

# A Parametrized Algorithm that Implements Sequential, Causal, and Cache Memory Consistencies

Ernesto Jiménez<sup>b</sup>, Antonio Fernández<sup>a</sup>, Vicent Cholvi<sup>c,\*</sup>

<sup>a</sup>*Universidad Rey Juan Carlos, 28933 Móstoles, Spain*

<sup>b</sup>*Universidad Politécnica de Madrid, 28031 Madrid, Spain*

<sup>c</sup>*Universitat Jaume I, 12071 Castellón, Spain*

---

## Abstract

In this paper we present an algorithm that can be used to implement sequential, causal, or cache consistency in distributed shared memory (DSM) systems. For this purpose it has a parameter that allows to choose the consistency model to be implemented. We can also use our algorithm such that not all processes have the same value in this parameter (we have shown the resulting consistency). This characteristic allows to choose a concrete consistency model but implementing it with the algorithm more efficient in each case (in function of the requirements of the applications). As far as we know, this is the first algorithm proposed that implements cache coherence.

In our algorithm, when implementing causal and cache consistency all read and write operations are executed locally (i.e., are *fast*). It is known that no sequential algorithm has only fast memory operations. However, in our algorithm, when implementing sequential consistency all write operations and some read operations are fast. The algorithm uses propagation and full replication, where values written by a process are propagated to the rest of processes. It works in a cyclic turn fashion, with each process of the DSM system broadcasting one message in its turn. The values written by the process are sent in the message (instead of sending one message for each write operation), but unnecessary values are excluded. All this allows to control the amount of message traffic due to the algorithm.

*Key words:* Distributed shared memory, causal consistency, sequential consistency, distributed systems, coherency models.

---

\* Corresponding author. Address: Departamento de Lenguajes y Sistemas Informáticos, Universitat Jaume I, Campus de Riu Sec, 12071 Castellón (Spain). Email: vcholvi@uji.es

## 1 Introduction

Distributed shared memory (DSM) is a well-known mechanism for inter-process communication in a distributed environment [1,2]. Roughly speaking, it consists of using read and write operations for inter-process communication thus hiding from the programmer the particular communication technique employed so that they do not need to be involved in the management of messages.

In general, in order to increase concurrency, most DSM protocols support *replication* of data. With replication, there are copies (replicas) of the same variables in the local memories of several processes of the system, which allows these processes to use the variables simultaneously. However, in order to guarantee the *consistency* of the shared memory, the system must control the replicas when the variables are updated. That control can be done by either *invalidating* outdated replicas or by *propagating* the new variable values to update the replicas. When propagation is used, a replica of the whole shared memory is usually kept in each process.

An interesting property of any algorithm implementing a consistency model is how long can a memory operation take. If a memory operation does not need to wait for any communication to finish, and can be completed based only on the local state of the process that issued it, it is said that the operation is *fast*, which is a very desirable feature. An algorithm is fast if all its operations are fast. For instance, one of the most widely known memory models, the causal consistency model [3], has been implemented in a fast way a number of times [3–5]. On the contrary, another of the widely known memory models, the sequential [6], has been shown that cannot be implemented in a fast manner [7]. This impossibility result restricts the efficiency of any implementation of sequential consistency[8,9].

### *Our Work*

In this paper, we introduce a parametrized algorithm that implements sequential, causal, and cache consistency [10] and allows us to change the model it implements on-line. The main reasons to choose these three models of consistency are the following. It has been shown that many practical distributed applications require competing operations (i.e., operations that need synchronization among them) ([11]). We have chosen to implement the sequential consistency model because it is the most popular proposed model that provides competing operations (beside atomic consistency model ([12]), which is more restrictive). However, it has also been shown that there are several classes of applications which when executed with algorithms that implement

causal consistency behave as sequentially consistent ([3,13]). Hence, we have also chosen to implement the causal consistency model with an algorithm where all memory operations are fast, and in this manner to avoid the efficiency problems of sequential consistency algorithms ([7]). The cache model is also included, even though it is not so popular, because it is extremely simply integrated in our algorithm and it has an interest (at least theoretical) for applications that require competing operations but only on the same variable. Furthermore, as far as we know, this is the first algorithm proposed to implement cache consistency.

In order to increase concurrency, most DSM algorithms support *replication* of data. With replication, there are copies (replicas) of the same variables in the local memories of several processes of the system, which allows these processes to use the variables simultaneously. However, in order to guarantee the consistency of the shared memory, the system must control the replicas when the variables are updated. That control can be done by either *invalidating* outdated replicas or by *propagating* the new variable values to update the replicas. When propagation is used, a replica of the whole shared memory is usually kept in each process. Our algorithm uses full propagation and uses broadcasts to perform such a task. It works as follows: a write operations is propagated from the process that issues it to the rest of processes so that they can apply it locally. However, write operations are not propagated immediately. The algorithm works on a cyclic turn fashion, with each process broadcasting one message in its turn. This scheme allows a very simple control of the load of messages in the network, since only one message is sent by each process at its turn. That makes several write operations to be grouped in a single propagation message, thus reducing the network load.

When implementing causal and cache consistency, all the operations in our algorithm are fast. Obviously, this is not the case for the sequential model (remember that, from the results in [7], it is derived the impossibility of having all the memory operations fast). However, even in the case of the sequential model, all write operations are always fast. In turn, this does not happen for read operations, but there is only one situation where read operations must be non-fast: when the process that issues that read operation has not issued, since its last turn, any write operations on the variable being read and has issued some write operation on another variable. In this case, the process must wait until reaching its turn.

### *Comparison with Previous Work*

From the set of algorithms that implement DSM, two of them have features similar to those presented in this paper. The first one, which has been pro-

posed by Afek et al. [14] (for sequential memory), they also ensure that write operations will be fast. Also, read operations are fast except for some situations. But those situations are more general than that in our algorithm, which makes our algorithm more "fast". Furthermore, we do not send each variable update in a single message as it is done in [14] and we can also bound the number of messages sent. Finally, in [14] it is assumed that there is a communication medium among all processes (and with the shared memory) that guarantees total order among concurrent write operations. In our case, we do not have such a restriction and enforce the order of the operations by using a cyclic turn technique.

On the other hand, in [13], the authors propose an algorithm that implements three consistency models (sequential, causal and a hybrid between both models). Such an algorithm can dynamically switch among these three consistency models. However, there are a number of differences with the algorithm we propose. First, their algorithm is designed separating the propagation mechanism from the consistency policy. On the contrary, in our algorithm the propagation mechanism is enough to maintain the consistency model. Furthermore, in their implementation they use an adaptation of vector clocks [15] (called *version vectors*) which results in a big waste of memory in each node, and bigger messages to send through the network. Finally, their implementation also forces several restrictions to achieve a total order: a) two update transactions can not be executed concurrently, and b) no update transaction is allowed whenever query transactions are ongoing.

The rest of the paper is organized as follows. In Section 2 we introduce basic definitions. In Section 3 we introduce the algorithm we propose. In Sections 4, 5, and 6 we prove the correctness of our algorithm. In Section 7 we provide an analysis of the complexity of our algorithm. In Section 8 we show the consistency when not all processes are executing our algorithm with the same parameter, and in Section 9 we present our concluding remarks.

## 2 Definitions

In this paper we assume a distributed system that consists of a set of  $n$  processes (each uniquely identified by a value in the range  $0 \dots n-1$ ). Processes do not fail and are connected by a reliable message passing subsystem. These processes use their local memory and the message passing system to implement a shared memory abstraction. This abstraction is accessed through read and write operations on variables of the memory. The execution of these memory operations must be consistent with the particular *memory consistency model*.

