

# Global State Estimates for Distributed Systems

Gabriel Kalyon<sup>1</sup>, Tristan Le Gall<sup>2</sup>, Hervé Marchand<sup>3</sup> and Thierry Massart<sup>1\*</sup>

<sup>1</sup> Université Libre de Bruxelles (U.L.B.), Campus de la Plaine, Bruxelles, Belgique

<sup>2</sup> CEA LIST, LMeASI, boîte 94, 91141 Gif-sur-Yvette, France

<sup>3</sup> INRIA, Rennes - Bretagne Atlantique, France

**Abstract.** We consider distributed systems modeled as *communicating finite state machines* with reliable unbounded FIFO channels. As an essential subroutine for control, monitoring and diagnosis applications, we provide an algorithm that computes, during the execution of the system, an estimate of the current global state of the distributed system for each local subsystem. This algorithm does not change the behavior of the system; each subsystem only computes and records a symbolic representation of the state estimates, and piggybacks some extra information to the messages sent to the other subsystems in order to refine their estimates. Our algorithm relies on the computation of reachable states. Since the reachability problem is undecidable in our model, we use abstract interpretation techniques to obtain regular overapproximations of the possible FIFO channel contents, and hence of the possible current global states. An implementation of this algorithm provides an empirical evaluation of our method.

## 1 Introduction

During the execution of a computer system, the knowledge of its global state may be crucial information, for instance to control which action can or must be done, to monitor its behavior or perform some diagnostic. Distributed systems, are generally divided into two classes, depending on whether the communication between subsystems is *synchronous* or not. When the synchrony hypothesis [1] can be made, each local subsystem can easily know, at each step of the execution, the global state of the system (assuming that there is no internal action). When considering *asynchronous* distributed systems, this knowledge is in general impossible, since the communication delays between the components of the system must be taken into account. Therefore, each local subsystem can *a priori* not immediately know either the local state of the other subsystems or the messages that are currently in transfer. In this paper, we consider the asynchronous framework where a system is composed of  $n$  subsystems that asynchronously communicate through reliable *unbounded* FIFO channels and modeled by *communicating finite state machines* (CFSM) [3]. This model appears to be essential for concurrent systems in which components cooperate via asynchronous message passing through unbounded buffers (they are e.g. widely used to model communication protocols). We thus assume that the distributed system is already built and the architecture of communication between the different subsystems is fixed. Our aim is to provide an algorithm that allows

---

\* This work has been done in the MoVES project (P6/39), part of the IAP-Phase VI Interuniversity Attraction Poles Programme funded by the Belgian State, Belgian Science Policy.

us to compute, in each subsystem of a distributed system  $\mathcal{T}$ , an estimate of the current state of  $\mathcal{T}$ . More precisely, each subsystem or a local associated *estimator* computes a set of possible global states, including the contents of the channels, in which the system  $\mathcal{T}$  can be; it can be seen as a particular case of monitoring with partial observation. We assume that the subsystems (or associated estimators) can record their own state estimate and that some extra information can be piggybacked to the messages normally exchanged by the subsystems. Without this additional information, since a local subsystem cannot observe the other subsystems nor the FIFO channel contents, the computed state estimates might be too rough. Our computation is based on the use of the reachability operator, which cannot always be done in the CFSM model for undecidability reasons. Therefore, we rely on the abstract interpretation techniques we presented previously in [13]. They ensure the termination of the computations by overapproximating in a symbolic way the possible FIFO channel contents (and hence the state estimates) by *regular languages*. Computing state estimates is useful in many applications. For example, this information can be used to control the system in order to prevent it from reaching some given forbidden global states [4], or to perform some diagnosis to detect some faults in the system [7, 19]. For these two potential applications, a more precise state estimate allows the controller or the diagnoser to take better decisions.

This problem differs from the synthesis problem (see e.g. [15, 8, 5]) which consists in synthesizing a distributed system (together with its architecture of communication) equivalent to a given specification. It also differs from the methodology described in [9] where the problem is to infer from a distributed observation of a distributed system (modeled by a High Level Message Sequence Chart) the set of sequences that explains this observation. It is also different from *model checking* techniques [2, 10] that proceed to a symbolic exploration of all the possible states of the system, without running it. We however use the same symbolic representation of queue contents as in [2, 10]. In [21], Kumar and Xu propose a distributed algorithm which computes an estimate of the current state of a system. Local estimators maintain and update local state estimates from their own observation of the system and information received from the other estimators. In their framework, the local estimators communicate between them through reliable FIFO channels with delays, whereas the system is monolithic and therefore in their case, a global state is simpler than for our distributed systems composed of several subsystems together with communicating FIFO channels. In [20], Tripakis studies the decidability of the existence of controllers such that a set of responsiveness properties is satisfied in a decentralized framework with communication delays between the controllers. This problem is undecidable when there is no communication or when the communication delays are unbounded. He conjectures that the problem is decidable when the communication delays are bounded. See [18, 14] for other works dealing with communication (with or without delay) between agents.

Below, in section 2, we define the formalism of *communicating finite state machines*, that we use. We formally define, in section 3, the state estimate mechanisms and the notion of state estimators. In section 4, we provide an algorithm to compute an estimate of the current state of a distributed system and prove its correctness. We explain, in section 5, how the termination of this algorithm is ensured by using abstract interpretation techniques. Section 6 gives some experimental results. Proofs can be found in [11].

## 2 Communicating Finite State Machines as a Model of the System

We model a distributed system by *communicating finite state machines* [3] which use reliable unbounded FIFO channels (also called *queues*) to communicate. A *global state* in this model is given by the local state of each subsystem together with the content of each FIFO queue. As no bound is given either in the transmission delay, or on the length of the queues, the state space of the system is *a priori* infinite.

**Definition 1 (Communicating Finite State Machines).** A communicating finite state machine (CFSM)  $\mathcal{T}$  is defined as a 6-tuple  $\langle L, \ell_0, Q, M, \Sigma, \Delta \rangle$ , where (i)  $L$  is a finite set of locations, (ii)  $\ell_0 \in L$  is the initial location, (iii)  $Q$  is a set of queues that  $\mathcal{T}$  can use, (iv)  $M$  is a finite set of messages, (v)  $\Sigma \subseteq Q \times \{!, ?\} \times M$  is a finite set of actions, that are either an output  $a!m$  to specify that the message  $m \in M$  is written on the queue  $a \in Q$  or an input  $a?m$  to specify that the message  $m \in M$  is read on the queue  $a \in Q$ , (vi)  $\Delta \subseteq L \times \Sigma \times L$  is a finite set of transitions.

A transition  $\langle \ell, i!m, \ell' \rangle$  indicates that when the system moves from the  $\ell$  to  $\ell'$ , a message  $m$  is added at the end of the queue  $i$ .  $\langle \ell, i?m, \ell' \rangle$  indicates that, when the system moves from  $\ell$  to  $\ell'$ , a message  $m$  must be present at the beginning of the queue  $i$  and is removed from this queue. Moreover, throughout this paper, we assume that  $\mathcal{T}$  is deterministic, meaning that for all  $\ell \in L$  and  $\sigma \in \Sigma$ , there exists at most one location  $\ell' \in L$  such that  $\langle \ell, \sigma, \ell' \rangle \in \Delta$ . For  $\sigma \in \Sigma$ ,  $\text{Trans}(\sigma)$  denotes the set of transitions of  $\mathcal{T}$  labeled by  $\sigma$ . The occurrence of a transition will be called an *event* and given an event  $e$ ,  $\delta_e$  denotes the corresponding transition. The semantics of a CFSM is defined as follows:

**Definition 2.** The semantics of a CFSM  $\mathcal{T} = \langle L, \ell_0, Q, M, \Sigma, \Delta \rangle$  is given by an infinite Labeled Transition System (LTS)  $\llbracket \mathcal{T} \rrbracket = \langle X, \mathbf{x}_0, \Sigma, \rightarrow \rangle$ , where (i)  $X \stackrel{\text{def}}{=} L \times (M^*)^{|Q|}$  is the set of states, (ii)  $\mathbf{x}_0 \stackrel{\text{def}}{=} \langle \ell_0, \epsilon, \dots, \epsilon \rangle$  is the initial state, (iii)  $\Sigma$  is the set of actions, and (iv)  $\rightarrow \stackrel{\text{def}}{=} \bigcup_{\delta \in \Delta} \delta \subseteq X \times \Sigma \times X$  is the transition relation where  $\delta \rightarrow$  is defined by:

$$\frac{\delta = \langle \ell, i!m, \ell' \rangle \in \Delta \quad w'_i = w_i \cdot m}{\langle \ell, w_1, \dots, w_i, \dots, w_{|Q|} \rangle \xrightarrow{\delta} \langle \ell', w_1, \dots, w'_i, \dots, w_{|Q|} \rangle}$$

$$\frac{\delta = \langle \ell, i?m, \ell' \rangle \in \Delta \quad w_i = m \cdot w'_i}{\langle \ell, w_1, \dots, w_i, \dots, w_{|Q|} \rangle \xrightarrow{\delta} \langle \ell', w_1, \dots, w'_i, \dots, w_{|Q|} \rangle}$$

A global state of a CFSM  $\mathcal{T}$  is thus a tuple  $\langle \ell, w_1, \dots, w_{|Q|} \rangle \in X = L \times (M^*)^{|Q|}$  where  $\ell$  is the current location of  $\mathcal{T}$  and  $w_1, \dots, w_{|Q|}$  are finite words on  $M^*$  which give the content of the queues in  $Q$ . At the beginning, all queues are empty, so the initial state is  $\mathbf{x}_0 = \langle \ell_0, \epsilon, \dots, \epsilon \rangle$ . Given a CFSM  $\mathcal{T}$ , two states  $\mathbf{x}, \mathbf{x}' \in X$  and an event  $e$ , to simplify the notations we sometimes denote  $\mathbf{x} \xrightarrow{\delta_e} \mathbf{x}'$  by  $\mathbf{x} \xrightarrow{e} \mathbf{x}'$ . An *execution* of  $\mathcal{T}$  is a sequence  $\mathbf{x}_0 \xrightarrow{e_1} \mathbf{x}_1 \xrightarrow{e_2} \dots \xrightarrow{e_m} \mathbf{x}_m$  where  $\mathbf{x}_i \xrightarrow{e_{i+1}} \mathbf{x}_{i+1} \in \rightarrow \forall i \in [0, m-1]$ . Given a set of states  $Y \subseteq X$ ,  $\text{Reach}_{\Delta'}^{\mathcal{T}}(Y)$  corresponds to the set of states that are reachable in  $\llbracket \mathcal{T} \rrbracket$  from  $Y$  only firing transitions of  $\Delta' \subseteq \Delta$  in  $\mathcal{T}$ . It is defined by  $\text{Reach}_{\Delta'}^{\mathcal{T}}(Y) \stackrel{\text{def}}{=} \bigcup_{n \geq 0} (\text{Post}_{\Delta'}^{\mathcal{T}}(Y))^n$  where  $(\text{Post}_{\Delta'}^{\mathcal{T}}(Y))^n$  is the  $n^{\text{th}}$  functional power of  $\text{Post}_{\Delta'}^{\mathcal{T}}(Y)$ ,

defined by:  $\text{Post}_{\Delta'}^{\mathcal{T}}(Y) \stackrel{\text{def}}{=} \{x' \in X \mid \exists x \in Y, \exists \delta \in \Delta' : x \xrightarrow{\delta} x'\}$ . Although there is no general algorithm that can exactly compute the reachability set in our setting [3], there exist some techniques that allow us to compute an overapproximation of this set (see section 5). Given a sequence of actions  $\bar{\sigma} = \sigma_1 \cdots \sigma_m \in \Sigma^*$  and two states  $x, x' \in X$ ,  $x \xrightarrow{\bar{\sigma}} x'$  denotes that the state  $x'$  is reachable from  $x$  by executing  $\bar{\sigma}$ .

**Asynchronous Product.** A distributed system  $\mathcal{T}$  is generally composed of several subsystems  $\mathcal{T}_i$  ( $\forall i \in [1, n]$ ) acting in parallel. In fact,  $\mathcal{T}$  is defined by a CFSM resulting from the asynchronous (interleaved) product of the  $n$  subsystems  $\mathcal{T}_i$ , also modeled by CFSMs. This can be defined through the asynchronous product of two subsystems.

**Definition 3.** Given 2 CFSMs  $\mathcal{T}_i = \langle L_i, \ell_{0,i}, Q_i, M_i, \Sigma_i, \Delta_i \rangle$  ( $i = 1, 2$ ), their asynchronous product, denoted by  $\mathcal{T}_1 \parallel \mathcal{T}_2$ , is defined by a CFSM  $\mathcal{T} = \langle L, \ell_0, Q, M, \Sigma, \Delta \rangle$ , where  $L \stackrel{\text{def}}{=} L_1 \times L_2$ ,  $\ell_0 \stackrel{\text{def}}{=} \ell_{0,1} \times \ell_{0,2}$ ,  $Q \stackrel{\text{def}}{=} Q_1 \cup Q_2$ ,  $M \stackrel{\text{def}}{=} M_1 \cup M_2$ ,  $\Sigma \stackrel{\text{def}}{=} \Sigma_1 \cup \Sigma_2$ , and  $\Delta \stackrel{\text{def}}{=} \{ \langle \langle \ell_1, \ell_2 \rangle, \sigma_1, \langle \ell'_1, \ell'_2 \rangle \rangle \mid (\langle \ell_1, \sigma_1, \ell'_1 \rangle \in \Delta_1) \wedge (\ell_2 \in L_2) \} \cup \{ \langle \langle \ell_1, \ell_2 \rangle, \sigma_2, \langle \ell'_1, \ell'_2 \rangle \rangle \mid (\langle \ell_2, \sigma_2, \ell'_2 \rangle \in \Delta_2) \wedge (\ell_1 \in L_1) \}$ .

Note that in the previous definition,  $Q_1$  and  $Q_2$  are not necessarily disjoint; this allows the subsystems to communicate between them via common queues. Composing the various subsystems  $\mathcal{T}_i$  ( $\forall i \in [1, n]$ ) two-by-two in any order gives the global distributed system  $\mathcal{T}$  whose semantics (up to state isomorphism) does not depend on the order.

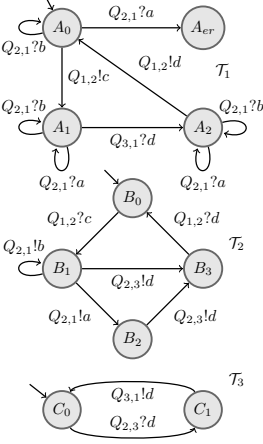
**Definition 4 (Distributed system).** A distributed system  $\mathcal{T} = \langle L, \ell_0, Q, M, \Sigma, \Delta \rangle$  is defined by the asynchronous product of  $n$  CFSMs  $\mathcal{T}_i = \langle L_i, \ell_{0,i}, Q_i, M, \Sigma_i, \Delta_i \rangle$  ( $\forall i \in [1, n]$ ) acting in parallel and exchanging information through FIFO channels.

Note that a distributed system is also modeled by a CFSM, since the asynchronous product of several CFSMs is a CFSM. In the sequel, a CFSM  $\mathcal{T}_i$  always denotes the model of a single process, and a distributed system  $\mathcal{T} = \langle L, \ell_0, Q, M, \Sigma, \Delta \rangle$  always denotes the model of the global system, as in Definition 4. Below, unless stated explicitly,  $\mathcal{T} = \mathcal{T}_1 \parallel \dots \parallel \mathcal{T}_n$  is the considered distributed system.

