

CAST – A Task-Level Concurrency Analysis Tool

Sander Stuijk, Jan Ypma and Twan Basten

Eindhoven University of Technology, PO Box 513, NL-5600 MB, Eindhoven, The Netherlands.
s.stuijk@tue.nl

Keywords: Kahn Process Networks, Embedded-Systems Design, Visualization, Streaming, Concurrency.

Abstract

CAST is a system-level software tool that supports the design of concurrent systems. It assumes that the functional model of a system is specified as a set of compute nodes that communicate with each other using point-to-point connections. CAST provides concurrency analysis, support for refinement of the concurrency in the model to implement, and guidance in the design-space exploration. Its main focus is on streaming applications. This paper provides a brief tutorial on CAST. It describes the implementation of the techniques used to model and analyze the concurrency in a specification. Visualization of the metrics is key for performing a guided design-space exploration, which is why visualization is an integral part of CAST. This paper presents therefore also our solutions for visualizing concurrency in a specification.

1 Introduction

Next generations of embedded multi-media systems require high compute power combined with a low energy consumption for use in mobile applications that provide streaming video, audio and/or graphics. To realize these systems, (single-chip) multi-processor systems are becoming a trend. These systems are inherently concurrent. To exploit this concurrency, the parallelism available in an application mapped onto the multi-processor system must be made visible in the mapping trajectory. See Fig. 1 for an overview of this trajectory. To extract parallelism, a specification should allow reasoning about the concurrency in the application and about how to exploit this concurrency. In this paper, we present a software implementation (CAST) of a concurrency model that allows task-level architecture-independent concurrency optimization in executable specifications (source code). The main focus is on streaming applications. The optimization leads to a specification that

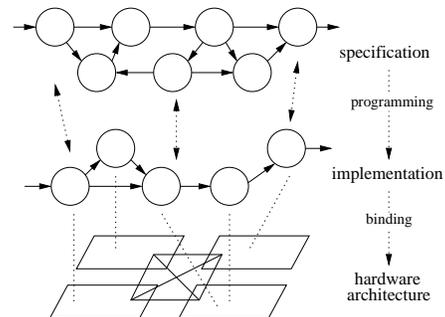


Figure 1: Mapping trajectory.

forms a good starting point for mapping the application onto a multi-processor system. The concurrency model is used in the programming step of the mapping trajectory. The optimization criteria used in our concurrency model are inspired by those used in performance analysis (see, e.g., [12, 11, 7]) but targeted towards streaming and concurrency. The novelty of our method is that we consider all aspects of concurrency including both computation and communication in an integrated framework, at a relatively high level of abstraction, allowing architecture-independent concurrency optimization. It allows the easy identification of concurrency bottlenecks. The result of the optimization process is a specification that can easily be further optimized for many different platforms. In other words, our techniques help in making re-usable specifications. The final architecture-dependent step is not covered in this paper.

Recently, a trend has emerged to use high-level models, written in languages such as C and C++, to introduce concurrency as well as communication and synchronization into originally sequential languages. SystemC [1], YAPI [5], SpecC [2] and VCC [9] are all good examples of relatively high-level models that are suitable for system-level design. They mostly focus on providing several communication models which can be refined down to the bus transaction level. Such refinement can be done, in an object-

oriented flavor, by providing different implementations of the same interface, where the fixed interface is all that the communicating processes care about.

CAST builds on the above trend toward executable communicating process models. Many theoretical models developed, mainly in computer science, are too abstract for proper concurrency optimization of resource-constrained embedded systems. Many source-code optimizations techniques are too low level to scale to future multi-processor systems. Hence, CAST focuses on source-code optimizations at the task-level. The extracted concurrency can be exploited in the mapping of the specification onto a multi-processor system. Our approach is related to the task-concurrency-management step in the system-synthesis methodology described in [14]. The task-concurrency-management step consists of task-level concurrency extraction, task scheduling and inter-task refinement. The extraction of task-level concurrency is limited to the number of independent tasks that can be found by a designer, based on the deterministic behavior of these tasks. The designer does not get support from the methodology in identifying these tasks. The approach used in CAST complements the task-concurrency-management approach since it provides support for extracting task-level concurrency.

