

A Monitoring-Aware Network-on-Chip Design Flow

Calin Ciordas^a Andreas Hansson^a Kees Goossens^b
Twan Basten^a

^a*Eindhoven University of Technology,
{c.ciordas,m.a.hansson,a.a.basten}@tue.nl*

^b*Philips Research Laboratories Eindhoven,
kees.goossens@philips.com*

Abstract

Networks-on-chip (NoC) are a scalable interconnect solution for systems on chip and are rapidly becoming reality. Monitoring is a key enabler for debugging or performance analysis and quality-of-service techniques. The NoC design problem and the NoC monitoring problem cannot be treated in isolation. We propose a monitoring-aware NoC design flow able to take into account the monitoring requirements in general. We illustrate our flow with a debug driven monitoring case study of transaction monitoring. By treating the NoC design and monitoring problems in synergy, the area cost of monitoring can be limited to 3-20% in general. We also investigate run-time configuration options for the NoC monitoring system resulting in acceptable configuration times.

Key words: run-time monitoring, network on chip, transaction monitoring, design flow

1 Introduction

Advances in semiconductor technology have enabled very complex large scale systems on a chip (SoCs) designs. Each new SoC generation integrates more processing elements (IPs) and offers increased functionality. As the number of IPs increases, traditional interconnects, such as busses, become a bottleneck.

Networks-on-chip (NoCs) are a modular, scalable interconnect solution [1, 9, 14, 17, 19]. Concrete examples of NoCs are *Æthereal* [13], *Xpipes* [8], *QNoC* [3] and *Mango* [2]. Currently, they tend to become the preferred interconnect solution for large scale inherently multiprocessor SoCs. However, NoCs require

sophisticated tools to aid in design-time decisions [3, 12, 15, 23]. Furthermore, with increasing complexity there is also a strong need for run-time NoC monitoring [4, 6, 24, 25], which must be accounted for in the design phase. This is in turn driven by debugging [4, 6] and performance monitoring/Quality of Service (QoS) [24, 25, 27].

With the introduction of NoCs, the on-chip communication becomes more sophisticated relying on run-time programmable solutions. In centralized bus-based systems a single bus monitor is enough to be able to track the whole history of the system. In NoC-based SoCs, due to the inherent parallel behavior of communications, where multiple pipelined parallel communications may exist between IPs, multiple monitors have to be employed. The problem of how many such monitors are needed, their automatic placement in the NoC-based SoC by means of a monitoring-aware NoC design flow and the associated area cost implications have not been previously investigated.

Monitors and the traffic they generate are traditionally added non-intrusively into the SoC by using a separate monitoring NoC [24]. The cost of such a solution is high however, and a more efficient solution is to use the same NoC for both monitor data and user data, as suggested in [4, 6, 25]. When monitoring traffic uses an interconnect of its own, it can be dimensioned after the user data NoC is designed. This merely adds an extra step in the design flow. However, when monitor and user data must share the same NoC, the overall design flow must be revised [6].

NoC design flows for ASIC type designs are normally split in several steps as topology selection, mapping, path selection and slot allocation [3, 12, 15, 23]. Some design flows may omit or combine various steps. Each step adheres to the decisions taken in the previous steps. As prerequisites for NoC design, communication requirements must be derived, and the set of IPs to be connected to the NoC must be specified. In the topology selection step, the router network together with the bordering NIs are generated, based on the previously derived communication requirements. Using this topology together with the IP specification, the binding of IP ports to NI ports is done in the mapping step. In the path selection step, paths are allocated for all the communication flows specified, and in the slot allocation step each of the flows gets its own TDMA time slots for the traversed NoC links.

We have two interdependent problems: the one of functional dimensioning of the NoC and mapping of cores while accounting for their communication requirements, and the other of monitor placement and monitoring bandwidth specification. If these two problems are solved sequentially, the monitoring communication requirements can be precomputed. However, if the communication requirements of the monitors do not fit directly on the generated application NoC, a new NoC must be generated, e.g., by increasing the topology

and repeating the process. However, by increasing the topology, the number of NoC routers increases. In turn, the mapping, path selection and allocation of resources may change and the number of required monitoring probes may increase as well (e.g. if probing all routers is required) and their communication requirements may change. In the mentioned cases the monitoring problem (whether driven by debugging or by run-time performance analysis) must be solved within or at least tightly coupled with the NoC design process. The task of placing the monitors must therefore be automated and integrated in the NoC design flow.

2 Contribution

We propose a monitoring-aware NoC design flow able to take into account the monitoring requirements at all steps in the NoC design flow. We illustrate this with a debug driven monitoring case study. Simple, area-efficient transaction monitors, attached to selectively chosen NoC routers, are used to enable debugging of the NoC-based SoC at transaction level. This is one of the most difficult cases, where the monitoring requirements are only known after the path selection step. In the context of application specific designs, the proposed flow is able to automatically insert transaction monitors, by determining the number and placement of these transaction monitors and accounting for their communication requirements. The smallest area NoC which satisfies the application requirements, as well as the monitoring requirements is generated as a result. The area implications are quantified and compared to original NoCs without monitoring. The efficiency of the flow is shown on several realistic examples. Several run-time configuration options for the monitoring service are also detailed and experimentally investigated. This paper is an extended version of [7].

3 Related Work

In [24], the use of end-to-end monitors collecting network interface statistics is proposed in order to assist the operating system controlling the NoC. The work focuses on the use of such performance monitors to optimize communication resource usage. The monitored data uses a separate NoC, called the control NoC instead of the application NoC.

Router-attached performance monitors are used in [25] to keep track of the network utilization. By means of a network manager this information is made useful to a QoS manager to increase/decrease the quality levels of running

applications. The monitored performance data uses the same NoC as the user data.

[28] uses link utilization monitors as support for a congestion-controlled best effort communication service, by means of run-time centralized model predictive control. It uses the existing NoC to transport the monitored link utilization at precomputed periods.

Embedded monitors in an FPGA environment are used to track end to end run-time behavior (queue utilization) as feedback for the design exploration phase [21]. The employed hardware monitors have dedicated wires for transport of their results multiplexed in front on an output port, showing an approach that is not scalable.

[4] proposes a generic NoC monitoring service comprising monitors attached to NoC components, routers or NIs, offered by the NoC. Targeted at debugging, it focuses on generic concepts of the service, architectural and general cost implications. The monitored data uses the same network as the user data.

[6] shows that using the same interconnect for the user traffic and monitoring traffic is area-efficient but may require modifications in the NoC design flow. However, it falls short on showing how to solve this problem in general and what are the associated cost implications.

All previous work assumes that: (1) the placement of the monitors is known, (2) the monitoring generated traffic or communication requirements are known in advance, (3) this traffic fits on top of the user traffic on the shared NoC or (4) on a separate NoC.

