

Fork-Consistent Constructions From Registers

Matthias Majuntke¹, Dan Dobre², Christian Cachin³, and Neeraj Suri¹

¹{majuntke, suri}@cs.tu-darmstadt.de

Technische Universität Darmstadt, Hochschulstraße 10, 64289 Darmstadt, Germany

Phone: +49 6151 16 3121 / Fax: +49 6151 16 4310

²dan.dobre@neclab.eu

³cca@zurich.ibm.com

NEC Laboratories Europe

IBM Research – Zurich

Technical Report

Abstract

Users increasingly execute services online at remote providers, but they may have security concerns and not always trust the providers. Fork-consistent emulations offer one way to protect the clients of a remote service, which is usually correct but may suffer from Byzantine faults. They feature linearizability as long as the service behaves correctly, and gracefully degrade to fork-consistent semantics in case the service becomes faulty. This guarantees data integrity and service consistency to the clients.

All currently known fork-consistent emulations require the execution of non-trivial computation steps by the service. From a theoretical viewpoint, such a service constitutes a *read-modify-write* object, representing the strongest object in Herlihy’s wait-free hierarchy [12]. A read-modify-write object is much more powerful than a shared memory made of so-called *registers*, which lie in the weakest class of all shared objects in this hierarchy. In practical terms, it is important to reduce the complexity and cost of a remote service implementation as computation resources are typically more expensive than storage resources.

In this paper, we address the fundamental structure of a fork-consistent emulation and ask the question: Can one provide a fork-consistent emulation in which the service does not execute computation steps, but can be realized only by a shared memory? Surprisingly, the answer is yes. Specifically, we provide two such algorithms that can be built only from registers: A fork-linearizable construction of a universal type, in which operations are allowed to abort under concurrency, and a weakly fork-linearizable emulation of a shared memory that ensures wait-freedom when the registers are correct.

Keywords: distributed system, shared memory, fork-consistency, universal object, atomic register, Byzantine faults

1 Introduction

The increasing trend of executing services online “in the cloud” [23] offers many economic advantages, but also raises the challenge of guaranteeing security and strong consistency to its users. As the service is provided by a remote entity that wants to retain its customers, the service usually acts as specified. But online services may fail for various reasons, ranging from simply closing down (corresponding to a crash fault) to deliberate and sometimes malicious behavior (corresponding to a Byzantine fault).

For some kinds of services, cryptographic techniques can prevent a malicious provider from forging responses or snooping on customer data. But other violations are still possible in the asynchronous model considered here: for instance, when multiple isolated clients interact only through a remote provider, the latter may send diverging and inconsistent replies to the clients. In this context, “forking” consistency conditions [22, 8] offer a gracefully degrading solution because they make it much easier for the clients to detect such violations. More precisely, they ensure that if a Byzantine provider only *once* sent a wrong response to some client, then this client becomes *forever isolated* or *forked* from those other clients to which the provider responded differently. With this notion, clients may easily detect service misbehavior from a single inconsistent operation, e.g. by out-of-band communication.

Forking consistency conditions are often encapsulated in the notion of a *Byzantine emulation* [8], which ensures graceful degradation of the service’s semantics: If the service is correct, then operations execute atomically. In any other case, the clients still observe operations according to the forking consistency notion. Fork-consistency represents a safety property — after all, a faulty service may simply stop. The liveness property in a Byzantine emulation refers to the good case when the service behaves correctly.

Fork-linearizability [22, 8] ensures that clients always observe linearizable [14] service behavior and that two clients, once forked, will never again see each other’s updates to the system (i.e. they share the same history prefix up to the forking point). However, it has been found that fork-linearizable Byzantine emulations of a shared memory *cannot* always provide *wait-free* operations [8], i.e., some clients may be blocked because of other clients that execute operations concurrently. An escape is offered by the weaker liveness property of abortable emulations, which allow client operations to *abort* under contention [20]. As another alternative, the notion of *weak fork-linearizability* relaxes fork-linearizability in order to allow wait-free client operations in Byzantine emulations [7]. *Weak fork-linearizability* [7] allows two clients, after being forked, to observe a single operation of the other one (at-most-one-join), and that the real-time order induced by linearizability may be violated by the last operation of each client (weak real-time order).

In this paper, we explore the fundamental assumptions required for building a Byzantine service emulation. Up to now, all fork-consistent emulation protocols have required the service to execute non-trivial computation steps, i.e., the service must be implemented by an object of *universal* type [12], capable of *read-modify-write* operations [16]. We show the surprising result that this requirement can be dropped, and implement fork-consistent emulation protocols only from memory objects, so-called *registers*. They provide simple read and write operations and represent one of the weakest forms of computational objects. A long tradition of research has already addressed how to realize powerful abstractions from weaker base objects (e.g., [12, 2]).

Specifically, we propose the *first fork-linearizable Byzantine emulation* of a universal object only from *registers*. Our algorithm necessarily offers abortable operations because a wait-free construction of a universal object from registers is not possible in an asynchronous system using only registers [12]. Moreover, we give an algorithm for a *weakly fork-linearizable Byzantine emulation* of a shared memory only from registers. It allows wait-free client operations when the underlying registers are correct.

Our two algorithms may directly replace the computation-based constructions in the existing respective emulations of shared memory on Byzantine servers [20, 7, 24]. For instance, our second construction, which yields a weakly fork-linearizable Byzantine emulation, allows to eliminate the server code from Venus [24]. Currently, Venus runs server code implemented by a *cloud computing* service, but our construction may realize it from a *cloud storage* service. For practical systems this can make a big difference in cost because full-fledged servers or virtual machines (e.g., Amazon EC2) are typically more expensive than simple disks or cloud-based key-value stores (e.g., Amazon S3).

Note that although our approach uses a collection of registers, we refrain from making more specific failure assumptions on them. Our remote service is comprised of registers, and as soon as one register is faulty, we consider the service to be faulty. It is conceivable to use fault-prone registers in our algorithms. Standard methods implementing robust shared registers from fault-prone base registers show how to *tolerate* up to a fraction of Byzantine base registers [21]. This extension, which is orthogonal to our work, would further refine our notion of graceful service degradation with faulty base objects.

Related Work The notion of fork-linearizability was introduced by Mazières and Shasha [22]. They implemented a fork-linearizable multi-user storage system called SUNDRA. An improved fork-linearizable storage protocol is described by Cachin *et al.* [8]; it reduces the communication complexity compared to SUNDRA from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$. More recently, fork-linearizable Byzantine emulations have been extended to *universal services* [5]. All fork-linearizable emulations are blocking and sometimes require one client to wait for another client to complete [8].

In order to circumvent blocking the clients, Majuntke *et al.* [20] propose the first *abortable* fork-linearizable storage implementations. Their work takes up the notion of an abortable object introduced by Aguilera *et al.* [1]. They demonstrated, for the first time, how an abortable (and, hence, obstruction-free [13]) universal

object can be constructed from abortable registers, which are base objects weaker than registers. In more recent work, it has been shown that abortable objects can be boosted to wait-free objects in a partially synchronous system [3]. This makes our Byzantine emulations of abortable objects very attractive in practical systems.

Actually implemented systems offering data storage integrity through forking consistency semantics include SUNDR (LKMS) [17], which realizes the protocol of Mazières and Shasha [22]. Furthermore, Cachin *et al.* [6] add fork-linearizable semantics to the Subversion revision control system, such that integrity and consistency of the server can be verified. The “blind stone tablet” of Williams *et al.* [25] provides fork-linearizable semantics for an untrusted database server; it may abort conflicting operations. Using a relaxation of fork-linearizability, called *fork-* consistency*, Feldman *et al.* [10] introduce a lock-free implementation for online collaboration that protects consistency and integrity of the service against a malicious provider.

Cachin *et al.* [7] present the storage service FAUST, which emulates a shared memory in a wait-free manner by exploiting the notion of *weak fork-linearizability*. It relaxes fork-linearizability in two fundamental ways: (1) after being forked, two clients may observe each others’ operations once more and (2) the real-time order of the last operation of each client is not preserved. FAUST incorporates client-to-client communication in a higher layer, which ensures that all operations become eventually consistent over time (or the server is detected to misbehave). The Venus system [24] implements the mechanisms behind FAUST and describes a practical solution for ensuring integrity and consistency to the users of cloud storage.

Li and Mazières [18] study storage systems, built from $3f + 1$ server replicas, where more than f replicas are Byzantine faulty. Their storage protocol ensures *fork-* consistency*. Similar to weak fork-linearizability, fork-* consistency allows that two forked clients observe again at most one common operation.

Contributions We present, for the first time, Byzantine emulations with forking consistency conditions only from *registers*, instead of more powerful computation objects. Any number of registers may be affected by Byzantine failures. Our constructions are linearizable provided that the base registers are correct. The constructions are:

- A register-based abortable Byzantine emulation of a fork-linearizable universal type.
- A register-based wait-free Byzantine emulation of weak fork-linearizable shared memory.

In Section 1, we discuss related work; Section 2 introduces the underlying system model. The two main constructions are given in Sections 3 and 4. The paper concludes in Section 5. The correctness proofs of the protocols have been moved to the Appendix (page 13).

2 System Model

We consider a distributed system consisting of $n > 1$ *clients* C_1, \dots, C_n that communicate through shared *objects*. Each such base object has a *type* which is given by a set of *invocations*, a set of *responses*, and by its *sequential specification*. The sequential specification defines the allowed sequences of invocations and responses. An *invocation* and the corresponding *response* constitute an *operation* of an object. A collection of base objects is used to implement high-level objects, where clients execute algorithm A , consisting of n state machines A_1, \dots, A_n (where C_i implements A_i). When client C_i receives an *invocation* of an operation to the high-level object, it takes steps of A_i , where it (1) either invokes an operation on some base object, (2) or receives the response to its previous invocation to a base object, (3) or it performs some local computation. At the end of a step, C_i changes its local state and possibly returns a response to the pending high-level operation.

