# Managing MPSoCs beyond their Thermal Design Power*

Luca Benini

*Università di Bologna & STMicroelectronics*
*Luca.benini@unibo.it*

# Thermal Power Wall

## Transistor count increases exponentially, but…
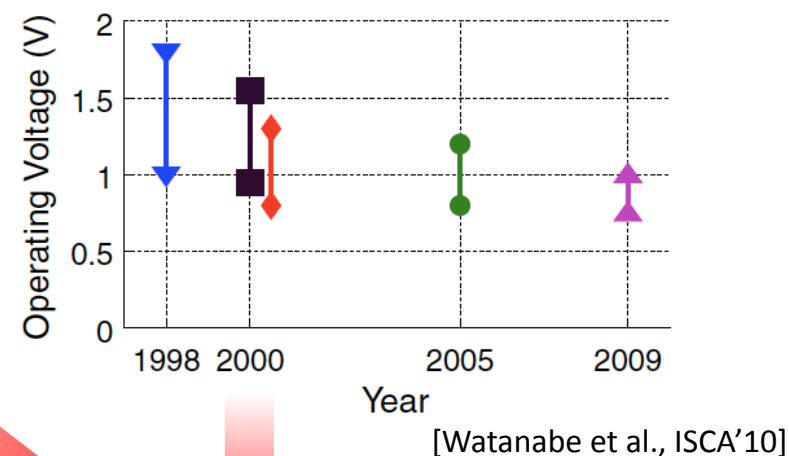
Can no longer power the entire chip (voltages, cooling do not scale)

Traditional HW power-aware techniques are insufficient (e.g., voltage-freq. scaling)



[Watanabe et al., ISCA'10]

**Dark Silicon !!!**

[Hardavellas11]

# Mobile SoCs are cool…right?

## Wrong!

So.. got myself a HTC (T-Mobile) HD2 […].
As i found out the problem is pretty common: damn thing restart itself - **thermally related - the old CPU overheat problem**. By searching the net I found out that it's pretty common with some HTC models. HD2 has it, Desire has it, Nexus One has it, hell even some xperia models have it.. about half of the devices powered by anything from the Snapdragon series could have it.

June 2011 - http://forum.xda-developers.com/showthread.php?t=982454
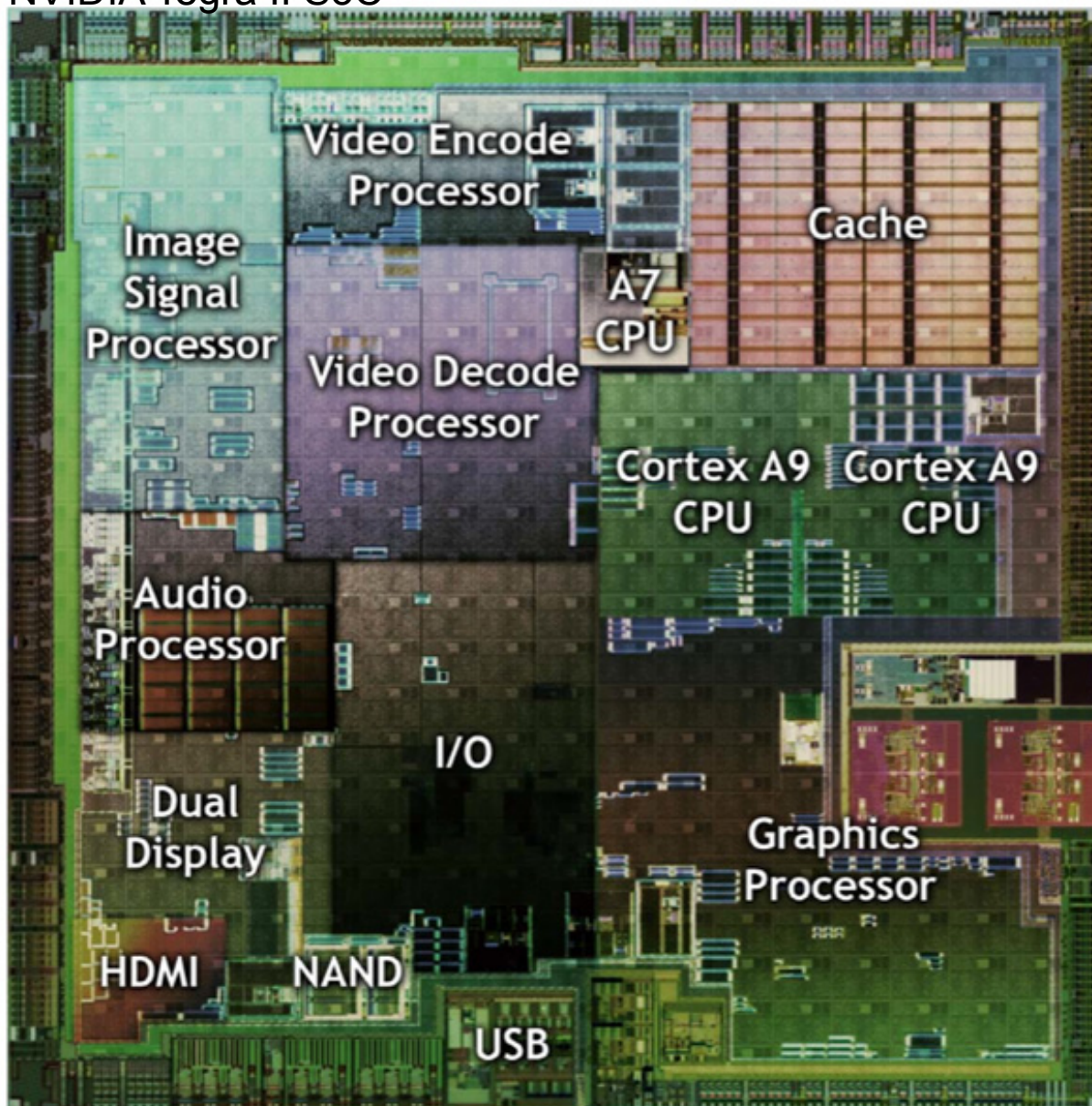
## Why?

ARM has unveiled its next generation "Eagle" (Cortex A15) processor, pitching it at **everything from smartphones to energy-efficient servers**. The A15 will be initially produced at 32nm or 28nm, although ARM claims the roadmap stretches down to 20nm. It will **deliver clock speeds of up to 2.5GHz**.
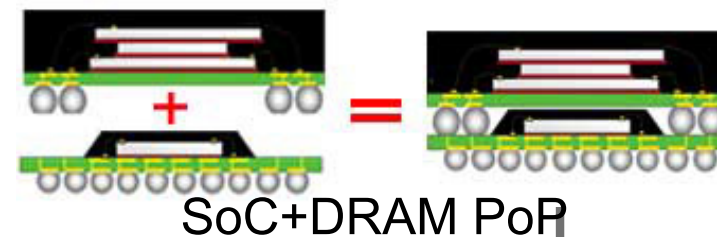
Aug. 2011 - http://www.pcpro.co.uk/news/360994/arm-preys-on-smartphones-and-servers-with-eagle
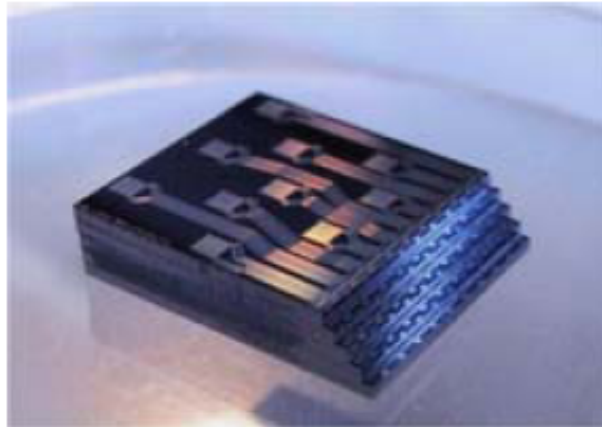
# A 2011 Mobile SoC

NVIDIA Tegra II SoC



- Tegra II
  - TSMC 40nm (LP/G)
  - A9 - 1GHz (G)
  - GPU, etc. - 330MHz (LP)
  - GEForce ULV (8 shaders)
  - 2 separate Vdd rails
  - 1MB L2$
  - 32b LPDDR2 (600MHz DR)

- Tegra II 3D
  - A9 – 1.2GHz
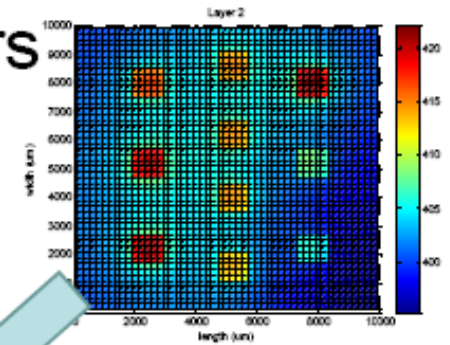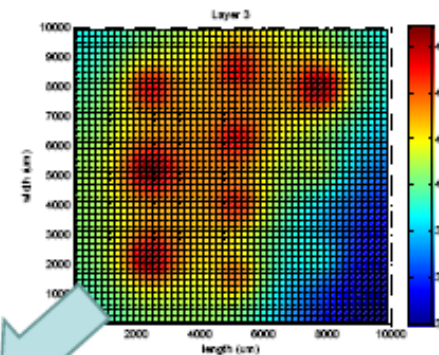  - GPU – 400Mhz



SoC+DRAM PoP

# 3D-SoCs are even worse
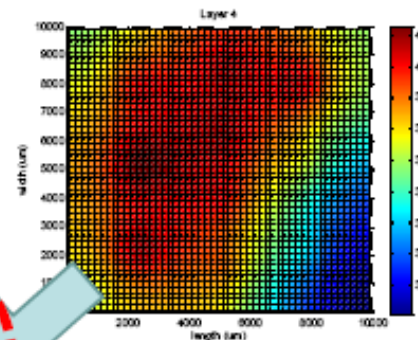
5-tier 3D stack: 10 heat sources and sensors



Inject between 4W – 1.5W

2nd Tier

3rd Tier

4th Tier

5th Tier

Large and non-uniform heat propagation!
(up to 130º C on top tier)

ESL/EPFL 2010

# Rushing to Many-Core

Hardware Trends → 1000+ core system
Software Trends → Concurrency (1000x +)



Massively Parallel
Large Scale SoCs
(1,000s of cores)

Intel 80 Core chip

SCC Tera Scale
Project (48 Cores)

Multi
Processors
cores on

UniProcessors
(single core)

**Research Challenge:**
**Power management for a 1,000 core**
**heterogeneous SoC  →Extreme MIMO!**

# Management Loop: Holistic view



Workload CPU utilization, queue status

MW: Job manager
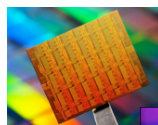OS: Scheduler

SW-introspection: system tracing

SW-control: scheduling, allocation

Resource management policy

SW

HW

HW-Introspection: monitors

HW-control: knobs (Vdd, Vb, f,on/off )

Multicore Platform

Reliability alarms

Temperature Power

Perf. Counters

Core 1

Core N

Proc. 1

Proc. 2

...

Proc. N

INTERCONNECT

Private Mem

Private Mem

...

