



저작자표시-비영리-변경금지 2.0 대한민국

이용자는 아래의 조건을 따르는 경우에 한하여 자유롭게

- 이 저작물을 복제, 배포, 전송, 전시, 공연 및 방송할 수 있습니다.

다음과 같은 조건을 따라야 합니다:



저작자표시. 귀하는 원저작자를 표시하여야 합니다.



비영리. 귀하는 이 저작물을 영리 목적으로 이용할 수 없습니다.



변경금지. 귀하는 이 저작물을 개작, 변형 또는 가공할 수 없습니다.

- 귀하는, 이 저작물의 재이용이나 배포의 경우, 이 저작물에 적용된 이용허락조건을 명확하게 나타내어야 합니다.
- 저작권자로부터 별도의 허가를 받으면 이러한 조건들은 적용되지 않습니다.

저작권법에 따른 이용자의 권리는 위의 내용에 의하여 영향을 받지 않습니다.

이것은 [이용허락규약\(Legal Code\)](#)을 이해하기 쉽게 요약한 것입니다.

[Disclaimer](#)

M.S. THESIS

NUMA-aware Hierarchical Power Management for Chip Multiprocessors

NUMA 구조를 인지한 칩 멀티프로세서를 위한 계층적 전력
관리

August 2017

DEPARTMENT OF COMPUTER SCIENCE ENGINEERING
COLLEGE OF ENGINEERING
SEOUL NATIONAL UNIVERSITY

Changmin Ahn

M.S. THESIS

NUMA-aware Hierarchical Power Management for Chip Multiprocessors

NUMA 구조를 인지한 칩 멀티프로세서를 위한 계층적 전력
관리

August 2017

DEPARTMENT OF COMPUTER SCIENCE ENGINEERING
COLLEGE OF ENGINEERING
SEOUL NATIONAL UNIVERSITY

Changmin Ahn

NUMA-aware Hierarchical Power Management for Chip
Multiprocessors

NUMA 구조를 인지한 칩 멀티프로세서를 위한 계층적
전력 관리

지도교수 Bernhard Egger

이 논문을 공학석사 학위논문으로 제출함

2017 년 04 월

서울대학교 대학원

컴퓨터 공학부

안 창 민

안 창 민의 공학석사 학위논문을 인준함

2017 년 06 월

위 원 장	<u>Srinivasa Rao Satti</u>	(인)
부위원장	<u>Bernhard Egger</u>	(인)
위 원	<u>허충길</u>	(인)

Abstract

Traditional approaches for cache-coherent shared-memory architectures running symmetric multiprocessing (SMP) operating systems are not adequate for future many-core chips where power management presents one of the most important challenges. In this thesis, we present a hierarchical power management framework for many-core systems. The framework does not require coherent shared memory and supports multiple-voltage/multiple-frequency (MVMF) architectures where several cores share the same voltage/frequency. We propose a hierarchical NUMA-aware power management technique that combines dynamic voltage and frequency scaling (DVFS) with workload migration. A greedy algorithm considers the conflicting goals of grouping workloads with similar utilization patterns in voltage domains and placing workloads as close as possible to their data. We implement the proposed scheme in software and evaluated it on existing hardware, a non-cache-coherent 48-core CMP. Compared to state-of-the-art power management techniques using DVFS-only and DVFS with NUMA-unaware migration, we achieve on average, a relative performance-per-watt improvement of 30 and 5 percent, respectively, for a wide range of datacenter workloads at no significant performance degradation.

Keywords: Many-core Architecture, NUMA, Scheduling, DVFS, Energy Efficiency

Student Number: 2015-22902

Contents

Abstract	i
Contents	ii
List of Figures	v
List of Tables	vii
Chapter 1 Introduction	1
Chapter 2 Motivation and Related Work	5
2.1 Characteristics of Chip Multiprocessors	5
2.2 Dynamic Voltage and Frequency Scaling	7
2.3 Power Management on CMPs	8
2.4 Related Work	10
Chapter 3 Cooperative Power Management	13
3.1 Cooperative Workload Migration	13
3.2 Hierarchical Organization	14
3.3 Domain Controllers	15
3.3.1 Core Controller	15

3.3.2	Frequency Controller	15
3.3.3	Voltage Controller	16
3.3.4	Chip Controller	16
3.3.5	Location of the Controllers	16
Chapter 4	DVFS and Workload Migration Policies	18
4.1	DVFS Policies	18
4.2	Phase Ordering and Frequency Considerations	19
4.3	Migration of Workloads	20
4.4	Scheduling Workload Migration	20
4.4.1	Schedule migration	21
4.4.2	Level migration	22
4.4.3	Assign target	25
4.4.4	Assign victim	27
4.5	Workload Migration Evaluation Model	27
Chapter 5	Implementation	29
5.1	The Intel Single-chip Cloud Computer	29
5.2	Implementing Workload Migration	31
5.2.1	Migration Steps	31
5.2.2	Networking	33
5.3	Domain Controller Implementation	33
Chapter 6	Experimental Setup	34
6.1	Hardware	34
6.2	Benchmark Scenarios	35
6.3	Comparison of Results	37

Chapter 7	Results	38
7.1	Synthetic Scenarios	38
7.2	Datacenter Scenarios	42
7.2.1	Varying Number of Workloads	42
7.2.2	Independent Workloads	45
7.3	Overall Results Comparison	46
Chapter 8	Discussion	48
8.1	Limitations	48
8.2	Extra Hardware Support	49
Chapter 9	Conclusion	50
Appendices		51
Chapter A	Benchmark Scenario Details	51
A.1	Synthetic Benchmark	53
A.2	Real World Benchmark	56
Bibliography		67
요약		73

List of Figures

Figure 2.1	Normalized memory bandwidth in dependence on the frequency and the distance from the memory controller.	6
Figure 2.2	Potential of DVFS.	7
Figure 2.3	Frequency and voltage domains in many-core CMPs.	9
Figure 4.1	Workload migration steps with f_{mid} example.	23
Figure 4.2	Workload migration steps with f_{mid} example result.	25
Figure 5.1	Intel SCC block diagram.	30
Figure 5.2	Workload migration.	32
Figure 6.1	G6 workload patterns.	36
Figure 7.1	Synthetic scenario PPW.	39
Figure 7.2	Synthetic scenario workload patterns.	40
Figure 7.3	PPW for a varying number of workloads.	42
Figure 7.4	PPW for scaled scenarios with a varying number of workloads.	43
Figure 7.5	Frequency map example for $G6$ and $Allhigh$	44
Figure 7.6	Experiment results for $G7$ to $G11$	45

Figure A.1	SCC core map	51
------------	------------------------	----

List of Tables

Table 4.1	Result of migration example	24
Table 6.1	Datacenter scenarios: distinct workloads patterns	35
Table 6.2	Average CPU and memory load	35
Table 7.1	Normalized PPW for synthetic scenarios	39
Table 7.2	Normalized performance per watt for Google cluster data scenarios	45
Table A.1	<i>SynMem</i> benchmark scenario	53
Table A.2	<i>SynCPU</i> benchmark scenario	53
Table A.3	<i>SynRand</i> benchmark scenario first half	54
Table A.4	<i>SynRand</i> benchmark scenario second half	55
Table A.5	Google cluster data benchmark scenario #1	56
Table A.6	Google cluster data benchmark scenario #2	57
Table A.7	Google cluster data benchmark scenario #3	58
Table A.8	Google cluster data benchmark scenario #4	59
Table A.9	Google cluster data benchmark scenario #5	60
Table A.10	Google cluster data benchmark scenario #6	61

Table A.11	Google cluster data benchmark scenario #7	62
Table A.12	Google cluster data benchmark scenario #8	63
Table A.13	Google cluster data benchmark scenario #9	64
Table A.14	Google cluster data benchmark scenario #10	65
Table A.15	Google cluster data benchmark scenario #11	66

Chapter 1

Introduction

The past decade has brought a shift from high-performance single- or dual-core processors to chip multiprocessors (CMPs) integrating from a few tens up to a thousand cores into one processor die [1,2,8,9,30,31]. Chip-level power and thermal constraints are now one of the primary design constraints and performance limiters [2]. Higher power consumption not only leads to increased energy cost but also causes higher die temperatures that adversely affect chip reliability and lifetime [7]. Even in commodity processors, such as Intel processors based on the Haswell, micro architecture power constraints result in reduced per-core performance when multiple cores are active [16].

To reduce overall chip energy consumption, modern processors provide hardware support to dynamically lower the operating voltage and clock frequency of clocked resources through dynamic voltage and frequency scaling (DVFS). Depending on the utilization of processor cores, for example, core voltages and frequencies are adjusted in order to minimize power consumption while, at the same time, meeting performance requirements [5]. On CMPs, the required logic for individually controlling the voltage for each core is becoming too costly [22]; instead, cores are physically clustered

into voltage and frequency domains leading to so-called *multiple-voltage/multiple-frequency* (MVMF) designs where all cores within a domain run at the same voltage or frequency [7, 14, 15, 34].

Managing power on CMPs has recently received considerable attention [6, 10, 11, 13, 18, 25, 28, 32, 33, 35]. Existing research on power management for CMPs foremost focuses on minimizing power consumption or optimizing performance for a given power budget [10, 18, 25, 28, 33]. Solutions for MVMF architectures combine DVFS with thread migration [6, 17, 19, 21, 33], because co-locating threads with similar performance requirements into the same domain allows for better tailored DVFS settings for that domain [17].

The integration of more and more cores into CMPs poses several other architectural challenges. First, to cope with the increased bandwidth requirements the cores of a CMP are typically connected to several memory controllers by a network-on-chip (NoC). Depending on the location of the core and the accessed memory controller large differences in memory access latency are observed, resulting in a non-uniform memory access (NUMA) architecture on a single chip. The second challenge for CMPs is that maintaining a coherent global view of shared memory in the presence of local caches is becoming difficult. While today's commercial CMPs typically still maintain cache coherency to support existing operating systems and parallel runtimes, the trend goes towards partial or no coherence [2, 15, 31].

In this thesis, we propose a hierarchical power management technique for MVMF CMPs that considers the (not necessarily cache-coherent) NUMA memory architecture. Existing techniques fall short for a variety of reasons. Many works assume per-core DVFS control which limits their applicability to MVMF designs. Researches employing thread migration assume symmetric multiprocessing with one centralized kernel and cannot easily cope with non-coherent memory architectures. Lastly, to the best of our knowledge, no work considers the NUMA properties of CMPs resulting in core

mappings that are not optimal with respect to the locality of the data accessed by individual threads.

The presented power management technique can be applied to monolithic kernels running on a cache-coherent SMP processor as well as non-coherent memory architectures running a distributed micro kernel. The hierarchical design naturally maps to the architecture with per-core utilization monitors, individual frequency controllers for the frequency domains, voltage controllers for the voltage domains, and a central migration controller. The solution is entirely implemented in software and does not require special hardware support. The migration controller computes and orchestrates migration of workloads based on a cost-benefit model. The individual frequency and voltage controllers regulate the frequency/voltage for the controlled domain. A working implementation is provided for the Intel Single-Chip Cloud Computer (SCC) [15] and evaluated with real-world workloads. All experiments and measurements are performed on the architecture itself and thus include overhead incurred by DVFS transitions, cold cache misses, workload migration, and the overhead caused by the different controllers. We compare the proposed technique to a DVFS-only approach [17] and a method with DVFS and migration [21]. On average, we achieve a 54, 33, and 5% higher performance-per-watt ratio over standard Linux, DVFS-only and DVFS with migration at no performance degradation.

We show that, even with complete separation and isolation of processes, it is possible to change the allocation of the physical cores to the applications with almost zero overhead on CMPs and very little support from the application-specific runtime environment. This technique allows us to group, on a global level, cores that exhibit similar load patterns onto voltage and frequency domains before applying DVFS.

In previous work [21], we proposed the design and implementation of a hierarchical power management technique using workload migration for multi-voltage/multi-frequency CMPs. The logical abstraction mimics the physical layout of the CMP (core,

frequency domain, voltage domain, and chip). In summary, the contributions of this work are as follows:

- we analyze the interplay of DVFS with workload migration and propose a data-locality-aware migration heuristic.
- we describe and evaluate a proof-of-concept software implementation for the Intel SCC [15] architecture running up to 40 different real-world workloads. All measurements are performed on a real system.

The remainder of this thesis is organized as follows: Chapter 2 discusses the problem formulation and related work. Chapter 3 describes the hierarchical power management framework, and Chapter 4 discusses the DVFS and migration policies in detail. Chapter 5 describes the implementation for the Intel SCC. The experimental setup and the results are presented in Chapters 6 and 7, respectively. In Chapter 8, we talk about the limitations and expected improvement of this work. Finally, Chapter 9 concludes this thesis.

Chapter 2

Motivation and Related Work

2.1 Characteristics of Chip Multiprocessors

Technology scaling, thermal limitations and the insight that doubling the logic in a processor core only delivers about 40% more performance, known as Pollack's Rule, have led to the introduction of chip multiprocessors with tens or hundreds of cores on one processor die [3, 4]. Architectural characteristics of today's and future many-core CMPs impose new restrictions on the design and implementation of operating systems in particular with respect to workload scheduling and power management.

The cores of a CMP are typically organized in a 2-D array. The Kilocore processor, for example, arranges its 1000 cores on a 32x32 grid [2]. A network-on-chip (NoC) interconnects the cores of a CMP and is used both for inter-core communication and accesses to memory and external devices such as network or storage controllers. The flow of data packets through the NoC is controlled by routers; this routing comes with a small delay. As a consequence, the distance between the source and the destination (number of hops in the NoC) has a significant effect on the access latency of indi-

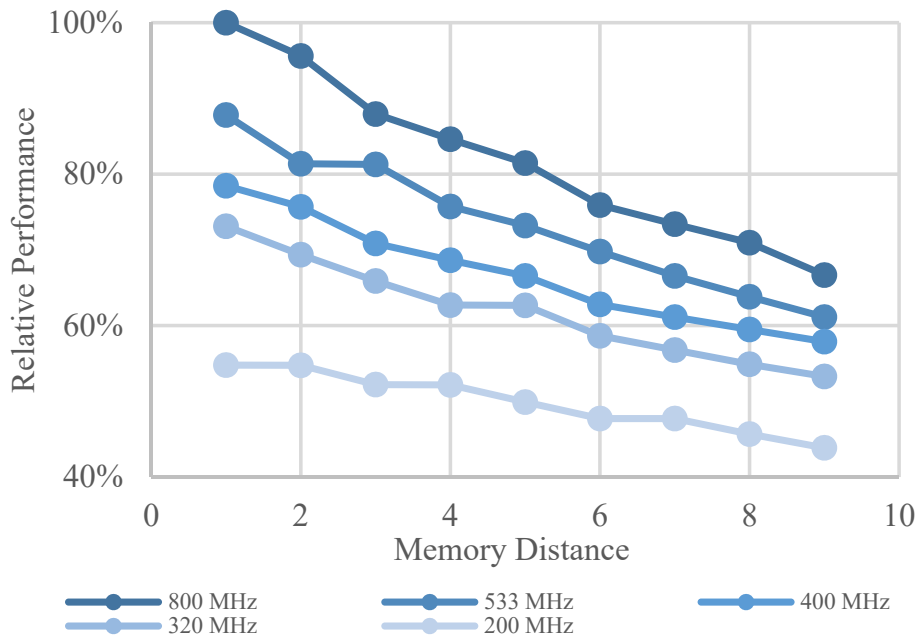
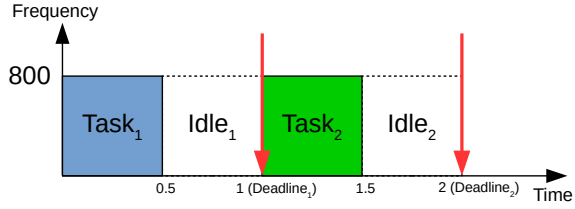


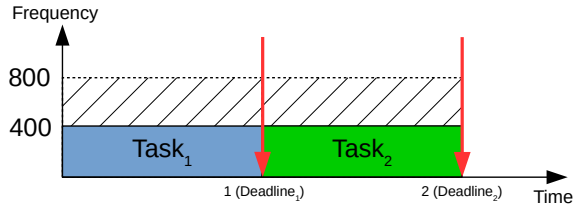
Figure 2.1: Normalized memory bandwidth in dependence on the frequency and the distance from the memory controller.

vidual cores to memory, resulting in a NUMA architecture on a single die. Figure 2.1 shows the relative memory bandwidth of a 48-core CMP, the Intel Single-Chip Cloud Computer (SCC) [15], in dependence the number of hops between the core and the memory controller where the accessed data is. Each criteria represents each frequency level which SCC supports. Y-axis is relative memory performance to the best performance which a core can achieve when it has the highest frequency and the closest location.