An *execution* of algorithm A is defined as the (interleaved) sequence of invocation and response events. Every execution induces a *history* which is the sequence of invocations and responses of the high-level operations. If σ is a history of an execution of algorithm A , then $\sigma|_{C_i}$ denotes the subsequence of σ containing all events of client C_i . For sequence σ and operation o , $\sigma|_o$ denotes the prefix of σ that ends with the last event of o . We say that a response *matches* an invocation, if both are events of the same operation. An operation is called *complete*, if there exists a matching response to its invocation, else *incomplete*. We assume that each

client invokes a new operation only after the previous operation has completed. A history consisting only of matching invocation/response pairs is called *well-formed*. Operation o *precedes* operation o' in a sequence of events σ ($o <_{\sigma} o'$) iff o is complete and the response of o happens before the invocation of o' . If o precedes o' we denote o and o' as *sequential*, if neither one precedes the other, then o and o' are said to be *concurrent*.

For the proposed *abortable* construction (Sec. 3), we introduce the special response ABORT. A complete operation o is called *unsuccessful* (“ o is aborted”), if it returns ABORT, else it is called *successful* (“ o successfully completes”). The formal definition of an *abortable* object comprises a non-triviality property which allows aborts only under concurrency [1].

Clients may fail by *crashing*, i.e. they stop taking steps and hence, the last operation of each client might be *incomplete*. Base objects may deviate arbitrarily from their specification exhibiting *non-responsive-arbitrary faults* [15] (called *Byzantine*). Clients have access to a digital signature scheme used by each client to *sign* its data such that any other client can determine the authenticity of a datum by *verifying* the corresponding signature. We assume that signatures cannot be forged.

All constructions appearing in this paper are based on *atomic registers*. An atomic register provides two operations, *read* and *write*¹. Operation $write(v)$ stores value v from domain *Values* into the register. A call of $read()$ returns the latest written value from the register or the special value \perp if no value has been written. As the register is atomic, its history satisfies linearizability [12], i.e. operations seem to appear as sequential, atomic events². Further, the atomic registers used allow single-writer-multiple-reader access (SWMR), i.e. to each register we assign a dedicated client that may call *write* and *read*, while all other clients may only call *read* to that register.

A sequence of operations π satisfies *weak real-time* order of σ if π , excluding the last operation of each client in π , satisfies real-time order of σ . *Causality* between two operations depends on the type of the implemented object³. For two operations of a shared memory o and o' in σ , o *causally precedes* o' ($o \rightarrow_{\sigma} o'$), if o, o' are called by the same client and o happens before o' , or if o' is a READ operation that returns the value written by WRITE operation o . The next definition formalizes the notion of *fork-linearizability* [8] and *weak fork-linearizability* [7]; for a formal definition of the term *possible view* as well as the above-mentioned notions we refer to the Appendix.

Definition 1. Let σ be a history of an object of type T and for each client C_i there exists a sequence of events π_i such that π_i is a possible view of σ at C_i with respect to T .

History σ is *fork-linearizable* with respect to object type T if for each client C_i :

1. π_i preserves the real-time order of σ , and
2. for every client C_j and for every $o \in \pi_i \cap \pi_j$, it holds $\pi_i|_o = \pi_j|_o$.

History σ is *weak fork-linearizable* with respect to object type T if for each client C_i :

1. π_i preserves the weak real-time order of σ , and
2. for every operation $o \in \pi_i$ and every operation $o' \in \sigma$ such that $o' \rightarrow_{\sigma} o$, it holds that $o' \in \pi_i$ and that $o' <_{\pi_i} o$, and
3. (At-most-one-join) for every client C_j and every two operations $o, o' \in \pi_i \cap \pi_j$ by the same client such that $o <_{\sigma} o'$, it holds $\pi_i|_o = \pi_j|_o$.

The notion of a *Byzantine emulation* [8] allows us to formally define the safety and liveness properties of our protocols. Note that the liveness condition of abortable operations is weaker than *wait-freedom* but still not weaker than *obstruction-freedom* [1].

¹We type operation calls to base registers in *italic* font and calls to constructed objects in CAPITALS.

²Hence, the “latest written value” is well-defined.

³As causality is needed to define *weak fork-linearizability*, here, we give causality for a *shared memory*, which is the type we implement with weak fork-linearizability.

Definition 2. An algorithm A *emulates* an object of type T on a set of Byzantine base objects B with $\{\text{fork}|\text{weak fork}\}$ -linearizability whenever the following conditions hold:

1. If all objects in set B are correct, the history of every fair⁴ and well-formed execution of A is linearizable with respect to type T , and
2. the history of every fair and well-formed execution of A is $\{\text{fork}|\text{weak fork}\}$ -linearizable with respect to type T .

Such an emulation is *wait-free* (*abortable* resp.), iff every fair and well-formed execution of the protocol with correct base objects is wait-free [12] (*abortable* [1] resp.).

3 A Fork-linearizable Universal Type

In this section we present as our first main contribution an abortable fork-linearizable Byzantine emulation of a universal type implemented from atomic registers. The shared object ensures fork-linearizability in the presence of any number of faulty base registers. High-level operations are *abortable* [1], i.e. under concurrency, the special response ABORT may be returned. The functionality of a universal type T is encoded in the procedure APPLY_T . For client C_i , state s and operation o , $\text{APPLY}_T(s, o, i)$ returns (s', res) , where s' is the new state of the universal object, res the computation result, and where the sequence of invoking $\text{APPLY}_T(s, o, i)$ and returning (s', res) is defined by the sequential specification of type T .

Our algorithm uses timestamp vectors called *versions* whose order reflects the real-time order in which operations are applied to the shared object. Each operation carries a version and the linearization of operations is achieved through the use of an INC&READ counter object C with two atomic operations INC&READ and READ. An invocation to $\text{INC\&READ}(C)$ advances the counter object C and returns a value which is higher than any value returned before, and $\text{READ}(C)$ returns the current value of the counter object. An implementation of the INC&READ counter is given in Algorithm 3 in Appendix B together with its formal properties. Our implementation uses wait-free atomic registers as base objects which makes it a wait-free variant of the abortable INC&READ counter described by Aguilera *et al.* [1].

3.1 Algorithm Ideas

Universal Type To implement universal type T , we use n SWMR registers R_1, \dots, R_n such that client C_i can read from all registers but may write only to R_i . The registers store states of the universal object. To implement high-level operations, client C_i reads from the register which holds the most current state, applies the relevant state transformation, and writes the new state to R_i . Note, that all information are digitally signed by the clients as base objects are untrusted. Thereby, operations “affect” each other which leads to the following relation on operations: Operation o of C_i *affects* operation o' of C_j , if during o' , C_j is able to verify the signature of C_i on state s that has been written during o and if C_j executes APPLY_T on s during o' ; further, an operation of C_i *affects* each later operation of C_i .

Concurrency detection We allow operations to abort under concurrency for two reasons: there is no wait-free construction of a universal type from registers, as shown by Herlihy [12], and no fork-linearizable protocol can be wait-free in all executions, as shown in a more recent work of Cachin *et al.* [8]. Cachin’s impossibility is based on two runs, indistinguishable for the reader: In the first run a READ operation does not return value v as it is concurrently written, while in the second run v has been previously written and is hidden by malicious registers. To avoid such a situation, our protocol implements a concurrency detection mechanism [1] using INC&READ counter object C . If concurrency is detected, a pending operation is aborted. At the invocation of a high-level operation o , our protocol calls $\text{INC\&READ}(C)$ and remembers the timestamp returned. At the end of o , $\text{READ}(C)$ is executed to check whether counter C still returns the same timestamp. If not, another operation o' was invoked during o — thus, o is aborted. Else, if at the end of o C has not been changed, all

⁴For a formal definition we refer to standard literature [19]

successful operations either terminated before o or will be invoked after o has terminated. This is because the timestamps, returned from INC&READ, are used to linearize operations: The current state is written together with the timestamp, and the timestamp is used to determine the most recent state. Hence, all other operations invoked so far write a state with a lower timestamp than o . Consequently, such operations are linearized before o and only the state written by o can be read by later operations.

Fork-Linearizability In addition to the timestamp from INC&READ counter C , each operation is assigned a vector of timestamps of length n , called *version*. The order relation \leq defined on versions respects real-time order and the "affected by" relation on operations. The idea is that each operation reads the most recent version from the storage, increments its own entry and writes the new version back to the storage. Thereby, each operation checks, if the version it reads, has been affected by the version of its own last successful operation, i.e. one which was not aborted. If the last successful operation of client C_i is hidden from C_j , then C_i does not accept operations of C_j as they have *not* been affected by the last successful operation of C_i . This ensures that the views of the clients after a forking attack are not rejoined. This principle is based on ideas of Mazières and Shasha [22], and Cachin *et al.* [8]. To apply it to this work, we have to add a specific handling for aborted operations: If operation o of client C_i is aborted, C_i cannot expect that o will affect later operations. However, it is still possible that some operation of C_j is affected by aborted o . In this case we call o *relevant* for C_j (Definition 8 in Appendix B).

3.2 Description of Algorithm 1

We now describe the steps performed by client C_i when executing high-level operation o . The algorithm is given as Algorithm 1, the variables used are collected in Variables 5 (see Appendix A).

The protocol is framed by INC&READ(C) and READ(C) calls to the counter object C implementing the concurrency detection mechanism (lines 1.2 and 1.14). If the returned timestamps are not equal, the operation is aborted in line 1.16. In lines 1.3–1.5, the client reads from all atomic registers R_1, \dots, R_n and determines by means of the assigned timestamps the index l of the register holding the latest written data $\langle ts_l, V_l, s_l, sig_l \rangle$, where ts_l is a timestamp, V_l is the version, s_l is the state and sig_l is a signature. If some data have been written to R_l , the signature of the content of R_l is verified (line 1.6). Then, client C_i checks whether the read version V_l is not smaller than V_{suc} the version of its own last successful operation (line 1.7). When the check is passed the new state of the universal object and the computation result is computed by calling $APPLY_T(s_l, o, i)$ (line 1.8). Finally the new version for operation o has to be computed. This is done by taking the per-entry maximum of version V , which is the local version of C_i , and V_l , and by incrementing the i th entry (lines 1.9–1.11). After signing the current timestamp, the new version V , and new state s in line 1.12, client C_i writes ts , V , s and the signature into register R_i (line 1.13). If operation o is successful, version V is stored as last successful version V_{suc} and the computation result is returned (lines 1.17–1.19).