Private Mem

# Outline

- Introduction
- Scalable Control
- Scalable model learning
- Experimental Environment
- Challenges ahead

# DRM - General Architecture

- System (Chip Scale)
- **Controller**
- Sensors
  - Performance counter
    - **Minimize energy**
  - - PMU
  - Core temperature
    - **Bound the CPU temperature**
- Actuator - Knobs
  - ACPI states
- **Layered approach**
  - P-State → DVFS
  - C-State → $P_{GATING}$
    - **Stack of controllers**
  - Task allocation
- Controller
    - **Energy controller**
  - Reactive
    - **Thermal controller**
    - Threshold/Heuristic
    - Controller theory
  - Proactive
    - Predictors

# DRM - General Architecture

- System (Chip Scale)
- Sensors
  - Performance counter - PMU
  - Core temperature
- Actuator - Knobs
  - ACPI states
    - P-State → DVFS
    - C-State → $P_{GATING}$
  - Task allocation
- Controller
  - Reactive
    - Threshold/Heuristic
    - Controller theory
  - Proactive
    - Predictors

# Energy Controller



**CPU BOUND TASK**

*High Frequency*

*Low Frequency*

- Performance Loss
- Power reduction
- Energy Efficiency Loss!

- Same Performance
- Power reduction
- Energy Efficiency Gain!

**MEMORY BOUND TASK**

*High Frequency*

*Low Frequency*

# Thermal & Power Modeling

Linear dynamic state space model

Non linear

Task$_j$ → **Power model** → P$_j$ → **Thermal model** → T$_j$

$$P = g(\text{task}, f)$$

$$x_{n+1} = A x_n + B P_{n,j}$$

$$T_{n,j} = C x_n$$

**task**   **task**   ..........   **task**

# Thermal transient

Thermal locality (Direct Fourier law implication)

- Continuous model:
  - Thermal neighborhood = Physical
- Discrete model:
  - Thermal neighborhood depends on sample time
- Hotspot simulation of 'Intel SCC like' 48core
  - Each core : Area = 11.82mm2, Pmax = 2.6W
  - We powered on only Core(5,3)
  - T neighborhood > +0.1°C
- Thermal transient – Model Order
  - Different materials reflects in different time constants [1]
    - Silicon die, heat spreader, heat sink
    - Second order model

O.S time scale Neighborhood

Ts < 2ms    Ts < 0.25s

Ts < 50ms

Ts < 75ms    Ts < 0.5s

Ts < 0.1s    Ts > 0.5s

1 2 3 4 5 6 7 8

1
2
3
4
5
6

[1] W. Huang Differentiating the roles of IR measurement and simulation for power and temperature-aware design 2009.

Core Temperature Response

core15    core20    core21    core48

Temperature [°C]

Time [s]

# Thermal Controller


[Intel®, ISSCC 2007]

**COMPLEXITY**

## Model Predictive Controller

- Internal prediction: avoid overshoot
- Optimization: maximizes performance



**MPC**

- Centralized
  - aware of neighbor cores thermal influence
  - All at once – MIMO controller
- Complexity !!!

## Classical feed-back controller

- PID controllers
- Better than threshold based approach
- Cannot prevent overshoot



## Threshold based controller

- $T > Tmax \rightarrow$ low freq
- $T < Tmin \rightarrow$ high freq
- cannot prevent overshoot
- thermal cycle

# MPC Scalability

**MPC *Complexity***

- **Implicit** - *a.k.a. on-line*
  - computational burden
- **Explicit** – *a.k.a. off-line*
  - high memory occupation

**Out of memory**

**# Explicit regions**

1296

36

4    8    16    **# CORES**

**Complexity grows *superlinearly* with number of cores!!**

# Addressing Scalability

**On a short time window, power has a local thermal effect!**

**One controller for each core**
**Controller uses:**
- local power & thermal model
- neighbor's temperatures

1296

**# Explicit regions**

36     8          16          32

4          8          16  **# CORES**

**Fully distributed**
**Complexity scales linearly with #cores**

# Distributed Control

## Distributed and hierarchical controllers:

- **Energy Controller (EC)**
  - **Output Frequency $f_{EC}$**
    - **Minimize power – CPI based**
    - **Performance degradation < 5%**
- **Temperature Controller (TC)**
  - **Distributed MPC**
  - **Inputs:**
    - $f_{EC}, T_{CORE}, T_{NEIGHBOURS}$
  - **Output**
    - **Core frequency ($f_{TC}$)**



$CPI_{N,k+1}$

controller node

$EC_N$

$f_{EC\ N,k+1}$

$CPI_{i,k+1}$

$T_{N-1,k}$

$TC_N$

$EC_i$

$f_{EC\ i,k+1}$

$T_{i+1,k}$

$T_{x,k}$

$f_{TC\ N,k+1}$

$Core_N$

$TC_i$

$T_{i,k}$

$T_{i-1,k}$

$Core_i$

$Core_1$

$T_{i,k}$

multicore

$T_{2,k}$

$f_{TC\ i,k+1}$

# High Level Architecture



**Energy Controller**

| | | |
|---|---|---|
| $f_{1,EC}$ | | $T_1 + T_{neigh}$ |
| CPI | Thermal Controller Core 1 | $f_{1,TC}$ |
| $f_{1,EC}$ | | $T_1 + T_{neigh}$ |
| CPI | Thermal Controller Core 2 | $f_{2,TC}$ |
| $f_{1,EC}$ | | $f_{3,TC}$ |
| CPI | Thermal Controller Core 3 | $T_3 + T_{neigh}$ |
| $f_{1,EC}$ | | $f_{4,TC}$ |
| CPI | Thermal Controller Core 4 | $T_4 + T_{neigh}$ |

**PLANT**

**Distributed Controller**

# Distributed Thermal Controller



MPC Controller Core 1

CPI
$f_{1,EC}$
$g(\cdot)$
$P_{1,EC}$

**MPC Controller**
Linear Model
QP Optimiz

$T_{ENV}$
$x_1$

CPI
$g^{-1}(\cdot)$
$P_{1,TC}$
$f_{1,TC}$

**Observer**

$T_1$

2 states per core

Nonlinear (Frequency to Power)

**Implicit formulation**

Linear (Power to Temperature)

$$\min \sum_{k=0}^{N-1} \| Q \cdot (f_{TC,1}(t+k) - f_{EC,1}(t+k)) \|_2$$

s.t

$$0 \le T_1(t+k|t) \le T_{MAX,1}$$

$$\hat{x}(k+1) = A_1 \cdot \hat{x}(k) + B_1 \cdot u(k) + L_1 \cdot (T_1 - C_1 \cdot \hat{x}(k))$$

Classic Luenberger state observer

# Explicit Distributed Controller

$$\begin{pmatrix} x_{11} = T_1 \\ x_{12} \end{pmatrix}$$

$x_1$

$x_{1SHIFTED}$

$H_i \cdot x(k) \leq K_i$ ?

↓

Region number

$[x_1 ; T_{ENV} ; P_{1,EC}]$

$A_1^{-1} \cdot B_1$

| REGION NUMBER | GAIN MATRICES |
|---|---|
| 1 | $F_1$, $G_1$ |
| 2 | $F_2$, $G_2$ |
| ... | ... |
| $nr$ | $F_{nr}$, $G_{nr}$ |

$u(k) = \Delta P_1$

$P_{1,TC}$

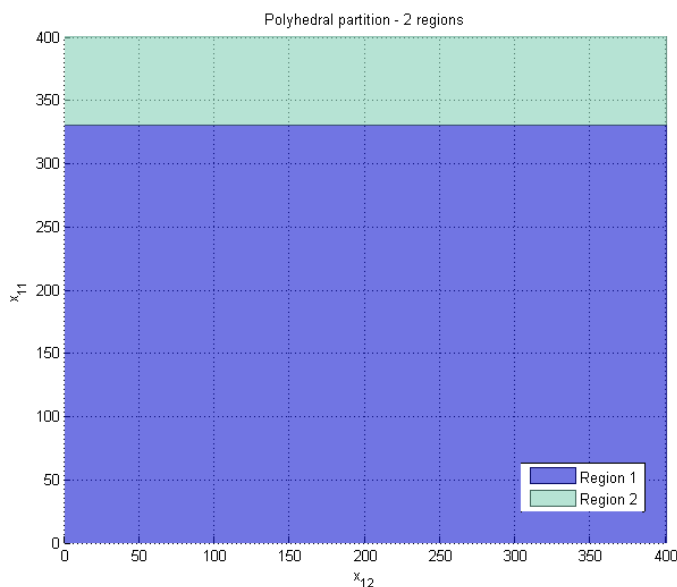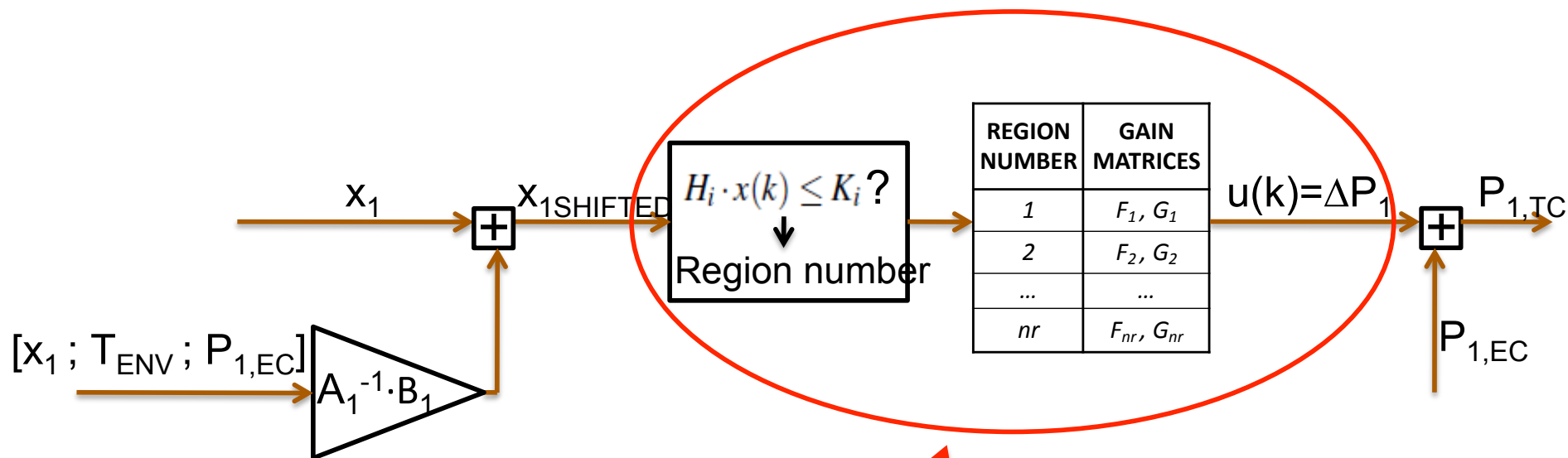$P_{1,EC}$

Our aim is to minimize the difference between the input $P_{1,TC}$ (also called manipulated variable MV) and the reference ($P_{1,EC}$). Our controller can only take in account a constant reference. To overcome this limitation we reformulate the tracking problem as a regulation problem consisting in taking the $\Delta P_1$ (the new MV) to 0. The regulated power $P_{1,TC}$ is:

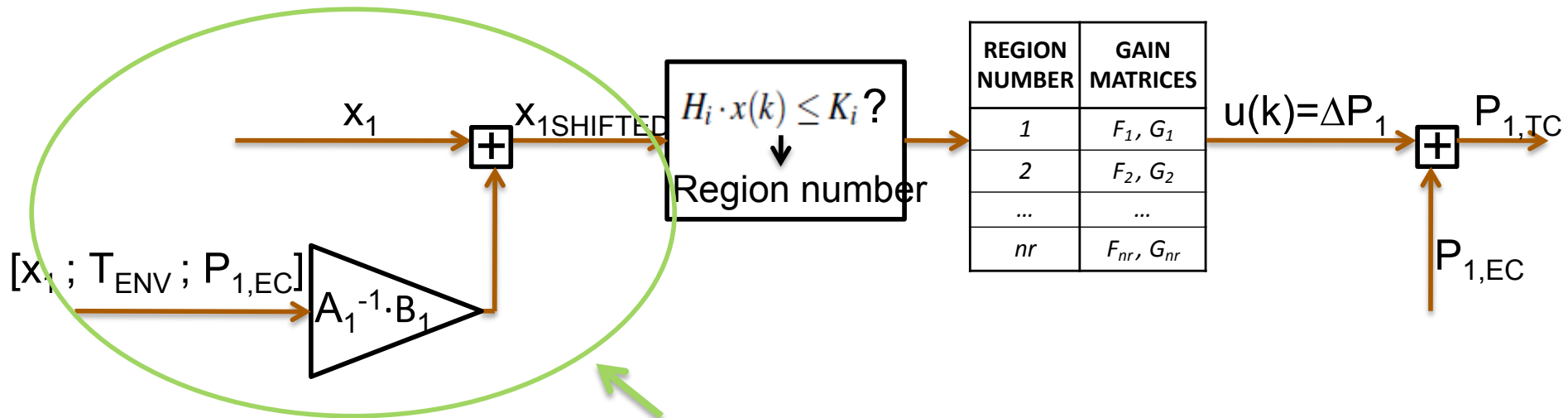$$P_{1,TC} = P_{1,EC} + \Delta P_{1,EC}$$

# Explicit Distributed Controller



At each time instant the system belongs to a region according with its current state. On each region the explicit controller executes the following linear control law:

$$u(x) = F_i \cdot x(t) + G_i$$

# Explicit Distributed Controller



The prediction evaluated by our explicit controller cannot take into account the measured disturbances ($u_{MD}$=[Tenv, P1, $T_{neigh}$]). Thus we exploit the superposition principle of linear systems:

$$x(k+1) = f(x(k), u_{MV}(k), u_{MD}(k)) \rightarrow x(k+1) = f(x(k), u_{MV}(k), 0) + f(x(k), 0, uMD(k)))$$

To remap the effect of these elements we exploit the model to modify the state ($x(k) \rightarrow x_{SHIFTED}(k)$) projecting one step forward the MDs effects.
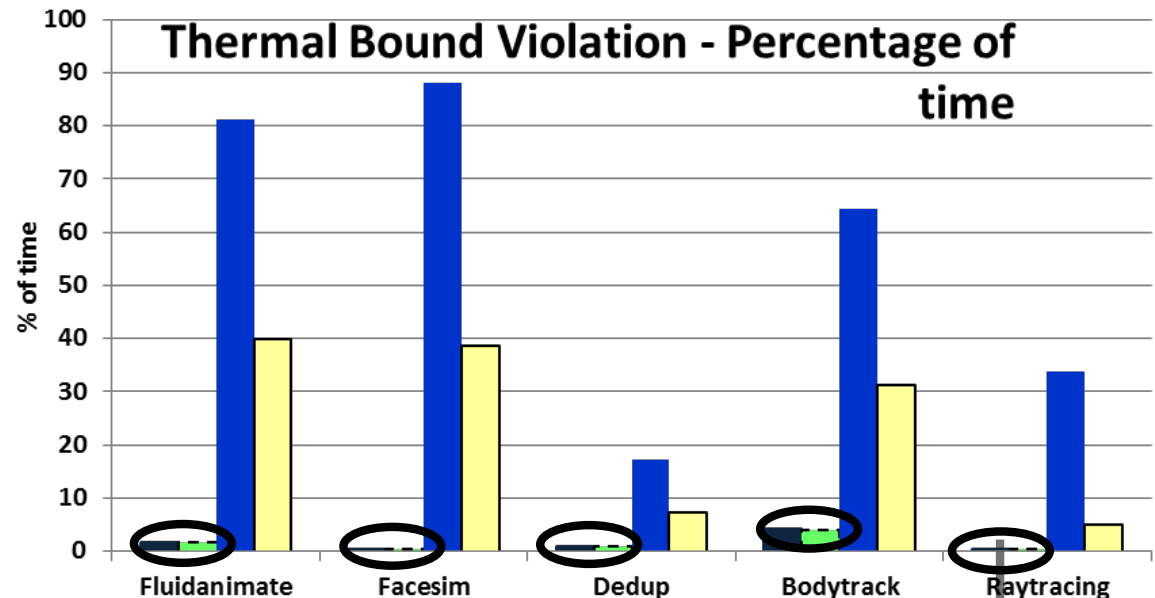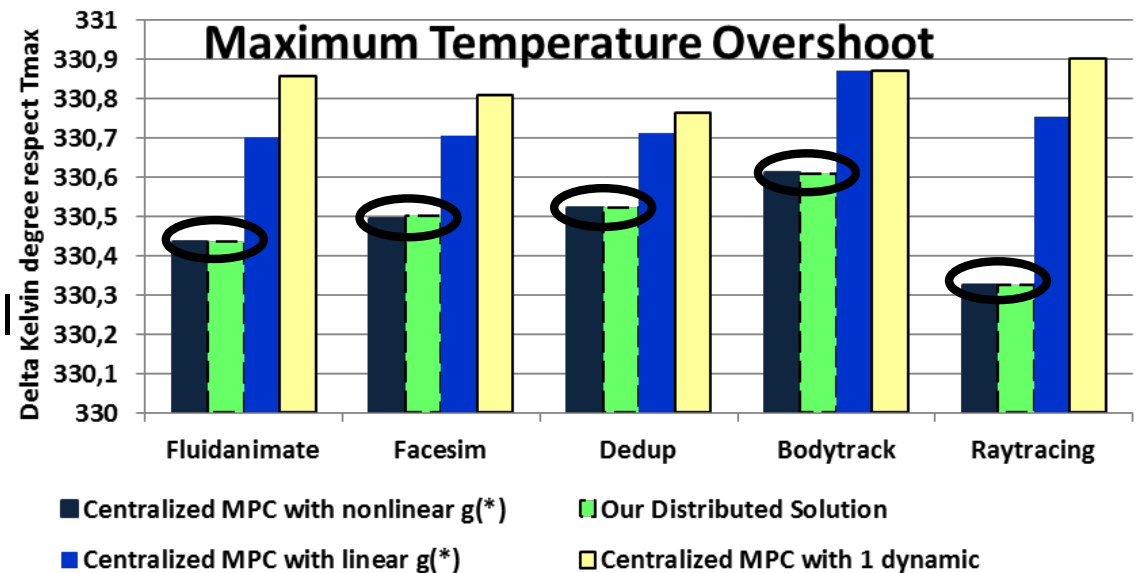
$$x(k+1) = A_1 \cdot x(k) + B_1' \cdot \Delta P_1 + B_1'' \cdot [P_{1,EC} \ T_{ENV} \ T_{neig}] \rightarrow x(k+1) = A_1 \cdot x_{SHIFTED}(k) + B_1' \cdot \Delta P_{1,EC}$$

$$A_1 \cdot x_{SHIFTED}(k) = A_1 \cdot x(k) + B_1'' \cdot [P_{1,EC} \ T_{ENV} \ T_{neig}] \rightarrow x_{SHIFTED}(k) = x(k) + A_1^{-1} \cdot B_1'' \cdot [P_{1,EC} \ T_{ENV} \ T_{neig}]$$

# MPC trade-off

## Trace Driven Simulation (Matlab) – gold model

- Parsec trace obtained on real HW

- Power Model: Nonlinear vs. linear

- Thermal Model: one vs. two
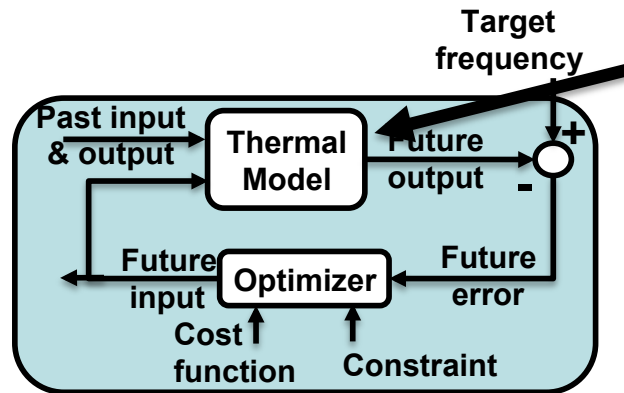
- Centralized vs. Distributed

# Outline

- Introduction
- Scalable Control
- Scalable model learning
- Experimental Environment
- Challenges Ahead

# Model Identification

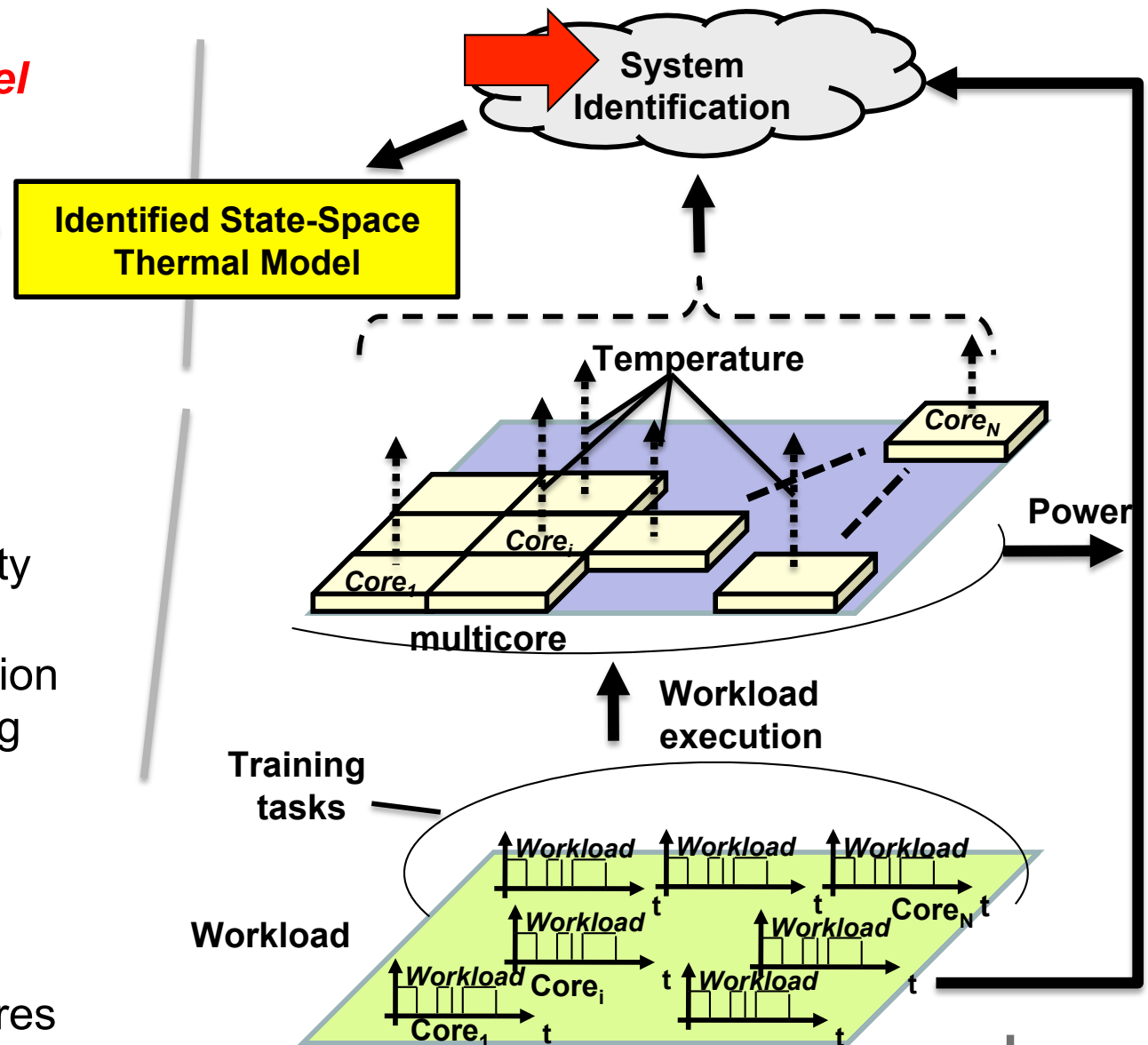**MPC needs a *Thermal Model***



**Identified State-Space Thermal Model**

- Accurate, with low complexity
- Must be known "at priori"
- Depends on user configuration
- Changes with system ageing

