

Labelled Transition System Approach to Distributed Computing Systems

Dang Van Hung
Institute of Informatics
National Centre for Scientific Research of Vietnam

1. Introduction

Labelled transition systems defined by Keller [15] in 1976 has been used widely for defining operational semantics for parallel programming languages as well as models for parallel systems [1,5,15]. One advantage of them is that they are very simple. However, semantics represented by them is interleaving one, i.e. they model independence by arbitrary interleaving. The situation that two events are independent is modeled by saying that the two events may occur in any order, which refers to a global ordering of the events. In order to overcome this shortcoming, Mazurkiewicz has introduced the concept of traces. Each Mazurkiewicz's trace contains sequences of events which are equivalent by an independence. One problem with traces is that the independence is fixed in the system, which restricts the expressive power of the models with traces. Vectors of firing sequences introduced by Shield [19] represents the behavior of concurrent and distributed systems more intuitively. This model has the same expressive power as traces [11], and is suitable for synchronous systems only.

In this paper we combined the mentioned concepts by considering an event of concurrent and distributed systems to contain atomic events executed in parallel by component processes. In this way we can define labelled transition systems, an intuitive model for distributed systems, which can express the concurrency explicitly and can be used in designing, implementing and verifying distributed algorithms.

In order to show the advantages of our model we relate it to partial orderings of events [5,6,16,18] and consider the problem of determining global states and the logical time of distributed systems. By keeping partial knowledge of the others processes in the states of each process in the systems, we propose a new algorithm for distributed snapshots. It is shown that the amount of global

information which can be maintained in the local states of the processes of the system by exchanging messages is not more than what derived from actions leading to them.

The paper is organized as follows. In the next section we introduce labelled transition systems for distributed computing systems and some behavioral properties. In the third section we relate our models to partial orderings of events, define the logical time in distributed systems and present a general algorithm for the maintenance of interesting global knowledge for solving some problems concerning the distributed systems.

2. Labelled Transition Systems for Distributed Computing

In this section we introduce a formal, general model for distributed computing systems and associate them with labelled transition systems in order to describe their behavior. In this model of the behavior, the semantics of the systems becomes easier to understand, and we can formulate precisely some problems with the systems. The distributed computing systems considered in this paper are taken from [2,10] as follows. A distributed computing system consists of multiple autonomous processors that do not share primary memory but cooperate by sending messages over a communication network. The communication network is assumed to be connected so that each processor can send message to any other, and to be reliable (i.e., messages are delivered error-free without ever being lost). Logically, between any two ordered processors there exists an one way communication channel. We consider in this paper two kinds of message passing: synchronous and asynchronous one. With synchronous message passing, the sender is blocked until the receiver has accepted the message. With asynchronous message passing, the sender does not wait for the receiver to be ready to accept its message. Conceptually, the sender continues immediately after sending the message. We also assume that each message sent or accepted involves an explicit action.

We start with the definition of labelled transition systems which are essentially those given in [6,15,17].

Definition 1 (Labelled transition systems). A labeled transition system is a triple $P = (Q, \rightarrow, A)$, where Q is the set of configurations, A is the set of labels and \rightarrow is the transition relation.

In the sequel, we shall write $p - a \rightarrow q$ instead of $(p, q, a) \in \rightarrow$. P is said to be deterministic iff

$$|\{q \in Q | p - a \rightarrow q\}| \leq 1, \forall p \in Q, a \in A.$$

Definition 2 (Distributed Computing Systems). A distributed computing system with N processes is a tuple $D = (P_1, P_2, \dots, P_N, C, \rightarrow, M)$, where

- (i) $P_i = (Q_i, \rightarrow_i, A_i)$ is a labelled transition system, called the i^{th} process of the system, Q_i is the set of local states of the i^{th} process, and A_i is the set of its local computation actions,
- (ii) C is the set of communication actions of the system, $C = \cup_{i,j=1}^N \{s_{ij}, r_{ij}\}$, r_{ij} is the message receive of the i^{th} process from j^{th} process, s_{ij} is the message send of the i^{th} process destinationed the j^{th} process,

(iii) M is a matrix, of which the i D row j D column element M_{ij} is a set of possible messages from the process i to the process j ,

(iv) $\rightarrow C$ is a transition relation labelled over C , $-s_{ij} \rightarrow C \subseteq Q_i \times (M_{ij} \times Q_i)$, and $-r_{ij} \rightarrow C \subseteq (Q_i \times M_{ji}) \times Q_i$, $i, j = 1, 2, \dots, N$

For the simplicity, we assume that the sets $C, A_i, Q_i, M_{ij}, i, j = 1, 2, \dots, N$ are pairwise disjoint.

D is said to be deterministic iff each labelled transition system in the system is deterministic.

In order to describe the behaviour of distributed systems we first define the set of actions of the system at each point of discrete time. An action of the system at each moment consists of actions of all process in the system including local actions, communicating actions and the idle action.

Definition 3 (Action of distributed systems). Let D be a distributed system given as above. $\Sigma_i = A_i \cup \{r_{ij}, s_{ij}, \lambda | j = 1, 2, \dots, N\}$ is called the set of actions executed by the process i , where λ denotes the idle action. The set $A = \Sigma_1 \times \Sigma_2 \times \dots \times \Sigma_N$ is called the set of possible actions of the system, while the set $A_s = \{\bar{a} = (a_1, a_2, \dots, a_N) | \bar{a} \in A \& a_i = s_{ij} \text{ iff } a_j = r_{ji}\}$ is called the set of its synchronous actions.

The behavior of the system now can be describe as a set of sequences of actions of the system, called vector of action sequences.

Definition 4 (Vector of action sequences). Let

$$VAS(A) = \Sigma_1^* \times \Sigma_2^* \times \dots \times \Sigma_N^*.$$

Each element of $VAS(A)$ is called a vector of action sequences of D .

For the simplicity of our notation, in this paper we adopt the following conventions:

For $\bar{w} \in VAS(A)$, w_i will denote its i^{th} component. For $\bar{w}, \bar{w}' \in VAS(A)$, the product $\bar{w}\bar{w}'$ (concatenation) of \bar{w} and \bar{w}' is defined as

$$\bar{w}\bar{w}' = (w_1 w'_1, w_2 w'_2, \dots, w_N w'_N),$$

For $B \subseteq \Sigma$ and $w = a_1 a_2 \dots a_n \in \Sigma^*$, $\#_B w$ will denote the number of occurrences of the elements of B in w ; n is called the length (number of letter occurrences) of w and will be denoted by $|w|$. The elements of the set $\{(i, a_i) | i = 1, 2, \dots, n\}$ is said to be letter occurrences of w . For a letter occurrence (i, a_i) , we call a_i its label. A single set will be identified with its element. $\bar{\lambda}$ will denotes the element $(\lambda, \lambda, \dots, \lambda) \in A$, while $\lambda_i(a)$ denotes the element in A of which all components is λ except that the i^{th} component is a .

In order to associate a labelled transition system to a given distributed system, we have to specify the state space for it. A state of the distributed system will consist of states of its processes and a state of the communication channels. The states of the channels will be the sequences of messages which have not been accepted.

Definition 5 (State space of distributed systems). Let D be a distributed system given as above. The state space for D is the set

$$S = Q_1 \times Q_2 \times \dots \times Q_N \times \Gamma,$$

where Γ is the set of $N \times N$ matrices γ' 's of which the i -Drow j -column element is a priori-queue of possible messages from the process i to the process j . ($\gamma_{ij} \in M_{ij}^*$).

Each element of S will be called a complete state of the system, while its n first components is called a state of the system.

In this paper we consider a state of a channel as a FIFO queue of possible messages. The empty queue will represent the fact that there is no message on the channel and will be denoted by λ . For FIFO queue γ_{ij} (content of the channel from the process i to the process j) the operation $put(\gamma_{ij}, m)$ means putting the message m into the channel, resulting $m\gamma_{ij}$ as its new content, while the operation $getout(\gamma_{ij})$ returns the element at the front of the queue and takes it out from the queue if the queue is not empty. Otherwise, the operation is undefined.

Labelled transition system for the distributed system D now is defined as follows.

Definition 6 (Labelled transition system for D). Let S and A be as above. The labelled transition system

$$T = (S, \rightarrow, A)$$

is said to be the labelled transition system for D , where \rightarrow is defined below.

Let

$$AL = (A_1 \cup \{\lambda\}) \times (A_2 \cup \{\lambda\}) \times \dots (A_N \cup \{\lambda\}).$$

$$AR = (R_1 \cup \{\lambda\}) \times (R_2 \cup \{\lambda\}) \times \dots (R_N \cup \{\lambda\}), \quad R_i = \bigcup_{j=1}^N \{r_{ij}\}, \quad i = 1, 2, \dots, N,$$

$$AS = (S_1 \cup \{\lambda\}) \times (S_2 \cup \{\lambda\}) \times \dots (S_N \cup \{\lambda\}), \quad S_i = \bigcup_{j=1}^N \{s_{ij}\}, \quad i = 1, 2, \dots, N.$$

In the sequel, when labelled transition relations can be distinguished one from another by the set of labels, we omit the subscript if it does not cause any confusion.

For $\bar{q} = (q_1, q_2, \dots, q_N, \gamma)$, $\bar{q}' = (q'_1, q'_2, \dots, q'_N, \gamma') \in S$, $\bar{a} = (a_1, a_2, \dots, a_N) \in A$, we say $\bar{q} - \bar{a} \rightarrow \bar{q}'$ iff

$$(i). \quad \bar{a} \in AL \& (\gamma = \gamma') \& (\forall i \leq N) : q_i - a_i \rightarrow q'_i,$$

$$(ii). \quad (\bar{a} \in AS) \& (\forall i \in N) :$$

$$((a_i = \lambda) \Rightarrow ((q_i = q'_i) \& (\forall j \leq N) : (\gamma_{ij} = \gamma'_{ij}))) \& ((a_i = s_{ik_i}) \Rightarrow ((q_i - s_{ik_i} \rightarrow (m, q'_i)) \&$$

$$(put(\gamma_{ik_i}, m) = \gamma'_{ik_i}) \&$$

$$((\forall j \leq N, j \neq k_i) : (\gamma'_{ij} = \gamma_{ij}))),$$

$$(iii). \quad (\bar{a} \in AR) \& (\forall i \leq N) :$$

$$((a_i = \lambda) \Rightarrow ((q_i = q'_i) \& (\forall j \leq N) : (\gamma_{ij} = \gamma'_{ij})))$$

$$\& ((a_i = r_{ik_i}) \Rightarrow ((m = getout(\gamma_{ik_i})) \& ((q_i, m) - r_{ik_i} \rightarrow q'_i) \& ((\forall j \leq N, j \neq k_i) : (\gamma'_{ij} = \gamma_{ij}))),$$

$$(iv).$$

$$E(\bar{a} = \bar{a}_1 \bar{a}_2 \bar{a}_3) \& (\bar{a}_1 \in AL, \bar{a}_2 \in AS, \bar{a}_3 \in AR) \& (\exists \bar{p}, \bar{p}' \in S) :$$

$$(\bar{q} - \bar{a}_1 \rightarrow \bar{p} - \bar{a}_2 \rightarrow \bar{p}' - \bar{a}_3 \rightarrow \bar{q}').$$

(Notice that for each $\bar{a} \in A$, there exist uniquely $\bar{a}_1 \in AL$, $\bar{a}_2 \in AS$, $\bar{a}_3 \in AR$ such that \bar{a} can be written as $\bar{a} = \bar{a}_1 \bar{a}_2 \bar{a}_3$.)

The definition achieves the general behaviour of distributed systems. Each process in a system can perform its local computations independently from the others. The message passing is synchronous if the specified channel is empty and the receiver is ready to accept message. The accept action can be delayed until the receiver becomes ready to perform it.

The relation \rightarrow is extended to the one labeled over $VAS(A)$ in the natural way as follows:

For $\bar{q}, \bar{q}' \in Q$, $\bar{w} \in VAS(A)$, $\bar{q} - \bar{w} \rightarrow \bar{q}'$ iff there exist $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_k \in A$, $\bar{q}_0, \bar{q}_1, \bar{q}_k \in S$, and s.t. $\bar{q}_0 = \bar{q}$, $\bar{q}_k = \bar{q}'$, and $\bar{q}_{i-1} - \bar{a}_i \rightarrow \bar{q}_i$, $i = 1, 2, \dots, k$.

Definition 7 (Transition system for synchronized computations). Let T be the labelled transition system for D . The restriction of T on the label set A_s

$$T_s = (S, \rightarrow, A_s)$$

is called the labelled transition system for synchronized computations of D , or the synchronized transition system for D .

Let $\bar{q}_0 = (q_1^0, q_2^0, \dots, q_N^0, \Phi)$ be a fixed state, called initial state of D . (Φ denotes the matrix having empty queues as its elements).

Definition 8 (Behaviour and reachable states of distributed systems). Let

$$B = \{\bar{w} \in VAS(A) \mid \bar{q}_0 - \bar{w} \rightarrow \bar{q}, \text{ for some } \bar{q} \in S\},$$

$$R = \{\bar{q} \in S \mid \bar{q}_0 - \bar{w} \rightarrow \bar{q} \text{ for some } \bar{w} \in VAS(A)\}.$$

We call B the behaviour, R the set of reachable states of D (with the initial state \bar{q}_0).

With the substitution of $VAS(A)$ by $VAS(A_s)$ in this definition, we get the definition of the synchronized behavior and the set of synchronized reachable states of D .

Now, we relate our model of behaviour to the net theory by associating to each element of B an occurrence net. The readers, who are not familiar with net theory, are referred to [18] for their details.

From the definition of T it follows that if $\bar{q}_0 - \bar{w} \rightarrow \bar{q}$, \bar{w} can be written as $\bar{a}_1 \bar{a}_2 \dots \bar{a}_K$ such that $\bar{q}_0 - \bar{a}_1 \rightarrow \bar{q}_1 \rightarrow \bar{a}_2 \dots \rightarrow \bar{a}_K \rightarrow \bar{q}_K$, $\bar{q}_K = \bar{q}$, where $\bar{a}_k, k \leq K$, is an element of A having only one component different from λ .

Let $\bar{q}_k = (q_1^k, q_2^k, \dots, q_N^k, \gamma_k)$, $\bar{w}_k = \bar{a}_1 \bar{a}_2 \dots \bar{a}_k, k \leq K, \bar{w} = (w_1, w_2, \dots, w_N)$. The occurrence net $\pi(\bar{w}) = (G, E, F)$ and the labelling function $l: G \cup E \rightarrow (\cup_{i=1}^N Q_i) \cup (\cup_{i,j=1}^N M_{ij})$ are constructed by induction on k as follows.

$k = 0$:

Let $\pi(\lambda) = (G_0, \emptyset, \emptyset), G_0 = \{g_1^0, g_2^0, \dots, g_N^0\}, l_0(g_j^0) = q_j^0, j \leq N, I_0 = \emptyset$ (R is introduced for the construction of $\pi(\bar{w}_{k+1})$).

$k \rightarrow k+1$:

Let $\pi(\bar{w}_k) = (G_k, E_k, F_k), l_k : G_k \cup E_k \rightarrow (\cup_{i=1}^N Q_i) \cup (\cup_{i,j=1}^N M_{ij})$ have been constructed, I_k be a partial function from G_k to $\{(i, j, n) | i, j \leq N, n \geq 0\}$ such that $I_k(g)$ is defined when $l(g) \in \cup_{i,j=1}^N M_{ij}$. Let $\bar{a}_k = \bar{\lambda}_i(a)$.

Then, $E_{k+1} = E_k \cup \{e\}, l_{k+1}(e) = a, l_{k+1}|_{E_k \cup G_k} = l_k$, where e is a new transition $\varphi|_\Omega$ denotes the restriction of φ on Ω . Depending on a , we have the following cases:

a). $a \in A_i$. By the definition of $\rightarrow, q_i^k - a \rightarrow q_i^{k+1}$. Then, $B_{k+1} = G_k \cup \{g\}$, where g is a new place, $l_{k+1}(g) = q_i^{k+1}$. Let g' be the place in B_k having no outgoing arc such that $l_k(g') = q_i^k$. (Such a place exists uniquely). Now, $F_{k+1} = F_k \cup \{(g', e), (e, g)\}$.

b). $a = s_{ij}$ for some $j \leq N$. By the definition of $\rightarrow, q_i^k - a \rightarrow (q_i^{k+1}, m)$. Then, $G_{k+1} = G_k \cup \{g, g'\}$, where g, g' is a new place, $l_{k+1}(g) = q_i^{k+1}, l_{k+1}(g') = m, I_{k+1}(g') = (i, j, n)$, where n is the number of the occurrences of s_{ij} in the i^{th} component of \bar{w}_k . Let g'' be the place in G_k having no outgoing arc such that $l_k(g'') = q_i^k$. (Such a place exists uniquely). Now, $F_{k+1} = F_k \cup \{(g'', e), (e, g), (e, g')\}$.

c). $a = r_{ij}$ for some $j \leq N$. By the definition of $\rightarrow, (q_i^k, m) - a \rightarrow q_i^{k+1}$, where $m = \text{getout}(\gamma_{ji}^k)$. Then, $G_{k+1} = G_k \cup \{g\}$, where g is a new place, $l_{k+1}(g) = q_i^{k+1}$. Let g' be the place with $I_k(g') = (j, i, n)$, where n is the smallest number such that g' has no outgoing arcs, g'' be the place in G_k having no outgoing arc such that $l_k(g'') = q_i^k$. (Such places exist uniquely). Now, $F_{k+1} = F_k \cup \{(g'', e), (e, g), (g', e)\}$.

Example 1. Let

$$D = (P_1, P_2, \rightarrow, C, M), M_{12} = \{1, 0\}$$

$$\begin{aligned} P_1 &= (\{q_0, q_1, q_2\}, \rightarrow, \{b\}), & P_2 &= (\{p_0, p_1, p_2\}, \rightarrow, \{a\}) \\ (q_0, 1) - r_{12} &\rightarrow q_1, & p_1 - s_{21} &\rightarrow (p_0, 1), \\ (q_0, 0) - r_{12} &\rightarrow q_2, & p_0 - s_{21} &\rightarrow (p_2, 0), \\ q_0 - a &\rightarrow q_1, & p_0 - b &\rightarrow p_1, \end{aligned}$$

Let

$$\bar{w} = (r_{12} a r_{12}, b s_{21} b s_{21} s_{21}) = (\lambda, b)(\lambda, s_{21})(r_{12}, \lambda)(\lambda, b)(a, \lambda)(\lambda, s_{21})(\lambda, s_{21})(r_{12}, \lambda), q_0 = (q_0, p_0, \Lambda), \text{ where}$$

$$\Lambda = \begin{bmatrix} \lambda & \lambda \\ \lambda & \lambda \end{bmatrix}$$

The occurrence net defined for \bar{w} is represented by the Fig. 1

Lemma 1. $\pi(\bar{w})$ is defined. Furthermore, let

$$S^0 = \{g \in G | \text{has no outgoing arc and } l(g) \in \cup_{i=1}^N Q_i\},$$

$$\Gamma_{ij}^0 = \{g \in B | \text{has no outgoing arc and } l(g) \in M_{ij}\},$$

$a_{ij} = g_1 g_2 \dots g_n$ be the sorting of its elements such that g precedes g' iff $n < n'$, where $I(g) = (i, j, n)$ and $I(g') = (i, j, n')$, $i, j \leq N$.

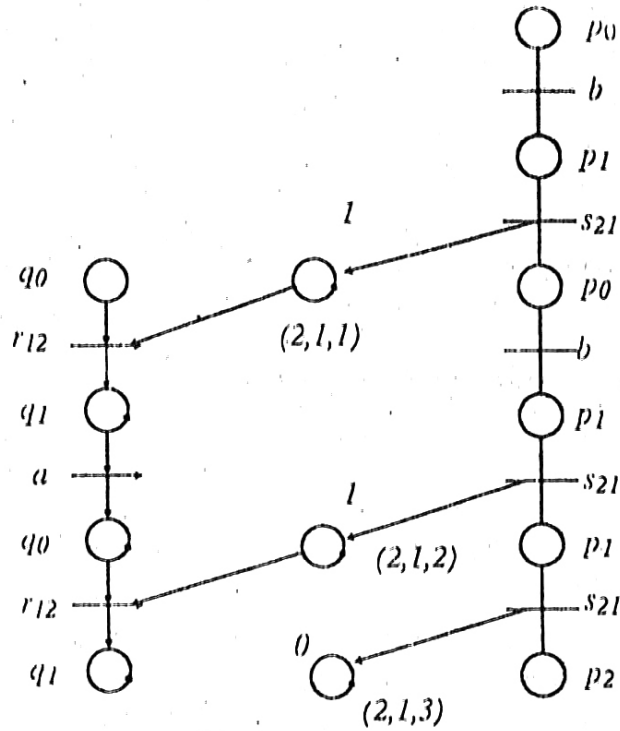


Fig. 1

Then, $q_i = I(S^0) \cap Q_i, \gamma_{ij} = I(\alpha_{ij})$

Lemma 2. There exists uniquely an one-to-one mapping τ between the set of letter occurrences of \bar{w} and the set E of transitions, preserving labels such that:

- for letter occurrences a, b in w_i , a immediately follows b in w_i if and only if there is a place g with label in Q_i such that $(\tau(b), g), (g, \tau(a)) \in F$,
- for letter occurrences a, b , a is the h^{th} occurrence of r_{ij} in w_i , b is the h^{th} occurrence of s_{ji} in w_j for some h if and only if there is a place g in G with label in M_{ji} such that $(\tau(b), g), (g, \tau(a)) \in F$.

The lemmas 1 and 2 can be verified straightforward by induction on k . Details of the proofs are left to the reader.

From Lemmas 1 and 2, it follows that many interesting behavioural aspects of distributed systems can be formulated in terms of vectors of action sequences.

Theorem 1. If distributed system is a deterministic, so is its transition system, i.e.

$$|\{q \in S | p - \bar{w} \rightarrow q\}| \leq 1, \forall p \in S, \bar{w} \in VAS(A).$$

Proof. Since each transition in E represents a transition labelled over A_i ($i = 1, 2, \dots, N$) or C , it follows from Lemma 2 that if the component transition systems in D are deterministic then all occurrence nets constructed for \bar{w} corresponding to its different factorizations into a sequence of

actions of D are isomorphic. Hence, by Lemma 1, the state \bar{q} is determined uniquely independently from the way of factoring \bar{w} .

It follows from Theorem 1 that the synchronized transition system for a deterministic distributed system is a deterministic one, and if $\bar{p} - \bar{w} \rightarrow \bar{q}$, $\bar{w} \in VAS(A_s)$, then \bar{q} does not depend on the way of factoring \bar{w} into a concatenation of synchronous actions.

In practice, parallel programs are frequently designed to be synchronized (e.g., OCCAM, CSP programs), and then are implemented by asynchronous message passing distributed systems. Theorem 1 says that such implementation is correct.

3. Maintaining the amount of global information in local states of process

We will formulate some basic notions such as causality relation, logical time, of the system, and then propose an algorithm for maintaining the amount of global information of the system from which some properties of interest can be derived. The implementations of the algorithm may be expensive, but it shows the upper bound of what can be maintained in the states of processes by sending messages.

For simplicity, we present the above mentioned notions and results only for the case of asynchronous message passing distributed systems only. The results, however, can be modified and extended to the case of synchronous and asynchronous message passing distributed systems as well.

Let, in the sequel, D be a distributed system, $T = (S, \rightarrow, A)$ its labelled transition system with the initial state \bar{q}_0 , B its behaviour.

For $\bar{w} \in B$, an initial part of \bar{w} is defined to be a $\bar{w}' \in B$, such that $\bar{w} = \bar{w}'\bar{u}$ with some $\bar{u} \in VAS(A)$. Let $init(\bar{w})$ denote the set of all initial parts of \bar{w} .

For $\bar{w}' = (w'_1, w'_2, \dots, w'_N), \bar{w}'' = (w''_1, w''_2, \dots, w''_N) \in init(\bar{w})$, denote:

$$\max(\bar{w}', \bar{w}'') = (\max(w'_1, w''_1), \max(w'_2, w''_2), \dots, \max(w'_N, w''_N)),$$

$$\min(\bar{w}', \bar{w}'') = (\min(w'_1, w''_1), \min(w'_2, w''_2), \dots, \min(w'_N, w''_N)),$$

where for any u, v over an alphabet Σ ,

$$\max(u, v) = \begin{cases} u, & \text{if } v \text{ is a prefix of } u, \\ v, & \text{if } u \text{ is a prefix of } v, \\ \text{undefined} & \text{otherwise,} \end{cases}$$

$$\min(u, v) = \begin{cases} u, & \text{if } u \text{ is a prefix of } v, \\ v, & \text{if } v \text{ is a prefix of } u, \\ \text{undefined} & \text{otherwise,} \end{cases}$$

Lemma 3. *With the operations \max, \min defined as above, $init(\bar{w})$ is a complete finite lattice.*

Proof. The lemma follows immediately from the fact that $\bar{w}' = (w'_1, w'_2, \dots, w'_N)$ is an initial part of $\bar{w} = (w_1, w_2, \dots, w_N)$ if and only if:

a) w'_i is a prefix of w_i , $i = 1, 2, \dots, N$ b) $\forall i, j \leq N$, the number of sends to j in w'_i is not greater than the number of receives from i in w'_j .

The details of the proof are omitted.

$init(\bar{w})$ is intended to represent all the possible past parts of the execution resulting \bar{w} of the system.

For a given $\bar{w} = (w_1, w_2, \dots, w_N) \in B$, an event structure can be derived to represent the actions performed by local processors and the causality relation between them. Each letter occurrence in w'_i 's will be called a event. Let $EVENT(\bar{w})$ denote all events in \bar{w} .

Definition 9. Let $a, b \in EVENT(\bar{w})$. $a \leq b$ iff for any $\bar{w}' \in init(\bar{w})$, $b \in EVENT(\bar{w}')$ implies $a \in EVENT(\bar{w}')$.

It follows from Lemma 3 that $a \leq b$ iff the minimal initial part (according to the partial ordering relation generated by the operation \min) of \bar{w} containing b also contains a . Thus, the minimal initial part of \bar{w} containing b contains all events which lead to the occurrence of b .

It can be seen easily that if $\bar{w}' \in init(\bar{w})$, then the poset $(EVENT(\bar{w}'), \leq)$ is a left closed subposet of $(EVENT(\bar{w}), \leq)$.

When a process i can perform an action b during an execution of the system resulting \bar{w} , what it can learn about the system is each predecessor of b has been occurred. Thus, the amount of global information that the processor i can have at this moment is the minimal initial part of \bar{w} containing b . It follows that we can maintain only the knowledge of the system which is derived from this amount of global information. Denote by $K_i(\nu a)$ the minimal initial part of \bar{w} having νa as its i^{th} component which is the minimal initial part of \bar{w} containing a for an event a performed by the process i (a occurs in \bar{w}). When a is a fixed letter occurrence in w_i , we write $K_i(a)$ instead of $K_i(\nu a)$ if it does not cause any confusion.

Proposition 1. Let $w_i = ua'av$, where a is a receive, and b the corresponding send of a in w_j (i.e. $w_j = xby$ and $\#_{r,i}, ua' = \#_{s,j}, x$). Then

$$K_i(a) = \max (K_i(a'), K_j(b)\bar{\lambda}_i(a))$$

where $\bar{\lambda}_i(a)$ denotes the element of A , of which all components is λ except that its i^{th} component is a .

Proof. $\max (K_i(a'), K_j(b))$, is an initial part of \bar{w} containing a' . Since b is the corresponding send of a , if $\bar{q}_0 = \max (K_i(a'), K_j(b)) \rightarrow \bar{q} = (q_1, \dots, q_N, \gamma)$ then $\gamma_{ji} = \lambda$ and the message m generated by b is the first element of γ_{ji} . Thus, $\bar{w}' = \max (K_i(a'), K_j(b))\bar{\lambda}_i(a)$ is in B , and hence, \bar{w}' must be the minimal initial part of \bar{w} containing a . To see why, we suppose in contrast that \bar{w}' is not so. This means that at least one component of $K_i(a)$ must be a proper prefix of \bar{w}' , which is the corresponding component of either $K_i(a')$ or $K_j(b)$. It follows that either $\min (K_i(a'), K_i(a))$ is a

proper initial part of $K_i(a')$ containing a' or $\min(K_i(a), K_j(b))$ is a proper initial part of $K_j(b)$ containing b . This is a contradiction to the definition of either $K_i(a')$ or $K_j(b)$. The proof is complete.

Proposition 1 is a foundation of the algorithm for maintaining the amount of global information of the system in the processes.

In order to keep our algorithm as general as possible, we introduce a function f for representing the knowledge of global properties that can be derived from $K_i(a)$'s. Let P be a partially ordered lattice, f be a function from $init(\bar{w})$ to P satisfying:

$$f(\max(\bar{w}', \bar{w}'')) = \max(f(\bar{w}'), f(\bar{w}''))$$

$$f(\bar{w}'\bar{\lambda})(a) = g_i(f(\bar{w}'), a),$$

where g_i is a given function from $init(\bar{w}) \times A_i$ to P , $i = 1, 2, \dots, N$.

The operation \max represents synthesizing global knowledge of the system of a process, and the function g_i represents the fact that if the process i has the amount of global knowledge K of the system before it performs the action a , then it has the amount of global knowledge $g_i(K, a)$ of the system after it performs the action a .

Let

$$\bar{q}_0 - \bar{w}' \rightarrow \bar{q}', \bar{q}_0 - \bar{w}'' \rightarrow \bar{q}'', \bar{q}' = (q'_1, q'_2, \dots, q'_N, \gamma'), \bar{q}'' = (q''_1, q''_2, \dots, q''_N, \gamma''), \bar{w}', \bar{w}'' \in init(\bar{w}).$$

a).

$$f(\bar{w}') = (|w'_1|, |w'_2|, \dots, |w'_N|), \bar{w}' \in init(\bar{w}), P = \mathbf{N}^N,$$

$$g_i(j_1, j_2, \dots, j_i, \dots, j_N) = (j_1, j_2, \dots, j_i + 1, \dots, j_N).$$

In this case, $f(\bar{w}')$ is considered as partially ordered logical time of the system after the execution represented by \bar{w}' . Then, $f(K_i(a))$ will be a timestamp assigned to a to ensure that $K_i(a) \leq K_j(b)$ iff a is a predecessor of b . An implementation of this can be found in [7].

b).

$$f(\bar{w}') = (|w'_1|, q'_1, C'_1), (|w'_2|, q'_2, C'_2), \dots, (|w'_N|, q'_N, C'_N), \bar{w}' \in init(\bar{w}),$$

where $C'_i = ((n'_{i1}, c'_{i2}), \dots, (n'_{iN}, c'_{iN}))$, n_{ij} is the number of sends from i to j in w'_i , $c'_{ij} = \gamma'_{ij}$ is the queue of messages having not received by j from i after the execution represented by \bar{w}' .

$$g_i(((j_1, q'_1, C'_1), \dots, (j_i, q'_i, C'_i), \dots, (j_N, q'_N, C'_N)), a)$$

$$((j_1, q'_1, C'_1), \dots, (j_i + 1, p_i, C''_i), \dots, (j_N, q'_N, C'_N)),$$

where

$$C''_j = ((n''_{j1}, c''_{j1}), \dots, (n''_{ji}, c''_{ji}), \dots, (n''_{jN}, c''_{jN})), j \leq N,$$

$$(q'_1, \dots, q'_i, \dots, q'_N, \gamma') - \bar{\lambda}_i(a) \rightarrow (q'_1, \dots, p_i, \dots, q'_N, \beta), \beta_{ij} = c''_{ij}, i, j \leq N,$$

$$n''_{ji} = n'_{ij} + 1 \text{ if } a \text{ is } s_{ij} \text{ and } n''_{ij} = n'_{ij} \text{ otherwise.}$$

$f(\bar{w}')$ is considered as the global state of the system after the execution represented by \bar{w}' which is added with message counters and local time clocks. Thus, $f(K_i(a))$ will give a possible global state when the processor i performs the event a . The operation \max in this interpretation is defined as

$$\max (f(\bar{w}'), f(\bar{w}'')) = ((\max (|\bar{w}'_1|, |w''_1|), q_1, C_1), \\ (\max (|w'_2|, |w''_2|), q_2, C_2), \dots, (\max (|w'_N|, |w''_N|), q_N, C_N)),$$

where

$$q_i = \begin{cases} q'_i & \text{if } \max (|w'_i|, |w''_i|) = |w'_i|, \\ q''_i & \text{if } \max (|w'_i|, |w''_i|) = |w''_i|, \end{cases}$$

if $\max (|w'_i|, |w''_i|) = |w'_i|$, $\max (|w'_j|, |w''_j|) = |w'_j|$, then

$$n_{ij} = n'_{ij}, \quad c_{ij} = c'_{ij},$$

if $\max (|w'_i|, |w''_i|) = |w'_i|$, $\max (|w'_j|, |w''_j|) = |w'_j|$, then $n_{ij} = n'_{ij}$, c_{ij} is the concatenation of the first $n'_{ij} - n''_{ij}$ elements of c'_{ij} and c''_{ij} .

The remaining cases of c_{ij} are obtained by changing the roles of \bar{w}' and \bar{w}'' in the above definition.

Now, we can give an algorithm for maintaining $f(K_i(a))$ in the local state of the process i after it has performed the event a . The algorithm is presented as a distributed system D' , which is behaviourally equivalent to D with $f(K_i(a))$ being kept as a component of the local state of the process i after it performs a .

$$D' = (P'_1 P'_2, \dots, P'_N, C, \rightarrow' C, M')$$

where

$$P'_i = (Q_i \times f(VAS(A)), \rightarrow' i', A_i), \quad M'_{ij} = M_{ij} \times f(VAS(A)),$$

$\rightarrow' C'$ and $\rightarrow' i'$ is defined as follows.

Let $q, q' \in Q_i$, $a \in \sum_i \bar{w}, \bar{w}' \in VAS(A)$.

If $a \in A_i$, $q - a \rightarrow \bar{q}'$, then $(q, f(\bar{w}) - a \rightarrow i'(q', g_i(f(\bar{w}, a)))$.

If $a = s_{ij}$, $q - a \rightarrow (q', m)$, then $(q, f(\bar{w}) - a \rightarrow ((q', g_i(f(\bar{w}, a))), (m, g_j(f(\bar{w}, a))))$.

If $a = r_{ij}$, $(m, q) - a \rightarrow q'$, then $((m, f(\bar{w}')), (q', f(\bar{w}')) - a \rightarrow ((q', \max (f(\bar{w}), f(\bar{w}'))))$.

Let $q'_0 = ((q_1^0, f(\bar{\lambda})), (q_2^0, f(\bar{\lambda})), \dots, (q_N^0, f(\bar{\lambda})), \Phi)$ be the initial state of D' .

Theorem 2. D and D' are behaviourally equivalent, i.e. $B = B'$. Furthermore, if $(q, d), q' \in Q_i$, $d \in f(VAS(A))$ is the local state of the process i of D' after it performs the event a in an execution resulting \bar{w} , then $d = f(K_i(a))$.

Proof. Let T and T' be labelled transition systems for D and D' respectively. The proof goes by induction on the length of derivation that

$$\bar{q}_0 - \bar{w} \rightarrow \bar{q} = (q_1, q_2, \dots, q_N, \gamma)$$

if and only if

$$\bar{q}'_0 - \bar{w} \rightarrow' \bar{q}' = ((q_1, f(K_1(\bar{w}_1))), (q_2, f(K_2(\bar{w}_2))), \dots, (q_N, f(K_N(\bar{w}_N))), \gamma'),$$

using Proposition 1 and the property of the function f , where $\gamma_{ij} = h(\gamma'_{ij})$ with h being the extended homomorphism of the projection on the first component.

By choosing a suitable interpretations of f , D' can be used in debugging distributed programs as in [14,15], studying invariant properties and reasoning of the systems [8], distributed snapshots [4], etc. With f being the global state function, a global state of D can be determined in only one local process. If the determined global state is too old to be considered, some of local states of some another processes may be accessed as well. Some of these uses of the algorithm will be presented in separate papers.

Example 2. Let D, \bar{w} be as in Example 1, f and $g_i, i = 1, 2$ be defined in b). Since only the process 2 sends messages to process 1, a channel state γ can be written as a queue of messages from the process 2 to the process 1 (matrix γ contains one element, and hence, $C_i, i = 1, 2$ contains one element). Then, the derivation represented by \bar{w} in D is as follows.

$$\begin{aligned}
& ((q_0, ((0, q_0, (0, \lambda)), (0, p_0, (0, \lambda))), (p_0, ((0, q_0, (0, \lambda)), (0, p_0, (0, \lambda))), \lambda) - (\lambda, b) \rightarrow \\
& ((q_0, ((0, q_0, (0, \lambda)), (0, p_0, (0, \lambda))), (p_1, ((0, q_0, (0, \lambda)), (1, p_1, (0, \lambda))), \lambda) - (\lambda, s_{21}) \rightarrow \\
& ((q_0, ((0, q_0, (0, \lambda)), (0, p_0, (0, \lambda))), (p_0, ((0, q_0, (0, \lambda)), (2, p_0, (1, 1))), \\
& (1, ((0, q_0, (0, \lambda)), (2, p_0, (1, 1)))) - (r_{12}, b) \rightarrow \\
& ((q_1, ((1, q_1, (0, \lambda)), (2, p_0, (1, \lambda))), (p_0, ((0, q_0, (0, \lambda)), (3, p_1, (1, 1))), \lambda) - (a, s_{21}) \rightarrow \\
& ((q_0, ((2, q_0, (0, \lambda)), (2, p_0, (1, \lambda))), (p_1, ((0, q_0, (0, \lambda)), (4, p_0, (2, 11))), \\
& (1, ((0, q_0, (0, \lambda)), (4, p_1, (2, 11)))) - (r_{12}, s_{12}) \rightarrow \\
& ((q_1, ((3, q_1, (0, \lambda)), (4, p_1, (2, \lambda))), (p_2, ((0, q_0, (0, \lambda)), (5, p_2, (3, 110))), \\
& (0, ((0, q_0, (0, \lambda)), (5, p_2, (3, 110))))).
\end{aligned}$$

From the state $((q_1, ((3, q_1, (0, \lambda)), (4, p_1, (2, \lambda))))$ of the process 1, we know that one global state of D is (q_1, p_1, λ) .

Conclusion.

The vectors of sequential words generated by distributed labelled transition systems are an intuitive model for describing the behaviour of concurrent and distributed systems. The notions of causality relation, knowledge of global information of processes, logical time are formulated accurately in term of this model of behavior. It has been proved to be a good model for proving the equivalence and verification of distributed systems. Thus, an extension of the model to the case of systems of the processes created dynamically and an algebra of processes in this model should be studied in support of design and understanding of concurrent and distributed systems.

References

1. America H.M. & Rutten J.J.M.M., A Parallel Object-Oriented Language: Design and Semantic Foundations, Languages for Parallel Architectures, J.W. Bakker, Ed., Wiley Series in Parallel Computing, 1989, 1-50.
2. Bal H.E.; Steiner J.G. & Tanebaum A.S., *Programming Languages for Distributed Computing Systems*, *ACM Computing Surveys*, Vol. 21, (1989), 261-322.
3. Casavant T.L. & Kuhl J.G., *A Communicating Finite Automata Approach to Modeling Distributed Computation and Its Application to Distributed Decision-Making*, *IEEE Transactions on Computers*, Vol. 39, (1990), 628-639.
4. Chandy K.M. & Lamport L., *Distributed Snapshots: Determining the Global State of Distributed Systems*, *Ann. Rev. Computer Science*, Vol. 2, (1987), 37-68.
5. Degano P. & Montanari U., Distributed Systems, Partial Orderings of Events, and Event Structures, Control Flow and Data Flow: Concepts of Distributed Programming, M. Broy, ed., NATO ASI Series, Vol. F14, 1985, 7-106.
6. Degano P., Gorrieri R. & Marchetti S., *An Exercise in Concurrency: A CSP Process as a Condition/Event System*, *Lecture Notes in Computer Science*, Vol. ?, 1989; 85-105.
7. Fridge C., *Logical Time in Distributed Computing Systems*, *Computer*, Vol. 24, (1991), 28-33.
8. Halpern J.I., *Using Reasoning about Knowledge to Analyze Distributed Systems*, *Ann. Rev. Computer Sci.* 2, (1987), 37-68.
9. Hoare C.A.R., *Communicating Sequential Processes*, *CACM*, 21, 1978, 666-677.
10. Hung D.V., *A Model for Analyzing Distributed Systems*, *Journal of Informatics and Cybernetics*, Vol. 1, No. 2, 1991, 15-23 (in Vietnamese).
11. Hung D.V., *Notes on Projection Products, Trace Languages and Synthesized Concurrent Computation Systems*; *MTA-SZTAKI Közlemenyek*, Vol. 32, (1985), 87-104.
12. Hung D.V. & Knuth E., *Semi-Commutations and Petri Nets*, *Theoretical Computer Science*, n^o. 64, (1989), 67-81.
13. Hung D.V. & Knuth E., *A Non-interleaving Semantics for Communicating Sequential Processes: A fixed point approach*, *Acta Cybernetica*, Szeged, Hungary, Tom. 8, f.3, (1988), 293-311.
14. Joseph T.A., Rauchle T. & Toueg S., *State Machines and Assertions: Integrated Approach to Modelling and Verification of Distributed Systems*, *Science of Computer Programming*, Vol. 7, (1987), 1-22.
15. Keller R., *Formal Verification of Parallel Programs*, *CACM*, 19, 1976, 561-572.
16. Mazurkiewicz A. Ochmanski E. & Peucek W., *Concurrent Systems and Inevitability*, *Theoretical Computer Science*, n^o. 64, (1989), 281-304.
17. De Nicola R., *Extensional Equivalences for Transition Systems*, *Acta Informatica*, 24, (1987), 211-237.
18. Nielsen M., Plotkin G. & Winskel G., *Petri Nets, Event Structures and Domains*, *Theoretical Computer Science*, 13, (1981), 85-108.

19. Shield M.W., Nonsequential Behaviour, Part I, Tech. Rept.CSR-120-82 Dept. of Computer Science, University of Edinburgh, 1982.

Abstract**Labelled Transition System Approach to the Distributed Computing Systems**

The present paper introduces the notion of distributed transition systems for modelling, designing and understanding distributed computing systems. The concurrency can be expressed explicitly in the model. Some of the global properties of the systems are discussed and determined. It is shown in the paper that by keeping knowledge of other processes in each process of a system, some of its global properties can be synthesized from only few local process states.