**Communication Architecture of the System.** We consider an architecture for the system  $\mathcal{T} = \mathcal{T}_1 \parallel \dots \parallel \mathcal{T}_n$  defined in Definition 4 with *point-to-point* communication i.e., any subsystem  $\mathcal{T}_i$  can send messages to any other subsystem  $\mathcal{T}_j$  through a queue  $Q_{i,j}$ . Thus, only  $\mathcal{T}_i$  can write a message  $m$  on  $Q_{i,j}$  (this is denoted by  $Q_{i,j}!m$ ) and only  $\mathcal{T}_j$  can read a message  $m$  on this queue (this is denoted by  $Q_{i,j}?m$ ). Moreover, we suppose that the queues are unbounded, that the message transfers between the subsystems are reliable and may suffer from arbitrary non-zero delays, and that no *global clock* or *perfectly synchronized local clocks* are available. With this architecture, the set  $Q_i$  of  $\mathcal{T}_i$  ( $\forall i \in [1, n]$ ) can be rewritten as  $Q_i = \{Q_{i,j}, Q_{j,i} \mid (1 \leq j \leq n) \wedge (j \neq i)\}$  and  $\forall j \neq i \in [1, n]$ ,  $\Sigma_i \cap \Sigma_j = \emptyset$ . Let  $\delta_i = \langle \ell_i, \sigma_i, \ell'_i \rangle \in \Delta_i$  be a transition of  $\mathcal{T}_i$ ,  $\text{global}(\delta_i) \stackrel{\text{def}}{=} \{ \langle \langle \ell_1, \dots, \ell_{i-1}, \ell_i, \ell_{i+1}, \dots, \ell_n \rangle, \sigma_i, \langle \ell'_1, \dots, \ell'_{i-1}, \ell'_i, \ell'_{i+1}, \dots, \ell'_n \rangle \rangle \in \Delta \mid \forall j \neq i \in [1, n] : \ell_j \in L_j \}$  is the set of transitions of  $\Delta$  that can be built from  $\delta_i$  in  $\mathcal{T}$ . We extend this definition to sets of transitions  $D \subseteq \Delta_i$  of the subsystem  $\mathcal{T}_i$ :  $\text{global}(D) \stackrel{\text{def}}{=} \bigcup_{\delta_i \in D} \text{global}(\delta_i)$ . We abuse notation and write  $\Delta \setminus \Delta_i$  instead of  $\Delta \setminus \text{global}(\Delta_i)$  to denote the set of transitions of  $\Delta$  that are not built from  $\Delta_i$ . Given the set  $\Sigma_i$  of  $\mathcal{T}_i$  ( $\forall i \in [1, n]$ ) and the set  $\Sigma$  of  $\mathcal{T}$ , the projection  $P_i$  of  $\Sigma$  onto

$\Sigma_i$  is standard:  $P_i(\varepsilon) = \varepsilon$  and  $\forall w \in \Sigma^*, \forall a \in \Sigma, P_i(wa) = P_i(w)a$  if  $a \in \Sigma_i$ , and  $P_i(w)$  otherwise. The inverse projection  $P_i^{-1}$  is defined, for each  $L \subseteq \Sigma_i^*$ , by  $P_i^{-1}(L) = \{w \in \Sigma^* \mid P_i(w) \in L\}$ .

*Example 1.* Let us illustrate the concepts of distributed system and CFSM with our running example depicted on the right hand side. It models a factory composed of three components  $\mathcal{T}_1$ ,  $\mathcal{T}_2$  and  $\mathcal{T}_3$ . The subsystem  $\mathcal{T}_2$  produces two kinds of items,  $a$  and  $b$ , and sends these items to  $\mathcal{T}_1$  to finish the job. At reception,  $\mathcal{T}_1$  must immediately terminate the process of each received item.  $\mathcal{T}_1$  can receive and process  $b$  items at any time, but must be in a *turbo mode* to receive and process  $a$  items. The subsystem  $\mathcal{T}_1$  can therefore be in *normal mode* modeled by the location  $A_0$  or in *turbo mode* (locations  $A_1$  and  $A_2$ ). In normal mode, if  $\mathcal{T}_1$  receives an item  $a$ , an error occurs (transition in location  $A_{er}$ ). Since  $\mathcal{T}_1$  cannot always be in turbo mode, a protocol between  $\mathcal{T}_1$  and  $\mathcal{T}_2$  is imagined. At the beginning,  $\mathcal{T}_1$  informs (*connect* action, modeled by  $\overset{Q_{1,2}!c}{\rightarrow}$ )  $\mathcal{T}_2$  that it goes in a turbo mode, then  $\mathcal{T}_2$  sends  $a$  and  $b$  items. At the end of a working session,  $\mathcal{T}_2$  informs  $\mathcal{T}_1$  (*disconnect* action, modeled by  $\overset{Q_{2,3}!d}{\rightarrow}$ ) that it has completed its session, so that  $\mathcal{T}_1$  can go back in normal mode. This information has to transit through  $\mathcal{T}_3$  via queues  $Q_{2,3}$  and  $Q_{3,1}$ , as  $\mathcal{T}_3$  must also record this end of session. Since  $d$  can be transmitted faster than some items  $a$  and  $b$ , one can easily find a scenario where  $\mathcal{T}_1$  decides to go back to  $A_0$  and ends up in the  $A_{er}$  location by reading the message  $a$ . Indeed, as  $\mathcal{T}_1$  cannot observe the content of the queues, it does not know whether there is a message  $a$  in queue  $Q_{2,1}$  when it arrives in  $A_0$ . This motivates the interest of computing good state estimates of the current state of the system. If each subsystem maintains good estimates of the current state of the system, then  $\mathcal{T}_1$  can know whether there is a message  $a$  in  $Q_{2,1}$ , and reach  $A_0$  only if it is not the case.



### 3 State Estimates of Distributed Systems

We introduce here the framework and the problem we are interested in.

**Local View of the Global System.** A global state of  $\mathcal{T} = \mathcal{T}_1 \parallel \dots \parallel \mathcal{T}_n$  is given by a tuple of locations (one for each subsystem) and the content of all the FIFO queues. Informally our problem consists in defining one local estimator per subsystem, knowing that each of them can only observe the occurrences of actions of its own local subsystem, such that these estimators compute *online* (i.e., during the execution of the system) estimates of the global state of  $\mathcal{T}$ . We assume that each local *estimator*  $\mathcal{E}_i$  has a precise observation of subsystem  $\mathcal{T}_i$ , and that the model of the global system is known by all the estimators (i.e., the structure of each subsystem and the architecture of the queues between them). Each estimator  $\mathcal{E}_i$  must determine online the *smallest possible set* of global states  $E_i$  that contains the actual current global state. Note that if  $\mathcal{E}_i$  observes that the location of  $\mathcal{T}_i$  is  $\ell_i$ , a very rough state estimate is  $L_1 \times \dots \times \{\ell_i\} \times \dots \times L_n \times (M^*)^{|Q|}$ . In other words all the global states of the system such that location of  $\mathcal{T}_i$  is  $\ell_i$ ; however, this rough estimate does not provide a very useful information.

**Online State Estimates.** The estimators must compute the state estimates online. Since each estimator  $\mathcal{E}_i$  is local to its subsystem, we suppose that  $\mathcal{E}_i$  synchronously observes the actions fired by its subsystem; hence since each subsystem is deterministic, each time an event occurs in the local subsystem, it can immediately infer the new location of  $\mathcal{T}_i$  and use this information to define its new state estimate. In order to have better state estimates, we also assume that the estimators can communicate with each other by adding some information (some timestamps and their state estimates) to the messages exchanged by the subsystems. Notice that, due to the communication delay, the estimators cannot communicate synchronously, and therefore the state estimate attached to a message might be out-of-date. A classical way to reduce this uncertainty is to timestamp the messages, e.g., by means of vector clocks (see section 4.1).

**Estimates Based on Reachability Sets.** Each local estimator maintains a symbolic representation of all global states of the distributed system that are compatible with its observation and with the information it received previously from the other estimators. In section 4.2, we detail the algorithms which update these symbolic representations whenever an event occurs. But first, let us explain the intuition behind the computation of an estimate. We consider the simplest case: the initial state estimate before the system begins its execution. Each FIFO channel is empty, and each subsystem  $\mathcal{T}_i$  is in its initial location  $\ell_{i,0}$ . So the initial global state is known by every estimator  $\mathcal{E}_i$ . A subsystem  $\mathcal{T}_j$  may however start its execution, while  $\mathcal{T}_i$  is still in its initial location, and therefore  $\mathcal{E}_i$  must thus take into account all the global states that are reachable by taking the transitions of the other subsystems  $\mathcal{T}_j$ . The initial estimate  $E_i$  is this set of reachable global states. This computation of reachable global states also occurs in the update algorithms which take into account any new local event occurred or message received (see section 4.2). The reachability problem is however undecidable for distributed FIFO systems. In section 5, we explain how we overcome this obstacle by using abstract interpretation techniques.

**Properties of the Estimators.** Estimators may have two important properties: soundness and completeness. Completeness refers to the fact that the current state of the global system is always included in the state estimates computed by each state estimator. Soundness refers to the fact that all states included in the state estimate of  $\mathcal{E}_i$  ( $\forall i \in [1, n]$ ) can be reached by one of the sequences of actions that are compatible with the observation of  $\mathcal{T}_i$  performed by  $\mathcal{E}_i$ .

**Definition 5 (Completeness and Soundness).** *The estimators  $(\mathcal{E}_i)_{i \leq n}$  are (i) complete if and only if, for any execution  $\mathbf{x}_0 \xrightarrow{e_1} \mathbf{x}_1 \xrightarrow{e_2} \dots \xrightarrow{e_m} \mathbf{x}_m$  of  $\mathcal{T}$ ,  $\mathbf{x}_m \in \bigcap_{i=1}^n E_i$ , and (ii) sound if and only if, for any execution  $\mathbf{x}_0 \xrightarrow{e_1} \mathbf{x}_1 \xrightarrow{e_2} \dots \xrightarrow{e_m} \mathbf{x}_m$  of  $\mathcal{T}$ ,  $E_i \subseteq \{x' \in X \mid \exists \bar{\sigma} \in P_i^{-1}(P_i(\sigma_{e_1} \cdot \sigma_{e_2} \dots \sigma_{e_m})) : \mathbf{x}_0 \xrightarrow{\bar{\sigma}} x'\}$  ( $\forall i \leq n$ ) where  $\sigma_{e_k}$  ( $\forall k \in [1, m]$ ) is the action that labels the transition corresponding to  $e_k$ .*

## 4 Algorithm to Compute the State Estimates

We now present our algorithm that computes estimates of the current state of a distributed system. But first, we recall the notion of *vector clocks* [12], a standard concept that we shall use to compute a more precise state estimates.

## 4.1 Vector Clocks

To allow the estimators to have a better understanding of the concurrent execution of the distributed system, it is important to determine the causal and temporal relationship between the events that occur in its execution. In a distributed system, events emitted by the same process are ordered, while events emitted by different processes are generally not. When the concurrent processes communicate, additional ordering information can however be obtained. In this case, the communication scheme can be used to obtain a partial order on the events of the system. In practice, vectors of logical clocks, called *vector clocks* [12], can be used to time-stamp the events of the distributed system. The order of two events can then be determined by comparing the value of their respective vector clocks. When these vector clocks are incomparable, the exact order in which the events occur cannot be determined. Vector clocks are formally defined as follows:

**Definition 6 (Vector Clocks).** Let  $\langle D, \sqsubseteq \rangle$  be a partially ordered set, a vector clock mapping of width  $n$  is a function  $V : D \rightarrow \mathbb{N}^n$  such that  $\forall d_1, d_2 \in D : (d_1 \sqsubseteq d_2) \Leftrightarrow (V(d_1) \leq V(d_2))$ .

In general, for a distributed system composed of  $n$  subsystems, the partial order on events is represented by a vector clock mapping of width  $n$ . The method for computing this vector clock mapping depends on the communication scheme of the distributed system. For CFSMs, this vector clock mapping can be computed by the Mattern's algorithm [16], which is based on the causal and thus temporal relationship between the sending and reception of any message transferred through any FIFO channel. This information is then used to determine a partial order, called *causality (or happened-before) relation*  $\prec_c$ , on the events of the distributed system. This relation is the smallest transitive relation satisfying the following conditions: (i) if the events  $e_i \neq e_j$  occur in the same subsystem  $\mathcal{T}_i$  and if  $e_i$  comes before  $e_j$  in the execution, then  $e_i \prec_c e_j$ , and (ii) if  $e_i$  is an output event occurring in  $\mathcal{T}_i$  and if  $e_j$  is the corresponding input event occurring in  $\mathcal{T}_j$ , then  $e_i \prec_c e_j$ . In Mattern's algorithm [16], each process  $\mathcal{T}_i$  ( $\forall i \in [1, n]$ ) has a vector clock  $V_i \in \mathbb{N}^n$  of width  $n$  and each element  $V_i[j]$  ( $\forall j \in [1, n]$ ) is a counter which represents the knowledge of  $\mathcal{T}_i$  regarding  $\mathcal{T}_j$  and which means that  $\mathcal{T}_i$  knows that  $\mathcal{T}_j$  has executed at least  $V_i[j]$  events. Each time an event occurs in a subsystem  $\mathcal{T}_i$ , the vector clock  $V_i$  is updated to take into account the occurrence of this event (see [16] for details). When  $\mathcal{T}_i$  sends a message to some subsystem  $\mathcal{T}_j$ , this vector clock is piggybacked and allows  $\mathcal{T}_j$ , after reception, to update its own vector clock. Our state estimate algorithm uses vector clocks and follows Mattern's algorithm, which ensures the correctness of the vector clocks that we use (see section 4.2).

## 4.2 Computation of State Estimates

Our state estimate algorithm computes, for each estimator  $\mathcal{E}_i$  and for each event occurring in the subsystem  $\mathcal{T}_i$ , a vector clock  $V_i$  and a state estimate  $E_i$  that contains the current state of  $\mathcal{T}$  and any future state that can be reached from this current state by firing actions that do not belong to  $\mathcal{T}_i$ . This computation obviously depends on the information that  $\mathcal{E}_i$  receives. As a reminder,  $\mathcal{E}_i$  observes the last action fired by  $\mathcal{T}_i$  and can infer the fired transition.  $\mathcal{T}_i$  also receives from the other estimators  $\mathcal{E}_j$  their state estimate  $E_j$  and their vector clock  $V_j$ . Our state estimate algorithm proceeds as follows :

---

**Algorithm 1:** initialization( $\mathcal{T}$ )

---

**input** :  $\mathcal{T} = \mathcal{T}_1 || \dots || \mathcal{T}_n$  .  
**output**: The initial state estimate  $E_i$  of the estimator  $\mathcal{E}_i$  ( $\forall i \in [1, n]$ ).  
**1 begin**  
**2** | **for**  $i \leftarrow 1$  **to**  $n$  **do** **for**  $j \leftarrow 1$  **to**  $n$  **do**  $V_i[j] \leftarrow 0$   
**3** | **for**  $i \leftarrow 1$  **to**  $n$  **do**  $E_i \leftarrow \text{Reach}_{\Delta \setminus \Delta_i}^{\mathcal{T}}(\langle \ell_{0,1}, \dots, \ell_{0,n}, \epsilon, \dots, \epsilon \rangle)$   
**4 end**

---

- When the subsystem  $\mathcal{T}_i$  sends a message  $m$  to  $\mathcal{T}_j$ ,  $\mathcal{T}_i$  attaches the vector clock  $V_i$  and the state estimate  $E_i$  of  $\mathcal{E}_i$  to this message. Next,  $\mathcal{E}_i$  receives the action fired by  $\mathcal{T}_i$ , and infers the fired transition. It then uses this information to update its state estimate  $E_i$ .
- When the subsystem  $\mathcal{T}_i$  receives a message  $m$  from  $\mathcal{T}_j$ ,  $\mathcal{E}_i$  receives the action fired by  $\mathcal{T}_i$  and the information sent by  $\mathcal{T}_j$  i.e., the state estimate  $E_j$  and the vector clock  $V_j$  of  $\mathcal{E}_j$ . It computes its new state estimate from these elements.

In both cases, the computation of the new state estimate  $E_i$  depends on the computation of reachable states. In this section, we assume that we have an operator that can compute an *approximation* of the reachable states (which is *undecidable* in the CFMS model). We will explain in section 5 how such an operator can be computed effectively.

**State Estimate Algorithm.** Our algorithm, called *SE-algorithm*, computes estimates of the current state of a distributed system. It is composed of three sub-algorithms: (i) the initialization algorithm, which is only used when the system starts its execution, computes, for each estimator, its initial state estimate (ii) the outputTransition algorithm computes online the new state estimate of  $\mathcal{E}_i$  after an output of  $\mathcal{T}_i$ , and (iii) the inputTransition algorithm computes online the new state estimate of  $\mathcal{E}_i$  after an input of  $\mathcal{T}_i$ .

*INITIALIZATION Algorithm:* According to the Mattern’s algorithm [16], each component of the vector  $V_i$  is set to 0. To take into account that, before the execution of the first action of  $\mathcal{T}_i$ , the other subsystems  $\mathcal{T}_j$  ( $\forall j \neq i \in [1, n]$ ) could perform inputs and outputs, the initial state estimate of  $\mathcal{E}_i$  is given by  $E_i = \text{Reach}_{\Delta \setminus \Delta_i}^{\mathcal{T}}(\langle \ell_{0,1}, \dots, \ell_{0,n}, \epsilon, \dots, \epsilon \rangle)$ .

---

**Algorithm 2:** outputTransition( $\mathcal{T}, V_i, E_i, \delta$ )

---

**input** :  $\mathcal{T} = \mathcal{T}_1 || \dots || \mathcal{T}_n$ , the vector clock  $V_i$  of  $\mathcal{E}_i$ , the current state estimate  $E_i$  of  $\mathcal{E}_i$ , and a transition  $\delta = \langle \ell_1, Q_{i,j}!m, \ell_2 \rangle \in \Delta_i$ .  
**output**: The state estimate  $E_i$  after the output transition  $\delta$ .  
**1 begin**  
**2** |  $V_i[i] \leftarrow V_i[i] + 1$   
**3** |  $\mathcal{T}_i$  tags message  $m$  with  $\langle E_i, V_i, \delta \rangle$  and it writes this tagged message on  $Q_{i,j}$   
**4** |  $E_i \leftarrow \text{Reach}_{\Delta \setminus \Delta_i}^{\mathcal{T}}(\text{Post}_{\delta}^{\mathcal{T}}(E_i))$   
**5 end**

---



OUTPUT *Algorithm*: Let  $E_i$  be the current state estimate of  $\mathcal{E}_i$ . When  $\mathcal{T}_i$  wants to execute a transition  $\delta = \langle \ell_1, Q_{i,j}!m, \ell_2 \rangle \in \Delta_i$  corresponding to an output on the queue  $Q_{i,j}$ , the following instructions are computed to update the state estimate  $E_i$ :

- according to the Mattern's algorithm [16],  $V_i[i]$  is incremented (i.e.,  $V_i[i] \leftarrow V_i[i] + 1$ ) to indicate that a new event has occurred in  $\mathcal{T}_i$ .
- $\mathcal{T}_i$  tags message  $m$  with  $\langle E_i, V_i, \delta \rangle$  and writes this information on  $Q_{i,j}$ . The estimate  $E_i$  tagging  $m$  contains the set of states in which  $\mathcal{T}$  can be *before* the execution of  $\delta$ . The additional information  $\langle E_i, V_i, \delta \rangle$  will be used by  $\mathcal{T}_j$  to refine its state estimate.
- $E_i$  is updated as follows, to contain the current state of  $\mathcal{T}$  and any future state that can be reached from this current state by firing actions that do not belong to  $\mathcal{T}_i$ :  $E_i \leftarrow \text{Reach}_{\Delta \setminus \Delta_i}^{\mathcal{T}}(\text{Post}_{\delta}^{\mathcal{T}}(E_i))$ . More precisely,  $\text{Post}_{\delta}^{\mathcal{T}}(E_i)$  gives the set of states in which  $\mathcal{T}$  can be after the execution of  $\delta$ . But, after the execution of this transition,  $\mathcal{T}_j$  ( $\forall j \neq i \in [1, n]$ ) could read and write on their queues. Therefore, we define the estimate  $E_i$  by  $\text{Reach}_{\Delta \setminus \Delta_i}^{\mathcal{T}}(\text{Post}_{\delta}^{\mathcal{T}}(E_i))$ .

---

**Algorithm 3:** inputTransition( $\mathcal{T}, V_i, E_i, \delta$ )

---

**input** :  $\mathcal{T} = \mathcal{T}_1 || \dots || \mathcal{T}_n$ , the vector clock  $V_i$  of  $\mathcal{E}_i$ , the current state estimate  $E_i$  of  $\mathcal{E}_i$  and a transition  $\delta = \langle \ell_1, Q_{j,i}?m, \ell_2 \rangle \in \Delta_i$ . Message  $m$  is tagged with the triple  $\langle E_j, V_j, \delta' \rangle$  where (i)  $E_j$  is the state estimate of  $\mathcal{E}_j$  before the execution of  $\delta'$  by  $\mathcal{T}_j$ , (ii)  $V_j$  is the vector clock of  $\mathcal{E}_j$  after the execution of  $\delta'$  by  $\mathcal{T}_j$ , and (iii)  $\delta' = \langle \ell'_1, Q_{j,i}!m, \ell'_2 \rangle \in \Delta_j$  is the output corresponding to  $\delta$ .

**output**: The state estimate  $E_i$  after the input transition  $\delta$ .

```

1 begin
2   \ \ We consider three cases to update  $E_j$ 
3   if  $V_j[i] = V_i[i]$  then  $Temp_1 \leftarrow \text{Post}_{\delta}^{\mathcal{T}}(\text{Reach}_{\Delta \setminus \Delta_i}^{\mathcal{T}}(\text{Post}_{\delta'}^{\mathcal{T}}(E_j)))$ 
4   else if  $V_j[j] > V_i[j]$  then
5      $Temp_1 \leftarrow \text{Post}_{\delta}^{\mathcal{T}}(\text{Reach}_{\Delta \setminus \Delta_i}^{\mathcal{T}}(\text{Reach}_{\Delta \setminus \Delta_j}^{\mathcal{T}}(\text{Post}_{\delta'}^{\mathcal{T}}(E_j))))$ 
6   else  $Temp_1 \leftarrow \text{Post}_{\delta}^{\mathcal{T}}(\text{Reach}_{\Delta}^{\mathcal{T}}(\text{Post}_{\delta'}^{\mathcal{T}}(E_j)))$ 
7    $E_i \leftarrow \text{Post}_{\delta}^{\mathcal{T}}(E_i)$  \ \ We update  $E_i$ 
8    $E_i \leftarrow E_i \cap Temp_1$  \ \  $E_i = \text{update of } E_i \cap \text{update of } E_j$  (i.e.,  $Temp_1$ )
9    $V_i[i] \leftarrow V_i[i] + 1$ 
10  for  $k \leftarrow 1$  to  $n$  do  $V_i[k] \leftarrow \max(V_i[k], V_j[k])$ 
11 end
```

---

INPUT *Algorithm*: Let  $E_i$  be the current state estimate of  $\mathcal{E}_i$ . When  $\mathcal{T}_i$  executes a transition  $\delta = \langle \ell_1, Q_{j,i}?m, \ell_2 \rangle \in \Delta_i$ , corresponding to an input on the queue  $Q_{j,i}$ , it also reads the information  $\langle E_j, V_j, \delta' \rangle$  (where  $E_j$  is the state estimate of  $\mathcal{E}_j$  before the execution of  $\delta'$  by  $\mathcal{T}_j$ ,  $V_j$  is the vector clock of  $\mathcal{E}_j$  after the execution of  $\delta'$  by  $\mathcal{T}_j$ , and  $\delta' = \langle \ell'_1, Q_{j,i}!m, \ell'_2 \rangle \in \Delta_j$  is the output corresponding to  $\delta$ ) tagging  $m$ , and the following operations are performed to update  $E_i$ :

- we update the state estimate  $E_j$  of  $\mathcal{E}_j$  (this update is denoted by  $Temp_1$ ) by using the vector clocks to guess the possible behaviors of  $\mathcal{T}$  between the execution of the transition  $\delta'$  and the execution of  $\delta$ . We consider three cases :

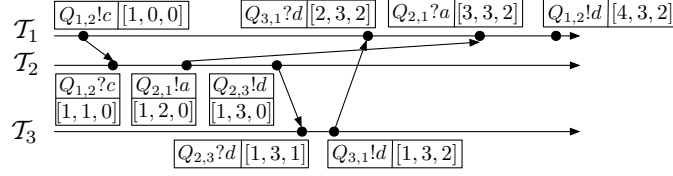


Fig. 1: An execution of the running example.

- if  $V_j[i] = V_i[i] : Temp_1 \leftarrow Post_\delta^T(Reach_{\Delta \setminus \Delta_i}^T(Post_{\delta'}^T(E_j)))$ . In this case, thanks to the vector clocks, we know that  $T_i$  has executed no transition between the execution of  $\delta'$  by  $T_j$  and the execution of  $\delta$  by  $T_i$ . Thus, only transitions in  $\Delta \setminus \Delta_i$  could have occurred during this period. We then update  $E_j$  as follows. We compute (i)  $Post_{\delta'}^T(E_j)$  to take into account the execution of  $\delta'$  by  $T_j$ , (ii)  $Reach_{\Delta \setminus \Delta_i}^T(Post_{\delta'}^T(E_j))$  to take into account the transitions that could occur between the execution of  $\delta'$  and the execution of  $\delta$ , and (iii)  $Post_\delta^T(Reach_{\Delta \setminus \Delta_i}^T(Post_{\delta'}^T(E_j)))$  to take into account the execution of  $\delta$ .
- else if  $V_j[j] > V_i[j] : Temp_1 \leftarrow Post_\delta^T(Reach_{\Delta \setminus \Delta_i}^T(Reach_{\Delta \setminus \Delta_j}^T(Post_{\delta'}^T(E_j))))$ . Indeed, in this case, we can prove (see Theorem 1) that if we reorder the transitions executed between the occurrence of  $\delta'$  and the occurrence of  $\delta$  in order to execute the transitions of  $\Delta_i$  before the ones of  $\Delta_j$ , we obtain a correct update of  $E_i$ . Intuitively, this reordering is possible, because there is no causal relation between the events of  $T_i$  and the events of  $T_j$ , that have occurred between  $\delta'$  and  $\delta$ . So, in this reordered sequence, we know that, after the execution of  $\delta$ , only transitions in  $\Delta \setminus \Delta_j$  could occur followed by transitions in  $\Delta \setminus \Delta_i$ .
- else  $Temp_1 \leftarrow Post_\delta^T(Reach_\Delta^T(Post_{\delta'}^T(E_j)))$ . Indeed, in this case, the vector clocks do not allow us to deduce information regarding the behavior of  $T$  between the execution of  $\delta'$  and the execution of  $\delta$ . Therefore, to have a correct state estimate, we update  $E_j$  by taking into account all the possible behaviors of  $T$  between the execution of  $\delta'$  and the execution of  $\delta$ .
- we update the estimate  $E_i$  to take into account the execution of  $\delta$ :  $E_i \leftarrow Post_\delta^T(E_i)$ .
- we intersect  $Temp_1$  and  $E_i$  to obtain a better state estimate:  $E_i \leftarrow E_i \cap Temp_1$ .
- according to the Mattern's algorithm [16], the vector clock  $V_i$  is incremented to take into account the execution of  $\delta$  and subsequently is set to the component-wise maximum of  $V_i$  and  $V_j$ . This last operation allows us to take into account the fact that any event that precedes the sending of  $m$  should also precede the occurrence of  $\delta$ .

*Example 2.* We illustrate SE-algorithm with a sequence of actions of our running example depicted in Figure 1 (the vector clocks are given in the figure). A state of the global system is denoted by  $\langle \ell_1, \ell_2, \ell_3, w_{1,2}, w_{2,1}, w_{2,3}, w_{3,1} \rangle$  where  $\ell_i$  is the location of  $T_i$  (for  $i = 1, 2, 3$ ) and  $w_{1,2}, w_{2,1}, w_{2,3}$  and  $w_{3,1}$  denote the content of the queues  $Q_{1,2}, Q_{2,1}, Q_{2,3}$  and  $Q_{3,1}$ . At the beginning of the execution, the state estimates of the three subsystems are (i)  $E_1 = \{ \langle A_0, B_0, C_0, \epsilon, \epsilon, \epsilon, \epsilon \rangle \}$ , (ii)  $E_2 = \{ \langle A_0, B_0, C_0, \epsilon, \epsilon, \epsilon, \epsilon \rangle \}$ , and (iii)  $E_3 = \{ \langle A_0, B_0, C_0, \epsilon, \epsilon, \epsilon, \epsilon \rangle, \langle A_1, B_0, C_0, c, \epsilon, \epsilon, \epsilon \rangle \}$ .

$\langle A_1, B_1, C_0, \epsilon, b^*, \epsilon, \epsilon \rangle, \langle A_1, B_2, C_0, \epsilon, b^*(a + \epsilon), \epsilon, \epsilon \rangle, \langle A_1, B_3, C_0, \epsilon, b^*(a + \epsilon), d, \epsilon \rangle$ .  
 After the first transition  $\langle A_0, Q_{1,2}!c, A_1 \rangle$ , the state estimate of  $\mathcal{E}_1$  is not really precise, because a lot of events may have happened without the estimator  $\mathcal{E}_1$  being informed:  $E_1 = \{ \langle A_1, B_0, C_0, c, \epsilon, \epsilon, \epsilon \rangle, \langle A_1, B_1, C_0, \epsilon, b^*, \epsilon, \epsilon \rangle, \langle A_1, B_2, C_0, \epsilon, b^*a, \epsilon, \epsilon \rangle, \langle A_1, B_3, C_0, \epsilon, b^*(a + \epsilon), d, \epsilon \rangle, \langle A_1, B_3, C_1, \epsilon, b^*(a + \epsilon), \epsilon, \epsilon \rangle, \langle A_1, B_3, C_0, \epsilon, b^*(a + \epsilon), \epsilon, d \rangle$ . But, after the second transition  $\langle B_0, Q_{1,2}?c, B_1 \rangle$ ,  $\mathcal{E}_2$  has an accurate state estimate:  $E_2 = \{ \langle A_1, B_1, C_0, \epsilon, \epsilon, \epsilon, \epsilon \rangle$ . We skip a few steps and consider the state estimates before the sixth transition  $\langle C_1, Q_{3,1}!d, C_0 \rangle$ :  $E_1$  is still the same, because  $\mathcal{T}_1$  did not perform any action,  $E_3 = \{ \langle A_1, B_3, C_1, \epsilon, b^*(a + \epsilon), \epsilon, \epsilon \rangle$ , and we do not indicate  $E_2$ , because  $\mathcal{T}_2$  is no longer involved. When  $\mathcal{T}_3$  sends the message  $d$  to  $\mathcal{T}_1$  (the transition  $\langle C_1, Q_{3,1}!d, C_0 \rangle$ ), it attaches  $E_3$  to this message. When  $\mathcal{T}_1$  reads this message, it computes  $E_1 = \{ \langle A_2, B_3, C_0, \epsilon, b^*(a + \epsilon), \epsilon, \epsilon \rangle$  and when it reads the message  $a$ , it updates  $E_1$ :  $E_1 = \{ \langle A_2, B_3, C_0, \epsilon, b^*, \epsilon, \epsilon \rangle$ . Thus,  $\mathcal{E}_1$  knows, after this action, that there is no  $a$  in  $Q_{2,1}$ , and that after writing  $d$  on  $Q_{1,2}$ , it cannot reach  $A_{er}$  from  $A_0$ . This example shows the importance of knowing the content of the queues as without this knowledge,  $\mathcal{E}_1$  may think that there is an  $a$  in  $Q_{2,1}$ , so an error might occur if  $\langle A_2, Q_{1,2}!d, A_0 \rangle$  is enabled.  $\diamond$

**Properties.** As explained above, we assume that we can compute an approximation of the reachable states. In this part, we present the properties of our state estimate algorithm w.r.t. the kind of approximations that we use.

**Theorem 1.** *SE-algorithm is complete, if the Reach operator computes an overapproximation of the reachable states.*

**Theorem 2.** *SE-algorithm is sound, if the Reach operator computes an underapproximation of the reachable states.*

The proofs of these theorems are given in [11]. If we compute an underapproximation of the reachable states, our state estimate algorithm is not complete. If we compute an overapproximation of the reachable states, our state estimate algorithm is not sound. So, depending on the approximations, our algorithm is either complete or sound. Completeness is a more important property, because it ensures that the computed state estimates always contain the current global state. Therefore, in section 5, we define an effective algorithm for the state estimate problem by computing overapproximations of the reachable states. Finally, note that our method proposes that we only add information to existing transmitted messages. We can show that increasing the information exchanged between the estimators (for example, each time an estimator computes a new state estimate, this estimate is sent to all the other estimators) improves their state estimate. This can be done only if the channels and the subsystems can handle this extra load.

## 5 Effective Computation of State Estimates by Means of Abstract Interpretation

The algorithm described in the previous section requires the computation of reachability operators, which cannot always be computed exactly in general. In this section, we

overcome this obstacle by using abstract interpretation techniques (see e.g. [6, 13]) to compute, in a finite number of steps, an overapproximation of the reachability operators and thus of the state estimates  $E_i$ .

**Computation of Reachability Sets by the Means of Abstract Interpretation.** For a given set of global states  $X' \subseteq X$  and a given set of transitions  $\Delta' \subseteq \Delta$ , the states reachable from  $X'$  can be characterized by the least fixpoint:  $\text{Reach}_{\Delta'}^{\mathcal{T}}(X') = \mu Y. X' \cup \text{Post}_{\Delta'}^{\mathcal{T}}(Y)$ . Abstract interpretation provides a theoretical framework to compute efficient overapproximation of such fixpoints. The concrete domain (i.e., the sets of states  $2^X$ ), is substituted by a simpler abstract domain  $\Lambda$ , linked by a *Galois Connection*  $2^X \xrightleftharpoons[\alpha]{\gamma} \Lambda$  [6], where  $\alpha$  (resp.  $\gamma$ ) is the abstraction (resp. concretization) function.

The fixpoint equation is transposed into the abstract domain:  $\lambda = F_{\Delta'}^{\#}(\lambda)$ , with  $\lambda \in \Lambda$  and  $F_{\Delta'}^{\#} \sqsupseteq \alpha \circ F_{\Delta'} \circ \gamma$ . In this setting, a standard way to ensure that the fixpoint computation converges after a finite number of steps to some overapproximation  $\lambda_{\infty}$ , is to use a *widening operator*  $\nabla$ . The concretization  $c_{\infty} = \gamma(\lambda_{\infty})$  is an overapproximation of the least fixpoint of the function  $F_{\Delta'}$ .

**Choice of the Abstract Domain.** In abstract interpretation based techniques, the quality of the approximation we obtain depends on the choice of the abstract domain  $\Lambda$ . In our case, the main issue is to abstract the content of the FIFO channels. As discussed in [13], a good abstract domain is the class of regular languages, which can be represented by finite automata. Let us recall the main ideas of this abstraction.

**Finite Automata as an Abstract Domain.** We first assume that there is only one queue in the distributed system  $\mathcal{T}$ ; we explain later how to handle a distributed system with several queues. With one queue, the concrete domain of the system  $\mathcal{T}$  is defined by  $X = 2^{L \times M^*}$ . A set of states  $Y \subseteq 2^{L \times M^*}$  can be viewed as a map  $Y : L \mapsto 2^{M^*}$  that associates a language  $Y(\ell)$  with each location  $\ell \in L$ ;  $Y(\ell)$  therefore represents the possible contents of the queue in the location  $\ell$ . To simplify the computation, we substitute the concrete domain  $\langle L \mapsto 2^{M^*}, \subseteq, \cup, \cap, L \times M^*, \emptyset \rangle$  by the abstract domain  $\langle L \mapsto \text{Reg}(M), \subseteq, \cup, \cap, L \times M^*, \emptyset \rangle$ , where  $\text{Reg}(M)$  is the set of *regular languages* over the alphabet  $M$ . Since regular languages have a canonical representation given by finite automata, each operation (union, intersection, left concatenation,...) in the abstract domain can be performed on finite automata.

**Widening Operator.** The widening operator is also performed on a finite automaton, and consists in quotienting the nodes<sup>4</sup> of the automaton by the *k-bounded bisimulation equivalence relation*  $\equiv_k$ ;  $k \in \mathbb{N}$  is a parameter which allows us to tune the precision, since increasing  $k$  improves the quality of the abstractions in general. Two nodes are equivalent w.r.t.  $\equiv_k$  if they have the same outgoing path (sequence of labeled transitions) up to length  $k$ . While we merge the equivalent nodes, we keep all transitions and we obtain an automaton recognizing a larger language. Note that for a fixed  $k$ , the class of automata which results from such a quotient operation from any original automaton, is finite and its cardinality is bounded by a number which is only function of  $k$ . So, when we apply this widening operator regularly, the fixpoint computation terminates (see [13] for more details).

<sup>4</sup> The states of an automaton representing the queue contents are called nodes to avoid the confusion with the states of a CFSM.

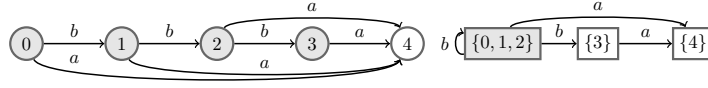


Fig. 2: Illustration of the 1-bounded bisimulation relation  $\equiv_1$  for  $\mathcal{A}$ .

*Example 3.* We consider the automaton  $\mathcal{A}$  depicted in Figure 2, whose recognized language is  $a + ba + bba + bbba$ . We consider the 1-bounded bisimulation relation i.e., two nodes of the automaton are equivalent if they have the same outgoing transitions. So, nodes 0, 1, 2 are equivalent, since they all have two transitions labeled by  $a$  and  $b$ . Nodes 3 and 4 are equivalent to no other node since 4 has no outgoing transition whereas only  $a$  is enabled in node 3. When we quotient  $\mathcal{A}$  by this equivalent relation, we obtain the automaton on the right side of Figure 2, whose recognized language is  $b^*a$ .  $\diamond$

When the system contains several queues  $Q = \{Q_1, \dots, Q_r\}$ , their content can be represented by a concatenated word  $w_1\# \dots \#w_r$  with one  $w_i$  for each queue  $Q_i$  and  $\#$ , a delimiter. With this encoding, we represent a set of queue contents by a finite automaton of a special kind, namely a QDD [2]. Since QDDs are finite automata, classical operations (union, intersection, left concatenation,...) in the abstract domain are performed as was done previously. We must only use a slightly different widening operator not to merge the different queue contents [13].

**Effective SE-algorithm.** The Reach operator is computed using those abstract interpretation techniques: we proceed to an iterative computation in the abstract domain of regular languages and the widening operator ensures that this computation terminates after a finite number of steps [6]. So the Reach operator always gives an overapproximation of the reachable states regardless the distributed system. The efficiency of these approximations is measured in the experiments of section 6. Because of Theorem 1, our SE-algorithm is complete.

## 6 Experiments

We have implemented the SE-algorithm as a new feature of the McScM tool [17], a model checker for distributed systems modeled by CFSM. Since it represents queue contents by QDDs, this software provides most of the functionalities needed by our algorithm, like effective computation of reachable states. We have also added a mechanism to manage vector clocks, and an interactive simulator. This simulator first computes and displays the initial state estimates. At each step, it asks the user to choose a possible transition.

We proceeded to an evaluation of our algorithm measuring the size of the state estimates. Note that this size is not the number of global states of the state estimate (which may be infinite) but the number of nodes of its QDD representation. We generated random sequences of transitions for our running example and some other examples of [10]. Table 1 shows the average execution time for a random sequence of 100 transitions, the memory required (heap size), the average and maximal size of the state estimates.

example	# subsystems	# channels	time [s]	memory [MB]	maximal size	average size
running example	3	4	7.13	5.09	143	73.0
c/d protocol	2	2	5.32	8.00	183	83.2
non-regular protocol	2	1	0.99	2.19	172	47.4
ABP	2	3	1.19	2.19	49	24.8
sliding window	2	2	3.26	4.12	21	10.1
POP3	2	2	3.08	4.12	22	8.5

Table 1: Experiments

Default value of the widening parameter is  $k = 1$ . Experiments were done on a standard MacBook Pro with a 2.4 GHz Intel core 2 duo CPU. These results show that the computation of state estimates takes about 50ms per transition and that the symbolic representation of state estimates we add to messages are automata with a few dozen nodes. A sensitive element in the method is the size of the computed and transmitted information. It can be improved by the use of compression techniques to reduce the size of this information. A more evolved technique would consist in the offline computation of the set of possible estimates. Estimates are indexed in a table, available at execution time to each local estimator. If we want to keep an online algorithm, we can use the memoization technique. When a state estimate is computed for the first time, it is associated with an index that is transmitted to the subsystem which records both values. If the same estimate must be transmitted, only its index can be transmitted and the receiver can find from its table the corresponding estimate. We also highlight that our method works better on the real-life communication protocols we have tested (alternating bit protocol, sliding window, POP3) than on the examples we introduced to test our tool.

## 7 Conclusion and Future Work

We have proposed an effective algorithm to compute online, locally to each subsystem, an estimate of the global state of a running distributed system, modeled as *communicating finite state machines* with reliable unbounded FIFO queues. With such a system, a global state is composed of the current location of each subsystem together with the channel contents. The principle is to add a local estimator to each subsystem such that most of the system is preserved; each local estimator is only able to compute information and in particular symbolic representations of state estimates and to piggyback some of this computed information to the transmitted messages. Since these estimates may be infinite, a crucial point of our work has been to propose and evaluate the use of regular languages to abstract sets of FIFO queues. In practice, we have used  $k$ -bisimilarity relations, which allows us to represent each (possibly infinite) set of queue contents by the minimal and canonical  $k$ -bisimilar finite automaton which gives an overapproximation of this set. Our algorithm transmits state estimates and vector clocks between subsystems to allow them to refine and preserve consistent state estimates. More elaborate examples must be taken to analyze the precision of our algorithm and see, in practice, if the estimates are sufficient to solve diagnosis or control problems. Anyway, it appears

important to study the possibility of reducing the size of the added communication while preserving or even increasing the precision in the transmitted state estimates.

## References

1. G. Berry and G. Gonthier. The estereel synchronous programming language: Design, semantics, implementation. *Sci. Comput. Program.*, 19(2):87–152, 1992.
2. B. Boigelot, P. Godefroid, B. Willems, and P. Wolper. The power of QDDs. In *SAS '97: Proceedings of the 4th International Symposium on Static Analysis*, pages 172–186, 1997.
3. D. Brand and P. Zafropulo. On communicating finite-state machines. *J. ACM*, 30(2):323–342, 1983.
4. C. Cassandras and S. Lafortune. *Introduction to Discrete Event Systems*. Kluwer Academic Publishers, 1999.
5. T. Chatain, P. Gastin, and N. Sznajder. Natural specifications yield decidability for distributed synthesis of asynchronous systems. In *SOFSEM, LNCS 5404*, pages 141–152, 2009.
6. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL'77*, pages 238–252, 1977.
7. R. Debouk, S. Lafortune, and D. Teneketzis. Coordinated decentralized protocols for failure diagnosis of discrete event systems. *Discrete Event Dynamical Systems: Theory and Applications*, 10:33–79, 2000.
8. B. Genest. On implementation of global concurrent systems with local asynchronous controllers. In *CONCUR, LNCS 3653*, pages 443–457, 2005.
9. L. Hérouet, T. Gazagnaire, and B. Genest. Diagnosis from scenarios. In *proc. of the 8th Int. Workshop on Discrete Events Systems, WODES'06*, pages 307–312, 2006.
10. A. Heußner, T. Le Gall, and G. Sutre. Extrapolation-Based Path Invariants for Abstraction Refinement of Fifo Systems. In *Proc. Model Checking Software, SPIN Workshop 2009, LNCS 5578*, pages 107–124. Springer, 2009.
11. G. Kalyon, T. Le Gall, H. Marchand, and T. Massart. Global state estimates for distributed systems (version with proofs). In *FORTE*, 2011. <http://hal.inria.fr/inria-00581259/>.
12. L. Lamport. Time, clocks, and the ordering of events in a distributed system. *Communications of the ACM*, 21(7):558–565, 1978.
13. T. Le Gall, B. Jeannot, and T. Jéron. Verification of communication protocols using abstract interpretation of fifo queues. In *AMAST '06, LNCS 4019*, July 2006.
14. F. Lin, K. Rudie, and S. Lafortune. Minimal communication for essential transitions in a distributed discrete-event system. *Trans. on Automatic Control*, 52(8):1495–1502, 2007.
15. T. Massart. A calculus to define correct transformations of lotos specifications. In *FORTE*, volume C-2 of *IFIP Transactions*, pages 281–296, 1991.
16. F. Mattern. Virtual time and global states of distributed systems. In *Proceedings of the Workshop on Parallel and Distributed Algorithms*, pages 215–226, 1989.
17. The McScM library, 2009. <http://altarica.labri.fr/forge/projects/mcscm/wiki/>.
18. L. Ricker and B. Caillaud. Mind the gap: Expanding communication options in decentralized discrete-event control. In *Conference on Decision and Control*, 2007.
19. M. Sampath, R. Sengupta, S. Lafortune, K. Sinaamohideen, and D. Teneketzis. Failure diagnosis using discrete event models. *IEEE Transactions on Control Systems Technology*, 4(2):105–124, March 1996.
20. S. Tripakis. Decentralized control of discrete event systems with bounded or unbounded delay communication. *IEEE Trans. on Automatic Control*, 49(9):1489–1501, 2004.
21. S. Xu and R. Kumar. Distributed state estimation in discrete event systems. In *ACC'09: Proc. of the 2009 conference on American Control Conference*, pages 4735–4740, 2009.