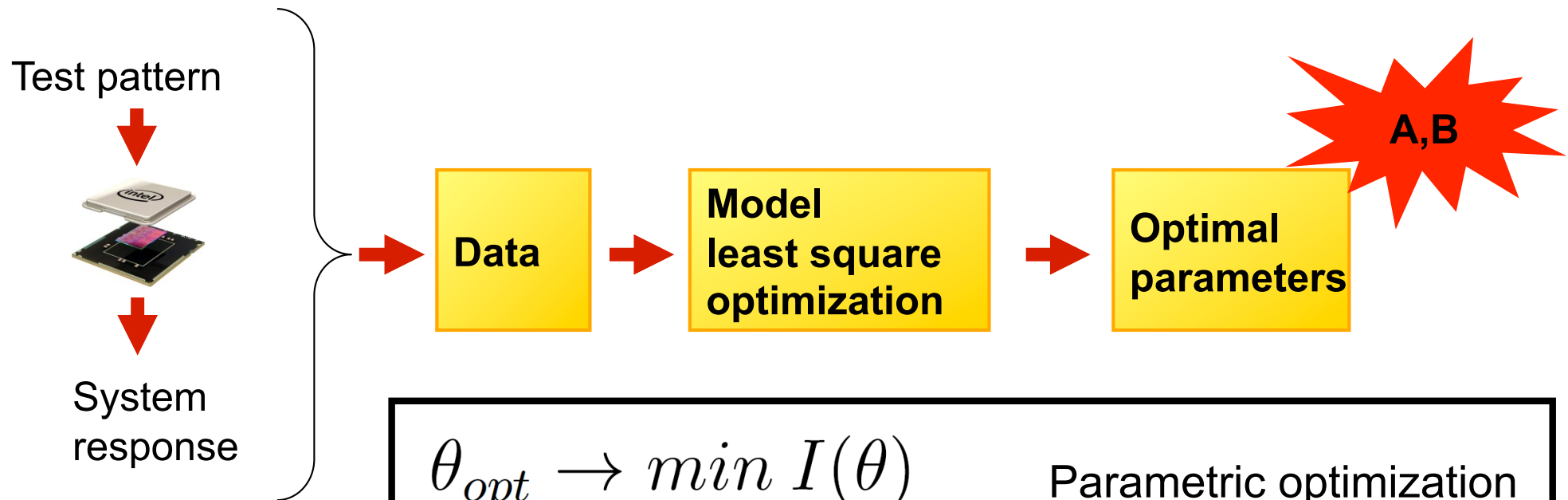**"In field" Self-Calibration**

- Force test workloads
- Measure cores temperatures
- System identification

# LS System Identification

Test pattern

System response

**Data** → **Model least square optimization** → **Optimal parameters**

**A,B**

$$\theta_{opt} \rightarrow \min_{\theta} I(\theta)$$  Parametric optimization

$$\theta = \{A, B\}$$  Parameters
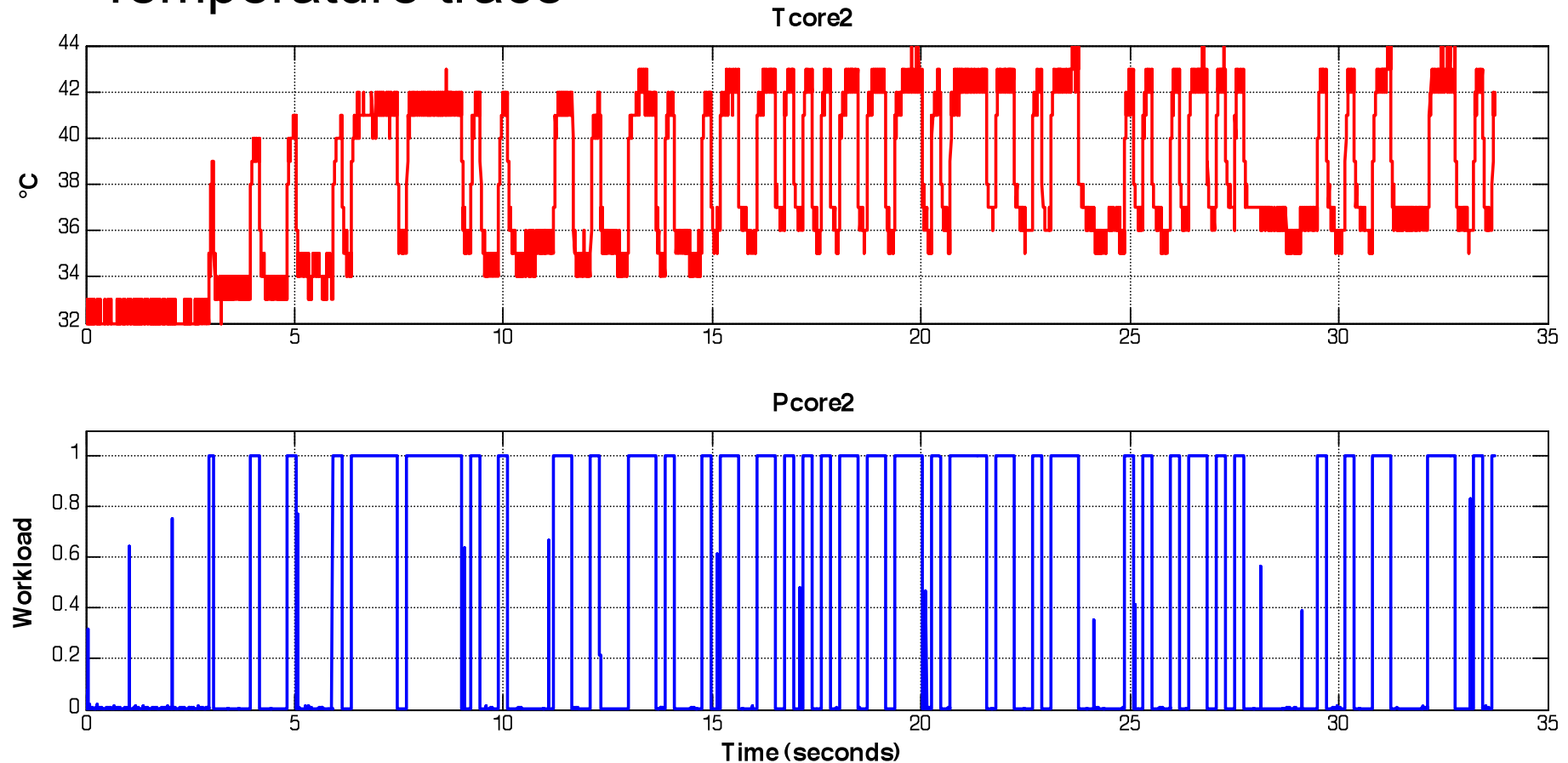
$$I(\theta) = \frac{1}{N_s} \sum_{i=1}^{N_s} e_i^2$$  Cost function

