

# About the efficiency of partial replication to implement Distributed Shared Memory

Jean-Michel Héлары  
IRISA  
Campus de Beaulieu, 35042 Rennes-cedex, France  
helary@irisa.fr

Alessia Milani  
DIS, Università di Roma  
La Sapienza  
Via Salaria 113, Roma, Italia  
Alessia.Milani@dis.uniroma1.it

## Abstract

*Distributed Shared Memory abstraction (DSM) is traditionally realized through a distributed memory consistency system (MCS) on top of a message passing system. In this paper we analyze the impossibility of efficient partial replication implementation of causally consistent DSM. Efficiency is discussed in terms of control information that processes have to propagate to maintain consistency. We introduce the notions of share graph and hoop to model variable distribution and the concept of dependency chain to characterize processes that have to manage information about a variable even though they do not read or write that variable. Then, we consider PRAM, a consistency criterion weaker enough to allow efficient partial replication implementations and strong enough to solve interesting problems. Finally, we illustrate the power of PRAM with the Bellman-Ford shortest path algorithm.*

## 1. Introduction

Distributed Shared Memory (DSM) is one of the most interesting abstraction providing data-centric communication among a set of application processes which are decoupled in time, space and flow. This abstraction allows programmers to design solutions by considering the well-known shared variables programming paradigm, independently of the system (centralized or distributed) that will run his program. Moreover, there are a lot of problems (in numerical analysis, image or signal processing, to cite just a few) that are easier to solve by using the shared variables paradigm rather than using the message passing one.

Distributed shared memory abstraction is traditionally realized through a distributed *memory consistency system* (MCS) on top of a message passing system providing a communication primitive with a certain quality of service in

terms of ordering and reliability [2]. Such a system consists of a collection of nodes. On each node there is an application process and a MCS process. An application process invokes an operation through its local MCS process which is in charge of the actual execution of the operation. To improve performance, the implementation of MCS is based on replication of variables at MCS processes and propagation of the variable updates [4]. As variables can be concurrently accessed (by read and write operations), users must be provided with a consistency criterion that precisely defines the semantics of the shared memory. Such a criterion defines the values returned by each read operation executed on the shared memory. Many consistency criteria have been considered, e.g., from more to less constraining ones: Atomic [7], Sequential [6], Causal [8] and PRAM (PipelinedRAM) [16]. Less constraining MCS are easier to implement, but, conversely, they offer a more restricted programming model. The Causal consistency model has gained interest because it offers a good tradeoff between memory access order constraints and the complexity of the programming model as well as of the complexity of the memory model itself. To improve performance, MCS enforcing Causal (or stronger) consistency have been usually implemented by protocols based on complete replication of memory locations [1, 9, 14], i.e. each MCS process manages a copy of each shared variable. It is easy to notice that in the case of complete replication, dealing with a large number of shared variables avoids scalability. Thus, in large scale systems, implementations based on partial replication, i.e. each process manages only a subset of shared variables, seems to be more reasonable. Since each process in the system could be justifiably interested only in a subset of shared variables, partial replication is intended to avoid a process to manage information it is not interested in. In this sense, partial replication loses its meaning if to provide consistent values to the corresponding application process, each MCS process has to consider information about vari-

ables that the corresponding application process will never read or write. Some implementations are based on partial replication [15, 12], but they suffer this drawback.

In this paper we study the problem of maintaining consistency in a partial replicated environment. More precisely, according to the variables distribution and to the consistency criterion chosen, we discuss the possibility of an *efficient partial replication implementation*, i.e., for each shared variable, only MCS processes owning a local copy have to manage information concerning this variable. Our study shows that MCS enforcing Causal consistency criterion (or stronger consistency criteria) have no efficient partial replication implementation. Then, it is shown that the PRAM consistency criterion is weak enough to allow efficient partial replication implementation.

The rest of the paper is organized as follows. In Section 2 we present the shared memory model. In Section 3, we discuss partial replication issues and we present our main result, namely a characterization for the possibility of efficient partial replication implementation. Finally, section 4 is devoted to the PRAM consistency criterion and section 5 to the solution of Bellman-Ford algorithm in such a MCS.

## 2. The Shared Memory Model

We consider a finite set of sequential *application* processes  $\Pi = \{ap_1, ap_2, \dots, ap_n\}$  interacting via a finite set of shared variables,  $\mathcal{X} = \{x_1, x_2, \dots, x_m\}$ . Each variable  $x_h$  can be accessed through *read* and *write* operations. A write operation invoked by an application process  $ap_i$ , denoted  $w_i(x_h)v$ , stores a new value  $v$  in variable  $x_h$ . A read operation invoked by an application process  $ap_i$ , denoted  $r_i(x_h)v$ , returns to  $ap_i$  the value  $v$  stored in variable  $x_h$ <sup>1</sup>. Each variable has an initial value  $\perp$ .

