

# The Impact of RDMA on Agreement

## [Extended Version]

Marcos K. Aguilera  
VMware  
maguilera@vmware.com

Naama Ben-David  
CMU  
nbendavi@cs.cmu.edu

Rachid Guerraoui  
EPFL  
rachid.guerraoui@epfl.ch

Virendra Marathe  
Oracle  
virendra.marathe@oracle.com

Igor Zablotchi  
EPFL  
igor.zablotchi@epfl.ch

### ABSTRACT

Remote Direct Memory Access (RDMA) is becoming widely available in data centers. This technology allows a process to directly read and write the memory of a remote host, with a mechanism to control access permissions. In this paper, we study the fundamental power of these capabilities. We consider the well-known problem of achieving consensus despite failures, and find that RDMA can improve the inherent trade-off in distributed computing between failure resilience and performance. Specifically, we show that RDMA allows algorithms that simultaneously achieve high resilience and high performance, while traditional algorithms had to choose one or another. With Byzantine failures, we give an algorithm that only requires  $n \geq 2fp + 1$  processes (where  $fp$  is the maximum number of faulty processes) and decides in two (network) delays in common executions. With crash failures, we give an algorithm that only requires  $n \geq fp + 1$  processes and also decides in two delays. Both algorithms tolerate a minority of memory failures inherent to RDMA, and they provide safety in asynchronous systems and liveness with standard additional assumptions.

### 1 INTRODUCTION

In recent years, a technology known as Remote Direct Memory Access (RDMA) has made its way into data centers, earning a spotlight in distributed systems research. RDMA provides the traditional send/receive communication primitives, but also allows a process to directly read/write remote memory. Research work shows that RDMA leads to some new and exciting distributed algorithms [3, 9, 25, 31, 45, 49].

RDMA provides a different interface from previous communication mechanisms, as it combines message-passing with shared-memory [3]. Furthermore, to safeguard the remote memory, RDMA provides *protection* mechanisms to grant and revoke access for reading and writing data. This mechanism is fine grained: an application can choose subsets of remote memory called *regions* to protect; it can choose whether a region can be read, written, or both; and it can choose individual processes to be given access, where different processes can have different accesses. Furthermore, protections are *dynamic*: they can be changed by the application over time. In this paper, we lay the groundwork for a theoretical understanding of these RDMA capabilities, and we show that they lead to distributed algorithms that are inherently more powerful than before.

While RDMA brings additional power, it also introduces some challenges. With RDMA, the remote memories are subject to failures that cause them to become unresponsive. This behavior differs from

traditional shared memory, which is often assumed to be reliable<sup>1</sup>. In this paper, we show that the additional power of RDMA more than compensates for these challenges.

Our main contribution is to show that RDMA improves on the fundamental trade-off in distributed systems between failure resilience and performance—specifically, we show how a consensus protocol can use RDMA to achieve *both* high resilience and high performance, while traditional algorithms had to choose one or another. We illustrate this on the fundamental problem of achieving consensus and capture the above RDMA capabilities as an M&M model [3], in which processes can use both message-passing and shared-memory. We consider asynchronous systems and require safety in all executions and liveness under standard additional assumptions (e.g., partial synchrony). We measure resiliency by the number of failures an algorithm tolerates, and performance by the number of (network) delays in common-case executions. Failure resilience and performance depend on whether processes fail by crashing or by being Byzantine, so we consider both.

With Byzantine failures, we consider the consensus problem called weak Byzantine agreement, defined by Lamport [37]. We give an algorithm that (a) requires only  $n \geq 2fp + 1$  processes (where  $fp$  is the maximum number of faulty processes) and (b) decides in two delays in the common case. With crash failures, we give the first algorithm for consensus that requires only  $n \geq fp + 1$  processes and decides in two delays in the common case. With both Byzantine or crash failures, our algorithms can also tolerate crashes of memory—only  $m \geq 2f_M + 1$  memories are required, where  $f_M$  is the maximum number of faulty memories. Furthermore, with crash failures, we improve resilience further, to tolerate crashes of a minority of the combined set of memories and processes.

Our algorithms appear to violate known impossibility results: it is known that with message-passing, Byzantine agreement requires  $n \geq 3fp + 1$  even if the system is synchronous [44], while consensus with crash failures require  $n \geq 2fp + 1$  if the system is partially synchronous [27]. There is no contradiction: our algorithms rely on the power of RDMA, not available in other systems.

RDMA's power comes from two features: (1) simultaneous access to message-passing and shared-memory, and (2) dynamic permissions. Intuitively, shared-memory helps resilience, message-passing helps performance, and dynamic permissions help both.

To see how shared-memory helps resilience, consider the Disk Paxos algorithm [29], which uses shared-memory (disks) but no messages. Disk Paxos requires only  $n \geq fp + 1$  processes, matching

<sup>1</sup>There are a few studies of failure-prone memory, as we discuss in related work.

the resilience of our algorithm. However, Disk Paxos is not as fast: it takes at least four delays. In fact, we show that no shared-memory consensus algorithm can decide in two delays (Section 6).

To see how message-passing helps performance, consider the Fast Paxos algorithm [39], which uses message-passing and no shared-memory. Fast Paxos decides in only two delays in common executions, but it requires  $n \geq 2fp + 1$  processes.

Of course, the challenge is achieving both high resilience and good performance in a single algorithm. This is where RDMA’s dynamic permissions shine. Clearly, dynamic permissions improve resilience against Byzantine failures, by preventing a Byzantine process from overwriting memory and making it useless. More surprising, perhaps, is that dynamic permissions help performance, by providing an uncontended instantaneous guarantee: if each process revokes the write permission of other processes before writing to a register, then a process that writes successfully knows that it executed uncontended, without having to take additional steps (e.g., to read the register). We use this technique in our algorithms for both Byzantine and crash failures.

In summary, our contributions are as follows:

- We consider distributed systems with RDMA, and we propose a model that captures some of its key properties while accounting for failures of processes and memories, with support of dynamic permissions.
- We show that the shared-memory part of our RDMA improves resilience: our Byzantine agreement algorithm requires only  $n \geq 2fp + 1$  processes.
- We show that the shared-memory by itself does not permit consensus algorithms that decide in two steps in common executions.
- We show that with dynamic permissions, we can improve the performance of our Byzantine Agreement algorithm, to decide in two steps in common executions.
- We give similar results for the case of crash failures: decision in two steps while requiring only  $n \geq fp + 1$  processes.
- Our algorithms can tolerate the failure of memories, up to a minority of them.

The rest of the paper is organized as follows. Section 2 gives an overview of related work. In Section 3 we formally define the RDMA-compliant M&M model that we use in the rest of the paper, and specify the agreement problems that we solve. We then proceed to present the main contributions of the paper. Section 4 presents our fast and resilient Byzantine agreement algorithm. In Section 5 we consider the special case of crash-only failures, and show an improvement of the algorithm and tolerance bounds for this setting. In Section 6 we briefly outline a lower bound that shows that the dynamic permissions of RDMA are necessary for achieving our results. Finally, in Section 7 we discuss the semantics of RDMA in practice, and how our model reflects these features.

To ease readability, most proofs have been deferred to the Appendices.

## 2 RELATED WORK

**RDMA.** Many high-performance systems were recently proposed using RDMA, such as distributed key-value stores [25, 31], communication primitives [25, 33], and shared address spaces across

clusters [25]. Kalia *et al.* [32] provide guidelines for designing systems using RDMA. RDMA has also been applied to solve consensus [9, 45, 49]. Our model shares similarities with DARE [45] and APUS [49], which modify queue-pair state at run time to prevent or allow access to memory regions, similar to our dynamic permissions. These systems perform better than TCP/IP-based solutions, by exploiting better raw performance of RDMA, without changing the fundamental communication complexity or failure-resilience of the consensus protocol. Similarly, Rüsç *et al.* [46] use RDMA as a replacement for TCP/IP in existing BFT protocols.

**M&M.** Message-and-memory (M&M) refers to a broad class of models that combine message-passing with shared-memory, introduced by Aguilera *et al.* in [3]. In that work, Aguilera *et al.* consider M&M models without memory permissions and failures, and show that such models lead to algorithms that are more robust to failures and asynchrony. In particular, they give a consensus algorithm that tolerates more crash failures than message-passing systems, but is more scalable than shared-memory systems, as well as a leader election algorithm that reduces the synchrony requirements. In this paper, our goal is to understand how memory permissions and failures in RDMA impact agreement.

**Byzantine Fault Tolerance.** Lamport, Shostak and Pease [40, 44] show that Byzantine agreement can be solved in synchronous systems iff  $n \geq 3fp + 1$ . With unforgeable signatures, Byzantine agreement can be solved iff  $n \geq 2fp + 1$ . In asynchronous systems subject to failures, consensus cannot be solved [28]. However, this result is circumvented by making additional assumptions for liveness, such as randomization [10] or partial synchrony [18, 27]. Many Byzantine agreement algorithms focus on safety and implicitly use the additional assumptions for liveness. Even with signatures, asynchronous Byzantine agreement can be solved only if  $n \geq 3fp + 1$  [15].

It is well known that the resilience of Byzantine agreement varies depending on various model assumptions like synchrony, signatures, equivocation, and the exact variant of the problem to be solved. A system that has non-equivocation is one that can prevent a Byzantine process from sending different values to different processes. Table 1 summarizes some known results that are relevant to this paper.

Work	Synchrony	Signatures	Non-Equiv	Strong Validity	Resiliency
[40]	✓	✓	✗	✓	$2f + 1$
[40]	✓	✗	✗	✓	$3f + 1$
[4, 41]	✗	✓	✓	✓	$3f + 1$
[21]	✗	✓	✗	✗	$3f + 1$
[21]	✗	✗	✓	✗	$3f + 1$
[21]	✗	✓	✓	✗	$2f + 1$
This paper	✗	✓	✗ (RDMA)	✗	$2f + 1$

**Table 1: Known fault tolerance results for Byzantine agreement.**

Our Byzantine agreement results share similarities with results for shared memory. Malkhi *et al.* [41] and Alon *et al.* [4] show bounds on the resilience of strong and weak consensus in a model with reliable memory but Byzantine processes. They also provide consensus protocols, using read-write registers enhanced with sticky bits (write-once memory) and access control lists not unlike our permissions. Bessani *et al.* [11] propose an alternative to sticky bits

and access control lists through Policy-Enforced Augmented Tuple Spaces. All these works handle Byzantine failures with powerful objects rather than registers. Bouzid *et al.* [13] show that  $3f_p + 1$  processes are necessary for strong Byzantine agreement with read-write registers.

Some prior work solves Byzantine agreement with  $2f_p + 1$  processes using specialized trusted components that Byzantine processes cannot control [19, 20, 22, 23, 34, 48]. Some schemes decide in two delays but require a large trusted component: a coordinator [19], reliable broadcast [23], or message ordering [34]. For us, permission checking in RDMA is a trusted component of sorts, but it is small and readily available.

At a high-level, our improved Byzantine fault tolerance is achieved by preventing equivocation by Byzantine processes, thereby effectively translating each Byzantine failure into a crash failure. Such translations from one type of failure into a less serious one have appeared extensively in the literature [8, 15, 21, 43]. Early work [8, 43] shows how to translate a crash tolerant algorithm into a Byzantine tolerant algorithm in the synchronous setting. Bracha [14] presents a similar translation for the asynchronous setting, in which  $n \geq 3f_p + 1$  processes are required to tolerate  $f_p$  Byzantine failures. Bracha’s translation relies on the definition and implementation of a reliable broadcast primitive, very similar to the one in this paper. However, we show that using the capabilities of RDMA, we can implement it with higher fault tolerance.

**Faulty memory.** Afek *et al.* [2] and Jayanti *et al.* [30] study the problem of masking the benign failures of shared memory or objects. We use their ideas of replicating data across memories. Abraham *et al.* [1] considers honest processes but malicious memory.

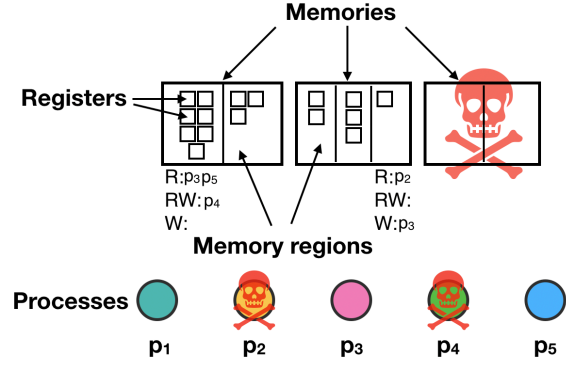
**Common-case executions.** Many systems and algorithms tolerate adversarial scheduling but optimize for common-case executions without failures, asynchrony, contention, etc (e.g., [12, 24, 26, 35, 36, 39, 42]). None of these match both the resilience and performance of our algorithms. Some algorithms decide in one delay but require  $n \geq 5f_p + 1$  for Byzantine failures [47] or  $n \geq 3f_p + 1$  for crash failures [16, 24].

### 3 MODEL AND PRELIMINARIES

We consider a message-and-memory (M&M) model, which allows processes to use both message-passing and shared-memory [3]. The system has  $n$  processes  $P = \{p_1, \dots, p_n\}$  and  $m$  (shared) memories  $M = \{\mu_1, \dots, \mu_m\}$ . Processes communicate by accessing memories or sending messages. Throughout the paper, memory refers to the shared memories, not the local state of processes.

The system is asynchronous in that it can experience arbitrary delays. We expect algorithms to satisfy the safety properties of the problems we consider, under this asynchronous system. For liveness, we require additional standard assumptions, such as partial synchrony, randomization, or failure detection.

**Memory permissions.** Each memory consists of a set of *registers*. To control access, an algorithm groups those registers into a set of (possibly overlapping) *memory regions*, and then defines permissions for those memory regions. Formally, a memory region  $mr$  of a memory  $\mu$  is a subset of the registers of  $\mu$ . We often refer to  $mr$  without specifying the memory  $\mu$  explicitly. Each memory region  $mr$  has a *permission*, which consists of three disjoint sets



**Figure 1: Our model with processes and memories, which may both fail. Processes can send messages to each other or access registers in the memories. Registers in a memory are grouped into memory regions that may overlap, but in our algorithms they do not. Each region has a permission indicating what processes can read, write, and read-write the registers in the region (shown for two regions).**

of processes  $R_{mr}$ ,  $W_{mr}$ ,  $RW_{mr}$  indicating whether each process can read, write, or read-write the registers in the region. We say that  $p$  has *read permission* on  $mr$  if  $p \in R_{mr}$  or  $p \in RW_{mr}$ ; we say that  $p$  has *write permission* on  $mr$  if  $p \in W_{mr}$  or  $p \in RW_{mr}$ . In the special case when  $R_{mr} = P \setminus \{p\}$ ,  $W_{mr} = \emptyset$ ,  $RW_{mr} = \{p\}$ , we say that  $mr$  is a Single-Writer Multi-Reader (SWMR) region—registers in  $mr$  correspond to the traditional notion of SWMR registers. Note that a register may belong to several regions, and a process may have access to the register on one region but not another—this models the existing RDMA behavior. Intuitively, when reading or writing data, a process specifies the region and the register, and the system uses the region to determine if access is allowed (we make this precise below).

**Permission change.** An algorithm indicates an initial permission for each memory region  $mr$ . Subsequently, the algorithm may wish to change the permission of  $mr$  during execution. For that, processes can invoke an operation  $changePermission(mr, new\_perm)$ , where  $new\_perm$  is a triple  $(R, W, RW)$ . This operation returns no results and it is intended to modify  $R_{mr}$ ,  $W_{mr}$ ,  $RW_{mr}$  to  $R, W, RW$ . To tolerate Byzantine processes, an algorithm can restrict processes from changing permissions. For that, the algorithm specifies a function  $legalChange(p, mr, old\_perm, new\_perm)$  which returns a boolean indicating whether process  $p$  can change the permission of  $mr$  to  $new\_perm$  when the current permissions are  $old\_perm$ . More precisely, when  $changePermission$  is invoked, the system evaluates  $legalChange$  to determine whether  $changePermission$  takes effect or becomes a no-op. When  $legalChange$  always returns false, we

say that the *permissions are static*; otherwise, the *permissions are dynamic*.

**Accessing memories.** Processes access the memories via operations  $write(mr, r, v)$  and  $read(mr, r)$  for memory region  $mr$ , register  $r$ , and value  $v$ . A  $write(mr, r, v)$  by process  $p$  changes register  $r$  to  $v$  and returns *ack* if  $r \in mr$  and  $p$  has write permission on  $mr$ ; otherwise, the operation returns *nak*. A  $read(mr, r)$  by process  $p$  returns the last value successfully written to  $r$  if  $r \in mr$  and  $p$  has read permission on  $mr$ ; otherwise, the operation returns *nak*. In our algorithms, a register belongs to exactly one region, so we omit the  $mr$  parameter from write and read operations.

**Sending messages.** Processes can also communicate by sending messages over a set of directed links. We assume messages are unique. If there is a link from process  $p$  to process  $q$ , then  $p$  can send messages to  $q$ . Links satisfy two properties: *integrity* and *no-loss*. Given two correct processes  $p$  and  $q$ , integrity requires that a message  $m$  be received by  $q$  from  $p$  at most once and only if  $m$  was previously sent by  $p$  to  $q$ . No-loss requires that a message  $m$  sent from  $p$  to  $q$  be eventually received by  $q$ . In our algorithms, we typically assume a fully connected network so that every pair of correct processes can communicate. We also consider the special case when there are no links (see below).

**Executions and steps.** An execution is as a sequence of process steps. In each step, a process does the following, according to its local state: (1) sends a message or invokes an operation on a memory (read, write, or changePermission), (2) tries to receive a message or a response from an outstanding operation, and (3) changes local state. We require a process to have at most one outstanding operation on each memory.

**Failures.** A memory  $m$  may fail by crashing, which causes subsequent operations on its registers to hang without returning a response. Because the system is asynchronous, a process cannot differentiate a crashed memory from a slow one. We assume there is an upper bound  $f_M$  on the maximum number of memories that may crash. Processes may fail by crashing or becoming Byzantine. If a process crashes, it stops taking steps forever. If a process becomes Byzantine, it can deviate arbitrarily from the algorithm. However, that process cannot operate on memories without the required permission. We assume there is an upper bound  $f_P$  on the maximum number of processes that may be faulty. Where the context is clear, we omit the  $P$  and  $M$  subscripts from the number of failures,  $f$ .

**Signatures.** Our algorithms assume unforgeable signatures: there are primitives  $sign(v)$  and  $sValid(p, v)$  which, respectively, signs a value  $v$  and determines if  $v$  is signed by process  $p$ .

**Messages and disks.** The model defined above includes two common models as special cases. In the *message-passing* model, there are no memories ( $m = 0$ ), so processes can communicate only by sending messages. In the *disk* model [29], there are no links, so processes can communicate only via memories; moreover, each memory has a single region which always permits all processes to read and write all registers.

## Consensus

In the consensus problem, processes propose an initial value and must make an irrevocable decision on a value. With crash failures, we require the following properties:

- **Uniform Agreement.** If processes  $p$  and  $q$  decide  $v_p$  and  $v_q$ , then  $v_p = v_q$ .
- **Validity.** If some process decides  $v$ , then  $v$  is the initial value proposed by some process.
- **Termination.** Eventually all correct processes decide.

We expect Agreement and Validity to hold in an asynchronous system, while Termination requires standard additional assumptions (partial synchrony, randomization, failure detection, etc). With Byzantine failures, we change these definitions so the problem can be solved. We consider weak Byzantine agreement [37], with the following properties:

- **Agreement.** If correct processes  $p$  and  $q$  decide  $v_p$  and  $v_q$ , then  $v_p = v_q$ .
- **Validity.** With no faulty processes, if some process decides  $v$ , then  $v$  is the input of some process.
- **Termination.** Eventually all correct processes decide.

**Complexity of algorithms.** We are interested in the performance of algorithms in *common-case executions*, when the system is synchronous and there are no failures. In those cases, we measure performance using the notion of *delays*, which extends message-delays to our model. Under this metric, computations are instantaneous, each message takes one delay, and each memory operation takes two delays. Intuitively, a delay represents the time incurred by the network to transmit a message; a memory operation takes two delays because its hardware implementation requires a round trip. We say that a consensus protocol is *k-deciding* if, in common-case executions, some process decides in  $k$  delays.

## 4 BYZANTINE FAILURES

We now consider Byzantine failures and give a 2-deciding algorithm for weak Byzantine agreement with  $n \geq 2f_P + 1$  processes and  $m \geq 2f_M + 1$  memories. The algorithm consists of the composition of two sub-algorithms: a slow one that always works, and a fast one that gives up under hard conditions.

The first sub-algorithm, called *Robust Backup*, is developed in two steps. We first implement a *reliable broadcast* primitive, which prevents Byzantine processes from sending different values to different processes. Then, we use the framework of Clement *et al.* [21] combined with this primitive to convert a message-passing consensus algorithm that tolerates crash failures into a consensus algorithm that tolerates Byzantine failures. This yields Robust Backup.<sup>2</sup> It uses only static permissions and assumes memories are split into SWMR regions. Therefore, this sub-algorithm works in the traditional shared-memory model with SWMR registers, and it may be of independent interest.

The second sub-algorithm is called *Cheap Quorum*. It uses dynamic permissions to decide in two delays using one signature in common executions. However, the sub-algorithm gives up if the system is not synchronous or there are Byzantine failures.

Finally, we combine both sub-algorithms using ideas from the Abstract framework of Aublin *et al.* [7]. More precisely, we start by running Cheap Quorum; if it aborts, we run Robust Backup. There

<sup>2</sup>The attentive reader may wonder why at this point we have not achieved a 2-deciding algorithm already: if we apply Clement *et al.* [21] to a 2-deciding crash-tolerant algorithm (such as Fast Paxos [12]), will the result not be a 2-deciding Byzantine-tolerant algorithm? The answer is no, because Clement *et al.* needs reliable broadcast, which incurs at least 6 delays.

is a subtlety: for this idea to work, Robust Backup must decide on a value  $v$  if Cheap Quorum decided  $v$  previously. To do that, Robust Backup decides on a *preferred value* if at least  $f + 1$  processes have this value as input. To do so, we use the classic crash-tolerant Paxos algorithm (run under the Robust Backup algorithm to ensure Byzantine tolerance) but with an initial set-up phase that ensures this safe decision. We call the protocol *Preferential Paxos*.

#### 4.1 The Robust Backup Sub-Algorithm

We develop Robust Backup using the construction by Clement *et al.* [21], which we now explain. Clement *et al.* show how to transform a message-passing algorithm  $\mathcal{A}$  that tolerates  $f_p$  crash failures into a message-passing algorithm that tolerates  $f_p$  Byzantine failures in a system where  $n \geq 2f_p + 1$  processes, assuming unforgeable signatures and a non-equivocation mechanism. They do so by implementing trusted message-passing primitives, *T-send* and *T-receive*, using non-equivocation and signature verification on every message. Processes include their full history with each message, and then verify locally whether a received message is consistent with the protocol. This restricts Byzantine behavior to crash failures.

To apply this construction in our model, we show that our model can implement non-equivocation and message passing. We first show that shared-memory with SWMR registers (and no memory failures) can implement these primitives, and then show how our model can implement shared-memory with SWMR registers.

**4.1.1 Reliable Broadcast.** Consider a shared-memory system. We present a way to prevent equivocation through a solution to the *reliable broadcast* problem, which we recall below. Note that our definition of reliable broadcast includes the sequence number  $k$  (as opposed to being single-shot) so as to facilitate the integration with the Clement *et al.* construction, as we explain in Section 4.1.2.

**DEFINITION 1.** Reliable broadcast is defined in terms of two primitives, *broadcast*( $k, m$ ) and *deliver*( $k, m, q$ ). When a process  $p$  invokes *broadcast*( $k, m$ ) we say that  $p$  broadcasts ( $k, m$ ). When a process  $p$  invokes *deliver*( $k, m, q$ ) we say that  $p$  delivers ( $k, m$ ) from  $q$ . Each correct process  $p$  must invoke *broadcast*( $k, *$ ) with  $k$  one higher than  $p$ 's previous invocation (and first invocation with  $k=1$ ). The following holds:

- (1) If a correct process  $p$  broadcasts ( $k, m$ ), then all correct processes eventually deliver ( $k, m$ ) from  $p$ .
- (2) If  $p$  and  $q$  are correct processes,  $p$  delivers ( $k, m$ ) from  $r$ , and  $q$  delivers ( $k, m'$ ) from  $r$ , then  $m=m'$ .
- (3) If a correct process delivers ( $k, m$ ) from a correct process  $p$ , then  $p$  must have broadcast ( $k, m$ ).
- (4) If a correct process delivers ( $k, m$ ) from  $p$ , then all correct processes eventually deliver ( $k, m'$ ) from  $p$  for some  $m'$ .

Algorithm 2 shows how to implement reliable broadcast that is tolerant to a minority of Byzantine failures in shared-memory using SWMR registers.

To broadcast its  $k$ -th message  $m$ ,  $p$  simply signs ( $k, m$ ) and writes it in slot  $Value[p, k, p]$  of its memory<sup>3</sup>.

#### Algorithm 2: Reliable Broadcast

```

1  SWMR Value[n,M,n]; initialized to  $\perp$ . Value[p] is
    $\hookrightarrow$  array of SWMR(p) registers.

3  SWMR L1Proof[n,M,n]; initialized to  $\perp$ . L1Proof[p] is
    $\hookrightarrow$  array of SWMR(p) registers.

5  SWMR L2Proof[n,M,n]; initialized to  $\perp$ . L2Proof[p] is
    $\hookrightarrow$  array of SWMR(p) registers.

7  Code for process p
8  last[n]: local array with last k delivered from each
    $\hookrightarrow$  process. Initially, last[q] = 0
9  state[n]: local array of registers. state[q]  $\in$  {
    $\hookrightarrow$  WaitForSender, WaitForL1Proof, WaitForL2Proof}.
    $\hookrightarrow$  Initially, state[q] = WaitForSender

11 broadcast (k,m){
12   Value[p,k,p].write(sign((k,m))); }

14 value checkL2proof(q,k) {
15   for i  $\in$   $\Pi$  {
16     proof = L2Proof[i,k,q].read();
17     if (proof !=  $\perp$  && check(proof)) {
18       L2Proof[p,k,q].write(proof);
19       return proof.msg; } }
20   return null; }

22 for q in  $\Pi$  in parallel {
23   while true {
24     try_deliver(q); } }

26 try_deliver(q) {
27   k = last[q];
28   val = checkL2Proof(q,k);
29   if (val != null) {
30     deliver(k, proof.msg, q);
31     last[q] += 1;
32     state = WaitForSender;
33     return;
34   }

36   if state == WaitForSender {
37     val = Value[q,k,q].read();
38     if (val== $\perp$  || !sValid(p, val) || key!=k)
39       return;
40     Value[p,k,q].write(sign(val));
41     state = WaitForL1Proof; }

43   if state == WaitForL1Proof {
44     checkedVals =  $\emptyset$ ;
45     for i  $\in$   $\Pi$ {
46       val = Value[i,k,q].read();
47       if (val!= $\perp$  && sValid(p,val) && key==k) {
48         add val to checkedVals; } }

50     if size(checkedVals)  $\geq$  majority and
        $\hookrightarrow$  checkedVals contains only one value {
51       llprf = sign(checkedVals);
52       L1Proof[p,k,q].write(llprf);
53       state = WaitForL2Proof; } }

55   if state == WaitForL2Proof{
56     checkedL1Proofs =  $\emptyset$ ;
57     for i in  $\Pi$ {
58       proof = L1Proof[i,k,q].read();
59       if ( checkL1Proof(proof) ) {
60         add proof to checkedL1Proofs; } }

62     if size(checkedL1Proofs)  $\geq$  majority {
63       l2prf = sign(checkedL1Proofs);
64       L2Proof[p,k,q].write(l2prf); } } }
```

<sup>3</sup>The indexing of the slots is as follows: the first index is the writer of the SWMR register, the second index is the sequence number of the message, and the third index is the sender of the message.

Delivering a message from another process is more involved, requiring verification steps to ensure that all correct processes will eventually deliver the same message and no other. The high-level idea is that before delivering a message  $(k, m)$  from  $q$ , each process  $p$  checks that no other process saw a different value from  $q$ , and waits to hear that “enough” other processes also saw the same value. More specifically, each process  $p$  has 3 slots per process per sequence number, that only  $p$  can write to, but all processes can read from. These slots are initialized to  $\perp$ , and  $p$  uses them to write the values that it has seen. The 3 slots represent 3 levels of ‘proofs’ that this value is correct; for each process  $q$  and sequence number  $k$ ,  $p$  has a slot to write (1) the initial value  $v$  it read from  $q$  for  $k$ , (2) a proof that at least  $f + 1$  processes saw the same value  $v$  from  $q$  for  $k$ , and (3) a proof that at least  $f + 1$  processes wrote a proof of seeing value  $v$  from  $q$  for  $k$  in their second slot. We call these slots the Value slot, the L1Proof slot, and the L2Proof slot, respectively.

We note that each such valid proof has signed copies of only one value for the message. Any proof that shows copies of two different values or a value that is not signed is not considered valid. If a proof has copies of only value  $v$ , we say that this proof *supports*  $v$ .

To deliver a value  $v$  from process  $q$  with sequence number  $k$ , process  $p$  must successfully write a valid proof-of-proofs in its L2Proof slot supporting value  $v$  (we call this an L2 proof). It has two options of how to do this; firstly, if it sees a valid L2 proof in some other process  $i$ ’s  $L2Proof[i, k, q]$  slot, it copies this proof over to its own L2 proof slot, and can then deliver the value that this proof supports. If  $p$  does not find a valid L2 proof in some other process’s slot, it must try to construct one itself. We now describe how this is done.

A correct process  $p$  goes through three stages when constructing a valid L2 proof for  $(k, m)$  from  $q$ . In the pseudocode, the three stages are denoted using states that  $p$  goes through: `WaitForSender`, `WaitForL1Proof`, and `WaitForL2Proof`.

In the first stage, `WaitForSender`,  $p$  reads  $q$ ’s  $Value[q, k, q]$  slot. If  $p$  finds a  $(k, m)$  pair,  $p$  signs and copies it to its  $Value[p, k, q]$  slot and enters the `WaitForL1Proof` state.

In the second stage, `WaitForL1Proof`,  $p$  reads all  $Value[i, k, q]$  slots, for  $i \in \Pi$ . If all the values  $p$  reads are correctly signed and equal to  $(k, m)$ , and if there are at least  $f + 1$  such values, then  $p$  compiles them into an L1 proof, which it signs and writes to  $L1Proof[p, k, q]$ ;  $p$  then enters the `WaitForL2Proof` state.

In the third stage, `WaitForL2Proof`,  $p$  reads all  $L1Proof[i, k, q]$  slots, for  $i \in \Pi$ . If  $p$  finds at least  $f + 1$  valid and signed L1 proofs for  $(k, m)$ , then  $p$  compiles them into an L2 proof, which it signs and writes to  $L2Proof[p, k, q]$ . The next time that  $p$  scans the  $L2Proof[\cdot, k, q]$  slots,  $p$  will see its own L2 proof (or some other valid proof for  $(k, m)$ ) and deliver  $(k, m)$ .

This three-stage validation process ensures the following crucial property: no two valid L2 proofs can support different values. Intuitively, this property is achieved because for both L1 and L2 proofs, at least  $f + 1$  values of the previous stage must be copied, meaning that at least one correct process was involved in the quorum needed to construct each proof. Because correct processes read the slots of *all* others at each stage before constructing the next proof, and because they never overwrite or delete values that they already wrote, it is guaranteed that no two correct processes will create valid L1 proofs for different values, since one must see the Value slot of the

other. Thus, no two processes, Byzantine or otherwise, can construct valid L2 proofs for different values.

Notably, a weaker version of broadcast, which does not require Property 4 (i.e., Byzantine consistent broadcast [17, Module 3.11]), can be solved with just the first stage of Algorithm 2, without the L1 and L2 proofs. The purpose of those proofs is to ensure the 4th property holds; that is, to enable all correct processes to deliver a value once some correct process delivered.

In the appendix, we formally prove the above intuition and arrive at the following lemma.

LEMMA 4.1. *Reliable broadcast can be solved in shared-memory with SWMR regular registers with  $n \geq 2f + 1$  processes.*

**4.1.2 Applying Clement et al’s Construction.** Clement et al. show that given unforgeable transferable signatures and non-equivocation, one can reduce Byzantine failures to crash failures in message passing systems [21]. They define non-equivocation as a predicate  $valid_p$  for each process  $p$ , which takes a sequence number and a value and evaluates to true for just one value per sequence number. All processes must be able to call the same  $valid_p$  predicate, which always terminates every time it is called.

We now show how to use reliable broadcast to implement messages with transferable signatures and non-equivocation as defined by Clement et al. [21]. Note that our reliable broadcast mechanism already involves the use of transferable signatures, so to send and receive signed messages, one can simply use broadcast and deliver those messages. However, simply using broadcast and deliver is not enough to satisfy the requirements of the  $valid_p$  predicate of Clement et al. The problem occurs when trying to validate nested messages recursively.

In particular, recall that in Clement et al’s construction, whenever a message is sent, the entire history of that process, including all messages it has sent and received, is attached. Consider two Byzantine processes  $q_1$  and  $q_2$ , and assume that  $q_1$  attempts to equivocate in its  $k$ th message, signing both  $(k, m)$  and  $(k, m')$ . Assume therefore that no correct process delivers any message from  $q_1$  in its  $k$ th round. However, since  $q_2$  is also Byzantine, it could claim to have delivered  $(k, m)$  from  $q_1$ . If  $q_2$  then sends a message that includes  $(q, k, m)$  as part of its history, a correct process  $p$  receiving  $q_2$ ’s message must recursively verify the history  $q_2$  sent. To do so,  $p$  can call `try_deliver` on  $(q_1, k)$ . However, since no correct process delivered any message from  $(q_1, k)$ , it is possible that this call never returns.

To solve this issue, we introduce a `validate` operation that can be used along with broadcast and deliver to validate the correctness of a given message. The `validate` operation is very simple: it takes in a process id, a sequence number, and a message value  $m$ , and simply runs the `checkL2proof` helper function. If the function returns a proof supporting  $m$ , `validate` returns true. Otherwise it returns false. The pseudocode is shown in Algorithm 3.

In this way, Algorithms 2 and 3 together provide signed messages and a non-equivocation primitive. Thus, combined with the construction of Clement et al. [21], we immediately get the following result.

THEOREM 4.2. *There exists an algorithm for weak Byzantine agreement in a shared-memory system with SWMR regular registers, signatures, and up to  $f_p$  process crashes where  $n \geq 2f_p + 1$ .*



### Algorithm 3: Validate Operation for Reliable Broadcast

```

1 bool validate(q, k, m) {
2   val = checkL2proof(q, k);
3   if (val == m) {
4     return true; }
5   return false; }

```

*Non-equivocation in our model.* To convert the above algorithm to our model, where memory may fail, we use the ideas in [2, 6, 30] to implement failure-free SWMR regular registers from the fail-prone memory, and then run weak Byzantine agreement using those regular registers. To implement an SWMR register, a process writes or reads all memories, and waits for a majority to respond. When reading, if  $p$  sees exactly one distinct non- $\perp$  value  $v$  across the memories, it returns  $v$ ; otherwise, it returns  $\perp$ .

**DEFINITION 2.** Let  $\mathcal{A}$  be a message-passing algorithm. Robust Backup( $\mathcal{A}$ ) is the algorithm  $\mathcal{A}$  in which all send and receive operations are replaced by T-send and T-receive operations (respectively) implemented with reliable broadcast.

Thus we get the following lemma, from the result of Clement et al. [21], Lemma 4.1, and the above handling of memory failures.

**LEMMA 4.3.** If  $\mathcal{A}$  is a consensus algorithm that is tolerant to  $f$  process crash failures, then Robust Backup( $\mathcal{A}$ ) is a weak Byzantine agreement algorithm that is tolerant to up to  $f_P$  Byzantine processes and  $f_M$  memory crashes, where  $n \geq 2f_P + 1$  and  $m \geq 2f_M + 1$  in the message-and-memory model.

The following theorem is an immediate corollary of the lemma.

**THEOREM 4.4.** There exists an algorithm for Weak Byzantine Agreement in a message-and-memory model with up to  $f_P$  Byzantine processes and  $f_M$  memory crashes, where  $n \geq 2f_P + 1$  and  $m \geq 2f_M + 1$ .

## 4.2 The Cheap Quorum Sub-Algorithm

We now give an algorithm that decides in two delays in common executions in which the system is synchronous and there are no failures. It requires only one signature for a fast decision, whereas the best prior algorithm requires  $6f_P + 2$  signatures and  $n \geq 3f_P + 1$  [7]. Our algorithm, called Cheap Quorum, is not in itself a complete consensus algorithm; it may abort in some executions. If Cheap Quorum aborts, it outputs an *abort value*, which is used to initialize the Robust Backup so that their composition preserves weak Byzantine agreement. This composition is inspired by the Abstract framework of Aublin et al. [7].

The algorithm has a special process  $\ell$ , say  $\ell = p_1$ , which serves both as a *leader* and a *follower*. Other processes act only as *followers*. The memory is partitioned into  $n + 1$  regions denoted  $Region[p]$  for each  $p \in \Pi$ , plus an extra one for  $p_1$ ,  $Region[\ell]$  in which it proposes a value. Initially,  $Region[p]$  is a regular SWMR region where  $p$  is the writer. Unlike in Algorithm 2, some of the permissions are dynamic; processes may remove  $p_1$ 's write permission to  $Region[\ell]$  (i.e., the *legalChange* function returns false to any permission change requests, except for ones revoking  $p_1$ 's permission to write on  $Region[\ell]$ ).

### Algorithm 4: Cheap Quorum normal operation—code for process $p$

```

1 Leader code
2 propose(v) {
3   sign(v);
4   status = Value[ℓ].write(v);
5   if (status == nak) Panic_mode();
6   else decide(v); }

8 Follower code
9 propose(w) {
10  do {v = read(Value[ℓ]);
11     for all q ∈ Π do pan[q] = read(Panic[q]);
12  } until (v ≠ ⊥ || pan[q] == true for some q ||
13         ↪ timeout);
14  if (v ≠ ⊥ && sValid(p1, v)) {
15    sign(v);
16    write(Value[p], v);
17    do {for all q ∈ Π do val[q] = read(Value[q]);
18       if (|{q : val[q] == v}| ≥ n then {
19         Proof[p].write(sign(val[1..n]));
20         for all q ∈ Π do prf[q] = read(Proof[q]);
21         if (|{q : verifyProof(prf[q]) == true}| ≥
22             ↪ n { decide(v); exit; } } }
23    for all q ∈ Π do pan[q] = read(Panic[q]);
24    } until (pan[q] == true for some q || timeout);
25    }
26  Panic_mode(); }

```

### Algorithm 5: Cheap Quorum panic mode—code for process $p$

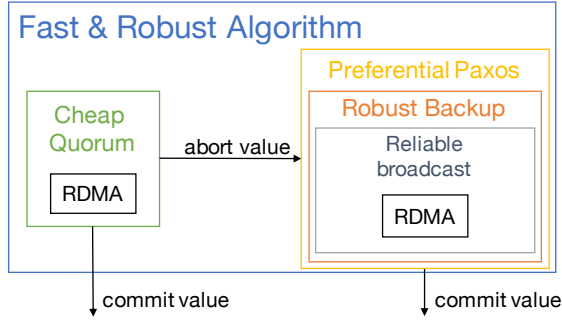
```

1 panic_mode() {
2   Panic[p] = true;
3   changePermission(Region[ℓ], R: Π, W: {}, RW: {});
4   ↪ // remove write permission
5   v = read(Value[p]);
6   prf = read(Proof[p]);
7   if (v ≠ ⊥) { Abort with ⟨v, prf⟩; return; }
8   LVal = read(Value[ℓ]);
9   if (LVal ≠ ⊥) { Abort with ⟨LVal, ⊥⟩; return; }
10  Abort with ⟨myInput, ⊥⟩; }

```

Processes initially execute under a *normal* mode in common-case executions, but may switch to *panic* mode if they intend to abort, as in [7]. The pseudo-code of the normal mode is in Algorithm 4.  $Region[p]$  contains three registers  $Value[p]$ ,  $Panic[p]$ ,  $Proof[p]$  initially set to  $\perp$ , *false*,  $\perp$ . To propose  $v$ , the leader  $p_1$  signs  $v$  and writes it to  $Value[\ell]$ . If the write is successful (it may fail because its write permission was removed), then  $p_1$  decides  $v$ ; otherwise  $p_1$  calls *Panic\_mode*(). Note that all processes, including  $p_1$ , continue their execution after deciding. However,  $p_1$  never decides again if it decided as the leader. A follower  $q$  checks if  $p_1$  wrote to  $Value[\ell]$  and, if so, whether the value is properly signed. If so,  $q$  signs  $v$ , writes it to  $Value[q]$ , and waits for other processes to write the same value to  $Value[*]$ . If  $q$  sees  $2f + 1$  copies of  $v$  signed by different processes,  $q$  assembles these copies in a *unanimity proof*, which it signs and writes to  $Proof[q]$ .  $q$  then waits for  $2f + 1$  unanimity proofs for  $v$  to appear in  $Proof[*]$ , and checks that they are valid, in which case  $q$  decides  $v$ . This waiting continues until a timeout expires<sup>4</sup>, at which time  $q$  calls *Panic\_mode*(). In *Panic\_mode*(), a

<sup>4</sup>The timeout is chosen to be an upper bound on the communication, processing and computation delays in the common case.



**Figure 6: Interactions of the components of the Fast & Robust Algorithm.**

process  $p$  sets  $\text{Panic}[p]$  to *true* to tell other processes it is panicking; other processes periodically check to see if they should panic too.  $p$  then removes write permission from  $\text{Region}[\ell]$ , and decides on a value to abort: either  $\text{Value}[p]$  if it is non- $\perp$ ,  $\text{Value}[\ell]$  if it is non- $\perp$ , or  $p$ 's input value. If  $p$  has a unanimity proof in  $\text{Proof}[p]$ , it adds it to the abort value.

In Appendix B, we prove the correctness of Cheap Quorum, and in particular we show the following two important agreement properties:

**LEMMA 4.5 (CHEAP QUORUM DECISION AGREEMENT).** *Let  $p$  and  $q$  be correct processes. If  $p$  decides  $v_1$  while  $q$  decides  $v_2$ , then  $v_1 = v_2$ .*

**LEMMA 4.6 (CHEAP QUORUM ABORT AGREEMENT).** *Let  $p$  and  $q$  be correct processes (possibly identical). If  $p$  decides  $v$  in Cheap Quorum while  $q$  aborts from Cheap Quorum, then  $v$  will be  $q$ 's abort value. Furthermore, if  $p$  is a follower,  $q$ 's abort proof is a correct unanimity proof.*

The above construction assumes a fail-free memory with regular registers, but we can extend it to tolerate memory failures using the approach of Section 4.1, noting that each register has a single writer process.

### 4.3 Putting it Together: the Fast & Robust Algorithm

The final algorithm, called Fast & Robust, combines Cheap Quorum (§4.2) and Robust Backup (§4.1), as we now explain. Recall that Robust Backup is parameterized by a message-passing consensus algorithm  $\mathcal{A}$  that tolerates crash-failures.  $\mathcal{A}$  can be any such algorithm (e.g., Paxos).

Roughly, in Fast & Robust, we run Cheap Quorum and, if it aborts, we use a process's abort value as its input value to Robust Backup. However, we must carefully glue the two algorithms together to ensure that if some correct process decided  $v$  in Cheap Quorum, then  $v$  is the only value that can be decided in Robust Backup.

For this purpose, we propose a simple wrapper for Robust Backup, called *Preferential Paxos*. Preferential Paxos first runs a set-up phase, in which processes may adopt new values, and then runs Robust Backup with the new values. More specifically, there are some *preferred* input values  $v_1 \dots v_k$ , ordered by priority. We guarantee that

every process adopts one of the top  $f + 1$  priority inputs. In particular, this means that if a majority of processes get the highest priority value,  $v_1$ , as input, then  $v_1$  is guaranteed to be the decision value. The set-up phase is simple; all processes send each other their input values. Each process  $p$  waits to receive  $n - f$  such messages, and adopts the value with the highest priority that it sees. This is the value that  $p$  uses as its input to Paxos. The pseudocode for Preferential Paxos is given in Algorithm 8 in Appendix C, where we also prove the following lemma about Preferential Paxos:

**LEMMA 4.7 (PREFERENTIAL PAXOS PRIORITY DECISION).** *Preferential Paxos implements weak Byzantine agreement with  $n \geq 2f_P + 1$  processes. Furthermore, let  $v_1, \dots, v_n$  be the input values of an instance  $C$  of Preferential Paxos, ordered by priority. The decision value of correct processes is always one of  $v_1, \dots, v_{f+1}$ .*

We can now describe Fast & Robust in detail. We start executing Cheap Quorum. If Cheap Quorum aborts, we execute Preferential Paxos, with each process receiving its abort value from Cheap Quorum as its input value to Preferential Paxos. We define the priorities of inputs to Preferential Paxos as follows.

**DEFINITION 3 (INPUT PRIORITIES FOR PREFERENTIAL PAXOS).** *The input values for Preferential Paxos as it is used in Fast & Robust are split into three sets (here,  $p_1$  is the leader of Cheap Quorum):*

- $T = \{v \mid v \text{ contains a correct unanimity proof}\}$
- $M = \{v \mid v \notin T \wedge v \text{ contains the signature of } p_1\}$
- $B = \{v \mid v \notin T \wedge v \notin M\}$

*The priority order of the input values is such that for all values  $v_T \in T$ ,  $v_M \in M$ , and  $v_B \in B$ ,  $\text{priority}(v_T) > \text{priority}(v_M) > \text{priority}(v_B)$ .*

Figure 6 shows how the various algorithms presented in this section come together to form the Fast & Robust algorithm. In Appendices B and C, we show that Fast & Robust is correct, with the following key lemma:

**LEMMA 4.8 (COMPOSITION LEMMA).** *If some correct process decides a value  $v$  in Cheap Quorum before an abort, then  $v$  is the only value that can be decided in Preferential Paxos with priorities as defined in Definition 3.*

**THEOREM 4.9.** *There exists a 2-deciding algorithm for Weak Byzantine Agreement in a message-and-memory model with up to  $f_P$  Byzantine processes and  $f_M$  memory crashes, where  $n \geq 2f_P + 1$  and  $m \geq 2f_M + 1$ .*

## 5 CRASH FAILURES

We now restrict ourselves to crash failures of processes and memories. Clearly, we can use the algorithms of Section 4 in this setting, to obtain a 2-deciding consensus algorithm with  $n \geq 2f_P + 1$  and  $m \geq 2f_M + 1$ . However, this is overkill since those algorithms use sophisticated mechanisms (signatures, non-equivocation) to guard against Byzantine behavior. With only crash failures, we now show it is possible to retain the efficiency of a 2-deciding algorithm while improving resiliency. In Section 5.1, we first give a 2-deciding algorithm that allows the crash of all but one process ( $n \geq f_P + 1$ ) and a minority of memories ( $m \geq 2f_M + 1$ ). In Section 5.2, we improve resiliency further by giving a 2-deciding algorithm that



### Algorithm 7: Protected Memory Paxos—code for process $p$

```

1  Registers: for  $i=1..m$ ,  $p \in \Pi$ ,
2  slot[i,p]: tuple (minProp, accProp, value) // in
   ↪ memory i
3   $\Omega$ : failure detector that returns current leader

5  startPhase2(i) {
6  add i to ListOfReady processes;
7  while (size(ListOfReady) < majority of memories) {}
8  Phase2Started = true; }

10 propose(v) {
11 repeat forever {
12 wait until  $\Omega == p$ ; // wait to become leader
13 propNr = a higher value than any proposal number
   ↪ seen before;
14 CurrentVal = v;
15 CurrentMaxProp = 0;
16 Phase2Started = false;
17 ListOfReady = 0;
18 for every memory i in parallel {
19 if ( $p \neq p_1$  || not first attempt) {
20 getPermission(i);
21 success = write(slot[i,p], (propNr, 1, 1));
22 if (not success) { abort(); }
23 vals = read all slots from i;
24 if (vals contains a non-null value) {
25 val = v  $\in$  vals with highest propNr;
26 if (val.propNr > propNr) { abort(); }
27 atomic {
28 if (val.propNr > CurrentMaxProp) {
29 if (Phase2Started) { abort(); }
30 CurrentVal = val.value;
31 CurrentMaxProp = val.propNr;
32 } } }
33 startPhase2(i); }
34 // done phase 1 or ( $p == p_1$  &&  $p_1$ 's first
   ↪ attempt)
35 success = write(slot[i,p], (propNr, propNr,
   ↪ CurrentVal));
36 if (not success) { abort(); }
37 } until this has been done at a majority of the
   ↪ memories, or until 'abort' has been
   ↪ called
38 if (loop completed without abort) {
39 decide CurrentVal; } } }
```

tolerates crashes of a minority of the combined set of memories and processes.

## 5.1 Protected Memory Paxos

Our starting point is the Disk Paxos algorithm [29], which works in a system with processes and memories where  $n \geq f_p + 1$  and  $m \geq 2f_M + 1$ . This is our resiliency goal, but Disk Paxos takes four delays in common executions. Our new algorithm, called Protected Memory Paxos, removes two delays; it retains the structure of Disk Paxos but uses permissions to skip steps. Initially some fixed leader  $\ell = p_1$  has exclusive write permission to all memories; if another process becomes leader, it takes the exclusive permission. Having exclusive permission permits a leader  $\ell$  to optimize execution, because  $\ell$  can do two things simultaneously: (1) write its consensus proposal and (2) determine whether another leader took over. Specifically, if  $\ell$  succeeds in (1), it knows no leader  $\ell'$  took over because  $\ell'$  would have taken the permission. Thus  $\ell$  avoids the last read in Disk Paxos, saving two delays. Of course, care must be taken to implement this without violating safety.

The pseudocode of Protected Memory Paxos is in Algorithm 7. Each memory has one memory region, and at any time exactly one process can write to the region. Each memory  $i$  holds a register  $slot[i,p]$  for each process  $p$ . Intuitively,  $slot[i,p]$  is intended for  $p$  to write, but  $p$  may not have write permission to do that if it is not the leader—in that case, no process writes  $slot[i,p]$ .

When a process  $p$  becomes the leader, it must execute a sequence of steps on a majority of the memories to successfully commit a value. It is important that  $p$  execute *all* of these steps on each of the memories that counts toward its majority; otherwise two leaders could miss each other's values and commit conflicting values. We therefore present the pseudocode for this algorithm in a *parallel-for* loop (lines 18–37), with one thread per memory that  $p$  accesses. The algorithm has two phases similar to Paxos, where the second phase may only begin after the first phase has been completed for a majority of the memories. We represent this in the code with a barrier that waits for a majority of the threads.

When a process  $p$  becomes leader, it executes the prepare phase (the first leader  $p_1$  can skip this phase in its first execution of the loop), where, for each memory,  $p$  attempts to (1) acquire exclusive write permission, (2) write a new proposal number in its slot, and (3) read all slots of that memory.  $p$  waits to succeed in executing these steps on a majority of the memories. If any of  $p$ 's writes fail or  $p$  finds a proposal with a higher proposal number, then  $p$  gives up. This is represented with an `abort` in the pseudocode; when an `abort` is executed, the `for` loop terminates. We assume that when the `for` loop terminates—either because some thread has aborted or because a majority of threads have reached the end of the loop—all threads of the `for` loop are terminated and control returns to the main loop (lines 11–37).

If  $p$  does not abort, it adopts the value with highest proposal number of all those it read in the memories. To make it clear that races should be avoided among parallel threads in the pseudocode, we wrap this part in an `atomic` environment.

In the next phase, each of  $p$ 's threads writes its value to its slot on its memory. If a write fails,  $p$  gives up. If  $p$  succeeds, this is where we optimize time:  $p$  can simply decide, whereas Disk Paxos must read the memories again.

Note that it is possible that some of the memories that made up the majority that passed the initial barrier may crash later on. To prevent  $p$  from stalling forever in such a situation, it is important that straggler threads that complete phase 1 later on be allowed to participate in phase 2. However, if such a straggler thread observes a more up-to-date value in its memory than the one adopted by  $p$  for phase 2, this must be taken into account. In this case, to avoid inconsistencies,  $p$  must abort its current attempt and restart the loop from scratch.

The code ensures that some correct process eventually decides, but it is easy to extend it so all correct processes decide [18], by having a decided process broadcast its decision. Also, the code shows one instance of consensus, with  $p_1$  as initial leader. With many consensus instances, the leader terminates one instance and becomes the default leader in the next.

**THEOREM 5.1.** *Consider a message-and-memory model with up to  $f_p$  process crashes and  $f_M$  memory crashes, where  $n \geq f_p + 1$  and  $m \geq 2f_M + 1$ . There exists a 2-deciding algorithm for consensus.*

## 5.2 Aligned Paxos

We now further enhance the failure resilience. We show that memories and processes are equivalent *agents*, in that it suffices for a majority of the agents (processes and memories together) to remain alive to solve consensus. Our new algorithm, *Aligned Paxos*, achieves this resiliency. To do so, the algorithm relies on the ability to use both the messages and the memories in our model; permissions are not needed. The key idea is to align a message-passing algorithm and a memory-based algorithm to use any majority of agents. We align Paxos [38] and Protected Memory Paxos so that their decisions are coordinated. More specifically, Protected Memory Paxos and Paxos have two phases. To align these algorithms, we factor out their differences and replace their steps with an abstraction that is implemented differently for each algorithm. The result is our *Aligned Paxos* algorithm, which has two phases, each with three steps: *communicate*, *hear back*, and *analyze*. Each step treats processes and memories separately, and translates the results of operations on different agents to a common language. We implement the steps using their analogues in Paxos and Protected Memory Paxos<sup>5</sup>. The pseudocode of Aligned Paxos is shown in Appendix E.

## 6 DYNAMIC PERMISSIONS ARE NECESSARY FOR EFFICIENT CONSENSUS

In §5.1, we showed how dynamic permissions can improve the performance of Disk Paxos. Are dynamic permissions necessary? We prove that with shared memory (or disks) alone, one cannot achieve 2-deciding consensus, even if the memory never fails, it has static permissions, processes may only fail by crashing, and the system is partially synchronous in the sense that eventually there is a known upper bound on the time it takes a correct process to take a step [27]. This result applies a fortiori to the Disk Paxos model [29].

**THEOREM 6.1.** *Consider a partially synchronous shared-memory model with registers, where registers can have arbitrary static permissions, memory never fails, and at most one process may fail by crashing. No consensus algorithm is 2-deciding.*

**PROOF.** Assume by contradiction that  $A$  is an algorithm in the stated model that is 2-deciding. That is, there is some execution  $E$  of  $A$  in which some process  $p$  decides a value  $v$  with 2 delays. We denote by  $R$  and  $W$  the set of objects which  $p$  reads and writes in  $E$  respectively. Note that since  $p$  decides in 2 delays in  $E$ ,  $R$  and  $W$  must be disjoint, by the definition of operation delay and the fact that a process has at most one outstanding operation per object. Furthermore,  $p$  must issue all of its read and writes without waiting for the response of any operation.

Consider an execution  $E'$  in which  $p$  reads from the same set  $R$  of objects and writes the same values as in  $E$  to the same set  $W$  of objects. All of the read operations that  $p$  issues return by some time  $t_0$ , but the write operations of  $p$  are delayed for a long time. Another process  $p'$  begins its proposal of a value  $v' \neq v$  after  $t_0$ . Since no process other than  $p'$  writes to any objects,  $E'$  is indistinguishable to  $p'$  from an execution in which it runs alone. Since  $A$  is a correct consensus algorithm that terminates if there is no contention,  $p'$  must

eventually decide value  $v'$ . Let  $t'$  be the time at which  $p'$  decides. All of  $p$ 's write operations terminate and are linearized in  $E'$  after time  $t'$ . Execution  $E'$  is indistinguishable to  $p$  from execution  $E$ , in which it ran alone. Therefore,  $p$  decides  $v \neq v'$ , violating agreement.  $\square$

Theorem 6.1, together with the Fast Paxos algorithm of Lamport [39], shows that an atomic read-write shared memory model is strictly weaker than the message passing model in its ability to solve consensus quickly. This result may be of independent interest, since often the classic shared memory and message passing models are seen as equivalent, because of the seminal computational equivalence result of Attiya, Bar-Noy, and Dolev [6]. Interestingly, it is known that shared memory can tolerate more failures when solving consensus (with randomization or partial synchrony) [5, 15], and therefore it seems that perhaps shared memory is strictly stronger than message passing for solving consensus. However, our result shows that there are aspects in which message passing is stronger than shared memory. In particular, message passing can solve consensus faster than shared memory in well-behaved executions.

## 7 RDMA IN PRACTICE

Our model is meant to reflect capabilities of RDMA, while providing a clean abstraction to reason about. We now give an overview of how RDMA works, and how features of our model can be implemented using RDMA.

RDMA enables a remote process to access local memory directly through the network interface card (NIC), without involving the CPU. For a piece of local memory to be accessible to a remote process  $p$ , the CPU has to *register* that memory region and associate it with the appropriate connection (called *Queue Pair*) for  $p$ . The association of a registered memory region and a queue pair is done indirectly through a *protection domain*: both memory regions and queue pairs are associated with a protection domain, and a queue pair  $q$  can be used to access a memory region  $r$  if  $q$  and  $r$  are in the same protection domain. The CPU must also specify what access level (read, write, read-write) is allowed to the memory region in each protection domain. A local memory area can thus be registered and associated with several queue pairs, with the same or different access levels, by associating it with one or more protection domains. Each RDMA connection can be used by the remote server to access registered memory regions using a unique region-specific key created as a part of the registration process.

As highlighted by previous work [45], failures of the CPU, NIC and DRAM can be seen as independent (e.g., arbitrary delays, too many bit errors, failed ECC checks, respectively). For instance, *zombie servers* in which the CPU is blocked but RDMA requests can still be served account for roughly half of all failures [45]. This motivates our choice to treat processes and memory separately in our model. In practice, if a CPU fails permanently, the memory will also become unreachable through RDMA eventually; however, in such cases memory may remain available long enough for ongoing operations to complete. Also, in practical settings it is possible for full-system crashes to occur (e.g., machine restarts), which correspond to a process and a memory failing at the same time—this is allowed by our model.

<sup>5</sup>We believe other implementations are possible. For example, replacing the Protected Memory Paxos implementation for memories with the Disk Paxos implementation yields an algorithm that does not use permissions.

Memory regions in our model correspond to RDMA memory regions. Static permissions can be implemented by making the appropriate memory region registration before the execution of the algorithm; these permissions then persist during execution without CPU involvement. Dynamic permissions require the host CPU to change the access levels; this should be done in the OS kernel: the kernel creates regions and controls their permissions, and then shares memory with user-space processes. In this way, Byzantine processes cannot change permissions illegally. The assumption is that the kernel is not Byzantine. Alternatively, future hardware support similar to SGX could even allow parts of the kernel to be Byzantine.

Using RDMA, a process  $p$  can grant permissions to a remote process  $q$  by registering memory regions with the appropriate access permissions (read, write, or read/write) and sending the corresponding key to  $q$ .  $p$  can revoke permissions dynamically by simply deregistering the memory region.

For our reliable broadcast algorithm, each process can register the two dimensional array of values in read-only mode with a protection domain. All the queue pairs used by that process are also created in the context of the same protection domain. Additionally, the process can preserve write access permission to its row via another registration of just that row with the protection domain, thus enabling single-writer multiple-reader access. Thereafter the reliable broadcast algorithm can be implemented trivially by using RDMA reads and writes by all processes. Reliable broadcast with unreliable memories is similarly straightforward since failure of the memory ensures that no process will be able to access the memory.

For Cheap Quorum, the static memory region registrations are straightforward as above. To revoke the leader's write permission, it suffices for a region's host process to deregister the memory region. Panic messages can be relayed using RDMA message sends.

In our crash-only consensus algorithm, we leverage the capability of registering overlapping memory regions in a protection domain. As in above algorithms, each process uses one protection domain for RDMA accesses. Queue pairs for connections to all other processes are associated with this protection domain. The process' entire slot array is registered with the protection domain in read-only mode. In addition, the same slot array can be dynamically registered (and deregistered) in write mode based on incoming write permission requests: A proposer requests write permission using an RDMA message send. In response, the acceptor first deregisters write permission for the immediate previous proposer. The acceptor thereafter registers the slot array in write mode and responds to the proposer with the new key associated with the newly registered slot array. Reads of the slot array are performed by the proposer using RDMA reads. Subsequent second phase RDMA write of the value can be performed on the slot array as long as the proposer continues to have write permission to the slot array. The RDMA write fails if the acceptor granted write permission to another proposer in the meantime.

**Acknowledgements.** We wish to thank the anonymous reviewers for their helpful comments on improving the paper. This work has been supported in part by the European Research Council (ERC) Grant 339539 (AOC), and a Microsoft PhD Fellowship.

## REFERENCES

- [1] Ittai Abraham, Gregory Chockler, Idit Keidar, and Dahlia Malkhi. Byzantine disk paxos: optimal resilience with byzantine shared memory. *Distributed computing (DIST)*, 18(5):387–408, 2006.
- [2] Yehuda Afek, David S. Greenberg, Michael Merritt, and Gadi Taubenfeld. Computing with faulty shared memory. In *ACM Symposium on Principles of Distributed Computing (PODC)*, pages 47–58, August 1992.
- [3] Marcos K Aguilera, Naama Ben-David, Irina Calciu, Rachid Guerraoui, Erez Petrank, and Sam Toueg. Passing messages while sharing memory. In *ACM Symposium on Principles of Distributed Computing (PODC)*, pages 51–60. ACM, 2018.
- [4] Noga Alon, Michael Merritt, Omer Reingold, Gadi Taubenfeld, and Rebecca N Wright. Tight bounds for shared memory systems accessed by byzantine processes. *Distributed computing (DIST)*, 18(2):99–109, 2005.
- [5] James Aspnes and Maurice Herlihy. Fast randomized consensus using shared memory. *Journal of algorithms*, 11(3):441–461, 1990.
- [6] Hagit Attiya, Amotz Bar-Noy, and Danny Dolev. Sharing memory robustly in message-passing systems. *Journal of the ACM (JACM)*, 42(1):124–142, 1995.
- [7] Pierre-Louis Aublin, Rachid Guerraoui, Nikola Knežević, Vivien Quéma, and Marko Vukolić. The next 700 BFT protocols. *ACM Transactions on Computer Systems (TOCS)*, 32(4):12, 2015.
- [8] Rida Bazzi and Gil Neiger. Optimally simulating crash failures in a byzantine environment. In *International Workshop on Distributed Algorithms (WDAG)*, pages 108–128. Springer, 1991.
- [9] Jonathan Behrens, Ken Birman, Sagar Jha, Matthew Milano, Edward Tremel, Eugene Bagdasaryan, Theo Gkountouvas, Weijia Song, and Robbert Van Renesse. Derecho: Group communication at the speed of light. Technical report, Technical Report. Cornell University, 2016.
- [10] Michael Ben-Or. Another advantage of free choice (extended abstract): Completely asynchronous agreement protocols. In *ACM Symposium on Principles of Distributed Computing (PODC)*, pages 27–30. ACM, 1983.
- [11] Alysson Neves Bessani, Miguel Correia, Joni da Silva Fraga, and Lau Cheuk Lung. Sharing memory between byzantine processes using policy-enforced tuple spaces. *IEEE Transactions on Parallel and Distributed Systems*, 20(3):419–432, 2009.
- [12] Romain Boichat, Partha Dutta, Svend Frolund, and Rachid Guerraoui. Reconstructing paxos. *SIGACT News*, 34(2):42–57, March 2003.
- [13] Zohir Bouzid, Damien Imbs, and Michel Raynal. A necessary condition for byzantine k-set agreement. *Information Processing Letters*, 116(12):757–759, 2016.
- [14] Gabriel Bracha. Asynchronous byzantine agreement protocols. *Information and Computation*, 75(2):130–143, 1987.
- [15] Gabriel Bracha and Sam Toueg. Asynchronous consensus and broadcast protocols. *Journal of the ACM (JACM)*, 32(4):824–840, 1985.
- [16] Francisco Brasileiro, Fabíola Greve, Achour Mostéfaoui, and Michel Raynal. Consensus in one communication step. In *International Conference on Parallel Computing Technologies*, pages 42–50. Springer, 2001.
- [17] Christian Cachin, Rachid Guerraoui, and Luís E. T. Rodrigues. *Introduction to Reliable and Secure Distributed Programming* (2. ed.). Springer, 2011.
- [18] Tushar Deepak Chandra and Sam Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM (JACM)*, 43(2):225–267, 1996.
- [19] Byung-Gon Chun, Petros Maniatis, and Scott Shenker. Diverse replication for single-machine byzantine-fault tolerance. In *USENIX Annual Technical Conference (ATC)*, pages 287–292, 2008.
- [20] Byung-Gon Chun, Petros Maniatis, Scott Shenker, and John Kubiawicz. Attested append-only memory: making adversaries stick to their word. In *ACM Symposium on Operating Systems Principles (SOSP)*, pages 189–204, 2007.
- [21] Allen Clement, Flavio Junqueira, Aniket Kate, and Rodrigo Rodrigues. On the (limited) power of non-equivocation. In *ACM Symposium on Principles of Distributed Computing (PODC)*, pages 301–308. ACM, 2012.
- [22] Miguel Correia, Nuno Ferreira Neves, and Paulo Veríssimo. How to tolerate half less one byzantine nodes in practical distributed systems. In *International Symposium on Reliable Distributed Systems (SRDS)*, pages 174–183, 2004.
- [23] Miguel Correia, Giuliana S Veronese, and Lau Cheuk Lung. Asynchronous byzantine consensus with  $2f+1$  processes. In *ACM symposium on applied computing (SAC)*, pages 475–480. ACM, 2010.
- [24] Dan Dobre and Neeraj Suri. One-step consensus with zero-degradation. In *International Conference on Dependable Systems and Networks (DSN)*, pages 137–146. IEEE Computer Society, 2006.
- [25] Aleksandar Dragojević, Dushyanth Narayanan, Miguel Castro, and Orion Hodson. FaRM: Fast remote memory. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 401–414, 2014.
- [26] Partha Dutta, Rachid Guerraoui, and Leslie Lamport. How fast can eventual synchrony lead to consensus? In *International Conference on Dependable Systems and Networks (DSN)*, pages 22–27. IEEE, 2005.
- [27] Cynthia Dwork, Nancy Lynch, and Larry Stockmeyer. Consensus in the presence of partial synchrony. *Journal of the ACM (JACM)*, 35(2):288–323, 1988.

- [28] Michael J Fischer, Nancy A Lynch, and Michael S Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM (JACM)*, 1985.
- [29] Eli Gafni and Leslie Lamport. Disk paxos. *Distributed computing (DIST)*, 16(1):1–20, 2003.
- [30] Prasad Jayanti, Tushar Deepak Chandra, and Sam Toueg. Fault-tolerant wait-free shared objects. *Journal of the ACM (JACM)*, 45(3):451–500, May 1998.
- [31] Anuj Kalia, Michael Kaminsky, and David G Andersen. Using RDMA efficiently for key-value services. *ACM SIGCOMM Computer Communication Review*, 44(4):295–306, 2015.
- [32] Anuj Kalia, Michael Kaminsky, and David G Andersen. Design guidelines for high performance RDMA systems. In *USENIX Annual Technical Conference (ATC)*, page 437, 2016.
- [33] Anuj Kalia, Michael Kaminsky, and David G Andersen. FaSST: Fast, scalable and simple distributed transactions with two-sided (RDMA) datagram RPCs. In *USENIX Symposium on Operating System Design and Implementation (OSDI)*, volume 16, pages 185–201, 2016.
- [34] Rüdiger Kapitza, Johannes Behl, Christian Cachin, Tobias Distler, Simon Kuhnle, Seyed Vahid Mohammadi, Wolfgang Schröder-Preikschat, and Klaus Stengel. Cheapbft: resource-efficient byzantine fault tolerance. In *European Conference on Computer Systems (EuroSys)*, pages 295–308, 2012.
- [35] Idit Keidar and Sergio Rajsbaum. On the cost of fault-tolerant consensus when there are no faults: preliminary version. *ACM SIGACT News*, 32(2):45–63, 2001.
- [36] Klaus Kursawe. Optimistic byzantine agreement. In *International Symposium on Reliable Distributed Systems (SRDS)*, pages 262–267, October 2002.
- [37] Leslie Lamport. The weak byzantine generals problem. *Journal of the ACM (JACM)*, 30(3):668–676, July 1983.
- [38] Leslie Lamport. The part-time parliament. *ACM Transactions on Computer Systems (TOCS)*, 16(2):133–169, 1998.
- [39] Leslie Lamport. Fast paxos. *Distributed computing (DIST)*, 19(2):79–103, 2006.
- [40] Leslie Lamport, Robert Shostak, and Marshall Pease. The byzantine generals problem. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 4(3):382–401, 1982.
- [41] Dahlia Malkhi, Michael Merritt, Michael K Reiter, and Gadi Taubenfeld. Objects shared by byzantine processes. *Distributed computing (DIST)*, 16(1):37–48, 2003.
- [42] Jean-Philippe Martin and Lorenzo Alvisi. Fast byzantine consensus. *IEEE Transactions on Dependable and Secure Computing (TDSC)*, 3(3):202–215, July 2006.
- [43] Gil Neiger and Sam Toueg. Automatically increasing the fault-tolerance of distributed algorithms. *Journal of Algorithms*, 11(3):374–419, 1990.
- [44] Marshall Pease, Robert Shostak, and Leslie Lamport. Reaching agreement in the presence of faults. *Journal of the ACM (JACM)*, 27(2):228–234, 1980.
- [45] Marius Poke and Torsten Hoefler. DARE: High-performance state machine replication on RDMA networks. In *Symposium on High-Performance Parallel and Distributed Computing (HPDC)*, pages 107–118. ACM, 2015.
- [46] Signe Rüschi, Ines Messadi, and Rüdiger Kapitza. Towards low-latency byzantine agreement protocols using RDMA. In *IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*, pages 146–151. IEEE, 2018.
- [47] Yee Jiun Song and Robbert van Renesse. Bosco: One-step byzantine asynchronous consensus. In *International Symposium on Distributed Computing (DISC)*, pages 438–450, September 2008.
- [48] Giuliana Santos Veronese, Miguel Correia, Alysson Neves Bessani, Lau Cheuk Lung, and Paulo Verissimo. Efficient byzantine fault-tolerance. *IEEE Trans. Computers*, 62(1):16–30, 2013.
- [49] Cheng Wang, Jianyu Jiang, Xusheng Chen, Ning Yi, and Heming Cui. APUS: Fast and scalable paxos on RDMA. In *Symposium on Cloud Computing (SoCC)*, pages 94–107. ACM, 2017.

## A CORRECTNESS OF RELIABLE BROADCAST

**OBSERVATION 1.** *In Algorithm 2, if  $p$  is a correct process, then no slot that belongs to  $p$  is written to more than once.*

**PROOF.** Since  $p$  is correct,  $p$  never writes on any slot more than once. Furthermore, since all slots are single-writer registers, no other process can write on these slots.  $\square$

**OBSERVATION 2.** *In Algorithm 2, correct processes invoke and return from  $\text{try\_deliver}(q)$  infinitely often, for all  $q \in \Pi$ .*

**PROOF.** The  $\text{try\_deliver}()$  function does not contain any blocking steps, loops or goto statements. Thus, if a correct process invokes  $\text{try\_deliver}()$ , it will eventually return. Therefore, for a fixed  $q$  the

infinite loop at line 24 will invoke and return  $\text{try\_deliver}(q)$  infinitely often. Since the parallel for loop at line 22 performs the infinite loop in parallel for each  $q \in \Pi$ , the Observation holds.  $\square$

**PROOF OF LEMMA 4.1.** We prove the lemma by showing that Algorithm 2 correctly implements reliable broadcast. That is, we need to show that Algorithm 2 satisfies the four properties of reliable broadcast.

**Property 1.** Let  $p$  be a correct process that broadcasts  $(k, m)$ . We show that all correct processes eventually deliver  $(k, m)$  from  $p$ . Assume by contradiction that there exists some correct process  $q$  which does not deliver  $(k, m)$ . Furthermore, assume without loss of generality that  $k$  is the smallest key for which Property 1 is broken. That is, all correct processes must eventually deliver all messages  $(k', m')$  from  $p$ , for  $k' < k$ . Thus, all correct processes must eventually increment  $\text{last}[p]$  to  $k$ .

We consider two cases, depending on whether or not some process eventually writes an L2 proof for some  $(k, m')$  message from  $p$  in its L2Proof slot.

First consider the case where no process ever writes an L2 proof of any value  $(k, m')$  from  $p$ . Since  $p$  is correct, upon broadcasting  $(k, m)$ ,  $p$  must sign and write  $(k, m)$  into  $\text{Values}[p, k, p]$  at line 12. By Observation 1,  $(k, m)$  will remain in that slot forever. Because of this, and because there is no L2 proof, all correct processes, after reaching  $\text{last}[p] = k$ , will eventually read  $(k, m)$  in line 37, write it into their own Value slot in line 40 and change their state to WaitForL1Proof.

Furthermore, since  $p$  is correct and we assume signatures are unforgeable, no process  $q$  can write any other valid value  $(k', m') \neq (k, m)$  into its  $\text{Values}[q, k, p]$  slot. Thus, eventually each correct process  $q$  will add at least  $f + 1$  copies of  $(k, m)$  to its checkedVals, write an L1proof consisting of these values into  $\text{L1Proof}[q, k, p]$  in line 52, and change their state to WaitForL2Proof.

Therefore, all correct processes will eventually read at least  $f + 1$  valid L1 Proofs for  $(k, m)$  in line 58 and construct and write valid L2 proofs for  $(k, m)$ . This contradicts the assumption that no L2 proof ever gets written.

In the case where there is some L2 proof, by the argument above, the only value it can prove is  $(k, m)$ . Therefore, all correct processes will see at least one valid L2 proof at deliver. This contradicts our assumption that  $q$  is correct but does not deliver  $(k, m)$  from  $p$ .

**Property 2.** We now prove the second property of reliable broadcast. Let  $p$  and  $p'$  be any two correct processes, and  $q$  be some process, such that  $p$  delivers  $(k, m)$  from  $q$  and  $p'$  delivers  $(k, m')$  from  $q$ . Assume by contradiction that  $m \neq m'$ .

Since  $p$  and  $p'$  are correct, they must have seen valid L2 proofs at line 16 before delivering  $(k, m)$  and  $(k, m')$  respectively. Let  $\mathcal{P}$  and  $\mathcal{P}'$  be those valid proofs for  $(k, m)$  and  $(k, m')$  respectively.  $\mathcal{P}$  (resp.  $\mathcal{P}'$ ) consists of at least  $f + 1$  valid L1 proofs; therefore, at least one of those proofs was created by some correct process  $r$  (resp.  $r'$ ). Since  $r$  (resp.  $r'$ ) is correct, it must have written  $(k, m)$  (resp.  $(k, m')$ ) to its Values slot in line 40. Note that after copying a value to their slot, in the WaitForL1Proof state, correct processes read all Value slots line 46. Thus, both  $r$  and  $r'$  read all Value slots before compiling their L1 proof for  $(k, m)$  (resp.  $(k, m')$ ).

Assume without loss of generality that  $r$  wrote  $(k, m)$  before  $r'$  wrote  $(k, m')$ ; by Observation 1, it must then be the case that  $r'$  later

saw both  $(k, m)$  and  $(k, m')$  when it read all Values slots (line 46). Since  $r'$  is correct, it cannot have then compiled an L1 proof for  $(k, m')$  (the check at line 50 failed). We have reached a contradiction.

**Property 3.** We show that if a correct process  $p$  delivers  $(k, m)$  from a correct process  $p'$ , then  $p'$  broadcast  $(k, m)$ . Correct processes only deliver values for which a valid L2 proof exists (lines 16—30). Therefore,  $p$  must have seen a valid L2 proof  $\mathcal{P}$  for  $(k, m)$ .  $\mathcal{P}$  consists of at least  $f + 1$  L1 proofs for  $(k, m)$  and each L1 proof consists of at least  $f + 1$  matching copies of  $(k, m)$ , signed by  $p'$ . Since  $p'$  is correct and we assume signatures are unforgeable,  $p'$  must have broadcast  $(k, m)$  (otherwise  $p'$  would not have attached its signature to  $(k, m)$ ).

**Property 4.** Let  $p$  be a correct process such that  $p$  delivers  $(k, m)$  from  $q$ . We show that all correct process must deliver  $(k, m')$  from  $q$ , for some  $m'$ .

By construction of the algorithm, if  $p$  delivers  $(k, m)$  from  $q$ , then for all  $i < k$  there exists  $m_i$  such that  $p$  delivered  $(i, m_i)$  from  $q$  before delivering  $(k, m)$  (this is because  $p$  can only deliver  $(k, m)$  if  $\text{last}[q] = k$  and  $\text{last}[q]$  is only incremented to  $k$  after  $p$  delivers  $(k - 1, m_{k-1})$ ).

Assume by contradiction that there exists some correct process  $r$  which does not deliver  $(k, m')$  from  $q$ , for any  $m'$ . Further assume without loss of generality that  $k$  is the smallest key for which  $r$  does not deliver any message from  $q$ . Thus,  $r$  must have delivered  $(i, m'_i)$  from  $q$  for all  $i < k$ ; thus,  $r$  must have incremented  $\text{last}[q]$  to  $k$ . Since  $r$  never delivers any message from  $q$  for key  $k$ ,  $r$ 's  $\text{last}[q]$  will never increase past  $k$ .

Since  $p$  delivers  $(k, m)$  from  $q$ , then  $p$  must have written a valid L2 proof  $\mathcal{P}$  of  $(k, m)$  in its L2Proof slot in line 18. By Observation 1,  $\mathcal{P}$  will remain in  $p$ 's L2Proof[p,k,q] slot forever. Thus, at least one of the slots  $\text{L2Proof}[:, k, q]$  will forever contain a valid L2 proof. Since  $r$ 's  $\text{last}[q]$  eventually reaches  $k$  and never increases past  $k$ ,  $r$  will eventually (by Observation 2) see a valid L2 proof in line 16 and deliver a message for key  $k$  from  $q$ . We have reached a contradiction.  $\square$

## B CORRECTNESS OF CHEAP QUORUM

We prove that Cheap Quorum satisfies certain useful properties that will help us show that it composes with Preferential Paxos to form a correct weak Byzantine agreement protocol. For the proofs, we first formalize some terminology. We say that a process *proposed* a value  $v$  by time  $t$  if it successfully executes line 4; that is,  $p$  receives the response *ack* in line 4 by  $t$ . When a process aborts, note that it outputs a tuple. We say that the first element of its tuple is its *abort value*, and the second is its *abort proof*. We sometimes say that a process  $p$  aborts with value  $v$  and proof  $pr$ , meaning that  $p$  outputs  $(v, pr)$  in its abort. Furthermore, the value in a process  $p$ 's Proof region is called a *correct unanimity proof* if it contains  $n$  copies of the same value, each correctly signed by a different process.

**OBSERVATION 3.** *In Cheap Quorum, no value written by a correct process is ever overwritten.*

**PROOF.** By inspecting the code, we can see that the correct behavior is for processes to never overwrite any values. Furthermore, since all regions are initially single-writer, and the legalChange function never allows another process to acquire write permission on a region

that they cannot write to initially, no other process can overwrite these values.  $\square$

**LEMMA B.1 (CHEAP QUORUM VALIDITY).** *In Cheap Quorum, if there are no faulty processes and some process decides  $v$ , then  $v$  is the input of some process.*

**PROOF.** If  $p = p_1$ , the lemma is trivially true, because  $p_1$  can only decide on its input value. If  $p \neq p_1$ ,  $p$  can only decide on a value  $v$  if it read that value from the leader's region. Since only the leader can write to its region, it follows that  $p$  can only decide on a value that was proposed by the leader ( $p_1$ ).  $\square$

**LEMMA B.2 (CHEAP QUORUM TERMINATION).** *If a correct process  $p$  proposes some value, every correct process  $q$  will decide a value or abort.*

**PROOF.** Clearly, if  $q = p_1$  proposes a value, then  $q$  decides. Now let  $q \neq p_1$  be a correct follower and assume  $p_1$  is a correct leader that proposes  $v$ . Since  $p_1$  proposed  $v$ ,  $p_1$  was able to write  $v$  in the leader region, where  $v$  remains forever by Observation 3. Clearly, if  $q$  eventually enters panic mode, then it eventually aborts; there is no waiting done in panic mode. If  $q$  never enters panic mode, then  $q$  eventually sees  $v$  on the leader region and eventually finds  $2f + 1$  copies of  $v$  on the regions of other followers (otherwise  $q$  would enter panic mode). Thus  $q$  eventually decides  $v$ .  $\square$

**LEMMA B.3 (CHEAP QUORUM PROGRESS).** *If the system is synchronous and all processes are correct, then no correct process aborts in Cheap Quorum.*

**PROOF.** Assume the contrary: there exists an execution in which the system is synchronous and all processes are correct, yet some process aborts. Processes can only abort after entering panic mode, so let  $t$  be the first time when a process enters panic mode and let  $p$  be that process. Since  $p$  cannot have seen any other process declare panic,  $p$  must have either timed out at line 12 or 22, or its checks failed on line 13. However, since the entire system is synchronous and  $p$  is correct,  $p$  could not have panicked because of a time-out at line 12. So,  $p_1$  must have written its value  $v$ , correctly signed, to  $p_1$ 's region at a time  $t' < t$ . Therefore,  $p$  also could not have panicked by failing its checks on line 13. Finally, since all processes are correct and the system is synchronous, all processes must have seen  $p_1$ 's value and copied it to their slot. Thus,  $p$  must have seen these values and decided on  $v$  at line 20, contradicting the assumption that  $p$  entered panic mode.  $\square$

**LEMMA B.4 (LEMMA 4.5: CHEAP QUORUM DECISION AGREEMENT).** *Let  $p$  and  $q$  be correct processes. If  $p$  decides  $v_1$  while  $q$  decides  $v_2$ , then  $v_1 = v_2$ .*

**PROOF.** Assume the property does not hold:  $p$  decided some value  $v_1$  and  $q$  decided some different value  $v_2$ . Since  $p$  decided  $v_1$ , then  $p$  must have seen a copy of  $v_1$  at  $2fp + 1$  replicas, including  $q$ . But then  $q$  cannot have decided  $v_2$ , because by Observation 3,  $v_1$  never gets overwritten from  $q$ 's region, and by the code,  $q$  only can decide a value written in its region.  $\square$

**LEMMA B.5 (LEMMA 4.6: CHEAP QUORUM ABORT AGREEMENT).** *Let  $p$  and  $q$  be correct processes (possibly identical). If  $p$  decides  $v$  in Cheap Quorum while  $q$  aborts from Cheap Quorum,*

then  $v$  will be  $q$ 's abort value. Furthermore, if  $p$  is a follower,  $q$ 's abort proof is a correct unanimity proof.

PROOF. If  $p = q$ , the property follows immediately, because of lines 4 through 6 of panic mode. If  $p \neq q$ , we consider two cases:

- If  $p$  is a follower, then for  $p$  to decide, all processes, and in particular,  $q$ , must have replicated both  $v$  and a correct proof of unanimity before  $p$  decided. Therefore, by Observation 3,  $v$  and the unanimity proof are still there when  $q$  executes the panic code in lines 4 through 6. Therefore  $q$  will abort with  $v$  as its value and a correct unanimity proof as its abort proof.
- If  $p$  is the leader, then first note that since  $p$  is correct, by Observation 3  $v$  remains the value written in the leader's Value region. There are two cases. If  $q$  has replicated a value into its Value region, then it must have read it from  $Value[p_1]$ , and therefore it must be  $v$ . Again by Observation 3,  $v$  must still be the value written in  $q$ 's Value region when  $q$  executes the panic code. Therefore  $q$  aborts with value  $v$ . Otherwise, if  $q$  has not replicated a value, then  $q$ 's Value region must be empty at the time of the panic, since the  $legalChange$  function disallows other processes from writing on that region. Therefore  $q$  reads  $v$  from  $Value[p_1]$  and aborts with  $v$ .  $\square$

LEMMA B.6. *Cheap Quorum is 2-deciding.*

PROOF. Consider an execution in which every process is correct and the system is synchronous. Then no process will enter panic mode (by Lemma B.3) and thus  $p_1$  will not have its permission revoked.  $p_1$  will therefore be able to write its input value to  $p_1$ 's region and decide after this single write (2 delays).  $\square$

## C CORRECTNESS OF THE FAST & ROBUST

The following is the pseudocode of Preferential Paxos. Recall that *T-send* and *T-receive* are the trusted message passing primitives that are implemented in [21] using non-equivocation and signatures.

**Algorithm 8: Preferential Paxos—code for process  $p$**

```

1 propose( $v$ , priorityTag){
2   T-send( $v$ , priorityTag) to all;
3   Wait to T-receive ( $val$ , priorityTag) from  $n - f_p$ 
   ↪ processes;
4   best = value with highest priority out of
   ↪ messages received;
5   RobustBackup(Paxos).propose(best); }
```

LEMMA C.1 (LEMMA 4.7: PREFERENTIAL PAXOS PRIORITY DECISION). *Preferential Paxos implements weak Byzantine agreement with  $n \geq 2f_p + 1$  processes. Furthermore, let  $v_1, \dots, v_n$  be the input values of an instance  $C$  of Preferential Paxos, ordered by priority. The decision value of correct processes is always one of  $v_1, \dots, v_{f_p+1}$ .*

PROOF. By Lemma 4.3, Robust Backup(Paxos) solves weak Byzantine agreement with  $n \geq 2f_p + 1$  processes. Note that before calling Robust Backup(Paxos), each process may change its input, but only to the input of another process. Thus, by the correctness and fault tolerance of Paxos, Preferential Paxos clearly solves

weak Byzantine agreement with  $n \geq 2f_p + 1$  processes. Thus we only need to show that Preferential Paxos satisfies the priority decision property with  $2f_p + 1$  processes that may only fail by crashing.

Since Robust Backup(Paxos) satisfies validity, if all processes call Robust Backup(Paxos) in line 5 with a value  $v$  that is one of the  $f_p + 1$  top priority values (that is,  $v \in \{v_1, \dots, v_{f_p+1}\}$ ), then the decision of correct processes will also be in  $\{v_1, \dots, v_{f_p+1}\}$ . So we just need to show that every process indeed adopts one of the top  $f_p + 1$  values. Note that each process  $p$  waits to see  $n - f_p$  values, and then picks the highest priority value that it saw. No process can lie or pick a different value, since we use T-send and T-receive throughout. Thus,  $p$  can miss at most  $f_p$  values that are higher priority than the one that it adopts.  $\square$

We now prove the following key composition property that shows that the composition of Cheap Quorum and Preferential Paxos is safe.

LEMMA C.2 (LEMMA 4.8: COMPOSITION LEMMA). *If some correct process decides a value  $v$  in Cheap Quorum before an abort, then  $v$  is the only value that can be decided in Preferential Paxos with priorities as defined in Definition 3.*

PROOF. To prove this lemma, we mainly rely on two properties: the Cheap Quorum Abort Agreement (Lemma 4.6) and Preferential Paxos Priority Decision (Lemma 4.7). We consider two cases.

*Case 1.* Some correct follower process  $p \neq p_1$  decided  $v$  in Cheap Quorum. Then note that by Lemma 4.6, all correct processes aborted with value  $v$  and a correct unanimity proof. Since  $n \geq 2f + 1$ , there are at least  $f + 1$  correct processes. Note that by the way we assign priorities to inputs of Preferential Paxos in the composition of the two algorithms, all correct processes have inputs with the highest priority. Therefore, by Lemma 4.7, the only decision value possible in Preferential Paxos is  $v$ . Furthermore, note that by Lemma 4.5, if any other correct process decided in Cheap Quorum, that process's decision value was also  $v$ .

*Case 2.* Only the leader,  $p_1$ , decides in Cheap Quorum, and  $p_1$  is correct. Then by Lemma 4.6, all correct processes aborted with value  $v$ . Since  $p_1$  is correct,  $v$  is signed by  $p_1$ . It is possible that some of the processes also had a correct unanimity proof as their abort proof. However, note that in this scenario, all correct processes (at least  $f + 1$  processes) had inputs with either the highest or second highest priorities, all with the same abort value. Therefore, by Lemma 4.7, the decision value must have been the value of one of these inputs. Since all these inputs had the same value  $v$ ,  $v$  must be the decision value of Preferential Paxos.  $\square$

THEOREM C.3 (END-TO-END VALIDITY). *In the Fast & Robust algorithm, if there are no faulty processes and some process decides  $v$ , then  $v$  is the input of some process.*

PROOF. Note that by Lemmas B.1 and 4.7, this holds for each of the algorithms individually. Furthermore, recall that the abort values of Cheap Quorum become the input values of Preferential Paxos, and the set-up phase does not invent new values. Therefore, we just have to show that if Cheap Quorum aborts, then all abort values are inputs of some process. Note that by the code in panic mode, if Cheap Quorum aborts, a process  $p$  can output an abort value from one of three sources: its own Value region, the leader's Value region,



or its own input value. Clearly, if its abort value is its input, then we are done. Furthermore note that a correct leader only writes its input in the Value region, and correct followers only write a copy of the leader's Value region in their own region. Since there are no faults, this means that only the input of the leader may be written in any Value region, and therefore all processes always abort with some processes input as their abort value.  $\square$

**THEOREM C.4 (END-TO-END AGREEMENT).** *In the Fast & Robust algorithm, if  $p$  and  $q$  are correct processes such that  $p$  decides  $v_1$  and  $q$  decides  $v_2$ , then  $v_1 = v_2$ .*

**PROOF.** First note that by Lemmas 4.5 and 4.7, each of the algorithms satisfy this individually. Thus Lemma 4.8 implies the theorem.  $\square$

**THEOREM C.5 (END-TO-END TERMINATION).** *In Fast & Robust algorithm, if some correct process is eventually the sole leader forever, then every correct process eventually decides.*

**PROOF.** Assume towards a contradiction that some correct process  $p$  is eventually the sole leader forever, and let  $t$  be the time when  $p$  last becomes leader. Now consider some process  $q$  that has not decided before  $t$ . We consider several cases:

- (1) If  $q$  is executing Preferential Paxos at time  $t$ , then  $q$  will eventually decide, by termination of Preferential Paxos (Lemma 4.7).
- (2) If  $q$  is executing Cheap Quorum at time  $t$ , we distinguish two sub-cases:
  - (a)  $p$  is also executing as the leader of Cheap Quorum at time  $t$ . Then  $p$  will eventually propose a value, so  $q$  will either decide in Cheap Quorum or abort from Cheap Quorum (by Lemma B.2) and decide in Preferential Paxos by Lemma 4.7.
  - (b)  $p$  is executing in Preferential Paxos. Then  $p$  must have panicked and aborted from Cheap Quorum. Thus,  $q$  will also abort from Cheap Quorum and decide in Preferential Paxos by Lemma 4.7.  $\square$

Note that to strengthen C.5 to general termination as stated in our model, we require the additional standard assumption [38] that some correct process  $p$  is eventually the sole leader forever. In practice, however,  $p$  does not need to be the sole leader forever, but rather *long enough* so that all correct processes decide.

## D CORRECTNESS OF PROTECTED MEMORY PAXOS

In this section, we present the proof of Theorem 5.1. We do so by showing the Algorithm 7 is an algorithm that satisfies all of the properties in the theorem.

We first show that Algorithm 7 correctly implements consensus, starting with validity. Intuitively, validity is preserved because each process that writes any value in a slot either writes its own value, or adopts a value that was previously written in a slot. We show that every value written in any slot must have been the input of some process.

**THEOREM D.1 (VALIDITY).** *In Algorithm 7, if a process  $p$  decides a value  $v$ , then  $v$  was the input to some process.*

**PROOF.** Assume by contradiction that some process  $p$  decides a value  $v$  and  $v$  is not the input of any process. Since  $v$  is not the input value of  $p$ , then  $p$  must have adopted  $v$  by reading it from some process  $p'$  at line 23. Note also that a process cannot adopt the initial value  $\perp$ , and thus,  $v$  must have been written in  $p'$ 's memory by some other process. Thus we can define a sequence of processes  $s_1, s_2, \dots, s_k$ , where  $s_i$  adopts  $v$  from the location where it was written by  $s_{i+1}$  and  $s_1 = p$ . This sequence is necessarily finite since there have been a finite number of steps taken up to the point when  $p$  decided  $v$ . Therefore, there must be a last element of the sequence,  $s_k$  who wrote  $v$  in line 35 without having adopted  $v$ . This implies  $v$  was  $s_k$ 's input value, a contradiction.  $\square$

We now focus on agreement.

**THEOREM D.2 (AGREEMENT).** *In Algorithm 7, for any processes  $p$  and  $q$ , if  $p$  and  $q$  decide values  $v_p$  and  $v_q$  respectively, then  $v_p = v_q$ .*

Before showing the proof of the theorem, we first introduce the following useful lemmas.

**LEMMA D.3.** *The values a leader accesses on remote memory cannot change between when it reads them and when it writes them.*

**PROOF.** Recall that each memory only allows write-access to the most recent process that acquired it. In particular, that means that each memory only gives access to one process at a time. Note that the only place at which a process acquires write-permissions on a memory is at the very beginning of its run, before reading the values written on the memory. In particular, for each memory  $d$  a process does not issue a read on  $d$  before its permission request on  $d$  successfully completes. Therefore, if a process  $p$  succeeds in writing on memory  $m$ , then no other process could have acquired  $d$ 's write-permission after  $p$  did, and therefore, no other process could have changed the values written on  $m$  after  $p$ 's read of  $m$ .  $\square$

**LEMMA D.4.** *If a leader writes values  $v_i$  and  $v_j$  at line 35 with the same proposal number to memories  $i$  and  $j$ , respectively, then  $v_i = v_j$ .*

**PROOF.** Assume by contradiction that a leader  $p$  writes different values  $v_1 \neq v_2$  with the same proposal number. Since each thread of  $p$  executes the phase 2 write (line 35) at most once per proposal number, it must be that different threads  $T_1$  and  $T_2$  of  $p$  wrote  $v_1$  and  $v_2$ , respectively. If  $p$  does not perform phase 1 (i.e., if  $p = p_1$  and this is  $p$ 's first attempt), then it is impossible for  $T_1$  and  $T_2$  to write different values, since CurrentVal was set to  $v$  at line 14 and was not changed afterwards. Otherwise (if  $p$  performs phase 1), then let  $t_1$  and  $t_2$  be the times when  $T_1$  and  $T_2$  executed line 8, respectively ( $T_1$  and  $T_2$  must have done so, since we assume that they both reached the phase 2 write at line 35). Assume wlog that  $t_1 \leq t_2$ . Due to the check and abort at line 29, CurrentVal cannot change after  $t_1$  while keeping the same proposal number. Thus, when  $T_1$  and  $T_2$  perform their phase 2 writes (after  $t_1$ ), CurrentVal has the same value as it did at  $t_1$ ; it is therefore impossible for  $T_1$  and  $T_2$  to write different values. We have reached a contradiction.  $\square$

**LEMMA D.5.** *If a process  $p$  performs phase 1 and then writes to some memory  $m$  with proposal number  $b$  at line 35, then  $p$  must have written  $b$  to  $m$  at line 21 and read from  $m$  at line 23.*

PROOF. Let  $T$  be the thread of  $p$  which writes to  $m$  at line 35. If phase 1 is performed (i.e., the condition at line 19 is satisfied), then by construction  $T$  cannot reach line 35 without first performing lines 21 and 23. Since  $T$  only communicates with  $m$ , it must be that lines 21 and 23 are performed on  $m$ .  $\square$

PROOF OF THEOREM D.2. Assume by contradiction that  $v_p \neq v_q$ . Let  $b_p$  and  $b_q$  be the proposal numbers with which  $v_p$  and  $v_q$  are decided, respectively. Let  $W_p$  (resp.  $W_q$ ) be the set of memories to which  $p$  (resp.  $q$ ) successfully wrote in phase 2 line 35 before deciding  $v_p$  (resp.  $v_q$ ). Since  $W_p$  and  $W_q$  are both majorities, their intersection must be non-empty. Let  $m$  be any memory in  $W_p \cap W_q$ .

We first consider the case in which one of the processes did not perform phase 1 before deciding (i.e., one of the processes is  $p_1$  and it decided on its first attempt). Let that process be  $p$  wlog. Further assume wlog that  $q$  is the first process to enter phase 2 with a value different from  $v_p$ .  $p$ 's phase 2 write on  $m$  must have preceded  $q$  obtaining permissions from  $m$  (otherwise,  $p$ 's write would have failed due to lack of permissions). Thus,  $q$  must have seen  $p$ 's value during its read on  $m$  at line 23, and thus  $q$  cannot have adopted its own value. Since  $q$  is the first process to enter phase 2 with a value different from  $v_p$ ,  $q$  cannot have adopted any other value than  $v_p$ , so  $q$  must have adopted  $v_p$ . Contradiction.

We now consider the remaining case: both  $p$  and  $q$  performed phase 1 before deciding. We assume wlog that  $b_p < b_q$  and that  $b_q$  is the smallest proposal number larger than  $b_p$  for which some process enters phase 2 with  $\text{CurrentVal} \neq v_p$ .

Since  $b_p < b_q$ ,  $p$ 's read at  $m$  must have preceded  $q$ 's phase 1 write at  $m$  (otherwise  $p$  would have seen  $q$ 's larger proposal number and aborted). This implies that  $p$ 's phase 2 write at  $m$  must have preceded  $q$ 's phase 1 write at  $m$  (by Lemma D.3). Thus  $q$  must have seen  $v_p$  during its read and cannot have adopted its own input value. However,  $q$  cannot have adopted  $v_p$ , so  $q$  must have adopted  $v_q$  from some other slot  $sl$  that  $q$  saw during its read. It must be that  $sl.\text{minProposal} < b_q$ , otherwise  $q$  would have aborted. Since  $sl.\text{minProposal} \geq sl.\text{accProposal}$  for any slot, it follows that  $sl.\text{accProposal} < b_q$ . If  $sl.\text{accProposal} < b_p$ ,  $q$  cannot have adopted  $sl.\text{value}$  in line 30 (it would have adopted  $v_p$  instead). Thus it must be that  $b_p \leq sl.\text{accProposal} < b_q$ ; however, this is impossible because we assumed that  $b_q$  is the smallest proposal number larger than  $b_p$  for which some process enters phase 2 with  $\text{CurrentVal} \neq v_p$ . We have reached a contradiction.  $\square$

Finally, we prove that the termination property holds.

THEOREM D.6 (TERMINATION). *Eventually, all correct processes decide.*

LEMMA D.7. *If a correct process  $p$  is executing the for loop at lines 18–37, then  $p$  will eventually exit from the loop.*

PROOF. The threads of the for loop perform the following potentially blocking steps: obtaining permission (line 20), writing (lines 21 and 35), reading (line 23), and waiting for other threads (the barrier at line 7 and the exit condition at line 37). By our assumption that a majority of memories are correct, a majority of the threads of the for loop must eventually obtain permission in line 20 and invoke the write at line 21. If one of these writes fails due to lack of permission, the loop aborts and we are done. Otherwise, a majority of threads

will perform the read at line 23. If some thread aborts at lines 22 and 26, then the loop aborts and we are done. Otherwise, a majority of threads must add themselves to `ListOfReady`, pass the barrier at line 7 and invoke the write at line 35. If some such write fails, the loop aborts; otherwise, a majority of threads will reach the check at line 37 and thus the loop will exit.  $\square$

PROOF OF THEOREM D.6. The  $\Omega$  failure detector guarantees that eventually, all processes trust the same correct process  $p$ . Let  $t$  be the time after which all correct processes trust  $p$  forever. By Lemma D.7, at some time  $t' \geq t$ , all correct processes except  $p$  will be blocked at line 12. Therefore, the `minProposal` values of all memories, on all slots except those of  $p$  stop increasing. Thus, eventually,  $p$  picks a `propNr` that is larger than all others written on any memory, and stops restarting at line 26. Furthermore, since no process other than  $p$  is executing any steps of the algorithm, and in particular, no process other than  $p$  ever acquires any memory after time  $t'$ ,  $p$  never loses its permission on any of the memories. So, all writes executed by  $p$  on any correct memory must return `ack`. Therefore,  $p$  will decide and broadcast its decision to all. All correct processes will receive  $p$ 's decision and decide as well.  $\square$

To complete the proof of Theorem 5.1, we now show that Algorithm 7 is 2-deciding.

THEOREM D.8. *Algorithm 7 is 2-deciding.*

PROOF. Consider an execution in which  $p_1$  is timely, and no process's failure detector ever suspects  $p_1$ . Then, since no process thinks itself the leader, and processes do not deviate from their protocols, no process calls `changePermission` on any memory. Furthermore,  $p_1$ 's does not perform phase 1 (lines 19–33), since it is  $p_1$ 's first attempt. Thus, since  $p_1$  initially has write permission on all memories, all of  $p_1$ 's phase 2 writes succeed. Therefore,  $p_1$  terminates, deciding its own proposed value  $v$ , after one write per memory.  $\square$

## E PSEUDOCODE FOR THE ALIGNED PAXOS

### Algorithm 9: Aligned Paxos

```

1  A=(P, M) is set of acceptors
2  propose(v){
3    resps = [] //prepare empty responses list
4    choose propNr bigger than any seen before
5    for all a in A{
6      cflag = communicate1(a, propNr)
7      resp = hearback1(a)
8      if (cflag){resps.append((a, resp)) } }
9    wait until resps has responses from a majority of A
10   next = analyze1(resps)
11   if (next == RESTART) restart;
12   resps = []
13   for all a in A{
14     cflag = communicate2(a, next)
15     resp = hearback2(a)
16     if (cflag){resps.append((a, resp)) } }
17   wait until resps has responses from a majority of A
18   next = analyze2(resps)
19   if (next == RESTART) restart;
20   decide next; }

```

**Algorithm 15: Analyze Phase 2—code for process  $p$** 

```

1 value analyze2(value v, responses resps){
2   if there are at least  $A/2 + 1$  resps such that resp.
3     ↳ response==ack{
4     return v }
5   return RESTART }

```

**Algorithm 10: Communicate Phase 1—code for process  $p$** 

```

1 bool communicate1(agent a, value propNr){
2   if (a is memory){
3     changePermission(a, {(R:II-{p}, W:0, RW: {p})}); //
4     ↳ acquire write permission
5     return write(a[p], {propNr, -, -}) }
6   else{
7     send prepare(propNr) to a
8     return true } }

```

**Algorithm 11: Hear Back Phase 1—code for process  $p$** 

```

1 value hearback1(agent a){
2   if (a is memory){
3     for all processes q{
4       localInfo[q] = read(a[q]) }
5     return = localInfo}
6   else{
7     return value received from a } }

```

**Algorithm 12: Analyze Phase 1—code for process  $p$** 

```

1 responses is a list of (agent, response) pairs
2 value analyze1(responses resps){
3   for resp in resps {
4     if (resp.agent is memory){
5       for all slots s in v.info {
6         if (s.minProposal > propNr) return RESTART }
7       (v, accProposal) = value and accProposal of slot
8         ↳ with highest
9         accProposal that had a value } }
10  return v where (v, accProposal) is the highest
11  ↳ accProposal seen in resps.response of all
12  ↳ agents }

```

**Algorithm 13: Communicate Phase 2—code for process  $p$** 

```

1 bool communicate2(agent a, value msg){
2   if (a is memory){
3     return write(a[p], {msg.propNr, msg.propNr, msg.
4       ↳ val})}
5   else{
6     send accepted(msg.propNr, msg.val) to a
7     return true } }

```

**Algorithm 14: Hear Back Phase 2—code for process  $p$** 

```

1 value hearback2(agent a){
2   if (a is memory){
3     return ack }
4   else{
5     return value received from a } }

```