# Evaluation and Validation

Peter Marwedel
TU Dortmund, Informatik 12
Germany
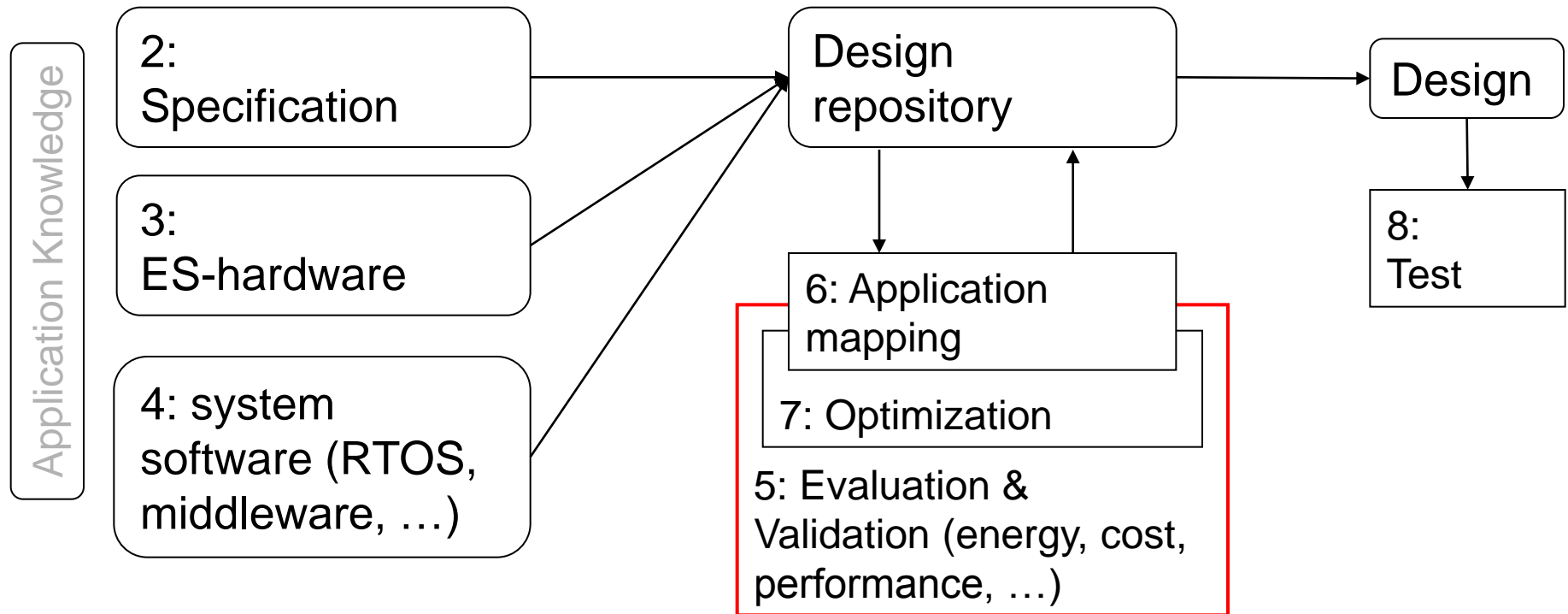
Embedded Systems

© Springer, 2010

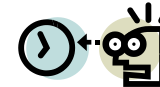2012 年 12 月 05 日

# Structure of this course



Numbers denote sequence of chapters

# How to evaluate designs according to multiple criteria?

Many different criteria are relevant for evaluating designs:

➡ - Average & worst case delay

- power/energy consumption

- thermal behavior

- reliability, safety, security

- cost, size

- weight

- EMC characteristics

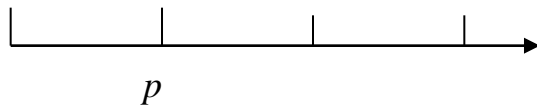- radiation hardness, environmental friendliness, ..

How to compare different designs?
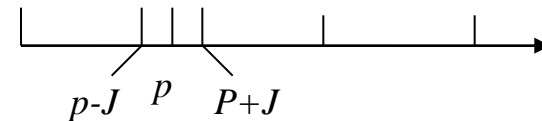(Some designs are "better" than others)

# Real-time calculus (RTC)/ Modular performance analysis (MPA)

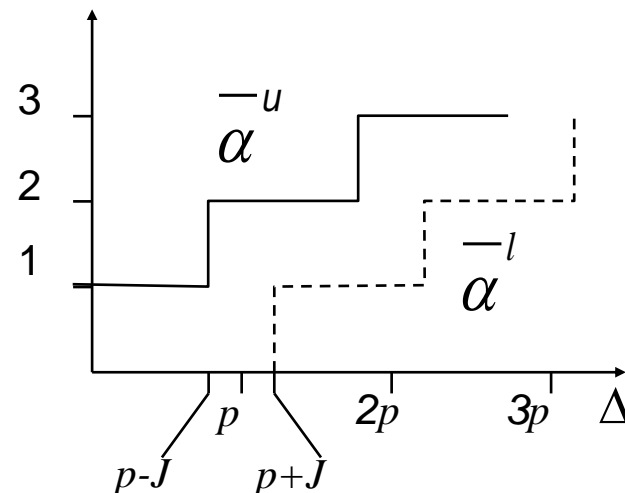## Streams of events important: Examples
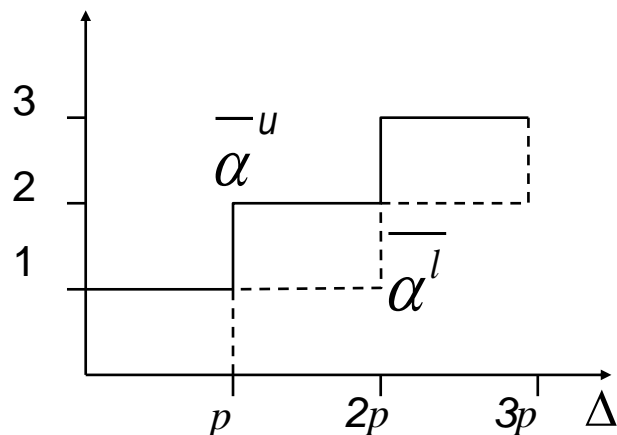
periodic event stream

periodic event stream with jitter
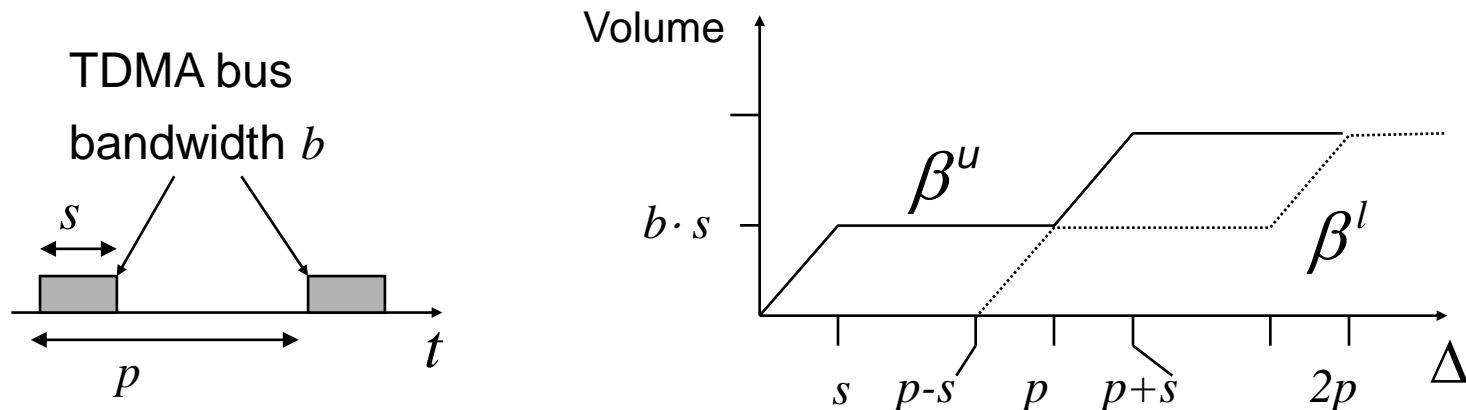


Thiele et al. (ETHZ): Extended **network calculus:**
**Arrival curves** describe the maximum and minimum number of events arriving in some time interval $\Delta$.

# RTC/MPA: Service curves

Service curves $\beta^u$ resp. $\beta^\ell$ describe the maximum and minimum service capacity available in some time interval $\Delta$
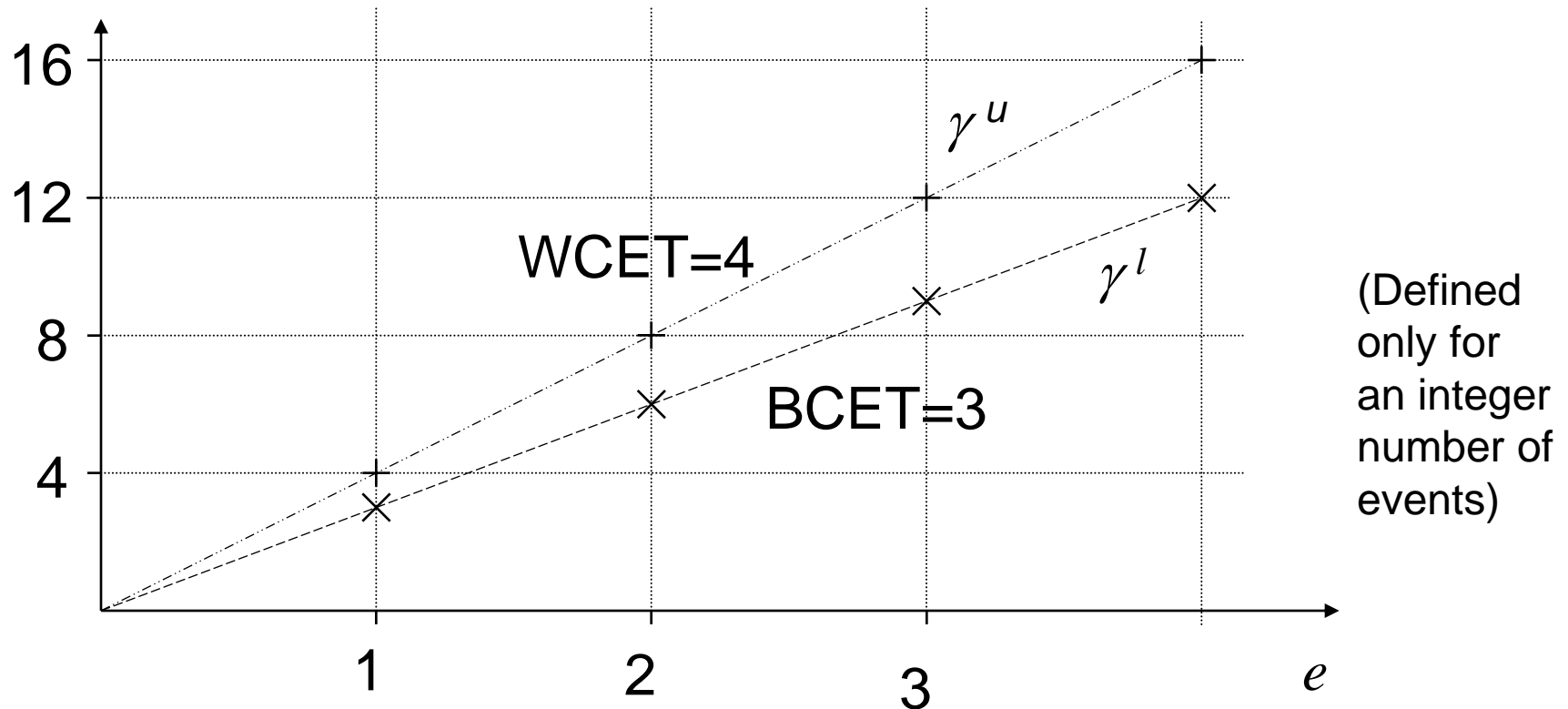
Example:

# RTC/MPA:
# Workload characterization

$\gamma^u$ resp. $\gamma^\ell$ describe the maximum and minimum service capacity required as a function of the number $e$ of events. Example:



(Defined only for an integer number of events)

technische universität dortmund

fakultät für informatik

# RTC/MPA:
# Workload required for incoming stream

Incoming workload

$$\alpha^u(\Delta) = \gamma^u\left(\overline{\alpha^u}(\Delta)\right) \qquad \alpha^l(\Delta) = \gamma^l\left(\overline{\alpha^l}(\Delta)\right)$$

Upper and lower bounds on the number of events

$$\overline{\beta}^u(\Delta) = \left(\gamma^l\right)^{-1}\left(\beta^u(\Delta)\right) \qquad \overline{\beta}^l(\Delta) = \left(\gamma^u\right)^{-1}\left(\beta^l(\Delta)\right)$$

# RTC/MPA:
# System of real time components

Incoming event streams and available capacity
are transformed by real-time components:

$$\left[\overline{\beta}^{l}, \overline{\beta}^{u}\right]$$

$$\left[\overline{\alpha}^{l}, \overline{\alpha}^{u}\right] \qquad \left[\overline{\alpha}^{\ell}{}', \overline{\alpha}^{u}{}'\right]$$

RTC          RTC"

$$\left[\overline{\beta}^{l}{}', \overline{\beta}^{u}{}'\right]$$

RTC'          …

Theoretical results
allow the computation
of properties of
outgoing streams ☞

© p. marwedel,
informatik 12, 2012

technische universität
dortmund

fakultät für
informatik

- 8 -

# RTC/MPA:
# Transformation of arrival and service curves

Resulting arrival curves:

$$\overline{\alpha}^u{}' = \min\left(\left\lfloor\left(\overline{\alpha}^u \underline{\otimes} \overline{\beta}^u\right)\overline{\oplus}\overline{\beta}^l\right\rfloor, \overline{\beta}^u\right)$$

$$\overline{\alpha}^l{}' = \min\left(\left\lfloor\left(\overline{\alpha}^l \overline{\oplus}\overline{\beta}^u\right)\underline{\otimes}\overline{\beta}^l\right\rfloor, \overline{\beta}^l\right)$$

Remaining service curves:

$$\overline{\beta}^u{}' = \left(\overline{\beta}^u - \overline{\alpha}^l\right)\underline{\oplus}0$$

$$\overline{\beta}^l{}' = \left(\overline{\beta}^l - \overline{\alpha}^u\right)\overline{\otimes}0$$

Where:

$$\left(f\underline{\otimes}g\right)(t) = \inf_{0\leq u\leq t}\left\{f(t-u) + g(u)\right\} \qquad \left(f\overline{\otimes}g\right)(t) = \sup_{0\leq u\leq t}\left\{f(t-u) + g(u)\right\}$$

$$\left(f\underline{\oplus}g\right)(t) = \inf_{u\geq 0}\left\{f(t+u) - g(u)\right\} \qquad \left(f\overline{\oplus}g\right)(t) = \sup_{u\geq 0}\left\{f(t+u) - g(u)\right\}$$

# RTC/MPA: Remarks

- Details of the proofs can be found in relevant references

- Results also include bounds on buffer sizes and on maximum latency.

- Theory has been extended into various directions, e.g. for computing remaining battery capacities

# Application: In-Car Navigation System

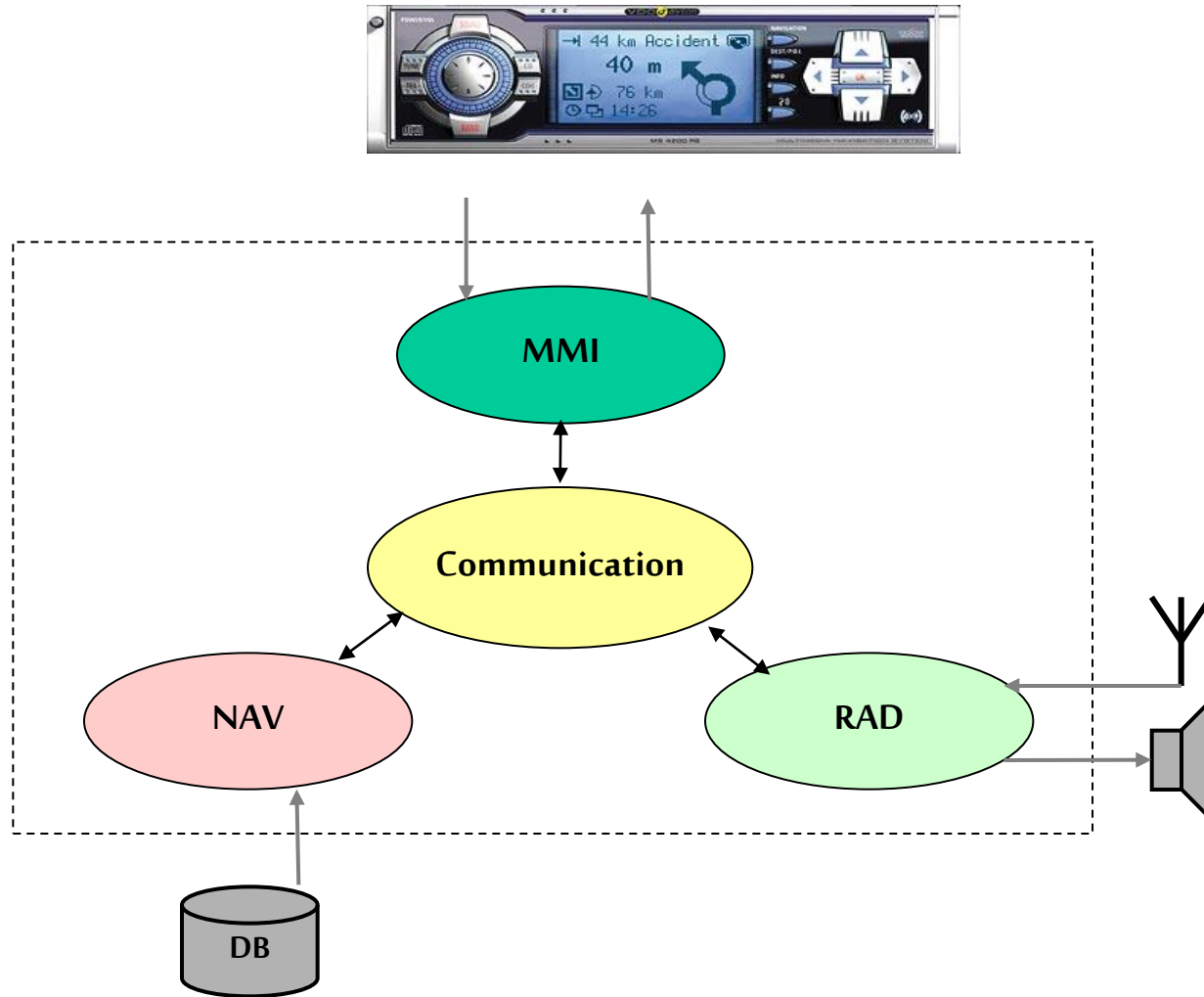Car radio with navigation system
User interface needs to be responsive
Traffic messages (TMC) must be processed in a timely way
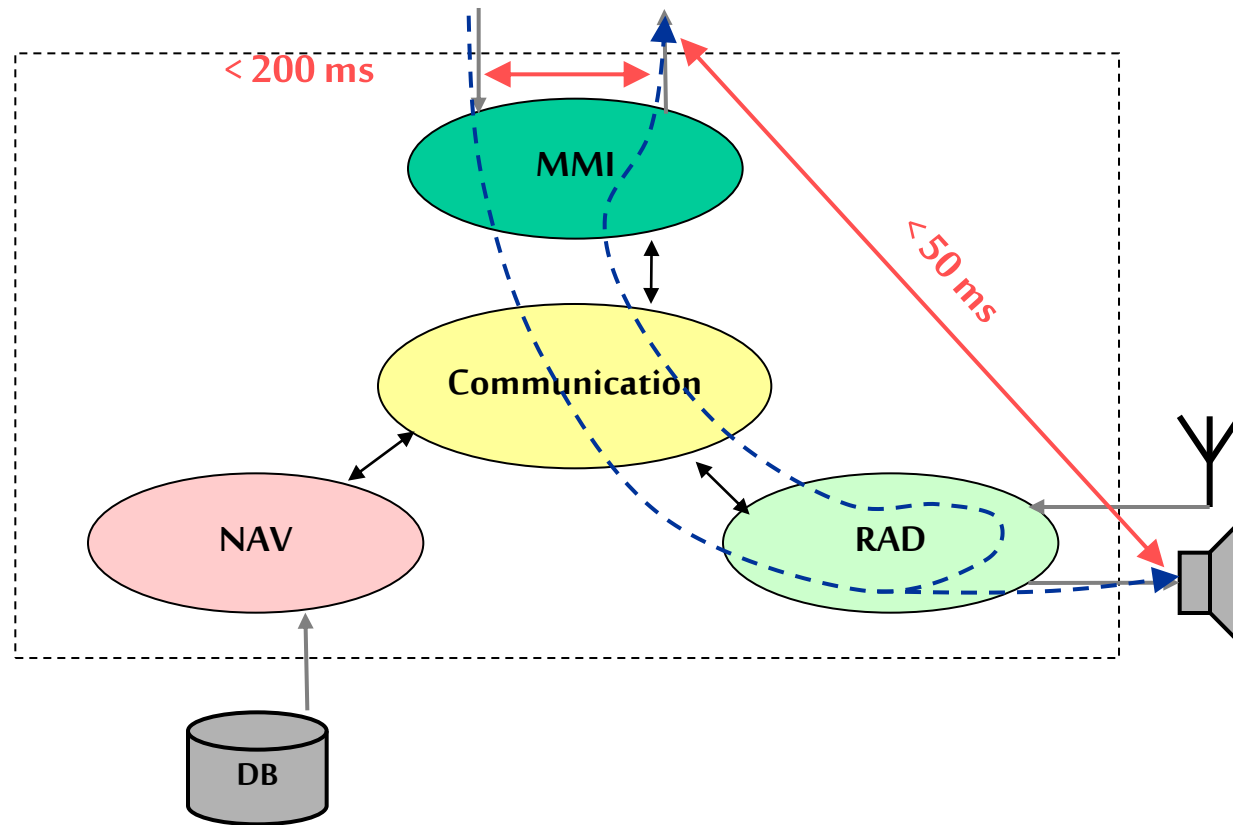Several applications may execute concurrently



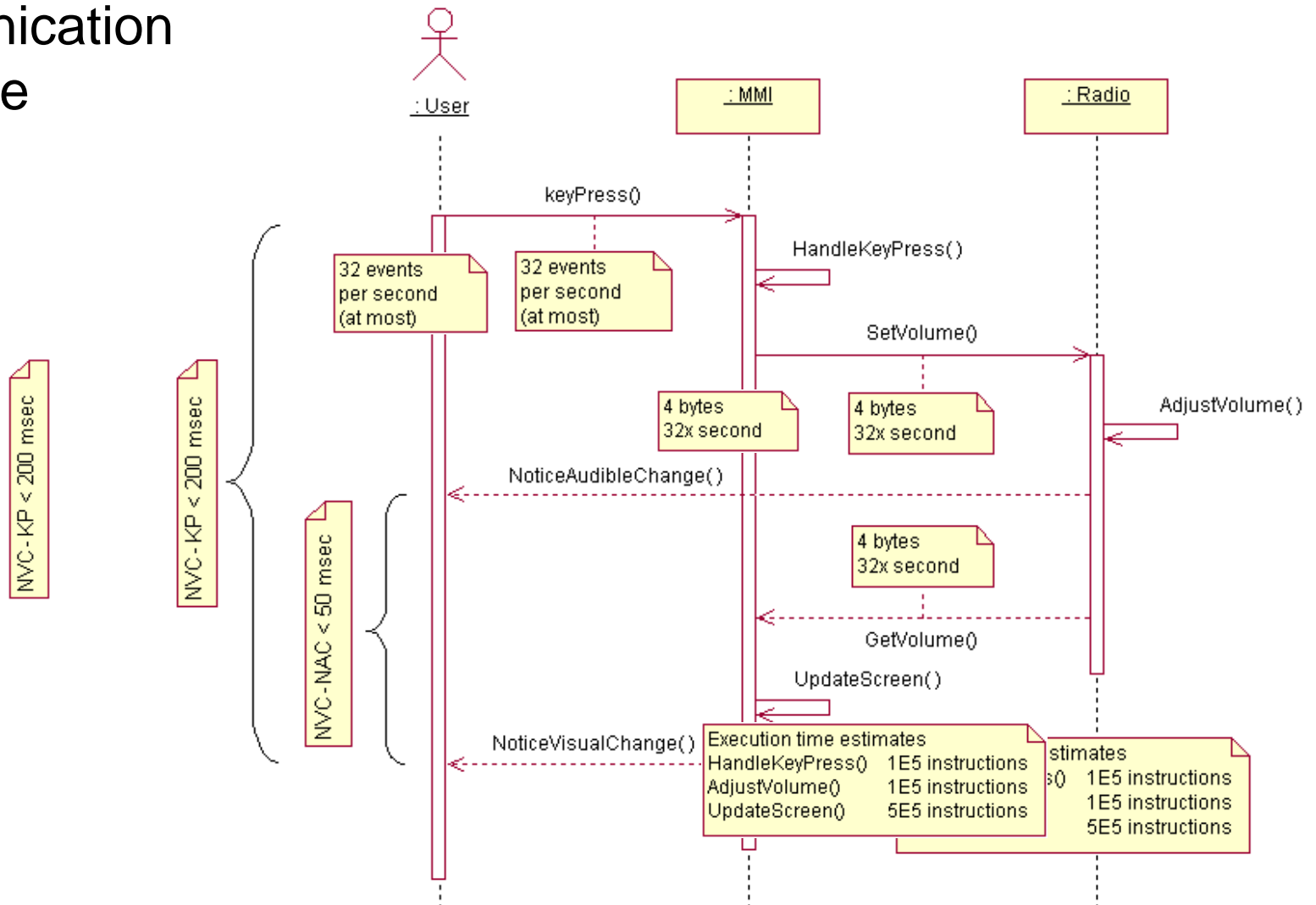© Thiele, ETHZ
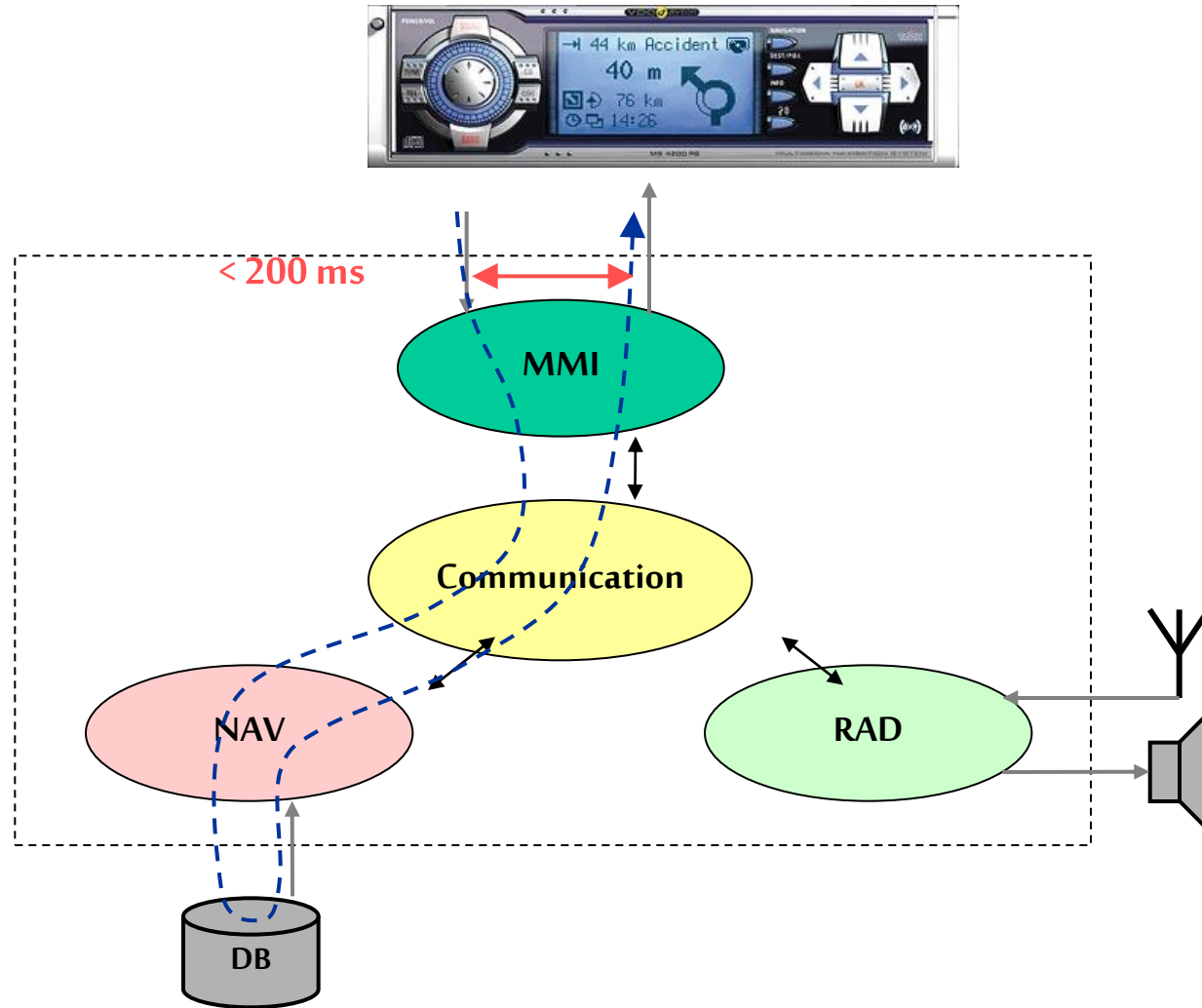
# System Overview