A *local history* of an application process  $ap_i$ , denoted  $h_i$ , is a sequence of read and write operations performed by  $ap_i$ . If an operation  $o_1$  precedes an operation  $o_2$  in  $h_i$ , we say that  $o_1$  precedes  $o_2$  in *program order*. This precedence relation, denoted by  $o_1 \mapsto_i o_2$ , is a total order. A *history*  $H = \langle h_1, h_2, \dots, h_n \rangle$  is the collection of local histories, one for each application process. The set of operations in a history  $H$  is denoted  $O_H$ .

Operations done by distinct application processes can be related by the *read-from order* relation. Given two operations  $o_1$  and  $o_2$  in  $O_H$ , the read-from order relation,  $\mapsto_{ro}$ , on some history  $H$  is any relation with the following properties [8]<sup>2</sup>:

- if  $o_1 \mapsto_{ro} o_2$ , then there are  $x$  and  $v$  such that  $o_1 = w(x)v$  and  $o_2 = r(x)v$ ;
- for any operation  $o_2$ , there is at most one operation  $o_1$  such that  $o_1 \mapsto_{ro} o_2$ ;
- if  $o_2 = r(x)v$  for some  $x$  and there is no operation  $o_1$  such that  $o_1 \mapsto_{ro} o_2$ , then  $v = \perp$ ; that is, a read with no write must read the initial value.

Finally, given a history  $H$ , the causality order  $\mapsto_{co}$ , [8], is a partial order that is the transitive closure of the union of the history's program order and the read-from order. Formally, given two operations  $o_1$  and  $o_2$  in  $O_H$ ,  $o_1 \mapsto_{co} o_2$  if and only if one of the following cases holds:

- $\exists ap_i$  s.t.  $o_1 \mapsto_i o_2$  (program order),
- $\exists ap_i, ap_j$  s.t.  $o_1$  is invoked by  $ap_i$ ,  $o_2$  is invoked by  $ap_j$  and  $o_1 \mapsto_{ro} o_2$  (read-from order),
- $\exists o_3 \in O_H$  s.t.  $o_1 \mapsto_{co} o_3$  and  $o_3 \mapsto_{co} o_2$  (transitive closure).

If  $o_1$  and  $o_2$  are two operations belonging to  $O_H$ , we say that  $o_1$  and  $o_2$  are *concurrent* w.r.t.  $\mapsto_{co}$ , denoted  $o_1 \parallel_{co} o_2$ , if and only if  $\neg(o_1 \mapsto_{co} o_2)$  and  $\neg(o_2 \mapsto_{co} o_1)$ .

### Properties of a history

**Definition 1 (Serialization).** Given a history  $H$ ,  $S$  is a *serialization* of  $H$  if  $S$  is a sequence containing exactly the operations of  $H$  such that each read operation of a variable  $x$  returns the value written by the most recent precedent write on  $x$  in  $S$ .

A serialization  $S$  respects a given order if, for any two operations  $o_1$  and  $o_2$  in  $S$ ,  $o_1$  precedes  $o_2$  in that order implies that  $o_1$  precedes  $o_2$  in  $S$ . Let  $H_{i+w}$  be the history containing all operation in  $h_i$  and all write operations of  $H$ .

**Definition 2 (Causally Consistent History [8]).** A history  $H$  is *causally consistent* if for each application process  $ap_i$  there is a serialization  $S_i$  of  $H_{i+w}$  that respects  $\mapsto_{co}$ .

A memory is causal if it admits only causally consistent histories.

## 3. The problem of efficient partial replication implementation of causal memories

In this section we analyze the efficiency of implementing causal memories when each application process  $ap_i$  accesses only a subset of the shared variables  $\mathcal{X}$ , denoted  $\mathcal{X}_i$ . Assuming a partial replicated environment means that each MCS process  $p_i$  manages a replica of a variable  $x$  iff  $x \in \mathcal{X}_i$ . Our aim is to determine which MCS processes are concerned by information on the occurrence of operations performed on the variable  $x$  in the system. More precisely, given a variable  $x$ , we will say that a MCS process

<sup>1</sup>Whenever we are not interested in pointing out the value or the variable or the process identifier, we omit it in the notation of the operation. For example  $w$  represents a generic write operation while  $w_i$  represents a write operation invoked by the application process  $ap_i$ , etc.

<sup>2</sup>It must be noted that the read-from order relation just introduced is the same as the writes-into relation defined in [8].

$p_i$  is  $x$ -relevant if, in at least one history, it has to transmit some information on the occurrence of operations performed on variable  $x$  in this history, to ensure a causally consistent shared memory. Of course, each process managing a replica of  $x$  is  $x$ -relevant. Ideally, we would like that only those processes are  $x$ -relevant. But unfortunately, as will be proved in this section, if the variable distribution is not known a priori, it is not possible for the MCS to ensure a causally consistent shared memory, if each MCS process  $p_i$  only manages information about  $\mathcal{X}_i$ . The main result of the section is a characterization of  $x$ -relevant processes.

