

Computationally Light “Multi-Speed” Atomic Memory

Antonio Fernández Anta¹, Theophanis Hadjistasi², and Nicolas Nicolaou¹

- 1 IMDEA Networks Institute
Madrid, Spain
antonio.fernandez@imdea.org, nicolas.nicolaou@imdea.org
- 2 University of Connecticut
Storrs CT, USA
theophanis.hadjistasi@uconn.edu

Abstract

Communication demands are usually the leading factor that defines the efficiency of operations on a read/write shared memory emulation in the message-passing environment. In the quest for minimizing the communication demands, the algorithms proposed either require restrictions in the system or incur high computation demands. As a result, such solutions may be not suitable to be used in practice.

In this paper we focus on the practicality of implementations of *atomic read/write shared memory* emulation in the message-passing environment. In particular we investigate implementations that reduce both *communication* and *computation* demands. We first examine the shortcomings of the best two (in terms of communication demands) known algorithms that implement atomic single-writer multiple-reader (SWMR) atomic memory, [4, 7]. The algorithm CCFast proposed in [4], achieves optimal communication by allowing each operation to complete in one round trip, with light computation requirements. Unfortunately, it relies on strict limitations on the number of readers. On the other hand, algorithm OHSAM [7], imposes no restrictions on the system, but provides operations that require one and a half communication rounds. In the light of these shortcomings, we present two algorithms that implement *multi-speed* operations with *light computation*, and *without imposing* any restriction on the system. In particular, algorithm CCHYBRID adopts the fast (one-round) writes presented in [4], and makes *clients* to switch to a slow (two-round) mode whenever the system is congested. On the other hand, algorithm OHFAST, pushes the responsibility of deciding for the speed switch to the *servers*. This allows the algorithm to utilize the fast operations presented in [4], and the slow one-and-a-half-rounds operations of [7], whenever is necessary. We prove that both new algorithms preserve atomicity. To evaluate the new algorithms we implement five different atomic memory algorithms in the NS3 simulator, and we compare their performance in terms of *operation latency*, and *ratio of slow over fast operations* performed. We test the algorithms over different: (i) topologies, and (ii) operation loads. Our results support that the newly presented algorithms increase the practicality of atomic read/write atomic shared memory implementations in the message-passing, asynchronous environment.

1998 ACM Subject Classification C.3.4 Distributed Systems, C.4 Performance of Systems

Keywords and phrases atomicity, read/write objects, shared memory, operation latency



© Antonio Fernández Anta, Theophanis Hadjistasi, and Nicolas Nicolaou;
licensed under Creative Commons License CC-BY

Leibniz International Proceedings in Informatics

LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

1 Introduction

Emulating atomic [9] (linearizable [8]) read/write objects in message-passing environments is one of the fundamental problems in distributed computing. The problem becomes more difficult when participants in the service may fail and the environment is asynchronous, i.e., it cannot provide any time guarantees on the delivery of the messages and the computation speeds. To cope with failures, traditional distributed object implementations like [2, 11], use *redundancy* by replicating the object to multiple (possibly geographically dispersed) locations (replica servers). Replication however raises the challenge of consistency, as multiple object copies can be accessed concurrently by multiple processes. Atomicity is the most intuitive consistency semantic, as it provides the illusion of a single-copy object that serializes all accesses: each read operation returns the value of the latest write operation.

Attiya, Bar-Noy, and Dolev [2] were the first to present an algorithm, known as ABD, to implement single-writer multi-reader (SWMR) atomic objects in message-passing, crash-prone, asynchronous environments. The authors associate logical *timestamps* to the values written, to impose an order on the write operations. The propagation of the latest timestamp (and its corresponding value) is based on the assumption that at least a majority of replica servers do not fail. In this setting, ABD has write operations that terminate with a single communication round-trip, and read operations that involve two round-trips. Based on basic value comparisons, ABD incurs almost no computational overhead to the service participants. Atomicity is guaranteed by the intersecting properties of two majorities and the second phase of a read operation. Following ABD, a folklore belief persisted that in asynchronous multi-reader (MR) atomic memory algorithms, “reads must write.”

The work by Dutta et al. [3] refuted this belief, by presenting atomic register algorithms in which every operation involves only a *single* round-trip. Such an algorithm is called *fast*. They showed that fast reads are possible only in the single-writer (SW) model, and given that the number of readers R is constrained with respect to the number of replicas S and the maximum number of failures f ; in particular, $R < \frac{S}{f} - 2$. A recent work by Fernández Anta, Nicolaou, and Popa [4], has shown that, although the result in [3] is efficient in terms of communication, it requires processes to evaluate a computationally hard (NP-hard) predicate. A new algorithm CCFast, with a new predicate, was proposed in that paper to allow operations terminate with a linear computation overhead. Despite improving the practicality of [3], the algorithm in [4] inherited the same system constraint as [3].

The idea of exploring “multi-speed” read operations is not new. An algorithm is said to be “multi-speed” when different read operations may perform different number of communication rounds before completing. Works like [5, 6] proposed implementations in the SWMR model with *two-speed* operations, in an attempt to relax the constraints proposed in [3], and to allow unbounded number of readers. In particular, the work in [6] presents algorithm SF, which applies a predicate similar to the one introduced in [3], but on *virtual nodes* (i.e., sets of readers) instead of individual reader processes. In [5], the authors introduced *quorum views*, which are client-side tools that examine the distribution of the latest value among the replicas, in order to enable fast read operations. Both [5, 6] trade communication for scalability. Under conditions of low concurrency, both algorithms allow most reads to complete in a single communication round-trip; otherwise a *two round-trip* operation (similar to ABD) is required. To determine the speed of an operation, both algorithms inflicted significant computational demands: (i) [6] exploited the same predicate as in [3], which is NP-hard [4], and (ii) [5] needed to examine the distribution of the object value within all the possible replica subsets. Thus, a trend appeared in the algorithms that

aimed for fast operations: algorithms with lower communication rounds demanded higher computation overhead at the processes.

Following the above findings, we say that an operation is *fast* if it completes in a single communication round trip, and *slow* if it completes in two round trips. A recent work by Hadjistasi, Nicolaou and Schwarzmann [7] redefines *slowness*, as they present an algorithm for the SWMR model, called OHSAM, where each operation takes *one and a half* round-trips to complete. As the number of readers is bounded when all operations are fast [3], the authors claim the optimality of their approach in terms of communication when no constraint is imposed. Furthermore their algorithm relies on basic comparisons, inflicting negligible computation overhead.

Contributions. In this paper, we focus in improving the practicality of SWMR atomic read/write register algorithms, by achieving low communication and computation costs on the atomic operations. We trade communication for scalability, by adopting the predicate presented in [4] and allowing some operations to be slow. Also, we combine ideas presented in both [4] and [7], to introduce implementations that allow only *single* and *one-and-a-half* round operations. Enumerated, our contributions are the following:

- We introduce a new “multi-speed” algorithm, CCHYBRID, that allows operations to terminate in *one* or *two* communication round-trips, and does not impose any bounds on the number of readers. CCHYBRID uses the predicate introduced in [4] to determine the speed of a read operation, and it requires at most one *complete* slow operation per written value. This is similar to the semifast algorithm SF [6]. However, in contrast to SF, in which processes have to decide NP-hard predicates, it incurs light (linear) computation.
- Next we examine whether we can combine the techniques presented in [4] and [7] to obtain a “multi-speed” algorithm that allows *one* and *one-and-a-half* round-trip operations. We present algorithm OHFAST, that achieves the targeted performance by moving the decision on whether a slow read operation is necessary to the servers. When servers determine that a slow read is necessary, they perform a *relay* phase to inform other servers before replying to the reader. It is interesting that in OHFAST not all the servers need to perform a relay for a single read operation. Some of the servers may reply directly to the read whereas some others may perform a relay phase for the same read. Thus a read operation may terminate before receiving a reply from a relaying server.
- We complement our algorithms with experimental results for five algorithms: ABD, OHSAM, CCHYBRID, OHFAST, and SF. ABD sets the threshold for the rest of the algorithms, while OHSAM sets the threshold on the operations that use one and a half rounds. Algorithm SF is used to demonstrate whether computation has an impact to the latency of operations. We test our algorithms under different scenarios by changing the number of participants, the frequency of operations, and using two network topologies: (i) a topology where servers are distributed evenly over the network, and (ii) a topology that resembles a datacenter where servers are concentrated in close proximity and communicate through high bandwidth links. Our results show that the proposed algorithms outperform the algorithms with “one speed” operations (i.e., ABD and OHSAM) in all scenarios, reducing the latency per operation to less than half in most cases. Compared with the semifast “multi-speed” algorithm SF, our algorithms achieve a similar read latency, even though the scenarios explored were extremely favorable for SF, since we observed that practically all its operations were fast and the NP-hard predicate evaluations were not heavy (mainly due to the good communication conditions). Finally, as expected, we observed that the topology has a great impact on the algorithms that use *one and a half* round operations.

2 Model

We assume a system consisting of three distinct sets of processes: a writer process with identifier w , a set \mathcal{R} of readers, and a set \mathcal{S} of replica servers. Let $\mathcal{I} = \{w\} \cup \mathcal{R} \cup \mathcal{S}$. In a read/write object implementation, we assume that the object may take a value from a set V . The writer is the sole process that is allowed to modify the value of the object, the readers are allowed to obtain the value of the object, and each server maintains a copy of the object to ensure the availability of the object in case of failures. We assume an *asynchronous* environment, where processes communicate by exchanging messages. The writer, any subset of readers, and up to $f < \frac{|\mathcal{S}|}{2}$ servers may *crash* without any notice.