For monitoring, in general, these assumptions are not valid. The number and placement of monitors and their associated monitoring communication requirements are usually not known beforehand, but only after the NoC to be probed has been fully designed, or at least some steps in the NoC design flow have been performed. For example, some requirements may be known only after topology generation, such as the number of routers employed in the NoC, which is relevant if all routers or a coverage of routers need to be probed e.g. with router monitors showing link utilization. In this case the number of routers determines the number of probes and their placement, while their communication requirements are fixed, depending only on the number of links being traced. Other communication requirements may be known after the path selection step in the design flow, e.g. router monitors able to trace a connection, e.g. the functional traffic for debug reasons (or for connection utilization). In this case, assuming a desired full coverage of the connections, the number of probes and their placement is given by the routers in the cover. Their communication requirements depend on the number of connections passing the probed router and their sizes. We propose a monitoring aware design flow

that fully integrates the design of the NoC and its monitoring service, solving all the above mentioned issues.

4 Architectural Platform

4.1 NoCs and *Æthereal*

NoCs comprise two components: routers (R) and network interfaces (NI), as depicted in Figure 1. The routers can be randomly connected among themselves and to the NIs (i.e., there are no topology constraints). Note that in principle there can be multiple links between routers. The routers transport packets of data from one NI to another.

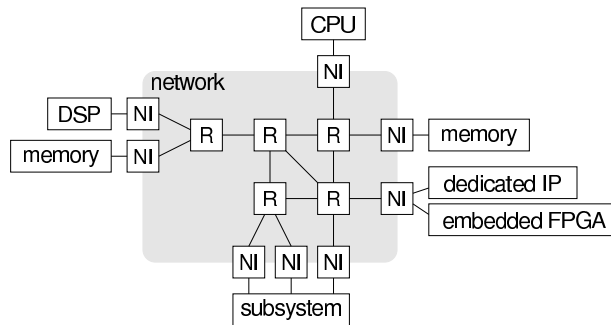


Fig. 1. Example NoC

The NIs enable end-to-end services [26] to the IP modules and are key in decoupling computation from communication [2, 29]. The NI allows the designer to simplify communication issues to local point-to-point transactions at IP module boundaries, using protocols natural to the IP [29]. They are responsible for (de-)packetization, for implementing the connections and services, and for offering a standard interface (e.g., AXI or OCP) to the IP modules.

We use the *Æthereal* NoC [12, 13] as an example for our work. The *Æthereal* NoC runs at 500 MHz and offers a raw link bandwidth of 2GB/s in a 0.13 μ m CMOS technology. It is supported by state of the art design time decision tools [12, 15]. *Æthereal* offers transport-layer communication services to IPs, in the form of connections, comprising best-effort (BE) and guaranteed-throughput (GT) services. Guarantees are obtained by means of TDMA slot reservations in NIs. *Æthereal* NoC instances are reconfigurable at run-time. This is achieved by programming the NIs using standard memory-mapped I/O ports. The current setup uses centralized programming of the NoC and source routing. The *Æthereal* NoC allows the mapping of potentially multiple IPs per NI and potentially multiple NIs per router with any topology.

The interconnected IPs interact with each other by means of transactions, which are read and write transactions from IPs. Transactions consist of one request message and one optional response message. E.g. a request message can be a write message. A response message is for example data coming back as a result of a read operation, or an acknowledgment as a result of a write operation. Transactions are performed on connections, consisting of one request and one response channel. The paths of request and response channels may be different.

The NIs convert these messages into packets, by chopping them into pieces of a maximum length and adding a header to each of these pieces, resulting in packets. Packets may be of different lengths. Packets are further split into flits, the minimum flow-control unit between hops. One flit corresponds to one TDMA slot.

4.2 Transaction Monitoring

4.2.1 The Transaction Monitoring Problem

To increase the operational speed of system-level debugging, the NoC debugging infrastructure must bring the abstraction level of the monitored data at transaction-level, and allow run-time transaction monitoring in particular, at a reasonable cost.

The problem of how many transaction monitors are needed relates to the desired coverage of the user communication flows. In general a full coverage is desired. However, it is prohibitively expensive to duplicate all traffic in the NoC; therefore the coverage may be full but has to be selective at certain moments in time. This means that the monitors must cover all channels, but not at the same time. At run-time, any (potentially more) of the desired channels can be selected to be monitored. The number of simultaneously active monitors in the system is bounded by the number of monitors deployed, as each monitor can only track a single channel.

The problem of the cost implications of the monitoring relates to the area of the monitors, the number of monitors involved and also to the area of the resulting NoC which supports both the application and monitoring communication requirements. The resulting NoC, potentially larger than the original NoC, accounts for the extra NIs, NI ports or enlarged topology to support monitoring in addition to the application communication.

4.2.2 NoC Monitoring Service

We use a monitoring service (NoCMS) as described in [4]. The NoCMS is offered by the NoC in addition to the communication services offered to the IPs. It consists of configurable probes attached to NoC components, see Figure 2 for details. The probe modular design comprises three parts: the sniffer (S), the event generator (EG) and the monitoring network interface (MNI). The MNI can be a separate NI or it can be merged with an existing NI. The monitoring service access point (MSA) is an IP which controls the configuration of the monitors at run-time and receives the monitored data from all monitors. E.g. the MSA can stream this data outside the chip through a debug port. The NoCMS is configured by means of probe programming via the NoC using memory-mapped I/O write transactions.

The generic NOCMS concepts must be instantiated for the monitoring task at hand, in our case transaction monitoring. This implies the replacement of the EGs with transaction monitors, the placement of transaction monitors to offer a full channel coverage of the system, the placement of the MSA, and the dimensioning of the communication requirements of the monitors (as this data should go to the MSA via the NoC). We use centralized monitoring with a single MSA.

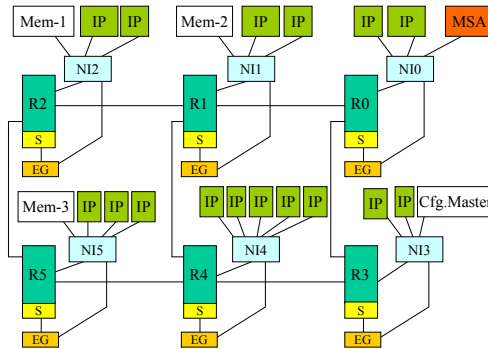


Fig. 2. NoC Monitoring Service

4.2.3 Transaction Monitors

The transaction monitors can be attached to routers or NIs. For simplicity we only consider them as attached to routers. They can ultimately track transactions over a single channel passing any of the router's links. The transaction monitors consist of a configuration block and a set of five pipelined filters, as illustrated in Figure 3.

The monitors can be (re-)programmed at run-time to track any channel. All run-time settings are done through the configuration block. The configuration path is marked CFG in the figure. As previously mentioned this is done by means of simple write transactions. The configuration time required for

configuring the NoC monitoring service together with possible configuration options and policies are discussed in Section 7.2.6.