3.3 Correctness Arguments

In this section we argue why Algorithm 1 satisfies fork-linearizability. The goal is to construct for each client C_i a view π_i of σ that satisfies the properties of fork-linearizability. To construct π_i , we simplify our argumentation by ignoring operations that are not relevant for C_i . Recall, any operation is *relevant* for client C_i that affects C_i 's last successful operation. Hence, operations that are not relevant for client C_i do not change the object's state from C_i 's point of view. Thus, we can order them arbitrarily among the operations in π_i and the resulting sequences still satisfy fork-linearizability.

The idea behind the construction of the π_i in the proof is that operations are ordered according to their assigned versions. The proof shows that this order respects the "affected by" relation, the sequential specification of a universal type, and the real-time order. As during an operation the new version is computed using the client's last version and the read version, proving "affected by" and real-time order is straightforward. The core of the proof is to show that the order of version also respects the sequential specification. We sketch the intuition behind this with the following argument leading to a contradiction:

Algorithm 1: Universal Object Implementation, Algorithm of Client i

```

1.1 EXECUTE( $o$ ) do
1.2    $ts \leftarrow \text{INC\&READ}(C)$                                 /* increment and read from counter */
1.3   for  $j = 1, \dots, n$  do
1.4      $\langle ts_j, V_j, s_j, sig_j \rangle \leftarrow \text{read}(R_j)$       /* low-level atomic read */
1.5   let  $l$  be such that  $ts_l = \max_{1 \leq j \leq n}(ts_j)$           /* find register with most recent data */
1.6   if  $V_l \neq [0 \dots 0] \wedge \neg \text{verify}_l(sig_l, \langle ts_l, V_l, s_l \rangle)$  then halt /* signature verified? */
1.7   if  $\exists k : V_{suc}[k] > V_l[k]$  then halt                    /* fork-linearizability check passed? */
1.8    $\langle s, res \rangle \leftarrow \text{APPLY}_T(s_l, o, i)$               /* compute new state + result */
1.9   for  $j = 1, \dots, n, j \neq i$  do
1.10     $V[j] \leftarrow \max(V[j], V_l[j])$                        /* determine
1.11     $V[i] \leftarrow V[i] + 1$                                   new version */
1.12     $sig \leftarrow \text{sign}_i(ts || V || s)$                    /* signature on ts, version, state */
1.13     $\text{write}(R_i, \langle ts, V, s, sig \rangle)$                  /* low-level atomic write */
1.14     $ts' \leftarrow \text{READ}(C)$                                 /* read from counter */
1.15    if  $ts \neq ts'$  then
1.16      return ABORT                                           /* concurrency detected */
1.17    else
1.18       $V_{suc} \leftarrow V$                                     /* reset last successful version */
1.19      return  $res$                                            /* return result */

```

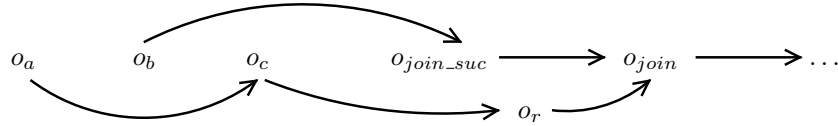


Figure 1: Correctness Idea of Algorithm 1. Arrows denote the “affected by” relation.

Assume that some operation o_c is not affected by the most recent state of the universal object, which has been written by relevant operation o_b , but is affected by an older state written by operation o_a . In this case, the clients of o_b and o_c are forked, and neither o_b nor o_c affect each other. We argue, that in such a situation, there is no relevant operation that has been affected by both o_b and o_c , as such an operation would join the two clients violating fork-consistency. We assume for contradiction, that a relevant operation o_{join} of client C_{join} , affected by o_b and o_c exists which is also the first among such operations (see Figure 1). Operation o_{join} is affected by o_{join_suc} , the last successful operation of C_{join} previous to o_{join} , and by o_r that wrote the state which is read during o_{join} . Hence, without loss of generality o_{join_suc} is affected by o_b while o_r is affected by o_c . During operation o_{join_suc} , client C_{join} raises its value in the version to $V[join]_{join_suc}$. This implies that o_{join} only accepts versions where the $join$ th entry is at least $V[join]_{join_suc}$ (line 1.7). As o_{join_suc} is not on the path of “affected by” relations from o_c to o_r , o_{join} would block while reading the state of o_r which is a contradiction. Thus, o_{join} does not exist.

Finally, it follows directly from the described construction, that sequences π_i satisfy the no-join property. To complete the correctness proof of the Byzantine emulation, we show that when all base objects are correct, no operation blocks and that no operation trivially aborts.

4 A Weak Fork-Linearizable Shared Memory

In this section we describe as our second contribution a wait-free, weak fork-linearizable Byzantine emulation of a shared memory implemented from atomic registers. The presented construction satisfies weak fork-linearizability in the presence of any number of faulty base objects. The implemented shared memory provides n atomic registers, such that each client can write to one dedicated register exclusively and may read from all registers. Operation $\text{WRITE}(v)$, called by client C_i , writes value v to C_i ’s register. Operation $\text{READ}(i)$ returns

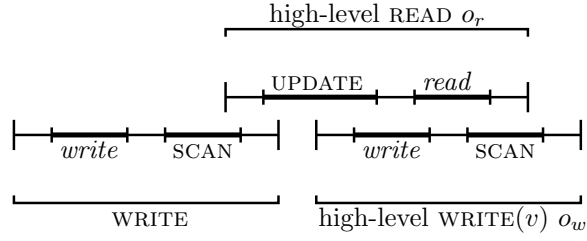


Figure 2: Basic principle implemented by Algorithm 2.

the last written value from C_i 's register, and may be called by any client. Our algorithm makes use of an atomic single-writer snapshot object S with n components [4, 11]. Snapshot object S provides two atomic operations: $\text{UPDATE}(d, S, i)$, that changes the state of component i of S to d , and $\text{SCAN}(S)$ that returns vector (d_1, \dots, d_n) such that d_i is the state of component i of S , $i = 1 \dots, n$. Formally, d_i is the state written by the last UPDATE to component i prior to SCAN . It has been shown, that such a shared snapshot object can be wait-free implemented only from registers [4, 11].

4.1 Algorithm Ideas

Each client locally maintains a timestamp that respects causality and real-time order of its *own* operations. As the basic principle, during each operation this timestamp is written to the shared memory and timestamps left by other operations are read. For each client C_i our implementation uses two registers only C_i may write to, but which can be read by all clients. The first one is needed to store value and timestamp written by C_i 's WRITE operations and is implemented by a SWMR atomic register W_i (i.e. registers W_1, \dots, W_n in total). The second “register” is required to store the latest timestamp of C_i 's READ operations. It is implemented as the i th component within the single-writer snapshot object with n components, S .

During $\text{READ}(j)$ operation of C_i , C_i 's current timestamp is written to S using UPDATE , thereafter, C_i reads a timestamp-value pair from register W_j (using low-level read). High-level $\text{WRITE}(v)$ of C_i proceeds analogously: C_i writes its current timestamp plus value v to register W_i using low-level write , thereafter, it reads *all* components from S using SCAN . By this, operations are able to observe each other, as expressed in the relation “seen”: We say that a WRITE operation o_w of C_j *sees* a READ operation o_r of C_i with timestamp ts if C_i digitally signed ts and updated the i th component of S by signed ts during o_r and, if during o_w , C_j scanned S and was able to verify the signature of C_i on ts ; READ operation o_r *sees* WRITE operation o_w if o_r returns the value written by o_w .

This construction guarantees the following property on interleaved high-level operations: Whenever high-level $\text{READ}(j)$ o_r of C_i and $\text{WRITE}(v)$ o_w of C_j appear in an execution such that o_r does not return v but a value written before v , then, by regularity of the atomic base registers, $o_w.\text{write}$ ⁵ does not precede $o_r.\text{read}$, i.e., $o_r.\text{read}$ has been invoked *before* $o_w.\text{write}$ finishes. Consequently, $o_r.\text{UPDATE}$ precedes $o_w.\text{SCAN}$ (see Figure 2). Thus, if o_r does not “see” o_w , then o_w “sees” o_r . A similar property on interleaving operations has also been leveraged in our previous work [9] as well as by Aguilera *et al.* [2].

We can expect that client C_j writes information during its next WRITE operation such that future operations of C_i may verify whether operation o_w actually has seen operation o_r . More concrete, if READ o_r has seen WRITE o_w then the client checks during o_r whether the next WRITE operation after o_w (of the same client as o_w), has seen READ operation o_r or a newer one. Else, the base objects are faulty, as shown in the following example: Let o_w and o'_w be two sequential WRITE operations of C_i , o'_w precedes READ operation o_r of C_j but it is hidden by the malicious base objects such that o_r sees only o_w . As o'_w precedes o_r , o'_w cannot see o_r . However, as o_r sees o_w , it expects that o'_w will see o_r . The next WRITE operation o''_w of C_i will write this information. If client C_j sees o''_w , which would violate weak fork-linearizability, the check, explained above, is not passed.

⁵The notation $x.y$ denotes the call of low-level operation y during high-level operation x .

4.2 Description of Algorithm 2

Algorithm 2: Weak Fork-Linearizable Memory for n Clients, Algorithm of Client C_i

```

2.1 READ( $j$ ) do
2.2    $ots \leftarrow ots + 1$                                 /* increment timestamp */
2.3    $sig \leftarrow \text{sign}_i(ots)$                         /* signature on timestamp */
2.4   UPDATE( $(ots, sig), S, i$ )                            /* update call to snapshot object */
2.5    $(wv, wts, r\_read\_seen, r\_write\_seen, sig) \leftarrow \text{read}(W_j)$  /* low-level atomic read */
2.6   if not verify $_j(sig)$  then halt                    /* signature verified? */
2.7    $read\_seen \leftarrow \text{merge}(read\_seen, r\_read\_seen)$  /* update read\_seen */
2.8    $read\_seen[i][j] \leftarrow read\_seen[i][j].\text{add}((ots, wts))$  /* add seen write */
2.9   check()                                              /* check passed? */
2.10   $write\_seen \leftarrow \text{merge}(write\_seen, r\_write\_seen)$  /* update write\_seen */
2.11  return  $wv$                                           /* return read value */

2.12 WRITE( $v$ ) do
2.13   $ots \leftarrow ots + 1$                                 /* increment timestamp */
2.14   $sig \leftarrow \text{sign}_i(v, ots, read\_seen, write\_seen)$  /* signature on timestamp */
2.15  write( $(v, ots, read\_seen, write\_seen, sig), W_i$ )      /* low-level atomic write */
2.16   $\langle (tmp_1, sig_1), \dots, (tmp_n, sig_n) \rangle \leftarrow \text{SCAN}(S)$  /* scan call to snapshot object */
2.17  for  $k = 1, \dots, n$  do
2.18    if not verify $_k(sig_k)$  then halt                /* signature verified? */
2.19     $write\_seen[i][k] \leftarrow write\_seen[i][k].\text{add}((tmp_k, ots))$  /* add all seen reads */
2.20  return OK                                           /* successfully return */

2.21 check() do
2.22  for  $k = 1, \dots, n$  do
2.23    forall  $(r, w) \in read\_seen[k][i]$  do
2.24      /* check if own writes have seen read operations reading my values */
2.25      if  $\exists (r', w') \in write\_seen[i][k]$  s.t.  $w' > w$  and  $w'$  minimal then
2.26        if  $r' < r$  then halt
2.27    forall  $(r, w) \in read\_seen[i][k]$  do
2.28      /* check if own reads have been seen by other's write operations */
2.29      if  $\exists (r', w') \in r\_write\_seen[k][i]$  s.t.  $w' > w$  and  $w'$  minimal then
2.30        if  $r' < r$  then halt

```

This section explains the steps taken by client C_i to implement high-level READ and WRITE operations. The algorithm is given as Algorithm 2, its variables in Variables 6 (see Appendix A).

