# Optimistic Erasure-Coded Distributed Storage[*]

Partha Dutta
IBM India Research Lab
Bangalore, India

Rachid Guerraoui
EPFL IC
Lausanne, Switzerland

Ron R. Levy
EPFL IC
Lausanne, Switzerland

### Abstract

We study erasure-coded atomic register implementations in an asynchronous crash-recovery model. Erasure coding provides a cheap and space-efficient way to tolerate failures in a distributed system. This paper presents ORCAS, Optimistic eRasure-Coded Atomic Storage, which consists of two separate implementations, ORCAS-A and ORCAS-B. In terms of storage space used, ORCAS-A is more efficient in systems where we expect large number of concurrent writes, whereas, ORCAS-B is more suitable if not many writes are invoked concurrently. Compared to replication based implementations, both ORCAS implementations significantly save on the storage space. The implementations are optimistic in the sense that the used storage is lower in synchronous periods, which are considered common in practice, as compared to asynchronous periods. Indirectly, we show that tolerating asynchronous periods does not increase storage overhead during synchronous periods.

## 1 Introduction

### 1.1 Motivation

Preventing data loss in storage devices is one of the most critical requirements in any storage system. Enterprise storage systems in particular have multiple levels of redundancy built in for fault tolerance. The cost of a specialized centralized storage server is very high and yet it does not offer protection against unforseen consequences such as fires and floods. Distributed storage systems based on commodity hardware, as alternatives to their centralized counterparts, have gained in popularity since they are cheaper, can be more reliable and offer better scalability. However, implementing such systems is more complicated due to their very distributed nature.

Most existing distributed storage systems rely on data replication to provide fault tolerance [14]. Recently however, it has been argued that erasure coding is a better alternative to data replication since it reduces the cost of ensuring fault tolerance [6, 7]. In erasure-coded storage systems, instead of keeping an identical version of a data $V$ on each server, $V$ is encoded into $n$ fragments such that $V$ can be reconstructed from any set of at least $k$ fragments (called $k$-of-$n$ encoding), where the size of each fragment is roughly $|V|/k$. A different encoded fragment is stored on each of the $n$ servers, and ideally such a system can tolerate the failure of $f = n - k$ servers.

The main advantage of erasure-coded storage over replicated storage is its *storage usage*, i.e., less storage space is used to provide fault tolerance. For instance, it is well-known that a replicated

---

[*]To appear in $22^{nd}$ International Symposium on Distributed Computing (DISC 2008).

storage system with 4 servers can tolerate at most 1 failure in an asynchronous environment. If each server has a storage capacity of 1 TB, the total capacity of the replicated storage system is still 1 TB. In this case the storage overhead (total capacity/useable capacity) is 4, i.e. only 1/4 of the total capacity is available. Erasure coding allows the reduction of this overhead to 2 in an asynchronous system, i.e., makes 2 TB useable. In a synchronous system (with 4 servers and at most 1 failure), it is even possible to further reduce this overhead and make 3 TB available to the user. Clearly, the synchronous erasure-coded storage is more desirable in terms of storage usage. Unfortunately, synchrony assumptions are often not realistic in practice. Even if we expect the system to be synchronous most of the time, it is good practice to tolerate asynchronous periods. The idea underlying our contribution is the common practice of designing distributed systems that can cope with worst case conditions (e.g., asynchrony and failures) but are optimized for best case situations (e.g., synchrony and no failures) that are considered common in practice.

## 1.2  Contributions

In this paper we investigate one of the fundamental building blocks of a fault-tolerant distributed storage $-$ multi-writer multi-reader atomic register implementations [12, 3, 14]. An atomic register is a distributed data-structure that can be concurrently accessed by multiple processes and yet provide an "illusion" of a sequential register. (A sequential register is a data-structure that is accessed by a single process with read and write operations, where a read always returns the last value written.) We consider implementations over a set of $n$ server processes in an asynchronous crash-recovery message-passing system where (1) each process may crash and recover but has access to a stable storage, (2) in a run, at most $f$ out of $n$ servers are faulty (i.e., eventually crash and never recover), and (3) channels are fair-lossy.

We present two wait-free atomic register implementations ORCAS-A and ORCAS-B (Optimistic eRasure-Coded Atomic Storage). Our implementations are the first wait-free atomic register implementations in a crash-recovery model that have an "optimistic" (stable) storage usage. Suppose that all possible write values are of a fixed size $\Delta$.[1] Then in both of our implementations, during synchronous periods with $q$ alive (non-crashed) servers and when there is no write operation in progress, the stable storage used at every alive server is $\frac{\Delta}{q-f}$, whereas during asynchronous periods when there is no write operation in progress, the storage used is $\frac{\Delta}{n-2f}$ at all but $f$ servers (in ORCAS-A, at most $f$ servers may use $\Delta$). However, the two implementations differ in their storage usage when there is a write in progress. In ORCAS-A, when one or more writes are in progress, the storage used at a server can be $\Delta$, but even in the worst-case the storage used is never higher than $\Delta$. In contrast, if there are $w$ concurrent writes in progress in ORCAS-B then, in the worst-case, the storage used at a server can be $\frac{w\Delta}{n-2f}$. Thus in terms of storage space used, ORCAS-A is more efficient in systems where we expect large number of concurrent writes, whereas, ORCAS-B is more suitable if not many writes are invoked concurrently. We also show how the number of messages exchanged in ORCAS-A can be significantly reduced by weakening the termination condition of the read from wait-free to Finite Write (FW) termination [1].[2]

Both ORCAS implementations are based on a simple but effective idea. The write first "gauges" the number of alive servers in the system by sending a message to all servers and counting the

---

[1]Through out the paper we assume that, other than the write value and its encoded fragments, all other values (e.g., timestamp) at a server are of negligible size.

[2]A similar idea was earlier used in [9] where the read satisfied obstruction-free termination.

number of replies received during a short waiting period. Depending on the number of replies, the write decides how to erasure code its value. Additionally, to limit the communication overhead, the ORCAS implementations ensure that the write value or the encoded fragments are sent to the servers in only one of the phases of a write; later, the servers can locally compute the final encoded fragments on receiving a small message that specifies how the value needs to be encoded (but the message does not contain the final encoded fragments).

In particular, in ORCAS-A, the write sends the unencoded write value to all servers and waits for replies. If it receives replies from $q$ servers, the write sends a message to the servers that requests them to locally encode the received value with $(q - f)$-of-$n$ encoding. (Note that $q \geq n - f$ because at most $f$ servers can be faulty.) Roughly speaking, a subsequent read can contact at least $q - f$ of the servers that reply to the write, and (1) either the read receives an unencoded value from one of those servers, or (2) it receives $q - f$ encoded fragments. In both cases, as the write does a $(q - f)$-of-$n$ encoding, the read can reconstruct the written value. Note that, as $q - f \geq n - 2f$, in the worst-case ORCAS-A does a $(n - 2f)$-of-$n$ encoding, and in synchronous periods with $q$ alive servers, it does a $(q - f)$-of-$n$ encoding.

On the other hand, ORCAS-B, like previous erasure-coded atomic register implementations, never sends an unencoded write value to the servers. Ideally, to obtain the same storage usage as ORCAS-A, in a write of ORCAS-B we would like to send an $(n-2f)$-of-$n$ encoded fragment to each server, and on receiving replies from $q$ servers, request the servers to keep a $(q - f)$-of-$n$ encoded fragment. However, in general, it is not possible to extract a particular fragment of a $(q - f)$-of-$n$ encoding from a *single* fragment of a $(n - 2f)$-of-$n$ encoding. Thus with this naive approach, either a write would need to send another set of fragments to the servers or the servers would need to exchange their fragments, resulting in significant increase in communication overhead. We solve this problem in ORCAS-B by a novel approach of storing multiple, much smaller fragments at each server (instead of a single large one). Suppose the write estimates that there are $q$ alive servers. Then it encodes the value with $x$-of-$nz$ encoding, where $x$ is any common multiple of $n - 2f$ and $q - f$, and $z = \frac{x}{n-2f}$, and sends $z$ fragments to each server. If the write receives replies from $q$ servers, then it requests each server to *trim* its stored fragments in the next phase, i.e., retain any $y = \frac{x}{q-f}$ out of its $z$ fragments and delete the rest. Roughly speaking, since a read can miss at most $f$ servers that replied to the write, if a subsequent read sees a trimmed server then it will eventually receive $y$ fragments from at least $q - f$ servers, and if the read does not see a trimmed server, then it will receive $z$ fragments from at least $n - 2f$ servers. In both cases, it receives $y(q - f) = z(n - 2f) = x$ fragments, and therefore, can reconstruct the written value. The advantage of our approach over the naive approach is that, our approach has the same storage usage as latter, but has lower communication overhead.