To this aim, we first introduce the notion of *share graph*, denoted  $SG$ , to characterize variable distribution and then we define the concepts of *hoop* and of *dependency chain* to highlight how particular variables distribution can impose global information propagation.

### 3.1. The share graph, hoops and dependency chains

The share graph is an undirected (symmetric) graph whose vertices are *processes*, and an edge  $(i, j)$  exists between  $p_i$  and  $p_j$  iff there exists a variable  $x$  replicated both on  $p_i$  and  $p_j$  (i.e.  $x \in \mathcal{X}_i \cap \mathcal{X}_j$ ). Possibly, each edge  $(i, j)$  is labelled with the set of variables replicated both on  $p_i$  and  $p_j$ . From this definition results that two processes can communicate (through a shared variable) if and only if they are linked by an edge in the share graph.

Figure 1 depicts an example of share graph representing a system of three processes  $p_i, p_j$  and  $p_k$  interacting through the following set of shared variables  $\mathcal{X} = \{x_1, x_2\}$ . In particular,  $\mathcal{X}_i = \{x_1, x_2\}$ ,  $\mathcal{X}_k = \{x_2\}$  and  $\mathcal{X}_j = \{x_1\}$ .

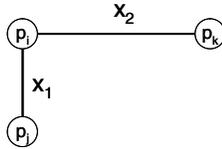


Figure 1. A share graph

It is simple to notice that each variable  $x$  defines a subgraph  $C(x)$  of  $SG$  spanned by the processes on which  $x$  is replicated (and the edges having  $x$  on their label). This subgraph  $C(x)$  is a clique, i.e. there is an edge between every pair of vertices. The "share graph" is the union of all cliques  $C(x)$ . Formally,  $SG = \bigcup_{x \in \mathcal{X}} C(x)$ .

In the example depicted in Figure 1, we have the following cliques:

- i)  $C(x_1) = (V_{x_1}, E_{x_1})$  where  $V_{x_1} = \{p_i, p_j\}$  and  $E_{x_1} = \{(i, j)\}$ ,
- ii)  $C(x_2) = (V_{x_2}, E_{x_2})$  where  $V_{x_2} = \{p_i, p_k\}$  and  $E_{x_2} = \{(i, k)\}$ .

Given a variable  $x$ , we call  $x$ -hoop, any path of  $SG$ , between two distinct processes in  $C(x)$ , whose intermediate vertices do not belong to  $C(x)$  (figure 2). Formally:

**Definition 3 (Hoop).** Given a variable  $x$  and two processes  $p_a$  and  $p_b$  in  $C(x)$ , we say that there is a  $x$ -hoop between  $p_a$  and  $p_b$  (or simply a hoop, if no confusion arises), if there exists a path  $[p_a = p_0, p_1, \dots, p_k = p_b]$  in  $SG$  such that:

- i)  $p_h \notin C(x)$  ( $1 \leq h \leq k-1$ ) and
- ii) each consecutive pair  $(p_{h-1}, p_h)$  shares a variable  $x_h$  such that  $x_h \neq x$  ( $1 \leq h \leq k$ )

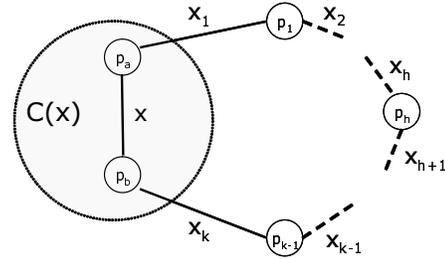


Figure 2. An  $x$ -hoop

Let us remark that the notion of hoop depends only on the distribution of variables on the processes, i.e. on the topology of the corresponding share graph. In particular, it is independent of any particular history.

**Definition 4 (Minimal Hoop).** A  $x$ -hoop  $[p_a = p_0, p_1, \dots, p_k = p_b]$  is said to be minimal, iff i) each edge of the hoop is labelled with a different variable and ii) none of the edge label is shared by processes  $p_a$  and  $p_b$ .

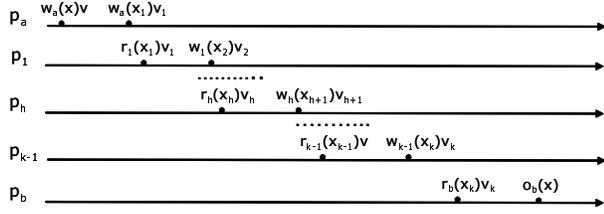
The following concept of *dependency chain along an hoop* captures the dependencies that can be created between operations occurring in a history, when these operations are performed by processes belonging to a hoop.