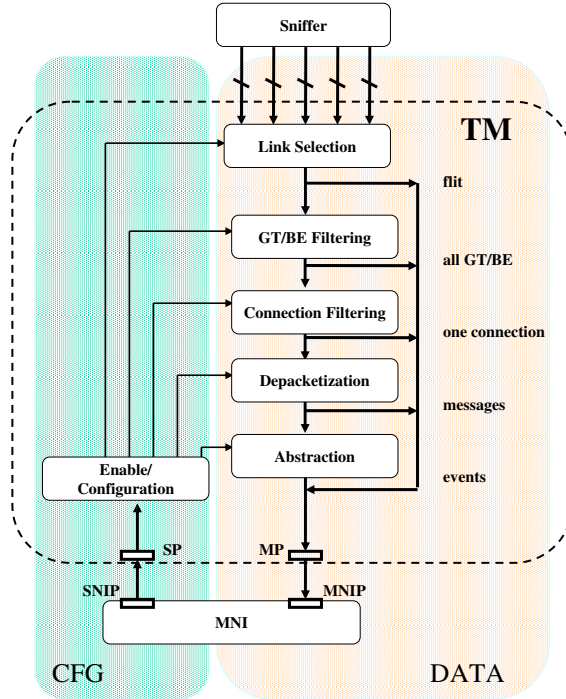


Fig. 3. Transaction Monitor

The monitoring data path starts at the sniffer and is marked DATA in Figure 3. The raw data is provided to the transaction monitor by the sniffer, which captures it from the router links. The link of interest can be selected at runtime by configuring the first filter. The flits can be further filtered as BE or GT in the second filtering block. Further filtering of flits is done by identifying a single connection from the set of connections sharing the same link, in the next filter.

Transactions are composed of messages. Message identification allows to see, from within the NoC, when a write or a read message has been issued and from where or to which of the IPs or memories. Messages are payload packed in packets. Therefore, message identification requires depacketization, a procedure usually done at the NI. For the fourth filter, which is the essential one to provide transaction monitoring, we reused available \AA threal hardware modules for depacketization. The fifth filter has abstraction capabilities and is not discussed here because the details are not important. The interface between the transaction monitor and the MNI consist of a slave port (SP) for configuration and a master port (MP) for sending the monitoring data. Their corresponding MNI ports are the SNIP and the MNIP.

Up to 64 bits of configuration data are required by a single transaction mon-

itor during the configuration. This can be done using either one 64-bit DTL-MMBD write operation or two 32-bit DTL-MMIO write operations. For the DTL details see [10]. In the first case the configuration data can be packed into a 32-bit (one word) command (C), 32-bit (one word) address (A), 64-bit (two words) payload (P) write message which would enter the MSA connected NI. In the second case the same data is packed into two (C,A,P) write messages. Note that in general the amount of configuration data required per monitor may differ with the monitor type, e.g. it may not be the same for a transaction monitor and a performance monitor, potentially resulting in more/less required write messages per monitor. However, the same run-time configuration policies and techniques may apply, see Section 7.2.6.

A $0.13\mu\text{m}$ CMOS technology implementation of a transaction monitor supporting the first four filtering stages shows an area cost of 0.026mm^2 . Assuming that no filtering/abstraction is done locally at the monitor, the bandwidth requirements of the transaction monitors are comparable with the bandwidth of the monitored connection. Further details on the (single) transaction monitor architecture and implementation with the associated problems (and their solutions) can be found in [5].

5 Application-aware placement

Since we are considering ASIC-like design, the application is known at design time. For the NoC-based SoC it means that also the set of connections (all request and response channels) is known at design time. The bandwidth and latency constraints of the channels are determined beforehand by means of static analysis or simulation.

At least one probe is required on the path of each channel, regardless whether it is a request or response channel. This means that any of the existing channels can be probed, achieving a full channel coverage. Furthermore, the concurrent observation of multiple channels is only limited by the number of probes in the NoC. We can simultaneously monitor one channel per probe. At run-time, the monitored channels per probe may change by means of programming the probes. This selectivity is acceptable as usually not all streams are required to be monitored at once (duplicating all traffic, even at a high abstraction level is prohibitive).

In ASIC design, a full coverage of routers with monitors may potentially be avoided, see for example the four monitors in Figure 4(a) covering each one of the four channels, versus the two monitors in Figure 4(b) covering each the two channels passing through. This leads to a reduction of the total monitoring solution area cost. Note that even assuming a full coverage of NoC routers with

transaction monitors the communication requirements of these monitors are not known before the path selection step in the NoC design flow, as we do not know earlier what channels will pass through each of the monitored routers. Therefore, the problem of modifying the design flow to support monitoring constraints cannot be avoided.

6 Design Flow

6.1 Design Flow vs. NoC Monitoring

When adding monitoring to an existing NoC two main architectural constructions can apply. One is the use of separate interconnects for the monitoring data and one is sharing the existing interconnect for both the user and the monitoring data. A hybrid version may also exist as a compromise between the two extremes, in which a single NoC is shared but physical resources for user and monitoring are disjoint. Note that although any interconnect may be used, as a separate monitoring interconnect it is logical to use a NoC, the monitoring NoC, because it is scalable. Using point to point wires or busses will eventually lead to scalability problems.

Separate NoCs. In this case, a separate NoC is chosen for monitoring. While Figures 5(a) and 6(a) show the original NoC and design flow without accounting for monitoring, Figure 5(b) shows the resulting system after adding the transaction monitors, and Figure 6(b) the corresponding design flow. The monitoring NoC is used for transporting the monitoring data from transac-

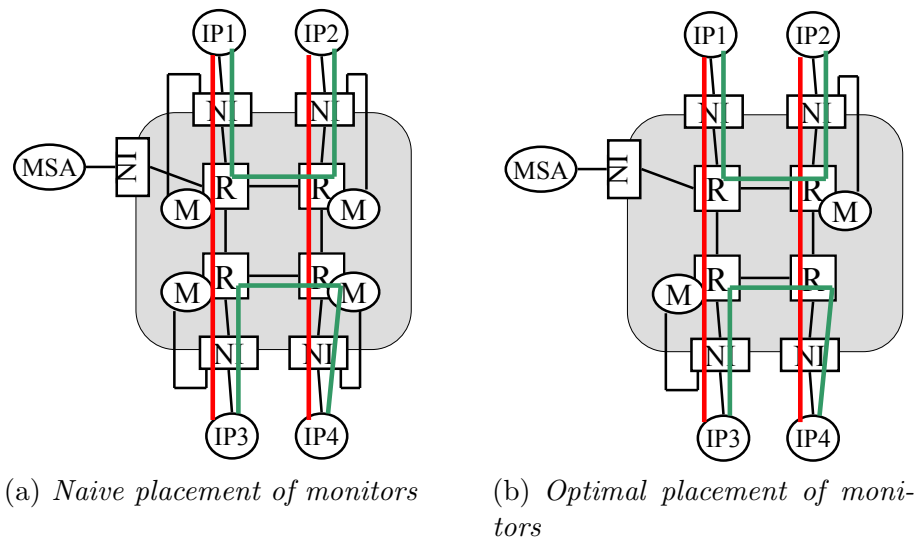


Fig. 4. Placement of Transaction Monitors

tion monitors to the MSA and for monitoring configuration traffic from MSA to the transaction monitors. In principle the monitoring NoC is similar in topology with the user NoC interconnect. For simplicity, we only show a fully probed NoC in Figure 5(b), with probes attached to all routers. A more advanced, selective NoC probe placement at routers is possible, e.g. ensuring a coverage of all NoC physical links or all NoC logical channels as shown in Section 5. For each of the probed routers, a new router and an NI are added. The NI is used by the probe to connect to the monitoring NoC. After the user NoC design process, the regular NoC design flow is applied also for the monitoring NoC, taking into account the monitoring communication requirements. Dimensioning of the monitoring communication requirements and of the number of debug IPs (e.g. router probes) required, which are dependent on the user NoC topology, mapping, and path selection, is simple as the user NoC design was done beforehand. While applying the NoC design flow for the monitoring NoC, the topology is already given by the original NoC, and mapping is given by the probe placement in the original NoC, as previously explained. Therefore only the path selection and slot allocation have to be done for the monitoring NoC. As the monitoring communication requirements are (should be) below the user communication requirements, in most cases these can be accommodated on the monitoring NoC; therefore the path selection and slot allocation would succeed. The major disadvantage of this solution is that it is expensive in area.