Each memory operation acts on a named variable and has an associated value.

A write operation by process  $p$ , denoted  $w_p(x)v$ , stores the value  $v$  in the variable  $x$ . Similarly, a read operation, denoted  $r_p(x)v$ , reports to the process  $p$  that issued it the value  $v$  stored in the variable  $x$  by write operation. When it doesn't matter the process that performs the operation, we simply will denote them as  $w(x)v$  and  $r(x)v$ . To simplify the analysis, we assume that a given value is written at most once in any given variable and that the initial values of the variables are set by using write operations.

In this paper we present an algorithm that uses replication and propagation. We assume each process holds a copy of the whole set of variables in the shared memory. When needed, we use  $x_p$  to denote the local copy of variable  $x$  in process  $p$ . Different copies of the same variable can hold different values at the same time.

We call  $\alpha$  to the set of read and write operations obtained in an execution of the memory algorithm.

Now we define an order among the operations observed by the processes in an execution of the memory algorithm.

**Definition 1 (Execution Order)** *Let  $op$  and  $op' \in \alpha$ . Then  $op$  precedes  $op'$  in the execution order, denoted  $op \prec op'$ , if:*

1.  $op$  and  $op'$  are operations from the same process and  $op$  is issued before  $op'$ ,
2.  $op = w(x)v$  and  $op' = r(x)v$ , or
3.  $\exists op'' \in \alpha : op \prec op'' \prec op'$

From this last definition, we also derive the non-transitive execution order (denoted as  $\prec_{nt}$ ) as a restriction of the execution order if the transitive closure (i.e., the third condition) is not applied. Note that if  $op \prec_{nt} op'$ , then  $op$  has been executed before  $op'$ . Hence, if  $op \prec op'$ , then  $op$  has also been executed before  $op'$ . If  $op \prec op'$ , we define by *related sequence* between  $op$  and  $op'$  a sequence of operations  $op^1, op^2, \dots, op^m$  such that  $op^1 = op$ ,  $op^m = op'$ , and  $op^j \prec_{nt} op^{j+1}$  for  $1 \leq j < m$ .

We say that  $\alpha_p$  is the set of operations obtained by removing from  $\alpha$  all read operations issued by processes other than  $p$ . We also say that  $\alpha(x)$  is the set of operations obtained by removing from  $\alpha$  all the operations on variables other than  $x$ .

**Definition 2 (View)** *We denote by system view  $\beta$ , process view  $\beta_p$  or variable view  $\beta(x)$  a sequence formed with all operations of  $\alpha$ ,  $\alpha_p$  or  $\alpha(x)$ , respectively, such that this sequence preserves the execution order  $\prec$ .*

Note that, due to the existence of operations that are not affected by the execution order, there can be a lot of sequences on  $\alpha$ ,  $\alpha_p$  or  $\alpha(x)$ , not only  $\beta$ ,

$\beta_p$  or  $\beta(x)$ , that preserve  $\prec$ .

We use  $op \rightarrow op'$  to denote that  $op$  precedes  $op'$  in a sequence of operations. Abusing the notation, we will also use  $set_1 \rightarrow set_2$ , where  $set_1$  and  $set_2$  are set of operations, to denote that all the operations in  $set_1$  precede all the operations in  $set_2$ .

**Definition 3 (Legal View)** *A view  $\beta$  on  $\alpha$  is legal if  $\forall r(x)v \in \alpha, \nexists w(x)u \in \alpha : w(x)v \rightarrow w(x)u \rightarrow r(x)v$  in  $\beta$ .*

**Definition 4 (Sequential, Causal or Cache Algorithm)**

- *An algorithm implements sequential consistency if for each execution  $\alpha$  there exists a legal view of it.*
- *An algorithm implements causal consistency if for each execution  $\alpha$  there exists a legal view  $\beta_p$  of each  $\alpha_p, \forall p$*
- *An algorithm implements cache consistency if for each execution  $\alpha$  there exists a legal view  $\beta(x)$  of each  $\alpha(x), \forall x$ .*

### 3 The Algorithm

In this section we present the parametrized algorithm  $A$  that implements causal, cache and sequential consistency. Figure 1 presents the algorithm in detail. As it can be noted, it is run with a parameter `model`, which defines the consistency model that the algorithm must implement. Hence, the parameter must take one of the values `causal`, `sequential`, or `cache`.

In Figure 1 it can be seen that all write operations are fast. When a process  $p$  issues a write operation  $w_p(x)v$ , the algorithm changes the local copy of variable  $x$  (which we denote by  $x_p$ ) to the value  $v$ , includes the pair  $(x, v)$  in a local set of variable updates (which we call `updatesp`), and returns control. This set `updatesp` will later be asynchronously propagated to the rest of processes. Note that, if a pair with the variable  $x$  was already in `updatesp`, it is removed before inserting the new pair, since it does not need to be propagated anymore.

Processes propagate their respective sets `updatesp` in a cyclic turn fashion, following the order of their identifiers. To maintain the turn, each process  $p$  uses a variable `turnp` which contains the identifier of the process whose set must be propagated next (from  $p$ 's view). When `turnp = p`, process  $p$  itself uses the communication channels among processes to send to the rest of processes its local set of updates `updatesp`. This is done in the algorithm

```

Initialization ::
begin
  turnp ← 0
  updatesp ← ∅
end

wp(x)v :: atomic function
begin
  xp ← v
  if ((x, ·) ∈ updatesp) then
    remove (x, ·) from updatesp
  include (x, v) in updatesp
end

rp(x) :: atomic function
begin
  if (model = sequential) and (updatesp ≠ ∅) and ((x, ·) ∉ updatesp)
  then
    wait until turnp = p
  return(xp)
end

send_updates() :: atomic task activated whenever turnp = p
begin
  /* send to all processes, except itself */
  broadcast(updatesp)
  updatesp ← ∅
  turnp ← (turnp + 1) mod n
end

apply_updates() :: atomic task activated whenever turnp = q, p ≠ q, and
the set updatesq from process q is in the receiving buffer of process p
begin
  take updatesq from the receiving buffer
  while updatesq ≠ ∅ do
    extract (x, v) from updatesq
    if (model = causal) or ((x, ·) ∉ updatesp) then
      xp ← v
    turnp ← (turnp + 1) mod n
  end
end

```

Fig. 1. The algorithm  $A(\text{model})$  for process  $p$ . It is invoked with the parameter  $\text{model}$ , which defines the consistency model that it must implement.

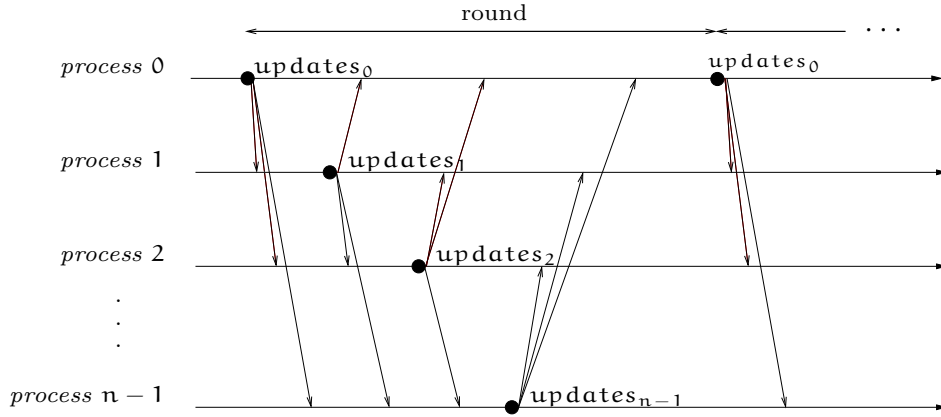


Fig. 2. Cyclic turn fashion.

with a generic broadcast call, which could be simply implemented by sending  $n-1$  point-to-point messages if the underlying message passing subsystem does not provide a more appropriate communication primitive. All this is done by the atomic task `send_updates()`, which also empties the set `updatesp`. The message sent implicitly passes the turn to the next process in order  $(\text{turn}_p + 1) \bmod n$  (see Figure 2).