**Definition 5 (Dependency chain).** Let  $[p_a, \dots, p_b]$  be a  $x$ -hoop in a share graph  $SG$ . Let  $H$  be a history. We say that  $H$  includes a  $x$ -dependency chain<sup>3</sup> along this hoop if the following three conditions are verified:

- $O_H$  includes  $w_a(x)v$ , and
- $O_H$  includes  $o_b(x)$ , where  $o_b$  can be a read or a write on  $x$ , and
- $O_H$  includes a pattern of operations, at least one for each process belonging to the hoop, that implies  $w_a(x)v \mapsto_{co} o_b(x)$ .

More precisely, we also say that  $w_a(x)v$  and  $o_b(x)$  are the *initial* and the *final* operations of the  $x$ -dependency chain from  $w_a(x)v$  to  $o_b(x)$ . Figure 3 depicts such a dependency chain.

<sup>3</sup>simply  $x$ -dependency chain when confusion cannot arise



**Figure 3. An  $x$ -dependency chain from  $w_a(x)v$  to  $o_b(x)$**

### 3.2. A characterization of $x$ -relevant processes

In this section, a characterization of  $x$ -relevant processes, where  $x$  is a variable, is given.

**Proposition 1.** *Given a variable  $x$ , a process  $p_i$  is  $x$ -relevant if it belongs to  $C(x)$  or to a minimal  $x$ -hoop.*

*Proof.* If  $p_i \in C(x)$ , then it is obviously  $x$ -relevant. Consider now a process  $p_i \notin C(x)$ , but belonging to a minimal  $x$ -hoop between two processes in  $C(x)$ , namely  $p_a$  and  $p_b$ . Let  $[p_a = p_0, p_1, \dots, p_k = p_b]$  be such minimal  $x$ -hoop, then the history  $H$  depicted in Figure 3 can be generated.  $H$  includes an  $x$ -dependency chain along this hoop from  $w_a(x)v$  to  $o_b(x)v$  and, by definition 5, it follows that  $w_a(x)v \mapsto_{co} o_b(x)$ . Thus if  $o_b(x)$  is a read operation, the value that can be returned is constrained by the operation  $w_a(x)v$ , i.e., to ensure causal consistency, it cannot return neither  $\perp$  nor any value written by a write operation belonging to the causal past of  $w_a(x)v$ . Similarly, if  $o_b(x)$  is a write operation, namely  $o_b(x) = w_b(x)v'$ , the dependency  $w_a(x)v \mapsto_{co} w_b(x)v'$  implies that, to ensure causal consistency, if a process  $p_c \in C(x)$  reads both values  $v$  and  $v'$  then it reads them in such a order.

In both cases, to ensure causal consistency, process  $p_b$  has to be aware that  $w_a(x)v$  is in the causal past of  $w_{k-1}(x_k)v_k$ . Since  $p_a$  cannot be aware of  $w_a(x)v$  causal future, this information has to arrive from  $p_{k-1}$ . Let us remark that  $w_a(x)v \mapsto_{co} o_b(x)$  since  $w_a(x)v \rightarrow_a w_a(x_1)v_1$ , and, for each  $h$ ,  $1 \leq h \leq k-1$ ,  $r_h(x_h)v_h \rightarrow_h w_h(x_{h+1})v_{h+1}$ , and  $r_b(x_k)v_k \rightarrow_b o_b(x)$ . This means that  $w_a(x)v$  is in the causal past of  $w_h(x_{h+1})v_{h+1}$  because  $w_a(x)v$  is in the causal past of  $w_{h-1}(x_h)v_h$ , for each  $h$  such that  $1 \leq h \leq k-1$ . Moreover, since the hoop is minimal, the only way for process  $p_h$  to be aware that  $w_a(x)v$  is in the causal past of  $w_{h-1}(x_h)v_h$  is through  $p_{h-1}$ , for each  $h$  s.t.  $1 \leq h \leq k-1$ . Then for each  $h$  such that  $1 \leq h \leq k-1$ ,  $p_h$  is  $x$ -relevant. In particular  $p_i$  is  $x$ -relevant.  $\square$

**Proposition 2.** *Given a variable  $x$ , if a process  $p_i$  is  $x$ -relevant then it belongs to  $C(x)$  or it belongs to a  $x$ -hoop.*

*Proof.* The analysis above shows that the purpose of transmitting control information concerning the variable  $x$  is to ensure causal consistency. In particular, if an operation  $o_1 = w_a(x)v$  is performed by a process  $p_a \in C(x)$ , then any operation  $o_2 = o_b(x)$  performed by another process  $p_b \in C(x)$  is constrained by  $o_1$  only if  $o_1 \mapsto_{co} o_2$ .