# Use case 1: Change Audio Volume

# Use case 1: Change Audio Volume
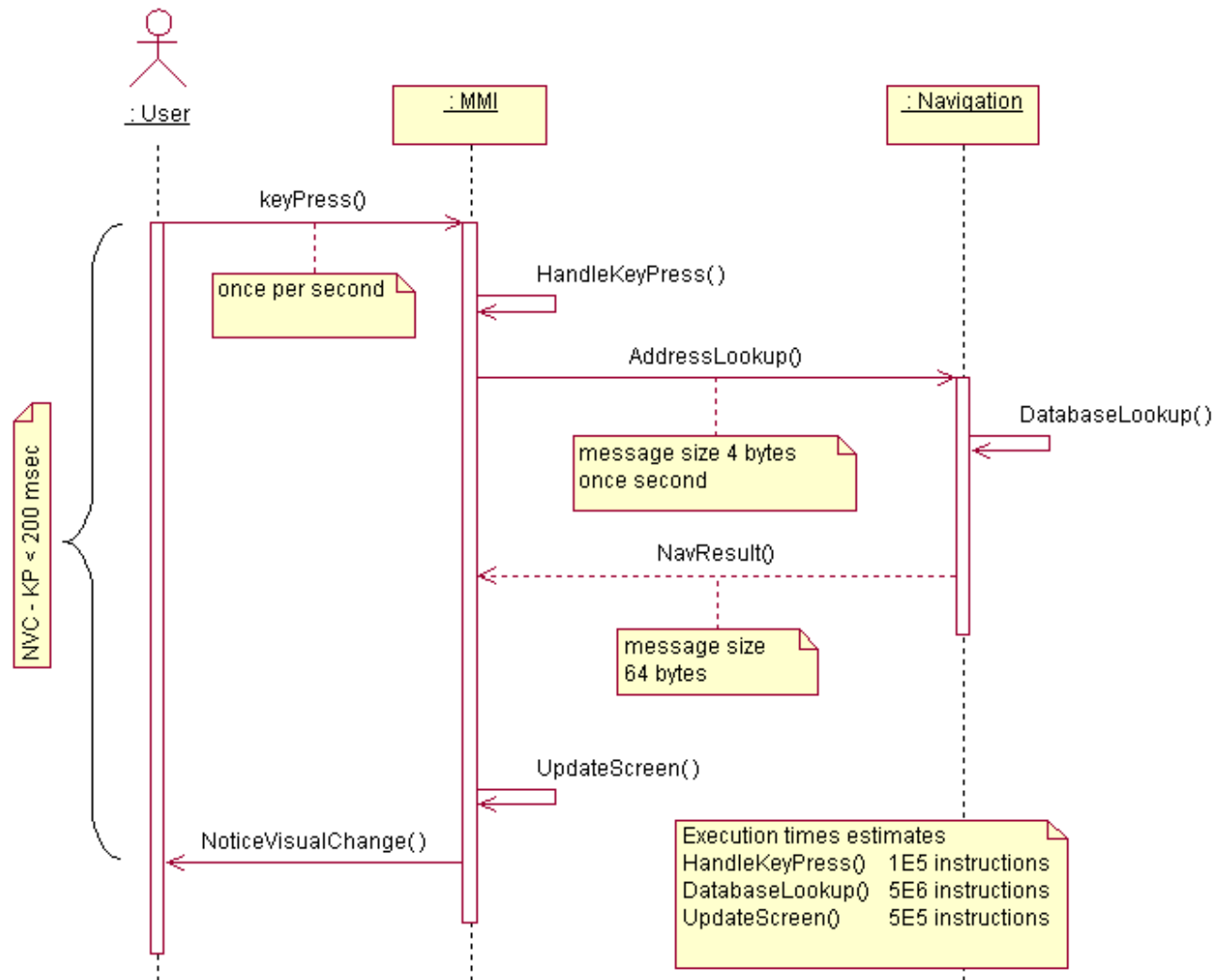
Communication
Resource
Demand

technische universität dortmund

fakultät für informatik
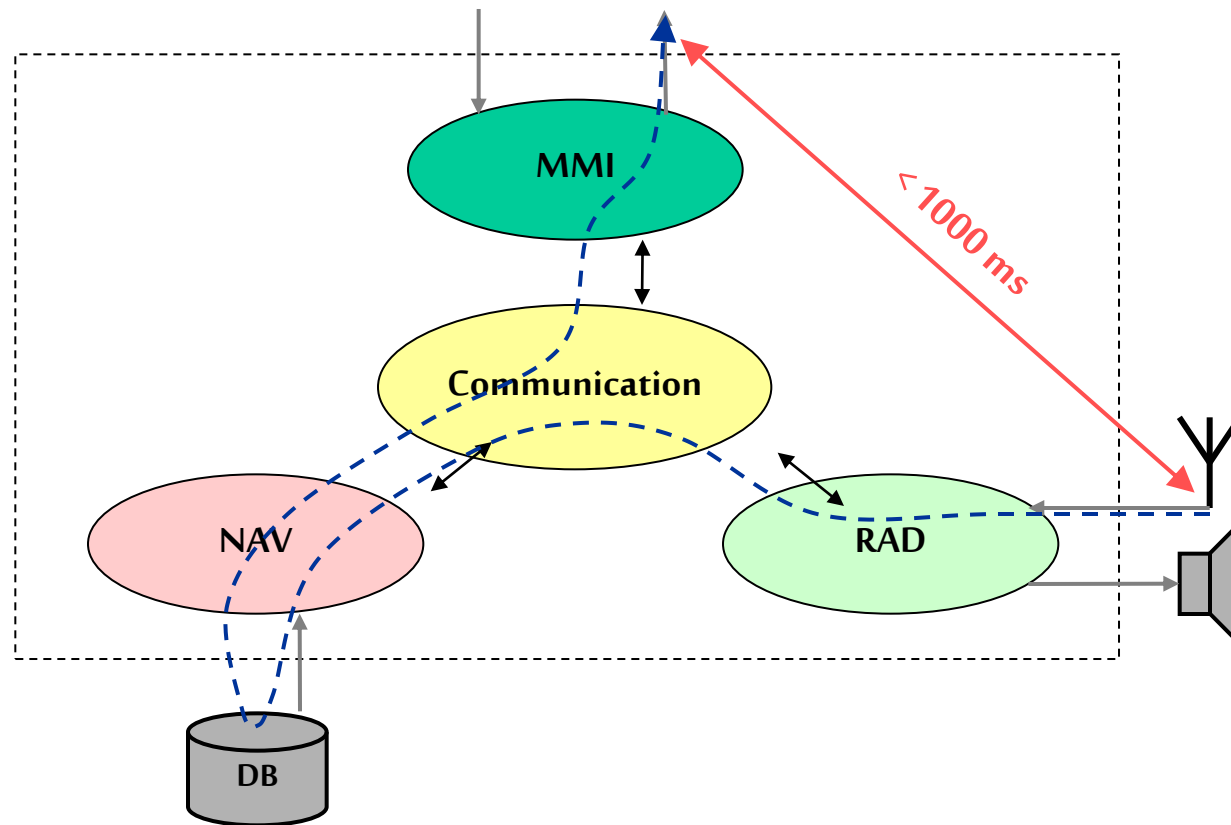
# Use case 2: Lookup Destination Address

# Use case 2: Lookup Destination Address

# Use case 3: Receive TMC Messages

# Use case 3: Receive TMC Messages

# Proposed Architecture Alternatives

technische universität
dortmund

fakultät für
informatik

© p. marwedel,
informatik 12, 2012

© Thiele, ETHZ    - 19 -
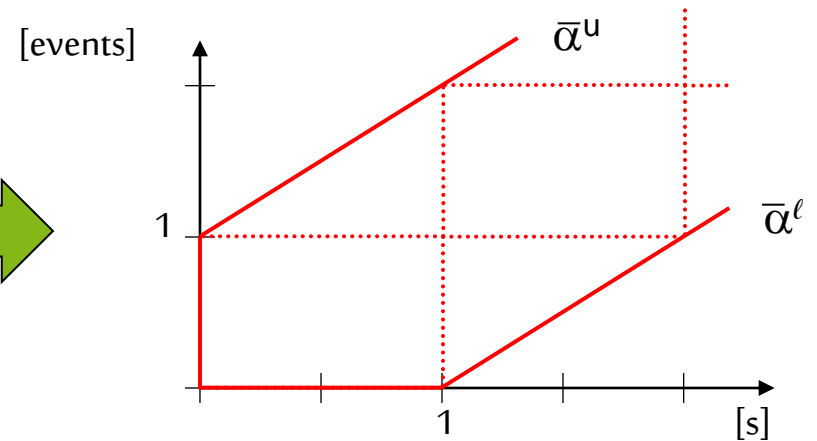
# Step 1: Environment (Event Steams)
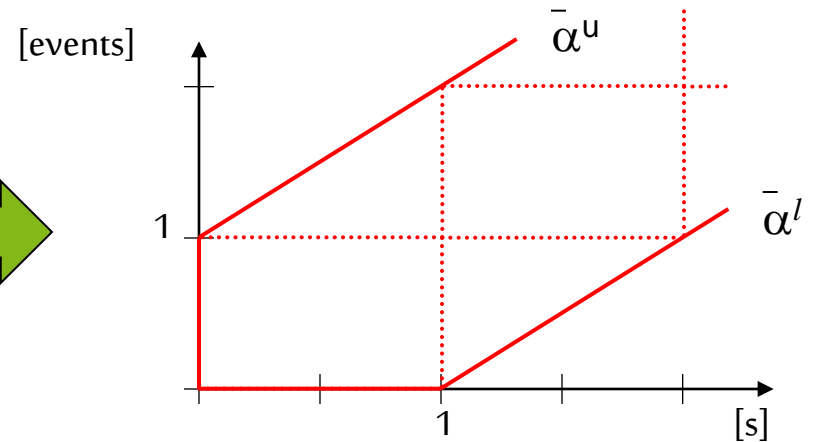
## Event Stream Model

e.g. Address Lookup
(1 event / sec)

# Step 2: Architectural Elements

## Event Stream Model
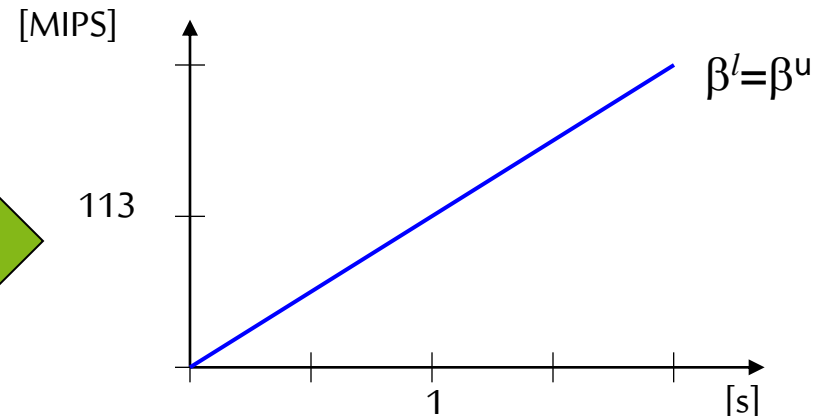
e.g. Address Lookup
(1 event / sec)



## Resource Model

e.g. unloaded RISC CPU
(113 MIPS)