At invocation of high-level READ(j), client C_i increments its local timestamp and generates a digital signature of it. The signed timestamp is stored to snapshot object S using operation UPDATE($(ots, sig), S, i$) (lines 2.2–2.4). Then, client C_i reads register W_j and verifies the signature (line 2.5–2.6). The content of register W_j contains the written value wv , the corresponding timestamp wts , as well as two matrices r_read_seen and r_write_seen . Both matrices are of size $n \times n$ where each entry holds a set of integer pairs (r, w) . Client C_i maintains a variable $read_seen$ of the same type, where a pair $(r, w) \in read_seen[i][j]$ denotes that READ of client C_i with timestamp r has seen WRITE of client C_j with timestamp w . Analogously, client C_i maintains a second matrix $write_seen$, where $(r, w) \in write_seen[i][j]$ denotes that WRITE of client C_i with timestamp w has seen READ of client C_j with timestamp r . In the next step (line 2.7), client C_i “merges” variables r_read_seen and $read_seen$. The merge procedure returns for each entry of two $n \times n$ set matrices A, B set $A[i][j] \cup B[i][j]$, $i, j = 1, \dots, n$. Then, C_i adds a pair consisting of its current timestamp and timestamp wts from W_j to $read_seen[i][j]$. To ensure weak fork-linearizability, client C_i calls procedure “check” (line 2.9). If all checks are passed, C_i merges r_write_seen and $write_seen$ and returns value wv (lines 2.10–2.11).

At invocation of WRITE(v), client C_i increments its timestamp (line 2.13). It digitally signs value v , its

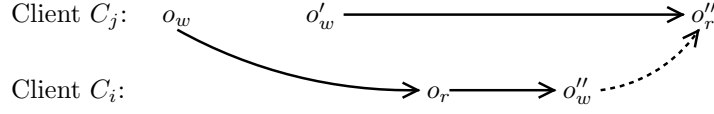


Figure 3: Correctness Ideas of Algorithm 2. Arrows denote the “seen” relation.

timestamp, and variables *read_seen* and *write_seen* to write to register W_i (lines 2.14–2.15). Next, it reads all timestamps of READS by calling SCAN to snapshot object S (line 2.16). All entries in S are digitally signed and thus client C_i verifies the signatures (line 2.18). Then, it adds to all sets $write_seen[i][k]$ ($k = 1, \dots, n$) a pair consisting of the timestamp of the k th component of S and C_i ’s current timestamp (line 2.19). Finally, client C_i successfully returns (line 2.20).

Procedure “check” implements the principle sketched in section 4.1 for n clients. It ensures that *weak fork-linearizability* is never violated. The procedure, called by C_i during $READ(j)$ (line 2.21), moves through a loop performing two checks: The first check (line 2.24–2.25) considers the information left by clients during $READ(i)$ operations (this information is stored in the i th column of *read_seen*). If $READ(i)$ with timestamp r of client C_k has seen WRITE of C_i with timestamp w , then it is tested whether the next WRITE of C_i has read (using SCAN) timestamp r or higher of client C_k . The check uses the local *write_seen* variable of C_i . The second check (line 2.27–2.28) reviews the information left by client C_i during any $READ(k)$ (which is kept in the i th row of *read_seen*). If $READ(k)$ with timestamp r of client C_i has seen WRITE of C_k with timestamp w , then we check whether the next WRITE of C_k has read (using SCAN) timestamp r or higher of client C_i . This check requires matrix r_write_seen , which has been fetched from W_j in line 2.5 before procedure “check” is called.

4.3 Correctness Arguments

In this section we give the intuition why Algorithm 2 satisfies the properties of a wait-free Byzantine emulation of a shared memory with weak fork-linearizability. Intuitively, the definition of weak fork-linearizability requires for each client C_i to construct a sequence π_i such that causality among operations, the sequential specification a shared memory, and weak real-time order is satisfied, and that two sequences π_i and π_j share the same prefix up to the second last common operation (at-most-one-join). The proof proceeds in steps, where in the first step all operations that have to be included in sequence π_i are causally ordered. Next, this order is extended such that it additionally respects the sequential specification. Intuitively, as all written values are digitally signed, the sequential specification never interferes with causality. The hardest step is to prove, that this order can be further refined such that it does not violate the weak real-time order. The intuition for this is given below as a proof by contradiction:

We assume that $READ(j)$ operation o_r of client C_i does not return the latest value, written by WRITE operation o'_w , but an older value written by operation o_w (see Figure 3). Further, let o_r be not the last operation of C_i . During operation o_r , the pair (r, w) ⁶ is added to set $read_seen[i][j]$. The data written by the next WRITE operation o''_w of C_i contains this information. Now, the algorithm prevents client C_j from reading the value written by o''_w which would violate weak real-time order (as o_r is ordered before o'_w according to the sequential specification). When during o''_w C_j sees operation o''_w , it finds the pair (r, w) in r_read_seen . As o'_w precedes o_r , it could not have seen o_r , thus $write_seen[j][i]$ contains a pair (r', w') such that $r' < r$ and the check in line 2.25 is not passed. Hence, operation o''_w of client C_j would block — a contradiction. This implies that such a situation does not appear and the constructed order of operations also satisfies weak real-time order.

As the last step, showing that the sequences π_i satisfy the at-most-one-join property follows directly from a simple construction argument. To prove liveness, as required in the definition of a Byzantine emulation (Definition 2), we show that no operation blocks when all base objects are correct, which follows from the principle sketched in section 4.1 as in this case all checks are passed.

⁶We assume that operation o_x is assigned timestamp x .

5 Analysis & Conclusions

The abortable construction in Algorithm 1 requires n atomic registers plus n additional ones to implement the INC&READ counter. The presented construction has an overall communication complexity of $O(n^2)$, as the size of the version vectors used in Algorithm 1 is linear in the number of clients n and as a linear number of such version vectors are exchanged per operation. In contrast, the *lock-step* protocol of Cachin *et al.* [8], also based on linear size version vectors, has an overall communication complexity of $O(n)$. This difference results from the fact that the server objects used by Cachin *et al.* are computationally strong enough to select the latest written version vector while in Algorithm 1 the client is required to read from *all* register objects to find the latest one by itself. For the implementation of Algorithm 2, we need n atomic registers plus $2n$ additional ones for the atomic snapshot object. Algorithm 2, uses matrices of size $n \times n$ where the size of each entry depends on the total number of operations N , resulting in a communication complexity of $O(N \cdot n^2)$. We leave for future research whether this complexity can be reduced by implementing a “garbage collection”. However, both of our algorithms require only a linear number of base registers.

We have shown by ways of two protocols as a first known result that fork-consistent semantics can be implemented only from registers. Our first protocol satisfies fork-linearizability and implements a shared object of universal type. Similar to non-fork-consistent universal constructions from registers, our protocol may abort operations under concurrency. Hence, fork-linearizability may be “added” to such protocols without making additional assumptions. Our second protocol implements a shared memory object that ensures *weak* fork-linearizability and where operations are wait-free as long as the base registers behave correctly. Weak fork-linearizability is the strongest known fork-consistency property that may be implemented in a wait-free manner. Although it weakens fork-linearizability, it has shown to be of practical relevance [7]. Moreover, our second algorithm shows for the first time that registers are sufficient to implement a fork-consistent shared memory. So far, all existing implementations are based on computationally stronger objects (featuring read-modify-write operations [16]). We leave as an open question whether there is a weak fork-linearizable construction of a universal type providing a stronger liveness condition than *abortable* in the fault-free case.

References

- [1] Marcos K. Aguilera, Svend Frolund, Vassos Hadzilacos, Stephanie L. Horn, and Sam Toueg. Abortable and Query-Abortable Objects and Their Efficient Implementation. In *PODC: Principles of distributed computing*, pages 23–32, New York, NY, USA, 2007. ACM.
- [2] Marcos K. Aguilera, Idit Keidar, Dahlia Malkhi, and Alexander Shraer. Dynamic Atomic Storage Without Consensus. *J. ACM*, 58:7:1–7:32, April 2011.
- [3] Marcos K. Aguilera and Sam Toueg. Timeliness-Based Wait-Freedom: A Gracefully Degrading Progress Condition. In *PODC '08: Proceedings of the twenty-seventh ACM symposium on Principles of distributed computing*, pages 305–314, New York, NY, USA, 2008. ACM.
- [4] Hagit Attiya, Rachid Guerraoui, and Eric Ruppert. Partial Snapshot Objects. In *Proc. SPAA*, pages 336–343, 2008.
- [5] Christian Cachin. Integrity and Consistency for Untrusted Services. In *Proceedings of the 37th international conference on Current trends in theory and practice of computer science*, SOFSEM'11, pages 1–14, Berlin, Heidelberg, 2011. Springer-Verlag.
- [6] Christian Cachin and Martin Geisler. Integrity Protection for Revision Control. In *Proceedings of the 7th International Conference on Applied Cryptography and Network Security*, ACNS '09, pages 382–399, Berlin, Heidelberg, 2009. Springer-Verlag.
- [7] Christian Cachin, Idit Keidar, and Alexander Shraer. Fail-Aware Untrusted Storage. *SIAM Journal on Computing*, 40(2):493–533, April 2011.

