



Scheduling with Bus Access Optimization for Distributed Embedded Systems

Eles, Petru; Doboli, Alex; Pop, Paul; Peng, Zebo

Published in:
IEEE Transactions on VLSI Systems

Link to article, DOI:
[10.1109/92.894152](https://doi.org/10.1109/92.894152)

Publication date:
2000

Document Version
Publisher's PDF, also known as Version of record

[Link back to DTU Orbit](#)

Citation (APA):
Eles, P., Doboli, A., Pop, P., & Peng, Z. (2000). Scheduling with Bus Access Optimization for Distributed Embedded Systems. *IEEE Transactions on VLSI Systems*, 8(5), 472-491. <https://doi.org/10.1109/92.894152>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Scheduling with Bus Access Optimization for Distributed Embedded Systems

Petru Eles, *Member, IEEE*, Alex Doboli, *Student Member, IEEE*, Paul Pop, and Zebo Peng, *Member, IEEE*

Abstract—In this paper, we concentrate on aspects related to the synthesis of distributed embedded systems consisting of programmable processors and application-specific hardware components. The approach is based on an abstract graph representation that captures, at process level, both dataflow and the flow of control. Our goal is to derive a worst case delay by which the system completes execution, such that this delay is as small as possible; to generate a logically and temporally deterministic schedule; and to optimize parameters of the communication protocol such that this delay is guaranteed. We have further investigated the impact of particular communication infrastructures and protocols on the overall performance and, specially, how the requirements of such an infrastructure have to be considered for process and communication scheduling. Not only do particularities of the underlying architecture have to be considered during scheduling but also the parameters of the communication protocol should be adapted to fit the particular embedded application. The optimization algorithm, which implies both process scheduling and optimization of the parameters related to the communication protocol, generates an efficient bus access scheme as well as the schedule tables for activation of processes and communications.

Index Terms—Communication synthesis, distributed embedded systems, process scheduling, real-time systems, system synthesis, time-triggered protocol.

I. INTRODUCTION

MANY embedded systems have to fulfill strict requirements in terms of performance and cost efficiency. Emerging designs are usually based on heterogeneous architectures that integrate multiple programmable processors and dedicated hardware components. New tools that extend design automation to system level have to support the integrated design of both the hardware and software components of such systems.

During synthesis of an embedded system the designer maps the functionality captured by the input specification on different architectures, trying to find the most cost-efficient solution that, at the same time, meets the design requirements. This design process implies the iterative execution of several allocation and partitioning steps before the hardware and software components of the final implementation are generated. The term “hardware/software cosynthesis” is often used to denote this

system-level synthesis process. Surveys on this topic can be found in [1]–[6].

An important characteristic of an embedded system is its performance in terms of timing behavior. In this paper, we concentrate on several aspects related to the synthesis of systems consisting of communicating processes, which are implemented on multiple processors and dedicated hardware components. In such a system, in which several processes communicate with each other and share resources like processors and buses, scheduling of processes and communications is a factor with a decisive influence on the performance of the system and on the way it meets its timing constraints. Thus, process scheduling has to be performed not only for the synthesis of the final system but also as part of the performance estimation task.

Optimal scheduling, in even simpler contexts than that presented above, has been proven to be an NP complete problem [7]. Thus, it is essential to develop heuristics that produce good quality results in a reasonable time. In our approach, we assume that some processes can only be activated if certain conditions, computed by previously executed processes, are fulfilled [8], [9]. Thus, process scheduling is further complicated since at a given activation of the system, only a certain subset of the total amount of processes is executed, and this subset differs from one activation to the other. This is an important contribution of our approach because we capture both the flow of data and that of control at the process level, which allows a more accurate and direct modeling of a wide range of applications.

Performance estimation at the process level has been well studied in the last years. Papers like [10]–[16] provide a good background for derivation of execution time (or worst case execution time) for a single process. Starting from estimated execution times of single processes, performance estimation and scheduling of a system consisting of several processes can be performed. Preemptive scheduling of independent processes with static priorities running on single-processor architectures has its roots in [17]. The approach has been later extended to accommodate more general computational models and has also been applied to distributed systems [18]. The reader is referred to [19] and [20] for surveys on this topic. In [21], performance estimation is based on a preemptive scheduling strategy with static priorities using rate monotonic analysis. In [22], an earlier deadline first strategy is used for nonpreemptive scheduling of processes with possible data dependencies. Preemptive and nonpreemptive static scheduling are combined in the cosynthesis environment described in [23] and [24].

Several research groups have considered hardware/software architectures consisting of a single programmable processor

Manuscript received August 15, 1999; revised February 18, 2000.

P. Eles, P. Pop, and Z. Pang are with the Department of Computer and Information Science, Linköping University, Sweden (e-mail: petel@ida.liu.se; paupo@ida.liu.se; zebpe@ida.liu.se).

A. Doboli is with the Department of Electrical and Computer Engineering and Computer Science, University of Cincinnati, Cincinnati, OH 45221 USA (e-mail: adoboli@ececs.uc.edu).

Publisher Item Identifier S 1063-8210(00)09504-4.

and an application-specific integrated circuit acting as a hardware coprocessor. Under these circumstances, deriving a static schedule for the software component is practically reduced to the linearization of a dataflow graph with nodes representing elementary operations or processes [25]. In the Vulcan system [26], software is implemented as a set of linear threads that are scheduled dynamically at execution. Linearization for thread generation can be performed both by exact, exponential complexity, algorithms and by faster urgency-based heuristics. Given an application specified as a collection of tasks, the tool presented in [27] automatically generates a scheduler consisting of two parts: a static scheduler that is implemented in hardware and a dynamic scheduler for the software tasks running on a microprocessor.

Static cyclic scheduling of a set of data-dependent software processes on a multiprocessor architecture has been intensively researched [28]. Several approaches are based on list scheduling heuristics using different priority criteria [29]–[32] or on branch-and-bound algorithms [33], [34]. In [35] and [36], static scheduling and partitioning of processes, and allocation of system components, are formulated as a mixed integer linear programming (MILP) problem. A disadvantage of this approach is the complexity of solving the MILP model. The size of such a model grows quickly with the number of processes and allocated resources. In [37], a formulation using constraint logic programming has been proposed for similar problems.

It is important to mention that in all the approaches discussed above, process interaction is only in terms of dataflow. This is the case also in [38], where a two-level internal representation is introduced: control-dataflow graphs for operation-level representation and pure dataflow graphs for representation at process level. The representation is used as a basis for derivation and validation of internal timing constraints for real-time embedded systems. In [39] and [40], an internal design representation is presented that is able to capture mixed data/control flow specifications. It combines dataflow properties with finite-state machine behavior. The scheduling algorithm discussed in [39] handles a subset of the proposed representation. Timing aspects are ignored and only software scheduling on a single processor system is considered.

In our approach, we consider embedded systems specified as interacting processes, which have been mapped on an architecture consisting of several programmable processors and dedicated hardware components interconnected by shared buses. Process interaction in our model is not only in terms of dataflow but also captures the flow of control. Considering a nonpreemptive execution environment, we statically generate a schedule table and derive a guaranteed worst case delay.

Currently, more and more real-time systems are used in physically distributed environments and have to be implemented on distributed architectures in order to meet reliability, functional, and performance constraints. However, researchers have often ignored or very much simplified aspects concerning the communication infrastructure. One typical approach is to consider communication processes as processes with a given execution time (depending on the amount of information exchanged) and schedule them as any other process, without considering issues like communication protocol, bus arbitration, packaging of mes-

sages, clock synchronization, etc. These aspects are, however, essential in the context of safety-critical distributed real-time applications, and one of our objectives is to develop a strategy that takes them into consideration for process scheduling.

Many efforts dedicated to communication synthesis have concentrated on the synthesis support for the communication infrastructure but without considering hard real-time constraints and system-level scheduling aspects [41]–[45].

We have to mention here some results obtained in extending real-time schedulability analysis so that network communication aspects can be handled. In [46], for example, the CAN protocol is investigated, while [47] considers systems based on the asynchronous transfer mode (ATM) protocol. These works, however, are restricted to software systems implemented with priority-based preemptive scheduling.

In the first part of this paper we consider a communication model based on simple bus sharing. There we concentrate on the aspects of scheduling with data and control dependencies, and such a simpler communication model allows us to focus on these issues. However, one of the goals of this paper is to highlight how communication and process scheduling strongly interact with each other and how system-level optimization can only be performed by taking into consideration both aspects. Therefore, in the second part of this paper, we introduce a particular communication model and execution environment. We take into consideration the overheads due to communications and to the execution environment and consider the requirements of the communication protocol during the scheduling process. Moreover, our algorithm performs an optimization of parameters defining the communication protocol, which is essential for reduction of the execution delay. Our system architecture is built on a communication model that is based on the time-triggered protocol (TTP) [48]. TTP is well suited for safety-critical distributed real-time control systems and represents one of the emerging standards for several application areas, such as automotive electronics [28], [49].

This paper is divided as follows. In Section II, we formulate our basic assumptions and set the specific goals of this work. Section III defines the formal graph-based model, which is used for system representation, introduces the schedule table, and creates the background needed for presentation of our scheduling technique. The scheduling algorithm for conditional process graphs is presented in Section IV. In Section V, we introduce the hardware and software aspects of the TTP-based system architecture. The mutual interaction between scheduling and the communication protocol as well as our strategy for scheduling with optimization of the bus access scheme are discussed in Section VI. Section VII describes the experimental evaluation, and Section VIII presents our conclusions.

II. PROBLEM FORMULATION

We consider a generic architecture consisting of *programmable processors* and application specific *hardware processors* (ASICs) connected through several *buses*. The buses can be shared by several communication channels connecting processes assigned to different processors. Only one process can be executed at a time by a programmable processor,

while a hardware processor can execute processes in parallel.¹ Processes on different processors can be executed in parallel. Only one data transfer can be performed by a bus at a given moment. Data transfer on buses and computation can overlap.

Each process in the specification can be, potentially, assigned to several programmable or hardware processors, which are able to execute that process. For each process estimated cost and execution time on each potential host processor are given [50]. We assume that the amount of data to be transferred during communication between two processes has been determined in advance. In [50], we presented algorithms for automatic hardware/software partitioning based on iterative improvement heuristics. The problem we are discussing in this paper concerns performance estimation of a given design alternative and scheduling of processes and communications. Thus, we assume that each process has been assigned to a (programmable or hardware) processor and that each communication channel, which connects processes assigned to different processors, has been assigned to a bus. Our goal is to derive a worst case delay by which the system completes execution such that this delay is as small as possible, to generate the static schedule and optimize parameters of the communication protocol, such that this delay is guaranteed.

For the beginning, we will consider an architecture based on a communication infrastructure in which communication tasks are scheduled on buses similar to the way processes are scheduled on programmable processors. The time needed for a given communication is estimated depending on the parameters of the bus to which the respective communication channel is assigned and the number of transferred bits. Communication time between processes assigned to the same processor is ignored. Based on this architectural model we introduce our approach to process scheduling in the context of both control and data dependencies.

In the second part of the paper we introduce an architectural model with a communication infrastructure suitable for safety critical hard real-time systems. This allows us to further investigate the scheduling problem and to explore the impact of the communication infrastructure on the overall system performance. The main goal is to determine the parameters of the communication protocol so that the overall system performance is optimized and, thus, the imposed time constraints can be satisfied. We show that system optimization and, in particular, scheduling cannot be efficiently performed without taking into consideration the underlying communication infrastructure.

III. PRELIMINARIES

A. The Conditional Process Graph

We consider that an application specified as a set of interacting processes is mapped to an abstract representation consisting of a directed acyclic polar graph $G(V, E_S, E_C)$ called

a *process graph*. Each node $P_i \in V$ represents one process. E_S and E_C are the sets of simple and conditional edges, respectively. $E_S \cap E_C = \emptyset$ and $E_S \cup E_C = E$, where E is the set of all edges. An edge $e_{ij} \in E$ from P_i to P_j indicates that the output of P_i is the input of P_j . The graph is polar, which means that there are two nodes, called *source* and *sink*, that conventionally represent the first and last task. These nodes are introduced as dummy processes, with zero execution time and no resources assigned, so that all other nodes in the graph are successors of the source and predecessors of the sink, respectively.

A mapped process graph $\Gamma(V^*, E_S^*, E_C^*, M)$ is generated from a process graph $G(V, E_S, E_C)$ by inserting additional processes (communication processes) on certain edges and by mapping each process to a given processing element. The mapping of processes $P_i \in V^*$ to processors and buses is given by a function $M: V^* \rightarrow PE$, where $PE = \{pe_1, pe_2, \dots, pe_{N_{pe}}\}$ is the set of processing elements. $PE = PP \cup HPP \cup B$, where PP is the set of programmable processors, HPP is the set of dedicated hardware components, and B is the set of allocated buses. In certain contexts, we will call both programmable processors and hardware components simply processors. For any process P_i , $M(P_i)$ is the processing element to which P_i is assigned for execution. In the rest of this paper, when we use the term *conditional process graph* (CPG), we consider a mapped process graph as defined here.

Each process P_i , assigned to a programmable or hardware processor $M(P_i)$, is characterized by an execution time t_{P_i} . In the CPG depicted in Fig. 1, P_0 and P_{32} are the source and sink nodes, respectively. For the rest of 31 nodes, 17, denoted P_1, P_2, \dots, P_{17} , are ordinary processes specified by the designer. They are assigned to one of the two programmable processors pe_1 and pe_2 or to the hardware component pe_3 . The rest of 14 nodes are so-called *communication processes* ($P_{18}, P_{19}, \dots, P_{31}$). They are represented in Fig. 1 as solid circles and are introduced during the mapping process for each connection, which links processes assigned to different processors. These processes model interprocessor communication and their execution time $t_{i,j}$ (where P_i is the sender and P_j the receiver process) is equal to the corresponding communication time. All communications in Fig. 1 are performed on bus b_1 . As discussed in the previous section, we treat, for the beginning, communication processes exactly as ordinary processes. Buses are similar to programmable processors in the sense that only one communication can take place on a bus at a given moment.

An edge $e_{ij} \in E_C$ is a *conditional edge* (represented with thick lines in Fig. 1) and has an associated condition value. Transmission on such an edge takes place only if the associated condition value is *true* and not, like on simple edges, for each activation of the input process P_i . In Fig. 1, processes P_2, P_{11} , and P_{12} have conditional edges at their output. Process P_2 , for example, communicates alternatively with P_4 and P_5 , or with P_6 . Process P_{12} , if activated (which occurs only if condition D in P_{11} has value *true*), always communicates with P_{16} but alternatively with P_{14} or P_{15} , depending on the value of condition K .

We call a node with conditional edges at its output a *disjunction node* (and the corresponding process a *disjunction process*). A disjunction process has one associated condition, the value of

¹In some designs certain processes implemented on the same hardware processor can share resources and, thus, cannot execute in parallel. This situation can easily be handled in our scheduling algorithm by considering such processes in a similar way as those allocated to programmable processors. For simplicity, here we consider that processes allocated to ASICs do not share resources.

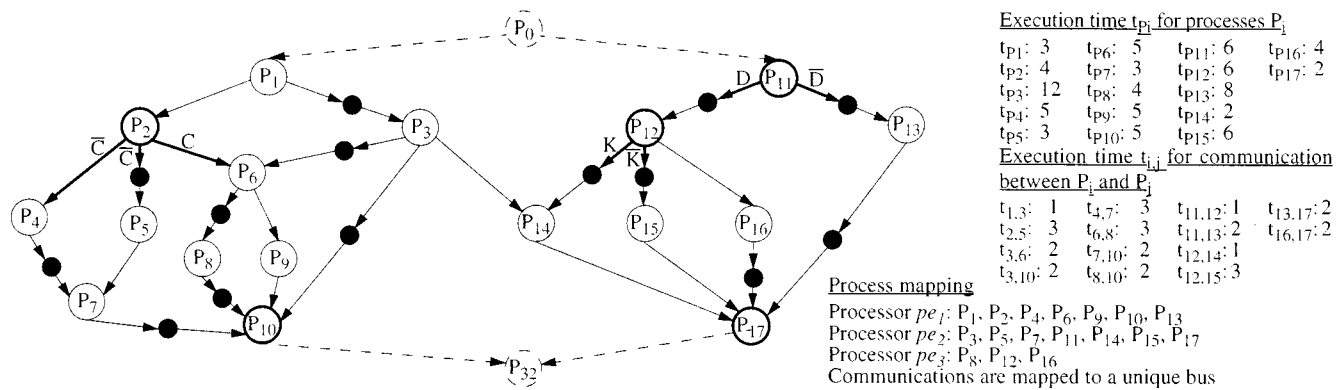


Fig. 1. Conditional process graph with execution times and mapping.

which it computes. Alternative paths starting from a disjunction node, which correspond to complementary values of the condition, are disjoint, and they meet in a so-called *conjunction node* (with the corresponding process called *conjunction process*).² In Fig. 1, circles representing conjunction and disjunction nodes are depicted with thick borders. The alternative paths starting from disjunction node P_2 , which computes condition C , meet in conjunction node P_{10} . Node P_{17} is the joining point for both the paths corresponding to condition K (starting from disjunction node P_{12}) and condition D (starting from disjunction node P_{11}). We assume that conditions are independent and alternatives starting from different processes cannot depend on the same condition.

A process that is not a conjunction process can be activated only after all its inputs have arrived. A conjunction process can be activated after messages coming on one of the alternative paths have arrived. All processes issue their outputs when they terminate. In Fig. 1, process P_7 can be activated after it receives messages sent by P_4 and P_5 ; process P_{10} waits for messages sent by P_3, P_8 , and P_9 or by P_3 and P_7 . If we consider the activation time of the source process as a reference, the activation time of the sink process is the delay of the system at a certain execution. This delay has to be, in the worst case, smaller than a certain imposed deadline. Release times of some processes as well as multiple deadlines can be easily modeled by inserting dummy nodes between certain processes and the source or the sink node, respectively. These dummy nodes represent processes with a certain execution time but that are not allocated to any processing element.

A Boolean expression X_{P_i} , called a *guard*, can be associated to each node P_i in the graph. It represents the necessary conditions for the respective process to be activated. In Fig. 1, for example, $X_{P_3} = \text{true}$, $X_{P_{14}} = D \wedge K$, $X_{P_{17}} = \text{true}$, and $X_{P_5} = \bar{C}$. X_{P_i} is not only necessary but also sufficient for process P_i to be activated during a given system execution. Thus, two nodes P_i and P_j , where P_j is not a conjunction node, are connected by an edge e_{ij} only if $X_{P_j} \Rightarrow X_{P_i}$ (which means that X_{P_i} is true whenever X_{P_j} is true). This avoids specifications

²If no process is specified on an alternative path, it is modeled by a conditional edge from the disjunction to the corresponding conjunction node (a communication process may be inserted on this edge at mapping).

in which a process is blocked even if its guard is true, because it waits for a message from a process that will not be activated. If P_j is a conjunction node, predecessor nodes P_i can be situated on alternative paths corresponding to a condition.

The above execution semantics is that of a so-called single rate system. It assumes that a node is executed at most once for each activation of the system. If processes with different periods have to be handled, this can be solved by generating several instances of the processes and building a CPG that corresponds to a set of processes as they occur within a time period that is equal to the least common multiple of the periods of the involved processes.

As mentioned, we consider execution times of processes, as well as the communication times, to be given. In the case of hard real-time systems this will, typically, be worst case execution times, and their estimation has been extensively discussed in the literature [13], [14]. For many applications, actual execution times of processes are depending on the current data and/or the internal state of the system. By explicitly capturing the control flow in our model, we allow for a more fine-tuned modeling and a tighter (less pessimistic) assignment of worst case execution times to processes, compared to traditional dataflow-based approaches.

B. The Schedule Table

For a given execution of the system, that subset of the processes is activated that corresponds to the actual track followed through the CPG. The actual track taken depends on the value of certain conditions. For each individual track there exists an optimal schedule of the processes that produces a minimal delay. Let us consider the CPG in Fig. 1. If all three conditions C , D , and K are true, the optimal schedule requires P_1 to be activated at time $t = 0$ on processor pe_1 and processor pe_2 to be kept idle until $t = 4$, in order to activate P_3 as soon as possible [see Fig. 2(a)]. However, if C and D are true but K is false, the optimal schedule requires starting both P_1 on pe_1 and P_{11} on pe_2 at $t = 0$; P_3 will be activated in this case at $t = 6$, after P_{11} has terminated and, thus, pe_2 becomes free [see Fig. 2(b)]. This example reveals one of the difficulties when generating a schedule for a system like that in Fig. 1. As the values of the conditions are unpredictable, the decision of on which process to ac-

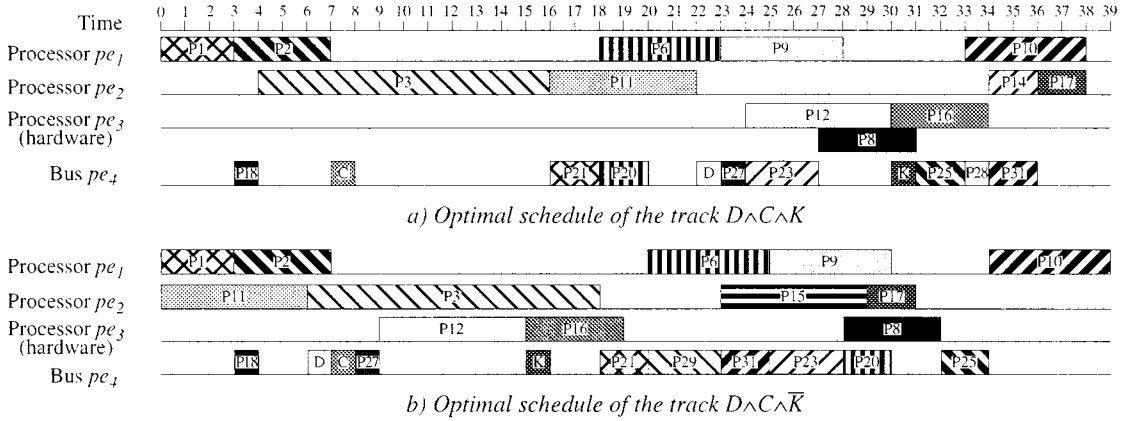


Fig. 2. Optimal schedules for two tracks extracted from the CPG in Fig. 1.

tivate on pe_2 and at which time has to be taken without knowing which values the conditions will later get. On the other side, at a certain moment during execution, when the values of some conditions are already known, they have to be used in order to take the best possible decisions on when and which process to activate. Heuristic algorithms have to be developed to produce a schedule of processes such that the worst case delay is as small as possible. The output of such an algorithm is a so-called *schedule table*. In this table, there is one row for each ordinary or communication process, which contains activation times for that process corresponding to different values of the conditions. Each column in the table is headed by a logical expression constructed as a conjunction of condition values. Activation times in a given column represent starting times of the processes when the respective expression is true.

Table I shows a possible schedule table corresponding to the system depicted in Fig. 1. According to this schedule processes P_1 , P_2 , P_{11} as well as the communication process P_{18} are activated unconditionally at the times given in the first column of the table. No condition value has yet been determined to select between alternative schedules. Process P_{14} , on the other hand, has to be activated at $t = 24$ if $D \wedge \bar{C} \wedge K = \text{true}$ and $t = 18$ if $D \wedge C \wedge K = \text{true}$. To determine the worst case delay δ_{\max} , we have to observe the rows corresponding to processes P_{10} and P_{17} : $\delta_{\max} = \max\{35 + t_{10}, 30 + t_{17}\} = 40$.

The schedule table contains all information needed by a distributed run-time scheduler to take decisions on activation of processes. We consider that, during execution, a very simple scheduler located on each processor decides on process and communication activation depending on actual values of conditions. Once activated, a process executes until it completes. Only one part of the table has to be stored on each processor, namely, the part concerning decisions that are taken by the corresponding scheduler.

Our goal is to derive a minimal worst case delay and to generate the corresponding schedule table given a conditional process graph $\Gamma(V^*, E_S^*, E_C^*, M)$ and estimated worst case execution times for each process $P_i \in V^*$.

At a certain execution of the system, one of the N_{alt} alternative tracks through the CPG will be executed. Each alternative track corresponds to one subgraph $G_k \in \Gamma$,

$k = 1, 2, \dots, N_{\text{alt}}$. For each subgraph, there is an associated logical expression L_k (the label of the track) that represents the necessary conditions for that subgraph to be executed. The subgraph G_k contains those nodes P_j of the conditional process graph Γ , for which $L_k \Rightarrow X_{P_j}$ (X_{P_j} is the guard of node P_j and has to be true whenever L_k is true). For the CPG in Fig. 1, we have six subgraphs (alternative tracks) corresponding to the following labels: $C \wedge D \wedge K$, $C \wedge D \wedge \bar{K}$, $C \wedge \bar{D}$, $\bar{C} \wedge D \wedge K$, $\bar{C} \wedge D \wedge \bar{K}$, and $\bar{C} \wedge \bar{D}$.

If at activation of the system all the condition values *would be* known in advance, the processes *could be* executed according to the (near) optimal schedule of the corresponding subgraph G_k . Under these circumstances, the worst case delay δ_{\max} would be

$$\delta_{\max} = \delta_M$$

with

$$\delta_M = \max\{\delta_k, k = 1, 2, \dots, N_{\text{alt}}\}$$

where δ_k is the delay corresponding to subgraph G_k .

However, this is not the case, as we do not assume any prediction of the condition values at the start of the system. Thus, what we can say is only that:³ $\delta_{\max} \geq \delta_M$.

A scheduling heuristic has to produce a schedule table for which the difference $\delta_{\max} - \delta_M$ is minimized. This means that the perturbation of the individual schedules, introduced by the fact that the actual track is not known in advance, should be as small as possible.

C. Conditions, Guards, and Influences

We first introduce some notations. If E is a logical expression, we use the notation $/E/$ to denote the set of conditions used in E . Thus, $/C \wedge \bar{D} \wedge K/ = /C \wedge D \wedge \bar{K}/ = \{C, D, K\}$; $/\text{true}/ = /false/ = \emptyset$. Similarly, if M is a set of condition values, $/M/$ is the set of conditions used in M . For example, if $M = \{C, \bar{D}, K\}$, then $/M/ = \{C, D, K\}$. For a set M of condition values, we denote with $\wedge M$ the logical expression consisting of the conjunction of that values. If $M = \{C, \bar{D}, K\}$, then $\wedge M = C \wedge \bar{D} \wedge K$.

³For this formula to be rigorously correct, δ_M has to be the maximum of the optimal delays for each subgraph.

