



# Predictive Models for Thermal Modelling



Andrea Bartolini DEI, Università di Bologna <a.bartolini@unibo.it>



ALMA MATER STUDIORUM - UNIVERSITÀ DI BOLOGNA

IL PRESENTE MATERIALE È RISERVATO AL PERSONALE DELL'UNIVERSITÀ DI BOLOGNA E NON PUÒ ESSERE UTILIZZATO AI TERMINI DI LEGGE DA ALTRE PERSONE O PER FINI NON ISTITUZIONALI



- Motivation
- Static Thermal Modelling
- Dynamic Thermal Modelling
- Thermal Management
- Datacenter automation



# Motivation

#### Intel Pentium - 1 core



ALMA MATE



### Datacentre: The problem scales up

#### Tihane-2

(most powerful supercomputer 2013-2015)

- 3120000 cores,
- 17.8MW IT only => 24MW w. cooling



Marconi @Cineca - 1512 nodes - 54432 cores



A MATER STUDIORUM ~ UNIVERSITÀ DI BOLOGNA



### Datacentre: The problem scales up



MATER STUDIORUM ~ UNIVERSITÀ DI BOLOGNA

### Datacentre: The problem scales up



# **Dynamic Power**



- Quadratic  $\downarrow$  with  $\downarrow$  V<sub>dd</sub>
- Cubic  $\downarrow$  with  $\downarrow$  both V<sub>dd</sub> and f

David H. Albonesi ACACES10

# Sub-threshold Leakage Current



- Exponential  $\downarrow$  with  $\downarrow$  V<sub>gs</sub> (~V<sub>dd</sub>)
- Exponential  $\downarrow$  with  $\uparrow V_{TH}$

David H. Albonesi ACACES10



Delay:  

$$D_{p} = \frac{C_{out} V_{dd}}{I_{ON}} = \frac{C_{out} V_{dd}}{\mu(T) [V_{dd} - V_{th}(T)]^{\alpha}}$$

Carrier Mobility:

$$\mu(\mathsf{T}) = \mu(\mathsf{T}_{o}) \left( \frac{\mathsf{T}_{o}}{\mathsf{T}} \right)^{m}$$

Threshold Voltage:

$$V_{th} = V_{th}(T_0) - k(T - T_0)$$
$$T \uparrow \qquad \mu \downarrow \qquad V_{th} \downarrow$$



- For wires, the resistivity is linearly dependent from T
  - Delay increases as T increases
- For Low VT (LVT) design ( $V_{dd} >> V_{th}$ )
  - $-\mu$  dominates w.r.t. V<sub>th</sub>
  - Delay Increases as T increases
- For High VT (HVT) design (V<sub>dd</sub> ≈ V<sub>th</sub>)
  - $-V_{th}$  dominates w.r.t  $\mu$
  - Delay decreases as T increases, Indirect Temperature Dependence (ITD)

# Thermal Behavior of CMOS gates



ALMA MATER STUDIORUM ~ UNIVERSITÀ DI BOLOGNA



# **TEST** Chip



A 60 GOPS/W, -1.8 V to 0.9 V body bias ULP cluster in 28 nm UTBB FD-SOI technology, D. Rossi, et al., Solid-State Electronics, 2016

Impact of body bias on maximum frequency and leakage power over the supported range of voltage supply.vb

#### **First Multi-core** ever in RVT FDSOI



Performance



http://iis-projects.ee.ethz.ch/index.php/PULP



## **Temperature Impact**



ALMA MATER STUDIORUM ~ UNIVERSITA DI BOLOGNA





# **Thermal Impact**

#### #1 Core Rotating Power Virus

Up to 20 C Temperature difference on DIE ~ 30 C Temperature difference in between sockets - Thermal neighbours exists!





Fan Speed [RPM] Fan #



## Haswell - PowerVirus #1





# Haswell - PowerVirus #1



ALMA MATER STUDIORUM - UNIVERSITÀ DI BOLOGNA



## Haswell - PowerVirus #1





- Motivation
- Static Thermal Modelling
- Dynamic Thermal Modelling
- Thermal Management
- Datacenter automation



# **Thermal Model**









#### Core thermal response

Dynamic Model  $T[n+1] = A \cdot T[n] + B \cdot P[n]$ Steady-state condition: T[n+1] = T[n] = TStatic thermal model:  $T = SG \cdot P$ Steady-state Gain matrix:  $SG = (I - A)^{-1} \cdot B$ 