We have that  $o_1 \mapsto_{co} o_2$  only if one of the two following cases holds:

- A "direct" relation:  $o_1 \mapsto_{ro} o_2$ . In this case, no third part process is involved in the transmission of information concerning the occurrence of the operation  $o_1$ .
- An "indirect" relation: there exists at least one  $o_h$  such that  $o_1 \mapsto_{co} o_h$  and  $o_h \mapsto_{co} o_2$ . Such an indirect relation involve a sequence  $\sigma$  of processes  $p_0 = p_a, \dots, p_h, \dots, p_k = p_b$  ( $k \geq 2$ ) such that two consecutive processes  $p_{h-1}$  and  $p_h$  ( $1 \leq h \leq k$ ) respectively perform operations  $o_{h-1}$  and  $o_h$  with  $o_{h-1} \mapsto_{ro} o_h$ . This implies that there exists a variable  $x_h$  such that  $o_{h-1} = w_{h-1}(x_h)v_h$  and  $o_h = r_h(x_h)v_h$ . Consequently,  $x_h$  is shared by  $p_{h-1}$  and  $p_h$ , i.e.,  $p_{h-1}$  and  $p_h$  are linked by an edge in the graph  $SG$ , meaning that the sequence  $\sigma$  is a path between  $p_a$  and  $p_b$  in the share graph  $SG$ . Such a path is either completely included in  $C(x)$ , or is a succession of  $x$ -hoops, and along each of them there is a  $x$ -dependency chain. Thus, a process  $p_i \notin C(x)$  and not belonging to any  $x$ -hoop cannot be involved in these dependency chains. The result follows from the fact that this reasoning can be applied to any variable  $x$ , then to any pair of processes  $p_a$  and  $p_b$  in  $C(x)$ , and finally to any  $x$ -dependency chain along any  $x$ -hoop between  $p_a$  and  $p_b$ .  $\square$

### 3.3. Impossibility of efficient partial replication

Shared memory is a powerful abstraction in large-scale systems spanning geographically distant sites; these environments are naturally appropriate for distributed applications supporting collaboration. Two fundamental requirements of large-scale systems are scalability and low-latency accesses:

- to be scalable a system should accommodate large number of processes and should allow applications to manage a great deal of data;
- in order to ensure low latency in accessing shared data, copy of interested data are replicated at each site.

According to this, causal consistency criterion has been introduced by Ahamad et al. [13], [8] in order to avoid large latencies and high communication costs that arise in implementing traditional stronger consistency criteria, e.g., atomic [7] and sequential consistency [6]. "Many applications can easily be programmed with shared data that is causally consistent, and there are efficient and scalable implementations of causal memory" [13]. In particular, low

latency is guaranteed by allowing processes to access local copy of shared data through wait-free operations. It means that causal consistency reduces the global synchronization between processes which is necessary to return consistent values.

This criterion is meaningful in systems in which complete replication is requested, i.e., when each process accesses all data in the system. On the other hand, considering large scale system with a huge and probably increasing number of processes and data, partial replication seems to be more reasonable: each process can directly access data it is interested in without introducing a heavy information flow in the network. From the results obtained in Section 3.2, several observations can be made, depending which is the *a priori* knowledge on variable distribution.

If a particular distribution of variables is assumed, it could be possible to build the share graph and analyze it off-line in order to enumerate, for each variable  $x$  not totally replicated, all the minimal  $x$ -hoops. It results from Proposition 1 that only processes belonging to one of these  $x$ -hoops will be concerned by the variable  $x$ . Thus, an ad-hoc implementation of causal DSM can be optimally designed. However, even under this assumption on variable distribution, enumerating all the hoops can be very long because it amounts to enumerate a set of paths in a graph that can be very big if there are many processes. In a more general setting, implementations of DSM cannot rely on a particular and static variable distribution, and, in that case, any process is likely to belong to any hoop. It results from Proposition 1 that each process in the system has to transmit control information regarding all the shared data, contradicting scalability.

Thus, causal consistency does not appear as the most appropriate consistency criterion for large-scale systems. For this reason, in the next section, we weaken the causal consistency in order to find a consistency criterion that allows efficient partial replication implementations of the shared memory, while being strong enough to solve interesting problems.

#### 4. Weakening the causal consistency criterion: PRAM

In this Section we consider three consistency criteria, obtained by successive relaxations of causal consistency, namely: *Lazy Causal*, *Lazy Semi-Causal* and *PRAM* consistency criteria. The *Lazy Causal* criterion relaxes the program order by observing that some operations performed by a process could be permuted without effect on the output of the program (e.g. two successive read operations on two different variables). The *Lazy Semi-Causal Consistency* is based on the previous lazy program order and on a relaxation of the read-from order relation, namely *weak writes-*

*before*, introduced by Ahamad et al. [10]. Due to the lack of space, these criteria are presented in details in the full paper [5], where it is also shown that they are still too strong to allow efficient partial replication.