The memory organization of CMPs combines several memory controllers to access off-chip data with on-chip local memory in the form of scratchpad memories or local or shared caches. As an example, the Intel SCC processor has four memory controllers and integrates both caches and user-managed scratchpad memory on the die [15]. It is foreseeable that maintaining cache coherency at the hardware level will become



(a) Run tasks as fast as possible.



(b) Run tasks as slow as possible.

Figure 2.2: Potential of DVFS.

extremely challenging with hundreds or thousands of local caches on one die, and most research prototypes integrating many cores do not provide a coherent cached view of global memory anymore [2, 15, 31].

2.2 Dynamic Voltage and Frequency Scaling

$$P_{chip} = P_{dyn} + P_{SC} + P_{leak} \quad (2.1)$$

The power consumption of a chip includes dynamic power consumption, short-circuit power consumption, and power loss of transistor leakage currents (Equation 2.1). In these factors, short-circuit power consumption and transistor leakage are steady. So, we can only save dynamic power consumption, $P_{dynamic}$, calculated by equation from [29]:

$$P_{dynamic} = ACV^2 f \quad (2.2)$$

where A is gate activity factor, C is capacitance, V is voltage, and f is frequency.

From Equation 2.2, because A and C are constants, we save power by lower voltage and frequency. This technique which changes voltage and frequency of a chip in runtime is dynamic voltage and frequency scaling, which called DVFS. But, when we use DVFS, a chip has a necessary condition when supports a frequency. For each frequency level, a chip has to support some minimum voltage matching with the frequency level.

The potential of DVFS is represented in the Figure 2.2. These two graphs are represent running time and frequency when running same tasks within deadlines at different frequencies. Then, with Equation 2.2 and voltage and frequency values from real system, 1.1 V to run frequency 800 MHz and 0.8 V for 400 MHz, we can calculate the energy for Figure (a) and Figure (b). For the Figure (a) which it runs tasks without DVFS, we can get $(1.1 \text{ V})^2 * 800 \text{ MHz} * 0.5 \text{ sec} * 2 = 48.400 \text{ V}^2 \text{ MHz sec}$, even if we assume 0 power consumption at idle stage, which is not true. For the Figure (b) which it runs tasks with DVFS, we can get $(0.8 \text{ V})^2 * 400 \text{ MHz} * 1 \text{ sec} * 2 = 25.600 \text{ V}^2 \text{ MHz sec}$, which save almost half of power when we run the tasks without DVFS.

2.3 Power Management on CMPs

In the past, DVFS has proved to be an effective technique to limit power dissipation, and an enormous body of research exists on that topic. On CMPs, equipping every single core with a voltage regulator is becoming too costly [22]; multiple-voltage/multiple-frequency (MVMF) designs are being proposed. In a MVMF design, cores in the same voltage or frequency domain share the same voltage or frequency, respectively. Since the range of valid frequencies depends on the supply voltage, one voltage domain typically contains one or more frequency domains (Figure 2.3).

The additional constraints imposed by the CMP architectures and in particular their

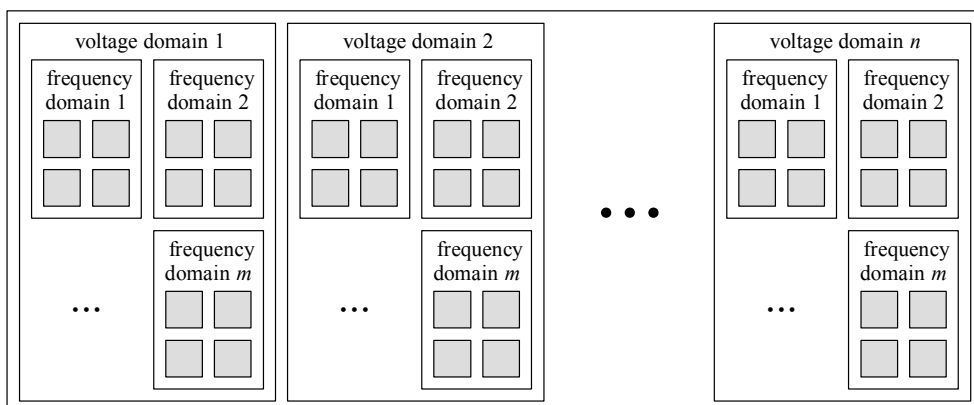


Figure 2.3: Frequency and voltage domains in many-core CMPs.

power management capabilities require new approaches power management. Existing designs that assume per-core DVFS control, cache-coherent global shared memory, and uniform access latency to the memory are not able to cope with today's characteristics of CMPs. What is required is a power management approach that considers all aspects of modern CMPs. In addition, power management has to hand-in-hand with workload scheduling. From a power management perspective only, cores with similar performance requirements need to be grouped together in voltage and frequency domains in order to achieve optimal power savings. On the other hand, the NUMA characteristics of the chip require the scheduler to place workloads as close as possible to the accessed memory controllers. In addition, the scheduler and the power manager may need to adhere to user-defined performance goals such as improving performance, maintaining Quality of Service (QoS), dissipate heat evenly, or minimize power for a given budget. In this paper, we describe our solution of a cooperative and hierarchical power management technique that balances the conflicting goals of scheduling and power management in order to achieve a better energy efficiency at no or minimal performance degradation.

2.4 Related Work

There is a significant amount of work focusing on the design and implementation of power management techniques for CMPs. One line of related work considers heterogeneous CMP designs with the goal of minimizing power consumption with no or minimal performance loss. Kumar *et al.* [23] propose heterogeneous CMPs composed of cores with an identical ISA but different power characteristics. Ghiasi [12] proposes CMPs with cores executing at different frequencies. Both works show that such systems offer improved power consumption and thermal management. Our work differs in that our approach modifies the voltage/frequency of cores dynamically, without being bound to certain hardware heterogeneity.

Another line of research has focused on exploiting idle periods. Meisner *et al.* propose PowerNap [26] and DreamWeaver [27]. Both assume hardware support for quick transitions between on- and off-states; the latter work batches wake-up events to increase the sleep periods. Our work is orthogonal to such approaches.

A number of researchers have proposed heterogeneous power management techniques for CMPs [6, 17, 18, 20, 24, 25, 28, 33].

Li *et al.* [24] provide an analytical model and experiments to show to what extent parallel applications can be parallelized given a power-budget. Isci *et al.* [18] apply different DVFS policies under a given power budget and show that their best policy performs almost as good as an oracle policy having limited knowledge of the future. Meng *et al.* [28] propose an adaptive power saving strategy that adheres to a global chip-wide power budget through run-time adaptation of configurable processor cores. They integrate multiple power optimization techniques (in this case DVFS and cache resizing) into single power management unit. To optimize performance of a CMP under given power budget by using power optimization techniques, they use a greedy search algorithm to select a technique. They introduce models to predict performance

of a core after apply power optimization techniques. Especially, for cache resizing, they achieve reasonable accuracy. But unlike us, they target per-core DVFS supported system which doesn't need thread migration.

Rangan *et al.* [33] propose ThreadMotion, a technique that moves threads around in order to improve power consumption. In a multi-core system which has homogeneous cores with heterogeneous power setting, they use thread migration to exploit fine variation in program behavior which DVFS cannot exploit because DVFS is too slow. They use a coarse-grained prediction-driven approach and a last-level cache miss driven approach to trigger thread migration. This technique requires hardware support to quickly move threads from one core to another. Our approach is similar but can be implemented on available CMPs without extra hardware support.

Cai *et al.* [6] propose Thread Shuffling, a technique that migrates hardware contexts around to exploit non-critical threads; non-critical threads can then be executed at reduced speed. They identify critical threads by using meeting point thread characterization. They assume a core has multiple hardware context and per-core DVFS which corresponds with per-tile DVFS. This work can be applied when a single parallel application is running on the system. In this thesis, we focus on independent OSes as opposed to threads within a parallel application.

Ma *et al.* [25] propose a scalable solution aiming at a mixed group of single-threaded and multi-threaded applications. They introduce hierarchical management to reduce the complexity of scheduling. Their framework periodically groups cores which running same applications. It partitions chip power budget between groups according to power efficiency. Then, it partitions quota again among cores in a group by analyzing thread criticality. They implemented and tested the power controller not only on a simulator but also on a real system. Unlike our approach which is best-effort, they aim at minimal performance reduction while maintaining a global power budget in per-core DVFS systems.

Jha *et al.* [20] propose a hierarchical power management on systems which have per-tile DVFS and shared a last-level cache. They classify threads with DVFS sensitivity and cache behavior. They migrate threads based on its' class. They call this DVFS and cache-aware thread migration (DCTM). This work aims getting best performance under power budget. But our focus is reducing power consumption under satisfying performance target. And our system doesn't have shared caches.

A previous work of our lab is a hierarchical power manager for the Intel SCC [17, 21]. While Ioannou *et al.* [17] apply DVFS to a static workload assignment, Kang *et al.* [21] demonstrate that adding workload migration can yield a significant improvement in the performance/watt ratio. The proposed buyer-seller algorithm used for workload migration, however, fails to consider data locality and thus results in sub-optimal core assignments. In this thesis, we follow the overall system architecture of the previous two works but present a new greedy workload placement algorithm that balances the conflicting goals of the scheduler and power management for MVMF CMPs.

Chapter 3

Cooperative Power Management

The proposed cooperative hierarchical power manager combines workload migration with DVFS to achieve optimal power efficiency. We target both CMPs with and without cache-coherent global shared memory. We employ a distributed OS design with small individual kernels running on each core.

3.1 Cooperative Workload Migration

Workloads with similar performance characteristics need to be grouped in frequency / voltage domains to allow for optimal DVFS and, as a consequence, improved power efficiency. As an example, consider two frequency domains with two cores each. In both domains, one core is running at 100% CPU load, the other one is only 10% loaded. To maintain throughput, both domains need to run at the highest frequency f_{max} in order to provide the computing power required by the busy core. Both lightly loaded cores will also run at f_{max} even though theoretically $\frac{1}{10}f_{max}$ would suffice. Workload migration allows us to group the heavily loaded cores into one and the lightly loaded

cores into the other domain. One domain can then run at f_{max} , the other one at $\frac{1}{10}f_{max}$ without sacrificing performance. Taking data locality into consideration complicates the situation. In the above example, one of the busy cores needs to be moved into the domain of the other busy core. This relocation may change the distance of the core from the memory controller holding its data and reduce the throughput on the newly assigned core.

With a distributed OS comprising microkernels that individually schedule the tasks assigned to them, task migration is more difficult to achieve than in global shared memory systems with a monolithic kernel. Since each kernel has its own network ID on the NoC, moving a task from one core to another would disconnect established communication channels. If properly orchestrated, however, the architecture of CMPs allows for dynamic workload re-allocation without side-effects. The idea is to migrate the entire microkernel from one core to another with its entire workload. Since we assume (non cache-coherent) global shared memory, migration of a microkernel only requires moving the volatile state of a core, i.e., its processor state, from one core to another. A greedy algorithm to deal with these potentially conflicting goals of core placement is presented in Chapter 4. Implementation details about microkernel migration on existing hardware are discussed in Chapter 5.

3.2 Hierarchical Organization

The logical structure of the hierarchical power manager reflects the structure of the CMP with separate frequency and voltage domains. At the lowest level in the hierarchy are the *core controllers* that represent a single core. The second level, the *frequency controllers*, represents a frequency domain with m individual cores all running at the same frequency. The *voltage controllers* at next level constitute a voltage domain with n number of *frequency controllers*. At this level, voltage changes are initiated. The

top level in the hierarchy, finally, is represented by the *chip controller* and models the entire chip.

3.3 Domain Controllers

Each domain, from core to frequency, voltage, and the global chip level, operates its own domain manager. Each level only communicates directly with the level above or below, i.e., the clock domain manager interacts with the voltage domain manager, the voltage domain manager interacts downstream with the clock domain, and upstream with the global domain manager. The functionality of the different domain managers is elaborated in more detail in the following sections and a possible implementation is discussed in Chapter 5.3.

3.3.1 Core Controller

The task of the *core controllers* is to monitor and predict the performance of the workload on the associated core. Each microkernel runs a core controller daemon monitoring the performance (load or instructions per clock (IPC)) and the number of memory accesses by periodically querying the performance monitoring unit (PMU). The core controllers also predict the required computational performance based on extrapolated measured data. At regular intervals, the core controllers communicate with their frequency controllers To report the required operating frequency and memory-boundness. The core controllers run on every kernel.

3.3.2 Frequency Controller

For each frequency domain, the *frequency controller* gathers data about the frequencies and workloads from core controllers within its domain, and processes and forwards that data to the voltage controller. The frequency controllers also compute and

set the operation frequency of the domain. The clock frequency is constrained by the current voltage level of the corresponding voltage domain and computed based on the requested frequency levels reported by the core controllers and the currently active DVFS policy (see Chapter 4.1).

3.3.3 Voltage Controller

The *voltage controllers* gather data from their frequency controllers and forward it to the chip controller. In addition, the voltage controllers also compute and set the operating voltage of their domains. Note that voltage changes must happen in close collaboration with the frequency controllers because the maximal operating frequency has a linear relationship to the supply voltage. This is particularly important if the voltage of a domain is to be lowered. In that case, the frequency controllers must first reduce the frequency to a value below or equal to the maximal operating frequency of the new supply voltage before the voltage change can occur.

3.3.4 Chip Controller

The *chip controller* uses the processed frequency and voltage requests from the subordinate controllers to compute a core assignment that allows more optimal DVFS settings at the voltage and frequency domain levels. The chip controller migrates the microkernels before signaling the voltage controllers to initiate DVFS adjustments.

3.3.5 Location of the Controllers

In a pure software implementation, one kernel per frequency and voltage domain needs to run the respective controller. Similarly, the chip controller also runs on one of the cores. Since we migrate entire kernels, it is impossible to designate the kernel for each of the controllers offline. Instead, we run an instance of each controller in *every* kernel. The frequency, domain, and chip controller in a kernel are activated and deactivated

depending on the physical location on the CMP. In other words, the functionality is pinned to the physical core and not the kernel. For example, if the frequency controller for frequency domain 1 is pinned to core 0, the kernel that is currently running on core 0 will activate its frequency controller. Such a scheme has the additional benefit that no discovery service is needed to find the controllers.

Chapter 4

DVFS and Workload Migration Policies

In this thesis, we focus on optimizing the *performance per watt* ratio of the overall chip. Other policies, such as, for example, even heat dissipation or adhering to a given power budget, can also be implemented within the framework of the presented collaborative hierarchical framework and are part of future work.

