the servers are Byzantine. Prime meets a new, performance-oriented correctness criterion, called BOUNDED-DELAY. Informally, BOUNDED-DELAY bounds the latency between a correct server receiving a client operation and the correct servers executing the operation. The bound is a function of the network delays between the correct servers in the system. This is a much stronger performance guarantee than the eventual execution promised by liveness.

Like many existing Byzantine fault-tolerant replication protocols, Prime is leader based. Unlike existing protocols, Prime bounds the amount of performance degradation that can be caused by the faulty servers, including by a malicious leader. Two main insights motivate Prime's design. First, most protocol steps do not need any messages from the faulty servers to complete. Faulty servers cannot delay these steps beyond the time it would take if only correct servers were participating in the protocol. Second, the leader should require a predictable amount of resources to fulfill its role as leader. In Prime, the resources required by the leader to do its job as leader are bounded as a function of the number of servers in the system and are independent of the offered load. The result is that the performance of the few protocol steps that do depend on the (potentially malicious) leader can be effectively monitored by the nonleader servers.

We present experimental results evaluating the performance of Prime in fault-free and under attack executions. Our results demonstrate that Prime performs competitively with existing Byzantine fault-tolerant replication protocols in fault-free configurations and that Prime performs an order of magnitude better in under attack executions in the configurations tested.

2 CASE STUDY: BFT UNDER ATTACK

This section presents an attack analysis of Castro and Liskov's BFT protocol [6], a leader-based Byzantine fault-tolerant replication protocol. We chose BFT because 1) it is a widely studied protocol to which other Byzantine-resilient protocols are often compared, 2) many of the attacks that can be applied to BFT (and the corresponding lessons learned) also apply to other leader-based protocols, and 3) its implementation was publicly available. BFT achieves high throughput in fault-free executions or when servers exhibit only benign faults. We first provide background on BFT and then describe two attacks that can be used to significantly degrade its performance when under attack.

BFT assigns a total order to client operations. The protocol requires $3f + 1$ servers, where f is the maximum number of servers that may be Byzantine. An elected leader coordinates the protocol. If a server suspects that the leader has failed, it votes to replace it. When $2f + 1$ servers vote to replace the leader, a view change occurs, in which a new leader is elected and servers collect information regarding pending operations so that progress can safely resume in a new view.

A client sends its operation directly to the leader. The leader proposes a sequence number for the operation by broadcasting a PRE-PREPARE message, which contains the view number, the proposed sequence number, and the operation itself. Upon receiving the PRE-PREPARE, a non-leader server accepts the proposed assignment by broadcasting a PREPARE message. When a server collects the PRE-PREPARE and $2f$ corresponding PREPARE messages, it broadcasts a COMMIT message. A server globally orders the operation when it collects $2f + 1$ COMMIT messages. Each server executes globally ordered operations according to sequence number.

2.1 Attack 1: Pre-Prepare Delay

A malicious leader can introduce latency into the global ordering path simply by waiting some amount of time after receiving a client operation before sending it in a PRE-PREPARE message. The amount of delay a leader can add without being detected as faulty is dependent on 1) the way in which nonleaders place time-outs on operations they have not yet executed and 2) the duration of these time-outs.

A malicious leader can ignore operations sent directly by clients. If a client's time-out expires before receiving a reply to its operation, it broadcasts the operation to all servers, which forward the operation to the leader. Each nonleader server maintains a FIFO queue of pending operations (i.e., those operations it has forwarded to the leader but has not yet executed). A server places a time-out on the execution of the first operation in its queue; that is, it expects to execute the operation within the time-out period. If the time-out expires, the server suspects the leader is faulty and votes to replace it. When a server executes the first operation in its queue, it restarts the timer if the queue is not empty. Note that a server does not stop the timer if it executes a pending operation that is not the first in its queue. The duration of the time-out is dependent on its initial value (which is implementation and configuration dependent) and the history of past view changes. Servers double the value of their time-out each time a view change occurs. BFT does not provide a mechanism for reducing time-out values.

To retain its role as leader, the leader must prevent $f + 1$ correct servers from voting to replace it. Thus, assuming a time-out value of T , a malicious leader can use the following attack to cause delay: 1) Choose a set, S , of $f + 1$ correct servers, 2) For each server $i \in S$, maintain a FIFO queue of the operations forwarded by i , and 3) For each such queue, send a PRE-PREPARE containing the first operation on the queue every $T - \epsilon$ time units. This guarantees that the $f + 1$ correct servers in S execute the first operation on their queue each time-out period. If these operations are all different, the fastest the leader would need to introduce operations is at a rate of $f + 1$ per time-out period. In the

worst case, the $f + 1$ servers would have identical queues, and the leader could introduce one operation per time-out.

This attack exploits the fact that nonleader servers place time-outs only on the first operation in their queues. To understand the ramifications of placing a time-out on *all* pending operations, we consider a hypothetical protocol that is identical to BFT except that nonleader servers place a time-out on all pending operations. Suppose nonleader server i simultaneously forwards n operations to the leader. If server i sets a time-out on all n operations, then i will suspect the leader if the system fails to execute n operations per time-out period. Since the system has a maximal throughput, if n is sufficiently large, i will suspect a correct leader. The fundamental problem is that correct servers have no way to assess the rate at which a correct leader can coordinate the global ordering.

Recent protocols [13], [14] attempt to mitigate the PRE-PREPARE attack by rotating the leader (an idea suggested in [15]). While these protocols allow good long-term throughput and avoid the scenario in which a faulty leader can degrade performance indefinitely, they do not guarantee that individual operations will be ordered in a timely manner. Prime takes a different approach, guaranteeing that the system eventually settles on a leader that is forced to propose an ordering on *all* operations in a timely manner. To meet this requirement, the leader needs only a bounded amount of incoming and outgoing bandwidth, which would not be the case if servers placed a time-out on all operations in BFT.

We note that BFT can be configured to use an optimization in which clients disseminate operations to all servers and the leader sends PRE-PREPARE messages containing hashes of the operations, thus limiting (but not bounding) the bandwidth requirements of the leader. However, when this optimization is used, faulty clients can repeatedly cause performance degradation by disseminating their operations only to $f + 1$ of the $2f + 1$ correct servers. This causes f correct servers to learn the ordering of an operation without having the operation itself, which prevents them from executing subsequent operations until they recover the missing operations.

2.2 Attack 2: Time-Out Manipulation

BFT ensures safety regardless of network synchrony. However, while attacks that impact message delay (e.g., denial-of-service attacks) cannot cause inconsistency, they can be used to increase the time-out value used to detect a faulty leader. During the attack, the time-out doubles with each view change. If the adversary stops the attack when a malicious server is the leader, then that leader will be able to slow the system down to a throughput of roughly $f + 1$ operations per time-out T , where T is potentially very large, using the attack described in the previous section. This vulnerability stems from BFT's inability to reduce the time-out and adapt to the network conditions after the system stabilizes.

One might try to overcome this problem in several ways, such as by resetting the time-out when the system reaches a view in which progress occurs, or by adapting the time-out using a multiplicative increase and additive decrease mechanism. In the former approach, if the time-out is set too low originally, then it will be reset just when it reaches

a large enough value. This may cause the system to experience long periods during which new operations cannot be executed, because leaders (even correct ones) continue to be suspected until the time-out becomes large enough again. The latter approach may be more effective but will be slow to adapt after periods of instability. As will be explained in Section 5.4, Prime adapts to changing network conditions and dynamically determines an acceptable level of timeliness based on the current latencies between correct servers.

3 SYSTEM MODEL AND SERVICE PROPERTIES

We consider a system consisting of N servers and M clients (collectively called *processors*), which communicate by passing messages. Each server is uniquely identified from the set $\mathcal{R} = \{1, 2, \dots, N\}$, and each client is uniquely identified from the set $\mathcal{S} = \{N + 1, N + 2, \dots, N + M\}$. We assume a Byzantine fault model in which processors are either *correct* or *faulty*; correct processors follow the protocol specification exactly, while faulty processors can deviate from the protocol specification arbitrarily by sending any message at any time, subject to the cryptographic assumptions stated below. We assume that $N \geq 3f + 1$, where f is an upper bound on the number of servers that may be faulty. For simplicity, we describe the protocol for the case when $N = 3f + 1$. Any number of clients may be faulty.

We assume an asynchronous network, in which message delay for any message is unbounded. The system meets our safety criteria in all executions in which f or fewer servers are faulty. The system guarantees our liveness and performance properties only in subsets of the executions in which message delay satisfies certain constraints. For some of our analysis, we will be interested in the subset of executions that model Diff-Serv [16] with two traffic classes. To facilitate this modeling, we allow each correct processor to designate each message that it sends as either *TIMELY* or *BOUNDED*.

All messages sent between processors are digitally signed. We denote a message, m , signed by processor i as $\langle m \rangle_{\sigma_i}$. We assume that digital signatures are unforgeable without knowing a processor's private key. We also make use of a collision-resistant hash function, D , for computing message digests. We denote the digest of message m as $D(m)$.