In Appendix B, we show lower bounds on storage space usage in atomic register implementation in synchronous and asynchronous systems, for a specific class of implementations (that include both ORCAS-A and ORCAS-B, and the implementation in [6]). Roughly speaking, we show that implementations $-$ (1) which at the end of a write store equal number of encoded fragments in the stable storage of the servers, and (2) do not use different encoding schemes in the same operation $-$ cannot have a stable storage usage better than ORCAS-A (and ORCAS-B) in either synchronous or asynchronous periods.

## 1.3 Related Work

Recently there has been lot of work on erasure-coded distributed storage [6, 7, 2, 5, 10, 9]. We discuss below three representative papers that are close to our work.

Frolund et al. [6] describe an erasure-coded distributed storage (called FAB) in the same system model as this paper, i.e., an asynchronous crash-recovery model. The primary algorithm in FAB implements an atomic register. Servers have stable-storage and keep a log containing old versions of the data, which is periodically garbage collected. The main difference with our approach is that in FAB the stable storage is not used optimistically. In particular, ORCAS-B has the same storage overhead as FAB during asynchronous periods (even when writes are in progress) but performs better during synchronous periods. Another difference is that FAB provides strict linearizability, which ensures that partial write operations appear to take effect before the crash or not at all. The price that is paid by FAB is to give up wait-freedom: concurrent operations may abort. ORCAS-B ensures that write operations are at worst completed upon recovery of the writer and guarantees wait-freedom: all operations invoked by correct clients eventually terminate despite the concurrent invocations of other clients.

Aguilera et al. [2] present an erasure-coded storage (that we call AJX) for *synchronous systems* that is optimized for $f << n$. AJX provides the same low storage overhead as ORCAS during failure-free synchronous periods, and performs better than ORCAS when there are failures. However, AJX provides consistency guarantees of only a regular register and puts a limit on the maximum number of client failures. Also, wait-freedom is not ensured since concurrent writes may abort.

Cachin et al. [5] propose a wait-free atomic register implementation for the byzantine model. It uses a reliable broadcast like primitive to disseminate the data fragments to all servers, thus guaranteeing that if one server receives a fragment, then all do. The storage required at the servers when there is no write in progress is $\frac{\Delta}{n-f}$. At first glance, one might be tempted to compare our implementations with a crash-failure restriction of the algorithm in [5], and conclude that our implementations have worse storage requirements in asynchronous periods ($\frac{\Delta}{n-2f}$). However, one of the implications of our lower bounds in Appendix B is that there is no obvious translation of the algorithm in [5] to a crash-recovery model while maintaining the same storage usage. (We discuss this comparison further in Appendix B.)

## 2 Model and Definitions

**Processes.** We consider an asynchronous message passing model, without any assumptions on communication delay or relative process speeds. For presentation simplicity, we assume the existence of a global clock. This clock however is inaccessible to the servers and clients.

The set of servers is denoted by $S$ and $|S| = n$. The $j^{th}$ server is denoted by $s_j$, $1 \leq j \leq n$. The set of clients is denoted by $C$ and it is bounded in size. Clients know all servers in $S$, but the set of clients is unknown to the servers. A client or a server is also called a *process*.

Every process executes a deterministic algorithm assigned to it, unless it *crashes*. (The process does not behave maliciously.) If it crashes, the process simply stops its execution, unless it possibly *recovers*, in which case the process executes a *recovery procedure* which is part of the algorithm assigned to it. (Note that in this case we assume that the process is aware that it had crashed and recovered.) A process is *faulty* if there is a time after which the process crashes and never recovers. A non-faulty process is also called a *correct* process. The set of faulty processes in a run

is not known in advance. In particular, any number of clients can fail in a run. However, there is an known upper bound $f \geq 1$ on the number of faulty servers in a run. We also assume $f < n/2$ which is necessary to implement a register in asynchronous model.

Every process has a volatile storage and a stable storage (e.g., hard disk). If a process crashes and recovers, the content of its volatile storage is lost but the content of its stable storage is unaffected. Whenever a process updates one of its variables, it does so in its volatile storage by default. If the process decides to store information in its stable storage, it uses a specific operation **store**: we also say that the process *logs* the information. The process retrieves the logged information using the operation **retrieve**.

**Fair-lossy channels.** We assume that any pair of processes, say $p_i$ and $p_j$, communicate using fair-lossy channels [13, 4], which satisfies the following three properties: (1) If $p_j$ receives a message $m$ from $p_i$ at time $t$ then $p_i$ sent $m$ to $p_j$ at time $t$, (2) if $p_i$ sends a message $m$ to $p_j$ a finite number of times, then $p_j$ receives the message a finite number of times, and (3) if $p_i$ sends a message $m$ to $p_j$ an infinite number of times and $p_j$ is correct, then $p_j$ receives $m$ from $p_i$ an infinite number of times.

On top of the fair-lossy channels we can implement more useful stubborn communication procedures (*s-send* and *s-receive*) which are used to send and receive messages reliably [4]. In addition to the first two properties of fair-lossy channels, stubborn procedures satisfy the following third property: If $p_i$ s-sends a message $m$ to a correct process $p_j$ at some time $t$, and $p_i$ does not crash after time $t$, then $p_j$ eventually s-receives $m$. We would like to note that stubborn primitives can be implemented without using stable storage [4].

**Registers.** A sequential register is a data structure accessed by a single process that provides two operations: write($v$), which stores $v$ in the register and returns OK, and read(), which returns the last value stored in the register. (We assume that the initial value of the register is $\perp$, which is not a valid input value for a write operation.) An atomic register is a distributed data-structure that can be concurrently accessed by multiple processes and yet provide an "illusion" of a sequential register to the accessing processes [12, 13]. An algorithm *implements* an atomic register if all runs of the algorithm satisfy the *atomicity* and *termination* properties. We follow the definition of atomicity for a crash-recovery model given in [8], which in turn extends the definition given in [11]. We recall the definition in Appendix A.

We use the following two termination conditions in this paper. (1) An implementation satisfies wait-free termination (for clients) if for every run where at most $f$ of the servers are faulty (and any number of clients are faulty), every operation invoked by a correct client completes. (2) An implementation satisfies Finite-Write (FW) termination [1] if for every run where at most $f$ of the servers are faulty (and any number of clients are faulty), every write invocation by a correct client is complete, and moreover, either every read invocation by a correct client is complete, or infinitely many writes are invoked in the run. (Note that wait-free termination implies FW-termination.)

**Erasure coding.** A $k$-of-$n$ erasure coding [16] is defined by the following two primitives:

- **encode**($V, k, n$) which returns a vector $[V[1], \ldots, V[n]]$, where $V[i]$ denotes $i^{th}$ encoded fragment. (For presentation simplicity, we will assume that encode returns a *set* of $n$ encoded fragments of $V$, where each fragment is tagged by its fragment number.)

- **decode**($X, k, n$) which given a set $X$ of at least $k$ fragments of $V$ (that were generated by encode($V, k, n$)), returns $V$.

For our algorithms, we make no assumption on the specific implementation of the primitives except the following one: each fragment in a $k$-of-$n$ encoding of $V$ is roughly of size $|V|/k$.

In the next two sections, we present two algorithms that implement erasure-coded, multi-writer multi-reader, atomic registers in a crash-recovery model, ORCAS-A and ORCAS-B. Both implementations have low storage overhead when no write operation is in progress. The implementations differ in the storage overhead during a write, and in their message sizes.