- [8] Christian Cachin, Abhi Shelat, and Alexander Shraer. Efficient Fork-Linearizable Access to Untrusted Shared Memory. In *PODC*, pages 129–138, New York, NY, USA, 2007. ACM.
- [9] Dan Dobre, Matthias Majuntke, and Neeraj Suri. On the time-complexity of robust and amnesic storage. In *OPODIS*, pages 197–216, 2008.
- [10] Ariel J. Feldman, William P. Zeller, Michael J. Freedman, and Edward W. Felten. SPORC: Group Collaboration on Untrusted Resources. In *Proc. 9th Symposium on Operating Systems Design and Implementation (OSDI 10)*, Vancouver, BC, October 2010.
- [11] Faith Ellen Fich. How Hard Is It to Take a Snapshot? In *Proc. SOFSEM*, pages 28–37, 2005.
- [12] Maurice Herlihy. Wait-Free Synchronization. *ACM Trans. Program. Lang. Syst.*, 13(1):124–149, 1991.
- [13] Maurice Herlihy, Victor Luchangco, and Mark Moir. Obstruction-Free Synchronization: Double-Ended Queues as an Example. In *ICDCS*, page 522, Washington, DC, USA, 2003. IEEE Computer Society.
- [14] Maurice P. Herlihy and Jeannette M. Wing. Linearizability: A Correctness Condition for Concurrent Objects. *ACM Trans. Program. Lang. Syst.*, 12(3):463–492, 1990.
- [15] Prasad Jayanti, Tushar Deepak Chandra, and Sam Toueg. Fault-tolerant Wait-free Shared Objects. *J. ACM*, 45(3):451–500, 1998.
- [16] Clyde P. Kruskal, Larry Rudolph, and Marc Snir. Efficient Synchronization of Multiprocessors with Shared Memory. *ACM Trans. Program. Lang. Syst.*, 10:579–601, October 1988.
- [17] Jinyuan Li, Maxwell Krohn, David Mazières, and Dennis Shasha. Secure Untrusted Data Repository (SUNDR). In *Proc. 6th Symp. Operating Systems Design and Implementation (OSDI 04)*, pages 121–136, 2004.
- [18] Jinyuan Li and David Mazières. Beyond One-Third Faulty Replicas in Byzantine Fault Tolerant Systems. In *Proc. NSDI*, 2007.
- [19] Nancy A. Lynch. *Distributed Algorithms*. Morgan Kaufmann, 1998.
- [20] Matthias Majuntke, Dan Dobre, Marco Serafini, and Neeraj Suri. Abortable Fork-Linearizable Storage. In *Proceedings of the 13th International Conference on Principles of Distributed Systems, OPODIS '09*, pages 255–269, Berlin, Heidelberg, 2009. Springer-Verlag.
- [21] Dahlia Malkhi and Michael K. Reiter. Byzantine Quorum Systems. *Distributed Computing*, 11(4):203–213, 1998.
- [22] David Mazières and Dennis Shasha. Building Secure File Systems out of Byzantine Storage. In *PODC*, pages 108–117, New York, NY, USA, 2002. ACM.
- [23] Peter Mell and Tim Grance. The NIST Definition of Cloud Computing. Report, National Institute of Standards and Technology (NIST), January 2011. Available online at http://csrc.nist.gov/publications/drafts/800-145/Draft-SP-800-145_cloud-definition.pdf.
- [24] Alexander Shraer, Christian Cachin, Asaf Cidon, Idit Keidar, Yan Michalevsky, and Dani Shaket. Venus: Verification for Untrusted Cloud Storage. In *Proceedings of the 2010 ACM Workshop on Cloud Computing Security, CCSW '10*, pages 19–30, New York, NY, USA, 2010. ACM.
- [25] Peter Williams, Radu Sion, and Dennis Shasha. The Blind Stone Tablet: Outsourcing Durability to Untrusted Parties. In *Proc. NDSS*, 2009.

Appendix

A Definitions

To complete the definition of (*weak*) *fork-linearizability* in Definition 1 on page 4, we first have to introduce the notion of a *possible view* [7] and the *weak real-time order* [7]. Definitions of the *causal precedence* relation between operations are given as Definition 18 in Appendix C for the weak fork-linearizable shared memory.

Definition 3. A sequence of events π is called a *possible view* of a history σ at a client C_i with respect to a type T if σ can be extended (by appending zero or more responses) to a history σ' such that:

1. π is a sequential permutation of some subsequence of $complete(\sigma')$,
2. $\pi|_{C_i} = complete(\sigma')|_{C_i}$, and
3. π satisfies the sequential specification of T .

Where for a sequence of events σ , $complete(\sigma)$ is the maximal subsequence of σ consisting only of complete operations.

Definition 4. Let π be a sequence of events and let $lastops(\pi)$ be a function of π returning the set containing the last operation from every client in π (if it exists), that is,

$$lastops(\pi) := \bigcup_{i=1, \dots, n} \{o \in \pi|_{C_i} \mid \nexists o' \in \pi|_{C_i} \text{ s.t. } o \text{ precedes } o' \text{ in } \pi\}$$

We say that π preserves the *weak real-time order* of a sequence of operations σ whenever π excluding all events belonging to operations in $lastops(\pi)$ preserves the real-time order⁷ of π .

Variables 5. Variables used in Algorithm 1:

Universal Object Implementation, Algorithm of Client i

C INC&READ counter object, initially 0	
R_1, \dots, R_n SWMR atomic register, initially $\langle 0, (0, \dots, 0), \perp, \perp \rangle$	/* ts+version+state+sig */
ts, ts', ts_i, cn integer, initially 0	/* timestamp & counter */
$V[1..n], V_i[1..n], V_{suc}[1..n]$ array of integers, initially $(0, \dots, 0)$	/* version */
s, s_i state, initially \perp	/* state */
res operation result, initially \perp	/* return value */
sig, sig_i signature, initially \perp	/* signature */

Variables 6. Variables used in Algorithm 2:

Weak Fork-Linearizable Memory for n Clients, Algorithm of Client C_i

S , atomic snapshot object with n componenets, initially $((0, \perp), \dots, (0, \perp))$	/* timestamp+sig */
W_1, \dots, W_n , SWMR atomic registers, initially $(\perp, 0, \emptyset, \perp)$	/* val+ts+rs+ws+sig */
v, wv value, initially \perp	/* value written to storage */
$wts, ots, i, k, r, r', w, w', tmp_1, \dots, tmp_n$ integer, initially 0	/* timestamps + temp. variables */
$read_seen[1..n][1..n], write_seen[1..n][1..n],$	/* matrices of seen
$r_write_seen[1..n][1..n]$, matrix of sets of pairs (integer, integer), initially \emptyset	operations */
sig, sig_1, \dots, sig_n signature, initially \perp	/* signatures */

⁷Sequence π preserves the real-time order of history σ if for each operations o, o' in π holds: if $o <_{\sigma} o'$ then $o <_{\pi} o'$.

B Proof of Correctness of Algorithm 1

This section formally proves that Algorithm 1 implements an abortable, fork-linearizable Byzantine emulation of a universal type.

The implementation uses an INC&READ counter object, given as Algorithm 3. An INC&READ counter object C provides two atomic operations INC&READ(C) and READ(C). An invocation to INC&READ(C) advances the counter object C and returns a value which is higher than any value returned before the invocation of INC&READ(C). An invocation to READ(C) returns the current value of the counter object. The INC&READ counter C has two properties:

- P1** If a client process runs in isolation and it first calls INC&READ(C) and then later READ(C), then the same value is returned by both invocations, and
- P2** the values returned by INC&READ(C) reflect the real-time order of invocations to INC&READ(C).

The counter object C is a wait-free variant of the abortable INC&READ counter described by Aguilera *et al.* [1]. For the implementation of the INC&READ counter, instead of abortable base registers [1], wait-free atomic registers are used here, hence the counter does not need to abort.

Algorithm 3: INC&READ Counter for n Clients, Algorithm of Client C_i

Variables:

R_1, \dots, R_n , SWMR atomic registers, initially $(0, \perp)$
 $cnt_1, \dots, cnt_n, k, c, id$, integers, initially 0

```

3.1 INC&READ() do
3.2   for  $k = 1, \dots, n$  do  $cnt_k \leftarrow read(R_k)$ 
3.3    $c \leftarrow \max_{1 \leq k \leq n} \{cnt_k\} + 1$ 
3.4   write( $c, R_i$ )
3.5   return  $n \cdot c + i$ 

3.6 READ() do
3.7   for  $k = 1, \dots, n$  do  $cnt_k \leftarrow read(R_k)$ 
3.8    $c \leftarrow \max_{1 \leq k \leq n} \{cnt_k\}$ 
3.9    $id \leftarrow \max_{1 \leq k \leq n} \{k \mid cnt_k = c\}$ 
3.10  return  $n \cdot c + id$ 

```

We further define the “affected by” relation of two (high-level) operations implemented by our protocol (Definition 7), the notion of relevant operations (Definitions 8 and 9), and the \leq order relation on versions (Definition 10).

Definition 7. For two operations o, o' in history σ of the universal type implemented by Algorithm 1 we say that o *affects* o' in σ (o' is *affected by* o) whenever one of the following conditions hold:

1. Operations o and o' are both invoked by the same client, o is successful and o finishes before o' is invoked.
2. Operation o' reads the state s_l (version V_l) written by o , successfully verifies the signature and executes APPLY_T to s_l (during o' in lines (1.5—1.10), V_l is the version, s_l the state, ts_l the timestamp and sig_l the signature written during o).
3. There exists an operation o'' such that o *affects* o'' and o'' *affects* o' .

The notion of a *relevant* operation is defined recursively.

Definition 8. An operation o is *relevant* if and only if

1. o is *successful* OR
2. there exists a relevant operation o' that has been affected by o .

Definition 9. An operation o is *relevant* for client C_i if and only if some successful operation of C_i has been affected by o .

Definition 10 (Order Relation). A *version* V is a vector of integers of length n , initially $(0, \dots, 0)$. For two versions V and V' holds $V \leq V'$ if and only if

$$\forall i : V[i] \leq V'[i].$$

It holds $V = V'$ if and only if V and V' are the same versions.