The power management policy is implemented in the global domain manager. The migration and DVFS algorithms are invoked at regular intervals by the scheduler. The DVFS and migration policy, though the former depends on the latter, are completely separated to be able to combine different migration and DVFS policies freely.

4.1 DVFS Policies

We implement two DVFS policies employed in the hierarchical power manager for CMPs proposed by Ioannou *et al.* [17] and employed by Kang *et al.* [21]. Both works have been implemented and evaluated on the same hardware and provide a good reference point.

- *Allhigh*: this DVFS policy runs all cores within a *voltage domain* at the highest frequency requested by the subordinate frequency domains. The supply voltage is set to the lowest voltage that supports the requested frequency.
- *Tile*: grants the requested frequency to each *frequency domain* and set the voltage accordingly. In [17] this policy is denoted *Simple*, we follow Kang’s nomenclature here.

In both policies, the supply voltage of the domain is set to the lowest voltage that supports the highest frequency of any of the subordinate frequency domains.

We have not implemented the *Allow* and *Allmean* policies since they sacrifice too much performance in return for power savings.

4.2 Phase Ordering and Frequency Considerations

In order to achieve maximum power savings, migration should occur before applying DVFS. The frequency of migration, and voltage/frequency changes is determined by the cost of the individual operations. The time required for migration is largely unaffected by the number of kernels that are migrated because the involved kernels can migrate in parallel. Kernels involved in a migration flush their caches and are briefly stopped, while the other kernels continue to run. Voltage changes incur a not insignificant overhead because all cores in the domain are stopped during the rather long voltage adjustment. Frequency changes, on the other hand, are almost instantaneous and can be performed often. On our specific target architecture, the Intel SCC, we have measured the following latencies: $\leq 3ms$ for migration, $\leq 10ms$ for voltage changes, and a few thousand cycles for frequency changes. We perform workload migration and DVFS at a 3 second interval, because of high latency for voltage changes. Besides, the SCC only supports one voltage change at a time; i.e., different domains cannot change the voltage in parallel. Nevertheless, in our experiment, workload migration and voltage

changes can be performed at every step. Chapter 7 discusses the benchmarks and results in more detail.

4.3 Migration of Workloads

As outlined in Chapter 3.1 workloads which have similar load pattern need to be grouped in order to achieve good power savings. A naïve algorithm would be to sort the workloads by their performance requirements and then assign them into the voltage and frequency domains. While the resulting migration of workload to domains is optimal for CPU-bound applications, the algorithm fails to consider the overhead of kernel migration. The migration of a kernel itself is very quick (measurements on a real system yield an overhead of $\leq 3ms$). However, each time a kernel is migrated to a different core, the workload running on the kernel will experience cold misses in the local caches that in turn lead to a loss in performance as well as increased memory traffic. To minimize this overhead, the number of migrated kernels should be kept as low as possible.

Due to the NUMA nature of CMPs, kernel migration can have a significant effect on the access latency and bandwidth of memory accesses. Since the data of a migrated workload is not moved, migrating a memory-intensive workload executing on a core close to the memory controller to a core far away from the memory controller can cause a significant drop in memory bandwidth and access latency (Figure 2.1).

4.4 Scheduling Workload Migration

We have two conflicting goals. One of the goals is optimizing power. Another is optimizing memory performance. To solve this problem, we develop a greedy algorithm to schedule workload migration which makes CMP optimize power consumption and memory distance with given performance requirement. When we make every voltage

domain always use the minimum voltage which can support maximum frequency request of workloads in it, the main concepts of our algorithm are as follows:

- for each target frequency level f_{target} , collect T , which are workloads with f_{target} , into minimum number of voltage domains.
- migrate each workload in T , in descending order of memory load, to a core which has minimum distance to memory controller.

In this algorithm, for each frequency level, we first generate voltage domain combination, $Comb$, which has every possible combination of voltage domains to place T . For every voltage domain set, set , in $Comb$, we assign each T , in descending of memory intensity, on a core which has minimum distance to memory in set .

To evaluate expected energy consumption of each set , we introduce a model to evaluate the power state of a chip. We discuss details about this model later in Chapter 4.5. Base on this evaluation model, we select the best voltage domain set among $Comb$, and repeat for next target frequency which one step lower than f_{target} until the minimum frequency. After schedule migration, we perform workload migration only if the chip's power status over a threshold rate.

4.4.1 Schedule migration

In this step, we get workload migration schedule for each f_{target} in descending order, as represented line 3 – 7 in Algorithm 1. If the schedule is better than the previous result, we save it and give the workload mapping to next f_{target} level (line 4 in Algorithm 1). Then, we use the evaluation model. If and only if the final migration result's power status is over the threshold, we perform workload migration (line 7 – 8 in Algorithm 1).

Algorithm 1 Decide Migration

EvalMigBenefit(*migMap*): returns a power rate between original state of chip and *migMap*

```
1: function DecideMigration(migThreshold)
2:   migMap  $\leftarrow$  Current core mapping
3:   for each  $f \in \text{Freq\_Range}$  do ▷ Decending order
4:     tmpMap  $\leftarrow$  LevelMig(migMap,  $f$ )
5:     if EvalMigBenefit(tmpMap) is better than before then
6:       migMap  $\leftarrow$  tmpMap
7:     end if
8:   end for
9:   if maxMigBenefit > migThreshold then
10:    return migMap
11:  end if
12:  return NULL
13: end function
```

4.4.2 Level migration

For given f_{target} , in this step, we try to minimize the number of voltage domains which have a core that has requested f_{target} . In Algorithm 2, this algorithm collect T (line 3 in Algorithm 2), and calculate how many voltage domains we need to allocate all T . To calculate this, we divide the voltage domains into two groups. When V_l represent a voltage which is minimum to support frequency level l , one group consist of voltage domains which have $V > V_{target}$. We call this group as $vDom_{used}$. Another is a group of voltage domains which have $V \leq V_{target}$, which called $vDom_{left}$.

In the context of minimizing N , the number of voltage domains which have V_{target} , we can achieve this by placing T in $vDom_{used}$. Because this migration won't increase N . From $vDom_{used}$, we collect candidate cores which have a workload with $f \leq f_{target}$, which called $Cand_{in_used}$. Then, we can calculate N , the minimum number of voltage domains without the number of $vDom_{used}$ by dividing the number of cores in a voltage

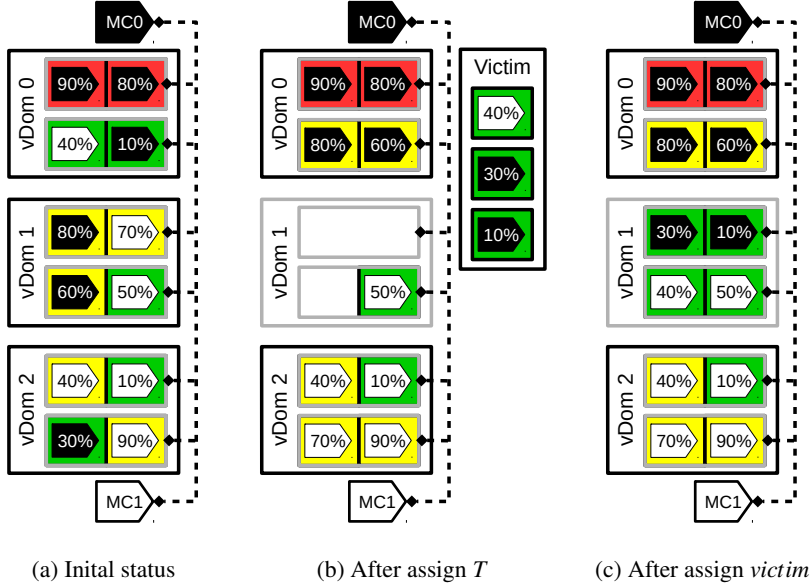


Figure 4.1: Workload migration steps with f_{mid} example.

domain with the ceiling function. This calculation is given as follows

$$N = \left\lceil \frac{|T| - |Cand_{in_used}|}{|CoresInVoltDom|} \right\rceil \quad (4.1)$$

After get the number of voltage domain we need, we can make combination with N number of voltage domains from $vDom_{left}$ (line 4 in Algorithm 2). Then, we can make complete *set* which T going to be by adding $vDom_{used}$ to each combination (line 6 in Algorithm 2).

For example, let's assume we have target frequency f_{mid} with core mapping like Figure 4.1 (a). In the figure colors red, yellow, and green boxes represent high, middle, and low workloads in a core, and the pentagons represent memory controller which the workload uses with color and amount of memory load with the number in it. We have voltage domain groups, $vDom_{used} = \{vDom0\}$ and $vDom_{left} = \{vDom1, vDom2\}$.

Algorithm 2 Level Migration

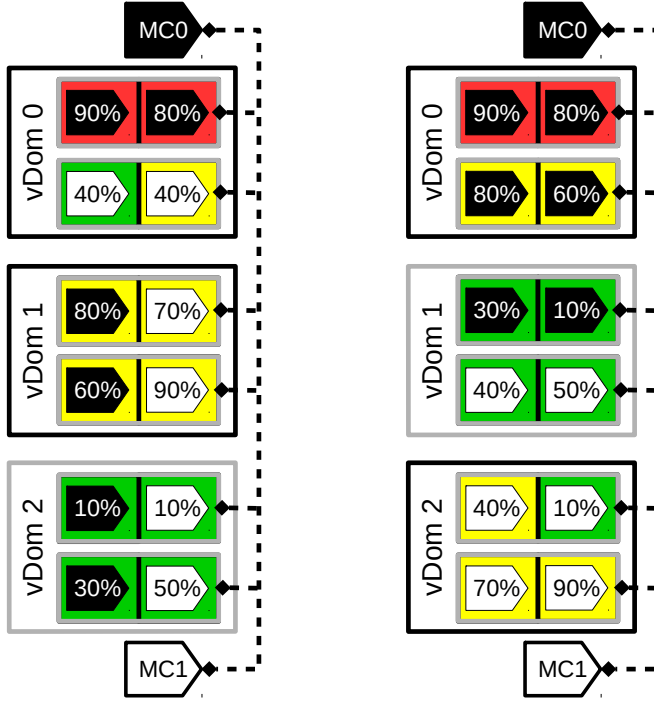
```
1: function LevelMig(migMap, ftarget)
2:   migResult  $\leftarrow$  migMap
3:   T  $\leftarrow$  GetWorkloadList(ftarget)
4:   Comb  $\leftarrow$  MakeVDomComb(ftarget, migMap)
5:   for each set  $\in$  Comb do
6:     set  $\leftarrow$  set + vDomused
7:     tmpMap  $\leftarrow$  PlaceTarget(T, set, ftarget)
8:     migMapList.add(tmpMap)
9:   end for
10:  migResult  $\leftarrow$  BestPowerState(migMapList)
11:  return migResult
12: end function
```

Table 4.1: Result of migration example

	{vDom0, vDom1}	{vDom0, vDom2}
# of migration	4	6
sum of weighted memory distance	19.8	11

Also, the number of T and $Cand_{in_used}$ are 5 and 2. Then $N = \lceil \frac{5-2}{4} \rceil = 1$ and we can make set by selecting 1 voltage domain from $vDom_{left}$ and add $vDom_{used}$. The result is $Comb = \{\{vDom0, vDom1\}, \{vDom0, vDom2\}\}$.

Then, we assign T in a for each set in $Comb$, evaluate the assignments, and return the best migration map from the migration mapping list (line 9 – 11 in Algorithm 2). The results look like Figure 4.2. We can calculate the number of migrations and sum of weighted memory distance like in Table 4.1 for each Figure (a) and Figure (b). How to place workloads will be discussed in Chapter 4.4.3. After getting the results, we evaluate each migration schedules' power status by using the evaluation model and choose the best in the $Comb$. Then, we calculate migration for next frequency level.



(a) Result for a set $\{vDom0, vDom1\}$. (b) Result for a set $\{vDom0, vDom2\}$.

Figure 4.2: Workload migration steps with f_{mid} example result.

4.4.3 Assign target

In this step, we show how assign T to $Cand$ for each set from $Comb$. First, we collect $Cand$ in a set (line 3 – 5 in Algorithm 3). But, at this time, the $Cand$ are cores which have a workload $< f_{target}$. Because, we will make workloads which are initially placed in set keep position (line 9 – 10 in Algorithm 3). Then, we allocate for each $workload$ in T , in descending order of memory intensity, on the core which has the shortest distance to the workload's memory controller in $Cand$ (line 8 in Algorithm 3).

After placing $workload$, there will be workloads which are kicked out from the own core. We shall call these workloads as *victim*. Moreover, cores which T have been

Algorithm 3 Place Target Workload

cand: Candidates of migration destination

```
1: function PlaceTarget( $T, vDomSet, f_{target}$ )
2:   SortByMemoryLoad( $T$ )
3:   for each  $w \in vDomSet$  do
4:     if  $w.requestFreq < f_{target}$  then
5:        $cand.add(workload)$ 
6:     end if
7:   end for
8:   for each  $w \in T$  do
9:     if  $vDomSet.contains(w.id)$  then
10:       $continue$  ▷ Keep position
11:    end if
12:     $dest \leftarrow GetMinDistCore(cand, w.memCntrl)$ 
13:     $victim.add(dest)$ 
14:     $empty.add(w.id)$ 
15:     $migMap.update(w, dest)$ 
16:  end for
17:   $migMap \leftarrow PlaceVictim(victim, empty, migMap)$ 
18:  return  $migMap$ 
19: end function
```

initially placed will be empty. We call the empty cores as *empty*. The list of victims, *victim*, and empty cores, *empty*, is updated at line 9 and 10 at Algorithm 3. And update *workload*'s placement in migration map (line 11 in Algorithm 3).

For example, we have $Comb = \{\{vDom0, vDom1\}, \{vDom0, vDom2\}\}$ in Figure 4.1 (a) with target frequency f_{mid} . For a set $\{vDom0, vDom2\}$, the *Cand* are four cores which have f_{high} and T which will move are three cores in $vDom1$. The workload which uses $MC0$ and the other workload which uses $MC1$ are assigned to two cores which on the bottom of $vDom_0$ and a workload on the bottom left of $vDom_2$ which places close to each memory controller. The result will be like Figure 4.1 (b). There are three empty cores in $vDom_1$ with three victim cores on the right-hand side of the figure. To handle *empty* and *victim*, we call *PlaceVictim* (line 12 in Algorithm 3).

4.4.4 Assign victim

Algorithm 4 Place Victim OS

victim: A victim OS is a OS originally placed at a target OS's detination.

```
1: function PlaceVictim(victim, empty, migMap)
2:   SortByMemoryLoad(victim)                                ▷ Decending order
3:   for os ∈ victim do
4:     memCntrl ← os.memCntrl
5:     dest ← GetMinDistCore(empty, memCntrl)
6:     migMap.update(os, dest)
7:   end for
8:   return migMap
9: end function
```

Because *victim* lost its core, we should allocate *victim* to *empty*. This step is same as assigning *T* with *victim* and *empty* correspond to *T* and *Cand* (see Algorithm 4). Also, because *victim* and *empty* don't have original place and allocated *workload*, it does not generate any *victim* or *empty*. The result of example case in Figure 4.1 is Figure 4.1 (c) because *MC0* and *MC1* are closer from the and the bottom at *vDom₁*.