## 3  ORCAS-A

We now present our first implementation which we call ORCAS-A. (The pseudocode is given in Figures 1 and 2.) The implementation is inspired by the well-known atomic register implementations in [3, 14]. Also, the registration process of a read at the servers is inspired by the listeners communication pattern in [15]. The first two phases of the write function are similar to that in [3, 14]— they store the unencoded values at $n - f$ (a majority) of servers with an appropriate timestamp. Additionally in ORCAS-A, depending on the number of servers from which the write receives a reply, it selects an encoding $r$-of-$n$. Then, the write performs another round trip where it requests the servers to encode the value using $r$-of-$n$ encoding and retain the fragment corresponding to its server id. The crucial parts of the implementation are choosing an encoding $r$-of-$n$ and the condition for waiting for fragments at a read, such that, any read can recover the written value without blocking permanently. We now describe the implementation in more detail.

```
 1: function initialization:
 2:     ts, wid, rid ← 0; r ← 1; T_c ← timer()    {at every
        client}
 3:     A_j ← ⊥; τ, δ ← 0; ρ ← 1   {at every server s_j}

 4: function write (V) at client c_i
 5:     wid ← wid + 1; ts ← 0
 6:     store(wid, ts)
 7:     repeat
 8:         send(⟨get_ts, wid⟩, S)
 9:     until s-receive ⟨ts_ack, *, wid⟩ from n − f servers
10:     ts ← 1+ max{ts_j : s-received ⟨ts_ack, ts_j, wid⟩}
11:     store(ts, V)
12:     trigger(T_c)
13:     repeat
14:         send(⟨write, ts, wid, 0, V⟩, S)
15:     until s-receive ⟨w_ack, ts, wid, 0⟩ from n−f servers and
        expired(T_c)
16:     r ← (number of servers from which s-received
        ⟨w_ack, ts, wid, 0⟩ messages) −f
17:     if r > 1 then
18:         repeat
19:             S' ← set of servers from which s-received
                ⟨w_ack, ts, wid, 0⟩ until now
20:             send (⟨encode, ts, wid, r⟩, S')
21:         until s-receive ⟨enc_ack, ts, wid, r⟩ from n−f servers
22:     return(OK)

23: upon receive ⟨get_ts, wid⟩ from client c_i at server s_j do
24:     s-send(⟨ts_ack, τ, wid⟩, {c_i})

25: upon receive ⟨write, ts', wid', rid', V'⟩ from client c_i at
    server s_j do
26:     if rid' > 0 then
27:         R ← R \ {[rid', *, *, i]}
28:     if V' ≠ ⊥ then
29:         if [ts', wid'] >_lex [τ, δ] then
30:             τ ← ts'; δ ← wid'; ρ ← 1; A_j ← V'
31:             store(τ, δ, ρ, A_j)
32:         for all [rid, ts, id, l] ∈ R do
33:             s-send(⟨r_ack, rid, ts, wid', 1, V'⟩, {c_l})
34:     s-send(⟨w_ack, ts', wid', rid'⟩, {c_i})

35: upon receive ⟨encode, ts', wid', r'⟩ from client c_i at server
    s_j do
36:     if [ts', id'] = [τ, δ] then
37:         A_j ← j^{th} fragment of encode(A_j, r', n)
38:         ρ ← r'
39:         store(ρ, A_j)
40:     s-send(⟨enc_ack, ts', wid', r'⟩, {c_i})

41: upon recovery() at server s_j do
42:     [τ, δ, ρ, A_j] ← retrieve()

43: upon recovery() at client c_i do
44:     [rid, ts, wid, r, V] ← retrieve()
45:     if ts ≠ 0 then
46:         repeat
47:             send(⟨write, ts, wid, 0, V⟩, S)
48:         until s-receive ⟨w_ack, ts, wid, 0⟩ from n − f servers
```

Figure 1: ORCAS-A: initialization, write and recovery procedures

```
1: function read() at client c_i
2:    rid ← rid + 1; Γ ← 0; M ← ∅; once ← false
3:    store(rid)
4:    repeat
5:       send(⟨read, rid⟩, S)
6:       M ← {msg = ⟨r_ack, rid, *, *, *, *⟩  :  s-received
             msg}
7:       TS ← max_lex{[ts, id]  :  ⟨r_ack, rid, ts, id, *, *⟩ ∈ M}
8:       if (M contains messages from at least n − f servers)
            and (once = false) then
9:          Γ ← TS; once ← true
10:            if TS = [0, 0] then return(⊥)
11:    until (once = true) and (∃ r', ts', id' such that
         ([ts', id'] ≥_lex Γ) and (|{A_j : ⟨r_ack, rid, ts', id', r', A_j⟩ ∈
         M}| ≥ r'))
12:    A ← set of A_j satisfying the condition in line 11
13:    if r' = 1 then
14:       V ← any A_j in A; V' ← V
15:    else
16:       V ← decode(A, r', n); V' ← ⊥
17:    repeat
18:       send(⟨write, ts', id', rid, V'⟩, S)
19:    until s-receive ⟨w_ack, ts', id', rid⟩ from n − f servers
20:    return(V)

21: upon receive ⟨read, rid⟩ from client c_i at server s_j do
22:    if R does not contain any [rid, *, *, i] then
23:       R ← R ∪ [rid, τ, δ, i]
24:       s-send(⟨r_ack, rid, τ, δ, ρ, A_j⟩, {c_i})
```

Figure 2: ORCAS-A: read procedure

## 3.1 Description

**Local variables.** The clients maintain the following local variables: (1) $ts$: part of the timestamp of the current write operation, and (2) $wid, rid$: the identifiers of write and read operations, respectively, which are used to distinguish between messages from different operations of the same client, and (3) a timer $T_c$ whose timeout duration is set to the round-trip time for contacting the servers in synchronous periods. The pair $[ts, wid]$ form the timestamp for the current write. The local variables at a server $s_j$ are as follows: (1) $A_j$: its share of the value stored in the register, which can either be the unencoded value or the $j^{th}$ encoded fragment, (2) $\tau, \delta$: the $ts$ and the $wid$, respectively, associated with the value in $A_j$, and (3) $\rho$: the encoding associated with the value in $A_j$, namely, $A_j$ is the $j^{th}$ fragment of a $\rho$-of-$n$ encoding of some value. (In particular, $\rho = 1$ implies that $A_j$ contains an unencoded value.)

**Write operation.** The write operation consists of three phases, where each phase is a round-trip of communication from the client to the servers. The first phase is used to compute the timestamp for the servers, the second phase to write the unencoded value at the servers, and the final phase is used to encode the value at the servers. On invoking a write($V$), the client first increments and logs its $wid$. This helps in distinguishing messages from different operations of the same server even across a crash-recovery. It also logs $ts = 0$ so as to detect an incomplete write across a crash-recovery. Next, the client sends $get\_ts$ messages to all servers and waits until it receives $ts$ from at least $n - f$ servers. (The notation $send(m, X)$ is a shorthand for the following: for every processes $p \in X$, send the message $m$ to $p$. It is not an atomic operation.) To overcome the effect of the fair-lossy channels, a client encloses the sending of its messages to the servers in a repeat-until loop, and the servers reply back using the s-send primitive. On receiving the $ts$ from at least $n - f$ servers, the client increments by one the maximum $ts$ received, to obtain the $ts$ for this write. It then logs $ts$ and $V$ so that in case of a crash during the write, the client can complete the write upon recovery. Next, it starts its timer, and sends a $write$ message with the timestamp $[ts, wid]$ and the value $V$, to the servers. (To distinguish this message from the $write$ message sent by a read operation, the message also contains a $rid$ field which is set to 0.) A server on receiving a $write$ message with a higher timestamp than its current timestamp $[\tau, wid]$, updates $A_j, \tau$ and $\delta$ to $V$, $ts$ and $wid$ of the message, respectively. It also updates the encoding $\rho$ to 1 (to denote that the

7

contents of $A_j$ is unencoded), and logs the updated variables. (The server also sends some message to the readers which we will discuss later.) The client waits until it receives $w\_ack$ messages from at least $n - f$ servers, and the timer expires. (Waiting for the timer to expire ensures that the client receives a reply from all non-crashed processes in synchronous periods.)

Next, the client select the encoding for the write to be $r = q - f$, where $q$ is the number of $w\_ack$ messages received by the client. Note $r \geq 1$ because $q \geq n - f$ and $f < n/2$. Then the client sends an *encode* message to all servers which have replied to the *write* message. A server $s_j$ on receiving this message encodes it value $A_j$ using $r$-of-$n$ encoding, and retains only the $j^{th}$ fragment in $A_j$. It also updates its encoding $\rho$ to $r$, logs $A_j$ and $\rho$, and replies to the client. The client returns from the write on receiving $n - f$ replies. (Note that the encode phase is skipped if $r = 1$, because 1-of-$n$ encoding is same as not encoding the value at all.)

**Read operation.** The read operation consists of two phases. The first phase gathers enough fragments to reconstruct a written value, and the second phase writes back the value at the servers to ensure that any subsequent reader does not read an older value.

On invoking a read, the client increments and logs its $rid$. It then sends a *read* message to the servers. On receiving a *read* message, a server *registers* the read[3] by appending it to a local list $\mathcal{R}$ with the following parameters: the $rid$ of the *read* message, and the timestamp $[\tau, \delta]$ at the server when the *read* message was received. (The client *de-registers* in the second phase of the read: line 27, Figure 1.) The server then replies with its current value of $A_j$ and its associated timestamp and encoding. In addition, whenever the server receives a new *write* message with a higher timestamp, it forwards it to its registered readers. The client on the other hand, first chooses a timestamp $\Gamma$ which is greater than or equal to the timestamp seen at $n - f$ processes,[4] and then waits for enough fragments to reconstruct a written value that has an associated timestamp greater than or equal to $\Gamma$: the condition in line 11 of Figure 2 simply requires that (1) the client receives $r\_ack$ from at least $n - f$ servers, and (2) there is an encoding $r'$ and timestamp $[ts', id']$ such that the client has received at least $r'$ fragments of the associated value, and $[ts', id']$ is greater than or equal to $\Gamma$. In Appendix C, we show that this condition is eventually satisfied for every read whose invoking client does not crash.

The second phase of a read is very similar to the second phase of a write except for the following case. If the read selects a value in the first phase that was encoded by the corresponding write ($r' > 1$), then the read does not need to write back the value to the servers because the write has already completed its second phase. In this case, the second phase of the read is only used to deregister the read at the servers.

**Recovery Procedures.** The recovery procedure at a server is straightforward: it retrieves all the logged values. The client, in addition to retrieving the logged values, also completes any incomplete write. (Note than, even if the last write invocation, before the crash at a client, is complete, $ts$ can be greater than 0. In this case, the recovery procedure tries to rewrite the same value with the same timestamp. It is easy to see that this attempt to rewrite the value is harmless.)