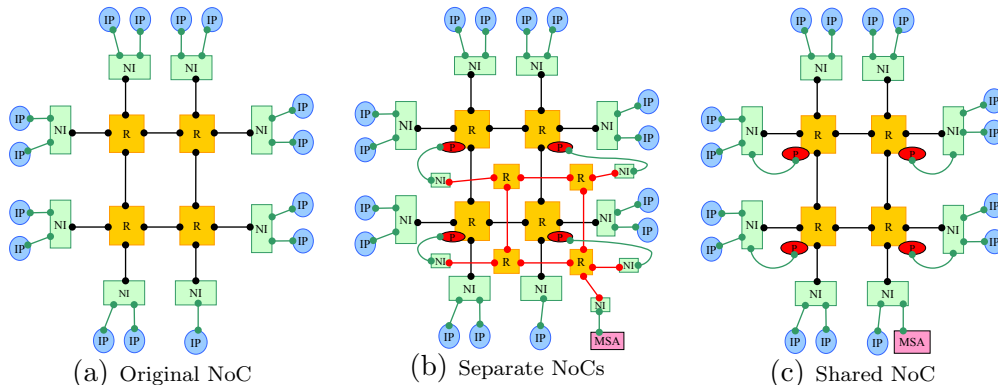


Fig. 5. Monitoring Transport Options

Shared NoC. In this case, we share all the NoC resources for user traffic or for monitoring traffic but we keep the NoC user traffic and the monitoring traffic separated, creating a virtual NoC for monitoring. After the user NoC is obtained, the monitoring communication requirements and debug IPs are computed and transaction monitors are added to the design. Figure 5(c) shows that transaction monitors are connected to the existing NoC by means of an extra port on the existing user NIs. Note that this can also be done by adding separate NIs for monitoring on the corresponding routers. The NoC is shared between the user and the monitoring traffic. The mapping of the probes to existing NIs is based on the closest available NI. We have the mapping of IPs

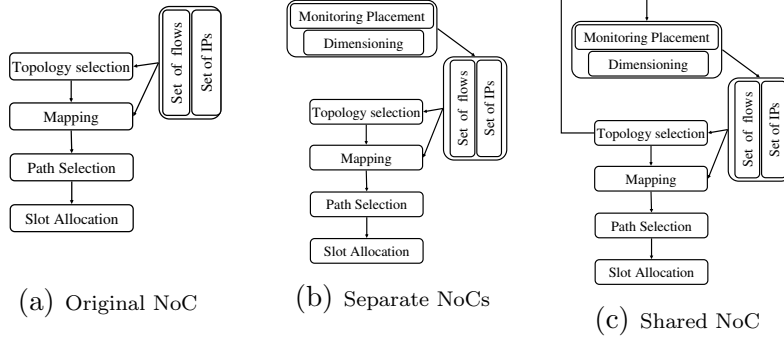


Fig. 6. Design Flows

to NIs and the NoC topology. Path selection and slot allocation is computed together for all the communication requirements: user and monitoring. Figure 6(c) shows the corresponding design flow. It is possible that everything fits on the existing user NoC. This means that the user NoC can accommodate the monitoring communication requirements on top of the existing user communication requirements. Topology of the NoC will therefore not change. This is exactly the situation shown in Figure 5(c). In this case, we have the lowest area cost, as no new NoC components, routers and NIs for monitoring, are added, except the new NI ports to connect the probes to the NoC. However, the combined communication requirements may not fit on the existing user NoC. In this case, it is clear that a new NoC must be generated, e.g. by increasing the topology and repeating the process. By increasing the topology, the number of user NoC routers increases and in turn the number of required transaction monitors may increase as well (e.g. if probing all routers is required). This leads to the recomputing of the monitoring communication requirements and monitoring IPs. *This means that the NoC monitoring flow must be revised, as illustrated in the remainder of this paper.* The reason for selecting and investigating this option is that the resulting NoC will in general be the cheapest one. For a more detailed analysis of the monitoring interconnect options, see [6].

6.2 UMARS

UMARS [15] is a QoS constrained NoC design algorithm. It unifies the three resource allocation phases: spatial mapping of cores, spatial routing of communication, and the restricted form of temporal mapping that assigns time-slots to these routes. UMARS considers the real-time communication requirements, and guarantees that application constraints on bandwidth and latency are met.

UMARS is a greedy algorithm, iterating over the monotonically decreasing set of unallocated channels until they are all accommodated in the NoC, or until allocation failed. The algorithm, as outlined in Algorithm 1, never back-

tracks to reevaluate an already allocated flow, enabling run-times in the order of milli-seconds.

Algorithm 1 Outer loop of UMARS

- (1) While there are unallocated channels
 - (a) Select the channel with highest bandwidth
 - (b) Find a mapping and a path
 - (c) Select slots on this path
-

An important property of UMARS that we exploit in this work is the fact that channels are allocated ordered on their bandwidth requirements. This is done as it: 1) helps in reducing bandwidth fragmentation [18], 2) is important from an energy consumption and resource conservation perspective since the benefits of a shorter path grow with communication demands [16], 3) gives precedence to flows with a more limited set of possible paths [16]. This ordering assures us that no channel succeeding the one currently being allocated has higher bandwidth requirements.

6.3 Monitoring-Awareness

The proposed monitoring aware NoC design flow is depicted in Figure 7. The coupling of mapping, path selection and time-slot allocation from the original UMARS is extended with the mapping of transaction monitors to routers such that a full coverage of user channels is achieved. Here, we do not discuss the original UMARS mapping, routing and slot allocation; for these refer to [15].

As a *preprocessing* step to the modified UMARS, transaction monitors are virtually added to all routers (as this would be the maximum set of transaction monitors that we consider). These virtual monitors are added to the set of IPs present in the system. They are connected to the closest local NI, attached to the router they monitor.

Due to the centralized monitoring used, a single MSA is further added to the set of IPs and it gets its own NI. A single GT connection is assumed from any monitor to the MSA although yet of unknown required bandwidth. We consider monitoring connections as latency insensitive, so no latency constraints are added to them.