4.5 Workload Migration Evaluation Model

The energy for the next time quantum *t* of the status quo is computed as

$$E_{status_quo} = P_{status_quo} \cdot t \quad (4.2)$$

where P_{status_quo} can be obtained from the on-chip sensors or, in the absence of such, from Equation 2.2. Constants are obtained offline for each frequency. The expected

energy consumption if the migration is performed is given as follows

$$E_{migrated} = P_{migrated} \cdot (t + O_{migration} + O_{memory}) \quad (4.3)$$

$$O_{migration} = t_{migration} + t_{cache_fill}(f_{target}) \quad (4.4)$$

$$O_{memory} = t \cdot \frac{throughput_{status_quo}}{throughput_{migrated}} \quad (4.5)$$

where $P_{migrated}$ is computed based on offline power consumption data for each frequency level. The migration overhead, $O_{migration}$ is the overhead incurred by the actual migration and the (worst-case) time required to re-fill the entire caches at the target frequency f_{target} . The memory overhead, O_{memory} captures the sensitivity of an application to the location of the assigned core(s) on the CMP. The maximum throughput at each frequency and core location is profiled once offline; the actually required throughput of an application based on the core's last-level cache misses (as obtained by the core controllers).

The migration plan is only executed if the following equation holds

$$E_{status_quo} > E_{migrated} \cdot (1 + \Delta m) \quad (4.6)$$

that is, the expected benefit of migration has to be above a certain threshold Δm .

Chapter 5

Implementation

This chapter describes the implementation of the proposed cooperative hierarchical power management on a concrete hardware platform, the Intel Single-Chip Cloud Computer [15]. We first provide a short overview of the SCC platform and its capabilities and then describe the implementation in detail. The implemented application-specific runtime is a modified version of the sccLinux provided by Intel.

5.1 The Intel Single-chip Cloud Computer

The Intel SCC is a concept vehicle created by Intel Labs as a platform for many-core research. It consists of 48 independent cores interconnected by a routed network-on-chip (NoC). The cores are Intel P54C Pentium[®] cores with bigger L1 caches (16KB) and additional support for managing the on-chip scratchpad memory, the so-called *message passing buffer* (MPB). The Intel SCC provides no cache coherence for the core-local L1 and L2 caches. Always two cores are grouped together to form a *tile*; the 24 tiles are organized on a 6 by 4 grid. Four memory controllers in the four corners of

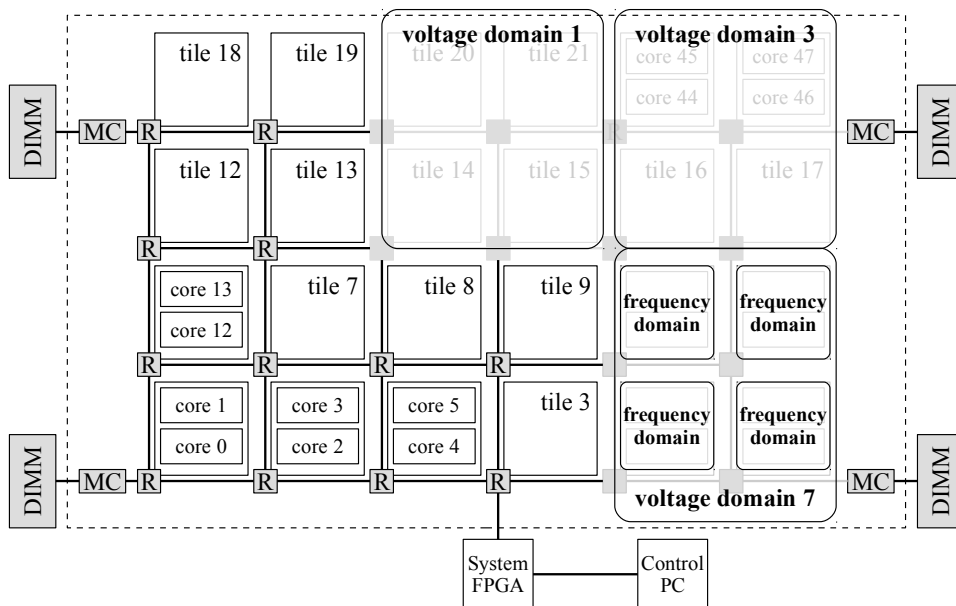


Figure 5.1: Intel SCC block diagram.

the chip provide access to up to 64 GB of memory. An FPGA provides the interface between the CMP and the management PC (MCPC). Figure 5.1 shows a block diagram of the SCC.

Memory Addressing. To support addressing up to 64GB of memory with 32-bit cores, the SCC implements the second level of indirection in the virtual-to-physical address translation. On the core, virtual-to-physical translation is performed as usual. The core-level physical addresses are then translated once again into system-level addresses through core-local lookup tables (LUT). With 64-bit cores, the LUT translation process will not be necessary anymore.

DVFS Capabilities. The SCC allows control over voltage and frequency for cores and the NoC. The frequency can be controlled *per tile*, that is, the two cores located at the same time always run at the same frequency and constitute a frequency domain (FD). The voltage can be regulated for a group of four tiles, i.e., a voltage domain

(VD) comprises a total of eight cores. The right upper hand of Figure 5.1 illustrates frequency and voltage domains on the SCC. In total, there are six voltage domains comprising four frequency domains à two cores each. The SCC supports seven different supply voltage levels. However, only four are of practical interest: 1.1V to run at a frequency of 800MHz, 0.9V to run at 533MHz, 0.8V for 400MHz, and 0.7V for frequencies between 320 and 100MHz.

Power Measurement. The SCC provides a number of voltage and ampere meters on-board. The total power consumed by the SCC chip is obtained by multiplying the (constant) supply voltage with the supply current for the entire SCC chip. The power consumption of individual voltage domains cannot be computed because only the per-domain supply voltage is available but not the current consumed by the domain. We thus always report the total chip power in our experiments in Chapter 7.

5.2 Implementing Workload Migration

The Intel SCC provides no means to read/write core-local registers from outside a core. A minimal level of cooperation is thus required by the application-specific runtime. Here, we first describe the logical steps necessary to re-assign a core to a new application container and then discuss concrete implementation details.

5.2.1 Migration Steps

Figure 5.2 illustrates the necessary steps to carry out a migration plan computed by the chip controller (Chapter 4.3). The migration manager first signals all kernels that are about to be migrated through an interrupt. Upon receipt of a migration interrupt, the OSes first save their complete volatile state of the core to a designated area in shared memory and set a flag to indicate completion of saving the state. They then flush the TLB and the caches, and then enter a busy loop, waiting for a flag set by the global mi-

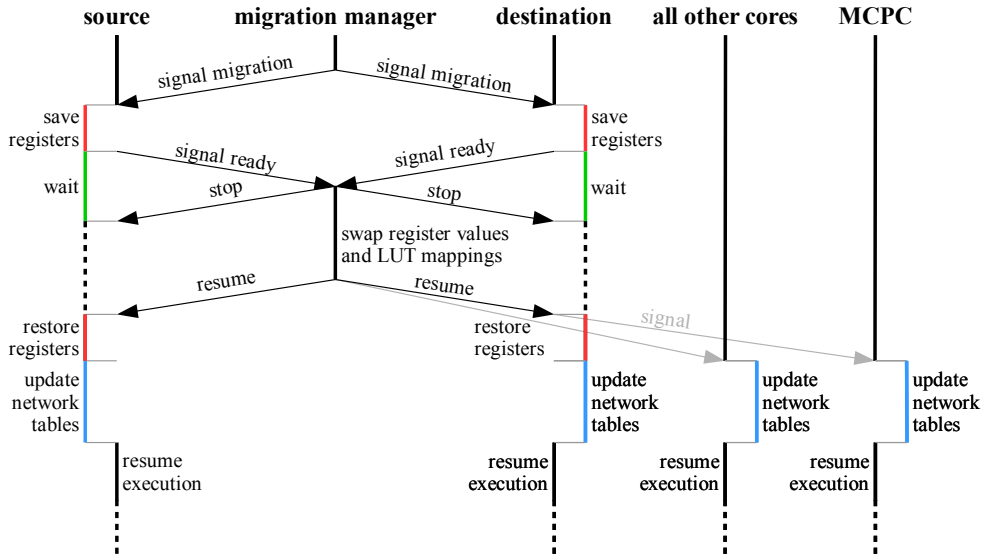


Figure 5.2: Workload migration.

gration manager to continue. Upon continuation, the volatile states are restored from the designated area, and the kernels return from the interrupt and continue execution. The migration manager waits for all kernels to save their volatile state and enter the busy loop. Before setting the completion flag, thereby allowing the cores to continue, the manager exchanges the contents of the saved volatile states of the migrated kernels with that of each target core. This means, in effect, that an entire kernel can be migrated from one core to another by copying a few hundred bytes of volatile state. This process is not much different from task switching with the difference that kernels are not scheduled in or out but rather swapped. To maintain consistent networking state, all cores, including the MCPC need to update internal network routing tables to reflect the new locations of the cores (Chapter 5.2.2).

5.2.2 Networking

The SCC provides two separate networks: one network for on-chip networking, and a subnet for communication with the MCPC. The target core of a network interrupt is identified by its physical core ID which corresponds to the x/y -coordinates of the core on the grid. In the original sccLinux the interrupt target ID is computed from the core ID. In order to support transparent migrations, we have added a table holding the current IP-to-coreID mappings in each kernel. After each migration, the migration manager notifies *all* cores about the changes to the IP-to-coreID mapping tables. The same method is used on the MCPC. These simple modifications are enough to keep networking, including open connections, alive across migrations. DMA is not supported, and no other devices exist on the SCC; input/output, including access to permanent storage, are routed through the network.

5.3 Domain Controller Implementation

The domain controllers (core, frequency, voltage, and chip) are implemented in C and are present on each kernel. As outlined in Chapter 3.3.5, the physical core ID determines which controllers are (de-)activated in a kernel. The reason is workload migration. We implement workload migration as a OS migration that makes OSes float around in frequency and voltage domains. If we assign a domain controller to a certain OS, the OS might control a domain even if it is not in the domain and every domain controllers can not distinguish where the information are come from. After migration and before returning from the migration interrupt, kernels check if the core they are running on requires activation/deactivation of one of the four controllers. Core controllers are active on every kernel. The 24 frequency controllers are activated on the cores with an odd core ID. The six voltage controllers run on the lower-left core of each domain (core IDs 0, 4, 8, 24, 28, and 32). The chip controller runs on core 30.

Chapter 6

Experimental Setup

6.1 Hardware

All experiments were conducted on the Intel SCC [15]. The chip controller and other services such as monitoring logging, run on dedicated cores in voltage domain 1. The microkernels run a modified version of sccLinux that supports kernel migration and dynamic IP-to-coreID mappings. We chose this separation in order to separate the power consumption of the core OS from the application containers, voltage domain 1 does not participate in workload migration. However, the SCC only allows measuring the *total chip power*; the power consumption of the OS services are therefore also included in all results. Power consumption is computed using the on-chip voltage and ampere meters. The meters are queried 10,000 times per second. Power is computed by multiplying the measured chip supply voltage by the current. This includes the power consumed by all 48 cores and the NoC. In particular, since the power manager is implemented entirely in software and runs on the cores of the SCC, the power measurements include all the overhead caused by the propose power manager.

Table 6.1: Datacenter scenarios: distinct workloads patterns

Scenario	G1	G2-G5	G6	G7-G11
# patterns	4	7	10	40

Table 6.2: Average CPU and memory load

Average	CPU load	Mem load
G1	34	14
G2	47	14
G3	40	13
G4	37	14
G5	39	16
G6	39	17
G7	42	17
G8	37	16
G9	41	17
G10	39	14
G11	39	18

6.2 Benchmark Scenarios

A benchmark scenario is defined by a mapping of a number of workload patterns to a number of cores. Depending on the scenario, we map from 2-40 different patterns onto 8-40 cores. Cores with no assigned workload only run the modified sccLinux kernel and domain managers depending on the core location (Section 5.3).

In this evaluation, we focus on workload scenarios occurring on the servers of a datacenter. We have created three synthetic benchmark scenarios composed of synthetic workloads in order to explore the best and worst cases and show the effect of NUMA awareness for the proposed technique. The workload patterns of the datacenter scenarios are based on the Google cluster data [36]. For the average CPU usage and memory intensity, we used the information of mean CPU usage rate and the memory accesses per instruction (MAI) from the data set. We add up the number of individual processes running on the same physical machine to obtain a real-world workload of a

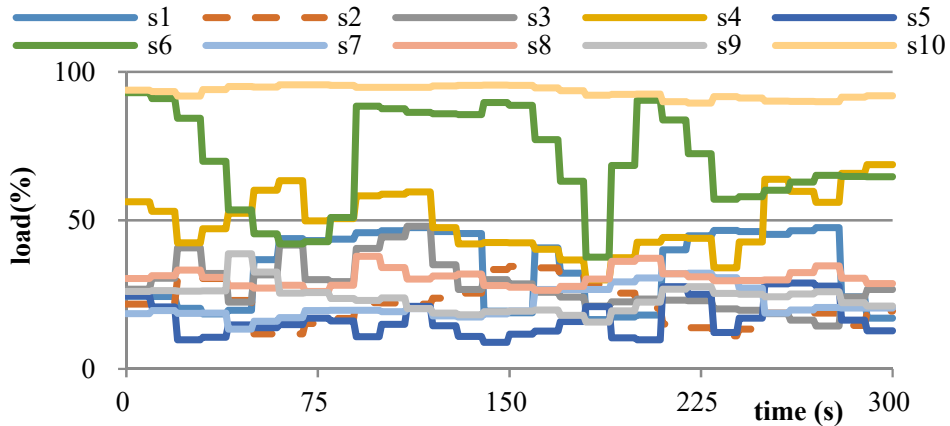


Figure 6.1: G6 workload patterns.

server over time. We scaled profiled time from average 5000 seconds to 300 seconds. We convert the numbers to a sequence of average utilization rates per 10 seconds. Then we assign the sequence as a workload pattern of a machine per a core. Each scenario has multiple workload patterns with different numbers. By this setup, we simulate not only multi-threaded applications but also multi-program environment. We have generated a total of 11 scenarios based on the Google cluster data, Table 6.1 lists the number of distinct patterns per scenario. The distinct patterns are assigned to a varying number of cores; details are given in Chapter 7. Figure 6.1 shows the 10 distinct patterns of the G6 scenario as an example. For simplicity we only display the CPU load; the memory load shows similar patterns. The average CPU and memory load of benchmark scenarios are in Table 6.2. Also, we have tested with a different number of containers which have workloads to show how it is hard to save power consumption with DVFS without workload migration according to the number of workloads increased.

6.3 Comparison of Results

The baseline of the experiments is result which is obtained by running the benchmark scenario on the SCC at full speed (800MHz) with no power management enabled. Unlike the work in [17] we do not use a phase-detector based on message passing since we are aiming at independent workloads running on a CMP. The workload of a kernel is estimated based on a weighted average of the past 10 measurements.

To show impact of workload migration and NUMA-awareness, we compare the presented NUMA-aware power management technique with the DVFS-only approach of Ioannou *et al.* [17] and the DVFS+migration technique with its locality-unaware buyer-seller algorithm described by Kang *et al.* [21]. The hierarchical framework and the DVFS policies for all three methods are identical. We evaluated the different core migration algorithms using the DVFS policies Allhigh and Tile (Section 4.1).

For all methods and benchmarks scenarios, the migration benefit threshold Δm is set to 10%. Because we want to keep overhead less than 1% of an epoch, with latencies $\leq 3ms$, $\leq 10ms$, and a few thousand cycles for migration, voltage change, and frequency changes, migrations are evaluated and performed once every 3 seconds. All benchmark scenarios are executed for 300 seconds. The reported results are the average of at least 3 runs executed at similar thermal conditions. Also, to reduce the effect of temperature, we use results which run in similar temperature for each scenario.