## 3.2 Correctness

The proof of the atomicity of ORCAS-A is similar to the implementations in [3, 14]. The only non-trivial argument in the proof of wait-free termination is proving that the waiting condition

---

[3]When there is no ambiguity, we also say that the server registers the client.

[4]The $\Gamma$ selected in this way is higher than or equal to the timestamp of all preceding writes because two server sets of size $n - f$ always has a non-empty intersection.

in line 11 in Figure 2 eventually becomes true in every run where the client does not crash after invoking the read. In this section, we give an intuition for this proof by considering a simple case where a (possibly incomplete) write is followed by a read, and there are no other operations.

Suppose there is a write($V$) that is followed by a read(). We claim that the read() can always reconstruct $V$ or the initial value of the register, and it can always reconstruct $V$ if the write is complete. The write() operation has two phases that modify the state of the servers: the write phase and the encode phase. Suppose that during the write phase, the writer receives replies from $q$ servers (denoted by set $Q$) such that $q \geq n - f > f$. If the writer fails without completing this phase, the read() can return the initial value of the register, which does not violate atomicity. In the encode phase, an $r$-of-$n$ encode message is sent to all servers, where $r = q - f \geq n - 2f > 0$. If the writer crashes, this message reaches an arbitrary subset of servers. Subsequently, the read() contacts a set $R$ containing at least $n - f$ servers. We denote the intersection of the read and write sets, by $U$, i.e. $U = Q \cap R$, and it follows that $|U| \geq q - f = r > 0$. There are two cases:

**Case 1:** There is at least one server in $U$ which still has the unencoded value $V$. The read can thus directly obtain $V$ from this server.

**Case 2:** All the servers in $U$ have received the encode message and encoded $V$. Since $|U| \geq r$ and an $r$-of-$n$ erasure code was used, there are enough fragments for the read to reconstruct $V$.

However, we must also consider the case where the read() is concurrent with multiple writes. If there is a series of consecutive writes, the write procedure ensures that all values are eventually encoded. If the read is slow, it could receive an encoded fragment of a different write from each server, making it impossible for the read to reconstruct any value. But the reader registration ensures that the servers will send all new fragments to the reader until the reader is able to reconstruct some written value. A detailed proof of wait-freedom is given in Appendix C.

## 3.3 Algorithm Complexity

In this section we discuss the theoretical performance of ORCAS-A.

**Timing guarantees.** For timing guarantees we consider periods of a run where links are timely, local computation time is negligible, at least $n - f$ servers are alive, and no process crashes or recovers. It is easy to show that a write operation completes in three round-trips (i.e., six communication steps), as compared to two round-trips in the implementation of [14]. (We discuss this comparison further in Section 5.) Also it is straightforward to show that a read can complete in two round-trips if there is no write in progress. In Appendix E, we show that even in the presence of concurrent writes, the read registration ensures that a read operation terminates within five communication steps.

**Messages.** Except the $r\_ack$ messages, the number of messages used by an operation is linear in the number of servers. In Appendix D we show how to circumvent the reader registration by slightly weakening the termination condition of the read. Message sizes in ORCAS-A are as large as those in the replication based register implementations of [3, 14]: the first phase of the write in ORCAS-A sends the unencoded value to all servers.

**Worst-case bound on storage.** Suppose that all possible write values are of a fixed size $\Delta$, and the size of variables, other than those containing a value of a write operation or an encoded fragment of such a value, is negligible. Consider a partial run $pr$ that has no incomplete write invocation. (An invocation that has no matching return event in the partial run is called incomplete.) The $r$ computed in line 16 of Figure 1 of every write is at least $n - 2f$. Thus, every encoded fragment is at most of size $\Delta/(n - 2f)$. Since, there are no incomplete write invocation in $pr$,

and every write encodes the value at $n - f$ processes before it returns, the size of (stable) storage at $n - f$ servers is at most $\Delta/(n - 2f)$ at the end of $pr$. In addition, note that the size of the storage at *all* servers is *always* bounded by $\Delta$. This is in contrast to the implementation in [6] and ORCAS-B implementation that we describe later, where the worst-case storage size is dependent on the maximum number of concurrent writes.

**Bound on storage in synchronous periods.** Consider a partial run $pr$ which has no incomplete write invocation. Let $wr$ be the write with the highest timestamp in $pr$. Let $t$ be the time when $wr$ was invoked. Now, assume that (1) the links were timely in $pr$ from time $t$ onwards, and (2) at least $n - f$ servers are alive at time $t$, and no process crashes or recovers from time $t$ onwards. Let $q \geq n - f$ be the set of servers that are alive at time $t$. Then, it is easy to see that the $r$ computed in line 16 of Figure 1 is $q - f$ in $wr$, and hence, the size of storage at *all* alive servers is at most $\Delta/(q - f)$ at the end of $pr$. It also follows that, if $pr$ is a synchronous failure-free partial run, then the size of storage at all servers is at most $\Delta/(n - f)$ at the end of $pr$.

**FW-termination.** Consider the case in the above implementation when a client invokes a read, registers at all the servers, and then crashes. If a server does not crash, its s-send module will send the $r\_ack$ message to the client forever. Since, these messages are of large sizes, it may significantly increase the load on the system. Following [1], we show in Appendix D that if we slightly weaken the waif-free termination condition of the read to Finite-Write (FW) termination, then such messages are not required.

# 4 ORCAS-B

Although the ORCAS-A implementation saves storage space in synchronous periods, it has two important drawbacks because it sends the unencoded values to the servers in the first phase of the write. First, it uses larger messages compared to implementations which never send any unencoded values to the servers. Second, if a client crashes before sending an *encode* message during a write, servers are left with an unencoded value in the stable storage.[5] In this section, we present our second implementation, ORCAS-B, which like most erasure-coded register implementations, never sends an unencoded value to the servers.

Due to lack of space, we discuss only those parts of ORCAS-B that significantly differ from ORCAS-A. (The pseudocode is presented in Figures 4 and 5 in the Appendix). The crucial difference between ORCAS-A and ORCAS-B is how the write value is encoded during a write and how it is reconstructed during a read. In ORCAS-B, the write consists of three phases. The first phase finds a suitable timestamp for the write, and tries to guess the number of alive servers, say $r'$. The write then encodes the value such that the following three conditions holds. (1) If the second phase of the write succeeds in contacting only $n - f$ servers (a worst-case scenario), a subsequent read can reconstruct the value. (2) If the second phase succeeds in contacting $r'$ servers (the optimistic case), then in the third phase, the write can "trim" (i.e., reduce the size of) the stored encoded value at the servers, and still a subsequent read can reconstruct the value. (3) The size of the stored encoded value at a server should be equal to the size of a fragment in $(n - 2f)$-of-$n$ encoding in the first case, and $(r' - f)$-of-$n$ encoding in the second case. The motivation behind these three conditions is to have the same optimistic storage requirements as in ORCAS-A.

---

[5]In practice, the second case might not cause a significant overhead because any subsequent complete write will erase such unencoded values.

It is not difficult to see that if the write uses a $(n - 2f)$-of-$n$ encoding, then a server cannot *locally* extract its trimmed fragment in the third phase from the encoded fragment it receives in the second phase, without making extra assumptions about $n$ or the erasure coding algorithm. Thus with $(n - 2f)$-of-$n$ encoding in the second phase, in the third phase of the write, either the write needs to send the trimmed fragment to each server, or the servers need to exchange their (second-phase) fragments. In ORCAS-B we avoid this issue by simply storing multiple fragments at a server, while still satisfying the three conditions above.

We define the following variables: (1) $r = r' - f$, (2) $x$ be the *lcm* (least common multiple) or $r$ and $n - 2f$, (3) $z = x/(n - 2f)$, and (4) $y = x/r$. Now the second phase of the write encodes the value using $x$-of-$(nz)$ encoding. It then tries to store $z$ fragments at each server. If the write succeed in storing the fragments at $r'$ servers, then in the next phase, it sends a *trim* message that requests the servers to retain $y$ out of its $z$ fragments (and delete the remaining fragments). Now it is easy to verify the above three conditions. If the second phase of the write stores the fragments at $n - f$ servers, a subsequent read can access at least $n - 2f$ of those servers, and thus receive at least $(n - 2f)z = x$ fragments. On the other hand, if the stored fragments at some server are trimmed, then at least $r'$ servers have at least $y$ fragments, and therefore a subsequent read receives $y$ fragments from at least $r' - f = r$ servers; i.e., $ry = x$ fragments in total. In both cases, since the write has used $x$-of-$(nz)$ encoding, the read can reconstruct the value. To see that the third condition is satisfied, notice that the total size of $z$ stored fragments at a server after the second phase of the write is $z(\Delta/x) = \Delta/(n-2f)$. After trimming, the size of the stored fragments become $y(\Delta/x) = \Delta/(r' - f)$.

