

# Correctness Proof for a Distributed Memory System

Vicente Cholvi-Juan<sup>1</sup> and José M. Bernabéu-Aubán<sup>2</sup>

<sup>1</sup> University Jaume I  
Castelló (Spain)

<sup>2</sup> Polytechnic University of Valencia  
Valencia (Spain)

**Abstract.** Distributed shared memory systems offer the possibility of implementing parallel algorithms using the familiar operations of shared memory access. In this paper we present the proof of correctness of an algorithm implementing such a shared memory system over a distributed environment. Such proof is based in the formal model for memory systems introduced by the authors in previous works.

## 1 Introduction

A typical (loosely coupled) distributed systems is composed of a collection of independent computers interconnected through some type of network. In order to cooperate, applications written to span several computers on such a system need to have some mechanism to allow each one of their parts to exchange information.

The *shared memory model* (SMM) provides a shared address space which can be used by processes in the same way as local memory, even if they are executed concurrently in different processors. Thus, every process can access any address by means of the basic operations  $data = read(address)$  and  $write(address, data)$ , where  $read$  returns the  $data$  in  $address$ , and  $write$  associates  $data$  with  $address$ .

When a SMM is built on top of a distributed system, we get what is known as a DSMM. In this paper we present the formal proof of correctness of a DSMM algorithm which implements a particular shared memory model. This model is a mixture of the *atomic* memory model [3] and the *sequential* model [4], thus presenting the developer with a familiar memory access semantics, while allowing a more efficient implementation. The proof is carried out using the formal model introduced by the authors in [1]. We base our framework on the I/O automata formalism [5] for discrete event systems. The I/O automata formalism has been successfully used to prove the correctness of algorithms as well as to specify properties of memory systems and transactional systems.

## 2 Correctness Proof

It is well known the inherent complexity of distributed algorithms. Indeed, there are several widely known algorithms that were considered correctly designed and

that later have been found incorrect. Therefore, it is very important to be able to ensure that implementations of models work correctly.

Most of the work on shared models has not paid enough attention to the way such models are specified. The availability of an adequate formal framework makes it possible to answer questions about models and algorithms implementing them. In another work [2], we have introduced such a formal framework. Here, we show, basing on the specifications of two given models, how to we verify the correctness of a proposed DSMM implementing algorithm.

Our framework takes into account two type of operations:  $\text{write}(i, x, u)$  which associates the value  $u$  to the variable  $x$ , and  $\text{read}(i, x, u)$  which informs of the value  $u$  associated with the variable  $x$  (being  $i$  the process that executes the operation). The way these operations are executed are described as actions of that formalism. A  $\text{write}(i, x, u)$  operation can be identified as a pair of  $\text{bwrite}(i, x, u) - \text{ewrite}(i, x)$  events, and a  $\text{read}(i, x, u)$  operation can be identified as a pair of  $\text{bread}(i, x) - \text{eread}(i, x, u)$  events (the first to begin the operation and the second to finish it).

In our formalism to specify a given memory model it is enough to specify the set of executions it allows. Moreover, the sequences of the model that we are going to implement (a mixture of the atomic and sequential) are characterized by the existence of atomics and sequential “views” of those sequences. Formally, from [2], “a sequence  $\alpha$  is *sequential* if  $\text{VIEWS}(\text{SEQ}, \alpha)$  is not an empty set (ditto for the *atomic* sequences)”.

The proposed DSMM implementing algorithm permits to take advantage of the larger degree of parallelism given by the sequential model, allowing at the same time the use of a stricter model (atomic) for a subset of the variables.

The way in which the DSMM implementing algorithm deals with atomic variables consists of invalidating all the replicas previously to each writing. Thus, each variable in isolation will be atomic [3]. Thereby, if we consider that the atomic model is compositional, we can assume (in the verification) that all the variables are sequential without losing generality.

Our solution contains three system’s descriptions:  $\text{OpSpec}$ , the operational requirement specification;  $\text{SysImpl}$ , the high level system implementation; and  $\text{DiscSysImpl}$ , the discrete system implementation.

The operational specification describes the system by using a single automaton. Next, we split it into two automata: one of them models the architectural platform and the other the DSMM implementing algorithm. Finally, we split the later automaton into the automata that model its distributed implementation.

Here, we present only those parts we consider necessary to understand the verification process. The details of the proofs are available in [2]. In what follows,  $\mathcal{P}$  will stand for the set of processes,  $\mathcal{M}$  for the set of variables (being  $\mathcal{M}_0$  the initial values) and  $\text{msig}(\mathcal{P}, \mathcal{M})$  for the set of memory operations.