Chapter 7

Results

We have conducted a wide range of experiments to evaluate our proposed power management technique which is NUMA-aware and hierarchical. To show the potential of our technique, we compare it with state-of-the-art methods which are using DVFS-only [17] and NUMA-unaware workload migration [21] on synthetic benchmark scenarios. The real-world server workloads obtained from Google cluster data [36] are then used to compare the three techniques in more realistic workload scenarios. For these workloads, we conducted experiments with different number of workload patterns and different number of workloads to show the effects of the differences. At last of this chapter, we conclude this section with the overall results overall benchmark scenarios.

7.1 Synthetic Scenarios

We first compare the *DVFS only* method [17] with a data-locality-unaware *Buyer-Seller* migration algorithm [21] and our NUMA-aware *Greedy* migration technique in terms of the performance per watt (PPW) at equal turnaround time. The goal is to

Table 7.1: Normalized PPW for synthetic scenarios

BM	DVFS only		Buyer-Seller		Greedy	
	AH	T	AH	T	AH	T
<i>SynMem</i>	1.12	1.23	1.51	1.52	1.67	1.68
<i>SynCpu</i>	1.30	1.33	1.59	1.58	1.61	1.60
<i>SynRandom</i>	1.00	1.01	1.00	1.00	1.02	1.04

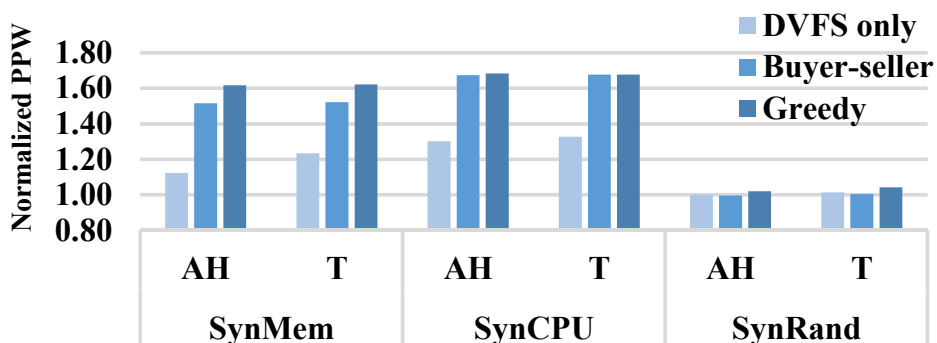
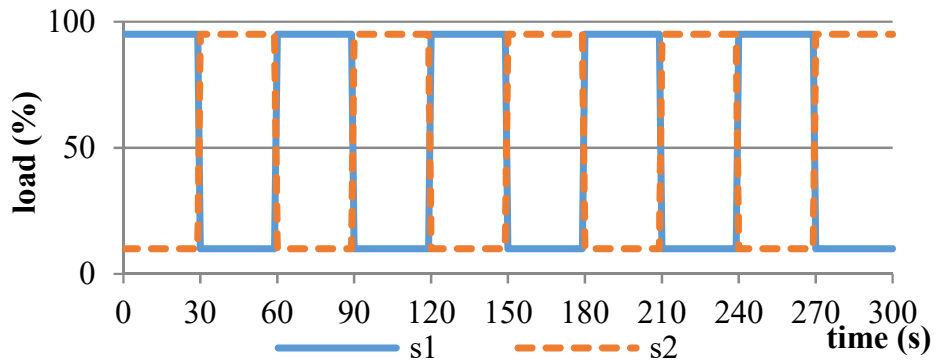
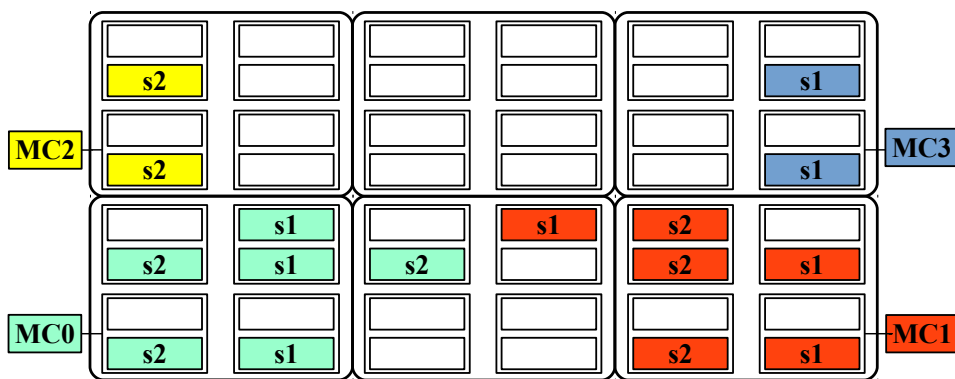


Figure 7.1: Synthetic scenario PPW.

show the necessity and the potential of NUMA-aware migration. For this, we have crafted three synthetic scenarios running synthetic workload patterns. Figure 7.2 (a) shows the workload patterns of *SynMem* and *SynCPU*. Two patterns denoted *s1* and *s2* show alternatively high and low utilizations. The patterns are crafted such that when *s1* shows a high utilization *s2* has a low load and vice-versa. The initial distribution of the workload patterns to the target architecture is shown in Figure 7.2 (b). The label refers to the workload pattern, and the coloring shows the affinity of the different workloads to the respective memory controllers. For example, *s2* running on core 0 in the left-bottom corner has its data located in memory controller *MC0*. This workload patterns and distribution make DVFS cannot lower voltage because *s1* and *s2* alternatively need high frequency. *SynCPU* and *SynMem* differ in that the former is CPU-bound (i.e., has very little memory accesses) whereas the latter is memory-bound. We expect that



(a) Load pattern



(b) Load Distribution

Figure 7.2: Synthetic scenario workload patterns.

migration outperforms DVFS on both scenarios and that especially for *SynMem* the proposed NUMA-aware algorithm achieves a better PPW. *SynRandom* is the worst-case scenario. It comprises forty distinct random workload patterns. With completely random utilizations and full occupation, migration is not expected to perform much better than DVFS only.

Table 7.1 shows the results of the synthetic scenarios for each of the three algorithms, we report the normalized PPW with respect to the baseline (no DVFS) for the *AH (Allhigh)* and *T(ile)* DVFS policy. We observe the expected behavior: for *Syn-*

CPU the migration-based algorithms outperform *DVFS only* by around 35%, and there is no significant difference between *Buyer-Seller* and *Greedy*. This result confirms the importance of workload migration on MVMF CMPs.

For *SynMem*, data locality comes into play. Even if both migration policies are effective, *Buyer-Seller* outperforms *DVFS only* by a similar margin as *SynCpu*. *Buyer-Seller* congregates high workloads to same voltage at a time, makes rest voltage domains lower power. In this case, however, the NUMA-aware *Greedy* algorithm is able to improve the PPW by 16% over *Buyer-Seller*, emphasizing the need to consider data-locality to achieve maximal power savings. This is because *Buyer-seller* algorithm does not consider memory locality.

Especially in case of *SynMem*, the result of *Buyer-seller*'s workload migration collect workload *s1* to left bottom voltage domain and *s2* to right top domain. This workload mapping makes containers apart from memory controllers, which lead to high memory access delay. On the contrary, our greedy algorithm assigns *s1* to right bottom and *s2* to middle bottom voltage domain which make lower memory access delay than the *Buyer-seller* case.

We can find a clue in the result of *SynCpu* which has same workload patterns and initial distribution with *SynMem* except it has CPU workload. *SynMem* shows the Greedy algorithm has an advantage over the *Buyer-Seller* algorithm, *SynCpu* shows similar results in both workload migration policies. In this context, we can infer the importance of memory locality.

For *SynRandom*, there is not much room for improvement for any algorithm. *DVFS only* and *Greedy* fail to improve the PPW compared to no power management, while the proposed *Greedy* algorithm improves the PPW by a few percent only. With 40 random workloads and no migration and the constraint of equal turnaround time, *DVFS only* is unable to apply DVFS.

The DVFS policies evaluated in this paper do not trade performance for power. As

a result there is no noticeable slowdown for any of the three algorithms, and the numbers have been omitted for brevity. Table 7.2 show the performance loss in numbers for the 11 datacenter scenarios.

7.2 Datacenter Scenarios

7.2.1 Varying Number of Workloads

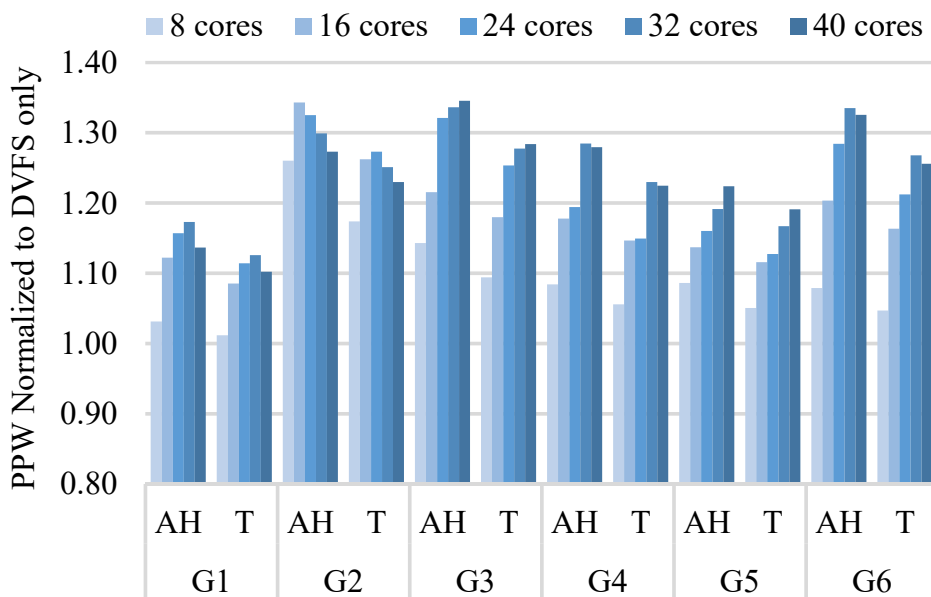


Figure 7.3: PPW for a varying number of workloads.

We first evaluate the real-world datacenter scenarios with respect to a varying number of assigned workloads from 8 to 40 in increments of 8. The number of distinct workload patterns for each scenario is given in Table 6.1; patterns are randomly assigned to the number of workloads (i.e., for the 8-workload case and *G1* we make 8 random selections from the pool containing the four workload patterns). The initial location of the workloads on the chip can affect the result; we create three different

random assignments and report the average of running each of the tree assignments three times. In other words, each individual result represents the average of nine runs.

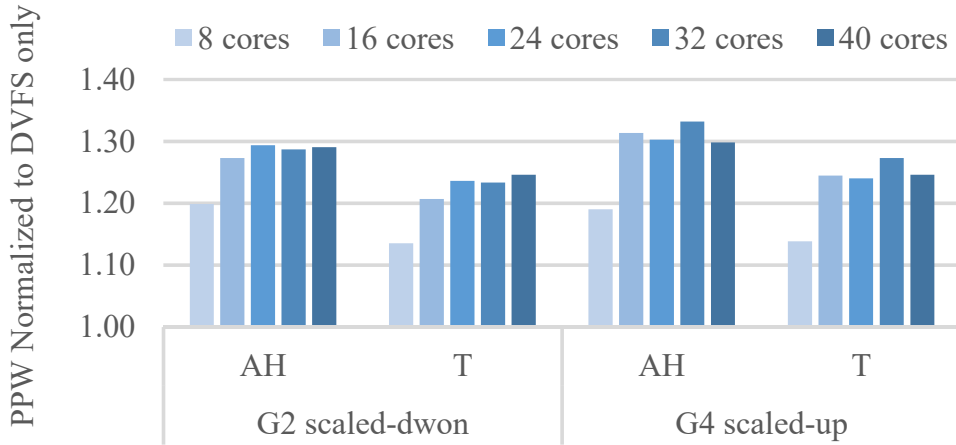


Figure 7.4: PPW for scaled scenarios with a varying number of workloads.

Figure 7.3 displays the results for the datacenter scenarios $G1$ to $G6$ with 8, 16, 24, 32, and 40 workloads running simultaneously. The Y-axis shows the PPW of the proposed greedy algorithm relative to *DVFS only*. We observe that *Greedy* shows better relative improvements if the number of active workloads (i.e., active cores) is 32 cores. In the case of 8 workloads, *DVFS only* manages to do quite a good job (50% over the baseline) despite its inability to migrate workloads because the low occupation still provides sufficient opportunities to apply DVFS. The reason is the number of workloads removes conflicting frequency needs. On the other end of the spectrum with 40 cores there are less opportunities for power savings with or without migration. The best case are moderately loaded CMPs where the proposed greedy algorithm outperforms *DVFS only* by 25% on average.

For $G2$, it shows different tendency with increasing number of workloads. The reason why it shows different result is average utilization. The other scenarios have average utilization 35%~41%, but $G2$ has 49% average utilization. High workload av-

erage reduces space which we can save power. To prove it, we have conducted experiments with *G2* which is 25% scaled down. Figure 7.4 shows the result of scaled-down *G2* which similar with the other scenarios. To show the dependency between PPW and average workload, we have conducted scenario *G4* with scaled up version. The result is showed in Figure 7.4 which decreased gap between the varying number of workloads.

Also, the workload migration shows less effective in scenario *G1*. That's because DVFS only policies also save significant power. It based on the minimum utilization (in this scenario group), and the number of workload patterns (see Table 6.1). Less number of workload patterns make a voltage domain be more likely to have similar workload patterns in it, not only for *DVFS only* but also for both migration policies. It is also shown on Table 7.2, which is the highest PPW for every policies.

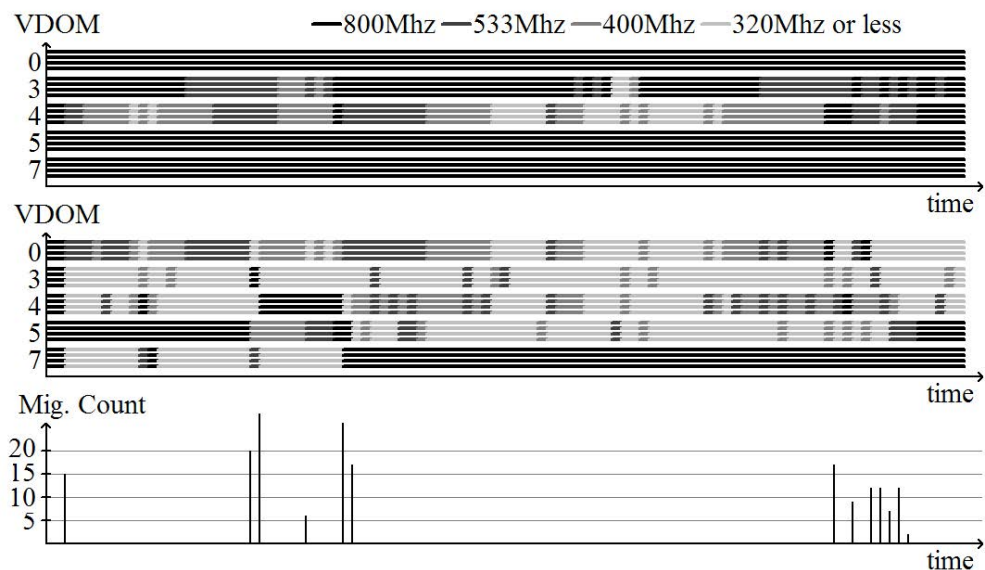


Figure 7.5: Frequency map example for *G6* and *Allhigh*.