Another significant difference between ORCAS-A and ORCAS-B is the condition for deleting an old value at a server. In ORCAS-A, whenever a server receives an unencoded value with a higher timestamp, the old fragment or the old unencoded value is overwritten. However in ORCAS-B, if the server receives fragments with a timestamp $ts$ that is higher than its current timestamp, the server adds the fragments to a set $L$ of received fragments. Subsequently, if it receives a trim message (i.e., a message from the third phase of a write) with timestamp $ts$, it deletes all fragments in $L$ with a lower timestamp. Also, the server sends the whole set $L$ in its $r\_ack$ reply messages to a read. (Thus the trim message also acts as a garbage collection message.) This modification is necessary in ORCAS-B because, until a sufficient number of encoded fragments are stored at the servers, the newly written value is not recoverable from the stored data obtained from any set of servers. The trim message acts as a confirmation that enough fragments of the new value have been stored. A similar garbage collection mechanism is also present in the implementation in [6]. On the other hand, since a server in ORCAS-A receives an unencoded value first, it can directly overwrite values with lower timestamps. An important consequence of this modification is that the worst-case storage size of ORCAS-B (and the implementation in [6]) is proportional to the number of concurrent writes, whereas, the storage requirement in ORCAS-A is never worse than that in replication (i.e., storing the unencoded value at all servers). We show the wait-free termination property of ORCAS-B in Appendix F.

## 5 Discussion and Future Work

There are two related disadvantages of ORCAS-A when compared to most replication based implementations. The write needs three phases to complete as compared to two phases in the latter. Also, the write needs four stable storage accesses (in its critical path) as compared to two such

accesses in replication based implementations. Both disadvantages primarily result from the last phase that is used for encoding the value at the servers, and which can be removed if we slightly relax the requirement on the storage space. ORCAS-A ensures that the stable storage is encoded whenever there is no write in progress. Instead, if we require that the stable storage is *eventually* encoded whenever there is no write in progress, then (with some minor modifications in ORCAS-A) the write operation can return without waiting for the last phase. The last phase can then be executed "lazily" by the client. (The two disadvantages and the above discussion hold for ORCAS-B as well.) On a similar note, in ORCAS-A, if a read selects a value in the first phase that is already encoded at some server, then it can return after the first phase, and lazily complete the second phase (which in this case is used only for deregistering at the servers, and not for writing back the value). It follows that, a read that has no concurrent write in ORCAS-A can return after the first phase.

An important open problem is to study storage lower bounds on register implementations in a crash-recovery model. In particular, it would be interesting to study if our lower bounds hold when some of the underlying assumptions are removed. Another interesting direction for investigation can be implementations that tolerate both process crash-recovery with fair-lossy channels and malicious processes.

# References

[1] Ittai Abraham, Gregory Chockler, Idit Keidar, and Dahlia Malkhi. Byzantine disk paxos: optimal resilience with byzantine shared memory. *Distributed Computing*, 18(5):387–408, 2006.

[2] M. K. Aguilera, R. Janakiraman, and L. Xu. Using erasure codes efficiently for storage in a distributed system. In *Proceedings of the International Conference on Dependable Systems and Networks (DSN)*, pages 336–345, 2005.

[3] H. Attiya, A. Bar-Noy, and D. Dolev. Sharing memory robustly in a message passing system. *Journal of the ACM*, 42(1):124–142, 1995.

[4] R. Boichat and R. Guerraoui. Reliable and total order broadcast in the crash-recovery model. *Journal of Parallel and Distributed Computing*, 65(4):397–413, 2005.

[5] C. Cachin and S. Tessaro. Optimal resilience for erasure-coded byzantine distributed storage. In *Proceedings of the International Conference on Dependable Systems and Networks (DSN)*, pages 115–124, 2006.

[6] S. Frolund, A. Merchant, Y. Saito, S. Spence, and A. Veitch. A decentralized algorithm for erasure-coded virtual disks. In *Proceedings of the International Conference on Dependable Systems and Networks (DSN)*, pages 125–134, 2004.

[7] G. R. Goodson, J. J. Wylie, G. R. Ganger, and M. K. Reiter. Efficient byzantine-tolerant erasure-coded storage. In *Proceedings of the International Conference on Dependable Systems and Networks (DSN)*, 2004.

[8] Rachid Guerraoui, Ron R. Levy, Bastian Pochon, and Jim Pugh. The collective memory of amnesic processes. *ACM Transactions on Algorithms*, 4(1), 2008.

[9] James Hendricks, Gregory R. Ganger, and Michael K. Reiter. Low-overhead byzantine fault-tolerant storage. In *Proceedings of the 21st ACM Symposium on Operating Systems Principles (SOSP)*, pages 73–86, 2007.

[10] James Hendricks, Gregory R. Ganger, and Michael K. Reiter. Verifying distributed erasure-coded data. In *Proceedings of the Twenty-Sixth Annual ACM Symposium on Principles of Distributed Computing (PODC)*, pages 139–146, 2007.

[11] Maurice Herlihy. Wait-free synchronization. *ACM Transactions on Programming Languages and Systems*, 13(1):124–149, 1991.

[12] L. Lamport. On interprocess communication - part i: Basic formalism, part ii: Algorithms. *DEC SRC Report*, 8, 1985. Also in Distributed Computing, 1, 1986, pages 77-101.

[13] N.A. Lynch. *Distributed Algorithms*. Morgan Kaufmann Publishers, San Mateo, CA, 1996.

[14] N.A. Lynch and A.A. Shvartsman. Robust emulation of shared memory using dynamic quorum-acknowledged broadcasts. *Proceedings of the International Symposium on Fault-Tolerant Computing Systems (FTCS)*, 1997.

[15] Jean-Philippe Martin, Lorenzo Alvisi, and Michael Dahlin. Minimal byzantine storage. In *Proceedings of the International Symposium on Distributed Computing (DISC)*, pages 311–325, 2002.

[16] I. S. Reed and G. Solomon. Polynomial codes over certain finite fields. *SIAM Journal of Applied Mathematics*, 8:300–304, 1960.

# A  Atomicity in a crash-recovery model (from [8])

A history of a run is a total order of events of four kinds in the run: invocations, replies, crashes and recoveries. An invocation with no matching reply in a history is said to be pending (or incomplete) in that history. (An invocation with a matching reply is complete, and the corresponding operation is said to be complete.) An operation $op$ is said to precede an operation $op'$ in a history if the reply of $op$ precedes the invocation of $op'$ in that history. An operation $op$ is said to immediately precede an operation $op'$ in a history if the reply of $op$ precedes the invocation of $op'$ in that history such that no operation $op''$ precedes $op'$ where $op$ precedes $op''$. A local history is a sequence of events associated with the same process. A local history is said to be well-formed if: (a) its first event is either an invocation or a crash, (b) a crash is followed by a matching recovery event or is not followed by any event, and (c) an invocation is followed by a crash or a reply event. A history is said to be well-formed if all its local histories are well-formed. Two histories $H$ and $H'$ are said to be equivalent if, for every process $p$, the local history $H$ restricted to $p$ has the same object events as the local history $H'$ restricted to $p$. To define atomicity, we reason about histories that are complete. These are histories without any crash or recovery events where every invocation has a matching reply. Given any well-formed history $H1$, we say that history $H2$ completes $H1$ if $H2$ does not contain any crash or recovery events and is made of the very same invocation and reply events in the same order as in $H1$, with one exception: any pending invocation in $H1$ is either absent in $H2$, or has a matching reply that appears in $H2$ before the subsequent invocation of the same process. A completed operation has a pending invocation in $H1$ that has a matching reply that appears in

$H2$ before the subsequent invocation of the same process. A history is said to be sequential if it is complete and every invocation is followed by a matching reply. A sequential history of an atomic register is said to be legal if it belongs to the sequential specification of register: a read returns the last written value. A history $H$ is said to be atomic if it can be completed to a history that is equivalent to some legal sequential history. An implementation satisfies the atomicity property if the history of every run of the implementation is atomic.

## B  Bounds on storage overhead

In this section, we show two lower bounds on the storage space usage in a specific kind of erasure coding based atomic register implementation. (The bounds also hold for safe and regular registers.) These lower bounds show that the storage requirements of ORCAS implementations are optimal if we assume a specific kind of implementations.

We make the following assumptions on the erasure coding based register implementations. (Recall that $n$ denotes the number of servers.)

- There is a value $V$ such that, for any $k' < n'$, $V$ cannot be reconstructed from less than $k'$ fragments of encode($V, k', n'$).

- Any given write or a read operation uses at most one encoding and that encoding is of the form $qk$-of-$qn$ for some $q \geq 1$ and $1 \leq k \leq n$.[6]

- If there are no pending writes at the end of some partial run $pr$ (and there is at least one complete write), then for at least $n - f$ servers the following condition holds at the end of $pr$. The only content of the stable storage of the server that is of non-negligible size are $q'$ fragments of some encode($V', q'k', q'n$), such that there is a write($V'$) in $pr$ that is $q'k'$-of-$q'n$ erasure-coded.