# Platform



#### SUN FIRE X4270

- Intel Nehalem 5500
- 8core/16thread
- 1.6÷2.9GHz
- 95W TDP
- IPMI

ALMA MATER STUDIORUM - UNIVERSITÀ DI BOLOGNA



# Test#1 - Power Mode

Power Data acquisition methodology



- Real multicore platform (4 core per CPU)
- System level power measurements only
- Benchmark based data profiling





Power Data acquisition methodology

#### Benchmark: Per core Workload allocation

- "1" = Core in fully busy
   state = Maximum core power consumption (PowerVirus process)
- "0" = Idle state = minimum core power consumption
- M<sub>coreON</sub> = #Active cores



## Test#1 - Power Mod





## Test#1 - Power Mod





# PM: breakthrough





## Test#2: Thermal Model





# TM: Data Fitting





### Results



ALMA MATER STUDIORUM - UNIVERSITÀ DI BOLOGNA



- Motivation
- Static Thermal Modelling
- Dynamic Thermal Modelling
- Scalable Thermal Modelling
- Thermal Management
- Datacenter automation



#### **Model Structure**



ALMA MATER STUDIORUM ~ UNIVERSITÀ DI BOLOGNA



#### **LS System Identification**











#### **Thermal Modeling**

i7 Server Platform – 4 cores

Ts = 1ms - Quantizzation noise

<u>Step response:</u>

$$y(t) = K_0 - \sum_{i=1}^{n} k_i \cdot e^{-p_i \cdot t}$$







#### **Thermal Modeling**

i7 Server Platform – 4 cores

Ts = 1ms - Quantizzation noise

<u>Step response:</u>  $y(t) = K_0 - \sum_{i=1}^{n} k_i \cdot e^{-p_i \cdot t}$ 

Black Box Model Estimation Performance

Order	1	2	3
SSE	1,02E+05	2,00E+04	2,59E+04
Time elapsed (s)	300,15	2434,24	12681,57



<u>Black Box</u>

<u>Our approach</u> LS + physical constraints


## **Thermal Model Identification**

Single chip cloud computer SCC – 48 cores – 24 tiles





## **SID Standard ARX**









# **Thermal Model Identification**

#### Why ARX is failing?

- Standard ARX considers white innovation
- Its impact on the output temperature get coloured by the model poles
- In our system we know we have a white additive noise in the temperature sensors
  - Two different noise sources
    - Trying to model it with just the innovation noise source biases the model parameters
      - Noise in the output cannot just be filtered out as it is unknown cannot power down a CPU and access the internal temperature sensors

w(t)

u(t)

 $\overline{A(z^{-1})}$ 

 $B(z^{-1})$ 

 $\overline{A(z^{-1})}$ 

T(t)



**ARX MO** 

Not realist



#### **Noisy Temperature Sensors**

•  $T(t) = \overline{T}(t) + v(t)$ 



















## **Thermal Model Identification**

Single chip cloud computer SCC – 48 cores – 24 tiles

#### Ts = 100ms - Measurment noise

#### Standard ARX:

- Designed only for process noise !
- Measurement noise induces biases !

#### Bias Compensated ARX

(Diversi et al., 2013a, 2014)

 Iterativelly estimate the noise variance and compensate it in the LS

#### <u>Dynamic Frisch</u>

(Diversi et al., 2013b)

- Searches solutions compatible with the covariance matrix of the noisy data
- Selection based on a set of low-order and high-order Yule-Walker equations



#### [TCAS 2014] [DATE13 BPA]



## **Thermal Model Identification**



#### [TCAS 2014] [DATE13 BPA]



## Galileo Modelling





- Motivation
- Static Thermal Modelling
- Dynamic Thermal Modelling
- Thermal Management
- Datacenter automation



# Multitherman Holistic Approach

FP7 ERC Advance MULTITHERMAN: Multiscale Thermal Management of Computing Systems PI: Prof Luca Benini

Ы





*Timescale: ~msec, ~sec* 

Built-in sensors and actuators for feedback control of P & T

- 1. Temperature sensors (core)
- 2. Architectural utilization (core)
- 3. Power monitors (cpu)
- 4. Clock Frequency (core)
- 5. Shutdown (core)

**Optimal Control Goal:** 

Maximize core's performance (i.e. frequency ) while constraining maximum temperature

Timescale: ~hours







#### **Chip Level Control**





### **Thermal Controller**





#### **MPC Scalability**





#### **Addressing Scalability**







Complexity

#### Model Predictive Distributed Control

(Bartolini et al., 2013) — (Tilli et al., 2012) — (Tilli et al., 2015) -



a)	Centralized MPC Complexity		b) <i>Di</i>	omplexity		
	MPC_explicit	MPC_implicit		()		
4	81	7,70	f2P		0,061	time (us)
8	6561	9,00	Obs	server	0,743	time (us)
16	OUT	24,20	MPO	C (Impl)	4,690	time (us)
48	OUT	85,50	MPO	C (Expl)	2	<i># regions</i>
	# regions	time (us)	P2f		1,188	time (us)