The effect of workload migration is visualized in Figure 7.5. The topmost graph shows the frequency map of *DVFS only* with the *Allhigh* policy for the different voltage domains. Darker colors represent higher frequencies. The middle graph shows the fre-

quency map for the same workload with the proposed *Greedy* algorithm. While *DVFS only* is required to run most domains at a high frequency for most of the time, we observe that *Greedy* is able to group workloads with similar utilization into a few domains and apply aggressive DVFS on the lightly loaded domains. The bottom graph in Figure 7.5, finally, shows the number of workload migrations over time.

7.2.2 Independent Workloads

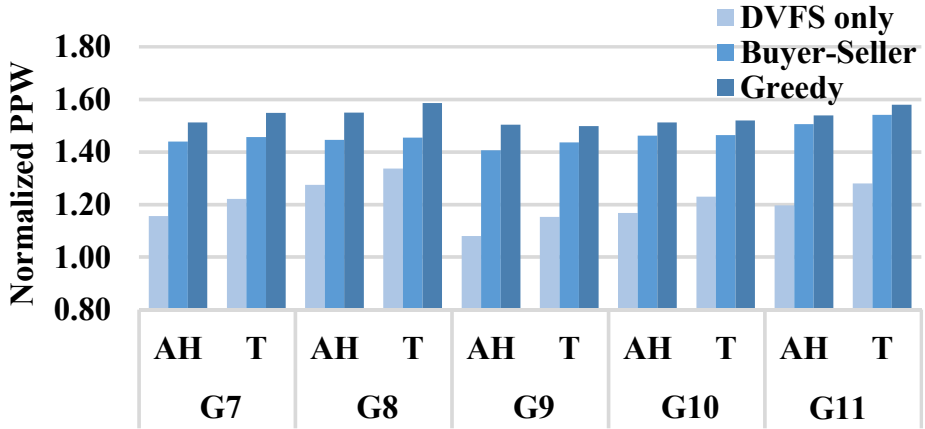


Figure 7.6: Experiment results for *G7* to *G11*.

Table 7.2: Normalized performance per watt for Google cluster data scenarios

BM	DVFS only				Buyer-Seller				Greedy			
	PPW	AH Perf Loss	PPW	T Perf Loss	PPW	AH Perf Loss	PPW	T Perf Loss	PPW	AH Perf Loss	PPW	T Perf Loss
G1	1.58	0.00%	1.61	0.00%	1.78	0.00%	1.76	0.00%	1.79	0.00%	1.77	0.01%
G2	1.00	0.00%	1.05	0.03%	1.26	0.02%	1.29	0.07%	1.28	0.02%	1.29	0.08%
G3	1.04	0.00%	1.10	0.00%	1.38	0.00%	1.39	0.04%	1.40	0.00%	1.41	0.06%
G4	1.23	0.00%	1.28	0.00%	1.55	0.00%	1.55	0.07%	1.57	0.01%	1.57	0.05%
G5	1.32	0.00%	1.37	0.00%	1.57	0.00%	1.59	0.30%	1.62	0.00%	1.64	0.04%
G6	1.16	0.00%	1.23	0.04%	1.51	0.51%	1.51	0.53%	1.54	0.20%	1.54	0.26%
G7	1.16	0.00%	1.22	0.00%	1.44	0.00%	1.46	0.00%	1.51	0.00%	1.55	0.00%
G8	1.27	0.00%	1.34	0.00%	1.45	0.00%	1.45	0.00%	1.55	0.00%	1.59	0.03%
G9	1.08	0.00%	1.15	0.00%	1.41	0.00%	1.44	0.01%	1.50	0.00%	1.50	0.04%
G10	1.17	0.00%	1.23	0.00%	1.46	0.00%	1.46	0.00%	1.51	0.00%	1.52	0.05%
G11	1.20	0.00%	1.28	0.00%	1.51	0.00%	1.54	0.00%	1.54	0.00%	1.58	0.05%
AVG	1.20	0.00%	1.26	0.01%	1.48	0.05%	1.49	0.09%	1.53	0.02%	1.54	0.06%

Figure 7.6 shows the normalized performance-per-watt over the baseline for data-center scenarios *G7-G11*. Each scenario is composed of 40 independent server workloads as recorded in Google’s datacenters which can be more realistic in real world. Overall, we observe that the proposed *Greedy* outperforms *DVFS only* by a large margin, once again emphasizing the importance of workload migration. Compared to *Buyer-Seller*, the NUMA-awareness of *Greedy* pays off in a 8% better energy efficiency. The reason for the larger gap is because as the number of workload patterns increase, the probability which a voltage domain get high frequency request increases. Also, more significant number of patterns needs more frequent workload migration. The *Buyer-Seller* algorithm which doesn’t consider memory locality keeps cause high memory access delay in contrast with the greedy algorithm which make cores need higher frequency. As a result, in scenario *G7* to *G11*, performance per watt gap between the greedy and the Buyer-seller are increased as 7% for *AH* and 8% for *T* from 2% in case of *G1* to *G6*.

7.3 Overall Results Comparison

Table 7.1 and Table 7.2, finally, displays the performance per watt (PPW) and the performance loss over the baseline policy, respectively, for the *Allhigh* and the *Tile* policy for *DVFS only*, *Buyer-Seller*, and the proposed *Greedy* algorithm. Each scenario is run with the number of workloads listed in Table 6.1. Every result is the average of evaluated at least three times. For the real world workloads, we have also tested three random initial placements of workload patterns.

From synthetic workloads, Table 7.1, *SynMem* and *SynCpu* are periodic workloads scenarios. For *SynMem*, we use long period length and large memory workloads. It shows large gap not only between DVFS-only policies and DVFS with migration policies but also *Buyer-Seller* and *Greedy* algorithm which we introduce. For *SynCpu*,

we use long period length and large CPU workloads. It shows, also, the large gap between *DVFS-only* and *DVFS with workload migration* policies, but there is little gap between *Buyer-Seller* and *Greedy*. The *SynRandom* is a scenario consist of randomly generated workload patterns for each 40 application containers. The workload patterns change every 3 seconds. For this workload scenario, we get minimum improvement at performance per watt, 2%-4%.

In case of the real world workloads, overall, the NUMA-aware *Greedy* algorithm outperforms *DVFS only* by about 30% and *Buyer-Seller* by 5%. We observe that *Tile* outperforms *Allhigh* without migration whereas with migration they achieve similar performance. The reason is that OS migration can group OSes with similar performance requirements into voltage domains such that the superior *Tile* DVFS policy has less effect. The performance degradation is negligible for all three policies. The algorithms that support migration suffer from a slightly higher performance degradation (0.05% on average over no degradation with *DVFS only*), but the slowdown is insignificant compared to the 30% improved energy efficiency.

Chapter 8

Discussion

In this section, we discuss about limitations and extra hardware support to get more improvements of this work.

8.1 Limitations

The first topic is the limitations of this work. One of the limitations of this work is on an assumption. We assume each core has single workload on it which is not true on real system. It make scheduling workload migration more difficult. Because, if there is multiple processes which uses different memory controller with each other in a core. Then, workload migration scheduling would be more complicate.

Anther limitation of this work is scalability of algorithm. This algorithm runs at M number of voltage levels traverses combination of voltage domains which can be represented as

$$O(M * N^{\min\{K, N-K\}}) \quad (8.1)$$

where N and K are the number of voltage domain which we can select and we need to

select. In Intel Single Chip Cloud Computer (SCC), our greedy algorithm take around 0.002 second, because, we have only four voltage levels and five voltage domains. But, if a system which have more fine granularity of voltage domain, the computation time would take longer.

The point we doesn't considered is temperature. In our framework, we collecting high workloads into same voltage domains. It could increase heat of the voltage domains which could make problems. This is a problem one of problems considered in future work.

8.2 Extra Hardware Support

In the system we use, Intel SCC, there is some hardware characteristics which make bottleneck. In this section, we discuss extra hardware support which can solve this problems and achieve more improvements from this work.

In the Intel SCC, dynamic voltage and frequency scaling (DVFS) too have to use more frequently ($\leq 10ms$). If we can have more light weight DVFS (less delay), we can response more quickly the cores' request changes.

Also, we implemented OS migration for workload migration. OS migration make more overhead than process migration because of swapping LUT and updating network tables. If we work with process migration, we can have less overhead on workload migration and better results.

Chapter 9

Conclusion

We have presented a NUMA-aware cooperative hierarchical power management technique for existing and future many-core systems. The technique employs workload migration. Without explicit hardware support, the microkernels running on the individual cores cooperate with the global power management by saving and restoring the volatile state of the core on demand. Combined with dynamic monitoring of each core's performance metrics this technique allows the power manager to group cores with similar performance requirement together so that traditional DVFS policies can apply DVF settings closer to the optimal value. In order to remain scalable, the power manager is implemented in a hierarchical fashion, logically re-creating the hierarchy imposed by the hardware through the different power management domains.

The cooperative power manager has been implemented and evaluated on a real system, the Intel Single-Chip Cloud Computer. Experiments with a wide range of real world workload benchmark scenarios show that, on average, the proposed technique outperforms existing DVFS policies by 30% and by 5% compared to a NUMA-unaware approach at the expense of little performance loss, less than 1%.

Appendix A

Benchmark Scenario Details

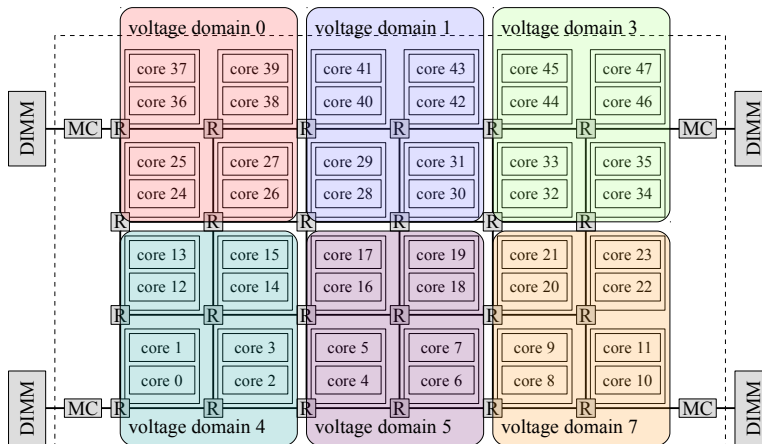


Figure A.1: SCC core map

In this appendix, we describe details of the benchmarks used for experiments. We have total 14 benchmark scenario, 3 for synthetic and 11 for real-world benchmarks. Each benchmark has workload patterns and distribution of workloads on cores. We

describe *synthetic benchmarks* and *real-world benchmarks* in Section A.1 and Section A.2.

Each workload pattern has CPU load and memory load. We described these in the tables below. The numbers represent ratio to maximum performance of a core at the highest frequency and the closest memory distance.

Benchmark has multiple workload patterns but the number of patterns is different. For benchmarks which have 40 workload patterns, each pattern initially placed at a core which has the same number which represented in Figure A.1 (a block diagram of Intel SCC). On the other hand, for benchmarks which have workload patterns less than 40, multiple cores can have same workload pattern. So, we describe initial workloads distribution in the tables followed by workload pattern table of each benchmark.

In our setup, voltage domain 1 is a control domain. So, it runs control processes and various measurement services. That's why cores in the domain don't get any workloads.

A.1 Synthetic Benchmark

Table A.1: *SynMem* benchmark scenario

(a) Workload pattern

WL		Epoch (1 epoch = 15 sec)																			
		00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19
S1	CPU	0	10	10	0	0	10	10	0	0	10	10	0	0	10	10	0	0	10	10	0
	Mem	95	0	0	95	95	0	0	95	95	0	0	95	95	0	0	95	95	0	0	95
S2	CPU	10	0	0	10	10	0	0	10	10	0	0	10	10	0	0	10	10	0	0	10
	Mem	0	95	95	0	0	95	95	0	0	95	95	0	0	95	95	0	0	95	95	0

(b) Workload distribution

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
n/a	n/a	n/a	n/a	n/a	s1	n/a	s1	n/a	s1	s2	n/a
s2	n/a	n/a	n/a	n/a	n/a	s2	s1	s2	n/a	s2	s1
n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
s2	n/a	n/a	n/a	n/a	s1	s2	s1	n/a	n/a	s2	s1

Table A.2: *SynCPU* benchmark scenario

(a) Workload pattern

WL		Epoch (1 epoch = 15 sec)																			
		00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19
S1	CPU	95	95	10	10	95	95	10	10	95	95	10	10	95	95	10	10	95	95	10	10
	Mem	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
S2	CPU	10	95	95	10	10	95	95	10	10	95	95	10	10	95	95	10	10	95	95	10
	Mem	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

(b) Workload distribution

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
s2	n/a	n/a	n/a	s2	n/a	s2	s2	s2	n/a	s2	n/a
n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
s1	s2	n/a	n/a	s1	s1	s1	s1	s1	s2	s1	s1

A.2 Real World Benchmark

Table A.5: Google cluster data benchmark scenario #1

(a) Workload pattern

Workload		Epoch (1 epoch = 10 sec)																													
		00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29
S1	CPU	1	3	3	33	41	41	41	44	47	47	47	44	44	44	43	40	37	42	46	51	51	51	61	69	69	69	76	79	79	79
	Mem	1	1	1	8	10	10	10	13	18	19	20	25	24	25	23	17	10	12	11	10	10	10	15	18	18	18	30	36	36	36
S2	CPU	36	36	36	36	19	16	20	16	18	33	36	36	36	35	29	23	21	24	26	30	38	38	38	38	38	44	44	41	34	34
	Mem	11	11	11	11	7	9	9	8	8	19	25	25	25	24	21	22	21	22	24	25	26	26	26	26	27	32	32	31	27	27
S3	CPU	35	35	35	35	35	35	31	27	27	27	27	27	27	27	27	27	27	27	27	27	23	15	13	13	13	13	13	13	13	13
	Mem	7	7	7	7	7	7	6	6	6	6	6	6	6	6	6	6	6	6	6	6	5	3	3	3	3	3	3	3	3	3
S4	CPU	32	32	33	33	33	33	33	33	33	33	33	33	33	33	34	34	35	34	32	32	32	33	33	33	33	33	33	33	32	31
	Mem	10	10	10	10	11	10	10	10	11	11	11	10	11	11	11	11	12	11	11	10	10	11	11	11	11	10	11	11	11	11

(b) Workload distribution #1

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s4	s4	n/a	n/a	s4	s4	s1	s1	s2	s4	s2	s3
s3	s3	n/a	n/a	s3	s4	s3	s3	s4	s1	s2	s2
s2	s3	n/a	n/a	s2	s2	s2	s4	s2	s2	s3	s4
s1	s1	n/a	n/a	s1	s3	s1	s4	s1	s3	s1	s1

(c) Workload distribution #2

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s1	s3	n/a	n/a	s2	s2	s3	s2	s3	s4	s4	s1
s2	s1	n/a	n/a	s4	s2	s1	s4	s4	s1	s1	s1
s1	s4	n/a	n/a	s2	s2	s3	s4	s4	s4	s3	s1
s2	s2	n/a	n/a	s3	s3	s2	s3	s1	s4	s3	s3

(d) Workload distribution #3

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s4	s4	n/a	n/a	s4	s2	s1	s2	s3	s1	s4	s3
s3	s2	n/a	n/a	s4	s4	s1	s4	s3	s3	s4	s3
s1	s2	n/a	n/a	s2	s2	s3	s4	s4	s1	s1	s2
s3	s1	n/a	n/a	s1	s2	s3	s2	s1	s2	s3	s1

