# Analysis of Shared Coprocessor Accesses in MPSoCs

## Overview

Simon Schliecker
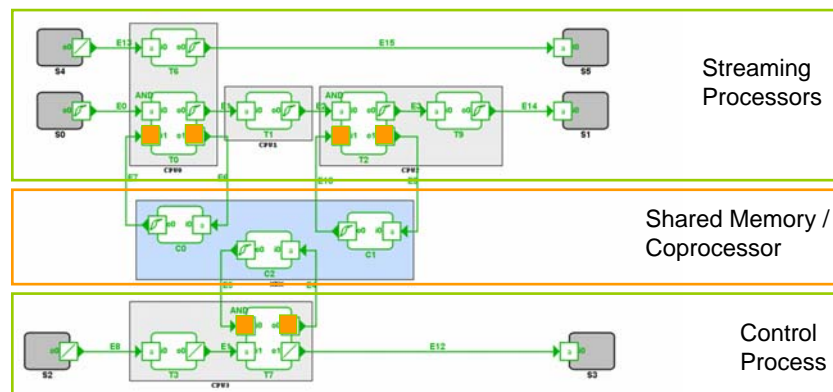Matthias Ivers
Rolf Ernst                    Bologna, 22.05.2006
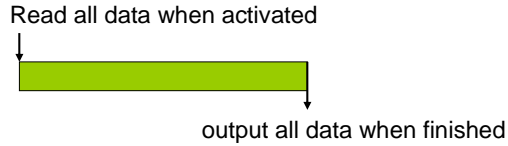
**IDA Institut für Datentechnik und Kommunikationetze**

---

# System Setup



Streaming Processors

Shared Memory / Coprocessor

Control Process

**IDA Institut für Datentechnik und Kommunikationetze**

# New Task Model

Read all data when activated

**Classical Task-model**

$C_i$ Core Execution Time

output all data when finished

Read *some* data when activated

**"Communicating Task"-model**

$C_i$ Core Execution Time

$Q_i$ Set of incurred *transactions*

Output *some* data when finished

Access shared resource multiple times during execution.

Synchronization still only at task activation.

**Institut für Datentechnik und Kommunikationetze**

3

---

# Effects on Coprocessor Accesses

Load from other system parts

Coprocessor (e.g. **load-dependant**)

Bus Transfer (e.g. **volume-dependant**)

Execution of Task A on the CPU **rescheduling effects**

WCRT

t

Executed Instructions

```
byte x = 10;
long y = 10;
z = COP(y);
```

```
x = x * z;
y = y / z;
z = COP(x);
```

```
x = x * z;
y = y / z;
z = x + y;
```

**Institut für Datentechnik und Kommunikationetze**

4

2

**Analysis Steps**

Streaming Processors

Shared Memory / Coprocessor

Control Process

Output Event Models     Local Response Times

Request Latencies

**Institut für Datentechnik und Kommunikationetze**

5



**Analysis Distribution**

System level      Local level

$E^{env}$

$E^{in}$

Task Analyses

Output Traffic Analysis

$E^{out}$

Response Time Analysis

$E^{in}$

$E^{in}$

Request Latency Analysis

Task Analyses

$\mathbf{Q}$

$S(\mathbf{Q})$

Output Traffic Analysis

$E^{out}$

Response Time Analysis

modifications for communicating tasks

**Institut für Datentechnik und Kommunikationetze**
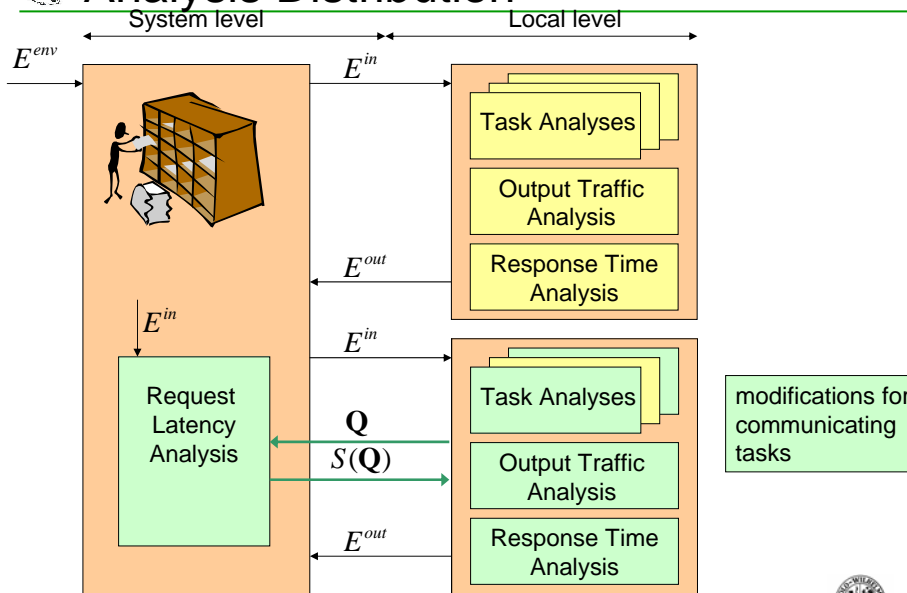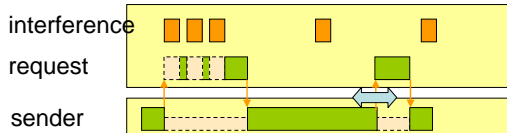
6

3

# Correlated Event Total Busy Times

- Not every request experiences the worst case in dynamic systems!