For two operations o and o' with versions V and V' holds $o \leq o'$ if and only if

$$V \leq V'$$

It holds $o = o'$ if and only if o and o' are the same operations.

It is easy to see that \leq relation on operations (versions) is *transitive*. The next definition introduces the notion of operations taking *effect*. Note, that the last operation of each client, when the client crashes, may be incomplete but may appear as a complete operation to others — i.e. it took effect.

Definition 11. An operation of client C_i *takes effect* if and only if the low-level *write* operation in line 1.13 successfully returns.

The next Corollary shows that \leq relation on operations respects the real-time order of sequential operations.

Corollary 12. If o and o' are two operations and $o \leq o'$ then o' does not precede o .

Proof. Let o and o' have associated versions V and V' respectively. Assume by contradiction that o' precedes o and that $o \leq o'$. During o , the entry $V[i]$ is incremented. As o' precedes o and as versions are digitally signed (line 1.12), it holds that $V'[i] < V[i]$. Hence, $V' \not\leq V$ and therefore $o \not\leq o'$. \square

The following two Corollaries show that operations which affect each other are ordered by \leq such that the “affected by” relation is respected. According to Definition 7, the successful operations of one client affect each other (Corollary 13) as well as an operation that is applied to the state updated by another operation (Corollary 14).

Corollary 13. All operations of the same client are totally ordered by \leq relation on operations.

Proof. We show that operation o' of client C_i is greater than its previous completed operation o . Let V and V' be the versions of operation o and o' respectively. Let V_l be the version read by operation o' . The entries $V'[k]$, $k \neq i$ is assigned the maximum of $V_l[k]$ and $V[k]$ (line 1.10), and $V'[i]$ is updated with a value larger than $V[i]$, as $V[i]$ is incremented with every invoked operation of C_i (line 1.11). Clearly, $V' > V$. By induction on C_i 's operations, it follows that o' is greater than any preceding operation of C_i . \square

Corollary 14. If o' is reading state s from some register R_i updated by operation o , then $o > o'$.

Proof. Let V and V' be the versions of o and o' respectively. When operation o' is applied to state s , then V is version V_l during operation o' . By lines 1.10–1.11 and analogously to the proof of Corollary 13 it follows directly that $V' > V$. \square

The following Lemma shows the main result of this section: The universal type, implemented in Algorithm 1 satisfies fork-linearizability (Definition 1 on page 4). The proof shows how for each client the subsequences π_i are constructed. Then, by proving two claims, we show that sequences π satisfy the properties of fork-linearizability. To ease the argumentation, operations which are not relevant at all or not relevant for client C_i are ignored.

Lemma 15. The history σ induced by any execution of Algorithm 1 satisfies fork-linearizability with respect to the universal object type T .

Proof. Let σ be the sequence of events observed by the clients in the protocol. We first remove all invocations of incomplete operations that do not take effect (Definition 11). We add the corresponding completion event of incomplete operations that take effect directly at the end of σ . Then we remove all operations that are not relevant. Note, that σ now contains only complete and relevant operations.

We construct a sequential permutation π by totally ordering all events in σ . To achieve this, we order the events in σ by the following rules:

1. Sort the operations in σ by \leq relation on operations.
2. Sort any yet unsorted operations by the real-time order of their completion event.

We construct the subsequences π_i (for $i = 1, \dots, n$) as required by the definition of fork-linearizability (Definition 1). We include in π_i all operations of client C_i in π . Then, for all $o \in \pi_i$ we include into π_i all operations o' in π such that $o' \leq o$. Finally, we remove all operations that are not *relevant* for C_i .

By Corollary 12, as \leq relation on operations respects real-time order, the following claim follows directly:

Claim 15.1 Let o and o' be two operations and o precedes o' in σ . Then, o precedes o' in π .

Claim 15.2 Let o_c be an operation of client C_i in sequence π_m of client C_m , $m \in 1, \dots, n$, that updates state s in register R_j ; state s was written by operation o_a of client C_j , $j \in 1, \dots, n$, into register R_j . Then:

1. Operation o_a is in π_m , and
2. in π_m there is no operation by client C_k , $k \in 1, \dots, n$, that is *relevant* for C_m , that is subsequent to o_a in π_m , and that completes before o_c is invoked.

By Corollary 14 holds that $o_c > o_a$. Hence, o_a is included in π_m by construction and the first statement of Claim 15.2 follows directly.

To prove the second statement of Claim 15.2, let us assume for contradiction that such an operation o_b of client C_k exists in π_m which is invoked after o_a completes and which completes before o_c is invoked. Hence, o_a, o_b, o_c are three sequential operations in that order in σ .

We first show that o_b and o_c do not affect each other, i.e. o_b is not affected by o_c and o_c is not affected by o_b :

- “ o_c is not affected by o_b ”: Operation o_c is affected by o_a and by o_{suc} , the last successful operation by client C_i previous to o_c , if it exists. If o_{suc} is affected by o_b then o_a precedes o_{suc} . Hence, the i th entry in the version of o_{suc} is greater then the one of o_a and therefore o_a could not affect o_c (as the check in line 1.7 is not passed) — a contradiction. If o_{suc} is not affected by o_b , then o_c is also not affected by o_b (as “affected by” relation is transitive; Definition 7).
- “ o_b is not affected by o_c ”: Follows directly as o_b precedes o_c .

Next, we derive a contradiction to the assumption that operation o_b exists. As operations o_b and o_c are both relevant for client C_m and o_c and o_b do not affect each other, there are successful operations o'_b and o'_c of client C_m such that o'_b is affected by o_b and o'_c is affected by o_c . Let o'_b and o'_c be the operations of C_m that are affected by o_b or o_c respectively with the smallest versions (they exists by Corollary 13). Note that, as o'_b and

o'_c are both successful operations of the same client C_m , they affect each other. Let us assume w.l.o.g. that o'_c is affected by o'_b . This means, there exists some operation o_{join} of client C_{join} , $join \in 1, \dots, n$, which is the operation with the smallest timestamp that is affected by both o_b and o_c . For operation o_{join} either holds

- (A) $o_{join} \leq o'_b$ and $o'_b (= o'_c)$ is affected by o_{join} , or
- (B) o_{join} is affected by o'_b and $o'_c (\neq o'_b)$ is affected by o_{join} .

To have o_{join} to be affected by two operations that do not affect each other, by Definition 7, (1) there must be some operation o_{join_suc} which is the last successful operation of client C_{join} previous to o_{join} that is affected either by o_b or o_c . W.l.o.g. we assume that o_{join_suc} is affected by o_b . Note, that $o_{join_suc} \geq o_b$. Further, o_{join_suc} is not affected by o_c , as otherwise o_{join} would not be the first operation affected by both o_b and o_c ($o_{join_suc} < o_{join}$ by Corollary 13). Further, (2) o_{join} reads the state written by some operation o_r that is affected by o_c and $o_r \geq o_c$. Analogously, o_r is not affected by o_b ($o_r < o_{join}$ by Corollary 14). Hence, as o_b and o_c do not affect each other, there are disjunct “affected by” paths⁸ from o_b to o_{join_suc} and from o_c to o_r (Figure 4).

During operation o_{join_suc} the *jointh* entry of its version is raised to $V_{join_suc}[join]$ (line and 1.11) and as o_{join_suc} is successful, also $V_{suc}[join] \geq V_{join_suc}[join]$ (line 1.18) from this point on at client C_{join} . Consequently, during o_{join} , as o_{join_suc} precedes o_{join} , client C_{join} does not accept version V_l such that $V_l[join] < V_{join_suc}[join]$ (check in line 1.7 would not be passed). Hence, there must be an operation o'_{join} on the path from o_c to o_r that raises the *jointh* entry in the versions to $V_{join_suc}[join]$ or higher. Note, as the *jointh* entry is only raised by an operation of C_{join} (line 1.11), o'_{join} has to be an operation of C_{join} as well. Operation o'_{join} cannot precede o_{join_suc} , as o_{join_suc} is the first operation to raise the *jointh* entry to $V_{join_suc}[join]$ (line 1.11). If o'_{join} follows o_{join_suc} , then o'_{join} does not accept versions V_l with $V_l[join] < V_{join_suc}[join]$ (check in line 1.7), as $V_{suc}[join] \geq V_{join_suc}[join]$ at C_{join} after o_{join_suc} has finished. Hence, o'_{join} does not exist on the path from o_c to o_r and o_{join} would block when reading the version of o_r . Thus, we have a contradiction and o_{join} does not exist.

This means that either o'_b does not exist (case A) which implies that o_b is not relevant for C_m or o'_c does not exist (case B) which implies that o_c is not relevant for C_m . Consequently, the assumption that there is an operation o_b between o_a and o_c in π_m is wrong and thus, all operations in π_m are totally ordered by \leq .

We now show that π_m for all $m = 1, \dots, n$ satisfy fork-linearizability as given in Definition 1. To show that π_m is a possible view of client C_m , properties 1. and 2. of Definition 3 follow directly from the construction of π_m given at the beginning of the Lemma. Claim 15.2 proves property 3. of Definition 3. Hence, π_m is a possible view of client C_m . Each sequence π_m satisfies real-time ordering as shown in claim 15.1. The no-join property (condition 3. in Definition 1) is also an easy consequence of the construction of π_m . The non-relevant operations that have been removed at the beginning of this proof, can be added to all π_m in real-time order of their completion event. As they are not relevant, they do not effect the sequential specification and thus, they do not violate fork-linearizability. \square

The next two Lemmas show that Algorithm 1 implements an abortable Byzantine emulation with fork-linearizability of a universal type (see Definition 2 on page 5). Lemma 16 shows that no operation blocks, and Lemma 17 proves that no operation is trivially aborted.

Lemma 16. If registers R_1, \dots, R_n and INC&READ object C are correct, and σ is the history induced by any execution of Algorithm 1, then no operation in σ halts in line 1.6 nor in line 1.7 of Algorithm 1.

Proof. We show that no operation in Algorithm 1 blocks: If the base objects are correct, and as clients are trusted, no signature is forged and thus no operation blocks in line 1.6.