The atomic task `apply_updates()` is the one in charge of applying the updates received from another process  $q$  in `updatesq`. This task is activated whenever  $\text{turn}_p = q$  and the set `updatesq` is in the receiving buffer of process  $p$ . Note that, when implementing sequential and cache consistency, after a local write operation has been performed in some variable, this task will stop applying the write operations on the same variable from other processes. That allows the system to “view” those writes as if they were overwritten with the write value issued by the local process.

Read operations are always fast with causal and cache consistencies. When implementing sequential consistency, a read operation  $r_p(x)u$  is fast unless `updatesp` contains a pair with a variable different from  $x$ . That is, the read operation is not fast only if, since the latest time it held the turn, process  $p$  has not issued write operations on  $x$  and has issued write operations on other variables. In this case, and only in this case, it is necessary to delay such a read operation until  $\text{turn}_p = p$  for the next time (see Fig. 3). Note that this condition is the same as the condition to execute the task `send_updates()`. We enforce a blocked read operation to have priority over the task `send_updates()`. Hence, when  $\text{turn}_p = p$ , a blocked read operation finished before `send_updates()` is executed.

We have labeled the code of the read operation as atomic because we do not want it to be executed while the variable `updatesp` is manipulated by some other task. However, if the read operation blocks, other tasks are free to access the algorithm variables. In particular, it is necessary that `apply_updates()`



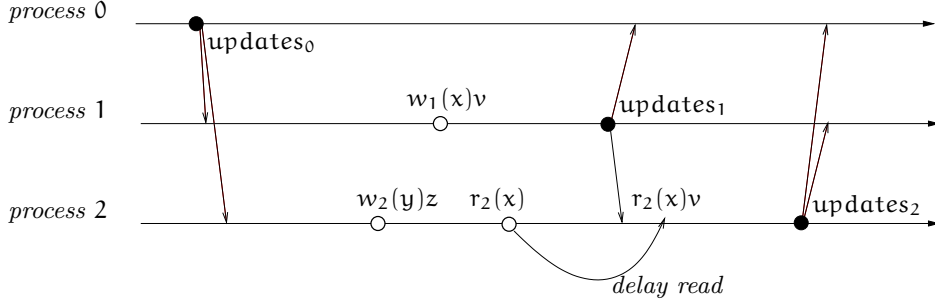


Fig. 3. An example of “non fast” read operation.

updates the variable  $\text{turn}_p$  for the operation to finish eventually.

#### 4 $A(\text{causal})$ Implements Causal Consistency

In this section, we show that the algorithm  $A$ , executed with the parameter  $\text{causal}$ , implements causal consistency. In the rest of this section we assume that  $\alpha$  is the set of operations obtained in the execution of the algorithm  $A(\text{causal})$ , and  $\alpha_p$  is the set of operations obtained by removing from  $\alpha$  all read operations issued by processes other than  $p$ .

**Definition 5** *The  $i^{\text{th}}$  writes of process  $q$ , denoted  $\text{writes}_q^i$ ,  $i > 0$ , is the sequence of all write operations of process  $q$  in  $\alpha_p$ , in the order they are issued, after  $\text{send\_updates}()$  is executed for the  $i^{\text{th}}$  time in this process  $q$ , and before it is executed for the  $i + 1^{\text{st}}$  time.*

For simplicity, we assume that no write operation is issued by any process before it executes  $\text{send\_updates}()$  for the first time. This allows us to consider  $\text{writes}_p^0$  as the empty sequence. Observe in  $A(\text{causal})$  that the  $i + 1^{\text{st}}$  set  $\text{updates}_q$  broadcasted by process  $q$  contains, for each variable, the last (if any) write operation in  $\text{writes}_q^i$  on that variable.

Then, we construct a sequence  $\beta_p$  that we will show in the following lemmas that preserves  $\prec$  and is legal.

**Definition 6** *We denote by  $\beta_p$  the sequence formed with all operations of  $\alpha_p$  as follows. Given the sequence of operations issued by  $p$ , in the order they are issued, we insert the sequence  $\text{writes}_q^i$  in the point of the sequence in which  $\text{apply\_updates}()$  is executed with the set  $\text{updates}_q$  for the  $i + 1^{\text{st}}$  time, for all  $q \neq p$  and  $i \geq 0$ .*

Since the execution of  $\text{apply\_updates}()$  is atomic, it does not overlap any of the operations issued by  $p$ , and the placement of every sequence  $\text{writes}_q^i$  can be easily found.

**Lemma 1** *Let  $\text{op}$  and  $\text{op}'$  be two write operations in  $\alpha_p$  issued by different processes. If  $\text{op} \prec \text{op}'$ , then  $\text{op} \rightarrow \text{op}'$  in  $\beta_p$ .*

**PROOF.** From Definition 1, we know there is a related sequence of operations  $\text{op}^1, \text{op}^2, \dots, \text{op}^m$  such that  $\text{op}^1 = \text{op}$ ,  $\text{op}^m = \text{op}'$ , and  $\text{op}^j \prec_{\text{nt}} \text{op}^{j+1}$  for  $1 \leq j < m$ . This sequence can be divided in  $r$  subsequences of consecutive operations  $s_1, \dots, s_r$ , such that the operations in each subsequence  $s_i$  are issued by the same process, the first operation of  $s_1$  is  $\text{op}$ , the last operation of  $s_r$  is  $\text{op}'$ , and for two consecutive subsequences  $s_i$  and  $s_{i+1}$ , the last operation of  $s_i$  writes the value read by the first operation of  $s_{i+1}$ . Then, from the algorithm  $\mathcal{A}(\text{causal})$  it can be seen that the last operations of two consecutive subsequences  $s_i$  and  $s_{i+1}$ ,  $i \leq r - 1$ , belong, from Definition 6, to sequences  $\text{write}_q^i$  and  $\text{write}_s^j$  such that either  $j > i$ , or  $j = i$  and  $s > q$ . Then,  $\text{op} \rightarrow \text{op}'$  in  $\beta_p$ .  $\square$

**Lemma 2**  *$\beta_p$  preserves the order  $\prec$*

**PROOF.** Let  $\text{op}$  and  $\text{op}'$  be two operations of  $\beta_p$  such that  $\text{op} \prec \text{op}'$ . Let us assume by way of contradiction that  $\text{op}' \rightarrow \text{op}$ .

**Case 1.**  $\text{op}$  and  $\text{op}'$  are issued by the same process. Let us suppose that they are issued by process  $q$ . Recall that  $\alpha_p$  only contains write operations of a process different of  $p$ . From Definition 5, if  $\text{op}' \rightarrow \text{op}$  is because  $\text{op}'$  is executed before  $\text{op}$ , but as we know, from Definition 1, if  $\text{op} \prec \text{op}'$ ,  $\text{op}$  must be executed before  $\text{op}'$ . Hence, we reach a contradiction. Similarly, now let us suppose that  $\text{op}$  and  $\text{op}'$  are issued by  $p$ . From Definition 6, if  $\text{op}' \rightarrow \text{op}$  is because  $\text{op}'$  is executed before  $\text{op}$ , but, from Definition 1, if  $\text{op} \prec \text{op}'$ , then  $\text{op}$  must be executed before  $\text{op}'$ . Hence, we also reach a contradiction.