The last possibility is to weaken the transitivity property such that two operations executed by different processes can be related only by the direct read-from relation. The PRAM consistency criterion [16] relaxes the transitivity due to intermediary processes [11]. In other words, it only requires that all processes observe the writes performed by a same process in the same order, while they may disagree on the order of writes by different processes. The PRAM consistency is based on a relation, denoted  $\mapsto_{pram}$ , weaker than  $\mapsto_{co}$ . Formally [11]<sup>4</sup>:

**Definition 6 (PRAM relation).** *Given two operations  $o_1$  and  $o_2$  in  $O_H$ ,  $o_1 \mapsto_{pram} o_2$  if, and only if, one of the following conditions holds:*

1.  $\exists p_i : o_1 \mapsto_i o_2$  (program order), or
2.  $\exists p_i \exists p_j \ i \neq j : o_1 = w_i(x)v$  and  $o_2 = r_j(x)v$ , i.e.  $o_1 \mapsto_{ro} o_2$  (read-from relation).

Note that  $\mapsto_{pram}$  is an acyclic relation, but is not a partial order due to the lack of transitivity.

**Definition 7 (PRAM consistent history).** *A history  $H$  is PRAM consistent if, for each application process  $ap_i$ , there exists a serialization  $H_{i+w}$  that respects  $\mapsto_{pram}$ .*

A memory is PRAM iff it allows only PRAM consistent histories.

The following result shows that PRAM memories allow efficient partial replication implementations.

**Theorem 1.** *In a PRAM consistent history, no dependency chain can be created along hoops.*

*Proof.* Let  $x$  be a variable and  $[p_a, \dots, p_b]$  be a  $x$ -hoop. A  $x$ -dependency chain along this hoop is created if  $H$  includes  $w_a(x)v$ ,  $o_b(x)$  and a pattern of operations, at least one for each process of the  $x$ -hoop, implying  $w_a(x)v \mapsto_{pram} o_b(x)$ . But the latter dependency can occur only if point 1 or point 2 of Definition 6 holds. Point 1 is excluded because  $a \neq b$ . Point 2 is possible only if  $o_b(x) = r_b(x)v$  and the dependency  $w_a(x)v \mapsto_{pram} r_b(x)v$  is  $w_a(x)v \mapsto_{ro} r_b(x)v$ , i.e., does not result from the operations performed by the intermediary processes of the hoop.  $\square$

As a consequence of this result, for each variable  $x$ , there is no  $x$ -pertinent process out of  $C(x)$ , and thus, PRAM memories allow efficient partial replication implementations.

Although being weaker than causal memories, Lipton and Sandberg show in [16] that PRAM memories are strong

<sup>4</sup>in [11] this relation is denoted  $\mapsto_{H'}$ .

enough to solve a large number of applications like FFT, matrix product, dynamic programming and more generally the class of oblivious computations<sup>5</sup>. In his PhD, Sinha [17] shows that totally asynchronous iterative methods to find fixed points can converge in Slow memories, which are still weaker than PRAM. In the next section, we illustrate the power of PRAM, together with the usefulness of partial replication, by showing how the Bellman-Ford shortest path algorithm can be solved by using PRAM memory.

## 5. Case study: Bellmann-Ford algorithm

A packet-switching network can be seen as a directed graph,  $G = (V, \Gamma)$ , where each packet-switching node is a vertex in  $V$  and each communication link between node is a pair of parallel edges in  $\Gamma$ , each carrying data in one direction. In such a network, a routing decision is necessary to transmit a packet from a source node to a destination node traversing several links and packet switches. This can be modelled as the problem of finding a path through the graph. Analogously for an Internet or an intranet network. In general, all packet-switching networks and all internets base their routing decision on some form of least-cost criterion, i.e minimize the number of hops that correspond in graph theory to finding the minimum path distance. Most least-cost routing algorithms widespread are a variations of one of the two common algorithms, Dijkstra's algorithm and the Bellman-Ford algorithm[3].

### 5.1. A distributed implementation of the Bellman-Ford algorithm exploiting partial replication

In the following we propose a distributed implementation of the Bellman-Ford algorithm to compute the minimum path from a source node to every other nodes in a system, pointing out the usefulness of partial replication to efficiently distribute the computation. In the following we refer to nodes as processes.

The system (network) is composed by  $N$  processes  $ap_1, \dots, ap_N$  and it is modelled with a graph  $G = (V, \Gamma)$ , where  $V$  is the set of vertex, one for each process in the system and  $\Gamma$  is the set of edges  $(i, j)$  such that  $ap_i, ap_j$  belong to  $V$  and there exists a link between  $i$  and  $j$ .

Let us use the following notation:

- $\Gamma^{-1}(i) = \{j \in V | (i, j) \in \Gamma\}$  is the set of predecessors of process  $ap_i$ ,
- $s$ =source process,

- $w(i, j)$ =link cost from process  $ap_i$  to process  $ap_j$ . In particular:
  - i)  $w(i, i) = 0$ ,
  - ii)  $w(i, j) = \infty$  if the two processes are not directly connected,
  - iii)  $w(i, j) \geq 0$  if the two processes are directly connected;
- $x_i^k$ = cost of the least-cost path from source process  $s$  to process  $n$  under the constraint of no more than  $k$  links traversed.