technische universität dortmund

fakultät für informatik
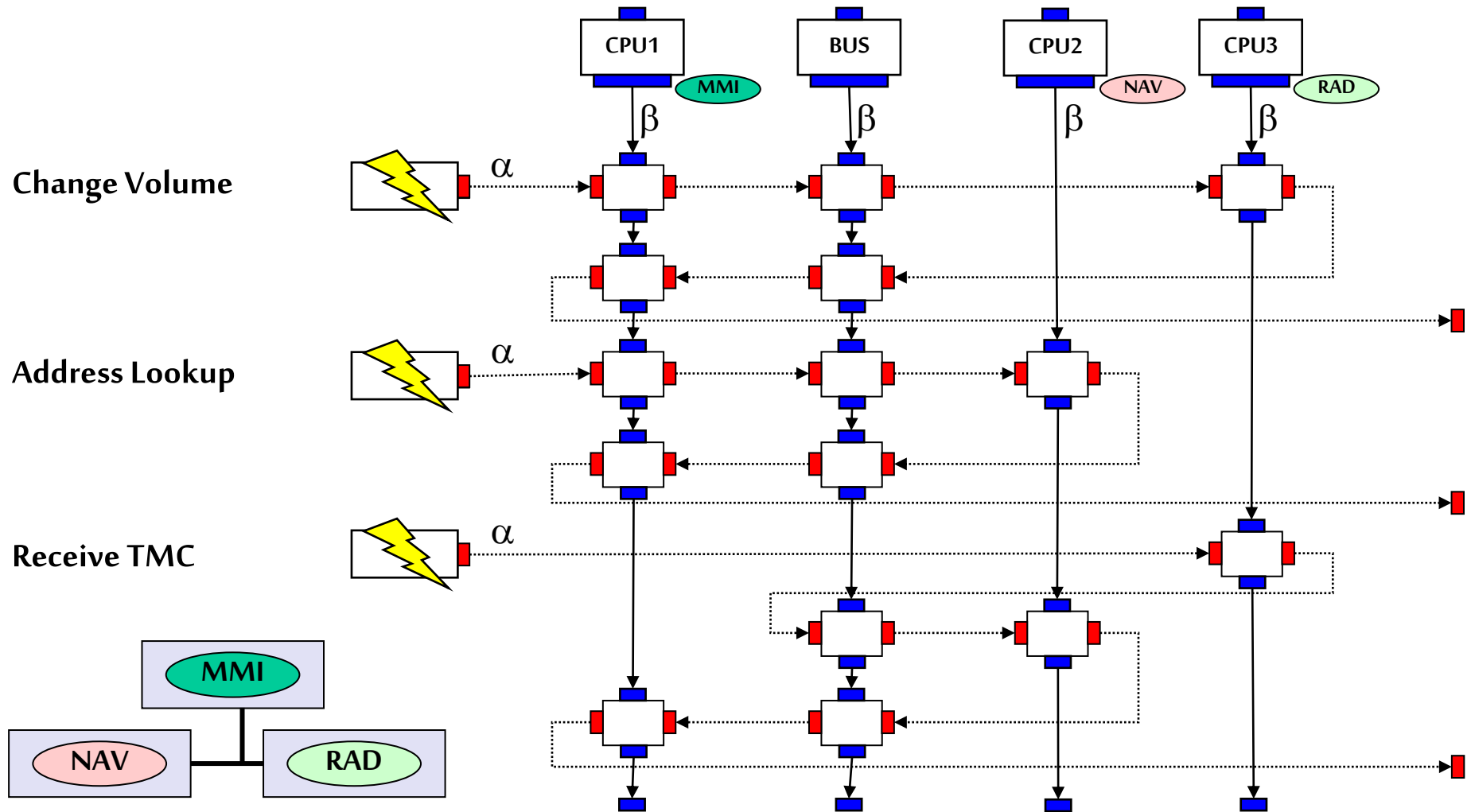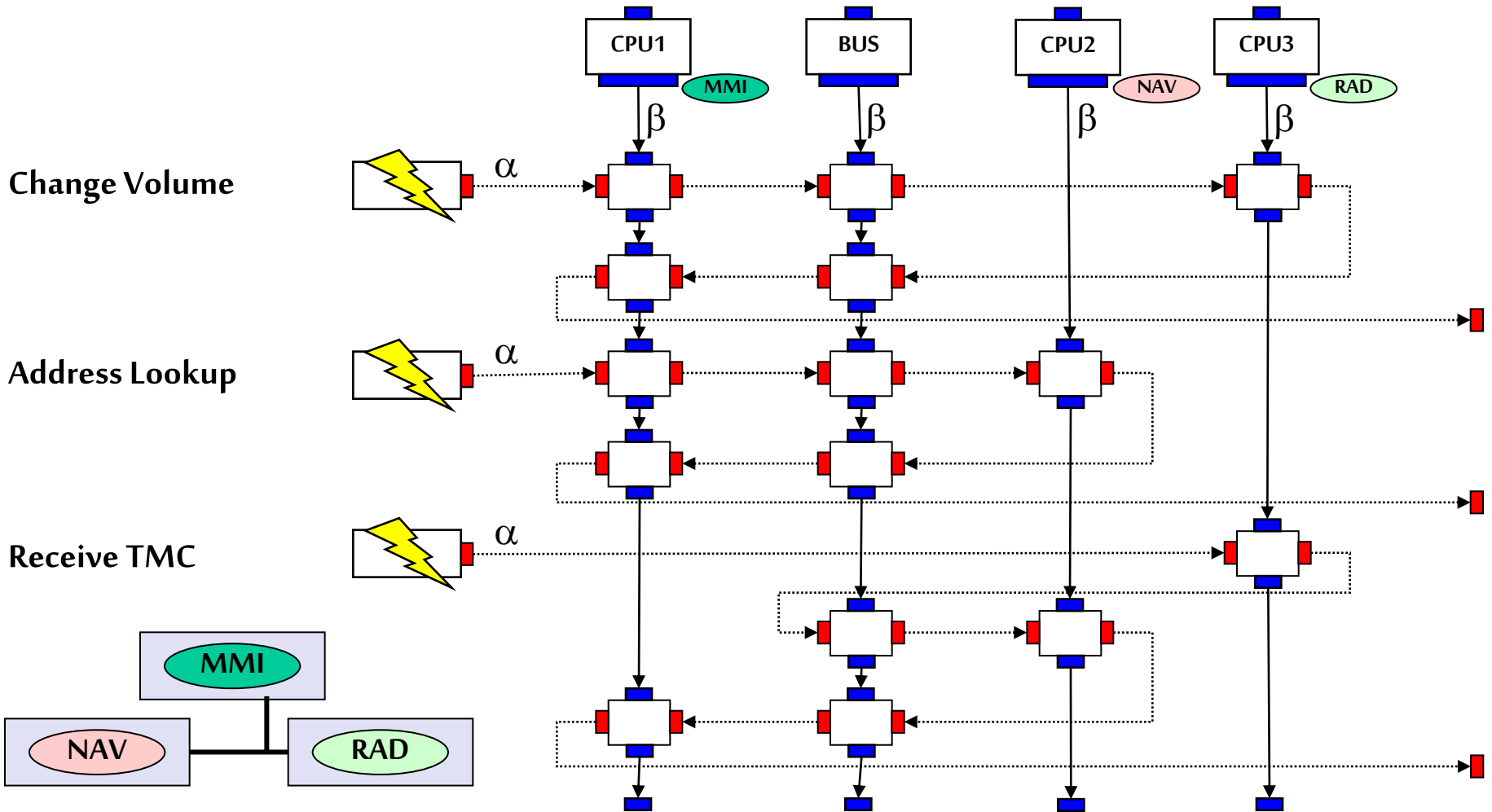
# Step 3: Mapping / Scheduling

Rate Monotonic Scheduling
(Pre-emptive fixed priority scheduling):

- Priority 1: Change Volume    (p=1/32 s)

- Priority 2: Address Lookup    (p=1 s)

- Priority 3: Receive TMC    (p=6 s)

# Step 4: Performance Model

© p. marwedel,
informatik 12, 2012

© Thiele, ETHZ

technische universität dortmund

fakultät für informatik

# Step 5: Analysis

# Analysis – Design Question 1

How do the proposed system architectures compare in respect to end-to-end delays?

technische universität dortmund

fakultät für informatik

# Analysis – Design Question 1

End-to-end delays:

# Analysis – Design Question 2

How robust is architecture A?

Where is the bottleneck of this architecture?

# Analysis – Design Question 2

TMC delay vs. MMI processor speed:

# Conclusions

- Easy to construct models (~ half day)

- Evaluation speed is fast and linear to model complexity (~ 1s per evaluation)

- Needs little information to construct early models (Fits early design cycle very well)

- Even though involved mathematics is very complex, the method is easy to use (Language of engineers)

# How to evaluate designs according to multiple criteria?

Many different criteria are relevant for evaluating designs:

- Average & worst case delay
- power/energy consumption
- thermal behavior
- reliability, safety, security
- cost, size
- weight
- EMC characteristics
- radiation hardness, environmental friendliness, ..

How to compare different designs?
(Some designs are "better" than others)

# Average vs. worst case energy consumption

- The **average energy consumption** $E_{AV}$ is based on the consumption for selected sets of input data (which?)

- The **worst case energy consumption** $E_{WC}$ is a safe upper bound on the energy consumption

- The worst case usage pattern for the battery is $\neq$ from the worst case for the overall energy consumption

- In general, the pattern for worst case energy consumption is $\neq$ from the worst case thermal pattern

# Evaluation of energy consumption: Challenges

- Energy consumption hardly predictable from the source code, due to difficult to predict impact of compiler & linker

- Small variations of the code can lead to large variations of the energy consumption

  - ex. notorious examples

  - Example: shifting code in memory by one byte

- Energy consumption must be predicted from executable code (like the WCET)

- The energy consumption might even depend very much on which instance of the hardware is used

# Energy models

- Measurements: (potentially) precise, fixed architecture

- Models: flexible architecture, less precise

- Combined models

In general, accuracy remains a problem

- Currents difficult to measure

# Steinke's model

E.g.: ATMEL board with ARM7TDMI and ext. SRAM

$$E_{total} = E_{cpu\_instr} + E_{cpu\_data} + E_{mem\_instr} + E_{mem\_data}$$



$V_{DD}$ mA mA

ARM7

| ALU | Register File | DAddr | | Data Memory |
| Barrel Shifter | Reg Value | Data | IAddr | |
| Multi-plier | Imm Reg# Instr Opcode | Instr | | Instruction Memory |

Instr. Decoder & Control Logic

# Example: Instruction dependent costs in the CPU

Cost for a sequence of $m$ instructions

$$E_{cpu\_instr} = \sum MinCostCPU(\textbf{\textit{Opcode}}_i) + FUCost(\textbf{\textit{Instr}}_{i-1}, \textbf{\textit{Instr}}_i) +$$

$$\alpha_1 * \sum w(\textbf{\textit{Imm}}_{i,j}) \quad + \beta_1 * \sum h(\textbf{\textit{Imm}}_{i-1,j}, \textbf{\textit{Imm}}_{i,j}) +$$

$$\alpha_2 * \sum w(\textbf{\textit{Reg}}_{i,k}) \quad + \beta_2 * \sum h(\textbf{\textit{Reg}}_{i-1,k}, \textbf{\textit{Reg}}_{i,k}) +$$

$$\alpha_3 * \sum w(\textbf{\textit{RegVal}}_{i,k}) + \beta_3 * \sum h(\textbf{\textit{RegVal}}_{i-1,k}, \textbf{\textit{RegVal}}_{i,k}) +$$

$$\alpha_4 * \sum w(\textbf{\textit{IAddr}}_i) \quad + \beta_4 * \sum h(\textbf{\textit{IAddr}}_{i-1}, \textbf{\textit{IAddr}}_i)$$

$w$:              number of ones;

$h$:              Hamming distance;

$FUCost$: cost of switching functional units;

$\alpha, \beta$:     determined through experiments.

# Hamming Distance between adjacent addresses is playing major role



h-costs, address bus, CPU + memory current

# Energy-efficient execution on graphics processor (GPU)



current clamp

| Energy per frame CPU | 3.26 J | 5.84 J | 10.52 J | Reduced to |
|---|---|---|---|---|
| Energy per frame GPU | 0.93 J | 1.56 J | 2.76 J | avg 27% |

C. Timm, A. Gelenberg, P. Marwedel, F. Weichert: Energy Considerations within the Integration of General Purpose GPUs in Embedded Systems. Intern. Conf. on Advances in Distributed and Parallel Computing, 2010

C. Timm, F. Weichert, P. Marwedel, H. Müller: Design Space Exploration Towards a Realtime and Energy-Aware GPGPU-based Analysis of Biosensor Data. Computer Science - Research and Development, ENA-HPC, 2011

# Measurements also used for SFB-B2 project
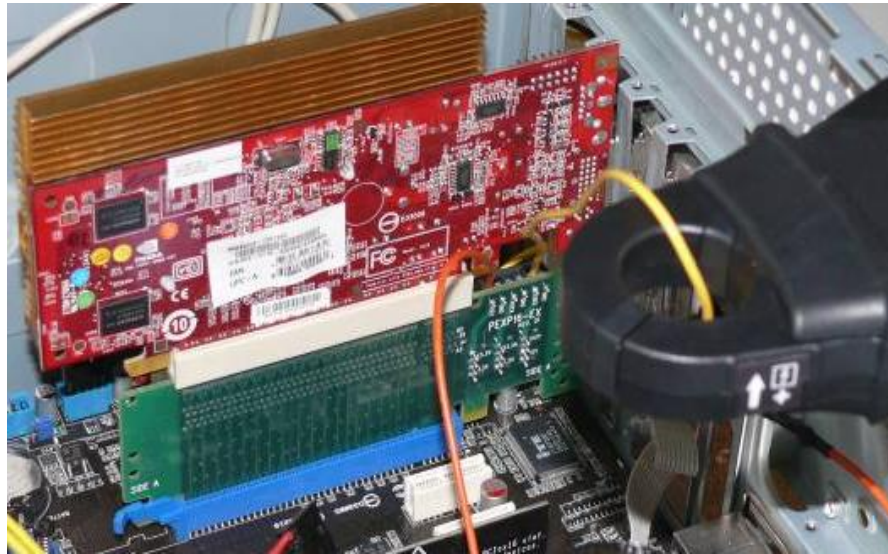
Used to estimate energy consumption of graphics card

Constantin Timm, Andrej Gelenberg, Peter Marwedel and Frank Weichert. Energy Considerations within the Integration of General Purpose GPUs in Embedded Systems. In Proceedings of the International Conference on Advances in Distributed and Parallel Computing, November 2010
Constantin Timm, Frank Weichert, Peter Marwedel and Heinrich Müller. Design Space Exploration Towards a Realtime and Energy-Aware GPGPU-based Analysis of Biosensor Data. Computer Science - Research and Development, Special Issue "International Conference on Energy-Aware High Performance Computing (ENA-HPC)", September 2011

# CACTI model



Cache model used

Comparison with SPICE

# Energy consumption of memories

Example: CACTI / high performance Scratchpad (SRAM):



**Energy (nJ) - read**

16 bit read; size in bytes; 65 nm technology

# Energy consumption of memories (2)

Example CACTI: Scratchpad (SRAM) vs. DRAM (DDR2):



16 bit read; size in bytes;
65 nm for SRAM, 80 nm for DRAM

# DRAM power

Complex DRAM models:

- http://www.micron.com/products/support/power-calc

- T. Vogelsang: Understanding the Energy Consumption of Dynamic Random Access Memories, Proceedings of the 2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture, pp. 363—374, http://dx.doi.org/10.1109/MICRO.2010.42

# Steinke's "combined" model

- Measured values for the processor

- Model-based values for memories
  (validated against existing measurements)

# Examples of energy models

- Measurements:
  - Tiwari (1994): Energy consumption within processors
  - Russell, Jacome (1998): Measurements for 2 fixed configurations
  - Simunic (1999): Values from data sheets. Not very precise.
  - Timm: measurements for graphics card

- Models:
  - CACTI [Jouppi, 1996]: Predicted energy consumption of caches
  - Wattch [Brooks, 2000]: Power estimation at the architectural level, without circuit or layout, known to be imprecise

- Combined models
  - Steinke et al., TU Dortmund (2001): mixed model

# Worst case energy consumption via worst case computing time?



- Computing the $E_{WC}$ using WCET estimations $WCET_{EST}$

$$E_{WC} = \int_{0}^{WCET_{EST}} P(t)\,dt$$

- Tight bounds if $P(t)$ has small variations & $WCET_{EST}$ is tight

- Little value if $P(t)$ varies too much.

# Battery models

- (Chemical) & physical models
  e.g. concentrated solution theory, partial differential eq.s
  many, frequently unknown parameters (50+); $xy$ hours simulation time

- Empirical models
  Simple equations, inaccurate
    - Peukert's law: lifetime= $C/I^\alpha$, with empirical $\alpha$
    - Weibull fit

- Abstract models
  ➡ - Electrical circuit models
    - Discrete time model (e.g. in VHDL)
    - Stochastic models (e.g. Markov processes)

- Mixed models
  e.g. electrical models with physical explanation

*frequently with fitting*

# Model proposed by Chen and Rincón-Mora



Source: M. Chen, G. A. Rincón-Mora: Accurate Electrical Battery Model Capable of Predicting Runtime and *I-V* Performance, *IEEE Trans. on Energy Conversion*, 2006, pp. 504

- Full charge capacitor: $C_{Capacity}$ = 3600 $\cdot$ *Capacity* $\cdot$ $f_1$(cycle) $\cdot$ $f_2$(Temp)
- Self-discharge resistor: $R_{Self\text{-}Discharge}$ (might depend on parameters)
- Current dependency of $V_{Batt}$: modeled by $R_{series}+R_{Transient\_S}+R_{Transient\_L}$
- $I_{Batt}$ charges and discharges $C_{Capacity}$
- Voltage controlled voltage source $V_{0C}$ captures nonlinear dependency between the state of charge and $V_{0C}$ (measurement can take days)
- $R_{Series}$: models immediate voltage drop at load change

# Battery capacity sufficient?

Question can be solved with adapted real-time calculus



Lipskoch, H., Albers, K. and Slomka, F.: Battery discharge aware energy feasibility analysis, Proceedings of the 4th international Conference on Hardware/Software Codesign and System Synthesis, CODES+ISSS '06, pp. 22-27, 2006.

# Energy models for communication:
# An Energy Model for Mobile WiMAX Devices

- How does the application data rate influence the energy efficiency?

- What is the impact of very small amounts of data on the efficiency?

- Relationship between submitted power and consumed energy?

# Traffic Dependent Energy Consumption



Channel Quality Reports

Idle State

Data Transmission

- Average Power in Idle State: 880 mW

- Channel Quality Reports every 300 ms lead to increased average power of 930 mW

- Transmission is costing significantly increased energy consumption

- Reception is not increasing the power consumption compared with idle state

# Modeling the Impact of $\underline{T_x}$-Power Variation



- Energy per Bit is constant for low $T_x$ power (below -25 dBm)
- For higher $T_x$ power, the consumed energy can be approximated be 2nd degree polynomial
- Significant energy savings can be achieved by using as low power as possible
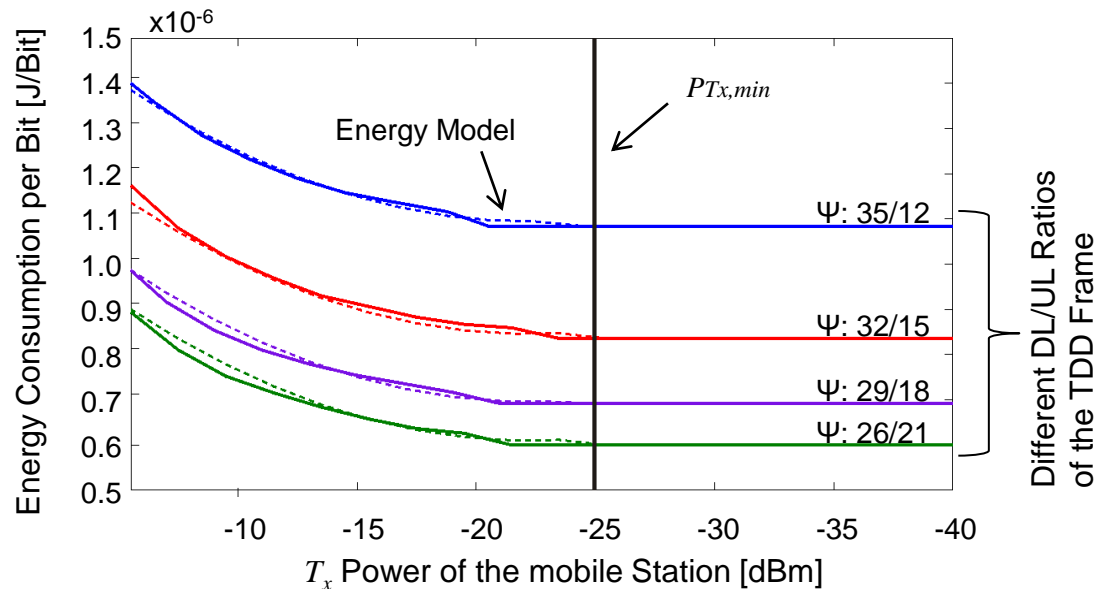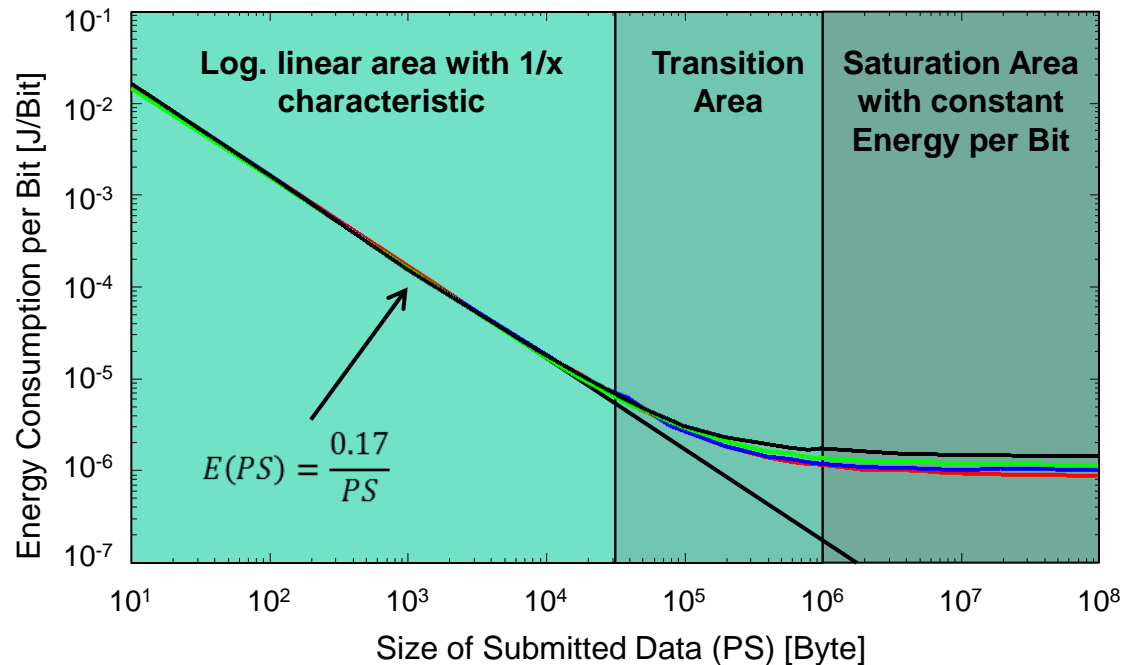
$$E(P_{Tx}) = \begin{cases} C(\Psi) & , P_{Tx} < P_{Tx,min} \\ \alpha \cdot P_{Tx}^2 + \beta \cdot P_{Tx} + \gamma + C(\Psi) & , P_{Tx} > P_{Tx,min} \end{cases}$$

| α | β | γ | C(Ψ=35/12) | C(Ψ=32/15) | C(Ψ=29/18) | C(Ψ=26/21) | $P_{Tx,min}$ |
|---|---|---|---|---|---|---|---|
| 1.0325e-6 | 4.6e-5 | 5.225e-4 | 1.1084e-3 | 8.3530e-4 | 6.916e-4 | 5.994e-4 | -25 dBm |