**Case 2.**  $\text{op}$  and  $\text{op}'$  are issued by different processes. First, let us suppose that  $\text{op}$  and  $\text{op}'$  are operations issued by processes other than  $p$ , and different between them. Then, as  $\alpha_p$  only contains read operations of process  $p$ ,  $\text{op}$  and  $\text{op}'$  must be write operations. Therefore, from Lemma 1,  $\text{op} \rightarrow \text{op}'$ , and we reach a contradiction.

Now, let us suppose that  $\text{op}$  is a read operation issued by  $p$ , and, as  $\alpha_p$  only contains write operations of processes other than  $p$ ,  $\text{op}'$  must be a write operation of a process different of  $p$ . We know that if  $\text{op} \prec \text{op}'$ , we have a related sequence  $\text{op} = \text{op}^1 \prec_{\text{nt}} \text{op}^2 \prec_{\text{nt}} \dots \prec_{\text{nt}} \text{op}^n = \text{op}' = w(y)v$ . If  $\text{op}^i$  is the first write operation after  $\text{op}$ , from Definition 1,  $\text{op}^i$  has to be issued by  $p$ , and  $\text{op}$  has to be issued before  $\text{op}^i$ . Therefore, from Definition 6,  $\text{op} \rightarrow \text{op}^i$ , and, from Lemma 1,  $\text{op}^i \rightarrow \text{op}'$ . Hence,  $\text{op} \rightarrow \text{op}'$ , and we reach a contradiction.

Finally, let us suppose that  $\text{op}'$  is a read operation issued by  $p$ , and, as  $\alpha_p$  only contains write operations of processes other than  $p$ ,  $\text{op}$  must be a write operation of a process different of  $p$ . We know that if  $\text{op} \prec \text{op}'$ , we have a related sequence  $\text{op} = \text{op}^1 = w_q(x)v \prec_{\text{nt}} \dots \prec_{\text{nt}} \text{op}^{n-1} \prec_{\text{nt}} \text{op}'$ . If  $\text{op}^j = w_q(y)v$  is the last write operation before  $\text{op}'$ , from Definition 1,  $\text{op}^j$  has to be executed before  $\text{op}'$ . This implies, with our algorithm  $A(\text{causal})$ , that  $\text{op}^j$  is propagated to process  $p$  before  $\text{op}'$  is issued. Then, this is because  $\text{apply\_updates}()$  is executed in process  $p$  with the set  $\text{updates}_q$  containing  $\text{op}^j$  before  $\text{op}'$  is issued. Therefore, from Definition 6,  $\text{op}^j \rightarrow \text{op}'$ , and, from Lemma 1,  $\text{op} \rightarrow \text{op}^j$ . Hence,  $\text{op} \rightarrow \text{op}'$ , and we reach a contradiction.  $\square$

**Lemma 3**  $\beta_p$  is legal.

**PROOF.** Let us assume, by way of contradiction, that  $\beta_p$  is illegal because there exists  $\text{op}' = w_q(x)u \rightarrow \text{op}'' = w_s(x)v \rightarrow \text{op} = r_p(x)u$  in  $\beta_p$ . From Definition 6, if  $\text{op}'$  precedes  $\text{op}''$  and  $\text{op}''$  precedes  $\text{op}$ , then we have in process  $p$  that: first,  $\text{op}'$  is issued (or applied if  $q \neq p$ ), next,  $\text{op}''$  is issued (or applied if  $s \neq p$ ), and finally,  $\text{op}$  is issued. From our algorithm  $A(\text{causal})$  we can see that, due to these write operations  $\text{op}'$  and  $\text{op}''$ , the local copy  $x_p$  of  $x$  will have the value  $u$  and later the value  $v$ . We can also see that in  $A(\text{causal})$  a read operation always returns the value of the local copy of a variable. Then, it is not possible to have  $\text{op}$  in  $\beta_p$  after  $\text{op}''$ , since it would mean that  $\text{op}$  would have found the value  $v$  in  $x_p$ , instead of the value  $u$ . Hence, we reach a contradiction, and  $\beta_p$  is legal.  $\square$

**Theorem 1** The algorithm  $A(\text{causal})$  implements causal consistency.

**PROOF.** From Lemma 2 and Lemma 3, every execution of the algorithm  $A(\text{causal})$  has a view  $\beta_p$  of  $\alpha_p$ ,  $\forall p$ , that preserves  $\prec$  and legal. Hence, from Definition 4, the algorithm  $A(\text{causal})$  is causal.  $\square$

## 5 $A(\text{sequential})$ Implements Sequential Consistency

In this section, we show that the algorithm  $A$ , executed with the parameter **sequential**, implements sequential consistency. In the rest of this section we assume that  $\alpha$  is the set of operations obtained in the execution of the algorithm  $A(\text{sequential})$ . Any time reference in this section has to do with the time in which the operations of  $\alpha$  are executed. Now we first introduce some definitions of subsets of  $\alpha$ .

**Definition 7** The  $i^{\text{th}}$  iteration of process  $p$ , denoted  $it_p^i$ ,  $i > 0$ , is the subset of  $\alpha$  that contains all the operations issued by process  $p$  after  $\text{send\_updates}()$  is executed for the  $i^{\text{th}}$  time, and before it is executed for the  $i + 1^{\text{st}}$  time.

Observe that any operation in  $it_p^i$  finishes before  $\text{send\_updates}()$  is executed for the  $i + 1^{\text{st}}$  time, since all write and most read operations are fast, and we assume that blocked read operations have priority over the execution of  $\text{send\_updates}()$ .

**Definition 8** The  $i^{\text{th}}$  iteration tail of process  $p$ , denoted  $\text{tail}_p^i$ , is the subset of  $it_p^i$  that includes all the operations from the first write operation (included) until the end of  $it_p^i$ . If  $it_p^i$  does not contain any write operation,  $\text{tail}_p^i$  is the empty sequence.

Observe that all write operations in  $it_p^i$  are in  $\text{tail}_p^i$ . Furthermore, it is easy to check in  $A(\text{sequential})$  that the  $i + 1^{\text{st}}$  set  $\text{updates}_p$  broadcasted by process  $p$  contains, for each variable, the last (if any) write operation in  $\text{tail}_p^i$ .

**Definition 9** The  $i^{\text{th}}$  iteration header of process  $p$ , denoted  $\text{head}_p^i$ , is the subset of  $it_p^i$  that contains all the operations in  $it_p^i$  that are not in  $\text{tail}_p^i$ .

It should be clear that all the operations in  $\text{head}_p^i$  precede all the operations in  $\text{tail}_p^i$  in the execution of  $A$ . We use now the time instants sets received from other processes are applied to partition the sequence  $\text{head}_p^i$ . Note that between the  $i^{\text{th}}$  and the  $i + 1^{\text{st}}$  execution of  $\text{send\_updates}()$  by  $p$  (which defines the operations that are in  $it_p^i$ , and hence in  $\text{head}_p^i$ ) the task  $\text{apply\_updates}()$  is executed  $n - 1$  times, with sets from processes  $(p + 1) \bmod n, \dots, n - 1, 0, \dots, (p - 1) \bmod n$  (in this order).

**Definition 10** The iteration subheader  $q$  of  $\text{head}_p^i$ , denoted  $\text{subhead}_{p,q}^i$ , is the subset of  $\text{head}_p^i$  that contains the following operations.

- If  $q = p$ , then  $\text{subhead}_{p,p}^i$  contains all the operations issued before  $\text{apply\_updates}()$  is executed with the set  $\text{updates}_{(p+1) \bmod n}$ .
- If  $q = (p - 1) \bmod n$ , then  $\text{subhead}_{p,q}^i$  contains all the operations issued after  $\text{apply\_updates}()$  is executed with the set  $\text{updates}_q$ .
- Otherwise,  $\text{subhead}_{p,q}^i$  contains all the operations issued after  $\text{apply\_updates}(\text{mess}_q)$  is executed with the set  $\text{updates}_q$  and before it is executed with the set  $\text{updates}_{(q+1) \bmod n}$ .

Clearly, if the first write operation in  $it_p^i$  is issued before  $\text{apply\_updates}()$  is executed with the set  $\text{updates}_q$ , then  $\text{subhead}_{p,q}^i$  is the empty sequence (see  $it_2^{i-1}$  in Fig. 4).

To simplify the notation and the analysis, we assume that no operation is issued by any process before it executes  $\text{send\_updates}()$  for the first time.

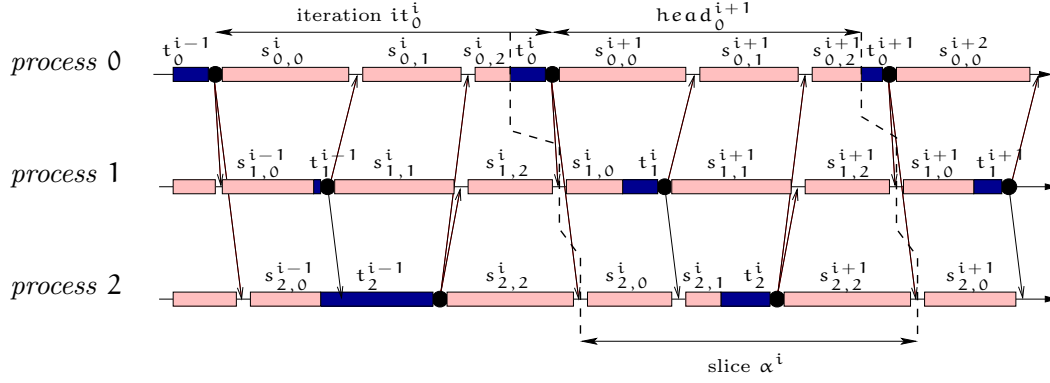


Fig. 4. Iterations and slices. We have abbreviated tail with  $t$  and subhead with  $s$ .

This allows us to define, for any  $p$  and  $q$ , the sequences  $it_p^0$ ,  $tail_p^0$ ,  $head_p^0$ , and  $subhead_{p,q}^0$  as empty sets of operations.

With these definitions, we divide now the set of operations  $\alpha$  in in slices. This division is done in such a way that preserves the order in the execution of  $\alpha$  (see Fig. 4).

**Definition 11** *The  $i^{\text{th}}$  slice of  $\alpha$ , denoted  $\alpha^i$ ,  $i \geq 0$ , is the subset of  $\alpha$  formed by the sets of operations  $tail_p^i, \forall p$ ,  $subhead_{p,q}^i, \forall p, q : p > q$ , and  $subhead_{p,q}^{i+1}, \forall p, q : p \leq q$ .*

Note that, if we consider  $\alpha^0$  the first slice, every operations in  $\alpha$  is in one and only one slice. There are subheaders of iteration 0 that are not assigned to any slice, but since by definition they are empty, they do not need further consideration.

The slice is the basic unit that we will use to define the sequential order that our algorithm enforces. We present now the sequential order for each slice separately. The order for the whole execution is obtained by simply concatenating the slices in their numerical order. However, to complete this sequential order, we yet need to define an order into each subset of operations in tails and subheads that constitute the slice  $\alpha^i$ . Then, from now to the rest of this section, we suppose that the operations into any  $tail_p^i$  and  $subhead_{q,p}^i$  are ordered among them as they were issued by process  $p$ . Hence, we define now, for each slice  $\alpha^i$ , the sequence  $\beta^i$  which contains all the operations of the slice in the sequential order.

**Definition 12** *The sequence  $\beta^i$  is obtained by ordering the operations into each  $tail_p^i$  and  $subhead_p^i$  of  $\alpha^i$  in the order as they were issued by process  $p$ ,*

and by concatenating the set of tails and subheads of  $\alpha^i$  as follows.

$$\begin{aligned}
& \text{tail}_0^i \rightarrow \text{subhead}_{0,0}^{i+1} \rightarrow \text{subhead}_{1,0}^i \rightarrow \text{subhead}_{2,0}^i \rightarrow \dots \rightarrow \text{subhead}_{n-1,0}^i \rightarrow \\
& \text{tail}_1^i \rightarrow \text{subhead}_{0,1}^{i+1} \rightarrow \text{subhead}_{1,1}^{i+1} \rightarrow \text{subhead}_{2,1}^i \rightarrow \dots \rightarrow \text{subhead}_{n-1,1}^i \rightarrow \\
& \dots \\
& \text{tail}_p^i \rightarrow \text{subhead}_{0,p}^{i+1} \rightarrow \dots \rightarrow \text{subhead}_{p,p}^{i+1} \rightarrow \text{subhead}_{p+1,p}^i \rightarrow \dots \rightarrow \text{subhead}_{n-1,p}^i \rightarrow \\
& \dots \\
& \text{tail}_{n-1}^i \rightarrow \text{subhead}_{0,n-1}^{i+1} \rightarrow \text{subhead}_{1,n-1}^{i+1} \rightarrow \text{subhead}_{2,n-1}^{i+1} \rightarrow \dots \rightarrow \text{subhead}_{n-1,n-1}^{i+1}
\end{aligned}$$

In fact, this is only one of many ways to order the sequences of the slice to obtain a sequential order. All the subheaders that appear above in the same line could be permuted in any possible way, since they only contain read operations and each contains operations from a different process. We choose the above order for simplicity.

We define now the sequence  $\beta$ .

**Definition 13**  $\beta$  is the sequence of  $\alpha$  obtained by the concatenation of all sequences  $\beta^i$  in order (i.e.,  $\beta^i \rightarrow \beta^{i+1}, \forall i \geq 0$ ).

From the above definitions, in  $\beta$ , we have that  $\text{tail}_p^i \rightarrow \text{tail}_q^j$  if and only if either  $i < j$  or  $i = j$  and  $p < q$ . This is exactly the order in which the sets associated with each tail are processed and applied in the algorithm.

We show in the following lemmas that  $\beta$  is in fact a view that preserves the order  $\prec$  and is legal.

**Lemma 4**  $\beta$  preserves the order  $\prec$ .

**PROOF.** Let  $\text{op}$  and  $\text{op}'$  be two operations of  $\beta$  such that  $\text{op} \prec \text{op}'$ . From Definition 1, we know that there is a related sequence of operations  $\text{op}^1, \text{op}^2, \dots, \text{op}^m$  such that  $\text{op}^1 = \text{op}$ ,  $\text{op}^m = \text{op}'$ , and  $\text{op}^j \prec_{\text{nt}} \text{op}^{j+1}$  for  $1 \leq j < m$ . If  $\beta$  preserves  $\prec$ , then  $\text{op}^j \rightarrow \text{op}^{j+1}, \forall j$ , and, hence,  $\text{op} \rightarrow \text{op}'$ . We consider several cases.