The centralized algorithm proceeds in steps.

#### 1. [Initialization]

$$x_i^0 = \infty, \forall n \neq s$$

$$x_s^k = 0, \text{ for all } k$$

#### 2. [Update] for each successive $k \geq 0$ :

$$\forall i \neq s, \text{ compute } x_i^{k+1} = \min_{j \in \Gamma^{-1}(i)} [x_j^k + w(j, i)]$$

It is well-known that, if there are no negative cost cycles, the algorithm converge in at most  $N$  steps.

The algorithm is distributively implemented as follows . Without loss of generality, we assume that process  $ap_1$  is the source node. We denote as  $x_i$  the current minimum value from node 1 to node  $i$ . Then, to compute all the minimum path from process  $ap_1$  to every other process in the system, processes cooperate reading and writing the following set of shared variables  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ . Moreover, since the algorithm is iterative, in order to ensure liveness we need to introduce synchronization points in the computation. In particular, we want to ensure that at the beginning of each iteration each process  $ap_i$  reads the new values written by his predecessors  $\Gamma^{-1}(i)$ . Thus each process knows that at most after  $N$  iterations, it has computed the shortest path. With this aim, we introduce the following set of shared variables  $\mathcal{S} = \{k_1, k_2, \dots, k_N\}$ .

Each application process  $ap_i$  only access a subset of shared variables. More precisely,  $ap_i$  accesses  $x_h \in \mathcal{X}$  and  $k_h \in \mathcal{S}$ , such that  $h = i$  or  $ap_h$  is a predecessor of  $ap_i$ .

Since each variable  $x_i$  and  $k_i$  is written only by one process, namely  $ap_i$ , it is simple to notice that the algorithm in Figure 4, correctly runs on a PRAM shared memory. Moreover, since each process has to access only a subset of the shared memory, we can assume a partial replication implementation of such memory. In particular, at each node where the process  $ap_i$  is running to compute the shortest path, there is also a MCS process that ensure Pram consistency in the access to the shared variables.

The algorithm proposed is deadlock-free. In fact, given two processes  $ap_i$  and  $ap_j$  such that  $ap_i$  is a predecessor of  $ap_j$  and viceversa, the corresponding barrier conditions

<sup>5</sup>"A computation is oblivious if its data motion and the operations it executes at a given step are independent of the actual values of data." [16]

```

MINIMUM PATH
1   $k_i := 0;$ 
2  if ( $i == 1$ )
3     $x_i := 0;$ 
4  else  $x_i := \infty;$ 
5  while ( $k_i < N$ ) do
6    while ( $\bigwedge_{h \in \Gamma^{-1}(i)} (k_h < k_i)$ ) do;
7     $x_i := \min([x_j + w(j, i)] \forall j \in \Gamma^{-1}(i));$ 
8     $k_i := k_i + 1$ 

```

Figure 4. pseudocode executed by process  $ap_i$

(line 6 of Figure 4) cannot be satisfied at the same time:  $k_i < k_j$  and  $k_j < k_i$ .

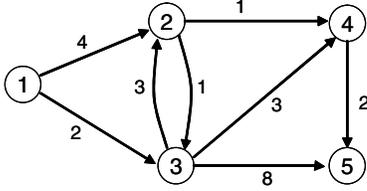


Figure 5. An example

As an example, let us consider the network depicted in Figure 5. We have the following set of processes  $\Pi = \{ap_1, ap_2, ap_3, ap_4, ap_5\}$  and the corresponding variable distribution:

$$\begin{aligned}
\mathcal{X}_1 &= \{x_1, k_1\}, \\
\mathcal{X}_2 &= \{x_1, x_2, x_3, k_1, k_2, k_3\}, \\
\mathcal{X}_3 &= \{x_1, x_2, x_3, k_1, k_2, k_3\}, \\
\mathcal{X}_4 &= \{x_2, x_3, x_4, k_2, k_3, k_4\}, \\
\mathcal{X}_5 &= \{x_3, x_4, x_5, k_3, k_4, k_5\}.
\end{aligned}$$

In Figure 6 we show the pattern of operations generated by each process at the  $k$ -th step of iteration, we only explicit value returned by operations of interest. In reality, in order to point out the sufficiency of PRAM shared memory to ensure the safety and the liveness of the algorithm, we start the scenario showing the two last write operations made by each process at  $(k - 1)$ -th step. In this sense, it must be noticed that the protocol correctly runs if each process reads the values written by each of its neighbors according to their program order.