TABLE I
SCHEDULE TABLE CORRESPONDING TO THE CPG SHOWN IN FIG. 1

| | <i>true</i> | <i>D</i> | $D \wedge C$ | $D \wedge C \wedge \bar{K}$ | $D \wedge C \wedge K$ | $D \wedge \bar{C}$ | $D \wedge \bar{C} \wedge \bar{K}$ | $D \wedge \bar{C} \wedge K$ | \bar{D} | $\bar{D} \wedge C$ | $\bar{D} \wedge \bar{C}$ |
|------------------------------|-------------|----------|--------------|-----------------------------|-----------------------|--------------------|-----------------------------------|-----------------------------|-----------|--------------------|--------------------------|
| P_1 | 0 | | | | | | | | | | |
| P_2 | 3 | | | | | | | | | | |
| P_3 | | 6 | | | | | | | 6 | | |
| P_4 | | | | | | 7 | | | | | 7 |
| P_5 | | | | | | | 18 | 18 | | | 18 |
| P_6 | | | | 21 | 20 | | | | | 20 | |
| P_7 | | | | | | | 21 | 21 | | | 21 |
| P_8 | | | | 29 | 28 | | | | | 28 | |
| P_9 | | | | 26 | 25 | | | | | 25 | |
| P_{10} | | | | 35 | 34 | | 27 | 26 | | 34 | 26 |
| P_{11} | 0 | | | | | | | | | | |
| P_{12} | | | 9 | | | 9 | | | | | |
| P_{13} | | | | | | | | | | 10 | 13 |
| P_{14} | | | | | 18 | | | 24 | | | |
| P_{15} | | | | 19 | | | 24 | | | | |
| P_{16} | | | | 15 | 15 | | 15 | 15 | | | |
| P_{17} | | | | 25 | 24 | | 30 | 26 | | 24 | 24 |
| $P_{18} (1 \rightarrow 3)$ | 3 | | | | | | | | | | |
| $P_{19} (2 \rightarrow 5)$ | | | | | | 9 | | | | | 8 |
| $P_{20} (3 \rightarrow 10)$ | | | | 21 | 20 | | 21 | 20 | | 20 | 20 |
| $P_{21} (3 \rightarrow 6)$ | | | | 19 | 18 | | | | | 18 | |
| $P_{22} (4 \rightarrow 7)$ | | | | | | 12 | | | | | 13 |
| $P_{23} (6 \rightarrow 8)$ | | | | 26 | 25 | | | | | 25 | |
| $P_{24} (7 \rightarrow 10)$ | | | | | | | 25 | 24 | | | 24 |
| $P_{25} (8 \rightarrow 10)$ | | | | 33 | 32 | | | | | 32 | |
| $P_{26} (11 \rightarrow 13)$ | | | | | | | | | | 8 | 11 |
| $P_{27} (11 \rightarrow 12)$ | | | 8 | | | 8 | | | | | |
| $P_{28} (12 \rightarrow 14)$ | | | | | 16 | | | 16 | | | |
| $P_{29} (12 \rightarrow 15)$ | | | | 16 | | | 16 | | | | |
| $P_{30} (13 \rightarrow 17)$ | | | | | | | | | | 22 | 22 |
| $P_{31} (16 \rightarrow 17)$ | | | | 23 | 22 | | 23 | 22 | | | |
| <i>D</i> | 6 | | | | | | | | | | |
| <i>C</i> | | 7 | | | | | | | 7 | | |
| <i>K</i> | | | 15 | | | 15 | | | | | |

Activation times in the schedule table are such that a process is started only after its predecessors, corresponding to the actually executed track, have terminated. This is a direct consequence of the dataflow dependencies as captured in the CPG model. However, in the context of the CPG semantics, there are more subtle requirements that the schedule table has to fulfill in order to produce a deterministic behavior that is correct for any combination of condition values.

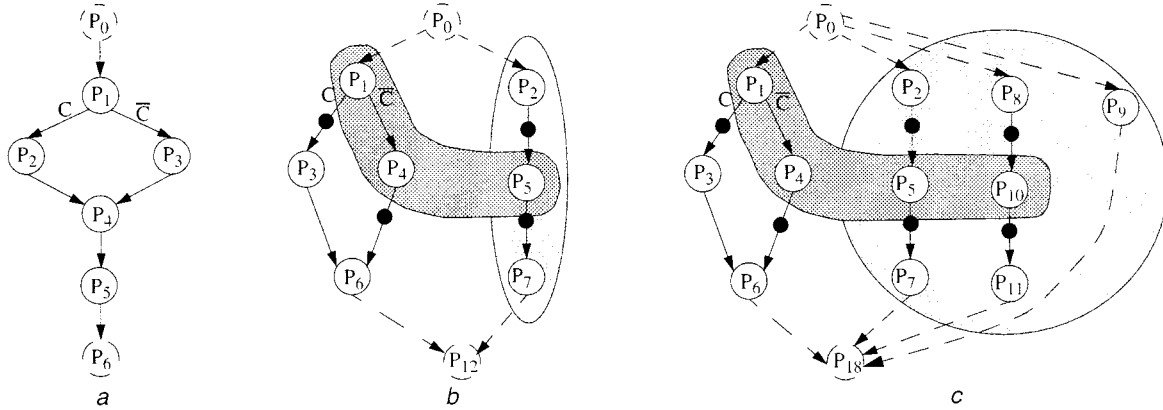
- R1) If for a certain execution of the system the guard X_{P_i} becomes true, then P_i has to be activated during that execution.
- R2) If for a certain process P_i , with guard X_{P_i} , there exists an activation time in the column headed by the logical expression E_k , then $E_k \Rightarrow X_{P_i}$; this means that no process will be activated if the conditions required for its execution are not fulfilled.
- R3) Activation times and the corresponding logical expressions have to be determined so that a process is activated not more than one single time for a given execution of the system. Thus, if for a certain execution a process P_i is scheduled to be executed at $\tau_{P_i}^{E_j}$ (placed in column headed by E_j), there is no other execution time $\tau_{P_i}^{E_k} \neq \tau_{P_i}^{E_j}$ for P_i (placed in column headed by E_k) so that E_k becomes true during the same execution.

- R4) Activation of a process P_i at a certain time t has to depend only on condition values that are determined at the respective moment t and are known to the processor that takes the decision on P_i to be executed.

The correct behavior, according to the semantics associated to conditions, is guaranteed by requirements R1 and R2. R3 guarantees the functionality as a single rate system. Requirement R4 states that decisions at any location of the distributed system are based on information available at the particular location and at the given time.

A scheduling algorithm has to traverse, in one way or another, all the alternative tracks in the CPG and to place activation times of processes and the corresponding logical expressions into the schedule table. For a given execution track, labeled L_k , and a process P_i that is included in that track, one single activation time $\tau_{P_i}^{L_k}$ has to be placed into the table (see R3 above); the corresponding column will be headed by an expression E . We say that $/E/$ is the set of conditions used for scheduling P_i . It is obvious that $/E/ \subseteq /L_k/$ and $L_k \Rightarrow E$.

Which conditions in set $/E/$ are to be used for scheduling a certain process? This is a question to be answered by the scheduling algorithm. A straightforward answer would be: those conditions that determine if the process has to be activated or that have an influence on the moment when this activation has to occur.



Process mapping to processors is illustrated by shading; all processors are programmable.

Fig. 3. Examples of conditional process graphs.

In the following section, we will give a more exact answer to the question above. As a prerequisite, we define the *guard set* GS_{P_i} and *influence set* IS_{P_i} of a process P_i .

The *guard set* GS_{P_i} of a process P_i is the set of conditions that decide whether or not the process is activated during a certain execution of the system. $GS_{P_i} = /X_{P_i}/$, where X_{P_i} is the guard of process P_i . In Fig. 1, for example, $X_{P_3} = \text{true}$, $X_{P_{14}} = D \wedge K$, and $X_{P_{15}} = D \wedge \bar{K}$; thus, $GS_{P_3} = \emptyset$, $GS_{P_{14}} = GS_{P_{15}} = \{D, K\}$.

However, not only the conditions in the guard set have an impact on the execution of a process. There are other conditions that do not decide if a process has or does not have to be activated during a given execution, but that influence the moment when the process is scheduled to start. In Fig. 3(a), for example, $GS_{P_4} = \emptyset$; process P_4 is executed regardless of whether C is true or false. However, the activation time of P_4 depends on the value of condition C : P_4 will be activated at $\tau_{P_4}^C = t_{P_1} + t_{P_2}$, if $C = \text{true}$, and at $\tau_{P_4}^{\bar{C}} = t_{P_1} + t_{P_3}$, if $C = \text{false}$. As a consequence, the activation time of P_5 is also dependent on the value of C . We say that both P_4 and P_5 are influenced by C , and this influence is induced to them (from their predecessors) by succession in the CPG.

However, not only by succession can the influence of a condition be induced to a process. This is illustrated in Fig. 3(b), where we assume that each communication is mapped to a distinct bus and that $t_{P_1} < t_{P_2} + t_{2,5} < t_{P_1} + t_{P_4}$. For $C = \text{true}$, we have $\tau_{P_5}^C = \max(t_{P_2} + t_{2,5}, t_{P_1}) = t_{P_2} + t_{2,5}$. For $C = \text{false}$, P_4 becomes ready and will be scheduled before P_5 ; thus, $\tau_{P_5}^{\bar{C}} = \max(t_{P_2} + t_{2,5}, t_{P_1} + t_{P_4}) = t_{P_1} + t_{P_4}$. The activation time of P_5 depends on condition C , because P_5 shares the same programmable processor with P_4 and the activation of P_4 is conditioned by C . In this case, the influence of condition C is induced on P_5 from P_4 by resource sharing. This influence is then further induced from P_5 , by succession, to the communication process between P_5 and P_7 and to P_7 . P_2 and the communication P_2 to P_5 as well as P_1 are not influenced by C .

Fig. 3(c) illustrates a more complex situation. We assume that communications are mapped to distinct buses and that, if

there is a conflict for the processor, P_8 is scheduled first and, with decreasing priority, P_9 , P_{11} , and P_2 . Concerning execution times, we assume: $t_{P_1} < t_{P_8} + t_{8,10} < t_{P_1} + t_{P_4}$, $t_{P_8} + t_{8,10} + t_{P_{10}} + t_{10,11} < t_{P_8} + t_9 < t_{P_1} + t_{P_4} + t_{P_{10}} + t_{10,11}$. Under these circumstances, if $C = \text{true}$, then P_{11} is ready when the processor is released by P_9 , and thus will be scheduled before P_2 ; we have $\tau_{P_2}^C = t_{P_8} + t_{P_9} + t_{P_{11}}$. For $C = \text{false}$, P_{10} is delayed because of the earlier activation of P_4 , which means that P_{11} will not be ready when P_9 has finished; thus, P_2 is scheduled before P_{11} and $\tau_{P_2}^{\bar{C}} = t_{P_8} + t_{P_9}$. In this case, the influence of C on P_2 has been induced by resource sharing from P_{11} , which is influenced from its predecessor P_{10} .

We can observe from the examples above that the influence of a certain condition on a process can only be determined dynamically, during the scheduling process. A particular scheduling policy, priorities, and execution times all have an impact on whether at a certain moment the scheduling of a process is influenced or not by a condition. In Fig. 3(c), for example, scheduling of P_2 is influenced by C only because the scheduler has considered P_8 with a higher priority, and this is also why P_8 is not influenced by C .

The set of conditions to be used for scheduling a certain process P_i at a moment t consists of the guard set GS_{P_i} corresponding to that process and of all the other conditions that influence the activation of P_i . An efficient heuristic has to be developed for the scheduler in order to determine these conditions so that a correct and close to optimal schedule table is produced. In order to solve this problem, we define the influence set IS_{P_i} corresponding to a process P_i . It is important to notice that this set of conditions can be determined statically for a given CPG, and it consists of all conditions that are not part of the guard set but potentially could influence the scheduling of process P_i .

A condition C is in the *influence set* IS_{P_i} of a process P_i if C is not in the guard set GS_{P_i} of P_i and if the influence of C can be induced to P_i by *succession* or by *resource sharing*.

The influence of condition C can be induced to a process P_i by *succession* if P_i is a successor of a process P_j and the following condition is satisfied: $C \in IS_{P_j}$ or $C \in GS_{P_j}$.

The influence of condition C can be induced to a process P_i by *resource sharing* if all the following three conditions are satisfied.

- 1) P_i is mapped on a programmable processor or a bus.
- 2) P_i is neither the process that computes condition C nor a predecessor of that process.
- 3) There exists a process P_j mapped to the same resource as P_i and:
 - a) P_j is neither a predecessor nor a successor of P_i ;
 - b) $C \in IS_{P_j}$ or $C \in GS_{P_j}$.

In Fig. 3(b), for example, $IS_{P_5} = IS_{P_6} = IS_{P_7} = IS_{5,7} = \{C\}$. For the other processes, the influence set is void. P_5 can be influenced by resource sharing and the other processes by succession. In Fig. 3(c), the influence of C can be induced by succession to P_6 and by resource sharing to P_5 and P_{10} . This is then transferred by succession to P_7 , P_{11} and to the communication processes preceding them. The influence can be induced by resource sharing to P_2 , P_8 , and P_9 and then by succession to the communications following P_2 and P_8 . For both examples, we considered that communications do not share any bus. If in Fig. 3(b) we consider that all communications are mapped to the same bus, then $IS_{2,5} = \{C\}$, as a result of resource sharing.

In the next section we will see how the influence set is utilized by a particular scheduling algorithm in order to determine the actual set of conditions used for scheduling a process.

The activation of a process P_i at a certain moment t depends on the conditions specified in the expression heading the respective column of the schedule table (the conditions used to schedule P_i). According to requirement R4, the values of all these conditions have to be determined at the respective moment t and have to be known to the processor $Q(P_i)$. $Q(P_i) = M(P_i)$, the processor that executes P_i , if P_i is not a communication process. If P_i is a communication process, then $Q(P_i) = M(P_j)$, where P_j is the process that initiates the communication.

The value of a condition is determined at the moment τ at which the corresponding disjunction process terminates. This means that at any moment t , $t \geq \tau$, the value of the condition is available for scheduling decisions on the processor that has executed the corresponding disjunction process. However, in order to be available on any other processor, the value has to arrive at that processor. The scheduling algorithm has to consider both the time and the resource (bus) needed for this communication. Transmissions of condition values are scheduled as communications on buses. These communication processes are not represented in the CPG and are the only activities that are not *a priori* mapped to a specific resource (bus). A condition communication will be mapped as part of the scheduling process to a bus. For broadcasting of condition values, only buses are considered to which all processors are connected, and we assume that at least one such bus exists.⁴ The time τ_0 needed for such a communication is the same for all conditions and depends on the features of the employed buses. The transmitted condition value is available for scheduling decisions on all other processors τ_0

⁴This assumption is made for simplification of the further discussion. If no bus is connected to all processors, communication tasks have to be scheduled on several buses according to the actual interconnection topology.

```

List_schedule ()
for each processing element  $pe_i$ ,  $i=1,2, \dots, N_{pe}$ , do
     $free_{pe_i}=0$ 
end for;
schedule  $P_0$  at  $t=0$ ; initialize ready process list  $Ls\_Ready$ 
with direct successors of  $P_0$ ;    --  $P_0$  is the source node
while  $Ls\_Ready$  is not empty do
     $p=Head(Ls\_Ready)$ ;
    if  $M(p) \in HP$  then    -- hardware supports several
        schedule  $p$  at  $t=p.t_{ready}$  -- processes at a time;
    else
         $p=Select(Ls\_Ready, M(p))$ ;
        schedule  $p$  at  $t=\max(p.t_{ready}, free_{M(p)})$ ;
         $free_{M(p)}=t+t_p$ 
    end if;
    delete  $p$  from  $Ls\_Ready$ ;
     $Ls\_Ready = Ls\_Ready +$  processes which become
        ready after  $p$  is executed
end while
end List_schedule;

```

Fig. 4. Basic list scheduling algorithm.

time units after initiation of the broadcast. For the example given in Table I, communication time for condition values has been considered $\tau_0 = 1$. In Table I, the last three rows indicate the schedule for communication of the values of conditions C , D , and K .

IV. SCHEDULING OF CONDITIONAL PROCESS GRAPHS

A. The Scheduling Algorithm

Our algorithm for scheduling of conditional process graphs relies on a list scheduling heuristic. List scheduling heuristics [29] are based on ordered lists from which processes are extracted to be scheduled at certain moments. In the algorithm presented in Fig. 4 we have such a list, Ls_Ready , in which processes are placed in increasing order of the time when they become ready for execution (this time is stored as an attribute $p.t_{ready}$ of each process p in the list, and is the moment when all the predecessors of p have terminated). If p is mapped to a hardware processor, it is scheduled, without any restriction, at the moment when it is ready. On a programmable processor or bus, however, a new process p can be scheduled only after the respective processing element $M(p)$ becomes free [$free_{M(p)}$ indicates this moment]. There can be several processes mapped on $M(p)$, so that, at a given moment, $p.t_{ready} \leq free_{M(p)}$. All of them will be ready when processing element $M(p)$ becomes free. From these processes, function *Select* will pick out the one that has the highest priority assigned to it, in order to be scheduled. This priority function is one of the essential features that differentiate various list scheduling approaches. We will later discuss the particular priority function used by our algorithm.

The algorithm presented in Fig. 4 is able to schedule, based on a certain priority function, process graphs without conditional control dependencies. In Fig. 5, we show the algorithm for scheduling of conditional process graphs. It is based on the general principles of list scheduling, using the list Ls_Ready to store processes ordered by the time when they become ready

```

List_schedule_CPG (Ls_Ready, free_pei | i=1... Npe, Cond_Set)
  p=empty_process;
  while Ls_Ready is not empty and p is not a disjunction process do
    p=Head(Ls_Ready);
    if  $M(p) \in HP$  then                                     -- hardware processors support several processes at a time;
       $\tau = p.t_{ready}$ 
    else                                                    -- programmable processors and buses support one process at a time
      p=Select(Ls_Ready, M(p));
       $\tau = \max(p.t_{ready}, free_{M(p)})$ 
    end if;
     $Z = \{c \mid c \in Cond\_Set \text{ and } c.t_{comp} \leq \tau\}$ ;          -- set of condition values computed at  $\tau$ 
     $CS = \{c \mid c \in Z \text{ and } (\{c\} \subset GS_p \text{ or } \{c\} \subset IS_p)\}$ ; -- set of condition values used to schedule p
     $Exp = \vee CS$ ;                                           -- the logical expression used to schedule p
     $\theta = \text{maximum among } c.t_{arr}$ , where  $c \in CS$  and c is computed on a processor different from Q(p);
    -- this is the earliest time when all condition values needed to schedule p are known on processor Q(p)
    schedule p at  $t = \max(\tau, \theta)$ , controlled by expression Exp;
     $free_{M(p)} = t + t_p$ ;
    delete p from Ls_Ready;
    if p is not a disjunction process then
      Ls_Ready = Ls_Ready + processes which become ready after p is executed
    end if
  end while
  if p is a disjunction process then                       -- we consider that C is the condition computed by p
     $C.t_{comp} = t + t_p$ ;
    Schedule the communication of the condition value and fix C.tarr accordingly;
    Ls_Ready_true = Ls_Ready + processes on the true alternative, which become ready after p is executed;
    List_schedule_CPG(Ls_Ready_true, free_pei | i=1... Npe, Cond_Set  $\cup$  {C});
    Ls_Ready_false = Ls_Ready + processes on the false alternative, which become ready after p is executed;
    List_schedule_CPG(Ls_Ready_false, free_pei | i=1... Npe, Cond_Set  $\cup$  { $\bar{C}$ });
  end List_schedule_CPG;

```

Fig. 5. List scheduling of conditional process graphs.

for execution. *Cond_Set* is the set of those condition values for which the corresponding disjunction process (which computes the value) has been already scheduled. Each condition value $c \in Cond_Set$ has two attributes attached: $c.t_{comp}$ is the moment when the value *c* has been computed (the disjunction process has terminated) and $c.t_{arr}$ is the time when *c* is available on processors different from the one that executes the disjunction process. The recursive procedure *List_schedule_CPG* traverses the CPG analyzing each possible alternative track and considering for each track only the processes executed for the respective condition values. The algorithm, thus, proceeds along a binary decision tree corresponding to the alternative tracks, which is explored in depth first order.

Fig. 6 shows the decision tree explored during generation of Table I for the CPG in Fig. 1. The nodes of the tree correspond to the states reached during scheduling when a disjunction process has been scheduled and the algorithm branches for the alternative condition values. Whenever a process P_i has been scheduled, it is eliminated from the ready list *Ls_Ready*, and all processes that become ready after P_i has been executed are added to *Ls_Ready*. If, however, P_i is a disjunction process computing condition *C*, the two possible alternatives are handled separately. First, the processes are added to *Ls_Ready* that become ready on the true alternative, and scheduling is continued with the condition set $Cond_Set \cup \{C\}$. Then, those processes are added to *Ls_Ready* that become ready on the false

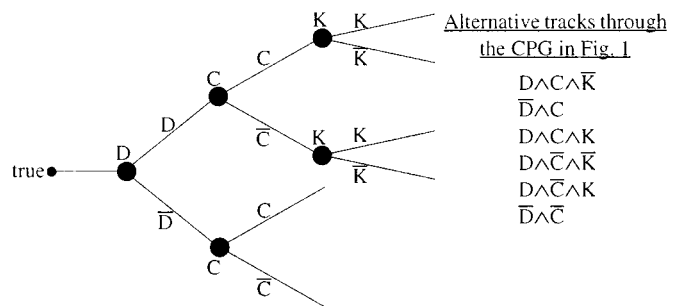


Fig. 6. The decision tree explored at generation of the schedule table for the CPG in Fig. 1.

alternative, and scheduling is continued with the condition set $Cond_Set \cup \{\bar{C}\}$.

Let us suppose that a process P_i has become ready and has been added to the list *Ls_Ready* before the branching for a certain condition *Z*. If the process is scheduled before the branching, it will be eliminated from *Ls_Ready* and, thus, not handled any more during further exploration of the tree below that branching point. This scheduling decision is the one and only decision valid for P_i on all tracks including that point. If P_i has not been scheduled before branching on *Z*, it will be inherited by *Ls_Ready* on both branches (these branches belong to different execution tracks through the CPG) and

will be eventually scheduled on both of them. There are two important conclusions to be drawn here.

- 1) If a process becomes ready on a certain execution track through the CPG, it will be scheduled for activation on that track. By this, requirement R1 is satisfied.
- 2) During exploration of a certain execution track, a given process that has to be activated on that track is considered one single time for scheduling.

We will now discuss three additional aspects concerning the scheduling algorithm in Fig. 5: the set of conditions used to schedule a certain process, the communication of condition values, and the priority assignment.

B. Conditions Used for Process Scheduling

When a process is ready to be scheduled at a moment τ , the set of conditions to be used to schedule it has to be determined:

A condition C will be used at moment τ for scheduling a process P_i if C is a member of the guard set or of the influence set of P_i and, according to the current schedule, the value of C has already been computed at time τ .

There are several aspects to be discussed in the context of this rule.

- 1) As discussed in Section III, the influence set IS_{P_i} consists of all conditions, except those in the guard set, which *potentially could have* an influence on the execution of P_i . The above rule selects from the influence set those conditions that *actually have* an influence on the execution of P_i , considering the actual track that is scheduled and the particular moment τ . Those conditions are eliminated that have not been computed at moment τ , according to the current schedule, and, thus, are excluded to have any influence (by sharing or succession) on the scheduling of P_i in the current context.
- 2) Every condition C that is a member of the guard set of a process P_i ($C \in GS_{P_i}$) is used to schedule P_i . By this, requirement R2 is satisfied.

We use the notation c for the value of condition C corresponding to the current track. C is a member of the guard set of P_i , which means that it is computed by a predecessor of P_i . Thus, $c \in Cond_Set$ and $c.t_{comp} \leq \tau$ (any predecessor of P_i has been scheduled and has terminated before P_i becomes ready). Consequently, $c \in Z$ and also $c \in CS$ (Fig. 5), which means that C will be used for scheduling P_i .

- 3) Consider a process P_i with guard X_{P_i} . For a given execution track, labeled L_k ($L_k \Rightarrow X_{P_i}$), P_i is scheduled at $\tau_{P_i}^E$, conditioned by expression E . Let us suppose that the same process P_i is scheduled also on another track, labeled L_q ($L_q \Rightarrow X_{P_i}$) at $\tau_{P_i}^{E'} \neq \tau_{P_i}^E$, conditioned by expression E' . We have shown in the previous subsection that P_i is scheduled one single time for a given execution track. In order to prove that requirement R3 is satisfied, we have to show that during execution of neither of the two tracks can both expressions E and E' become true.

We have $E = X_{P_i} \wedge Y$ and $E' = X_{P_i} \wedge Y'$, where Y and Y' are conjunctions of condition values selected

from the influence set of P_i ; we know that $L_k \Rightarrow Y$ and $L_q \Rightarrow Y'$. Let us consider I_1, I_2, \dots, I_N , the conditions that appear both in L_k and in L_q , but with complementary values. If we take, for example, the tracks labeled $D \wedge \overline{C} \wedge \overline{K}$ and $D \wedge C \wedge K$ in Fig. 6, then I_1 is C and I_2 is K . For any pair of tracks there has to be at least one such condition (the one on which the two tracks branch). If none of the conditions I_1, \dots, I_N is a member of the influence set of P_i , then P_i is scheduled at identical times on the two tracks (there is no influence on the process from any condition which differentiates the tracks) and, thus, there is no conflict ($\tau_{P_i}^{E'} = \tau_{P_i}^E$).

Let us suppose that I_1 is a member of the influence set of P_i . If P_i is considered for scheduling before the algorithm has branched on condition I_1 , then the respective scheduling time is the only one for both tracks L_k and L_q (see previous subsection) and there is no conflict.

If P_i is considered for scheduling after the algorithm has branched on condition I_1 , but before $I_1.t_{calc}$, we have $\tau_{P_i}^{E'} = \tau_{P_i}^E$ and, thus, no conflict (before $I_1.t_{calc}$, which means that before the disjunction process has terminated and I_1 has been computed, the two branches act identically). If P_i is scheduled after $I_1.t_{calc}$, then I_1 will be used to schedule P_i on both tracks ($I_1 \in /E/$ and $I_1 \in /E'/$) but with complementary values. Consequently, E and E' cannot both become true during the same execution.

C. Communication of Condition Values

In Section III-C, we have shown that the values of all conditions needed to schedule a certain process P_i have to be known on processor $Q(P_i)$ at the activation moment of that process. In Fig. 5, we have used θ to denote the earliest moment when all these conditions are available, and process P_i will not be scheduled before θ . Hence, correctness requirement R4 is also satisfied by the algorithm in Fig. 5.

Transmission of a value for condition C is scheduled as soon as possible on the first bus that becomes available after termination of the disjunction process that computes C . At the same time, the moment when the respective value becomes available on processors different from the one running the disjunction process has to be determined: $C.t_{arr} = t + \tau_0$, where t is the time at which transmission of the condition value has been scheduled. We mention that communication of a condition value will only be scheduled if there exists at least one process P_j so that $C \in GS_{P_j}$ or $C \in IS_{P_j}$, and $Q(P_j)$ is different from the processor running the disjunction process.

D. Priority Assignment

In this section, we first introduce the priority function used by our list scheduling algorithm and then show how priorities are assigned in the particular context of a CPG.

List scheduling algorithms rely on priorities in order to solve conflicts between several processes that are ready to be executed on the same programmable processor or bus. Such priorities very often are based on the critical path (CP) from the respective process to the sink node. In this case (we call it CP

scheduling), the priority assigned to a process P_i will be the maximal total execution time on a path from the current node to the sink

$$l_{Pi} = \max_k \sum_{P_j \in \pi_{ik}} t_{P_j}$$

where π_{ik} is the k th path from node P_i to the sink node.

Considering the concrete definition of our problem, significant improvements of the resulting schedule can be obtained, without any penalty in scheduling time, by making use of the available information on process mapping.

Let us consider the graph in Fig. 7 and suppose that the list scheduling algorithm has to decide between scheduling process P_A or P_B , which are both ready to be scheduled on the same programmable processor or bus pe_i . In Fig. 7, we have depicted only the critical path from P_A and P_B to the sink node. Let us consider that P_X is the last successor of P_A on the critical path such that all processes from P_A to P_X are assigned to the same processing element pe_i . The same holds for P_Y relative to P_B . t_A and t_B are the total execution time of the chain of processes from P_A to P_X and from P_B to P_Y , respectively, following the critical paths. λ_A and λ_B are the total execution times of the processes on the rest of the two critical paths. Thus, we have

$$l_{PA} = t_A + \lambda_A \quad \text{and} \quad l_{PB} = t_B + \lambda_B.$$

However, we will not use the length of these critical paths as a *priority*. Our policy, called partial critical path scheduling (PCP), is based on the estimation of a lower bound L on the total delay, taking into consideration that the two chains of processes $P_A - P_X$ and $P_B - P_Y$ are executed on the same processor. L_{PA} and L_{PB} are the lower bounds if P_A and P_B , respectively, are scheduled first

$$\begin{aligned} L_{PA} &= \max(T_{current} + t_A + \lambda_A, \\ &\quad T_{current} + t_A + t_B + \lambda_B) \\ L_{PB} &= \max(T_{current} + t_B + \lambda_B, \\ &\quad T_{current} + t_B + t_A + \lambda_A). \end{aligned}$$

We select the alternative that offers the perspective of the shorter delay $L = \min(L_{PA}, L_{PB})$. It can be observed that if $\lambda_A > \lambda_B$, then $L_{PA} < L_{PB}$, which means that we have to schedule P_A first so that $L = L_{PA}$; similarly, if $\lambda_B > \lambda_A$, then $L_{PB} < L_{PA}$, and we have to schedule P_B first in order to get $L = L_{PB}$.

As a conclusion, for PCP scheduling we use the value of λ_{Pi} as a priority criterion instead of the length l_{pi} of the whole critical path. Thus, we take into consideration only that part of the critical path corresponding to a process P_i that starts with the first successor of P_i that is assigned to a processor different from $M(P_i)$.

For evaluation of the PCP policy, we used 1250 graphs generated for experimental purpose. Two-hundred fifty graphs have been generated for each dimension of 20, 40, 75, 130, and 200 nodes. We considered architectures consisting of one ASIC, one to 11 processors, and one to eight buses. We have evaluated the percentage deviations of the schedule lengths produced by CP,

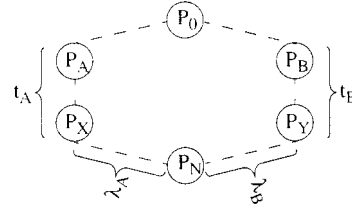


Fig. 7. Delay estimation for PCP scheduling.

UB,⁵ and PCP-based list scheduling from the lengths of the optimal schedules. The optimal schedules were produced running a branch-and-bound based algorithm. The average deviation for all graphs is 4.73% with UB, 4.69% with CP, and 2.35%, two times smaller, with PCP. Further details concerning the PCP policy and its experimental evaluation can be found in [33].

How is a PCP priority assignment applied to conditional process graphs? The problem is that in a CPG a process potentially belongs to several tracks, and relative to each of these tracks it has a different PCP priority. Thus, a process is characterized not by a single priority but by a set of priorities, one for each track to which the process belongs. P_{11} in Fig. 1, for example, has priority 14 corresponding to all tracks that imply condition value D , priority 15 corresponding to the tracks that imply $D \wedge K$, and 18 for those implying $D \wedge \bar{K}$.

When a certain process is considered for scheduling, the scheduler does not know the particular track to be followed, and thus it is not obvious which priority to use for the given process. We have used a very simple and efficient heuristic for priority assignment to processes in a CPG: for each process P_i , that PCP priority is considered that corresponds to the most critical track to which P_i belongs. This means that P_i gets the maximal priority among those assigned for each particular track. In the case of P_{11} (Fig. 1), the chosen priority will be 18.

From a practical point of view, the above heuristic means that priorities are assigned to processes by simply applying PCP priority assignment to the simple (unconditional) process graph that is obtained from the CPG by ignoring conditional dependencies and considering conditional edges as simple edges.

By this priority assignment policy, we try to enforce schedules so that the execution of longer tracks is as similar as possible to the best schedule we could generate if we could know in advance that the particular track is going to be selected. By introducing perturbations (possible delays) into the shorter tracks and letting the longer ones proceed as similar as possible to their (near) optimal schedule, we hope to produce a schedule table with a minimized worst case execution time. This is in line with the objective formulated in Section III-B.

V. A TTP-BASED SYSTEM ARCHITECTURE

In the previous sections, we have considered a general system architecture with a relatively simple communication infrastructure. Such a model is representative for a large class of applications, and this is the reason why it has been used in order to develop our scheduling strategy for CPGs.

⁵Urgency-based (UB) scheduling uses the difference between the as late as possible (ALAP) schedule of a process and the current time as a priority.

When designing a particular system, however, very specific architectural solutions and particular communication protocols often have to be used. In such a situation, particulars of the communication infrastructure have to be considered for system scheduling. At the same time, several parameters of the communication protocol have to be optimized in order to meet time constraints at the lowest possible cost. These two aspects are strongly interrelated and have a strong impact on the quality of the final design.

We consider a system architecture built around a communication model that is based on the TTP [48]. TTP is well suited for the category of systems targeted by our synthesis environment. These are systems that include both dataflow and control aspects and that are implemented on distributed architectures. They have to fulfill hard real-time constraints and are often safety-critical. TTP is also perfectly suited for systems implemented with static nonpreemptive scheduling, and thus represents an ideal target architecture for the scheduling approach presented in the previous sections.

While in the previous sections we emphasized the impact of conditional dependencies on system scheduling, in the rest of this paper we mainly concentrate on how the communication infrastructure and protocol have to be considered during scheduling and how they can be optimized in order to reduce the system delay. In this section, we describe our hardware and software architecture based on the TTP. In Section VI, we then show how scheduling and optimization of the communication protocol interact with each other and how they can be considered together in order to improve the overall system performance.

A. Hardware Architecture

We consider architectures consisting of nodes connected by a broadcast communication channel [Fig. 8(a)]. Every node consists of a TTP controller [48], a CPU, a RAM, a ROM, and an I/O interface to sensors and actuators. A node can also have an ASIC in order to accelerate parts of its functionality.

Communication between nodes is based on the TTP. TTP was designed for distributed real-time applications that require predictability and reliability. The communication is performed on a broadcast channel, so a message sent by a node is received by all the other nodes. The bus access scheme is time-division multiple-access (TDMA) [Fig. 8(b)]. Each node N_i can transmit only during a predetermined time interval, the so-called TDMA slot S_i . In such a slot, a node can send several messages packaged in a frame. We consider that a slot S_i is at least large enough to accommodate the largest message generated by any process assigned to node N_i , so the messages do not have to be split in order to be sent. A sequence of slots corresponding to all the nodes in the architecture is called a TDMA round. A node can have only one slot in a TDMA round. Several TDMA rounds can be combined together in a cycle that is repeated periodically. The sequence and length of the slots are the same for all the TDMA rounds. However, the length and contents of the frames may differ.

Every node has a TTP controller that implements the protocol services and runs independently of the node's CPU. Communication with the CPU is performed through a message base in-

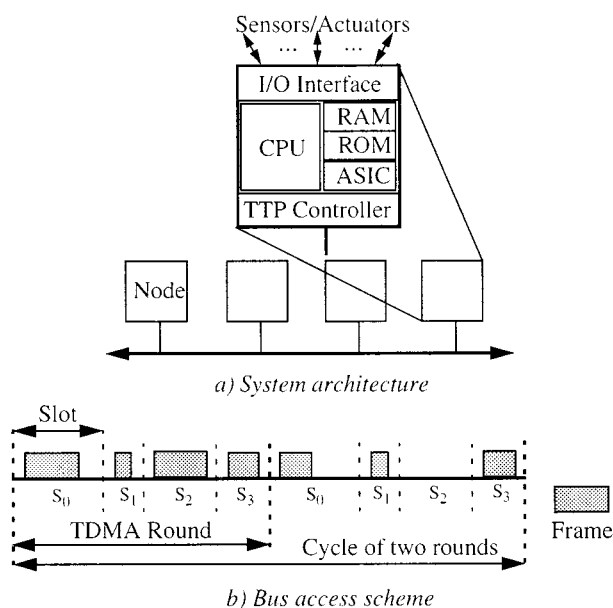


Fig. 8. TTP-based system architecture

terface (MBI), which is usually implemented as a dual-ported RAM (see Fig. 9).

The TDMA access scheme is imposed by a message descriptor list (MEDL) that is located in every TTP controller. The MEDL basically contains the time when a frame has to be sent or received, the address of the frame in the MBI, and the length of the frame. MEDL serves as a schedule table for the TTP controller, which has to know when to send or receive a frame to or from the communication channel.

The TTP controller provides each CPU with a timer interrupt based on a local clock, synchronized with the local clocks of the other nodes. The clock synchronization is done by comparing the *a priori* known time of arrival of a frame with the observed arrival time. By applying a clock synchronization algorithm, TTP provides a global time-base of known precision, without any overhead on the communication.

B. Software Architecture

We have designed a software architecture that runs on the CPU in each node and that has a real-time kernel as its main component. Each kernel has a schedule table that contains all the information needed to take decisions on activation of processes and transmission of messages, based on the values of conditions.

The message passing mechanism is illustrated in Fig. 9, where we have three processes P_1 to P_3 . P_1 and P_2 are mapped to node N_0 that transmits in slot S_0 , and P_3 is mapped to node N_1 that transmits in slot S_1 . Message m_1 is transmitted between P_1 and P_2 that are on the same node, while message m_2 is transmitted from P_1 to P_3 between the two nodes. We consider that each process has its own memory locations for the messages it sends or receives and that the addresses of the memory locations are known to the kernel through the schedule table.

P_1 is activated according to the schedule table, and when it finishes it calls the *send* kernel function in order to send m_1

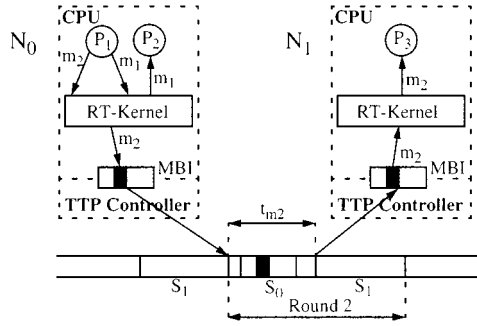


Fig. 9. Message passing mechanism.

and then m_2 . Based on the schedule table, the kernel copies m_1 from the corresponding memory location in P_1 to the memory location in P_2 . When P_2 will be activated, it finds the message in the right location. According to our scheduling policy, whenever a receiving process needs a message, the message is already placed in the corresponding memory location. Thus, there is no overhead on the receiving side for messages exchanged on the same node.

Message m_2 has to be sent from node N_0 to node N_1 . At a certain time, known from the schedule table, the kernel transfers m_2 to the TTP controller by packaging it into a frame in the MBI. Later on, the TTP controller knows from its MEDL when it has to take the frame from the MBI in order to broadcast it on the bus. In our example the timing information in the schedule table of the kernel and the MEDL is determined in such a way that the broadcasting of the frame is done in slot S_0 of Round 2. The TTP controller of node N_1 knows from its MEDL that it has to read a frame from slot S_0 of Round 2 and to transfer it into the MBI. The kernel in node N_1 will read message m_2 from the MBI. When P_3 will be activated based on the local schedule table of node N_1 , it will already have m_2 in its right memory location.

In [51], we presented a detailed discussion concerning the overheads due to the kernel and to every system call. We also presented formulas for derivation of the worst case execution delay of a process, taking into account the overhead of the timer interrupt, the worst case overhead of the process activation, and message passing functions.

VI. SCHEDULING WITH BUS ACCESS OPTIMIZATION

A. Problem Formulation

We consider a system captured as a CPG. The target architecture is as described in Section V. Each process is mapped on a CPU or an ASIC of a node. We are interested in deriving a delay on the system execution time so that this delay is as small as possible, and in synthesizing the local schedule tables for each node, as well as the MEDL for the TTP controllers, which guarantee this delay.

For each message, its length b_{mi} is given. If the message is exchanged by two processes mapped on the same node, the message communication time is completely accounted for in the worst case execution delay of the two processes as shown in Section V-B. Thus, from the scheduling point of view, commu-

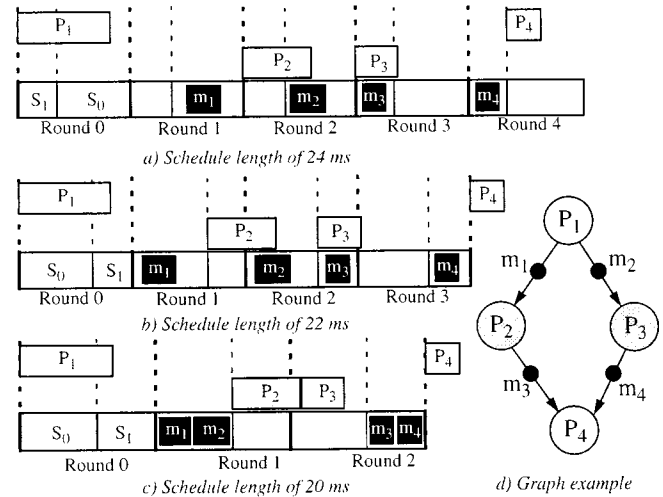


Fig. 10. Scheduling example.

nication of such messages is instantaneous. This is in line with our CPG representation where no communication processes are introduced between processes mapped to the same processor. However, if the message is sent between two processes mapped onto different nodes, the message has to be scheduled according to the TTP protocol. Several messages can be packaged together in the data field of a frame. The number of messages that can be packaged depends on the slot length corresponding to the node. The effective time spent by a message m_i on the bus is $t_{mi} = b_{Si}/T$, where b_{Si} is the length of slot S_i and T is the transmission speed of the channel. In Fig. 9, t_{m2} depicts the time spent by m_2 on the bus. The previous equation shows that the communication time t_{mi} does not depend on the bit length b_{mi} of the message m_i but on the slot length corresponding to the node sending m_i .

The important impact of the communication parameters on the performance of the application is illustrated in Fig. 10. In Fig. 10(d), we have a graph consisting of four processes P_1 to P_4 and four messages m_1 to m_4 . The architecture consists of two nodes interconnected by a TTP channel. The first node, N_0 , transmits on slot S_0 of the TDMA round, and the second node, N_1 , transmits on slot S_1 . Processes P_1 and P_4 are mapped on node N_0 , while processes P_2 and P_3 are mapped on node N_1 . With the TDMA configuration in Fig. 10(a), where slot S_1 is scheduled first and slot S_0 is second, we have a resulting schedule length of 24 ms. However, if we swap the two slots inside the TDMA round without changing their lengths, we can improve the schedule by 2 ms, as seen in Fig. 10(b). Furthermore, if we have the TDMA configuration in Fig. 10(c) where slot S_0 is first, slot S_1 is second, and we increase the slot lengths so that the slots can accommodate both of the messages generated on the same node, we obtain a schedule length of 20 ms, which is optimal. However, increasing the length of slots is not necessarily improving a schedule, as it delays the communication of messages generated by other nodes.

Our optimization strategy, described in the following sections, determines the sequence and length of the slots in a TDMA round with the goal of reducing the delay on the execution time of the system. Before discussing this optimization of

```

Plan_message (time_ready, b_m, Node_m)
  Slot=the slot assigned to Node_m;
  Round=[time_ready/round_length];
  if time_ready - Round*round_length > start_Slot then
    Round = Round+1
  end if;
  while b_m > b_Slot - b_occupied do
    Round = Round+1
  end while;
  return (Round, Slot)
end Plan_message;

```

-- the slot in which the message has to be sent
-- the first round which could be a candidate
-- is the right slot in this round already gone?
-- if yes, take the next round
-- is enough space left in the slot for the message?
-- if not, take the next round
-- return the right round and slot

Fig. 11. Message communication planning.

the bus access scheme, we first analyze how the particularities of the TTP protocol have to be taken into consideration at scheduling of CPGs.

B. Scheduling with a Given Bus Access Scheme

Given a certain bus access scheme, which means a given ordering of the slots in the TDMA round and fixed slot lengths, the CPG has to be scheduled with the goal to minimize the worst case execution delay. This can be performed using our algorithm *List_schedule_CPG* (Fig. 5) presented in Section IV. Two aspects have to be discussed here: the planning of messages in pre-determined slots and the impact of this communication strategy on the priority assignment.

The function *Plan_message* in Fig. 11 is called in order to plan the communication of a message m , with length b_m , generated on $Node_m$, which is ready to be transmitted at $time_ready$. *Plan_message* returns the first round and corresponding slot (the slot corresponding to $Node_m$) that can host the message. In Fig. 11, $round_length$ is the length of a TDMA round expressed in time units (in Fig. 12, for example, $round_length = 18$ ms). The first round after $time_ready$ is the initial candidate to be considered. For this round, however, it can be too late to catch the right slot, in which case the next round is selected. When a candidate round is selected, we have to check that there is enough space left in the slot for our message ($b_{occupied}$ represents the total number of bits occupied by messages already scheduled in the respective slot of that round). If no space is left, the communication has to be delayed for another round.

With this message planning scheme, the algorithm in Fig. 5 will generate correct schedules for a TTP-based architecture, with guaranteed worst case execution delays. However, the quality of the schedules can be much improved by adapting the priority assignment scheme so that particularities of the communication protocol are taken into consideration.

Let us consider the graph in Fig. 12(c), and suppose that the list scheduling algorithm has to decide between scheduling process P_1 or P_2 , which are both ready to be scheduled on the same programmable processor. The worst case execution time of the processes is depicted on the right side of the respective node and is expressed in milliseconds. The architecture consists of two nodes interconnected by a TTP channel. Processes P_1 and P_2 are mapped on node N_1 , while processes P_3 and P_4 are mapped on node N_0 . Node N_0 transmits on slot S_0 of the TDMA round, and N_1 transmits on slot S_1 . Slot S_0 has a length of 10 ms, while slot S_1 has a length of 8 ms. For simplicity, we

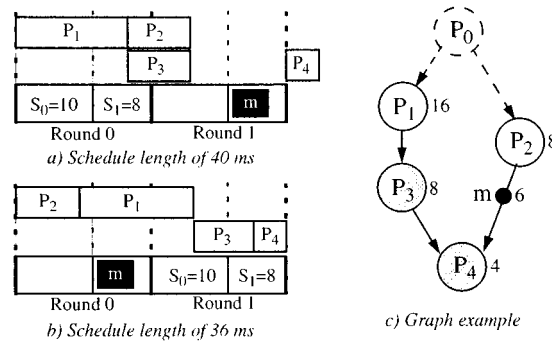


Fig. 12. Priority function example.

suppose that there is no message transferred between P_1 and P_3 . PCP (see Section IV-D) assigns a higher priority to P_1 because it has a partial critical path of 12, starting from P_3 , longer than the partial critical path of P_2 which is ten and starts from m . This results in a schedule length of 40 ms depicted in Fig. 12(a). On the other hand, if we schedule P_2 first, the resulting schedule, depicted in Fig. 12(b), is only 36 ms.

This apparent anomaly is due to the fact that the way we have computed PCP priorities, considering message communication as a simple activity of delay 6 ms, is not realistic in the context of a TDMA protocol. Let us consider the particular TDMA configuration in Fig. 12 and suppose that the scheduler has to decide at $t = 0$ which one of the processes P_1 or P_2 to schedule. If P_2 is scheduled, the message is ready to be transmitted at $t' = 8$. Based on a computation similar to that used in Fig. 11, it follows that message m will be placed in round $\lceil 8/18 \rceil = 0$, and it arrives in time to get slot S_1 of that round ($time_ready = 8 < start_{S_1} = 10$). Thus, m arrives at $t_{arr} = 18$, which means a delay relative to $t' = 8$ (when the message was ready) of $\delta = 10$. This is the delay that should be considered for computing the partial critical path of P_2 , which now results in $\delta + t_{P_4} = 14$ (longer than the one corresponding to P_1).

The obvious conclusion is that priority estimation has to be based on message planning with the TDMA scheme. Such an estimation, however, cannot be performed statically before scheduling. If we take the same example in Fig. 12, but consider that the priority-based decision is taken by the scheduler at $t = 5$, m will be ready at $t' = 13$. This is too late for m to get into slot S_1 of round 0. The message arrives with round 1 at $t_{arr} = 36$. This leads to a delay due to the message passing of $\delta = 36 - 13 = 23$, different from the one computed above.

```

Optimize_access ()
  for i:= 0 to Nr_slot-1 do          -- creates the initial, straightforward solution
    NodeSi=Ni; lengthSi=min_lengthSi
  end for;
  for i:= 0 to Nr_slot-1 do          -- over all slots
    for j:= i to Nr_slot-1 do        -- over all slots which have not yet been allocated a node and slot length
      swap values (NodeSi, lengthSi) with (NodeSj, lengthSj);
      for all slot lengths λ, larger than lengthSi do -- initially, lengthSi has the minimal allowed value
        lengthSi=λ;
        List_schedule_CPG ( ... );
        remember best_solution=(NodeSi, lengthSi), with the smallest δmax produced by List_schedule_CPG
      end for;
      swap back values (NodeSi, lengthSi) with (NodeSj, lengthSj) to the state before entering the for cycle
    end for;
    bind (NodeSi, lengthSi)=best_solution; -- slot Si gets a node allocated and a length fixed
  end for;
end Optimize_access;

```

Fig. 13. Optimization of the bus access scheme.

We introduce a new priority function, the modified PCP (MPCP), which is computed during scheduling, whenever several processes are in competition to be scheduled on the same resource. Similar to PCP, the priority metric is the length of that portion of the critical path corresponding to a process P_i , which starts with the first successor of P_i that is assigned to a processor different from $M(P_i)$. The critical path estimation starts with time t at which the processes in competition are ready to be scheduled on the available resource. During the partial traversal of the graph, the delay introduced by a certain node P_j is estimated as follows:

$$\delta_{P_j} = \begin{cases} t_{P_j}, & \text{if } P_j \text{ is not a message passing} \\ t_{arr} - t', & \text{if } P_j \text{ is a message passing.} \end{cases}$$

t' is the time when the node generating the message terminates (and the message is ready) and t_{arr} is the time when the slot to which the message is supposed to be assigned has arrived. The slot is determined as in Fig. 11, but without taking into consideration space limitations in slots. As the experimental results (Section VII) show, using MPCP instead of PCP for the TTP-based architecture results in an important improvement of the quality of generated schedules, with a slight increase in scheduling time.

C. Optimization of the Bus Access Scheme

In the previous section, we have shown how our algorithm *List_schedule_CPG* (Fig. 5) can produce an efficient schedule for a CPG, given a certain TDMA bus access scheme. However, as discussed in Section VI-A, both the ordering of slots and the slot lengths strongly influence the worst case execution delay of the system. In Fig. 13, we show a heuristic that, based on a greedy approach, determines an ordering of slots and their lengths so that the worst case delay corresponding to a certain CPG is as small as possible.

The initial solution, the “straightforward” one, assigns in order nodes to the slots ($Node_{S_i} = N_i$) and fixes the slot length $length_{S_i}$ to the minimal allowed value, which is equal to the length of the largest message generated by a process assigned to $Node_{S_i}$. The algorithm starts with the first slot and

tries to find the node that, when transmitting in this slot, will minimize the worst case delay of the system, as produced by *List_Schedule_CPG*. Simultaneously with searching for the right node to be assigned to the slot, the algorithm looks for the optimal slot length. Once a node is selected for the first slot and a slot length fixed, the algorithm continues with the next slots, trying to assign nodes (and fix slot lengths) from those nodes that have not yet been assigned.

When calculating the length of a certain slot, a first alternative could be to try all the slot lengths l allowed by the protocol. Such an approach starts with the minimum slot length determined by the largest message to be sent from the candidate node, and it continues incrementing with the smallest data unit (e.g., 2 bits) up to the largest slot length determined by the maximum allowed data field in a TTP frame (e.g., 32 bits, depending on the controller implementation). We call this alternative *Optimize_access_1*. A second alternative, *Optimize_access_2*, is based on feedback from the scheduling algorithm, which recommends slot sizes to be tried out. Before starting the actual optimization process for the bus access scheme, a scheduling of the straightforward solution is performed that generates the recommended slot lengths. These lengths are produced by the *Plan_message* function (Fig. 11) whenever a new round has to be selected because of lack of space in the current slot. In such a case, the slot length that would be needed in order to accommodate the new message is added to the list of recommended lengths for the respective slot. With this alternative, the optimization algorithm in Fig. 13 only selects among the recommended lengths when searching for the right dimension of a certain slot.

As the experimental results show (see Section VII), optimization of the bus access scheme can produce huge improvements in terms of performance. As expected, the *Optimize_access_2* alternative is much faster than the first one, while the quality of the results produced is only slightly lower.

VII. EXPERIMENTAL RESULTS

In the following two sections, we show a series of experiments that demonstrate the effectiveness of the proposed algo-

gorithms. The first set of results is related to the scheduling of conditional process graphs, while the second set targets the problem of scheduling with optimization of the bus access scheme. As a general strategy, we have evaluated our algorithms performing experiments on a large number of test cases generated for experimental purpose. We then have validated the proposed approaches using real-life examples. All experiments were run on SPARCstation 20.

A. Evaluation of the Scheduling Algorithm

As discussed in Section III-B, there are N_{alt} alternative tracks through a CPG, and a schedule could be generated for each one separately. Suppose that δ_M is the longest of these schedules. This, however, does not mean that the worst case delay δ_{max} , corresponding to the CPG, is guaranteed to be δ_M . Such a delay cannot be guaranteed in theory, as the values of conditions, and thus the actual track to be followed, cannot be predicted. The objective of our scheduling heuristic is to generate a schedule table so that the difference $\delta_{max} - \delta_M$ is minimized.

For evaluation of the scheduling algorithm, we used 1080 conditional process graphs generated for experimental purpose; 360 graphs have been generated for each dimension of 60, 80, and 120 processes. The number of alternative tracks through the graphs is 10, 12, 18, 24, or 32. Execution times were assigned randomly using both uniform and exponential distribution. We considered architectures consisting of one ASIC, one to 11 processors, and one to eight buses.

Fig. 14(a) presents the percentage increase of the worst case delay δ_{max} over the delay δ_M of the longest track. The delay δ_M has been obtained by scheduling separately each alternative track through the respective CPG, using PCP list scheduling, and selecting the delay that corresponds to the longest track. The average increase is between 0.1% and 8.1% and, practically, does not depend on the number of processes in the graph but only on the number of alternative tracks. It is worth mentioning that a zero increase ($\delta_{max} = \delta_M$) was produced for 90% of the graphs with ten alternative tracks, 82% with 12 tracks, 57% with 18 tracks, 46% with 24 tracks, and 33% with 32 tracks.

Concerning execution time, the interesting aspect is how the algorithm scales with the number of alternative tracks and that of processes. The worst case complexity of the scheduling algorithm depends on the number of tracks, which theoretically can grow exponentially. However, such an explosion is unlikely for practically significant applications. Fig. 14(b) shows the average execution time for the scheduling algorithm as a function of the number of alternative tracks. We observe very small execution times for even large graphs and very good scaling with the number of alternative tracks. The increase of execution time with the number of processes, for a given number of alternative tracks, is practically linear, which corresponds to the theoretical complexity of list scheduling algorithms [52], [53].

One of the very important applications of our scheduling algorithm is for performance estimation during design space exploration. We have performed such an experiment as part of a project aiming to implement the operation and maintenance (OAM) functions corresponding to the F4 level of the ATM protocol layer [54]. Fig. 15(a) shows an abstract model of the

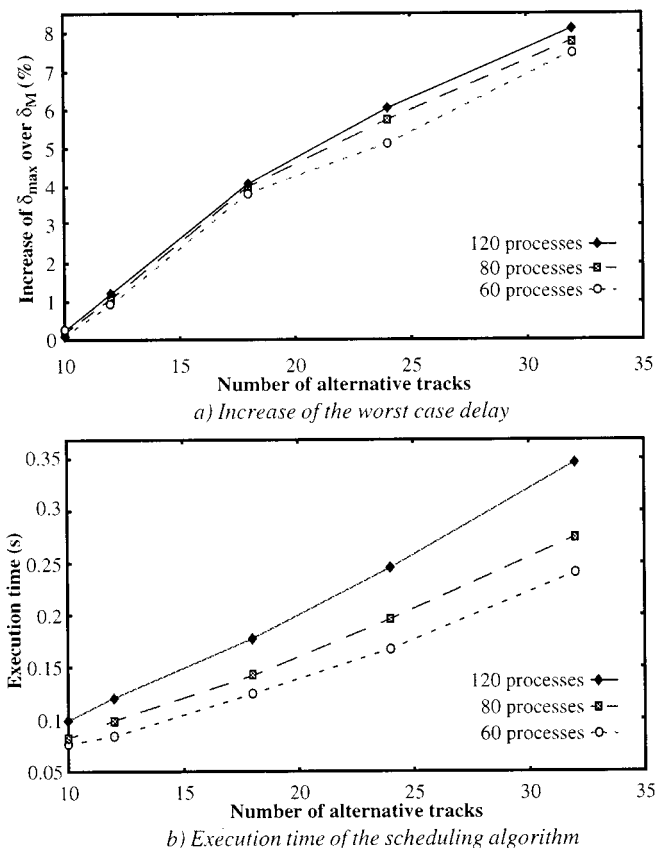


Fig. 14. Evaluation of the scheduling algorithm for CPGs.

ATM switch. Through the switching network, cells are routed between the n input and q output lines. In addition, the ATM switch also performs several OAM related tasks.

In [50], we discussed hardware/software partitioning of the OAM functions corresponding to the F4 level. We concluded that filtering of the input cells and redirecting of the OAM cells toward the OAM block have to be performed in hardware as part of the line interfaces (LI). The other functions are performed by the *OAM block* and can be implemented in software.

We have identified three independent modes in the functionality of the OAM block. Depending on the content of the input buffers [Fig. 15(b)], the OAM block switches between these three modes. Execution in each mode is controlled by a statically generated schedule table for the respective mode. We specified the functionality corresponding to each mode as a set of interacting VHDL processes. These specifications have then been translated to the corresponding CPGs. Table II shows the characteristics of the resulting CPGs. The main objective of this experiment was to estimate, using our scheduling algorithm, the worst case delays in each mode for different alternative architectures of the OAM block. Based on these estimations as well as on the particular features of the environment in which the switch will be used, an appropriate architecture can be selected and the dimensions of the buffers can be determined.

Fig. 15(b) shows a possible implementation architecture of the OAM block, using one processor and one memory module (1P/1M). Our experiments included also architecture models with two processors and one memory module (2P/1M), as well

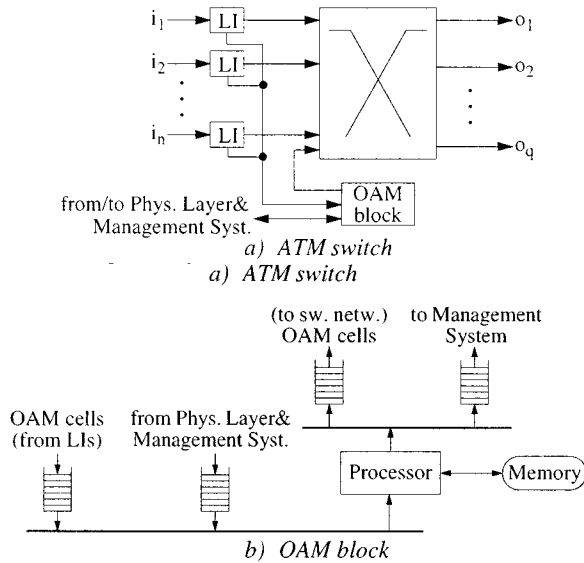


Fig. 15. ATM switch with OAM block

as structures consisting of one (respectively, two) processors and two memory modules (1P/2M, 2P/2M). The inclusion of alternatives with two memory modules is justified by the fact that the information processed by the OAM block is organized in two main tables that potentially could be accessed in parallel. The target architectures are based on two types of processors: one ($\mu p1$) running at 80 MHz and another ($\mu p2$) running at 120 MHz. For each architecture, processes have been assigned to processors taking into consideration the potential parallelism of the CPGs and the amount of communication between processes. The worst case delays that resulted after generation of the schedule table for each of the three modes are given in Table II. As expected, using a faster processor reduces the delay in each of the three modes. Introducing an additional processor, however, has no effect on the execution delay in mode 2, which does not present any potential parallelism. In mode 3, the delay is reduced by using two $\mu p1$ processors instead of one. For the faster, $\mu p2$ processor, however, the worst case delay cannot be improved by introducing an additional processor. Using two processors will always improve the worst case delay in *mode 1*. As for the additional memory module, only in mode 1 does the model contain memory accesses that are potentially executed in parallel. Table II shows that only for the architecture consisting of two $\mu p2$ processors, providing an additional memory module pays back by a reduction of the worst case delay in mode 1. The reason is that, with the process execution times corresponding to this processor and the 2P/1M architecture, the track containing parallel memory accesses is the one that dictates the worst case execution time. Thus, adding a second memory module results in a reduced worst case delay of the system in mode 1.

B. Evaluation of Scheduling with Bus Access Optimization

In this set of experiments, we were interested to investigate the efficiency of our scheduling algorithm in the context of the TDMA-based protocol and the potential of our optimization strategies for the bus access scheme. We considered TTP-based architectures consisting of two, four, six, eight, and ten nodes.

Forty processes were assigned to each node, resulting in graphs of 80, 160, 240, 320, and 400 processes. Thirty CPGs were generated for each graph dimension; thus a total of 150 CPGs were used for experimental evaluation. Execution times and message lengths were assigned randomly using both uniform and exponential distribution. For the communication channel, we considered a transmission speed of 256 kbps and a length below 20 m. The maximum length of the data field was 8 bytes, and the frequency of the TTP controller was chosen to be 20 MHz.

The first results concern the improvement of the schedules produced by our algorithm when using the MPCP priority function instead of the one based on the general PCP priority. In order to compare the two priority functions, the 150 CPGs were scheduled, considering the TTP-based architectures presented above, using first PCP for priority assignment and then MPCP. We calculated the average percentage deviations of the schedule lengths produced with MPCP and PCP for each graph, from the length of the best schedule among the two. The results are depicted in Fig. 16(a). The diagram shows an important improvement of the resulted schedules if the TDMA-specific priority function MPCP is used. On average, the deviation with MPCP is 11.34 times smaller than with PCP. However, due to its dynamic nature, MPCP implies a slightly larger execution time, as shown in Fig. 16(b).

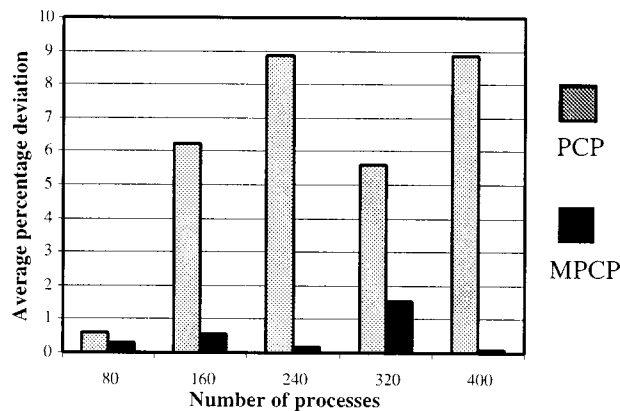
In the following experiments, we were interested to check the potential of the algorithm presented in Section VI-C to improve the generated schedules by optimizing the bus access scheme. We compared schedule lengths obtained for the 150 CPGs considering four different bus access schemes: the straightforward solution, the optimized schemes generated with the two alternatives of our greedy algorithm (*Optimize_access_1* and *Optimize_access_2*), and a near-optimal scheme. The near-optimal scheme was produced using a simulated annealing (SA)-based algorithm for bus access optimization, which is presented in [51]. Very long and extensive runs have been performed with the SA algorithm for each graph, and the best ever solution produced has been considered as the near optimum for that graph.

Table III presents the average and maximum percentage deviation of the schedule lengths obtained with the straightforward solution and with the two optimized schemes from the length obtained with the near-optimal scheme. For each of the graph dimensions, the average optimization time, expressed in seconds, is also given. The first conclusion is that by considering the optimization of the bus access scheme, the results improve significantly compared to the straightforward solution. The greedy heuristic performs well for all the graph dimensions. As expected, the alternative *Optimize_access_1* (which considers all allowed slot lengths) produces slightly better results, on average, than *Optimize_access_2*. However, the execution times are much smaller for *Optimize_access_2*. It is interesting to mention that the average execution times for the SA algorithm, needed to find the near-optimal solutions, are between 5 min for the CPGs with 80 processes and 275 min for 400 processes [51].

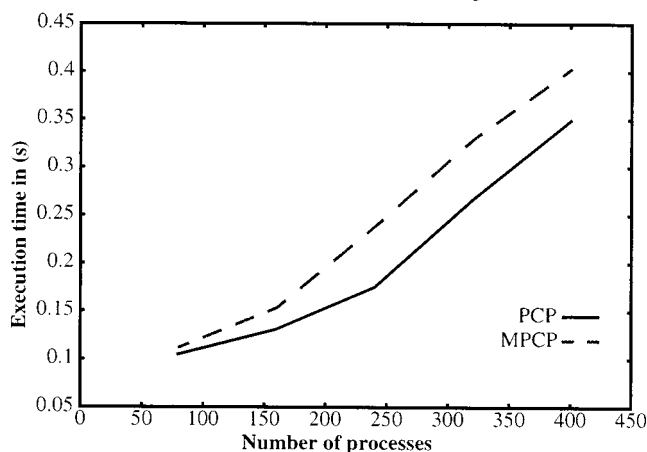
A typical safety-critical application with hard real-time constraints, to be implemented on a TTP-based architecture, is a vehicle cruise controller (CC). We have considered a CC system

TABLE II
WORST CASE DELAYS FOR THE OAM BLOCK

| Mode | System model | | δ_{\max} (ns) | | | | | | | | | |
|--------|------------------|---------------|----------------------|----------|----------|----------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| | nr. of processes | nr. of tracks | 1P/1M | | 1P/2M | | 2P/1M | | | 2P/2M | | |
| | | | $\mu p1$ | $\mu p2$ | $\mu p1$ | $\mu p2$ | $2\times\mu p1$ | $2\times\mu p2$ | $\mu p1+\mu p2$ | $2\times\mu p1$ | $2\times\mu p2$ | $\mu p1+\mu p2$ |
| Mode 1 | 32 | 6 | 4471 | 2701 | 4471 | 2701 | 2932 | 2131 | 2532 | 2932 | 1932 | 2532 |
| Mode 2 | 23 | 3 | 1732 | 1167 | 1732 | 1167 | 1732 | 1167 | 1167 | 1732 | 1167 | 1167 |
| Mode 3 | 42 | 8 | 5852 | 3548 | 5852 | 3548 | 5033 | 3548 | 3548 | 5033 | 3548 | 3548 |



a) Percentage deviation of schedule lengths



b) Execution time of the scheduling algorithm

Fig. 16. Scheduling with PCP and MPCP for TTP-based architectures.

derived from a requirement specification provided by the industry. The CC described in this specification delivers the following functionality: it maintains a constant speed for speeds over 35 km/h and under 200 km/h, offers an interface (buttons) to increase or decrease the reference speed, and is able to resume its operation at the previous reference speed. The CC operation is suspended when the driver presses the brake pedal.

The specification assumes that the CC will operate in an environment consisting of several nodes interconnected by a TTP channel. There are five nodes that functionally interact with the CC system: the antiblocking system, the transmission control module, the engine control module, the electronic throttle module, and the central electronic module. It has been decided

to distribute the functionality (processes) of the CC over these five nodes.

The CPG corresponding to the CC system consists of 32 processes and includes two alternative tracks. The maximum allowed delay is 110 ms. For our model, the straightforward solution for bus access resulted in a schedule corresponding to a maximal delay of 114 ms (which does not meet the deadline) when PCP was used as a priority function, while using MPCP we obtained a schedule length of 109 ms. Both of the greedy heuristics for bus access optimization produced solutions so that the worst case delay was reduced to 103 ms. The near-optimal solution (produced with the SA-based approach) results in a delay of 97 ms.

VIII. CONCLUSIONS

We have presented an approach to process scheduling for the synthesis of embedded systems implemented on architectures consisting of several programmable processors and application-specific hardware components. The approach is based on an abstract graph representation that captures, at process level, both dataflow and the flow of control. The scheduling problem has been considered in strong interrelation with the problem of communication in distributed systems.

We first presented a general approach to process scheduling with control and data dependencies, considering a generic bus-based distributed architecture. The proposed algorithm is based on a list scheduling approach and statically generates a schedule table that contains activation times for processes and communications. The main problems that have been solved in this context are the minimization of the worst case delay and the generation of a logically and temporally deterministic table, taking into consideration communication times and the sharing of the communication support.

We have further investigated the impact of particular communication infrastructures and protocols on the overall performance and, specially, how the requirements of such an infrastructure have to be considered for process and communication scheduling. Considering a TTP-based system architecture, we have shown that the general scheduling algorithm for conditional process graphs can be successfully applied if the strategy for message planning is adapted to the requirements of the TDMA protocol. At the same time, the quality of generated schedules has been much improved after adjusting the priority function used by the scheduling algorithm to the particular communication protocol.

TABLE III
EVALUATION OF THE BUS ACCESS OPTIMIZATION ALGORITHM

| Nr. of proc. | Straightforward solution | | Optimize_access_1 | | | Optimize_access_2 | | |
|--------------|--------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| | average deviation | maximal deviation | average deviation | maximal deviation | optimization time | average deviation | maximal deviation | optimization time |
| 80 | 3.16% | 21% | 0.02% | 0.5% | 0.25s | 1.8% | 19.7% | 0.04s |
| 160 | 14.4% | 53.4% | 2.5% | 9.5% | 2.07s | 4.9% | 26.3% | 0.28s |
| 240 | 37.6% | 110% | 7.4% | 24.8% | 10.46s | 9.3% | 31.4% | 1.34s |
| 320 | 51.5% | 135% | 8.5% | 31.9% | 34.69s | 12.1% | 37.1% | 4.8s |
| 400 | 48% | 135% | 10.5% | 32.9% | 56.04s | 11.8% | 31.6% | 8.2s |

However, not only do particulars of the underlying architecture have to be considered during scheduling but also the parameters of the communication protocol should also be adapted to fit the particular embedded application. We have shown that important performance gains can be obtained, without any additional cost, by optimizing the bus access scheme. The optimization algorithm, which now implies both process scheduling and optimization of the parameters related to the communication protocol, generates an efficient bus access scheme as well as the schedule tables for activation of processes and communications.

The algorithms have been evaluated based on extensive experiments using a large number of graphs generated for experimental purpose as well as real-life examples.

There are several aspects that were omitted from the discussion in this paper. In [55], we have analyzed the optimization of a TDMA bus access scheme in the context of priority-based preemptive scheduling. There we also considered the possibility of messages being split over several successive frames. We neither insist here on the relatively simple procedure for postprocessing of the schedule table, during which the table can be simplified for certain situations in which identical activation times are scheduled for a given process on different columns. During postprocessing, the table is also split into subtables containing the particular activities to be performed by a certain processor.

REFERENCES

- [1] G. De Micheli and M. G. Sami, Eds., *Hardware/Software Co-Design*. Norwell, MA: NATO ASI 1995, Kluwer Academic, 1996.
- [2] G. De Micheli and R. K. Gupta, "Hardware/software co-design," *Proc. IEEE*, vol. 85, no. 3, pp. 349–365, 1997.
- [3] R. Ernst, "Codesign of embedded systems: Status and trends," *IEEE Design Test Comput.*, pp. 45–54, Apr.–June 1998.
- [4] D. D. Gajski and F. Vahid, "Specification and design of embedded hardware-software systems," *IEEE Design Test Comput.*, pp. 53–67, Spring 1995.
- [5] J. Staunstrup and W. Wolf, Eds., *Hardware/Software Co-Design: Principles and Practice*. Norwell, MA: Kluwer Academic, 1997.
- [6] W. Wolf, "Hardware-software co-design of embedded systems," *Proc. IEEE*, vol. 82, no. 7, pp. 967–989, 1994.
- [7] J. D. Ullman, "NP-complete scheduling problems," *J. Comput. Syst. Sci.*, vol. 10, pp. 384–393, 1975.
- [8] A. Doboli and P. Eles, "Scheduling under control dependencies for heterogeneous architectures," in *Proc. Int. Conf. Computer Design (ICCD)*, 1998, pp. 602–608.
- [9] P. Eles, K. Kuchcinski, Z. Peng, A. Doboli, and P. Pop, "Scheduling of conditional process graphs for the synthesis of embedded systems," in *Proc. Design Aut. Test Eur.*, 1998, pp. 132–138.
- [10] R. Ernst and W. Ye, "Embedded program timing analysis based on path clustering and architecture classification," in *Proc. Int. Conf. CAD*, 1997, pp. 598–604.
- [11] J. Gong, D. D. Gajski, and S. Narayan, "Software estimation using a generic-processor model," in *Proc. Eur. Design Test Conf.*, 1995, pp. 498–502.
- [12] J. Henkel and R. Ernst, "A path-based technique for estimating hardware run-time in Hw/Sw-cosynthesis," in *Proc. Int. Symp. Syst. Synthesis*, 1995, pp. 116–121.
- [13] Y. S. Li and S. Malik, "Performance analysis of embedded software using implicit path enumeration," in *Proc. ACM/IEEE DAC*, 1995, pp. 456–461.
- [14] T. Lundqvist and P. Stenström, "An integrated path and timing analysis method based on cycle-level symbolic execution," *Real-Time Syst.*, vol. 17, no. 2/3, pp. 183–207, 1999.
- [15] S. Malik, M. Martonosi, and Y. S. Li, "Static timing analysis of embedded software," in *Proc. ACM/IEEE DAC*, 1997, pp. 147–152.
- [16] K. Suzuki and A. Sangiovanni-Vincentelli, "Efficient software performance estimation methods for hardware/software codesign," in *Proc. ACM/IEEE DAC*, 1996, pp. 605–610.
- [17] C. L. Liu and J. W. Layland, "Scheduling algorithms for multiprogramming in a hard real-time environment," *J. ACM*, vol. 20, no. 1, pp. 46–61, 1973.
- [18] K. Tindell and J. Clark, "Holistic schedulability analysis for distributed hard real-time systems," *Microprocess. Microprogram.*, vol. 40, pp. 117–134, 1994.
- [19] N. C. Audsley, A. Burns, R. I. Davis, K. Tindell, and A. J. Wellings, "Fixed priority pre-emptive scheduling: An historical perspective," *Real-Time Syst.*, vol. 8, no. 2/3, pp. 173–198, 1995.
- [20] F. Balarin, L. Lavagno, P. Murthy, and A. Sangiovanni-Vincentelli, "Scheduling for embedded real-time systems," *IEEE Design Test Comput.*, pp. 71–82, Jan.–Mar. 1998.
- [21] T. Y. Yen and W. Wolf, *Hardware-Software Co-Synthesis of Distributed Embedded Systems*. Norwell, MA: Kluwer Academic, 1997.
- [22] C. Lee, M. Potkonjak, and W. Wolf, "Synthesis of hard real-time application specific systems," *Design Automat. Embedded Syst.*, vol. 4, no. 4, pp. 215–241, 1999.
- [23] B. P. Dave and N. K. Jha, "COHRA: Hardware-software cosynthesis of hierarchical heterogeneous distributed systems," *IEEE Trans. Computer-Aided Design*, vol. 17, no. 10, pp. 900–919, 1998.
- [24] B. P. Dave, G. Lakshminarayana, and N. J. Jha, "COSYN: Hardware-software co-synthesis of heterogeneous distributed embedded systems," *IEEE Trans. VLSI Syst.*, vol. 7, no. 1, pp. 92–104, 1999.
- [25] P. Chou and G. Borriello, "Interval scheduling: Fine-grained code scheduling for embedded systems," in *Proc. ACM/IEEE DAC*, 1995, pp. 462–467.
- [26] R. K. Gupta, *Co-Synthesis of Hardware and Software for Digital Embedded Systems*. Boston, MA: Kluwer Academic, 1995.
- [27] V. Mooney, T. Sakamoto, and G. De Micheli, "Run-time scheduler synthesis for hardware-software systems and application to robot control design," in *Proc. Int. Workshop Hardware-Software Co-Design*, 1997, pp. 95–99.
- [28] H. Kopetz, *Real-Time Systems-Design Principles for Distributed Embedded Applications*. Norwell, MA: Kluwer Academic, 1997.
- [29] E. G. Coffman Jr and R. L. Graham, "Optimal scheduling for two processor systems," *Acta Inform.*, vol. 1, pp. 200–213, 1972.
- [30] P. B. Jorgensen and J. Madsen, "Critical path driven cosynthesis for heterogeneous target architectures," in *Proc. Int. Workshop Hardware-Software Co-Design*, 1997, pp. 15–19.

- [31] Y. K. Kwok and I. Ahmad, "Dynamic critical-path scheduling: An effective technique for allocating task graphs to multiprocessors," *IEEE Trans. Parallel Distrib. Syst.*, vol. 7, no. 5, pp. 506–521, 1996.
- [32] M. Y. Wu and D. D. Gajski, "Hypertool: A programming aid for message-passing systems," *IEEE Trans. Parallel Distrib. Syst.*, vol. 1, no. 3, pp. 330–343, 1990.
- [33] P. Eles, K. Kuchcinski, Z. Peng, A. Doboli, and P. Pop, "Process scheduling for performance estimation and synthesis of hardware/software systems," in *Proc. Euromicro Conf.*, 1998, pp. 168–175.
- [34] H. Kasahara and S. Narita, "Practical multiprocessor scheduling algorithms for efficient parallel processing," *IEEE Trans. Comput.*, vol. C-33, no. 11, pp. 1023–1029, 1984.
- [35] A. Bender, "Design of an optimal loosely coupled heterogeneous multiprocessor system," in *Proc. ED&TC*, 1996, pp. 275–281.
- [36] S. Prakash and A. Parker, "SOS: Synthesis of application-specific heterogeneous multiprocessor systems," *J. Parallel Distrib. Comput.*, vol. 16, pp. 338–351, 1992.
- [37] K. Kuchcinski, "Embedded system synthesis by timing constraint solving," in *Proc. Int. Symp. Syst. Synth.*, 1997, pp. 50–57.
- [38] A. Dasdan, D. Ramanathan, and R. K. Gupta, "A timing-driven design and validation methodology for embedded real-time systems," *ACM Trans. Des. Aut. Electron. Syst.*, vol. 3, no. 4, pp. 533–553, 1998.
- [39] K. Strehl, L. Thiele, D. Ziegenbein, R. Ernst, and J. Teich, "Scheduling hardware/software systems using symbolic techniques," in *Proc. Int. Workshop Hardware-Software Co-Design*, 1999, pp. 173–177.
- [40] D. Ziegenbein, K. Richter, R. Ernst, J. Teich, and L. Thiele, "Representation of process model correlation for scheduling," in *Proc. Int. Conf. CAD*, 1998, pp. 54–61.
- [41] P. H. Chou, R. B. Ortega, and G. Borriello, "The Chinook hardware/software co-synthesis system," in *Proc. Int. Symp. Syst. Synthesis*, 1995, pp. 22–27.
- [42] J. M. Daveau, T. B. Ismail, and A. A. Jerraya, "Synthesis of system-level communication by an allocation-based approach," in *Proc. Int. Symp. Syst. Synthesis*, 1995, pp. 150–155.
- [43] P. V. Knudsen and J. Madsen, "Integrating communication protocol selection with hardware/software codesign," *IEEE Trans. Computer-Aided Design*, vol. 18, no. 8, pp. 1077–1095, 1999.
- [44] S. Narayan and D. D. Gajski, "Synthesis of system-level bus interfaces," in *Proc. Eur. Design Test Conf.*, 1994, pp. 395–399.
- [45] R. B. Ortega and G. Borriello, "Communication synthesis for distributed embedded systems," in *Proc. Int. Conf. CAD*, 1998, pp. 437–444.
- [46] K. Tindell, A. Burns, and A. J. Wellings, "Calculating controller area network (CAN) message response times," *Contr. Eng. Practice*, vol. 3, no. 8, pp. 1163–1169, 1995.
- [47] H. Ermedahl, H. Hansson, and M. Sjödin, "Response-time guarantees in ATM networks," in *Proc. IEEE Real-Time Systems Symp.*, 1997, pp. 274–284.
- [48] H. Kopetz and G. Grünsteidl, "TTP—A protocol for fault-tolerant real-time systems," *IEEE Comput.*, vol. 27, no. 1, pp. 14–23, 1997.
- [49] "X-by-wire consortium," URL: <http://www.vmars.tuwien.ac.at/projects/xbywire/>.
- [50] P. Eles, Z. Peng, K. Kuchcinski, and A. Doboli, "System level hardware/software partitioning based on simulated annealing and tabu search," *Design Automat. Embedded Syst.*, vol. 2, no. 1, pp. 5–32, 1997.
- [51] P. Pop, P. Eles, and Z. Peng, "Scheduling with optimized communication for time-triggered embedded systems," in *Proc. Int. Workshop Hardware-Software Co-Design*, 1999, pp. 178–182.
- [52] G. De Micheli, *Synthesis and Optimization of Digital Circuits*: McGraw-Hill, 1994.
- [53] S. H. Gerez, *Algorithms for VLSI Design Automation*. New York: Wiley, 1999.
- [54] T. M. Chen and S. S. Liu, *ATM Switching Systems*. Norwood, MA: Artech House, 1995.
- [55] P. Pop, P. Eles, and Z. Peng, "Bus access optimization for distributed embedded systems based on schedulability analysis," in *Proc. Design Aut. Test Eur.*, 2000.



Petru Eles (M'99) received the M.S. degree in computer science from the Politehnica University Timisoara, Romania, in 1979 and the Ph.D. degree in computer science from the Politehnica University Bucuresti, Romania, in 1993.

He is currently an Associate Professor with the Department of Computer and Information Science at Linköping University, Sweden. His research interests include design of embedded systems, hardware/software codesign, real-time systems, system specification and testing, and computer-aided design for digital systems. He has published extensively in these areas and has coauthored several books, among them *System Synthesis with VHDL* (Norwell, MA: Kluwer Academic, 1997).

Dr. Eles was a corecipient of best paper awards at the 1992 and 1994 European Design Automation Conference.



Alex Doboli (S'99) received the M.S. and Ph.D. degrees in computer science from Politehnica University Timisoara, Romania, in 1990 and 1997, respectively. He is currently pursuing the Ph.D. degree in computer engineering at the University of Cincinnati, Cincinnati, OH.

His research interest is in VLSI design automation, with special interest in mixed-signal CAD, hardware–software codesign, and CAD for reconfigurable computing.

Mr. Doboli is a member of Sigma XI and ACM.



Paul Pop (S'99) received the M.S. degree in computer science from the Politehnica University Timisoara, Romania, in 1997. He is currently pursuing the Ph.D. degree in computer science at Linköping University, Sweden.

His research interests include hardware/software codesign, systems engineering, and real-time systems.



Zebo Peng (M'91) received the Ph.D. degree in computer science from Linköping University, Sweden, in 1987.

He is Professor and Chair of Computer Systems and Director of the Embedded Systems Laboratory (ESLAB) at Linköping University. His current research interests include design and test of embedded systems, electronic design automation, design for testability, hardware/software codesign, and real-time systems. He has published more than 90 technical papers in these areas and is coauthor of

System Synthesis with VHDL (Norwell, MA: Kluwer Academic, 1997).

Dr. Peng was corecipient of two best paper awards at the European Design Automation Conferences in 1992 and 1994.