## 2.1 Operational Specification

The operational specification,  $\text{OpSpec}$ , consists of a single automaton,  $\text{BLACK}(\mathcal{P}, \mathcal{M}, \mathcal{M}_0)$  (Fig. 1), which models all required properties in order to

solve the sequential model. Moreover, it incorporates the fairness condition that guarantees liveness of the implementation. The state transitions are described by specifying the “preconditions” under which each action can occur and the “effect” of each action. This constitutes an abstract representation harder to understand than the axiomatic specification of the sequential model. However, it is easier to use in the proofs.

**Actions:**

- Inputs:  $\text{bread}(i, x)$ ,  $\text{bwrite}(i, x, u)$
  - Outputs:  $\text{eread}(i, x, u)$ ,  $\text{ewrite}(i, x)$
- where  $i \in \mathcal{P}$ ,  $x \in \mathcal{M}$  and  $u \in \text{range}(x)$

**State:** a set of elements,  $S$ , in  $\text{msig}(\mathcal{P}, \mathcal{M})$ , initially empty

**Transitions:**

- $\text{bwrite}(i, x, u)$

Effect:

$$S = S \cdot \text{bwrite}(i, x, u)$$

- $\text{bread}(i, x)$

Effect:

$$S = S \cdot \text{bread}(i, x)$$

**Partition:**  $\{\{\text{ewrite}(i, \cdot)\}, \{\text{eread}(i, \cdot, \cdot)\}\}$

- $\text{ewrite}(i, x)$

Precondition:

$$\text{last}(S \mid i) = \text{bwrite}(i, x, \cdot)$$

Effect:

$$S = S \cdot \text{ewrite}(i, x)$$

- $\text{eread}(i, x, u)$

Precondition:

$$\text{last}(S \mid i) = \text{bread}(i, x)$$

$$\text{VIEW}(\text{SEQ}, T) \neq \emptyset$$

$$\text{(where } T = S \cdot \text{eread}(i, x, u)\text{)}$$

Effect:

$$S = S \cdot \text{eread}(i, x, u)$$

**Fig. 1.** Specification of  $BLACK(\mathcal{P}, \mathcal{M}, \mathcal{M}_0)$ .

## 2.2 System Implementation

To implement a model, memory operations have to be translated into operations of the underlying system. The nature of such operations depends on the architecture of the system on which the model is to be implemented. The architecture gives us only part of the full implementation of the model. The other part is the algorithm which, based on the architecture, implements the model ( $MMS$ ).

In this section we provide the specification of the components that compose the system by using two automata. The first implements the architecture and the other one the  $MMS$ .

We consider a distributed architecture consisting of a set of locations  $\mathcal{D}$  connected by a fault-free, ordered communication channel. We identify the set of variables at each location as  $\mathcal{W}$ , being  $\mathcal{W}_0$  their initial values. Then, the architecture is globally modeled by using a single automaton  $DST(\mathcal{D}, \mathcal{W}, \mathcal{W}_0)$

(Fig. 2). Roughly speaking, each memory operation is modeled by using three actions (the first to start it, the second to “physically perform” it and the later to finish it). Moreover, data at each location is modeled by an array of storing values. Accesses to those data are made through three FIFO operations: *HEAD* returns the queue’s bottom value, *DEQ* dequeues that value and *ENQ* puts the value on the top of the queue.

**Actions:**

- Inputs:  $br(j, y)$ ,  $bw(j, y, v)$
  - Outputs:  $er(j, y, v)$ ,  $ew(j, y)$
  - Internals:  $xr(j, y, v)$ ,  $xw(j, y, v)$
- where  $j \in \mathcal{D}$ ,  $y \in \mathcal{W}$  and  $v \in range(y)$

**State:**

- $F_{in}[j]$ ,  $F_{out}[j]$ : array of FIFO queues of elements in  $in(sig(DST(\mathcal{D}, \mathcal{W}, \mathcal{W}_0)))$  and  $out(sig(DST(\mathcal{D}, \mathcal{W}, \mathcal{W}_0)))$  respectively, initially empties
- $Store[j, y]$ : array of variables in the range of  $y$ , initially  $y_0$