Monitor Placement. The loop of Algorithm 1 is extended with a fourth step, after a channel is allocated. This step is described in Algorithm 2. First, we check whether we need to insert additional monitoring. If the channel passes through a router that is monitored, we know, as channels are traversed in decreasing bandwidth order, that the monitor is able to monitor also this channel. Hence, nothing changes in this case. However, if none of the routers

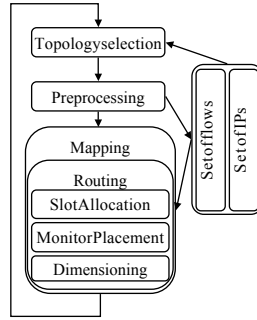


Fig. 7. Monitoring-aware design flow

that the channel passes through are yet monitored, we select one in Step 1a of Algorithm 2. We select a router with the highest arity on the channel path, because it maximizes the number of potential observed channels for this monitor. Once we select the router to be probed we are sure that the router will stay in the final set of transaction monitors. Therefore, the virtually probed router is added to the set of probed routers.

In Step 1b of Algorithm 2 we then add a channel from the now monitored router (and its associated NI) to the MSA. This channel is added to the set of unallocated flows.

Dimensioning. The requirement in terms of bandwidth is derived as a function of the channel that mandated the insertion of the probe. Note that the way in which the communication requirements are dimensioned does not impact the overall proposed design flow. For the transaction monitoring example we set the traffic numbers for the monitoring channels equal to the bandwidth required by the monitored channel. The next channel to be monitored by the same monitor, whose monitoring channel has been allocated, is guaranteed to require a lower bandwidth. As one monitor can only monitor one channel at a time, the previously allocated monitoring channel would be reused. The same holds if the monitoring channels would require, e.g. 10% of the monitored connection bandwidth, due to a higher abstraction power of the monitors.

Algorithm 2 Step four

- (1) If the path does not pass a monitored router
 - (a) Select a router on the path
 - (b) Add a channel from this router to the MSA
-

The newly added channel is a monitoring channel. The only difference between a genuine user channel and a monitoring channel is that we only want to monitor the user channels and not the monitoring channels themselves. Besides allocating the user and monitoring channels we also take care not to monitor the monitoring channels. Therefore, Algorithm 2 is only executed for user channels.

Results. If UMARS completes the allocation successfully, we have as results the

mapping, routing, slot allocation, monitor placement and monitoring dimensioning. After UMARS completes the allocation for all flows, all the routers in the set of probed routers have monitors attached. All the rest of virtual monitors are removed, as well as all the unallocated monitoring flows.

Iterations. If an allocation was not found by varying the slot table size till some predefined upper limit, the topology can be increased and the process repeated.

7 Experiments

7.1 Application Examples

Real Examples. We have used two real applications. (*mpeg*) an mpeg2 encoder/decoder using the main profile (4:2:0 chroma sampling) at main level (720x480 resolution with 15Mb/s) supporting interlaced video up to 30 frames per second. This application consists of 15 processing cores and an external SDRAM, and has 42 channels (with an aggregated bandwidth of 3GB/s), all configured to use guaranteed throughput, as presented in [12].

(*audio*) this application performs sample rate conversion, MP3, audio-postprocessing and radio. It closely resembles the chip presented in [20]. The application consists of 18 cores and has 66 channels all configured to use guaranteed throughput.

We have combined the two applications into four cases to be used as examples: mpeg (*Design1*), mpeg + audio (*Design2*), $2 \times$ mpeg + audio (*Design3*), $4 \times$ mpeg + audio (*Design4*).

Synthetic Examples. We have also generated synthetic application benchmarks for testing our proposed design flow. These benchmarks are structured to follow the application patterns of real SoCs. We have generated applications into two classes of such benchmarks, as presented in [22]: (i) Spread communication benchmarks (*Spread*), where each core communicates to a few other cores. These benchmarks characterize designs such as the TV processor that has many small local memories with communication evenly spread in the design. (ii) Bottleneck communication benchmarks (*Bottleneck*) where there are one or multiple bottleneck vertices to which the core communication takes place. These benchmarks resemble designs using shared memory/external devices such as the set-top boxes.

We have used spread communication of 12 IPs, in which every IP commu-

nicates with three others. We have used bottleneck communication with two converging points and 12 IPs. We have generated 500 synthetic application examples with spread and bottleneck communication.

7.2 Results

7.2.1 Setup

For both the real and synthetic application examples we have investigated what the original UMARS vs. monitoring-aware UMARS output is. The original UMARS generates the minimal NoC on which only the application requirements fit, while the monitoring-aware UMARS generates the minimal NoC on which both the application and monitoring requirements fit. To evaluate the performance of our approach, we looked at: (i) required number of transaction monitors, (ii) resulting topology size, (iii) resulting slot table size and (iv) resulting area.

For each application we have evaluated all possible meshes, from one by one up to seven by seven. For each of these topologies we have added one, two and three NIs per router, as depicted in Figure 8 and evaluated slot table sizes up to 65 TDMA slots. A larger slot table size mitigates overprovisioning due to granularity, but is often associated with a growth in buffer sizes as network consumption tends to become more bursty. Out of all the configurations for which UMARS finds an allocation, we present the one with lowest total area cost.

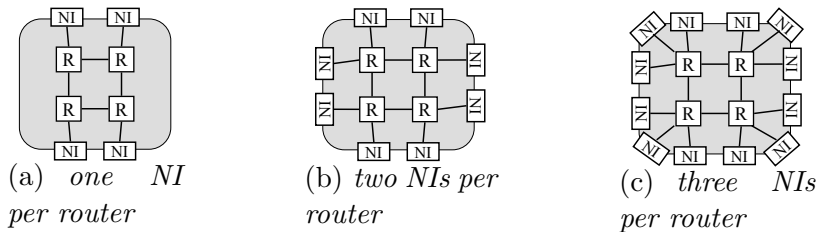


Fig. 8. NIs per router

Table 1 summarizes the results for the real examples, when one, two or three NIs per routers are tried. Due to the large communication demands, and given the constraints on topology and slot-table size we set for our experiments, *Design4* only fits on a topology using three NIs per router. For the synthetic examples, Figures 9 and 10 summarize the results for bottleneck and spread communication respectively. Each of the four aspects is discussed in detail in the following subsections. In Section 7.2.6, we briefly look at run-time reconfiguration of the monitoring system.

Table 1
Real Examples

<i>Designs</i>	<i>area</i>	<i>inc</i>	<i>size</i>	<i>mon</i>	<i>slot table</i>
<i>1NI/R</i>					
Design1	5.15	-	2x4	-	21
Design1+M	5.43	+5.5%	2x4	5	27
Design2	8.75	-	3x3	-	30
Design2+M	10.16	+16.1%	3x4	10	27
Design3	12.03	-	3x4	-	44
Design3+M	13.95	+16%	3x4	9	60
<i>2NIs/R</i>					
Design1	4.03	-	1x4	-	21
Design1+M	4.12	+2.2%	2x2	3	20
Design2	7.88	-	2x3	-	20
Design2+M	8.2	+3.9%	2x3	6	20
Design3	10.82	-	3x3	-	22
Design3+M	11.64	+7.6 %	2x4	8	29
<i>3NIs/R</i>					
Design1	3.62	-	1x2	-	30
Design1+M	3.85	+6.3%	1x3	3	18
Design2	6.97	-	1x3	-	27
Design2+M	7.16	+5.4%	1x3	3	30
Design3	10.26	-	2x3	-	21
Design3+M	10.78	+5%	2x3	6	22
Design4	18.45	-	3x4	-	21
Design4+M	19.07	+3.4%	2x4	8	36