$$e = T_{model} - T_{real}$$  Error Function

# Workload & Temperature

Temperature trace
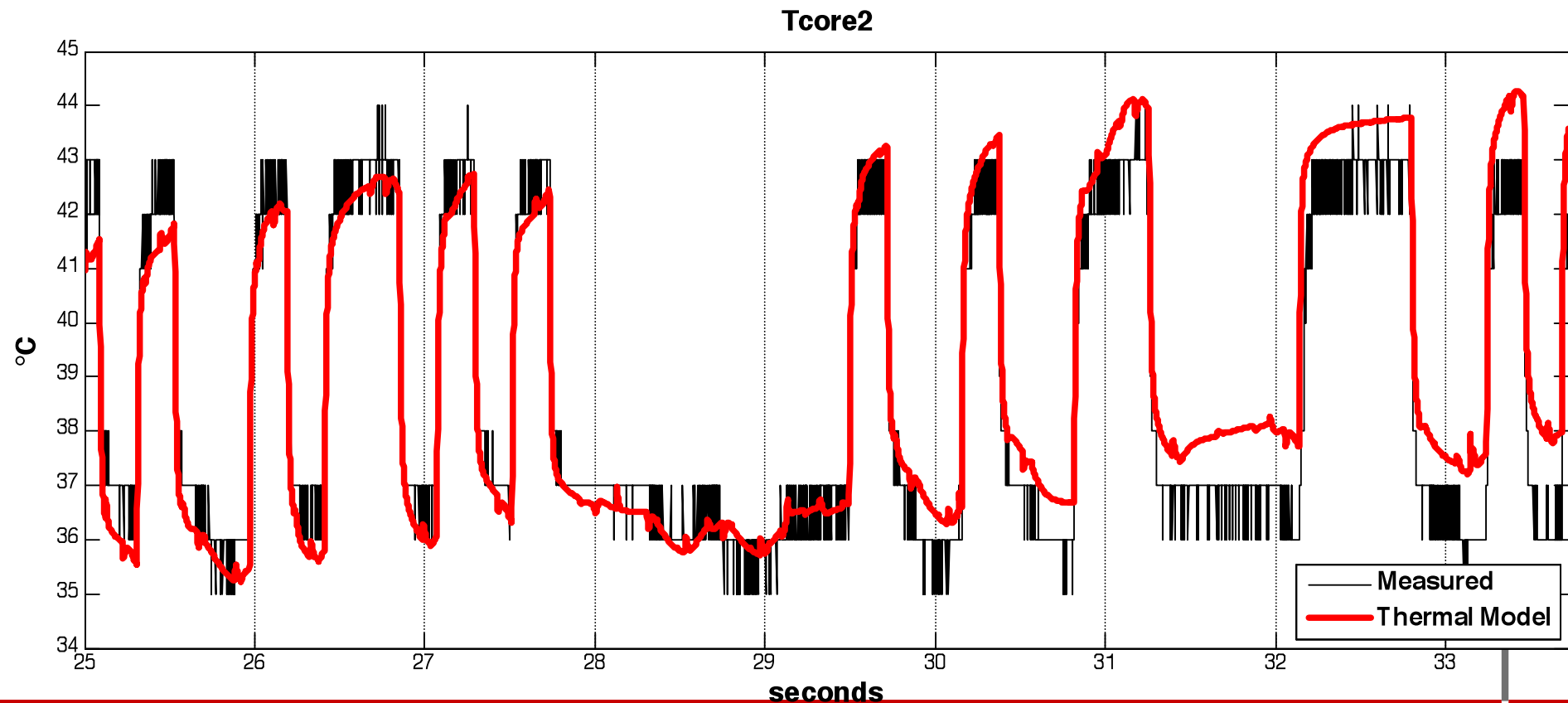


Pseudorandom workload pattern

# Black-box Identification

**Identification based on pure LS fitting**

## MEASURED vs. SIMULATED TEMPERATURE

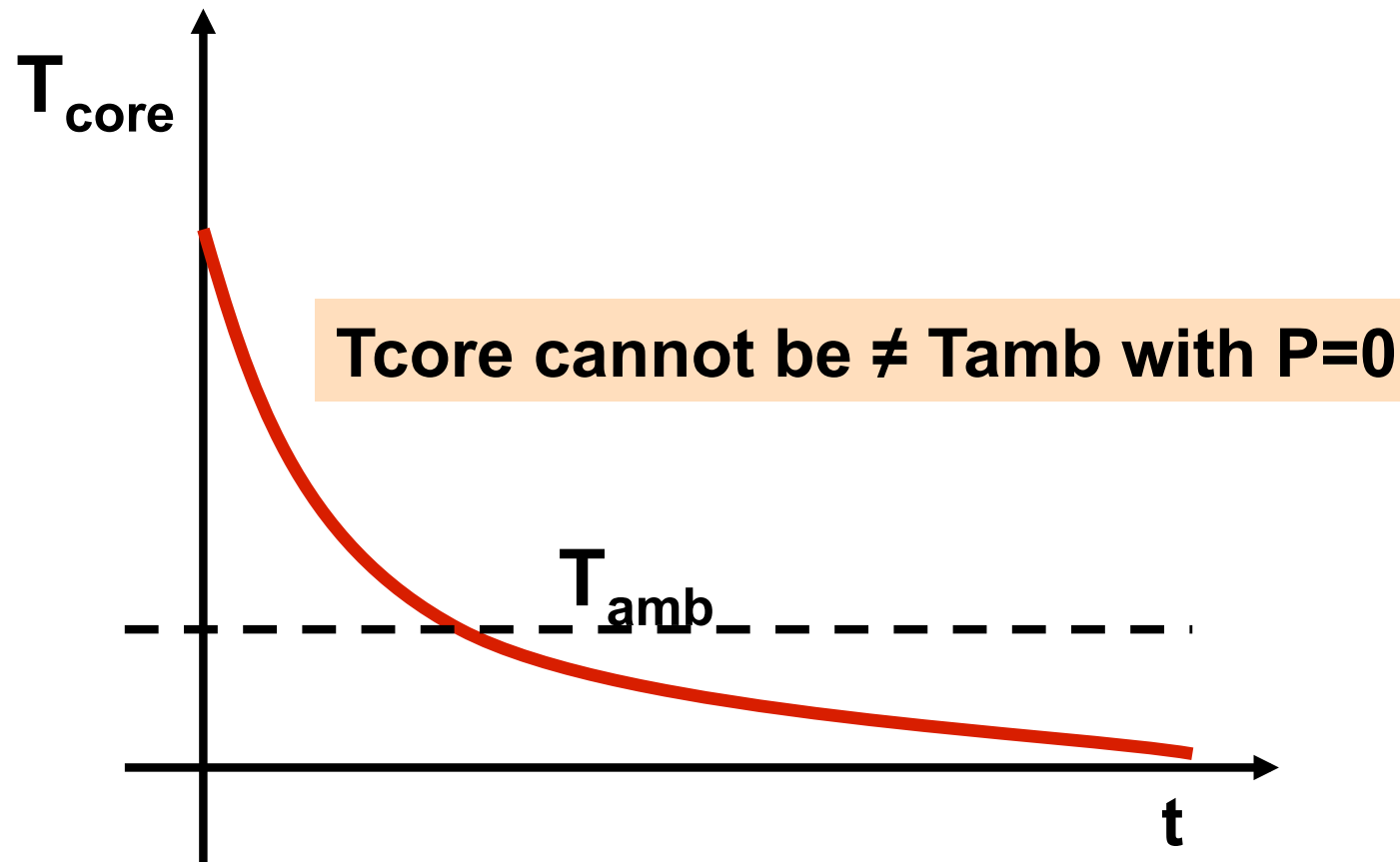**Tcore2**

# Validation



**Problem**: unstable Model

# Gray Box identification

LS model must be constrained by physical properties to avoid over-fitting



Tcore cannot be ≠ Tamb with P=0

# Physical Constraints

# Model learning Scalability

**MPC Weaknesses – 2nd**
*Internal Thermal Model*



**Identified State-Space Thermal Model**

**System Identification**

**MPC**

- Accurate, with low complexity
- Must be known "at priori"
- Depends on user configuration
- Changes with system ageing

**"In field" Self-Calibration**

- Force test workloads
- Measure cores temperatures
- System identification

# Model learning Scalability

**Complexity issue**
- State-of-the-art is centralized MIMO
- Least square based – is based on matrix inversion (cubic with #cores)



**1230**

**105**

**Distributed approach: each core identifies its local thermal model Complexity scales linearly with #cores**

- Accurate, with low complexity
- Must be known "at p..."
- Depends on user configuration
- Changes with system being
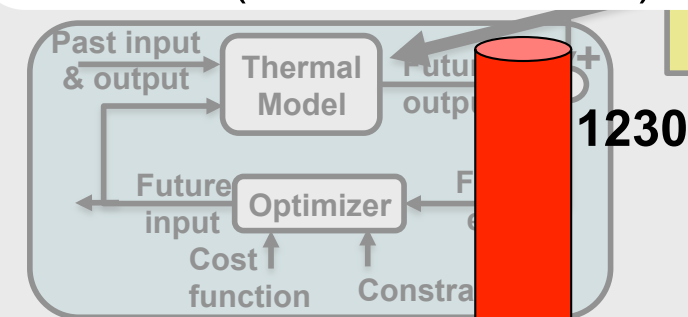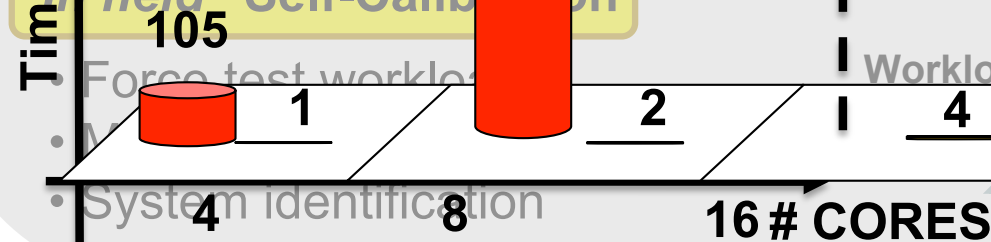
"in field" Self-Calibration

- Force test workload
- System identification

Time - s

1     2     4

4     8     16 # CORES

# Distributed Model learning

Distributed MISO identification

• ARX model

**State Space Model**

$$\begin{bmatrix} T_i(k+1) \\ x_i(k+1) \end{bmatrix} = A \cdot \begin{bmatrix} T_i(k) \\ x_i(k) \end{bmatrix} + B \cdot \begin{bmatrix} P_i(k) \\ T_{neig,i}(k) \end{bmatrix}$$

**ARX Model**

$T_{neigW,i}(k-s)$    $T_{neigE,i}(k-s)$

**Model$_i$**

$P_i(k-s)$

$T_i(k-s)$

$Core_i$

$T_{neigN,i}(k-s)$

$T_i(k) = \alpha_2 \cdot T_i(k-1) + \alpha_1 \cdot T_i(k-2) +$    **A**uto**R**egressive term

$+ \beta_{1,1} \cdot P_i(k-2) + \beta_{1,2} \cdot P_i(k-1) +$    Power input

Neightbours temperatures $+ \beta_{2,1} \cdot T_{neig1,i}(k-2) + \beta_{2,2} \cdot T_{neigE,i}(k-1) + ...$

$+ e(t)$    Errors (disturbances, measurement errors)

# Distributed model-learning

Distributed identification
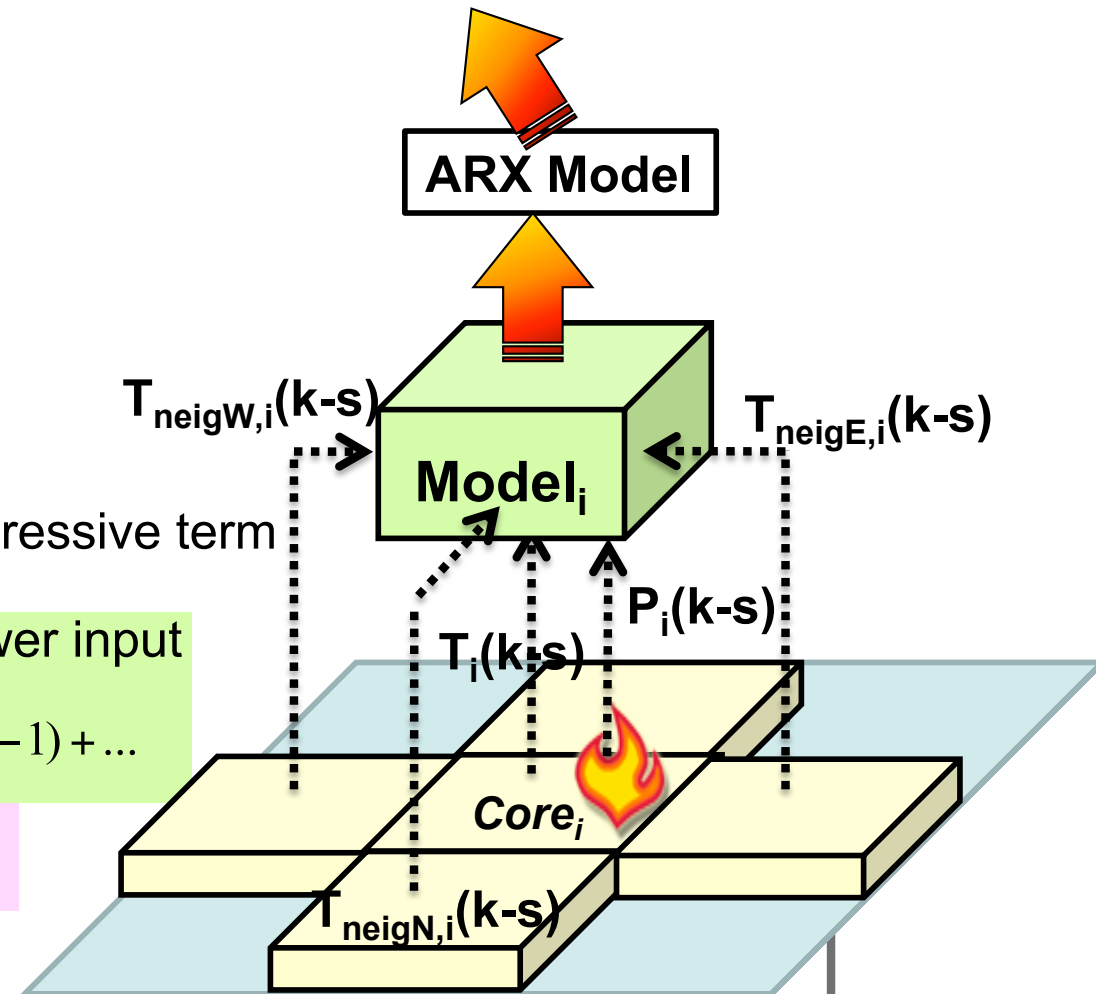
- ARX model
- Computation Algorithm
- Results

- Objective: find $a_i$ and $b_{i,j}$

- System data collection:
  - input: PRBS signals to all cores (persistently exciting inpt sequence)
  - output: Temperatures of all cores ($T_i^o$ con  i= # core)

- Parameters computation:
  - $T^p(a_i, b_{i,j})$=T(k+1|k)  computed with previous equation
  - $T^o$ output temperature (measured)
  - Least square algorithm:

$$\min \frac{1}{L} \sum_{k=3}^{L+2} \left( T^0 - T^p \right)^2 = \min \frac{\left\| T^0 - T^p \right\|_2^2}{L}$$

**1) Mean Absolute Error between original and identified models**

**2) Temperature response of core 1**

# Outline

- Introduction
- Scalable Control
- Scalable model learning
- Experimental Environment
- Challenges ahead