It remains to show that no operation of client C_i blocks in line 1.7. Assume by contradiction that during operation o of client $C_i \exists k : V_{suc}[k] > V_l[k]$. Let o_{suc} be the successful operation of C_i that wrote some

⁸An “affected by” path from operation o_1 to o_x is a sequence of operations o_1, o_2, \dots, o_x such that for $i = 1, \dots, x - 1$, o_i affects o_{i+1} .

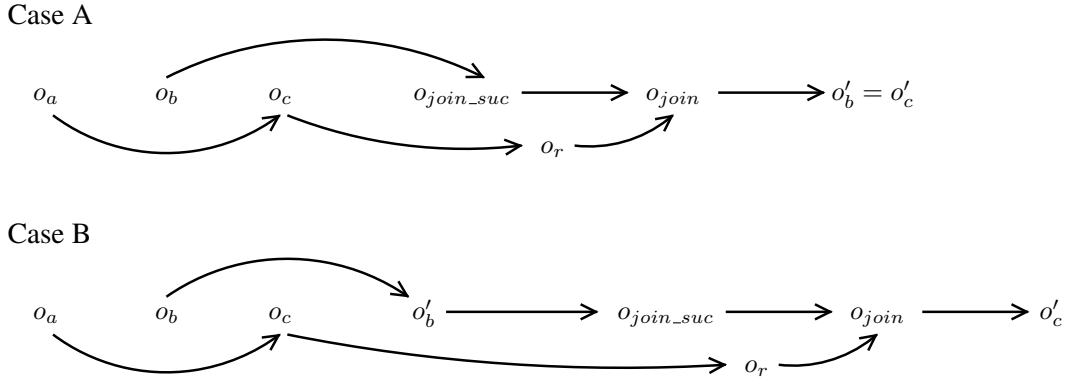


Figure 4: Proof of Lemma 15. The arrows denote the “affected by” relation between operations.

state with version V_{suc} to register R_i . As o_{suc} is not aborted, client C_i has read the same timestamp from the INC&READ object C in line 1.2 and 1.14 during o_{suc} . This means, as the INC&READ object C is correct, that no other operation executed INC&READ(C) between INC&READ(C) and READ(C) of o_{suc} and that o_{suc} has written the highest timestamp so far (line 1.13). Hence, INC&READ(C) of operation o_l , that wrote version V_l , happened either (1) before INC&READ(C) of o_{suc} or (2) after READ(C) of o_{suc} . In case (1) o would not find l as the highest index in line 1.5 and thus it would not read o_l as o_{suc} holds a higher timestamp than o_l (INC&READ object is correct) — a contradiction. For case (2), o_l would read a version $\geq V_{suc}$ and thus $V_l \geq V_{suc}$ (line 1.10)— a contradiction. Concluding, during operation o of client C_i no such $k : V_{suc}[k] > V_l[k]$ exists and thus, the protocol does not block in line 1.7. \square

Lemma 17. If registers R_1, \dots, R_n and INC&READ object C are correct then, if operation o in an execution of Algorithm 1 returns ABORT, then o is concurrent with some other operation.

Proof. The correctness follows directly from the properties of INC&READ object C : Operation o is only aborted if the condition in line 1.15 is satisfied. This is the case when object C returns a different value to call in line 1.14 than in line 1.2. The properties of C imply, that this happens only if some other operations calls INC&READ(C) in the meanwhile. This means, some other operation is concurrent with o (according to the definition in Section 2 on page 3) and thus, we are done. \square

Note, that Lemma 17 is sufficient to show that Algorithm 1 implements an *abortable* object. It is easy to see that in every situation where an operation of Algorithm 1 aborts, Aguilera’s universal type construction ([1], Algorithm 2) would abort as well.

Finally, the correctness of Algorithm 1 has been shown in Lemma 15 (Fork-Linearizability), Lemma 16 (No Blocking), and Lemma 17 (Nontriviality).

C Proof of Correctnes of Algorithm 2

Now we prove that Algorithm 2 implements a wait-free, weakly fork-linearizable Byzantine emulation of a shared memory. To complete the definition of *weak fork-linearizability* in Definition 1 on page 4, we first add the formal definition of the *causal precedence* relation between READ and WRITE operation of Algorithm 2.

Definition 18. For two operations o, o' in history σ of the shared memory implemented by Algorithm 2 we say that o *causally precedes* o' in σ (o' *causally depends on* o), denoted $o \rightarrow_{\sigma} o'$ whenever one of the following conditions hold:

1. Operations o and o' are both invoked by the same client and o finishes before o' is invoked.
2. operation o' is a READ operation, o is a WRITE operation, and o' reads the value written by o .
3. There exists an operation o'' such that $o \rightarrow_{\sigma} o''$ and $o'' \rightarrow_{\sigma} o'$.

Now, we can proceed with the proof that Algorithm 2 satisfies the properties of weak fork-linearizability as given in Definition 1.

To show the existence of sequential permutations π_i for every client C_i that satisfy weak fork-linearizability, let $o_{i1}, \dots, o_{p_{il_i}}$ be the operations of C_i ordered by their timestamps, $i \in 1, \dots, n, l_i \in \mathbb{N}$. The timestamp of an operation is variable ots that is assigned to a READ operation in line 2.4 and to a WRITE operation in line 2.15, respectively. For every client C_i we define a directed graph G_{π_i} , where the set of operations $o_{i1}, \dots, o_{p_{il_i}}$ are the vertices. For all $k \in 1, \dots, l_i - 1$, we draw an edge from o_{ik} to $o_{i(k+1)}$.

Next, we construct a directed graph G_{π} as $G_{\pi_1} \cup G_{\pi_2} \cup \dots \cup G_{\pi_n}$. We add an edge from o_{iw_i} to o_{jr_j} to graph G_{π} if o_{iw_i} is a WRITE operation of $C_i, i \in 1, \dots, n$, and o_{jr_j} is a READ operation of $C_j, j \in 1, \dots, n$, that reads the value written by o_{iw_i} .

The purpose of the next Corollary is to show that a partial order of the operations can be defined according to the ordering of the vertices of graph G_{π}

Corollary 19. Graph G_{π} does not contain directed cycles.

Proof. Let us assume there exists the following directed cycle which is also the shortest possible one: $(o_{ir}, o_{iw}, o_{jr}, o_{jw})$, where o_{ir} is a READ and o_{iw} is a WRITE operation of client C_i , and o_{jr} is a READ and o_{jw} is a WRITE operation of client C_j . Further, let o_{ir} have a lower timestamp than o_{iw} , let o_{jr} read the value written by o_{iw} , let o_{jw} have a lower timestamp than o_{jr} , and let o_{ir} read the value written by o_{jw} . We now can deduce the following statements:

1. o_{ir} precedes o_{iw} , as both are operations of C_i and o_{ir} has a lower timestamp than o_{iw} (line 2.13).
2. o_{jr} returns after o_{iw} has been invoked, as otherwise o_{jr} cannot read the value written by o_{iw} (written values are digitally signed).
3. o_{jr} precedes o_{jw} , as both are operations of C_j and o_{jr} has a lower timestamp than o_{jw} (line 2.13).
4. o_{jw} is invoked after o_{iw} has been invoked, by 2. and 3.
5. o_{ir} returns after o_{jw} has been invoked, as otherwise o_{ir} cannot read the value written by o_{jw} (written values are digitally signed).

By statements 1.–5. we have the contradiction that o_{ir} returns before and after o_{iw} is invoked. The analogous arguments hold, if the circle is extended between o_{jw} and o_{ir} to a circle $(o_{ir}, o_{iw}, o_{jr}, o_{jw}, \dots)$ of arbitrary length.

Hence, graph G_{π} does not contain directed cycles. □

For each client C_i we recursively define the subgraph $T(o_{il_i})$ that contains o_{il_i} as a vertex, and if o is a vertex of $T(o_{il_i})$, and (o', o) is an edge of G_π , then vertex o' and edge (o', o) is added to $T(o_{il_i})$ until no more edges can be added.

Corollary 20. The set of operations represented by the vertices of $T(o_{il_i})$ contains all operations of client C_i .

Proof. By construction of graph G_{π_i} , there is a path from any operation of client C_i to o_{il_i} . Thus, all operation of C_i are contained in $T(o_{il_i})$. \square

Now, we start constructing for each client C_i a subsequence π_i of the history σ induced by any execution of Algorithm 2 that satisfies the properties of weak fork-linearizability (Definition 1). The next corollary constructs an order relation among operations in π_i .

Corollary 21. There is a sequential permutation π_i of the set of operations represented by the vertices of $T(o_{il_i})$ and an order relation $<_{\pi_i}$ that satisfies the following condition: For every operation $o \in \pi_i$ and every WRITE operation $o' \in V(G_\pi)$ s.t. o' causally precedes o , it holds that $o' \in \pi_i$ and that $o' <_{\pi_i} o$.

Proof. By Corollary 20, we know that every operation of C_i is contained in π_i . Further, by the construction of graphs $G(\pi)$ and $T(o_{il_i})$ every operation that causally precedes an operation in π_i (Definition 18) is contained in π_i . As $T(o_{il_i})$ contains no cycles (Corollary 19), for $o, o' \in \pi_i$ we order o before o' ($o <_{\pi_i} o'$) if there is an edge from o to o' in $T(o_{il_i})$. \square

The order relation constructed in the proof of Corollary 21 does not necessarily respect the sequential specification of a shared memory. The next corollary shows how this can be achieved.

Corollary 22. The order relation $<_{\pi_i}$, constructed in Corollary 21 can be extended such that π_i satisfies the sequential specification of shared memory: If o_{kr} is a READ operation of some client C_k , $k \in 1, \dots, n$, in π_i that reads the value written by WRITE operation o_{lw} from client C_l , $l \in 1, \dots, n$, then we additionally order o_{kr} before $o_{l(w+1)}$ where $o_{l(w+1)}$ is the next WRITE operation of C_l in π_i (if it exists in π_i).

Proof. By causality o_{lw} is ordered before o_{kr} . There is no WRITE operation of C_l between o_{lw} and o_{kr} in π_i , as $o_{l(w+1)}$ is the next WRITE operation of C_l in π_i after o_{lw} , and o_{kr} can be ordered before it. \square

Now we define how the remaining operations have to be ordered such that π_i satisfies weak real-time ordering. The proof of the following lemma distinguishes two cases to show that when a READ reads a value written by some WRITE, then the WRITE is the last one that precedes the READ. The two cases correspond to the fact that a READ operation of client C_i may appear in its own sequence π_i (case B) as well as in sequence π_j of C_j .