We presented in [13] a model of computation for specifying concurrency and an accompanying concurrency model. These models are briefly introduced in the next section. The implementation of these models is presented in Section 3. The implementation of these models forms the basis of CAST, as the output of the concurrency analysis is a set of values for all concurrency properties of the specified application. To present this information in an intuitive way to a designer, we need an intelligent user-interface. A first prototype of such an interface is presented in Section 4. It is discussed how we are able to guide a designer in an intuitive way through the design space. This requires techniques to handle the information explosion that occurs when large applications (e.g. MPEG) are analyzed.

2 Model of Computation

The model of computation and the accompanying concurrency model, which are the basis of CAST, are introduced in this section. For detailed information on these models, we refer to [13].

2.1 Computational Networks

The computational-network model assumes that a parallel computation is organized as a hierarchical collection of autonomous compute nodes that are connected to each other by means of point-to-point connections. The set of compute nodes and their con-

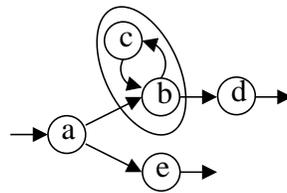


Figure 2: A computational network.

nections form the computational network. The connections are the only way of communication between the nodes. A given node computes on data it receives along its input connections to produce output on some or all of its output connections. Thus, each compute node performs a sequence of actions (e.g., C/C++ statements in an executable specification) which are modeled as a totally ordered sequence of events. Events that read from input connections and write to output connections are respectively called *read events* and *write events*. All other events are called *internal events*.

Hierarchy might be useful in a larger computational network, as it allows abstraction from the primitive operations taking place in the system. We introduce hierarchy in the computational network by allowing a node in a computational network to be a network of compute nodes. The outside world sees this node as an indivisible compute node with a set of input ports, a set of output ports and a defined behavior, while the node is in fact a set of compute nodes. Figure 2 shows a computational network that consists of the compute nodes *a*, *d*, *e* and a network that contains the compute nodes *b* and *c*.

The computational-network model is used to model applications for image, video and graphics processing (e.g., an MPEG decoder or a still-texture decoder). It captures the core of parallel (streaming) applications. It specifies only those aspects that are necessary for concurrency analysis. In this way, it allows for many instantiations. The model is, for example, sufficiently abstract to comprise a number of data-flow models like Kahn Process Networks [3, 4] and Synchronous Dataflow [10]. Our concurrency analysis can be applied to executable specifications in all these models.

Lamport [8] has shown that the events resulting from a concurrent computation form a partial order, referred to as the *causality relation* or *happened before relation*. Lamport also introduces *logical clocks* that can be used to create an ordering that is consistent with causality for all events that occur during a computation. We use an adapted version of Lamport's clocks to obtain an abstract notion of time for use in our concurrency model. To reason accurately about timing aspects without referring to concrete imple-

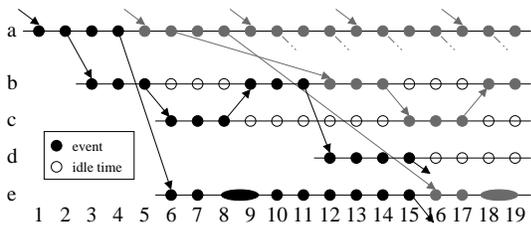


Figure 3: A partial event diagram.

mentations, we associate a delay, using a *delay function*, with the events that take place in the compute nodes and with the communication over the connections. Of course, these delays must be in some way reasonable for actual system implementations.

An execution can be displayed graphically in a *partial event diagram*, such as the one shown in Fig. 3. This figure shows an example ordering of all events that take place during an execution in the computational network of Fig. 2.

Node *a* starts with reading input from the environment. At the end, *d* and *e* produce values for the environment. The gray nodes represent repetition of the events caused by a single input to the network. As we are targeting streaming applications, it is common to have this repetition. In practice, one should possibly use more than one input to derive an event diagram that can be reliably used by CAST. It can be seen from Fig. 3 that the connection between the nodes *a* and *e* has a delay of one logical clock value, all other connections have a delay of zero logical clock values. Node *e* executes one event that takes two logical clock values; all other events require one logical clock value. Note that the diagram is kept simple for illustrative purposes.

2.2 Concurrency Model

Our concurrency model aims at performing a target-architecture-independent concurrency optimization. Its concurrency measures abstract from the environment in which a computational network operates, and are calculated from the computational-network structure and an event diagram of an execution. The event diagram is used to compute timing information such as run-times of nodes, idle times, etc. The concurrency measures are used in conjunction with a design method that consists of four steps. Each step tries to optimize one different aspect of task-level concurrency; it optimizes one of the concurrency measures, while the other measures are used to balance the overall concurrency optimization. The measures are computed for the network and for the individual nodes in the network. The measures for the compute nodes provide insight into the concurrency bottlenecks. The measures for the network can provide global guidance when optimizing concurrency.