The last assumption ensures that the servers only store the encoded fragments when there are no write in progress. This requirement is restricted to $n - f$ servers because in an asynchronous system, a write operation may miss $f$ servers.

We now present bounds on the storage requirements in synchronous and asynchronous systems. Theorem 1 considers a synchronous system [13]: a system where message delays and process speeds are bounded and known. We say that a run of a synchronous system is failure-free if there are no crashes in that run.

**Theorem 1** *Consider an erasure-coded storage that implements an atomic register in a synchronous system with $n$ servers out of which $f$ can be faulty. Then the implementation has a failure-free run in which some write is $qk$-of-$qn$ erasure-coded, such that $k \leq n - f$.*

PROOF: Suppose by contradiction that in every failure-free run of the implementation, every write has an encoding $qk$-of-$qn$ such that $k > n-f$. (Different operation may use different values of $q$ and $k$.) First consider a failure-free run in which a client completes $qk$-of-$qn$ erasure-coded write($V$) and there are no other operations. Let $pr1$ be a partial run of the failure-free run till the step where the write returns.

---

[6]We say that a write or a read operation uses a $qk$-of-$qn$ encoding (or is $qk$-of-$qn$ erasure-coded) if the operation ever invokes the encode($*, qk, qn$). Different operations may use different values of $q$ and $k$.

From our assumptions, there are at least $n - f$ servers such that each of these servers contains at most $q$ fragments from encode($V, qk, qn$) or the initial value of the register. Let $X$ be the set of the remaining $f$ servers. Now consider another partial run $pr2$ of the implementation that extends $pr1$. Partial run $pr2$ extends $pr1$ as follows. After the write returns the following events occur: (1) the servers in set $X$ and the writer crash but never recover, (2) all other processes crash and all messages that are in transit in the communication channels are lost, and then (3) all processes, other than the writer and those in set $X$, recover. (Recall that the channels are fair-lossy.) Now another client invokes and completes a read().

From the definition of a register, the read must return $V$. However, the servers that are not in $X$ have at most $q(n - f)$ fragments out of the $qn$ fragments returned by encode($V, qk, qn$). Since $k > n - f$, the read() cannot recover $V$ from the fragments available at the alive servers, a contradiction. □

The above theorem implies that there is a run in which some server stores $q$ fragments of some write($V$) that is $q(n - f)$-of-$qn$ erasure-coded. Thus its worst-case storage usage is $\frac{|V|}{q(n-f)}q = \frac{|V|}{n-f}$.

Next we consider a lower bound for an asynchronous system.

**Theorem 2** *Consider an erasure-coded storage that implements an atomic register in an asynchronous system with $n$ servers out of which $f$ can be faulty, and $n > 2f$. Then the implementation has a run in which some write is $qk$-of-$qn$ erasure-coded, such that $k \leq n - 2f$.*

PROOF: Suppose by contradiction that in every run of the implementation, every write has an encoding $qk$-of-$qn$ such that $k > n - 2f$. Consider the following partial run $pr$ of the implementation. A set $X$ of $f$ servers crash at the beginning of the run and do not recover. Then a client $c$ completes $qk$-of-$qn$ erasure-coded write($V$). From our assumptions, there is a set $Y$ of $n - f$ servers, each of which contain at most $q$ fragments from encode($V, qk, qn$) or the initial value of the register. Now construct a set $Z$ of $f$ servers as follows. First choose all servers in $S - X - Y$. If $|S - X - Y| < f$, choose any $f - |S - X - Y|$ servers from $Y - X$. (Since $n > 2f$, we can always choose such a set $Z$.) Note $X$ and $Z$ are disjoint sets of size $f$ each, and hence, $|S - X - Z| = n - 2f$.

Now consider another run which is exactly the same as $pr$ till write($V$) returns, except that the servers in $X$ do not crash, but do not take any steps till write ($V$) returns. (Recall that the system is asynchronous.) Now, the following events occur: (1) servers in $Z$ and the writer crash but never recover, (2) all other processes crash and all messages that are in transit in the communication channels are lost, and (3) all processes, other than the writer and those in set $Z$, recover. Now another client invokes and completes a read().

From the definition of a register, the read must return $V$. Note that, when the read is invoked, no server in $X$ has any fragment of $V$ because none of them took a step before write($V$) returned, and all messages in communication channels are lost before they took any step. Among the remaining servers, the ones in $Z$ do not recover from the crash, and the $n - 2f$ servers in $S - X - Z \subseteq Y$ have at most $q(n - 2f)$ fragments out of the $qn$ fragments returned by encode($V, qk, qn$). Since $k > n - 2f$, the read() cannot reconstruct $V$, a contradiction. □

The above theorem implies that there is a run in which some server stores $q$ fragments of some write($V$) that is $q(n-2f)$-of-$qn$ erasure-coded. Thus its worst-case storage usage is $\frac{|V|}{q(n-2f)}q = \frac{|V|}{n-2f}$.

It follows that any translation of Cachin and Tessaro's algorithm [5] to a crash-recovery model with fair-lossy channels cannot preserve the storage overhead of $n/(n - f)$. We now explain this in more detail. For simplicity of discussion, assume that the algorithm in [5] is implemented in the

15

crash-recovery model with $k$-of-$n$ encoding, where, $k = n - f$. Then the *Disperse* primitive becomes a reliable broadcast of the encoded fragments, where each server eventually *rdelivers* its own fragment. However, in a crash-recovery model with fair-lossy channels, implementing this primitive will require servers to log more than one $k$-of-$n$ encoded fragments. The proof for this is essentially same as our lower bound proof. Consider a run in which a set $X$ of $f$ servers crash initially. Now some client $c$ *rbcasts* $n$ fragments, and $n - f$ servers rdelivers their respective fragments, and then the client crashes. Suppose that each of these $n - f$ servers, say set $Y$, store only one encoded $k$-of-$n$ fragment. Now suppose $f$ of the servers in $Y$ crash permanently, the remaining servers in $Y$ crash and recover, the set $X$ containing $f$ initially crashed servers recover, and all messages in transit in the channels are lost (fair-lossy channels). Then there are only $n - 2f < k$ fragments remaining in the system, and therefore, the servers in $X$ cannot reconstruct and rdeliver their fragments.

## C   ORCAS-A Wait-free Termination

**Lemma 1** *The ORCAS-A implementation is* wait-free.

(*Sketch.*) We sketch the only non-trivial argument in the proof of wait-freedom: the condition in line 11 of Figure 2 eventually becomes true if the client does not crash after invoking the read. Suppose by contradiction that a client $c$ invokes a read operation $rd$, and then does not crash, and the condition in line 11 of Figure 2 is never satisfied in that operation. Thus $c$ executes the first repeat-until loop in the operation forever. Consider the two cases: (1) there is a time after which the read does not receive any new timestamp $[ts, id]$, and (2) otherwise. We now show that in both the cases the condition in line 11 is eventually satisfied.
Let $\gamma$ denote the timestamp assigned to $\Gamma$ in line 9, and $t$ be the time when the assignment was executed. (Note that line 9 is executed exactly once in $rd$.)
**Case 1.** Let $TS$ be the highest timestamp seen by the read in the run. We note that no correct server ever stores a timestamp higher than $TS$ in the run. Otherwise, since both $c$ and the server are correct, and timestamp at a server is non-decreasing with time, $c$ will eventually receive a timestamp higher than $TS$, which violates the definition of $TS$.
By definition, $TS \geq_{lex} \gamma$. Let $t_1 > t$ be a time such that no correct process crashes after time $t_1$. Let $wr$ be the write operation that corresponds to timestamp $TS$, and suppose it writes a value $V$. Let $m$ be an $r\_ack$ message received by $rd$ with timestamp $TS$. Also, suppose $m$ is received at time $t_2$. If $m$ contains an unencoded value $V$, then the condition in line 11 is satisfied with $r' = 1$, and $[ts', id'] = TS$. Suppose otherwise. Then, $m$ contains an encoded fragment of $V$, say with encoding $x$-of-$n$. Since, a server encodes a value only in the third phase of a write, $wr$ has completed its second phase, and has contacted at least $x + f$ servers in the second phase. As there are at most $f$ faulty servers, and no correct server ever stores a timestamp higher than $TS$, $wr$ has stored the unencoded value at $x$ correct servers by time $t_2$. Let us denote this set of $x$ servers by $X$.
Let $t_3 = 1 + Max\{t_1, t_2\}$. For every server $s$ in $X$, consider *any $r\_ack$ messages* $m'$ such that (1) $s$ sends $m'$ in line 24 of Figure 2 after time $t_3$, and (2) $c$ receives $m'$. (Note that, such a message exists because $c$ executes the repeat-until loop an infinite number times, and therefore, $s$ replies back in line 24 an infinite number of times.) We have already showed that no server in $X$ stores a timestamp greater than $TS$. Thus, the timestamp of $m'$ is $TS$, and therefore, $m'$ either contains the unencoded value $V$ or an encoded fragment of $V$. Let $M_x$ be the set of such $x$ messages from the servers in $X$. Now, if any message in $M_x$ contains the unencoded $V$, then the condition in

line 11 is satisfied with $r' = 1$, and $[ts', id'] = TS$. Otherwise, if all message in $M_x$ contains an encoded fragment of $V$, then the condition in line 11 is satisfied with $r' = x$, and $[ts', id'] = TS$.

**Case 2.** Let $t_1 > t$ be a time such that (1) no correct server crashes after time $t_1$, and (2) all correct servers have received the *read* message from $rd$ (and therefore, all correct servers have registered the read). Since there are bounded number of clients, if the read keeps on seeing new timestamps then there is at least one client $c'$ that invokes an infinite number of writes. (As a client cannot invoke a new operation without completing the last invoked operation, $c'$ completes an infinite number of writes.) Thus there is a complete write operation $wr$ such that (1) $wr$ has a timestamp $TS$ higher than $\gamma$, and (2) $wr$ is invoked after time $t_1$. Let $V$ be the value with which $wr$ is invoked. Consider any correct server that replies to the *write* message of $wr$. Before replying, the server s-sends $V$ to all its registered readers. Since, $wr$ is invoked after time $t_1$, the server s-sends $V$ to $c$. Thus, $rd$ eventually s-receives $V$, and the condition in line 11 becomes true with $r' = 1$, and $[ts', id'] = TS$. □

# D    FW-Termination of ORCAS-A

The FW-Termination property requires that in every run where at most $f$ of the servers are faulty (and any number of clients are faulty), every write invocation by a correct process is complete, and moreover, either every read invocation by a correct process is complete, or infinitely many writes are invoked in the run. Observe that in a run with finite number of writes, a read cannot see new values and new timestamps forever. This case is same as the Case 1 in the proof of Lemma 1, whose proof does not rely on the fact that the read registers at the servers. Thus, if we want to satisfy FW-Termination, we can safely remove the registration. This change is presented in Figure 3. We would like to point out another slightly subtle change from the wait-free version. The s-send in line 24 of Figure 2 can now be replaced by a simple send.[7]

   With these above changes, the implementation has the desirable property that, if there are no new operations then the processes will eventually stop sending large messages (where a large message is a message containing either a written value or an encoded fragment of a written value).

# E    Timing guarantee of reads in ORCAS-A

In this section we study the timing performance of the read in the presence of concurrent writes. We make the following additional assumptions: (1) the time required for local computation is negligible,[8] (2) the links are timely, i.e., a message sent from a correct process to a correct process is received within a known constant delay, and (3) we consider a period in which at least $n - f$ server are alive and no alive process crashes and no crashed process recovers.

---

[7]In runs with finite number of writes, the servers eventually stop updating the variables whose values are sent in the $r\_ack$ messages. Therefore, if the condition in line 11 is never satisfied, the read operation sends an infinite number of *read* messages, and each server replies back with the same $r\_ack$ message an infinite number of times. Thus, from the property of fair-lossy channel, the reader eventually s-receives those $r\_ack$ messages. This modification does not work when there are infinite number of writes because in those runs the contents of the $r\_ack$ messages keep changing with new writes, and (if send replaces s-send then) no $r\_ack$ message is sent an infinite number of times.

[8]In fact, the time required for stable storage access and erasure coding the data can be significant, investigating which is outside the scope of this paper. We refer the readers to [8] for registers implementations that focus on reducing the number of stable storage accesses.

```
1:  In Figure 1 replace lines 25 to 34 by the following lines:
2:  upon receive ⟨write, ts′, wid′, rid′, V′⟩ from client c_i at server s_j do
3:     if ([ts′, wid′] >_{lex} [τ, δ]) and (V′ ≠ ⊥) then
4:        τ ← ts′; δ ← wid′; ρ ← 1; A_j ← V′
5:        store(τ, δ, ρ, A_j)
6:     s-send(⟨w_ack, ts′, wid′, rid′⟩, {c_i})


7:  In Figure 2 replace lines 21 to 24 by the following lines:
8:  upon receive ⟨read, rid⟩ from client c_i at server s_j do
9:     send(⟨r_ack, rid, τ, δ, ρ, A_j⟩, {c_i})
```

Figure 3: Modifications of ORCAS-A for Finite-Write termination

We now show that even in the presence of concurrent writes, the reader registration ensures that a read operation terminates within 5 communication steps. (The argument is similar to the to the proof of Lemma 1.) Suppose that some client $c$ invokes a read at time $t$. Then, within two communication steps of $t$, $c$ is registered at $n-f$ servers and it has computed $\Gamma$ in line 9 in Figure 2, say at time $t_1$. Now consider the $r\_ack$ message $m$ which was received by $c$, and which contained the timestamp $\Gamma$. If that message contains an unencoded value, then on receiving that message, the condition in line 11 is satisfied. Otherwise, if $m$ contains an encoded fragment, then the write with timestamp $\Gamma$ has completed its second phase by time $t_1$, and within one communication step of $t_1$, $c$ receives either (1) an unencoded value of that write, or enough encoded fragments of that write, or (2) an unencoded value of some write with a timestamp higher than $\Gamma$. (Recall that, from the algorithm, values with a smaller timestamp cannot overwrite values with a larger timestamp.) Thus, within one communication step of $t_1$, the condition in line 11 is satisfied, and within three communication steps of $t_1$, $c$ completes the read operation.

# F  Wait-free termination of ORCAS-B

**Lemma 2** *The ORCAS-B implementation is* wait-free.

(*Sketch.*) We sketch the only non-trivial argument in the proof of wait-freedom: the condition in line 13 of Figure 5 eventually becomes true if the client does not crash after invoking the read.
Suppose by contradiction that a client $c$ invokes a read operation $rd$, and then does not crash, and the condition in line 13 of Figure 5 is never satisfied in that operation. Thus $c$ executes the first repeat-until loop in the operation forever. Consider the two cases: (1) there is a time after which the read does not receive any new timestamp $[ts, id]$, and (2) otherwise. We now show that in both the cases the condition in line 13 is eventually satisfied.
Let $\gamma$ denote the timestamp assigned to $\Gamma$ in line 10, and $t$ be the time when the assignment was executed. (Note that line 10 is executed exactly once in $rd$.)
**Case 1.** Let $TS$ be the highest timestamp ever computed in line 8 of the read $rd$ in the run. By definition, $TS \geq_{lex} \gamma$. Let $wr$ be the write operation that corresponds to timestamp $TS$, and suppose it is invoked with a value $V$.
Let us call a timestamp $[ts, id]$, a *checked* timestamp if the read $rd$ ever receives an $r\_ack$ message with $L$ containing an element with timestamp $[ts, id]$ and $A = \emptyset$. We note that from line 44 of Figure 4, for an element $[ts, id, r, A, B]$ of $L$, a server sets $A = \emptyset$ only if it receives a trim message