## 6. Conclusion

This paper has focused on the pertinence of implementing distributed shared memories by using partial replication of variables. It has introduced the notions of share graph

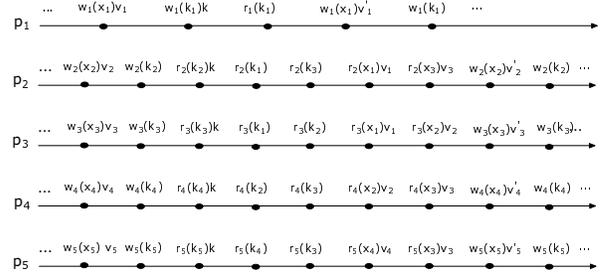


Figure 6. A step of the protocol in Figure 4 for the network in Figure 5

and hoops to model the distribution of variables on the processes, and the notion of dependency chain along hoops to characterize processes that have to transmit information on variables that they don't manage. As a consequence, it has been shown that, in general, distributed shared memories enforcing consistency criteria stronger than causality do not allow efficient implementation based on partial replication. It has also been shown that distributed shared memories enforcing consistency criteria weaker than PRAM are prone to efficient implementation based on partial replication. The power of PRAM memories has been illustrated with the particular example of Bellman-Ford shortest path algorithm.

This paper opens the way for future work. First, the design of an efficient implementation of PRAM memories based on partial replication. Second, on a more theoretical side, the "optimality" of the PRAM consistency criterion, with respect to efficient implementation based on partial replication. In other words, the existence of a consistency criterion stronger than PRAM, and allowing efficient partial replication implementation, remains open. Finally and more subtly, we point out the necessity of understanding and modelling consistency requirements in new distributed system paradigms characterized by large-scale and dynamism properties. In such systems a big amount of processes can dynamically enter and leave the system. Thus traditional consistency guarantees, that is the ones thought

for static systems in which the number of processes  $n$  is fixed and known by each process, seem to be not reasonable.

**Acknowledgement** We like to thank Michel Raynal for suggesting this subject of research and for insightful discussions on this work.

## References

- [1] M. S. A.D. Kshemkalyani. Necessary and sufficient conditions on the information for causal message ordering and their optimal implementation. *Distributed Computing*, 11:91–111, 1988.
- [2] H. Attiya and J. Welch. *Distributed Computing (second edition)*. Wiley, 2004.
- [3] R. E. Bellman. On a routing problem. *Quarterly Applied Mathematics*, XVI(1):87–90, 1958.
- [4] V. C. E. Jimenez, A. Fernández. On the interconnection of causal memory systems. *In: press in Journal of Parallel and Distributed Computing*, 2004.
- [5] A. M. J.M. Hlary. About the efficiency of partial replication to implement distributed shared memory. Technical Report **PI-1727**, IRISA, Campus de Beaulieu, 2005.
- [6] L. Lamport. How to make a multiprocessor computer that correctly executes multiprocess programs. *IEEE Transactions on Computers*, 28(9):690–691, January 1979.
- [7] L. Lamport. On interprocess communication; part i: Basic formalism. *Distributed Computing*, 1(2):77–85, 1986.
- [8] J. B.-P. K. M. Ahamad, G. Neiger and P. Hutto. Causal memory: Definitions, implementation and programming. *Distributed Computing*, 9(1):37–49, 1995.
- [9] M. R. M. Ahamad and G. Thia-Kime. An adaptive architecture for causally consistent distributed services. *Distributed System Engineering*, 6(2):63–70, 1999.
- [10] P. K. M. Ahamad, R.A. Bazy and G. Neiger. The power of processor consistency. *ACM*, 1993.
- [11] A. S. M. Raynal. A suite of formal definitions for consistency criteria in distributed shared memories. Proc. 9-th Int. IEEE Conference on Parallel and Distributed Computing Systems (PDCS96), Dijon, France, pages 125–131, 1996.
- [12] M. A. M. Raynal. Exploiting write semantics in implementing partially replicated causal objects. Proc. 6th Euromicro Conference on Parallel and Distributed Systems, pages 164–175, 1998.
- [13] P. K. M. Ahamad, R. John and G. Neiger. Causal memory meets the consistency and performance needs of distributed application! EW 6: Proceedings of the 6th workshop on ACM SIGOPS European workshop, pages 45–50, 1994.
- [14] S. T.-P. R. Baldoni, A. Milani. Optimal propagation-based protocols implementing causal memories. *Distributed Computing*, 2006.
- [15] S. T.-P. R. Baldoni, C. Spaziani and D. Tulone. An implementation of causal memories using the writing semantics. Proc. 6th Int. Conf. on Principles of Distributed Systems, pages 43–52, 2002.
- [16] J. S. R. Lipton. Pram: a scalable shared memory. *Technical Report CS-TR-180-88*, Princeton University, 1988.
- [17] H. S. Sinha. Mermera: non-coherent distributed shared memory for parallel computing. *Technical Report BU-CS-93-005*, Boston University, 1993.