# Simulation Strategy

## Trace driven Simulator [1]:

- Not suitable for full system simulation (How to simulate O.S.?)
- looses information on cross-dependencies
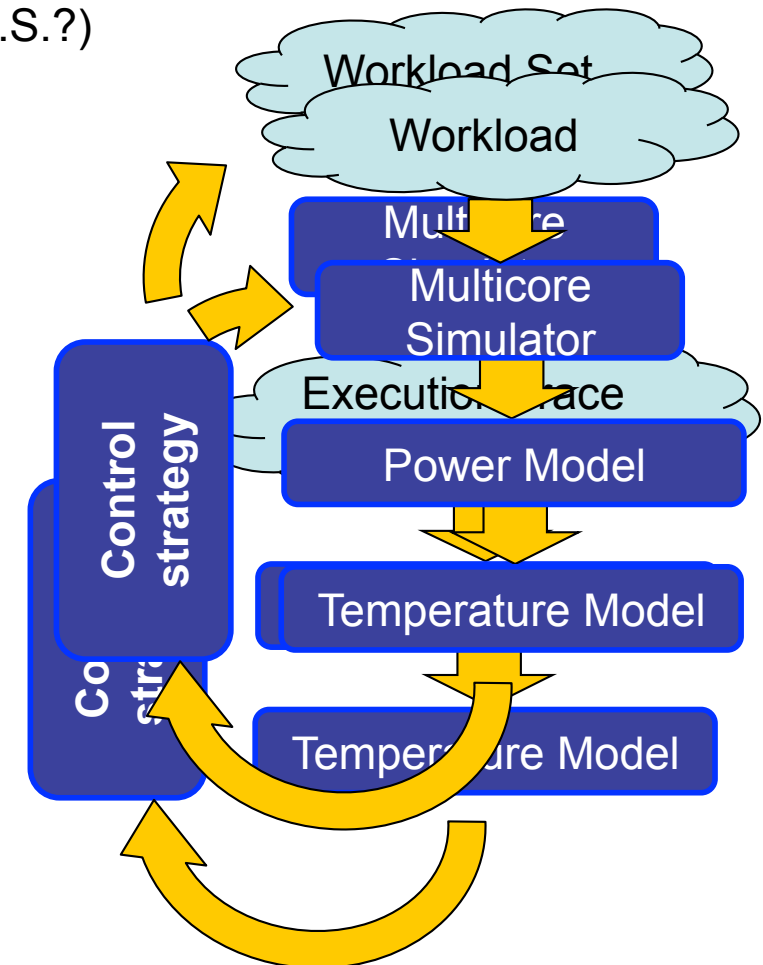    → resulting in degraded simulation accuracy

## Closed loop simulator:

- Cycle accurate simulators [2] :
  - High modeling accuracy
  - Suppost well-established power and temperature co-simulation based on analytical models and system micro-architectural knowledge
  - Low simulation speed
  - Not suitable for full-system simulation
- Functional and instruction set simulators:
  - Allow full system simulation
  - Lower internal precision
  - Less detailed data → no micro-architectural model
  - Introduces the challenge of having accurate power and temperature physical models

[1] P Chaparro et al. Understanding the thermal implications of multi-core architectures. 2007
[2] Benini L. et al. MPARM: Exploring the multi-processor SoC design space with SystemC 2005

# Virtual Platform



**Simics by Virtutech:**

- full system functional simulator
- models the entire system: peripherals, BIOS, network interfaces, cores, memories
- allows booting full OS, such as Linux SMP
- supports different target CPU (arm, sparc, x86)
- x86 model:
  - in-order
  - all instruction are retired in 1 cycle
  - does not account for memory latency
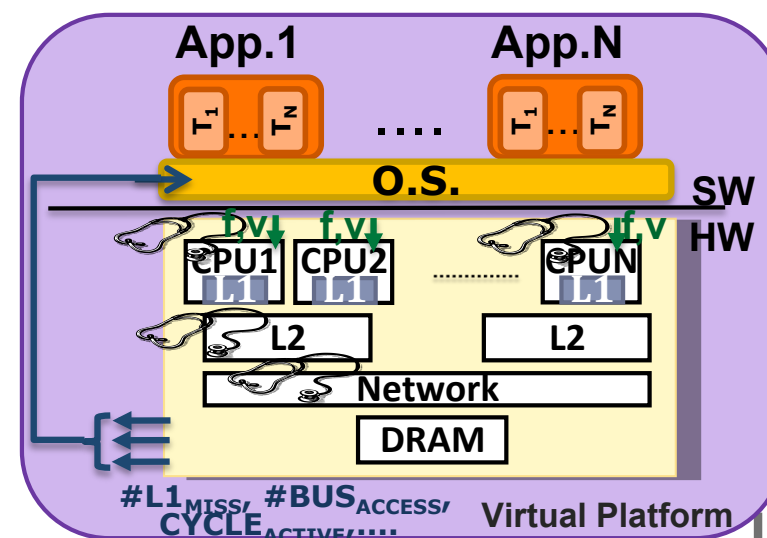
**Memory timing model**

- RUBY – GEMS (University of Wisconsin)[1]
  - Public cycle-accurate memory timing model
  - Different target memory architectures
  - fully integrated with Virtutech Simics
  - written in C++
  - we use it as skeleton to apply our add-ons (as C++ object)

*[1] Martin Milo M. K. et al. Multifacet's general execution-driven multiprocessor simulator (GEMS) toolset 2005*

# Virtual Platform

**Performance counter (DVFS) module:**

- Needed by Simics per control policy change at run-time
- RUBY does now Performance Counter module to support it
  - exports to O.S. an application defined frequency:
- We add a new DVFS module to support clock cycles and stall cycles expired,
  - ensure L2 structure and DRAM to have a constant clock frequency
  - L1 latency scale with Simics processor clock frequency

# Virtual Platform

## Power model module:

- At run-time estimate the power consumption of the target architecture
- Core model $P_T = [P_D(f, CPI) + P_S(T, VDD)] * (1 - \text{idleness}) + \text{idleness} * (P_{IDLE})$
- $P_D$ experimentally calibrated analytical power model
- Cache and memory power – access cost estimated with CACTI [1]



[1] Thoziyoor Shyamkumar et al. A comprehensive memory modeling tool and its application to the design and analysis of future memory hierarchies. 2008

# Power Model

- Power model interface

i-th CPU
# Istruction retired
# Stall Cycle
# HLT Istruction

i-th L1
# Line & WD Read
# Line & WD Write

$f, V_{dd}, T$

$f_1 = k_1 * f_{nom}$   $f_2 = k_2 * f_{nom}$   $f_N = k_N * f_{nom}$

CPU1   CPU2   CPU N
L1     L1     L1

L2     L2

Network

DRAM

Ruby & simics

i-th L2
# Line Read
# Line Write

DRAM
# Burst Read
# Burst Write

POWER MODEL

POWER & ENERGY

# Modeling Real Platform – Power

**Real Power Meas**

- *Intel server system*
  - 16 cores - 4 qu
  - 16GB FBDIMM
  - Intel® Core™
- *At the wall Power co*
  - test:
    - set of synth
    - forcing all t
    - for each be
      correlate it v



**CPI vs freq.**

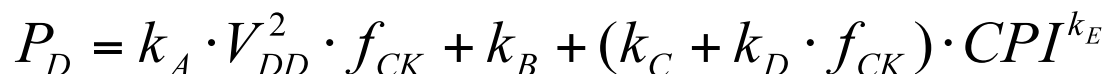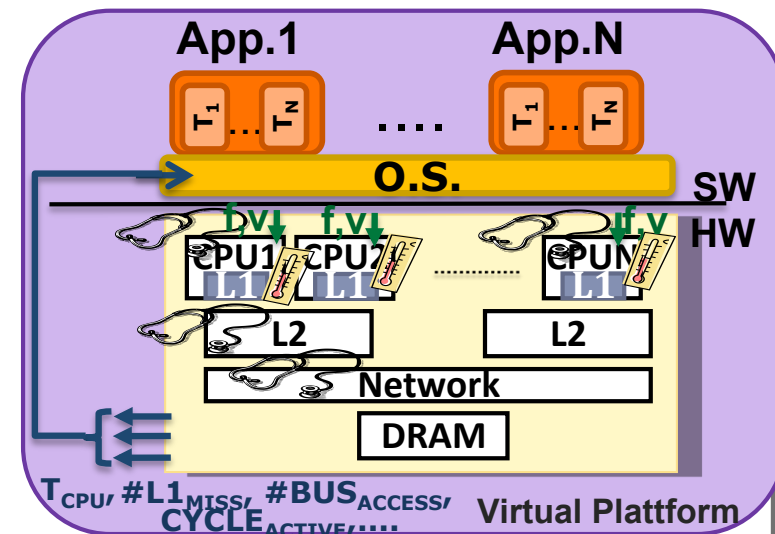| | | | |
|---|---|---|---|
| – – – 2970MHz – Fi | ✳ 2970MHz – Real | –··– 2630MHz – Fi | ✳ 2630MHz – Real |
| —— 2400MHz – Fi | ✳ 2400MHz – Real | ······ 2140MHz – Fi | ✳ 2140MHz – Real |
| –·–· 1800MHz – Fi | ✳ 1800MHz – Real | –··–· 1600MHz – Fi | ✳ 1600MHz – Real |

High accuracy at high and low CPI

Dynamic Power – [W]

Clock per Istruction (CPI)

$$P_D = k_A \cdot V_{DD}^2 \cdot f_{CK} + k_B + (k_C + k_D \cdot f_{CK}) \cdot CPI^{k_E}$$

- We relate the static power with the operating point by using an analytical model
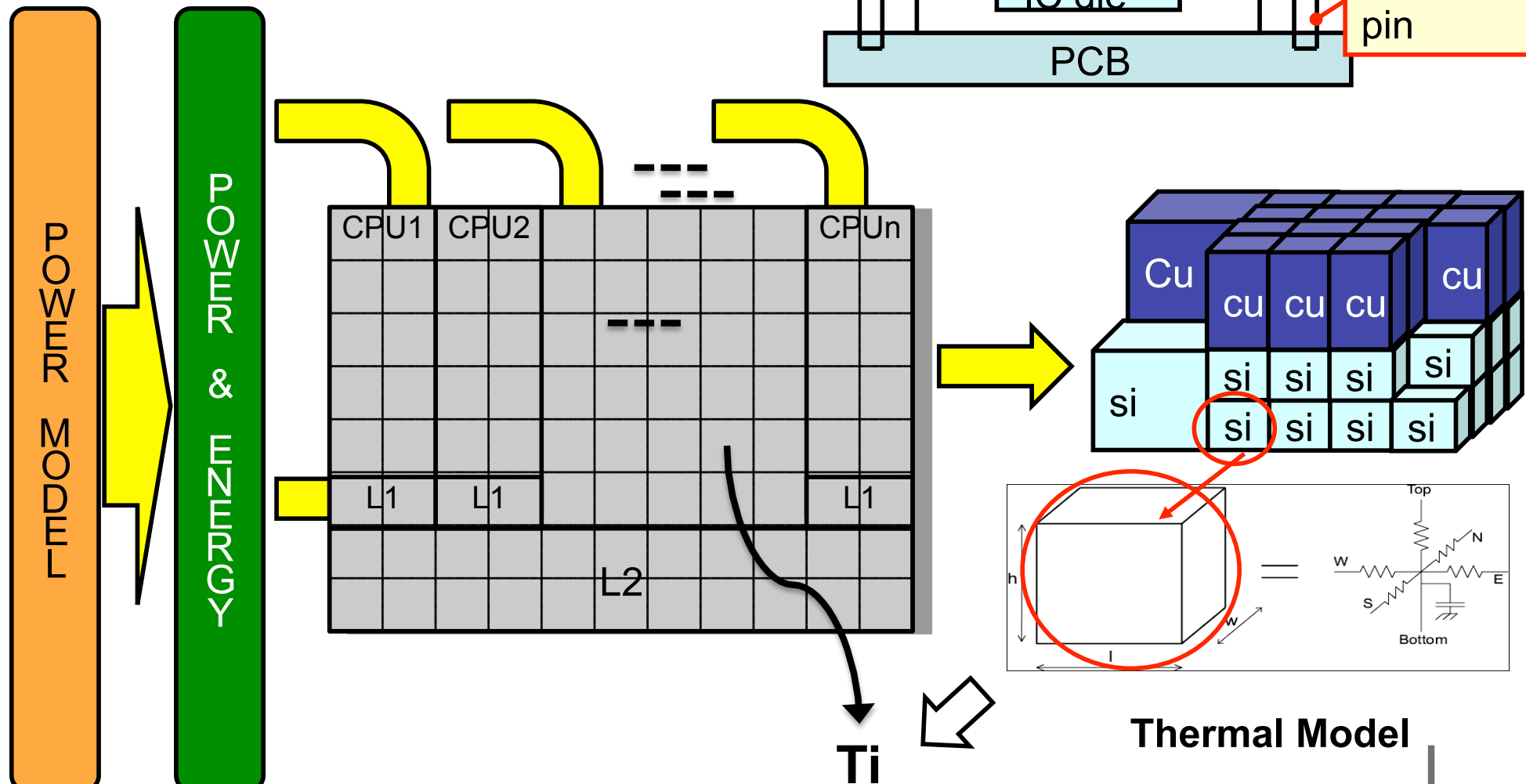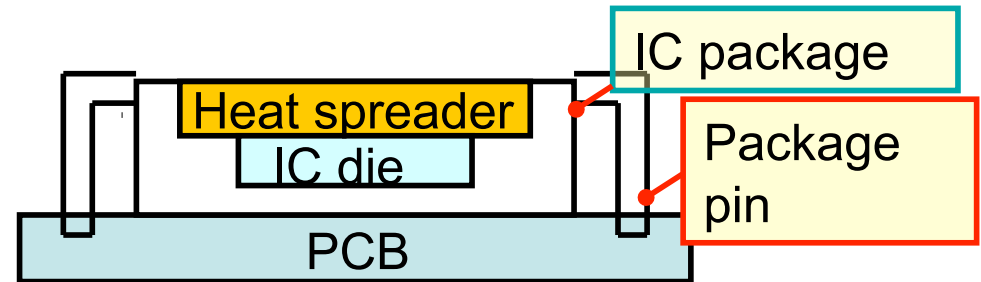
# Virtual Platform