7.2.2 Number of transaction monitors

For the synthetic cases with bottleneck communication, we see that the number of routers needed to be probed for full coverage varies between 50% and 100% with an average of 75%. Figure 9(a) displays the distribution. For spread communication Figure 10(a) displays the distribution. We see that the number of routers requiring a probe is higher compared to the bottleneck cases, but that is no surprise as the communication is more balanced (spread out)

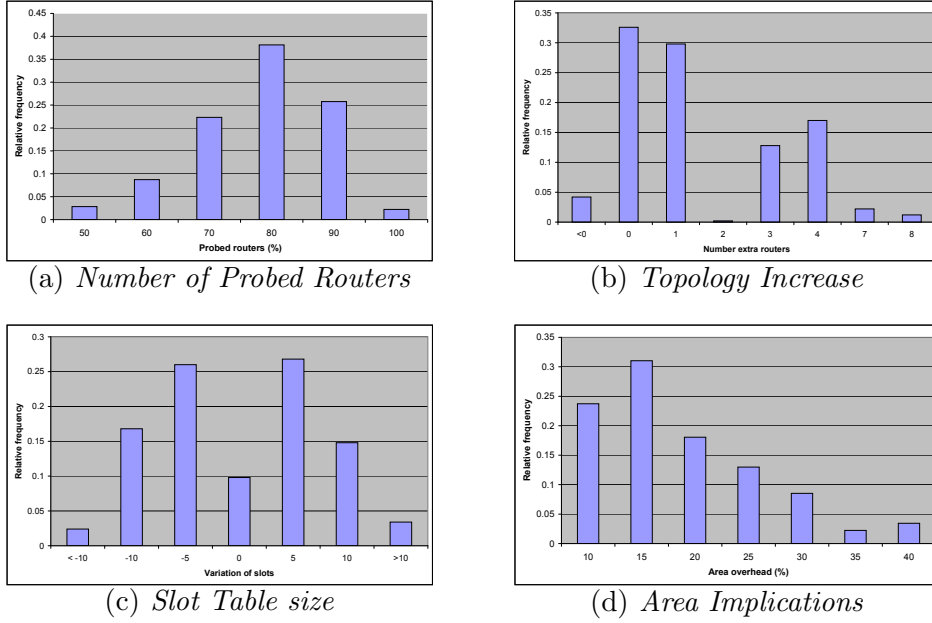


Fig. 9. Bottleneck

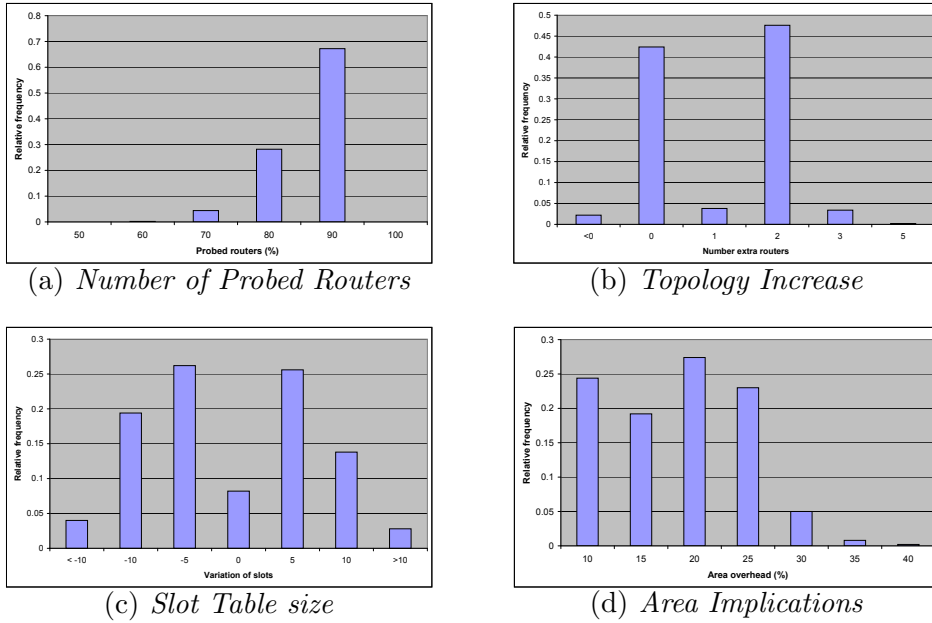


Fig. 10. Spread

over the routers. The minimum is 60% while the maximum is 90%. Hence, the interval is narrower than with bottleneck communication, the maximum is actually lower, and coverage of all routers was never required. Looking at the diagrams it is obvious that the number of routers needing probes is focused around the 80-90% bins.

Please note that the number of transaction monitors required is high because the \AA thetical NoC allows multiple IPs to be connected to the same NI and

multiple NIs to be connected to the same router. Therefore, channels can be very short, e.g. a channel between a master and a slave connected to the same NI will go through the NI starting from the master, then through one router and back to the same NI to the slave. All routers having at least one channel like this passing through will require one transaction monitor. Other NoCs may require a channel to pass through two different NIs, potentially lowering the number of transaction monitors being required.

For the real examples, see column *mon* in Table 1, showing the number of monitors and compare it to column *size* showing the mesh size. On average 87% of the routers need to be probed, but full coverage of routers with probes was required in 60% of the cases. Relating this with the area numbers from the same table, it is interesting to observe that the most area-efficient solutions required all routers probed. Therefore probing all routers must not be associated with area-inefficient solutions, the number of monitoring probes (in our case transaction monitors) being just one component which influences the total area cost of the monitoring solution.

7.2.3 Topology size

For the topology size we looked at the total number of routers employed. Figures 9(b) and 10(b) display the distribution for the synthetic examples. On average, topology stays the same (no extra routers required) or one or two extra routers are required. Increases in topology size with more than two routers, but with a maximum of 8, are still required in other cases, especially in the bottleneck applications. This can be explained because in bottleneck designs it is harder to accommodate the new monitoring channels due to the existing bottleneck vertices. Interesting is the fact that in 3-4% of the cases the number of routers actually decreased. This we can attribute to the heuristic nature of UMARS and to the higher number of slots used in the NoCs with monitoring.

For the real examples we see the number of routers kept constant in six cases, and both an increase and a decrease in two cases. The latter is accountable to an allocation found with a higher slot table size, see column *slot table* in Table 1.

In both real and synthetic examples we see that there is a good chance(30-60%) to find a solution on the same NoC topology, without requiring extra routers.

7.2.4 Slot Table size

Figures 9(c) and 10(c) display the distribution of the slot table size variation for the synthetic examples. In general we can see a similar shape for both bottleneck and spread communication examples. In a small number of cases (up to 10%) the slot table size is constant. It varies within a limit of ± 5 slots on a cumulated 50% of the cases. In roughly 30% of the cases the variation is between 5 and 10 slots, either in the negative or positive part. Higher variations than 10 slots are least visible in the figures.

For the real examples, we can observe the slot table size being constant in one case, bigger in six cases and smaller in three. Clearly, there exists a relation between the NoC topology and the slot table size.

