

# Distributed Peak Power Management for Many-core Architectures

John Sartori and Rakesh Kumar

Coordinated Science Laboratory

1308 West Main St

Urbana, IL 61801

## Abstract

Recently proposed techniques for peak power management [4] involve centralized decision-making and assume quick evaluation of the various power management states. These techniques do not prevent instantaneous power from exceeding the peak power budget, but instead trigger corrective action when the budget has been exceeded. Similarly, they are not suitable for many-core architectures (processors with tens or possibly hundreds of cores on the same die) due to an exponential explosion in the number of global power management states.

In this paper, we look at a hierarchical and a gradient ascent-based technique for decentralized peak power management for many-core architectures. The proposed techniques prevent power from exceeding the peak power budget and enable the placement of several more cores on a die than what the power budget would normally allow. We show up to 47% (33% on average) improvements in throughput for a given power budget. Our techniques outperform the static oracle by 22%.

## 1 Introduction

While power has long been a well-studied problem, most dynamic power reduction techniques, e.g., V/f scaling, clock gating, etc., exploit slack in the execution behavior of programs to reduce average power. Peak power is often left untouched. This is in spite of the fact that peak power plays a large role in determining the characteristics and hence the cost of the power supply, thermal budgeting for the chip, as well as the reliability qualification of the processor [2].

Recently proposed techniques for peak power management [4] involve centralized decision-making and assume quick evaluation of the various power management states. While these policies work well for multi-core processors (with a relatively small number of cores), they may not be suitable for many-core architectures (processors with tens or possibly hundreds of cores on the same die) due to an exponential explosion in the number of global power management states. For example, an 80-core processor like Intel's recent announcement [3], with two power states (full-power and half-power) can have over  $1.2 \times 10^{24}$  possibilities!

The second limitation of previously proposed techniques is that while these techniques attempt to limit the global power consumption, power overshoots are still allowed (overshoots trigger corrective action). The techniques that we propose in this paper attempt to guarantee that the power of a multi-core processor does not exceed a threshold. The guarantee is provided by dynamically selecting a subset of cores and scaling down their voltages.

## 2 Throughput Benefits of Peak Power Management

There is often a sizable gap between the average power and peak power of a multi-core processor. Figure 1 shows the distribution of power consumption for a 9-core chip

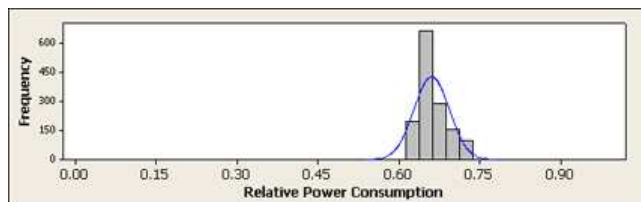


Figure 1. On average, a multicore processor consumes only a fraction of its maximum rated power.

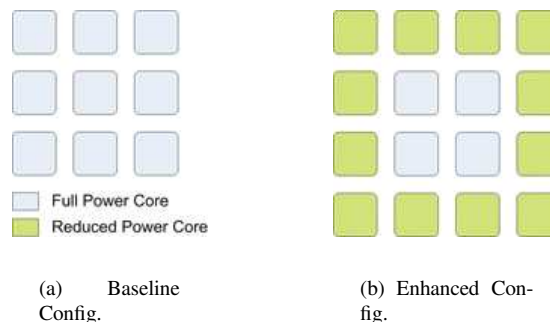


Figure 2. Power-Equivalent Configurations.

multi-processor (CMP) as a percentage of peak power for a set of workloads that we studied (details in Section 4). On average, the processor consumes only 66% of its maximum rated power. However, the processor and the power supply must be designed to supply the peak power and rated to handle this load. In theory, we should be able to add approximately 50% more cores *running at full power* and still remain below the peak power, on average.

Now consider an architecture (Figure 2) with peak power management that has several more cores on the die than the baseline processor. The average power of the new architecture is just below the peak power of the baseline processor while a peak power management mechanism prevents the new architecture from exceeding the peak power of the baseline (even though the processor contains several more cores than the power budget would normally allow). The peak power management mechanism bounds the processor power by intelligently scaling down the power for a subset of cores. Power can be scaled down through the application of V/f scaling, clock gating, or power gating. The throughput of this architecture is higher than the baseline processor, due to the increased number of cores.

In this paper, we employ V/f scaling to limit the maximum power consumption on a core.

### 2.1 Intelligent Core to Power State Mapping to Maximize Throughput

In this section, we describe the various baselines that we constructed to compare against our proposed approaches.

### 2.1.1 Static Mapping

We considered two static baselines – *random static* and *static oracle* – that prevent power overshoots and, therefore, provide the throughput benefits mentioned previously.

For a given processor configuration, the *random static* scheduler arbitrarily selects the cores that will be scaled.

Static oracular scheduling [5] requires foreknowledge of the behavior of applications in various power states and employs a metric called weighted speedup (WS) [7]. WS expresses the throughput of an application running at full power relative to the throughput of the same application running in a reduced power state. A large WS value indicates that the performance of an application deteriorates rapidly as power is decreased. The static oracular mapping is produced by allocating power to applications in order of decreasing WS.

We also evaluated a static configuration in which all cores are scaled to the lowest power state. This configuration – referred to as *all scaled* – maximizes core integration for a given power budget.

### 2.1.2 Dynamic Mapping

We used the previously proposed MaxBIPS policy [4] as our centralized dynamic baseline.

MaxBIPS, proposed by Isci *et al* [4], aims to optimize system throughput by predicting and choosing in the power scheduling stage the power mode combination that should maximize the throughput while obeying the power budget.

MaxBIPS was presented as a peak power management solution for a four-core multi-core processor and relies on a global arbiter/scheduler to make power management decisions for each core on the processor. As confirmed in Section 5), the technique is not scalable for a large number of cores. Secondly, while MaxBIPS attempts to limit the global power consumption, it allows power overshoots and triggers reactive correction for these overshoots.

The two dynamic techniques proposed in the next section do not allow power overshoots and are targeted towards architectures with a large number of cores.

## 3 Non-centralized Peak Power Management

In a many-core architecture, the responsibility of power management must be shifted away from a central arbiter and distributed to multiple locations around the processor.

### 3.1 A Hierarchical Approach

Consider a processor with 64 cores in which half of the cores run at full power and the other half run at reduced power to meet the peak power budget. The task of choosing the optimal power state boils down to choosing the 32 cores that can best utilize the full power state. Essentially, the size of the search space contains  $C(64,32)$  configurations, where  $C(x,y)$  denotes the number of  $y$ -combinations from an  $x$ -set.

Now, consider the same processor, divided into 4 clusters. For each cluster, the search space contains  $C(16,8)$  configurations when power is divided evenly between the clusters. Thus, the search space for the entire processor consists of 4 parallel decisions between  $C(16,8)$  states – 14 orders of magnitude fewer than the number required for a global decision.

Hierarchical Power Management involves dividing the processor into several clusters of cores and performing local scheduling within each cluster. The next level up in the hierarchy provides power management between clusters.

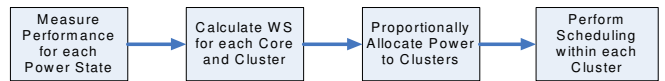


Figure 3. Hierarchical scheduling performs power management within clusters to reduce arbitration overheads.

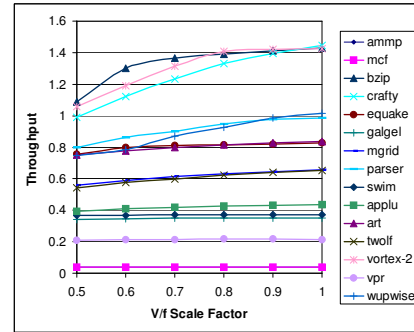


Figure 4. Different applications exhibit varying performance gradients with respect to power scaling.

Power is allocated to each cluster proportionally, based on WS. During allocation, the cluster-level arbiter ensures that no cluster receives more power than it can possibly use or less power than it needs to run each core at the most reduced power state. The algorithmic flow of hierarchical scheduling is depicted in figure 3. Note that hierarchical scheduling may result in decisions that are globally suboptimal due to inefficient power distribution between clusters.

In this paper, we consider a hierarchical scheduler that is based on hybrid a global trigger – i.e., re-allocation and local scheduling is triggered when the change in aggregate weighted speedup of the chip exceeds a threshold.

### 3.2 A Gradient Ascent-based Approach

Gradient ascent is motivated by the fact that different applications have different rates of performance decrease when voltage/frequency is scaled down and allocates power preferentially to applications with steep gradients. Figure 4 shows the performance gradients for some SPEC benchmarks. Since all the information needed to choose an efficient power state is present at the core level, the decision-making process can be executed locally, thereby eliminating the global arbiter and creating a scalable power management policy, independent of the number of cores on the processor. The algorithm in figure 5 describes the gradient ascent approach to peak power management.

To provide a global peak power guarantee, a single power balance counter is used to track the number of requests to ascend and descend one power level. If the power balance is not zero after each core asserts its desire, either some cores must accept a lower power state or some cores must assume a higher power state to keep the power requirement of the processor constant. A natural tradeoff exists here between the locality and scalability of decision-making and the amount of global information used, which in turn affects the optimality of the solution.

### 3.3 Decentralized MaxBIPS

Finally, for completeness, we also consider a decentralized version of MaxBIPS. The decentralized version of MaxBIPS is same as hierarchical scheduling where the local scheduling policy is MaxBIPS. Determination of when to do mapping is still based on a hybrid global trigger.

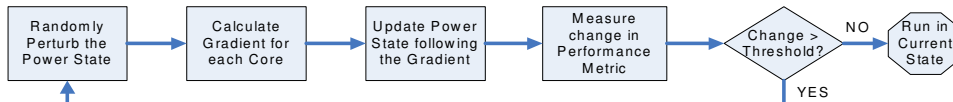


Figure 5. Gradient ascent performs power management locally rather than relying on a global arbiter.

Note that decentralized MaxBIPS will not prevent power overshoots because of the reactive nature of MaxBIPS.

#### 4 Methodology

The processors evaluated in our study are chip multiprocessors (CMPs) with homogeneous cores. All cores are modeled with 65 nm process parameters. The frequency and supply voltage of each core are 3 GHz and 1.5 Volts, respectively, at full power. All cores are connected to the L2 cache banks through a matrix crossbar interconnect. To account for increased area due to additional cores and interconnections, the L2 cache size of an enhanced configuration is reduced by half with respect to the corresponding baseline configuration.

Power estimates reported by Wattch [1] were used to calculate the peak power consumption of each core in various power states. Wattch reports the maximum dynamic power consumption for a core at a given supply voltage. To calculate the total peak power for a core in a given power state, an assumption is made that the peak dynamic power consumption represents 75% of the total processor power. Thus, the dynamic power values from Wattch are scaled up to represent the dynamic and static contributions to peak power. Table 1 gives the peak power figures for several V/f scaling factors, assuming 65 nm process technology.

Scale	1.0	0.9	0.8	0.7	0.6	0.5
Power (W)	18.289	15.549	13.098	10.888	8.899	7.107

Table 1. Peak Power for V/f Scaling Factors.

The supply voltage and frequency of each core in our modeled CMPs can be controlled independently. When a core switches V/f domains, we do not assume an instantaneous change. Instead, we model a gradual transition from one V/f domain to another at a rate of 10 mV/ $\mu$ s. When a transition between V/f domains occurs, the cores are halted until the transition is complete for all cores. During this time, the processor is assumed to still consume power, but no performance gains are registered. Modeling a V/f transition in this manner represents a very conservative approach.

Workloads are constructed from a set of 16 SPEC2000 benchmarks. Table 2 lists the benchmarks used in workload construction along with fast-forward distances in billions of instructions.

ampp	mcf	bzip	crafty	eon	equake	galgel	mgrid
2.0	12.6	0.4	0.7	0.1	3.5	5.0	2.1
parser	swim	applu	art	twolf	vpr	vortex	wupwise
0.4	0.3	0.3	7.5	0.9	36.1	6.0	1.1

Table 2. Benchmarks and Fast-Forward Distances.

For the 16 core results that we present in this paper, we average the results over five kinds of workloads (each consisting of 16 threads). One workload contains *all* benchmarks in equal proportions. Other workloads contain benchmarks with *high*, *low*, and *varied* sensitivity to V/f scaling. Finally, the *dynamic* workload contains benchmarks that exhibit more dynamic behavior than others.

Simulations are performed using SMTSIM [8] to simulate our various CMP configurations. Wattch is integrated

into SMTSIM to gather power statistics. SMTSIM executes statically-linked Alpha binaries. After fast forwarding each benchmark for an appropriate time [6], all simulations run for 1 Billion cycles.

We performed our evaluations with respect to weighted speedup (WS) and aggregate throughput and found no significant difference in trends or analysis.

#### 5 Analysis and Results

In this section, we demonstrate the performance benefits of peak power management for many-cores.

##### 5.1 Improving Throughput through Peak Power Management

$X \times V_{dd}$	Core Count	Full Power Cores	Reduced Cores	Peak Power (W)
1.0	8	8	0	146.0 W
0.8	9	5	4	143.8 W
0.7	10	5	5	145.9 W
0.5	11	6	5	145.3 W

Table 3. Alternative processor configurations with the same peak power bound.

To quantify the throughput benefits of peak power management, we consider an 8-core baseline processor, all cores running at full power, and compare it against peak power-equivalent processors with larger number of cores, where some cores are running at reduced power. Table 3 describes four different processor configurations which have nearly the same peak power requirements but differ in the number of cores. Figure 6 shows the corresponding throughputs for different static mapping policies. Results are shown relative to the average performance of a benchmark running on a reduced power core. As the results show, even static subsetting of cores can result in 16% throughput improvement when voltage is halved and by 5% when voltage is scaled by a factor of 0.8.

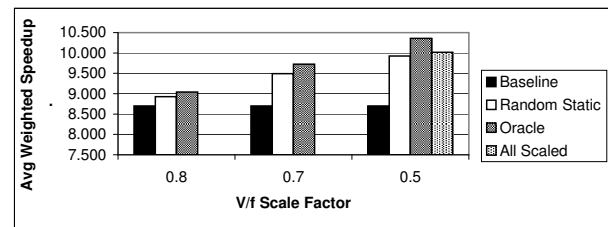


Figure 6. Performance comparison against the baseline processor for various peak power management techniques.

##### 5.2 Overhead of Managing Peak Power for Many-cores

The previous section shows the potential benefits of static peak power management. Performance improvements for dynamic peak power management depend on the overhead of making dynamic power scheduling decisions. Figure 7 shows the overhead of making peak power management decisions for the various techniques for various number of cores. MaxBIPS, which has been shown as an effective global power management technique in [4], is unsuitable for peak power management for many-core architectures (with

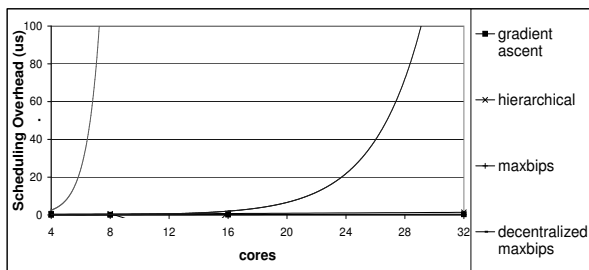


Figure 7. Timing Overhead of Making Dynamic Peak Power Management Decisions

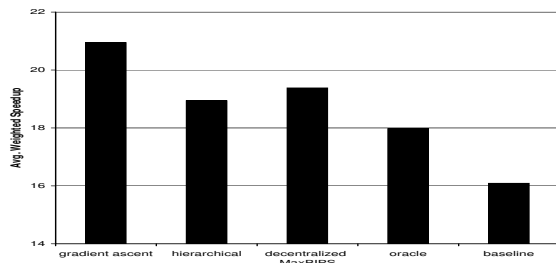


Figure 8. Performance of many-core power management policies against the baseline and the static oracle.

unacceptably large overheads even for 8 cores) due to the exhaustive nature of its state space search. Even the decentralized version of *MaxBIPS* has unacceptable (exponential) overheads for more than 16 cores due to the exhaustive nature of its local search. The decentralized policies (*gradient ascent* and *hierarchical*) have the least overhead. In fact, the overhead of *gradient ascent* shows little change with increasing number of cores till at least 1000 cores and is, therefore, the most suitable for many-core peak power management. Overheads for other techniques increase linearly and may be prohibitive for a large number of cores.

### 5.3 Power Management using Non-centralized Techniques

The previous section demonstrated that naive, centralized peak power management policies may be inadequate for peak power management for chips with a large number of cores (more than 16 or 32). There are certain policies that are inadequate even for 8 core processors. In this section, we consider a chip with a small number (16) of cores and compare the various strategies in terms of throughput improvement when their timing overheads are still acceptable. Figure 8 shows the results. The graph contains results for all policies except *centralized MaxBIPS* (which has unacceptably high overhead for 16 cores).

Results show that the proposed distributed techniques fare quite well at peak power management even for 16-core processors (for which the naive centralized techniques have acceptable timing overheads) due to their ability to reduce the state search space. With the hierarchical scheduling approach, for example, dividing the arbitration decision between four clusters in a 16-core processor reduces the complexity of the task by 99.8%. The use of hierarchical power management generated average improvements of 18.7% and 7.8%, respectively, over the baseline and static oracular models.

Note that hierarchical scheduling produces a power mapping that is potentially globally and also locally suboptimal. This is why worse performance is observed for 16 cores compared to decentralized *MaxBIPS* which is locally optimal. However, as Figure 7 shows, the overhead of decentralized *maxBIPS* becomes prohibitive for higher number of cores, so hierarchical scheduling will easily win out.

The gradient ascent approach to peak power management produced the best results of any policy, demonstrating a 33.3% improvement over the baseline processor and a 22.5% increase over the static oracle. While most other policies were designed with two possible power states for each core, the gradient ascent algorithm was given more freedom in the application of *V/f* scaling. Because the task of arbitration is distributed to each core, this additional control does not represent a substantial increase in the search space for power configurations. However, since our modified gradient ascent approach only shifts scale factors in discrete quanta, more power states will result in longer time to convergence when the optimal global mapping is far from the initial state of the processor.

## 6 Conclusion

Over our set of diverse workloads, our enhanced architectures (using the proposed techniques) averaged 30% better performance than comparable CMPs with equivalent area and power budgets. Also, the policies that we devised specifically for many-core architectures performed, on average, at least 22% better than our policies for multi-core architectures, even for 16 cores. As the number of cores on a processor die rapidly increases, the effectiveness of our techniques will only continue to increase.

## References

- [1] D. Brooks, V. Tiwari, and M. Martonosi. Wattch: a framework for architectural-level power analysis and optimizations. In *ISCA '00: Proceedings of the 27th annual international symposium on Computer architecture*, pages 83–94, 2000.
- [2] D. M. Brooks, P. Bose, S. E. Schuster, H. Jacobson, P. N. Kudva, A. Buyuktosunoglu, J.-D. Wellman, V. Zyuban, M. Gupta, and P. W. Cook. Power-aware microarchitecture: Design and modeling challenges for next-generation microprocessors. *IEEE Micro*, 20(6):26–44, 2000.
- [3] Intel Corp. *Intel's Teraflops Research Chip*.
- [4] C. Isci, A. Buyuktosunoglu, C.-Y. Cher, P. Bose, and M. Martonosi. An analysis of efficient multi-core global power management policies: Maximizing performance for a given power budget. In *MICRO 39: Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture*, 2006.
- [5] R. Kumar, D. M. Tullsen, P. Ranganathan, N. P. Jouppi, and K. I. Farkas. Single-ISA Heterogeneous Multi-core Architectures for Multithreaded Workload Performance. In *International Symposium on Computer Architecture*, June 2004.
- [6] T. Sherwood, E. Perelman, G. Hamerly, S. Sair, and B. Calder. Discovering and exploiting program phases. In *IEEE Micro: Micro's Top Picks from Computer Architecture Conferences*, Dec. 2003.
- [7] A. Snively and D. Tullsen. Symbiotic jobscheduling for a simultaneous multithreading architecture. In *Eighth International Conference on Architectural Support for Programming Languages and Operating Systems*, Nov. 2000.
- [8] D. Tullsen. Simulation and modeling of a simultaneous multithreading processor. In *22nd Annual Computer Measurement Group Conference*, Dec. 1996.