Corollary 23. The order relation $<_{\pi_i}$, constructed in Corollaries 21 and 22 does not violate weak real-time order.

Proof. We distinguish the following cases:

Case A: Let w and w' be two WRITE operations of client C_i and let r be a READ operation of client C_j that reads the value written by w . Let w precede w' in π_i . Let r be not the last operation of C_j in π_i . Then w' does not happen before r , i.e. $r.UPDATE$ happens before $w'.SCAN$ ⁹.

Proof. We assume by contradiction that WRITE operation w' happens before READ operation r , i.e. $w'.SCAN$ precedes $r.UPDATE$ (Ass. A). By assumption, r reads the value written by operation w and thus by line 2.8 set $read_seen_j[j][i]$ at client C_j contains the couple (r, w) ¹⁰. Let $w_{\min} \leq w'$ be the WRITE

⁹In the following $x.UPDATE$ ($x.SCAN$) denotes a call of procedure UPDATE (SCAN) during READ (WRITE) operation x in line 2.4 (line 2.16). The analogous notation holds for $x.write$ ($x.read$) in line 2.5 (line 2.15) during Lemma 28.

¹⁰To simplify the presentation, let r denote the timestamp assigned to operation r and w the timestamp assigned to w .

operations of client C_i directly following w . Then, during operation w_{\min} , by line 2.17, the couple (x_{\min}, w_{\min}) is added to set $write_seen_i[i][j]$ at client C_i . By Ass. A and as all written timestamps are digitally signed, it holds that $x_{\min} < r$. As r is not the last operation of C_j in π_i , there exists a read operation r'' of C_i that happens after w' , a write operation w'' of C_j that happens after r , and there exists a path in graph $T(op_{i,i})$ from r to r'' that contains w'' (see Figure 5). During operation w'' of client C_j , variable $read_seen_j$ is written by line 2.15. As there is the causal path from w'' to r'' , by lines 2.7 and 2.15, during operation r'' of client C_i , $r_read_seen[j][i]$ contains couple (r, w) . By line, 2.7 it is also contained in $read_seen_i[j][i]$ during operation r'' . We further know that $write_seen_i[i][j]$ at client C_i contains (x_{\min}, w_{\min}) . As operation w_{\min} is the minimal operation larger then w , the check in line 2.25 is not passed as $x_{\min} < r$ and operation r'' blocks. This means that r is not in π_i — a contradiction. \square

Case B: Let w and w' be two WRITE operations of client C_j and let r be a READ operation of client C_i that reads the value written by w . Let w precede w' in π_i . Let w' be not the last operation of C_j in π_i . Then w' does not happen before r , i.e. $r.UPDATE$ happens before $w'.SCAN$.

Proof. We assume by contradiction that WRITE operation w' happens before READ operation r , i.e. $w'.SCAN$ precedes $r.UPDATE$ (Ass. B). By assumption, r reads the value written by operation w and thus by line 2.8 set $read_seen[i][j]$ at client C_i contains the couple (r, w) . Let $w_{\min} \leq w'$ be the WRITE operations of client C_j directly following w . Then, during operation w_{\min} , by line 2.17, the couple (x_{\min}, w_{\min}) is added to set $write_seen_j[j][i]$ at client C_j . By Ass. B and as all written timestamps are digitally signed, it holds that $x_{\min} < r$. As w' is not the last operation of C_j in π_i , there exists a read operation r'' of C_i that happens after r , a write operation w'' of C_j that happens after w' , and there exists a path in graph $T(op_{i,i})$ from w to r'' that contains w' and w'' (see Figure 5). During operation w'' of client C_j , variable $write_seen_j$ is written by line 2.15. As there is the causal path from w'' to r'' , by lines 2.10 and 2.15, during operation r'' of client C_i , $r_write_seen[j][i]$ contains couple (x_{\min}, w_{\min}) . We further know that $read_seen_i[i][j]$ at client C_i contains (r, w) . As operation w_{\min} is the minimal operation larger then w , the check in line 2.28 is not passed as $x_{\min} < r$ and operation r'' blocks. This means that w' is not in π_i — a contradiction. \square

It remains to show that in π_i READ operations of client C_j that read values written by operations of client C_k can be ordered to satisfy weak real-time order. The proof is obvious as weak real-time order holds for π_j when case B from above is applied.

Hence, the order induced in Corollary 22 does not violate weak real-time order — i.e. the not yet ordered operations can be ordered in real-time order or in any deterministic order if they are concurrent. \square

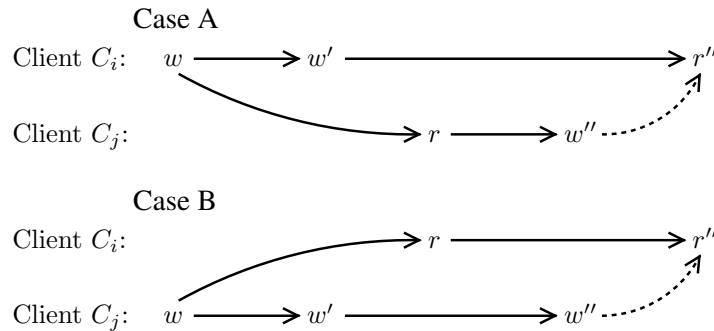


Figure 5: Proof of Corollary 23. Arrows denote the causality relation.

Corollary 24. If all operations in π_i which have not yet been ordered in Corollary 21 or 22 are ordered according to their real-time order if they are sequential and by the real-time order of their completion event else, then order relation $<_{\pi_i}$ of π_i satisfies weak real-time ordering.

Proof. The proof follows directly from construction and Corollary 23. \square

To complete the correctness proof, we have to add operations to each π_i such that the join-at-most-once property is satisfied. This is because π_i may contain operations of π_j but none of π_k , but π_j might have common operations with π_k . To ensure that π_i and π_j share a common prefix such operations have to be added to π_i . Thus, we define a merge operation on totally ordered command sequences π_i .

Definition 25. Let π_i and π_j be two totally ordered command sequences such that there are at least two operation o, o' for which holds $o \in \pi_i \cap \pi_j$ and $o' \in \pi_i \cap \pi_j$ ¹¹. Let $\pi|_x$ denote the prefix of an operation sequence π that ends with operation x . To merge π_i and π_j we perform the following steps: Let $o_{2\text{ndlast}}$ be the second last operation in $\pi_i \cap \pi_j$. In π_i and π_j we replace the prefix $\pi_i|_{o_{2\text{ndlast}}}$ and $\pi_j|_{o_{2\text{ndlast}}}$ by $\pi_{\text{merge}_{ij}}$:

- $\pi_{\text{merge}_{ij}}$ contains all operations from $\pi_i|_{o_{2\text{ndlast}}} \cup \pi_j|_{o_{2\text{ndlast}}}$
- If for two operations o, o' in $\pi_{\text{merge}_{ij}}$ holds $o <_{\pi_i} o'$ or $o <_{\pi_j} o'$, then we order o before o' in $\pi_{\text{merge}_{ij}}$, i.e. $o <_{\pi_{\text{merge}_{ij}}} o'$.
- If for two operations o, o' in $\pi_{\text{merge}_{ij}}$ neither $o <_{\pi_{\text{merge}_{ij}}} o'$ nor $o >_{\pi_{\text{merge}_{ij}}} o'$ holds, then we order o and o' in $\pi_{\text{merge}_{ij}}$ according to their real-time ordering or by the real-time order of their completion event if they are concurrent.

Corollary 26. For all pairs, $i, j \in 1, \dots, n$, if we merge π_i and π_j whenever they have two or more operations in common until no more changes appear. Then sequences π_i, \dots, π_n satisfy the *At-most-one-join* property (Definition 1).

Proof. Correctness follows directly from the construction given in Definition 25. \square

Lemma 27. The history σ induced by any execution of Algorithm 2 satisfies weak fork-linearizability with respect to a shared memory object with n registers.

Proof. The correctness follows from Corollaries 20 and 22 which ensures that π_i is a view of σ with respect to the functionality of a shared memory object, from Corollary 21 that guarantees that causality is respected, from Corollary 24 that ensures weak real-time ordering, and Corollary 26 that guarantees the *At-most-one-join* property. \square

Lemma 28. If all base registers and the snapshot object S is correct, and σ is the history induced by any execution of Algorithm 2, then no READ operation blocks in line 2.6, 2.25, nor 2.28 and no WRITE operation blocks in line 2.18.

Proof. We show that no operation of the shared memory implemented in Algorithm 2 blocks when the base registers behave correctly. As the clients behave correctly and registers do not forge signatures, it is easy to see that WRITE operations do not block. The same argument holds for the check in line 2.6 during READ operations. Thus, it remains to show that READ operations do not block in line 2.25 and 2.28.

Let us assume for contradiction that there is a READ operation of client C_i that blocks in line 2.25 or 2.28:

Line 2.25 There exists a READ operation r of client C_k that has read from WRITE operation w of client C_i (line 2.8). By assumption, the minimal WRITE operation of C_i after w , called w' has seen READ operation r' of client C_k (line 2.17). As $r' < r$, and as all registers are correct, $r.\text{UPDATE}$ does not precede $w'.\text{SCAN}$. Thus, $w'.\text{write}$ precedes $r.\text{read}$. However, as w precedes w' , we conclude that r reads the value written by w' — a contradiction.

¹¹By construction, for all such operations hold $o <_{\pi_i} o'$ and $o <_{\pi_j} o'$ or $o >_{\pi_i} o'$ and $o >_{\pi_j} o'$.

Line 2.28 There exists a READ operation r of client C_i that has read from WRITE operation w of client C_k (line 2.8). By assumption, the minimal WRITE operation of C_k after w , called w' has seen READ operation r' of client C_i (line 2.17). As $r' < r$, and as all registers are correct, r .UPDATE does not precede w' .SCAN. Thus, w' .write precedes r .read. However, as w precedes w' , we conclude that r reads the value written by w' — a contradiction.

Hence, no operation in σ blocks. □

Finally, it has been shown in Lemma 27 (Weak Fork-Linearizability), and Lemma 28 (No Blocking), that Algorithm 2 correctly implements a wait-free, weak fork-linearizable Byzantine emulation of a shared memory.