# Modeling the Impact of Different File Sizes



The Energy Model for different file sizes can be divided into three parts

**Log Linear Area:**
Rapidly decreasing energy consumption per Bit for packet sizes below 20 kByte

**Transition Area:**
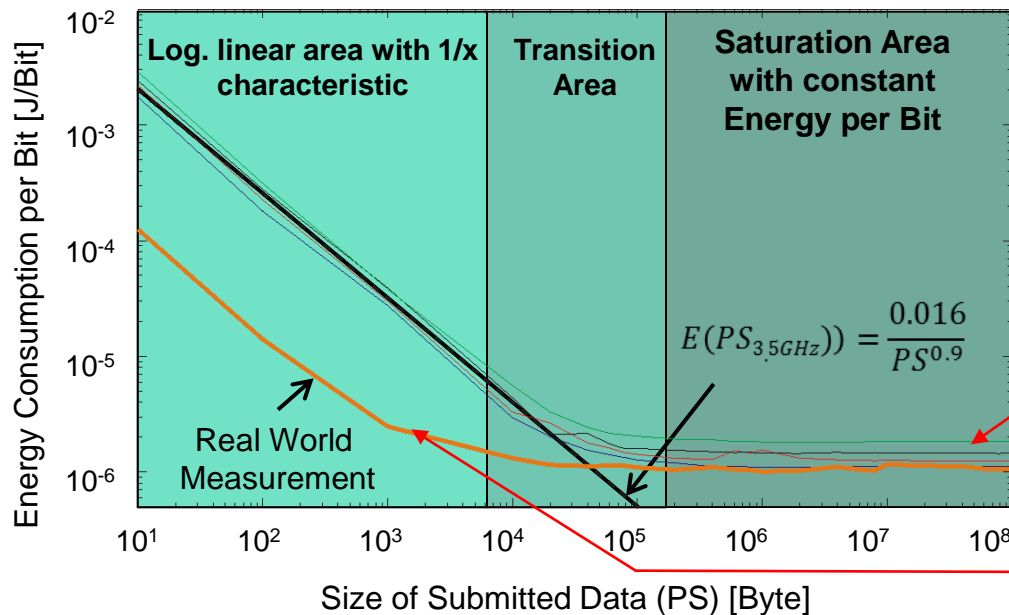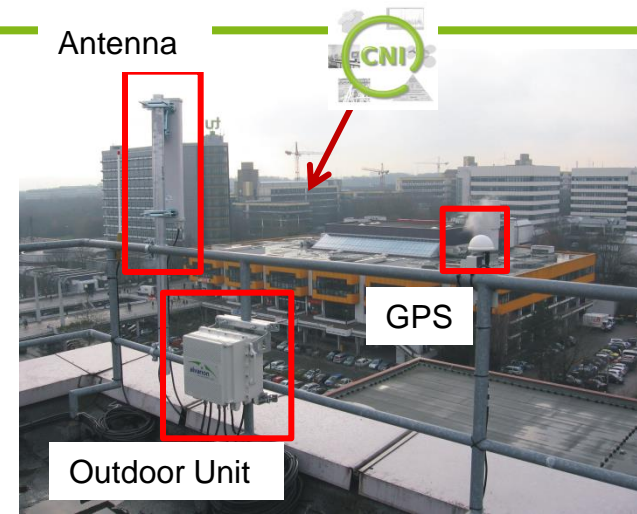Transition to constant energy consumption

**Saturation Area**
For packet sizes above 900 kByte, collecting more data does not make sense from an energy efficiency perspective

Figure labels:
- Log. linear area with 1/x characteristic
- Transition Area
- Saturation Area with constant Energy per Bit
- Energy Consumption per Bit [J/Bit]
- Size of Submitted Data (PS) [Byte]

$$E(PS) = \frac{0.17}{PS}$$

# Validation for Different Devices


Antenna

GPS

Outdoor Unit



Energy Consumption per Bit [J/Bit]

| Log. linear area with 1/x characteristic | Transition Area | Saturation Area with constant Energy per Bit |

$$E(PS_{3.5GHz})) = \frac{0.016}{PS^{0.9}}$$

Real World Measurement
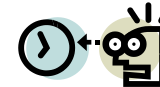
Size of Submitted Data (PS) [Byte]

- The overall model is valid for different devices, and different frequency bands
- Chipset specific offsets have to be applied for the log linear area
- Lab equipment consumes some energy after bits have been sent.
- Real (more modern) BS instead of lab equipment is covered by the model (offset has to be applied for low PS)

# How to evaluate designs according to multiple criteria?

Many different criteria are relevant for evaluating designs:

- Average & worst case delay
- power/energy consumption
- thermal behavior
- reliability, safety, security
- cost, size
- weight
- EMC characteristics
- radiation hardness, environmental friendliness, ..

How to compare different designs?
(Some designs are "better" than others)

# Thermal models

Thermal models becoming increasingly important

- since temperatures become more relevant due to increased performance, and

- since temperatures affect

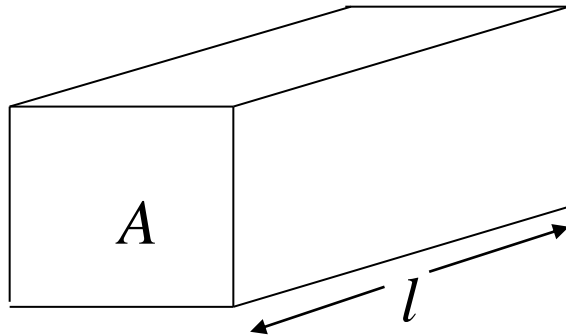  - usability and

  - dependability.

# Thermal conductivity

$$P_{th} = \kappa \frac{\Delta T \cdot A}{l} \qquad (1)$$



Where

| | |
|---|---|
| $P_{th}$ | : thermal power transferred |
| $\kappa$ | : thermal conductivity |
| $\Delta T$ | : temperature difference |
| $A$ | : area |
| $l$ | : length |

**Thermal conductivity** $\kappa$ reflects the amount of thermal energy per unit of time transferred through a plate made of some material of area $A$ and thickness $l$ when the temperatures at the opposite sides differ by one temperature unit (e.g. Kelvin)

# Examples of thermal conductivity

| Material | Thermal conductivity [W/(m K)] |
|---|---|
| Copper | 240-401 |
| Aluminum (95.5%) | 236 |
| Silicon | 148 |
| Wood (perpendicular to fibre) | 0.09-0.19 |
| Concrete | 0.08-0.25 |
| Air (21% oxygen) | 0.0262 |

http://de.wikipedia.org/wiki/Wärmeleitfähigkeit

# Thermal conductance & resistance

- **Thermal conductance** = amount of thermal energy which passes through a plate per unit of time if the temperatures at the two ends of the plate differ by one unit of temperature (e.g. Kelvin).

- The reciprocal of thermal conductance is called **thermal resistance** $R_{th}$.

$$P_{th} = \kappa \frac{\Delta T \cdot A}{l} \qquad (1)$$
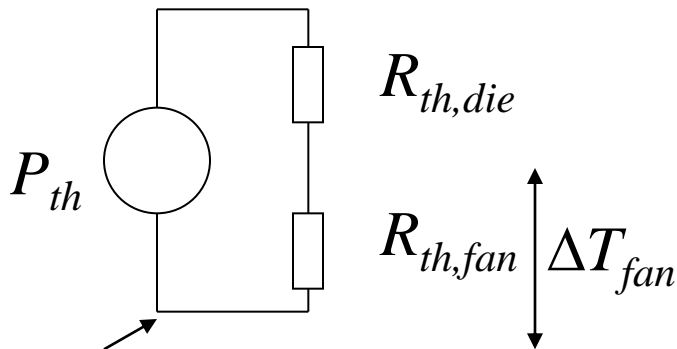
$$\downarrow$$

$$R_{th} = \frac{\Delta T}{P_{th}} = \frac{l}{\kappa \cdot A} \qquad (2)$$

# Equivalent thermal circuits

- Thermal resistances add up like electrical resistances

☞ Thermal modeling mapped to circuit modeling

e.g.: microprocessor:

$$\Delta T = R_{th} \cdot P_{th} \qquad (3)$$

$$R_{th} = R_{th,die} + R_{th,fan} \quad (4)$$

$R_{th,die}$

$P_{th}$

$R_{th,fan} \Big| \Delta T_{fan}$

Ground ≈ Reference temperature

For $R_{th,die}$=0.4 [W/K],
$\quad R_{th,fan}$=0.3 [W/K],
$\quad P_{th} \quad$ =10  [W]:
☞ $\Delta T \quad$ = 7 [K],
$\quad \Delta T_{fan} \quad$ = 3 [K]
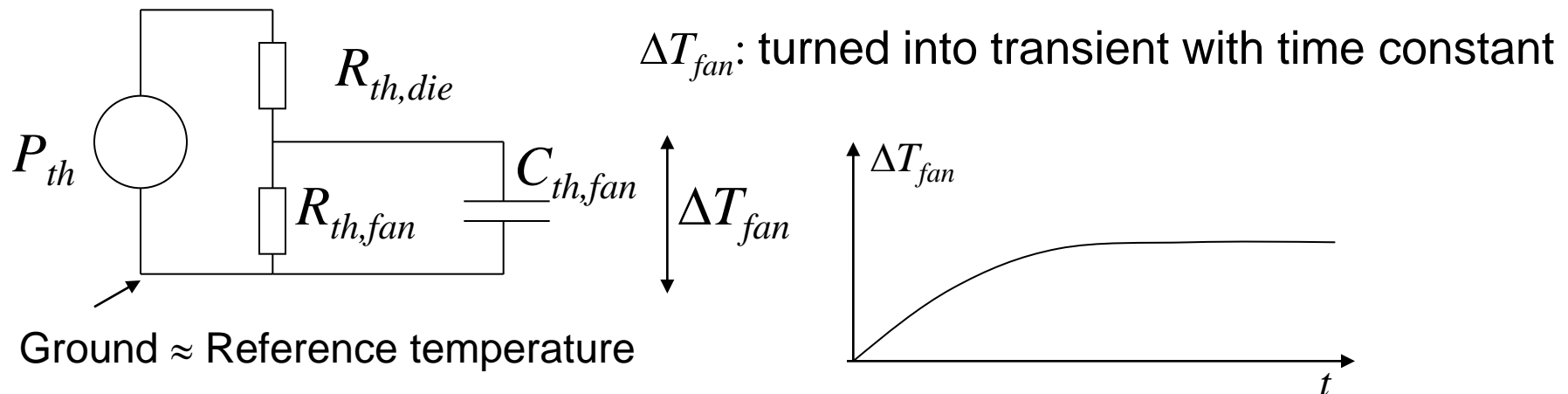
So far, we have just considered the steady state.