Table A.6: Google cluster data benchmark scenario #2

(a) Workload pattern

Workload	Epoch (1 epoch = 10 sec)																														
	00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	
S1	CPU	1	3	3	33	41	41	44	47	47	47	44	44	44	43	40	37	42	46	51	51	51	61	69	69	69	76	79	79	79	
	Mem	1	1	1	8	10	10	10	13	18	19	20	25	24	25	23	17	10	12	11	10	10	10	15	18	18	18	30	36	36	36
S2	CPU	38	48	75	86	82	92	85	87	91	93	91	93	89	73	36	25	26	26	28	28	25	24	26	21	19	21	23	31	18	31
	Mem	15	18	21	27	28	32	30	30	31	34	31	31	31	29	19	17	16	15	16	17	16	15	14	13	10	10	10	15	11	22
S3	CPU	16	9	5	4	7	47	80	92	78	83	91	86	89	87	79	91	91	93	93	90	91	91	84	86	66	59	69	79	62	47
	Mem	6	4	3	2	2	9	15	16	15	16	18	14	15	15	16	16	16	17	15	14	13	12	14	18	16	19	18	16	12	12
S4	CPU	85	82	83	69	67	84	49	63	81	92	90	76	70	58	18	14	14	15	14	14	12	10	15	24	24	25	28	27	69	
	Mem	16	14	16	23	29	35	22	18	35	51	58	39	24	21	10	6	7	8	7	7	7	6	5	5	7	7	9	8	17	
S5	CPU	27	25	60	79	82	85	84	84	85	87	83	60	62	77	80	84	82	84	85	89	82	69	19	13	12	13	13	15	15	25
	Mem	14	12	19	23	25	29	25	26	26	28	26	21	22	25	27	28	28	28	29	30	30	28	13	10	9	9	11	10	17	
S6	CPU	28	28	26	25	26	25	27	30	29	28	25	24	24	25	27	26	31	28	29	29	28	25	24	23	23	26	25	25	28	26
	Mem	18	20	20	17	15	14	17	18	19	17	15	15	15	19	21	18	21	20	22	22	21	19	17	13	15	16	16	17	17	18
S7	CPU	35	55	40	38	49	66	38	23	19	25	24	20	22	34	31	37	35	28	32	28	27	22	35	43	24	27	21	25	30	28
	Mem	14	19	20	17	21	29	19	12	10	14	12	10	10	13	16	17	14	15	18	14	13	14	13	22	10	12	10	10	11	13

(b) Workload distribution #1

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s4	s6	n/a	n/a	s4	s5	s3	s5	s5	s5	s5	s7
s3	s5	n/a	n/a	s3	s4	s6	s5	s4	s4	s2	s1
s2	s3	n/a	n/a	s2	s4	s2	s4	s3	s6	s3	s7
s1	s1	n/a	n/a	s1	s2	s1	s3	s1	s2	s1	s2

(c) Workload distribution #2

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s3	s7	n/a	n/a	s4	s3	s5	s4	s7	s3	s7	s7
s5	s5	n/a	n/a	s6	s7	s1	s2	s1	s4	s5	s3
s3	s6	n/a	n/a	s2	s2	s6	s6	s6	s2	s1	s2
s4	s5	n/a	n/a	s5	s1	s3	s1	s1	s4	s6	s4

(d) Workload distribution #3

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s2	s4	n/a	n/a	s5	s4	s4	s6	s4	s3	s1	s1
s7	s6	n/a	n/a	s2	s7	s1	s2	s6	s5	s1	s2
s1	s5	n/a	n/a	s7	s5	s7	s3	s5	s2	s7	s6
s6	s1	n/a	n/a	s2	s3	s4	s3	s3	s5	s6	s3

Table A.7: Google cluster data benchmark scenario #3

(a) Workload pattern

Workload	Epoch (1 epoch = 10 sec)																														
	00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	
S1	CPU	1	3	3	33	41	41	44	47	47	47	44	44	44	43	40	37	42	46	51	51	51	61	69	69	69	76	79	79	79	
	Mem	1	1	1	8	10	10	10	13	18	19	20	25	24	25	23	17	10	12	11	10	10	10	15	18	18	18	30	36	36	36
S2	CPU	38	48	75	86	82	92	85	87	91	93	91	93	89	73	36	25	26	26	28	28	25	24	26	21	19	21	23	31	18	31
	Mem	15	18	21	27	28	32	30	30	31	34	31	31	31	29	19	17	16	15	16	17	16	15	14	13	10	10	10	15	11	22
S3	CPU	51	51	44	45	49	51	58	40	39	41	40	38	30	25	22	23	24	23	22	35	44	43	43	54	38	19	39	57	52	50
	Mem	9	9	10	9	8	9	10	6	6	8	7	5	7	7	7	7	6	6	6	6	6	7	8	10	7	5	9	12	9	8
S4	CPU	91	91	90	87	84	86	87	86	73	12	3	0	0	0	0	2	31	1	3	4	2	2	21	26	28	28	28	28	28	
	Mem	38	38	39	38	37	37	38	36	30	1	0	0	0	0	0	1	43	1	1	1	1	1	3	3	3	3	3	3	3	
S5	CPU	27	25	60	79	82	85	84	84	85	87	83	60	62	77	80	84	82	84	85	89	82	69	19	13	12	13	13	15	15	25
	Mem	14	12	19	23	25	29	25	26	26	28	26	21	22	25	27	28	28	28	29	30	30	28	13	10	9	9	9	11	10	17
S6	CPU	1	1	1	10	30	36	38	40	40	39	40	41	40	40	40	40	41	35	25	8	1	1	1	1	1	1	1	1	1	1
	Mem	1	1	1	2	3	6	7	5	5	5	5	5	7	8	5	5	5	4	3	2	2	1	1	1	1	1	2	1	1	
S7	CPU	1	1	1	26	58	67	57	37	35	31	31	32	35	39	44	49	39	29	20	12	10	8	14	16	16	14	14	11	10	
	Mem	1	1	1	7	15	16	13	9	9	9	9	9	10	10	13	17	13	9	6	4	4	3	5	6	6	6	5	5	3	3

(b) Workload distribution #1

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s4	s6	n/a	n/a	s4	s5	s3	s5	s5	s5	s5	s7
s3	s5	n/a	n/a	s3	s4	s6	s5	s4	s4	s2	s1
s2	s3	n/a	n/a	s2	s4	s2	s4	s3	s6	s3	s7
s1	s1	n/a	n/a	s1	s2	s1	s3	s1	s2	s1	s2

(c) Workload distribution #2

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s4	s6	n/a	n/a	s3	s7	s3	s4	s3	s5	s6	s7
s6	s1	n/a	n/a	s1	s4	s2	s1	s2	s6	s2	s5
s7	s6	n/a	n/a	s3	s6	s1	s5	s1	s5	s3	s2
s4	s7	n/a	n/a	s2	s5	s5	s3	s4	s2	s1	s4

(d) Workload distribution #3

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s2	s5	n/a	n/a	s7	s3	s6	s3	s1	s6	s2	s7
s1	s4	n/a	n/a	s3	s1	s1	s3	s5	s6	s4	s5
s6	s2	n/a	n/a	s1	s4	s4	s7	s3	s4	s7	s6
s4	s2	n/a	n/a	s6	s5	s7	s2	s5	s3	s5	s2

Table A.8: Google cluster data benchmark scenario #4

(a) Workload pattern

Workload	Epoch (1 epoch = 10 sec)																														
	00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	
S1	CPU	33	39	59	54	47	66	74	83	83	82	69	51	47	45	34	32	31	33	36	34	33	23	23	45	49	58	58	57	47	43
	Mem	18	19	24	23	21	21	21	30	30	27	26	24	25	25	20	20	17	18	19	19	18	17	17	22	23	25	26	27	23	22
S2	CPU	26	27	28	27	25	27	28	26	26	27	28	29	26	25	28	28	27	28	31	31	33	36	31	29	29	31	31	31	22	20
	Mem	7	12	6	6	6	6	7	6	6	7	7	7	7	7	7	7	7	8	9	9	9	9	9	8	8	8	8	8	7	6
S3	CPU	70	67	56	49	40	52	57	67	63	45	58	71	74	76	64	64	79	76	71	56	45	59	62	71	65	54	61	68	64	63
	Mem	25	23	18	16	14	18	20	20	19	12	16	20	21	22	18	18	19	18	17	17	16	21	22	24	22	17	21	24	24	24
S4	CPU	26	27	28	27	25	27	28	26	26	27	28	29	26	25	28	28	27	28	31	31	33	36	31	29	29	31	31	31	22	20
	Mem	7	12	6	6	6	6	7	6	6	7	7	7	7	7	7	7	7	8	9	9	9	9	9	8	8	8	8	8	7	6
S5	CPU	56	55	54	48	40	47	49	44	44	45	54	53	39	41	38	37	37	50	48	59	59	49	29	37	32	22	21	41	27	22
	Mem	11	11	11	10	10	8	7	5	5	6	8	10	8	8	9	9	8	9	9	10	10	9	7	8	7	4	5	8	8	8
S6	CPU	48	47	57	48	37	40	41	34	33	32	28	23	25	26	32	33	25	27	31	28	26	30	31	38	42	53	48	43	37	28
	Mem	18	19	22	20	17	20	21	15	15	19	17	15	13	12	15	20	15	13	14	15	16	17	18	23	22	20	19	19	16	15
S7	CPU	19	21	26	25	23	18	16	23	22	18	17	14	16	17	15	15	15	15	16	16	15	19	20	17	17	18	20	22	18	16
	Mem	15	15	15	15	15	16	16	20	20	17	17	16	17	18	18	18	14	15	17	18	18	20	21	17	17	18	18	17	14	13

(b) Workload distribution #1

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s4	s6	n/a	n/a	s4	s6	s3	s5	s5	s5	s5	s7
s3	s5	n/a	n/a	s3	s6	s6	s5	s4	s4	s2	s1
s2	s3	n/a	n/a	s2	s4	s2	s4	s3	s6	s3	s7
s1	s1	n/a	n/a	s7	s2	s1	s3	s1	s2	s1	s2

(c) Workload distribution #2

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s6	s5	n/a	n/a	s7	s2	s2	s1	s6	s7	s4	s2
s6	s7	n/a	n/a	s7	s3	s4	s7	s4	s4	s2	s5
s2	s3	n/a	n/a	s3	s3	s1	s3	s6	s1	s6	s6
s5	s3	n/a	n/a	s1	s2	s4	s7	s1	s5	s5	s1

(d) Workload distribution #3

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s3	s1	n/a	n/a	s3	s6	s6	s5	s3	s2	s6	s5
s2	s6	n/a	n/a	s4	s5	s7	s7	s7	s1	s1	s6
s3	s1	n/a	n/a	s4	s5	s6	s7	s4	s3	s2	s7
s5	s4	n/a	n/a	s1	s2	s4	s5	s4	s2	s1	s2

Table A.9: Google cluster data benchmark scenario #5

(a) Workload patterns

Workload		Epoch (1 epoch = 10 sec)																													
		00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29
S1	CPU	46	44	35	39	43	44	44	44	45	53	48	40	41	42	40	40	42	42	43	43	45	45	46	57	57	57	56	56	40	35
	Mem	31	29	23	25	26	28	29	29	30	39	32	22	25	27	25	25	26	27	29	30	33	35	35	41	41	40	41	41	28	25
S2	CPU	59	57	52	53	53	55	56	60	61	66	70	75	66	61	50	49	47	47	48	50	51	69	73	70	69	68	68	68	60	32
	Mem	21	21	21	20	19	20	21	23	23	25	26	26	26	25	19	18	18	18	18	19	20	25	27	26	25	24	25	25	25	21
S3	CPU	27	26	22	21	20	18	17	20	21	22	19	19	16	14	20	21	21	21	22	23	23	26	27	28	28	27	22	17	13	12
	Mem	11	11	11	11	6	7	8	8	8	8	7	6	5	5	5	5	5	4	4	4	4	5	6	7	8	8	7	7	5	4
S4	CPU	42	41	37	37	37	37	37	26	28	40	38	40	42	43	35	37	35	36	39	42	44	42	42	40	39	38	34	30	29	36
	Mem	22	22	20	20	19	19	11	12	19	20	22	23	24	20	19	16	18	22	23	24	22	21	22	22	20	19	18	18	21	
S5	CPU	20	20	19	17	38	53	60	45	44	44	44	50	56	57	44	44	44	39	18	15	18	14	16	15	16	12	12	31	15	10
	Mem	7	7	8	10	9	11	13	8	8	8	8	11	11	11	8	8	9	9	8	6	6	5	5	8	8	6	7	10	7	7
S6	CPU	43	40	38	39	39	39	40	47	49	53	53	52	51	50	13	13	46	47	49	50	51	54	53	52	51	51	51	52	53	53
	Mem	8	7	6	6	6	7	7	7	9	9	9	9	8	5	5	10	10	10	10	10	10	10	10	7	7	7	7	7	7	
S7	CPU	73	69	54	52	49	47	46	40	42	54	47	37	27	21	19	18	18	18	19	20	20	24	25	28	27	25	24	22	21	20
	Mem	27	25	19	18	17	18	18	15	16	19	17	13	11	9	8	8	7	7	8	9	9	12	13	18	16	12	12	12	10	9

(b) Workload distribution #1

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s2	s5	n/a	n/a	s4	s6	s2	s7	s1	s7	s7	s2
s1	s3	n/a	n/a	s2	s6	s7	s5	s5	s1	s4	s4
s4	s3	n/a	n/a	s5	s1	s4	s1	s4	s2	s5	s3
s1	s2	n/a	n/a	s6	s2	s4	s3	s7	s3	s6	s3

(c) Workload distribution #2

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s7	s4	n/a	n/a	s3	s5	s2	s6	s1	s1	s3	s7
s2	s4	n/a	n/a	s3	s6	s5	s7	s4	s5	s7	s2
s7	s2	n/a	n/a	s7	s5	s1	s6	s4	s1	s6	s5
s6	s5	n/a	n/a	s3	s3	s6	s2	s3	s1	s4	s1

(d) Workload distribution #3

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s3	s2	n/a	n/a	s5	s3	s6	s7	s4	s5	s5	s6
s1	s3	n/a	n/a	s4	s5	s6	s2	s3	s7	s6	s7
s7	s4	n/a	n/a	s2	s1	s1	s2	s3	s1	s7	s2
s5	s4	n/a	n/a	s4	s7	s1	s6	s6	s4	s5	s1