The idea is to first extract as much parallelism as possible, then choose the granularity of communication, and finally put together tasks to arrive at a network with a balanced workload over compute nodes. We omit the formal definitions of all measures because that would require the introduction of a lot of formalism and focus on their intuitive explanation. The reader interested in more detail is referred to [13].

Task-splitting. The compute node with the longest run-time is determining the rate at which new computations can be started in the network. In other words, this node determines the throughput of the network. The throughput is an important property when a system designer is designing a streaming application. The *restart* measure provides an abstract notion of it. To optimize the restart, the slowest compute node must be split in a set of compute nodes with better values for the restart. The slowest node in Fig. 3 is node *e* which has a run-time of 10, determined by the events from logical time 6 to logical time 15.

Good values for the restart can be obtained through very fine-grained compute nodes. However, this gives communication overhead (and possibly scheduling overhead). The restart measure should therefore be balanced with other measures.

Data-splitting. The structure of a network reveals the chains of compute nodes that belong to different parts of the computation taking place in the network. In other words, it reveals the different data-streams that are processed in the network. The simple network of Fig. 2 processes two data streams, namely one through *a*, *b*, *c*, *d*, and one through *a* and *e*. The more data-streams that can be distinguished in a network, the more data-parallelism is present. However, if many different data-streams go through one node, then this node may be a synchronization bottleneck for those data-streams. Assuming that input node *a* in Fig. 2 cannot be split, there is not really an obvious bottleneck. The *structure* measure is used to quantify this concurrency property. The (task-level) data-parallelism that is present in the specification should be made explicit to optimize this concurrency property.

Communication granularity. In a parallel execution, we want to minimize the overhead of communicating data between nodes. The nodes should spend as much time as possible on computation and not on communication, as computation, i.e., data transformation, is the main goal of every computational network. The ratio between time spent on computation and time spent on both computation and communication is expressed in the *computation load*. This ratio can be calculated for the network as a whole and for individual nodes and should be as high as possible. In Fig. 3, node *a* has a computation load of $1/4$ (only 1 in four events is an internal computation event), where

as e.g. node d has a computation load of $4/5$; this suggests node a as a bottleneck. The granularity of communication must balance time spent on communication and time that nodes have to wait for input data.

Merging. During an execution, a compute node is either busy, performing events, or it is idle. It can be idle because it is waiting for data or because it has finished its execution while other nodes have not yet finished. To get a balanced workload over nodes, we must balance the execution times (computation plus communication time) and run-times (execution time plus idle time) of the different nodes. This is important to optimize streaming behavior. To get a notion of the workload balance, the *execution load* considers the ratio between the execution time and the run-time. Nodes that have a low execution load must be merged with each other to get a better overall execution load for the network. In our running example, node b has the lowest execution load because of the idle time in the middle of its computation.

A parallel computation will in most cases be faster than a sequential implementation of that computation. This is often referred to as speed-up. The realized speed-up for a computational network depends on the synchronization that is required between the different nodes in the network, the introduced communication overhead, and the balance of the computation over the different nodes. The second and third aspect are covered by the computation load and execution load respectively. The influence of synchronization is not yet fully captured in these measures, although a poor synchronization does affect the execution load. Synchronization is important when considering concurrency, because synchronization is limiting the execution of compute nodes and with that the number of nodes that can run in parallel. Synchronization constraints may impose the restriction that two nodes can only execute in sequence. A typical example of such a problem can be seen in Fig. 3 for nodes b , c and d . This concurrency property is captured in the final measure of our model, the *synchronization* measure. The design method does not contain a special step in which this property is optimized. It must be taken into account in all steps.

This section introduced the five concurrency measures that form our concurrency model. In [13], it is shown in more detail how these measures are defined and that all five are meaningful and do not (fully) overlap, i.e., all five are needed to allow a balanced optimization of a computational network.

3 Implementation

This section presents the implementation of the computational-network model and the concurrency

model, as introduced in the previous section. An important part of the computational-network model is the delay function that determines the format (timing) of the event diagram. Its implementation is discussed in detail in the next sub-section. The implementation of the concurrency model, which forms the basis of CAST, is presented in Section 3.2.