**Transitions:**

- |  |   |  |
|--|---|--|
| <ul style="list-style-type: none"> <li>• <math>bw(j, y, v)</math></li> </ul> <p>Effect:<br/><math>ENQ(F_{in}[j], bw(j, y, v))</math></p> | <ul style="list-style-type: none"> <li>• <math>ew(j, y)</math></li> </ul> <p>Precondition:<br/><math>HEAD(F_{out}[j]) = ew(j, y)</math></p> <p>Effect:<br/><math>DEQ(F_{out}[j])</math></p>       | <ul style="list-style-type: none"> <li>• <math>xw(j, y, v)</math></li> </ul> <p>Precondition:<br/><math>HEAD(F_{in}[j]) = bw(j, y, v)</math></p> <p>Effect:<br/><math>Store[j, y] = v</math><br/><math>DEQ(F_{in}[j])</math><br/><math>ENQ(F_{out}[j], ew(j, y))</math> </p> |
| <ul style="list-style-type: none"> <li>• <math>br(j, y)</math></li> </ul> <p>Effect:<br/><math>ENQ(F_{in}[j], br(j, y))</math></p>       | <ul style="list-style-type: none"> <li>• <math>er(j, y, v)</math></li> </ul> <p>Precondition:<br/><math>HEAD(F_{out}[j]) = er(j, y, v)</math></p> <p>Effect:<br/><math>DEQ(F_{out}[j])</math></p> | <ul style="list-style-type: none"> <li>• <math>xr(j, y, v)</math></li> </ul> <p>Precondition:<br/><math>HEAD(F_{in}[j]) = br(j, y)</math><br/><math>v = Store[j, y]</math></p> <p>Effect:<br/><math>DEQ(F_{in}[j])</math><br/><math>ENQ(F_{out}[j], er(j, y, v))</math></p>  |

**Partition:**  $\{\{ew(j, \cdot)\}, \{er(j, \cdot, \cdot)\}, \{xw(j, \cdot, \cdot)\}, \{xr(j, \cdot, \cdot)\}\}$

**Fig. 2.** Specification of  $DST(\mathcal{D}, \mathcal{W}, \mathcal{W}_0)$ .

The MMS is modeled by using a single automaton  $RED(\mathcal{P}, \mathcal{M})$  (Fig. 3). That continues being a rather abstract representation of the system which is still far away from a realistic representation of the proposed distributed memory system. However, it is intended to be used in the discrete system specification correctness proof.

At this stage it has to be proved that the interaction of the MMS with the architecture provides sequential executions.

**Lemma 1.** *Let  $SysImpl = hide_{\Sigma}(RED(\mathcal{P}, \mathcal{M}) \amalg DST(\mathcal{D}, \mathcal{W}, \mathcal{W}_0))$  (where  $\Sigma = extsig(DST(\mathcal{D}, \mathcal{W}, \mathcal{W}_0))$ ). Then  $SysImpl$  solves  $OpSpec$ .*

**Actions:**

- Inputs:  $\text{bread}(i, x)$ ,  $\text{bwrite}(i, x, u)$ ,  $\text{er}(j, y, v)$ ,  $\text{ew}(j, y)$
  - Outputs:  $\text{eread}(i, x, u)$ ,  $\text{ewrite}(i, x)$ ,  $\text{br}(j, y)$ ,  $\text{bw}(j, y, v)$
- where  $i \in \mathcal{P}$ ,  $x \in \mathcal{M}$ ,  $u \in \text{range}(x)$ ,  $j \in \mathcal{D}$ ,  $y \in \mathcal{W}$  and  $v \in \text{range}(y)$

**State:** a set of elements,  $S$ , in  $\text{sig}(\text{RED}(\mathcal{P}, \mathcal{M}))$ , initially empty

**Transitions:**

- |  |  |
|--|--|
| <ul style="list-style-type: none"> <li>• <math>\text{bwrite}(i, x, u)</math></li> </ul> <p>Effect:<br/><math>S = S \cdot \text{bwrite}(i, x, u)</math></p> <ul style="list-style-type: none"> <li>• <math>\text{bread}(i, x)</math></li> </ul> <p>Effect:<br/><math>S = S \cdot \text{bread}(i, x)</math> <ul style="list-style-type: none"> <li>• <math>\text{bw}(j, y, v)</math></li> </ul> <p>Precondition:<br/><math>\text{last}(S \mid j) = \text{bwrite}(j, y, v)</math></p> <p>Effect:<br/><math>S = S \cdot \text{bw}(j, y, v)</math> <ul style="list-style-type: none"> <li>• <math>\text{br}(j, y)</math></li> </ul> <p>Precondition:<br/><math>\text{last}(S \mid j) = \text{bread}(j, y)</math></p> <p>Effect:<br/><math>S = S \cdot \text{br}(j, y)</math></p> </p></p> | <ul style="list-style-type: none"> <li>• <math>\text{ewrite}(i, x)</math></li> </ul> <p>Precondition:<br/><math>\text{last}(S \mid i) = \text{ew}(i, x)</math></p> <p>Effect:<br/><math>S = S \cdot \text{ewrite}(i, x)</math> <ul style="list-style-type: none"> <li>• <math>\text{eread}(i, x, u)</math></li> </ul> <p>Precondition:<br/><math>\text{last}(S \mid i) = \text{er}(i, x, u)</math><br/><math>\text{VIEW}(\text{SEQ}, T) \neq \emptyset</math><br/>(where <math>T = (S \cdot \text{eread}(i, x, u)) \mid \text{msig}(\mathcal{P}, \mathcal{M})</math>)</p> <p>Effect:<br/><math>S = S \cdot \text{eread}(i, x, u)</math> <ul style="list-style-type: none"> <li>• <math>\text{ew}(j, y)</math></li> </ul> <p>Effect:<br/><math>S = S \cdot \text{ew}(j, y)</math> <ul style="list-style-type: none"> <li>• <math>\text{er}(j, y, v)</math></li> </ul> <p>Effect:<br/><math>S = S \cdot \text{er}(j, y, v)</math></p> </p></p></p> |
|--|--|