Problem: Can not predict exact request times!

interference
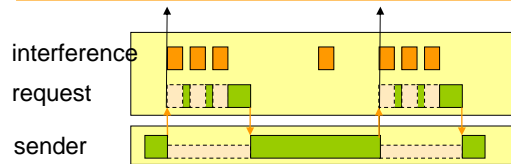
request

sender

variable times of interference
+ variable times of request
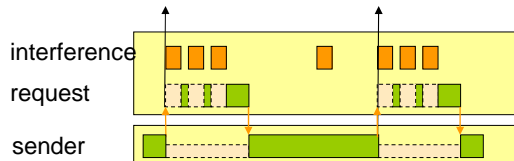___
highly *variable* and *sensitive* request latencies!

Possible solution: assume "worst case for every request"

interference

request

sender

---

# Correlated Event Total Busy Times
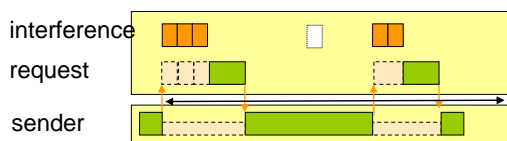
- Deriving individual bounds is difficult
- Using individual bounds for analysis is difficult
- but ignoring it is easy!

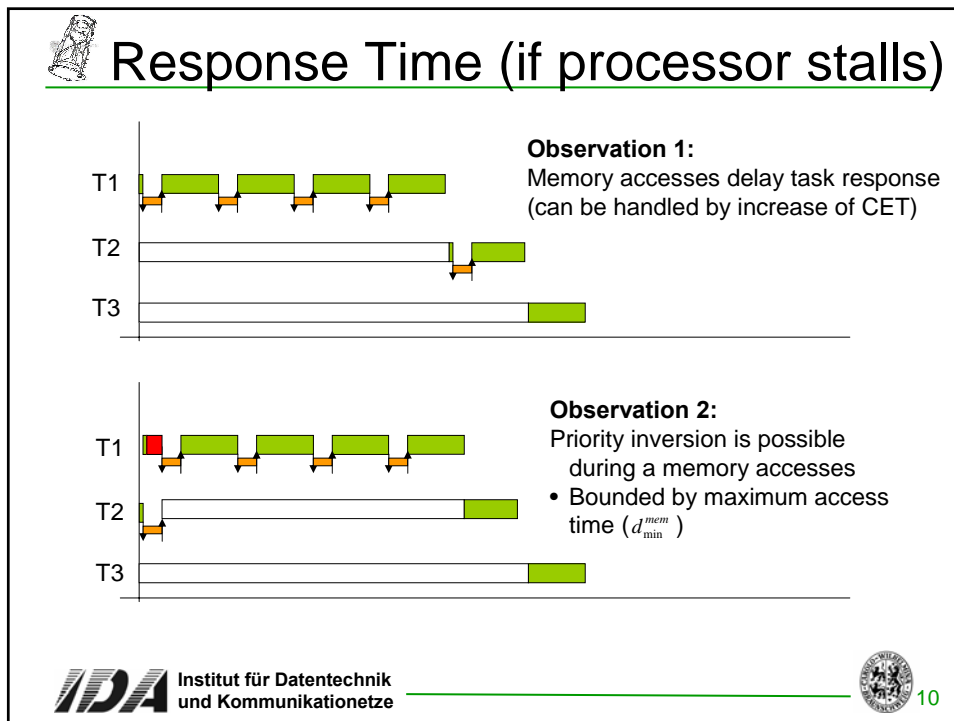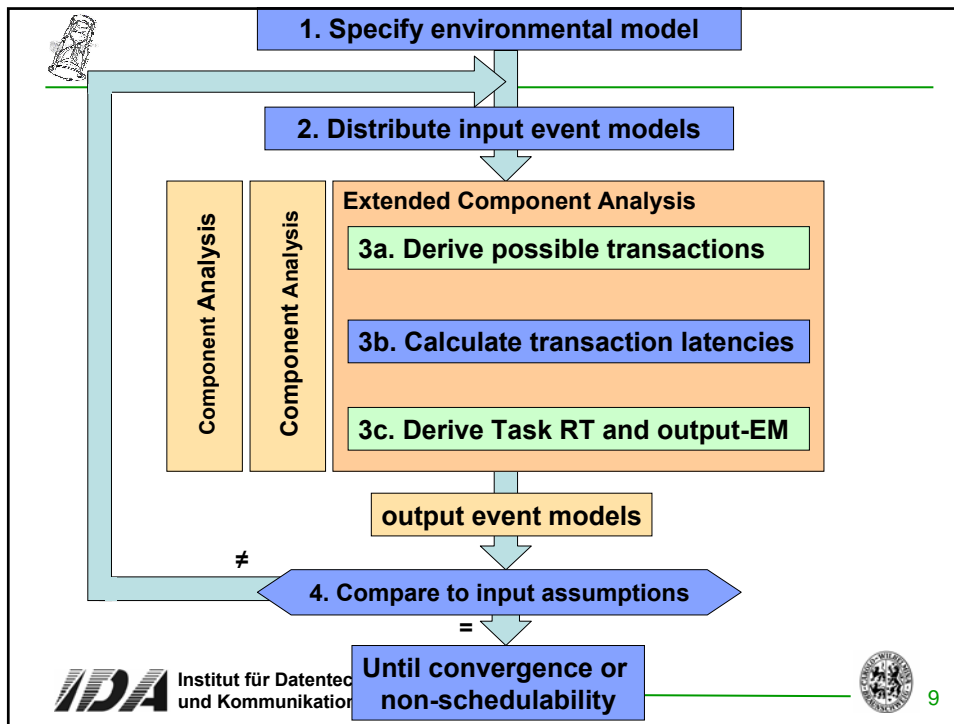Possible solution: assume "worst case for every request"