**Case 1.**  $\text{op}^j$  and  $\text{op}^{j+1}$  are operations issued by the same process. If  $\text{op}^j \prec \text{op}^{j+1}$ , from Definition 1,  $\text{op}^j$  must be issued before  $\text{op}^{j+1}$ . Then, it is easy to check from the above definitions of  $\beta$  and  $\beta^i$  that operations from the same process appear in the same order in  $\beta$  as they were issued. Then,  $\text{op}^j \rightarrow \text{op}^{j+1}, \forall j$ . Hence,  $\text{op} \rightarrow \text{op}'$ .

**Case 2.**  $\text{op}^j$  and  $\text{op}^{j+1}$  are a write operation and a read operation, respectively, issued by different processes. Let us suppose that, from Definition 1,  $\text{op}^j = w_q(x)u$  and  $\text{op}^{j+1} = r_s(x)u$ . We know that if  $\text{op}^j \prec \text{op}^{j+1}$ , then  $\text{op}^j$  must be executed before  $\text{op}^{j+1}$ . From our algorithm  $A(\text{sequential})$  and from the above definitions of  $\beta$  and  $\beta^i$ , we can see that  $\text{op}^j$  always belongs to  $\text{tail}_q^i$ , and in the case of  $\text{op}^{j+1}$  we have two possibilities: a)  $\text{op}^{j+1}$  belongs to  $\text{subhead}_{s,l}^j$ ,  $i < j$ , or if  $i = j$   $q \leq l$ ; or b)  $\text{op}^{j+1}$  belongs to  $\text{tail}_p^j$ ,  $i < j$ . In both cases,  $\text{op}^j \rightarrow \text{op}^{j+1}$ ,  $\forall j$ . Hence,  $\text{op} \rightarrow \text{op}'$ .  $\square$

**Lemma 5**  $\beta$  is legal.

**PROOF.** Let us suppose that there exists  $\text{op}' = w(x)v \rightarrow \text{op} = r(x)v$  in  $\beta$ . From definition Definition 3,  $\beta$  is legal if  $\nexists \text{op}'' = w(x)u \in \alpha$  such that  $\text{op}' \rightarrow \text{op}'' \rightarrow \text{op}$  in  $\beta$ . Then, this is equivalent to say that  $\beta$  is legal if for every read operation  $\text{op} = r(x)v$  in  $\beta$ , the nearest previous write operation in  $\beta$  on the variable  $x$  is  $\text{op}' = w(x)v$ .

Let us assume  $\text{op}$  is issued by process  $p$ . Note first that the order in which the iteration tails appear in  $\beta$  is exactly the order imposed by the token passing procedure. Then, in  $p$ , the order in  $\beta$  reflects exactly the order in which the sets  $\text{updates}_q$  are applied in the local memory of  $p$ . The only exceptions are the sets  $\text{updates}_p$ , since the write operations of  $p$  itself, are applied in its local memory immediately, and do not wait until  $p$  holds the token. However, note that any update from other processes on a variable written locally is not applied (see `apply_updates()`). This gives the illusion that the local write operations have in fact been applied at the time of  $p$ 's token possession. Then we consider several cases.

**Case 1.** Both  $\text{op}$  and  $\text{op}'$  belong to the same iteration tail  $\text{tail}_p^i$ . When issued by  $p$ ,  $\text{op}'$  sets the value of the local local copy  $x_p$  of  $x$ . After  $\text{op}'$  is executed,  $(x, \cdot) \in \text{updates}_p$ , and no update applied from other process changes this value (see `apply_updates()`). Hence, if  $\text{op}$  returns the value  $v$  is because  $\text{op}'$  wrote the value  $v$  in  $x$ .

**Case 2.**  $\text{op}$  belongs to an iteration subheader  $\text{subheader}_{p,q}^i$ . The value  $v$  returned by  $\text{op}$  is the value of  $x_p$  after locally applying the write operations in the following tails.

- If  $p > q$ ,  $\text{tail}_r^j$  for each  $j < i$ , and for each  $r \leq q$  when  $j = i$ .
- If  $p \leq q$ ,  $\text{tail}_r^j$  for each  $j < i - 1$ , and for each  $r \leq q$  when  $j = i - 1$ .

These are the tails that precede  $\text{subheader}_{p,q}^i$  in  $\beta$ . As we said above, these tails are applied in the order they appear in  $\beta$ . Then,  $v$  has to be the value written by the nearest write operation on  $x$  that precedes  $\text{op}$  in  $\beta$ , which by definition is  $\text{op}'$ .

**Case 3.**  $\text{op}$  belongs to an iteration tail  $\text{tail}_p^i$ , while  $\text{op}$  belongs to a different

iteration tail. Then the read operation  $\text{op}$  was issued when  $p$  had already issued a write operation (since it belong to a tail) on a variable different from  $x$  (by definition of  $\text{op}'$ ). Then,  $\text{op}$  was blocked until the token was assigned to  $p$ . The value  $v$  returned by  $\text{op}$  is the value of  $x_p$  after locally applying the write operations in the tails  $\text{tail}_q^i$  for each  $j < i$  and for each  $q < p$  when  $j = i$ , which are the tails that precede  $\text{tail}_p^i$  in  $\beta$ . As we said above, these tails are applied in the order they appear in  $\beta$ . Then,  $v$  has to be the value written by the nearest write operation on  $x$  that precedes  $\text{op}$  in  $\beta$ , which by definition is  $\text{op}'$ .

Thus, in the above three cases we have shown that  $\text{op}' = w(x)v$  is the nearest write operation on variable  $x$  previous to  $\text{op} = r(x)v$  in  $\beta$ . Hence,  $\beta$  is legal.  $\square$

**Theorem 2** *The algorithm  $A(\text{sequential})$  implements sequential consistency.*

**PROOF.** From Lemma 4 and From Lemma 5, every execution of the algorithm  $A(\text{sequential})$  has a view  $\beta$  of  $\alpha$  that preserves the order  $\prec$  and legal. Hence, from Definition 4, the algorithm  $A(\text{sequential})$  is sequential.  $\square$

## 6 $A(\text{cache})$ Implements Cache Consistency

In this section, we show that the algorithm  $A$ , executed with the parameter `cache` in each process, implements cache consistency. In the rest of this section we assume that  $\alpha$  is a set of operations produced in the execution of the algorithm  $A(\text{cache})$ , and  $\alpha(x)$  is a set of operations formed by all operations in  $\alpha$  on the variable  $x$ .

The proof of correctness follows the same lines as the proof of correctness for  $A(\text{sequential})$ , but on  $\alpha(x)$  instead of  $\alpha$ . First we define the sequences  $\text{it}(x)_p^i$ ,  $\text{tail}(x)_p^i$ ,  $\text{head}(x)_p^i$ ,  $\text{subhead}(x)_{p,q}^i$ , and the slice  $\alpha(x)^i$  of  $\alpha(x)$ . Then we construct the sequence  $\beta(x)$  from these sequences in a similar way as the sequence  $\beta$  was defined in Section 5. A version for  $\beta(x)$  of Lemma 4 is directly derived. In a version for  $\beta(x)$  of Lemma 5 with the above sequences the case 3 disappears. Hence we have that  $\beta(x)$  is a view of  $\alpha(x)$  that preserves the order  $\alpha$  and legal. Since this is true for any variable  $x$ , we have the following theorem.