18

```
 1: function initialization()                          28: upon receive ⟨get_ts, wid⟩ from client c_i at server s_j do
 2:    ts, wid, rid, x, y, z ← 0; r ← 1; A, B ← ∅; T_c ← timer()   29:    s-send(⟨ts_ack, wid, τ⟩, {c_i})
       {at every client}
 3:    τ, δ, τ_1, δ_1 ← 0; L ← {[0,0,1,∅,∅]}    {at every server    30: upon receive ⟨write, ts', wid', rid', r', A', B'⟩ from client
       s_j}                                                  c_i at server s_j do
                                                      31:    if rid' > 0 then
                                                      32:        R ← R \ {[rid', *, *, i]}
 4: function write (V) at client c_i                  33:    if [ts', wid'] >_lex [τ_1, δ_1] then
 5:    wid ← wid + 1; ts ← 0                          34:        if [ts', wid'] >_lex [τ, δ] then [τ, δ] ← [ts', wid']
 6:    store(wid, ts)                                 35:        L ← L ∪ {[ts', wid', r', A', B']}
 7:    trigger(T_c)                                   36:        store(L)
 8:    repeat                                         37:    for all [rid, ts, id, l] ∈ R do
 9:        send (⟨get_ts, wid⟩, S)                    38:        s-send(⟨r_ack, rid, L⟩, {c_l})
10:    until s-receive ⟨ts_ack, wid, *⟩ from n − f servers and   39:    s-send(⟨w_ack, ts', wid', rid'⟩, {c_i})
       expired(T_c)
11:    r ← (number of processes from which received
       ⟨ts_ack, wid, *⟩) − f                          40: upon s-received ⟨trim, isTrim, ts', id', rid'⟩ from client c_i
12:    ts ← 1 + max{ts_j : received ⟨ts_ack, wid, ts_j⟩}      at server s_j do
13:    store(ts, r, V)                                41:    modify L as follows
14:    x ← lcm(n − 2f, r); z ← x/(n−2f); y ← x/r      42:        if there is an element [ts'', id'', *, A'', B''] ∈ L such
15:    [C_1, ..., C_zn] ← encode(V, x, nz)                       that [ts'', id''] = [ts', id'] then
16:    1 ≤ l ≤ n, B_l ← [C_(l−1)z+1, ..., C_(l−1)z+y]  43:            if isTrim = false then B'' ← B'' ∪ A''
17:    1 ≤ l ≤ n, A_l ← [C_(l−1)z+y+1, ..., C_lz]      44:            A'' ← ∅; [τ_1, δ_1] ← [ts', wid']
18:    trigger(T_c)                                   45:        remove all elements [ts'', id'', *, *, *] ∈ L such
19:    repeat                                                    that [ts'', id''] <_lex [ts', id']
20:        ∀p_j ∈ S, send(⟨write, ts, wid, 0, r, A_j, B_j⟩, {p_j})   46:    store(L)
21:    until s-receive ⟨w_ack, ts, wid, 0⟩ from n − f servers and   47:    s-send(⟨trim_ack, ts', id', rid'⟩, {c_i})
       expired(T_c)
22:    if s-received ⟨w_ack, ts, wid, 0⟩ from r + f servers then   48: upon recovery() at server s_i do
       isTrim = true else isTrim = false              49:    L ← retrieve()
23:    repeat                                         50:    [τ, δ] ← highest timestamp in L
24:        S' ← set of servers from which s-received   51:    [τ_1, δ_1] ← lowest timestamp in L
       ⟨w_ack, ts, wid, 0⟩ until now
25:        send (⟨trim, isTrim, ts, wid, 0⟩, S')      52: upon recovery() at client c_i do
26:    until s-receive ⟨trim_ack, ts, wid, 0⟩ from n − f servers   53:    [rid, ts, wid, r, V] ← retrieve()
                                                      54:    if ts ≠ 0 then
27:    return(OK)                                     55:        execute lines 14 to 26 of this figure
```
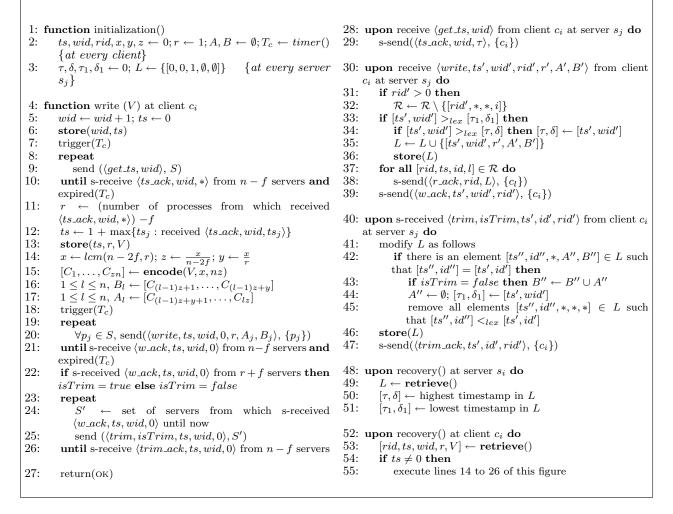
Figure 4: ORCAS-B: initialization, write and recovery procedures

from a write with timestamp $[ts, id]$. Thus, a checked timestamp corresponds to a write that has sent a trim message.

Since $TS$ is a checked timestamp, $wr$ has sent a trim message, and therefore, has completed its second phase. Let $r_1, x_1, z_1, y_1$ be the values of $r, x, z, y$, respectively in $wr$. Thus at least $r_1 + f$ servers have received $z$ fragments of $wr$, and at least $r_1$ among them are correct.

Next we note that all timestamps selected in line 8 of Figure 5 are checked timestamp. Thus, by definition, $TS$ is the highest checked timestamp received by the read $rd$. Also, it follows that no correct process ever receives a trim message with timestamp $TS' >_{lex} TS$ (otherwise, since $c$ is a correct process, the read $rd$ will eventually receive an $r\_ack$ message with checked timestamp $TS'$, thus violating the definition of $TS$). Thus, $[\tau_1, \delta_1]$ at correct processes is never higher than $TS$, and whenever a correct process receives a $write$ message with timestamp $TS$ for the first time, it adds the corresponding fragments in $L$. Also, some of those fragments may be later deleted (i.e., their $A$ set to $\emptyset$), but at least $y_1$ fragments are always retained in $L$.

Now there are two cases depending on whether write $wr$ received at least $r_1 + f$ $w\_ack$ replies in line 21 of Figure 4 or not. In the first case, at least $r_1$ of the replying servers are correct and have

```
 1: function read() at client c_i                         17:    1 ≤ l ≤ n, B_l ← [C_{(l-1)z+1}, ..., C_{(l-1)z+y}]
 2:    rid ← rid + 1; Γ ← 0; once ← false               18:    1 ≤ l ≤ n, A_l ← [C_{(l-1)z+y+1}, ..., C_{lz}]
 3:    store(rid)                                          19:    trigger(T_c)
 4:    repeat                                              20:    repeat
 5:       send(⟨read, rid⟩, S)                             21:       send(⟨write, ts, wid, rid, r, A_j, B_j⟩, S)
 6:       M ← {msg = ⟨r_ack, rid, *⟩ : s-received msg}    22:    until s-receive ⟨w_ack, ts, wid, rid⟩ from n − f servers
 7:       L ← {[*,*,*,*,*] ∈ L : ⟨r_ack, rid, L⟩ ∈ M}           and expired(T_c)
 8:       TS ← max_{lex}{[ts,id] : [ts,id,*,∅,*] ∈ L}    23:    if s-received ⟨w_ack, ts, wid, rid⟩ from r + f servers
 9:       if (M contains messages from at least n − f servers)    then isTrim = true else isTrim = false
          and (once = false) then                         24:    repeat
10:          Γ ← TS; once ← true                           25:       S' ← set of servers from which s-received
11:          if TS = [0,0] then return(⊥)                            ⟨w_ack, ts, wid, rid⟩ until now
12:       ∀(r,ts,id), Fragments(r,ts,id) ← union of all sets D    26:       send (⟨trim, isTrim, ts, wid, rid⟩, S')
          such that ([ts,id,r,D,*] ∈ L) ∨ ([ts,id,r,*,D] ∈ L)    27:    until s-receive ⟨trim_ack, ts, wid, rid⟩ from n − f
13:    until (once = true) and (∃(r,ts,id) such that              servers
          ([ts,id] ≥_{lex} Γ) and (|Fragments(r,ts,id)| ≥ lcm(n −    28:    return(V)
          2f, r)))
14:    x ← lcm(n − 2f, r); z ← x/(n-2f); y ← x/r          29: upon receive ⟨read, rid⟩ from c_i at server s_j do
15:    V ← decode((set of fragments satisfying condition in   30:    if R does not contain any [rid,*,*,i] then
       line 13), x, nz)                                   31:       R ← R ∪ [rid, τ, δ, i]
16:    [C_1, ..., C_{zn}] ← encode(V, x, nz)              32:       s-send(⟨r_ack, rid, L⟩, {c_i})
```

Figure 5: ORCAS-B: read procedure

stored the corresponding fragments. Eventually, $rd$ will receive at least $y_1$ of these fragments from each of these $r_1$ servers, and the condition in line 13 of Figure 5 will be satisfied with $r = r_1$, and $[ts, id] = TS$. In the second case, at least $n - f$ servers have replies with $w\_ack$ to $wr$, and at least $n - 2f$ of them are correct and have stored the corresponding fragments. Also in this case, no fragment in the corresponding element of $L$ is deleted − the contents of $A$ is added to $B$. Thus eventually, $wr$ will receive at least $z_1$ of these fragments from each of these $n - 2f$ servers, and the condition in line 13 of Figure 5 is satisfied with $r = r_1$, and $[ts, id] = TS$.

**Case 2.** Let $t_1 > t$ be a time such that (1) no correct server crashes after time $t_1$, and (2) all correct servers have received the $read$ message from $rd$ (and therefore, all correct servers have registered the read). Since there are bounded number of clients, if the read keeps on seeing new timestamps then there is at least one client $c'$ that invokes an infinite number of writes. As a client cannot invoke a new operation without completing the last invoked operation, $c'$ completes an infinite number of writes. Thus there is a complete write operation $wr$ such that (1) $wr$ has a timestamp $TS$ higher than $\gamma$, and (2) $wr$ is invoked after time $t_1$. Let $V$ be the value with which $wr$ is invoked and let $r_1$ be the value of $r$ in the write. Consider a correct server that replies to the $write$ message of $wr$. Before replying to the write, the server adds all $\frac{lcm(n-2f, r_1)}{n-2f}$ received fragments of $V$ to its set $L$, and then s-sends these fragments to all registered reads. Since there are at least $n - f$ correct servers, the read eventually s-receives at least $(n - f)\frac{lcm(n-2f, r_1)}{n-2f} > lcm(n - 2f, r_1)$ fragments of $V$, and the condition in line 13 is satisfied with $r = r_1$, and $[ts, id] = TS$. □