A client submits an *operation* to the system by sending it to one or more servers. Operations are classified as read-only (*queries*) and read/write (*updates*). Each client operation is signed. Each server produces a sequence of operations, $\{o_1, o_2, \dots\}$, as its output. The output reflects the order in which the server executes client operations. When a server outputs an operation, it sends a reply containing the result of the operation to the client.

Safety. Server replies for operations submitted by correct clients are correct according to linearizability [17], as modified to cope with faulty clients [18]. Prime establishes a total order on client operations, as encapsulated in the following safety property:

Definition 3.1 (Safety-S1). *In all executions in which f or fewer servers are faulty, the output sequences of two correct servers are identical, or one is a prefix of the other.*

Liveness and performance properties. Before defining Prime’s liveness and performance properties, which we call PRIME-LIVENESS and BOUNDED-DELAY, we first require several definitions:

Definition 3.2. A stable set is a set of correct servers, $Stable$, such that $|Stable| \geq 2f + 1$. We refer to the members of $Stable$ as the stable servers.

Definition 3.3 (Bounded-Variance(T, S, K)). For each pair of servers, s and r , in S , there exists a value, $Min_Latency(s, r)$, unknown to the servers, such that if s sends a message in traffic class T to r , it will arrive with delay $\Delta_{s,r}$, where $Min_Latency(s, r) \leq \Delta_{s,r} \leq Min_Latency(s, r) * K$.

Definition 3.4 (Eventual-Synchrony(T, S)). Any message in traffic class T sent from server $s \in S$ to server $r \in S$ will arrive within some unknown bounded time.

We now specify the degrees of network stability needed for Prime to meet PRIME-LIVENESS and BOUNDED-DELAY, respectively:

Definition 3.5 (Stability-S1). Let T_{timely} be a traffic class containing all messages designated as TIMELY. Then, there exists a stable set, S , a known network-specific constant, K_{Lat} , and a time, t , after which $Bounded-Variance(T_{timely}, S, K_{Lat})$ holds.

Definition 3.6 (Stability-S2). Let T_{timely} and $T_{bounded}$ be traffic classes containing messages designated as TIMELY and BOUNDED, respectively. Then, there exists a stable set, S , a known network-specific constant, K_{Lat} , and a time, t , after which $Bounded-Variance(T_{timely}, S, K_{Lat})$ and $Eventual-Synchrony(T_{bounded}, S)$ hold.

We now state Prime’s liveness and performance properties:

Definition 3.7 (Prime-Liveness). If Stability-S1 holds for a stable set, S , and no more than f servers are faulty, then if a server in S receives an operation from a correct client, the operation will eventually be executed by all servers in S .

Definition 3.8 (Bounded-Delay). If Stability-S2 holds for a stable set, S , and no more than f servers are faulty, then there exists a time after which the latency between a server in S receiving a client operation and all servers in S executing that operation is upper bounded.

Discussion: The degree of stability needed for liveness in Prime (i.e., *Stability-S1*) is incomparable with the degree of stability needed in BFT and similar protocols. BFT requires *Eventual-Synchrony* [3] to hold for all protocol messages, while Prime requires *Bounded-Variance* (a stronger degree of synchrony) but only for messages in the TIMELY traffic class. As described in Section 5, the TIMELY messages account for a small fraction of the total traffic, are only sent periodically, and have a small, bounded size. Prime is live even when messages in the BOUNDED traffic class arrive completely asynchronously.

It is possible for a strong network adversary capable of controlling the network variance to construct scenarios in which BFT is live and Prime is not. These scenarios occur when the variance for TIMELY messages becomes greater than K_{Lat} , yet the delay is still bounded. This can be made less likely to occur in practice by increasing K_{Lat} , although

at the cost of giving a faulty leader more leeway to cause delay (see Section 5.4).

In practice, we believe both *Stability-S1* and *Stability-S2* can be made to hold. In well-provisioned local-area networks, network delay is often predictable and queuing is unlikely to occur. On wide-area networks, one could use a quality-of-service mechanism such as Diff-Serv [16], with one low-volume class for TIMELY messages and a second class for BOUNDED messages, to give *Bounded-Variance* sufficient coverage, provided enough bandwidth is available to pass the TIMELY messages without queuing. The required level of bandwidth is tunable and independent of the offered load; it is based only on the number of servers in the system and the rate at which the periodic messages are sent. Thus, in a well-engineered system, *Bounded-Variance* should hold for TIMELY messages, regardless of the offered load.

Finally, we remark that resource exhaustion denial-of-service attacks may cause *Stability-S2* to be violated for the duration of the attack. Such attacks fundamentally differ from the attacks that are the focus of this paper, where malicious leaders can slow down the system without triggering defense mechanisms. Handling resource exhaustion attacks at the systemwide level is a difficult problem that is orthogonal and complementary to the solution strategies considered in this work.

4 PRIME: DESIGN AND OVERVIEW

In any execution where the network is well behaved, Prime eventually bounds the time between when a client submits an operation to a stable server and when all of the stable servers execute the operation. The nonleader servers monitor the performance of the leader and replace any leader that is performing too slowly.

In existing leader-based protocols, variation in workload can give rise to legitimate variations in a leader’s performance, which, in turn, make it hard or impossible for nonleaders to judge the leader’s performance. In contrast, for any fixed system size, Prime assigns the leader a fixed amount of work (to do in its role as leader), independent of the load on the system. Thus, it is much easier for the nonleaders to judge the performance of the leader. A leader in Prime is also assigned tasks unrelated to its role as leader. It must give priority to its leader tasks so that its performance as leader will be independent of system workload.

Prime reduces the amount of work assigned to the leader by offloading to other servers most of the tasks required to establish an order on client operations. Most offloaded tasks are done by all of the servers as a group, with the leader participating but playing no special role. In any execution where the network is well behaved, the progress of these tasks has bounded delay regardless of the behavior of Byzantine servers (including a Byzantine leader).

One offloaded task is handled differently. Much of the initial work to order and propagate each client operation is done by a subprotocol coordinated by the server to which the operation was submitted. If the coordinator of this subprotocol is correct, the subprotocol will have bounded delay in executions where the network is well behaved. If instead the coordinator of the subprotocol is Byzantine, there is no requirement on whether or when the subprotocol completes (and, in turn, whether or when the submitted operation is ordered).

The view change protocol in Prime is also designed so that the leader has a fixed-size task and hence can have its performance accurately assessed by its peers.

4.1 Operation Ordering

A client requests the ordering of an operation by submitting it to a server. A server incrementally constructs a server-specific ordering of those operations that clients submit directly to it, and it assumes responsibility for disseminating the operations to the other servers. Each server periodically broadcasts a bounded-size *summary message* that indicates how much of each server's server-specific ordering this server has learned about.

The only thing that the current leader of the main protocol must do to build the global ordering of client operations is to incrementally construct an interleaving of the server-specific orderings. A correct leader periodically sends an *ordering message* containing the most recent summary message from each server. The ordering messages are of fixed size, and the extension to the global order is implicit in the set of summary messages included in the ordering messages sent by the leader. The ordering message describes for each server a (possibly empty) window of additional operations from that server's server-specific order to add to the global order. The specified window always starts with the earliest operation from each server that has not yet been added to the global order. The window adds only those operations known widely enough among the correct servers so that eventually all correct servers will be able to learn what the operations are.

Because the leader's job of extending the global order requires a small, bounded amount of work, the nonleader servers can judge the leader's performance effectively. The nonleaders measure the round-trip times to each other to determine how long it should take between sending a summary to the leader and receiving a corresponding ordering message; we call this the *turnaround time* provided by the leader. Nonleaders also monitor that the leader is sending current summary messages. When a nonleader server sends a summary message to the leader, it can expect the leader's next ordering message to reflect at least as much information about the server-specific orderings as is contained in the summary.

4.2 Overview of Subprotocols

We now list and describe the subprotocols in Prime.

Client subprotocol. The Client subprotocol defines how a client injects an operation into the system and collects replies from servers once the operation has been executed.

Preordering subprotocol. The Preordering subprotocol implements the server-specific orderings that are later interleaved by the leader to construct the global ordering. At all times a separate instance of the subprotocol, coordinated by each server in the system, runs in order to process operations submitted to that server. The subprotocol has three main functions. First, it is used to disseminate to $2f + 1$ servers each client operation. Second, it binds each operation to a unique *preorder identifier*; we say that a server *preorders* an operation when it learns the operation's unique binding. Third, it summarizes each server's knowledge of the server-specific orderings by generating summary messages.

A summary generated by server i contains a value, x , for each server j such that x is the longest gap-free prefix of the server-specific ordering generated by j that is known to i .

Global Ordering subprotocol. The Global Ordering subprotocol runs periodically and is used to incrementally extend the global order. The subprotocol is coordinated by the current leader and, like BFT [6], establishes a total order on PRE-PREPARE messages. Instead of sending a PRE-PREPARE message containing client operations (or even operation identifiers) like in BFT, the leader in Prime sends a PRE-PREPARE message that contains a vector of at most $3f + 1$ summary messages, each from a different server. The summaries contained in the totally ordered sequence of PRE-PREPARE messages induce a total order on the pre-ordered operations.