**Theorem 3** *The algorithm  $A(\text{cache})$  implements cache consistency.*

**PROOF.** From Lemma 4 (but only with operations of  $\alpha(x)$ ) and From Lemma 5 (only with operations of  $\alpha(x)$ ), every execution of the algorithm



$A(\text{cache})$  has a view  $\beta(x)$  of  $\alpha(x)$ ,  $\forall x$ , that preserves the order  $\prec$  and legal. Hence, from Definition 4, the algorithm  $A(\text{cache})$  is cache.  $\square$

## 7 Complexity Measures

### 7.1 Worst-Case Response Time

In this section we consider that local operations are executed instantaneously (i.e., in 0 time units) and that any communication takes  $d$  time units. In the algorithm  $A$  executed with parameter `causal` or `cache` all operations are executed locally, while when executed with parameter `sequential` all write and some read operations are also executed locally. Therefore, the response time for them is always 0.

Let us now consider a read operation that is blocked in algorithm  $A(\text{sequential})$ . To obtain the maximum response time for such a read operation, we will consider the worst case. This can happen if the operation blocks (almost) immediately after the process that issued it sent a message. Then, the read operation will be blocked until the turn of this process again, which can take up to  $n$  message transmissions. Therefore, in the worst case, a process will have to wait  $nd$  time units.

The previous analysis assumes that the messages are never delayed at the processes. However, the protocol allows the processes to control when to send the messages. For instance, it is possible for a process  $p$ , when  $\text{turn}_p = p$ , to wait a time  $T$  before executing its task `send_updates()` (see Fig. 1). Thus, we can reduce the number of messages sent by this process per unit of time. Obviously, this can increase the response time, since in this case the delay time of a message sent by  $p$ , in the worst case, will be  $T + d$ .

### 7.2 Message Size

It is easy to check in Fig. 1 that the size of the list `updatesp` of process  $p$  depends on the number of write operations performed by  $p$  during each round, which can be very high. However, the number of pairs  $(x, v)$  in `updatesp` will be, at most, the same as the number of shared variables, since we only hold at most one pair for each variable.

The bound obtained may seem extremely bad. However, note that the real number of pairs in a set `updatesp` really depends on the frequency  $f$  of write operations and the rotation time  $nd$ . Hence if we have a write operation on a

variable every milisecond, in a system with 100 processes and 1 milisecond of delay, we will have at most 100 pairs in the set  $\text{updates}_p$  broadcasted, which is a reasonable number.

Furthermore, note that most algorithms that implement propagation and full replication send a message for *every* write operation that is performed. This would mean that 100 messages would have to be sent. With our algorithm, only one pair per variable is sent, and all of them are grouped into one single message. With the overhead per message in current networks, this implies a significant saving in bandwidth.

### 7.3 Memory Space

Finally, note that we do not require the communication channels among processes to deliver messages in order. Hence, a process could have received messages that are held until the message from the appropriate process arrives. It is easy to check that the maximum number of messages that will ever be held is  $n - 2$ .

## 8 Consistency in $A$ with different parameters

In this section we show that the algorithm  $A$ , executed in some processes with the parameter `sequential` and with `causal` in the rest of processes, implements causal consistency. Also we prove in this section that if there are processes executig  $A(\text{sequential})$  and others  $A(\text{cache})$ , the algorithm  $A$  implements cache consistency.

In this part of this section we assume that  $\alpha$  is the set of operations obtained in the execution of the algorithm  $A(\text{sequential})$  and  $A(\text{causal})$ , and  $\alpha_p$  is the set of operations obtained by removing from  $\alpha$  all read operations issued by processes other than  $p$ .

**Theorem 4** *The algorithm  $A$  implements causal consistency when there exist processes executing  $A(\text{sequential})$  and others  $A(\text{causal})$ .*

**PROOF.** We can see in the Figure 1 that  $A(\text{causal})$  and  $A(\text{sequential})$  have the same mechanism to propagate the write operations. Besides, in both algorithms read operations are always done locally (without generating messages) in the process where they were issued.

Recall that a sequence  $\beta_p$  is formed with every write operation issued by any process and with every read operation issued only by process  $p$ . Hence, from the point of view of each process executing  $A(\text{causal})$ , the existence of processes executing  $A(\text{sequential})$  neither include nor modify the set of operations  $\alpha_p$  to order in  $\beta_p$  different from other process executing  $A(\text{causal})$ . Therefore, if we now construct  $\beta_p$  for each process  $p$  executing  $A(\text{causal})$  as we did in Definition 6,  $\beta_p$  will remain legal and preserving the order  $\prec$ .

Now, in the following two definitions, we redefine  $\beta_p$  for each process  $p$  executing  $A(\text{sequential})$ . After that, we will show that this new sequence  $\beta_p$  preserves  $\prec$  and is legal.

For construct  $\beta_p$  we use, from Section 5, the notation of slice, but now on  $\alpha_p$  instead of  $\alpha$ , to see how the set of operations of  $\alpha_p$  is divided. We also use  $\text{writes}_p^i$  and  $\text{subhead}_{p,q}^i$  as we defined in Section 4 and Section 5, respectively, to see how the set of operations of  $\alpha_p^i$  is ordered.

**Definition 14** *The  $i^{\text{th}}$  slice of  $\alpha_p$ , denoted  $\alpha_p^i$ ,  $i \geq 0$ , is the subset of  $\alpha_p$  formed by the sets of operations  $\text{writes}_q^i, \forall q$ ,  $\text{subhead}_{q,p}^i, \forall q : q > p$ , and  $\text{subhead}_{q,p}^{i+1}, \forall q : q \leq p$ .*

We denote by  $\beta_p^i$  the sequence of all operations of the slice  $\alpha_p^i$ , and by  $\beta_p$  the sequence of  $\alpha_p$  formed by the concatenation of  $\beta_p^i$  in increasing numerical order

**Definition 15** *The sequence  $\beta_p^i, \forall p$  executing  $A(\text{sequential})$ , is obtained by ordering the operations into each  $\text{writes}_p^i$  and  $\text{subhead}_p^i$  of  $\alpha_p^i$  in the order as they were issued by process  $p$ , and by concatenating the set of writes and subheads of  $\alpha_p^i$  as follows.*

$$\begin{aligned} & \text{writes}_0^i \rightarrow \text{subhead}_{p,0}^i \rightarrow \\ & \text{writes}_1^i \rightarrow \text{subhead}_{p,1}^i \rightarrow \\ & \dots \\ & \text{tail}_p^i \rightarrow \text{subhead}_{p,p}^{i+1} \rightarrow \\ & \dots \\ & \text{writes}_{n-1}^i \rightarrow \text{subhead}_{p,n-1}^{i+1} \end{aligned}$$

**Definition 16** *The sequence  $\beta_p, \forall p$  executing  $A(\text{sequential})$ , is the sequence of  $\alpha$  obtained by the concatenation of all sequences  $\beta_p^i$  in order (i.e.,  $\beta_p^i \rightarrow \beta_p^{i+1}, \forall i \geq 0$ ).*

As we can see comparing Definition 12 and Definition 15,  $\beta_p^i$  is the same sequence that  $\beta^i$  (and, therefore,  $\beta_p$  and  $\beta$ ) but with the following two differences:

- $\text{writes}_q^i, q \neq p$ , is the same sequence that  $\text{it}_q^i$  where we have removed every

- read operation issued by  $q$ .
- We have removed every subhead with operations other than process  $p$ . That is to say, every  $\text{subhead}_{q,n}^j$  such that  $q \neq p$ ,  $j = i$  or  $j = i + 1$ , and  $\forall n$ .

As we can see, we have constructed  $\beta_p$  in a similar way as the sequence  $\beta$  was defined in Section 5. Then, all operations belonging to  $\beta_p$  are in the same order as in the definition of  $\beta$ , and a version for  $\beta_p$  of Lemma 4 is directly derived. Hence,  $\beta_p$  preserves the order  $\prec$ .