# Dynamic thermal properties

In general, transients and thermal capacity to be considered:

$$C_{th} = m \cdot c$$

where $C_{th}$: thermal capacity, $m$: mass, $c$: specific heat

☞ Networks comprising resistances and capacities



$\Delta T_{fan}$: turned into transient with time constant

Ground ≈ Reference temperature

Extra voltage source can make reference temperature explicit

# Equivalences

| Electrical model | | Thermal model | |
|---|---|---|---|
| Current | $I$ | Thermal flow, flow of "power" | $P_{th}=\dot{Q}$ |
| Total charge | $Q = \int I\, dt$ | Thermal energy | $E_{th} = \int P_{th}\, dt$ |
| Resistance | $R$ | Thermal resistance | $R_{th}$ |
| Potential | $\varphi$ | Temperature | $T$ |
| Voltage = potential difference | $U$ | Temperature difference | $\Delta T$ |
| Capacitance | $C$ | Thermal capacitance | $C_{th}$ |
| Ohms law | $U = R\,I$ | $\Delta$Temperature at $R_{th}$ | $\Delta T = R_{th} \cdot P_{th}$ |

# Examples of thermal resistance of P-TO263-7-3

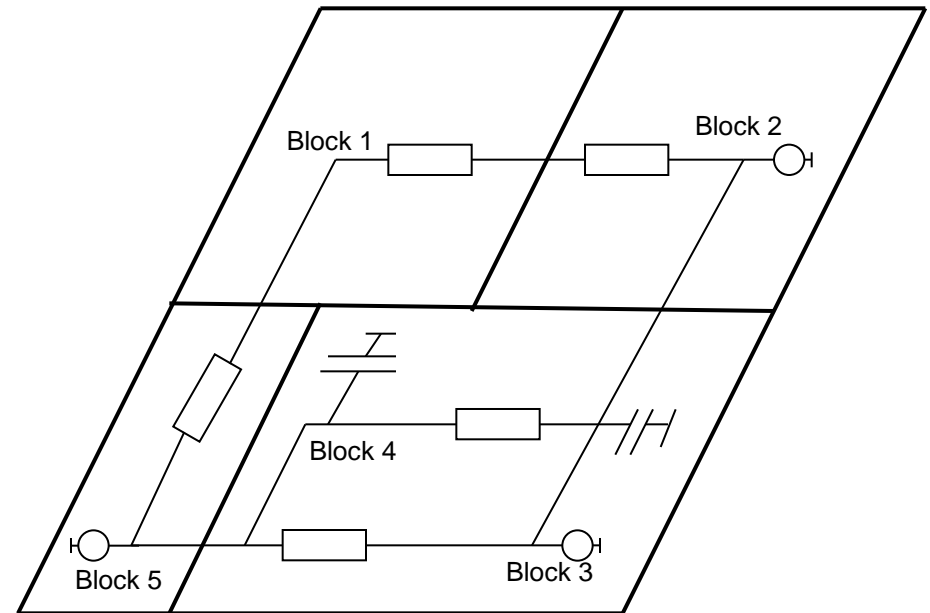| Component | Value & Dimension |
|---|---|
| Thermal resistance of chip | 0.48 [K/W] |
| Thermal time constant of chip | ≈1.5 ms |
| Thermal capacity of chip | ≈3 [mWs/K] |
| Thermal resistance of heat slug | 0.24 [K/W] |
| Thermal capacity of heat slug | 310 [mWs/K] |
| Thermal time constant of heat slug | 70 [ms] |

# Hotspot – A popular thermal simulator for processors

- Localized heating much faster than chip-wide (millisec time scale)

- Chip-wide treatment is inaccurate (neglects hot spots)

- Temperature is sensitive to chip layout (floorplan)

☞ Fine-grained, dynamic model of temperature
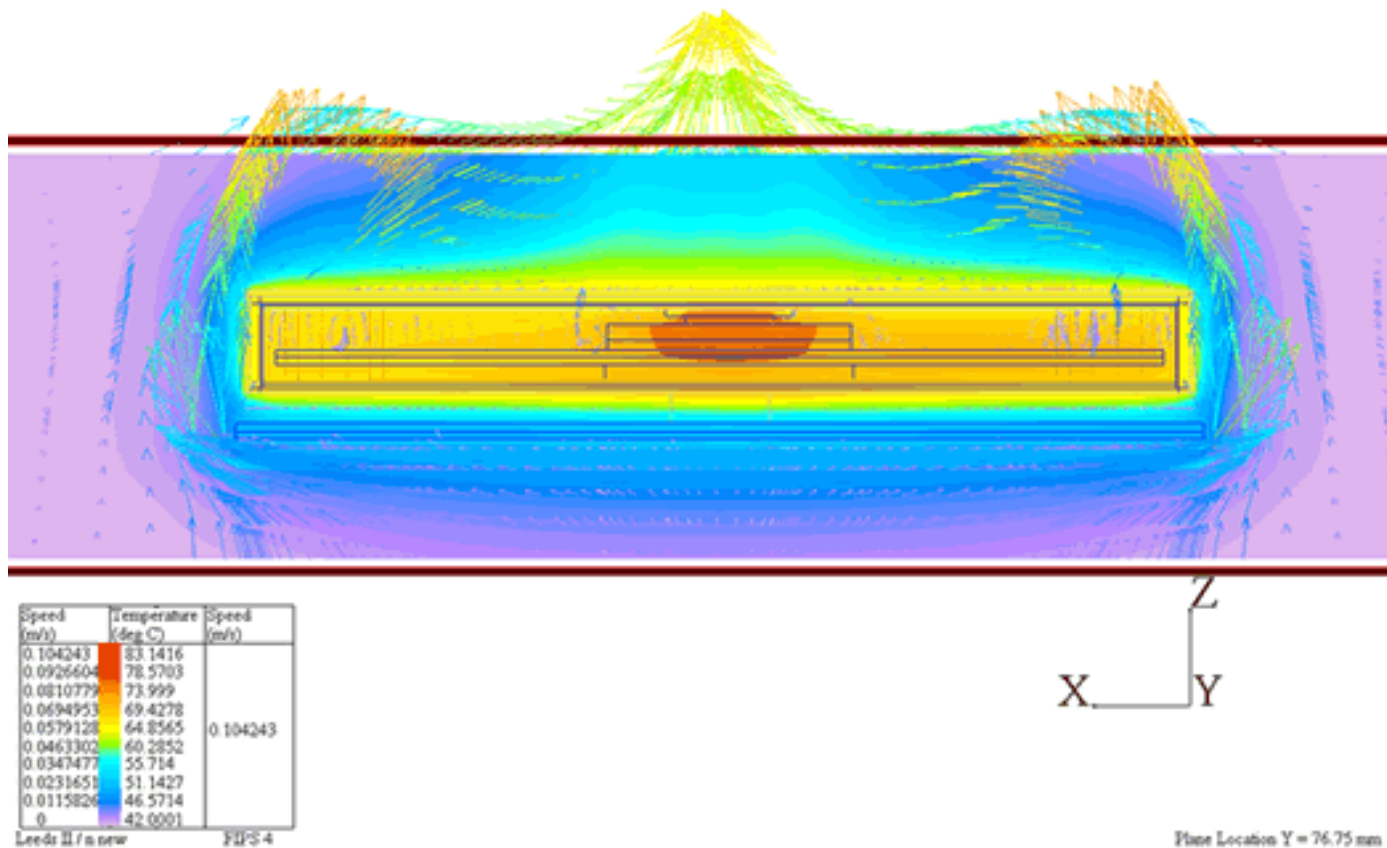
- Authors say: Validated against FEM models



(2D model, 2.5 D exists)

# Results of simulations based on thermal models (1)

Encapsulated cryptographic coprocessor:



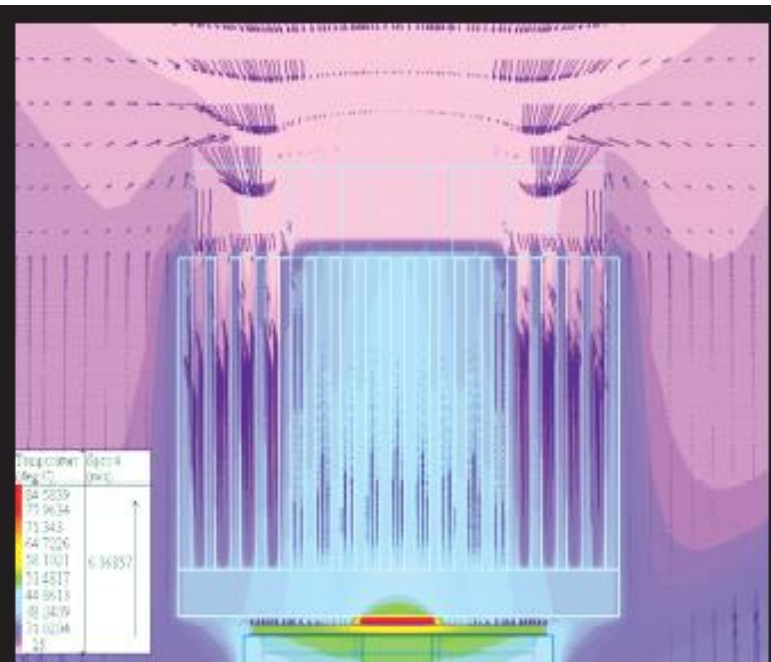Source: http://www.coolingzone.com/Guest/News/ NL_JUN_2001/Campi/Jun_Campi_2001.html

# Results of simulations based on thermal models (2)

Microprocessor



▲ Flomerics image showing the thermal solution with a metal lid.

▲ Flomerics image showing the thermal solution without a metal lid.

# Summary

- Thiele's real-time calculus (RTC)/MPA
  - Using bounds on the number of events in input streams
  - Using bounds on available processing capability
  - Derives bounds on the number of events in output streams
  - Derives bound on remaining processing capability, buffer sizes, …
  - Examples demonstrate design procedure based on RTC
- Energy and power consumption
  - Measurements
  - Models (with calibration)
  - Mixed approaches
  - Energy for computation, storage and communication
- Thermal behavior
  - Mapping to thermal circuit model

# Reserve