

Hot Spots and Core-to-Core Thermal Coupling in Future Multi-Core Architectures

M. Janicki¹, J.H. Collet², A. Louri³ and A. Napieralski¹

¹ Department of Electronics and Computer Science, Technical University of Łódź
Wólczajska 221/223, 90-924 Łódź, Poland

² Laboratoire d'Analyse et d'Architecture des Systèmes - CNRS, Université de Toulouse, UPS, INSA, INP, ISAE
7 avenue du colonel Roche, F-31077 Toulouse, France

³ Department of Electrical and Computer Engineering, University of Arizona, Tucson, AZ - 85721, USA
E-mails: janicki@dmcs.pl, jacques.collet@laas.fr, louri@ece.arizona.edu, napier@dmcs.pl

Abstract

This paper studies hot spot and thermal coupling problems in future multicore architectures as CMOS technology scales from 65 nm feature size to 15 nm. We demonstrate that the thermal coupling between neighboring cores will dramatically increase as the technology scales to smaller feature sizes. The simulation studies were based on solving the heat equation using the analytical Green's function method. Our simulations indicate that the thermal coupling in the 15 nm feature size just after 100 ms of operation will increase from 20 % to 42 % and in the steady state might reach even 65 %. This finding uncovers a major challenge for the design of future multi-core architectures as the technology keeps scaling down. This will require a holistic approach to the design of future multi-core architectures encompassing low power computing, thermal management technologies and workload distribution.

1. Introduction

Unsustainable power consumption and ever-increasing design complexity have pushed the microprocessor industry to move away from designing single complex monolithic processing core to multiple cores on a single chip. Today there is a wide consensus, both in industry and academia that multi-core chips, also called chip multiprocessors (CMPs) are the only efficient way for utilizing the billions of transistors resulting from the continued scaling of technology. These CMP architectures are expected to rely on the full exploitation of parallelism to achieve further increases in performance. Almost all major microprocessor vendors, including Intel, IBM, HP, Sun, are currently offering a family of multicore chips. It is further predicted that we will move from multi-core to many-core chips with more than thousand cores per chip [1]. As chip geometry shrinks, transistor density increases, and clock frequencies rise, the transistor leakage current increases, leading to excessive power consumption and heat generation. It is widely agreed upon that power consumption and heat will be a major challenge for the future of CMPs. Extensive research efforts are under way to determine the power-performance relationship in light of future technology scaling.

While many investigations have studied so far the thermal implications of multi-core architectures in the current existing technologies [2]-[5], to the best of our knowledge, this is the first study that investigates the core-to-core thermal coupling in future technologies, scaling the feature size from

the current 65 nm down to 15 nm. From the time- and space-resolved we show several important results in the small (15 nm) feature size, namely:

- 1) in the 15-nm process the dynamic thermal coupling between adjacent cores occurs within shorter time intervals, typically after 10 ms of operation when one starts two adjacent cores simultaneously at full power,
- 2) the temperature overhead of each core may reach even 65% (with respect to the uncoupled thermal regime) typically after 1 second of operation,
- 3) the leakage power is the major factor responsible for most of the temperature rise,
- 4) the advances in cooling techniques may lower temperature, but they may not be sufficient to reduce the dynamic thermal coupling between the cores.

Throughout this work, we consider multi-core chips based on the replication of an Alpha-like architecture. Section 2 provides a brief description of the floorplan and geometry and discusses the significant influence of leakage power on chip temperature. Section 3 details different thermal problems associated with technology scaling in multicore chips from 65 down to 15 nm, especially the increase of the core-to-core thermal coupling. Thermal simulations were carried out solving the heat equation using the analytical Green's function method [6].

2. Benchmark Geometry

The multi-core chips considered in the paper are built from the replication of the one-core Alpha microarchitecture shown in Figure 1, which is often used in literature as a standard benchmark [2]-[3]. This basic architecture contains Alpha-like cores and two levels of cache memory. We initially start with the 65 nm technology node where the core area is 9 mm² and L2 cache is assumed to be 4 MB with an area slightly larger than 50 mm². The core layout visible in Figure 1 shows only the components dissipating the most power, which are the Floating Point Unit (1 adder, 1 multiplier, 1 register), the INTegeR Unit (4 ALUs, 1 register), the Branch Target Buffer (BTB), and finally the Data and Instruction L1 caches, which exchange 32-byte blocks with the L2 cache.

The power consumption data necessary for the simulations were taken from [3]. Thermal simulations were performed for a case when the core operates at the full load and the L2 cache dissipates the static power. The full load is understood here as the maximal active power (both static and

dynamic), which is specified in [3] for the supply voltage 0.9 V and the clock frequency 5 GHz in temperature 100 °C. Then, the total chip power is 20 W, with only one third of it dissipated in the core. This might be surprising, though quite justified because, due to the very thin gate dielectric layer, the 65 nm technology had the highest leakage currents.

These currents are certainly very important in the L2 cache whereas the dynamic power dissipated in it (i.e., the frequency of READ/WRITE operations) critically depends on the miss rate (MR) in the L1 cache. Obviously, the MR depends on the application specificity (especially on the locality of data) and on the internal structure of L1 caches [7]. However, previous studies have shown that the MR is typically of the order of a few percent (for the 64-kB L1 caches considered here) when executing the SPEC92 benchmark suite on a DEC Station [8]. Thus, the statement that the L2 cache dissipates mainly static power is justified.

The thermal model assumed that the chip is mainly cooled from the back side and conduction cooled on the other one, which is typical for flip-chip processor assemblies. The heat transfer rate at the back surface was adjusted to model the presence of a heat sink having thermal resistance of 0.6 K/W for a 12 mm die. Owing to the use of the Green's function solution method, it was possible to compute the temperature map in 10,000 locations in less than 1 minute.

This method is quite flexible and for simple geometries, e.g. for multilayered slabs, allows the computation of thermal influence coefficients linking power dissipation to temperature rise only for selected time instants and locations in a structure. The main advantage over numerical methods is that it makes it possible to compute the temperature map for different heat source configurations without re-solving the heat equation. Evidently, for complex geometries numerical methods remain the only possible choice.

The steady state temperature map obtained for the layout from Figure 1 is presented in Figure 2. Here, the steady state is understood as the situation when the chip is operated at full load and it is powered long enough so that all the thermal processes are already in equilibrium. The exact core location is marked by the black square. The vertical bar on the right-

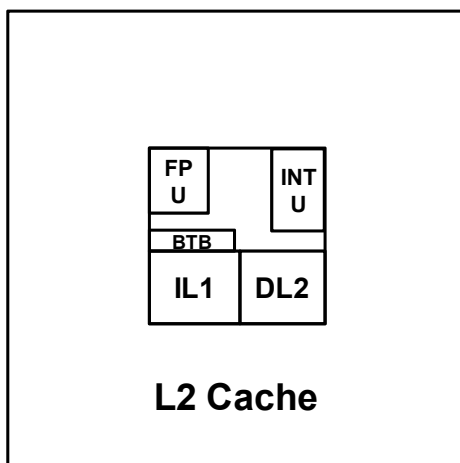


Figure 1: Layout of the Alpha architecture.

hand side of the figure shows the mapping between colors and local temperatures. As expected, the hottest components are the floating point units and the integer ALUs. To illustrate the possible advantages coming from the use of multiple cores, the simulations were repeated for the architecture with 4 cores operating at their maximal processing speed and located in the middle of each quarter of chip area.

The temperature rise map obtained in this case is presented in Figure 3. The core layout used in this architecture is not only optimal from the thermal point of view but it corresponds also to the one which is expected to be replicated in future architectures containing hundreds of cores. The total power dissipation grew to 40 W and the chip area increased by 45 % due to the introduction of three additional cores. This caused the chip temperature increase by 12-14 K, but this architecture is capable of much better performance than the one-core one discussed in the beginning.

Analyzing the results in more detail, one can say that the cache leakage is the dominant factor responsible for 23-26 K of the total temperature rise. Thus, we consider that data given in [3] are significantly overestimated and having in view the analysis of circuit scaling throughout future technologies, we decided to reduce the leakage power 10 times. As a result, the total power dropped to 25.9 W, of which only 1.4 W is the leakage power. These values will be used and scaled in all the subsequent simulations, as explained in the following section.

The proposed leakage power reduction in fact corresponds to the real situation, when the 65 nm technology was replaced by the high-dielectric dual metal gate 45 nm technology [9]. Then, as reported in [10], the subthreshold leakage and the gate leakage power were reduced by more than 5 and 10 times respectively. The steady state temperature rise map obtained for the reduced leakage power is presented in Figure 4. Now, the dynamic power is not any longer dominated by the L2 cache leakage and the leakage contributes only 4% of the total

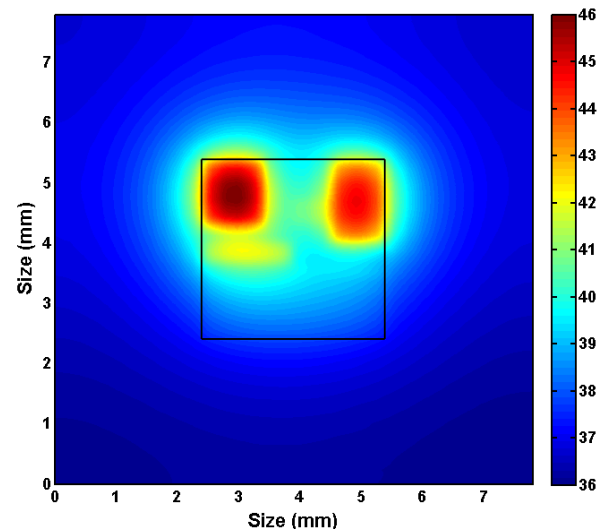


Figure 2: Temperature rise map for the layout from Figure 1.

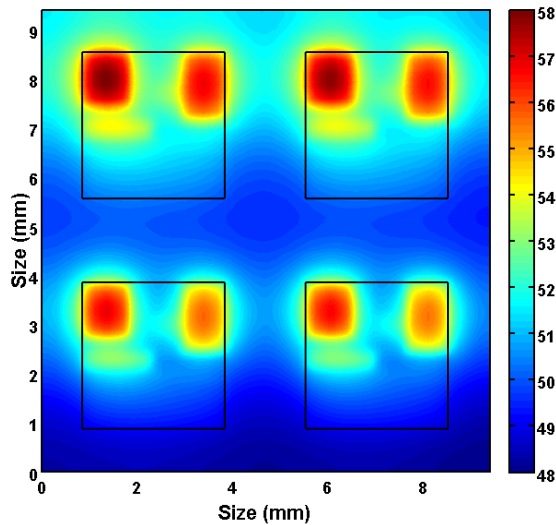


Figure 3: Temperature rise map for 4 cores with high leakage.

temperature rise. It is important to stress that the reduction of leakage power does not affect the subsequent discussion on the thermal core coupling, but in our opinion it brings the proportions between the dynamic and leakage power closer to reality when high- k dielectrics are used.

3. Technology Migration

This section is devoted to the investigation of the thermal problems induced by the migration of the multicore design to future technology nodes. First, we study the evolution of the steady state temperature map throughout the technologies. In particular, the contribution of each power component to the total temperature rise is analyzed. Then, we study the core-to-core dynamic coupling effects. Additionally, we include the discussion of chip cooling influence on the coupling.

In the experiment it was assumed that all the dimensions of the Alpha architecture (see Figure 1) in the next technology are scaled by the standard factor of 0.7 and consequently all the feature areas are halved. We will present results obtained for the following technology scales: 65 nm, 45 nm, 32 nm, 22 nm and 15 nm. Equally, or even more important, issue for the technology migration is the choice of a particular power scaling method. Obviously, moving to a next technology node one would like to improve circuit performance and keep chip temperature at the same level. Theoretically, the increase of temperature can be prevented if constant power density, i.e. the heat flux, is maintained during the scaling. However, this is rarely possible or desired.

Actually, dissipated power has two major components: the static (leakage) power and the dynamic (active) power. The first component, consisting of the subthreshold leakage and the gate leakage, is strongly technology related and designers have little influence on it. The other one is linked to the chip activity and can be shaped according to particular needs, but usually it is scaled aggressively so that to improve processor performance.

Technology node (nm)	Max. temp. rise (K)	Min. temp. rise (K)	Av. temp. rise (K)
65	42.5	27.9	33.0
45	41.3	31.7	35.1
32	43.4	37.1	39.4
22	50.6	46.6	48.0
15	67.2	64.9	65.6

Table 1: Evolution of steady state temperatures.

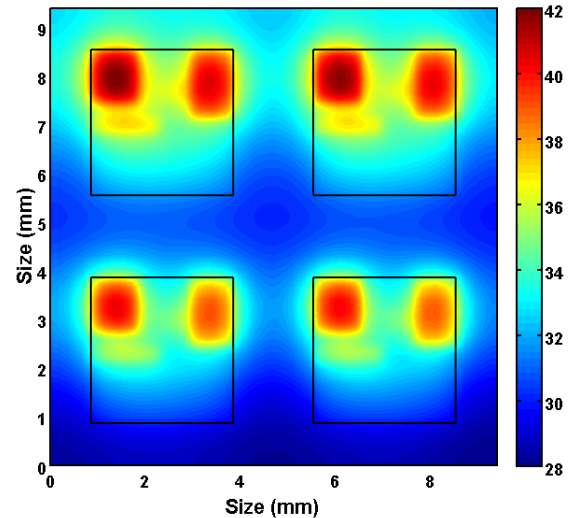


Figure 4: Temperature rise map for 4 cores with low leakage.

According to the ITRS report [11], leakage power tends to scale exponentially unless some technological breakthrough is achieved or a novel device developed. Thus, we believe that it is reasonable when this power component is scaled at the constant power rate. This assumption is quite optimistic and will probably require the introduction of fully depleted devices with gate dielectrics of the relative permittivity higher than 100, but it is a feasible task. The scaling scenario adopted here for active power is more optimistic because we assume that owing to the combined influence of lower power supply and capacitances as well as moderate frequency scaling, it should be possible to achieve the scaling at constant power density.

The above-described strategy for power scaling resulted in the decrease of total power dissipated in the chip from 26 W in the 65 nm technology to only 3 W in the 15 nm technology; always with the cores processing at full load as it was defined in Section 2. The maximal hot spot temperature (the floating point unit of the top left core in Figure 4), the minimal surface temperature (L2 cache at the bottom) and the average steady state surface temperature rise computed for all the considered technologies are presented in Table 1. As can be seen from the table, the temperature profile gets ever flatter throughout the technologies and the average surface temperature rise almost doubles. The particularly dramatic increase of all presented temperature values is observed for the last two technologies. The cause for this is clearly visible in Table 2 presenting the individual contributions of different

power components to the total temperature rise in the hot spot location. It is evident that the temperature rise again becomes dominated by the leakage, which is scaled at constant power.

Technology node (nm)	L2 cache leakage (%)	Core self-heating (%)	Core mutual heating (%)
65	3.5	57.4	39.1
45	7.8	46.5	45.7
32	15.4	36.9	47.7
22	27.4	27.7	44.9
15	42.4	19.4	38.2

Table 2: Contribution to the temperature rise.

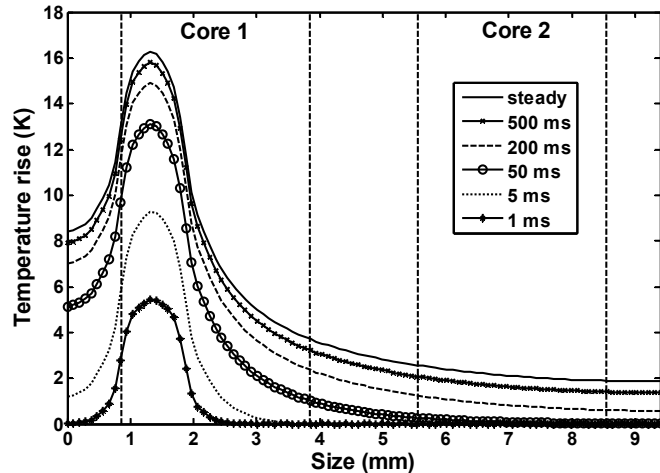


Figure 5: Heat diffusion – 65 nm technology.

However, the most striking fact is that already in the 45 nm technology less than half of the total temperature rise in the hot spot is due to the power dissipation in a core itself and the remaining part is generated by the cache leakage and heating by other cores. This situation becomes ever worse and in the 15 nm technology 80 % of the temperature rise is due to the external factors. In particular, the heating by other cores grows drastically, what leads to the situation when in each core the neighboring ones contribute twice as much of the temperature rise.

The analysis presented so far was based only on the steady state temperature maps, which were influenced by leakage. However, more information on the problem of core thermal coupling can be gained from the analysis of heat diffusion processes when only dynamic power is dissipated in the cores. Thus, in order to emphasize the dynamics of the heat transfer, in what follows we consider that the processes related to the leakage are already in thermal equilibrium. In this context, we investigate the impact of technology scaling on the heat diffusion times and the dynamic thermal coupling between adjacent cores. The simulated kinetics of the heat diffusion process when the floating point unit of a core suddenly starts to operate at the full load are shown in Figures 5-6 for the two limit technologies, i.e. 65 nm and 15 nm.

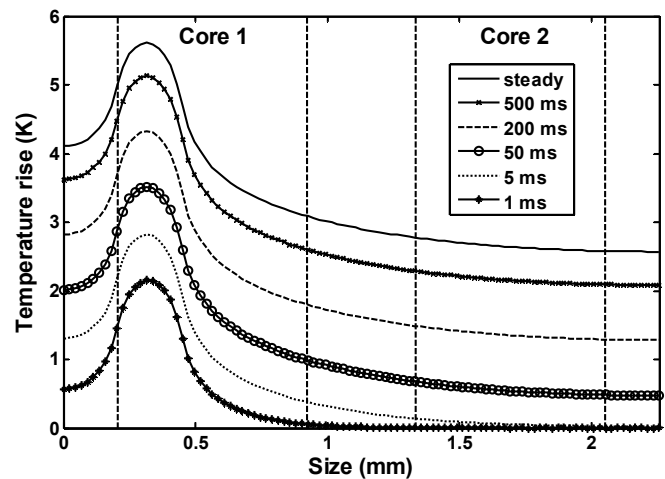


Figure 6: Heat diffusion – 15 nm technology.

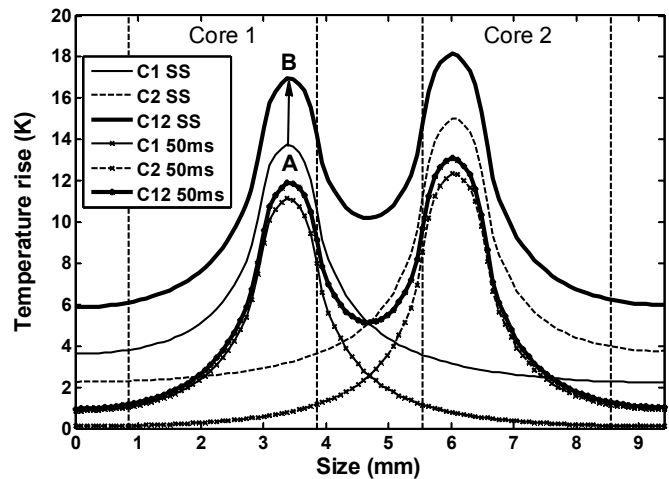


Figure 7: Core thermal coupling – 65 nm technology.

The dashed vertical lines mark the locations of Core 1, where power is dissipated, and the closest neighboring Core 2. These two figures differ substantially and directly show the increase of the thermal coupling between cores. Indeed, in the 65-nm process the power dissipation in Core 1 has virtually no influence on the temperature of Core 2 during the first 50 ms, whereas in the 15-nm technology the temperature coupling appears much faster, already after 5 ms. This is caused mainly by the fact that the core spacing decreased substantially from 1.7 mm to 0.4 mm.

The problem of the dynamic thermal coupling was further investigated in the simulations where the closest units in the adjacent cores, i.e. the integer unit in Core 1 and the floating point unit in Core 2, were started simultaneously at their full processing power. The simulated temperature profiles in this case for the two earlier mentioned technologies are visualized in Figures 7-8. The figures show the individual contributions of each core and the total temperature rise (thick lines) for the steady state and after 50 ms of operation. For the purposes of the core dynamic thermal analysis, the measure of the coupling was defined here as the ratio of the difference of the total temperatures rise of Core 1, T_B , and the temperature rise of Core 1 caused by self-heating, T_A , related

to the latter value, i.e. the core coupling is expressed as $(T_B - T_A) / T_A$ (see the figures).

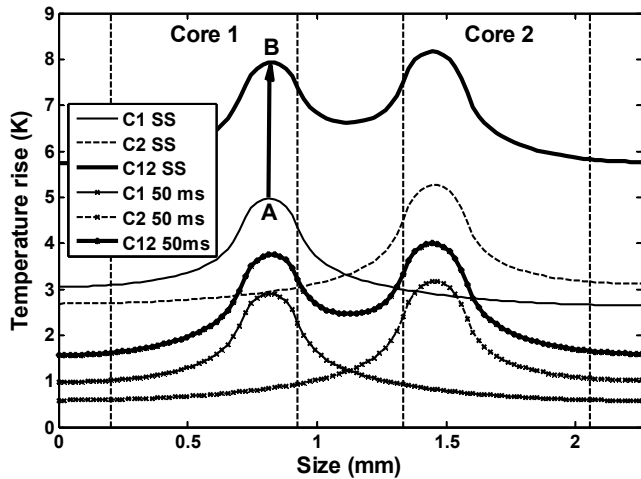


Figure 8: Core thermal coupling – 15 nm technology.

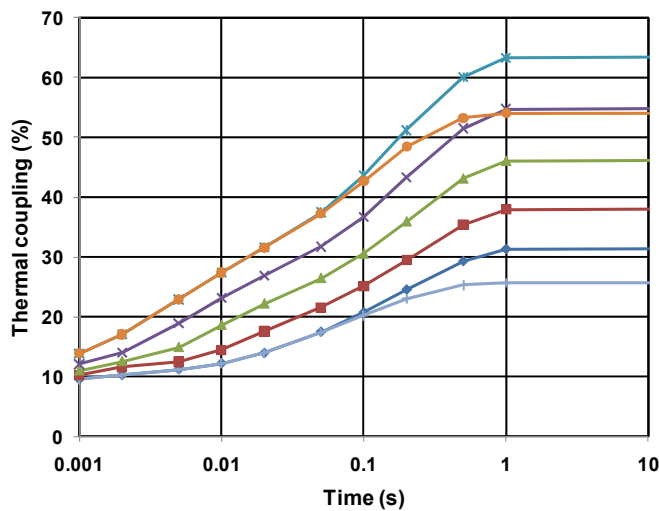


Figure 9: Kinetics of core-to-core thermal coupling.
 ◆ 65 nm; □ 65 nm cool; ■ 45 nm; ▲ 32 nm,
 ◆ 22 nm; * 15 nm; ● 15 nm cool.

The core coupling values computed using this formula for all the technologies and different instants are plotted against the logarithmic time in Figure 9. As can be seen in the figure, the core thermal coupling will increase in future technologies and in the technology range considered here it will at least double for all the heat diffusion times larger than a few milliseconds. For the steady state, the core thermal coupling increased from 32 % to 64 %. However, the most important increase of the coupling is observed for the diffusion times between 10 ms and 200 ms, which is important from the parallel computing point of view.

Theoretically, the reduction of the core thermal coupling could be achieved through the improvement of chip cooling, because this would force the heat to flow more vertically and consequently make the temperature profiles steeper. Thus, the last simulations presented in this paper consisted in investigating of the influence of chip cooling on the core coupling. For this purpose, the heat transfer coefficient at the back side of the chip was doubled, what

corresponds to the thermal resistance of 0.3 K/W for a 12 mm die. The results obtained for the two limit technologies are compared in Figure 9 with the previous simulations. As can be seen, the improvement of chip cooling indeed lowered the coupling in the steady state, but the values for the time intervals shorter than 100 ms remained the same. This is caused by the fact that diffusing heat before it reaches the heat sink does not ‘sense’ the change of cooling conditions and as a result the dynamic coupling is not affected by the cooling.

The last comment on the presented results concerns the decrease of the core temperature during scaling observed in Figures 5-8. Theoretically, when scaling with the constant power density, the maximal temperature rise should remain the same, but here the vertical dimensions are not scaled, so the heat spreading angle changes and the situation approaches the one-dimensional heat flow, hence the temperature rise drops. When scaling is done at constant chip area and the number of cores is increased, the core temperature rise should not decrease.

4. Conclusions

This paper presented a study of the thermal issues that will occur due to the reduction of dimensions in the forthcoming silicon technologies. The significant increase of the thermal coupling between neighboring cores revealed by the presented simulations will aggravate the occurrence of hot spots in the future nanoscale technologies. Moreover, the advances in the cooling technology will lower chip temperature, but they may not alleviate the problems due to the dynamic core coupling.

Acknowledgments

The participation in the conference was supported by the EU European Social Fund project POKL.04.01.01-00-179/08-00. This work was also partially supported by the US NSF Grant number CFF-0915537.

References

1. Borkar, “Thousand Core Chips - Technology Perspective”, Design Automation Conference, pp. 746-749, 2007.
2. Skadron, K., M. Stan, K. Sankaranarayanan, W Huang, S. Velusamy, and D. Tarjan, “Temperature-aware Micro-architecture: Modeling and Implementation”, ACM Trans. on Arch. and Code Optim., Vol. 1, pp. 94-125, 2004.
3. Liao, W., L. He, K.M. Lepak, “Temperature and Supply Voltage Aware Performance & Power Modeling at Micro-architecture level”, IEEE Trans. CAD Integrated Circuits and Systems, Vol. 24, pp. 1042-1053, 2005.
4. Monchiero, R. Canal, A. Gonzalez, “Power/Performance/Thermal Design Space Exploration for Multicore Architectures”, IEEE Trans. on Parallel and Distributed Systems, Vol. 19, No. 5, pp. 666-681, 2008.
5. Li J., J. F. Martínez, “Power-Performance Considerations of Parallel Computing on Chip Multiprocessors”, ACM Trans. Arch. and Code Optim., Vol. 2 , pp. 397-422, 2005.

6. Janicki, G. De Mey, A. Napieralski, "Thermal Analysis of Layered Electronic Circuits with Green's Functions", *Microelectronics Journal*, Vol. 38, pp. 177-184, 2007.
7. Patterson, D.A., J.L. Hennessy, "Computer Architecture: A Quantitative Approach", Morgan Kaufmann Publishers, see chapter 5 of Second Edition, 1996.
8. Gee, J.D., M.D. Hill, D.N. Pnevmatikakos, and A.J. Smith, "Cache performance on the SPEC92 benchmark suite", *IEEE Micro*, pp. 278-283, 1993.
9. Chau, R., J. Brask, S. Datta, G. Dewey, M. Doczy, B. Doyle, J. Kavalieros, B. Jin, M. Metz, A. Majumdar, and M. Radosavljevic, "Application of High-k Gate Dielectrics and Metal Gate Electrodes to Enable Silicon and Non-silicon Logic Nanotechnology", *Microelectronic Engineering*, Vol. 80, pp. 1-6, 2005.
10. Varghese, S. Jahagirdar, C. Tong, K. Smits, S. Damaraju, S. Siers, V. Naydenov, T. Khondker, S. Sarkar, P. Singh, "Penryn: 45nm Next Generation Intel® Core™ 2 Processor" in *Proc. of IEEE Asian Solid-State Circuits Conference*, pp. 14-17, 2007.
11. International Roadmap for Semiconductors available at: <http://www.itrs.net/reports.html>