**A General Framework for Average-Case** Performance Analysis of Shared Resources

Sahar Foroutan<sup>1</sup>, Benny Akesson<sup>2</sup>, Kees Goossens<sup>2</sup>, and Frederic Petrot<sup>1</sup> <sup>1</sup>TIMA Laboratory, Grenoble, France <sup>2</sup>Eindhoven University of Technology, Eindhoven, the Netherlands {sahar.foroutan, frederic.petrot}@imag.fr {k.b.akesson, k.g.w.goossens}@tue.nl



**EUROMICRO DSD/SEAA** SANTANDER, SPAIN SEPTEMBER 4-6, 2013



# Soft Real-Time Requirements in Shared Memory MPSoCs

- Embedded systems are growing in complexity as more and more applications are integrated into MPSoCs.
- Applications share resources (e.g. processors , interconnect, and memories ) to reduce cost, resulting in temporal interference.
- Some applications, e.g. video decoders, have soft real-time (SRT) requirements, where deadlines must be satisfied with high probability.
- Cost-efficient performance analysis requires the average execution time of the SRT applications to be accurately estimated.

2) Application

### **Average Performance Analysis**

- To determine the average execution time of application:  $\bullet$ 
  - $\blacktriangleright$  Average latency of shared resources is required.
  - > Probabilistic models, e.g. **<u>queuing theory</u>**, can be used. 1) Resource Modeling

Any performance analysis needs to:  $\bullet$ 

- 1. Model shared resources (i.e. arbitration mechanism)
- 2. Characterize the **application** (i.e. traffic)

Existing queuing models: 

> Are either based on an <u>exponential arrival traffic</u> or address only a single arbitration mechanism, thus failing to recognize the diversity of arbiters in complex systems.

In an MPSoC, access to the shared resources may be provided by different arbiters, e.g. :

- Time-Division Multiplexing (**TDM**),
- Static-Priority (SP), or
- Round-Robin (**RR**)

### **Contribution 1:**

A high-level model for resource sharing that can be used with <u>different arbiters.</u>

TDM RR SP Arbiter  $\mu$  ,  $\sigma_{s}$ S  $\lambda_p, \sigma_p$ 

Models assuming exponentially distributed traffic do not cover dynamic applications executing on MPSoCs. Characterization

E.g. histograms of request intervals of JPEG and H.263 clearly show a non-exponential distribution.

### **Contribution 2:**

Our high-level model is based on general distributions of request interval and service times (i.e. G/G/1).



## **General Resource Model Based on Queuing Theory**

• The <u>analytical</u> model aims to estimate the average waiting time a request arriving to a multiple-queue resource with *p* queues, spends in any queue *i*. This is called the *queuing delay* (*W<sub>i</sub>*) :



(1) is the waiting time due to requests in the <u>same queue</u>, depends on:

- *n<sub>i</sub>*: average number of requests already waiting in queue *i*
- **Ts**<sub>*i*</sub>: Average service time of queue *i*

( $\Pi$ ) represents the interference from <u>other queues</u>. At the arrival of the considered request to queue *i*, there is a number of requests  $(n_i)$  in any other queue **j** that depending on <u>arbitration</u> may be served earlier.

•  $A(n_{i}, n_{j}) = Arbitration Indicator$ : determines the interference of

# **Analytical vs. Simulation Performance Results**

#### **Simulation Platform**

- A cycle-accurate SystemC model of a real-time memory controller, supporting a variety of memories and arbiters.
- Experiments consider a 32-bit **<u>SRAM</u>** with a peak bandwidth of 2 GB/s.
- **Four** requestors share the bandwidth of the memory.
- Traffic (read and write requests toward the SRAM) is injected by traffic generators issuing either **synthetic** traffic or **real application traces**.
- **<u>Request size</u>** is set to 64 B (16 words), resulting in uniform service times of 16 cycles for both reads and writes.

#### **Synthetic Traffic Results**

- Average <u>memory latency</u> (queuing + architectural delays), plotted against <u>bandwidth</u>.
- Synthetic traffic: Request intervals follow a <u>normal distribution</u> with:
  - <u>**Mean</u>** request intervals = 1/bandwidth</u>
  - **Standard deviation** *σ*=30, 60, 90, 120 ns for TDM and RR, and *σ*=120 for SP

#### • For **SP**, 4 requestors with different priority levels.

• <u>Analytical</u> vs. <u>simulation</u> **average error**: 4.1% for TDM, 12.9% for SP, and 14.2% for RR.

other queues on queue *i* and thus models the arbitration.

(III)  $\underline{R}$  is the average <u>residual times</u> (remaining service times) of all queues in multiple G/G/1 models.

 $A(n_i, n_j)^{TDM} = n_i \quad for \ 1 \le j \le p$  $A(n_{i}, n_{j})^{SP} = \begin{cases} n_{j}, & \text{if } \boldsymbol{j} \text{ has a higher priority than } \boldsymbol{i} \\ 0, & otherwise \end{cases}$  $A(n_{i}, n_{j})^{RR} = Min(n_{i}, n_{j}) \text{ for all } \boldsymbol{j}$ 