3.1 Computational Networks

The computational-network model is presented in Section 2. It is our belief that the model is suitable for modeling embedded streaming application, which are often signal-processing applications. These are nowadays often modeled using Kahn Process Networks or Dataflow graphs, which can be seen as realizations of our computational-network model. For a first experimental evaluation of the concurrency model, we have therefore chosen to use an existing Kahn Process Network (KPN) implementation, namely YAPI [5]. The processes in the YAPI/KPN model become our compute nodes, the process networks become our computational networks, and the channels of a YAPI/KPN are the connections in our computational-network model.

The previous section introduced the time-stamping mechanism based on Lamport's logical clocks that is a part of the model of computation. This time-stamping mechanism must allow reasoning about causality and some timing aspects on a relatively high level of abstraction without referring to implementations/physical time. To implement this time-stamping mechanism, we need an implementation for the delay functions of the events and connections. The delay function for events is implemented as two different functions. One function associates a delay with each internal event that occurs in the compute nodes and one function associates a delay with each read/write event.

A single internal event in a compute node is defined to be equal to the execution of a single C++ statement. The delay function for internal events must map each C++ statement on a delay value associated with that statement. The value of the delay must represent the amount of work associated with performing the event. Therefore, we relate internal events to the number of instructions needed to execute these events on a processor using a standard compiler. We assume that the influence of a specific instruction set does not have too much influence on the results. Our first experiments show that the proposed notion of time is both accurate and abstract enough to perform target-architecture-independent optimization.

When the design-space exploration is started, little information about the communication medium is known. As a first strategy to get some abstract notion of the amount of work related to read and write

events, we could use the same approach as with the internal events. However, the read and write functions used in our software library as communication primitives will in general not be used in the actual implementation. They will be mapped onto lower-level more efficient communication primitives. Using the number of instructions needed to execute a read or write function call is thus not a very good measure for the delay associated with the read and write events. The functions used in the implementation will have a different delay. Similar problems arise when a delay must be associated with a connection. Although little information about the implementation is known in the early stages of design, we observe two important properties for the communication. First, in general, the more data-elements that are communicated, the more time that a communication will take. The number of data-elements communicated should thus be a part of the delay associated with the read or write event and the connection. Secondly, each function that interacts with the communication medium (e.g., read, write event and implementation of the connection) will have to get access to the communication medium. The time needed for this does in general not depend on the number of data-elements communicated. Based on these two observations, we propose to implement the delay functions for read/write events and for connections using the linear functions $ax + b$ and $cy + d$ of respectively the number of data-elements x communicated and the size of the data-elements going through a connection y . The constants b and d approximate respectively the time needed to call the communication primitives and the access time to the communication medium. The constants a and c approximate respectively the time needed to read/write one data-element and the time needed to transport one data-element.

3.2 CAST

CAST is a software tool for computing the concurrency measures of a computational network. The overview of CAST is shown in Figure 4. The current version operates on YAPI/KPNs. The core of CAST consists of three steps (i.e., parser, simulator and analyzer). These steps are described below in some detail.

The actual computation of the concurrency measures for a given network is performed in the analyzer. The analyzer uses for that a trace of all events that have occurred during a simulation of the computational network with a given input. The analyzer needs also a description of the network structure, the network graph. Using these two and the settings for the delay functions, it maps the events onto the appropriate delay and then orders these events according to the causality relations; it creates a (partial) event dia-

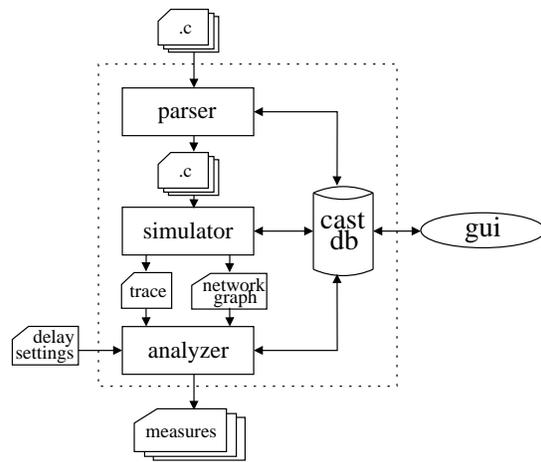


Figure 4: Overview of CAST.

gram. After that, it has enough information to compute the values of the different concurrency measures. The information needed to compute the event diagram is obtained through simulation. The simulator takes C++ source code as an input. It simulates the network described in these files with a given input and traces all events that occur in the network. To be able to do this, the original C++ files containing the network are modified in the parser. The parser adds functions to the computational network that log the execution of individual events to a file when the network is simulated. The parser adds also a CAST run-time environment, which is used to govern the event tracing and extraction of the network structure during the simulation. The network structure is thus extracted at run-time and not through a static analysis. It would of course be possible to do this using static analysis of the code, but it is more convenient to do it at run-time.