To ensure that client operations known only to faulty processors will not be globally ordered, we define an operation as *eligible for execution* when the collection of summaries in a PRE-PREPARE message indicate that the operation has been preordered by at least $2f + 1$ servers. An operation that is eligible for execution is known to enough correct servers so that all correct servers will eventually be able to execute it, regardless of the behavior of faulty servers and clients. Totally ordering a PRE-PREPARE extends the global order to include those operations that become eligible for the first time.

Reconciliation subprotocol. The Reconciliation subprotocol proactively recovers globally ordered operations known to some servers but not others. Because correct servers can only execute the gap-free prefix of globally ordered operations, this prevents faulty servers from blocking execution at some correct servers by intentionally failing to disseminate operations to them.

Suspect-Leader subprotocol. The Suspect-Leader subprotocol determines whether the leader is extending the global ordering in a timely manner. The servers measure the round-trip times to each other in order to compute two values. The first is an acceptable turnaround time that the leader should provide, computed as a function of the latencies between the correct servers in the system. The second is a measure of the turnaround time actually being provided by the leader since its election. Suspect-Leader guarantees that a leader will be replaced unless it provides an acceptable turnaround time to at least one correct server, and that the threshold for turnaround time is high enough so that at least $f + 1$ correct servers will not be suspected.

Leader Election subprotocol. When the current leader is suspected to be faulty by enough servers, the nonleader servers vote to elect a new leader. Each leader election is associated with a unique *view number*; the resulting configuration, in which one server is the leader and the rest are nonleaders, is called a *view*. The view number is increased by one and the new leader is chosen by rotation.

View Change subprotocol. When a new leader is elected, the View Change subprotocol wraps up in-progress activity from prior views by deciding which updates to the global order to process and which can be safely abandoned.

TABLE 1
Traffic Class of Each Prime Message Type

Subprotocol	Message Type	Traffic Class
Client	CLIENT-OP	BOUNDED
	CLIENT-REPLY	BOUNDED
Preordering	PO-REQUEST	BOUNDED
	PO-ACK	BOUNDED
	PO-SUMMARY	BOUNDED
Global Ordering	PRE-PREPARE (from leader only)	TIMELY
	PRE-PREPARE (flooded)	BOUNDED
	PREPARE	BOUNDED
	COMMIT	BOUNDED
Reconciliation	RECON	BOUNDED
	INQUIRY	BOUNDED
	CORRUPTION-PROOF	BOUNDED
Suspect-Leader	SUMMARY-MATRIX	TIMELY
	RTT-PING	TIMELY
	RTT-PONG	TIMELY
	RTT-MEASURE	BOUNDED
	TAT-UB	BOUNDED
	TAT-MEASURE	BOUNDED
Leader Election	NEW-LEADER	BOUNDED
	NEW-LEADER-PROOF	BOUNDED
View Change	REPORT	BOUNDED
	PC-SET	BOUNDED
	VC-LIST	BOUNDED
	VC-PARTIAL-SIG	BOUNDED
	REPLAY-PREPARE	BOUNDED
	REPLAY-COMMIT	BOUNDED
	VC-PROOF	TIMELY
	REPLAY	TIMELY

5 PRIME: TECHNICAL DETAILS

This section describes the technical details of Prime. Due to space limitations, we only present the details of the Preordering, Global Ordering, Reconciliation, and Suspect-Leader subprotocols. Details of the remaining subprotocols can be found in Appendix A, on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TDSC.2010.70>. Table 1 lists the types and traffic classes of all of Prime's messages.

5.1 The Preordering Subprotocol

The Preordering subprotocol binds each client operation to a unique *preorder identifier*. The preorder identifier consists of a pair of integers, (i, seq) , where i is the identifier of the server that introduces the operation for preordering, and seq is a *preorder sequence number*, a local variable at i incremented each time it introduces an operation for preordering. Note that the preorder sequence number corresponds to server i 's server-specific ordering.

Operation dissemination and binding. Upon receiving a client operation, o , server i broadcasts a $\langle \text{PO-REQUEST}, seq_i, o, i \rangle_{\sigma_i}$ message. The PO-REQUEST disseminates the client's operation and proposes that it be bound to the preorder identifier (i, seq_i) . When a server, j , receives the PO-REQUEST, it broadcasts a $\langle \text{PO-ACK}, i, seq_i, D(o), j \rangle_{\sigma_j}$ message if it has not previously received a PO-REQUEST from i with preorder sequence number seq_i .

A set consisting of a PO-REQUEST and $2f$ matching PO-ACK messages from different servers constitutes a *preorder certificate*. The preorder certificate proves that the preorder identifier (i, seq_i) is uniquely bound to client operation o . We say that a server that collects a preorder certificate *preorders* the corresponding operation. Prime guarantees that if two servers bind operations o and o' to preorder identifier (i, seq_i) , then $o = o'$.

Let $m_1 = \langle \text{PO-SUMMARY}, vec_1, i \rangle_{\sigma_i}$
Let $m_2 = \langle \text{PO-SUMMARY}, vec_2, i \rangle_{\sigma_i}$

1. m_1 is at least as up-to-date as m_2 when
 - $(\forall j \in \mathcal{R})[vec_1[j] \geq vec_2[j]]$.
2. m_1 is more up-to-date than m_2 when
 - m_1 is at least as up to date as m_2 , and
 - $(\exists j \in \mathcal{R})[vec_1[j] > vec_2[j]]$.
3. m_1 and m_2 are consistent when
 - m_1 is at least as up to date as m_2 , or
 - m_2 is at least as up to date as m_1 .

Fig. 1. Terminology used by the Preordering subprotocol.

Summary generation and exchange. Each correct server maintains a local vector, *PreorderSummary*[], indexed by server identifier. At correct server j , *PreorderSummary*[i] contains the maximum sequence number, n , such that j has preordered all operations bound to preorder identifiers (i, seq) , with $1 \leq seq \leq n$. For example, if server 1 has *PreorderSummary*[] = {2, 1, 3, 0}, then server 1 has preordered the client operations bound to preorder identifiers (1, 1) and (1, 2) from server 1; (2, 1) from server 2; (3, 1), (3, 2), and (3, 3) from server 3; and the server has not yet preordered any operations introduced by server 4.

Each correct server, i , periodically broadcasts the current state of its *PreorderSummary* vector by sending a $\langle \text{PO-SUMMARY}, vec, i \rangle_{\sigma_i}$ message. Note that the PO-SUMMARY message serves as a cumulative acknowledgment for preordered operations and is a short representation of every operation the sender has contiguously preordered (i.e., with no holes) from each server.

A key property of the Preordering subprotocol is that if an operation is introduced for preordering by a correct server, the faulty servers cannot delay the time at which the operation is cumulatively acknowledged (in PO-SUMMARY messages) by at least $2f + 1$ correct servers. This property holds because the rounds are driven by message exchanges between correct servers.

Each correct server stores the most *up-to-date* and *consistent* PO-SUMMARY messages that it has received from each server; these terms are defined formally in Fig. 1. Intuitively, two PO-SUMMARY messages from server i , containing vectors vec and vec' , are consistent if either all of the entries in vec are greater than or equal to the corresponding entries in vec' , or vice versa. Note that correct servers will never send inconsistent PO-SUMMARY messages, since entries in the *PreorderSummary* vector never decrease. Therefore, a pair of inconsistent PO-SUMMARY messages from the same server constitutes proof that the server is malicious. Each correct server, i , maintains a *Blacklist* data structure that stores the server identifiers of any servers from which i has collected inconsistent PO-SUMMARY messages.

The collected PO-SUMMARY messages are stored in a local vector, *LastPreorderSummaries*[], indexed by server identifier. In Section 5.2, we show how the leader uses the contents of its own *LastPreorderSummaries* vector to propose an ordering on preordered operations. In Section 5.4, we show how the nonleaders' *LastPreorderSummaries* vectors determine what they expect to see in the leader's ordering messages, thus allowing them to monitor the leader's performance.

5.2 The Global Ordering Subprotocol

Like BFT, the Global Ordering subprotocol uses three message rounds: PRE-PREPARE, PREPARE, and COMMIT. In Prime, the PRE-PREPARE messages contain and propose a global order on *summary matrices*, not client operations. Each summary matrix is a vector of $3f + 1$ PO-SUMMARY messages. The term “matrix” is used because each PO-SUMMARY message sent by a correct server itself contains a vector, with each entry reflecting the operations that have been preordered from each server. Thus, row i in summary matrix sm (denoted $sm[i]$) either contains a PO-SUMMARY message generated and signed by server i , or a special *empty* PO-SUMMARY message, containing a null vector of length $3f + 1$, indicating that the server has not yet collected a PO-SUMMARY from server i . When indexing into a summary matrix, we let $sm[i][j]$ serve as shorthand for $sm[i].vec[j]$.

A correct leader, l , of view v periodically broadcasts a $\langle \text{PRE-PREPARE}, v, seq, sm, l \rangle_{\sigma_l}$ message, where seq is a global sequence number (analogous to the one assigned to PRE-PREPARE messages in BFT) and sm is the leader’s *LastPreorderSummaries* vector. When correct server i receives a $\langle \text{PRE-PREPARE}, v, seq, sm, l \rangle_{\sigma_l}$ message, it takes the following steps. First, server i checks each PO-SUMMARY message in the summary matrix to see if it is consistent with what i has in its *LastPreorderSummaries* vector. Any server whose PO-SUMMARY is inconsistent is added to i ’s *Blacklist*. Second, i decides if it will respond to the PRE-PREPARE message using similar logic to the corresponding round in BFT. Specifically, i responds to the message if 1) v is the current view number and 2) i has not already processed a PRE-PREPARE in view v with the same sequence number but different content.