**Temperature model module:**

- we integrate our virtual platform with a thermal simulator [1]
- Input: power dissipated by the main functional units composing the target platform
- Output: Provides the temperature distribution along the simulated multicore die area as output
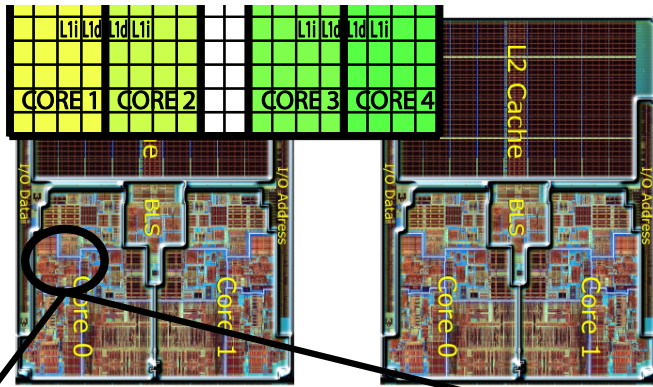
[1] Paci G. et al. Exploring "temperature-aware" design in low-power MPSoCs

# Thermal Model

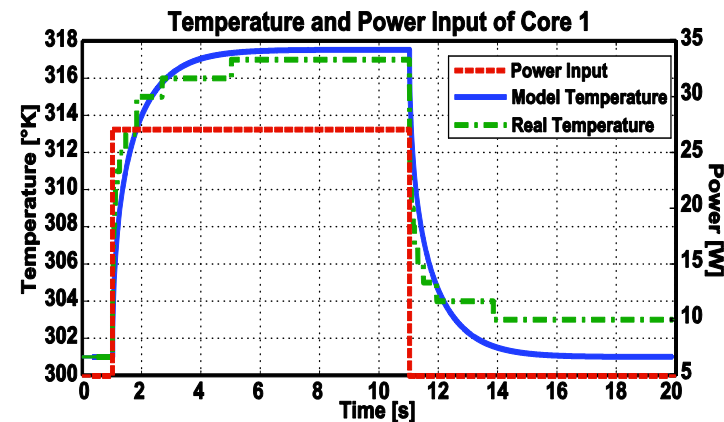- Methods to solve temperature



Thermal Model

# Modeling Real Platform– Thermal

- Thermal Model Calibration :

  - Derived from Intel® Core™ 2 Duo layout

  - We calibrate the model parameter to simulate real HW transient

  - High accuracy (error < 1%) and same transient behavior



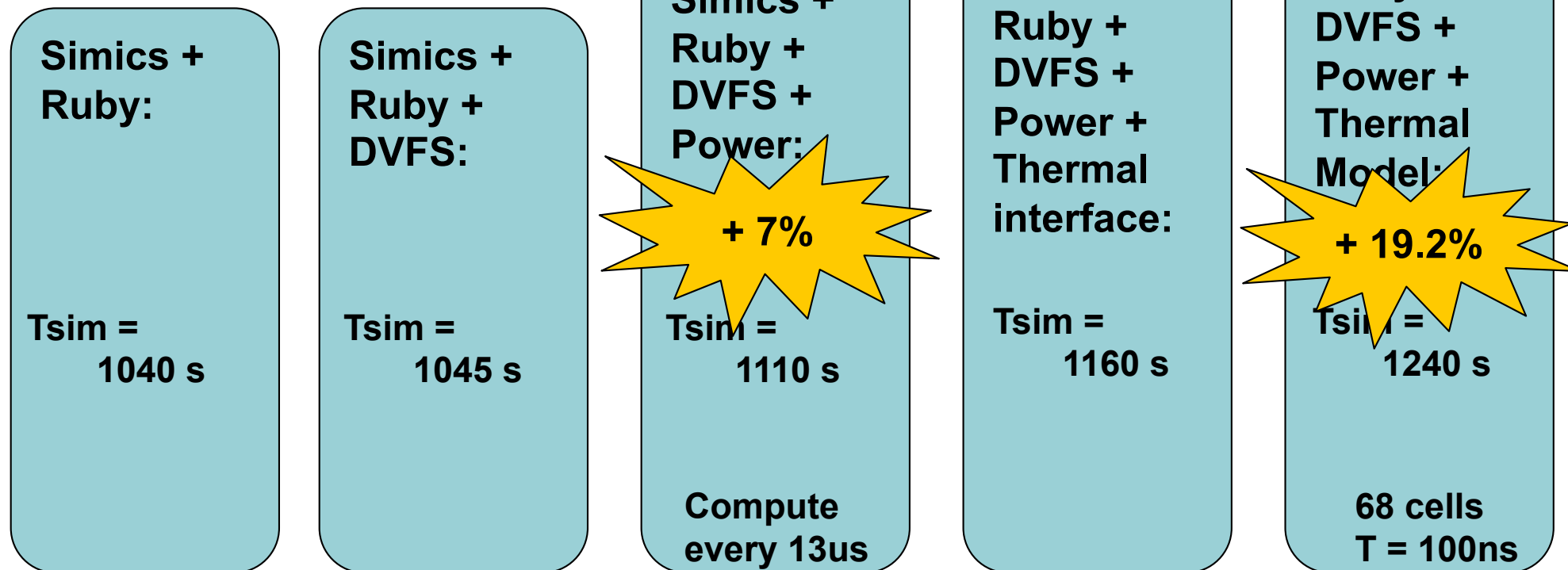| silicon thermal conductivity | $150 \cdot \left(\frac{300}{T}\right)^{4/3}$ W/mK |
|---|---|
| silicon specific heat | $1.628e^{-12}$ J/um |
| silicon thickness | 350um |
| copper thermal conductivity | 400W/mK |
| copper specific heat | $3.55e^{-12}$ J/um 3K |
| copper thickness | 2057um |
| elementary cell length | 1312um |
| package-to-air conductivity | 0.4K/W |



Temperature and Power Input of Core 1

# Virtual Platform Performance

- Target:
  - 4 core Pentium® 4
  - 2GB RAM
  - 32 KB private L1 cache
  - 4 MB shared L2 cache
  - Linux OS

- Host:
  - Intel® Core™ 2 Duo
  - 2.4 Ghz
  - 2GB RAM

1 Billion instruction

**Simics + Ruby:**

**Tsim = 1040 s**

**Simics + Ruby + DVFS:**

**Tsim = 1045 s**

**Simics + Ruby + DVFS + Power:**

**Tsim = 1110 s**

**Compute every 13us**

**+ 7%**

**Simics + Ruby + DVFS + Power + Thermal interface:**

**Tsim = 1160 s**

**Simics + Ruby + DVFS + Power + Thermal Model:**

**Tsim = 1240 s**

**68 cells T = 100ns**

**+ 19.2%**

# Mathworks Matlab/Simulink

- **Numerical computing environment** developed to design, implement and test numerical algorithms

- Mathworks Simulink – for **simulation of dynamic systems**: simplifies and speedups the development cycle of control systems

- **Can be called** as **a computational engine by** writing C and Fortran **programs** that use Mathworks Matlab's engine library

- Controller design - two steps:

  - developing the control algorithm that optimizes the system performance

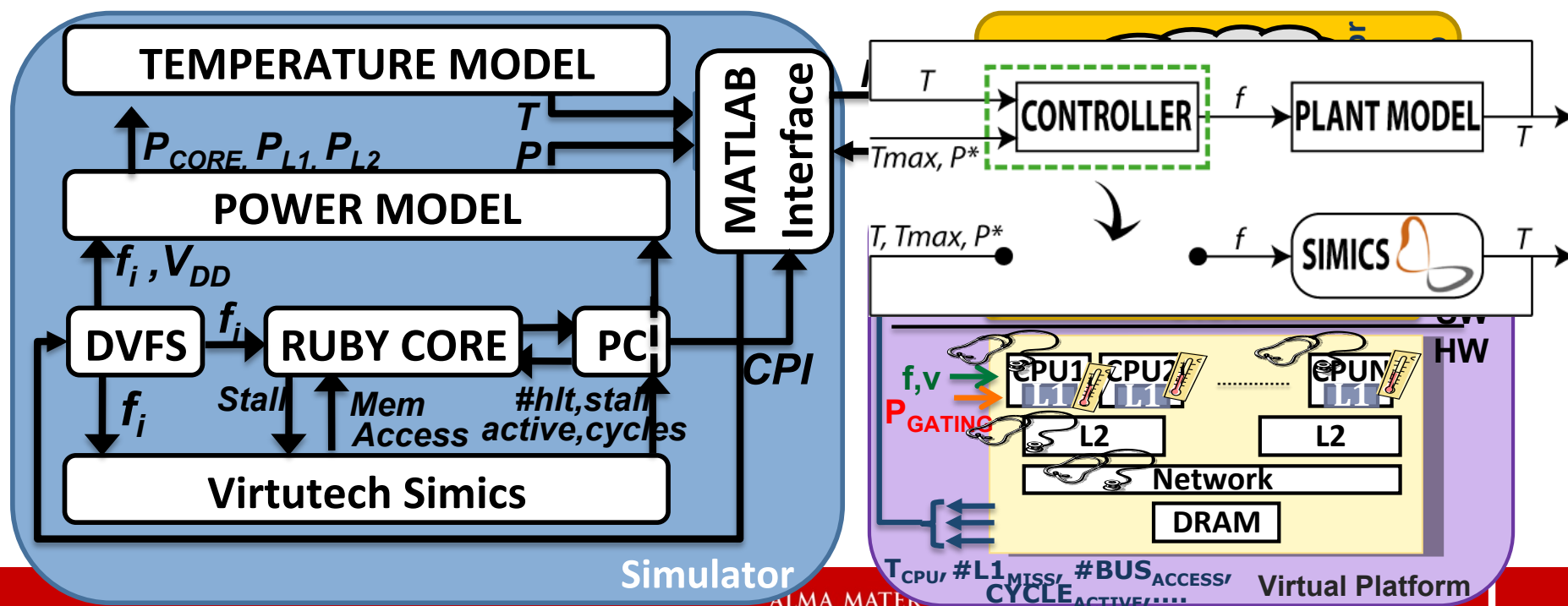  - implementing it in the system

We allow a Mathworks Matlab/Simulink description of the controller to directly drive at run-time the performance knobs of the emulated system

