

# A framework for modeling the distributed deployment of synchronous designs\*

Luca P. Carloni · Alberto L. Sangiovanni-Vincentelli

Published online: 3 May 2006  
© Springer Science + Business Media, LLC 2006

**Abstract** Synchronous specifications are appealing in the design of large scale hardware and software systems because of their properties that facilitate verification *and* synthesis. When the target architecture is a *distributed system*, implementing a synchronous specification as a synchronous design may be inefficient in terms of both size (memory for software implementations or area for hardware implementations) and performance. A more elaborate implementation style where the basic synchronous paradigm is adapted to distributed architectures by introducing elements of asynchrony is, hence, highly desirable. Building on the tagged-signal model, we present a modeling for the distributed deployment of synchronous design. We offer a comparative exposition of various design approaches (synchronous, asynchronous, GALS, latency-insensitive, and synchronous programming) and we provide some insight on the role of signal absence in modeling synchronization in distributed concurrent systems. Finally, we compare two distinct methodologies, desynchronization and latency-insensitive design, and we elaborate on possible options to combine their results.

**Keywords** Desynchronization · GALS · Distributed systems · Latency-insensitive design

## 1. Introduction

The synchronous design paradigm is pervasive in electronic system engineering. It is used in discrete-time dynamical control systems, it is the basis of digital integrated circuit design, and it is the foundation of programming languages and design environments used for

---

\*This research was supported in part by the NSF under the project ITR (CCR-0225610), the Gigascale System Research Center (GSRC), and by the European Commission under the projects COLUMBUS, IST-2002-38314, and ARTIST, IST-2001-34820.

L. P. Carloni (✉)  
Department of Computer Science, Columbia University, New York, NY 10027-6902, USA  
e-mail: luca@cs.columbia.edu

A. L. Sangiovanni-Vincentelli  
EECS Department, University of California at Berkeley, Berkeley, CA 94720-1770, USA  
e-mail: alberto@eecs.berkeley.edu

software development for real-time embedded systems. In this paradigm, a complex system is represented as a collection of interacting modules whose state is updated collectively in one *zero-time* step. A synchronous specification is naturally simpler than specifying the same system as the interaction of components whose state is updated following an intricate set of time-based interdependency relations. However, for an increasing number of important applications, e.g., transportation systems, sensor networks, and industrial control, the implementation architecture is distributed. In addition, the advent of deep-submicron (DSM) technologies for IC design, where hundreds of millions of transistors can be integrated on a single die, is making the synchronous paradigm very expensive to implement since the chip becomes a distributed system with interconnect delays that are up to an order of magnitude larger than the switching delays of the gates and that are very difficult to estimate [13]. In this scenario, we believe that new methodologies that combine specification simplicity with implementation constraints will take center place in the design stage.

The *desynchronization problem* can be informally described as the task of deploying a synchronous design on a distributed architecture in a *correct-by-construction* (and mostly automatic) fashion. The relevance of this problem follows naturally from the desire of leveraging the well-known tools and practices of synchronous design for the specification and optimization of a system [5], while targeting efficient final implementations that are distributed in nature.

We present a modeling framework for distributed deployment that focuses on the synchronization aspects of system design. Our framework encompasses different design styles from the “strong assumptions” of synchronous design and asynchronous design, to more “relaxed and realistic” models for distributed design, like globally asynchronous locally synchronous (GALS) systems [15]. We argue for the importance of the notions of absence to distinguish (and relate) these systems, and we illustrate their interplay in modeling the desynchronization problem. Finally, we revisit previous work on distributed embedded code generation (desynchronization) and latency-insensitive design. These are two distinct approaches that share the goal of conjugating the theoretical properties of synchronous designs with the efficiency of implementations where the constraints imposed by synchrony are relaxed:

- *Desynchronization* was motivated by the problem of reaching a correct-by-construction modular deployment of embedded software on distributed architectures [3, 4, 24]. In a desynchronized implementation processes that compose the large scale system are implemented synchronously while their communication is implemented in an asynchronous style. This approach allows also running each process at its own “speed”.
- *Latency-insensitive design* was motivated by the problem of deriving hardware implementations in the presence of long interconnect delays [11, 12]. In fact, with DSM technologies, the long paths between the design components may introduce delays that force the overall clock of the system to run too slow in order to maintain synchronous behavior. By enabling automatic pipelining of long interconnect, latency-insensitive design allows the implementation to avoid slowing down the clock and to “recover” at least part of the throughput that could have been achieved with communication delays of the same order of the clock of the components.

Besides comparing the two approaches, we show how the causality analysis of the desynchronization approach can be applied to deriving more efficient implementations of latency-insensitive protocols.

## 2. Related work

Benveniste et al. formally define the desynchronization problem in the context of embedded system applications [2, 3]. Their main motivation is to address the issue of compositionality of synchronous languages and enable modular code generation for distributed architectures. In particular, they advocate a methodology centered on the use of the synchronous paradigm for system specification and validation followed by a provably correct desynchronization step to derive a distributed implementation (e.g. on GALS architectures). Building on the first works, in [24] they study the relationship between synchrony and asynchrony and identify the key property that must be satisfied by the communication architecture in a desynchronized implementation: each of its channels must act as a lossless queue or, in other words, “messages shall not be lost and shall be delivered in order”. Interestingly enough, this property is also at the core of the theory of latency-insensitive protocols.

Recently, desynchronization has been the object of investigation of several projects in the areas of embedded systems design and integrated circuit design. A mathematical framework to support the composition of heterogeneous reactive systems is presented in [4] together with a set of theorems supporting the automatic generation of correct-by-construction adapters between heterogeneous designs. The idea is applied to the deployment of synchronous design on GALS architectures and LTTA architectures [6]. Desynchronization approaches targeting hardware design have been presented both by Jacobson et al. [20] and Cortadella et al. [16]: the basic idea is to start from a fully synchronous synthesized (or manually designed) integrated circuit, and then replace the global clock network with a set of local handshaking circuits. The advantage of the desynchronized circuit with respect to the corresponding synchronous circuit is the removal of the clock tree. This gives important benefits in terms of power dissipation, electro-magnetic interference control, and robustness to temperature-related and manufacturing-related variations.

The Polychrony project aims to support design refinement from the early stages of requirement specification to the later stages of synthesis and deployment [18]. The term *polychrony* denotes the capability of describing circuits and systems using the synchronous assumption together with multiple clocks. This can be applied to abstracting the key properties of a system (while completing a high-level specification) as well as to describing the characteristics of the components that can be used to implement it. The concept of polychrony is used in [22, 25] to address the formal validation of the refinement of synchronous multi-clocked designs into GALS architectures.

Latency-insensitive protocols were proposed in [11] and, then, applied to synchronous hardware design in [10, 12]. A complete presentation of latency-insensitive design is given in [9], which includes a detailed discussion of the analysis and optimization of latency-insensitive systems.

## 3. The tagged-signal model

We summarize here the main concepts of Lee and Sangiovanni-Vincentelli’s (LSV) tagged-signal model [21], the basis of our formal framework.

Given a set of *values*  $\mathcal{D}$  and a set of *tags*  $\mathcal{T}$ , an *event* is a member of  $\mathcal{D} \times \mathcal{T}$ . A *signal*  $s$  is a set of events. The set of all  $M$ -tuples of signals is denoted  $S^M$  and a *process*  $P$  is a subset of  $S^M$ . An  $M$ -tuple  $b = (s_1, \dots, s_M) \in S^M$  is called a *behavior* of process  $P$  when it *satisfies* the process, i.e. when  $b \in P$ . Thus, a process is a set of possible behaviors. A *tagged system* is a composition of processes  $\{P_1, \dots, P_I\}$ , i.e. a new process  $P$  that is defined as the

	<i>tag</i>	:		$t_0$		$t_1$		$t_2$		$t_3$		$t_4$		$t_5$		$t_6$		$t_7$		...
$P$ :	$u$	:		4		-2		5		$\perp$		-1		3		4		2		...
	$v$	:		1		$\perp$		1		$\perp$		$\perp$		1		1		1		...

Fig. 1 Presence and absence of events in two signals

intersection of their behaviors  $P = \bigcap_{i=1}^I P_i$ . To distinguish signals, we assume an underlying set  $\mathcal{V}$  of variables with domain  $\mathcal{D}$ . We denote a tagged system as a triple  $P = (V, \mathcal{T}, \mathcal{B})$ , where  $V \subset \mathcal{V}$  is a finite set of variables,  $\mathcal{T}$  is a tag set, and  $\mathcal{B}$  a set of behaviors with domain  $V$ . The *composition* of two systems  $P_1$  and  $P_2$  is given by the intersection of their behaviors:

$$P_1 \cap P_2 =_{\text{def}} (V_1 \cup V_2, \mathcal{T}_1 \cup \mathcal{T}_2, \mathcal{B}_1 \cap \mathcal{B}_2), \text{ where}$$

$$\mathcal{B}_1 \cap \mathcal{B}_2 =_{\text{def}} \{b \mid b|_{V_i} \in \mathcal{B}_i, i = 1, 2\},$$

and  $b|_W$  denotes the restriction of  $b$  to a subset  $W$  of variables. We denote with  $\mathcal{T}(s)$  the tag set of signal  $s$  (and, similarly, for a behavior and a process).

In some models of computation the set  $\mathcal{D}$  includes a special value  $\perp$ , which indicates the absence of a value. For any tag  $t \in \mathcal{T}$ , we call  $(t, \perp)$  the *absent-value event*, or simply,  $\perp$  event. We say that a signal  $s$  is *present* at a given tag  $t$  when  $(\exists e = (t, d) \in s, (d \neq \perp))$ ; otherwise, we say that  $s$  is *absent* at  $t$  (we further elaborate on this point in the sequel).

*Example 1.* The diagram of Fig. 1 represents the unique<sup>1</sup> behavior of a system that has two signals with names  $u, v$ . At any given tag, signal  $u$  is present with unit value if and only if signal  $v$  is present and carries a positive integer value.

*Ordering among signal tags.* We assume that for any tag  $t \in \mathcal{T}$ , each signal  $s$  in the system has at most one event, i.e.:

$$\forall b \in \mathcal{B}, \forall s \in b, \neg[\exists e_1 \in s, \exists e_2 \in s, (tag(e_1) = tag(e_2))]$$

This assumption naturally leads to the definition of a total order  $<$  among the tags of a *signal*. Then, the total order over the tag set  $\mathcal{T}(s)$  of signal  $s$  induces a total order among its events. Therefore, a signal can be seen as a sequence of events. We use  $t_i$  to denote the  $i$ -th tag of a signal and, naturally, we rely on the fact that  $t_i < t_j \Leftrightarrow i < j$ . Further, we can use tags to identify an event of a signal (somewhat like the indexes of an array) as well as its values. Given a signal  $s$  and a tag  $t$ , we write  $e = eve(s, t)$  to denote the event of  $s$  whose tag is  $t$  and we write  $d = val(s, t)$  to denote the value of  $eve(s, t)$ .

$\Sigma$  denotes the set of all sequences of elements in  $\mathcal{D} \cup \{\perp\}$ . Function  $\sigma : \mathcal{S}^1 \times \mathcal{T}^2 \rightarrow \Sigma$  takes a signal  $s = \{(d_0, t_0), (d_1, t_1), \dots\}$  and an ordered tag pair  $(t_i, t_j), i \leq j$ , and returns a sequence  $\sigma_{[t_i, t_j]} \in \Sigma$  s.t.  $\sigma_{[t_i, t_j]}(s) = d_i, d_{i+1}, \dots, d_j$ . The sequence of values of a signal is denoted  $\sigma(s)$ . The empty sequence is denoted as  $\epsilon$ . To manipulate sequences of values we

<sup>1</sup> We used the *unique* attribute to express the fact that this system has exactly one behavior, i.e. the one illustrated in the figure.

define the filtering operator  $\mathcal{F}_\perp : \Sigma \rightarrow \mathcal{D}$  that returns a sequence  $\sigma' = \mathcal{F}_\perp[\sigma]$  s.t.

$$\sigma'_i = \begin{cases} \sigma_{[t_i, t_i]}(s) & \text{if } \sigma_{[t_i, t_i]}(s) \in \mathcal{D} \\ \epsilon & \text{if } \sigma_{[t_i, t_i]}(s) = \perp \end{cases}$$

*Ordering among process tags.* In general, the tag set  $\mathcal{T}$  of a process is not totally ordered, but only partially ordered. When the partial order is the identity relation, the tag set is an unordered set. When tags are used to express causality relations among signals, it is common to assume that  $\mathcal{T}$  is partially ordered. In this case,  $\leq$  is used to denote the partial order on  $\mathcal{T}$  by writing  $t < t'$  when  $t \leq t'$  and  $t \neq t'$ . Finally, a tag system is *timed* if  $\mathcal{T}$  is a totally ordered set, i.e. for each pair of distinct tags  $t, t'$  either  $t < t'$  or  $t' < t$ .

Often tags are used as a mechanism to express time (as in this in paper). This may be useful, for instance, to move across the various representations of a design at different levels of abstraction from initial specification, where *logical time* is central, to final implementation, where each event occurs at a given instant of the *physical*, or *real*, time. However, the reader should notice that tags are essentially a tool to express constraints, like coordination constraints, among events of different signals (and, transitively, among signals and among processes). In [4], morphisms among tag sets are used to handle semantic-preserving transformations of synchronous designs.

#### 4. Models of computation

We use models of computation to specify the mathematical behavior of the systems under design [17]. The models of computation addressed in this paper fall under the category of synchronous, asynchronous, and *in between* to indicate models that are neither.

*Synchronous systems.* Two events  $e_1, e_2$  are *synchronous* ( $e_1 \approx e_2$ ) when they have the same tag, i.e.  $e_1 \approx e_2 \Leftrightarrow \text{tag}(e_1) = \text{tag}(e_2)$ . Two signals  $s_1, s_2$  are synchronous ( $s_1 \approx s_2$ ) when they share the same tag set, i.e.:

$$s_1 \approx s_2 \Leftrightarrow \mathcal{T}(s_1) = \mathcal{T}(s_2)$$

The definitions of two synchronous behaviors  $b_1, b_2$  and two synchronous processes  $P_1 = (V_1, \mathcal{T}_1, \mathcal{B}_1), P_2 = (V_2, \mathcal{T}_2, \mathcal{B}_2)$  naturally follow:

$$b_1 \approx b_2 \Leftrightarrow \forall s_i \in b_1, \forall s_j \in b_2 (s_i \approx s_j)$$

$$P_1 \approx P_2 \Leftrightarrow \forall b_i \in \mathcal{B}_1, \forall b_j \in \mathcal{B}_2, (b_i \approx b_j)$$

A stand-alone behavior  $b$  is synchronous when  $b \approx b$ . A stand-alone process  $P$  is synchronous when  $P \approx P$ . In a behavior of a synchronous system, every signal is synchronous with every other signal and, equivalently, for each tag a signal has *exactly one* corresponding event:

$$\forall b \in \mathcal{B}, \forall s \in b, \forall t \in \mathcal{T} (\exists! e \in s (\text{tag}(e) = t)).$$

tag :		t <sub>0</sub>	t <sub>1</sub>	t <sub>2</sub>	t <sub>3</sub>	t <sub>4</sub>	t <sub>5</sub>	t <sub>6</sub>	t <sub>7</sub>	...
P :	w :	1	0	1	0	1	0	1	0	...
	y :	0	2	2	6	6	10	10	14	...
	z :	0	0	4	0	8	4	12	8	...
Q :	w :	1	0	1	0	1	0	1	0	...
	x :	1	3	5	7	9	11	13	15	...
	y :	0	2	2	6	6	10	10	14	...
R :	w :	1	0	1	0	1	0	1	0	...
	x :	1	3	5	7	9	11	13	15	...
	z :	0	0	4	0	8	4	12	8	...

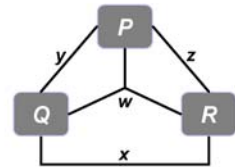


Fig. 2 The synchronous system of Example 2 and its behavior

Example 2. The diagram of Fig. 2 represents the unique behavior of a synchronous system that is the result of the composition of three processes *P*, *Q*, and *R*. Signal *w*, a binary, is shared by all processes, while the remaining signals, integers *x*, *y*, and *z*, are shared in pairwise manner. In Fig. 2, the signals are purposely represented by simple lines and not arrows. In fact, by observing only the event sequences we can not say which input/output relations exist among the system processes. However, in the sequel, we focus our attention on functional systems [21] and we use this example assuming that signal *w* is produced by process *P*, signals *x*, *y* by process *Q*, and signal *z* by process *R*.

### 4.1. Synchronous languages

Synchronous programming languages like ESTEREL, LUSTRE, and SIGNAL represent powerful tools for the specification of complex real-time embedded systems because they allow to combine the simplicity of the synchronous assumption with the power of concurrency in functional specification [7, 8, 19]. They are synchronous systems with particular properties and for this reason, they are often considered a model of computation in addition to the generic synchronous model. The *synchronous programming model* can be expressed by the following “pseudo-mathematical” statements [3, 5]:

$$P \equiv R^\omega$$

$$P_1 || P_2 \equiv (R_1 \wedge R_2)^\omega$$

where *P*, *P*<sub>1</sub>, *P*<sub>2</sub> denote synchronous programs, *R*, *R*<sub>1</sub>, *R*<sub>2</sub> denote the sets of all the possible reactions of the corresponding programs, and the superscript  $\omega$  indicates non-terminating iterations. The first expression interprets the essence of the synchronous assumption: a synchronous program *P* evolves according to an infinite sequence of successive atomic reactions. At each reaction, the program variables may or may not present a value. The second expression defines the parallel composition of two components as the conjunction of the reactions for each component. This implies that communication among components is performed via instantaneous broadcast. To cast the synchronous programming model into the LSV formalism, we naturally associate signals to variables and use tags to index the program reactions. An important feature offered by the synchronous programming model is the ability of taking decisions based on the absence of a value for a variable at a given reaction, i.e., *in synchronous systems absence can be sensed*. This is perfectly in line with the definition of the

absent-value event since processes react to events and hence can also react to the particular absent-value event. The absent-value event plays an important role in synchronous models of computation. In fact, the essence of the model is that all computation processes awake simultaneously when any of them posts an event for communication. Some of the signals that connect the processes may be not present. The synchronous model requires that these signals be read with the absent-value event posted. If indeed the information on the presence of an absent-value event does not cause a process to react to it, then reading this event is an unnecessary complication. We shall see later that recognizing this situation and eliminating the associated steps are key in deriving a more effective deployment that, while formally giving up the synchronous model, maintains behavior equivalence with the original synchronous specifications.

The notion of *clock of a variable* is introduced as a *Boolean meta-variable* that it is implicitly defined in order to track the absence/presence of a value for each corresponding variable.<sup>2</sup> Variables that are always present simultaneously are said to have the same clock, so that clocks can be seen as *equivalence classes of simultaneously-present variables*. In the sequel, we focus our attention on SIGNAL, which is a declarative language [7]. Besides parallel composition, SIGNAL's main operators are the followings:

- statement “ $c := a \text{ op } b$ ”, where *op* denotes a generic logic or arithmetic operator, defines not only that the values of *c* are function of those of *a* and *b*, but also that the three variables have the same clock;
- statement “ $c := a \S k$ ”, where *k* is a positive integer constant, specifies both that *c* and *a* have the same clock and that at the *n*-th reaction when the two signals are present, the value of *c* is equal to the value held by *a* at the (*n* − *k*)-th reaction;
- statement “ $c := a \text{ default } b$ ” specifies that variable *c* is present at every reaction where either *a* or *b* is present while taking the value of *b* only if *a* is not present (*oversampling*);
- statement “ $c := a \text{ when } b$ ” specifies that variable *c* is present (taking the value of *a*) only when both *a* is present and the Boolean condition expressed by variable *b* is true (*undersampling*).

While the first two statements are *single-clock*, the last two are *multi-clock*. Additional operators are available to directly relate the variable clocks: for instance, statement  $c \wedge = a$  constraints variables *c* and *a* to have the same clock, without relating the values that they assume. The SIGNAL compiler uses *clock calculus* to statically analyze every program statement, identify the structure of the clock of each variable, and schedule the overall computation. The compiler rejects the program when it detects that the collection of its statements as a whole contains clock constraint violations.

*Example 3.* Figure 3 reports the code of a SIGNAL program that is structured as a main process with three sub-processes *P*, *Q*, and *R*. These processes communicate via signals *w*, *x*, *y*, *z* that are constrained to be synchronized (first statement of the main process). Hence, using SIGNAL jargon, these signals belong to the same *clock equivalence class* [7], which is also the class of signal *tag*. Signal *tag* is an external input whose values evolve as an infinite alternating sequence of 0s and 1s. Under this assumption, a run of program MAIN returns deterministically a tuple of sequences of values for variables *w*, *x*, *y*, *z* that coincide with

<sup>2</sup> Notice that despite its name the clock of a variable is not necessarily a periodic signal in hardware terms. Rather, it is more like an enable signal that gates a single global clock.

<pre> <b>process</b> MAIN (   ? <b>boolean</b> tag;   ! <b>boolean</b> w, x, y, z;) (      x <math>\wedge</math> = y <math>\wedge</math> = z <math>\wedge</math> = w <math>\wedge</math> = tag     w := P(tag, y, z)     (x,y) := Q(tag, w)     z := R(tag, w, x)    ) ;         </pre>	<pre> <b>process</b> P (   ? <b>boolean</b> tag; <b>integer</b> y, z;)   (    ! <b>integer</b> w;) (     i <math>\wedge</math> = tag     i := (i\$1 <b>init</b> (-1)) + 1     iW := true <b>when</b> (i &lt; 1)     w := iW <b>default</b> (y\$1 &gt; z\$1)    ) <b>where</b> <b>integer</b> i, iW; <b>end</b>;         </pre>
<pre> <b>process</b> Q (   ? <b>boolean</b> tag; <b>integer</b> w;   ! <b>integer</b> x, y;) (      i <math>\wedge</math> = tag     i := (i\$1 <b>init</b> (-1)) + 1     iY := 0 <b>when</b> (i &lt; 1)     y := iY <b>default</b> (     <b>if</b> w\$1 <b>then</b> (x\$1+1) <b>else</b>     (x\$1-1))     x := (x + 2) \$1 <b>init</b> 1    ) <b>where</b> <b>integer</b> i, iY; <b>end</b>;         </pre>	<pre> <b>process</b> R (   ? <b>boolean</b> tag; <b>integer</b> w, x;   ! <b>integer</b> z;) (      i <math>\wedge</math> = tag     i := (i\$1 <b>init</b> (-1)) + 1     iZ := 0 <b>when</b> (i &lt; 2)     z := iZ <b>default</b> (     <b>if</b> w\$1 <b>then</b> x\$2 - 3 <b>else</b> x\$2     + 3)    ) <b>where</b> <b>integer</b> i, iZ;   <b>end</b>;         </pre>

Fig. 3 SIGNAL program with a deterministic behavior as in Example 2

the behavior of the synchronous system of Example 2. By analyzing the program we derive the functional relationships between its signals: e.g., we see that  $y$  and  $z$  are input signals for process  $P$ , which produces output signal  $w$ . Also, we learn causality dependencies among signals like, for instance, that every event of signal  $w$ , besides the first, depends on the events of  $y$  and  $z$  occurred at the previous reaction. Similarly, while the first two events of  $z$  carry values equal to 0, each subsequent event depends on the event occurred on  $w$  at the previous reaction as well as on the event occurred on  $x$  two reactions earlier. Hence, events of  $w$  depend on events of  $z$  and vice versa. In fact, cyclic causality dependencies across signals of a synchronous program are quite common and may be problematic only in the presence of a *combinational cycle*, i.e. when two events with the same tag depend on each other. The discussion of methods to handle this issue goes beyond the scope of this paper (see [5]).

*Asynchronous systems.* The definition of asynchrony as used in the literature is vague: some use the term to indicate any systems that is *not* synchronous, others are more restrictive. According to [21], two events  $e_1, e_2$  are *asynchronous* ( $e_1 \simeq e_2$ ) if they have different tags, i.e.  $e_1 \simeq e_2 \Leftrightarrow tag(e_1) \neq tag(e_2)$ . Two signals  $s_1, s_2$  are asynchronous ( $s_1 \simeq s_2$ ) when they have disjoint tag sets, i.e.:

$$s_1 \simeq s_2 \Leftrightarrow (\mathcal{T}(s_1) \cap \mathcal{T}(s_2) = \emptyset)$$

The definitions of asynchronous behaviors  $b_1, b_2$  and asynchronous processes  $P_1 = (V_1, \mathcal{T}_1, \mathcal{B}_1), P_2 = (V_2, \mathcal{T}_2, \mathcal{B}_2)$  follow:

$$b_1 \simeq b_2 \Leftrightarrow \forall s_i \in b_1, \forall s_j \in b_2, s_i \neq s_j, (s_i \simeq s_j)$$

$$P_1 \simeq P_2 \Leftrightarrow \forall b_i \in \mathcal{B}_1, \forall b_j \in \mathcal{B}_2, b_i \neq b_j, (b_i \simeq b_j)$$



	tag :	t <sub>0</sub>	t <sub>1</sub>	t <sub>2</sub>	t <sub>3</sub>	t <sub>4</sub>	t <sub>5</sub>	t <sub>6</sub>	t <sub>7</sub>	t <sub>8</sub>	t <sub>9</sub>	t <sub>10</sub>	t <sub>11</sub>	t <sub>12</sub>	t <sub>13</sub>	...
P <sub>a</sub> :	w <sub>a</sub> :	1				0			2		1				0	...
	y <sub>a</sub> :			0								2				...
	z <sub>a</sub> :				0				0				4			...
Q <sub>a</sub> :	w <sub>a</sub> :	1				0				1					0	...
	x <sub>a</sub> :		1				3				5					7
	y <sub>a</sub> :			0				2				2				...
R <sub>a</sub> :	w <sub>a</sub> :	1				0				1					0	...
	x <sub>a</sub> :		1				3				5					7
	z <sub>a</sub> :				0				0				4			...

Fig. 4 The behavior of an asynchronous system

A stand-alone behavior  $b$  is asynchronous when  $b \simeq b$ . A stand-alone process  $P$  is asynchronous when  $P \simeq P$ . In a behavior of an asynchronous system, every signal is asynchronous with every other signal and, equivalently, for each tag there is one and only one event across all signals:

$$\forall b = (s_1, \dots, s_M) \in \mathcal{B}, \forall t \in \mathcal{T}, \left( \exists! e \in \bigcup_i s_i, (tag(e) = t) \right)$$

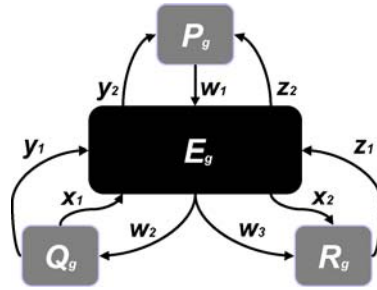
Example 4. Figure 4 illustrates the unique behavior of the asynchronous system  $S_a = P_a \cap Q_a \cap R_a$ . Processes  $P_a$ ,  $Q_a$ , and  $R_a$  communicate by sharing signals (as in the case for synchronous systems), but signals do not share tags.

*Between synchronous and asynchronous: Globally-asynchronous locally-synchronous (GALS) systems.* Formally, the set of asynchronous systems is *not* the complement of the set of synchronous systems. In fact, there is a set of systems that sits *in between* these two sets and whose elements are useful to model heterogeneous systems and distributed architectures. An element of this *in-between set* is a process with a behavior that has both at least a pair of synchronous events (hence, it is not asynchronous) and at least a tag for which a signal does not present a corresponding event while another does (hence, it is not synchronous). A relevant subset of this set is the class of GALS systems. GALS systems are of particular interest because they represent a compromise that allows designers to leverage the traditional practices and tools of synchronous design for implementations of synchronous processes on distributed architectures. In a GALS system, computation occurs in synchronous clusters exchanging data asynchronously via a set of communication media. Each cluster runs with its own clock that controls also the sampling of new values for its input signals. At each sampling period, some of these new values may or may not be present, depending on the transferring latencies in the asynchronous communication media. Since we want to focus on the communication mechanisms at the interface between synchronous and asynchronous, our LSV definition for GALS systems assumes, without loss of generality, that all asynchronous communications can be modeled as occurring within a single media process. A GALS system  $S_g = \bigcap_{P_i \in \mathcal{P}} P_i \cap E$  is the composition of a collection  $\mathcal{P}$  of *computation processes* and one *communication, or media, process*  $E = (V_e, \mathcal{T}_e, \mathcal{B}_e)$  s.t.:

$$\forall P_i, P_j \in \mathcal{P}, ((i = j \Rightarrow P_i \approx P_j) \wedge (i \neq j \Rightarrow P_i \simeq P_j)), \quad \text{and}$$

$$\forall P_i = (V_i, \mathcal{T}_i, \mathcal{B}_i), \forall b \in \mathcal{B}_e, (b|_{V_i} \in \mathcal{B}_i)$$

**Fig. 5** GALS system for Example 5



Each computation process is a stand-alone synchronous process because it runs with its own logical clock whose occurrences are represented by tags. In the general case, we assume that no relation exists between the clocks of distinct computation processes leaving total freedom in the implementation process. This is captured by saying that the intersection of their tag sets is empty (i.e., they are pairwise asynchronous processes). Instead, a media process is not synchronous (because it models the communication latency and the sharing of communication resources among processes that are pairwise asynchronous) nor asynchronous (because each subset of its signals that interfaces a specific computation process is a synchronous sub-process). Hence, from a LSV perspective, the name “GALS” is justified when considering the system from the viewpoint of any of its computation processes.

*Example 5.* The diagram of Fig. 6 reports the unique behavior of the GALS system  $S_g$  of Fig. 5, which is the result of the composition of processes  $P_g$ ,  $Q_g$ ,  $R_g$ , and  $E_g$ . Process  $P_g$  is synchronous because all its signals are synchronous. The same is true for processes  $Q_g$  and  $R_g$  separately. However, the composition of processes  $P_g$ ,  $Q_g$ , and  $R_g$  is not a synchronous process (no tag is shared across signals of different processes). Observe that  $P_g$  has a period twice as fast as those of  $Q_g$  and  $R_g$ . Processes  $P_g$ ,  $Q_g$ , and  $R_g$  communicate only via the media process  $E_g$ . Process  $E_g$ , acting as the communication environment, has subsets of its signals synchronized with signals of the other processes, but, as a stand-alone process,

	tag :	$t_0$	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$	$t_7$	$t_8$	$t_9$	$t_{10}$	$t_{11}$	$t_{12}$	$t_{13}$	...
$P_g$ :	$w_1$ :	1		⊥		⊥		0		⊥		1		⊥		...
	$y_2$ :	⊥		0		⊥		2		⊥		⊥		...		...
	$z_2$ :	⊥		⊥		0		⊥		0		⊥		4		...
$Q_g$ :	$w_2$ :		1				⊥				0				1	...
	$x_1$ :		1				3				⊥				5	...
	$y_1$ :		0				2				⊥				2	...
$R_g$ :	$w_3$ :				1				0					1		...
	$x_2$ :				1				3				⊥	...		...
	$z_1$ :				0				0				4			...
$E_g$ :	$w_1$ :	1		⊥		⊥		0		⊥		1		⊥		...
	$w_2$ :		1				⊥				0				1	...
	$w_3$ :				1				0				1			...
	$x_1$ :		1				3			⊥					5	...
	$x_2$ :				1				3				⊥			...
	$y_1$ :		0				2			⊥			⊥		2	...
	$y_2$ :	⊥		0		⊥		2		⊥		⊥		...		...
	$z_1$ :			0		0			0				4			...
$z_2$ :	⊥		⊥		0		⊥		0		⊥		4		...	

**Fig. 6** The behavior of a GALS system

it is not synchronous. The signals of  $E_g$  can be partitioned in equivalence classes, whose members carry the same sequence of values at different tags (e.g., signal  $x_2$  is a “delayed” version of signal  $x_1$ ). We call this relation semantic equivalence.

Our proposal for modeling a distributed system with the LSV model is to use more than one signal to capture each communication thread between two processes. For instance, if process  $P_g$  sends data to process  $Q_g$ , we must be able to distinguish between the sending event and the receiving event. To do so, we need at least two signals, e.g.  $w_1$  and  $w_2$ . Each new event of  $w_1$  is created by  $P_g$ , whose overall activity of reading input events and computing output events proceeds according to its tag set  $\mathcal{T}(P_g)$ . Then, a new event of  $w_1$  causes at least a corresponding event of  $w_2$  within the media process (more events could be necessary to model arbitrary latencies or the sharing of communication resources). Finally, event  $w_2$  is consumed by  $Q_g$ , whose activity is controlled by tag set  $\mathcal{T}(Q_g)$  that has empty intersection with the tag set of every other synchronous processes, including  $\mathcal{T}(P_g)$ . In synchronous systems, signal decoupling is not necessary thanks to the power of the synchronous abstraction: all processes create and sample events at the same tags and a unique signal  $w$  is sufficient to express the *instantaneous* communication,<sup>3</sup> between process  $P$  and process  $Q$  (see Example 2). Strictly asynchronous systems rely on an abstraction that is equally powerful: there is no notion of global (i.e., system) or local (i.e., process) tag set and two processes communicate by sharing signals that are produced and sampled independently from the rest of the communication and computation activities in the system. The sharing of a signal in asynchronous systems represents the presence of an *ad-hoc* handshaking communication protocol, which is dedicated to the particular communication of, say,  $w_a$  from  $P_a$  to  $Q_a$ : a new event for  $w_a$  is created by  $P_a$  only when  $Q_a$  is ready to sample it (see Example 4).

*The role of event absence.* In the GALS system  $S_g = (V_g, \mathcal{T}_g, \mathcal{B}_g)$  of Example 5, for any  $t \in \mathcal{T}_g$  and  $s \in V_g$ , one of three things can happen:

1.  $t \notin \mathcal{T}(s)$  (event absence)
2.  $t \in \mathcal{T}(s) \wedge \exists e = (t, d) \in s (d = \perp)$  (value absence)
3.  $t \in \mathcal{T}(s) \wedge \exists e = (t, d) \in s (d \neq \perp)$  (presence)

From a global viewpoint, a GALS system is a system with multiple tag sets (a multi-clock system), each representing a dimension that is familiar to a synchronous process and extraneous to all the remaining synchronous processes. For instance, the signals of process  $P_g$  do not have events at tag  $t_1$  and the signals of processes  $Q_g$  and  $R_g$  do not have events at tag  $t_0$ . However, process  $P_g$  does not “expect” an event at tag  $t_1$  nor at tags  $t_3, t_5, \dots$  because these are tags that do not belong to the tag set of  $P_g$ : they represent instants of a time dimension that is completely extraneous to  $P_g$ . The meaning of the absence that  $P_g$  detects on signal  $y_2$  at  $t_4$ , which is a tag belonging to  $\mathcal{T}(P_g)$ , is different. In fact, at tag  $t_4$ ,  $P_g$  is looking for a significant event, but it ends up sampling the absent-value event (the awaited event will arrive only at tag  $t_6$ , after being created by process  $Q_g$ ) because of the latency introduced by the communication medium. This case is typical of a GALS system, because, in a purely synchronous model, the absent value event is always an “expected” event. In this case, the computation can be affected in a substantial way since  $P_g$  can compute on the absence value event and produce an output that is different from the one originally expected.

<sup>3</sup> Instantaneous has to be interpreted in the sense of a process that is not “seen” by the computation part of the system. Communication and computation in synchronous systems never overlap.

Deploying automatically a synchronous design on a distributed architecture entails the development of techniques for making each process robust with respect to absence. In other words, we ask under which conditions we can guarantee that sensing the absent value event when a different value was expected does not produce incorrect behavior or that *not* sensing an absent value event when one is expected, does not change the behavior of the system. Section 5 discusses methods to achieve this result.

*Remark .* Consider again the case of asynchronous design (see Example 4). By definition, the tag sets of any two asynchronous signals  $x_a, y_a$  are disjoint. Thus, for each tag in  $\mathcal{T}(x_a)$  there is no corresponding event in signal  $y_a$  and vice versa. If we consider an asynchronous system with several signals, we have that for each event that is present in one of its signal, there are absences in all the remaining signals. This phenomenon is so inherent to the notion of asynchronous system that here the notion of event absence loses its meaning: when events are always systematically absent, there is no point in looking for their presence! In fact, in asynchronous systems no common references exist across processes and processes cannot (and do not attempt) to detect event absence: inter-process communication occurs according to handshake protocols that don't leave space for this kind of uncertainty.

### 5. Deploying synchronous design on distributed architectures

In this section, we revisit previous research that targets the implementation of a synchronous design on a distributed architecture both in software and in hardware. We introduce the definition of semantic equivalence, which provides a formal ground to establish when the behavior of a distributed implementation conforms to the one of the synchronous specification. Then, we summarize the theory of latency-insensitive design and we present the main results on the desynchronization of synchronous programs for distributed embedded code generation. Finally, we compare these two approaches and we sketch a possible research avenues to combine them.

*Semantic equivalence.* With the notion of semantic equivalence between processes we capture the fact that their operations produce exactly the same result from the viewpoint of an observer watching the sequences of values of their signals.

Two signals are *semantic equivalent* if they have the same sequence of values after discarding the  $\perp$  events. This is written:  $s \equiv s' \Leftrightarrow \mathcal{F}_\perp[\sigma(s)] = \mathcal{F}_\perp[\sigma(s')]$ . Two behaviors  $b = (s_1, \dots, s_M), b' = (s'_1, \dots, s'_M)$  are semantic equivalent  $b \equiv b'$  when there exists a bijective map  $\psi : b \mapsto b'$  s.t.  $\forall i, (s_i \equiv \psi(s'_i))$ . Finally, for two processes  $P = (V, \mathcal{T}, \mathcal{B}), P' = (V, \mathcal{T}', \mathcal{B}')$  we have:  $P \equiv P' \Leftrightarrow [(\forall b \in \mathcal{B}, \exists b' \in \mathcal{B}', (b \equiv b')) \wedge (\forall b' \in \mathcal{B}', \exists b \in \mathcal{B}, (b' \equiv b))]$ .

Similarly to *flow equivalence* [18], semantic equivalence indicates that two behaviors have exactly the same sequence of present events, which, however, may be interleaved by a different number of  $\perp$  events. Hence, we can use it to model implicitly the communication latency between processes: e.g., the communication of events from  $P$  to  $Q$  occurs via media process  $E$  and involves several signals  $u_p, u_{e_1}, u_{e_2}, \dots, u_q$  that belong all to the same class of semantic equivalence. Observe that semantic equivalence doesn't say anything about tags: processes  $P, P'$  can be semantic equivalent even if  $\mathcal{T}(P) \cap \mathcal{T}(P') = \emptyset$ . Finally, it is a compositional property: if two pairs of processes are semantic equivalent, their pairwise intersections give semantic equivalence processes, i.e.  $(P \equiv P' \wedge Q \equiv Q') \Rightarrow (P \cap Q \equiv P' \cap Q')$ .

*Example 6.* Reconsidering the systems of Examples 2 and 4, we have that  $\psi(w) = w_a$ ,  $\psi(x) = x_a$ ,  $\psi(y) = y_a$ ,  $\psi(z) = z_a$  and, consequently:  $P \equiv P_a$ ,  $Q \equiv Q_a$ ,  $R \equiv R_a$  and, finally,  $P \cap Q \cap R \equiv P_a \cap Q_a \cap R_a$ . Now, consider the GALs system of Example 5. Semantic equivalence models the communication among computation processes  $P_g$ ,  $Q_g$  and  $R_g$  via media process  $E_g$ : e.g.,  $\{w_1, w_2, w_3\}$  is a semantic equivalence class representing the communication from  $P_g$  to both  $Q_g$  and  $R_g$ . The other equivalence classes are  $\{x_1, x_2\}$ ,  $\{y_1, y_2\}$ , and  $\{z_1, z_2\}$ . Let  $V^* = \{w_1, x_1, y_1, z_1\}$  be the set of representative variables for each equivalence class and for all behaviors  $b \in P_g \cap Q_g \cap R_g$  let  $b^* = b|_{V^*}$ . Then,  $\psi'(w_1) = w$ ,  $\psi'(x_1) = x$ ,  $\psi'(y_1) = y$ ,  $\psi'(z_1) = z$ , and:  $(P_g \cap Q_g \cap R_g)|_{V^*} \equiv P \cap Q \cap R$ .

*Latency-insensitive systems.* Latency-insensitive protocols [11, 12] were originally proposed to address the interconnect delay problem in synchronous hardware design. The latency-insensitive design methodology takes a *strict synchronous* system specification and automatically derives a *latency-equivalent synchronous* implementation. This implementation formally operates as a synchronous system, but, practically, does it in a looser fashion that makes it robust with respect to arbitrary, but discrete, latency variations between its processes.

A key element of a latency-insensitive protocols is the concept of *stalling event* (or,  $\tau$  event), i.e. an event carrying as value the special symbol  $\tau$ , as opposed to an *informative event*, i.e. an event carrying a value  $d$  in accordance with the original specification. A stalling events is the *modeling unit* to represent explicitly latency variations in inter-process communication, while remaining within the boundaries of the synchronous model. Hence, we augment domain  $\mathcal{D}$  with  $\tau$ , while  $\Sigma_{lat}$  denotes the set of all sequences of elements in  $\mathcal{D} \cup \{\tau\}$ . A *strict* signal  $s$  is always present and contains only informative events:  $\forall t \in \mathcal{T}(s), (val(s, t) \notin \{\perp, \tau\})$ . A *stalled* signal  $s$  contains at least one  $\tau$  event:  $\exists t \in \mathcal{T}(s) (val(s, t) = \tau)$ . Similarly to the definition of operator  $\mathcal{F}_\perp$ , we define operator  $\mathcal{F}_\tau : \Sigma_{lat} \rightarrow \Sigma$  as follows:

$$\sigma'_i = \begin{cases} \sigma_{[t_i, t_i]}(s) & \text{if } \sigma_{[t_i, t_i]}(s) \neq \tau \\ \epsilon & \text{otherwise} \end{cases}$$

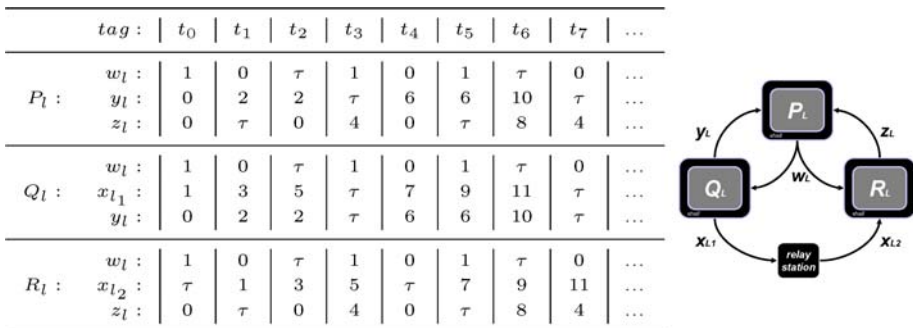
Two signals  $s, s'$  are *latency-equivalent* if they have the same sequence of values after discarding the  $\tau$  events, i.e.  $s \equiv_\tau s' \Leftrightarrow \mathcal{F}_\tau[\sigma(s)] = \mathcal{F}_\tau[\sigma(s')]$ . As for semantic-equivalence, the definition extends to behaviors and processes. The main result of latency-insensitive design follows:

**Theorem 1 ([12]).** *If  $S = \bigcap_i P_i$  is a strict synchronous system and  $S' = \bigcap_i P'_i$  is a system of patient processes s.t.  $\forall i (P'_i \equiv_\tau P_i)$  then  $S' \equiv_\tau S$ .*

Informally, a patient process is able to wait an arbitrary amount of reactions for an informative event to occur at any of its inputs and, when this occurs, it reacts as if no wait happened. Hence, patience, which is compositional, enables the compensation of any arbitrary latency in inter-process communication. While patience is generally a strong requirement to demand, each *stallable* component (*core* process) can be made patient by generating proper interface logic (*shell* process) implementing the latency-insensitive protocol [12]. A component is stallable if it can freeze its internal state for an arbitrary amount of time. In hardware systems, this property can be obtained through a *gated-clock mechanism*. At each reaction, the presence of

a  $\tau$  event at one of the input ports of a shell/core pair means that the expected informative event is not ready yet and will be delayed for at least another reaction. Hence, the shell logic reacts by stalling the internal core, while emitting new stalling events on the outputs and saving the informative events of its other input signals on internal buffers.<sup>4</sup> Once all the informative events for that reaction finally arrive, the interface reactivates the internal core, which produces new informative events. It is important to notice that until *all* informative events for a given reaction arrive, the shell logic does *not* reactivate the core process, regardless of its internal computational structure. In fact, the shell logic ignores this structure and sees the core simply as a *black box* component. This conservative approach is a consequence of the assumption, influenced by single-clock hardware design, that the original system specification is strictly synchronous (a signal never presents a  $\perp$  event at any reaction). In the sequel, we discuss how to lift this assumption to extend the application of latency-insensitive design to multi-clock and software systems.

*Example 7.* The application of latency-insensitive design to integrated circuits provides two main advantages [10]: (a) it enables the *a-posteriori* automatic pipelining of long wires by insertion of special patient processes called *relay stations*; (b) it facilitates the assembly of pre-designed components, that, as long as they are stallable, can be interfaced to the communication protocol without changing their internal structure. Assume that each process of Example 2 is implemented as a distinct finite state machine (FSM) on an integrated circuit and that the wire carrying signal  $x$  from  $Q$  to  $R$  is the only one that has been pipelined by introducing one relay station. Figure 7 reports the structure and the behavior of the resulting latency-insensitive system after each process has been made patient by wrapping it within a shell. The system is strictly synchronous (no absent-value  $\perp$  events occur at any tag). At the system initialization, the only stalling event is the one at the output of the relay station because a relay station represents a “design correction” that is extraneous to the original system specification, while each FSM is properly initialized according to it. As the behavior evolves, new stalling events are generated whenever a process must wait for a new informative event at its inputs. In fact, since the system is cyclic [9],  $\tau$  events occur periodically on each signal at the rate of 1/4.



**Fig. 7** The latency-insensitive system of Example 7 and its behavior

<sup>4</sup> A discussion on how to detect possible buffer overflow (as well as how to avoid it by introducing *back-pressure* in the protocol together with optimum buffer sizing) is given in [9].

*Desynchronization of synchronous programs.* The behavior of a synchronous system can be seen as a sequence of tuples of events with each tuple indexed by a tag. This is not the case for an asynchronous or a GALS system: the most one can say is that a behavior, being a tuple of signals, is a tuple of sequences of events. In [2, 3, 24], *desynchronization* is defined as the act of discarding the synchronization barriers that delimit successive reactions in a synchronous program. Since this corresponds to filtering away the absent-value events, desynchronization amounts to mapping a sequence of tuples of values in domains extended with the extra value  $\perp$  into a tuple of sequences of present values, one sequence per each variable. The desynchronization abstraction is relevant because it provides a formal way to characterize those synchronous programs that can be deployed on a distributed architecture without losing their semantics. As proven in [3], the notions of endochrony and isochrony are sufficient for this characterization.

Let  $clk(s, t)$  be a Boolean function denoting whether signal  $s$  presents an event at tag  $t$  or no, i.e.  $(clk(s, t) = 1 \Leftrightarrow val(s, t) \neq \perp)$ . For any process  $P = (V, \mathcal{T}, \mathcal{B})$ , any tag  $t \in \mathcal{T}$ , and any pair of subsets  $W, W'$  s.t.  $W \subset W' \subseteq V$ , we say that  $W$  *tag-determines*  $W'$  at  $t$ , written  $W \rightarrow_t W'$ , when:

$$\forall b \in \mathcal{B}, \forall s \in b_{|W'}, (\exists \phi : D'_W \rightarrow \{1, 0\}, (clk(s, t) = \phi(D'_W)))$$

where  $D'_W$  is the set of values  $val(s, t)$  for  $s \in b_{|W}$ . Since relation  $W \rightarrow_t W'$  is stable by union over  $W'$  sets, there is a maximal  $W'$  that is tag-determined by  $W$  at  $t$ . Thus, for any tag  $t \in \mathcal{T}$ , if  $V$  is a finite set then there is a maximal chain  $\emptyset = W_0 \rightarrow_t W_1 \rightarrow_t W_2 \rightarrow_t \dots \rightarrow_t W_{\max}$ . A process  $P = (V, \mathcal{T}, \mathcal{B})$  is *endochronous* when  $\forall t \in \mathcal{T}, (W_{\max} = V)$ . In other words: a process is endochronous when for each tag of its behaviors the presence/absence of event on all its signals can be inferred incrementally from the values carried by a subset of them that are guaranteed to be present at this tag.

Two processes  $P = (V, \mathcal{T}, \mathcal{B}), P' = (V', \mathcal{T}', \mathcal{B}')$  are *isochronous* when for each behaviors  $b \in \mathcal{B}$  there is a behavior  $b' \in \mathcal{B}'$  (and vice versa) s.t.:

$$\forall t \in \mathcal{T}, (W_t^* \neq \emptyset \Rightarrow W_t^* = W)$$

where  $W = V \cap V'$  and  $W_t^* = \{s \in W \mid val(b, s, t) = val(b', s, t) \neq \perp\}$ . In other words: two processes are isochronous when, at each tag, if there is a pair of shared signals that are present and agree on the event value, then for each other pair of shared signals, either they are present and agree on their value or they are absent.

Endochrony and isochrony allow the derivation of a key result for the automatic generation of distributed embedded code [3]: *if each process of a synchronous program is endochronous and all processes are pairwise isochronous, then deploying the program on a GALS architecture gives a semantic-equivalent implementation.* For any process  $P(V, \mathcal{T}, \mathcal{B})$ , let  $\alpha(P)$  denote the semantic equivalent asynchronous process that is constructed from  $P$  by: (1) applying transformation  $\mathcal{F}_\perp[\sigma(s)]$  to each signal  $s \in b$ , for all  $b \in \mathcal{B}$ , and (2) properly adding tags to each event of  $\alpha(P)$  s.t.  $\forall s, s' \in b, (\mathcal{T}(s) \cap \mathcal{T}(s') = \emptyset)$ .

**Theorem 2 ([3]).** *Let  $\bigcap_i P_i$  be a synchronous system s.t. each  $P_i$  is endochronous and each pair  $(P_i, P_j)$  is isochronous. Then  $\bigcap_i \alpha(P_i) = \alpha(\bigcap_i P_i)$ .*

Endochrony and isochrony can be expressed in terms of transition relations, are model checkable and synthesizable.<sup>5</sup> The SIGNAL compiler handles the parallelism specified using synchronous composition by organizing the computation of signals as a collection of hierarchical trees (the *clock hierarchy forest*) based on the relations among their clocks. Hence, in practice, a synchronous program can be made endo-isochronous by adding a set of “Boolean guard” variables and a master clock to transform this forest into a tree. This approach is somewhat equivalent to synthesize an inter-process communication protocol and carries a drawback: there is not unique solution, or, in other words, there are many possible communication protocols [3].

*Latency-insensitive design and endo-isochrony.* Several commonalities between the work on synchronous program desynchronization and latency-insensitive design have been already pointed out in [1]. Here, we return on the topic to understand how the two approaches could be combined.

**Theorem 3 ([14]).** *The processes of a latency-insensitive system satisfy the properties of endochrony and isochrony.*

The previous result should not come as a surprise if we reflect on the intrinsic differences between  $\tau$  events and  $\perp$  events: the former is used to maintain the semblance of synchronous behavior while the latter represents the lack of it. In other words, a  $\tau$  event tells the process that “the awaited value is not ready yet,” whereas a  $\perp$  event says “there is no value to wait for.” Hence, the  $\tau$  mechanism accounts for the arbitrary latency of interprocess communication while enforcing a synchronous behavior for the distributed latency-insensitive system. Consequently, as illustrated by Example 7,  $\tau$  events never leave the systems (unless for the particular case of acyclic systems) and a price in performance may be ultimately paid [9]. Instead, one may wonder whether it is possible to derive an alternative endo-isochronous implementation that doesn’t rely on latency-insensitive design. In theory, this is certainly possible but in practice it may be challenging and not necessarily advantageous. This depends on knowing the inner structure of each process in the system. In the case of Example 7 the result would be positive because the analysis of the causality dependencies among the events shows that the first two events of output signal  $z$  do not depend on the first event of input signal  $x_{l_2}$ , the delayed event (see also Example 3). Hence, the endo-isochronous implementation would be able to *absorb* the  $\perp$  event (which it would see instead of the  $\tau$  event seen in the corresponding latency-insensitive system) and the behavior would progress without further event absences.

## 6. Concluding remarks

We used the tagged-signal model together with a simple “running example” to provide a comparative view of various design approaches: synchronous, asynchronous, GALS, latency-insensitive, and synchronous programming. In particular, we gave a new formalization of GALS systems and we studied the interplay among the concepts of event absence, event sampling, and communication latency in modeling distributed concurrent systems. Finally, we

<sup>5</sup> A problem with endochrony and isochrony, i.e. the lack of compositionality, is discussed in [23] where new “weak” versions of these concepts are proposed in order to address it.



presented a new comparison of latency-insensitive design and synchronous program desynchronization. The main operational difference between these approaches can be expressed as follows: the former knows how to handle *black box* processes but does not know how to analyze/exploit *white box* ones (that are treated uniformly as if they were black box processes); the latter does not know how to handle black box processes and must analyze the inner structure of each white box process in the system (as well as the properties of each communicating pair), but it is clever in exploiting the information resulting from this analysis. These reflections naturally lead to consider a new research avenue for extending our work on latency-insensitive design: by analyzing the internal structure of each process that comes as a white box module, we can learn its input/output causality dependencies and see if they allow us to *absorb* some  $\tau$  events at given states.

**Acknowledgments** The collaboration with Albert Benveniste and Paul Caspi has been essential in deriving the considerations made in this paper. We gratefully acknowledge their contributions to our on-going discussions on specification and implementation of distributed systems. We gratefully thank Luciano Lavagno for his precious feedback.

## References

1. Benveniste A (2001) Some synchronization issues when designing embedded systems from components. In: Henzinger T, Kirsch C (eds) Embedded Software. Proc. of the First Intl. Workshop, EMSOFT 2001, vol. 2211 of LNCS, pp. 32–49
2. Benveniste A, Caillaud B, Guernic PL (1999) From synchrony to asynchrony. In: Baeten J, Mauw S. (eds) CONCUR'99, vol. 1664 of LNCS, pp. 162–177
3. Benveniste A, Caillaud B, Guernic PL (2000) Compositionality in dataflow synchronous languages: Specification & distributed code generation. Inform Comput 163:125–171
4. Benveniste A, Carloni LP, Caspi P, Sangiovanni-Vincentelli AL (2003a) Heterogeneous reactive systems modeling and correct-by-construction deployment. In: Alur R, Lee I (eds) Proc. of the Third Intl. Conf. on Embedded Software (EMSOFT), vol. 2855 of LNCS, pp. 35–50
5. Benveniste A, Caspi P, Edwards S, Halbwachs N, Guernic PL, de Simone R (2003b) The synchronous language twelve years later. Proc. of the IEEE 91(1):64–83
6. Benveniste A, Caspi P, Guernic PL, Marchand H, Talpin JP, Tripakis S (2002) A protocol for loosely time-triggered architectures. In: Sifakis J, Sangiovanni-Vincentelli A (eds) Embedded Software. Proc. of the Second Intl. Workshop, EMSOFT 2002., vol. 2491 of LNCS, pp. 32–49
7. Benveniste A, Guernic PL (1990) Hybrid dynamical systems theory and the signal language. IEEE Trans Automatic Control 5:535–546
8. Berry G (2000) The foundations of estereel. MIT Press
9. Carloni LP (2004) Latency-insensitive design. Ph.D. thesis, University of California Berkeley, Electronics Research Laboratory, College of Engineering, Berkeley, CA 94720. Memorandum No. UCB/ERL M04/29
10. Carloni LP, McMillan KL, Saldanha A, Sangiovanni-Vincentelli AL (1999a) A methodology for “correct-by-construction” latency insensitive design. In: Proc. Intl. Conf. on Computer-Aided Design, pp. 309–315
11. Carloni LP, McMillan KL, Sangiovanni-Vincentelli AL (1999b) Latency insensitive protocols. In: Halbwachs N, Peled D. (eds) Proc. of the 11th Intl. Conf. on Computer-Aided Verification, vol. 1633. Trento, Italy, pp. 123–133
12. Carloni LP, McMillan KL, Sangiovanni-Vincentelli AL (2001) Theory of latency-insensitive design. IEEE Trans Comp-Aided Design of Integ Cir Syst 20(9):1059–1076
13. Carloni LP, Sangiovanni-Vincentelli AL (2002) Coping with latency in SOC design. IEEE Micro 22(5):24–35
14. Carloni LP, Sangiovanni-Vincentelli AL (2003) A formal modeling framework for deploying synchronous designs on distributed architectures. In: FMGALS 2003: First Intl. Workshop on Formal Methods for Globally Asynchronous Locally Synchronous Architectures, pp. 11–31
15. Chapiro DM (1984) Globally-asynchronous locally-synchronous systems. Ph.D. thesis, Stanford University
16. Cortadella J, Kondratyev A, Lavagno L, Sotiriou C (2003) A concurrent model for de-synchronization. In: Proc. Intl. Workshop on Logic Synthesis, pp. 294–301

17. Edwards S, Lavagno L, Lee E, Sangiovanni-Vincentelli A (1997) Design of embedded systems: Formal methods, validation and synthesis. *Proc. of the IEEE* 85(3):266–290
18. Guernic PL, Talpin JP, Lann JCL (2003) Polychrony for system design. *J Cir, Syst Comp* 12(3):261–303
19. Halbwachs N, Caspi P, Raymond P, Pilaud D (1991) The synchronous data flow programming language LUSTRE. *Proc. of the IEEE* 79(9):1305–1320
20. Jacobson H, Kudva P, Bose P, Cook P, Schuster S, Mercer E, Myers C (2002) Synchronous interlocked pipelines. In: 8th Intl. Symp. on Asynchronous Circuits and Systems
21. Lee EA, Sangiovanni-Vincentelli A (1998) A framework for comparing models of computation. *IEEE Trans Comp-Aided Design Integ Cir Syst* 17(12):1217–1229
22. Mousavi M, Guernic PL, Talpin J, Shukla SK, Basten T (2004) Modeling and validating globally asynchronous design in synchronous frameworks. In: *Proc. European Design and Test Conf.*, pp. 384–389
23. Potop-Butucaru D, Caillaud B, Benveniste A (2004) Concurrency in synchronous systems. In: *Fourth Intl. Conf. on Application of Concurrency to System Design*
24. Talpin JP, Benveniste A, Caillaud B, Guernic PL (1999) Hierarchical normal form for desynchronization. Technical Report 3822, IRISA
25. Talpin JP, Guernic PL, Shukla SK, Gupta R, Doucet F (2004) Formal refinement-checking in a system-level design methodology. *Fundamenta Informaticae*, pp. 243–273