Similarly, we know that all write operations are in  $\beta_p$  in the same order as in the definition of  $\beta$ . Then, a version for  $\beta_p$  of Lemma 5 is also directly derived. Hence,  $\beta_p$  is legal.

Thus, we have for each process  $p$  executing  $A(\text{causal})$  a  $\beta_p$  formed as we described in Section 4. We also have for each process  $q$  executing  $A(\text{sequential})$  a  $\beta_q$  formed as Definition 16. As in both cases  $\beta_p$  is legal and preserves  $\prec$ , we can affirm that  $A$  implements causal consistency when there are processes executing  $A(\text{sequential})$  and others  $A(\text{causal})$ .  $\square$

In this part of this section we assume that  $\alpha$  is the set of operations obtained in the execution of the algorithm  $A(\text{sequential})$  and  $A(\text{cache})$ , and  $\alpha(x)$  is the set of operations of  $\alpha$  on variable  $x$ .

**Theorem 5** *The algorithm  $A$  implements cache consistency when there exist processes executing  $A(\text{sequential})$  and others  $A(\text{cache})$ .*

**PROOF.** We can see in the Figure 1 that  $A(\text{cache})$  and  $A(\text{sequential})$  are the same algorithm except in one case: when a read operation is not fast. Therefore, they use the same mechanism to propagate the write operations, and read operations do not add new messages because they are done locally in the process where it is invoked.

Hence, from the point of view of each process executing  $A(\text{cache})$ , the existence of processes executing  $A(\text{sequential})$  neither include nor modify the set of operations to order in  $\beta(x)$  different from other process executing  $A(\text{cache})$ . Therefore, if we now construct  $\beta(x)$  for each process  $p$  executing  $A(\text{cache})$  as we did in Section 6,  $\beta(x)$  will remain legal and preserving  $\prec$ .

Now, we are going to use the same definition and way of construction of  $\beta(x)$  from Section 6 for all processes that executes  $A(\text{sequential})$ . Then, the difference between the sequence  $\beta(x)$  for processes executing  $A(\text{sequential})$  with respect to processes executing  $A(\text{cache})$  is what happen when a non-fast read operation occurs. To analyse this case, we use the same notation from Section 6 to define a slice  $\alpha^i(x)$ , a tail  $\text{tail}_p^i(x)$ , and a subhead  $\text{subhead}_{q,p}^j(x)$ .

Then, let us suppose that a non-fast read operation happens in a process  $p$  executing  $A(\text{sequential})$  in the slice  $\alpha^i(x)$ . We know, by definition of  $\beta(x)$ , that this read operation belongs to  $\text{tail}_p^i(x)$ . We are also going to suppose, that the first operation of  $\text{tail}_p^i(x)$  is done just after  $\text{subhead}_{q,p}^j(x)$ ,  $i = j$  if  $q > p$ , or  $j = i + 1$  if  $q \leq p$ . Hence, in this case, the unique difference with respect to an execution in a process  $A(\text{cache})$  is that each subhead subsequent to  $\text{subhead}_{q,p}^j(x)$  in  $\alpha^i(x)$  is always empty. Therefore, as the sequence  $\beta(x)$  of a process executing  $A(\text{cache})$  from Section 6, A version of Lemma 4 and Lema 5 for  $\beta(x)$  of a process executing  $A(\text{sequential})$  are directly derived. Hence,  $\beta(x)$  of a process executing  $A(\text{sequential})$  is legal and preserves the order  $\prec$ .

Thus, we have shown for each process  $p$  executing  $A(\text{cache})$  or  $A(\text{sequential})$  that there is a  $\beta(x)$ ,  $\forall x$ , formed as we described in Section 6, such that preserves the order  $\prec$  and legal. Therefore,  $A$  implements cache consistency when there exist processes executing  $A(\text{sequential})$  and others  $A(\text{cache})$ .  $\square$

## 9 Conclusions and Future Work

In this paper, we have presented a parametrized algorithm that implements sequential, causal, and cache consistency on a distributed system. To our knowledge, this is the first algorithm that implements cache consistency.

The algorithm presented in this paper guarantees fast operations in its causal and cache executions. It is proven in [7] that it is imposible to have a sequential algorithm with all operations fast. The algorithm presented in this paper guarantees in its sequential execution fast writes and reduces to only one case the reads that can not be executed locally.

Considering possible extensions of this work for the sequential version, we would like to know how many read operations are fast in real applications with several system parameters. Our belief is that most read operations will be fast. A second line of work has to do with the scalability of the protocol. The worst case response time is linear on the number of processes. Hence, it will not scale well, since it may become high when the system has a large number of processes. It would be nice to remove this dependency. Finally, the protocol works in a token passing fashion, which can be very risky in an environment with failures, since a single failure can block the whole system. It would be interesting to extend the protocol with fault tolerance features.

## References

- [1] R. C. Steinke, G. J. Nutt, A unified theory of shared memory consistency, *J. ACM* 51 (5) (2004) 800–849.
- [2] C. Manovit, S. Hangal, Efficient algorithms for verifying memory consistency, in: *SPAA'05: Proceedings of the 17th annual ACM symposium on Parallelism in algorithms and architectures*, ACM Press, New York, NY, USA, 2005, pp. 245–252.
- [3] M. Ahamad, G. Neiger, J. Burns, P. Kohli, P. Hutto, Causal memory: Definitions, implementation and programming, *Distributed Computing* 9 (1) (1995) 37–49.
- [4] R. Prakash, M. Raynal, M. Singhal, An adaptive causal ordering algorithm suited to mobile computing environments, *Journal of Parallel and Distributed Computing* 41 (1997) 190–204.
- [5] M. Raynal, M. Ahamad, Exploiting write semantics in implementing partially replicated causal objects, in: *Proceedings of the 6th EUROMICRO Conference on Parallel and Distributed Computing*, 1998, pp. 157–163.
- [6] L. Lamport, How to make a multiprocessor computer that correctly executes multiprocess programs, *IEEE Transactions on Computers* 28 (9) (1979) 690–691.
- [7] H. Attiya, J. Welch, Sequential consistency versus linearizability, *ACM Transactions on Computer Systems* 12 (2) (1994) 91–122.
- [8] M. Raynal, Token-based sequential consistency., *Comput. Syst. Sci. Eng.* 17 (6) (2002) 359–365.
- [9] M. Raynal, K. Vidyasankar, A distributed implementation of sequential consistency with multi-object operations., in: *ICDCS*, 2004, pp. 544–551.
- [10] V. Cholvi, J. Bernabéu, Relationships between memory models., *Inf. Process. Lett.* 90 (2) (2004) 53–58.
- [11] H. Attiya, R. Friedman, A correctness condition for high-performance multiprocessors, in: *Proceedings of the 24th Annual ACM Symposium on the Theory of Computing*, 1992, pp. 679–690.
- [12] J. Misra, Axioms for memory access in asynchronous hardware systems, *ACM Transactions on Programming Languages and Systems* 8 (1) (1986) 142–153.
- [13] M. Raynal, A. Schiper, From causal consistency to sequential consistency in shared memory systems, *Tech. Rep. 926*, IRISA (May 1996).
- [14] G. B. Yehuda Afek, M. Merritt, Lazy caching, *ACM Transactions on Programming Languages and Systems* 15 (1) (1993) 182–205.
- [15] A. Singh, Bounded timestamps in process networks, *Parallel Processing Letters* 6 (2) (1996) 259–264.