- Motivation
- Static Thermal Modelling
- Dynamic Thermal Modelling
- Thermal Management
- Datacenter automation



#### **A New Trend: Datacentre Automation**





#### The Big Data & DL backbone

Computing clusters

- Not only the computing engine of Big Data solution
- Also a complex industrial plant and a growing industry in ER
- A compute nodes can produces ~ 100/1000 metrics/s \* "peta/exa scale" = Big Data!

Datacenter automation – improve energy/cost efficiency and effectiveness – Industry 4.0 thanks to:

- Live collection and processing of large telemetry data (>100GB/day x cluster)
- On-line generation of "plant models" a.k.a. digital tweens", security break detection









EXA

# **ExaMon: an Industy 4.0 approach to datacenter automation**



#### **Front-end**

- Host: management
  node
- Docker containers
- ~45KS/s

#### **Back-end**

 MQTT–enabled monitoring agents (e.g. Dig)



# **ExaMon: an Industy 4.0 approach to datacenter automation**





# **ExaMon: an Industy 4.0 approach to datacenter automation**





#### Broker:

 Forward data to the listeners (e.g. kairosDB)

#### Mqtt2kairosdb:

- Interface between MQTT and KairosDB
- KairosDB is a front-end to handle time series in Cassandra

#### Cassandra:

- NoSQL database
- Highly scalable
- Optimized to balance the load on multiple nodes



# ExaMon: an Industy 4.0 approach to datacenter automation



#### **Application layer:**

- Grafana, Apache Spark, etc …
- Aggregate metrics for
  Data Visualization, ML
  Analysis, Post Processing,
  etc ...





# **ExaMon: Scalable Data Collection and Analytics**



{Key,Value} = TS, Measurement Topic = /davide/node1/Metric





#### **ExaMon: Batch & Streaming**





## Use Case: Online Power Model Learning

 Idea: Coupling the monitoring framework with Apache Spark to calculate CPU power model parameters

Power Model:

$$\mathbf{P}_{pkg0+1} = \sum_{i=0}^{N_c-1} (a_i + b_i \mathbf{IPS}_i) \mathbf{f}_i$$

Spark MLlib: Streaming linear regression with SGD  $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$   $\mathbf{w} = \begin{bmatrix} a_0, ..., a_{N_c}, b_0, ..., b_{N_c} \end{bmatrix}$   $\mathbf{w}^{(t+\bar{1})} := \mathbf{w}^{(t)} - \gamma f'_{\mathbf{w}, \mathbf{i}}$ Spark





18 September 2018

## se Case: Online Power Model Learning



# Use Case: Galileo (CINECA) setup



- 516 nodes
- 430 metrics/node



- 221880 total metrics
- •144 Mb/s
  - @ 2 seconds sampling time


# Use Case: Results



- Actual Power vs Model prediction (CPU0 + CPU1)
  - Transition between two different workload phases running on the node
  - Online algorithm promptly recovers, learning a new model
  - 5-6 samples recovery time)
- Average model accuracy:

– ±25mW

- Model update time (iteration):
  - 600ms

Algorithm returned satisfactory performance!

# Thermal modeling at scale





# Problem with training





## How to evaluate goodness



#### Experts features not good!

Deep Learning approaches beats expert user in determining the «good» windows



### D.A.V.I.D.E. (#18 Green500 Nov'17)

E4 COMPUTER ENGINEERING

D.A.V.I.D.E. SUPERCOMPUTER (Development of an Added Value Infrastructure Designed in Europe)













#### DiG: High Resolution Out-of-band Power Monitoring

- Out-of-band  $\rightarrow$  Zero overhead
- Collect more than 1.5 kS/s, 7/7d, 24/24h, for all users
- Architecture independent (i.e. tested on Intel, ARM and IBM)
- Fine grain  $\rightarrow$  down to ms scale (sampling @800 kS/s + avg)
- IoT communication technology (MQTT)  $\rightarrow$  scalable
- Edge Computing!





- Machine learning and DL is a key tool for solving engineering problems!
- We are now able to design predictive control loop on supercomputing processors
- Future works: Datacenter automation and anomaly detection!

# Pubblications – Naif Modelling

- Davide Rossi, Antonio Pullini, Igor Loi, Michael Gautschi, Frank K. Gürkaynak, Andrea Bartolini, Philippe Flatresse, Luca Benini, A 60 GOPS/W, -1.8 V to 0.9 V body bias ULP cluster in 28 nm UTBB FD-SOI technology, JSSE15
- Conficoni C.; Bartolini A.; Tilli A.; Tecchiolli G.; Benini L.; Energy-Aware Cooling for Hot-Water Cooled Supercomputers, DATE 15
- Conficoni, Christian, et al. "HPC Cooling: A Flexible Modeling Tool for Effective Design and Management." IEEE Transactions on Sustainable Computing (2018).
- Bartolini, Andrea, et al. "The DAVIDE Big-Data-Powered Fine-Grain Power and Performance Monitoring Support." (2018).
- Beneventi, Francesco, et al. "Continuous learning of HPC infrastructure models using big data analytics and in-memory processing tools." 2017 Design, Automation & Test in Europe Conference & Exhibition (DATE). IEEE, 2017.
- Bartolini A.; Cacciari M.; Cavazzoni C.; Tecchiolli G.; Benini L.; **Unveiling Eurora Thermal and** power characterization of the most energy-efficient supercomputer in the world, DATE14
- Fraternali, Francesco, et al. "Quantifying the impact of variability and heterogeneity on the energy efficiency for a next-generation ultra-green supercomputer." IEEE Transactions on Parallel and Distributed Systems 29.7 (2018): 1575-1588.
- F. Fraternali, A. Bartolini, C. Cavazzoni, G. Tecchiolli, L. Benini, **Quantifying the impact of** variability on the energy efficiency for a next-generation ultra-green supercomputer; ISLPED14



# **Pubblications Thermal**

- Beneventi F. ; Bartolini A. ; Tilli A. ; Benini L., An Effective Gray-Box Identification Procedure for Multicore Thermal Modeling TC14
- Diversi R.; Tilli A.; Bartolini A.; Beneventi F.; Benini L.; Bias-Compensated Least Squares Identification of Distributed Thermal Models for Many-Core Systems-on-Chip, TCAS 14
- Diversi, R; Bartolini, A.; Tilli, A.; Beneventi, F.; Benini, L.; ,"SCC thermal model identification via advanced bias-compensated least-squares", DATE13 BPA
- Bartolini, A., Diversi, R., Cesarini, D., & Beneventi, F.. Self-Aware Thermal Management for High Performance Computing Processors. IEEE Design & Test.
- Diversi, Roberto, et al. "Thermal model identification of supercomputing nodes in production environment." Industrial Electronics Society, IECON 2016-42nd Annual Conference of the IEEE. IEEE, 2016.
- Bartolini, Andrea, et al. "Multiscale Thermal Management of Computing Systems-The MULTITHERMAN approach." IFAC-PapersOnLine 50.1 (2017): 6709-6716.
- Beneventi, Francesco, et al. "Static thermal model learning for high-performance multicore servers." Computer Communications and Networks (ICCCN), 2011 Proceedings of 20th International Conference on. IEEE, 2011.
- Beneventi, Francesco, et al. "Cooling-aware node-level task allocation for next-generation green HPC systems." High Performance Computing & Simulation (HPCS), 2016 International Conference on. IEEE, 2016.
- Bartolini, Andrea, et al. "Thermal and energy management of high-performance multicores: Distributed and self-calibrating model-predictive controller." IEEE Transactions on Parallel and Distributed Systems 24.1 (2013): 170-183.



# **Pubblications - NN**

- Netti, A., Kiziltan, Z., Babaoglu, O., Sirbu, A., Bartolini, A., & Borghesi, A. (2018). FINJ: A Fault Injection Tool for HPC Systems. arXiv preprint arXiv:1807.10056.
- Lombardi, Michele, Michela Milano, and Andrea Bartolini. "Empirical decision model learning." Artificial Intelligence 244 (2017): 343-367.
- Bartolini, Andrea, et al. "Optimization and Controlled Systems: A Case Study on Thermal Aware Workload Dispatching." AAAI. 2012.



# Pubblications – big data

- Bartolini, Andrea, et al. "The DAVIDE Big-Data-Powered Fine-Grain Power and Performance Monitoring Support." (2018).
- Beneventi, Francesco, et al. "Continuous learning of HPC infrastructure models using big data analytics and in-memory processing tools." 2017 Design, Automation & Test in Europe Conference & Exhibition (DATE). IEEE, 2017.



• HPC team:

 Luca Benini, Francesco Beneventi, Antonio Libri, Andrea Borghesi, Federico Pittino, Alessandro Petrella, Daniele Cesarini, Christian Conficoni, Roberto Diversi, Andrea Tilli, Michele Lombardi, Michela Milano