An algorithm A is a collection of processes, where process A_p is assigned to processor $p \in \mathcal{I}$. Each processor p has a *state* which is determined over a set of state variables. The state of A is a vector that contains the state of each process. Algorithm A performs a *step*, when some process p atomically: (i) receives a message, (ii) performs local computation, (iii) sends a message. Each such step causes the state at p to change from a pre-state σ_p to a post-state σ'_p . Hence, the state of A changes from σ to σ' where σ contains state σ_p for p and σ' contains state σ'_p , while the state of every $p' \neq p$ is the same in both σ and σ' . An *execution fragment* is an alternating sequence of states and actions of A ending in a state. An *execution* is an execution fragment that starts with the initial state. An execution fragment ξ' extends an execution fragment ξ if the last state of ξ is the first state of ξ' . A process p *crashes* in an execution if it stops taking steps; otherwise p is *correct*. Each process may perform a read or write operation, and each operation has *invocation* and *response* steps. An operation π is *complete* in an execution ξ , if ξ contains both the invocation and the *matching* response step for π ; otherwise π is *incomplete*. An execution ξ is *well formed* if any process p that invokes an operation π in ξ does not invoke any other operation π' before the matching response step of π appears in ξ . An operation π *precedes* an operation π' in an execution ξ , denoted by $\pi \rightarrow \pi'$, if the response step of π appears before the invocation step of π' in ξ . Two operations are *concurrent* if none precedes the other.

Correctness of an implementation of an atomic read/write object is defined in terms of the *atomicity* and *termination* properties. The termination property requires that any operation invoked by a correct process eventually completes. For atomicity we use the definition of [10, Lemma 13.16].

Efficiency Metrics. We measure the complexity of an operation π in terms of: (i) *message complexity*, i.e. the worst-case number of messages exchanged during π , and (ii) *operation latency*, i.e. the *computation time* and the *communication delays* incurred by π . Computation time accounts the computation steps the algorithm performs in each operation. Communication delays are measured in *communication exchanges*, as defined in [7].

In particular, a protocol requires each operation to involve a sequence of sends (or broadcasts) of typed messages and the corresponding receives. A *communication exchange* during an operation π in an execution ξ , is defined as the collection of send and receive actions for a specific typed message (as required by the protocol) between the invocation and response of π in ξ . Using this definition, implementations, such as ABD, are structured in terms of *rounds*, where each round consists of two message exchanges: a *broadcast*, initiated by the process executing an operation, and a *convergecast* of responses to the initiator. A *fast* operation as in [6, 3] consists of two communication exchanges (or one round), and a *slow* operation as used in [2, 5, 6] consists of four communication exchanges (or two rounds). A read operation as in [7] consists of three communication exchanges (or 1.5 rounds). The number of messages that a process expects during a convergecast depends on the implementation.

Algorithm	WE	RE	WC	RC	WM	RM	Bounds
ABD [2]	2	4	$O(1)$	$O(S)$	$2 S $	$4 S $	Unbounded
FAST [3]	2	2	$O(1)$	$O(S ^2 \cdot 2^{ \mathcal{S} })$	$2 S $	$2 S $	$ \mathcal{R} < \frac{ S }{f} - 2$
ccFAST [4]	2	2	$O(1)$	$O(S)$	$2 S $	$2 S $	$ \mathcal{R} < \frac{ S }{f} - 2$
SF [6]	2	2 or 4	$O(1)$	$O(S ^2 \cdot 2^{ \mathcal{S} })$	$2 S $	$O(4 S)$	$ \mathcal{V} < \frac{ S }{f} - 1$
OHSAM [7]	2	3	$O(1)$	$O(S)$	$2 S $	$2 S + S ^2$	Unbounded
ccHYBRID (here)	2	2 or 4	$O(1)$	$O(S)$	$2 S $	$O(4 S)$	Unbounded
OhFAST (here)	2	2 or 3	$O(1)$	$O(S)$	$2 S $	$O(S ^2)$	Unbounded

■ **Table 1** Communication, Computation, Message Complexities and Participation Bounds. (WE/RE: write/read-communication exchanges, WC/RC: write/read-computation, WM/RM: write/read-number of messages). \mathcal{V} is the set of virtual nodes.

3 State-of-the-Art Performance of Atomic Memory Implementations

The algorithm by Dutta et al. in 2004 [3], we refer to it as FAST, was the first to present atomic register implementations where all operations take a *single* communication round before completing. To allow fast reads, FAST deploys a recording mechanism at each server and evaluates a predicate at each reader. It was shown that fast reads are possible only if $|\mathcal{R}| < \frac{|S|}{f} - 2$ readers participate in the service. To avoid the bound on the number of readers, Georgiou, Nicolaou and Shvartsman [6], grouped the readers under logical sets, they called *virtual groups*, and allowed some of the read operations to perform two rounds (or 4 communication exchanges). The predicate of [3] was applied on the virtual groups instead of individual readers, where each group could have an arbitrary size. As expected, the use of the predicate imposed a bound on the number of virtual nodes, for atomicity to be preserved; that is $|\mathcal{V}| < \frac{|S|}{f} - 1$.

Fernández Anta, Nicolaou and Popa [4], showed that the predicate used by both [3] and [6], is computationally hard. This was due to the fact that the original predicate was searching among all the subsets of servers to identify if there is some subset of servers that replied to a “large enough” subset of readers. To avoid this computational overhead, they investigate whether it is possible to use *how many* instead of *which* readers obtained the latest value, and still be able to preserve atomicity. Thus, the paper introduced a new algorithm, called cCFast, that was using the following predicate at the readers:

$$\exists \alpha \in [1, |\mathcal{R}| + 1] \text{ s.t. } MS = \{s : (s, m) \in \text{maxAck} \wedge m.\text{views} \geq \alpha\} \text{ and } |MS| \geq |S| - \alpha f$$

Essentially, each server records the readers that observed its local timestamp in a set *seen*, and whenever requested, it reports the *cardinality* of that set to the requesting process. A reader collects the replies from the servers in each read operation, detects the replies that contain the maximum timestamp (set *maxAck*), and checks the cardinalities reported in those replies (*m.views*). If there are “enough” replies with “sufficiently” large cardinalities, the predicate holds and the reader returns the value associated with the maximum timestamp; otherwise the value associated with the previous timestamp is returned. The evaluation of the predicate can be done in linear time with respect to the number of servers in the system. Their algorithm inherited the necessary bound presented in [3] on the number of readers participants, $|\mathcal{R}| < \frac{|S|}{f} - 2$.

Finally, Hadjistasi, Nicolaou and Schwarzmann [7] closed the gap of the communication of read/write operations by presenting algorithm OHSAM, where writes take just one round (or 2 communication exchanges) and reads always take one and a half round (or 3 communication exchanges) to complete. The main idea of the algorithm is to allow servers to exchange information about the operations, before replying to the invoking process. OHSAM uses

Algorithm 1 Write, Read and Server protocols of algorithm CCHYBRID

```

1: at the writer  $w$ 
2: Components:
3:  $ts \in \mathbb{N}^+$ ;  $v, vp \in V$ ;  $wcounter \in \mathbb{N}^+$ 
4: Initialization:
5:  $ts \leftarrow 0$ ;  $v \leftarrow \perp$ ;  $vp \leftarrow \perp$ ;  $wcounter \leftarrow 0$ 
6: function WRITE( $val$ )
7:    $vp \leftarrow v$ ;  $v \leftarrow val$ 
8:    $ts \leftarrow ts + 1$ 
9:    $wcounter \leftarrow wcounter + 1$ 
10:  send(( $ts, v, vp$ ),  $w, wcounter$ ) to all servers
11:  wait until  $|\mathcal{S}| - f$  servers reply
12:  return(OK)
13: end function

14: at each reader  $r_i$ 
15: Components:
16:  $ts \in \mathbb{N}^+$ ;  $maxTS \in \mathbb{N}^+$ ;  $v, vp \in V$ ;  $rcounter \in \mathbb{N}^+$ 
17:  $srvAck \subseteq \mathcal{S} \times M$ 
18: Initialization:
19:  $ts \leftarrow 0$ ;  $maxTS \leftarrow 0$ ;  $v \leftarrow \perp$ ;  $vp \leftarrow \perp$ ;  $rcounter \leftarrow 0$ 
20: function READ()
21:   $rcounter \leftarrow rcounter + 1$ 
22:  send(( $ts, v, vp$ ),  $r_i, rcounter$ ) to all servers
23:  wait until  $|\mathcal{S}| - f$  servers reply
24:   $\triangleright$  Collect ( $sid, \langle (ts', v', vp'), views, prop \rangle$ ) msgs in  $srvAck$ 
25:   $maxTS \leftarrow \max(\{m.ts' \mid (s, m) \in srvAck\})$ 
26:   $maxAck \leftarrow \{(s, m) \mid (s, m) \in srvAck \wedge m.ts' = maxTS\}$ 
27:   $\langle ts, v, vp \rangle \leftarrow m.\langle ts', v', vp' \rangle$  for  $(s, m) \in maxAck$ 
28:   $maxViews \leftarrow \max(\{m.views \mid (s, m) \in maxAck\})$ 
29:   $propSet \leftarrow \{s \mid (s, m) \in maxAck \wedge m.prop = True\}$ 
30:  if  $maxViews > \frac{|\mathcal{S}|}{f} - 2 \vee propSet \neq \emptyset$  then
31:    send(( $ts, v, vp$ ),  $r_i, rcounter$ ) to all servers
32:    wait until  $|\mathcal{S}| - f$  servers reply
33:    end if
34:    return( $v$ )
35:  else
36:    if  $\exists \alpha \in [1, \frac{|\mathcal{S}|}{f} - 2]$  s.t.
37:       $MS = \{s : (s, m) \in maxAck \wedge m.views \geq \alpha\}$  and
38:       $|MS| \geq |\mathcal{S}| - \alpha f$  then
39:        return( $v$ )
40:      else
41:        return( $vp$ )
42:      end if
43:    end if
44:  end function

45: at each server  $s_i$ 
46: Components:
47:  $ts \in \mathbb{N}^+$ ;  $seen \subseteq \mathcal{R} \cup \{w\}$ ;  $v, vp \in V$ ;  $prop \in \{True, False\}$ 
48:  $Counter[|\mathcal{R}| + 1] \in \mathbb{N}^+$ 
49: Initialization:
50:  $ts \leftarrow 0$ ;  $seen \leftarrow \emptyset$ ;  $v, vp \leftarrow \perp$ ;  $prop \leftarrow False$ 
51:  $Counter[i] \leftarrow 0$  for  $i \in \mathcal{R} \cup \{w\}$ 
52: function RCV( $\langle (ts', v', vp'), q, counter \rangle$ )
53:    $\triangleright$  Called upon reception of a message
54:   if  $Counter[q] < counter$  then
55:     if  $ts' > ts$  then
56:        $\langle ts, v, vp \rangle \leftarrow \langle ts', v', vp' \rangle$ 
57:        $seen \leftarrow \{q\}$ 
58:        $prop \leftarrow False$ 
59:     else
60:        $seen \leftarrow seen \cup \{q\}$ 
61:     end if
62:     if  $ts' = ts \wedge q \in \mathcal{R}$  then
63:        $prop \leftarrow True$ 
64:     end if
65:     send(( $ts, v, vp$ ),  $|seen|, prop$ ) to  $q$ 
66:   end if
67: end function

```

negligible computation at the processors, as each operation performs only basic comparisons. However, the server communication in *every* operation makes the algorithm suitable for environments where server communication is being carried out by high capacity links.

Table 1, summarizes the efficiency of each of the algorithms in different efficiency metrics. It also presents any bounds that an algorithm may impose on the participation of the service in order to be able to provide atomic guarantees. The last two algorithms are the ones we present in this paper. Notice that the goal is to minimize communication without inflicting high computation overheads, or participation bounds in the system.

4 Algorithm CCHYBRID: Switching from One to Two Rounds

As discussed in Section 3, algorithm CCFast guarantees correctness only when the number of readers is bounded with respect to the ratio of the number of servers and the number of failures in the system, i.e. $|\mathcal{R}| < \frac{|\mathcal{S}|}{f} - 1$. In this section we propose a modification to CCFast that removes the bound on the number of readers. To unleash the number of readers, the new algorithm CCHYBRID, allows some read operations to complete in two rounds. In particular, CCHYBRID combines ideas from CCFast and ABD: (i) it exploits *timestamp-value* pairs to order the write operations, (ii) it uses the predicate proposed by CCFast to determine the value returned by a *fast* read, and (iii) it propagates the maximum *timestamp-value* pair to a majority of servers during a *slow* read. The biggest challenge in CCHYBRID is to determine *when* a second phase is necessary, and ensure that such a strategy does not violate atomicity. The idea of CCHYBRID is to have the reader examine if the number of processes that observed the latest value is over the bound $\frac{|\mathcal{S}|}{f} - 1$. If not, then CCHYBRID evaluates the predicate proposed in CCFast over the replies, to determine the value to return. Otherwise, it proceeds to a propagation phase to send the latest value to a majority of servers. To prevent readers from propagating an already propagated value, servers maintain a flag that indicates whether a timestamp has been propagated.

Algorithm 1 provides the formal pseudocode of CCHYBRID. The write protocol remains

the same as in both CCFast and ABD: the writer increments its local timestamp (L8) and propagates the *timestamp-value* pair to a majority of servers (L10-11). The server protocol is more involved. In addition to the replica state (timestamp and value), a server s maintains a set *seen* to record the processes that requested this replica, and a flag *prop* that, as we explain later, optimizes read operations. A server s waits for read and write requests. When a request is received, s updates its local *timestamp-value* pair (L51-57) if the *timestamp* attached in the received message is greater than its local timestamp. In addition, it initializes its *seen* set to contain the sender process, and sets the *prop* flag to *False*. In case the timestamp of the message is not greater than the local timestamp of s , then the server records the sender in its *seen* set (L59). The server s sets *prop* = *True* when it receives a message from a reader that contained a *timestamp-value* pair equal to the one that is locally stored in s . Notice that a reader propagates a *timestamp-value* pair in every phase. So, s may set *prop* during the first or second phase of a read.

The main departure of CCHYBRID from CCFast lies in the read protocol. A reader behaves as in CCFast as long as the maximum number of *views* reported by the servers remains below $\frac{|\mathcal{S}|}{f} - 2$. In particular, a reader sends read messages to all the servers and waits from $|\mathcal{S}| - f$ to reply (L22). When those replies are received, the reader discovers the maximum timestamp (*maxTS*) among the replies (L24), the set of messages that contained *maxTS* (L25), and the maximum views reported in those messages (L27). If the maximum views are less than $\frac{|\mathcal{S}|}{f} - 2$ and no reader propagated *maxTS* (L29), then the reader evaluates the predicate as in CCFast to decide which value to return; otherwise the reader returns the value associated with the *maxTS*. If at least $f + 1$ of the messages that contain *maxTS*, also contain *prop* = *True*, the reader returns without further action. If this is not the case then the reader performs a second phase propagating the maximum *timestamp-value* pair to $|\mathcal{S}| - f$ servers (L30-33). Notice that CCHYBRID performs equally to CCFast when the number of readers that return the same value (not necessarily the same readers for each value) satisfies the bound required by CCFast. In any other case, a *single complete, slow* read operation (similar to [6]) is necessary per write operation. The use of the *prop* flag allows any read that succeeds a slow read, and returns the same value, to be *fast*, as: (i) The slow read propagates the *maxTS* to $|\mathcal{S}| - f$ servers, (ii) a succeeding read receives replies from $|\mathcal{S}| - f$ servers, and (iii) the read discovers *prop* = *True* for *maxTS* in more than $|\mathcal{S}| - 2f > f + 1$ servers.

4.1 Algorithm Correctness

Our algorithm is correct if it can satisfy Termination (liveness condition) and Atomicity (safety condition). It is trivial to see that termination is satisfied given that the system respects our failure model. To proof atomicity we are going to express atomicity in terms of timestamps written and returned in a SWMR model, as also presented in [4]:

- A1.** For each process p the ts variable is non-negative and monotonically nondecreasing.
- A2.** If a read ρ succeeds a write operation $\omega(ts)$ and returns a timestamp ts' , then $ts' \geq ts$.
- A3.** If a read ρ returns ts' , then either a write $\omega(ts')$ precedes ρ , i.e. $\omega(ts') \rightarrow \rho$, or $\omega(ts')$ is concurrent with ρ .
- A4.** If ρ_1 and ρ_2 are two read operations such that $\rho_1 \rightarrow \rho_2$ and ρ_1 returns ts_1 , then ρ_2 returns $ts_2 \geq ts_1$.

Due to space limitations and due to the similarity of the writer and server protocols to the ones used in CCFast, we omit some of the proofs and we present them in the optional Appendix, or we refer the reader to specific lemmas presented in [4]. Properties **A1** and **A3**

can be extracted easily from the algorithm. Now let us prove an important lemma about the timestamp returned by a server process:

► **Lemma 1.** *In any execution ξ of the algorithm, if a server s receives a timestamp ts at time T from a process p , then s replies with a timestamp $ts' \geq ts$ at any time $T' > T$.*

To following lemma shows **A2**:

► **Lemma 2.** *In any execution ξ of the algorithm, if a read ρ from r_1 succeeds a write operation ω that writes timestamp ts from the writer w , i.e. $\omega \rightarrow \rho$, and returns a timestamp ts' , then $ts' \geq ts$.*

So it remains to investigate if property **A4** holds.

► **Lemma 3.** *In any execution ξ of CCHYBRID, if ρ_1 and ρ_2 are two read operations such that $\rho_1 \rightarrow \rho_2$, ρ_1 is fast satisfying the predicate for $maxTS = ts_1$, then ρ_2 receives a $maxTS = ts_2$ s.t. $ts_2 \geq ts_1$.*

► **Lemma 4.** *In any execution ξ of CCHYBRID, if ρ_1 and ρ_2 are two read operations such that $\rho_1 \rightarrow \rho_2$, and ρ_1 returns ts_1 , then ρ_2 returns $ts_2 \geq ts_1$.*

Proof. A read operation has two modes: fast and slow. Thus, we need to examine all the possible combinations of the speeds of ρ_1 and ρ_2 . There are four cases to investigate: (a) ρ_1 is fast, and ρ_2 is fast, (b) ρ_1 is fast, and ρ_2 is slow, (c) ρ_1 is slow, and ρ_2 is slow, and (d) ρ_1 is slow, and ρ_2 is fast. Let $maxTS_i$ be the maximum timestamp observed by a read ρ_i , for $i \in \{1, 2\}$, during its first phase.

Case a: In case both operations are fast then, according to CCHYBRID, either they observe $maxViews \leq \frac{|\mathcal{S}|}{f} - 2$ and $propSet = \emptyset$, or they observe an $|propSet| \geq f + 1$. If both observe $maxViews \leq \frac{|\mathcal{S}|}{f} - 2$ and check the predicate, then with the same reasoning as in [4, Lemma 8], it follows that $ts_2 \geq ts_1$.

If ρ_1 observes $|propSet| \geq f + 1$ then since ρ_2 receives replies from $|\mathcal{S}_2| = |\mathcal{S}| - f$ servers, then there exists a server $s \in propSet \cap \mathcal{S}_2$ such that s replies to both ρ_1 and ρ_2 . Since $\rho_1 \rightarrow \rho_2$, then s replies to ρ_1 before replying to ρ_2 . Since s replies with $maxTS_1$ to ρ_1 , then by Lemma 1, s replies with a timestamp $ts_s \geq maxTS_1$ to ρ_2 . So $maxTS_2 \geq ts_s$ and hence $maxTS_2 \geq maxTS_1$. If $maxTS_2 = maxTS_1$ then s will reply with $ts_s = maxTS_1$ and $prop = True$. In this case ρ_2 will return $ts_2 = maxTS_1 = ts_1$. If $maxTS_2 > maxTS_1$ then ρ_2 returns either $maxTS_2$ or $maxTS_2 - 1$ and thus $ts_2 \geq ts_1$.

It remains to examine the case where ρ_1 observes $maxViews \leq \frac{|\mathcal{S}|}{f} - 2$ and $propSet = \emptyset$, and ρ_2 observes $|propSet| \geq f + 1$. If the predicate holds for ρ_1 then by Lemma 3, ρ_2 observes $maxTS_2 \geq maxTS_1$. Since ρ_2 observes $|propSet| \geq f + 1$ then it returns $ts_2 = maxTS_2$, and thus $ts_2 \geq ts_1$. If the predicate does not hold for ρ_1 then we know that the write operation propagating $maxTS_1 - 1$ completed before or during ρ_1 . Since $\rho_1 \rightarrow \rho_2$ then this write completed before ρ_2 as well. Thus, by **A2**, ρ_2 observes $maxTS_2 \geq maxTS_1 - 1$. Since ρ_2 observes $|propSet| \geq f + 1$, then it returns $ts_2 = maxTS_2 \Rightarrow ts_2 \geq maxTS_1 - 1 \Rightarrow ts_2 \geq ts_1$.

Case b: Since ρ_1 in this case is fast then ρ_1 returns either: (i) $maxTS_1 - 1$, or (ii) $maxTS_1$.

In (i), since ρ_1 observed $maxTS_1$ and since we have a single writer, it follows that the write operation that wrote timestamp $maxTS_1 - 1$, say ω_1 , proceeds or is concurrent to ρ_1 , and completes before the response step of ρ_1 . Since $\rho_1 \rightarrow \rho_2$, then $\omega_1 \rightarrow \rho_2$. Since ρ_2 is slow, then it returns the maximum timestamp it observes, i.e. $ts_2 = maxTS_2$. Moreover, since $\omega_1 \rightarrow \rho_2$, and since both operations wait for $|\mathcal{S}| - f$ replies, then according to our failure model, there exist at least a single server s that replies to both operations, first to ω_1 and

then to ρ_2 . According to Lemma 1, s sends a timestamp $ts_s \geq \max TS_1 - 1$ to ρ_2 . Thus, $\max TS_2 \geq \max TS_1 - 1$, and therefore $ts_2 \geq ts_1$.

In (ii) it follows that either the predicate holds for ρ_1 , or ρ_1 observes $|propSet| \geq f + 1$. Since ρ_2 is slow and returns $ts_2 = \max TS_2$, then by Lemma 3 and with similar reasoning as in Case (a) for when ρ_1 observes $|propSet| \geq f + 1$, we can show that $\max TS_2 \geq \max TS_1$ and hence $ts_2 \geq ts_1$.

Case c: The case where both reads are slow is simple and resembles the behavior of the reads in ABD [2]. Here each read ρ_i , for $i \in [1, 2]$, returns $\max TS_i$ and before completing it propagates $\max TS_i$ to $|\mathcal{S}| - f$ servers. Thus, ρ_1 returns $ts_1 = \max TS_1$, and before completing propagates $\max TS_1$ to $|P_1| = |\mathcal{S}| - f$ servers. Since $\rho_1 \rightarrow \rho_2$, and since ρ_2 receives $|S_2| = |\mathcal{S}| - f$ replies, then it is going to receive a timestamp $ts_s \geq \max TS_1$ from at least a single server $s \in P_1 \cap S_2$. Thus, ρ_2 returns $ts_2 = \max TS_2 \geq \max TS_1$, and $ts_2 \geq ts_1$.

Case d: So it remains to investigate the case where ρ_1 is *slow* and ρ_2 is *fast*. Observe that this case is possible when a server s is “saturated” by concurrent reads (more than $\frac{|\mathcal{S}|}{f} - 2$) and s replies to ρ_1 but does not reply to ρ_2 . Now we have two cases to investigate: either ρ_2 observes $\max TS_2 \geq \max TS_1$, or $\max TS_2 = \max TS_1 - 1$. If ρ_2 observes a $\max TS_2 \geq \max TS_1$, it may either return $ts_2 = \max TS_2$ or $ts_2 = \max TS_2 - 1$. In either case $ts_2 \geq \max TS_1 - 1 \Rightarrow ts_2 \geq ts_1$.

Let us examine now the case where $\max TS_2 = \max TS_1 - 1$. Since ρ_1 is slow and returns $\max TS_1 - 1$, then before completing it propagates $\max TS_1 - 1$ to $|\mathcal{S}| - f$ servers. Let P_1 be the set of servers that received the messages and replied to the second phase of ρ_1 . Moreover, $|S_2| = |\mathcal{S}| - f$ are the servers that received messages and replied to ρ_2 . So by Lemma 1, every server $s \in P_1 \cap S_2$ replies to both ρ_1 and then to ρ_2 , with a timestamp $ts_s \geq \max TS_1 - 1$. In addition s sets $prop = True$ before replying to ρ_1 . Since $\max TS_2 = \max TS_1 - 1$, then s replies with $ts_s = \max TS_1 - 1$ to ρ_2 , and thus the $propSet$ contains at least s in ρ_2 . According to the algorithm ρ_2 returns $ts_2 = \max TS_2$ in this case and hence $ts_2 \geq ts_1$. ◀

► **Theorem 5.** *Algorithm CCHYBRID implements a SWMR atomic read/write register.*

5 Algorithm OHFAST: Switching from One to One and a Half Rounds

Similar to algorithm CCHYBRID, OHFAST aims to allow unbounded number of readers to participate in the service while allowing operations to complete in one round. In contrast to the classic approach of the two rounds per read operation, OHFAST tries to further reduce the communication required by *slow* reads. Thus OHFAST combines ideas from CCFast and the *one and a half round* approach suggested by OHSAM. With server to server communication, OHFAST is expected to perform better in environments where the servers communicate via high capacity links, e.g., data centers.

Like in OHSAM, servers assume the responsibility of propagating the value of the timestamp instead of the reader. Similarly, in OHFAST we move the decision on a slow read to the servers. In particular, the servers record the processes that requested their timestamp. If the recording set becomes “large” then a server relays a read to the other servers before replying to the reader. However, there is a major departure from OHSAM: the servers that receive relay messages do not broadcast relays to all the servers but just to the servers that send them a relay. So, only a single server may relay for a read operation keeping the message complexity of the algorithm low in cases of low contention. When a server that relays a timestamp gets appropriate relays from the other servers, it marks the timestamp as *secured*, and sends a reply to the reader. When now the reader receives the replies from $|\mathcal{S}| - f$ servers it collects the messages with the highest timestamp. If there is a server that

declares this timestamp as *secured* then the read immediately returns the value associated with this timestamp; otherwise the reader evaluates the predicate of CCFast on the replies to determine the value to return.

Algorithm 2 Read protocol of algorithm OHFAST

```

1: at each reader  $r_i$ 
2: Components:
3:  $ts \in \mathbb{N}^+$ ;  $maxTS \in \mathbb{N}^+$ ;  $v, vp \in V$ ;  $rcounter \in \mathbb{N}^+$ 
4:  $srvAck \subseteq \mathcal{S} \times M$ 
5: Initialization:
6:  $ts \leftarrow 0$ ,  $maxTS \leftarrow 0$ ,  $v \leftarrow \perp$ ,  $vp \leftarrow \perp$ ;  $rcounter \leftarrow 0$ 
7: function READ( )
8:    $rcounter \leftarrow rcounter + 1$ 
9:   send(( $ts, v, vp$ ),  $r_i$ ,  $rcounter$ ) to all servers
10:  wait until  $|\mathcal{S}| - f$  servers reply
    ▷ Collect the ( $sid, (ts', v', vp')$ ,  $views, secured$ ) msgs in  $srvAck$ 
11:   $maxTS \leftarrow \max\{m.ts' \mid (s, m) \in srvAck\}$ 
12:   $maxAck \leftarrow \{(s, m) \mid (s, m) \in srvAck \wedge m.ts' = maxTS\}$ 
13:   $(ts, v, vp) \leftarrow m.(ts', v', vp')$  for  $(s, m) \in maxAck$ 
14:   $maxViews \leftarrow \max\{m.views \mid (s, m) \in maxAck\}$ 
15:  if  $\exists (s, m) \in maxAck$  s.t.  $m.secured = True$  then
16:    return( $v$ )
17:  else if  $\exists \alpha \in [1, \frac{|\mathcal{S}|}{f} - 2]$  s.t.
18:     $MS = \{s : (s, m) \in maxAck \wedge m.views \geq \alpha\}$  and
19:     $|MS| \geq |\mathcal{S}| - \alpha f$  then
20:    return( $v$ )
21:  else
22:    return( $vp$ )
23:  end if
24: end function

```

Algorithms 2 and 3 provide the formal pseudocode of OHFAST. We omit the write protocol as it is the same to the one presented for CCHYBRID. The read protocol in OHFAST (Algorithm 2) is simpler than the read of CCHYBRID. The reader sends messages to all the servers and waits for $|\mathcal{S}| - f$ of them to reply (L9). Once those replies are received the reader discovers the maximum timestamp $maxTS$ among the replies (L11), and collects the messages that contain $maxTS$ (L12)¹ in the set $maxAck$. If some message in $maxAck$ indicates that $maxTS$ is secured, i.e. it contains $secured = True$ (L15), then the reader returns $maxTS$. Otherwise, it evaluates the predicate on the messages in $maxAck$ (L19) to determine which timestamp to return.

Algorithm 3 Server protocol of algorithm OHFAST

```

1: at each server  $s_j$ 
2: Components:
3:  $ts \in \mathbb{N}^+$ ;  $seen \subseteq \mathcal{R} \cup \{w\}$ ;  $v, vp \in V$ ;  $Counter[|\mathcal{R}|+1] \in \mathbb{N}^+$ 
4:  $scounter \in \mathbb{N}^+$ ;  $securedts \in \{True, False\}$ 
5:  $Relays[|\mathcal{R}|] \in \mathbb{N}^+$ 
6: Initialization:
7:  $ts \leftarrow 0$ ;  $seen \leftarrow \emptyset$ ;  $v, vp \leftarrow \perp$ ;  $prop \leftarrow False$ 
8:  $Counter[i] \leftarrow 0$  for  $i \in \mathcal{R} \cup \{w\}$ ;  $scounter \leftarrow 0$ 
9:  $Relays[i] \leftarrow 0$ ;  $securedts \leftarrow False$ 
10: function RCV(( $ts', v', vp'$ ),  $q$ ,  $counter$ )
    ▷ Called upon reception of a READ/WRITE message
11: if  $Counter[q] < counter$  then
12:   if  $ts' > ts$  then
13:      $(ts, v, vp) \leftarrow (ts', v', vp')$ ;  $seen \leftarrow \{q\}$ 
14:      $securedts \leftarrow False$ 
15:   else
16:      $seen \leftarrow seen \cup \{q\}$ 
17:   end if
18:   if  $q \in \mathcal{R}$  and  $|seen| > \frac{|\mathcal{S}|}{f} - 2$  and
19:      $securedts = False$  and  $Relays[q] < ts$  then
20:      $scounter \leftarrow scounter + 1$ 
21:     sendRelay(( $ts, v, vp$ ),  $q, s_j, counter, scounter$ )
22:     to all the servers
23:      $Relays[q] \leftarrow ts$ ;  $srvRelay \leftarrow \emptyset$ 
24:   else
25:     send(( $ts, v, vp$ ),  $|seen|$ ,  $counter, securedts$ ) to  $q$ 
26:   end if
27: end if
28: end function
29: function RCVRELAY(( $ts', v', vp'$ ),  $q, s, c1, c2$ )
    ▷ Called upon reception of a RELAY message
30: if  $Counter[s] < c2$  then
31:   if  $ts' > ts$  then
32:      $(ts, v, vp) \leftarrow (ts', v', vp')$ 
33:      $seen \leftarrow \{q\}$ 
34:   else if  $ts = ts'$  then
35:      $seen \leftarrow seen \cup \{q\}$ 
36:   end if
37:   if  $Relays[q] = ts'$  then
38:      $srvRelay \leftarrow srvRelay \cup \{s\}$ 
39:     if  $|srvRelay| = |\mathcal{S}| - f$  then
40:       if  $ts = ts'$  then
41:          $securedts \leftarrow True$ 
42:       end if
43:       send(( $ts', v', vp'$ ),  $0, c1, True$ ) to  $q$ 
44:     end if
45:   else
46:      $scounter \leftarrow scounter + 1$ 
47:     sendRelay(( $ts', v', vp'$ ),  $q, s_j, scounter$ ) to  $s$ 
48:   end if
49: end if
50: end function

```

The server protocol (Algorithm 3) is the most involved in OHFAST. The server’s state is composed of the state of the replica, the recording set $seen$, a flag $securedts$ which indicates whether a timestamp has been relayed to a majority of servers, and a $Relays$ list storing the latest timestamp the server relayed for each reader. A server s waits for read/write and relay requests. When s receives a read/write request it updates its local replica state and $seen$ set appropriately (L13-14). In case the timestamp in the request is higher than its

¹ Notice that this is another departure from OHSAM as each reader in OHSAM returns the smallest discovered timestamp.

local timestamp it also sets *securedts* flag to *False*. Then, s decides whether to relay the received timestamp or not. In particular, s relays a timestamp if (L19): (i) the sender is a reader, (ii) it sent this timestamp to more than $\frac{|\mathcal{S}|}{f} - 2$ processes, (iii) the timestamp has not already being relayed (i.e. *securedts* = *False*) and (iv) the server has not yet relayed this timestamp for the same reader. If some of these conditions does not hold then s just replies to the sender with its local timestamp (L25). Notice here that servers only relay for the readers and do not relay for the writer, as the sole writer always has the latest timestamp. In a relay message s includes its local replica state, the id of the reader that initiated the relay, and its own id. When a server s' receives a relay message from s , it first updates its local replica and *seen* set appropriately (L32-33, L35). Then s' checks if it also sent a relay with the same timestamp for the same reader (L37). If not then s' bounces the relay to s and completes (L47); otherwise s' adds s in the servers that received its relay (38). When it receives $|\mathcal{S}| - f$ relays, s' replies to the reader that initiated the relay along with the timestamp that it initially relayed (not its local timestamp) (L43). Finally, if its local timestamp is the same as the relayed timestamp, then s' also sets *securedts* = *True* (L41).

5.1 Algorithm Correctness

In order to show that OHFAST is correct we have to prove that it satisfies both termination (liveness) and atomicity (safety) properties. The omitted proofs are found in the Appendix. Termination of the write operation is easy to see as according to our failure model $|\mathcal{S}| - f$ servers do not fail and can receive and reply to the write request. However, termination of the read protocol is not straightforward: a server may communicate with other servers before responding to a reader. The next lemma shows that all the read operations terminate.

► **Lemma 6.** *In any execution ξ of OHFAST, every read operation ρ invoked by a correct process r eventually terminates.*

Next it remains to show that atomicity is preserved. To prove atomicity we are going to use the four properties that express atomicity in terms of timestamps written and returned, as presented in Section 4.1. It is easy to see from the algorithm, that every process updates its local replica only when a value with a higher timestamp is received. Thus, it can be easily seen that the algorithm satisfies properties **A1** and **A3**. Notice also that when a server receives a timestamp ts then it attaches a timestamp $ts_s \geq ts$ to any message it sends from that point onward. This can be shown with similar statements as in Lemma 1. We need to show that when a server receives a *relay* that contains a timestamp ts then it sends a timestamp $ts_s \geq ts$ from that point onward.

► **Lemma 7.** *In any execution ξ of OHFAST, if a server s receives a relay with a timestamp ts at time T from a server s' , then s attaches a timestamp $ts' \geq ts$ to any message it sends at any time $T' > T$.*

Now we can show that if a read operation succeeds a write operation, then it returns a value at least as recent as the one written. This shows the validity of property **A2**.

► **Lemma 8.** *In any execution ξ of the algorithm, if a read ρ from r succeeds a write operation ω that writes timestamp ts_ω from the writer w , i.e. $\omega \rightarrow \rho$, and returns a timestamp ts_ρ , then $ts_\rho \geq ts_\omega$.*

Finally, it remains to investigate if property **A4** holds. Before we do so, we prove a lemma showing that if a timestamp ts is secured from a server s , then at least $|\mathcal{S}| - f$ servers have a timestamp $ts' > ts$.

► **Lemma 9.** *In any execution ξ of OHFAST, if a server s sets $\text{securedts} = \text{True}$ for a timestamp ts at time T then $\exists \mathcal{S}' \subseteq \mathcal{S}$ at T , s.t. $|\mathcal{S}'| \geq |\mathcal{S}| - f$ and $\forall s' \in \mathcal{S}'$, the local timestamp of s' is $ts' \geq ts$.*

► **Lemma 10.** *In any execution ξ of OHFAST, if ρ_1 and ρ_2 are two read operations such that $\rho_1 \rightarrow \rho_2$, and ρ_1 returns ts_{ρ_1} , then ρ_2 returns $ts_{\rho_2} \geq ts_{\rho_1}$.*

Proof. A read operation may decide on the value to return in two ways in OHFAST: (i) it receives a secured timestamp, or (ii) it evaluates the predicate. Let us first examine what happens when the two reads are invoked by the same reader (i.e. $r_1 = r_2$). During ρ_2 , r_1 includes a timestamp $ts_{r_1} \geq ts_{\rho_1}$ in every message it sends to servers. According to Lemma 1 every server s replies with a timestamp $ts_s \geq ts_{\rho_1}$. Thus, $\text{maxTS}_2 \geq ts_{\rho_1}$. If $\text{maxTS}_2 > ts_{\rho_1}$ then since $ts_{\rho_2} = \text{maxTS}_2$ or $ts_{\rho_2} = \text{maxTS}_2 - 1$ it follows that $ts_{\rho_2} \geq ts_{\rho_1}$ in either case. If $\text{maxTS}_2 = ts_{\rho_1}$ then every server adds r_1 in their *seen* set before replying to ρ_2 . So the predicate is valid for $|\mathcal{MS}| \geq |\mathcal{S}| - f$ and $\alpha = 1$. Hence, ρ_2 returns $ts_{\rho_2} = \text{maxTS}_2 = ts_{\rho_1}$ in any case (i) or (ii).

So we need now to examine all the possible combinations for the two reads ρ_1 and ρ_2 when $r_1 \neq r_2$. If both read operations examine the predicate to decide on the value to return (i.e., they do not receive a secured timestamp), then with same reasoning as in [4, Lemma 8] we can show that atomicity is preserved. So it remains to examine the following three cases: (1) ρ_1 evaluates the predicate, and ρ_2 receives a secured maxTS_2 , (2) ρ_1 receives a secured maxTS_1 , and ρ_2 evaluates the predicate, and (3) ρ_1 receives a secured maxTS_1 , and ρ_2 receives a secured maxTS_2 .

Case 1: In this case, ρ_1 evaluates the predicate, and ρ_2 returns $ts_{\rho_2} = \text{maxTS}_2$ as it received a reply with maxTS_2 and $\text{secured} = \text{True}$. There are two subcases to examine: (a) ρ_1 returns maxTS_1 , and (b) ρ_1 returns $\text{maxTS}_1 - 1$.

Case 1a: If ρ_1 returns maxTS_1 it follows that the predicate is valid for ρ_1 . Hence:

$$\exists \alpha \in [1, \frac{|\mathcal{S}|}{f} - 2] \text{ and } \mathcal{MS} \subseteq \mathcal{S} \text{ s.t. } \mathcal{MS} = \{s : s.ts = \text{maxTS}_1 \wedge s.views \geq \alpha\} \wedge |\mathcal{MS}| \geq |\mathcal{S}| - \alpha f$$

Moreover, since ρ_1 examines the predicate, then none of the servers that replied with maxTS_1 sends $\text{secured} = \text{True}$. Therefore, $\forall s \in \mathcal{MS}$, it must be true that $s.views \leq \frac{\mathcal{S}}{f} - 2$ before replying to ρ_1 (L16), otherwise s would proceed to relay and secure maxTS_1 . Since every $s.views \leq \frac{\mathcal{S}}{f} - 2$, then it must be the case that $\alpha \leq \frac{\mathcal{S}}{f} - 2$ as well. Thus substituting:

$$|\mathcal{MS}| \geq |\mathcal{S}| - \alpha f \Rightarrow |\mathcal{MS}| \geq |\mathcal{S}| - (\frac{\mathcal{S}}{f} - 2)f \Rightarrow |\mathcal{MS}| > f$$

Since ρ_2 receives replies from $|\mathcal{S}_2| = |\mathcal{S}| - f$ servers then $\mathcal{S}_2 \cap \mathcal{MS} \neq \emptyset$. Also notice that since $\rho_1 \rightarrow \rho_2$, then a server $s \in \mathcal{S}_2 \cap \mathcal{MS}$ replies to ρ_1 with maxTS_1 before replying to ρ_2 . By Lemma 1, s replies to ρ_2 with a timestamp $ts_s \geq \text{maxTS}_1$. Thus, $\text{maxTS}_2 \geq ts_s \Rightarrow \text{maxTS}_2 \geq \text{maxTS}_1$ and ρ_2 returns $ts_{\rho_2} \geq \text{maxTS}_1 \Rightarrow ts_{\rho_2} \geq ts_{\rho_1}$.

Case 1b: Assume now the case where ρ_1 returns $\text{maxTS}_1 - 1$. Since ρ_1 received maxTS_1 , and since the sole writer invokes one operation at a time, then it follows that the write operation that wrote $\text{maxTS}_1 - 1$, say ω , completed during or before ρ_1 . Since though $\rho_1 \rightarrow \rho_2$, then it follows that $\omega \rightarrow \rho_2$. Since ω communicates with $|\mathcal{S}| - f$ servers before completing, and since ρ_2 waits for $|\mathcal{S}| - f$ replies, then there is a server s that replies to ω before replying to ρ_2 . By Lemma 1, s replies with a timestamp $ts_s \geq \text{maxTS}_1 - 1$ to ρ_2 . Thus ρ_2 observes a $\text{maxTS}_2 \geq \text{maxTS}_1 - 1$, and hence $ts_{\rho_2} \geq \text{maxTS}_1 - 1 \Rightarrow ts_{\rho_2} \geq ts_{\rho_1}$ in this case as well.

Case 2: Here, ρ_1 returns $ts_{\rho_1} = \max TS_1$ as it received a message that contained $\max TS_1$ and $\text{secured} = \text{True}$. Read ρ_2 evaluates the predicate to decide on the value to return. We have two subcases to examine again: (a) ρ_2 returns $\max TS_2$, or (b) ρ_2 returns $\max TS_2 - 1$. Since ρ_1 returned a secured timestamp, then it received $\max TS_1$ and $\text{secured} = \text{True}$ from some server s . By Lemma 9, a set $|\mathcal{S}'| \geq |\mathcal{S}| - f$ of servers have a timestamp $ts' \geq \max TS_1$ before s replies to ρ_1 . Since ρ_2 receives replies from $|\mathcal{S}_2| = |\mathcal{S}| - f$ servers, then $\mathcal{S}' \cap \mathcal{S}_2 \neq \emptyset$. Then by Lemmas 1 and 7, any server in $s' \in \mathcal{S}' \cap \mathcal{S}_2$ replies to ρ_2 with a timestamp $ts_{s'} \geq \max TS_1$. Thus, ρ_2 observes a $\max TS_2 \geq \max TS_1$. If $\max TS_2 > \max TS_1$ and since ρ_2 returns either $\max TS_2$ or $\max TS_2 - 1$, then in either case $ts_{\rho_2} \geq ts_{\rho_1}$.

So it remains to examine what happens when $\max TS_2 = \max TS_1$. If ρ_2 returns $ts_{\rho_2} = \max TS_2$ then $ts_{\rho_2} \geq ts_{\rho_1}$. Let us examine now if ρ_2 may return $\max TS_2 - 1$. As we said before every server s' in $\mathcal{S}' \cap \mathcal{S}_2$ replies with $ts_{s'} \geq \max TS_1$ to ρ_2 . Since $|\mathcal{S}'| \geq |\mathcal{S}| - f$ and $|\mathcal{S}_2| \geq |\mathcal{S}| - f$ then $|\mathcal{S}' \cap \mathcal{S}_2| \geq |\mathcal{S}| - 2f$. Also by the algorithm, every server in \mathcal{S}' adds r_1 in its *seen* set before replying to the relay message from s (L39). Furthermore, every server in \mathcal{S}_2 adds r_2 in its *seen* set before replying to ρ_2 . So every server $s' \in \mathcal{S}' \cap \mathcal{S}_2$ replies with a $s.\text{views} \geq 2$. Thus, the predicate holds for at least $|\mathcal{MS}| = |\mathcal{S}' \cap \mathcal{S}_2| \geq |\mathcal{S}| - 2f$ and $\alpha = 2$. Hence ρ_2 will return $\max TS_2$ contradicting our assumption that returns $\max TS_2 - 1$. So returning $\max TS_2 - 1$ is not possible.

Case 3: In this case both ρ_1 and ρ_2 return a secured timestamp. Let s_1 be the server that send $\max TS_1$ and $\text{secured} = \text{True}$ to ρ_1 , and s_2 (not necessarily different than s_1) be the server that sent $\max TS_2$ and $\text{secured} = \text{True}$ to ρ_2 . By Lemma 9, there exists a set \mathcal{S}' s.t. every server $s \in \mathcal{S}'$ has a timestamp $ts_s \geq \max TS_1$ before s_1 replies to ρ_1 . As explained in Case 2, $\mathcal{S}' \cap \mathcal{S}_2 \neq \emptyset$. Hence there exists a server that replied both to the relay message of s_1 and to ρ_2 . By Lemma 7, each server $s' \in \mathcal{S}' \cap \mathcal{S}_2$ replies to ρ_2 with a timestamp $ts_{s'} \geq \max TS_1$. Hence, $\max TS_2 \geq \max TS_1$. Since ρ_2 returns a secured timestamp, then it returns $\max TS_2$. Therefore, $ts_{\rho_2} = \max TS_2 \Rightarrow ts_{\rho_2} \geq \max TS_1 \Rightarrow ts_{\rho_2} \geq ts_{\rho_1}$. ◀

► **Theorem 11.** *Algorithm OHFAST implements a SWMR atomic read/write register.*

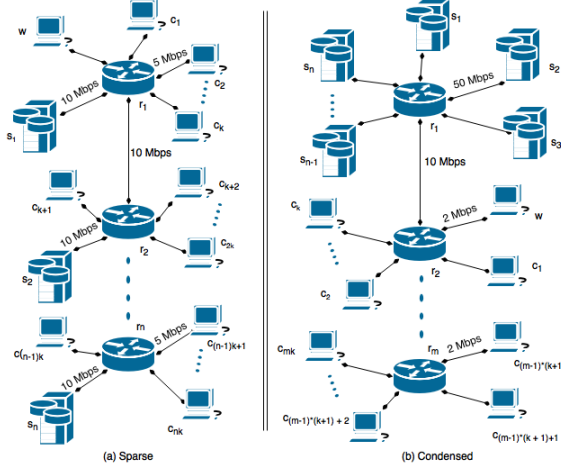
6 Empirical Results

In this section, we present empirical results that we obtained by implementing algorithms ABD [2], OHSAM [7], SF [6], CCHYBRID, and OHFAST, using the NS3 discrete event simulator [1]. NS3 is a highly customizable and extensible simulator that allows us to gain full control over the event scheduler and the deployment environment. Thus, it allows us to investigate the exact parameters that may affect the performance of our algorithms.

Experimentation Platform. The general testbed of our experiments consists of a single writer, a set of readers, and a set servers. We assume that $f = 1$ servers may fail. This assumption was chosen so as every operation would wait for all but one servers to reply, inflicting that way high concurrency and potentially inconsistency in our system. Communication between the nodes is established via point to point bidirectional links implemented with a DropTail queue. For the purpose of the experimental evaluation, we developed simulations representing two different topologies, *Sparse* and *Condensed*, which mainly differ on the deployment of server nodes.

Figure 1 presents the two topologies. In both topologies the clients are divided evenly and are connected on a series of router nodes. Clients are connected to the routers with 5Mbps links and 2ms delay, and routers are connected with 10Mbps links and 4ms delay. In the *Sparse* topology (Figure 1(a)), a server is connected to each router with 10Mbps bandwidth and 2ms delay. This topology demonstrates a network where servers are separated and

appear to be in different networks. In the *Condensed* topology (Figure 1(b)) all the servers are connected to a single router with 50Mbps links and 2ms delay, simulating a network where servers are connected in close proximity and with high bandwidth links (e.g., a datacenter).



■ **Figure 1** Simulated topologies.

(a) the simulation clock, and (b) a real time clock. The simulation clock was able to estimate the communication time, while the real clock allowed us obtain the time taken by the computation at each operation. The latency is calculated adding both times.

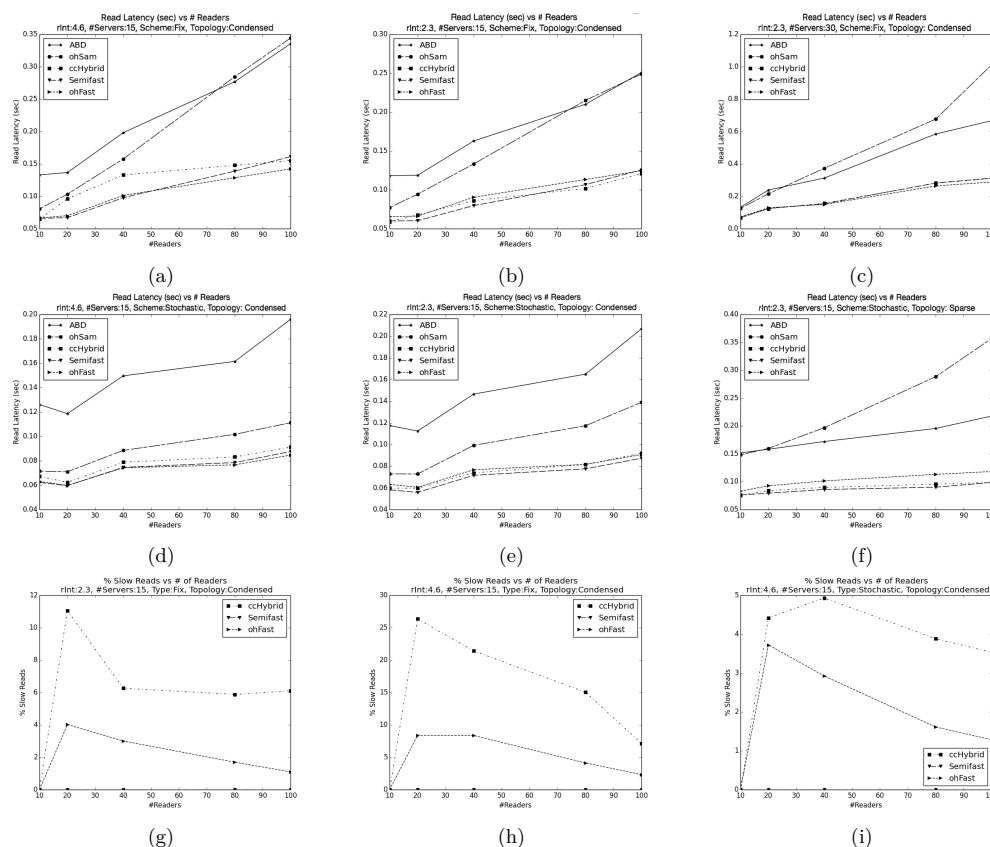
Scenarios. Measurements of the performance involves multiple execution scenarios. The scenarios were designed to test (i) the scalability of the algorithms as the number of readers and servers increases; (ii) the contention effect on efficiency, by running different concurrency scenarios; and (iii) the relation of the efficiency with the topology of the network that we use. To test scalability we range the number of readers $|\mathcal{R}| \in [10, 20, 40, 80, 100]$ and the number of servers $|\mathcal{S}| \in [10, 15, 20, 25, 30]$. To test contention we specify the frequency of read operation and we run our algorithm for different read intervals ($rInt \in [2.3, 4.6, 6.9]$ seconds). We issue write operations every 4 seconds. In addition, we define two read invocation schemes: (i) *fix* and (ii) *stochastic*. In the fix scheme all the read operations are scheduled periodically at the read interval. In the stochastic scheme each operation is scheduled at random between 1s and $rInt$ seconds in each read interval. Finally, to test topological effects we run our algorithms using both the *Sparse* and *Condensed* topologies.

Results. As a general observation, the new algorithms outperform all the other algorithms in most scenarios. In particular, it is clear that CCHYBRID and OHFAST outperform algorithms ABD and OHSAM. In addition, the two algorithms appear to achieve similar operation latencies as SF. A closer examination reveals that in many scenarios SF does not perform any slow reads, whereas in the same executions both CCHYBRID and OHFAST require some slow reads. The fact that the two algorithms perform the same as SF, despite the slow reads, demonstrates that the computation overhead of the two presented algorithms is much less than the computation needed by SF. Thus, in executions where SF will perform more slow operations, clearly this will result in even worse operation latencies. More in detail, taking our tests one by one we conclude to the following observations:

Scalability: As can be seen in Figures 2(b) and (c), the increasing number of readers and the servers have a negative impact on all the algorithms. The impact is higher on ABD and OHSAM, and lower for the rest of the algorithms.

We ran NS3 on a Macintosh machine running OS X El Capitan, with 2.5Ghz Intel Core i7 processor and 16GB of RAM. The average of 5 samples per scenario provided the stated operation latencies.

Performance. The performance of the algorithms is measured in terms of the ratio of the number of fast over slow R/W operations - *communication burden*; and the total time it takes for an operation to complete - *operation latency*. Operation latency is affected by both communication and computation latencies. As NS3 only provides simulated time events and omits any computation, we combined two clocks:



■ **Figure 2** Experimental Results from NS3 Simulation

Contention: Contention is generated by: (i) operation frequencies, and (ii) concurrency schemes. We observe that *operation frequency* affects the latency of the operations in the *fix* scheme. This can be seen in Figure 2(a) and (b). Algorithms ABD and OHSAM are not affected (as all of their reads are slow), but the multi-speed algorithms SF, CCHYBRID and OHFAST, are affected negatively. This behaviour is due to the fact that these algorithms perform a slow read operation per write operation. When the read interval is close to the write interval, e.g., $rInt = 4.6$, most of the reads are concurrent to the write and thus more reads are slow (Figure 2(h)). This is not the case when $rInt = 2.3$ (Figure 2(g)). Notice that the same behavior is not being observed when a *stochastic* scheme is used, as randomness prevents the operations to be invoked at exactly the same time (Figure 2(d) and (e)). Hence, a slow read operation may complete before any read operations that return the same value are invoked. Therefore, according to the multi-speed algorithms, once a slow read is completed, any read operation that succeeds such a read will be fast. This results in a low percentage of slow reads, as shown in Figure 2(i).

Finally, when the operation frequency is constant, it appears that in the *stochastic* scheme each operation completes almost two times faster than in the *fix* scheme (Figure 2(b) and 2(e)). Algorithms, ABD and OHSAM, can be used as points of reference as they have the same computation and communication requirements in both *fix* and *stochastic* scenarios. The difference can be explained due to the congestion that the *fix* scheme introduces in the network. On the contrary, a *stochastic* scheme distributes the invocation time intervals of the read operations uniformly, reducing the network congestion, and hence operation latency.

Topology: Plots 2(e) and 2(f) show that topology has an impact on the performance and the efficiency of all the algorithms. Most importantly, we can observe that OHSAM and OHFAST are the two algorithms that are affected the most. In particular, while in (e) OHSAM performs better than ABD and OHFAST performs similar to CCHYBRID and SF we notice that in (f) OHSAM performs worse than ABD and OHFAST worse than the 2 others. This behaviour is expected as both OHSAM and OHFAST need to exchange messages between the servers during a relay phase. However, notice that OHFAST performs much better since operation relays are not performed for every read operation.

7 Conclusions

In this paper we present two new algorithms CCHYBRID and OHFAST that implement atomic SWMR register in a message-passing, asynchronous environment. Both algorithms use the predicate introduced in [4], to achieve *single round* reads with small computational footprint. However, to avoid constraints in reader participation both algorithms allow some reads to be *slow*. In CCHYBRID the reader decides on the speed of its read operation, resulting in operations that perform *1 or 2 rounds*. OHFAST moves the decision of slow operations to the servers, enabling *1 or 1.5 round* operations. Simulation results show that our algorithms outperform all slow operation algorithms, as well as “multi-speed” implementations that have high computation demands. We claim that our developments take us closer to *practical* implementations of atomic read/write objects in the message-passing environment.

References

- 1 NS3 network simulator. <https://www.nsnam.org/>.
- 2 H. Attiya, A. Bar-Noy, and D. Dolev. Sharing memory robustly in message passing systems. *Journal of the ACM*, 42(1):124–142, 1996.
- 3 P. Dutta, R. Guerraoui, R. R. Levy, and A. Chakraborty. How fast can a distributed atomic read be? In *PODC’04*, pages 236–245, 2004.
- 4 A. Fernández Anta, N. Nicolaou, and A. Popa. Making “fast” atomic operations computationally tractable. In *OPODIS’15*, 2015.
- 5 C. Georgiou, N. C. Nicolaou, and A. A. Shvartsman. On the robustness of (semi) fast quorum-based implementations of atomic shared memory. In *DISC’08*, pp. 289–304, 2008.
- 6 C. Georgiou, N. C. Nicolaou, and A. A. Shvartsman. Fault-tolerant semifast implementations of atomic read/write registers. *Journal of Parallel and Distributed Computing*, 69(1):62–79, 2009.
- 7 T. Hadjistasi, N. Nicolaou, and A. A. Schwarzmann. Brief announcement: Oh-ram! one and a half round read/write atomic memory. In *PODC ’16*, pages 353–355, NY, USA, 2016.
- 8 M. P. Herlihy and J. M. Wing. Linearizability: a correctness condition for concurrent objects. *ACM Trans. on Progr. Languages and Systems (TOPLAS)*, 12(3):463–492, 1990.
- 9 L. Lamport. How to make a multiprocessor computer that correctly executes multiprocess program. *IEEE Transactions on Computers*, 28(9):690–691, 1979.
- 10 N. A. Lynch. *Distributed Algorithms*. Morgan Kaufmann Publishers, 1996.
- 11 N. A. Lynch and A. A. Shvartsman. Robust emulation of shared memory using dynamic quorum-acknowledged broadcasts. In *Proceedings of Symposium on Fault-Tolerant Computing*, pages 272–281, 1997.

A Missing Proofs

Proof of Lemma 1. If the local timestamp of the server s , ts_s , is smaller than ts , then $ts_s = ts$. Otherwise ts_s does not change and remains $ts_s \geq ts$. In any case s replies with a timestamp $ts_s \geq ts$ to π . Since the timestamp of s is monotonically incrementing, then s attaches a timestamp $ts' \geq ts_s$, and hence $ts' \geq ts$, to any subsequent reply. ◀

Proof of Lemma 2. If a read operation ρ is fast due to the predicate validation and succeeds a write operation ω , i.e. $\omega \rightarrow \rho$, then by [4, Lemma 7], the read operation returns a timestamp $ts' \geq ts$, where ts is the timestamp written by ω . On the other hand, if ρ is slow or observed $prop = True$ in more than $f + 1$ servers, then it returns $maxTS$. Since $|\mathcal{S}| - f$ servers received ω , and since ρ contacts $|\mathcal{S}| - f$ servers during its first phase, then there is at least a single server, say s that received the message for ω before replying to ρ . According to Lemma 1, s replies to ρ with a timestamp $ts_s \geq ts_\omega$, the timestamp it received from ω . Thus, ρ observes a $maxTS \geq ts_s \geq ts_\omega$, and hence returns $ts' = maxTS \geq ts_\omega$. ◀

Proof of Lemma 3. Since the predicate holds for ρ_1 , hence there exists an $\alpha \in [1, \frac{|\mathcal{S}|}{f} - 2]$, and $MS_1 \subseteq \mathcal{S}$ s.t. $|MS_1| = |\mathcal{S}| - \alpha f$, and $\forall s \in MS_1, s.ts = ts_1$ and $s.views \geq \alpha$. Performing the substitutions follows that:

$$|MS_1| \geq |\mathcal{S}| - \left(\frac{|\mathcal{S}|}{f} - 2\right)f \Rightarrow |MS_1| > f$$

Since ρ_2 receives replies from $|S_2| = |\mathcal{S}| - f$ servers, then there exists a server $s \in MS_1 \cap S_2$ that replies to both ρ_1 and ρ_2 . Since $\rho_1 \rightarrow \rho_2$ then s replies to ρ_1 before replying to ρ_2 . Since s replies with ts_1 to ρ_1 , then according to Lemma 1, it replies with a timestamp $ts_s \geq ts_1$ to ρ_2 . Thus, ρ_2 observes a timestamp $maxTS \geq ts_1$ and hence $ts_2 \geq ts_1$. ◀

Proof of Lemma 6. Each operation ρ sends messages to all the servers and waits for $|\mathcal{S}| - f$ replies before terminating. Thus termination of such process is prevented if less than $|\mathcal{S}| - f$ servers reply to r for operation ρ .

When a server receives a message for a read operation it may perform one of two actions: (i) replies with its local timestamp-value pair, or (ii) sends relay messages to other servers and replies to the reader when it collects $|\mathcal{S}| - f$ relays that contain its local-timestamp. So a read operation terminates if a correct server is guaranteed to send a reply to the reader in both cases. Notice that when a server s' receives a relay message from s with a timestamp ts it either, (a) sends a relay to s (L43), or (b) appends its local $srvRelay$ set with the sender if $Relays[q] = ts$ (L38). In (a) it is clear that s' replies to s with a relay that contains ts . However it is not clear if s' sends a relay message to s in (b). Notice that (b) is only possible if $Relays[q] = ts$, where ts the timestamp enclosed in the relay message. Server s' sets $Relays[q] = ts$ only when it sends relay messages for q for timestamp ts to all the servers (L22). So in that line s' sends relay message to s as well. Therefore, in any case (a) or (b), a relay message is sent by s' to s with timestamp ts . So s eventually receives $|\mathcal{S}| - f$ relays that contain ts and thus the check in Line 39 is satisfied and replies to the read operation. Thus, the reader collects replies from a server in both cases (i) and (ii). Hence, the reader receives at least $|\mathcal{S}| - f$ replies. ◀

Proof of Lemma 7. If the local timestamp of the server s , ts_s , is smaller than ts when it receives a relay message, then s updates $ts_s = ts$ (L32). Otherwise ts_s does not change and remains $ts_s \geq ts$. In any case s replies with a timestamp $ts_s \geq ts$ to s' . By monotonicity of the timestamp, s attaches a timestamp $ts' \geq ts_s$, and hence $ts' \geq ts$ to any subsequent message. ◀

Proof of Lemma 8. There are two cases to investigate: (i) ρ returns after examining the predicate, or (ii) ρ returns because it received a secured timestamp. In case (i) we can show with similar arguments with [4, Lemma 7], that $ts_\rho \geq ts_\omega$. So it remains to examine case (ii). In this case ρ received a reply that contained a timestamp $ts_s = \text{maxTS}$ and a *secured* flag equal to *True*. According to the algorithm ρ returns $ts_\rho = \text{maxTS}$. Since $|\mathcal{S}| - f$ servers received ω , and since ρ contacts $|\mathcal{S}| - f$ servers during its first phase, with $f < \frac{|\mathcal{S}|}{2}$, then there is at least a single server, say s , that received the message for ω before replying to ρ . According to Lemmas 1 and 7, s replies to ρ with a timestamp $ts_s \geq ts_\omega$, the timestamp it received from ω . Thus, ρ observes a $\text{maxTS} \geq ts_s \geq ts_\omega$, and hence returns $ts_\rho = \text{maxTS} \geq ts_\omega$. ◀

Proof of Lemma 9. This lemma follows from the way that a relay round is implemented by a server. In particular, when a server s relays a timestamp ts , it sends a message to all the servers. Each server $srvr'$ that receives such a relay replies by Lemma 6, with a timestamp $ts' = ts$. Before replying, s' either sets its timestamp to ts or has a larger timestamp. So when s sets *securedts* = *True* has received a set $|\mathcal{S}'| \geq |\mathcal{S}| - f$ of replies, and every server $s' \in \mathcal{S}'$ has a timestamp $ts' \geq ts$, by Lemma 7. Thus the lemma follows. ◀