In general a higher number of slots corresponds to adding the monitoring communication requirements on the same (or eventually smaller) NoC topology than the one used for user only communication requirements. A lower number of slots corresponds to a bigger NoC topology in the resulting shared NoC. The adapted UMARS design flow tries to balance these aspects.

7.2.5 Area

The total NoC area is derived according to the model in [11] extended with the area of the transaction monitors, $0.026mm^2$ per monitor in $0.13\mu m$ CMOS technology. Note that the total area presented includes NIs, routers and probes (transaction monitors). The area of NIs also accounts for buffer sizing in the NIs/NI ports corresponding to the real communication requirements of the users and monitors. The area numbers do not include the area of other IPs in the SoC, but refer to the NoC together with the complete monitoring service.

For bottleneck communication, area wise the cost is continuously below 50% with an average of 15%. Figure 9(d) shows the distribution of area overhead over the test cases and it is obvious that most lie in the left half of the span.

For spread communication Figure 10(d) shows the distribution. From an area point of view the overhead is between 10% and 40%, which again is a narrower interval than for bottleneck communication. In all it ends up on an average of 15% also for uniform traffic. No major difference in the area overhead is noticeable between uniform and bottleneck communication.

For the real examples the total area increase, see column *inc* in Table 1, amounts to between 2.2% and 16.1%. The area overhead is between 3% and 7% in the most area efficient case of three NIs per router which succeeded for all four designs. The resulting four designs we consider the end results of the monitoring-aware NoC design flow.

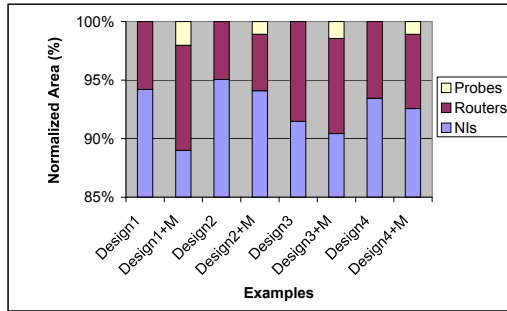


Fig. 11. Overall distribution of area

It is also interesting to see the overall distribution of this area between NIs, routers and monitors. This is presented in Figure 11 for the four designs in their most area-efficient case using $3NIs/R$. For the original designs the distribution of area between NIs and routers is shown. The main remark is that in all cases area of the transaction monitors is insignificant relative to the total area of the designs, dominated by the area of the NIs. Furthermore, in all designs the area of the monitors is even several times lower than the area of the routers involved.

7.2.6 Run-time reconfiguration

It has been previously mentioned that transaction monitors can be (re-)configured at run-time by means of write transactions. As a separate experiment we have looked at complete monitoring service configuration and evaluated the configuration options and the resulting configuration times. For this we have used the example MPEG design case using a 2x3 mesh topology, which was fully probed, resulting in six transaction monitors. We have used a centralized monitoring service with one MSA. We have used a slot table size of 128. Each transaction monitor uses a dedicated connection to the MSA. We have investigated both using the existing GT and BE communication services for monitoring system configuration. When using GT connections we have reserved a single slot for each monitoring connection.

We have tried two monitoring system-wide policies for configuration. One policy is based on simple write messages, which are not acknowledged by the transaction monitors. The total configuration time in this case is the time elapsed from the sending of the first message from the MSA to the first transaction monitor to be configured until the last received message at any of the transaction monitors. Note that the last message sent from the MSA may not be the last received message at the transaction monitors.

A second monitoring system configuration policy is based on acknowledgements. In this case, a 32-bit acknowledge is sent back from each of the transaction monitors upon reception of a configuration message and completion

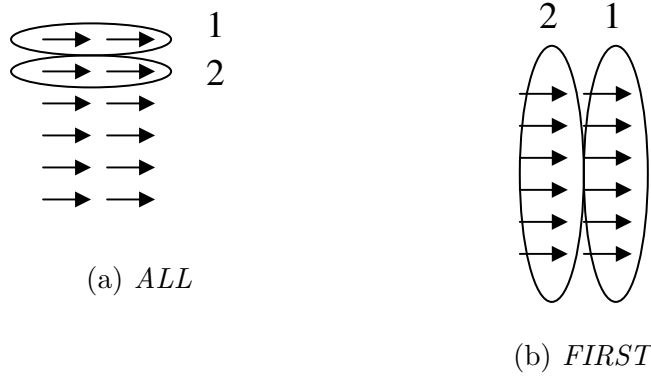


Fig. 12. Multiple configuration messages

Table 2
ALL-FIRST

<i>mpeg</i>	<i>ALL GT(ns)</i>	<i>ALL BE(ns)</i>	<i>FIRST GT(ns)</i>	<i>FIRST BE(ns)</i>
WR64	1212	78	1212	78
2xWR32	3378	174	1980	174
WR_ACK64	1832	152	1832	152
2xWR_ACK32	4136	224	2600	224

of the local transaction monitor configuration. The advantage of the second method is that the MSA knows when the monitoring system is configured. In this case the configuration time is the time elapsed between the time when the first message is sent from the MSA and the last acknowledgement is received at the MSA. Note that in general the acknowledgements are not received at the MSA in the sending order of the configuration messages from the MSA.

In the case of multiple configuration messages required for the same transaction monitor (e.g. two write messages) we have used two options. One option is to send all the messages for the same transaction monitor first then followed by all the messages for the second transaction monitor and so on. This is graphically depicted in Figure 12(a) and further referred as the ALL case. A second option is to send the first message to the first transaction monitor followed by the first message to the second transaction monitor, and so on, and only send the second message to all the transaction monitors when all the first messages for all transaction monitors have been sent, and so on. This is illustrated in Figure 12(b) and further referred to as the FIRST case.

Table 2 show the configuration time experimental results. On the first column we show use of write messages as described in Section 4.2.3, where WR64 and WR32 corresponds to a write with 64 bits or 32 bits of payload; 2xWR32 shows that two write messages are used for the configuration, while the presence of ACK shows the presence of a 32-bit acknowledgement in the configuration

process for a single transaction monitor. The table shows that using BE for configuration is several times faster than using GT as the configuration data does not have to wait for the reserved slot. This is expected because as soon as there is an empty slot or reserved but not used slot the BE configuration would sneak on the link. When using multiple configuration messages over GT monitoring connections for the same probe, it is more efficient to do it the FIRST way than to do it the ALL way. This is because in the ALL way the second configuration message for the first probe cannot be sent to the corresponding NI queue until there is space in the queue, thus delaying the first configuration message for the second probe. Table 2 finally shows the expected result that the use of acknowledgements increases the configuration time. Note that when using GT connections for configuration, the results in Table 2 do not account for the time required to set up these connections.

The results show that the run-time configuration is feasible for realistic cases, and the configuration time required for it is acceptable.

8 Conclusion

We propose a NoC design flow in which monitoring is taken into account at design time and is fully integrated in the flow. It automates the insertion of the monitors whenever their communication requirements are known, leading to a monitoring aware NoC design flow. Our flow was exemplified with the concrete case of transaction monitoring, in the context of the *Æthereal* NoC and UMARS design flow.

We are the first to quantify the complete cost of the complete monitoring solution accounting for the monitors, extra NIs, NI ports or enlarged topology needed to support monitoring in addition to the original application communication. Results show an area efficient solution for integrating monitoring in NoC designs. Monitors alone do not add much to the overall area numbers as the designs remain dominated by the area of NIs. We also considered the run-time reconfiguration of the monitoring system, showing acceptable reconfiguration times.

As future work we will look at more intelligent algorithms for application specific placement of transaction monitors and at the detailed trade-offs between the monitoring capabilities of single monitors and these placement strategies. We will also investigate whether an application independent monitor placement strategy can be developed.

References

- [1] L. Benini, G. De Micheli, Networks on chips: A new SoC paradigm, *IEEE Computer* 35 (1) (2002) 70–80.
- [2] T. Bjerregaard, J. Sparsø, A router architecture for connection-oriented service guarantees in the MANGO clockless network-on-chip, in: *Proc. Design, Automation and Test in Europe Conference and Exhibition (DATE)*, 2005.
- [3] E. Bolotin, I. Cidon, R. Ginosar, A. Kolodny, QNoC: QoS architecture and design process for network on chip, *Journal of Systems Architecture* 50 (2–3) (2004) 105–128, special issue on Networks on Chip.
- [4] C. Ciordas, T. Basten, A. Rădulescu, K. Goossens, J. van Meerbergen, An event-based monitoring service for networks on chip, *ACM Transactions on Design Automation of Electronic Systems* 10 (4) (2005) 702–723, HLDVT’04 Special Issue on Validation of Large Systems.
- [5] C. Ciordas, K. Goossens, T. Basten, A. Rădulescu, A. Boon, Transaction monitoring in networks on chip: The on-chip run-time perspective, in: *Proc. Symposium on Industrial Embedded Systems (IES)*, 2006.
- [6] C. Ciordas, K. Goossens, A. Rădulescu, T. Basten, NoC monitoring: Impact on the design flow, in: *Proc. Int’l Symposium on Circuits and Systems (ISCAS)*, 2006.
- [7] C. Ciordas, A. Hansson, K. Goossens, T. Basten, A Monitoring-aware NoC Design Flow, in: *Proc. Euromicro Symposium on Digital System Design (DSD)*, 2006.
- [8] M. Dall’Osso, G. Biccari, L. Giovannini, D. Bertozzi, L. Benini, xpipes: A latency insensitive parameterized network-on-chip architecture for multi-processor socs, in: *Proc. Int’l Conference on Computer Design (ICCD)*, 2003.
- [9] W. J. Dally, B. Towles, Route packets, not wires: on-chip interconnection networks, in: *Proc. Design Automation Conference (DAC)*, 2001.
- [10] Device Transaction Level (DTL) Protocol Specification. Version 2.2, in: Philips Semiconductors, 2002.
- [11] S. González Pestana, E. Rijpkema, A. Rădulescu, K. Goossens, O. P. Gangwal, Cost-performance trade-offs in networks on chip: A simulation-based approach, in: *Proc. Design, Automation and Test in Europe Conference and Exhibition (DATE)*, 2004.
- [12] K. Goossens, J. Dielissen, O. P. Gangwal, S. González Pestana, A. Rădulescu, E. Rijpkema, A design flow for application-specific networks on chip with guaranteed performance to accelerate SOC design and verification, in: *Proc. Design, Automation and Test in Europe Conference and Exhibition (DATE)*, 2005.

- [13] K. Goossens, J. Dielissen, A. Rădulescu, The *Æthereal* network on chip: Concepts, architectures, and implementations, *IEEE Design and Test of Computers* 22 (5) (2005) 21–31.
- [14] P. Guerrier, A. Greiner, A generic architecture for on-chip packet-switched interconnections, in: *Proc. Design, Automation and Test in Europe Conference and Exhibition (DATE)*, 2000.
- [15] A. Hansson, K. Goossens, A. Rădulescu, A unified approach to constrained mapping and routing on network-on-chip architectures, in: *Int’l Conf. on Hardware/Software Codesign and System Synthesis (CODES+ISSS)*, 2005.
- [16] J. Hu, R. Marculescu, Exploiting the routing flexibility for energy/performance aware mapping of regular NoC architectures, in: *DATE ’03: Proceedings of the conference on Design, Automation and Test in Europe*, IEEE Computer Society, 2003.
- [17] F. Karim, A. Nguyen, S. Dey, An interconnect architecture for networking systems on chips, *IEEE Micro* 22 (5) (2002) 36–45.
- [18] I. Matta, A. Bestavros, A load profiling approach to routing guaranteed bandwidth flows, in: *IEEE INFOCOM*, vol. 3, 1998.
- [19] M. Millberg, E. Nilsson, R. Thid, S. Kumar, A. Jantsch, The *Nostrum* backbone - a communication protocol stack for networks on chip, in: *Proc. Int’l Conference on VLSI Design*, 2004.
- [20] A. Moonen, R. van den Berg, M. Bekooij, H. Bhullar, J. van Meerbergen, A multi-core architecture for in-car digital entertainment, in: *GSPx*, 2005.
- [21] R. B. Mouhoub, O. Hammami, Noc monitoring hardware support for fast noc design space exploration and potential noc partial dynamic reconfiguration, in: *IES*, 2006.
- [22] S. Murali, M. Coenen, A. Rădulescu, K. Goossens, G. De Micheli, A methodology for mapping multiple use-cases on to networks on chip, in: *Proc. Design, Automation and Test in Europe Conference and Exhibition (DATE)*, 2006.
- [23] S. Murali, G. De Micheli, Bandwidth-constrained mapping of cores onto NoC architectures, in: *Proc. Design, Automation and Test in Europe Conference and Exhibition (DATE)*, 2004.
- [24] V. Nollet, T. Marescaux, P. Avasare, D. Verkest, J.-Y. Mignolet, Operating-system controlled network on chip, in: *Proc. Design Automation Conference (DAC)*, 2004.
- [25] M. Pastrnak, et al., Combined reservation and adaptation QoS for improving picture quality and resource usage of multimedia (NoC) chips, in: *International Symposium on Consumer Electronics*, 2006.

- [26] A. Rădulescu, J. Dielissen, S. González Pestana, O. P. Gangwal, E. Rijpkema, P. Wielage, K. Goossens, An efficient on-chip network interface offering guaranteed services, shared-memory abstraction, and flexible network programming, *IEEE Transactions on CAD of Integrated Circuits and Systems* 24 (1) (2005) 4–17.
- [27] J. W. van den Brand, Runtime networks-on-chip performance monitoring, Technical Report 2006/00218, Philips Research (Mar. 2006).
- [28] J. W. van den Brand, C. Ciordas, K. Goossens, T. Basten, Congestion-controlled best-effort communication for networks-on-chip, in: *Proc. Design, Automation and Test in Europe Conference and Exhibition (DATE)*, 2007.
- [29] D. Wingard, Socket-based design using decoupled interconnects, in: J. Nurmi, H. Tenhunen, J. Isoaho, A. Jantsch (eds.), *Interconnect-Centric Design for Advanced SoC and NoC*, chap. 15, Kluwer, 2004, pp. 367–396.