**Partition:**  $\{\{\text{eread}(i, \cdot, \cdot)\}, \{\text{ewrite}(i, \cdot)\}, \{\text{br}(j, \cdot)\}, \{\text{bw}(j, \cdot, \cdot)\}\}$

**Fig. 3.** Specification of  $\text{RED}(\mathcal{P}, \mathcal{M})$ .

### 2.3 Discrete System Implementation

In the previous sections we have presented two levels of description of our solution, which, by their very nature, are quite abstract and do not correspond directly with an implementation over a distributed system.

We now proceed to describe the discrete system model of the algorithm. In this model, the distributed nature of the underlying system is made explicit by using FIFO communication channels and explicit message passing between nodes in the distributed system.

The communications system is modeled by an automaton  $\text{FIFO}_M$ , which, for our purposes, is in charge of sending and making receive the messages used by our protocol.

To obtain  $\text{DiscSysImpl}$  we compose the automata that will model the modules responsible of requesting variables as well as providing them ( $\{WRKR_i(\mathcal{P}, \mathcal{M})\}_{i \in \mathcal{P}}$ ) with the module taking care of arbitrating conflicting requests for variables, ordering them and keeping track of circulating variables ( $SEQ(\mathcal{P}, \mathcal{M})$ ) and that which models the *communication channel* module ( $FIFOM$ ).

At this stage it has to be proved that the interaction of the automata that model the discrete system implementation can be safely used for any task for which  $RED(\mathcal{P}, \mathcal{M})$  is satisfactory.

**Lemma 2.**

Let  $\text{DiscSysImpl} = \text{hide}_\Sigma(FIFOM \amalg SEQ(\mathcal{P}, \mathcal{M}) \amalg \{WRKR_i(\mathcal{P}, \mathcal{M})\}_{i \in \mathcal{P}})$  (where  $\Sigma = \text{extsig}(FIFOM)$ ). Then  $\text{DiscSysImpl}$  solves  $RED(\mathcal{P}, \mathcal{M})$ .

### 3 Conclusion

In this paper we have shown how to model, using our formalization of memory coherency models, a particular algorithm implementing a distributed shared memory system. First we modeled the problem using an I/O automata. The possible execution sequences of this automata define the “problem” we claim our algorithm “solves”. Then we propose a more detailed I/O automaton which is shown to solve the original problem. This automaton, while still abstract, and not making use of the distributed nature of the underlying system, contains some details which facilitate its further break down into a fully distributed specification.

The method used in this paper would be appropriate in general to prove the correctness of any other MMS implementation. In our approach we use the full power of the I/O automata formalism. A similar approach can also be used to prove the correctness of algorithms built upon any particular memory model.

### References

1. Bernabéu-Aubán, J.M., Cholvi-Juan, V.: Formalizing memory coherency models. *Journal of Computing and Information* **1** (1994) 653–672
2. Cholvi-Juan, V.: Formalizing memory models. PhD thesis, Department of Computer Science, Polytechnic University of Valencia (1994)
3. Herlihy, M.P., Wing, J.M.: Linearizability: A correctness condition for concurrent objects **12** 3 (1990) 463–492
4. Lamport, L.: How to make a multiprocessor computer that correctly executes multiprocess programs, *IEEE Transactions on Computers* **28** 9 (1979) 690–691
5. Lynch, N.: I/O Automata: A model for discrete event system. Massachusetts Institute of Technology, Technical Report MIT/LCS/TM-351 (1988)