interference

request

sender

no remaining overestimation given interference with large jitter!

Better solution: approximate "maximum *total* busy time"

interference

request

sender

total interference and all requests in total time window ⇒ "total busy time"

## Slide 9

**1. Specify environmental model**

**2. Distribute input event models**

**Extended Component Analysis**

**3a. Derive possible transactions**

**3b. Calculate transaction latencies**

**3c. Derive Task RT and output-EM**

Component Analysis

Component Analysis

**output event models**

≠

**4. Compare to input assumptions**

=

**Until convergence or non-schedulability**

Institut für Datentec... und Kommunikatio...

9

## Slide 10

# Response Time (if processor stalls)

T1

T2

T3

**Observation 1:**
Memory accesses delay task response (can be handled by increase of CET)

T1

T2

T3

**Observation 2:**
Priority inversion is possible during a memory accesses
• Bounded by maximum access time $(d_{\min}^{mem})$

Institut für Datentechnik
und Kommunikationetze

10

# Scheduling Analysis

- If processor stalls during Memory requests:
  - Processor is NOT released, this extends CET.
  - Higher priority tasks can be *blocked* by maximum memory access time.
  - Buffer is always empty, because previous requests finished.

$$R_i = B_i + C_i + S(Q_i) + \sum_{j=1}^{i+1} \eta_j(R_i) \cdot (C_j + S(Q_j))$$

$$B_i = d_{min}^{mem}$$

CET and. CoP times

activations of hp task

CET and CoP times

Total blocking time

---

# Scheduling Analysis (2)

- The more requests are considered together, the smaller the overestimation!
  - Collect all requests that can lead to delay and add maximum total busy time
  - Perfect match for improved path latencies

$$R_i = (B_i + C_i + \sum_{j \in hp(i)} \eta_j(R_i) \cdot C_j) + S(\underset{j \in hp(i)}{Y} Q_j))$$

$$B_i = d_{min}^{mem}$$

Total CoP times that can lead to delay of task i

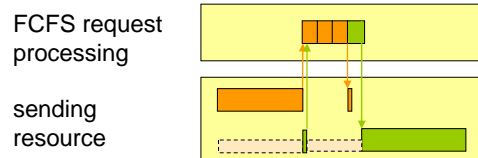# Multithreading from real time perspective

Generally:
  – stalling decreases and
  – processor utilization increases
But: Additional interference for requests
  – interference from previous requests can completely compensate the gain of reduced stalling!
  – A task can be fully delayed by higher priority execution **and** requests
  – FCFS ordering along request chain counters priorities on sending resource

⇒ **no gain for response time** under given task assumptions

FCFS request processing

sending resource

---

# Additional Critical Sections

- How much blocking time to take into account?
  – Blocking Memory Accesses can be "nested" into Critical Sections (not the other way around)
  – Assume a virtual semaphore "memory":
    • All tasks require "memory" to be free to start executing
    • Some tasks spend no time accessing "memory", but still must wait until it is free
    • Other tasks access "memory", and may enter the critical section multiple times
  – Memory Accesses are "automatically" protected with highest priority!
  – Problem mapped to "nested critical sections problem" (Sha, Rajkumar)
    • Depends on utilized protocol PIP, PCP

high priority task blocked!

T1

T2

Task in critical section