Table A.10: Google cluster data benchmark scenario #6

(a) Workload patterns

Workload	Epoch (1 epoch = 10 sec)																														
	00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	
S1	CPU	24	24	20	19	20	37	44	43	44	46	47	48	46	46	19	19	41	32	17	17	18	40	45	47	46	45	46	48	24	17
	Mem	10	10	8	7	7	6	6	8	8	10	11	13	11	9	7	6	5	6	6	8	9	8	8	10	10	10	9	9	8	7
S2	CPU	22	22	30	30	23	12	12	15	17	22	22	22	24	26	34	35	34	31	28	25	21	15	14	11	13	19	19	15	19	
	Mem	7	7	8	7	6	5	5	5	5	6	6	7	7	7	8	8	8	8	7	6	5	5	6	5	6	9	10	11	7	6
S3	CPU	27	31	41	32	22	32	42	30	29	41	44	48	35	27	30	29	26	24	21	23	24	23	23	20	20	19	16	14	23	26
	Mem	24	22	12	10	7	11	15	9	9	16	29	15	13	11	16	16	11	11	11	11	11	9	9	10	10	9	9	9	24	27
S4	CPU	56	53	42	47	52	60	63	50	51	58	59	60	48	42	43	42	40	37	30	37	43	44	44	34	43	64	60	56	66	69
	Mem	30	28	23	24	26	38	43	27	26	25	26	29	27	26	25	24	18	18	17	23	27	30	30	30	29	26	27	28	29	30
S5	CPU	25	21	10	10	11	14	15	17	16	11	15	21	14	11	9	12	13	16	21	10	10	28	25	12	17	29	29	28	16	13
	Mem	11	10	7	11	15	11	9	5	5	6	7	8	7	6	6	7	5	7	10	7	7	11	12	7	8	10	10	14	8	6
S6	CPU	93	91	84	70	54	45	42	43	51	88	88	86	86	86	90	89	77	63	38	68	90	84	72	57	58	60	63	65	65	65
	Mem	23	20	10	9	9	6	5	7	9	16	16	15	19	21	14	13	10	9	6	12	16	16	15	7	10	15	16	17	18	19
S7	CPU	19	20	19	19	13	16	17	20	20	20	19	20	18	17	18	20	27	27	27	29	31	32	32	31	27	19	20	21	20	
	Mem	8	8	7	6	4	6	6	5	4	5	8	8	8	5	5	8	8	9	9	9	8	8	7	8	8	7	6	7	7	
S8	CPU	30	31	33	31	28	27	28	26	28	38	34	30	31	32	28	28	26	28	30	36	37	32	31	30	30	30	32	35	31	29
	Mem	20	21	24	23	20	19	18	18	18	18	19	20	20	21	19	18	16	17	19	22	24	21	20	20	19	19	21	22	21	20
S9	CPU	26	26	26	26	39	33	26	26	24	23	24	20	19	18	19	19	20	18	15	20	22	27	28	25	25	24	25	26	22	21
	Mem	20	20	22	21	21	19	17	18	17	16	16	16	17	17	16	16	16	16	16	17	17	21	22	21	20	19	21	23	19	18
S10	CPU	94	93	92	94	95	95	96	96	95	95	95	95	95	96	96	95	95	94	92	92	93	90	90	92	91	90	90	90	92	92
	Mem	53	54	58	54	49	45	44	48	49	50	51	52	54	55	46	46	49	52	57	56	57	47	46	48	47	45	51	56	51	49

(b) Distribution #1

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s7	s6	n/a	n/a	s1	s6	s4	s5	s3	s10	s5	s9
s8	s4	n/a	n/a	s2	s1	s4	s2	s3	s10	s6	s9
s8	s3	n/a	n/a	s5	s2	s1	s3	s7	s8	s9	s8
s7	s10	n/a	n/a	s4	s5	s1	s2	s6	s7	s10	s9

(c) Distribution #2

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s5	s2	n/a	n/a	s1	s4	s2	s4	s8	s10	s9	s5
s10	s2	n/a	n/a	s2	s7	s6	s6	s9	s1	s5	s4
s10	s1	n/a	n/a	s3	s8	s6	s10	s8	s7	s7	s3
s6	s5	n/a	n/a	s3	s9	s8	s7	s9	s1	s3	s4

(d) Distribution #3

vDom0		vDom1		vDom3		vDom4		vDom5		vDom7	
s6	s6	n/a	n/a	s7	s7	s5	s9	s3	s5	s6	s1
s10	s10	n/a	n/a	s5	s3	s4	s8	s9	s2	s5	s7
s9	s1	n/a	n/a	s7	s8	s4	s9	s1	s8	s10	s8
s4	s4	n/a	n/a	s2	s3	s3	s1	s2	s10	s6	s2

Bibliography

- [1] Luiz André Barroso, Kourosh Gharachorloo, Robert McNamara, Andreas Nowatzky, Shaz Qadeer, Barton Sano, Scott Smith, Robert Stets, and Ben Verghese. Piranha: A scalable architecture based on single-chip multiprocessing. In *Proceedings of the 27th Annual International Symposium on Computer Architecture*, ISCA '00, pages 282–293, New York, NY, USA, 2000. ACM.
- [2] Brent Bohnenstiehl, Aaron Stillmaker, Jon J. Pimentel, Timothy Andreas, Bin Liu, Anh T. Tran, Emmanuel Adeagbo, and Bevan M. Baas. KiloCore: A 32-nm 1001-Processor Computational Array. *IEEE Journal of Solid-State Circuits*, PP(99):1–12, 2017.
- [3] Shekhar Borkar. Thousand core chips: A technology perspective. In *Proceedings of the 44th Annual Design Automation Conference*, DAC '07, pages 746–749, New York, NY, USA, 2007. ACM.
- [4] Shekhar Borkar and Andrew A. Chien. The future of microprocessors. *Communications of the ACM*, 54(5):67–77, May 2011.
- [5] Thomas D. Burd and Robert W. Brodersen. Energy efficient cmos microprocessor design. In *Proceedings of the Twenty-Eighth Hawaii International Conference on System Sciences*, volume 1, pages 288–297, Jan 1995.

- [6] Qiong Cai, José González, Grigorios Magklis, Pedro Chaparro, and Antonio González. Thread shuffling: Combining dvfs and thread migration to reduce energy consumptions for multi-core systems. In *Proceedings of the 17th IEEE/ACM International Symposium on Low-power Electronics and Design, ISLPED '11*, pages 379–384, Piscataway, NJ, USA, 2011. IEEE Press.
- [7] Saurabh Dighe, Sriram R. Vangal, Paolo Aseron, Shasi Kumar, Tiju Jacob, Keith A. Bowman, Jason Howard, James Tschanz, Vasantha Erraguntla, Nitin Borkar, Vivek K. De, and Shekhar Borkar. Within-die variation-aware dynamic-voltage-frequency-scaling with optimal core allocation and thread hopping for the 80-core teraflops processor. *Solid-State Circuits, IEEE Journal of*, 46(1):184–193, Jan 2011.
- [8] Alejandro Duran and Michael Klemm. The intel many integrated core architecture. In *High Performance Computing and Simulation (HPCS), 2012 International Conference on*, pages 365–366, July 2012.
- [9] EZchip. TILE-MX Multicore Processor. <http://www.tilera.com/products/?ezchip=585&spage=686>. Online, accessed March 2015.
- [10] Xing Fu and Xiaouri Wang. Utilization-controlled task consolidation for power optimization in multi-core real-time systems. In *2011 IEEE 17th International Conference on Embedded and Real-Time Computing Systems and Applications*, volume 1, pages 73–82, Aug 2011.
- [11] Mohammad Ghasemazar, Hadi Goudarzi, and Massoud Pedram. Robust optimization of a chip multiprocessor’s performance under power and thermal constraints. In *2012 IEEE 30th International Conference on Computer Design (ICCD)*, pages 108–114, Sept 2012.

- [12] Soraya Ghiasi. *Aide De Camp: Asymmetric Multi-core Design for Dynamic Thermal Management*. PhD thesis, University of Colorado at Boulder Boulder, CO, USA, Boulder, CO, USA, 2004. AAI3136618.
- [13] Vinay Hanumaiah and Sarma Vrudhula. Energy-efficient operation of multicore processors by dvfs, task migration, and active cooling. *IEEE Transactions on Computers*, 63(2):349–360, Feb 2014.
- [14] Sebastian Herbert and Diana Marculescu. Analysis of dynamic voltage/frequency scaling in chip-multiprocessors. In *2007 ACM/IEEE International Symposium on Low Power Electronics and Design (ISLPED)*, pages 38–43, Aug 2007.
- [15] Jason Howard, Saurabh Dighe, Yatin Hoskote, Sriram Vangal, David Finan, Gregory Ruhl, David Jenkins, Howard Wilson, Nitin Borkar, Gerhard Schrom, Fabrice Paillet, Shailendra Jain, Tiju Jacob, Satish Yada, Sraven Marella, Praveen Salihundam, Vasantha Erraguntla, Michael Konow, Michael Riepen, Guido Droege, Joerg Lindemann, Matthias Gries, Thomas Apel, Kersten Henriss, Tor Lund-Larsen, Sebastian Steibl, Shekhar Borkar, Vivek De, Rob Van der Wijngaart, and Timothy Mattson. A 48-core ia-32 message-passing processor with dvfs in 45nm cmos. In *Solid-State Circuits Conference Digest of Technical Papers (ISSCC), 2010 IEEE International*, pages 108–109, Feb 2010.
- [16] Intel. Intel turbo boost technology 2.0. <http://www.intel.com/content/www/us/en/architecture-and-technology/turbo-boost/turbo-boost-technology.html>, 2013. Online, accessed March 2015.
- [17] Nikolas Ioannou, Michael Kauschke, Matthias Gries, and Marcelo Cintra. Phase-based application-driven hierarchical power management on the single-chip cloud computer. In *Proceedings of the 2011 International Conference on Parallel*

- Architectures and Compilation Techniques*, PACT '11, pages 131–142, Washington, DC, USA, 2011. IEEE Computer Society.
- [18] Canturk Isci, Alper Buyuktosunoglu, Chen-Yong Cher, Pradip Bose, and Margaret Martonosi. An analysis of efficient multi-core global power management policies: Maximizing performance for a given power budget. In *Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO 39, pages 347–358, Washington, DC, USA, 2006. IEEE Computer Society.
- [19] Vaibhav Jain. Fast Process Migration on Intel SCC using Lookup Tables (LUTs). Technical Report Masters Thesis, Arizona State University, May 2013.
- [20] Sudhanshu Shekhar Jha, Wim Heirman, Ayose Falcón, Jordi Tubella, Antonio González, and Lieven Eeckhout. Shared resource aware scheduling on power-constrained tiled many-core processors. *Journal of Parallel and Distributed Computing*, 100:30–41, 2017.
- [21] Chanseok Kang, Seungyul Lee, Yong-Jun Lee, Jaejin Lee, and Bernhard Egger. Scheduling for better energy efficiency on many-core chips. In *19th Workshop on Job Scheduling Strategies for Parallel Processing (JSSPP) In Conjunction with IPDPS 2015*, 19th JSSPP. Springer-Verlag, Hyderabad, India, May 2015.
- [22] Wonyoung Kim, Meeta S. Gupta, Gu-Yeon Wei, and D. Brooks. System level analysis of fast, per-core dvfs using on-chip switching regulators. In *IEEE 14th International Symposium on High Performance Computer Architecture (HPCA 2008)*, pages 123–134, Feb 2008.
- [23] Rakesh Kumar, Dean M. Tullsen, Parthasarathy Ranganathan, Norman P. Jouppi, and Keith I. Farkas. Single-isa heterogeneous multi-core architectures for multi-threaded workload performance. In *Proceedings of the 31st Annual International*

- Symposium on Computer Architecture*, ISCA '04, pages 64–, Washington, DC, USA, 2004. IEEE Computer Society.
- [24] Jian Li and José F. Martínez. Power-performance implications of thread-level parallelism on chip multiprocessors. In *IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS 2005)*, pages 124–134, March 2005.
- [25] Kai Ma, Xue Li, Ming Chen, and Xiaorui Wang. Scalable power control for many-core architectures running multi-threaded applications. In *Proceedings of the 38th Annual International Symposium on Computer Architecture*, ISCA '11, pages 449–460, New York, NY, USA, 2011. ACM.
- [26] David Meisner, Brian T. Gold, and Thomas F. Wenisch. Powernap: Eliminating server idle power. In *Proceedings of the 14th International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS XIV, pages 205–216, New York, NY, USA, 2009. ACM.
- [27] David Meisner and Thomas F. Wenisch. Dreamweaver: Architectural support for deep sleep. In *Proceedings of the Seventeenth International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS XVII, pages 313–324, New York, NY, USA, 2012. ACM.
- [28] Ke Meng, Russ Joseph, Robert P. Dick, and Li Shang. Multi-optimization power management for chip multiprocessors. In *Proceedings of the 17th International Conference on Parallel Architectures and Compilation Techniques*, PACT '08, pages 177–186, New York, NY, USA, 2008. ACM.
- [29] Trevor Mudge. Power: A first-class architectural design constraint. *Computer*, 34(4):52–58, 2001.

- [30] NVIDIA. GeForce GTX TITAN X. <http://www.geforce.com/hardware/desktop-gpus/geforce-gtx-titan-x>. Online, accessed March 2015.
- [31] Andreas Olofsson. Epiphany-V: A 1024 processor 64-bit RISC System-On-Chip. <https://arxiv.org/abs/1610.01832>, Oct 2016.
- [32] Jean-Marc Pierson and Henri Casanova. On the utility of dvfs for power-aware job placement in clusters. In *Proceedings of the 17th International Conference on Parallel Processing - Volume Part I, Euro-Par'11*, pages 255–266, Berlin, Heidelberg, 2011. Springer-Verlag.
- [33] Krishna K. Rangan, Gu-Yeon Wei, and David Brooks. Thread motion: Fine-grained power management for multi-core systems. In *Proceedings of the 36th Annual International Symposium on Computer Architecture, ISCA '09*, pages 302–313, New York, NY, USA, 2009. ACM.
- [34] Efraim Rotem, Avi Mendelson, Ran Ginosar, and Uri Weiser. Multiple clock and voltage domains for chip multi processors. In *Proceedings of the 42Nd Annual IEEE/ACM International Symposium on Microarchitecture, MICRO 42*, pages 459–468, New York, NY, USA, 2009. ACM.
- [35] Niraj Tolia, Zhikui Wang, Manish Marwah, Cullen Bash, Parthasarathy Ranganathan, and Xiaoyun Zhu. Delivering energy proportionality with non energy-proportional systems: Optimizing the ensemble. In *Proceedings of the 2008 Conference on Power Aware Computing and Systems, HotPower'08*, pages 2–2, Berkeley, CA, USA, 2008. USENIX Association.
- [36] John Wilkes. More Google cluster data. Google research blog, November 2011. Posted at <http://googleresearch.blogspot.com/2011/11/more-google-cluster-data.html>.

요약

대칭형 다중 처리 운영체제를 실행 시키는 캐쉬 일관성을 가지는 공유 메모리 아키텍처를 위한 전통적인 접근 방법은 전력관리가 가장 중요한 문제 중 하나로 존재하는 미래의 매니코어 시스템에는 적합하지 않다. 본 논문에서는 매니코어 시스템을 위한 계층적 전력관리 프레임워크를 소개한다. 제안한 프레임워크는 캐쉬 일관성을 가지는 공유 메모리가 필요 없으며, 다수의 코어들이 전압/주파수를 공유하고 다중 전압/다중 주파수를 지원하는 아키텍처에서 사용 가능하다. 이 프레임워크는 NUMA-인지 계층적 전력관리 기술로 동적 전압 및 주파수 교환(DVFS)과 워크로드 마이그레이션을 사용한다. 여기서 워크로드 마이그레이션 계획을 위해 사용된 탐욕 알고리즘은 서로 상충하는 비슷한 작업량의 패턴을 가진 작업을 같은 전압 영역으로 모으는 목표와 작업을 데이터가 있는 위치와 가까운 곳으로 이동하는 목표를 고려한다. 제안된 프레임워크는 소프트웨어로 구현되어 캐쉬 일관성이 없는 48 코어의 칩 레벨 멀티프로세서 하드웨어에서 평가되었다. 본 논문의 프레임워크를 데이터 센터 작업 패턴으로 광범위에 걸친 실험을 수행한 결과 최첨단의 DVFS 기술과 DVFS와 NUMA-비인지 워크로드 마이그레이션을 같이 사용한 전력관리 기술에 비해 상대적으로 각각 30%와 5%의 전력소모당 처리 작업량 향상을 큰 성능손실 없이 이루었다.

주요어: 매니코어 아키텍처, 불균일 기억 장치 접근, 스케줄링, 동적 전압 및 주파수 변경, 에너지 효율

학번: 2015-22902