The three steps of CAST, the parser, simulator and analyzer communicate with a database in which they store relevant information. The analyzer stores for instance the values of the concurrency measures in this database. The data in the database is available for use by the graphical user interface.

For typical applications (e.g., JPEG or MPEG encoders or decoders), it is practically impossible to analyze the network for all possible inputs. Therefore, CAST contains a statistical analysis module. With this module, not shown explicitly in Fig. 4, it is possible to select analysis results of different simulations and compute the average, variance, minimum and maximum for each individual compute node and computational network. In this way, we can minimize effects from a specific simulation input, but also effects of a specific instruction set or setting for the delay functions.

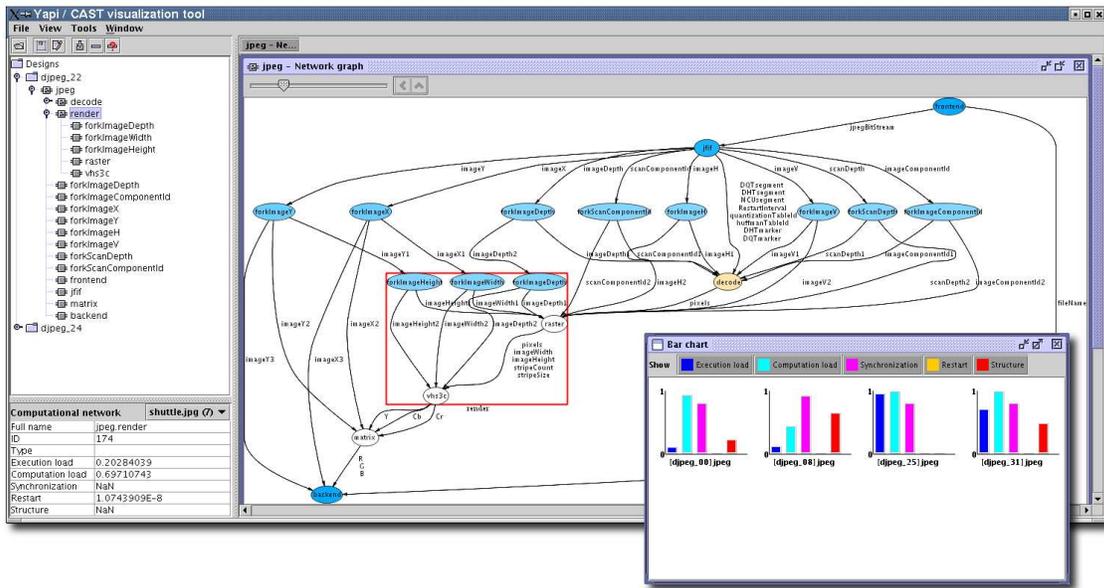


Figure 5: Screen-shot of CAST.

4 Visualization

CAST contains a visualization software package that helps the designer in understanding the concurrent behavior of the network and identifying potential concurrency bottlenecks by providing a direct method of feedback on the network analysis results. Figure 5 shows a screen-shot of CAST (operating on a JPEG decoder). The hierarchy that is present in computational networks is used to present the designer with a visual graph representation of manageable magnitude, along with a common tree view. Starting from the top level of the hierarchy, lower level details can be expanded. This approach makes it possible to work with networks consisting of hundreds of nodes, while still being able to reach every node without losing the overview.

We discuss in the remainder of this section how a designer can perform a design-space exploration using CAST. An important step in this exploration is the identification of bottlenecks. In Sec. 4.2 is discussed how CAST can be used for this. To identify the cause of a bottleneck, it might be necessary to look at the event diagram of the bottleneck node. The support that CAST offers in this situation is discussed in Sec. 4.3.

4.1 Design-Space Exploration

The goal of our task-level concurrency optimization is to transform a computational network in a structured way into a computational network that has a balanced workload and good communication behavior. A generally applicable design exploration method consisting of four steps is used in conjunction with the

concurrency model to realize this, as explained in Sec. 2.2.