# Virtual Platform

**Mathworks Matlab interface:**

- New module named Controller in RUBY
- Initialization: starts the Mathworks Matlab engine concurrent process,
- Every N cycle - wake-up:
  - send the current performance monitor output to the Mathworks Simulink model
  - execute one step of the controller Mathworks Simulink model
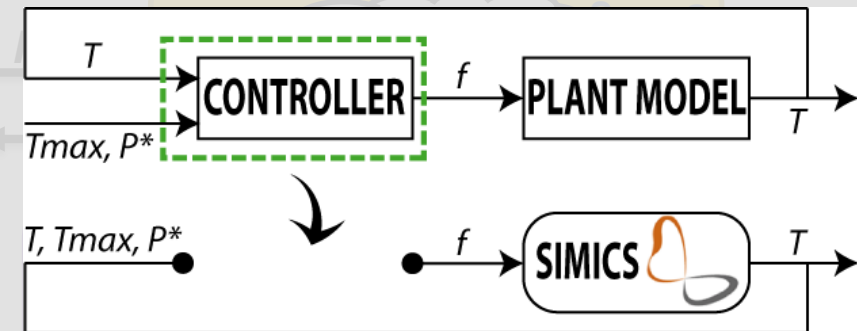  - propagate the Mathworks Simulink controller decision to the DVFS module

# Virtual Platform

## CONTROL-STRATEGIES DEVELOPMENT CYCLE

1. Controller design in Mathworks Matlab/Simulink framework
   • system represented by a simplified model
   • obtained by physical considerations and identification techniques
2. Set of simulation tests and design adjustments done in Simulink
3. Tuned controller evaluation with an accurate model of the plant done in the virtual platform
4. Performance analysis, by simulating the overall system

# Results

## Energy Controller (EC)

- Performance Loss < 5%
- Energy minimization

## Temperature Controller (TC)

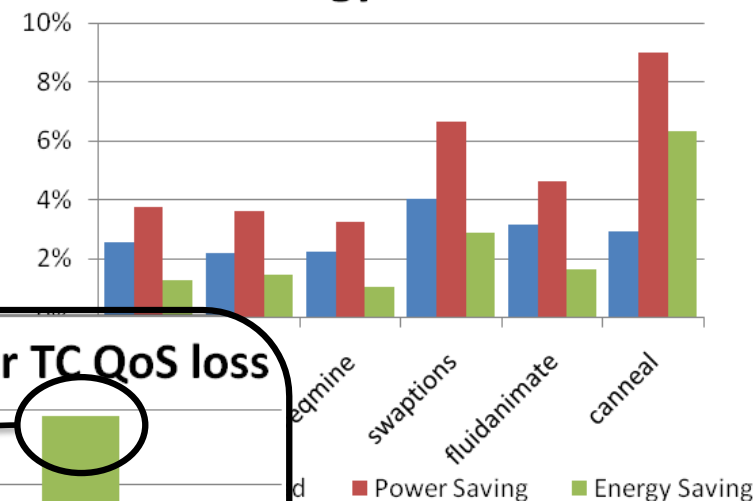- Complexity reduction
  - 2 explicit region for controller
- Performs as the centralized <3%
  - Thermal capping

# Working on Real Chips (Intel)
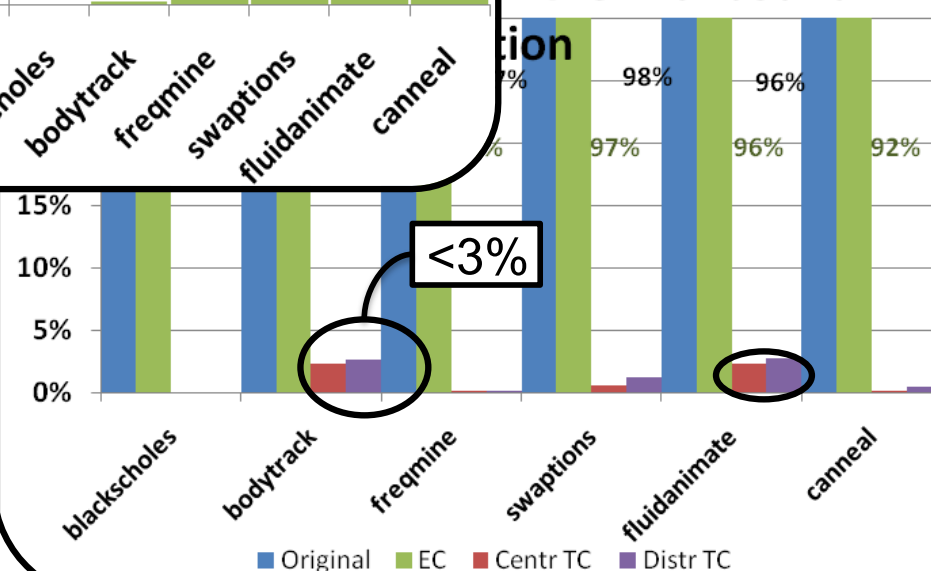


MODEL

LS

PREPROCESSING

Matlab System Identification
- N4SID
- PEM
- **LS (Levenberg-Marquardt)**

- **Temperature**
- Frequency
- Workload

..idle,idle,run,run,run,…

..0,0,1,1,1,…

Pattern Generator

PRBS.csv

DATA.csv

PREPROCESSING

XTS                    (Ts=1/10ms)

core0    core1    ........    coreN    **HW**

Workloader

# Working on Real Chips (Intel)



SUN FIRE X4270

Chipset

RAM CPU2 CPU1 RAM

Air flow

Fan board

Storage drives

- **Intel Nehalem**
- **8core/16thread**
- **2.9GHz**
- **95W TDP**
- **IPMI**

# Working on Real Chips (Intel)
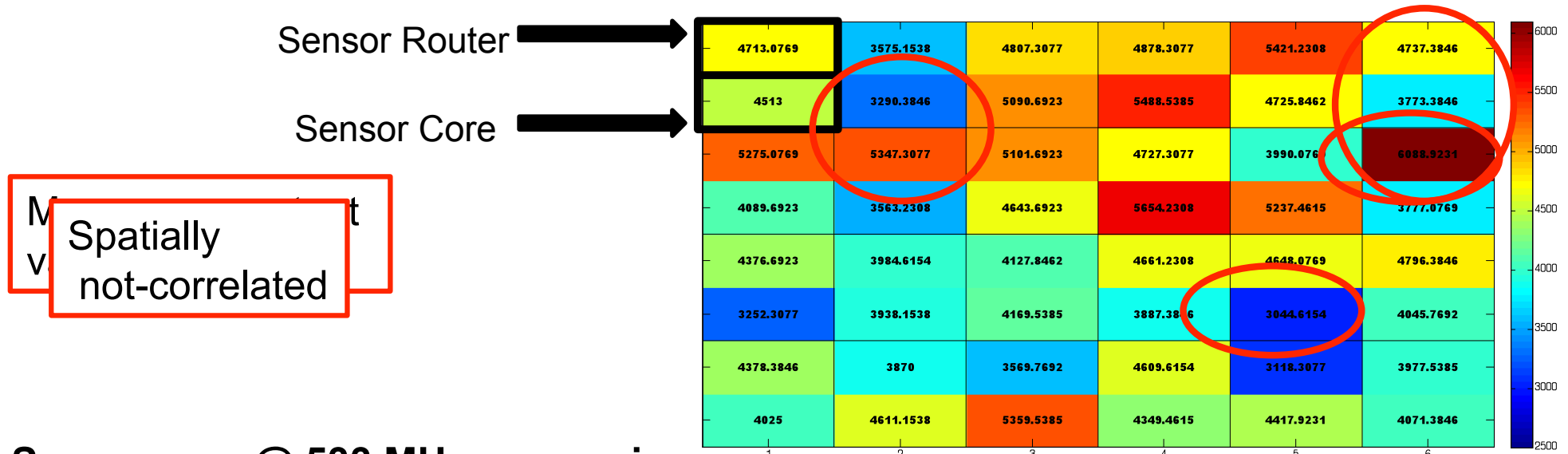
## Single Chip Cloud (45nm)



- 567.1 mm2
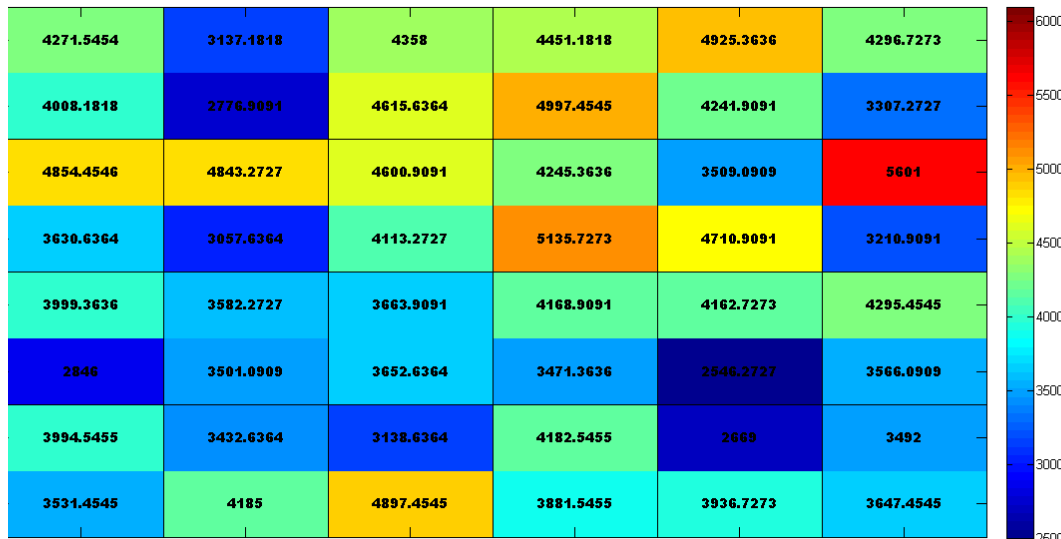- 48cores @1GHz
- 2GHz NoC
- 25-125W
- 27 (f), 8 (V) islands

# Thermal Sensor Variability



Sensor Router →

Sensor Core →

Spatially not-correlated

**Sensor map @ 533 MHz power virus**

**Sensor map @ 100 MHz idle**

# Outline

- Introduction
- Scalable Control
- Scalable model learning
- Experimental Environment
- Challenges ahead

# The 1,000 Cores Chip

- STM-CEA Platform 2012 project
  - Die with 4 16-cores tiles with L1 & L2 → few tens of mm2 (28nm)
  - SCC die → 20 of these dies: 1,280 cores
  - Thousands of Vdd, f domains
- 3D stacking currently the only technology which can provide sufficient L3 bandwidth
  - Vertical thermal dissipation!
  - Heterogeneous requirements (DRAM≠LOGIC)
- Major static and dynamic variability

# Power management Challenges

- Truly scalable algorithms → O(NlogN)

- Hardware support needed (e.g. DPM NoC)

- Cross-layer algos are needed
  – Real-time intra+inter layer communication
  – Abstraction and filtering
  – Multi-scale

- The threat of non-linearity
  – Hybrid control complexity  (MILP is NP-HARD)
  – Lack of robustness (Ill-conditioning) and stability proofs

**SoCs as complex systems (societies/markets)
→ DPM as political sociology/finance?**

# Thank you!

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA