# Cost-efficient Network-on-Chip Design Using Traffic Monitoring System

Kwanho Kim, Donghyun Kim, Kangmin Lee, and Hoi-Jun Yoo

Dept. of EE&CS, Korea Advanced Institute of Science and Technology (KAIST)
373-1, Guseong-dong, Yuseong-gu
Daejeon, 305-701, Republic of Korea
E-mail : kkh82@eeinfo.kaist.ac.kr

## Abstract

An in-depth NoC traffic monitoring system is presented for accurate evaluation and refinement of the NoC for the target application. It measures various real-time traffic parameters such as an end-to-end latency, queuing buffer usage and congestion level. A NoC-based portable multimedia system is implemented as a prototype on multiple FPGAs in order to demonstrate the effectiveness of the proposed traffic monitoring system. Utilizing the traffic monitoring system, the implemented NoC is diagnosed and two NoC parameters are modified: buffer size assignment and run-time routing path modification. As a result, a 42% reduction of buffer size and a 28% latency reduction are obtained without performance degradation in the mesh-connected NoC.

## Introduction

NoC involves a complex design process to select a topology, protocol, buffer size, routing and scheduling algorithm that are suitable for the target application [1]. The selection of the NoC design parameters should be made based on exact awareness of the application-specific on-chip real traffic patterns. A NoC traffic monitoring system is essential for the optimized NoC design because optimization requires knowledge of real traffic patterns and their effects on NoC performance.

In this paper, a real implementation of a NoC traffic monitoring system and its application to cost-efficient NoC design are reported. The traffic monitoring system observes the real-time internal traffic behavior of the NoC core to obtain cycle-accurate traffic information. We verify the effectiveness of the traffic monitoring system by applying it to a prototype of a NoC-based portable multimedia SoC implemented on multiple FPGAs. Through an analysis of the monitoring results, potential performance bottlenecks of the NoC can be identified and the design can be refined by cutting overestimated NoC resources. As a result, a cost-efficient and higher-performance NoC design is possible due to the exact information provided by the traffic monitoring system. The details of the traffic monitoring system and various experimental results in a multimedia application will be described in the following sections.

## Traffic Monitoring System

Fig. 1 shows the overall structure of the traffic monitoring system. It consists of three sub-systems: the Host Interface, the Central Controller, and Monitor Units.

The Host Interface, a bridge between the central controller and a host PC, transfers traffic monitoring results to the host PC via an Ethernet line. The central controller enables/disables each monitoring unit based on the requested monitoring regional scope and time interval.

The Monitor Unit consists of a traffic probe, a traffic manager, and a traffic memory. A traffic probe module is connected to a switch or a network interface in order to trace the real-time traffic
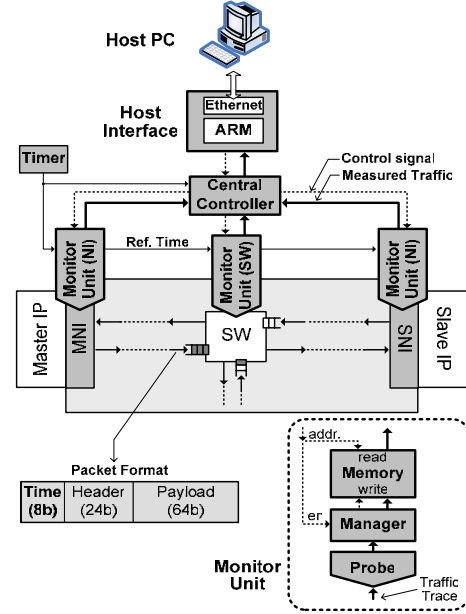


Figure 1. Proposed traffic monitoring system

parameters such as end-to-end latency, communication bandwidth, queuing buffer backlog and link utilization. The traffic manager then stores the traces in its local trace memory after attaching a time-stamp to each trace using a global timer connected to all monitor units. During the operation of the specific application, all monitoring results are stored in the corresponding local memory, which is accessible by the Host Interface's CPU via the central controller.

The traffic monitoring system has a modular architecture. Thus, the monitor unit can be attached to any NoC component at design-time.

## Application to Portable Multimedia System

### A. System Operation

We implement a portable multimedia system to demonstrate the effectiveness of the proposed traffic monitoring system. The implemented system includes various IPs: five masters (RISC CPU
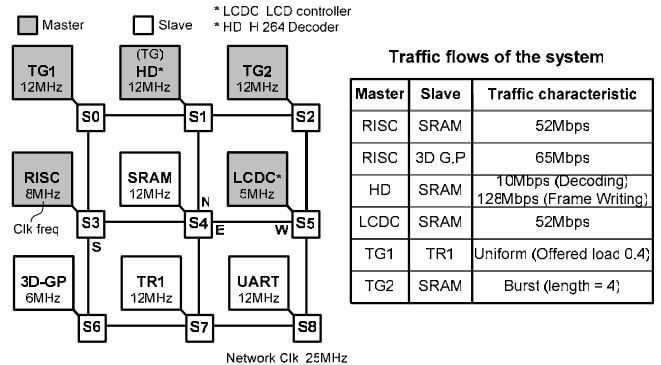


Figure 2. Portable multimedia system in Mesh topology

core, LCD controller, H.264 decoder, and two traffic generators (TG)) and four slaves (3-D graphics processor [3], SRAM, UART, and traffic receptor (TR)) as shown in Fig. 2. The H.264 decoder (HD) is replaced with a traffic generator since it requires too much logic-element resources to be implemented on a FPGA board. The traffic generator produces the real HD traces assuming that it decodes the CIF (352x288) H.264 baseline profile at level 2 with 30frames/sec.

## B. NoC Design Refinement

### 1) Buffer Size Assignment

The input queuing buffers in a switch occupy a significant portion of the chip area, and thus the buffer size should be minimized without significant throughput degradation. The initial NoC design has uniformly assigned input buffers of 5-packet capacity. Although the uniform choice of the input buffer size is quite straightforward and widely used in current NoC designs, it may lead to wasted silicon area because the size of all input buffers is uniformly assigned for the worst-case traffic pattern of the target application. Based on the backlog monitoring results, the proper buffer size can be decided as shown in Fig. 3. The buffer that has the maximum backlog is selected as the bottleneck buffer and the selected buffer size is increased by one until the execution time of the application reaches that of a uniform buffer configuration. As a result of the buffer size assignment, 42% total buffer size reduction is obtained without performance degradation.

### 2) Run-time Routing Path Modification

In NoCs, a deterministic routing scheme such as source routing is widely used [2] because it is a cost-effective scheme for network and transport layer design. In conventional source routing, a packet is transferred to a destination through a fixed routing path that is defined at design time. In the case of an application with time-varying traffic flows, however, fixed routing path selection may lead to performance degradation. To overcome the limitation on fixed source routing, in this work, a packet routing path is dynamically selected based on the traffic monitoring results.
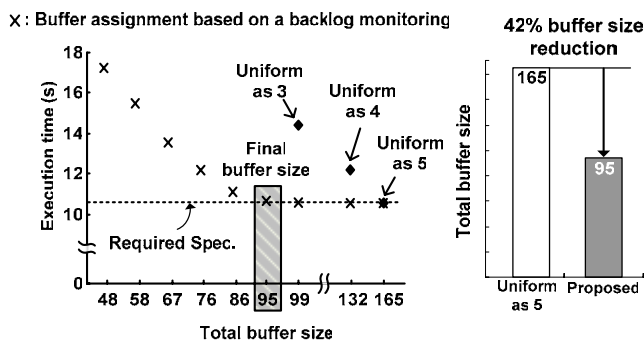


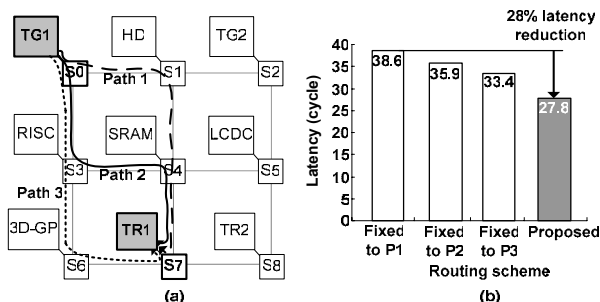Figure 3. Buffer size assignment based on monitoring results



Figure 4. (a) Alternate routing paths from TG1 to TR1 and (b) Latency comparison

There are three alternate routing paths from TG1 to TR1 as shown in Fig. 4(a). Before the routing path modification, three test packets are generated periodically at the source network interface and the latencies of the test packets are measured at the destination network interface using the latency monitoring framework. As a result, the routing path that has the shortest latency is selected as the next routing path, and the routing table at the source network interface is updated. Fig 4(b) shows that 28% average latency reduction is obtained compared to the worst fixed source routing and the variation is also diminished significantly.

## FPGA Prototype Implementation

Fig. 5 shows the overall system implemented on three Altera Stratix EP1S60 series FPGAs. The RISC and the 3D-GP are integrated on the left and right FPGA, respectively. The central FPGA integrates the total NoC with the proposed traffic monitoring system and other IPs such as the UART, the LCD controller, the H.264 traffic generator and the user interface logic. All implemented modules are designed in RTL level and 6 different clocks are used for plesiochronous communications.
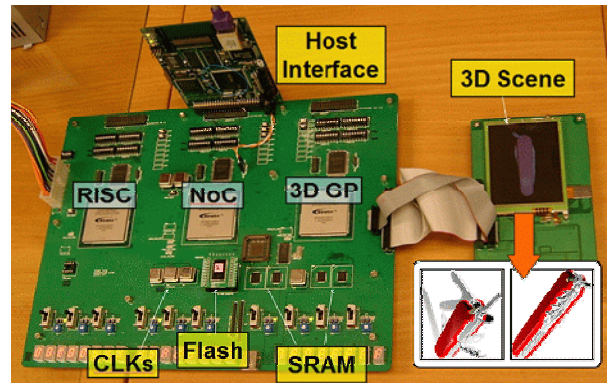


Figure 5. NoC-based system implemented on FPGA board

## Conclusion

A NoC traffic monitoring system is proposed to probe the internal traffic of the NoC cores at run-time for accurate performance evaluation and in-depth refinement of application-specific NoCs. A portable multimedia system is implemented on FPGAs to demonstrate the effectiveness of the traffic monitoring system. Based on the monitoring diagnosis, the target system is refined in two ways: buffer size assignment and run-time routing path modification. Through utilization of the traffic monitoring system, buffering cost and latency reduction is obtained up to 42% and 28%, respectively in a mesh topology. Thus, it is demonstrated that the traffic monitoring system can be effectively used to implement cost-efficient and high performance application-specific NoC-based systems.

## References

[1] N. Genko, et al., "A complete network-on-chip emulation framework", in *Proc. Design, Automation and Test in Europe Conf.* Mar. 2005, pp. 246-251.
[2] Se-Joong Lee, et al., "An 800MHz Star-Connected On-Chip Network for Application to Systems on a Chip," *IEEE Int. Solid-State Circuits Conf. Dig. Tech. Papers*, Feb. 2003, pp. 468-469.
[3] Jeong-Ho Woo, et al., "A 1.2Mpixels/s/mW 3-D Rendering Processor For Portable Multimedia Application", in *IEEE Asian Solid-State Circuits Conf.*, Nov. 2005, pp. 297-300.