The visualization software supports this design exploration method through its visualization of the concurrency measures and the offered support for finding concurrency bottlenecks. When a potential bottleneck is found, the designer will modify the computational network to resolve this bottleneck. The results of these actions can be evaluated using CAST by simulating the new computational network and comparing the concurrency measures for the two designs. The impact of the design changes can be visualized using bar charts. For convenience, all measures are normalized in the interval $[0, 1]$, with low values corresponding to poor concurrency and high values to good concurrency. This helps the designer to very quickly get a good overview of the quality of a design, and it supports an easy comparison of different designs. This visual feedback combined with the design methodology guides the designer in an intuitive way through the design space.

4.2 Identifying Bottlenecks

The bottleneck nodes in a computational network are those nodes that have low values for the concurrency measures. To help the designer in finding these nodes, CAST can map these measures onto the node size and colors of the graph representation of the network. This makes it very easy for a designer to identify potential concurrency bottlenecks. It requires a simple two-step approach. First, the designer must open a graph representation of the computational network and identify the node that needs attention because of its low concurrency figures. If the identi-

a JPEG decoder manually optimized for this platform [6]. The results illustrate that the approach followed by CAST works. More details can be found in [13]. Currently, the concurrency model helps the designer in finding bottlenecks in the system. It does not give support for situations in which a compute node must be subdivided into multiple nodes. An extension of the concurrency model is planned that provides this support. We further plan to add semi-automatic pattern recognition of patterns in event traces. This must help the designer in finding causes for bottleneck nodes. Further, the computational-network model and concurrency model will be extended to support reactive behavior and control applications. We also plan to extend the concurrency model to take architecture information into account for the architecture-dependent step of the design process.

References

- [1] Grotker, T., et al., *System design with SystemC*. Dordrecht, The Netherlands: Kluwer Academic Publishers, 2002.
- [2] Gajski, D.D., et al., *SpecC: specification language and methodology*. Dordrecht, The Netherlands: Kluwer Academic Publishers, 2000.
- [3] Kahn, G., *The semantics of a simple language for parallel programming*. In: Information Processing 74, Proc., Stockholm, Sweden, August 1974. Ed. by J.L. Rosenfeld. Amsterdam, The Netherlands: North-Holland, 1974. p. 471-475.
- [4] Kahn, G., and D.B. MacQueen, *Coroutines and networks of parallel processes*. In: Information Processing 77, Proc., Toronto, Canada, August 1977. Ed. by B. Gilchrist. Amsterdam, The Netherlands: North-Holland, 1977. p. 993-998.
- [5] Kock, E.A. de, et al., *YAPI: application modeling for signal processing systems*. In: Design Automation Conference, Proc. 37th Int. Symp., Los Angeles, USA, June 2000. IEEE, 2000. p. 402-405.
- [6] Kock, E.A. de, *Multiprocessor mapping of process networks: a JPEG decoding case study*. In: System Synthesis, Proc. 15th Int. Symp., Kyoto, Japan, October 2002. IEEE, 2002. p. 68-73.
- [7] Kung, S.Y., *VLSI array processors*. London, UK: Prentice Hall. 1998.
- [8] Lamport, L., *Time, clocks, and the ordering of events in a distributed system*. In: Communications of the ACM, 21(7), 1978. p. 558-565.
- [9] LaReu, W., et al., *Functional and performance modeling of concurrency in VCC*. In: Concurrency in Hardware Design, LNCS 2549. Berlin, Germany: Springer-Verlag, 2002. p. 191-227.
- [10] Lee, E.A., and T.M. Parks, *Dataflow process networks*. In: Proc. of the IEEE, 83(5), May 1995. IEEE, 1995. p.773-801.
- [11] Mazzeo, A., and N. Mazzocca, U. Villano, *Efficiency measurements in heterogeneous distributed computing systems: From theory to practice*. In: Concurrency: Practice and experience, 10(4), May 1998. p. 285-313.
- [12] Raynal, M., and M. Mizuno, M. Neilsen, *Synchronization and concurrency measures for distributed computations*. In: Distributed Computing Systems. Proc. 12th Int. Conf., Yokohama, Japan, June 1992. IEEE, 1992. p. 700-707.
- [13] Stuijk, S., *Concurrency in computational networks*. Master's thesis, TU Eindhoven, Eindhoven, The Netherlands, 2002.
- [14] Thoen, F., and F. Catthoor, *Modeling, verification and exploration of task-level concurrency in real-time embedded systems*. Dordrecht, The Netherlands: Kluwer Academic Publishers, 2000.