If i decides to respond to the PRE-PREPARE, it broadcasts a $\langle \text{PREPARE}, v, seq, D(sm), i \rangle_{\sigma_i}$ message, where v and seq correspond to the fields in the PRE-PREPARE and $D(sm)$ is a digest of the summary matrix found in the PRE-PREPARE. A set consisting of a PRE-PREPARE and $2f$ matching PREPARE messages constitutes a *prepare certificate*. Upon collecting a prepare certificate, server i broadcasts a $\langle \text{COMMIT}, v, seq, D(sm), i \rangle$ message. We say that a server *globally orders* a PRE-PREPARE when it collects $2f + 1$ COMMIT messages that match the PRE-PREPARE.

Obtaining a global order on client operations. At any time at any correct server, the current outcome of the Global Ordering subprotocol is a totally ordered stream of PRE-PREPARE messages: $T = \langle T_1, T_2, \dots, T_x \rangle$. The stream at one correct server may be a prefix of the stream at another correct server, but correct servers do not have inconsistent streams.

We now explain how a correct server obtains a total order on client operations from its current local value of T . Let mat be a function that takes a PRE-PREPARE message and returns the summary matrix that it contains. Let M , a function from PRE-PREPARE messages to sets of preorder identifiers, be defined as

$$M(T_y) = \{(i, seq) : i \in \mathcal{R} \wedge seq \in \mathcal{N} \wedge |\{j : j \in \mathcal{R} \wedge mat(T_y)[j][i] \geq seq\}| \geq 2f + 1\}.$$

Observe that any preorder identifier in $M(T_y)$ has been associated with a specific operation by at least $2f + 1$ servers, of which at least $f + 1$ are correct. The Preordering subprotocol guarantees that this association is unique. The

Reconciliation subprotocol guarantees that any operation in $M(T_y)$ is known to enough correct servers so that in any sufficiently stable execution, any correct server that does not yet have the operation will eventually receive it (see Section 5.3). Note also that since PO-SUMMARY messages are cumulative acknowledgments, if $M(T_y)$ contains a preorder identifier (i, seq) , then $M(T_y)$ also contains all preorder identifiers of the form (i, seq') for $1 \leq seq' < seq$.

Let L be a function that takes as input a set of preorder identifiers, P , and outputs the elements of P ordered lexicographically by their preorder identifiers, with the first element of the preorder identifier having the higher significance. Letting \parallel denote concatenation and \setminus denote set difference, the final total order on clients operations is obtained by

$$\begin{aligned} C_1 &= L(M(T_1)), \\ C_q &= L(M(T_q) \setminus M(T_{q-1})), \\ C &= C_1 \parallel C_2 \parallel \dots \parallel C_x. \end{aligned}$$

Intuitively, when a PRE-PREPARE is globally ordered, it expands the set of preordered operations that are *eligible for execution* to include all operations o for which the summary matrix in the PRE-PREPARE proves that at least $2f + 1$ servers have preordered o . Thus, the set difference operation in the definition of the C_q components causes only those operations that have not already become eligible for execution to be executed.

Pre-Prepare flooding. We now make a key observation about the Global Ordering subprotocol: If all correct servers receive a copy of a PRE-PREPARE message, then there is nothing the faulty servers can do to prevent the PRE-PREPARE from being globally ordered in a timely manner. Progress in the PREPARE and COMMIT rounds is based on collecting sets of $2f + 1$ messages. Therefore, since there are at least $2f + 1$ correct servers, the correct servers are not dependent on messages from the faulty servers to complete the global ordering.

We leverage this property by having a correct server broadcast a PRE-PREPARE upon receiving it for the first time. This guarantees that all correct servers receive the PRE-PREPARE within one round from the time that the first correct server receives it, after which no faulty server can delay the correct servers from globally ordering it. The benefit of this approach is that it forces a malicious leader to delay sending a PRE-PREPARE to *all* correct servers in order to add unbounded delay to the Global Ordering subprotocol. As described in Section 5.4, the Suspect-Leader subprotocol results in the replacement of any leader that fails to send a timely PRE-PREPARE to at least one correct server. This property, combined with PRE-PREPARE flooding, will be used to ensure timely ordering.

In the course of PRE-PREPARE flooding, if a correct server, i , receives two PRE-PREPARE messages with the same view number and sequence number but different summary matrices, server i adds the leader to its *Blacklist* and invokes the Leader Election subprotocol (see Appendix A, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TDSC.2010.70>). By periodically broadcasting the conflicting PRE-PREPARE messages, i will cause all correct servers to blacklist and suspect the leader, ensuring that a new leader is eventually elected.

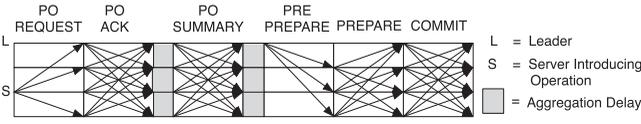


Fig. 2. Fault-free operation of Prime ($f = 1$).

Summary of normal-case operation. To summarize the Preordering and Global Ordering subprotocols, Fig. 2 follows the path of a client operation through the system during normal-case operation. The operation is first preordered in two rounds (PO-REQUEST and PO-ACK), after which its preordering is cumulatively acknowledged (PO-SUMMARY). When the leader is correct, it includes, in its next PRE-PREPARE, the set of at least $2f + 1$ PO-SUMMARY messages that prove that at least $2f + 1$ servers have preordered the operation. The PRE-PREPARE flooding step (not shown) runs in parallel with the PREPARE step. The client operation will be executed once the PRE-PREPARE is globally ordered. Note that in general, many operations are being preordered in parallel, and globally ordering a PRE-PREPARE will make many operations eligible for execution.

5.3 The Reconciliation Subprotocol

The Reconciliation subprotocol ensures that all correct servers will eventually receive any operation that becomes eligible for execution. Conceptually, the protocol operates on the totally ordered sequence of operations defined by the total order $C = C_1 \| C_2 \| \dots \| C_x$. Recall that each C_j is a sequence of preordered operations that became eligible for execution with the global ordering of pp_j , the PRE-PREPARE globally ordered with global sequence number j . From the way C_j is created, for each preordered operation (i, seq) in C_j , there exists a set, $R_{i,seq}$, of at least $2f + 1$ servers whose PO-SUMMARY messages cumulatively acknowledged (i, seq) in pp_j . The protocol operates by having $2f + 1$ deterministically chosen servers in $R_{i,seq}$ send *erasure encoded parts* of the PO-REQUEST containing (i, seq) to those servers that have not cumulatively acknowledged preordering it.

Prime uses a Maximum Distance Separable erasure-resilient coding scheme [19], in which the PO-REQUEST is encoded into $2f + 1$ parts, each $1/(f + 1)$ the size of the original message, such that any $f + 1$ parts are sufficient to decode. Each of the $2f + 1$ servers in $R_{i,seq}$ sends one part. Since at most f servers are faulty, this guarantees that a correct server will receive enough parts to be able to decode the PO-REQUEST. The servers run the reconciliation procedure speculatively, when they first receive a PRE-PREPARE message, rather than when they globally order it. This proactive approach allows operations to be recovered in parallel with the remainder of the Global Ordering subprotocol.

Analysis. Since a correct server will not send a reconciliation message unless at least $2f + 1$ servers have cumulatively acknowledged the corresponding PO-REQUEST, reconciliation messages for a given operation are sent to a maximum of f servers. Assuming an operation size of s_{op} , the $2f + 1$ erasure encoded parts have a total size of $(2f + 1)s_{op}/(f + 1)$. Since these parts are sent to at most f servers, the amount of reconciliation data sent per

operation across all links is at most $f(2f + 1)s_{op}/(f + 1) < (2f + 1)s_{op}$. During the Preordering subprotocol, an operation is sent to between $2f$ and $3f$ servers, which requires at least $2fs_{op}$. Therefore, reconciliation uses approximately the same amount of aggregate bandwidth as operation dissemination. Note that a single server needs to send at most one reconciliation part per operation, which guarantees that at least $f + 1$ correct servers share the cost of reconciliation.

5.4 The Suspect-Leader Subprotocol

There are two types of performance attacks that can be mounted by a malicious leader. First, it can send PRE-PREPARE messages at a rate slower than the one specified by the protocol. Second, even if the leader sends PRE-PREPARE messages at the correct rate, it can intentionally include a summary matrix that does not contain the most up-to-date PO-SUMMARY messages that it has received. This can prevent or delay preordered operations from becoming eligible for execution.

The Suspect-Leader subprotocol is designed to defend against these attacks. The protocol consists of three mechanisms that work together to enforce timely behavior from the leader. The first provides a means by which nonleader servers can tell the leader which PO-SUMMARY messages they expect the leader to include in a subsequent PRE-PREPARE message. The second mechanism allows the nonleader servers to periodically measure how long it takes for the leader to send a PRE-PREPARE containing PO-SUMMARY messages at least as up-to-date as those being reported. We call this time the *turnaround time* provided by the leader. The third mechanism is a *distributed monitoring protocol* in which the nonleader servers can dynamically determine, based on the current network conditions, how quickly the leader should be sending up-to-date PRE-PREPARE messages and decide, based on each server's measurements of the leader's performance, whether to suspect the leader. We now describe each mechanism in more detail.

5.4.1 Reporting the Latest PO-SUMMARY Messages

If the leader is to be expected to send PRE-PREPARE messages with the most up-to-date PO-SUMMARY messages, then each correct server must tell the leader which PO-SUMMARY messages it believes are the most up-to-date. This explicit notification is necessary because the reception of a particular PO-SUMMARY by a correct server does not imply that the leader will receive the same message—the server that originally sent the message may be faulty. Therefore, each correct server, i , periodically sends the leader the complete contents of its *LastPreorderSummaries* vector in a $\langle \text{SUMMARY-MATRIX}, sm, i \rangle_{\sigma_i}$ message.

Upon receiving a SUMMARY-MATRIX, a correct leader updates its *LastPreorderSummaries* vector by adopting any of the PO-SUMMARY messages in the SUMMARY-MATRIX that are more up-to-date than what the leader currently has in its data structure. Since SUMMARY-MATRIX messages have a bounded size-dependent only on the number of servers in the system, the leader requires a small, bounded amount of incoming bandwidth and processing resources to learn about the most up-to-date PO-SUMMARY messages in the system. Furthermore, since PRE-PREPARE messages also

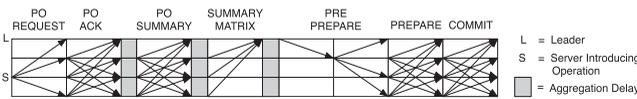


Fig. 3. Operation of Prime with a malicious leader that performs well enough to avoid being replaced ($f = 1$).

have a bounded size independent of the offered load, the leader requires a bounded amount of outgoing bandwidth to send timely, up-to-date PRE-PREPARE messages.

5.4.2 Measuring the Turnaround Time

The preceding discussion suggests a way for nonleader servers to monitor the leader's performance effectively. Given that a correct leader can send timely, up-to-date PRE-PREPARE messages, a nonleader server can measure the time between sending a SUMMARY-MATRIX message, SM , to the leader and receiving a PRE-PREPARE that contains PO-SUMMARY messages that are at least as up-to-date as those in SM . This is the turnaround time provided by the leader. As described below, Suspect-Leader's distributed monitoring protocol forces any server that retains its role as leader to provide a timely turnaround time to at least one correct server. Combined with PRE-PREPARE flooding (see Section 5.2), this ensures that all eligible client operations will be globally ordered in a timely manner.

Fig. 3 depicts the maximum amount of delay that can be added by a malicious leader that performs well enough to avoid being replaced. The leader ignores PO-SUMMARY messages and sends its PRE-PREPARE to only one correct server. PRE-PREPARE flooding ensures that all correct servers receive the PRE-PREPARE within one round of the first correct server receiving it.

In order to define the notion of turnaround time more formally, we first define the *covers* predicate:

Let $pp = \langle \text{PRE-PREPARE}, *, *, sm, * \rangle_{\sigma_*}$
 Let $SM = \langle \text{SUMMARY-MATRIX}, sm', * \rangle_{\sigma_*}$

Then $\text{covers}(pp, SM, i)$ is true at server i iff:

- $\forall j \in (\mathcal{R} \setminus \text{Blacklist}_i), sm[j]$ is at least as up-to-date as $sm'[j]$.

Thus, server i is satisfied that a PRE-PREPARE *covers* a SUMMARY-MATRIX, SM , if, for all servers not in i 's blacklist, each PO-SUMMARY in the PRE-PREPARE is at least as up-to-date (see Fig. 1) as the corresponding PO-SUMMARY in SM .

We can now define the turnaround time provided by the leader for a SUMMARY-MATRIX message, SM , sent by server i , as the time between i sending SM to the leader and i receiving a PRE-PREPARE that 1) covers SM , and 2) is for the next global sequence number for which i expects to receive a PRE-PREPARE. Note that the second condition establishes a connection between receiving an up-to-date PRE-PREPARE and actually being able to execute client operations once the PRE-PREPARE is globally ordered. Without this condition, a leader could provide fast turnaround times without this translating into fast global ordering.

5.4.3 The Distributed Monitoring Protocol

Before describing Suspect-Leader's distributed monitoring protocol, we first define what it means for a turnaround time to be timely. Timeliness is defined in terms of the current network conditions and the rate at which a correct

leader would send PRE-PREPARE messages. In the definition that follows, we let L_{timely}^* denote the maximum latency for a TIMELY message sent between any two correct servers; Δ_{pp} denote a value greater than the maximum time between a correct server sending successive PRE-PREPARE messages; and K_{Lat} be a known network-specific constant accounting for latency variability.

Property 5.1. *If Stability-S1 holds, then any server that retains a role as leader must provide a turnaround time to at least one correct server that is no more than $B = 2K_{Lat}L_{timely}^* + \Delta_{pp}$.*

Property 5.1 ensures that a faulty leader will be suspected unless it provides a timely turnaround time to at least one correct server. We consider a turnaround time, $t \leq B$, to be timely because B is within a constant factor of the turnaround time that the slowest correct server might provide. The factor is a function of the latency variability that Suspect-Leader is configured to tolerate. Note that malicious servers cannot affect the value of B , and that increasing the value of K_{Lat} gives the leader more power to cause delay.

Of course, it is important to make sure that Suspect-Leader is not overly aggressive in the timeliness it requires from the leader. The following property ensures that this is the case:

Property 5.2. *If Stability-S1 holds, then there exists a set of at least $f + 1$ correct servers that will not be suspected by any correct server if elected leader.*

Property 5.2 ensures that when the network is sufficiently stable, view changes cannot occur indefinitely. Prime does not guarantee that the slowest f correct servers will not be suspected because slow faulty leaders cannot be distinguished from slow correct leaders.

We now present Suspect-Leader's distributed monitoring protocol, which allows nonleader servers to dynamically determine how fast a turnaround time the leader should provide and to suspect the leader if it is not providing a fast enough turnaround time to at least one correct server. Pseudocode is contained in Fig. 4.

The protocol is organized as several tasks that run in parallel, with the outcome being that each server decides whether or not to suspect the current leader. This decision is encapsulated in the comparison of two values: TAT_{leader} and $TAT_{acceptable}$ (see Fig. 4, lines E1-E3). TAT_{leader} is a measure of the leader's performance in the current view and is computed as a function of the turnaround times measured by the nonleader servers. $TAT_{acceptable}$ is a standard against which the server judges the current leader and is computed as a function of the round-trip times between correct servers. A server decides to suspect the leader if $TAT_{leader} > TAT_{acceptable}$.

As seen in Fig. 4, lines A1-A4, the data structures used in the distributed monitoring protocol are reinitialized at the beginning of each new view. Thus, a newly elected leader is judged using fresh measurements, both of what turnaround time it is providing and what turnaround time is acceptable given the current network conditions. The following two sections describe how TAT_{leader} and $TAT_{acceptable}$ are computed.

Computing TAT_{leader} . Each server keeps track of the maximum turnaround time provided by the leader in the

```

/* Initialization, run at start of new view */
A1. For  $i = 1$  to  $N$ ,  $TATs\_If\_Leader[i] \leftarrow \infty$ 
A2. For  $i = 1$  to  $N$ ,  $TAT\_Leader\_UBs[i] \leftarrow \infty$ 
A3. For  $i = 1$  to  $N$ ,  $Reported\_TATs[i] \leftarrow 0$ 
A4.  $ping\_seq \leftarrow 0$ 

/* TAT Measurement Task, run at server  $i$  */
B1. Periodically:
B2.  $max\_tat \leftarrow$  Maximum TAT measured this view
B3. BROADCAST:  $\langle TAT-MEASURE, view, max\_tat, i \rangle_{\sigma_i}$ 
B4. Upon receiving  $\langle TAT-MEASURE, view, tat, j \rangle_{\sigma_j}$ 
B5. if  $tat > Reported\_TATs[j]$ 
B6.    $Reported\_TATs[j] \leftarrow tat$ 
B7.  $TAT_{leader} \leftarrow (f+1)^{st}$  lowest val in  $Reported\_TATs$ 

/* RTT Measurement Task, run at server  $i$  */
C1. Periodically:
C2. BROADCAST:  $\langle RTT-PING, view, ping\_seq, i \rangle_{\sigma_i}$ 
C3.  $ping\_seq++$ 
C4. Upon receiving  $\langle RTT-PING, view, seq, j \rangle_{\sigma_j}$ :
C5. SEND to server  $j$ :  $\langle RTT-PONG, view, seq, i \rangle_{\sigma_i}$ 
C6. Upon receiving  $\langle RTT-PONG, view, seq, j \rangle_{\sigma_j}$ :
C7.  $rtt \leftarrow$  Measured RTT for pong message
C8. SEND to server  $j$ :  $\langle RTT-MEASURE, view, rtt, i \rangle_{\sigma_i}$ 
C9. Upon receiving  $\langle RTT-MEASURE, view, rtt, j \rangle_{\sigma_j}$ :
C10.  $t \leftarrow rtt * K_{Lat} + \Delta_{pp}$ 
C11. if  $t < TATs\_If\_Leader[j]$ 
C12.    $TATs\_If\_Leader[j] \leftarrow t$ 

/* TAT Leader Upper Bound Task, run at server  $i$  */
D1. Periodically:
D2.  $\alpha \leftarrow (f+1)^{st}$  highest val in  $TATs\_If\_Leader$ 
D3. BROADCAST:  $\langle TAT-UB, view, \alpha, i \rangle_{\sigma_i}$ 
D4. Upon receiving  $\langle TAT-UB, view, tat\_ub, j \rangle_{\sigma_j}$ :
D5. if  $tat\_ub < TAT\_Leader\_UBs[j]$ 
D6.    $TAT\_Leader\_UBs[j] \leftarrow tat\_ub$ 
D7.  $TAT_{acceptable} \leftarrow (f+1)^{st}$  highest val in  $TAT\_Leader\_UBs$ 

/* Suspect Leader Task */
E1. Periodically:
E2. if  $TAT_{leader} > TAT_{acceptable}$ 
E3.   Suspect Leader

```

Fig. 4. The Suspect-Leader distributed monitoring protocol.

current view and periodically broadcasts the value in a TAT-MEASURE message (Fig. 4, lines B1-B3). The values reported by other servers are stored in a vector, $Reported_TATs$, indexed by server identifier. TAT_{leader} is computed as the $(f+1)$ st lowest value in $Reported_TATs$ (line B7). Since at most f servers are faulty, TAT_{leader} is therefore a value v such that the leader is providing a turnaround time $t \leq v$ to at least one correct server.

Computing $TAT_{acceptable}$. Each server periodically runs a ping protocol to measure the RTT to every other server (Fig. 4, lines C1-C5). Upon computing the RTT to server j , server i sends the RTT measurement to j in an RTT-MEASURE message (line C8). When j receives the RTT measurement, it can compute the maximum turnaround time, t , that i would compute if j were the leader (line C10). Note that t is a function of the latency variability constant, K_{Lat} , as well as the rate at which a correct leader would send PRE-PREPARE messages. Server j stores the minimum such t in $TATs_If_Leader[i]$.

Each server, i , can use the values stored in $TATs_If_Leader$ to compute an upper bound, α , on the value of TAT_{leader} that any correct server will compute for i if it were leader. This upper bound is computed as the $(f+1)$ st highest value in $TATs_If_Leader$ (line D2). The servers periodically exchange their α values by broadcasting TAT-UB messages, storing the values in TAT_Leader_UBs (lines D5-D6). $TAT_{acceptable}$ is computed as the $(f+1)$ st highest value in TAT_Leader_UBs .

We prove that Suspect-Leader meets Properties 5.1 and 5.2 in Appendix B, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TDSC.2010.70>.

6 ANALYSIS

In this section, we show that in those executions in which *Stability-S2* holds, Prime provides BOUNDED-DELAY (see Definition 3.8). As before, we let L_{timely}^* and $L_{bounded}^*$ denote the maximum message delay between correct servers for TIMELY and BOUNDED messages, respectively, and we let $B = 2K_{Lat}L_{timely}^* + \Delta_{pp}$. We also let Δ_{agg} denote a value greater than the maximum time between a correct server sending any of the following messages successively: PO-SUMMARY, SUMMARY-MATRIX, and PRE-PREPARE.

We first consider the maximum amount of delay that can be added by a malicious leader that performs well enough to avoid being replaced. The time between a server receiving and introducing a client operation, o , for preordering and all correct servers sending SUMMARY-MATRIX messages containing at least $2f+1$ PO-SUMMARY messages that cumulatively acknowledge the preordering of o is at most three bounded rounds plus $2\Delta_{agg}$. The malicious servers cannot increase this time beyond what it would take if only correct servers were participating. By Property 5.1, a leader that retains its role as leader must provide a TAT, $t \leq B$, to at least one correct server. By definition, $\Delta_{agg} \geq \Delta_{pp}$. Thus, $B \leq 2K_{Lat}L_{timely}^* + \Delta_{agg}$. Since correct servers flood PRE-PREPARE messages, all correct servers receive the PRE-PREPARE within three bounded rounds and one aggregation delay of when the SUMMARY-MATRIX messages are sent. All correct servers globally order the PRE-PREPARE in two bounded rounds from the time, t , the last correct server receives it. Reconciliation guarantees that all correct servers receive the PO-REQUEST containing the operation within one bounded round of time t . Summing the total delays yields a maximum latency of $\beta = 6L_{bounded}^* + 2K_{Lat}L_{timely}^* + 3\Delta_{agg}$.

If a malicious leader delays proposing an ordering, by more than B , on a summary matrix that proves that at least $2f+1$ servers preordered operation o , it will be suspected and a view change will occur. View changes require a finite (and, in practice, small) amount of state to be exchanged among correct servers, and thus they complete in finite time. As described in Section A.3 in Appendix A, a faulty leader will be suspected if it does not terminate the view change in a timely manner. Property 5.2 of Suspect-Leader guarantees that at most $2f$ view changes can occur before the system settles on a leader that will not be replaced. Therefore, there is a time after which the bound of β holds for any client operation received and introduced by a stable server.

7 PERFORMANCE EVALUATION

To evaluate the performance of Prime, we implemented the protocol and compared its performance to that of an available implementation of BFT. We show results evaluating the protocols in an emulated wide-area setting with seven servers ($f=2$). More extensive performance results, including results in a local-area setting, are presented in Appendix C, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TDSC.2010.70>.

7.1 Testbed and Network Setup

We used a system consisting of seven servers, each running on a 3.2 GHz, 64-bit Intel Xeon computer. RSA signatures [20]

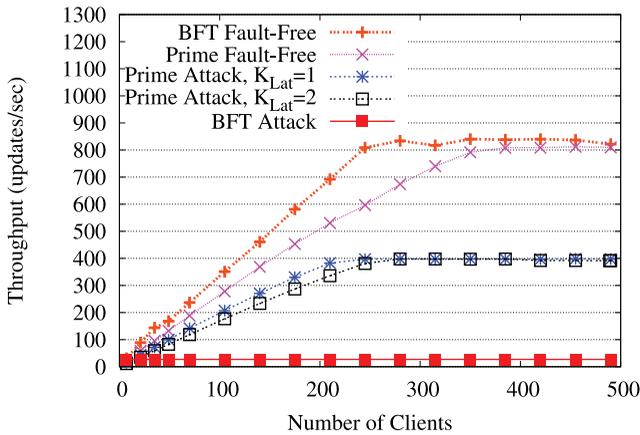


Fig. 5. Throughput of Prime and BFT as a function of the number of clients in a seven-server configuration.

provided authentication and nonrepudiation. We used the netem utility [21] to place delay and bandwidth constraints on the links between the servers. We added 50 ms delay (emulating a US-wide deployment) to each link and limited the aggregate outgoing bandwidth of each server to 10 Mbps. Clients were evenly distributed among the servers, and no delay or bandwidth constraints were set between the client and its server.

Clients submit one update operation to their local server, wait for proof that the update has been ordered, and then submit their next update. Updates contained 512 bytes of data. BFT uses an optimization where clients send updates directly to all of the servers and the BFT PRE-PREPARE message contains batches of update digests. Messages in BFT use message authentication codes for authentication.

7.2 Attack Strategies

Our experimental results during attack show the minimum performance that must be achieved in order for a malicious leader to avoid being replaced. Our measurements do not reflect the time required for view changes, during which a new leader is installed. Since a view change takes a finite and, in practice, relatively small amount of time, malicious leaders must cause performance degradation without being detected in order to have a prolonged effect on throughput. Therefore, we focus on the attack scenario where a malicious leader retains its role as leader indefinitely while degrading performance.

To attack Prime, the leader adds as much delay as possible (without being suspected) to the protocol, and faulty servers force as much reconciliation as possible. As described in Section 5.4, a malicious leader can add approximately two rounds of delay to the Global Ordering subprotocol, plus an aggregation delay. The malicious servers force reconciliation by not sending their PO-REQUEST messages to f of the correct servers. Therefore, all PO-REQUEST messages originating from the faulty servers must be sent to these f correct servers using the Reconciliation subprotocol (see Section 5.3). Moreover, the malicious servers only acknowledge each other's PO-REQUEST messages, forcing the correct servers to send reconciliation messages to them for all PO-REQUEST messages introduced by correct servers. Thus, all PO-REQUEST messages undergo

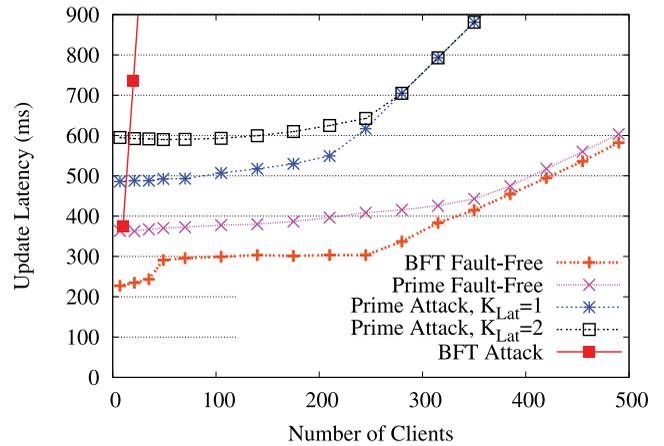


Fig. 6. Latency of Prime and BFT as a function of the number of clients in a seven-server configuration.

a reconciliation step, which consumes approximately the same outgoing bandwidth as the dissemination of the PO-REQUEST messages during the Preordering subprotocol.

To attack BFT, we use the attack described in Section 2.1. We present results for a very aggressive yet possible timeout (300 ms). This yields the most favorable performance for BFT under attack.

7.3 Results

Fig. 5 shows system throughput, measured in updates/sec, as a function of the number of clients in the emulated wide-area deployment. Fig. 6 shows the corresponding update latency, measured at the client. In the fault-free scenario, the throughput of BFT increases at a faster rate than the throughput of Prime because BFT has fewer protocol rounds. BFT's performance plateaus due to bandwidth constraints at slightly fewer than 850 updates/sec, with about 250 clients. Prime reaches a similar plateau with about 350 clients. As seen in Fig. 6, BFT has a lower latency than Prime when the protocols are not under attack, due to the differences in the number of protocol rounds. The latency of both protocols increases at different points before the plateau due to overhead associated with aggregation. The latency begins to climb steeply when the throughput plateaus due to update queuing at the servers.

The throughput results are different when the two protocols are attacked. With an aggressive time-out of 300 ms, BFT can order fewer than 30 updates/sec. With the default time-out of five seconds, BFT can only order two updates/sec (not shown). Prime plateaus at about 400 updates/sec due to the bandwidth overhead incurred by the Reconciliation subprotocol. Prime's throughput continues to increase until it becomes bandwidth constrained. BFT reaches its maximum throughput when there is one client per server. This throughput limitation, which occurs when only a small amount of the available bandwidth is used, is a consequence of judging the leader conservatively.

The slope of the curve corresponding to Prime under attack is less steep than when it is not under attack due to the delay added by the malicious leader. We include results with $K_{Lat} = 1$ and $K_{Lat} = 2$. K_{Lat} accounts for variability in latency (see Section 3). As K_{Lat} increases, a malicious leader can add more delay to the turnaround time without being

detected. The amount of delay that can be added by a malicious leader is directly proportional to K_{Lat} . For example, if K_{Lat} were set to 10, the leader could add roughly 10 round-trip times of delay without being suspected. When under attack, the latency of Prime increases due to the two extra protocol rounds added by the leader. When $K_{Lat} = 2$, the leader can add approximately 100 ms more delay than when $K_{Lat} = 1$. The latency of BFT under attack climbs as soon as more than one client is added to each server because the leader can order one update per server per time-out without being suspected.

8 RELATED WORK

This paper focused on leader-based Byzantine fault-tolerant protocols [6], [7], [8], [9], [10], [11], [12] that achieve replication via the state machine approach [22], [23]. The consistency of these systems does not rely on synchrony assumptions, while liveness is guaranteed assuming the network meets certain stability properties. To ensure that the stability properties are eventually met in practice, they use exponentially growing time-outs during view changes. This makes these systems vulnerable to the type of performance degradation when under attack described in Section 2.2. In contrast, Prime uses the Suspect-Leader subprotocol to allow correct servers to collectively decide whether the leader is performing fast enough by adapting to the network conditions once the system stabilizes. Aiyer et al. [15] first noted the problems that can be caused by a faulty primary and suggested rotating the primary to mitigate its attacks. Prime takes a different approach, enforcing timely behavior from any leader that remains in power and eventually settling on a leader that provides good performance. Singh et al. [24] demonstrate how the performance of different protocols can degrade under unfavorable network conditions.

More recently, the Aardvark system of Clement et al. [13] proposed building *robust* Byzantine replication systems that sacrifice some normal-case performance in order to ensure that performance remains acceptably high when the system exhibits Byzantine failures. The approaches taken by Prime and Aardvark are quite different. Prime aims to guarantee that *every* request known to correct servers will be executed in a timely manner, limiting the leader's responsibilities in order to enforce timeliness exactly where it is needed. Aardvark aims to guarantee that over sufficiently long periods, system throughput remains within a constant factor of what it would be if only correct servers were participating in the protocol. It achieves this by gradually increasing the level of work expected from the leader, which ensures that view changes take place. Aardvark guarantees high throughput when the system is saturated, but individual requests may take longer to execute (e.g., if they are introduced during the grace period that begins any view with a faulty primary). The Spinning protocol of Veronese et al. [14] constantly rotates the primary to reduce the impact of faulty servers.

Rampart [25] implements Byzantine atomic multicast over a reliable group multicast protocol. This is similar to how Prime uses preordering followed by global ordering. Both protocols disseminate requests to $2f + 1$ servers before a coordinator assigns the global order. Drabkin et al. [26]

observe the difficulty of setting time-outs in the context of group communication in malicious settings. Prime's Reconciliation subprotocol uses erasure codes for efficient data dissemination. A similar approach was taken by Cachin and Tessaro [27] and Fitzi and Hirt [28].

Other Byzantine fault-tolerant protocols [4], [5], [29], [30] use randomization to circumvent the FLP impossibility result, guaranteeing termination with probability one. These protocols incur a high number of communication rounds during normal-case operation (even those that terminate in an expected constant number of rounds). However, they do not rely on a leader to coordinate the ordering protocol and thus may not suffer the same kinds of performance vulnerabilities when under attack.

Byzantine quorum systems [31], [32], [33], [34] can also be used for replication. While early work in this area was restricted to a read/write interface, recent work uses quorum systems to provide state machine replication. The Q/U protocol [33] requires $5f + 1$ replicas for this purpose and suffers performance degradation when write contention occurs. The HQ protocol [34] showed how to mitigate this cost by reducing the number of replicas to $3f + 1$. Since HQ uses BFT to resolve contention when it arises, it is vulnerable to the same types of performance degradation as BFT.

A different approach to state machine replication is to use a hybrid architecture in which different parts of the system rely on different fault and/or timing assumptions [35], [36], [37]. The different components are therefore resilient to different types of attacks. We believe leveraging stronger timing assumptions may allow for more aggressive performance monitoring.

9 CONCLUSIONS

In this paper, we pointed out the vulnerability of current leader-based Byzantine fault-tolerant state machine replication protocols to performance degradation when under attack. We proposed the BOUNDED-DELAY correctness criterion to require consistent performance in all executions, even when the system exhibits Byzantine faults. We presented Prime, a new Byzantine fault-tolerant state machine replication protocol, which meets BOUNDED-DELAY and is an important step toward making Byzantine fault-tolerant replication resilient to performance attacks in malicious environments. Our experimental results show that Prime performs competitively with existing protocols in fault-free configurations and an order of magnitude better when under attack in the configurations tested.

ACKNOWLEDGMENTS

A preliminary version of this paper appeared in the Proceedings of the 38th IEEE International Conference on Dependable Systems and Networks, 2008 [1]. This publication was supported by Grants 0430271 and 0716620 from the US National Science Foundation. Its contents are solely the responsibility of the authors and do not necessarily represent the official view of Johns Hopkins University or the US National Science Foundation.

REFERENCES

- [1] Y. Amir, B. Coan, J. Kirsch, and J. Lane, "Byzantine Replication under Attack," *Proc. 38th IEEE Int'l Conf. Dependable Systems and Networks*, pp. 197-206, 2008.
- [2] M.J. Fischer, N.A. Lynch, and M.S. Paterson, "Impossibility of Distributed Consensus with One Faulty Process," *J. ACM*, vol. 32, no. 2, pp. 374-382, 1985.
- [3] C. Dwork, N. Lynch, and L. Stockmeyer, "Consensus in the Presence of Partial Synchrony," *J. ACM*, vol. 35, no. 2, pp. 288-323, 1988.
- [4] M. Ben-Or, "Another Advantage of Free Choice (Extended Abstract): Completely Asynchronous Agreement Protocols," *Proc. Second Ann. ACM Symp. Principles of Distributed Computing*, pp. 27-30, 1983.
- [5] M.O. Rabin, "Randomized Byzantine Generals," *Proc. 24th Ann. Symp. Foundations of Computer Science*, pp. 403-409, 1983.
- [6] M. Castro and B. Liskov, "Practical Byzantine Fault Tolerance and Proactive Recovery," *ACM Trans. Computer Systems*, vol. 20, no. 4, pp. 398-461, 2002.
- [7] R. Kotla, L. Alvisi, M. Dahlin, A. Clement, and E. Wong, "Zyzzyva: Speculative Byzantine Fault Tolerance," *ACM Trans. Computer Systems*, vol. 27, no. 4, pp. 7:1-7:39, 2009.
- [8] Y. Amir, C. Danilov, D. Dolev, J. Kirsch, J. Lane, C. Nita-Rotaru, J. Olsen, and D. Zage, "Steward: Scaling Byzantine Fault-Tolerant Replication to Wide Area Networks," *IEEE Trans. Dependable and Secure Computing*, vol. 7, no. 1, pp. 80-93, Jan-Mar. 2010.
- [9] J. Yin, J.-P. Martin, A. Venkataramani, L. Alvisi, and M. Dahlin, "Separating Agreement from Execution for Byzantine Fault-Tolerant Services," *Proc. 19th ACM Symp. Operating Systems Principles*, pp. 253-267, 2003.
- [10] J.-P. Martin and L. Alvisi, "Fast Byzantine Consensus," *IEEE Trans. Dependable and Secure Computing*, vol. 3, no. 3, pp. 202-215, July-Sept. 2006.
- [11] Y. Amir, B. Coan, J. Kirsch, and J. Lane, "Customizable Fault Tolerance for Wide-Area Replication," *Proc. 26th IEEE Int'l Symp. Reliable Distributed Systems*, pp. 66-80, 2007.
- [12] J. Li and D. Mazieres, "Beyond One-Third Faulty Replicas in Byzantine Fault Tolerant Systems," *Proc. Fourth USENIX Symp. Networked Systems Design and Implementation*, pp. 131-144, 2007.
- [13] A. Clement, E. Wong, L. Alvisi, M. Dahlin, and M. Marchetti, "Making Byzantine Fault Tolerant Systems Tolerate Byzantine Faults," *Proc. Sixth USENIX Symp. Networked Systems Design and Implementation*, pp. 153-168, 2009.
- [14] G.S. Veronese, M. Correia, A.N. Bessani, and L.C. Lung, "Spin One's Wheels? Byzantine Fault Tolerance with a Spinning Primary," *Proc. 28th IEEE Int'l Symp. Reliable Distributed Systems*, pp. 135-144, 2009.
- [15] A.S. Aiyer, L. Alvisi, A. Clement, M. Dahlin, J.-P. Martin, and C. Porth, "BAR Fault Tolerance for Cooperative Services," *Proc. 20th ACM Symp. Operating Systems Principles*, pp. 45-58, 2005.
- [16] S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, and W. Weiss, "An Architecture for Differentiated Services," RFC 2475, 1998.
- [17] M.P. Herlihy and J.M. Wing, "Linearizability: A Correctness Condition for Concurrent Objects," *ACM Trans. Programming and Languages and Systems*, vol. 12, no. 3, pp. 463-492, 1990.
- [18] M. Castro, "Practical Byzantine Fault Tolerance," PhD dissertation, Massachusetts Inst. of Technology, pp. 29-31, 2001.
- [19] F.J. MacWilliams and N.J.A. Sloane, *The Theory of Error-Correcting Codes*, North-Holland, 1988.
- [20] R.L. Rivest, A. Shamir, and L. Adleman, "A Method for Obtaining Digital Signatures and Public-Key Cryptosystems," *Comm. ACM*, vol. 21, no. 2, pp. 120-126, 1978.
- [21] "The netem Utility," <http://www.linuxfoundation.org/collaborate/workgroups/networking/netem>, 2010.
- [22] L. Lamport, "Time, Clocks, and the Ordering of Events in a Distributed System," *Comm. ACM*, vol. 21, no. 7, pp. 558-565, 1978.
- [23] F.B. Schneider, "Implementing Fault-Tolerant Services Using the State Machine Approach: A Tutorial," *ACM Computing Surveys*, vol. 22, no. 4, pp. 299-319, 1990.
- [24] A. Singh, T. Das, P. Maniatis, P. Druschel, and T. Roscoe, "BFT Protocols under Fire," *Proc. Fifth USENIX Symp. Networked Systems Design and Implementation*, pp. 189-204, 2008.
- [25] M.K. Reiter, "The Rampart Toolkit for Building High-Integrity Services," *Proc. Int'l Workshop Theory and Practice in Distributed Systems*, pp. 99-110, 1995.
- [26] V. Drabkin, R. Friedman, and A. Kama, "Practical Byzantine Group Communication," *Proc. 26th IEEE Int'l Conf. Distributed Computing Systems*, p. 36, 2006.
- [27] C. Cachin and S. Tessaro, "Asynchronous Verifiable Information Dispersal," *Proc. 24th IEEE Symp. Reliable Distributed Systems*, pp. 191-202, 2005.
- [28] M. Fitzi and M. Hirt, "Optimally Efficient Multi-Valued Byzantine Agreement," *Proc. 25th Ann. ACM Symp. Principles of Distributed Computing*, pp. 163-168, 2006.
- [29] C. Cachin and J.A. Portiz, "Secure Intrusion-Tolerant Replication on the Internet," *Proc. IEEE Int'l Conf. Dependable Systems and Networks*, pp. 167-176, 2002.
- [30] H. Moniz, N.F. Neves, M. Correia, and P. Verissimo, "Randomized Intrusion-Tolerant Asynchronous Services," *Proc. IEEE Int'l Conf. Dependable Systems and Networks*, pp. 568-577, 2006.
- [31] D. Malkhi and M. Reiter, "Byzantine Quorum Systems," *Distributed Computing*, vol. 11, no. 4, pp. 203-213, 1998.
- [32] D. Malkhi and M.K. Reiter, "Secure and Scalable Replication in Phalanx," *Proc. 17th IEEE Symp. Reliable Distributed Systems*, pp. 51-58, 1998.
- [33] M. Abd-El-Malek, G.R. Ganger, G.R. Goodson, M.K. Reiter, and J.J. Wylie, "Fault-Scalable Byzantine Fault-Tolerant Services," *Proc. 20th ACM Symp. Operating Systems Principles*, pp. 59-74, 2005.
- [34] J. Cowling, D. Myers, B. Liskov, R. Rodrigues, and L. Shrira, "HQ Replication: A Hybrid Quorum Protocol for Byzantine Fault Tolerance," *Proc. Seventh USENIX Symp. Operating Systems Design and Implementation*, pp. 177-190, 2006.
- [35] P.E. Verissimo, N.F. Neves, C. Cachin, J. Portiz, D. Powell, Y. Deswarte, R. Stroud, and I. Welch, "Intrusion-Tolerant Middleware: The Road to Automatic Security," *IEEE Security & Privacy*, vol. 4, no. 4, pp. 54-62, July-Aug. 2006.
- [36] M. Correia, N.F. Neves, and P. Verissimo, "How to Tolerate Half Less One Byzantine Nodes in Practical Distributed Systems," *Proc. 23rd IEEE Int'l Symp. Reliable Distributed Systems*, pp. 174-183, 2004.
- [37] M. Serafini and N. Suri, "The Fail-Heterogeneous Architectural Model," *Proc. 26th IEEE Int'l Symp. Reliable Distributed Systems*, pp. 103-113, 2007.
- [38] G. Bracha, "An Asynchronous $[(n-1)/3]$ -Resilient Consensus Protocol," *Proc. Third Ann. ACM Symp. Principles of Distributed Computing*, pp. 154-162, 1984.
- [39] "The BFT Project Homepage," <http://www.pmg.csail.mit.edu/bft>, 2010.
- [40] R.C. Merkle, "Secrecy, Authentication, and Public Key Systems," PhD dissertation, Stanford Univ. 1979.



Yair Amir received the BS and MS degrees from the Technion, Israel Institute of Technology, in 1985 and 1990, respectively, and the PhD degree from the Hebrew University of Jerusalem, Israel, in 1995. He is a professor in the Department of Computer Science, Johns Hopkins University, where he served as an assistant professor since 1995, an associate professor since 2000, and professor since 2004. Prior to his PhD, he gained extensive experience building C3I systems. He is a creator of the Spread and Secure Spread messaging toolkits, the Backhand and Wackamole clustering projects, the Spines overlay network platform, and the SMesh wireless mesh network. He has been a member of the program committees of the IEEE International Conference on Distributed Computing Systems (1999, 2002, and 2005-07), the ACM Conference on Principles of Distributed Computing (2001), and the International Conference on Dependable Systems and Networks (2001, 2003, and 2005). He is a member of the ACM and the IEEE Computer Society.



Brian Coan received the BSE degree from Princeton University in 1977, the MS degree from Stanford in 1979, and the PhD degree in computer science from MIT in the Theory of Distributed Systems Group in 1987. He is director of the Distributed Computing group at Telcordia, where he has worked since 1978. He works on providing reliable networking and information services in adverse network environments. He is a member of the ACM.



John Lane received the BA degree in biology from Cornell University in 1992, the MSE degree in computer science from Johns Hopkins University in 2006, and the PhD degree in computer science from Johns Hopkins University in 2008. He is a senior research scientist at LiveTimeNet, where he has worked since 2008. His research interests include distributed systems, replication, and byzantine fault tolerance. He is a member of the ACM, the IEEE, and the IEEE Computer Society.



Jonathan Kirsch received the BSc degree from Yale University in 2004, the MSE degree from Johns Hopkins University in 2007, and the PhD degree in computer science from Johns Hopkins University in 2010. He is currently a research scientist at the Siemens Technology to Business Center in Berkeley, California. His research interests include fault-tolerant replication and survivable systems. He is a member of the IEEE.

▷ **For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.**